**Faculty of Science and Engineering**
**Department of Chemical Engineering**

# Development of an Intelligent Dynamic Modelling System for the Diagnosis of Wastewater Treatment Processes

**Muhammad Imran Khalid**

**This thesis is presented for the Degree of**
**Master of Philosophy (Chemical Engineering)**
**of**
**Curtin University of Technology**

**March 2010**

**Declaration**


To the best of my knowledge and belief this thesis contains no material previously published by any other person except where due acknowledgment has been made.


This thesis contains no material which has been accepted for the award of any other degree or diploma in any university.


Signature: .. *ImranKhalid.* .........................


Date:         11/10/2010

# Abstract

In the 21$^{st}$ Century, water is already a limited and valuable resource, in particular the limited availability of fresh water sources. The projected increase in global population from 6 billion people in 2010 to 9 billion in 2050 will only increase the need for additional water sources to be identified and used. This situation is common in many countries and is frequently exacerbated by drought conditions. Water management planning requires both the efficient use of water sources and, increasingly, the re-use of domestic and industrial wastewaters. A large body of published research spanning several decades is available, and this research study looks specifically at ways of improving the operation of wastewater treatment processes.

Process fault diagnosis is a major challenge for the chemical and process industries, and is also important for wastewater treatment processes. Significant economic and environmental losses can be attributed to inappropriate Abnormal Event Management (AEM) in a chemical/processing operation, and this has been the focus of many researchers. Many researchers are now focusing on the application of several fault diagnosis techniques simultaneously in order to improve and overcome the limitations experienced by the individual techniques. This approach requires resolution of the conflicts ascribed to the individual methods, and incurs additional costs and resources when employing more than one technique. The research study presented in this thesis details a new method of using the available techniques. The proposal is to use different techniques in different roles within the diagnostic approach based upon their inherent individual strengths. The techniques that are excellent for the detection of a fault should be employed in the fault detection, and those best applied to diagnosis are used in the diagnosis section of a diagnostic system.

Two different techniques are used here, namely a mathematical model and data mining are used for detection and diagnosis respectively. A mathematical model is used which is based upon the principal of analytical redundancy in order to establish the presence of a fault in a process (the fault detection), and data mining is

used to produce production rules derived from the historical data for the diagnosis. A dataset from an industrial wastewater treatment facility is used in this study.

A diagnostic algorithm has been developed that employs the techniques identified above. An application in Java was constructed which allows the algorithm to be applied, eventually producing an intelligent modelling agent. Thus the focus of this research work was to develop an intelligent dynamic modelling system (using components such as mathematical model, data mining, diagnostic algorithm, and the dataset) for simulation of, and diagnosis of faults in, a wastewater treatment process where different techniques will be assigned different roles in the diagnostic system.

Results presented in Chapter 5 (section 5.5) show that the application of this combined technique yields better results for detection and diagnosis of faults in a process. Furthermore, the dynamic update of the set value for any process variable (presented in Chapter 5, section 5.2.1) makes possible the detection of any process disturbance for the algorithm, thereby mitigating the issue of false alarms. The successful embedding of both a detection and a diagnostic technique in a single algorithm is a key achievement of this work, thus reducing the time taken to detect and diagnose a fault. In addition, the implementation of the algorithm in the purpose-built software platform proved its practical application and potential to be used in the chemical and processing industries.

# Table of contents

# List of Figures

# List of Tables

# Acknowledgements

I would like to thank my academic supervisor Dr. Nicoleta Maynard and co-supervisor Dr. Martyn Ray of Curtin University of Technology for their guidance and help in accomplishing this research task.

Muhammad Imran Khalid

# Chapter 1

# Research Overview

## 1.1 Background

Water sources, and their use, cannot be considered as an unlimited resource. The treatment and reuse of wastewater is now accepted practice worldwide and within Australia. Wastewater treatment is a very active area of research yielding many publications, however there are many challenges and significant aspects requiring improvement within the wastewater industries. In broad terms, there are two types of wastewater, namely municipal and industrial. Although different operations produce different wastewaters, the schematic process applied to water treatment is often very similar.

Wastewater treatment facilities are mostly government or state operated and, unlike other chemical and process industries, there is no revenue directly generated from the end product (the reclaimed water), except for reducing the use and cost of freshwater. Reviewing the treatment process highlights the need to increase the efficiency and effectiveness of operation by using a robust fault detection and diagnostic system, optimization of utilities, work force, assets, etc. Models have been developed to assist the engineers and operators understanding of the process, and hence apply effective fault detection and diagnostics and optimization techniques, e.g. Activated Sludge Model (ASM) series (Gujer et al., 1999) and Anaerobic Digestion Model (Copp et al., 2005).

Process fault diagnosis is one of the major challenges facing the chemical and process industries and significant economic and environmental losses can be attributed to poor Abnormal Event Management (AEM). There are two key approaches to fault diagnosis, namely process model based diagnosis and data driven diagnosis. Process model based approaches assume that a fault will cause changes to certain physical parameters, which in turn will lead to changes in some of the model parameters or states. It is then possible to detect and diagnose these faults by monitoring the estimated model parameters or states. This technique uses a

diagnostic-driven mathematical model of a process (Simani et al., 2003) and generally exhibits good accuracy, since the process models are developed from the underlying fundamental principles. However, comprehensive theoretical models of complex processes are extremely challenging to develop due to their often inherent non-linear nature. Alternatively the use of computer aided systems means that data driven models are relatively easy to develop, and analysis of the data enables fault identification and diagnosis (Li, 2003).

A comprehensive review of the techniques available for fault detection and diagnostics, the strengths and challenges of each technique, and the key attribute required by an ideal diagnostic system, is given by Venkatasubramanian et al. (2003b, 2003c, 2003a). They determined that no single method has all the ideal desirable features for a diagnostic system, but some of these methods can complement one another resulting in an improved overall diagnostic system. The use of a hybrid system is a promising future research direction leading to developments in diagnostic systems. Integrating these complementary features is one way to develop hybrid systems that may overcome the limitations of an individual strategy. The drawbacks of single-method-based diagnostic systems are serious enough to limit their applications to small case studies and render them unsuitable for large-scale industrial situations. This makes the design and development of hybrid systems important (Venkatasubramanian et al., 2003c).

Few papers have been published on wastewater treatment in relation to process monitoring and control and fault diagnostics (Lee et al. 2006). Since 1998, 29 journal papers were published concerning the application of hybrid fault diagnostic systems. There are only two publications for fault detection and mathematical models (in 2003 and 2008).

## 1.2   Thesis Objectives

The main objective of this research work was to develop an intelligent modelling system to be used for dynamic simulation, and the diagnosis of wastewater processes using a hybrid of model based (mathematical model) and data driven (inductive data mining) to address the limitations exhibited by the individual approaches. Particular

emphasis was placed upon the development of a hybrid diagnostic system for its reliability, flexibility, robustness, relative low-cost of development, and ease of utilization in wastewater treatment processes. The overall objective of this thesis was to develop an effective fault detection and diagnostic platform which employs a hybrid system of a mathematical model and a data driven technique to overcome the limitations of the individual fault detection and diagnostic methodologies. This was achieved by the following approach:

- Develop a software platform using Java for dynamic simulation and fault detection on the wastewater treatment process; this includes the implementation of the model libraries of full-scale wastewater treatment processes.

- Select the most suitable data driven technique for fault detection and then implement on the wastewater treatment data.

- Implement the hybrid algorithm in the Java-based platform; this includes the development of the hybrid algorithm.

- Validate the results against wastewater treatment process data.

## 1.3 Thesis Outline

The organization of the thesis is as follows:

Chapter 2 discusses in detail the process diagnostics including an overview, the techniques available, and a critical literature review of current developments identifying the limitation of process diagnostics.

Chapter 3 focuses on the wastewater process. It describes briefly the types of operating units in a wastewater process before presenting a mathematical model to be used for the fault detection. It also includes the rational which supports implementation of the derived model over other available process models.

Chapter 4 provides a comparative study of available data driven techniques leading to the selection of the preferred data driven technique (inductive data mining), the wastewater data is then used to generate results. Discussion on the output format of results and manipulation of the format for use in the diagnostic system is also included.

Chapter 5 contains the development of the software platform in Java and discusses its flexibility, ease of use, and customization. The wastewater model is imported and initial simulation results are discussed. It includes the development of the algorithm for a hybrid approach of the model plus data driven technique. The results and validation of the wastewater data are presented together with the accuracy and the computational efficiency.

Chapter 6 presents the research conclusions, and identifies areas of research recommended for further study.

# Chapter 2

# Literature Review

## 2.1   Introduction

This chapter provides a general overview of process diagnosis and wastewater treatment, in addition to a critical review of selected topics. This chapter includes details of hybrid diagnosis systems and the software platform used in wastewater treatment and process diagnosis.

## 2.2   Process Diagnosis

Diagnosis is defined as "the process of identifying the nature and causes of certain phenomenon". It has been observed that the word "diagnosis" is always attributed to abnormal phenomenon rather than the normal situation. Similarly, process diagnosis is the study used to identify the nature and causes of an abnormal phenomenon occurring in a process. The abnormal phenomenon in the processing and manufacturing industries is termed a "fault". Himmelblau (1978) described a fault as a departure from an acceptable range of an observed variable or a calculated parameter associated with a process. Hence, fault is a process abnormality or symptom, such as high temperature in a reactor or low product quality, etc. The underlying causes of this abnormality are called basic events, or root causes. The basic event is also referred to as a malfunction or a failure.

In broad terms, a fault is generally related to one (or more) of three main classes of failures or malfunctions. These are:

*Gross parameter changes in a model*: Parameter changes arise when there is a disturbance entering the process from the environment. An example is the change in the heat transfer coefficient due to fouling of a heat exchanger.

*Structural changes*: Structural changes refer to changes in the model itself. They occur due to the equipment hardware failures. An example is a controller failure

which would imply that the manipulated variable is no longer functionally dependent on the controlled variable.

*Malfunctioning sensors and actuators*: Gross errors usually occur with actuators and sensors. These could be due to a fixed failure, a constant bias (positive or negative) or an out-of-range failure.

Process diagnosis is the study undertaken to identify the nature and causes of a fault in a process. This process usually includes the following four steps (Chiang et al., 2001):

*Fault Detection* - determining if a fault has occurred by observing the values of process variables or by user-defined fault indices.

*Fault Identification* - identifying the variable most relevant to the fault. The task of this step is to focus the attention of the plant operator on the particular subsystem which is most pertinent to the fault.

*Fault Diagnosis* - isolating the cause, type, location, and time of the fault.

*Process Recovery* - removing the fault and bringing the process back to normal conditions, or taking optimal steps to minimize loss of production.

Process diagnosis was traditionally performed by the process operators but this task has became more difficult due to the variety, uncertainty and time delay of malfunctions which can be very complex in their nature,. In the last decade, computer-aided systems have been investigated, implemented and proved to be a successful tool for fault diagnosis. A fault diagnostic system has two main components:

(i)    type of knowledge used, and

(ii)   type of diagnostic search strategy.

Diagnostic search strategy is usually strongly dependent on the knowledge representation scheme, which in turn is largely influenced by the type of prior knowledge available. Hence, the type of prior knowledge used is the most important distinguishing feature in diagnostic systems. The prior domain knowledge may be developed from a fundamental understanding of the process using a first principles approach. Such knowledge is referred to as deep, causal, or model-based knowledge (Simani et al., 2003). However it may also be obtained from past experience of the

process and is then referred to as shallow, compiled, evidential, or process-history-based knowledge.

The process diagnostic techniques can be broadly classified into two categories, i.e. process-model-based methods and process-history-based (data-driven) methods as discussed below.

## 2.2.1 Process Model Based Methods

These methods use qualitative knowledge and quantitative models extracted from an understanding of the process principles. The models present the interacting relationships between the process variables based upon the assumption that a fault will cause changes to certain physical parameters which in turn lead to changes in some of the model parameters or states. It is then possible to detect and diagnose these faults by monitoring the estimated model parameters or states.

The process model based methods can be further sub-divided into qualitative causal models and quantitative methods (Dash and Venkatasubramanian, 2000, Venkatasubramanian et al., 2003a, 2003b).

### *Quantitative model based methods*

Relying on an explicit model, all model based fault detection and isolation (FDI) methods require two steps. The first step is to generate inconsistencies between the actual and expected behavior, these are called residuals and reflect the potential faults of the system. The second step chooses a decision rule for diagnosis.

Some form of redundancy is required to check for the inconsistencies. There are two types of redundancies, namely hardware and analytical redundancy. Hardware redundancy uses redundant sensors and has been utilized in the control of safety-critical systems such as aircraft, space vehicles and nuclear power plants, but its applicability is limited due to the additional costs and space required. Alternatively, analytical redundancy is achieved from the functional dependence among the process variables and is usually provided by a set of algebraic or differential relationships among the states, the inputs, and the outputs of the system (Lou et al., 1986, Michle, 1988).

The essential element of analytical redundancy is to check the actual system behavior against the system model for consistency. Any inconsistency, expressed as residuals, can be used for detection and isolation purposes. The residuals should be close to zero when no fault occurs, but show a significant value when the underlying system changes. To generate the diagnostic residuals requires an explicit mathematical model of the system. The model may be obtained either analytically using first principles or empirically as a black-box model.

Most of the work on model-based diagnostic systems reported to date was mainly in the aerospace, mechanical or electrical engineering literature. There has not been much published research on its application for fault diagnosis in chemical process systems. There are some serious limitations that apply for its application in the chemical industries. One issue is the lack of availability, and the complexity, of models for chemical processes and their inherent non-linear nature. In addition to the modelling challenges, the model based qualitative methods do not include an explanation and descriptive facility. Furthermore, an estimation of classification errors cannot be provided when using these methods. Another disadvantage with these methods is that if a fault is not specifically and appropriately modelled, then there is no guarantee that the residuals will be able to detect it.

### Qualitative model based methods

Qualitative models are usually developed based on the fundamental understanding of the physics and chemistry of the process. Various forms of qualitative models such as causal models and abstraction hierarchies have been developed. The strategy employed in qualitative models is causal-effect reasoning related to the system behavior. The most popular methods are fault-trees and signed digraphs (SDG). Fault trees (Lapp and Powers, 1977) use backward chaining until a primary event is found that presents a possible root cause for the observed process deviation from normal operation. SDG (Iri et al., 1979) is another representation of the causal information in which the process variables are represented as graph nodes and causal relations by directed arcs. Causal model-based methods mimic human reasoning and so generation of explanation is relatively straightforward making them more interactive. There were drawbacks of the early SDG methods mainly because their expressive capability is often limited (Hunag and Wang, 1999a). In an SDG, a node or a branch

8

can often only take three values. i.e., -, 0, and + representing for example low, normal and high values for a node. This over simplified expression could create ambiguous solutions.

The two main concerns with qualitative model based methods are ambiguities and spurious/inauthentic solutions. Considerable research has been done in relation to the reduction of spurious solutions while reasoning with qualitative models.

### 2.2.2 Process History Based Methods

In contrast to the model-based approaches where a priori knowledge (either quantitative or qualitative) about the process is needed, in process history based methods, only the availability of a large amount of historical process data is needed. There are different ways in which this data can be transformed and presented as a priori knowledge to a diagnostic system. This is known as feature extraction, and this extraction process can be either qualitative or quantitative in nature. Two of the major methods that extract qualitative history information are expert systems and trend modelling methods. Methods that extract quantitative information can be broadly classified as non-statistical or statistical methods. Neural networks are an important class of non-statistical classifiers. Principal component analysis (PCA)/partial least squares (PLS), data mining and statistical pattern classifiers form a major component of statistical feature extraction methods. The knowledge can be available as rules and formulations.

## 2.3   Hybrid Diagnosis Systems

Venkatasubramanian and co-workers (see Venkatasubramanian et al., 2003b, 2003c, 2003a) have provided a very comprehensive review of the methods available for process diagnosis. A set of desirable characteristics that a diagnostic system should possess were also identified and listed.  Different approaches were evaluated against a common set of requirements or standards. From their evaluation, it was revealed that no single method has all the desirable features stipulated for a diagnostic system. It was postulated that some of these methods can complement one another, resulting in better diagnostic systems. Integrating these complementary features is one way to develop hybrid methods that could overcome the limitations of individual solution strategies. Hence, hybrid approaches are attractive where different methods work in

conjunction to solve parts of the problem. Although all the methods possess limitations, in the sense that they are only as good as the quality of information provided, it was shown that some methods might be better adapted to the knowledge available than others. For example, fault explanation through a causal chain is best done through the use of digraphs, whereas fault isolation might be very difficult using digraphs due to the qualitative ambiguity and then analytical model-based methods might be superior. It is expected that hybrid methods will provide a general, powerful problem-solving platform.

## 2.4   Software Platform

One of the challenges in the implementation of a diagnostic system is the software architecture (Venkatasubramanian et al., 2003c). The diagnostic task can be performed either off-line or on-line. In an off-line diagnostic task, the process behavior is recorded in the form of data or graphical trends. This data is then analyzed off-line using a suitable diagnostic method. This type of diagnostic is often used as a preventive action to save the process from repeating the same malfunction. Due to the complexity of chemical processes and the significant losses attributed to the poor management of faults, the online diagnostic is increasingly employed in the chemical industries for corrective actions. The online diagnostic system helps operators and engineers to manage an abnormal event (fault) as soon as it is encountered in the process. Using its knowledge base, the online diagnostic system will search for the most likely culprit for a specific fault. This vital information, along with the prior knowledge that the operators and engineers have about the process, will be critical in the isolation and diagnosis of a fault, thus reducing the losses and downtime and increasing process and personnel safety associated with the management of an abnormal event (fault).

From the literature review (Venkatasubramanian et al., 2003c), a software platform intended to be used for online diagnosis should posses the following characteristics.

i)      Flexibility

Flexibility refers to the ability of software to accommodate different configurations of plant items and an ability to use different types of equipment for any unit

operation or process. In relation to wastewater treatment where different configurations exist in different facilities due to variations of influents within a region (state or country), the software should have the timely ability to accommodate for such changes in configuration.

ii)     Detection algorithm

The software should incorporate an algorithm that is capable of detecting any abnormality in the process. The algorithm will be highly dependent on the type of benchmarking information used to detect faults, i.e. if a mathematical model is used then the algorithm should be able to: (a) obtain values from the mathematical model; (b) obtain information from real processes; and (c) should have the threshold value for noise/bias in real data. After the comparison of values of a process variable from real operations and benchmarking, the algorithm will be able to detect a fault and draw the attention of operators to that variable.

iii)     Data/Information import

Based on the detection algorithm requirements, it is essential that the software should be able to obtain information from the control system of a process for use in fault detection. The integration of a diagnostic software platform with the control module in a plant is a major challenge that requires precision skills in software architecture, and the availability of resources (e.g. a distributed control system of a pilot scale plant) for trials to validate the effectiveness, efficiency and robustness to the software platform.

iv)     Diagnostic algorithm

After the detection of a fault, the most important aspect of the diagnostic system is to then diagnose the fault. This necessitates a diagnostic algorithm that will diagnose the fault using the knowledge-base used for diagnosis purposes.

v)     Results display

An essential ingredient that a diagnostic software platform should possess is the display of diagnostic results in an easily understood format. This may be production rules, fault trees, decision trees, or any other acceptable form of results.

vi)     Cost effectiveness

Although management of a fault in a process is essential, and is associated with significant economic benefits, it would be an optimal solution to have a cost effective diagnosis system without compromising the effectiveness of the system. Commercial software such as Gensym G2, which is widely used in the chemical industries for

better management of abnormal process conditions, costs between $100,000 and $1 million. Although it may be cost effective, this is a critical investment decision for small-scale and service industries.

## 2.5   Wastewater Treatment

Water recycling and reuse is now generally accepted and adopted worldwide as an essential water resource. The average composition of municipal wastewater is 99.94% $H_2O$ and 0.06% dissolved and suspended solids, which indicates the potential for the reuse of waste water.

Wastewater treatment facilities can process wastewater using either an aerobic, anaerobic or anoxic process. Each process has its own advantages and disadvantages. Wastewater treatment facilities are mostly government or state operated and, unlike other chemical and process industries, there is no direct revenue generated from the end product, i.e. the reclaimed water (although it can be used for irrigation or drinking purposes). Hence there is more attention given to reducing the running costs by optimization of utilities, workforce, assets, etc. Models have been developed to assist the engineers and operators understanding of the process, e.g. Activated Sludge Model (ASM) series (Gujer et al., 1999) and Anaerobic Digestion Model (Copp et al., 2005). Wastewater processes have been widely researched in recent years. Lee et al. (2006) addressed fault diagnosis of sensors in wastewater treatment processes but the methodology proposed is unable to identify the faulty senor which causes process transitions. Kim et al. (2002) calibrated ASM1 using genetic algorithms but the components of ASM1 were not calibrated in detail. Gernaey et al. (2004) used artificial intelligence and white-box based modelling and simulation for wastewater treatment processes, and demonstrated how different methodologies can complement and support the process knowledge included in white-box activated sludge models. Puteh et al. (1999) present a mathematical model of the aeration tank and the secondary settler of a wastewater treatment facility, proposing that the performance of wastewater treatment processes consisting of an incomplete mixing reactor described by tanks-in-series model is better than that of a completely mixed aeration tank. Rigger et al. (2006) proposed a model for the response time of the aeration systems concluding that if more calibrated applications of the aeration system model are available, it should be possible to develop a classification system for design

purposes. Wintgens et al. (2003) presented the modelling of a membrane bioreactor replacing the conventional aeration system in wastewater treatment processes, which is more efficient than the traditional aeration tank in terms of investment cost but the operating cost is higher. Van Hulle and Vanrolleghem (2004) showed that model-based optimization is an efficient and cost-effective way to ensure that an industrial wastewater treatment plant functions well, but a more holistic evaluation is required before the proposed methodology can be applied.

Much research on wastewater is related to the treatment of industrial wastewaters with very specific requirements. Acharya et al. (2009b,2009a) used activated carbon prepared from Tamarind wood with zinc chloride activation for the removal of lead (II) and chromium (VI) from industrial wastewater. Hunag et al. (2009) developed an integrated neural-fuzzy process controller to control aeration in an aerated submerged biofilm wastewater treatment process (ASBWTP) which saved 33% of the operating costs during the time when the controller was used. Pai (2008) employed grey models (GM) to predict the effluent quality of a wastewater and compared the results with the use of artificial neural networks (ANN). The results indicated that GM can be used for effluent prediction while using less data than required in ANN. The amount of research work reported on wastewater is significant, and wastewater is indeed one research area that has attracted major attention worldwide.

## 2.6   Concluding Remarks

Process diagnosis is a key research area receiving significant attention from both academia and industry. In recent years this has lead to significant process improvements in abnormal event management with proposals for numerous new techniques. With the development of these techniques, the focus of research should now be on the development of a hybrid methodology that can address the weaknesses of individual techniques in order to further enhance the effectiveness of process diagnostics. So far, this area has not been studied in depth.

By comparison with other chemical processes, the wastewater treatment process is quite small in terms of the processes involved and the physical size of a facility. Furthermore, as discussed in Section 2.5, a municipal wastewater facility does not

generate any direct revenue from its process, unless it is involved in the treatment of an industrial wastewater. Hence, the optimization of capital and operating costs must be a key focus. A purpose-built, cheap and effective software platform that can be used for process diagnosis in wastewater treatment, or any other small chemical industry, is one application for budget optimization. A significant contribution of this research study is the development of a flexible and low-cost software platform.

# Chapter 3

# Process Model-Based Methods for Wastewater Treatment Processes

## 3.1   Introduction

In this chapter, Sections 3.1 to 3.7 present the necessary background understanding of the typical wastewater treatment process. There is a significant amount of published literature on wastewater treatment processes and this large volume of information can pose a problem for someone new to this field needing an overview of the essential processes involved in water treatment. Therefore Sections 3.1-3.7 present a comprehensive, yet concise, review of wastewater and its treatment.

Process Model-Based (PMB) methods/techniques for Fault Detection and Diagnosis (FDD) are then discussed from Section 3.8. PMB methods are generally classified as quantitative or qualitative. An introduction and brief discussion of the quantitative and qualitative PMB techniques are presented, followed by explanation of how mathematical modelling will be used in this work for FDD. A mathematical model of a municipal wastewater treatment plant is presented in Section 3.10. Results obtained from the process simulations are reported and discussed in Section 3.11.

## 3.2   Wastewater

Wastewater is a general term for any water that has been used for either domestic or industrial purposes, and hence becomes contaminated by various waste materials. It can contain human excreta, food waste, and industrial toxins, along with many other pollutants in the form of dissolved and suspended material, and hence is unfit for human consumption and can damage aquatic systems. The composition of wastewater in terms of the water itself and the contained wastes varies widely

between countries, and within each country, but the average composition of a wastewater is around 99% water and 1% solid.

Water shortages and deterioration in quality are major challenges faced by many countries, especially those pursuing economic and social development. Due to population growth and the associated increased use of water for agriculture, industry and recreation, the human consumption of natural waters has steadily increased over several centuries thus making it an increasingly valuable resource. This situation had lead many researchers to think about whether there will be enough water to accommodate the needs of future generations (Kumar, 2004). Water management strategies and specific techniques are being adopted in many countries for optimum water utilization based on defined water policies. Wastewater treatment compliments water management in two ways: (1) it increases the total available water resource by converting the wastewater into useable water; and (2) when the treated wastewater is discharged into a receiving body (lakes or rivers) after treatment, it does not deteriorate the overall water quality by being free from pollutants that may affect the aquatic system.

The composition of wastewater indicates the great potential for the application of wastewater treatment processes. Wastewater treatment is being increasingly adopted worldwide for optimizing water management systems.

## 3.3   Wastewater Treatment

Wastewater treatment is defined as the processing of wastewater for the removal or reduction of all undesirable constituents. Three types of processes are involved in the treatment of wastewater:

a)   Mechanical,

b)   Biological, and

c)   Chemical.

Mechanical or physical treatment processes are used for the removal of large objects, heavy inorganic matter, oil, greases and particulate solids from the raw wastewater before it is treated using biological processes.

The biological process then converts the influent from mechanical treatment unit, to be almost free from dissolved and suspended solids (pollutants). This goal is achieved by the action of microorganisms that thrive on the pollutants in wastewater

for their nutrients. Oxygen is required in this treatment stage as it is essential for the survival of living organisms.

The main purpose of chemical treatment is to achieve the final required water quality (specification) before it can be sent to a receiving body. Chemical agents, such as chlorine gas, are mainly used for disinfection of water in order to kill any microorganisms remaining after the biological treatment stage.

In any type of wastewater treatment plant, there are essentially the following stages:

a)    Preliminary treatment,

b)    Primary treatment,

c)    Secondary treatment,

d)    Tertiary treatment and

e)    Advanced treatment systems.

The last stage, i.e. tertiary treatment, depends on the specified required quality for the effluent water and, unless a high quality of water is required, this stage is seldom used in treatment plants.

A typical preliminary treatment stage consists of screens, comminutors and grits. Only mechanical treatment is carried out in the preliminary stage. In the primary treatment stage, heavy solids and oils or greases are separated from the wastewater using gravity action in a sedimentation tank often called a primary settler/clarifier. Effluent from the primary settler/clarifier is treated in a bio-reactor where microorganisms eliminate, or significantly reduce, the amount of dissolved and suspended matter (pollutants). The effluent is then passed to a secondary settler where the clean water as supernatant is either released into the receiving body or is passed to tertiary treatment for further purification. Settled solids at the bottom of the settler are either recycled back to the reactor or disposed of after further treatment. The bio-reactor and the secondary settler comprise the secondary stage for wastewater treatment. An example of a typical wastewater treatment plant is given in the schematic of Figure 3.1.

**Figure 3.1** Structure of a wastewater treatment plant (Wang et al., 2004)

Wastewater has to pass through these three treatment stages before it can be released back into the environment. The following section provides further details of these stages and the equipment used in each stage.

## 3.4   Preliminary Treatment Stage

The preliminary stage consists of screens, comminutors and grits. A screen is a device with openings of uniform size that is used to retain solids present in the influent wastewater. Its principal role is to remove coarse material from the influent that could damage subsequent process equipment (Tchobanoglous et al., 2003). Screens are followed by Comminutors which are used to reduce the particle size of wastewater solids as large, stringy solids can easily plug pump impellers (Degremont., 1991, Forster, 2003). The last unit in the primary treatment is Grit chamber which is used to remove the heavy inorganic material present in wastewater, such as sand, eggshells, gravel and cinders, which have settling velocities and specific gravities substantially greater than the organic solids in the wastewater. (Tchobanoglous et al., 2003). Only mechanical treatment of wastewater is carried out with the aim to eliminate the large sized materials and heavy inorganic material from the wastewater before it can enter the primary settler for further purification.

## 3.5 Primary Treatment Stage

After the preliminary treatment of raw wastewater, the remaining solids are extracted by gravity in large sedimentation tanks in the primary stage of treatment. Sedimentation tanks further slow the influent flow of wastewater so that organic and inorganic suspended solids can settle to the bottom of the clarifiers. Floatable solids and grease are skimmed off by a rotating arm and deposited as a scum. A clarified supernatant liquid leaves from the top of the sedimentation tank, while concentrated sludge exits from the bottom of the sedimentation tank. The primary clarifiers remove about 60% of the Total Suspended Solids and about 30% of the Biochemical Oxygen Demand in the incoming wastewater. Two types of sedimentation tanks/clarifier are used in wastewater treatment plants, namely rectangular tanks and circular tanks (Tchobanoglous et al., 2003).

## 3.6 Secondary Treatment Stage

The secondary treatment stage of a wastewater treatment process is the most important as the removal of organic material (which represents the pollutants in wastewater) takes place in this stage. It comprises a biological unit that facilitates the growth of microorganisms requiring a sufficient supply of oxygen. The microorganisms use the organic material in the wastewater stream as food, thus reducing the pollutant contents in the wastewater. Treated wastewater from the biological unit is then passed to a secondary clarifier/tank where the solids settle to the tank bottom. The treated effluent stream is either released into the receiving body, or is subjected to tertiary treatment if a high quality of water purity is required. A brief overview of the principal types of reactors used in this stage is given below.

**a) Batch Reactors**

A batch reactor used for wastewater treatment will incorporate the following operational phases (Buchanan and Seabloom, 2004):

i. *Fill:* Raw wastewater that has been through primary treatment is added to the reactor. During this phase, aeration may or may not be supplied in order to provide alternating periods of high or low dissolved oxygen. This mode may occupy 25% of the total cycle time.

ii. *React:* Aeration is provided in an effort to obtain rapid biodegradation of organic compounds. This mode will typically require about 35% of the total cycle time.

iii. *Settle:* Aeration is shut off to allow the wastewater to become anoxic (for denitrification) and to allow for quiescent conditions that allow very effective liquid-solid separation. Clarification will usually take about 20% of the overall cycle time.

iv. *Draw:* Clarified raw water is removed as the supernatant liquid. The decanting is accomplished using adjustable or floating weirs. Periodically the excess biosolids must be removed. Decanting generally takes about 15% of the total cycle time.

An important requirement in the batch reactor process is that a tank is never completely emptied, and a portion of the settled solids are left to seed the next cycle. This allows the establishment of a population of organisms uniquely suited to treating the wastewater. By subjecting the organisms to periods of high and low oxygen levels, and to high and low food availability, the population of organisms becomes very efficient at treating wastewater. A typical hydraulic retention time for a batch reactor varies between 20 to 40 hours (Tchobanoglous et al., 2003).

**b) Complete Mix Reactors**

It is assumed that a complete mixing occurs instantaneously and uniformly throughout a complete mix reactor as fluid particles enters the reactor. Fluid particles leave the reactor in proportion to their statistical population. The actual time required for complete mixed conditions depend on the reactor geometry and the power input.

**c) Plug-Flow Reactors**

Fluid particles pass through the reactor with a little or no longitudinal mixing and exit from rector in the same sequence in which they enter. The particles retain their identity and remain in the reactor for a time equal to the theoretical detention time.

**d) Packed Bed Reactors**

A packed bed reactor is filled with some type of packing material, such as rock, slag, ceramic or plastic with plastic being the most commonly used. With respect to the flow, a packed bed reactor can be operated in either down or up flow mode. The

input to the reactor can be continuous or intermittent. The packing material can also be continuous or arranged in multiple stages with flow from one stage to another.

**e)  Fluidized Bed Reactors**

The fluidized reactor is similar to the packed bed reactor in many aspects, but the packing material is expanded by the upward movement of fluid through the bed.

## 3.6.1 Biological Treatment

Biological methods of wastewater treatment are based upon induced contact with microorganisms, which feed on the organic materials in the wastewater, thereby reducing the Biological Oxygen Demand (BOD) content of wastewater. BOD in wastewater is used as an indicator of pollutant level, where the greater the BOD, the greater the degree of pollution (Green-Ideas, 2009).

The basic principle behind biological treatment lies in the microorganisms consuming the suspended organic material present in the wastewater as their food source. The organic material is transformed into cellular mass by the metabolic process which is no longer suspended and hard to separate from the water, but can be precipitated by gravity at the bottom of a settling tank. Thus, the water exiting the biological system (biological treatment unit and clarifier) is much clearer than the entering water. Biological treatment based on the metabolic action of the microorganisms can be classified as anaerobic, anoxic and aerobic treatments.

The anaerobic treatment, carried out in the absence of oxygen, utilizes anaerobic bacteria to decompose suspended organic substances. Wastewater or sludge is introduced into a closed tank which is kept under anaerobic conditions and the retention time in the tank is from several days to several weeks. Anaerobic treatment is generally suitable for the treatment of wastes containing high concentrations of organic substances (often used for sludge treatment).

Anoxia is defined as a condition where water is without, or has very low levels of, dissolved oxygen (U.S-EPA, 2006) . Anoxic treatment refers to the growth of microorganisms in anoxic conditions. In wastewater treatment process, this treatment is carried out in anoxic tanks that ultimately reduce the concentration of

nitrate in wastewater. Although this is not the primary metabolic reaction in aerobic treatment, anoxic treatment exists with aerobic treatment under favorable conditions.

Aerobic treatment is a means of oxidizing and decomposing organic substances in wastewater using aerobic microorganisms. Suspended organic substances are oxidized and decomposed by metabolic reactions of microorganisms, which also produces energy. Microorganisms multiply using a portion of this energy and the organic substances present, any excess of microorganisms grown must be separated and disposed of as excess sludge.

## 3.7   Tertiary Treatment Stage

The removal of nutrients such as phosphorous and nitrogen from the treated water using chemicals is considered as tertiary treatment of wastewater, although recently this has been performed using biological mass. Therefore, the removal of nutrients is performed close to the biological unit. The plant configuration can be such that the nutrient removal is either before or after the biological unit. Recent practice has used only the disinfection of treated water in the tertiary treatment stage (Norweco, 2006). The purpose of disinfection in the treatment of wastewater is to substantially reduce the number of microorganisms in the water that will be discharged back into the environment. Common means of disinfection include ozone, chlorine, or ultraviolet light.

## 3.8   PMB Techniques

Process Model-Based (PMB) methods/techniques for Fault Detection and Diagnosis (FDD) are discussed in this section. PMB methods are generally classified as quantitative or qualitative, an introduction and brief discussion is presented below. A detailed, but concise, review of these different methodologies is presented below.

### 3.8.1  Qualitative Model-Based Techniques

For qualitative model-based techniques, the relationships developed are based upon a fundamental understanding of the physical phenomena controlling the process that are expressed in terms of qualitative functions. The qualitative models can be

developed either as qualitative causal models or abstraction hierarchies. Diagraphs, fault trees and qualitative physics are the most popular techniques that belong to the class of casual models. The abstraction hierarchy used for the FDD in a process can be further classified as structural or functional hierarchy. The following is a brief explanation of some of the most commonly used techniques used in qualitative model-based FDD (Venkatasubramanian et al., 2003b).

**a)  Digraphs based causal models**

Cause and effect relations, or models, can be represented in the form of signed digraphs (SDG). Digraph is a graph with directed arcs between the nodes, and SDG is a graph in which the directed arcs have a positive or negative sign attached to them. The directed arcs lead from the 'cause' nodes to the 'effect' nodes. Each node in the SDG corresponds to the deviation from the steady-state value of a variable. SDGs have been the most widely used form of causal knowledge for process fault diagnosis and safety.

**b)  Fault Trees**

Fault tree is a logic tree that propagates primary events or faults to the top level event or a hazard. The tree usually has layers of nodes. At each node different logic operations such as AND and OR are performed for propagation. Fault trees have been used in a wide range of risk assessment and reliability analysis studies. Before the construction of the fault tree, the analyst should possess a complete understanding of the system. The fault tree is constructed by asking questions such as: "What could cause a top level event?" In answering this question, one generates other events connected by logic nodes. Fault trees provide a computational means for combining logic in order to analyze system faults. The attraction of using a fault tree stems from the fact that different logic nodes can be used (OR, AND, XOR) instead of the predominantly OR node used in the digraphs. This helps in eliminating spurious solutions and representing the system in a concise manner. The biggest problem with fault trees is that the development is prone to mistakes at different stages.  It is of primary importance that the underlying logic of the fault tree construction is correct, otherwise the entire model is faulty from the outset.

**c) Qualitative Physics**

Qualitative physics or "common sense" reasoning about physical systems has been an area of major interest in the artificial intelligence community. An important approach in qualitative physics is the derivation of qualitative behavior from the ordinary differential equations (ODEs). These qualitative behaviors for different failures can be used as a knowledge source. The goals of these methodologies are to reason from qualitative physical and equation-based descriptions. The advantage of these methods is their ability to yield partial conclusions from incomplete and often uncertain knowledge of the process.

**d) Abstraction Hierarchy**

Another form of model knowledge is through the development of abstraction hierarchies based on decomposition. There are two-dimensions along which abstraction at different levels is possible, i.e. structural and functional. The structural hierarchy represents the connectivity information of the system and its subsystems. The functional abstraction hierarchy represents the means-end relationships between a system and its subsystems. The majority of the work on fault diagnosis in chemical engineering depends on the development of functional decomposition, and the reason for its popularity is due to the complex functionalities of various units that cannot be expressed in terms of structure.

## 3.8.2 Quantitative Model-Based Techniques

Relying on an explicit model of the monitored plant, all model-based FDD methods require two steps. The first step generates inconsistencies between the actual and expected behavior known as "residuals". The second step chooses a decision rule for diagnosis. Some form of redundancy is always required in order to generate residuals to evaluate the inconsistency. There are two types of redundancies, hardware redundancy and analytical redundancy. Hardware redundancy has been utilized in the control of safety-critical systems such as aircraft, space vehicles and nuclear power plants, and requires redundant sensors. Its applicability in the chemical and process industry sector has been limited due to the additional costs and space required for the extra sensors. However, analytical redundancy is achieved from the functional dependence among the process variables and is usually provided by a set of algebraic or differential relationships among the states, and the inputs and the outputs of the

system. The main concept used in the quantitative model-based FDD techniques is analytical redundancy based upon checking the actual system behavior against the system model for consistency. Any inconsistency, expressed as residuals, can be used for detection and isolation purposes. The residuals should be close to zero when no fault occurs, but show 'significant' values when the underlying system changes.

The generation of the residuals requires an explicit mathematical model of the system, either a model derived analytically using first principles or a black-box model obtained empirically. The first principles models are obtained based on a physical understanding of the process. In a chemical engineering process, mass, energy and momentum balances are used in the development of model equations. Historical models developed from first principles were seldom used in process control and fault diagnosis mainly because of their complexity. In addition, chemical engineering processes are often nonlinear which makes the design of fault diagnosis procedures more difficult. However, this is changing due to increased computational power and speed and better understanding of nonlinear controller design and synthesis(Venkatasubramanian et al., 2003a,b,c).

## 3.9   Mathematical Modelling

A mathematical model describes the fundamental physical phenomena controlling the process expressed in terms of mathematical functional relationships between the inputs and outputs of the system. Most of the work on quantitative model-based approaches has been based on general input-output and state-space models. However, there are a wide variety of quantitative model types that have been considered in fault diagnosis such as first-principles models, frequency-response models, etc. The first-principles models have not been very popular in fault diagnosis studies because of the computational complexity in utilizing these models in real time fault diagnostic systems, and the difficulty in developing accurate models. The most important class of models that have been heavily investigated in fault diagnosis studies are the input-output or state-space models. In this work, the equations describing first principals are used rather than the inputs-output or state-space model of the secondary stage of wastewater treatment. This dynamic model will provide transient values of the process measurements which can be used for the generation of

residuals using the actual values from the process operation for the detection of a fault in the process.

## 3.10  Wastewater Treatment Mathematical Model

In order to promote development, and facilitate the application of practical models for design and operation of biological wastewater treatment systems, the International Association on Water Quality (IAWQ) formed a task group in 1983. The first goal was to review existing models and the second goal was to reach a consensus concerning the simplest mathematical model having the capability of realistically predicting the performance of single-sludge systems. The final result was presented in 1987 as the IAWQ Activated Sludge Model No. 1(ASM1). Although the model has been extended since then, for example to incorporate more fractions of COD (i.e. chemical oxygen demand), to describe growth and population dynamics of floc forming and filamentous bacteria and to include new processes for describing enhanced biological phosphorus removal, the original model is probably still the most widely used for describing WWT (wastewater treatment) processes all over the world (Jeppsson, 2003). Since then ASM1 has been the core of numerous models with a number of supplementary details added in almost every case. The model has grown more complex over the years, from ASM1 to ASM2, including biological phosphorus removal processes and to ASM2d including denitrifying PAOs. In 1998 the task group decided to develop a new modelling platform, the ASM3, in order to create a tool for use in the next generation of activated sludge models. The ASM3 is based on recent developments in the understanding of the activated sludge processes, among which are the possibilities of following internal storage compounds, which have an important role in the metabolism of the organisms (Henez et al., 2000).

The use of ASM1 as the core of recent models was the source of inspiration for adopting ASM1 as the basis for the mathematical modelling of wastewater in this research along with the fact that the data available for the validation (dated back to 1991) of the proposed model only covers the basic processes in the water treatment. The concepts of Monod's kinetics, population balance and Activated Sludge Model 1 (ASM1) are used to derive a mathematical model for the activated sludge treatment

of wastewater which can then be used for fault detection. The major biological and chemical processes occurring in the activated sludge system for the treatment of wastewater are:

1) Production and decay of microorganisms (under different conditions)
2) Utilization of suspended organic material (substrate, i.e. food for microorganisms)
3) Oxygen consumption
4) Production of volatile suspended solids

### 1) Production and decay of microorganisms

The change in microorganism population due to production is given (Morley, 1979) as:

$$\left(\frac{dX}{dt}\right)_p = \mu X \tag{3.1}$$

where

X = concentration of microorganisms. This is the g/m$^3$ of microorganisms present in the activated sludge system for the conversion of wastewater into treated water.

μ = specific growth rate

The specific growth rate in equation 3.1 can be modelled using Monod's kinetics (Shuler and Kargi, 2002) as given by:

$$\mu = \mu_{max}\left(\frac{S}{Ks+S}\right) \tag{3.2}$$

where

$\mu_{max}$ = maximal specific growth rate

Ks = half saturation coefficient

S = substrate concentration (i.e. the food for micro-organisms). Substrate is the suspended solids present in the wastewater. As it is a source of food for the

microorganisms to live on, the suspended solids are called "substrate" in this modelling exercise.

On substitution of equation 3.2 in 3.1, the change in microorganism population due to production is:

$$\left(\frac{dX}{dt}\right)_p = \mu_{max}\left(\frac{S}{K_s+S}\right)X \qquad (3.3)$$

Now, the decay rate of microorganisms due to endogenous metabolism (Morley, 1979):

$$\left(\frac{dX}{dt}\right)_d = -k_d X \qquad (3.4)$$

Hence, the net change in microorganism population can be modelled by:

$$\left(\frac{dX}{dt}\right) = \left[\mu_{max}\left(\frac{S}{Ks+S}\right)-k_d\right]X \qquad (3.5)$$

### 2) Utilization of suspended organic material

This section describes the dynamics of suspended solids (SS) in the wastewater treatment. The rate of substrate utilization (consumption) due to microorganisms (Morley, 1979) is given by:

$$(3.6)$$

where
$$\frac{dS}{dt} = -\frac{\mu X}{Y_{X/S}}$$

S = concentration of substrate (suspended organic material in wastewater)

$Y_{X/S}$ = yield factor (indicates how many units of microorganisms are produced per unit of substrate. Similar to $Y_{X/S}$ are $Y_{X/O2}$ that is the unit of $O_2$ consumed by a unit of microorganisms)

On substitution of equation 3.2 in 3.6, the equation will become:

$$\frac{dS}{dt} = -\mu_{max}\left(\frac{S}{K_s+S}\right)\frac{X}{Y_{X/S}} \qquad (3.7)$$

28

### 3) Oxygen Concentration

The oxygen concentration in the wastewater system is reduced by aerobic growth of microorganisms. Similar to the consumption of suspended organic material (equation 3.7), the rate of oxygen depletion from the process can be derived as:

$$\frac{dS_o}{dt} = -\left[ \mu_{max} \left( \frac{S}{Ks+S} \right) \left( \frac{So}{K_O+So} \right) \frac{X}{Y_{X/O_2}} \right] \tag{3.8}$$

Another important consideration is that the model equations are derived assuming the excess oxygen demand in the process. If the oxygen supply is limited or accounted for, then the factor (So/Ko+So) is to be incorporated in equations 3.5, and 3.7. The resulting models equations are given below by equations 3.9 and 3.10.

Net microorganism growth (equation 3.5):

$$\left( \frac{dX}{dt} \right) = \left[ \mu_{max} \left( \frac{S}{Ks+S} \right) \left( \frac{S_O}{K_o+S_O} \right) - k_d \right] X \tag{3.9}$$

and the rate of substrate utilization (equation 3.7):

$$\frac{dS}{dt} = -\mu_{max} \left( \frac{S}{K_s+S} \right) \left( \frac{S_O}{K_o+S_O} \right) \frac{X}{Y_{X/S}} \tag{3.10}$$

### 4) Production of Volatile Suspended Solids

The decay of microorganisms contributes towards the volatile suspended solids. A portion of the dead microorganisms is recycled as a source of food for other microorganisms, while the particulate part of the dead microorganisms contributes to the concentration of volatile solids. The decay rate can be modelled using death-regeneration hypothesis and the model equation (Jeppsson, 2003) is given as:

$$\frac{dS_{VS}}{dt} = f_P \left( k_d X \right) \tag{3.11}$$

$k_d$ = decay rate of microorganisms

$f_P$ = fraction of biomass yielding particulate products

Equation 3.11 above describes the dynamics of Volatile Suspended Solids (VSS).

The derivation of an equation for the production of "Substrate" from equation 3.11 is based upon a fraction of biomass yielding particulate products, and the remainder is utilized as a food source. The equation then derived for the production of substrate is:

$$\frac{dS}{dt} = (1\text{-}f_P)\left(k_d X\right) \tag{3.12}$$

This equation affects equation 3.7 such that the net rate of substrate utilization is:

$$\frac{dS}{dt} = -\mu_{max}\left(\frac{S}{K_s + S}\right)\frac{X}{Y_{X/S}} + (1\text{-}f_P)\left(k_d X\right) \tag{3.13}$$

**Mass Balance around Activate Sludge System**

If F is the influent flow rate in the activated sludge process, a mass balance will yield:

$$\text{In-Out+Generation-Consumption=Accumulation} \tag{3.14}$$

Substrate (suspended organic solids) mass balance:

$$FSo\text{-}FS+0\text{-}(\mu X/Y_{X/S})V = V\frac{dS}{dt} \tag{3.15}$$

Simplifying the above equation yields:

$$\frac{F}{V}So\text{-}\frac{F}{V}S\text{-}(\mu X/Y_{X/S}) = \frac{dS}{dt} \tag{3.16}$$

where

V = volume of reactor

S = concentration of substrate in reactor

So = initial concentration of substrate in influent. This is the concentration of the suspended solids in influent wastewater whereas S is the concentration of suspended solids in effluent (treated wastewater).

F/V = D (dilution rate, i.e. inverse of residence time) (Shuler and Kargi, 2002)

Equation 3.16, using D will become:

$$DSo\text{-}DS\text{-}(\mu X/Y_{X/S}) = \frac{dS}{dt} \tag{3.17}$$

Equation 3.17 predicts the dynamic behavior of the concentration of suspended solids in the activated sludge system.

**Cell mass balance**

Similarly to the substrate analysis, the microorganisms (cell) balance on the activated sludge process is modelled as:

$$FXo\text{-}FX+\mu XV - k_d XV = V\frac{dX}{dt} \tag{3.18}$$

Simplifying equation 3.18 by using D for F/V:

$$DXo\text{-}DX+\mu X\text{-}k_d X = \frac{dX}{dt} \tag{3.19}$$

Equation 3.19 is similar to equation 3.17 and predicts the dynamics of cell mass in the activated sludge process.

**Volatile Suspended Solids Balance**

Using equation 3.14, a mass balance for the volatile solids is:

$$FS_{VS}\text{-}FS_{VS}+f_P\left(k_d X\right)V = V\frac{dS_{VS}}{dt} \tag{3.20}$$

Simplifying the above equation and using D for F/V:

$$DS_{VS}\text{-}DS_{VS}+f_P\left(k_H X_H + k_A X_A\right) = V\frac{dS_{VS}}{dt} \tag{3.21}$$

## 3.11  Results from Model Simulations

The following simulation results were obtained by using the model equations for the activated sludge system. The biological treatment of wastewater can be carried out using a batch or continuous process, the later is more widely used.

### 3.11.1 Batch Process

The simulation study on a batch reactor is further divided into three steps in order to obtain a better understanding of the process and the effects of different parameters.

**1.** Microorganisms and organic material

The first simulation is for aerobic growth of the heterotrophic microorganism, and their effect on the concentration of the organic material in the wastewater. It is assumed that an excess of dissolved oxygen (DO) exists in the process (i.e. no effect of oxygen concentration on microorganisms). The equations used are given below.

*Microorganism concentration*

The net change in the microorganism's concentration can be predicted from equation 3.5:

$$\frac{dX}{dt} = \left[ \mu_{max} \left( \frac{S}{Ks+S} \right) - k_d \right] X$$

*Organic material concentration*

Using Equation 3.13 for the net concentration of the organic material (i.e. substrate) in the wastewater:

$$\frac{dS}{dt} = -\mu_{max} \left( \frac{S}{Ks+S} \right) \frac{X}{Y_{X/S}} + \left(1 - f_P\right) k_d X$$

Simulation results from the above equations show how the concentration of microorganisms and organic material will change with time in a batch process.

*Simulation Results and Discussion*

**Table 3.1** Parameters used for the simulation in Figure 3.1

| Parameter | Value | Units |
|---|---|---|
| Substrate initial concentration | 1000 | $g/m^3$ |
| Cell initial concentration | 1.5 | $g/m^3$ |
| Residence time | 24 | hr |

**Figure 3.2** Simulation result from batch process assuming excess DO

Figure 3.2 illustrates the dynamics of the wastewater reclamation process using activated sludge in batch conditions. The concentration of organic materials (i.e. pollutant or substrate) reduces to near zero at the end of batch operation.

Figure 3.3 is an extension of Figure 3.2 using the parameters in Table 3.1, it shows how the activated sludge (microorganisms or cell) concentration starts decreasing as the organic material in wastewater is completely consumed by cells.



**Figure 3.3** Simulation result from batch process assuming excess DO and a residence time of 30hr

**2.** Microorganisms, organic material and dissolved oxygen

For this simulation, the earlier assumption of excess DO is rejected. Now the microorganism growth (concentration) depends on the oxygen concentration at any time in the process. If enough DO is provided, the concentration of microorganism will increase resulting in the decrease of organic material concentration. The rate of growth will be reduced if DO is less than the oxygen demand of the microorganisms. The equations used are given below.

*Microorganism concentration*

Using equation 3.9 for the net growth of microorganism including oxygen concentration is:

$$\frac{dX}{dt} = \left[ \mu_{max} \left( \frac{S}{Ks+S} \right) \left( \frac{So}{K_{OH}+So} \right) - k_d \right] X$$

*Organic material concentration*

The model equation for substrate utilization is then derived from equation 3.13 to include the effect of oxygen as:

$$\frac{dS}{dt} = -\left[ \mu_{max} \left( \frac{S}{Ks+S} \right) \left( \frac{So}{K_{OH}+So} \right) \frac{X}{Y_{X/S}} \right] + (1-f_P) k_d X \tag{3.22}$$

Although oxygen is consumed in the batch process, it is provided continuously using turbines. Oxygen supply will fulfill the oxygen demand of microorganisms and also produce turbulence in the reactor that will help to keep the solution of organic material and microorganism suspended in the reactor, thus improving the contact between both phases and resulting in efficient removal of pollutants from wastewater. The amount of DO depends upon the Biological Oxygen Demand (BOD) of the wastewater and, in general, a minimum residual of 1mg of DO per liter of wastewater must be maintained (Buchanan and Seabloom, 2004).

The equation for oxygen consumption derived from equation 3.8 is:

$$\frac{dS_o}{dt} = DSo_{in} - \left[ \mu_{max} \left( \frac{S}{Ks+S} \right) \left( \frac{So}{K_{OH}+So} \right) \frac{X}{Y_{X/O_2}} \right]$$  (3.23)

Simulation from the model equations above will predict the microorganisms and organic material concentration under the influence of the oxygen supplied.

*Simulation Results and Discussion*

Oxygen is always used in excess in the activated sludge treatment of wastewater treatment. For illustration purposes, Figure 3.4 shows the dynamics of activated sludge system using $2g/m^3$ of air.



**Figure 3.4** Dynamics of activated sludge system using oxygen concentration as $2g/m^3$

## 3.11.2 Completely Mixed Reactor

Assume that complete mixing occurs instantaneously and uniformly throughout the reactor as the fluid-particles enter.

1)  Microorganisms and organic material

For the activated sludge process, the model equations assuming an excess of dissolved oxygen (DO) are:

Microorganism concentration predicted using equations 3.5 and 3.14:

$$\frac{dX}{dt} = DX_{in} - DX_{out} + \mu_{max}\left(\frac{S}{Ks+S}\right)X - k_d X \qquad (3.24)$$

Similarly, the organic material concentration from equations 3.13 and 3.15:

$$\frac{dS}{dt} = DS_{in} - DS_{out} - \mu_{max}\left(\frac{S}{Ks+S}\right)\frac{X}{Y_{X/S}} + (1-f_P)k_d X \qquad (3.25)$$

Simulation results from the above equations will predict how the concentration of microorganisms and organic material will change with time in the activated sludge process.

*Simulation Results and Discussion*

It is desirable to study the effect of continuous operation on the wastewater process. This dynamic behavior is predicted by simulations using the parameters in Table 3.2, as shown in Figure 3.5.



**Table 3.2** Simulation parameters for continuous process

| Parameter | Value | Units |
|---|---|---|
| Substrate initial concn. | 1000 | $g/m^3$ |
| Cell initial concn. | 500 | $g/m^3$ |
| Residence time | 4 | hr |
| Substrate influent | 250 | $g/m^3$ |
| Substrate effluent | 10 | $g/m^3$ |
| Cell influent | 0 | $g/m^3$ |
| Cell effluent | 400 | $g/m^3$ |

**Figure 3.5** Simulation dynamics for a continuous process using excess DO

2)  Microorganisms, organic material and dissolved oxygen

The model equations for microorganisms and organic material considering the DO concentration are:

Microorganism concentration using equations 3.9 and 3.14:

$$\frac{dX}{dt} = DX_{in} - DX_{out} + \mu_{max}\left(\frac{S}{Ks+S}\right)\left(\frac{So}{K_{OH}+So}\right)X - k_d X \qquad (3.26)$$

Similarly for the organic material:

$$\frac{dS}{dt} = DS_{in} - DS_{out} - \mu_{max}\left(\frac{S}{Ks+S}\right)\left(\frac{So}{K_{OH}+So}\right)\frac{X}{Y_{X/S}} + (1-f_P)k_d X \qquad (3.27)$$

Oxygen concentration is predicted using equation 3.23:

$$\frac{dS_o}{dt} = DSo - \left[\mu_{max}\left(\frac{S}{Ks+S}\right)\left(\frac{So}{K_{OH}+So}\right)\frac{X}{Y_{X/O_2}}\right]$$

Simulations obtained using these model equations predict the microorganisms and organic material concentrations under the influence of the oxygen supplied.

*Simulation, results and discussion*

Simulation results used to study the dynamic behavior of the processes discussed above are shown in Figure 3.6, and using the simulation parameters given in Table 3.3.



**Table 3.3** Parameters used to study the effect of $O_2$ consumption

| Parameter | Value | Units |
|---|---|---|
| Substrate initial concn. | 1000 | g/ m3 |
| Cell initial concn. | 500 | g/ m3 |
| Residence time | 4 | hr |
| Substrate influent | 250 | g/ m3 |
| Substrate effluent | 10 | g/ m3 |
| Cell influent | 0 | g/ m3 |
| Cell effluent | 400 | g/ m3 |
| DO | 500 | g/ m3 |

**Figure 3.6** Wastewater treatment dynamics for controlled oxygen supply

37

Table 3.4 shows all parameters that are used in the simulations discussed above. The parameters are adopted from the ASM1 model described by Jeppsson (2003).

Table 3.4 Parameters used in modelling and simulation with literature values

| Notation | Explanation | Value |
|---|---|---|
| $\mu_{max}$ | Maximum specific growth rate | 0.25 hr$^{-1}$ |
| $Y_{X/S}$ | Yield factor (unit of microorganisms produced per unit of substrate) | 0.4 |
| $Y_{X/O2}$ | Unit of microorganisms per unit of O2 | 0.9 - 1.4 |
| $Y_{X/NO2}$ | Unit of microorganisms per unit of NO | 0.24 |
| $K_s$ | Half saturation coefficient | 23 m$^{-3}$ |
| $K_{OH}$ | Oxygen half saturation constant | 0.2 m$^{-3}$ |
| $K_{NO}$ | Nitrate half saturation constant | 0.5 m$^{-3}$ |
| $\eta_g$ | Correction factor for anoxic growth | 0.8 |
| $f_P$ | Fraction of biomass yielding VSS | 0.15 |
| $k_d$ | Decay rate of microorganisms | 0.005 hr$^{-1}$ |
| F | Influent flow rate | 936 hr$^{-1}$ m$^3$ |
| V | Volume of reactor | 4021 m$^3$ |
| D | F/V | 0.23 hr$^{-1}$ |

## 3.12  Model Validation

The model developed is then validated before it can be used for the diagnostic purposes. It is essential to establish that the model does actually predict the behaviour of activated sludge process for the treatment of wastewater. Though it is evident from the simulation results (Figure 3.2 to 3.6) that the model does follow the expected behaviour of an activated sludge system, it is still to be established that how accurate the proposed model is. The wastewater treatment plant (WWTP) data (discussed in detail in section 4.5) is used for the validation of the model.

The value of an input parameter from the WWTP data (Appendix A) is used as the initial value for the simulations study of the model to obtain the output value of the parameter. This output value is then compared to the actual output given in the

data set to validate the proposed model. Table 3.5 summarizes the results obtained during the validation of proposed model.

**Table 3.5** Results for model validation

| Date | Input Parameter | Input Value (WWTP data set) | Output Parameter | Output Value (WWTP data set) | Output value (simulations) | Difference (%) |
|------|-----------------|------------------------------|------------------|-------------------------------|-----------------------------|-----------------|
| 11/1/90 | SS-E | 192 | SS-D | 100 | 107 | 7 |
| 1/3/90 | SS-E | 166 | SS-D | 94 | 98 | 4 |
| 3/5/90 | SS-D | 88 | SS-S | 49 | 53 | 8 |
| 29/7/90 | SS-D | 90 | SS-S | 34 | 37 | 9 |
| 28/09/90 | SSV-E | 57 | SSV-D | 77 | 81 | 5 |
| 30/11/90 | SSV-E | 75 | SSV-D | 83 | 89 | 7 |
| 15/1/90 | SSV-D | 71 | SSV-S | 76 | 81 | 7 |
| 8/3/91 | SSV-D | 64 | SSV-S | 85 | 88 | 4 |
| 21/5/91 | DBO-E | 238 | DBO-D | 101 | 105 | 4 |
| 31/7/91 | DBO-E | 170 | DBO-D | 101 | 100 | 1 |
| 6/8/91 | DBO-D | 90 | DBO-S | 16 | 18 | 12 |
| 16/10/91 | DBO-D | 121 | DBO-S | 33 | 36 | 9 |

Table 3.5 summarizes the results obtained from the validation of the activated sludge model proposed. The discrepancy in the model and actual values ranges from 1 – 12% with an average of 6.5%.

## 3.13 Summary Comments

A brief review of the processes involved in the wastewater treatment facility is presented, and a mathematical model derived from the ASM1 model. This model is used in this research work for the detection step of the fault detection and diagnosis system. The ideology used for fault detection is the analytical redundancy where the value of an observed parameter in a process is compared to the value obtained from mathematical model simulations, in order to decide whether the parameter under observation is out of the normal operating limit.

# Chapter 4

# Process History Based Method

## 4.1   Introduction

Process History Based (PHB) methods for fault detection and diagnosis are discussed in this chapter.   This includes an introduction and explanation of the different techniques available including a detailed review of data-mining, a technique used in this research work. Data mining is then applied to the dataset of a wastewater treatment plant, and the results are presented, analyzed and discussed. Some modifications are suggested and implemented in the selected technique of data mining that further improves its efficiency and effectiveness in order to interpret data for fault detection and diagnosis. An introduction and explanation of the Wastewater Treatment Plant (WWTP) dataset used in this work is included.

## 4.2   PHB Techniques

In contrast to the depth of knowledge required for fault detection and diagnosis when using Process Model Based methods, PHB methods only require access to the historical and/or operational data. This data is then used to extract knowledge for input to a diagnostic system. This process of knowledge extraction is known as feature extraction. The PHB methods are classified as either qualitative or quantitative on the basis of the type of knowledge extracted from the database. Expert System (ES) and Qualitative Trend Analysis (QTA) are the most important and widely applied techniques from the class of qualitative PHB approaches. Quantitative techniques are further divided into statistical and non-statistical approaches. Principal Component Analysis (PCA), Independent Component Analysis (ICA) and Partial Least Squares (PLS) form the majority of the statistical quantitative PHB approaches, while the best known techniques in non-statistical quantitative PHB approach is Artificial Neural Networks (ANN). Each of these techniques has its own strengths and weaknesses. During the diagnosis stage of FDD,

the onsite operator or engineer has to identify the symptoms, analyze the symptomatic information, interpret the various error messages and indications, and decide upon the correct diagnosis for the situation. Due to the inherent complexity of chemical processes, the diagnosis requires extensive technical skills and process experience, in addition to a complete understating of the process and some general concepts of diagnosis, in order to carry out the diagnostic operation when a fault is identified in the process. This requires a very experienced engineer with the deep domain-specific knowledge and the knowledge of the "ins-and-outs" of the system (Sun et al., 2007). Many diagnosis methods have been proposed to help operators and engineers perform diagnostic fault analysis. For example, expert systems, neural networks, and genetic algorithms are the most popular approaches among the PHB techniques. The application of an expert system is limited due to the knowledge acquisition because the knowledge-based systems developed from expert rules are very system specific, their representational ability is quite limited, and they are difficult to update. The complexity of a process makes the diagnosis methods based on ANN and Genetic Algorithm (GA) difficult, in addition the inherent nature of ANN approach means they lack the explanation and adaptability properties of a diagnostic system (Chen and Mo, 2004, Venkatasubramanian et al., 2003c, Sun et al., 2007, Yang et al., 2005).

The most helpful presentation of a domain-specific knowledge for a non-expert is the cause-symptom relationship that enables the quick comprehension of the situation. Production rules are the knowledge formalized into "rules" containing an If part and a Then part that explains the cause-symptom relationship in a process. Production rules are one of the most popular and widely used knowledge representation languages.

## 4.3   Data Mining

In modern processes, the computer control and data logging systems are able to easily collect large amounts of data. This data can be used for process monitoring and fault diagnosis, as well as in other decision making activities - if properly analyzed. Data mining is a powerful new technique with the potential to help engineers explore and focus on the most important information available from the

analysis of the historical and/or operational data available from a process. Data mining is defined as "The nontrivial process of extracting implicit, previously unknown, and potentially useful, information from data" (Fayyed et al., 1996). It uses machine learning, statistical and visualization techniques to discover and present discovered knowledge in a form which can be easily understood. It allows users to analyze large databases to solve decision problems encountered in an industry or business sector. The primary goal of data mining is the extraction of knowledge from the available data and is often known as Knowledge Discovery and Data Mining (KDDM).

## 4.3.1 Data Mining Process Description

Data Mining is a complex process which typically involves the following procedures (Fayyed et al., 1996).

*Understanding:* Developing an understanding of the application domain, the relevant prior knowledge, and the goals of mining.

*Creating a target data set:* Selecting a data set, or focus on a subset of variables or data samples, on which discovery is to be performed.

*Data pre-processing and cleaning:* This is frequently time consuming; data pre-processing is needed because most large databases were created for a different purpose from their current applications. Therefore, the data within these databases are not immediately ready to use in knowledge discovery algorithms or other information processing techniques. For example, the data may contain information that is not uniform, the data may be blank or inconsistent, certain data may be continuous while others are categorical, some data may contain sensitive information and require encryption, and finally, some data may contain uncertainties. Hence data pre-processing and cleaning involves basic operations such as the removal of noise or outliers if appropriate, collecting the necessary information to model, deciding on strategies for handling missing data fields, and accounting for noise, time sequence information and known changes.

***Data reduction and projection:*** Finding useful features to represent the data depending on the aim of the task, using dimensionality reduction or transformation, or finding invariant representation from data.

***Choosing the data mining task:*** Depends mainly on the application domain and on the interest of the miner. Decide whether the goal of the KDDM process is summarization, clustering, classification and regression, etc., and identification of several types of data mining tasks for which data mining offers possible answers.

***Choosing the data analysis algorithm(s):*** Selection of the methods to be used to search for patterns in the data. This includes deciding which models and parameters may be appropriate (e.g. models of categorical data are different from models on vectors over the real data) and matching a particular data mining method with the overall criteria of KDDM process.

***Data mining:*** searching for patterns of interest in a particular representational form, or a set of such representations, including clustering, dependency modelling, analysis, visualization, etc. The following steps can significantly aid the data mining process.

*a) Interpretation* - interpreting mined patterns, and possible return to any of the previous steps. This step can also involve visualization of the data given the extracted models.

*b) Using discovered knowledge* - this step involves acting directly on discovered knowledge, incorporating the knowledge into another system for further action, or documenting and reporting the knowledge. It also includes checking and resolving potential conflicts with previously believed or extracted knowledge.

*c) Evaluation of KDDM purpose* - newly discovered knowledge is often used to formulate new hypotheses; also new questions may be posed using the enlarged knowledge base. In this step, the KDDM process is evaluated for possible further use in both refinement and expansion.

The overall process, representing the move from data to information, and ultimately knowledge, is shown in Figure 4.1. Sometimes the analysis step itself is referred to as data mining, although this term is better applied to the whole process.



**Figure 4.1** An overview of the knowledge discovery process, showing the move from data to knowledge or information, through various steps.(Buontempo, 2005)

## 4.3.2 Applications for Data-Mining

The technique of data mining has attracted much interest, not only from information technology companies but also from the industrial and business sectors. The following is the list of domains where this technique is being, or potentially can be, applied (Wang, 1999).

*Manufacturing Process Analysis -* identifying the causes of faults in manufacturing processes.

*Production Design -* developing a system which will give product designers access to data and information from a range of corporate databases deemed essential to their function, in particular, customer complaints, product material features, R&D testing. Access to this data may point to fundamental design anomalies or inefficiencies which would not have been otherwise apparent.

*Scientific Data Analysis -* cataloguing in surveys, the basic processing needed before high-level scientific analysis can occur, scientific discovery over a large data set, e.g. the SKICAT system from JPL/Caltech was used to automatically identify stars and galaxies in a large-scale sky survey for cataloguing and scientific analysis. In the global climate area, spatio-temporal patterns such as cyclones were predicted from large simulated and observational datasets.

*Experimental Results Analysis -* summarizing the experimental results and the predictive models.

*Marketing and Sales Data Analysis -* identifying potential customers, establishing the effectiveness of a sales campaign.

*Investment Analysis -* predict a portfolio return on investment.

*Intelligent agents and World Wide Web (WWW) navigation -* model user preferences from data, collaborative filtering, advertising, etc.

*Fraud detection -* identify fraudulent transactions.

*Loan approval -* establishing the credit worthiness of a customer requesting a loan.

*Portfolio Trading -* trade a portfolio of financial investments by maximizing returns and minimizing risks.

## 4.3.3  Data Mining Approaches

Data mining approaches can be broadly categorized as either descriptive or predictive. Descriptive approaches aim to discover patterns that characterize the data, whereas predictive approaches aim to construct models to predict the outcome of a future event by learning from the observed parameters (Charaniya et al., 2008).

### 4.3.3.1 Descriptive approaches

The descriptive approaches fall into two categories; (a) identifying interesting patterns in the data; and (b) clustering the data into meaningful groups.

**a) Pattern discovery**

Algorithms for finding patterns in very large datasets are one of the key success stories of data mining research. These methods aim to analyze the parameters of various runs to identify a pattern that is observed in a large number of runs. Patterns discovered from process data can provide insights into the relationship between different parameters, and can also be used to discover association rules. Various algorithms have been developed that can mine process data to discover relationships between the parameters of the different runs that satisfy certain properties. The most efficient approaches for finding these patterns are FPgrowth and LPminer (Han et al., 2004).

**b) Clustering**

Clustering methods can be used to group different process runs into subsets (groups) according to the similarity in the behavior of some parameters. Clustering methods can be differentiated along multiple dimensions, one of them being the top-down (partitional) or bottom-up (agglomerative) nature of the algorithm. Partitional methods commence with all process runs (or object/record) belonging to one cluster and they are divided into designated number of clusters. *K*-means, Partitioning Around Medoids (PAM), Self-Organizing Maps (SOM), and graph-based clustering methods are popular examples of partitional algorithms.

By contrast, agglomerative methods start with each run belonging to a separate cluster and the clusters are merged, based on the similarities of their parameter profiles, until the runs have been grouped into a pre-specified number of clusters. Hierarchical agglomerative clustering is the most commonly used agglomerative method (Jain et al., 1999). Most statistical packages, such as S-Plus and R Project provide a range of clustering methods (R-Foundation, 2008, TIBCO-Software-Inc, 2008)

**4.3.3.2 Predictive approaches**

Predictive approaches can be used to analyze a set of process runs that exhibit different outcomes (e.g. final product concentration) to identify the relationship between process parameters and the outcome. The discovered relationships (called model or classifier) can be used to predict the process outcome and provide key

insights into how the predicted outcome might affect other parameters of the run, thereby allowing for an intelligent outcome-driven refinement of the process parameters. Commonly used predictive methods include regression, Support Vector Machines (SVM), Artificial Neural Networks (ANN) and Decision Trees (DT). These methods have been designed for problems that arise when process runs are divided into classes. Two of the commonly used predictive methods are discussed below.

## Artificial Neural Networks (ANN)

ANN models attempt to imitate the signal processing events that occur in the interconnected network of neurons in the brain. An ANN consists of several nodes that are organized into two or more layers. The first layer serves as input for process parameters and the final layer determines the run outcome. Any intermediate layers are referred to as hidden layers. Every node of a hidden layer receives all inputs from the previous layer, performs a weighted average of the inputs, and sends its output to the next layer after a threshold transformation. This process is continued until the final output layer is reached. The weighting factors and threshold parameters are learnt from the training runs in an attempt to minimize the error in classifying the runs (Krogh, 2008).

## Decision Trees (DT)

DT-based methods classify runs recursively based on chosen thresholds for one or more parameters. The process parameter that provides most information about the classes is used to split the runs into two or more branches. Splitting thus results in 'child' nodes that are most separated from each other in terms of the class. Thus, selecting a parameter and its threshold for the split is a key exercise for DT classifiers. This division is repeated until all the runs at a particular node belong to a single class (terminal node) or one or more stopping rules are satisfied. A top-down interpretation of a decision tree is intuitive and it also allows ranking of process parameters according to their relevance (Quinlan, 1990).

## 4.4 Inductive Data Mining

Inductive data mining refers to the technique used for the generation of the decision tree and production rules from a dataset. It is also an effective approach for automated acquisition of expert knowledge to be built into the knowledge base of an expert system. Classification problem is the current research focus in the area of data mining, and decision tree is one of the most widely used classification methods.

The appeal of decision trees for data analysis and as classifier systems originate primarily from three inherent properties: their ability to model non linear relationships; their ease of interpretability; and their nonmetric nature. Decision trees have been found to be able to handle large-scale problems due to their computational efficiency, to provide interpretable results and, in particular, to identify the most representative attributes for a given task. The traditional approach to inducing decision trees based upon given training data involves recursive partitioning which selects partitioning variables and their value in a greedy (indiscriminate acquisition?) manner to optimize a given measure of purity. A greedy algorithm makes each choice in a locally optimized manner and progresses making one greedy choice after another and reduces the problem to a smaller one this way. This methodology has numerous benefits including classifier interpretability and the capability of modelling non linear relationships.

While capable of modelling nonlinear relationships, decision trees retain a high level of interpretability. The typical structure of a decision tree consists of a root node linked to two or more child nodes which may or may not link to further child nodes. Each nominal node within the tree represents a point of decision or data splitting based upon the data (DeLisle and Dixon, 2004)

Most inductive data mining methods for decision tree generation use supervised learning, i.e. learning from a set of pre-classified cases. Many algorithms have been proposed for decision tree generation, e.g. ID3, C4.5, See5.0, CART, SLIQ, SPRINT and BOAT. The best known algorithms used are CART and See5.0 (with earlier versions as ID3 and C4.5). The decision tree created by CART is a binary tree in which each split generates exactly two branches. In the decision tree created by See5.0, each split can generate more than two branches; also See5.0 can solve the

classification problem with continuous-valued attributes. It was developed by Quinlan (1993, 1986, 1990, 1996), and is similar to most other decision tree algorithms in that See5.0 consists of two phases, a building (growing) phase followed by a pruning phase.

## a)  Building Phase

In the building phase, a subset of the training set called the window is chosen at random and a decision tree is formed from it; this tree correctly classifies all objects in the window. All other objects in the training set are then classified using the tree. If the tree gives the correct answer for all these objects, then it is correct for the entire training set and the process terminates. If not, then a selection of incorrectly classified objects is added to the window and the process continues. A decision tree is then constructed in a top-down fashion by iteratively selecting the most informative attribute at the current node in the tree. The most informative attribute for the current node is determined by the splitting criterion, i.e. Gain Ratio (Quinlan, 1993). The Gain Ratio is calculated in the following manner.

*Step 1:* Calculate Info(S) to identify the class in the training set S.

$$Info(S) = -\sum_{i=1}^{n}\left[\left\{\frac{freq(Ci,S)}{|S|}\right\}log_2\left\{\frac{freq(Ci,S)}{|S|}\right\}\right] \qquad (4.1)$$

where,

n = number of classes

$C_i$ = a class

|S| = total number of cases in the training set S

freq (Ci, S) = $|S_i|$ = number of cases in S belonging to the class Ci

**Step 2:** Calculate the expected information value, $Info_X(S)$ for test X to partition S:

$$Info_X(S) = -\sum_{i=1}^{m}\left[\left(\frac{|S_i|}{|S|}\right)Info(S_i)\right]$$  (4.2)

when,

m = number of outputs from test X

**Step 3:** Calculate the information gain after partition according to test X:

$$Gain\ (X) = Info(S) - Info_X(S)$$  (4.3)

**Step 4:** Calculate the partition information value SplitInfo(X) acquiring for S partitioned into m subsets:

$$SplitInfo = -\frac{1}{2}\sum_{i=1}^{m}\left[\left(\frac{|S_i|}{|S|}\right)log_2\left(\frac{|S_i|}{|S|}\right)+\left\{1-\left(\frac{|S_i|}{|S|}\right)\right\}log_2\left\{1-\left(\frac{|S_i|}{|S|}\right)\right\}\right]$$  (4.4)

**Step 5:** Calculate the Gain Ratio of Gain(X) over SplitInfo(X):

$$GainRatio(X) = \frac{Gain(X)}{SplitInfo(X)}$$  (4.5)

The Gain Ratio (X) compensates for the weak point of Gain(X) which represents the quantity of information provided by X in the training set. Therefore, an attribute with the highest Gain Ratio (X) is taken as the root of the decision tree.

## b) Pruning phase

A large decision tree constructed from a training set usually does not retain its accuracy over the whole sample space for over-training or over-fitting. Therefore, a fully grown decision tree needs to be pruned by removing the less reliable branches to obtain better classification performance over the whole instance space, even though it may have a higher error over the training set. A number of empirical methods have been proposed for pruning a decision tree and they can be divided into two types: construction-time pruning (or pre-pruning) and pruning after building a fully grown tree (or post pruning). Pre-pruning methods (e.g. threshold method and X2 test method) are used to decide when to stop expanding a decision tree. A serious limitation in the pre-pruning method is that the criterion to stop a tree is often based

on local information. In contrast, the post-pruning methods (e.g. cost-complexity, critical value and reduced error) use global information. The See5.0 algorithm applies an error-based post-pruning strategy to deal with the over-training problem, which is a pessimistic error pruning method. In practice for each classification node, See5.0 calculates a predicted error rate based on the total aggregate of misclassifications at that particular node (See5.0, 2008).

## 4.5   Wastewater Treatment Plant (WWTP) dataset

A wastewater treatment plant database containing 527 cases representing 527 days of operation is used in this study for the generation of production rules (to be used as the fault libraries for process fault diagnosis). It was collected by Poch and made publicly available by Bejar and Corts of the University of Catalonia, Spain (Sanchez et al., 1997). Each data case is represented by 38 attributes, i.e. process parameters/variables. Out of the 38 attributes, 7 are the output variables, 9 are related to process operational performance, and the rest are the variables related to the influent into the biological unit of a wastewater treatment plant. The 38 attributes are listed in Table 4.1. All attributes are numeric and have continuous values. The units for all parameters are $g/m^3$ except the input flow to plant (Q-E) which has the unit $m^3/day$. The WWTP used in this study consists of Sequential Batch Reactor (SBR).

**Table 4.1** Process variables of the wastewater treatment plant (Wastewaterdatabase, 2006)

| No. | Attribute | |
|---|---|---|
| *Pre Treatment* | | |
| 1 | Q-E | (input flow to plant) |
| 2 | ZN-E | (input zinc to plant) |
| 3 | PH-E | (input pH to plant) |
| 4 | DBO-E | (input biological demand of oxygen to plant) |
| 5 | DQO-E | (input chemical demand of oxygen to plant) |
| 6 | SS-E | (input suspended solids to plant) |
| 7 | SSV-E | (input volatile suspended solids to plant) |
| 8 | SED-E | (input sediments to plant) |
| 9 | COND-E | (input conductivity to plant) |
| *Primary Treatment* | | |
| 10 | PH-P | (input pH to primary settler) |
| 11 | DBO-P | (input biological demand of oxygen to primary settler) |
| 12 | SS-P | (input suspended solids to primary settler) |
| 13 | SSV-P | (input volatile suspended solids to primary settler) |
| 14 | SED-P | (input sediments to primary settler) |
| 15 | COND-P | (input conductivity to primary settler) |
| *Secondary Treatment* | | |
| 16 | PH-D | (input pH to secondary settler) |
| 17 | DBO-D | (input biological demand of oxygen to secondary settler) |
| 18 | DQO-D | (input chemical demand of oxygen to secondary settler) |
| 19 | SS-D | (input suspended solids to secondary settler) |
| 20 | SSV-D | (input volatile suspended solids to secondary settler) |
| 21 | SED-D | (input sediments to secondary settler) |
| 22 | COND-D | (input conductivity to secondary settler) |
| *Performance Inputs* | | |
| 23 | RD-DBO-P | (performance input biological demand of oxygen in primary settler) |
| 24 | RD-SS-P | (performance input suspended solids to primary settler) |
| 25 | RD-SED-P | (performance input sediments to primary settler) |
| 26 | RD-DBO-S | (performance input biological demand of oxygen to secondary settler) |
| 27 | RD-DQO-S | (performance input chemical demand of oxygen to secondary settler) |
| 28 | RD-DBO-G | (global performance input biological demand of oxygen) |
| 29 | RD-DQO-G | (global performance input chemical demand of oxygen) |
| 30 | RD-SS-G | (global performance input suspended solids) |
| 31 | RD-SED-G | (global performance input sediments) |
| *Outputs* | | |
| 32 | PH-S | (output pH) |
| 33 | DBO-S | (output biological demand of oxygen) |
| 34 | DQO-S | (output chemical demand of oxygen) |
| 35 | SS-S | (output suspended solids) |
| 36 | SSV-S | (output volatile suspended solids) |
| 37 | SED-S | (output sediments) |
| 38 | COND-S | (output conductivity) |

## The Process Details

The plant is an activated sludge process located in Manresa, a town near Barcelona (Catalonia, Spain) population of 100,000 inhabitants. It treats a daily flow of approx. 35,000 $m^3$ comprising mainly domestic wastewater, although other wastewaters from industries located near the town are also received in the plant. The plant consists of three main treatment sections (Albazzaz et al., 2005):

(i)      Pre-treatment,

(ii)     Primary treatment and

(iii)    Secondary treatment by means of activated sludge.

The database has been used for studies in classification by Sanchez  et al., (1997) where two methods, the K-means clustering method and Linneo+ methodology, a knowledge acquisition tool with unsupervised learning strategy, were investigated. ( Sanguesa and Cortes (1997) used the data to study a possibilistic network. Hunag and Wang (1999) used the data in developing fuzzy casual networks. Wang et al. (2004) used this data set to present an approach for multidimensional visualization of multiple principal coordinates using a technique called parallel coordinates for the purpose of process monitoring. This data was also used for the historical data analysis and an empirical comparison was made between multidimensional visualization using parallel coordinates, PCA based multivariable statistical process control charts, the T2 and SPE charts, and a clustering approach (Albazzaz et al., 2005).  Ma and Wang (2009) used the data for a new approach to data mining using Genetic Programming.  Dellana and West (2009) used the data in their work for the predictive modelling of wastewater. West and Mangiameli (2000) employed the WWTP data for the identification of process conditions.

There is an adequate amount of research work in literature where people have used different dataset for the fault detection and diagnosis systems but the data set they used, is normally not available publically. Researchers usually use dataset from industry with which they have research ties. Furthermore, all the researchers who have used this WWTP data-set which in non-linear in nature, only rely on it and

didn't use any other data-set for any validation etc. There is only one publication in the writer's knowledge that use another dataset for validation but the dataset used there was from their own source and that dataset is not available.

Out of the 38 attributes, the database has some missing values for about 144 days out of the 527. Missing data is a common problem for data mining, because in many of these situations, the missing data cannot be re-collected or reproduced. Albazzaz et al. (2005) used eight different methods available in the commercial statistics software system SPSS (ver. 11.5) to deal with the missing values in the database. Only three approaches, i.e., linear trend at point, series mean, and series median were able to give estimations for all the missing values. It was found that the estimations for the missing values using these three methods were reasonably close. A further comparison was made between series mean and linear trend. Both approaches were used to fill in all the missing values, and then calculated such statistics as minimum, maximum, mean, median, and standard deviation for all the 38 attributes in the wastewater treatment plant dataset. It was found that the differences in these statistics between the two approaches were negligible. Eventually the series mean method was used to fill in the missing values. The cleaned data (data after the missing values were filled) is used in this work.

## 4.6   Development of Production Rules from WWTP data

The data discussed above is used for the development of production rules via decision trees that are intended to be used in the diagnostic section of the hybrid system.

*Step 1*

To classify the data into normal, high or low operating conditions, either a descriptive approach of data mining (i.e. pattern recognition or clustering) or a prior knowledge about the process parameter under observation can be used. In this work, the normal value of a parameter using prior knowledge is used rather than using pattern recognition and clustering. The rational is that it is mandatory to provide the normal value (value that can be used to classify the data) of a parameter under study

for Inductive Data Mining using See5.0. Furthermore, this information is readily available from the process. Initially the parameter SS-S, i.e. Suspended Solids out from Secondary stage of wastewater is used for the study and the normal value is 20mg/L.

*Step 2*

The final product value (concentration, flow rate, etc.) in the chemical or process industries is not always fixed to a single numeric value. There exists a range (or limits) for a parameter such that if the output value for the parameter falls between the limits, it is assumed to be normal; otherwise it is high or low if the value is above or below the range for normal operation. From the diagnostic point of view, this is known as the Tolerance Limit of a parameter. As an example, tolerance limit of 15% is used below. This 15% tolerance limit means that if the value of the parameter under observation changes more than 15% of its set value, it indicates the presence of a fault.

*Step 3*

By applying the tolerance limit on the normal value, the range for normal operation of SS-S is calculated to be 17-23mg/L. The entire data is then split into normal, high and low classes using this limit with each class representing the corresponding operating condition. Although the low value of output product (concentration of pollutant in this case) in a wastewater is desired, the data that corresponds to a low value of SS-S is classified as at low operating condition.

*Step 4*

This data is then segregated into training and test data. Training data is used by the algorithm for its learning process, whereas the test data is unseen (by the algorithm) data which it uses for the validation of the model and its results. In this exercise, 75% of the total data is used for training and 25% for test.

**Results from WWTP Dataset**

After the processing of data as described above, it is then used in See5.0 to obtain production rules via decision tree. As we intend to use these results in fault diagnosis, the set of results that explains the cause-consequence relationship for process parameters while the process is in normal operating condition is rejected. This means that the production rules belonging to the normal class are not used any further, although this information may be important and helpful for other decision-making activities. The number of cases occurring in different classes for the WWTP data is reported in Table 4.2.

**Table 4.2** WWTP data in different classes

| Class | Total Data | Training | Test |
|---|---|---|---|
| *Normal* | 199 | 149 | 50 |
| *Low* | 193 | 143 | 50 |
| *High* | 135 | 101 | 34 |

The production rules explaining only low and high classes are listed in Table 4.3.

**Table 4.3** Production rules for SS-S

*CLASS HIGH*

**RULE NO 1:**
IF RD-SS-G ≤ 86.9
& DBO-E > 89.0
& SS-E > 168.0
**(53)**

**RULE NO 2:**
IF 86.9 < RD-SS-G ≤ 90.2
& SS-E > 214.0
& DBO-E ≤ 283.0
**(17)**

**RULE NO 3:**
IF RD-SS-G ≤ 83.3
& DBO-E > 89.0
& SS-E ≤ 168.0
**(20)**

**RULE NO 4:**
IF RD-SS-G > 90.2
& SED-S > 0.03
& SSV-E ≤ 71.7
& Q-E ≤ 33999
& SS-E > 238
**(7)**

*CLASS LOW*

**RULE NO 1:**
IF 86.9 < RD-SS-G ≤ 90.2
& 144.0 < SS-E ≤ 214.0
& DQO-E ≤ 340.0
**(2)**

**RULE NO 2:**
IF 90.2 < RD-SS-G ≤ 90.8
& SED-S ≤ 0.03
& SS-E ≤ 176.0
**(15/2)**

**RULE NO 3:**
IF RD-SS-G > 92.7
& SED-S ≤ 0.03
& SS-E > 194.0
& RD-DQO-S ≤ 66.7
**(27)**

**RULE NO 4:**
IF RD-SS-G > 90.2
& SED-S > 0.03
& SSV-E > 71.7
**(3)**

**RULE NO 5:**
IF RD-SS-G > 90.8
& SED-S ≤ 0.03
& SS-E ≤ 194.0

**(61)**

**RULE NO 6:**
IF RD-SS-G > 92.7
& SED-S ≤ 0.03
& 194.0 < SS-E ≤ 262.0
& RD-DQO-S > 66.7
**(23/1)**

**RULE NO 7:**
IF RD-SS-G > 92.7
& SED-S ≤ 0.03
& SS-E > 262.0
& DQO-S ≤ 17.0
& RD-DQO-S > 81.8
& PH-D ≤ 7.9
**(7)**

**RULE NO 8:**
IF RD-SS-G > 92.7
& SED-S ≤ 0.03
& SS-E > 262.0
& DQO-S ≤ 17.0
& RD-DQO-S > 81.8
& PH-D > 7.9
& Q-E ≤ 41073
**(2)**

These derived rules are then used in fault libraries for the diagnosis of an abnormal event. Fault libraries are a knowledge-base that the diagnostic algorithm uses to determine the reason for, and solution to, an abnormal event in a process. Rule No.1 from the high class corresponding to high operating condition of the process is used for illustration and explanation as follows.

***Rule No 1:***

If RD-SS-G $\leq$ 86.9

& DBO-E > 89.0

& SS-E > 168.0

**(53)**

Production Rule No. 1 implies that if at any stage of operation the value of global performance input for suspended solids RD-SS-G starts decreasing, and at that instance the biological demand of oxygen and suspended solids inputs to the plants are greater than 89.0 and 168.0 respectively, there is a probability of at least 53% that the output value for suspended solids from the plant will be higher than 23 mg/l (i.e. the maximum allowable limit for SS-S). The number (53) is the number of cases that have been classified to be at high operating condition (w.r.t. SS-S) by the information obtained from Rule No. 1. This probability can then determine if the total numbers of cases that belong to the high operating condition are known (101 in this case, see Table 4.2).

## 4.7   Weighting of Production Rules

A modification is suggested here for the enhancement of the production rules to be used in a diagnosis system. A methodology for weighting the production rules is presented and the production rules that correspond to the same class (for example high or low class) are arranged according to their weighting to further simplify and enhance the process of diagnosis.

There are a total of four production rules that describe the reason for a shift of the process to high operating condition. If a fault has been detected in the process

indicating high operating condition and the plant operator/engineer is provided with the four possible cause-symptom relationships in the form of production rules for diagnosis, then the first essential question to be answered is: "Which of the four production rules is most important and requires attention first?" As a timely response is very critical at that stage, information that can provide an operator/engineer with guidance in decision making regarding the most promising production rules is very helpful. As an example, there are 8 production rules that give comprehensive information about the reason why the process is in low operating condition. The number of production rules normally increases with an increase in the amount of data, therefore a large number of production rules are possible if there is much historical/operational data available from a process. In this case, the number of production rules can itself pose a challenge in the diagnosis of abnormal event. A simple concept of "weighting function" is introduced here that will help prioritize the knowledge obtained from production rules in the diagnostic process. Weighting function can be obtained from the relationship described in equation 4.6 below:

$$wtF_i = \frac{W_i}{\sum_{i=1}^{n} W_i} \tag{4.6}$$

where

$wtF_i$ = weight function of production rule "i"

$\sum_{i=1}^{n} W_i$ = total weight for all production rules (belonging to one class)

$W_i$ = weight of production rule "i" (calculated from equation 4.7) given as:

$$W_i = \frac{\left(\frac{C_c - C_{mc}}{T_c}\right)}{n} \tag{4.7}$$

when

$C_c$ = number of cases correctly classified by the production rule "i"

$C_{mc}$ = number of cases misclassified by the production rule "i"

$T_c$ = total number of cases that belong to a class "i" (e.g. high or low)

$n$ = number of leaf nodes/variables in production rules

If the production rules are arranged according to their weight function then this will help a plant operator/engineer to prioritize the knowledge during the decision-making process. The weight function for the production rules is reported in Table 4.4. The production rules reported in Table 4.3 are then rearranged using this weighting concept. This concept will help develop a hierarchy of production rules on the basis of their ability to diagnose a fault more accurately. The production having the highest weight function becomes the first candidate out of all those possible in the event of a fault diagnosis. The revised production rules are reported in Table 4.5.

**Table 4.4** Production rules for low and high class by their weight function (ascending order)

| Class Low | | Class High | |
|---|---|---|---|
| *Production Rule No.* | *Weight Function* | *Production Rule No.* | *Weight Function* |
| 3 | 0.143 | 1 | 0.173 |
| 5 | 0.048 | 3 | 0.067 |
| 6 | 0.038 | 2 | 0.056 |
| 4 | 0.030 | 4 | 0.014 |
| 7 | 0.010 | | |
| 2 | 0.006 | | |
| 1 | 0.003 | | |
| 8 | 0.002 | | |

**Table 4.5** Production rules for SS-S arranged by weight function

## *CLASS HIGH*

**RULE NO 1***(RANK UNCHANGED):*
IF RD-SS-G ≤ 86.9
& DBO-E > 89.0
& SS-E > 168.0
**(53)**

**RULE NO 2***(FORMALY RULE NO 3):*
IF RD-SS-G ≤ 83.3
& DBO-E > 89.0
& SS-E ≤ 168.0
**(20)**

**RULE NO 3***(FORMALY RULE NO 2):*
IF 86.9 < RD-SS-G ≤ 90.2
& SS-E > 214.0
& DBO-E ≤ 283.0
**(17)**

**RULE NO 4***(FORMALY RULE NO 4):*
IF RD-SS-G > 90.2
& SED-S > 0.03
& SSV-E ≤ 71.7
& Q-E ≤ 33999
& SS-E > 238
**(7)**

## *CLASS LOW*

**RULE NO 1***(FORMALY RULE NO 3):*
IF RD-SS-G > 92.7
& SED-S ≤ 0.03
& SS-E > 194.0
& RD-DQO-S ≤ 66.7
**(27)**

**RULE NO 2***(FORMALY RULE NO 5):*
IF RD-SS-G > 90.8
& SED-S ≤ 0.03
& SS-E ≤ 194.0
 **(61)**

**RULE NO 3***(FORMALY RULE NO 6):*
IF RD-SS-G > 92.7
& SED-S ≤ 0.03
& 194.0 < SS-E ≤ 262.0
& RD-DQO-S > 66.7
**(23/1)**

**RULE NO 4***(FORMALY RULE NO 4):*
IF RD-SS-G > 90.2
& SED-S > 0.03
& SSV-E > 71.7
**(3)**

**RULE NO 5***(FORMALY RULE NO 7):*
IF RD-SS-G > 92.7
& SED-S ≤ 0.03
& SS-E > 262.0
& DQO-S ≤ 17.0
& RD-DQO-S > 81.8
& PH-D ≤ 7.9
**(7)**

**RULE NO 6***(FORMALY RULE NO 2):*
IF 90.2 < RD-SS-G ≤ 90.8
& SED-S ≤ 0.03
& SS-E ≤ 176.0
**(15/2)**

**RULE NO 7***(FORMALY RULE NO 1):*
IF 86.9 < RD-SS-G ≤ 90.2
& 144.0 < SS-E ≤ 214.0
& DQO-E ≤ 340.0
**(2)**

**RULE NO 8***(FORMALY RULE NO 8):*
IF RD-SS-G > 92.7
& SED-S ≤ 0.03
& SS-E > 262.0
& DQO-S ≤ 17.0
& RD-DQO-S > 81.8
& PH-D > 7.9
& Q-E ≤ 41073
**(2)**

Similar to the development of production rules for the one output variable discussed above (i.e. Suspended Solids, SS-S), the production rules for other output variables namely Volatile Suspended Solids(VSS-S), Sediments (SED-S), Biological Oxygen Demand (DBO-S), and Chemical Demand of Oxygen (DQO-S) are developed and presented in Tables 4.6 to 4.9 respectively.

**Table 4.6** Production rules for VSS-S arranged by weight function

### *CLASS HIGH*

**RULE NO 1***(FORMALY RULE NO 2)***:**
IF COND – P > 921
& COND – D > 51
& SS-D > 2550
(6)

**RULE NO 2***(FORMALY RULE NO 1)***:**
IF COND – P > 921
& COND – D > 51
& DBO – D ≤ 78.9
& PH – D ≤ 104
& SS – D > 1846
& SED – D > 11
& ZN – E ≤ 1.8
& RD – DBO – P ≤ 8.0
& 235 < RD – SED – P ≤ 274
(3)

### *CLASS LOW*

**RULE NO 1***(RANK UNCHANGED)***:**
IF COND – P ≤ 921
& SS – D ≤ 863
**(5)**

**RULE NO 2***(FORMALY RULE NO 4)***:**
IF COND – P ≤ 921
& SS – D > 863
& COND – E > 829
& RD-DBO-G > 90.4
**(3)**

**RULE NO 3***(FORMALY RULE NO 2)***:**
IF COND – P > 921
& COND – D > 51
& DBO – D ≤ 78.9
& SS – D ≤ 921
**(3)**

**RULE NO 4***(FORMALY RULE NO 3)***:**
IF COND – P ≤ 921
& SS – D > 863
& COND – D > 829
& RD – DBO – G ≤ 90.4
& SED – D ≤ 12
**(3)**

**Table 4.7** Production rules for SED-S arranged by weight function

## *CLASS HIGH*

**RULE NO 1***(RANK UNCHANGED):*
IF RD-SS – G > 79.5
& RD – DQO –G > 81.7
& COND – D ≤ 39
& RD – SED – P > 204
 **(3)**

**RULE NO 2***(RANK UNCHANGED):*
IF RD-SS – G > 79.5
& SSV – P > 63.6
& PH – E > 7.6
& 39 < COND – D ≤ 52
& RD – DQO –G > 86.8
 **(4)**

## *CLASS LOW*

**RULE NO 1***(FORMALY RULE NO 2):*
IF RD – DQO – G ≤ 55.6
& COND – P > 1165
 **(3)**

**RULE NO 2***(FORMALY RULE NO 3):*
IF COND – D > 103
& 55.6 < RD – DQO – G ≤ 66
& RD – SED – P ≤ 220.0
 **(3)**

**RULE NO 3***(FORMALY RULE NO 1):*
IF RD – DQO – G ≤ 55.6
& COND – P ≤ 1165
& RD – DQO – S ≤ 74.4
**(3)**

**RULE NO 4***(FORMALY RULE NO 4):*
IF COND – D > 103
& 55.6 < RD – DQO – G ≤ 66
& RD – SED – P > 220.0
& DBQ – E > 488
**(3)**

**Table 4.8** Production rules for DBO-S arranged by weight function

## CLASS HIGH

**RULE NO 1***(FORMALY RULE NO 2):*
IF COND – P > 2340
& SS –D > 2550
 **(6)**

**RULE NO 2***(FORMALY RULE NO 1):*
IF COND – P > 2340
& 921 < SS –D ≤ 2550
& RD – SED – P ≤ 295
 **(2)**

## CLASS LOW

**RULE NO 1***(RANK UNCHANGED):*
IF SS-D ≤ 921
& DQO-D ≤ 0.1
 **(6)**

**RULE NO 2***(FORMALY RULE NO 3):*
IF SS-D ≤ 921
& SS-P ≤ 278
& DQO-D > 0.2
**(5)**

**RULE NO 3***(FORMALY RULE NO 2):*
IF SS-D ≤ 921
& SS-P ≤ 278
& 0.1 < DQO-D ≤ 0.2
& COND-D > 64
 **(3)**

**Table 4.9** Production rules for DQO-S arranged by weight function

## *CLASS HIGH*

**RULE NO 1***(FORMALY RULE NO 2):*
IF SED-E ≤ 7.0
& RD-DBO-G > 88.7
& RD-SS-P ≤ 64
& DBO-P > 134
 **(4)**

**RULE NO 2***(FORMALY RULE NO 1):*
IF SED-E ≤ 7.0
& RD-DBO-G > 88.7
& RD-SS-P > 64
& RD-SED-G ≤ 99.7
& DBO-D > 64.5
& DBQ-E > 297
& DBO-P > 296
& ZN-E > 0.4
 **(2)**

## *CLASS LOW*

**RULE NO 1***(FORMALY RULE NO 4):*
IF SED-E ≤ 7.0
& RD-DBO-G ≤ 88.7
& DBO-P ≤ 145
& RD-SED-P ≤ 317
& RD-DBO-P ≤ 7.8
& RD-DBO-S > 16
& DBQ-E > 319
 **(6)**

**RULE NO 2***(FORMALY RULE NO 3):*
IF SED-E ≤ 7.0
& RD-SS-P > 64
& RD-SED-G > 99.7
& RD-SS-G > 90.7
& DBO-P ≤ 146
& RD-DBO-G > 92.7
 **(3)**

**RULE NO 3***(FORMALY RULE NO 2):*
IF SED-E ≤7.0
& DBQ-E > 297
& DBO-P > 145
& DBO-D > 64.5
& 84.1< RD-DBO-G ≤ 88.7
& COND-D > 100
& 159 < SS-P ≤ 228
& RD-DBO-P ≤7.8
& RD-DBO-S ≤31
 **(4)**

## 4.8　Validation of Production Rules

Validation of a model is a very critical step in model development (either a mathematical model or model in the form of production rules). A model may initially produce very encouraging results but its accuracy on the new process or data set actually decides if the model is fit for the intended purposes. In this case study, the WWTP data set was divided into two portions i.e. 75% and 25%. The major portion of the data was used as training data-set for the learning of See 5.0 and the development of production rules. The remaining 25% is then used as a test data-set to validate the model (to measure the accuracy of the production rules See 5.0 developed on the training data-set). The following table summarizes the results obtained on the output variables used in this study.

**Table 4.10** Validation results of Production Rules developed by See 5.0

| | See 5.0 Results (% accuracy) | |
|---|---|---|
| **Variable** | **Training data-set** | **Test data-set** |
| SS-S | 97.5 | 96.9 |
| VSS-S | 97.3 | 95.5 |
| SED-S | 98.1 | 96.6 |
| DBO-S | 99.7 | 94.7 |
| DQO-S | 96.4 | 91.7 |

## 4.9　Concluding Remarks

A review of different data driven techniques is presented. The underlying principles of Inductive Data Mining, a technique selected for the research work, are discussed. Production Rules are developed using Inductive Data Mining on the Spanish data set of a wastewater treatment plant in order to build the knowledge base (fault libraries) of the diagnostic system. A new concept of Weight Factor is introduced and implemented on the production rules generated in the chapter. This concept helps arrange the production rules for easy retrieval.

# Chapter 5

# Diagnostic Algorithm and Java Application

## 5.1    Introduction

An algorithm is presented that uses the mathematical model from Chapter 3 for the detection of faults (as discussed in chapter 3) and the fault libraries that are made up of the production rules (derived from the WWTP dataset in chapter 4) for the diagnosis purposes.  These component parts provide a complete fault detection and diagnosis system. The ideology behind development of this diagnostic system, together with the assumptions made, is discussed below, and is followed by the development of a software platform (a Java application). This Java application is used for simulation of the process using the mathematical model for the detection of faults. The implementation of the diagnostic algorithm enables the Java application to be used for the diagnosis of faults in the process. A discussion of the development and architecture of the application is presented, followed by the results obtained.

## 5.2    Diagnostic Algorithm

As discussed in Section 2.4, it has been shown by Venkatasubramanian and co-workers (2003b, 2003c, 2003a) that a diagnostic system is a valuable area of research in fault detection and diagnosis. This section presents a new diagnostic system for the detection and diagnosis of fault in a process.

Analogous to the medical profession, where graduates have the opportunity to apply their expertise in many fields of medicine as a physician or a surgeon, it is proposed here that instead of using available techniques to perform simultaneously detection and diagnosis of faults, these techniques should be classified into Detection and Diagnosis techniques on the basis of their strengths. Furthermore, any future research should be focused on improving the specialty of a technique (either detection or diagnosis of faults). Based on this idea, a diagnostic algorithm is

presented here that uses analytical redundancy based upon the mathematical model for the detection of a fault and production rules for the diagnosis.

While developing the algorithm, the following assumptions were made:

a) Due to the complexity, multiple faults are not considered in this study.
b) For the detection of faults, the data is considered to be noise free.
c) The term: "multi faults" used in this work (later in Table 5.2) refers to the two faults that occur simultaneously with distinct characteristics and no interconnecting relationship exists between the two.

The structure of a typical fault detection and diagnosis system is shown in Figure 5.1.



**Figure 5.1** A schematic of fault detection and diagnosis

**Fault detection** is the first step in a diagnosis system. Fault detection is determining whether a fault has occurred in the process under observation. It helps engineers/operators identify if a fault exists.

**Fault diagnosis** is the determination of the cause of the observed out-of-control status of a process variable. In other words, fault diagnosis is determining the type, location, magnitude and time of the fault. This step is essential in counteraction or elimination of a fault.

Some researchers split diagnosis further into the fault isolation and its identification - which when combined have the same meaning as that defined above. In this work, rather than fault isolation and identification, the term fault diagnosis is used. Another term, "diagnosis system", when used in this work refers to both the fault detection and the diagnosis system rather than only the diagnosis system.

## 5.2.1 The Algorithm

A comprehensive survey of available and most commonly used techniques for the detection and diagnosis of faults is given by Venkatasubramanian et al. (2003b, 2003c, 2003a). In their extensive review, it was shown that despite all the research for improvement in detection and diagnosis of process faults, it is a widely accepted

by most researchers that none of the techniques introduced so far have the potential to meet the key requirements of a practical diagnosis system (as defined by the ten key characteristics of a diagnostic system in Venkatasubramanian et al. (2003a)). One possibility for enhancing the effectiveness of a diagnosis system is to adopt the hybrid methodology. This is based on the assumption that these methods can complement each other's limitations in different areas, thus resulting in an overall improved diagnostic system. As an example, fault explanation through a causal chain is best done through the use of digraphs, whereas fault isolation might be difficult using digraphs due to the qualitative ambiguity (Venkatasubramanian et al., 2003a).

Building upon the detailed review by Venkatasubramanian et al. (2003b, 2003c, 2003a), this work presents a new approach for fault detection and diagnosis employing both mathematical models and historical data from a wastewater treatment plant (for detection and diagnosis respectively). The knowledge base of the diagnostic system is built up of the historical data using inductive data mining and application of the commercial software See5.0. The mathematical model is used for the detection of faults in the process. A step by step overview of the system is given below.

### a)    *Fault Detection*

Assume the system has inputs $u_{(i)}$ and outputs $y_{(i)}$ ( i = 1 to n). Under the fault free operational mode, the output $y_{(i)}$ will conform to the input $u_{(i)}$ . Now consider the system with a fault $f_{(i)}$ and/or a process disturbance $d_{(i)}$ as shown below:



### Step 1:  Residual Generation

The first step in the system is the residual generation where residual (*r*) is defined as an error in a result. In the context of fault detection and diagnosis, it is the difference between the actual and the desired value of any process parameter. The residual in the output parameter can be calculated from the following equation:

$$r_{(i)} = y_{(i)} - y'_{(i)}$$

where $y_{(i)}$ is the output variable and $y'_{(i)}$ is the desired (set) value for parameter $y_{(i)}$. For the purpose of its implementation in computer code, although the value of $r_{(i)}$ is calculated above but $|r_{(i)}|$ will be used in next steps (the Java application need only positive value to decide if the process is performing above or under the control limits).



The desired value $y'_{(i)}$ can be obtained in a number of ways. It includes historical data, prior knowledge and principal models of the process. A statistical method can be applied on historical data to determine a value of the output that is considered to be in the normal operating state. Prior knowledge from plant engineers and operators can be useful in deciding the normal operating value for any output variable. Another way to obtain the normal operating value for a process variable as used in this work is the use of mathematical models. A mathematical model can be obtained by carrying out a simple mass and energy balance on the process under consideration. For a given input **u**, it is possible to find out the expected normal value of the output via mathematical models by the aid of computer simulations. In this study, the wastewater mathematical model (chapter 3) is used for the simulation. And before it can be employed for the residual generation in this step, it was validated against the wastewater data-set.

**Step 2: Residual Evaluation**

Once the residual is generated in Step 1, the next critical step is to evaluate this residual to establish if any inconsistency exists in the system. This is known as residual evaluation. The inputs to this step are the outputs from residual generation, i.e. $|r_{(i)}|$, and $\tau'_{(i)}$ which is the tolerance for each residual generated. Tolerance is similar to the desired value and can be obtained in different ways. Tolerance helps to identify if the process, although deviated from the expected value, is still within the

normal operating state. If not, then determine the direction of the drift, i.e. +ve (operating state above the allowable maxima; i.e. high) or −ve (operating state below the allowable minima; i.e. low). The final calculation in this step yields the value of error **e** (which potentially indicates a fault), and it can be calculated from the following relationships:

$$e_{(i)} = 0, \ |r_{(i)}| \leq \tau'_{(i)}$$

$$and$$

$$e_{(i)} = 1, \ |r_{(i)}| > \tau'_{(i)}$$



A value of $e_{(i)} = 0$ indicates that the system is in normal operating state, where a value of $e_{(i)} = 1$ indicates the presence of a potential fault. Further analysis is carried out once the value of $e_{(i)}$ has been assigned.

**Step 3:  Disturbance Detection**

The decoupling of a disturbance from a fault is a key research area in fault detection and diagnosis. Recursive use of mathematical models is employed in this diagnosis system which can isolate a disturbance from a process fault.  A mathematical model is employed using process input $u_{(i)}$ to calculate the expected value of the output variable, i.e. $y''_{(i)}$.  With the new expected value of the output variable and output value from the process $y_{(i)}$, then Step 2 and Step 3 are repeated for the detection of disturbance in this step. The value of **e** now determines whether this unwanted incident in the process is an uncontrollable input (process disturbance) or a fault;

where $e_{(i)} = 0$ indicates a disturbance in input $u_{(i)}$; and $e_{(i)} = 1$ indicates a fault in the process.



Depending on the nature of the disturbance, i.e. step, ramp or continuous (reconfiguration of input), the new calculated value of $y''_{(i)}$ can be passed onto Step 2 for either temporary or permanent use as the set value of output $y_{(i)}$.

## Step 4:  Fault Detection

If Step 3 produces $e_{(i)} = 1$, then it indicates the presence of a fault in the process which effects the value of the output variable $y_{(i)}$. The results are then passed to the diagnosis section of the system along with the information about the variable, i.e. $y_{(i)}$, and the direction of fault, i.e. (+ve or −ve drift).

An integrated overview of the fault detection section for this diagnostic methodology is given in Figure 5.2.



**Figure 5.2** Integrated view of Fault detection

### b) *Fault Diagnosis*

After the detection of a fault, the output variable away from the normal operating condition and its direction (+ve drift or −ve drift, where +ve refers to the process performing above the control limit and −ve refers to its performance under the control limit) is then passed onto the diagnosis section of this diagnostic system. For diagnosis of a fault, the major requirement is a knowledge base that can be used to compare symptoms and produce a corrective action for the fault under observation. In this diagnostic methodology, the knowledge base is produced from the historical data. Valuable information and hidden interactions between the process variables can

be obtained by exploring the data. The commercial software See5.0 is used in this work to conduct inductive data mining on the historical data available from a wastewater plant. The results can be obtained in the form of a decision tree or production rules.

**Step 1:  Inductive data mining**

Two inputs are required in this step in order to produce a meaningful relationship between the measured variables of the process, i.e. data, and the classification basis for the data. The data is classified using the numeric values from the tolerance $\tau'_{(i)}$ used in the detection section of this work. It should be noted that classification of data needs to be carried out on every measurable output $y'_{(i)}$ using its respective tolerance $\tau'_{(i)}$. The production rules for each of the variables used are obtained here. These production rules explain the reason for the deviation of a variable from its normal operating state. There will be essentially two sets of production rules, one that explains the +ve drift of the variable and other for −ve drift.



**Step 2:  Fault Libraries**

From the production rules obtained in the above step, fault libraries are built in order to compare the fault symptoms and hence correctly diagnose the fault.  The structure of a fault library used in this work is as follows;

**Table 5.1** The fault libraries used in the knowledge base of the diagnosis system

| Fault Library | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Parameter** | $y_{(1)}$ | | $y_{(2)}$ | | $y_{(3)}$ | | $y_{(n)}$ | |
| **Drift** | + | - | + | - | + | - | + | - |
| **Production Rules** | If…. | If…. | If…. | If…. | If…. | If…. | If…. | If…. |
| | and… | and… | and… | and… | and… | and… | and… | and… |
| | and… | and… | and… | and… | and… | and… | and… | and… |
| | and… | and… | and… | and… | and… | and… | and… | and… |
| | and… | and… | and… | and… | and… | and… | and… | and… |
| | and… | and… | and… | and… | and… | and… | and… | and… |
| | and… | and… | and… | and… | and… | and… | and… | and… |
| | then…. | then…. | then…. | then…. | then…. | then…. | then…. | then…. |

The production rules belonging to each parameter need to be organized in the fault library for access by the system when required.

$P.R\ (y_{(i)})$

Separate and Collate

0 Drift

+ve Drift

-ve Drift

Fault Libraries

**Step 3: Diagnosis Search**

This is the final step for the fault diagnosis system before it can present results for the mitigation of a fault, and hence recovery of the process back to its normal operating state. Using the information obtained from the fault detection system, i.e. $y_{(i)}$ and (+,- drift), the diagnosis system searches for the best candidate out of the fault library. The initial step of the search methodology is to match the variable to from those available in fault library. When the variable is identified; the search method will display the result (production rules) from the available drift (+ve or -ve).



An integrated overview of the complete diagnosis section for the proposed diagnostic system is presented in Figure 5.3.

**Figure 5.3** An integrated view of the components in the diagnosis section

## 5.3 Software Platform

After the successful development of a new diagnostic algorithm, efforts were made to develop a software platform that can be used for the simulation (to be used in the detection of faults using mathematical models), and diagnosis of faults using production rules from the fault libraries. Java Development Kit (**JDK** version 6) and NetBeans IDE (Integrated Development Environment version 6.1) is used for the development of the application. A brief introduction to Java and NetBeans is given below before presenting the development of the Java application.

**Java**

Java is the most influential and widely used programming languages of the 21$^{st}$ century (TIBCO-Software-Inc, 2008). The language derives much of its syntax from languages C and C++ (once the most widely used languages and now second in application). Java applications can run on any Java Virtual Machine (JVM) regardless of computer architecture. Java is general-purpose, object-oriented and is specifically designed to have as few implementation dependencies as possible. Java is used for application software through to web applications. It was originally designed for use on digital mobile devices, such as cell phones. However, when Java 1.0 was released to the public in 1996, its main focus had shifted to use on the Internet. Since 1996 it has evolved as a successful language for use both Internet and other uses. A decade later, it is still an extremely popular language used by over 6.5million developers worldwide (Palmer, 2003).

The following reasons lead to the choice of Java rather than other languages in this work, relating to a few key principles from the original Java design (Leahy, 2010):

- *Easy to Use:* The fundamentals of Java came from a programming language called C++. Although a powerful language, it was felt to be too complex in its syntax, and inadequate for all of Java's requirements. Java is built using improved ideas, to provide a programming language that is powerful and simple to use.
- *Reliability:* Java needed to reduce the likelihood of fatal errors from programmer mistakes. With this in mind, object-oriented programming was introduced. Once data and its manipulation were packaged together in one place, it increased Java's robustness.
- *Security:* As Java was originally targeting mobile devices that would be exchanging data over networks, it was built to include a high level of security. Java is probably the most secure programming language to date.
- *Platform Independent:* Programs needed to work regardless of the machine upon which they were being executed. Java was written to be a portable language that does depend upon a particular operating system, or the computer hardware.

78

In addition to these points, Java is available under the GNU General Public License (GPL), thus making it free software.

**NetBeans**

The NetBeans IDE is an open-source integrated development environment. NetBeans IDE supports development of all Java application types. It provides the "plumbing" for the Graphical User Interface (GUI) in any Java application that conventionally every developer had to write themselves otherwise. NetBeans provides these entire straight "out of the box" thus saving a developer a significant amount of time and work (NetBeans, 2010).

NetBeans is available as open source free software. It can be run on most operating systems including Windows, Linux, Mac OS X and Solaris (Fears, 2008).

## 5.4 Architecture of the Application

The GUI of the application "WWTP-Simulation and Diagnosis" is presented in Figure 5.4. The application consists of the following five screens and a graphical applet.

1. U.D. Inputs
2. Flowsheet
3. Results
4. Online Data
5. Faults



**Figure 5.4** The Graphical User Interface (GUI) of the application

Figure 5.4 shows the main screen of the "WWTP-Simulation and Diagnostic" application. The first screen to be discussed is the "U.D. Inputs (used defined inputs)".

The inputs to the WWTP are declared in the field marked as "1" in figure 5.4.

Field "2" points to the process parameters used in the different stages of the WWTP. Most of the process parameters are discussed in Chapter 3.

This application can be used for the simulation of different stages of the wastewater in a sequential and cumulative manner by selecting the number of stages (section "3" on figure 5.4). A point of emphasis here is that if Stage 2 is selected, then this will give the simulation of stages 1 and 2; whereas selecting stage 3 will simulate stage 1, 2 and 3.



**Figure 5.5** Flowsheet screen of the WWTP application

Equipment selection, simulation command, and selection of an individual stage is developed in the "Flowsheet" tab of WWTP application. In figure 5.5, "3" indicates the simulation command button of the application that will initiate the simulation of required stages (as discussed above under figure 5.4) using the inputs on the "U.D. Inputs" screen and a file "WWTP.java" that includes the model

equations. The file "WWTP.java" is solved using the Runge-Kutta 4 (RK4) method (also coded in Java).

The next screen next to Flowsheet is the "Results" of this application. The numeric results obtained from the simulation of the model are displayed in the screen on the execution of the software. The "Results" screen is used mainly to display the results from the simulation. Although graphical presentation is used to display the simulation results, it was also considered useful to display the results in numeric form in the WWTP application. The reasons were the use of numeric data for the detection of faults, and the assumption that the simulated data may be used in future research work. For example, the data may be required for the validation of some experimental results.

The detection of faults occurred in the screen labeled as "Online Data".



**Figure 5.6** Detection and Diagnosis Algorithm command window

The check boxes (No. "1" in figure 5.6) were designed so that that once selected, it will import the online data when coupled to the control system of a wastewater treatment facility. However, it was not possible to obtain technical support from any local wastewater treatment facility (preventing activation and use of this screen). For the WWTP application, the data for the detection of a fault is entered manually into the file. After the acquisition of data (in this application,

manual entry of data in the file), the detection and diagnosis of a fault can be initiated using the command box labeled "Initialize FDD" (labeled as "2" in figure 5.6).

The diagnostic results are the displayed in the last tab labeled as "Faults". A screenshot of this screen is given in figure 5.7 below.



**Figure 5.7** Screen for the display of the diagnostic results

The variable suspected of faulty behavior (for example, SS, VSS, etc.) is displayed in the field labeled as "1" in figure 5.7, whereas the drift (i.e. high or low) is displayed in the field below labeled as "2". The production rules (developed in Chapter 4) are displayed in the field labeled "3" with the title "Diagnosis".

## 5.5   Results and Discussion

The diagnostic algorithm proposed earlier in this chapter, and the WWTP application developed, along with the WWTP dataset were used to generate results. From the WWTP dataset, input values were selected and used in the WWTP Java application, the known output values (from WWTP dataset) were fed to the software for detection purposes. The production rules for the diagnostic suggestion were used for the development of the fault libraries as the knowledge base of the FDD platform. When

the suspended solids reduction was inhibited due to the poor production of microorganisms, the algorithm was able to detect and eventually provide suggestions for the diagnostics. Figure 5.8 below shows the inputs used for the fault scenario when the suspended solids concentration in effluent eventually increased. Figures 5.9 and 5.10 present the graphical results from the simulation and the diagnostic results after the detection of fault, respectively.



**Figure 5.8** Input of process parameters for simulation



**Figure 5.9** Graphical presentations of the simulation results

**Figure 5.10** Diagnostic results for high SS value

Similarly when a low concentration of suspended solids in effluent was detected (although in practice, the low concentration of suspended solids is not a fault, but is considered as a fault in this study), the WWTP software was able to detect and diagnose this fault.

The graphical presentation of the results and the diagnostic suggestions from the fault scenario discussed above are presented in figures 5.11 and 5.12 respectively.

**Figure 5.11** Simulation indicating low SS value



**Figure 5.12** Set of diagnostic results for low value of SS

Further scenarios using other output variables were studied and results are reported in Table 5.2.

**Table 5.2** Results of the diagnostic methodology

|  | Variable | Direction | Detection | Diagnosis | Time(Sec) |
|---|---|---|---|---|---|
| **Single Fault** | SS | + | ✓ | ✓ | 1.5 |
|  | SS | - | ✓ | ✓ | 1.5 |
|  | SSV | + | ✓ | ✓ | 1.5 |
|  | SSV | - | ✓ | ✓ | 1.5 |
| **Disturbance** | SS | + | ✓ | ✓ | 1 |
|  | SS | - | ✓ | ✓ | 1 |
|  | SSV | + | ✓ | ✓ | 1 |
|  | SSV | - | ✓ | ✓ | 1 |
| **Multi Fault** | SS & SSV | + & + | ✓ | No | 2 |
|  | SS & SSV | - & - | ✓ | No | 2 |
|  | SS & SSV | + & - | ✓ | No | 2 |
|  | SS & SSV | - & + | ✓ | No | 2 |

The results show that the most troublesome scenario for the proposed methodology is if two faults occur simultaneously. It is determined that there is a need to expand the knowledge base obtained from data mining in order to accommodate multiple faults. Recommendations to enhance this methodology are discussed in Chapter 6.

## 5.6   Concluding Remarks

An algorithm based on a new approach is presented in this chapter, and the results obtained confirmed its suitability. A software application is also developed to implement the algorithm and yielded reliable simulation results. Although there remain some challenges to be overcome with the proposed software application (see Chapter 6: Conclusions and Future Work) and also with the proposed algorithm (also

see Chapter 6), this study serves its objective to develop an initial platform for future researchers intending to use different techniques to perform different roles in process diagnostics.

# Chapter 6
# Conclusions and Recommendations

## 6.1 Conclusions

An algorithm is proposed which has been applied to a data set from a wastewater treatment plant. A mathematical model was used for detection, whereas the production rules were employed for the diagnosis, of process faults. The accuracy of detection of faults using the mathematical model was calculated to be 93% (Table No. 3.5) whereas, the production rules exhibited 95% accuracy on the validation step (Table No. 4.10). It was not possible to compare the response time for the detection and diagnosis of faults using this technique with other systems because that data is not available in the literature. A Java application was designed for the implementation of the algorithm and it performed successfully, although some aspects still require improvement.

A new technique of ranking the production rules is proposed and was applied in this research to reduce the large number of production rules and improve the efficiency of the predictions. In this study, the production rules are used to build the knowledge base of the algorithm, and are subsequently used in the diagnosis. The production rules were obtained by using a data-set of WWTP, and an inductive data mining technique using See 5.0 was applied to yield the production rules.

A mathematical model was derived, mainly based upon the ASM1 model, to be used for the detection of faults. The model was validated against the WWTP data-set. An accuracy of 93% on the validation was considered acceptable in this work, but it is expected that calibration of the proposed model would improve the effectiveness of this technique.

One major issue encountered in this research work was the lack of technical input and data sets from any local wastewater treatment facilities. Therefore, a data-set from 1991 was used, which mainly covers the dynamics of suspended solids removal from an activated sludge system.

Despite the lack of extensive technical data, the average response time of 1.5 sec for the proposed software (for the detection of fault and/or disturbance, search of the best possible diagnostic strategy from the fault libraries, and the display of results on screen), 93% model accuracy, 95% production rules accuracy demonstrate the validity and represent the key achievements from this work.

## 6.2   Recommendations for Future Work

The following areas are identified for continuation of this research work:

a. Modeling errors are a common problem especially when a model is to be used for process fault detection and diagnosis. It is suggested that if the mathematical model and operational data for a process is available, then genetic algorithm or programming (a technique well established and used for optimization) can be used to fit that model to the operational data thus improving the output (predictions) obtained from the model.

b. It is expected that the "WWTP-Simulation and Diagnostic" has the potential to be used in the chemical and processing industries. Thus it is recommended to apply this technique on a more complex process where: (a) an elaborated and exhaustive mathematical model of the process is available; and (b) where abundant historical and process data covering almost all aspects and variables of the process are available, in order to test the robustness of this technique.

c. Detection of multiple faults in a process is a challenging task. It is proposed that if sufficient operational or historical data containing multiple faults scenarios is available, then data mining should be employed for determination of the root causes and the process behavior when multiple faults are present.

# References

Acharya, J., Sahu, J. N., Mohanty, C. R. & Meikap, B. C. (2009a) Removal of lead(II) from wastewater by activated carbon developed from Tamarind wood by zinc chloride activation. *Chemical Engineering Journal,* vol. 149**,** no. 1-3, pp. 249-262.

Acharya, J., Sahu, J. N., Sahoo, B. K., Mohanty, C. R. & Meikap, B. C. (2009b) Removal of chromium(VI) from wastewater by activated carbon developed from Tamarind wood activated with zinc chloride. *Chemical Engineering Journal,* vol. 150**,** no. 1, pp. 25-39.

Aguado, D., Ferrer, A., Ferrer, J. & Seco, A. (2007) Multivariate SPC of a sequencing batch reactor for wastewater treatment. *Chemometrics and Intelligent Laboratory Systems,* vol. 85**,** no. 1, pp. 82-93.

Albazzaz, H., Wang, X. Z. & Marhoon, F. (2005) Multidimensional visualisation for process historical data analysis: a comparative study with multivariate statistical process control. *Journal of Process Control,* vol. 15**,** no., pp. 285-294.

Buchanan, J. R. & Seabloom, R. W. (2004) Aerobic Treatment of Wastewater and Aerobic Treatment Units. The University of Washington.

Buontempo, F. V. (2005) *Rapid toxicity prediction of organic chemicals using data mining techniques and QSAR based on genetic programming for decision tree generation* Department of Chemical Engineering, University of Leeds, Leeds, PhD.

Charaniya, S., Hu, W. S. & Karypis, G. (2008) Mining bioprocess data: opportunities and challenges. *Trends in Biotechnology,* vol. 26**,** no. 12, pp. 690-699.

Chen, C. Z. & Mo, C. T. (2004) A method for intelligent fault diagnosis of rotating machinery. *Digital Signal Processing,* vol. 14**,** no. 3, pp. 203-217.

Chiang, L. H., Russell, E. & Braatz, R. D. (2001) *Fault detection and diagnosis in industrial systems,* Springer, London.

Copp, J. B., Belia, E., Snowling, S. & Schraa, O. (2005) Anaerobic digestion: a new model for plant-wide wastewater treatment process modelling. *Water Science and Technology,* vol. 52**,** no. 10-11, pp. 1-11.

Dash, S. & Venkatasubramanian, V. (2000) Challanges in the industrial applications of fault diagnostic systems. *Computers and Chemical Engineering,* vol. 24**,** no. 2-7, pp. 785-791.

Degremont, G. (1991) *Water treatment handbook, Lavoisier Publishing,* Paris.

DeLisle, R. K. & Dixon, S. L. (2004) Induction of decision trees via evolutionary programming. *Journal of Chemical Information and Computer Sciences,* vol. 44**,** no. 3, pp. 862-870.

Dellana, S. A. & West, D. (2009) Predictive modeling for wastewater applications: Linear and nonlinear approaches. *Environmental Modelling & Software,* vol. 24**,** no. 1, pp. 96-106.

Fayyed, U. M., Shapiro, G., Smith, P. & Uthursuamy, R. (1996) *Advances in Knowledge Discovery and Data Mining,* MIT Press.

Fears, N. (2008) What Is Netbeans and What Can It Do For You? http://www.brighthub.com/internet/web- development/articles/8472.aspx#ixzz0fHVxjHBo

Forster, C. F. (2003) *Wastewater treatment and technology,* Thomas Telford, London.

Gernaey, K. V., van Loosdrecht, M. C. M., Henze, M., Lind, M. & Jorgensen, S. B. (2004) Activated sludge wastewater treatment plant modelling and simulation: state of the art. *Environmental Modelling & Software,* vol. 19**,** no. 9, pp. 763-783.

Green-Ideas (2009) Green Ideas, Inc. http://www.egreenideas.com/glossary.php?group=b

Gujer, W., Henze, M., Mino, T. & van Loosdrecht, M. (1999) Activated Sludge Model No. 3. *Water Science and Technology,* vol. 39**,** no. 1, pp. 183-193.

Han, J. W., Pei, J., Yin, Y. W. & Mao, R. Y. (2004) Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data Mining and Knowledge Discovery,* vol. 8**,** no. 1, pp. 53-87.

Himmelblau, D. M. (1978) *Fault detection and diagnosis in chemical and petrochemical processes,* Elsevier Scientific Publishing , Amsterdam.

Henze, M., Gujer, W., Mino, T. & van Loosdrecht, M. (2000) *Activated sludge models : ASM1, ASM2, ASM2d and ASM3,* IWA Publishing, London.

Huang, M. Z., Wan, J. Q., Ma, Y. W., Wang, Y., Li, W. J. & Sun, X. F. (2009) Control rules of aeration in a submerged biofilm wastewater treatment process using fuzzy neural networks. *Expert Systems with Applications,* vol. 36**,** no. 7, pp. 10428-10437.

Hunag, Y. C. & Wang, X. Z. (1999) Application of fuzzy casual networks to wastewater treatment plants. *Chemical Engineering Science,* vol. 54**,** pp. 2731-2738.

Iri, M., Aoki, K., O'Shima & Matsuyama, H. (1979) An algorithm for diagnosis of system failures in the chemical processes. *Computers and Chemical Engineering,* vol. 3**,** no. 1-4, pp. 489-493.

Jain, A. K., Murty, M. N. & Flynn, P. J. (1999) Data clustering: A review. *Association for Computing Machinery Computing Surveys,* vol. 31**,** no. 3, pp. 264-323.

Jeppsson, U. (2003) A General Description of the IAWQ Activated Sludge Model No. 1. Lund Institute of Technology, Lund, Sweden.

Khalid, M. I. (2005) *Automatic Generation of Decision Trees from Process Historical Data using Genetic Programming* Institute of Particle Science and Engineering, The University of Leeds, Leeds, MSc (Eng).

Kim, S., Lee, H., Kim, J., Kim, C., Ko, J., Woo, H. & Kim, S. (2002) Genetic algorithms for the application of Activated Sludge Model No. 1. *Water Science and Technology,* vol. 45**,** no. 4-5, pp. 405-411.

Krogh, A. (2008) What are artificial neural networks? *Natural Biotechnology,* vol. 26**,** no. 2, pp. 195-197.

Kumar, A. (2004) State-of-the-art report on Water Management, Documentation and Research training centre, Indian Statistical Institute, Bangalore.

Lapp, S. A. & Powers, G. A. (1977) Computer-aided synthesis of fault trees. *IEEE Transactions and Reliability,* vol. 37**,** no., pp. 2-13.

Leahy, P. (2010) What is Java? http://java.about.com/od/gettingstarted/a/whatisjava.htm

Lee, C., Choi, S. W. & Lee, I. B. (2006) Sensor fault diagnosis in a wastewater treatment process. *Water Science and Technology,* vol. 53**,** no. 1, pp. 251-257.

Lee, D. S., Vanrolleghem, P. A. & Park, J. M. (2005) Parallel hybrid modeling methods for a full-scale cokes wastewater treatment plant. *Journal of Biotechnology,* vol. 115**,** no. 3, pp. 317-328.

Li, R. F. (2003) *Advanced process monitoring & control using principal & independent component analysis* Department of Chemical Engineeirng, University of Leeds, Leeds, PhD.

Lou, X. C., Willsky, A. S. & Verghese, G. C. (1986) Optimally robust redundancy relations for failure detection in uncertain systems. *Automatica,* vol. 22**,** no. 3, pp. 333-344.

Ma, C. Y. & Wang, X. Z. (2009) Inductive data mining based on genetic programming: Automatic generation of decision trees from data for process historical data analysis. *Computers & Chemical Engineering,* vol. 33**,** no. 10, pp. 1602-1616.

Michle, B. (1988) Detecting changes in signals and systems\&mdash;a survey. *Automatica,* vol. 24**,** no. 3, pp. 309-326.

Morley, D. A. (1979) *Mathematical modelling in water and wastewater treatment,* Applied Science, London.

NetBeans (2010) NetBeans IDE - Get Tomorrow Today http://netbeans.org/features/index.html

Norweco (2006) Disinfection of Water and Wastewater
http://www.norweco.com/html/lab/Disinfection.htm

Pai, T. Y. (2008) Gray and neural network prediction of effluent from the wastewater treatment plant of industrial park using influent quality. *Environmental Engineering Science,* vol. 25**,** no. 5, pp. 757-766.

Palmer, G. (2003) *Technical Java : developing scientific and engineering applications,* Prentice Hall, Upper Saddle River, NJ.

Puteh, M., Minekawa, K., Hashimoto, N. & Kawase, Y. (1999) Modeling of activated sludge wastewater treatment processes. *Bioprocess Engineering,* vol. 21**,** no. 3, pp. 249-254.

Quinlan, J. R. (1986) Induction of decision trees. *Machine Learning,* vol. 1**,** pp. 81-106.

Quinlan, J. R. (1990) Decision trees and decision-making. *IEEE transactions on systems, man, and cybernetics,* vol. 20**,** no. 2, pp. 339.

Quinlan, J. R. (1993) *C4.5 : programs for machine learning,* Morgan Kaufmann, Calif.

Quinlan, J. R. (1996) Improved use of continuous attributes in C4.5. *Journal of Artificial Intelligence Research,* vol. 4**,** no., pp. 77-90.

R-Foundation (2008) The R Project for Statistical Computing http://www.r-project.org/

Radcliffe, J. C. (2004) Water Recycling in Australia, The Australian Academy of Technological Sciences and Engineering, Melbourne.

Rieger, L., Alex, J., Gujer, W. & Siegrist, H. (2006) Modelling of aeration systems at wastewater treatment plants. *Water Science and Technology,* vol. 53**,** no. 4-5, pp. 439-447.

Sanchez , M., Cortes, U., Bejar, J., Gracia, J. D., Lafuente , J. & Poch, M. (1997) Concept formation in WWTP by means of classification techniques:A compared study. *Applied Intelligence,* vol. 7**,** no., pp. 147-165.

Sanguesa, R. & Cortes, U. (1997) Learning causal networks from data:a survey and a new algorithm for recovering possibilistic causal networks. *A.I. Communications.*

See5.0 (2008) See 5.0. 2.06 ed. RuleQuest Research Pty Ltd, NSW, Australia.

Shuler, M. L. & Kargi, F. (2002) *Bioprocess engineering,* Prentice Hall, NJ.

Simani, S., Fantuzzi, C. & Patten, R. J. (2003) *Model based fault diagnosis in dynamic systems using identification techniques,* Springer, NY.

Sun, W. X., Chen, J. & Li, J. Q. (2007) Decision tree and PCA-based fault diagnosis of rotating machinery. *Mechanical Systems and Signal Processing,* vol. 21**,** no. 3, pp. 1300-1317.

Tchobanoglous, G., Burton, F. L., Stensel, H. D. & Metcalf & Eddy. (2003) *Wastewater engineering : treatment and reuse,* McGraw-Hill, Boston.

TIBCO-Software-Inc (2008) TIBCO:The Power of Now http://www.insightful.com/

U.S-EPA (2006) Lake Erie http://www.epa.gov/lakeerie/glossary.html

Van Hulle, S. W. H. & Vanrolleghem, P. A. (2004) Modelling and optimisation of a chemical industry wastewater treatment plant subjected to varying production schedules. *Journal of Chemical Technology and Biotechnology,* vol. 79**,** no. 10, pp. 1084-1091.

Venkatasubramanian, V., Rengaswamy, R. & Kavuri, S. N. (2003a) A review of process fault detection and diagnosis: Part I: Quantitative model based methods. *Computers and Chemical Engineering,* vol. 27**,** no. 3, pp. 293-311.

Venkatasubramanian, V., Rengaswamy, R. & Kavuri, S. N. (2003b) A review of process fault detection and diagnosis: Part II: Qualitative models and search strategies. *Computers and Chemical Engineering,* vol. 27**,** no. 3, pp. 313-326.

Venkatasubramanian, V., Rengaswamy, R. & Kavuri, S. N. (2003c) A review of process fault detection and diagnosis: Part III: Process history based methods. *Computers and Chemical Engineering,* vol. 27**,** no. 3, pp. 327-346.

Wang, X. Z. (1999) *Data mining and knowledge discovery for process monitorning and control,* Springer-Verlag, London.

Wang, X. Z., Medasani, S., Marhoon, F. & Albazzaz, H. (2004) Multidimensional visualization of principal component scores for process historical data analysis. *Industrial & Engineering Chemistry Research,* vol. 43**,** no., pp. 7036-7048.

West, D. & Mangiameli, P. (2000) Identifying process conditions in an urban wastewater treatment plant. *International Journal of Operations & Production Management,* vol. 20**,** no. 5-6, pp. 573-590.

Wimberger, D. & Verde, C. (2008) Fault diagnosticability for an aerobic batch wastewater treatment process. *Control Engineering Practice,* vol. 16**,** no. 11, pp. 1344-1353.

Wintgens, T., Rosen, J., Melin, T., Brepols, C., Drensla, K. & Engelhardt, N. (2003) Modelling of a membrane bioreactor system for municipal wastewater treatment. *Journal of Membrane Science,* vol. 216**,** no. 1-2, pp. 55-65.

Yang, B. S., Lim, D. S. & Tan, A. C. C. (2005) VIBEX: an expert system for vibration fault diagnosis of rotating machinery using decision tree and decision table. *Expert Systems with Applications,* vol. 28**,** no. 4, pp. 735-742.

# Appendix

## Data Statistics:-

| No. | Attrib. | Min | Max | Mean | St-dev |
|-----|---------|------|------|----------|---------|
| 1 | Q-E | 10000 | 60081 | 37226.56 | 6571.46 |
| 2 | ZN-E | 0.1 | 33.5 | 2.36 | 2.74 |
| 3 | PH-E | 6.9 | 8.7 | 7.81 | 0.24 |
| 4 | DBO-E | 31 | 438 | 188.71 | 60.69 |
| 5 | DQO-E | 81 | 941 | 406.89 | 119.67 |
| 6 | SS-E | 98 | 2008 | 227.44 | 135.81 |
| 7 | SSV-E | 13.2 | 85.0 | 61.39 | 12.28 |
| 8 | SED-E | 0.4 | 36 | 4.59 | 2.67 |
| 9 | COND-E | 651 | 3230 | 1478.62 | 394.89 |
| 10 | PH-P | 7.3 | 8.5 | 7.83 | 0.22 |
| 11 | DBO-P | 32 | 517 | 206.20 | 71.92 |
| 12 | SS-P | 104 | 1692 | 253.95 | 147.45 |
| 13 | SSV-P | 7.1 | 93.5 | 60.37 | 12.26 |
| 14 | SED-P | 1.0 | 46.0 | 5.03 | 3.27 |
| 15 | COND-P | 646 | 3170 | 1496.03 | 402.58 |
| 16 | PH-D | 7.1 | 8.4 | 7.81 | 0.19 |
| 17 | DBO-D | 26 | 285 | 122.34 | 36.02 |
| 18 | DQO-D | 80 | 511 | 274.04 | 73.48 |
| 19 | SS-D | 49 | 244 | 94.22 | 23.94 |
| 20 | SSV-D | 20.2 | 100 | 72.96 | 10.34 |
| 21 | SED-D | 0.0 | 3.5 | 0.41 | 0.37 |
| 22 | COND-D | 85 | 3690 | 1490.56 | 399.99 |
| 23 | PH-S | 7.0 | 9.7 | 7.70 | 0.18 |
| 24 | DBO-S | 3 | 320 | 19.98 | 17.20 |
| 25 | DQO-S | 9 | 350 | 87.29 | 38.35 |
| 26 | SS-S | 6 | 238 | 22.23 | 16.25 |
| 27 | SSV-S | 29.2 | 100 | 80.15 | 9.00 |
| 28 | SED-S | 0.0 | 3.5 | 0.03 | 0.19 |
| 29 | COND-S | 683 | 3950 | 1494.81 | 387.53 |
| 30 | RD-DBO-P | 0.6 | 79.1 | 39.08 | 13.89 |
| 31 | RD-SS-P | 5.3 | 96.1 | 58.51 | 12.75 |
| 32 | RD-SED-P | 7.7 | 100 | 90.55 | 8.71 |
| 33 | RD-DBO-S | 8.2 | 94.7 | 83.44 | 8.4 |
| 34 | RD-DQO-S | 1.4 | 96.8 | 67.67 | 11.61 |
| 35 | RD-DBO-G | 19.6 | 97 | 89.01 | 6.78 |
| 36 | RD-DQO-G | 19.2 | 98.1 | 77.85 | 8.67 |
| 37 | RD-SS-G | 10.3 | 99.4 | 88.96 | 8.15 |
| 38 | RD-SED-G | 36.4 | 100 | 99.08 | 4.32 |

## WWTP Data:-

The WWTP data used in this research work is included here.

| Date | Q-E | ZN-E | PH-E | DBO-E | DBQ-E | SS-E | SSV-E | SED-E | COND-E | PH-P | DBO-P | SS-P | SSV-P | SED-P | COND-P | PH-D | DBO-D | DQO-D | SS-D | SSV-D | SED-D | COND-D | PH-S | DBO-SS | DQO-S | SS-S | SSV-S | SED-S | COND-S | RD-DBO-P | RD-SS-P | RD-SED-P | RD-DBO-S | RD-DQO-S | RD-DBO-G | RD-DQO-G | RD-SS-G | RD-SED-G |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-1/3/90 | 44101 | 1.5 | 7.8 | 183 | 407 | 166 | 66 | 4.5 | 2110 | 7.9 | 197 | 228 | 70 | 5.5 | 2120 | 7.9 | 119 | 280 | 94 | 72 | 0.3 | 2010 | 7.3 | 18 | 84 | 21 | 81 | 0 | 2000 | 40 | 59 | 96 | 85 | 70 | 90 | 79 | 87 | 100 |
| D-2/3/90 | 39024 | 3 | 7.7 | 183 | 443 | 214 | 69 | 6.5 | 2660 | 7.7 | 197 | 244 | 75 | 7.7 | 2570 | 7.6 | 119 | 474 | 96 | 79 | 0.4 | 2700 | 7.5 | 18 | 91 | 17 | 94 | 0 | 2590 | 40 | 61 | 95 | 85 | 81 | 90 | 80 | 92 | 100 |
| D-4/3/90 | 32229 | 5 | 7.6 | 183 | 528 | 186 | 70 | 3.4 | 1666 | 7.7 | 197 | 220 | 73 | 4.5 | 1594 | 7.7 | 119 | 272 | 92 | 78 | 0.2 | 1742 | 7.6 | 18 | 128 | 21 | 81 | 0.1 | 1888 | 40 | 58 | 96 | 85 | 53 | 90 | 76 | 89 | 99 |
| D-5/3/90 | 35023 | 3.5 | 7.9 | 205 | 588 | 192 | 66 | 4.5 | 2430 | 7.8 | 236 | 268 | 73 | 8.5 | 2280 | 7.8 | 158 | 376 | 96 | 77 | 0.4 | 2060 | 7.6 | 20 | 104 | 20 | 97 | 0 | 1840 | 33 | 64 | 95 | 87 | 72 | 90 | 82 | 90 | 100 |
| D-6/3/90 | 36924 | 1.5 | 8 | 242 | 496 | 176 | 65 | 4 | 2110 | 7.9 | 197 | 236 | 58 | 4.5 | 2020 | 7.8 | 119 | 372 | 88 | 68 | 0.2 | 2250 | 7.6 | 19 | 108 | 22 | 66 | 0 | 2120 | 40 | 63 | 96 | 85 | 71 | 92 | 78 | 88 | 100 |
| D-7/3/90 | 38572 | 3 | 7.8 | 202 | 372 | 186 | 69 | 4.5 | 1644 | 7.8 | 197 | 248 | 66 | 8.5 | 1762 | 7.7 | 150 | 460 | 100 | 76 | 0.3 | 1768 | 7.5 | 20 | 100 | 28 | 82 | 0 | 1764 | 40 | 60 | 97 | 87 | 78 | 90 | 73 | 85 | 100 |
| D-8/3/90 | 41115 | 6 | 7.8 | 183 | 552 | 262 | 64 | 5 | 1603 | 7.8 | 197 | 320 | 68 | 6.5 | 1608 | 7.8 | 192 | 376 | 122 | 72 | 0.4 | 1668 | 7.5 | 21 | 76 | 26 | 85 | 0.1 | 1703 | 40 | 62 | 94 | 89 | 80 | 90 | 86 | 90 | 99 |
| D-9/3/90 | 36107 | 5 | 7.7 | 215 | 489 | 334 | 41 | 6 | 1613 | 7.6 | 197 | 304 | 54 | 8 | 1557 | 7.6 | 181 | 350 | 90 | 71 | 0.4 | 1596 | 7.5 | 17 | 162 | 18 | 67 | 0 | 1606 | 40 | 70 | 96 | 91 | 54 | 92 | 67 | 95 | 100 |
| D-11/3/90 | 29156 | 2.5 | 7.7 | 206 | 451 | 194 | 69 | 4.5 | 1249 | 7.7 | 206 | 220 | 62 | 4 | 1219 | 7.7 | 111 | 282 | 124 | 77 | 0.3 | 1233 | 7.5 | 16 | 118 | 19 | 84 | 0 | 1338 | 46 | 44 | 93 | 86 | 58 | 92 | 74 | 90 | 99 |
| D-12/3/90 | 39246 | 2 | 7.8 | 172 | 506 | 200 | 69 | 5 | 1865 | 7.8 | 208 | 248 | 66 | 6.5 | 1929 | 7.8 | 164 | 463 | 100 | 78 | 0.6 | 1825 | 7.6 | 19 | 157 | 27 | 87 | 0 | 1616 | 21 | 60 | 91 | 88 | 66 | 89 | 69 | 87 | 100 |
| D-13/3/90 | 42393 | 0.7 | 7.9 | 189 | 478 | 230 | 67 | 5.5 | 1410 | 8.1 | 173 | 192 | 63 | 5 | 1406 | 7.7 | 172 | 412 | 104 | 71 | 0.4 | 1562 | 7.6 | 152 | 306 | 131 | 80 | 3.5 | 1575 | 0.6 | 46 | 92 | 12 | 26 | 20 | 36 | 43 | 36 |
| D-14/3/90 | 42857 | 1.5 | 7.7 | 238 | 319 | 292 | 34 | 3.5 | 1261 | 7.6 | 170 | 268 | 31 | 4.2 | 1204 | 7.6 | 116 | 276 | 104 | 52 | 0.3 | 1261 | 7.4 | 320 | 350 | 238 | 74 | 2 | 1304 | 32 | 61 | 93 | 85 | 70 | 90 | 79 | 19 | 43 |
| D-15/3/90 | 42911 | 0.7 | 7.6 | 114 | 252 | 116 | 59 | 1.2 | 1238 | 7.9 | 148 | 136 | 65 | 3 | 1208 | 7.7 | 79 | 216 | 70 | 83 | 0.3 | 1177 | 7.5 | 84 | 172 | 104 | 79 | 0.1 | 1221 | 47 | 49 | 92 | 85 | 20 | 26 | 32 | 10 | 95 |
| D-16/3/90 | 40376 | 1.5 | 8.1 | 204 | 333 | 174 | 68 | 3 | 2390 | 7.8 | 231 | 156 | 74 | 2.5 | 2540 | 7.8 | 136 | 325 | 78 | 80 | 0.4 | 2580 | 7.6 | 32 | 153 | 98 | 88 | 0 | 2550 | 41 | 50 | 84 | 77 | 53 | 84 | 54 | 44 | 100 |
| D-18/3/90 | 40923 | 3.5 | 7.6 | 146 | 329 | 188 | 57 | 2.5 | 1300 | 7.6 | 162 | 132 | 64 | 2 | 1324 | 7.6 | 109 | 243 | 88 | 82 | 0.2 | 1467 | 7.5 | 19 | 94 | 41 | 83 | 0 | 1545 | 33 | 33 | 90 | 83 | 61 | 87 | 71 | 78 | 99 |
| D-19/3/90 | 43830 | 1.5 | 7.8 | 177 | 512 | 214 | 59 | 5.5 | 1605 | 7.7 | 164 | 256 | 72 | 5.5 | 1599 | 7.7 | 118 | 320 | 70 | 89 | 0.4 | 1401 | 7.6 | 25 | 203 | 20 | 85 | 0 | 1110 | 28 | 73 | 93 | 79 | 37 | 86 | 60 | 91 | 100 |
| D-20/3/90 | 39165 | 1.2 | 7.4 | 250 | 447 | 252 | 61 | 7 | 1533 | 7.4 | 275 | 216 | 57 | 6.5 | 1501 | 7.4 | 138 | 269 | 90 | 73 | 0.5 | 1458 | 7.3 | 14 | 9 | 20 | 83 | 0 | 1402 | 50 | 58 | 92 | 90 | 97 | 94 | 98 | 92 | 100 |
| D-21/3/90 | 35791 | 1.2 | 7.8 | 277 | 466 | 246 | 63 | 4 | 1556 | 7.7 | 197 | 288 | 65 | 6 | 1846 | 7.7 | 166 | 419 | 174 | 81 | 1.3 | 1664 | 7.5 | 24 | 124 | 26 | 83 | 0 | 1606 | 40 | 40 | 78 | 86 | 70 | 91 | 73 | 89 | 99 |
| D-22/3/90 | 37419 | 1.2 | 7.6 | 219 | 446 | 222 | 61 | 5.5 | 1600 | 7.7 | 266 | 240 | 70 | 5 | 1645 | 7.6 | 172 | 345 | 102 | 84 | 0.4 | 1670 | 7.5 | 42 | 175 | 53 | 84 | 0 | 1780 | 35 | 58 | 92 | 76 | 49 | 81 | 61 | 76 | 100 |
| D-23/3/90 | 40983 | 3 | 7.6 | 182 | 431 | 214 | 57 | 7 | 1591 | 7.5 | 219 | 248 | 58 | 5.5 | 1473 | 7.5 | 175 | 376 | 88 | 66 | 0.4 | 1537 | 7.5 | 23 | 120 | 25 | 68 | 0 | 1597 | 20 | 65 | 94 | 87 | 68 | 87 | 72 | 88 | 100 |
| D-25/3/90 | 42217 | 8.5 | 7.5 | 138 | 333 | 240 | 55 | 3.8 | 1087 | 7.5 | 153 | 184 | 67 | 4 | 1109 | 7.5 | 108 | 194 | 82 | 85 | 0.4 | 1136 | 7.1 | 16 | 62 | 17 | 94 | 0 | 1223 | 29 | 55 | 91 | 85 | 68 | 88 | 81 | 93 | 100 |
| D-26/3/90 | 47665 | 1.2 | 7.7 | 156 | 405 | 200 | 74 | 4 | 1856 | 7.6 | 178 | 184 | 72 | 3.5 | 1976 | 7.5 | 126 | 302 | 92 | 78 | 0.3 | 1920 | 7.6 | 19 | 71 | 23 | 78 | 0 | 1706 | 28 | 50 | 91 | 85 | 77 | 88 | 83 | 89 | 100 |
| D-27/3/90 | 44314 | 3 | 7.8 | 155 | 389 | 308 | 49 | 6 | 1927 | 7.7 | 252 | 308 | 49 | 6.5 | 2150 | 7.7 | 121 | 302 | 108 | 72 | 0.6 | 1950 | 7.6 | 15 | 87 | 23 | 70 | 0 | 1869 | 52 | 65 | 91 | 88 | 71 | 90 | 78 | 93 | 100 |
| D-28/3/90 | 40841 | 1 | 7.6 | 179 | 389 | 168 | 69 | 3.5 | 1240 | 7.8 | 202 | 272 | 72 | 6 | 1381 | 7.8 | 148 | 302 | 92 | 78 | 0.3 | 1425 | 7.9 | 16 | 83 | 20 | 85 | 0 | 1416 | 27 | 66 | 95 | 89 | 73 | 91 | 79 | 88 | 100 |
| D-29/3/90 | 41157 | 3 | 8 | 145 | 398 | 192 | 67 | 4.5 | 2240 | 8 | 213 | 240 | 62 | 6 | 2010 | 8 | 140 | 287 | 84 | 79 | 0.4 | 2270 | 7.8 | 15 | 87 | 21 | 81 | 0 | 2290 | 34 | 65 | 94 | 89 | 70 | 90 | 78 | 89 | 100 |
| D-30/3/90 | 40078 | 1.4 | 7.9 | 198 | 464 | 228 | 65 | 4.6 | 1431 | 7.6 | 243 | 272 | 65 | 7.5 | 1606 | 7.8 | 177 | 319 | 88 | 82 | 0.2 | 1556 | 7.8 | 17 | 102 | 22 | 82 | 0 | 1475 | 27 | 68 | 97 | 90 | 68 | 91 | 78 | 90 | 100 |
| D-1/2/90 | 44365 | 7.5 | 7.9 | 183 | 365 | 212 | 62 | 3.5 | 1339 | 7.9 | 197 | 184 | 65 | 4.7 | 1380 | 7.8 | 119 | 321 | 92 | 74 | 0.5 | 1386 | 7.5 | 18 | 75 | 20 | 75 | 0.1 | 1377 | 40 | 50 | 89 | 85 | 77 | 90 | 80 | 91 | 99 |
| D-2/2/90 | 43080 | 4.3 | 7.8 | 95 | 349 | 136 | 77 | 2.5 | 1063 | 7.8 | 132 | 188 | 75 | 2 | 1139 | 7.8 | 123 | 317 | 98 | 69 | 0.4 | 1218 | 7.5 | 19 | 67 | 24 | 83 | 0 | 1220 | 6.8 | 48 | 80 | 85 | 79 | 80 | 81 | 82 | 100 |
| D-4/2/90 | 29414 | 3 | 7.6 | 160 | 374 | 168 | 69 | 3.1 | 1042 | 7.6 | 220 | 246 | 70 | 4.6 | 1057 | 7.6 | 126 | 299 | 112 | 75 | 0.2 | 1085 | 7.4 | 19 | 79 | 28 | 82 | 0 | 1087 | 43 | 55 | 96 | 85 | 74 | 88 | 79 | 83 | 100 |
| D-5/2/90 | 37312 | 1 | 8.1 | 205 | 492 | 192 | 71 | 4 | 1454 | 8.1 | 197 | 200 | 72 | 5.5 | 1489 | 7.9 | 217 | 433 | 134 | 79 | 0.3 | 1423 | 7.7 | 32 | 114 | 37 | 84 | 0 | 1275 | 40 | 33 | 95 | 85 | 74 | 84 | 77 | 81 | 100 |
| D-6/2/90 | 38568 | 0.7 | 8.2 | 233 | 506 | 204 | 67 | 6.7 | 1692 | 8.3 | 218 | 212 | 66 | 11 | 1614 | 7.9 | 188 | 355 | 88 | 82 | 0.2 | 1516 | 7.8 | 47 | 153 | 59 | 81 | 0.1 | 1483 | 14 | 59 | 99 | 75 | 67 | 80 | 77 | 77 | 99 |
| D-7/2/90 | 38655 | 1.5 | 7.9 | 179 | 344 | 172 | 65 | 3.8 | 1379 | 8 | 148 | 156 | 74 | 4 | 1412 | 7.8 | 155 | 301 | 86 | 81 | 0.2 | 1426 | 7.5 | 26 | 97 | 35 | 83 | 0 | 1470 | 40 | 45 | 95 | 83 | 68 | 86 | 72 | 80 | 100 |
| D-8/2/90 | 34193 | 2 | 8 | 166 | 396 | 176 | 71 | 4 | 1265 | 8 | 178 | 188 | 70 | 5.5 | 1380 | 7.8 | 165 | 368 | 90 | 78 | 0.2 | 1434 | 7.5 | 26 | 106 | 31 | 84 | 0 | 1442 | 7.3 | 52 | 96 | 84 | 71 | 84 | 73 | 82 | 100 |
| D-9/2/90 | 36332 | 3.5 | 7.9 | 120 | 455 | 184 | 67 | 4 | 1224 | 8.1 | 205 | 188 | 68 | 5.5 | 1217 | 7.7 | 168 | 333 | 90 | 78 | 0.2 | 1353 | 7.6 | 24 | 98 | 32 | 81 | 0 | 1420 | 18 | 52 | 96 | 86 | 71 | 80 | 79 | 83 | 100 |
| D-11/2/90 | 32484 | 0.9 | 7.5 | 183 | 388 | 170 | 77 | 3.5 | 1130 | 7.6 | 197 | 178 | 75 | 4 | 1149 | 7.7 | 164 | 310 | 102 | 82 | 0.2 | 1212 | 7.5 | 22 | 89 | 33 | 91 | 0.1 | 1274 | 40 | 43 | 96 | 87 | 71 | 90 | 77 | 81 | 99 |
| D-12/2/90 | 37724 | 1 | 7.9 | 183 | 526 | 206 | 71 | 5.5 | 1422 | 7.9 | 197 | 218 | 72 | 6 | 1461 | 7.8 | 175 | 382 | 108 | 82 | 0.2 | 1595 | 7.5 | 34 | 128 | 40 | 84 | 0 | 1342 | 40 | 51 | 97 | 81 | 67 | 90 | 76 | 81 | 100 |
| D-13/2/90 | 36446 | 1 | 7.7 | 183 | 710 | 366 | 56 | 6.5 | 2400 | 7.8 | 197 | 256 | 63 | 6 | 2450 | 7.8 | 192 | 450 | 120 | 67 | 0.5 | 2330 | 7.6 | 18 | 295 | 88 | 76 | 0.3 | 2390 | 40 | 53 | 93 | 85 | 34 | 90 | 59 | 76 | 96 |
| D-14/2/90 | 35636 | 1.2 | 8 | 203 | 469 | 264 | 65 | 5.2 | 1489 | 8.1 | 197 | 304 | 66 | 8.5 | 1690 | 7.9 | 155 | 361 | 100 | 80 | 0.4 | 1718 | 7.3 | 27 |  | 38 | 97 | 0.1 | 1716 | 40 | 67 | 95 | 85 | 73 | 89 | 79 | 86 | 99 |
| D-15/2/90 | 34746 | 1 | 7.7 | 208 | 427 | 192 | 75 | 4.5 | 1426 | 7.7 | 195 | 236 | 75 | 7 | 1375 | 7.7 | 186 | 334 | 104 | 81 | 0.4 | 1518 | 7.4 | 27 | 78 | 33 | 85 | 0 | 1636 | 4.6 | 56 | 94 | 86 | 77 | 87 | 82 | 83 | 100 |
| D-16/2/90 | 34893 | 1.2 | 8 | 235 | 400 | 228 | 75 | 7 | 1532 | 8 | 232 | 252 | 78 | 8 | 1532 | 7.8 | 165 | 345 | 92 | 80 | 0.3 | 1478 | 7.6 | 25 | 125 | 25 | 91 | 0 | 1445 | 29 | 64 | 97 | 85 | 64 | 89 | 79 | 89 | 100 |
| D-18/2/90 | 37102 | 2 | 7.8 | 196 | 353 | 174 | 68 | 4 | 1315 | 7.8 | 152 | 162 | 77 | 3 | 1322 | 7.7 | 127 | 270 | 100 | 74 | 0.8 | 1337 | 7.6 | 24 | 71 | 24 | 92 | 0 | 1509 | 16 | 38 | 73 | 81 | 74 | 88 | 80 | 86 | 100 |
| D-19/2/90 | 41598 | 1.2 | 8.2 | 194 | 419 | 186 | 72 | 0.4 | 1310 | 8 | 210 | 208 | 71 | 4.5 | 1333 | 7.9 | 157 | 341 | 118 | 73 | 1 | 1474 | 7.6 | 23 | 71 | 33 | 79 | 0 | 1340 | 25 | 43 | 78 | 85 | 79 | 88 | 83 | 82 | 94 |
| D-21/2/90 | 38058 | 1 | 7.8 | 193 | 424 | 170 | 74 | 4 | 1406 | 7.7 | 226 | 356 | 72 | 4.5 | 1324 | 7.7 | 187 | 352 | 118 | 78 | 0.4 | 1360 | 7.5 | 24 | 88 | 29 | 87 | 0 | 1445 | 17 | 67 | 90 | 87 | 75 | 88 | 79 | 83 | 100 |
| D-22/2/90 | 40716 | 3.5 | 8.1 | 183 | 524 | 222 | 68 | 5.8 | 1597 | 8.1 | 230 | 248 | 66 | 7.5 | 1512 | 7.9 | 154 | 300 | 96 | 77 | 0.5 | 1521 | 7.4 | 29 | 76 | 25 | 84 | 0 | 1422 | 33 | 61 | 93 | 81 | 75 | 90 | 86 | 89 | 100 |
| D-23/2/90 | 40868 | 1.5 | 8.1 | 206 | 490 | 190 | 68 | 5.2 | 1392 | 8 | 220 | 224 | 70 | 6 | 1505 | 8.1 | 178 | 363 | 92 | 80 | 0.4 | 1532 | 7.7 | 16 | 86 | 26 | 82 | 0 | 1532 | 19 | 59 | 93 | 91 | 76 | 92 | 82 | 88 | 100 |
| D-25/2/90 | 36358 | 2 | 7.7 | 192 | 298 | 162 | 68 | 4 | 1241 | 7.7 | 160 | 188 | 68 | 4.5 | 1243 | 7.7 | 118 | 278 | 110 | 75 | 0.7 | 1285 | 7.5 | 30 | 98 | 45 | 74 | 0.1 | 1399 | 26 | 42 | 84 | 75 | 65 | 84 | 67 | 72 | 99 |
| D-26/2/90 | 40879 | 1.2 | 7.6 | 183 | 435 | 196 | 68 | 4.5 | 1421 | 7.7 | 197 | 264 | 70 | 7 | 1469 | 7.7 | 119 | 408 | 120 | 73 | 0.6 | 1532 | 7.6 | 18 | 102 | 23 | 84 | 0 | 1354 | 40 | 55 | 91 | 85 | 75 | 90 | 77 | 88 | 100 |
| D-27/2/90 | 44150 | 1 | 8.1 | 183 | 516 | 164 | 76 | 3.5 | 1548 | 8.1 | 197 | 232 | 74 | 5.5 | 1545 | 7.9 | 119 | 326 | 94 | 94 | 0.6 | 1415 | 7.6 | 18 | 113 | 37 | 96 | 0 | 1409 | 40 | 50 | 90 | 85 |  | 78 |  | 77 | 99 |
| D-28/2/90 | 45779 | 3 | 7.8 | 183 | 376 | 194 | 69 | 5 | 2020 | 7.8 | 197 | 276 | 62 | 7.5 | 2390 | 7.7 | 119 | 326 | 82 | 68 | 0.4 | 2260 | 7.7 | 18 | 66 | 15 | 93 | 0 | 2400 | 40 | 70 | 95 | 85 | 80 | 90 | 82 | 92 | 100 |
| D-1/1/90 | 41230 | 0.4 | 7.6 | 120 | 344 | 136 | 54 | 4.5 | 993 | 7.5 | 197 | 188 | 55 | 3 | 972 | 7.6 | 119 | 259 | 70 | 49 | 0.2 | 921 | 7.5 | 16 | 97 | 17 | 52 | 0 | 903 | 40 | 63 | 93 | 85 | 63 | 87 | 72 | 88 | 99 |
| D-2/1/90 | 37386 | 1.4 | 7.9 | 165 | 470 | 170 | 77 | 4 | 1365 | 7.9 | 197 | 192 | 71 | 4.5 | 1399 | 7.9 | 156 | 368 | 96 | 73 | 0.3 | 1338 | 7.6 | 22 | 97 | 18 | 81 | 0 | 1481 | 40 | 50 | 94 | 86 | 74 | 87 | 79 | 89 | 100 |
| D-3/1/90 | 34535 | 1 | 7.8 | 232 | 518 | 220 | 66 | 5.5 | 1617 | 7.9 | 230 | 202 | 71 | 4 | 1593 | 7.8 | 155 | 364 | 76 | 82 | 0.2 | 1594 | 7.5 | 29 | 146 | 31 | 77 | 0 | 1492 | 33 | 62 | 95 | 81 | 60 | 88 | 72 | 86 | 100 |
| D-4/1/90 | 32527 | 3 | 7.8 | 187 | 460 | 180 | 68 | 5.2 | 1832 | 7.9 | 219 | 236 | 66 | 5.5 | 1920 | 7.8 | 190 | 355 | 100 | 80 | 0.3 | 1646 | 7.5 | 28 | 105 | 30 | 82 | 0 | 1590 | 13 | 58 | 96 | 85 | 70 | 85 | 77 | 83 | 100 |

99

| ID | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-7/1/90 | 27760 | 1.2 | 7.6 | 199 | 466 | 186 | 74 | 4.5 | 1220 | 7.5 | 225 | 176 | 82 | 4 | 1208 | 7.5 | 139 | 314 | 94 | 87 | 0.2 | 1315 | 7.4 | 21 | 122 | 25 | 84 | 0 | 1411 | 38 | 47 | 95 | 85 | 61 | 89 | 74 | 87 | 100 |
| D-8/1/90 | 36281 | 2 | 7.8 | 183 | 612 | 226 | 71 | 8 | 1544 | 7.9 | 197 | 268 | 66 | 8 | 1503 | 7.8 | 158 | 259 | 100 | 80 | 0.4 | 1443 | 7.5 | 38 | 106 | 34 | 87 | 0 | 1239 | 40 | 63 | 96 | 76 | 59 | 90 | 83 | 85 | 100 |
| D-9/1/90 | 38055 | 3.5 | 7.8 | 221 | 524 | 188 | 72 | 5 | 1540 | 7.9 | 197 | 252 | 78 | 4.5 | 1477 | 7.8 | 128 | 299 | 82 | 90 | 0.2 | 1506 | 7.5 | 29 | 136 | 39 | 84 | 0 | 1503 | 40 | 68 | 96 | 77 | 55 | 87 | 74 | 79 | 100 |
| D-10/1/90 | 34064 | 1 | 8.1 | 230 | 535 | 242 | 67 | 6.5 | 1652 | 7.9 | 197 | 264 | 70 | 5.4 | 1700 | 7.9 | 174 | 322 | 100 | 72 | 0.3 | 1577 | 7.5 | 28 | 101 | 36 | 83 | 0.1 | 1552 | 40 | 62 | 94 | 84 | 69 | 88 | 81 | 85 | 99 |
| D-11/1/90 | 31447 | 3.5 | 7.9 | 190 | 374 | 192 | 71 | 6.5 | 1494 | 7.8 | 204 | 184 | 71 | 4.5 | 1462 | 7.7 | 121 | 259 | 100 | 82 | 0.5 | 1562 | 7.5 | 21 | 108 | 22 | 84 | 0.1 | 1596 | 41 | 46 | 89 | 83 | 58 | 89 | 71 | 89 | 99 |
| D-12/1/90 | 32127 | 17 | 7.7 | 183 | 526 | 292 | 64 | 7.5 | 2240 | 7.6 | 244 | 344 | 63 | 9 | 2220 | 7.7 | 193 | 450 | 134 | 70 | 1.3 | 2450 | 7.5 | 28 | 92 | 36 | 81 | 0 | 2580 | 21 | 61 | 86 | 86 | 80 | 90 | 83 | 88 | 100 |
| D-14/1/90 | 31059 | 3.5 | 7.8 | 202 | 431 | 200 | 74 | 5 | 1302 | 7.7 | 199 | 184 | 78 | 4.7 | 1307 | 7.6 | 124 | 269 | 96 | 75 | 0.4 | 1334 | 7.4 | 17 | 63 | 20 | 85 | 0.1 | 1473 | 38 | 48 | 92 | 86 | 77 | 92 | 85 | 90 | 99 |
| D-15/1/90 | 36470 | 4.5 | 7.8 | 227 | 526 | 212 | 69 | 4.5 | 1542 | 7.8 | 232 | 240 | 75 | 5.5 | 1583 | 7.8 | 172 | 411 | 106 | 83 | 0.2 | 1607 | 7.4 | 30 | 99 | 41 | 85 | 0.1 | 1395 | 26 | 56 | 96 | 83 | 76 | 87 | 81 | 81 | 99 |
| D-16/1/90 | 47449 | 1.7 | 7.8 | 170 | 401 | 158 | 67 | 4 | 1292 | 7.8 | 184 | 172 | 70 | 3.5 | 1413 | 7.8 | 167 | 345 | 102 | 75 | 0.7 | 1568 | 7.5 | 26 | 87 | 27 | 78 | 0.1 | 1635 | 9.2 | 41 | 80 | 84 | 75 | 85 | 78 | 83 | 99 |
| D-17/1/90 | 43940 | 3.5 | 7.8 | 149 | 361 | 186 | 62 | 3.2 | 1651 | 7.8 | 204 | 204 | 67 | 5 | 1751 | 7.9 | 155 | 345 | 106 | 70 | 0.5 | 1623 | 7.5 | 27 | 103 | 38 | 82 | 0.1 | 1597 | 24 | 48 | 90 | 83 | 70 | 82 | 72 | 80 | 98 |
| D-18/1/90 | 40347 | 1.8 | 7.7 | 155 | 338 | 132 | 70 | 2.7 | 1332 | 7.7 | 180 | 160 | 69 | 2.5 | 1366 | 7.7 | 152 | 319 | 94 | 72 | 0.6 | 1442 | 7.4 | 22 | 87 | 27 | 78 | 0.1 | 1482 | 16 | 41 | 76 | 86 | 73 | 86 | 74 | 80 | 98 |
| D-19/1/90 | 40267 | 1.8 | 7.9 | 180 | 433 | 186 | 72 | 4 | 1729 | 7.9 | 200 | 174 | 76 | 3.7 | 1820 | 7.8 | 127 | 354 | 86 | 79 | 0.4 | 1856 | 7.5 | 19 | 110 | 26 | 83 | 0 | 1861 | 37 | 51 | 89 | 85 | 69 | 89 | 75 | 86 | 99 |
| D-21/1/90 | 37976 | 1 | 7.7 | 148 | 345 | 162 | 78 | 3.2 | 1432 | 7.6 | 145 | 166 | 74 | 3.8 | 1415 | 7.6 | 137 | 302 | 94 | 85 | 0.5 | 1461 | 7.4 | 17 | 60 | 21 | 91 | 0.1 | 1572 | 5.5 | 43 | 87 | 88 | 80 | 89 | 83 | 87 | 98 |
| D-22/1/90 | 47368 | 2 | 7.9 | 156 | 417 | 152 | 74 | 4.5 | 1608 | 7.9 | 207 | 164 | 67 | 4.5 | 1798 | 7.8 | 147 | 321 | 76 | 76 | 0.4 | 1627 | 7.4 | 29 | 99 | 26 | 85 | 0 | 1490 | 29 | 54 | 91 | 80 | 69 | 81 | 76 | 83 | 100 |
| D-23/1/90 | 48086 | 5 | 8 | 247 | 444 | 166 | 74 | 3.5 | 1700 | 8.1 | 207 | 182 | 66 | 3.5 | 1724 | 7.9 | 176 | 360 | 112 | 77 | 0.4 | 1768 | 7.6 | 27 | 52 | 29 | 86 | 0 | 1764 | 15 | 39 | 89 | 85 | 86 | 89 | 88 | 83 | 100 |
| D-24/1/90 | 47642 | 5 | 7.9 | 157 | 428 | 204 | 54 | 4 | 1989 | 7.7 | 187 | 250 | 53 | 4.5 | 1869 | 7.8 | 219 | 303 | 90 | 69 | 0.4 | 1735 | 7.5 | 19 | 73 | 18 | 75 | 0 | 1800 | 40 | 64 | 92 | 91 | 76 | 88 | 83 | 91 | 100 |
| D-25/1/90 | 43174 | 4.5 | 7.7 | 179 | 420 | 158 | 68 | 2.5 | 1260 | 7.7 | 200 | 160 | 69 | 3.5 | 1256 | 7.8 | 183 | 376 | 110 | 71 | 1 | 1369 | 7.4 | 23 | 136 | 17 | 88 | 0 | 1365 | 8.5 | 31 | 71 | 87 | 64 | 87 | 68 | 89 | 100 |
| D-26/1/90 | 39891 | 2 | 7.6 | 178 | 416 | 188 | 61 | 4.5 | 1301 | 7.6 | 174 | 166 | 70 | 3.5 | 1267 | 7.6 | 164 | 368 | 106 | 70 | 0.4 | 1280 | 7.3 | 22 | 140 | 20 | 70 | 0 | 1262 | 5.7 | 36 | 89 | 87 | 62 | 88 | 66 | 89 | 100 |
| D-28/1/90 | 32257 | 3.5 | 7.5 | 246 | 583 | 504 | 45 | 7.3 | 1016 | 7.5 | 228 | 436 | 42 | 7.5 | 1079 | 7.5 | 85 | 236 | 76 | 68 | 0.5 | 1088 | 7.3 | 21 | 75 | 25 | 72 | 0.1 | 1164 | 63 | 83 | 94 | 75 | 68 | 92 | 87 | 95 | 99 |
| D-29/1/90 | 40498 | 10 | 8.1 | 202 | 476 | 300 | 49 | 3.7 | 1636 | 8 | 206 | 252 | 51 | 3.5 | 1579 | 7.9 | 186 | 394 | 108 | 72 | 1.3 | 1413 | 7.5 | 27 | 75 | 22 | 91 | 0.1 | 1291 | 9.7 | 57 | 63 | 86 | 81 | 87 | 84 | 93 | 99 |
| D-30/1/90 | 40221 | 2 | 8.1 | 177 | 407 | 172 | 58 | 2.5 | 1379 | 8 | 231 | 248 | 55 | 4.6 | 1454 | 8 | 188 | 379 | 108 | 72 | 0.7 | 1529 | 7.5 | 27 | 95 | 26 | 77 | 0 | 1542 | 19 | 57 | 85 | 86 | 75 | 85 | 77 | 85 | 100 |
| D-31/1/90 | 46669 | 1.8 | 7.8 | 183 | 340 | 168 | 71 | 2.3 | 1477 | 7.8 | 197 | 228 | 65 | 4.5 | 1379 | 7.8 | 119 | 368 | 110 | 73 | 0.6 | 1550 | 7.6 | 18 | 84 | 19 | 81 | 0 | 1445 | 40 | 52 | 87 | 85 | 70 | 90 | 79 | 91 | 100 |
| D-1/6/90 | 34669 | 1.2 | 7.8 | 198 | 381 | 216 | 52 | 3.5 | 1415 | 7.8 | 183 | 220 | 56 | 3.5 | 1453 | 7.7 | 123 | 246 | 76 | 79 | 0.2 | 1432 | 7.7 | 29 | 91 | 26 | 81 | 0.1 | 1390 | 33 | 66 | 94 | 76 | 63 | 85 | 76 | 88 | 99 |
| D-3/6/90 | 41824 | 1.2 | 7.8 | 161 | 281 | 164 | 59 | 2.3 | 1075 | 7.8 | 144 | 192 | 53 | 2.5 | 1068 | 7.8 | 118 | 233 | 88 | 71 | 0.4 | 1121 | 7.8 | 14 | 63 | 19 | 74 | 0 | 1246 | 18 | 54 | 84 | 88 | 73 | 91 | 78 | 88 | 100 |
| D-4/6/90 | 51520 | 2 | 7.3 | 156 | 336 | 192 | 63 | 5.5 | 1320 | 7.8 | 158 | 184 | 57 | 4 | 1327 | 7.6 | 79 | 198 | 92 | 63 | 0.3 | 1175 | 7.9 | 15 | 59 | 19 | 79 | 0 | 1054 | 50 | 50 | 93 | 81 | 70 | 90 | 82 | 90 | 100 |
| D-5/6/90 | 39421 | 1 | 7.9 | 189 | 457 | 1004 | 26 | 24 | 1218 | 7.8 | 234 | 1384 | 25 | 35 | 1257 | 7.7 | 156 | 323 | 140 | 66 | 0.3 | 1308 | 7.8 | 19 | 79 | 21 | 81 | 0.1 | 1172 | 33 | 90 | 99 | 88 | 76 | 90 | 83 | 98 | 100 |
| D-6/6/90 | 36131 | 1 | 7.9 | 215 | 500 | 252 | 62 | 4.7 | 1512 | 7.8 | 233 | 348 | 49 | 7.5 | 1427 | 7.8 | 147 | 327 | 102 | 69 | 0.3 | 1436 | 7.9 | 25 | 75 | 23 | 78 | 0 | 1409 | 37 | 71 | 96 | 83 | 77 | 88 | 85 | 91 | 100 |
| D-7/6/90 | 33251 | 1 | 7.6 | 225 | 578 | 256 | 66 | 5.5 | 1510 | 7.6 | 224 | 276 | 54 | 6.5 | 1486 | 7.7 | 119 | 319 | 102 | 65 | 0.3 | 1492 | 7.6 | 15 | 151 | 25 | 77 | 0.1 | 1461 | 47 | 63 | 95 | 87 | 53 | 93 | 74 | 90 | 99 |
| D-8/6/90 | 35789 | 1.5 | 7.4 | 316 | 533 | 264 | 55 | 5.5 | 1361 | 7.4 | 352 | 344 | 47 | 6.5 | 1453 | 7.5 | 190 | 361 | 160 | 51 | 1.5 | 1518 | 7.5 | 39 | 150 | 54 | 74 | 0 | 1506 | 46 | 54 | 77 | 80 | 58 | 88 | 72 | 80 | 98 |
| D-10/6/90 | 40106 | 0.6 | 7.8 | 238 | 504 | 292 | 59 | 6.5 | 1109 | 7.8 | 361 | 352 | 51 | 7 | 1113 | 7.8 | 147 | 256 | 128 | 66 | 0.8 | 1113 | 8 | 19 | 80 | 20 | 70 | 0.1 | 1238 | 59 | 64 | 89 | 87 | 69 | 92 | 84 | 93 | 99 |
| D-11/6/90 | 45191 | 2 | 8 | 125 | 324 | 362 | 36 | 5 | 1093 | 7.7 | 297 | 804 | 33 | 13 | 1086 | 7.9 | 84 | 204 | 124 | 53 | 0.6 | 1105 | 7.9 | 28 | 128 | 14 | 74 | 0 | 1008 | 72 | 85 | 95 | 67 | 37 | 78 | 61 | 96 | 100 |
| D-12/6/90 | 43308 | 1.4 | 7.9 | 265 | 330 | 562 | 27 | 7.5 | 1866 | 7.9 | 242 | 680 | 30 | 13 | 1858 | 7.9 | 133 | 206 | 120 | 48 | 0.6 | 1850 | 7.8 | 24 | 70 | 11 | 62 | 0 | 1800 | 45 | 82 | 95 | 82 | 66 | 91 | 79 | 98 | 100 |
| D-13/6/90 | 37615 | 1.2 | 7.8 | 199 | 404 | 232 | 53 | 5 | 1310 | 7.8 | 416 | 544 | 43 | 11 | 1366 | 7.7 | 143 | 299 | 114 | 61 | 0.4 | 1466 | 7.6 | 25 | 85 | 19 | 65 | 0 | 1404 | 66 | 79 | 96 | 83 | 72 | 87 | 79 | 92 | 100 |
| D-14/6/90 | 42596 | 3 | 7.7 | 138 | 259 | 456 | 23 | 4 | 1007 | 7.6 | 160 | 584 | 27 | 7.5 | 1039 | 7.6 | 101 | 176 | 126 | 48 | 0.5 | 1113 | 7.6 | 22 | 59 | 16 | 69 | 0 | 1194 | 37 | 78 | 93 | 78 | 67 | 84 | 77 | 97 | 100 |
| D-15/6/90 | 41948 | 1.5 | 7.7 | 198 | 396 | 216 | 53 | 3 | 1282 | 7.8 | 245 | 328 | 48 | 5 | 1321 | 7.7 | 109 | 244 | 120 | 57 | 0.3 | 1304 | 7.8 | 23 | 105 | 32 | 71 | 0 | 1295 | 56 | 63 | 94 | 79 | 57 | 88 | 74 | 85 | 99 |
| D-17/6/90 | 34647 | 1 | 7.5 | 193 | 342 | 260 | 51 | 3 | 985 | 7.6 | 230 | 424 | 44 | 6 | 1008 | 7.6 | 132 | 223 | 174 | 51 | 0.4 | 992 | 7.3 | 18 | 84 | 22 | 69 | 0 | 1218 | 43 | 59 | 93 | 85 | 70 | 90 | 79 | 92 | 100 |
| D-18/6/90 | 36967 | 1 | 7.6 | 202 | 426 | 248 | 81 | 5.5 | 2310 | 7.7 | 326 | 404 | 56 | 7 | 2180 | 7.7 | 142 | 280 | 146 | 66 | 0.6 | 1909 | 7.7 | 33 | 92 | 23 | 73 | 0 | 1813 | 56 | 64 | 92 | 77 | 67 | 84 | 78 | 91 | 100 |
| D-19/6/90 | 34879 | 1 | 7.5 | 319 | 465 | 214 | 64 | 5.5 | 1308 | 7.6 | 364 | 388 | 51 | 7.5 | 1344 | 7.7 | 143 | 261 | 100 | 70 | 0.4 | 1378 | 7.7 | 18 | 104 | 11 | 87 | 0 | 1345 | 61 | 74 | 95 | 87 | 60 | 94 | 78 | 95 | 100 |
| D-20/6/90 | 34365 | 6 | 7.6 | 236 | 444 | 236 | 63 | 4.9 | 1400 | 7.7 | 259 | 440 | 50 | 10 | 1439 | 7.6 | 140 | 246 | 122 | 59 | 0.3 | 1458 | 7.7 | 24 | 115 | 17 | 71 | 0.1 | 1480 | 46 | 72 | 97 | 83 | 53 | 90 | 74 | 93 | 99 |
| D-21/6/90 | 34291 | 8 | 7.9 | 192 | 433 | 300 | 57 | 7 | 1395 | 7.9 | 269 | 436 | 51 | 12 | 1335 | 7.9 | 138 | 294 | 108 | 69 | 0.3 | 1378 | 7.8 | 15 | 60 | 12 | 92 | 0 | 1423 | 49 | 75 | 98 | 89 | 80 | 92 | 86 | 96 | 100 |
| D-22/6/90 | 34886 | 8 | 7.7 | 211 | 488 | 268 | 54 | 5.8 | 1212 | 7.9 | 265 | 348 | 58 | 8.5 | 1274 | 7.7 | 151 | 306 | 88 | 80 | 0.2 | 1309 | 7.7 | 15 | 63 | 12 | 83 | 0 | 1320 | 43 | 75 | 98 | 90 | 79 | 93 | 87 | 96 | 100 |
| D-24/6/90 | 38731 | 1.2 | 7.5 | 200 | 402 | 184 | 64 | 2.8 | 1127 | 7.5 | 188 | 218 | 63 | 3 | 1140 | 7.5 | 122 | 233 | 90 | 78 | 0.2 | 1200 | 7.7 | 17 | 64 | 17 | 77 | 0 | 1317 | 35 | 59 | 95 | 86 | 73 | 92 | 84 | 91 | 100 |
| D-25/6/90 | 39308 | 3 | 7.8 | 217 | 349 | 172 | 70 | 3.5 | 1454 | 7.8 | 174 | 208 | 65 | 4 | 1487 | 7.8 | 98 | 261 | 78 | 77 | 0.2 | 1360 | 7.9 | 18 | 92 | 20 | 80 | 0 | 1230 | 44 | 63 | 95 | 85 | 65 | 90 | 74 | 88 | 100 |
| D-26/6/90 | 44198 | 7 | 7.7 | 257 | 667 | 1016 | 32 | 22 | 1478 | 7.8 | 212 | 572 | 36 | 13 | 1422 | 7.9 | 135 | 210 | 108 | 57 | 0.2 | 1358 | 8 | 16 | 56 | 15 | 73 | 0 | 1378 | 36 | 81 | 98 | 88 | 73 | 94 | 92 | 99 | 100 |
| D-27/6/90 | 39003 | 1.2 | 7.8 | 183 | 456 | 232 | 66 | 5 | 1262 | 7.8 | 198 | 216 | 56 | 4.5 | 1247 | 7.9 | 107 | 266 | 84 | 57 | 0.2 | 1305 | 7.8 | 19 | 60 | 12 | 75 | 0 | 1234 | 46 | 61 | 97 | 82 | 77 | 90 | 87 | 95 | 100 |
| D-28/6/90 | 34487 | 0.7 | 7.9 | 183 | 380 | 192 | 63 | 4.5 | 1339 | 7.8 | 196 | 216 | 56 | 6 | 1403 | 7.8 | 128 | 270 | 68 | 71 | 0.1 | 1409 | 8 | 13 | 63 | 15 | 67 | 0 | 1344 | 35 | 69 | 98 | 90 | 77 | 93 | 83 | 92 | 100 |
| D-29/6/90 | 35198 | 0.8 | 7.7 | 185 | 372 | 164 | 61 | 3 | 1623 | 7.8 | 210 | 192 | 65 | 3 | 1508 | 7.8 | 124 | 278 | 74 | 73 | 0.1 | 1482 | 7.9 | 18 | 78 | 15 | 67 | 0 | 1491 | 41 | 62 | 97 | 86 | 72 | 90 | 79 | 91 | 100 |
| D-1/5/90 | 27617 | 1 | 7.6 | 285 | 436 | 218 | 68 | 6 | 1095 | 7.6 | 292 | 238 | 67 | 5.5 | 1149 | 7.5 | 160 | 284 | 80 | 80 | 0.2 | 1139 | 7.5 | 26 | 104 | 41 | 73 | 0.5 | 1146 | 45 | 66 | 96 | 84 | 63 | 91 | 76 | 81 | 92 |
| D-2/5/90 | 37881 | 3 | 7.7 | 257 | 588 | 328 | 60 | 7 | 1392 | 7.7 | 213 | 272 | 60 | 6.5 | 1432 | 7.8 | 144 | 344 | 116 | 70 | 0.8 | 1356 | 7.6 | 26 | 92 | 44 | 74 | 0.3 | 1213 | 18 | 57 | 88 | 85 | 73 | 90 | 84 | 87 | 96 |
| D-3/5/90 | 39024 | 1.2 | 7.9 | 268 | 467 | 224 | 66 | 6.5 | 1409 | 7.9 | 328 | 312 | 63 | 9.5 | 1376 | 7.9 | 152 | 338 | 88 | 77 | 0.2 | 1336 | 7.7 | 29 | 148 | 49 | 80 | 0.1 | 1334 | 54 | 72 | 98 | 81 | 56 | 89 | 68 | 78 | 99 |
| D-4/5/90 | 38990 | 1.4 | 7.8 | 189 | 357 | 172 | 67 | 4 | 1160 | 7.7 | 213 | 212 | 62 | 5.5 | 1232 | 7.7 | 123 | 255 | 80 | 80 | 0.3 | 1237 | 7.6 | 30 | 122 | 42 | 81 | 0.1 | 1304 | 42 | 62 | 95 | 76 | 52 | 84 | 66 | 76 | 99 |
| D-6/5/90 | 37710 | 3 | 7.5 | 312 | 388 | 204 | 62 | 4 | 1026 | 7.5 | 211 | 214 | 65 | 5 | 1020 | 7.5 | 111 | 222 | 86 | 77 | 0.2 | 1016 | 7.4 | 42 | 154 | 69 | 83 | 0.1 | 1092 | 47 | 60 | 96 | 62 | 31 | 87 | 60 | 80 | 98 |
| D-7/5/90 | 25957 | 0.6 | 8.1 | 404 | 455 | 448 | 38 | 5 | 1229 | 7.8 | 491 | 692 | 39 | 12 | 1279 | 7.8 | 285 | 274 | 110 | 71 | 0.4 | 1156 | 7.8 | 20 | 83 | 25 | 91 | 0.1 | 1021 | 42 | 84 | 97 | 93 | 70 | 95 | 82 | 94 | 99 |
| D-8/5/90 | 38623 | 1.5 | 8.1 | 243 | 299 | 180 | 54 | 2.5 | 1615 | 8 | 291 | 324 | 65 | 3.5 | 1630 | 7.9 | 166 | 265 | 96 | 60 | 0.2 | 1485 | 7.8 | 26 | 82 | 36 | 74 | 0 | 1345 | 43 | 70 | 94 | 84 | 69 | 89 | 73 | 80 | 99 |
| D-9/5/90 | 41746 | 1 | 8 | 352 | 471 | 208 | 65 | 5.5 | 2150 | 7.9 | 386 | 216 | 65 | 5.5 | 2320 | 7.6 | 198 | 348 | 118 | 75 | 0.5 | 2290 | 7.5 | 24 | 91 | 32 | 84 | 0 | 2290 | 49 | 45 | 92 | 88 | 74 | 93 | 81 | 85 | 100 |
| D-10/5/90 | 43291 | 0.8 | 7.7 | 215 | 447 | 164 | 61 | 5.5 | 1177 | 7.7 | 258 | 256 | 53 | 5 | 1207 | 7.7 | 141 | 285 | 86 | 70 | 0.1 | 1324 | 7.6 | 19 | 75 | 32 | 77 | 0 | 1390 | 45 | 66 | 96 | 87 | 74 | 91 | 83 | 81 | 100 |
| D-11/5/90 | 41436 | 1.2 | 7.9 | 191 | 356 | 192 | 63 | 7 | 1434 | 7.9 | 226 | 288 | 56 | 10 | 1427 | 7.7 | 130 | 257 | 84 | 71 | 0.1 | 1467 | 7.6 | 28 | 103 | 34 | 77 | 0 | 1469 | 43 | 71 | 99 | 79 | 60 | 85 | 71 | 82 | 100 |
| D-13/5/90 | 39402 | 1.2 | 7.9 | 283 | 274 | 162 | 59 | 3 | 937 | 7.9 | 190 | 178 | 58 | 2.8 | 944 | 7.8 | 117 | 212 | 78 | 77 | 0.1 | 970 | 7.7 | 28 | 74 | 25 | 88 | 0 | 1060 | 38 | 56 | 98 | 76 | 65 | 90 | 73 | 85 | 100 |
| D-14/5/90 | 39383 | 0.6 | 8 | 216 | 529 | 248 | 61 | 6.3 | 1179 | 8 | 429 | 220 | 66 | 5 | 1294 | 7.8 | 140 | 337 | 138 | 64 | 1.5 | 1155 | 7.7 | 23 | 94 | 22 | 84 | 0 | 1107 | 67 | 37 | 70 | 84 | 72 | 89 | 82 | 91 | 100 |

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-15/5/90 | 37106 | 0.7 | 7.8 | 163 | 468 | 202 | 57 | 4.5 | 1525 | 8 | 319 | 292 | 59 | 7 | 1587 | 7.8 | 146 | 353 | 230 | 56 | 0.2 | 1328 | 7.7 | 14 | 100 | 19 | 87 | 0 | 1309 | 54 | 21 | 97 | 90 | 72 | 91 | 79 | 91 | 100 |
| D-16/5/90 | 36591 | 1.2 | 7.8 | 183 | 499 | 248 | 57 | 6.5 | 1213 | 7.8 | 197 | 304 | 58 | 7.2 | 1264 | 7.8 | 119 | 334 | 98 | 78 | 0.2 | 1420 | 7.6 | 18 | 77 | 20 | 75 | 0 | 1395 | 40 | 68 | 97 | 85 | 77 | 90 | 85 | 92 | 100 |
| D-17/5/90 | 33711 | 2.5 | 8 | 345 | 457 | 288 | 60 | 7.5 | 1272 | 7.9 | 197 | 348 | 63 | 9 | 1336 | 7.9 | 119 | 372 | 146 | 73 | 2 | 1369 | 7.8 | 19 | 88 | 19 | 79 | 0 | 1437 | 40 | 58 | 78 | 85 | 76 | 95 | 81 | 93 | 100 |
| D-18/5/90 | 35081 | 0.8 | 7.7 | 431 | 532 | 210 | 66 | 5 | 1535 | 7.6 | 409 | 224 | 59 | 6 | 1561 | 7.8 | 210 | 360 | 98 | 65 | 0.2 | 1472 | 7.7 | 13 | 132 | 13 | 50 | 0 | 1516 | 49 | 56 | 97 | 94 | 63 | 97 | 75 | 94 | 100 |
| D-20/5/90 | 32372 | 1.5 | 7.8 | 174 | 404 | 192 | 64 | 3.5 | 1056 | 7.7 | 182 | 210 | 61 | 4.5 | 1064 | 7.7 | 140 | 243 | 104 | 77 | 0.4 | 1078 | 7.7 | 13 | 55 | 18 | 83 | 0 | 1165 | 23 | 51 | 91 | 91 | 77 | 93 | 86 | 91 | 100 |
| D-21/5/90 | 37283 | 2 | 7.7 | 327 | 376 | 184 | 58 | 4 | 1226 | 7.7 | 287 | 236 | 46 | 3.5 | 1223 | 7.8 | 102 | 255 | 84 | 74 | 0.2 | 1199 | 7.8 | 15 | 47 | 13 | 69 | 0 | 1136 | 65 | 64 | 94 | 85 | 82 | 95 | 88 | 93 | 100 |
| D-22/5/90 | 42202 | 1 | 7.6 | 184 | 238 | 316 | 30 | 3 | 1079 | 7.7 | 163 | 344 | 29 | 4 | 1113 | 7.6 | 122 | 214 | 92 | 63 | 0.2 | 1139 | 7.6 | 27 | 95 | 39 | 59 | 0 | 1176 | 25 | 73 | 95 | 78 | 56 | 85 | 60 | 88 | 100 |
| D-23/5/90 | 50942 | 3 | 7.8 | 159 | 234 | 292 | 33 | 3 | 1140 | 7.8 | 182 | 344 | 33 | 4 | 1125 | 7.9 | 111 | 218 | 100 | 52 | 0.3 | 1231 | 7.8 | 12 | 36 | 14 | 77 | 0 | 1276 | 39 | 71 | 94 | 89 | 84 | 93 | 85 | 95 | 100 |
| D-24/5/90 | 44040 | 1 | 7.9 | 275 | 330 | 180 | 56 | 3 | 1376 | 7.9 | 190 | 200 | 60 | 2.5 | 1392 | 7.9 | 118 | 260 | 76 | 76 | 0.3 | 1435 | 7.7 | 25 | 109 | 24 | 75 | 0 | 1399 | 38 | 62 | 90 | 79 | 58 | 91 | 67 | 87 | 100 |
| D-25/5/90 | 43117 | 1.2 | 8.4 | 134 | 302 | 208 | 58 | 3.3 | 1160 | 8 | 159 | 276 | 61 | 3.5 | 1201 | 8 | 116 | 233 | 80 | 78 | 0.2 | 1264 | 7.9 | 16 | 66 | 20 | 80 | 0 | 1333 | 27 | 71 | 94 | 86 | 72 | 88 | 78 | 90 | 100 |
| D-27/5/90 | 48333 | 0.7 | 7.6 | 132 | 188 | 172 | 47 | 2.5 | 940 | 7.7 | 210 | 192 | 47 | 2 | 924 | 7.7 | 102 | 176 | 76 | 79 | 0.1 | 1001 | 7.7 | 17 | 78 | 21 | 88 | 0 | 1107 | 51 | 60 | 95 | 83 | 56 | 87 | 59 | 88 | 100 |
| D-28/5/90 | 46540 | 1.2 | 7.8 | 132 | 297 | 176 | 46 | 1.5 | 1140 | 8 | 133 | 172 | 44 | 1.3 | 1072 | 7.9 | 98 | 246 | 84 | 62 | 0.3 | 1111 | 7.6 | 24 | 63 | 16 | 69 | 0 | 1014 | 26 | 51 | 81 | 76 | 74 | 82 | 79 | 91 | 100 |
| D-29/5/90 | 46057 | 1 | 8.1 | 133 | 288 | 162 | 57 | 2.3 | 1050 | 8.1 | 133 | 152 | 63 | 2.5 | 1116 | 8 | 97 | 220 | 70 | 77 | 0.2 | 1142 | 7.9 | 32 | 128 | 40 | 83 | 0.1 | 1187 | 27 | 54 | 94 | 67 | 42 | 76 | 56 | 75 | 98 |
| D-30/5/90 | 45018 | 2 | 7.9 | 224 | 488 | 456 | 48 | 11 | 1147 | 7.9 | 153 | 252 | 49 | 4 | 1224 | 8 | 110 | 248 | 92 | 67 | 0.4 | 1238 | 8 | 15 | 36 | 17 | 82 | 0 | 1217 | 28 | 64 | 91 | 86 | 86 | 93 | 93 | 96 | 100 |
| D-1/4/90 | 40552 | 2 | 7.9 | 200 | 395 | 178 | 72 | 3.5 | 1038 | 7.8 | 190 | 156 | 72 | 3.5 | 1063 | 7.7 | 127 | 281 | 102 | 73 | 0.7 | 1023 | 7.6 | 17 | 75 | 19 | 90 | 0 | 1148 | 33 | 35 | 80 | 87 | 73 | 92 | 81 | 89 | 99 |
| D-2/4/90 | 53210 | 5 | 7.8 | 87 | 241 | 236 | 38 | 4 | 1179 | 8 | 97 | 248 | 34 | 3.5 | 1231 | 8 | 78 | 190 | 70 | 57 | 0 | 902 | 7.7 | 10 | 79 | 16 | 78 | 0 | 899 | 20 | 72 | 100 | 87 | 58 | 89 | 67 | 93 | 100 |
| D-3/4/90 | 53530 | 1.5 | 8.1 | 132 | 336 | 330 | 37 | 6 | 1234 | 8 | 161 | 264 | 38 | 5.5 | 1327 | 7.9 | 126 | 316 | 80 | 75 | 0.7 | 1257 | 7.7 | 13 | 68 | 13 | 69 | 0 | 1229 | 22 | 70 | 88 | 90 | 79 | 90 | 80 | 96 | 100 |
| D-4/4/90 | 46659 | 2 | 7.7 | 175 | 321 | 176 | 58 | 3.5 | 1335 | 7.8 | 198 | 192 | 58 | 4 | 1418 | 7.7 | 141 | 281 | 76 | 79 | 0.3 | 1316 | 7.5 | 12 | 67 | 11 | 100 | 0 | 1308 | 29 | 60 | 93 | 92 | 76 | 93 | 79 | 94 | 100 |
| D-5/4/90 | 45772 | 2 | 7.9 | 162 | 348 | 156 | 64 | 3.1 | 1340 | 7.9 | 228 | 252 | 65 | 5.9 | 1398 | 7.9 | 135 | 329 | 110 | 78 | 0.7 | 1377 | 7.7 | 17 | 127 | 14 | 93 | 0 | 1423 | 41 | 56 | 88 | 87 | 61 | 90 | 64 | 91 | 100 |
| D-6/4/90 | 52933 | 1.8 | 7.9 | 135 | 317 | 180 | 54 | 3.5 | 1362 | 7.9 | 181 | 268 | 60 | 4.6 | 1417 | 7.6 | 114 | 285 | 104 | 71 | 0.5 | 1459 | 7.8 | 14 | 91 | 22 | 86 | 0 | 1555 | 37 | 61 | 89 | 88 | 68 | 90 | 71 | 88 | 100 |
| D-8/4/90 | 36510 | 1.5 | 7.9 | 91 | 325 | 122 | 77 | 3.5 | 1026 | 7.9 | 109 | 124 | 94 | 3.7 | 1063 | 7.8 | 81 | 227 | 64 | 97 | 0.4 | 983 | 7.7 | 15 | 74 | 11 | 100 | 0 | 1039 | 26 | 48 | 89 | 82 | 67 | 84 | 77 | 91 | 100 |
| D-9/4/90 | 34299 | 3 | 7.8 | 210 | 725 | 350 | 42 | 5 | 1265 | 7.8 | 238 | 428 | 42 | 7.5 | 1262 | 7.8 | 136 | 306 | 106 | 47 | 0.4 | 1287 | 7.6 | 21 | 94 | 18 | 67 | 0 | 1309 | 43 | 75 | 95 | 85 | 69 | 90 | 87 | 95 | 100 |
| D-10/4/90 | 41073 | 0.8 | 8.1 | 166 | 422 | 184 | 63 | 3 | 1450 | 8 | 118 | 280 | 66 | 3.5 | 1429 | 8.1 | 119 | 323 | 114 | 70 | 0.4 | 1419 | 7.8 | 26 | 81 | 22 | 76 | 0 | 1357 | 40 | 59 | 89 | 78 | 75 | 84 | 81 | 88 | 99 |
| D-11/4/90 | 43536 | 2.5 | 7.8 | 267 | 342 | 202 | 58 | 2 | 1327 | 7.7 | 254 | 172 | 70 | 2 | 1306 | 7.8 | 130 | 349 | 120 | 67 | 0.4 | 1450 | 7.7 | 12 | 77 | 16 | 88 | 0 | 1454 | 49 | 30 | 80 | 91 | 78 | 96 | 78 | 92 | 100 |
| D-13/4/90 | 34667 | 6 | 7.2 | 165 | 315 | 170 | 65 | 4.5 | 1125 | 7.3 | 219 | 180 | 68 | 4.5 | 1151 | 7.3 | 121 | 235 | 86 | 81 | 0.3 | 1204 | 7.4 | 11 | 53 | 19 | 74 | 0 | 1306 | 45 | 52 | 93 | 91 | 77 | 93 | 83 | 89 | 100 |
| D-16/4/90 | 29624 | 1.4 | 7.5 | 184 | 219 | 148 | 62 | 3 | 1530 | 7.5 | 189 | 146 | 60 | 2 | 1553 | 7.5 | 92 | 192 | 78 | 77 | 0.3 | 1560 | 7.5 | 18 | 37 | 22 | 73 | 0 | 1100 | 51 | 47 | 85 | 80 | 81 | 90 | 83 | 85 | 99 |
| D-17/4/90 | 34069 | 1.2 | 7.5 | 89 | 298 | 116 | 69 | 2 | 1119 | 7.7 | 89 | 120 | 70 | 2.3 | 1103 | 7.7 | 79 | 250 | 74 | 81 | 0.2 | 1188 | 7.6 | 16 | 42 | 22 | 84 | 0 | 1149 | 11 | 38 | 91 | 80 | 83 | 82 | 86 | 81 | 99 |
| D-18/4/90 | 37782 | 1.2 | 8.1 | 155 | 382 | 174 | 67 | 7 | 1205 | 7.9 | 295 | 212 | 74 | 5.5 | 1319 | 7.8 | 147 | 312 | 90 | 76 | 0.3 | 1246 | 7.6 | 18 | 87 | 24 | 77 | 0.1 | 1244 | 50 | 58 | 95 | 88 | 72 | 88 | 77 | 86 | 99 |
| D-19/4/90 | 42109 | 0.7 | 7.8 | 159 | 350 | 150 | 69 | 3 | 1304 | 7.8 | 192 | 160 | 75 | 2.5 | 1206 | 7.7 | 166 | 312 | 82 | 76 | 0.2 | 1276 | 7.7 | 41 | 152 | 41 | 79 | 0 | 1341 | 14 | 49 | 94 | 75 | 51 | 74 | 57 | 73 | 99 |
| D-20/4/90 | 40871 | 0.7 | 7.7 | 193 | 298 | 120 | 60 | 3 | 1068 | 7.8 | 160 | 200 | 38 | 2.5 | 1167 | 7.8 | 135 | 222 | 90 | 75 | 0.3 | 1211 | 7.8 | 16 | 69 | 19 | 81 | 0 | 1199 | 16 | 60 | 88 | 88 | 69 | 92 | 77 | 91 | 99 |
| D-22/4/90 | 36088 | 15 | 7.5 | 75 | 133 | 142 | 39 | 2 | 860 | 7.6 | 127 | 172 | 47 | 1.8 | 838 | 7.6 | 73 | 124 | 62 | 68 | 0.3 | 871 | 7.3 | 9 | 32 | 12 | 80 | 0 | 889 | 43 | 64 | 86 | 88 | 74 | 88 | 76 | 92 | 100 |
| D-23/4/90 | 47255 | 0.6 | 7.7 | 181 | 264 | 124 | 61 | 2.6 | 1087 | 7.8 | 148 | 150 | 52 | 3 | 1047 | 7.8 | 122 | 232 | 68 | 62 | 0.3 | 1103 | 7.7 | 19 | 72 | 28 | 64 | 0.1 | 1008 | 18 | 55 | 92 | 84 | 69 | 90 | 73 | 77 | 96 |
| D-24/4/90 | 55300 | 1.5 | 7.9 | 254 | 356 | 206 | 56 | 4 | 1300 | 7.9 | 207 | 222 | 53 | 3.5 | 1347 | 7.9 | 91 | 184 | 66 | 70 | 0.3 | 1188 | 7.6 | 18 | 100 | 20 | 76 | 0 | 1202 | 56 | 70 | 91 | 80 | 46 | 93 | 72 | 90 | 100 |
| D-25/4/90 | 37646 | 1 | 7.8 | 355 | 556 | 184 | 50 | 3.3 | 1278 | 7.8 | 308 | 194 | 52 | 3 | 1376 | 7.8 | 122 | 216 | 78 | 72 | 0.2 | 1395 | 7.7 | 42 | 124 | 48 | 74 | 0.1 | 1310 | 60 | 60 | 93 | 66 | 43 | 88 | 78 | 74 | 99 |
| D-26/4/90 | 34528 | 1 | 7.9 | 262 | 444 | 324 | 48 | 6.5 | 1319 | 7.9 | 279 | 380 | 43 | 6 | 1362 | 7.9 | 161 | 347 | 90 | 78 | 0.4 | 1476 | 7.6 | 31 | 174 | 67 | 74 | 0.5 | 1398 | 42 | 76 | 94 | 81 | 50 | 88 | 61 | 79 | 92 |
| D-27/4/90 | 31417 | 3 | 7.8 | 237 | 549 | 236 | 71 | 4 | 1397 | 7.9 | 286 | 280 | 67 | 8 | 1368 | 7.8 | 176 | 416 | 112 | 70 | 0.4 | 1469 | 7.7 | 39 | 287 | 84 | 82 | 0.2 | 1480 | 39 | 60 | 96 | 78 | 31 | 84 | 48 | 64 | 95 |
| D-29/4/90 | 27333 | 2 | 7.6 | 238 | 348 | 174 | 64 | 3.5 | 1110 | 7.6 | 372 | 124 | 76 | 2.5 | 1105 | 7.4 | 172 | 364 | 104 | 79 | 1 | 1171 | 7.4 | 44 | 210 | 73 | 81 | 1.5 | 1256 | 54 | 16 | 60 | 74 | 42 | 82 | 40 | 58 | 57 |
| D-1/7/90 | 30201 | 1 | 7.3 | 137 | 398 | 188 | 57 | 4 | 1179 | 7.3 | 164 | 204 | 59 | 4.3 | 1164 | 7.3 | 81 | 204 | 62 | 65 | 0.1 | 1165 | 7.4 | 9 | 66 | 10 | 70 | 0 | 1264 | 51 | 70 | 98 | 89 | 68 | 93 | 93 | 95 | 100 |
| D-2/7/90 | 39445 | 0.8 | 8.1 | 187 | 488 | 216 | 56 | 6.5 | 2440 | 8 | 191 | 252 | 49 | 5 | 2340 | 7.9 | 130 | 304 | 66 | 73 | 0.2 | 2230 | 8 | 28 | 80 | 21 | 71 | 0 | 2190 | 32 | 74 | 96 | 79 | 74 | 85 | 84 | 90 | 100 |
| D-3/7/90 | 37252 | 3.5 | 7.6 | 131 | 436 | 476 | 28 | 4.5 | 1603 | 7.8 | 151 | 416 | 29 | 5 | 1479 | 7.8 | 114 | 269 | 78 | 56 | 0.1 | 1591 | 7.9 | 13 | 90 | 16 | 69 | 0 | 1625 | 25 | 81 | 98 | 89 | 67 | 90 | 79 | 97 | 100 |
| D-4/7/90 | 37643 | 1.3 | 7.7 | 145 | 382 | 168 | 55 | 3.5 | 1814 | 7.9 | 180 | 184 | 57 | 3 | 1820 | 8.1 | 113 | 271 | 70 | 66 | 0.2 | 1893 | 8.1 | 15 | 88 | 15 | 73 | 0 | 1879 | 37 | 62 | 95 | 87 | 68 | 90 | 77 | 91 | 100 |
| D-5/7/90 | 36389 | 0.9 | 7.7 | 156 | 391 | 156 | 54 | 3.5 | 1358 | 7.8 | 161 | 172 | 58 | 3 | 1412 | 7.8 | 116 | 285 | 66 | 73 | 0.3 | 1460 | 7.9 | 11 | 86 | 14 | 50 | 0 | 1425 | 28 | 62 | 92 | 91 | 70 | 93 | 78 | 91 | 100 |
| D-6/7/90 | 33020 | 0.8 | 7.7 | 176 | 422 | 176 | 61 | 4 | 2200 | 7.7 | 176 | 196 | 59 | 3.5 | 2180 | 7.7 | 120 | 274 | 72 | 64 | 0.3 | 2230 | 7.8 | 14 | 82 | 16 | 69 | 0 | 2140 | 32 | 63 | 91 | 88 | 70 | 92 | 81 | 91 | 100 |
| D-8/7/90 | 36095 | 4.5 | 7.5 | 112 | 341 | 256 | 47 | 4.2 | 2070 | 7.4 | 101 | 238 | 48 | 3 | 2430 | 7.5 | 95 | 251 | 82 | 76 | 0.2 | 1930 | 7.6 | 18 | 114 | 27 | 74 | 0 | 2700 | 5.9 | 66 | 93 | 81 | 55 | 84 | 67 | 90 | 100 |
| D-9/7/90 | 39590 | 5 | 7.8 | 144 | 361 | 220 | 54 | 4.5 | 1790 | 7.8 | 112 | 236 | 54 | 3.3 | 1760 | 7.8 | 123 | 239 | 78 | 64 | 0.2 | 1710 | 7.6 | 14 | 74 | 20 | 70 | 0 | 1660 | 53 | 67 | 94 | 74 | 69 | 90 | 80 | 91 | 100 |
| D-10/7/90 | 42859 | 13 | 7.8 | 161 | 326 | 268 | 45 | 4.5 | 1570 | 7.9 | 152 | 208 | 55 | 5 | 1640 | 7.9 | 76 | 173 | 74 | 65 | 0.2 | 1690 | 7.8 | 23 | 111 | 26 | 77 | 0 | 1720 | 50 | 64 | 97 | 70 | 36 | 86 | 66 | 90 | 100 |
| D-11/7/90 | 36325 | 2.5 | 7.7 | 135 | 304 | 222 | 47 | 3.5 | 1087 | 7.7 | 135 | 250 | 55 | 3.5 | 1143 | 7.8 | 107 | 260 | 96 | 73 | 0.2 | 1162 | 7.9 | 14 | 96 | 18 | 69 | 0 | 1169 | 21 | 62 | 94 | 87 | 63 | 90 | 68 | 92 | 100 |
| D-12/7/90 | 33522 | 2.5 | 7.7 | 189 | 416 | 248 | 58 | 5.5 | 1330 | 7.7 | 194 | 252 | 64 | 4.5 | 1931 | 7.8 | 114 | 304 | 104 | 71 | 0.2 | 1426 | 7.8 | 27 | 108 | 26 | 77 | 0 | 1416 | 41 | 59 | 96 | 76 | 65 | 86 | 74 | 90 | 100 |
| D-13/7/90 | 33680 | 1.8 | 7.6 | 215 | 532 | 282 | 64 | 5 | 1444 | 7.6 | 253 | 296 | 62 | 5 | 1437 | 7.8 | 120 | 320 | 90 | 76 | 0.2 | 1418 | 7.8 | 17 | 92 | 15 | 75 | 0 | 1433 | 53 | 70 | 96 | 86 | 71 | 92 | 83 | 95 | 100 |
| D-15/7/90 | 31293 | 1 | 7.7 | 150 | 337 | 182 | 64 | 3 | 1189 | 7.7 | 117 | 190 | 68 | 2.5 | 1199 | 7.7 | 108 | 256 | 62 | 75 | 0.5 | 1282 | 7.8 | 18 | 85 | 20 | 76 | 0 | 1378 | 7.7 | 67 | 80 | 83 | 67 | 88 | 75 | 89 | 100 |
| D-16/7/90 | 36010 | 1 | 7.9 | 189 | 424 | 264 | 55 | 4 | 2270 | 7.8 | 190 | 318 | 55 | 5.5 | 2310 | 8 | 119 | 250 | 90 | 67 | 0.4 | 2130 | 8 | 16 | 68 | 21 | 78 | 0 | 1848 | 37 | 72 | 94 | 87 | 73 | 92 | 84 | 92 | 100 |
| D-17/7/90 | 35445 | 2.5 | 7.9 | 182 | 388 | 272 | 43 | 3.5 | 1704 | 7.9 | 166 | 262 | 48 | 4 | 1731 | 7.9 | 110 | 274 | 82 | 51 | 0.2 | 1734 | 8 | 18 | 71 | 18 | 53 | 0 | 1788 | 34 | 69 | 96 | 84 | 74 | 90 | 82 | 93 | 100 |
| D-18/7/90 | 32540 | 1 | 7.7 | 290 | 566 | 352 | 53 | 8 | 2540 | 7.7 | 314 | 358 | 54 | 8 | 2730 | 7.7 | 122 | 304 | 78 | 90 | 0.2 | 2750 | 7.9 | 18 | 116 | 15 | 93 | 0 | 2580 | 61 | 78 | 98 | 85 | 62 | 94 | 80 | 96 | 100 |
| D-19/7/90 | 32929 | 1.2 | 7.8 | 246 | 641 | 376 | 56 | 6 | 1683 | 7.8 | 265 | 378 | 58 | 8 | 1724 | 7.9 | 132 | 338 | 170 | 64 | 1.7 | 1710 | 7.9 | 31 | 100 | 34 | 79 | 0 | 1742 | 50 | 55 | 79 | 77 | 70 | 87 | 84 | 91 | 100 |
| D-20/7/90 | 32575 | 2 | 7.7 | 244 | 607 | 380 | 54 | 8 | 2160 | 7.8 | 242 | 436 | 54 | 8 | 2220 | 7.9 | 139 | 353 | 164 | 46 | 1.8 | 2120 | 8 | 18 | 104 | 26 | 72 | 0 | 2270 | 43 | 62 | 78 | 87 | 71 | 93 | 83 | 93 | 100 |
| D-22/7/90 | 30019 | 4.5 | 7.6 | 175 | 416 | 224 | 52 | 4.5 | 1463 | 7.7 | 183 | 214 | 57 | 3.5 | 2120 | 7.7 | 105 | 274 | 104 | 63 | 0.8 | 2250 | 7.9 | 18 | 84 | 21 | 67 | 0.1 | 2610 | 43 | 51 | 77 | 83 | 70 | 90 | 79 | 91 | 99 |
| D-23/7/90 | 27711 | 1.5 | 7.4 | 342 | 400 | 232 | 65 | 4.5 | 1858 | 7.6 | 259 | 284 | 61 | 6 | 1907 | 7.6 | 133 | 278 | 86 | 72 | 0.2 | 1835 | 7.6 | 23 | 106 | 38 | 77 | 0 | 1511 | 49 | 70 | 97 | 83 | 62 | 93 | 74 | 84 | 100 |
| D-24/7/90 | 33999 | 1.5 | 7.5 | 191 | 488 | 506 | 36 | 6.3 | 2590 | 7.4 | 172 | 580 | 41 | 7.8 | 2530 | 7.6 | 102 | 267 | 84 | 71 | 0.2 | 2670 | 7.7 | 22 | 110 | 33 | 75 | 0.1 | 2850 | 41 | 86 | 97 | 78 | 59 | 89 | 78 | 94 | 99 |
| D-25/7/90 | 33959 | 2 | 7.7 | 235 | 614 | 692 | 30 | 8 | 1589 | 7.7 | 242 | 716 | 32 | 9.5 | 1632 | 7.7 | 94 | 246 | 244 | 24 | 0.2 | 1588 | 7.9 | 16 | 90 | 52 | 29 | 0 | 1625 | 61 | 66 | 98 | 83 | 63 | 93 | 85 | 93 | 100 |

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-26/7/90 | 33290 | 1 | 7.7 | 128 | 392 | 236 | 59 | 3.2 | 2100 | 7.6 | 130 | 176 | 57 | 4 | 2200 | 7.7 | 118 | 272 | 88 | 59 | 0.2 | 2010 | 7.8 | 18 | 164 | 49 | 67 | 0.1 | 2070 | 9.2 | 50 | 95 | 85 | 40 | 90 | 58 | 79 | 98 |
| D-27/7/90 | 33877 | 1.5 | 7.6 | 166 | 452 | 320 | 46 | 4.5 | 1160 | 7.7 | 221 | 388 | 47 | 5.5 | 1137 | 7.7 | 87 | 216 | 76 | 74 | 0.2 | 1147 | 7.8 | 19 | 56 | 28 | 75 | 0.1 | 1148 | 61 | 80 | 97 | 78 | 74 | 89 | 88 | 91 | 98 |
| D-29/7/90 | 26871 | 2 | 7.3 | 139 | 369 | 284 | 49 | 4 | 991 | 7.5 | 161 | 312 | 45 | 4 | 986 | 7.3 | 87 | 198 | 90 | 62 | 0.1 | 978 | 7.5 | 19 | 116 | 34 | 73 | 0 | 1070 | 46 | 71 | 98 | 78 | 41 | 86 | 69 | 88 | 100 |
| D-30/7/90 | 37634 | 2 | 8.1 | 151 | 400 | 268 | 48 | 3.5 | 1732 | 8 | 160 | 248 | 60 | 3 | 1803 | 8 | 85 | 274 | 84 | 69 | 0.2 | 1721 | 8.2 | 17 | 84 | 17 | 78 | 0 | 1643 | 47 | 66 | 93 | 86 | 70 | 92 | 79 | 94 | 100 |
| D-31/7/90 | 32909 | 0.4 | 7.5 | 221 | 391 | 260 | 46 | 4.3 | 2140 | 7.5 | 173 | 244 | 53 | 4.5 | 2090 | 7.5 | 112 | 261 | 74 | 78 | 0.2 | 2000 | 7.6 | 17 | 127 | 20 | 78 | 0 | 1950 | 35 | 70 | 97 | 85 | 51 | 92 | 68 | 92 | 100 |
| D-2/9/90 | 44601 | 1.7 | 7.5 | 140 | 231 | 360 | 64 | 4 | 806 | 7.5 | 146 | 328 | 63 | 3.5 | 811 | 7.5 | 79 | 145 | 64 | 75 | 0.2 | 834 | 7.6 | 12 | 39 | 11 | 81 | 0 | 905 | 46 | 81 | 94 | 85 | 73 | 91 | 83 | 97 | 100 |
| D-3/9/90 | 43614 | 6.5 | 7.8 | 143 | 400 | 472 | 64 | 6.5 | 2230 | 7.8 | 95 | 226 | 63 | 4 | 2160 | 7.8 | 97 | 235 | 74 | 75 | 0.3 | 2170 | 7.8 | 14 | 82 | 12 | 81 | 0 | 1910 | 40 | 67 | 94 | 86 | 65 | 90 | 80 | 98 | 100 |
| D-4/9/90 | 40529 | 4.5 | 7.7 | 192 | 253 | 186 | 64 | 5 | 1818 | 7.6 | 250 | 194 | 63 | 4.5 | 1745 | 7.5 | 108 | 197 | 80 | 75 | 0.3 | 1625 | 7.5 | 14 | 48 | 13 | 81 | 0 | 1726 | 57 | 59 | 94 | 87 | 76 | 93 | 81 | 93 | 100 |
| D-5/9/90 | 38231 | 0.8 | 7.7 | 124 | 234 | 120 | 64 | 1.8 | 1550 | 7.9 | 123 | 220 | 63 | 4 | 1651 | 7.8 | 87 | 156 | 76 | 75 | 0.2 | 1586 | 7.9 | 18 | 84 | 74 | 81 | 0.2 | 1661 | 29 | 66 | 95 | 79 | 46 | 86 | 64 | 38 | 89 |
| D-6/9/90 | 36909 | 1.3 | 7.6 | 172 | 394 | 226 | 51 | 5.5 | 1324 | 7.6 | 197 | 302 | 47 | 4 | 1306 | 7.7 | 119 | 257 | 92 | 65 | 0.2 | 1368 | 7.8 | 14 | 74 | 13 | 71 | 0 | 1418 | 40 | 70 | 96 | 88 | 71 | 92 | 81 | 94 | 100 |
| D-7/9/90 | 37002 | 1.7 | 7.6 | 268 | 324 | 226 | 53 | 4 | 1390 | 7.6 | 265 | 236 | 53 | 4.5 | 1389 | 7.6 | 110 | 179 | 98 | 69 | 0.2 | 1352 | 7.7 | 13 | 97 | 15 | 69 | 0 | 1358 | 59 | 59 | 96 | 88 | 46 | 95 | 70 | 93 | 100 |
| D-9/9/90 | 43377 | 6.5 | 7.6 | 176 | 200 | 176 | 50 | 3 | 859 | 7.6 | 149 | 172 | 51 | 2 | 860 | 7.6 | 90 | 178 | 76 | 66 | 0.2 | 926 | 7.7 | 13 | 45 | 18 | 69 | 0 | 1020 | 40 | 56 | 90 | 86 | 75 | 93 | 78 | 90 | 99 |
| D-11/9/90 | 37862 | 5 | 7.6 | 146 | 243 | 164 | 57 | 3.2 | 1095 | 7.6 | 193 | 208 | 56 | 3 | 1089 | 7.6 | 68 | 165 | 86 | 77 | 0.3 | 1095 | 7.7 | 6 | 59 | 14 | 83 | 0 | 1095 | 65 | 59 | 92 | 91 | 64 | 96 | 76 | 92 | 100 |
| D-12/9/90 | 35809 | 5 | 7.8 | 139 | 380 | 256 | 53 | 5.5 | 2040 | 7.8 | 201 | 340 | 48 | 8.5 | 2000 | 7.9 | 131 | 278 | 104 | 73 | 0.3 | 1950 | 7.9 | 11 | 58 | 8 | 58 | 0 | 1690 | 35 | 69 | 98 | 92 | 82 | 92 | 87 | 97 | 100 |
| D-13/9/90 | 35729 | 2.5 | 7.9 | 205 | 529 | 264 | 61 | 6.5 | 2750 | 7.9 | 220 | 324 | 61 | 6 | 2850 | 7.8 | 111 | 325 | 90 | 73 | 0.3 | 2550 | 7.8 | 11 | 84 | 9 | 62 | 0 | 2500 | 50 | 72 | 95 | 90 | 70 | 95 | 79 | 97 | 100 |
| D-14/9/90 | 41206 | 3.3 | 7.8 | 117 | 366 | 500 | 27 | 4.2 | 3230 | 7.8 | 135 | 620 | 32 | 7 | 3170 | 7.7 | 81 | 181 | 106 | 49 | 0.2 | 3690 | 7.8 | 12 | 67 | 12 | 75 | 0 | 3950 | 40 | 83 | 97 | 85 | 63 | 90 | 82 | 98 | 100 |
| D-16/9/90 | 15519 | 3.6 | 7.6 | 148 | 400 | 564 | 32 | 7 | 1554 | 7.6 | 235 | 798 | 29 | 11 | 1515 | 7.6 | 49 | 151 | 120 | 48 | 0.4 | 1259 | 7.7 | 11 | 54 | 9 | 80 | 0 | 1267 | 79 | 85 | 97 | 78 | 64 | 93 | 87 | 98 | 100 |
| D-17/9/90 | 49986 | 2.7 | 7.8 | 158 | 256 | 194 | 47 | 3 | 1831 | 7.8 | 218 | 240 | 47 | 4 | 1860 | 7.8 | 94 | 198 | 98 | 59 | 0.3 | 1862 | 7.7 | 10 | 66 | 9 | 76 | 0 | 1813 | 57 | 59 | 93 | 89 | 67 | 94 | 74 | 95 | 100 |
| D-18/9/90 | 51575 | 2.4 | 7.7 | 123 | 246 | 186 | 40 | 2.5 | 1288 | 7.7 | 241 | 300 | 25 | 4.5 | 1328 | 7.7 | 101 | 202 | 86 | 37 | 0.2 | 1287 | 7.7 | 12 | 70 | 15 | 56 | 0 | 1317 | 58 | 71 | 96 | 88 | 65 | 90 | 72 | 92 | 99 |
| D-19/9/90 | 44869 | 0.9 | 7.7 | 133 | 410 | 204 | 63 | 4.5 | 1490 | 7.7 | 161 | 256 | 58 | 4.5 | 1421 | 7.7 | 69 | 251 | 88 | 68 | 0.3 | 1386 | 7.7 | 8 | 144 | 11 | 80 | 0 | 1361 | 57 | 66 | 94 | 88 | 43 | 94 | 65 | 95 | 100 |
| D-20/9/90 | 43491 | 1 | 7.5 | 237 | 388 | 178 | 57 | 4 | 1620 | 7.6 | 300 | 240 | 53 | 4.5 | 1571 | 7.7 | 113 | 243 | 84 | 69 | 0.2 | 1631 | 7.8 | 12 | 87 | 6 | 87 | 0 | 1590 | 62 | 65 | 96 | 89 | 69 | 95 | 80 | 93 | 100 |
| D-21/9/90 | 45453 | 1.6 | 7.7 | 203 | 357 | 222 | 59 | 3 | 1828 | 7.6 | 215 | 208 | 73 | 5 | 1817 | 7.6 | 98 | 169 | 86 | 65 | 0.2 | 1890 | 7.6 | 9 | 73 | 13 | 81 | 0 | 1868 | 54 | 59 | 96 | 91 | 57 | 96 | 80 | 94 | 100 |
| D-23/9/90 | 37190 | 1.1 | 7.8 | 199 | 380 | 244 | 56 | 5 | 1453 | 7.8 | 110 | 196 | 60 | 3.5 | 1430 | 7.8 | 104 | 204 | 116 | 72 | 1.5 | 1436 | 7.9 | 13 | 43 | 15 | 75 | 0 | 1458 | 5.5 | 41 | 57 | 88 | 79 | 94 | 89 | 94 | 100 |
| D-24/9/90 | 40067 | 3.4 | 7.8 | 214 | 941 | 304 | 53 | 6.5 | 1454 | 7.8 | 396 | 320 | 58 | 7 | 1528 | 7.8 | 107 | 243 | 106 | 68 | 0.2 | 1456 | 7.9 | 11 | 55 | 14 | 74 | 0 | 1324 | 73 | 67 | 98 | 90 | 77 | 95 | 94 | 95 | 100 |
| D-25/9/90 | 57606 | 5.4 | 7.7 | 208 | 298 | 260 | 38 | 4 | 1236 | 7.8 | 245 | 284 | 47 | 5 | 1285 | 7.8 | 90 | 220 | 124 | 48 | 0.3 | 1416 | 7.9 | 19 | 84 | 14 | 83 | 0 | 1611 | 63 | 56 | 94 | 79 | 70 | 91 | 79 | 95 | 100 |
| D-26/9/90 | 46791 | 2.2 | 7.9 | 206 | 305 | 188 | 56 | 4 | 1200 | 7.9 | 197 | 272 | 50 | 6 | 1228 | 7.9 | 113 | 210 | 88 | 66 | 0.2 | 1237 | 7.9 | 14 | 55 | 10 | 92 | 0 | 1084 | 40 | 68 | 98 | 88 | 74 | 93 | 82 | 95 | 100 |
| D-27/9/90 | 46852 | 1.5 | 7.7 | 247 | 361 | 286 | 51 | 4.5 | 1180 | 7.7 | 146 | 212 | 56 | 2.5 | 1251 | 7.7 | 135 | 248 | 134 | 52 | 0.6 | 1144 | 7.9 | 11 | 56 | 10 | 78 | 0 | 1222 | 7.5 | 37 | 76 | 92 | 77 | 96 | 85 | 97 | 100 |
| D-28/9/90 | 38761 | 4 | 7.5 | 438 | 681 | 370 | 57 | 6 | 1396 | 7.6 | 212 | 348 | 64 | 6 | 1446 | 7.6 | 148 | 350 | 88 | 77 | 0.2 | 1356 | 7.6 | 15 | 120 | 11 | 71 | 0 | 1311 | 30 | 75 | 97 | 90 | 66 | 97 | 82 | 97 | 100 |
| D-30/9/90 | 42046 | 2 | 7.8 | 255 | 282 | 166 | 68 | 2 | 950 | 7.8 | 284 | 172 | 65 | 2 | 965 | 7.8 | 164 | 259 | 112 | 77 | 0.5 | 971 | 7.9 | 13 | 56 | 20 | 80 | 0 | 1093 | 42 | 35 | 75 | 92 | 78 | 95 | 80 | 88 | 99 |
| D-1/8/90 | 33322 | 0.4 | 7.6 | 217 | 315 | 400 | 59 | 3.5 | 2080 | 7.5 | 225 | 376 | 65 | 4 | 2090 | 7.4 | 131 | 215 | 67 | 78 | 0.1 | 2190 | 7.5 | 22 | 81 | 28 | 71 | 0.1 | 2170 | 42 | 82 | 98 | 83 | 62 | 90 | 74 | 93 | 99 |
| D-2/8/90 | 10050 | 0.4 | 7.6 | 208 | 556 | 210 | 52 | 4 | 2340 | 7.5 | 244 | 148 | 66 | 2.5 | 2080 | 7.4 | 95 | 205 | 73 | 75 | 0.1 | 2070 | 7.6 | 18 | 75 | 34 | 72 | 0 | 2080 | 61 | 51 | 96 | 81 | 63 | 91 | 87 | 84 | 100 |
| D-3/8/90 | 55930 | 1.2 | 7.8 | 223 | 459 | 364 | 46 | 5.5 | 2220 | 8 | 220 | 346 | 41 | 4.5 | 2300 | 7.9 | 85 | 199 | 63 | 75 | 0.1 | 2240 | 8 | 19 | 56 | 23 | 78 | 0.1 | 2350 | 61 | 82 | 98 | 78 | 72 | 92 | 88 | 94 | 99 |
| D-5/8/90 | 52851 | 0.3 | 8.1 | 64 | 161 | 172 | 41 | 2 | 1350 | 8 | 59 | 146 | 45 | 1.1 | 1315 | 7.9 | 44 | 114 | 63 | 64 | 0.1 | 292 | 8.1 | 8 | 31 | 14 | 80 | 0 | 1400 | 25 | 57 | 91 | 82 | 73 | 88 | 81 | 92 | 100 |
| D-6/8/90 | 40585 | 0.4 | 8 | 66 | 152 | 364 | 20 | 1.6 | 1403 | 8.1 | 86 | 412 | 20 | 2.5 | 1383 | 8 | 49 | 114 | 80 | 46 | 0.2 | 1407 | 8 | 11 | 24 | 12 | 77 | 0 | 1433 | 43 | 81 | 92 | 78 | 79 | 83 | 84 | 97 | 100 |
| D-7/8/90 | 45027 | 2.5 | 7.7 | 48 | 156 | 242 | 26 | 1.4 | 1400 | 7.9 | 63 | 266 | 25 | 1.9 | 1403 | 7.8 | 37 | 112 | 84 | 48 | 0.2 | 1363 | 8 | 5 | 60 | 12 | 83 | 0 | 1407 | 41 | 68 | 90 | 87 | 46 | 90 | 62 | 95 | 100 |
| D-8/8/90 | 47338 | 2 | 8.1 | 69 | 126 | 138 | 39 | 1.4 | 882 | 8.1 | 150 | 180 | 39 | 1.9 | 902 | 8.2 | 53 | 100 | 62 | 60 | 0.1 | 920 | 8 | 8 | 52 | 9 | 89 | 0 | 902 | 65 | 66 | 95 | 85 | 48 | 88 | 59 | 94 | 100 |
| D-9/8/90 | 44207 | 1.2 | 7.8 | 70 | 188 | 112 | 48 | 1.3 | 1031 | 7.9 | 125 | 194 | 41 | 2.8 | 1074 | 7.9 | 44 | 101 | 54 | 61 | 0.2 | 1064 | 8 | 12 | 53 | 16 | 63 | 0 | 1056 | 65 | 72 | 95 | 73 | 48 | 83 | 72 | 86 | 100 |
| D-10/8/90 | 43563 | 1 | 7.8 | 95 | 206 | 132 | 44 | 1.5 | 828 | 7.8 | 74 | 112 | 57 | 1.2 | 856 | 7.9 | 47 | 158 | 62 | 65 | 0.1 | 85 | 7.9 | 21 | 99 | 39 | 67 | 0.1 | 891 | 37 | 45 | 92 | 75 | 37 | 78 | 52 | 71 | 97 |
| D-12/8/90 | 47718 | 0.7 | 7.8 | 31 | 81 | 208 | 20 | 0.8 | 715 | 7.8 | 32 | 246 | 25 | 1.5 | 714 | 7.8 | 119 | 80 | 233 | 20 | 0.7 | 712 | 7.9 | 3 | 25 | 12 | 73 | 0 | 831 | 40 | 5.3 | 53 | 85 | 69 | 90 | 69 | 94 | 100 |
| D-13/8/90 | 42587 | 2 | 7.8 | 77 | 256 | 248 | 31 | 2.5 | 910 | 7.9 | 71 | 270 | 29 | 3 | 931 | 7.9 | 26 | 120 | 71 | 45 | 0.3 | 917 | 8 | 7 | 62 | 9 | 93 | 0 | 919 | 63 | 74 | 90 | 73 | 48 | 91 | 76 | 96 | 100 |
| D-15/8/90 | 35098 | 0.8 | 7.8 | 100 | 256 | 234 | 38 | 2.5 | 993 | 7.8 | 97 | 168 | 50 | 1.3 | 953 | 7.8 | 119 | 140 | 95 | 53 | 0.4 | 1034 | 8 | 11 | 58 | 10 | 76 | 0 | 985 | 40 | 44 | 68 | 85 | 59 | 89 | 77 | 96 | 100 |
| D-16/8/90 | 38052 | 0.8 | 7.6 | 94 | 409 | 194 | 58 | 2.5 | 997 | 7.6 | 197 | 192 | 54 | 2.5 | 996 | 7.7 | 119 | 194 | 103 | 57 | 0.3 | 988 | 7.9 | 11 | 65 | 11 | 69 | 0 | 990 | 40 | 46 | 72 | 85 | 67 | 85 | 84 | 94 | 99 |
| D-17/8/90 | 31404 | 0.8 | 8 | 183 | 321 | 160 | 68 | 2.5 | 1096 | 7.9 | 197 | 164 | 61 | 2.5 | 1136 | 7.9 | 119 | 179 | 85 | 66 | 0.3 | 1082 | 8.1 | 11 | 98 | 15 | 67 | 0 | 1061 | 40 | 48 | 88 | 85 | 45 | 90 | 70 | 91 | 100 |
| D-19/8/90 | 38905 | 0.3 | 7.7 | 58 | 197 | 130 | 65 | 2 | 1135 | 7.7 | 85 | 146 | 62 | 2.5 | 1169 | 7.8 | 44 | 123 | 54 | 80 | 0.2 | 1141 | 7.9 | 20 | 42 | 16 | 68 | 0 | 1143 | 48 | 63 | 94 | 55 | 66 | 66 | 79 | 88 | 99 |
| D-20/8/90 | 38620 | 0.7 | 7.5 | 95 | 302 | 176 | 59 | 3 | 1120 | 7.6 | 94 | 152 | 59 | 2.5 | 1145 | 7.6 | 48 | 316 | 107 | 65 | 1.2 | 1155 | 7.7 | 19 | 59 | 15 | 77 | 0.1 | 1040 | 49 | 30 | 52 | 60 | 81 | 80 | 81 | 92 | 98 |
| D-21/8/90 | 34352 | 0.3 | 7.4 | 112 | 470 | 172 | 65 | 4.5 | 1207 | 7.4 | 100 | 192 | 65 | 4 | 1208 | 7.4 | 74 | 213 | 69 | 73 | 0.2 | 1175 | 7.6 | 14 | 97 | 15 | 72 | 0 | 1162 | 26 | 64 | 95 | 81 | 55 | 88 | 79 | 91 | 100 |
| D-22/8/90 | 34785 | 1 | 7.5 | 126 | 397 | 188 | 67 | 5 | 1950 | 7.5 | 121 | 190 | 63 | 3.5 | 1760 | 7.6 | 75 | 190 | 68 | 72 | 0.2 | 1720 | 7.7 | 16 | 101 | 15 | 61 | 0 | 1570 | 38 | 64 | 96 | 79 | 47 | 87 | 75 | 92 | 100 |
| D-23/8/90 | 27109 | 0.4 | 7.6 | 158 | 276 | 142 | 65 | 1.5 | 1939 | 7.6 | 205 | 278 | 56 | 7.5 | 1878 | 7.6 | 102 | 269 | 123 | 66 | 1 | 1920 | 7.7 | 24 | 77 | 21 | 59 | 0 | 1959 | 50 | 56 | 87 | 77 | 71 | 85 | 72 | 85 | 98 |
| D-24/8/90 | 32802 | 0.7 | 7.3 | 203 | 405 | 212 | 68 | 4.5 | 1922 | 7.3 | 197 | 320 | 61 | 8.5 | 1874 | 7.4 | 97 | 225 | 74 | 78 | 0.3 | 1820 | 7.4 | 22 | 114 | 24 | 71 | 0.1 | 1796 | 40 | 77 | 97 | 77 | 49 | 89 | 72 | 89 | 98 |
| D-25/8/90 | 35876 | 0.5 | 7.9 | 81 | 448 | 296 | 64 | 4.5 | 1315 | 7.7 | 128 | 380 | 63 | 4.5 | 1355 | 7.8 | 65 | 319 | 70 | 75 | 0.3 | 1401 | 7.9 | 27 | 123 | 33 | 81 | 0 | 1525 | 49 | 82 | 93 | 59 | 61 | 67 | 73 | 89 | 100 |
| D-27/8/90 | 41410 | 0.6 | 7.9 | 85 | 269 | 110 | 64 | 1.5 | 1342 | 7.8 | 139 | 198 | 63 | 4 | 1325 | 7.8 | 65 | 168 | 65 | 75 | 0.2 | 1362 | 7.9 | 15 | 84 | 23 | 81 | 0 | 1375 | 53 | 67 | 95 | 77 | 70 | 82 | 79 | 79 | 100 |
| D-28/8/90 | 40933 | 1.5 | 7.8 | 120 | 303 | 290 | 64 | 4.5 | 1818 | 7.8 | 125 | 250 | 63 | 4 | 1800 | 7.8 | 78 | 192 | 68 | 75 | 0.2 | 1846 | 7.6 | 10 | 76 | 10 | 81 | 0 | 1723 | 38 | 73 | 95 | 87 | 60 | 92 | 75 | 97 | 100 |
| D-29/8/90 | 34764 | 0.9 | 7.5 | 127 | 284 | 188 | 64 | 2.8 | 2260 | 7.6 | 143 | 222 | 63 | 3.5 | 2500 | 7.6 | 114 | 192 | 92 | 75 | 0.4 | 2140 | 7.7 | 18 | 62 | 14 | 81 | 0 | 2100 | 20 | 59 | 90 | 85 | 68 | 90 | 78 | 93 | 100 |
| D-30/8/90 | 39489 | 0.9 | 7.8 | 131 | 320 | 166 | 64 | 2.5 | 1680 | 7.8 | 135 | 196 | 63 | 2.5 | 1690 | 7.8 | 90 | 214 | 88 | 75 | 0.2 | 1551 | 7.8 | 15 | 74 | 14 | 81 | 0 | 1672 | 33 | 55 | 92 | 83 | 65 | 89 | 77 | 92 | 99 |
| D-31/8/90 | 42230 | 0.7 | 8.1 | 132 | 288 | 144 | 64 | 1.3 | 1581 | 8.3 | 143 | 196 | 63 | 3 | 1705 | 8.3 | 80 | 184 | 74 | 75 | 0.2 | 1717 | 8.3 | 11 | 64 | 13 | 81 | 0 | 1720 | 44 | 62 | 93 | 86 | 65 | 92 | 78 | 91 | 100 |
| D-2/12/90 | 29388 | 0.8 | 8.2 | 339 | 713 | 356 | 73 | 10 | 2170 | 8 | 334 | 584 | 84 | 8.5 | 2110 | 7.9 | 223 | 372 | 116 | 85 | 1.3 | 2140 | 7.6 | 17 | 83 | 12 | 72 | 0 | 2330 | 33 | 80 | 85 | 92 | 78 | 95 | 88 | 97 | 100 |
| D-3/12/90 | 30935 | 2 | 8.1 | 255 | 550 | 214 | 69 | 5.5 | 1919 | 8.2 | 227 | 246 | 68 | 8 | 1917 | 8.1 | 161 | 372 | 100 | 78 | 0.2 | 1815 | 7.8 | 18 | 87 | 22 | 82 | 0 | 1690 | 29 | 59 | 98 | 89 | 77 | 93 | 84 | 90 | 100 |
| D-4/12/90 | 26348 | 1 | 8.1 | 253 | 473 | 212 | 79 | 7 | 1990 | 8.1 | 250 | 246 | 75 | 6 | 1950 | 8.1 | 202 | 384 | 112 | 82 | 0.2 | 2220 | 7.7 | 23 | 105 | 22 | 89 | 0 | 2130 | 19 | 55 | 97 | 89 | 73 | 91 | 78 | 90 | 100 |
| D-6/12/90 | 28680 | 1.7 | 8.1 | 256 | 539 | 188 | 75 | 4.7 | 1803 | 8.1 | 337 | 218 | 75 | 5.2 | 1732 | 8.1 | 186 | 365 | 102 | 84 | 0.5 | 1804 | 7.8 | 27 | 116 | 23 | 87 | 0 | 1932 | 45 | 53 | 90 | 86 | 68 | 90 | 79 | 88 | 100 |

| D-7/12/90 | 32799 | 5.9 | 8.4 | 215 | 440 | 190 | 70 | 4.5 | 1380 | 8.3 | 228 | 228 | 68 | 4.5 | 1408 | 8.1 | 113 | 213 | 78 | 77 | 0.3 | 1362 | 7.8 | 9 | 60 | 13 | 77 | 0 | 1379 | 50 | 66 | 93 | 92 | 72 | 96 | 86 | 93 | 100 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-9/12/90 | 33545 | 1.3 | 8.2 | 145 | 747 | 310 | 56 | 6 | 1059 | 8.2 | 181 | 264 | 53 | 4.7 | 1056 | 8.1 | 118 | 353 | 110 | 66 | 1 | 1082 | 7.7 | 15 | 100 | 18 | 78 | 0 | 1210 | 35 | 58 | 79 | 87 | 72 | 90 | 87 | 94 | 100 |
| D-10/12/90 | 28791 | 1.1 | 8.7 | 354 | 539 | 232 | 68 | 5.5 | 1620 | 8.4 | 197 | 178 | 71 | 2.5 | 1677 | 8.4 | 114 | 356 | 130 | 79 | 1.5 | 1529 | 7.9 | 14 | 67 | 17 | 82 | 0 | 1359 | 40 | 27 | 40 | 88 | 81 | 96 | 88 | 93 | 100 |
| D-11/12/90 | 27219 | 1.5 | 8.3 | 302 | 566 | 212 | 70 | 5.7 | 2270 | 8.3 | 319 | 216 | 69 | 6.3 | 2240 | 8.1 | 119 | 475 | 144 | 74 | 1.7 | 2200 | 7.7 | 18 | 84 | 19 | 81 | 0 | 1432 | 40 | 33 | 73 | 85 | 70 | 90 | 79 | 91 | 100 |
| D-12/12/90 | 31849 | 5.5 | 8 | 330 | 511 | 184 | 79 | 4 | 2110 | 8 | 372 | 212 | 75 | 5 | 2100 | 8.1 | 204 | 396 | 128 | 78 | 1 | 2050 | 7.7 | 34 | 119 | 33 | 81 | 0 | 2270 | 45 | 40 | 80 | 83 | 70 | 90 | 77 | 82 | 100 |
| D-13/12/90 | 30352 | 5.9 | 8.2 | 324 | 539 | 196 | 78 | 4.5 | 2090 | 8.2 | 282 | 190 | 73 | 4.5 | 2180 | 8.1 | 203 | 396 | 110 | 86 | 1 | 2130 | 7.8 | 21 | 89 | 19 | 90 | 0 | 2210 | 28 | 42 | 78 | 90 | 78 | 94 | 84 | 90 | 100 |
| D-14/12/90 | 32009 | 2.9 | 8.1 | 288 | 459 | 216 | 73 | 4.5 | 1755 | 8 | 399 | 258 | 71 | 4.5 | 1690 | 7.9 | 183 | 364 | 104 | 77 | 0 | 1729 | 7.6 | 18 | 75 | 18 | 84 | 0 | 1750 | 54 | 60 | 93 | 90 | 79 | 94 | 84 | 92 | 100 |
| D-16/12/90 | 34492 | 0.4 | 7.9 | 253 | 354 | 162 | 77 | 3.5 | 1443 | 8 | 232 | 202 | 71 | 4 | 1455 | 8.2 | 162 | 287 | 96 | 81 | 0.5 | 1520 | 7.9 | 16 | 69 | 12 | 90 | 0 | 1689 | 30 | 53 | 88 | 90 | 76 | 94 | 81 | 93 | 100 |
| D-17/12/90 | 36452 | 0.5 | 8.2 | 154 | 310 | 128 | 72 | 3 | 1520 | 8.2 | 164 | 148 | 70 | 3 | 1350 | 8.1 | 126 | 237 | 82 | 81 | 0.2 | 1403 | 7.7 | 15 | 50 | 11 | 91 | 0 | 1280 | 23 | 45 | 93 | 88 | 79 | 90 | 84 | 91 | 100 |
| D-18/12/90 | 34361 | 1.2 | 8.2 | 172 | 345 | 152 | 74 | 4.3 | 1552 | 8.2 | 286 | 156 | 74 | 3.5 | 1500 | 8.2 | 140 | 310 | 74 | 84 | 0.2 | 1736 | 7.8 | 16 | 85 | 14 | 94 | 0 | 1679 | 51 | 53 | 94 | 89 | 73 | 91 | 75 | 91 | 100 |
| D-19/12/90 | 36432 | 18 | 8.1 | 168 | 334 | 196 | 65 | 7 | 2230 | 8.1 | 188 | 188 | 66 | 4.5 | 2110 | 8.2 | 133 | 210 | 74 | 78 | 0.2 | 2190 | 7.8 | 20 | 74 | 20 | 85 | 0 | 2320 | 29 | 61 | 96 | 85 | 65 | 88 | 78 | 90 | 100 |
| D-20/12/90 | 37009 | 1.5 | 8.3 | 181 | 348 | 184 | 85 | 4.5 | 1337 | 8.3 | 278 | 200 | 78 | 5.5 | 1479 | 8.3 | 124 | 264 | 76 | 90 | 0.2 | 1544 | 7.8 | 21 | 48 | 19 | 91 | 0 | 1605 | 55 | 62 | 96 | 83 | 82 | 88 | 86 | 90 | 100 |
| D-21/12/90 | 37281 | 1.5 | 8.1 | 287 | 484 | 308 | 77 | 5.8 | 1653 | 8.1 | 392 | 308 | 75 | 5.5 | 1640 | 8.1 | 155 | 332 | 118 | 83 | 0.3 | 1777 | 7.8 | 22 | 80 | 16 | 93 | 0.1 | 1880 | 61 | 62 | 95 | 86 | 76 | 92 | 84 | 95 | 99 |
| D-23/12/90 | 28437 | 1 | 7.9 | 176 | 387 | 162 | 75 | 3.8 | 1344 | 7.9 | 178 | 156 | 80 | 3 | 1331 | 7.9 | 155 | 312 | 100 | 84 | 0.6 | 1420 | 7.7 | 17 | 51 | 14 | 71 | 0 | 1468 | 13 | 36 | 80 | 89 | 84 | 90 | 87 | 91 | 100 |
| D-24/12/90 | 29955 | 0.9 | 7.6 | 203 | 301 | 146 | 69 | 4.5 | 1299 | 7.8 | 229 | 176 | 65 | 4.5 | 1319 | 7.9 | 135 | 242 | 108 | 72 | 0.3 | 1375 | 7.6 | 18 | 51 | 32 | 79 | 0 | 1430 | 41 | 39 | 93 | 85 | 79 | 90 | 83 | 78 | 100 |
| D-26/12/90 | 35263 | 0.3 | 8 | 201 | 434 | 118 | 73 | 5 | 1727 | 8.1 | 300 | 214 | 73 | 4.5 | 1700 | 8.1 | 131 | 317 | 114 | 67 | 0.9 | 1749 | 7.7 | 12 | 55 | 18 | 67 | 0 | 1750 | 56 | 47 | 80 | 91 | 83 | 94 | 87 | 85 | 100 |
| D-27/12/90 | 34319 | 0.5 | 8.1 | 236 | 448 | 178 | 76 | 5 | 1325 | 8 | 241 | 180 | 77 | 4.5 | 1259 | 7.9 | 150 | 226 | 90 | 78 | 0.5 | 1344 | 7.8 | 16 | 79 | 17 | 85 | 0 | 1276 | 38 | 50 | 89 | 89 | 65 | 93 | 82 | 90 | 100 |
| D-28/12/90 | 32730 | 1 | 7.8 | 290 | 376 | 146 | 84 | 4.5 | 1441 | 7.8 | 222 | 154 | 86 | 4.5 | 1460 | 7.7 | 150 | 242 | 72 | 83 | 0 | 1446 | 7.5 | 18 | 89 | 19 | 93 | 0 | 1491 | 32 | 53 | 93 | 85 | 63 | 90 | 76 | 87 | 100 |
| D-30/12/90 | 30164 | 0.5 | 7.9 | 232 | 503 | 194 | 78 | 6 | 1200 | 7.9 | 407 | 220 | 80 | 5.5 | 1222 | 8 | 149 | 297 | 112 | 80 | 0.4 | 1167 | 7.8 | 18 | 84 | 14 | 94 | 0 | 1399 | 63 | 49 | 94 | 85 | 70 | 90 | 79 | 93 | 100 |
| D-1/11/90 | 45006 | 5.2 | 8 | 183 | 182 | 134 | 61 | 2.5 | 1007 | 8 | 197 | 190 | 51 | 3 | 1041 | 8 | 119 | 137 | 76 | 84 | 0.2 | 1106 | 7.8 | 18 | 42 | 11 | 76 | 0 | 1117 | 40 | 60 | 93 | 85 | 69 | 90 | 77 | 92 | 100 |
| D-2/11/90 | 44158 | 8.1 | 8 | 124 | 463 | 230 | 48 | 4.5 | 1599 | 7.9 | 123 | 230 | 49 | 3 | 1586 | 7.9 | 94 | 251 | 88 | 64 | 0.3 | 1631 | 7.6 | 27 | 90 | 20 | 70 | 0 | 1554 | 24 | 62 | 93 | 71 | 64 | 78 | 81 | 91 | 100 |
| D-4/11/90 | 39223 | 1.8 | 8.1 | 151 | 294 | 186 | 57 | 4 | 1006 | 8 | 148 | 154 | 61 | 2.8 | 1000 | 7.9 | 94 | 204 | 90 | 71 | 0.2 | 988 | 7.8 | 9 | 39 | 11 | 87 | 0 | 1038 | 37 | 42 | 93 | 90 | 81 | 94 | 87 | 94 | 100 |
| D-5/11/90 | 42394 | 1.5 | 8.2 | 241 | 688 | 270 | 53 | 4.5 | 1566 | 8.2 | 331 | 244 | 56 | 4.2 | 1571 | 8.2 | 125 | 294 | 108 | 61 | 0.2 | 1550 | 7.8 | 26 | 84 | 21 | 65 | 0 | 1409 | 62 | 56 | 95 | 79 | 70 | 89 | 79 | 92 | 100 |
| D-6/11/90 | 44364 | 2.5 | 8 | 182 | 364 | 154 | 69 | 3 | 1672 | 8.1 | 156 | 224 | 49 | 3.5 | 1750 | 8 | 136 | 261 | 86 | 65 | 0.2 | 1694 | 7.8 | 20 | 83 | 17 | 78 | 0 | 1683 | 13 | 62 | 94 | 85 | 68 | 89 | 77 | 89 | 100 |
| D-7/11/90 | 44235 | 1.9 | 7.9 | 195 | 428 | 210 | 70 | 6.5 | 1835 | 8.2 | 194 | 250 | 57 | 3.5 | 1873 | 8.2 | 121 | 289 | 104 | 67 | 0.2 | 1855 | 8 | 15 | 75 | 22 | 73 | 0 | 1820 | 38 | 58 | 94 | 88 | 74 | 92 | 83 | 90 | 100 |
| D-8/11/90 | 45151 | 1.7 | 8.3 | 185 | 457 | 184 | 70 | 4 | 1944 | 8.3 | 161 | 228 | 61 | 5.5 | 2050 | 8.1 | 109 | 325 | 98 | 76 | 0.4 | 2090 | 7.8 | 18 | 98 | 17 | 75 | 0.1 | 1824 | 32 | 57 | 93 | 84 | 70 | 90 | 79 | 91 | 99 |
| D-9/11/90 | 47032 | 1.5 | 8.2 | 139 | 294 | 256 | 41 | 3 | 1456 | 8.1 | 140 | 260 | 39 | 2.5 | 1361 | 8 | 81 | 172 | 90 | 56 | 0.2 | 1504 | 7.7 | 11 | 55 | 14 | 71 | 0 | 1751 | 42 | 65 | 92 | 86 | 68 | 92 | 81 | 95 | 100 |
| D-11/11/90 | 41372 | 1.6 | 7.8 | 119 | 220 | 128 | 63 | 1.5 | 1444 | 7.9 | 117 | 138 | 55 | 1.2 | 1493 | 7.9 | 91 | 161 | 76 | 68 | 0.3 | 1391 | 7.9 | 7 | 45 | 9 | 84 | 0 | 1431 | 22 | 45 | 75 | 92 | 72 | 94 | 80 | 93 | 99 |
| D-12/11/90 | 45729 | 4.5 | 8.2 | 139 | 263 | 144 | 61 | 1.5 | 1665 | 8.3 | 159 | 204 | 51 | 2.5 | 1717 | 8 | 85 | 208 | 80 | 83 | 0.2 | 1807 | 7.8 | 9 | 51 | 10 | 90 | 0 | 1698 | 47 | 61 | 92 | 89 | 76 | 94 | 81 | 93 | 99 |
| D-13/11/90 | 49314 | 2.3 | 8.1 | 166 | 318 | 118 | 56 | 2.5 | 1820 | 8 | 198 | 134 | 63 | 2 | 1779 | 8.1 | 95 | 213 | 60 | 70 | 0.2 | 1768 | 7.9 | 15 | 74 | 12 | 77 | 0 | 1810 | 52 | 55 | 90 | 84 | 65 | 91 | 77 | 90 | 100 |
| D-14/11/90 | 44038 | 0.4 | 8 | 138 | 304 | 136 | 66 | 2.5 | 2050 | 8.1 | 144 | 168 | 61 | 3 | 2000 | 8.2 | 107 | 243 | 68 | 90 | 0.2 | 1852 | 7.8 | 17 | 106 | 20 | 85 | 0 | 1854 | 26 | 60 | 93 | 84 | 56 | 88 | 65 | 85 | 100 |
| D-15/11/90 | 29816 | 0.5 | 8.5 | 251 | 447 | 152 | 66 | 3.3 | 1770 | 8.4 | 312 | 204 | 63 | 3.8 | 2880 | 8.2 | 192 | 439 | 108 | 72 | 0.9 | 2490 | 7.9 | 26 | 158 | 30 | 83 | 0 | 2240 | 39 | 47 | 76 | 87 | 64 | 90 | 65 | 80 | 100 |
| D-16/11/90 | 29448 | 0.1 | 8.2 | 192 | 357 | 286 | 79 | 4.5 | 2950 | 7.9 | 127 | 230 | 70 | 4.5 | 2980 | 7.7 | 158 | 304 | 92 | 89 | 0.3 | 2780 | 7.6 | 27 | 77 | 31 | 88 | 0 | 2790 | 40 | 60 | 93 | 83 | 75 | 86 | 78 | 89 | 100 |
| D-18/11/90 | 35825 | 1.4 | 8 | 111 | 371 | 142 | 73 | 4.5 | 1440 | 8 | 157 | 154 | 71 | 4 | 1437 | 7.9 | 109 | 197 | 80 | 90 | 0.2 | 1560 | 7.7 | 14 | 131 | 24 | 83 | 0 | 1757 | 31 | 48 | 96 | 87 | 34 | 87 | 65 | 83 | 100 |
| D-19/11/90 | 47384 | 0.8 | 8.2 | 147 | 224 | 146 | 66 | 4 | 1758 | 8.2 | 118 | 138 | 67 | 1.4 | 1740 | 8 | 85 | 106 | 74 | 70 | 0.2 | 1712 | 7.7 | 15 | 23 | 18 | 82 | 0 | 1707 | 28 | 46 | 86 | 82 | 78 | 90 | 90 | 88 | 100 |
| D-20/11/90 | 47306 | 0.1 | 8.2 | 206 | 324 | 166 | 64 | 3.5 | 1782 | 8 | 261 | 146 | 69 | 3 | 1830 | 8 | 105 | 220 | 84 | 71 | 0.2 | 1879 | 7.8 | 12 | 66 | 23 | 66 | 0.1 | 1833 | 60 | 43 | 93 | 89 | 70 | 94 | 80 | 86 | 99 |
| D-21/11/90 | 40127 | 7.9 | 8.1 | 233 | 456 | 214 | 67 | 6 | 1586 | 8.1 | 268 | 238 | 70 | 4.5 | 1672 | 8 | 114 | 260 | 74 | 81 | 0.2 | 1784 | 7.9 | 33 | 140 | 32 | 78 | 0 | 1839 | 58 | 69 | 96 | 71 | 46 | 86 | 69 | 85 | 100 |
| D-22/11/90 | 28005 | 3.5 | 8.2 | 226 | 419 | 178 | 71 | 4.5 | 2210 | 8.3 | 197 | 218 | 71 | 4 | 2150 | 8.2 | 141 | 318 | 100 | 80 | 0.2 | 2170 | 7.8 | 19 | 100 | 28 | 76 | 0 | 1914 | 28 | 54 | 91 | 87 | 69 | 92 | 76 | 84 | 99 |
| D-23/11/90 | 28819 | 2.2 | 8.4 | 195 | 392 | 188 | 70 | 4 | 2790 | 8.2 | 218 | 210 | 67 | 4 | 2680 | 8 | 130 | 267 | 82 | 78 | 0.3 | 2550 | 7.9 | 18 | 86 | 19 | 79 | 0 | 2520 | 40 | 61 | 94 | 86 | 68 | 91 | 78 | 90 | 100 |
| D-25/11/90 | 27098 | 2.1 | 8.2 | 197 | 497 | 254 | 73 | 6.5 | 1685 | 8.2 | 197 | 274 | 66 | 7 | 1681 | 8.2 | 130 | 307 | 132 | 73 | 1.3 | 1697 | 7.7 | 10 | 70 | 19 | 74 | 0 | 1894 | 40 | 52 | 82 | 92 | 77 | 95 | 86 | 93 | 100 |
| D-26/11/90 | 42667 | 2.7 | 8.3 | 173 | 427 | 208 | 71 | 5.5 | 1979 | 8.4 | 221 | 204 | 71 | 5 | 1955 | 8.3 | 113 | 291 | 82 | 83 | 0.2 | 1862 | 7.7 | 16 | 105 | 20 | 86 | 0 | 1789 | 49 | 60 | 97 | 86 | 64 | 91 | 75 | 90 | 100 |
| D-27/11/90 | 47222 | 2.2 | 8.1 | 148 | 321 | 142 | 73 | 3 | 1681 | 8.1 | 194 | 242 | 67 | 7 | 1704 | 8 | 130 | 274 | 104 | 75 | 0.8 | 1653 | 7.8 | 12 | 78 | 14 | 83 | 0 | 1705 | 33 | 57 | 89 | 91 | 72 | 92 | 76 | 90 | 100 |
| D-28/11/90 | 32157 | 0.6 | 8.1 | 243 | 572 | 200 | 78 | 4.5 | 2280 | 8.2 | 289 | 240 | 74 | 6.5 | 2340 | 8 | 166 | 372 | 98 | 82 | 0.2 | 2320 | 7.8 | 13 | 86 | 13 | 89 | 0 | 2230 | 43 | 59 | 98 | 92 | 77 | 95 | 85 | 94 | 100 |
| D-29/11/90 | 25687 | 1.2 | 8.2 | 316 | 743 | 280 | 72 | 6.5 | 2390 | 7.8 | 363 | 260 | 71 | 7 | 2460 | 8.1 | 169 | 408 | 112 | 77 | 0.3 | 2380 | 7.7 | 16 | 98 | 16 | 78 | 0 | 2210 | 53 | 57 | 96 | 91 | 76 | 95 | 87 | 94 | 100 |
| D-30/11/90 | 26040 | 2.1 | 8.1 | 302 | 702 | 244 | 75 | 8.5 | 2480 | 8 | 357 | 304 | 71 | 12 | 2250 | 8 | 194 | 412 | 114 | 83 | 0.2 | 2640 | 7.8 | 19 | 118 | 18 | 83 | 0 | 2590 | 46 | 63 | 98 | 90 | 71 | 94 | 83 | 93 | 100 |
| D-1/10/90 | 47623 | 3.4 | 7.7 | 283 | 310 | 170 | 61 | 2.5 | 1065 | 7.8 | 235 | 270 | 51 | 5.5 | 1100 | 7.8 | 100 | 227 | 108 | 67 | 0.5 | 1090 | 7.9 | 16 | 85 | 20 | 80 | 0 | 993 | 53 | 60 | 92 | 86 | 63 | 94 | 73 | 88 | 99 |
| D-2/10/90 | 54578 | 3.6 | 7.9 | 313 | 341 | 512 | 34 | 6 | 915 | 7.9 | 205 | 520 | 34 | 7.5 | 976 | 7.9 | 107 | 200 | 106 | 64 | 0.2 | 1002 | 7.9 | 13 | 74 | 16 | 80 | 0 | 1021 | 48 | 80 | 97 | 88 | 63 | 96 | 78 | 97 | 100 |
| D-3/10/90 | 36911 | 6 | 7.7 | 300 | 610 | 452 | 45 | 8 | 1313 | 7.7 | 375 | 556 | 44 | 10 | 1312 | 7.7 | 150 | 323 | 212 | 49 | 2.5 | 1238 | 7.7 | 21 | 116 | 14 | 86 | 0 | 1133 | 60 | 62 | 75 | 86 | 64 | 93 | 81 | 97 | 100 |
| D-4/10/90 | 35244 | 6 | 7.8 | 177 | 412 | 196 | 69 | 4.5 | 2190 | 7.8 | 330 | 476 | 52 | 7 | 2330 | 7.8 | 151 | 274 | 124 | 61 | 0.3 | 2350 | 7.7 | 8 | 90 | 17 | 77 | 0 | 2220 | 54 | 74 | 97 | 95 | 67 | 96 | 78 | 91 | 100 |
| D-5/10/90 | 39566 | 5.8 | 7.7 | 192 | 416 | 236 | 58 | 4.5 | 1447 | 7.6 | 296 | 384 | 54 | 7 | 1445 | 7.7 | 148 | 270 | 108 | 67 | 0.4 | 1380 | 7.8 | 13 | 78 | 14 | 71 | 0 | 1408 | 50 | 72 | 94 | 91 | 71 | 93 | 81 | 94 | 100 |
| D-7/10/90 | 45469 | 2 | 7.3 | 129 | 237 | 234 | 39 | 1.5 | 1190 | 7.4 | 123 | 228 | 42 | 2 | 1202 | 7.4 | 73 | 175 | 110 | 35 | 0.2 | 1309 | 7.5 | 7 | 66 | 19 | 63 | 0 | 1631 | 41 | 52 | 90 | 90 | 62 | 95 | 72 | 92 | 99 |
| D-8/10/90 | 46240 | 2.7 | 7.9 | 122 | 287 | 168 | 50 | 2 | 1906 | 7.8 | 157 | 188 | 53 | 2 | 2130 | 7.8 | 99 | 198 | 92 | 59 | 0.2 | 810 | 7.8 | 12 | 43 | 16 | 78 | 0.1 | 1425 | 37 | 51 | 93 | 88 | 78 | 90 | 85 | 91 | 98 |
| D-9/10/90 | 45903 | 1.3 | 8 | 90 | 271 | 144 | 67 | 4.5 | 1221 | 8 | 110 | 164 | 71 | 3 | 1209 | 7.8 | 78 | 179 | 74 | 73 | 0.2 | 1266 | 7.7 | 10 | 41 | 14 | 74 | 0.1 | 1283 | 29 | 55 | 95 | 87 | 77 | 89 | 85 | 90 | 99 |
| D-10/10/90 | 44343 | 1.3 | 8 | 169 | 327 | 126 | 68 | 2.5 | 1940 | 8.1 | 172 | 180 | 78 | 2.5 | 1950 | 7.8 | 106 | 272 | 78 | 80 | 0.2 | 1886 | 7.9 | 15 | 93 | 17 | 73 | 0 | 1891 | 38 | 57 | 92 | 86 | 66 | 91 | 72 | 87 | 99 |
| D-12/10/90 | 39343 | 4.8 | 8 | 140 | 555 | 282 | 45 | 6.5 | 850 | 8 | 134 | 208 | 48 | 3.5 | 890 | 7.9 | 44 | 160 | 76 | 61 | 0.2 | 906 | 7.8 | 7 | 75 | 10 | 72 | 0 | 882 | 67 | 64 | 94 | 84 | 53 | 95 | 87 | 97 | 100 |
| D-14/10/90 | 34347 | 1.8 | 7.9 | 155 | 243 | 210 | 48 | 2.5 | 848 | 7.9 | 170 | 138 | 52 | 1.5 | 858 | 9 | 89 | 133 | 74 | 76 | 0.3 | 902 | 8 | 9 | 35 | 12 | 58 | 0 | 898 | 48 | 46 | 83 | 90 | 74 | 94 | 86 | 94 | 99 |
| D-15/10/90 | 53012 | 3.5 | 8.2 | 157 | 361 | 208 | 51 | 5 | 1015 | 8.2 | 197 | 188 | 53 | 4 | 931 | 8.1 | 93 | 204 | 86 | 61 | 0.3 | 1019 | 7.9 | 12 | 51 | 17 | 59 | 0 | 967 | 53 | 54 | 94 | 87 | 75 | 92 | 86 | 92 | 100 |
| D-16/10/90 | 52258 | 1.5 | 8.4 | 195 | 246 | 172 | 61 | 3.5 | 1631 | 8.3 | 184 | 156 | 69 | 3.8 | 1564 | 8.2 | 92 | 165 | 86 | 72 | 0.4 | 1626 | 8 | 11 | 42 | 18 | 78 | 0 | 1635 | 50 | 45 | 89 | 88 | 75 | 94 | 83 | 90 | 100 |
| D-17/10/90 | 49493 | 3.7 | 8.1 | 102 | 269 | 156 | 53 | 2.5 | 1082 | 7.9 | 177 | 152 | 53 | 2.5 | 1180 | 8 | 68 | 165 | 76 | 71 | 0.2 | 1041 | 7.8 | 10 | 54 | 13 | 83 | 0 | 1079 | 62 | 50 | 92 | 85 | 67 | 90 | 80 | 92 | 99 |

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-18/10/90 | 46200 | 2.5 | 8.3 | 192 | 294 | 172 | 56 | 3.5 | 1983 | 8.3 | 157 | 150 | 61 | 2 | 1934 | 8.1 | 107 | 207 | 104 | 62 | 0.4 | 1886 | 7.9 | 15 | 69 | 20 | 84 | 0 | 1880 | 32 | 31 | 83 | 86 | 67 | 92 | 77 | 88 | 99 |
| D-19/10/90 | 46069 | 6 | 8.2 | 185 | 334 | 156 | 60 | 4 | 1987 | 8.3 | 217 | 162 | 64 | 4 | 2100 | 8.2 | 121 | 299 | 114 | 65 | 0.5 | 2050 | 7.9 | 16 | 65 | 11 | 80 | 0 | 2040 | 44 | 30 | 88 | 87 | 78 | 91 | 81 | 93 | 99 |
| D-21/10/90 | 44324 | 1.6 | 7.8 | 174 | 173 | 134 | 69 | 2.5 | 967 | 7.8 | 165 | 114 | 68 | 2 | 980 | 7.8 | 104 | 202 | 100 | 70 | 0.3 | 1013 | 7.8 | 21 | 81 | 28 | 76 | 0.1 | 1105 | 37 | 12 | 88 | 80 | 60 | 88 | 53 | 79 | 96 |
| D-22/10/90 | 48950 | 2.5 | 8.1 | 109 | 211 | 880 | 13 | 5.5 | 1745 | 7.9 | 154 | 1108 | 14 | 7 | 1534 | 8 | 46 | 111 | 118 | 44 | 0.4 | 1556 | 7.9 | 7 | 35 | 13 | 77 | 0 | 1495 | 70 | 89 | 95 | 85 | 69 | 94 | 83 | 99 | 100 |
| D-23/10/90 | 60017 | 3.5 | 7.9 | 120 | 284 | 518 | 18 | 4 | 1048 | 7.9 | 115 | 460 | 20 | 4 | 1046 | 8 | 92 | 215 | 134 | 43 | 0.5 | 1135 | 8 | 11 | 88 | 20 | 54 | 0 | 1171 | 20 | 71 | 89 | 88 | 59 | 91 | 69 | 96 | 100 |
| D-24/10/90 | 41569 | 3.3 | 7.9 | 115 | 357 | 334 | 34 | 3.5 | 1760 | 8 | 186 | 464 | 38 | 6 | 1858 | 8 | 92 | 230 | 110 | 55 | 0.5 | 1785 | 7.9 | 8 | 61 | 18 | 58 | 0 | 1595 | 51 | 76 | 93 | 91 | 74 | 93 | 83 | 95 | 99 |
| D-25/10/90 | 40915 | 2 | 8 | 244 | 400 | 194 | 51 | 3.2 | 1885 | 8 | 197 | 284 | 42 | 2.8 | 2080 | 7.9 | 141 | 251 | 140 | 53 | 0.6 | 2020 | 8 | 13 | 63 | 16 | 63 | 0 | 2030 | 40 | 51 | 79 | 91 | 75 | 95 | 79 | 92 | 99 |
| D-26/10/90 | 44858 | 1.4 | 8.1 | 380 | 318 | 194 | 56 | 2.7 | 1970 | 8.1 | 137 | 188 | 55 | 1.8 | 2080 | 8 | 99 | 239 | 104 | 64 | 0.4 | 2130 | 7.9 | 16 | 71 | 22 | 59 | 0 | 2190 | 28 | 45 | 78 | 84 | 70 | 96 | 78 | 89 | 99 |
| D-28/10/90 | 47576 | 6 | 8.1 | 183 | 188 | 128 | 67 | 2 | 1465 | 8.1 | 197 | 158 | 56 | 2 | 1489 | 7.9 | 119 | 176 | 80 | 80 | 0.3 | 1435 | 7.8 | 18 | 59 | 24 | 73 | 0 | 1566 | 40 | 49 | 85 | 85 | 67 | 90 | 69 | 81 | 99 |
| D-29/10/90 | 47501 | 1 | 8.2 | 183 | 384 | 180 | 64 | 3.5 | 1567 | 8.2 | 197 | 224 | 55 | 2.5 | 1512 | 8.1 | 119 | 270 | 110 | 64 | 0.3 | 1570 | 7.8 | 18 | 122 | 22 | 68 | 0 | 1463 | 40 | 51 | 88 | 85 | 55 | 90 | 68 | 88 | 100 |
| D-30/10/90 | 47506 | 0.5 | 8 | 183 | 307 | 196 | 60 | 4 | 1557 | 8 | 197 | 228 | 49 | 5 | 1575 | 8 | 119 | 233 | 96 | 63 | 0.3 | 1620 | 7.9 | 18 | 74 | 26 | 60 | 0 | 1639 | 40 | 58 | 94 | 85 | 68 | 90 | 76 | 87 | 100 |
| D-1/3/91 | 26343 | 2.1 | 7.7 | 275 | 553 | 528 | 44 | 7 | 1105 | 7.8 | 166 | 516 | 45 | 8.5 | 1174 | 7.9 | 84 | 265 | 158 | 60 | 1.4 | 1422 | 7.7 | 24 | 87 | 29 | 83 | 0.1 | 1782 | 49 | 69 | 84 | 71 | 67 | 91 | 84 | 95 | 99 |
| D-3/3/91 | 32884 | 0.8 | 7.8 | 169 | 337 | 132 | 76 | 2.5 | 1568 | 7.9 | 192 | 204 | 69 | 4.5 | 1526 | 7.9 | 137 | 332 | 126 | 78 | 1.8 | 1516 | 7.7 | 14 | 71 | 14 | 94 | 0 | 1450 | 29 | 38 | 60 | 90 | 79 | 92 | 79 | 89 | 99 |
| D-4/3/91 | 40745 | 1.5 | 7.7 | 156 | 547 | 166 | 66 | 6 | 1667 | 7.9 | 135 | 174 | 68 | 4 | 1730 | 7.9 | 118 | 333 | 100 | 74 | 1 | 1619 | 7.7 | 15 | 120 | 17 | 82 | 0 | 1576 | 13 | 43 | 75 | 87 | 64 | 90 | 78 | 90 | 100 |
| D-5/3/91 | 39804 | 4.9 | 7.9 | 209 | 416 | 196 | 68 | 5 | 1711 | 7.9 | 277 | 250 | 70 | 7 | 1774 | 7.9 | 138 | 348 | 110 | 78 | 0.6 | 1845 | 7.7 | 16 | 107 | 19 | 78 | 0 | 1813 | 50 | 56 | 91 | 88 | 69 | 92 | 74 | 90 | 100 |
| D-6/3/91 | 45804 | 2.2 | 7.9 | 170 | 380 | 146 | 80 | 4.5 | 1949 | 7.9 | 236 | 160 | 76 | 4.5 | 1957 | 7.9 | 161 | 328 | 82 | 90 | 0.6 | 2110 | 7.8 | 23 | 127 | 21 | 93 | 0 | 2130 | 32 | 49 | 87 | 86 | 61 | 87 | 67 | 86 | 100 |
| D-7/3/91 | 42289 | 3.3 | 7.7 | 150 | 289 | 182 | 74 | 2.5 | 1698 | 8 | 130 | 134 | 69 | 2 | 1718 | 8 | 76 | 202 | 68 | 74 | 0.5 | 1604 | 7.7 | 25 | 87 | 21 | 84 | 0 | 1694 | 42 | 49 | 75 | 67 | 57 | 83 | 70 | 89 | 100 |
| D-8/3/91 | 44548 | 1.4 | 7.9 | 218 | 272 | 216 | 41 | 4.5 | 1360 | 7.9 | 261 | 236 | 44 | 4.5 | 1328 | 7.8 | 154 | 213 | 100 | 64 | 0.3 | 1427 | 7.7 | 16 | 101 | 18 | 84 | 0 | 1514 | 41 | 58 | 93 | 90 | 53 | 93 | 63 | 92 | 100 |
| D-10/3/91 | 36792 | 0.7 | 7.9 | 168 | 303 | 148 | 66 | 2.5 | 1069 | 8 | 158 | 158 | 66 | 3 | 1070 | 7.9 | 92 | 248 | 74 | 87 | 0.4 | 1150 | 7.7 | 16 | 78 | 12 | 83 | 0 | 1244 | 42 | 53 | 88 | 83 | 69 | 91 | 74 | 92 | 100 |
| D-11/3/91 | 36410 | 1.6 | 8 | 173 | 396 | 242 | 57 | 6.8 | 1879 | 8.1 | 253 | 210 | 64 | 6.3 | 1916 | 7.9 | 85 | 341 | 108 | 67 | 0.7 | 1761 | 7.7 | 21 | 70 | 18 | 89 | 0 | 1605 | 66 | 49 | 89 | 75 | 80 | 88 | 82 | 93 | 100 |
| D-12/3/91 | 33988 | 2.5 | 7.8 | 161 | 391 | 148 | 73 | 4.5 | 1551 | 7.9 | 183 | 170 | 69 | 4 | 1505 | 7.9 | 151 | 313 | 86 | 89 | 0.9 | 1512 | 7.8 | 19 | 94 | 14 | 97 | 0 | 1439 | 18 | 49 | 78 | 87 | 70 | 88 | 76 | 91 | 100 |
| D-13/3/91 | 36479 | 3.8 | 8.1 | 192 | 422 | 176 | 68 | 4.5 | 2080 | 8 | 202 | 220 | 64 | 5.7 | 2100 | 7.9 | 132 | 367 | 118 | 80 | 1.2 | 2140 | 7.8 | 20 | 152 | 16 | 88 | 0 | 2290 | 35 | 46 | 79 | 85 | 59 | 90 | 64 | 91 | 100 |
| D-14/3/91 | 31592 | 4.7 | 8 | 217 | 549 | 248 | 68 | 7 | 1370 | 8.1 | 255 | 224 | 68 | 5 | 1483 | 8.1 | 154 | 406 | 102 | 80 | 0.6 | 1410 | 7.8 | 15 | 74 | 13 | 86 | 0 | 1406 | 40 | 55 | 88 | 90 | 82 | 93 | 87 | 95 | 100 |
| D-15/3/91 | 31789 | 2.1 | 7.9 | 215 | 549 | 234 | 72 | 5.5 | 2370 | 7.9 | 328 | 272 | 72 | 5 | 2410 | 7.9 | 204 | 428 | 112 | 79 | 1 | 2360 | 7.8 | 18 | 101 | 18 | 89 | 0 | 2300 | 38 | 59 | 80 | 91 | 76 | 92 | 82 | 92 | 100 |
| D-17/3/91 | 27968 | 1.6 | 8 | 202 | 471 | 242 | 64 | 5.5 | 1158 | 8 | 278 | 340 | 61 | 6.5 | 1160 | 7.9 | 118 | 277 | 106 | 74 | 0.4 | 1205 | 7.7 | 25 | 83 | 21 | 76 | 0 | 1285 | 58 | 69 | 95 | 79 | 70 | 88 | 82 | 91 | 100 |
| D-18/3/91 | 30853 | 1.2 | 8 | 179 | 455 | 168 | 70 | 4.3 | 1410 | 8 | 179 | 166 | 70 | 3 | 1347 | 7.9 | 119 | 376 | 104 | 79 | 0.1 | 1371 | 7.6 | 24 | 99 | 22 | 86 | 0 | 1203 | 40 | 37 | 97 | 85 | 74 | 87 | 78 | 87 | 100 |
| D-19/3/91 | 29815 | 1 | 8 | 224 | 400 | 160 | 69 | 4.5 | 2230 | 8.2 | 244 | 178 | 64 | 5 | 2270 | 8 | 168 | 392 | 96 | 77 | 0.6 | 1763 | 7.7 | 20 | 94 | 22 | 84 | 0 | 1754 | 31 | 46 | 88 | 88 | 76 | 91 | 77 | 86 | 100 |
| D-20/3/91 | 32578 | 2.5 | 7.8 | 223 | 443 | 184 | 70 | 4 | 1553 | 7.9 | 214 | 186 | 71 | 3.5 | 1670 | 7.9 | 181 | 345 | 86 | 84 | 0.1 | 1643 | 7.8 | 29 | 86 | 22 | 91 | 0 | 1665 | 15 | 54 | 97 | 84 | 75 | 87 | 81 | 88 | 100 |
| D-21/3/91 | 33784 | 1 | 7.9 | 166 | 392 | 186 | 65 | 7.5 | 2510 | 8 | 213 | 198 | 65 | 5 | 2470 | 8 | 149 | 357 | 110 | 75 | 1.2 | 2390 | 7.8 | 20 | 94 | 20 | 80 | 0 | 2420 | 30 | 44 | 76 | 87 | 74 | 88 | 76 | 89 | 100 |
| D-22/3/91 | 33029 | 3.4 | 7.9 | 226 | 624 | 264 | 75 | 4.5 | 2150 | 7.8 | 336 | 226 | 75 | 4.5 | 2070 | 7.7 | 164 | 348 | 82 | 95 | 0.3 | 2240 | 7.6 | 27 | 60 | 22 | 86 | 0 | 2210 | 51 | 64 | 93 | 84 | 83 | 88 | 90 | 92 | 100 |
| D-24/3/91 | 48657 | 1.7 | 7.7 | 72 | 180 | 122 | 49 | 2.5 | 990 | 7.7 | 98 | 122 | 46 | 1.8 | 986 | 7.6 | 61 | 180 | 70 | 63 | 0.2 | 923 | 7.6 | 21 | 72 | 18 | 83 | 0 | 1109 | 38 | 43 | 89 | 66 | 60 | 71 | 60 | 85 | 100 |
| D-25/3/91 | 45512 | 2.6 | 7.9 | 89 | 176 | 166 | 36 | 1.2 | 1059 | 8 | 97 | 218 | 32 | 1.9 | 1145 | 8 | 72 | 161 | 86 | 49 | 0.2 | 1060 | 7.7 | 15 | 84 | 14 | 86 | 0 | 1046 | 26 | 61 | 92 | 79 | 70 | 83 | 79 | 92 | 100 |
| D-26/3/91 | 44085 | 4.2 | 7.4 | 183 | 325 | 228 | 47 | 3.5 | 1229 | 7.7 | 125 | 218 | 51 | 3 | 1254 | 7.8 | 97 | 368 | 136 | 57 | 1.3 | 1343 | 7.7 | 15 | 94 | 18 | 78 | 0 | 1382 | 22 | 38 | 57 | 85 | 75 | 92 | 71 | 92 | 100 |
| D-27/3/91 | 40578 | 3.3 | 7.9 | 174 | 400 | 160 | 68 | 3.1 | 1433 | 7.9 | 203 | 146 | 66 | 3 | 1388 | 7.9 | 133 | 274 | 82 | 81 | 0.2 | 1472 | 7.8 | 17 | 84 | 18 | 80 | 0 | 1459 | 35 | 44 | 93 | 87 | 70 | 90 | 79 | 89 | 100 |
| D-29/3/91 | 34917 | 8.5 | 7.5 | 156 | 311 | 226 | 67 | 4.5 | 1075 | 7.5 | 256 | 204 | 65 | 4.5 | 1116 | 7.5 | 112 | 182 | 92 | 76 | 0.3 | 1141 | 7.4 | 12 | 29 | 15 | 80 | 0 | 1170 | 56 | 55 | 93 | 89 | 84 | 92 | 91 | 93 | 100 |
| D-31/3/91 | 32217 | 2 | 7.7 | 365 | 370 | 172 | 63 | 4.5 | 928 | 7.6 | 296 | 196 | 55 | 4.5 | 936 | 7.7 | 93 | 135 | 80 | 78 | 0.3 | 940 | 7.6 | 13 | 70 | 16 | 88 | 0 | 940 | 69 | 59 | 93 | 86 | 48 | 96 | 81 | 91 | 100 |
| D-1/2/91 | 38105 | 1.6 | 8.4 | 230 | 517 | 218 | 75 | 5 | 1645 | 8.5 | 269 | 212 | 74 | 4.5 | 1676 | 8.3 | 154 | 349 | 102 | 80 | 0.4 | 1658 | 8 | 16 | 82 | 13 | 92 | 0 | 1658 | 43 | 52 | 91 | 90 | 77 | 93 | 84 | 94 | 100 |
| D-3/2/91 | 30701 | 0.6 | 8.1 | 136 | 447 | 196 | 60 | 4.5 | 1105 | 8.3 | 159 | 188 | 61 | 3.5 | 1120 | 7.8 | 114 | 247 | 74 | 78 | 0.4 | 1155 | 7.6 | 29 | 106 | 27 | 85 | 0 | 1312 | 28 | 61 | 89 | 75 | 57 | 79 | 76 | 86 | 100 |
| D-4/2/91 | 34290 | 3 | 8.1 | 194 | 435 | 166 | 72 | 4 | 1770 | 8.1 | 213 | 260 | 70 | 5.5 | 1930 | 8.1 | 64 | 345 | 100 | 76 | 0.8 | 1880 | 7.7 | 14 | 51 | 13 | 85 | 0 | 1722 | 70 | 62 | 86 | 78 | 85 | 93 | 88 | 92 | 100 |
| D-5/2/91 | 35338 | 3.2 | 8 | 195 | 574 | 166 | 70 | 4.9 | 1563 | 8 | 156 | 188 | 68 | 4 | 1580 | 8 | 131 | 337 | 76 | 79 | 0.5 | 1695 | 7.7 | 17 | 83 | 16 | 88 | 0 | 1736 | 16 | 60 | 88 | 87 | 75 | 91 | 86 | 90 | 100 |
| D-6/2/91 | 37120 | 2.6 | 8 | 212 | 546 | 184 | 73 | 4.5 | 2240 | 8 | 312 | 182 | 75 | 4.3 | 2160 | 8 | 169 | 329 | 94 | 75 | 0.5 | 1940 | 7.8 | 16 | 91 | 14 | 93 | 0 | 1860 | 46 | 48 | 89 | 91 | 72 | 93 | 83 | 92 | 100 |
| D-7/2/91 | 35190 | 10 | 8.1 | 227 | 812 | 350 | 59 | 7.3 | 1512 | 8.1 | 323 | 358 | 60 | 7.3 | 1592 | 8.1 | 175 | 412 | 150 | 71 | 1.2 | 1524 | 7.8 | 21 | 91 | 20 | 90 | 0.1 | 1565 | 46 | 58 | 83 | 88 | 78 | 91 | 89 | 94 | 86 |
| D-8/2/91 | 33714 | 2.7 | 8.4 | 228 | 473 | 200 | 74 | 4.5 | 1542 | 8.3 | 224 | 280 | 69 | 4.5 | 1546 | 8.1 | 146 | 335 | 106 | 81 | 0.3 | 1658 | 7.7 | 31 | 117 | 22 | 91 | 0 | 1645 | 35 | 62 | 93 | 79 | 65 | 86 | 75 | 89 | 100 |
| D-10/2/91 | 29660 | 2 | 8 | 223 | 521 | 208 | 74 | 5.3 | 1585 | 8 | 230 | 200 | 75 | 5 | 1555 | 7.9 | 158 | 327 | 120 | 75 | 1.1 | 1597 | 9.7 | 21 | 84 | 17 | 88 | 0 | 1708 | 31 | 40 | 78 | 87 | 70 | 91 | 79 | 92 | 100 |
| D-11/2/91 | 31749 | 11 | 8.1 | 247 | 525 | 264 | 71 | 6.3 | 1606 | 8.2 | 230 | 262 | 68 | 7 | 1598 | 8.1 | 148 | 339 | 126 | 75 | 0.8 | 1564 | 7.8 | 24 | 81 | 21 | 86 | 0.1 | 1455 | 36 | 52 | 89 | 84 | 70 | 90 | 79 | 92 | 99 |
| D-12/2/91 | 32736 | 5.3 | 8.2 | 257 | 629 | 312 | 72 | 7.5 | 1823 | 8.2 | 259 | 264 | 72 | 6.3 | 1767 | 8.2 | 199 | 374 | 104 | 79 | 0.6 | 1871 | 7.8 | 28 | 108 | 22 | 91 | 0 | 1850 | 23 | 61 | 90 | 86 | 71 | 89 | 83 | 93 | 100 |
| D-13/2/91 | 34441 | 8.2 | 8.1 | 199 | 486 | 240 | 68 | 9 | 1605 | 8.1 | 259 | 352 | 68 | 13 | 1684 | 8.1 | 168 | 394 | 116 | 76 | 0.4 | 1679 | 7.7 | 31 | 147 | 28 | 86 | 0 | 1770 | 35 | 67 | 97 | 82 | 63 | 84 | 70 | 88 | 100 |
| D-14/2/91 | 32888 | 6.3 | 8.1 | 185 | 401 | 156 | 80 | 4 | 1513 | 8.2 | 170 | 224 | 73 | 5.5 | 1593 | 8.1 | 149 | 296 | 86 | 91 | 0.4 | 1653 | 7.8 | 24 | 97 | 22 | 91 | 0 | 1682 | 12 | 62 | 93 | 84 | 67 | 87 | 76 | 86 | 100 |
| D-15/2/91 | 34461 | 6.2 | 8.1 | 222 | 662 | 252 | 68 | 6.5 | 1734 | 8.2 | 250 | 310 | 71 | 7.1 | 1733 | 8.2 | 175 | 366 | 104 | 79 | 0.6 | 1789 | 7.8 | 26 | 85 | 25 | 96 | 0.1 | 1821 | 30 | 67 | 92 | 86 | 75 | 89 | 84 | 92 | 99 |
| D-17/2/91 | 34045 | 1.3 | 8 | 129 | 506 | 178 | 63 | 4 | 1180 | 8 | 109 | 186 | 62 | 3.8 | 1148 | 7.9 | 126 | 269 | 90 | 78 | 0.4 | 1255 | 7.6 | 22 | 54 | 19 | 90 | 0 | 1406 | 40 | 52 | 91 | 83 | 80 | 83 | 89 | 89 | 100 |
| D-18/2/91 | 36421 | 3 | 8.2 | 240 | 507 | 356 | 61 | 9 | 1672 | 8.2 | 225 | 300 | 61 | 8 | 1663 | 8.1 | 135 | 293 | 94 | 77 | 0.4 | 1622 | 7.7 | 22 | 78 | 20 | 86 | 0 | 1454 | 40 | 69 | 95 | 84 | 73 | 91 | 85 | 94 | 100 |
| D-19/2/91 | 37662 | 2.6 | 8.1 | 238 | 490 | 170 | 71 | 5 | 1652 | 8.2 | 305 | 244 | 67 | 5.5 | 1643 | 8.2 | 135 | 371 | 102 | 71 | 1 | 1772 | 8 | 23 | 95 | 20 | 85 | 0 | 1800 | 56 | 58 | 82 | 83 | 74 | 90 | 81 | 88 | 100 |
| D-21/2/91 | 29990 | 7.6 | 7.9 | 204 | 469 | 212 | 76 | 5 | 1145 | 7.9 | 252 | 246 | 71 | 6.5 | 1186 | 8.1 | 154 | 257 | 90 | 84 | 0.3 | 1172 | 7.8 | 18 | 20 | 14 | 89 | 0 | 1249 | 39 | 63 | 96 | 88 | 92 | 91 | 96 | 93 | 100 |
| D-22/2/91 | 37561 | 13 | 7.5 | 183 | 400 | 196 | 73 | 4.5 | 1245 | 7.6 | 370 | 194 | 75 | 4.5 | 1185 | 7.7 | 212 | 230 | 110 | 75 | 0.3 | 1163 | 7.6 | 15 | 99 | 14 | 91 | 0 | 1207 | 43 | 43 | 93 | 93 | 57 | 90 | 79 | 93 | 100 |
| D-24/2/91 | 27340 | 4.2 | 7.7 | 254 | 411 | 174 | 74 | 4.5 | 1339 | 7.9 | 375 | 292 | 67 | 8 | 1350 | 7.9 | 121 | 319 | 102 | 78 | 0.5 | 1423 | 7.7 | 17 | 96 | 19 | 86 | 0 | 1431 | 68 | 65 | 94 | 86 | 70 | 93 | 77 | 89 | 100 |
| D-25/2/91 | 30055 | 5.7 | 7.8 | 230 | 422 | 222 | 71 | 5 | 1733 | 7.9 | 379 | 400 | 65 | 10 | 1790 | 8 | 194 | 346 | 122 | 82 | 0.9 | 1743 | 7.7 | 27 | 106 | 24 | 85 | 0 | 1629 | 49 | 70 | 91 | 86 | 69 | 88 | 75 | 89 | 100 |
| D-26/2/91 | 31494 | 2.5 | 8 | 290 | 459 | 228 | 72 | 7.5 | 1684 | 8 | 449 | 298 | 71 | 9 | 1746 | 8 | 187 | 352 | 112 | 91 | 0.6 | 1828 | 7.8 | 28 | 57 | 22 | 87 | 0 | 1740 | 58 | 62 | 93 | 85 | 84 | 90 | 88 | 90 | 100 |
| D-27/2/91 | 31765 | 2 | 7.9 | 240 | 887 | 314 | 76 | 7 | 2120 | 7.9 | 158 | 292 | 74 | 7.2 | 2200 | 7.9 | 167 | 511 | 136 | 75 | 0.7 | 2190 | 7.7 | 26 | 95 | 23 | 92 | 0 | 2170 | 40 | 53 | 91 | 84 | 81 | 89 | 89 | 93 | 100 |
| D-28/2/91 | 25342 | 2 | 8.3 | 275 | 494 | 212 | 78 | 4.5 | 2070 | 8.3 | 411 | 322 | 73 | 10 | 2200 | 8.2 | 167 | 371 | 118 | 88 | 0.5 | 2240 | 7.8 | 18 | 138 | 22 | 81 | 0 | 2230 | 40 | 60 | 93 | 85 | 70 | 90 | 79 | 91 | 100 |

| ID | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-1/1/91 | 32441 | 0.8 | 8.1 | 198 | 351 | 134 | 76 | 2.5 | 1341 | 8.1 | 264 | 142 | 73 | 2.5 | 1300 | 8.1 | 138 | 234 | 98 | 84 | 0.3 | 1355 | 7.9 | 15 | 77 | 17 | 85 | 0 | 1394 | 48 | 31 | 88 | 89 | 67 | 92 | 78 | 87 | 100 |
| D-2/1/91 | 40740 | 1.2 | 8.2 | 127 | 444 | 172 | 73 | 3 | 1658 | 8.2 | 117 | 184 | 72 | 4 | 1629 | 8.1 | 66 | 315 | 94 | 83 | 0.5 | 1644 | 7.8 | 15 | 85 | 15 | 90 | 0 | 1491 | 44 | 49 | 88 | 77 | 73 | 88 | 81 | 91 | 100 |
| D-3/1/91 | 34637 | 3.5 | 8.1 | 231 | 368 | 162 | 77 | 4 | 1854 | 8.1 | 273 | 158 | 75 | 4.5 | 1846 | 8.1 | 158 | 308 | 76 | 87 | 0.4 | 1822 | 7.8 | 18 | 92 | 15 | 83 | 0 | 1734 | 42 | 52 | 92 | 89 | 70 | 92 | 75 | 91 | 100 |
| D-4/1/91 | 34322 | 3 | 8.4 | 249 | 600 | 216 | 78 | 5 | 1566 | 8.3 | 220 | 204 | 75 | 5.3 | 1583 | 8.1 | 161 | 328 | 86 | 77 | 0.4 | 1643 | 7.7 | 18 | 88 | 12 | 83 | 0.1 | 1655 | 27 | 58 | 92 | 89 | 73 | 93 | 85 | 94 | 99 |
| D-6/1/91 | 35111 | 0.5 | 7.9 | 228 | 461 | 126 | 75 | 2.5 | 1113 | 8.1 | 243 | 142 | 76 | 2.5 | 1104 | 8 | 122 | 279 | 92 | 78 | 0.5 | 1061 | 7.8 | 14 | 89 | 13 | 65 | 0 | 1142 | 50 | 35 | 80 | 89 | 68 | 94 | 81 | 90 | 99 |
| D-7/1/91 | 40585 | 0.8 | 8.3 | 230 | 524 | 158 | 71 | 5 | 1178 | 8.2 | 285 | 186 | 67 | 5.5 | 1217 | 8.1 | 117 | 304 | 82 | 71 | 0.3 | 1178 | 7.8 | 13 | 92 | 13 | 73 | 0 | 1114 | 59 | 56 | 95 | 89 | 70 | 94 | 82 | 92 | 100 |
| D-8/1/91 | 37735 | 1.7 | 7.9 | 203 | 487 | 168 | 77 | 3.5 | 1353 | 8.1 | 231 | 196 | 70 | 5.5 | 1338 | 8.1 | 107 | 341 | 60 | 87 | 0.2 | 1250 | 7.8 | 17 | 166 | 15 | 90 | 0 | 1268 | 54 | 69 | 96 | 84 | 51 | 92 | 66 | 91 | 100 |
| D-9/1/91 | 34277 | 1.5 | 8.2 | 133 | 531 | 220 | 74 | 7 | 1514 | 8.2 | 207 | 288 | 65 | 11 | 1535 | 8 | 92 | 257 | 72 | 72 | 0.4 | 1476 | 7.7 | 11 | 83 | 17 | 74 | 0 | 1429 | 56 | 75 | 97 | 88 | 68 | 92 | 84 | 92 | 100 |
| D-10/1/91 | 41451 | 0.5 | 8 | 131 | 384 | 154 | 74 | 4.5 | 1260 | 8.1 | 189 | 298 | 75 | 6.5 | 1300 | 8 | 114 | 273 | 80 | 80 | 0.3 | 1367 | 7.8 | 18 | 95 | 19 | 84 | 0 | 1443 | 40 | 73 | 96 | 84 | 65 | 86 | 75 | 88 | 100 |
| D-11/1/91 | 45183 | 0.4 | 7.8 | 205 | 347 | 142 | 68 | 4.5 | 1373 | 7.7 | 339 | 148 | 69 | 4.5 | 1312 | 7.8 | 101 | 219 | 72 | 82 | 0.3 | 1409 | 7.7 | 19 | 100 | 19 | 85 | 0 | 1473 | 70 | 51 | 93 | 81 | 54 | 91 | 71 | 87 | 100 |
| D-13/1/91 | 27415 | 0.9 | 7.7 | 241 | 332 | 158 | 66 | 4 | 1133 | 7.8 | 182 | 194 | 67 | 3.3 | 1148 | 7.9 | 118 | 256 | 90 | 80 | 0.3 | 1190 | 7.7 | 12 | 156 | 11 | 71 | 0.1 | 1236 | 35 | 54 | 91 | 90 | 39 | 95 | 53 | 93 | 99 |
| D-14/1/91 | 42614 | 0.4 | 8.3 | 113 | 347 | 120 | 72 | 2.5 | 1385 | 8.2 | 106 | 134 | 75 | 2.5 | 1377 | 8.1 | 98 | 279 | 68 | 88 | 0.4 | 1380 | 7.8 | 7 | 84 | 11 | 91 | 0 | 1355 | 7.5 | 49 | 84 | 93 | 70 | 94 | 76 | 91 | 100 |
| D-15/1/91 | 48914 | 0.4 | 8.1 | 203 | 434 | 128 | 70 | 4 | 1496 | 8.3 | 222 | 114 | 72 | 2.5 | 1533 | 8.2 | 134 | 271 | 84 | 71 | 1 | 1489 | 7.7 | 13 | 60 | 11 | 76 | 0 | 1520 | 40 | 26 | 60 | 90 | 78 | 94 | 86 | 91 | 100 |
| D-16/1/91 | 49174 | 0.7 | 7.9 | 188 | 434 | 126 | 67 | 2.5 | 2450 | 8.1 | 159 | 156 | 68 | 3 | 2680 | 8.2 | 139 | 422 | 90 | 71 | 1 | 2950 | 7.8 | 22 | 155 | 28 | 79 | 0.1 | 2740 | 13 | 42 | 67 | 84 | 63 | 88 | 64 | 78 | 98 |
| D-17/1/91 | 45151 | 0.5 | 8.2 | 166 | 307 | 118 | 80 | 4 | 1697 | 8.2 | 214 | 130 | 80 | 3.5 | 1661 | 8.2 | 148 | 295 | 96 | 77 | 1.2 | 1624 | 7.7 | 19 | 116 | 27 | 88 | 0 | 1821 | 31 | 26 | 66 | 87 | 61 | 89 | 62 | 77 | 100 |
| D-18/1/91 | 37143 | 0.7 | 8.2 | 222 | 493 | 162 | 75 | 5.3 | 1620 | 8.3 | 211 | 162 | 73 | 4 | 1746 | 8.2 | 140 | 330 | 86 | 84 | 0.6 | 1834 | 7.7 | 16 | 105 | 11 | 82 | 0 | 1734 | 34 | 47 | 85 | 89 | 68 | 93 | 79 | 93 | 100 |
| D-20/1/91 | 30244 | 0.9 | 7.6 | 224 | 337 | 146 | 70 | 3.5 | 1146 | 7.8 | 240 | 168 | 71 | 4 | 1144 | 7.9 | 134 | 302 | 92 | 78 | 0.4 | 1205 | 7.7 | 17 | 78 | 8 | 80 | 0 | 1283 | 44 | 45 | 90 | 87 | 74 | 92 | 77 | 95 | 100 |
| D-21/1/91 | 34032 | 1 | 8.2 | 223 | 470 | 166 | 72 | 5.5 | 1343 | 8.3 | 273 | 196 | 76 | 5 | 1362 | 8.2 | 161 | 337 | 92 | 74 | 0.5 | 1419 | 7.8 | 16 | 86 | 13 | 39 | 0 | 1311 | 41 | 53 | 90 | 90 | 75 | 93 | 82 | 92 | 100 |
| D-22/1/91 | 34904 | 1.2 | 8 | 272 | 665 | 188 | 73 | 6 | 2200 | 8.1 | 293 | 192 | 74 | 5 | 2200 | 8 | 164 | 304 | 96 | 88 | 0.6 | 2130 | 7.7 | 20 | 76 | 15 | 80 | 0 | 2150 | 44 | 50 | 88 | 88 | 75 | 93 | 89 | 92 | 100 |
| D-23/1/91 | 36063 | 1 | 8.2 | 346 | 350 | 144 | 72 | 5 | 1233 | 8.3 | 257 | 152 | 75 | 4 | 1230 | 8.1 | 174 | 300 | 86 | 98 | 0.6 | 1266 | 7.8 | 20 | 61 | 10 | 70 | 0 | 1297 | 32 | 43 | 85 | 89 | 80 | 94 | 83 | 93 | 100 |
| D-24/1/91 | 35500 | 2.7 | 8.2 | 148 | 545 | 188 | 77 | 5.3 | 1329 | 8.2 | 221 | 230 | 74 | 5 | 1354 | 8.2 | 152 | 337 | 94 | 81 | 0.4 | 1407 | 7.7 | 14 | 71 | 11 | 86 | 0 | 1387 | 31 | 59 | 92 | 91 | 79 | 91 | 87 | 94 | 100 |
| D-25/1/91 | 37730 | 0.7 | 8 | 427 | 815 | 204 | 71 | 4.5 | 1455 | 7.9 | 517 | 186 | 73 | 4.5 | 1558 | 7.9 | 119 | 423 | 118 | 78 | 0.3 | 1602 | 7.4 | 23 | 101 | 20 | 88 | 0 | 1590 | 40 | 37 | 93 | 85 | 76 | 95 | 88 | 90 | 100 |
| D-27/1/91 | 28209 | 1.2 | 7.9 | 319 | 532 | 250 | 71 | 6.5 | 1379 | 8 | 331 | 254 | 69 | 6 | 1453 | 8 | 150 | 299 | 92 | 74 | 0.4 | 1431 | 7.7 | 15 | 60 | 14 | 86 | 0 | 1508 | 55 | 64 | 93 | 90 | 80 | 95 | 89 | 94 | 100 |
| D-28/1/91 | 32680 | 2.6 | 8.2 | 334 | 594 | 256 | 67 | 7 | 1380 | 8.3 | 268 | 238 | 67 | 5.5 | 1400 | 8.2 | 180 | 372 | 118 | 66 | 0.9 | 1352 | 7.9 | 20 | 62 | 11 | 80 | 0 | 1167 | 33 | 50 | 84 | 89 | 83 | 94 | 90 | 96 | 100 |
| D-29/1/91 | 32974 | 2.6 | 8.1 | 311 | 420 | 208 | 69 | 6.5 | 1474 | 8.2 | 258 | 198 | 73 | 3.3 | 1559 | 8.1 | 189 | 309 | 152 | 71 | 3 | 1402 | 7.8 | 18 | 41 | 12 | 83 | 0 | 1466 | 27 | 23 | 7.7 | 91 | 87 | 94 | 90 | 94 | 100 |
| D-30/1/91 | 33189 | 1 | 8.1 | 238 | 591 | 202 | 84 | 6 | 1380 | 8 | 282 | 210 | 80 | 5.5 | 1372 | 8.2 | 153 | 290 | 106 | 91 | 0.5 | 1413 | 7.9 | 14 | 63 | 17 | 94 | 0 | 1417 | 46 | 50 | 91 | 90 | 78 | 93 | 89 | 92 | 100 |
| D-31/1/91 | 34579 | 2.5 | 8.3 | 261 | 592 | 216 | 71 | 7.5 | 1663 | 8.4 | 244 | 256 | 73 | 8 | 1685 | 8.2 | 163 | 400 | 170 | 78 | 3.5 | 1680 | 7.8 | 24 | 47 | 13 | 93 | 0 | 1848 | 33 | 34 | 56 | 85 | 88 | 91 | 92 | 94 | 100 |
| D-1/5/91 | 46126 | 1.3 | 7.5 | 122 | 289 | 114 | 74 | 2.5 | 1103 | 7.6 | 122 | 146 | 64 | 2 | 1060 | 7.6 | 97 | 242 | 70 | 77 | 0.3 | 1094 | 7.6 | 12 | 99 | 10 | 96 | 0 | 1131 | 21 | 52 | 88 | 88 | 59 | 90 | 66 | 91 | 100 |
| D-2/5/91 | 43445 | 4.1 | 7.8 | 133 | 295 | 158 | 56 | 3 | 1436 | 7.9 | 140 | 158 | 54 | 2.5 | 1448 | 7.9 | 92 | 206 | 74 | 73 | 0.3 | 1370 | 7.7 | 14 | 66 | 12 | 75 | 0 | 1203 | 34 | 53 | 88 | 85 | 68 | 90 | 78 | 92 | 100 |
| D-3/5/91 | 35990 | 1.7 | 7.9 | 142 | 272 | 160 | 56 | 2 | 1543 | 7.9 | 154 | 174 | 58 | 3 | 1485 | 7.8 | 114 | 295 | 86 | 74 | 0.3 | 1560 | 7.6 | 18 | 74 | 14 | 86 | 0 | 1565 | 26 | 51 | 90 | 84 | 75 | 87 | 73 | 91 | 100 |
| D-5/5/91 | 36976 | 0.9 | 7.9 | 152 | 510 | 136 | 68 | 3 | 1235 | 7.7 | 145 | 148 | 66 | 2.5 | 1266 | 7.7 | 75 | 204 | 72 | 69 | 0.3 | 1247 | 7.6 | 16 | 65 | 11 | 91 | 0 | 1279 | 48 | 51 | 88 | 79 | 68 | 90 | 87 | 92 | 99 |
| D-6/5/91 | 33085 | 0.8 | 7.7 | 185 | 518 | 202 | 66 | 3 | 2350 | 7.8 | 345 | 210 | 68 | 3.5 | 2120 | 7.7 | 156 | 355 | 98 | 84 | 0.5 | 2080 | 7.7 | 14 | 65 | 13 | 92 | 0 | 1568 | 55 | 53 | 87 | 91 | 82 | 92 | 88 | 94 | 99 |
| D-7/5/91 | 34150 | 3.1 | 7.9 | 228 | 486 | 184 | 70 | 4.5 | 1912 | 8 | 209 | 214 | 69 | 4 | 1867 | 8 | 145 | 297 | 100 | 78 | 0.5 | 1889 | 7.8 | 24 | 86 | 16 | 81 | 0 | 1888 | 31 | 53 | 89 | 83 | 71 | 90 | 82 | 91 | 100 |
| D-8/5/91 | 60081 | 3.5 | 7.6 | 100 | 212 | 280 | 34 | 3.5 | 651 | 7.6 | 121 | 288 | 36 | 3 | 646 | 7.5 | 73 | 149 | 98 | 45 | 0.5 | 697 | 7.4 | 15 | 59 | 14 | 86 | 0 | 937 | 40 | 66 | 83 | 80 | 60 | 85 | 72 | 95 | 100 |
| D-9/5/91 | 57629 | 3.7 | 7.3 | 80 | 204 | 180 | 32 | 1.8 | 940 | 7.4 | 84 | 198 | 34 | 1.5 | 912 | 7.5 | 45 | 125 | 74 | 46 | 0.3 | 863 | 7.4 | 10 | 51 | 10 | 80 | 0 | 683 | 46 | 63 | 80 | 78 | 59 | 88 | 75 | 94 | 100 |
| D-10/5/91 | 48110 | 1 | 7.7 | 179 | 340 | 150 | 53 | 4.5 | 1509 | 7.7 | 168 | 198 | 49 | 4.5 | 1517 | 7.6 | 108 | 163 | 60 | 67 | 0.3 | 1568 | 7.2 | 14 | 96 | 10 | 72 | 0 | 1548 | 36 | 70 | 93 | 87 | 41 | 92 | 72 | 93 | 100 |
| D-12/5/91 | 59184 | 0.6 | 7.5 | 94 | 189 | 120 | 57 | 1.5 | 1200 | 7.5 | 122 | 144 | 51 | 1.3 | 1206 | 7.6 | 71 | 203 | 54 | 67 | 0.2 | 1260 | 7.5 | 13 | 33 | 10 | 68 | 0 | 1353 | 42 | 63 | 84 | 82 | 84 | 86 | 83 | 92 | 100 |
| D-13/5/91 | 47489 | 0.2 | 7.6 | 135 | 297 | 164 | 57 | 2.5 | 1000 | 7.7 | 175 | 196 | 53 | 2.5 | 1040 | 7.8 | 99 | 250 | 98 | 59 | 0.5 | 968 | 7.7 | 20 | 43 | 18 | 83 | 0 | 906 | 43 | 50 | 80 | 80 | 83 | 85 | 86 | 89 | 99 |
| D-14/5/91 | 35374 | 4.4 | 7.9 | 175 | 566 | 292 | 58 | 7.5 | 1268 | 8 | 196 | 406 | 50 | 9.2 | 1358 | 7.9 | 139 | 327 | 110 | 73 | 0.5 | 1385 | 7.8 | 20 | 80 | 16 | 88 | 0.1 | 1234 | 29 | 73 | 95 | 86 | 76 | 89 | 86 | 95 | 99 |
| D-15/5/91 | 33434 | 3.2 | 7.9 | 223 | 538 | 284 | 63 | 7.5 | 1425 | 8 | 218 | 256 | 63 | 6.5 | 1441 | 7.9 | 153 | 355 | 112 | 79 | 0.5 | 1403 | 7.9 | 18 | 104 | 19 | 84 | 0 | 1391 | 30 | 56 | 92 | 88 | 71 | 92 | 81 | 93 | 100 |
| D-16/5/91 | 31967 | 3.3 | 7.9 | 222 | 516 | 456 | 47 | 8.5 | 1335 | 7.9 | 302 | 374 | 53 | 9 | 1393 | 7.8 | 161 | 352 | 110 | 75 | 0.4 | 1450 | 7.7 | 21 | 96 | 19 | 82 | 0 | 1493 | 47 | 71 | 96 | 87 | 73 | 91 | 81 | 96 | 100 |
| D-17/5/91 | 32835 | 1.7 | 7.8 | 159 | 516 | 248 | 61 | 5 | 1405 | 7.7 | 201 | 314 | 55 | 5 | 1340 | 7.8 | 127 | 314 | 104 | 71 | 0.4 | 1381 | 7.7 | 17 | 72 | 15 | 87 | 0 | 1390 | 37 | 67 | 92 | 87 | 77 | 89 | 86 | 94 | 100 |
| D-19/5/91 | 33000 | 1.5 | 7.7 | 153 | 404 | 238 | 56 | 4.5 | 1049 | 7.7 | 138 | 184 | 63 | 2.7 | 1073 | 7.7 | 114 | 265 | 116 | 69 | 0.5 | 1061 | 7.7 | 33 | 111 | 41 | 81 | 0 | 1126 | 17 | 37 | 82 | 71 | 58 | 78 | 73 | 83 | 99 |
| D-20/5/91 | 47243 | 0.8 | 7.7 | 168 | 376 | 272 | 46 | 5.3 | 1052 | 7.8 | 136 | 296 | 45 | 5.6 | 1063 | 7.7 | 97 | 242 | 80 | 73 | 0.2 | 1083 | 7.7 | 14 | 87 | 15 | 100 | 0 | 1070 | 29 | 73 | 96 | 86 | 64 | 92 | 77 | 95 | 100 |
| D-21/5/91 | 40295 | 0.9 | 7.7 | 238 | 327 | 194 | 61 | 5 | 1725 | 7.8 | 210 | 252 | 55 | 6.3 | 1724 | 7.8 | 101 | 233 | 82 | 68 | 0.3 | 1736 | 7.7 | 12 | 79 | 14 | 91 | 0 | 1673 | 52 | 68 | 95 | 88 | 66 | 95 | 76 | 93 | 100 |
| D-22/5/91 | 38792 | 1.9 | 7.7 | 250 | 431 | 196 | 64 | 5.5 | 1219 | 7.9 | 399 | 238 | 62 | 6 | 1232 | 7.9 | 133 | 310 | 86 | 74 | 0.2 | 1212 | 7.7 | 15 | 90 | 17 | 77 | 0.1 | 1197 | 67 | 64 | 97 | 89 | 71 | 94 | 79 | 91 | 99 |
| D-23/5/91 | 36162 | 2.5 | 7.4 | 224 | 421 | 204 | 69 | 5 | 1328 | 7.7 | 293 | 248 | 65 | 6.5 | 1341 | 7.7 | 119 | 282 | 94 | 72 | 0.2 | 1331 | 7.7 | 18 | 98 | 20 | 80 | 0.1 | 1265 | 59 | 62 | 97 | 85 | 65 | 92 | 77 | 90 | 98 |
| D-24/5/91 | 36495 | 0.1 | 7.7 | 213 | 627 | 2008 | 18 | 4.5 | 1257 | 7.6 | 308 | 1692 | 18 | 4.5 | 1335 | 7.5 | 97 | 226 | 66 | 70 | 0.3 | 1255 | 7.6 | 16 | 119 | 13 | 77 | 0 | 1289 | 69 | 96 | 93 | 84 | 47 | 93 | 81 | 99 | 100 |
| D-26/5/91 | 36922 | 0.5 | 7.6 | 122 | 338 | 174 | 62 | 5 | 1035 | 7.7 | 135 | 216 | 57 | 5 | 1030 | 7.7 | 108 | 244 | 68 | 79 | 0.2 | 1099 | 7.6 | 14 | 44 | 12 | 82 | 0 | 1140 | 20 | 69 | 97 | 87 | 82 | 89 | 73 | 91 | 100 |
| D-27/5/91 | 43497 | 2.1 | 7.8 | 134 | 323 | 190 | 61 | 4.5 | 2070 | 7.8 | 126 | 206 | 54 | 3 | 2050 | 7.9 | 85 | 229 | 118 | 61 | 0.3 | 1784 | 7.7 | 15 | 56 | 13 | 85 | 0 | 1680 | 33 | 43 | 90 | 82 | 76 | 89 | 83 | 93 | 100 |
| D-28/5/91 | 38809 | 0.8 | 7.6 | 179 | 432 | 1228 | 23 | 36 | 1889 | 7.6 | 174 | 1692 | 21 | 46 | 1906 | 7.7 | 99 | 227 | 90 | 67 | 0.5 | 1962 | 7.6 | 18 | 71 | 17 | 82 | 0 | 1932 | 43 | 95 | 99 | 82 | 69 | 90 | 84 | 99 | 100 |
| D-29/5/91 | 34301 | 3.9 | 7.9 | 243 | 459 | 286 | 55 | 6.5 | 1174 | 7.9 | 216 | 342 | 46 | 7.5 | 1202 | 7.9 | 123 | 253 | 90 | 58 | 0.3 | 1190 | 7.8 | 15 | 103 | 14 | 86 | 0 | 1167 | 43 | 74 | 96 | 88 | 59 | 94 | 78 | 95 | 100 |
| D-30/5/91 | 33968 | 1.5 | 7.7 | 198 | 546 | 308 | 62 | 13 | 1869 | 7.8 | 239 | 420 | 53 | 19 | 1893 | 7.9 | 81 | 222 | 76 | 84 | 0.4 | 1804 | 7.7 | 17 | 79 | 27 | 78 | 0.1 | 1792 | 66 | 82 | 98 | 79 | 64 | 91 | 86 | 91 | 100 |
| D-31/5/91 | 34094 | 1 | 7.8 | 156 | 483 | 964 | 24 | 18 | 2120 | 7.8 | 196 | 764 | 30 | 17 | 2110 | 7.8 | 97 | 170 | 92 | 61 | 0.4 | 1930 | 7.6 | 15 | 84 | 20 | 75 | 0 | 1966 | 51 | 88 | 98 | 85 | 70 | 90 | 79 | 98 | 100 |
| D-1/4/91 | 34573 | 0.7 | 7.7 | 156 | 276 | 146 | 71 | 3.3 | 1265 | 7.7 | 166 | 206 | 66 | 4.5 | 1270 | 7.7 | 114 | 176 | 124 | 73 | 0.9 | 1260 | 7.7 | 30 | 43 | 44 | 80 | 0.5 | 1270 | 31 | 40 | 80 | 74 | 76 | 81 | 84 | 70 | 85 |
| D-2/4/91 | 35395 | 0.6 | 7.8 | 273 | 473 | 210 | 73 | 4.5 | 1232 | 7.9 | 213 | 224 | 69 | 6 | 1257 | 7.9 | 170 | 310 | 116 | 79 | 0.1 | 1214 | 7.7 | 22 | 85 | 22 | 79 | 0 | 1116 | 20 | 48 | 98 | 87 | 73 | 92 | 82 | 90 | 100 |
| D-3/4/91 | 34525 | 0.8 | 7.8 | 312 | 576 | 224 | 68 | 5.5 | 1300 | 7.9 | 324 | 268 | 72 | 5 | 1280 | 7.9 | 157 | 306 | 92 | 74 | 0.4 | 1248 | 7.8 | 23 | 74 | 23 | 81 | 0 | 1251 | 52 | 66 | 93 | 85 | 76 | 93 | 87 | 90 | 100 |
| D-4/4/91 | 35861 | 0.8 | 8.1 | 242 | 492 | 176 | 75 | 5.5 | 1530 | 8.2 | 397 | 238 | 68 | 5 | 1612 | 8.1 | 148 | 312 | 88 | 75 | 0.1 | 1589 | 7.9 | 20 | 84 | 19 | 78 | 0 | 1566 | 63 | 63 | 98 | 87 | 73 | 92 | 83 | 89 | 100 |
| D-5/4/91 | 43082 | 0.7 | 7.8 | 173 | 496 | 178 | 66 | 4.5 | 1329 | 7.9 | 365 | 212 | 63 | 5 | 1303 | 7.9 | 124 | 304 | 100 | 78 | 0.4 | 1338 | 7.7 | 23 | 88 | 23 | 87 | 0.1 | 1408 | 66 | 53 | 92 | 82 | 71 | 87 | 82 | 87 | 99 |

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-7/4/91 | 27931 | 2.2 | 7.5 | 296 | 455 | 278 | 78 | 4.3 | 1439 | 7.6 | 255 | 182 | 69 | 3.9 | 1436 | 7.6 | 131 | 318 | 84 | 83 | 0.3 | 1421 | 7.6 | 15 | 106 | 13 | 96 | 0 | 1423 | 49 | 54 | 92 | 89 | 67 | 95 | 77 | 95 | 100 |
| D-8/4/91 | 32954 | 0.7 | 7.7 | 269 | 423 | 192 | 68 | 6 | 1164 | 7.9 | 360 | 266 | 59 | 7 | 1170 | 7.9 | 145 | 314 | 94 | 70 | 0.6 | 1222 | 7.7 | 16 | 110 | 15 | 98 | 0 | 1149 | 60 | 65 | 91 | 89 | 65 | 94 | 74 | 92 | 100 |
| D-9/4/91 | 33773 | 4.1 | 7.9 | 233 | 506 | 222 | 72 | 4.5 | 1410 | 7.9 | 328 | 242 | 69 | 5.8 | 1366 | 7.9 | 201 | 388 | 102 | 84 | 0.4 | 1589 | 7.7 | 18 | 90 | 16 | 94 | 0 | 1461 | 39 | 58 | 94 | 91 | 77 | 92 | 82 | 93 | 100 |
| D-10/4/91 | 33666 | 2.7 | 7.6 | 237 | 494 | 204 | 72 | 5.5 | 1187 | 7.7 | 241 | 224 | 65 | 5 | 1248 | 7.6 | 164 | 341 | 102 | 84 | 0.7 | 1350 | 7.6 | 24 | 71 | 57 | 88 | 0.1 | 1312 | 32 | 55 | 86 | 85 | 79 | 90 | 86 | 72 | 99 |
| D-11/4/91 | 39715 | 0.6 | 7.8 | 241 | 539 | 212 | 71 | 5.5 | 1338 | 7.9 | 232 | 218 | 67 | 5.5 | 1403 | 7.9 | 160 | 353 | 92 | 80 | 0.2 | 1427 | 7.8 | 20 | 81 | 13 | 92 | 0 | 1364 | 31 | 58 | 96 | 88 | 77 | 92 | 85 | 94 | 100 |
| D-12/4/91 | 28923 | 0.7 | 7.2 | 133 | 467 | 186 | 66 | 4.5 | 1398 | 7.4 | 141 | 216 | 65 | 4.5 | 1425 | 7.5 | 130 | 393 | 76 | 84 | 0.3 | 1586 | 7.5 | 20 | 147 | 31 | 84 | 0 | 1543 | 7.8 | 65 | 93 | 85 | 63 | 85 | 69 | 83 | 100 |
| D-14/4/91 | 32317 | 0.7 | 7.5 | 166 | 393 | 242 | 55 | 5.3 | 958 | 7.7 | 206 | 256 | 55 | 5.5 | 982 | 7.7 | 135 | 260 | 96 | 79 | 0.4 | 1038 | 7.6 | 13 | 27 | 10 | 90 | 0 | 1082 | 35 | 63 | 94 | 90 | 90 | 92 | 93 | 96 | 100 |
| D-15/4/91 | 33090 | 0.4 | 7.9 | 205 | 453 | 198 | 71 | 5 | 1479 | 7.9 | 199 | 214 | 70 | 5.5 | 1530 | 7.9 | 130 | 301 | 92 | 76 | 0.3 | 1401 | 7.8 | 15 | 70 | 22 | 64 | 0 | 1233 | 35 | 57 | 95 | 89 | 77 | 93 | 85 | 89 | 100 |
| D-16/4/91 | 33371 | 0.8 | 7.8 | 290 | 448 | 218 | 67 | 5.5 | 1842 | 8 | 254 | 200 | 62 | 5 | 1829 | 7.9 | 137 | 312 | 90 | 76 | 0.3 | 1856 | 7.8 | 19 | 92 | 22 | 80 | 0.1 | 1843 | 46 | 55 | 94 | 86 | 71 | 93 | 80 | 90 | 99 |
| D-17/4/91 | 33813 | 2.9 | 7.7 | 212 | 468 | 224 | 66 | 4.5 | 1475 | 7.8 | 337 | 260 | 62 | 6 | 1506 | 7.9 | 149 | 320 | 114 | 75 | 0.3 | 1612 | 7.8 | 18 | 92 | 14 | 86 | 0 | 1654 | 56 | 56 | 96 | 88 | 71 | 92 | 80 | 94 | 100 |
| D-18/4/91 | 35456 | 6.4 | 8.1 | 184 | 412 | 216 | 62 | 5 | 1800 | 8.1 | 194 | 234 | 62 | 5 | 1770 | 8 | 116 | 276 | 106 | 66 | 0.5 | 1817 | 7.7 | 14 | 84 | 12 | 92 | 0 | 1760 | 40 | 55 | 90 | 88 | 70 | 92 | 80 | 94 | 100 |
| D-19/4/91 | 38045 | 9.6 | 8.1 | 177 | 428 | 220 | 62 | 4.5 | 1372 | 8 | 160 | 244 | 56 | 5.3 | 1335 | 7.9 | 138 | 372 | 90 | 73 | 0.4 | 1340 | 7.7 | 20 | 104 | 18 | 89 | 0 | 1323 | 14 | 63 | 93 | 86 | 72 | 89 | 76 | 92 | 100 |
| D-21/4/91 | 31191 | 2 | 7.9 | 270 | 321 | 168 | 67 | 4 | 1026 | 7.9 | 226 | 200 | 63 | 3.5 | 971 | 7.8 | 100 | 226 | 82 | 81 | 0.3 | 1009 | 7.7 | 14 | 51 | 16 | 83 | 0 | 1086 | 56 | 59 | 91 | 86 | 77 | 95 | 84 | 91 | 100 |
| D-22/4/91 | 36215 | 6.5 | 7.5 | 161 | 392 | 266 | 48 | 6 | 1156 | 7.8 | 252 | 208 | 47 | 3.5 | 1105 | 7.9 | 110 | 238 | 82 | 78 | 0.3 | 1108 | 7.7 | 15 | 59 | 13 | 83 | 0 | 1035 | 56 | 61 | 89 | 86 | 75 | 91 | 85 | 95 | 100 |
| D-23/4/91 | 34719 | 9 | 7.7 | 173 | 388 | 322 | 50 | 5.5 | 1196 | 7.7 | 167 | 210 | 62 | 4.5 | 1171 | 7.8 | 149 | 319 | 110 | 76 | 0.7 | 1234 | 7.6 | 20 | 108 | 24 | 79 | 0 | 1250 | 11 | 48 | 84 | 87 | 66 | 88 | 72 | 93 | 100 |
| D-24/4/91 | 35729 | 2.9 | 7.7 | 334 | 841 | 616 | 60 | 14 | 1285 | 7.8 | 357 | 572 | 56 | 16 | 1378 | 7.9 | 170 | 274 | 102 | 80 | 0.4 | 1420 | 7.8 | 23 | 100 | 20 | 97 | 0 | 1433 | 52 | 82 | 97 | 87 | 70 | 93 | 88 | 97 | 100 |
| D-25/4/91 | 36395 | 6.5 | 7.7 | 183 | 449 | 380 | 53 | 6.5 | 1306 | 7.8 | 258 | 464 | 51 | 7.5 | 1410 | 7.9 | 125 | 292 | 88 | 77 | 0.3 | 1498 | 7.7 | 23 | 54 | 19 | 90 | 0 | 1512 | 52 | 81 | 96 | 82 | 82 | 87 | 88 | 95 | 100 |
| D-26/4/91 | 41503 | 8.7 | 7.5 | 133 | 346 | 274 | 46 | 4.5 | 1186 | 7.3 | 125 | 146 | 59 | 4.5 | 1203 | 7.1 | 113 | 196 | 74 | 73 | 0.3 | 1229 | 7 | 17 | 65 | 18 | 83 | 0 | 1272 | 9.6 | 49 | 93 | 85 | 67 | 87 | 81 | 93 | 100 |
| D-28/4/91 | 27642 | 1.8 | 7.5 | 69 | 170 | 180 | 40 | 1.4 | 810 | 7.5 | 130 | 310 | 40 | 3 | 827 | 7.5 | 83 | 124 | 80 | 65 | 0.3 | 866 | 7.5 | 13 | 39 | 14 | 93 | 0 | 949 | 36 | 74 | 90 | 84 | 69 | 81 | 77 | 92 | 99 |
| D-29/4/91 | 35760 | 1.8 | 7.6 | 115 | 295 | 182 | 52 | 25 | 1400 | 7.7 | 125 | 166 | 55 | 19 | 1418 | 7.8 | 98 | 225 | 84 | 81 | 1 | 1396 | 7.7 | 19 | 58 | 19 | 84 | 0.1 | 1316 | 22 | 49 | 95 | 81 | 74 | 84 | 80 | 90 | 100 |
| D-1/7/91 | 33416 | 1.7 | 7.4 | 167 | 333 | 242 | 66 | 4.5 | 1960 | 7.6 | 211 | 202 | 67 | 4.5 | 2090 | 7.7 | 106 | 274 | 80 | 88 | 0.1 | 1942 | 7.7 | 12 | 73 | 14 | 79 | 0 | 1788 | 50 | 60 | 98 | 89 | 73 | 93 | 78 | 94 | 100 |
| D-2/7/91 | 35518 | 4.2 | 7.8 | 133 | 105 | 208 | 55 | 3.5 | 1293 | 7.8 | 138 | 236 | 53 | 4.5 | 1347 | 7.8 | 91 | 125 | 58 | 83 | 0.1 | 1323 | 7.7 | 16 | 20 | 16 | 94 | 0 | 1318 | 34 | 75 | 98 | 82 | 84 | 88 | 81 | 92 | 100 |
| D-3/7/91 | 35623 | 4.4 | 7.6 | 151 | 404 | 204 | 68 | 3.5 | 1565 | 7.6 | 137 | 232 | 63 | 4 | 1629 | 7.6 | 88 | 277 | 80 | 70 | 0.1 | 1575 | 7.4 | 12 | 84 | 13 | 92 | 0 | 1467 | 36 | 66 | 98 | 86 | 70 | 92 | 79 | 94 | 100 |
| D-4/7/91 | 32815 | 6.6 | 7.8 | 151 | 485 | 198 | 72 | 3.5 | 1535 | 7.7 | 140 | 156 | 62 | 4 | 1528 | 7.7 | 102 | 283 | 78 | 77 | 0.1 | 1571 | 7.6 | 13 | 101 | 14 | 77 | 0 | 1605 | 27 | 50 | 98 | 87 | 64 | 91 | 79 | 93 | 100 |
| D-5/7/91 | 32454 | 3.4 | 7.4 | 148 | 545 | 202 | 72 | 4.5 | 1337 | 7.4 | 138 | 272 | 63 | 4.5 | 1334 | 7.4 | 67 | 200 | 76 | 89 | 0.3 | 1283 | 7.4 | 11 | 163 | 11 | 91 | 0 | 1365 | 51 | 72 | 93 | 84 | 19 | 93 | 70 | 95 | 100 |
| D-7/7/91 | 26590 | 2.9 | 7.5 | 134 | 351 | 108 | 82 | 3 | 1135 | 7.5 | 154 | 182 | 65 | 3 | 1115 | 7.5 | 123 | 351 | 124 | 76 | 0.9 | 1117 | 7.4 | 12 | 107 | 15 | 79 | 0 | 1220 | 20 | 32 | 70 | 90 | 70 | 91 | 70 | 86 | 100 |
| D-8/7/91 | 33636 | 34 | 7.5 | 166 | 481 | 368 | 51 | 6.3 | 1355 | 7.6 | 155 | 302 | 54 | 5.5 | 1359 | 7.9 | 110 | 283 | 96 | 75 | 0.3 | 1283 | 8 | 29 | 113 | 19 | 81 | 0 | 1240 | 29 | 68 | 95 | 74 | 60 | 83 | 77 | 91 | 100 |
| D-9/7/91 | 32334 | 19 | 7.6 | 179 | 461 | 298 | 56 | 5 | 1340 | 7.6 | 225 | 212 | 63 | 5.5 | 1395 | 7.5 | 111 | 307 | 84 | 76 | 0.2 | 1292 | 7.4 | 20 | 125 | 24 | 78 | 0 | 1300 | 51 | 60 | 96 | 82 | 59 | 89 | 73 | 92 | 100 |
| D-10/7/91 | 35178 | 3.7 | 7.7 | 159 | 396 | 154 | 69 | 3 | 1440 | 7.8 | 159 | 184 | 67 | 3.1 | 1520 | 7.8 | 105 | 295 | 82 | 83 | 0.2 | 1415 | 7.7 | 23 | 109 | 30 | 87 | 0 | 1479 | 34 | 55 | 94 | 78 | 63 | 86 | 73 | 81 | 100 |
| D-11/7/91 | 35990 | 6.3 | 7.7 | 146 | 375 | 244 | 57 | 4 | 1413 | 7.8 | 158 | 280 | 56 | 6.5 | 1486 | 7.7 | 116 | 282 | 96 | 81 | 0.7 | 1449 | 7.6 | 23 | 112 | 32 | 83 | 0 | 1410 | 27 | 66 | 89 | 80 | 60 | 84 | 70 | 87 | 100 |
| D-12/7/91 | 35990 | 2.6 | 7.8 | 198 | 604 | 288 | 58 | 4.7 | 1235 | 7.8 | 229 | 222 | 60 | 4.5 | 1330 | 7.7 | 144 | 391 | 104 | 73 | 0.9 | 1320 | 7.7 | 31 | 159 | 41 | 83 | 0 | 1399 | 37 | 53 | 80 | 79 | 59 | 84 | 74 | 86 | 100 |
| D-14/7/91 | 35990 | 0.8 | 7.6 | 177 | 388 | 170 | 69 | 4 | 1147 | 7.7 | 166 | 198 | 67 | 3 | 1126 | 7.6 | 89 | 208 | 68 | 82 | 0.2 | 1105 | 7.7 | 14 | 74 | 16 | 88 | 0 | 1170 | 46 | 66 | 93 | 84 | 64 | 92 | 81 | 91 | 100 |
| D-15/7/91 | 35990 | 5.7 | 7.7 | 197 | 545 | 254 | 60 | 6 | 1202 | 7.7 | 182 | 248 | 61 | 4 | 1218 | 7.7 | 120 | 325 | 124 | 77 | 0.6 | 1168 | 7.6 | 20 | 216 | 18 | 100 | 0 | 1205 | 34 | 50 | 85 | 83 | 34 | 90 | 60 | 93 | 100 |
| D-16/7/91 | 35990 | 1.1 | 7.6 | 149 | 412 | 208 | 58 | 4.3 | 1593 | 7.7 | 127 | 194 | 61 | 4 | 1710 | 7.7 | 96 | 298 | 128 | 66 | 0.7 | 1606 | 7.8 | 17 | 110 | 22 | 84 | 0 | 1602 | 24 | 34 | 83 | 82 | 63 | 89 | 73 | 89 | 100 |
| D-17/7/91 | 35990 | 0.7 | 7.6 | 149 | 359 | 242 | 61 | 4 | 1315 | 7.7 | 221 | 280 | 63 | 7 | 1240 | 7.7 | 121 | 294 | 138 | 62 | 0.7 | 1240 | 7.7 | 105 | 290 | 104 | 87 | 0 | 1434 | 45 | 51 | 90 | 13 | 1.4 | 30 | 19 | 57 | 100 |
| D-18/7/91 | 35990 | 1 | 7.6 | 186 | 495 | 222 | 66 | 5.5 | 1518 | 7.6 | 168 | 222 | 64 | 4.5 | 1496 | 7.6 | 110 | 274 | 112 | 70 | 1 | 1505 | 7.6 | 101 | 292 | 74 | 84 | 0.3 | 1642 | 35 | 50 | 78 | 8.2 | 70 | 46 | 41 | 67 | 96 |
| D-19/7/91 | 35990 | 1.6 | 6.9 | 233 | 472 | 242 | 65 | 5 | 1183 | 7.3 | 192 | 236 | 59 | 4.5 | 1165 | 7.3 | 124 | 357 | 116 | 66 | 1.4 | 1253 | 7.5 | 101 | 236 | 78 | 74 | 0 | 1374 | 35 | 51 | 69 | 19 | 34 | 57 | 50 | 68 | 100 |
| D-21/7/91 | 35990 | 0.9 | 7.3 | 185 | 395 | 216 | 62 | 3 | 1367 | 7.4 | 238 | 232 | 60 | 3 | 1383 | 7.6 | 93 | 263 | 102 | 59 | 0.3 | 1356 | 7.6 | 31 | 117 | 39 | 72 | 0 | 1444 | 61 | 56 | 90 | 67 | 56 | 83 | 70 | 82 | 100 |
| D-22/7/91 | 35990 | 4 | 7.5 | 182 | 605 | 208 | 75 | 5 | 1672 | 7.5 | 155 | 230 | 62 | 5 | 1734 | 7.5 | 129 | 308 | 82 | 76 | 0.1 | 1693 | 7.6 | 26 | 86 | 21 | 71 | 0 | 1603 | 17 | 64 | 98 | 80 | 72 | 86 | 86 | 90 | 100 |
| D-23/7/91 | 35990 | 1.8 | 7.6 | 198 | 432 | 224 | 64 | 4 | 2150 | 7.7 | 195 | 208 | 67 | 3.5 | 2100 | 7.6 | 106 | 290 | 92 | 74 | 0.2 | 2140 | 7.6 | 19 | 90 | 33 | 82 | 0 | 2240 | 46 | 56 | 96 | 82 | 69 | 90 | 79 | 85 | 100 |
| D-24/7/91 | 35990 | 1.5 | 8.1 | 129 | 318 | 140 | 67 | 3 | 1319 | 8 | 124 | 136 | 69 | 2 | 1320 | 7.6 | 101 | 277 | 70 | 86 | 0.1 | 1270 | 7.5 | 15 | 86 | 22 | 86 | 0 | 1289 | 19 | 49 | 95 | 85 | 69 | 88 | 73 | 84 | 100 |
| D-25/7/91 | 35990 | 5.7 | 7.7 | 115 | 408 | 290 | 51 | 4 | 2150 | 7.7 | 139 | 198 | 58 | 2.5 | 2050 | 7.6 | 96 | 290 | 98 | 76 | 0.1 | 1924 | 7.5 | 20 | 118 | 46 | 76 | 0 | 2060 | 31 | 51 | 96 | 79 | 59 | 83 | 71 | 84 | 100 |
| D-26/7/91 | 35990 | 2.4 | 7.6 | 132 | 253 | 200 | 49 | 2.8 | 2110 | 7.6 | 129 | 256 | 50 | 2.5 | 2060 | 7.6 | 106 | 245 | 94 | 83 | 0.1 | 1974 | 7.6 | 15 | 90 | 28 | 86 | 0 | 1959 | 18 | 63 | 96 | 86 | 63 | 89 | 64 | 86 | 99 |
| D-28/7/91 | 35990 | 1.4 | 7.5 | 188 | 263 | 148 | 68 | 2 | 1645 | 7.5 | 148 | 134 | 78 | 1.7 | 1668 | 7.6 | 98 | 198 | 56 | 89 | 0.1 | 1671 | 7.7 | 14 | 85 | 21 | 86 | 0 | 1760 | 34 | 58 | 94 | 86 | 57 | 93 | 68 | 86 | 100 |
| D-29/7/91 | 35990 | 0.2 | 7.8 | 114 | 273 | 138 | 74 | 3.8 | 1870 | 7.8 | 128 | 172 | 76 | 3.5 | 1929 | 7.8 | 80 | 192 | 68 | 82 | 0.2 | 1835 | 7.5 | 14 | 73 | 14 | 89 | 0 | 1809 | 38 | 61 | 96 | 83 | 62 | 88 | 73 | 90 | 100 |
| D-30/7/91 | 35990 | 0.6 | 7.4 | 159 | 394 | 170 | 69 | 3 | 1450 | 7.4 | 155 | 196 | 60 | 3 | 1556 | 7.4 | 111 | 274 | 78 | 72 | 0.1 | 1690 | 7.6 | 16 | 110 | 26 | 79 | 0 | 1693 | 28 | 60 | 98 | 86 | 70 | 90 | 72 | 85 | 100 |
| D-31/7/91 | 35990 | 1.6 | 7.5 | 170 | 336 | 168 | 69 | 2.6 | 1531 | 7.5 | 192 | 200 | 67 | 3 | 1485 | 7.6 | 101 | 265 | 80 | 80 | 0.1 | 1642 | 7.7 | 19 | 99 | 34 | 79 | 0 | 1648 | 47 | 60 | 98 | 81 | 63 | 89 | 71 | 80 | 99 |
| D-2/6/91 | 32308 | 1.8 | 7.7 | 118 | 295 | 178 | 63 | 3 | 1459 | 7.7 | 137 | 236 | 53 | 4 | 1442 | 7.7 | 103 | 231 | 104 | 67 | 0.3 | 1474 | 7.6 | 12 | 72 | 11 | 82 | 0 | 1607 | 25 | 56 | 93 | 88 | 69 | 90 | 76 | 94 | 99 |
| D-3/6/91 | 31114 | 1 | 7.8 | 181 | 462 | 216 | 62 | 5 | 1315 | 7.8 | 197 | 354 | 54 | 8 | 1377 | 7.8 | 154 | 307 | 118 | 64 | 0.3 | 1270 | 7.7 | 22 | 68 | 20 | 75 | 0.1 | 1142 | 37 | 67 | 96 | 82 | 78 | 88 | 85 | 91 | 99 |
| D-4/6/91 | 31205 | 0.6 | 7.7 | 214 | 467 | 242 | 61 | 4.5 | 1171 | 7.8 | 218 | 318 | 55 | 6 | 1205 | 7.7 | 132 | 300 | 116 | 67 | 0.2 | 1312 | 7.7 | 22 | 86 | 22 | 82 | 0.1 | 1333 | 39 | 64 | 97 | 83 | 71 | 90 | 82 | 91 | 99 |
| D-5/6/91 | 35509 | 1.1 | 7.8 | 181 | 358 | 228 | 64 | 4.3 | 1224 | 7.9 | 243 | 456 | 50 | 9.5 | 1263 | 7.9 | 107 | 272 | 112 | 67 | 0.3 | 1265 | 7.8 | 16 | 105 | 16 | 75 | 0 | 1268 | 56 | 75 | 97 | 85 | 61 | 91 | 71 | 93 | 100 |
| D-6/6/91 | 34903 | 2.2 | 7.7 | 227 | 416 | 218 | 62 | 5.4 | 1335 | 7.7 | 257 | 416 | 50 | 8 | 1350 | 7.7 | 156 | 311 | 134 | 60 | 0.5 | 1368 | 7.7 | 25 | 86 | 16 | 75 | 0 | 1379 | 39 | 68 | 94 | 84 | 72 | 93 | 89 | 95 | 100 |
| D-7/6/91 | 34294 | 5.9 | 7.6 | 146 | 438 | 174 | 68 | 4.5 | 1102 | 7.6 | 182 | 332 | 55 | 4.5 | 1093 | 7.5 | 123 | 340 | 112 | 64 | 0.3 | 1239 | 7.5 | 23 | 78 | 16 | 88 | 0 | 1213 | 32 | 66 | 93 | 81 | 77 | 84 | 82 | 91 | 100 |
| D-9/6/91 | 30614 | 0.9 | 7.7 | 146 | 313 | 168 | 64 | 4.3 | 1019 | 7.7 | 173 | 246 | 59 | 5.9 | 1020 | 7.7 | 116 | 183 | 124 | 65 | 0.2 | 1034 | 7.6 | 16 | 59 | 15 | 80 | 0 | 1056 | 33 | 50 | 97 | 86 | 68 | 89 | 81 | 91 | 100 |
| D-10/6/91 | 33239 | 2.5 | 7.8 | 217 | 591 | 264 | 67 | 7.5 | 1234 | 7.7 | 227 | 500 | 53 | 8.5 | 1219 | 7.8 | 135 | 338 | 132 | 67 | 0.3 | 1294 | 7.7 | 17 | 84 | 27 | 78 | 0 | 1247 | 41 | 74 | 97 | 87 | 75 | 92 | 86 | 90 | 100 |
| D-12/6/91 | 32100 | 1.5 | 7.5 | 277 | 523 | 324 | 68 | 7 | 1817 | 7.8 | 315 | 348 | 56 | 9 | 1816 | 7.8 | 154 | 297 | 84 | 67 | 0.2 | 1829 | 7.7 | 21 | 53 | 19 | 78 | 0 | 1819 | 51 | 76 | 98 | 86 | 82 | 92 | 90 | 94 | 100 |
| D-13/6/91 | 32538 | 1 | 7.6 | 219 | 511 | 286 | 58 | 6.5 | 1326 | 7.8 | 243 | 382 | 57 | 7 | 1420 | 7.8 | 142 | 323 | 106 | 72 | 0.2 | 1359 | 7.7 | 22 | 98 | 21 | 74 | 0 | 1361 | 42 | 72 | 97 | 85 | 70 | 90 | 81 | 93 | 100 |
| D-14/6/91 | 35571 | 1.8 | 7.2 | 166 | 549 | 138 | 64 | 4.5 | 1769 | 7.4 | 197 | 426 | 50 | 4.5 | 1882 | 7.5 | 111 | 287 | 74 | 70 | 0.3 | 2100 | 7.6 | 15 | 85 | 9 | 89 | 0 | 2030 | 44 | 83 | 93 | 87 | 70 | 91 | 85 | 94 | 100 |
| D-16/6/91 | 33210 | 1.3 | 7.5 | 164 | 353 | 218 | 62 | 4.5 | 1535 | 7.6 | 185 | 244 | 58 | 5 | 1489 | 7.7 | 107 | 229 | 80 | 75 | 0.2 | 1573 | 7.6 | 15 | 59 | 14 | 71 | 0 | 1679 | 42 | 67 | 97 | 86 | 74 | 91 | 83 | 94 | 100 |

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-17/6/91 | 34097 | 0.7 | 7.8 | 192 | 330 | 190 | 58 | 5.5 | 1775 | 7.8 | 206 | 218 | 55 | 5.5 | 1782 | 7.8 | 117 | 275 | 82 | 73 | 0.2 | 1716 | 7.7 | 10 | 70 | 19 | 67 | 0 | 1605 | 43 | 62 | 96 | 92 | 75 | 95 | 79 | 90 | 100 |
| D-18/6/91 | 31224 | 0.7 | 7.8 | 231 | 473 | 222 | 62 | 4.5 | 1309 | 7.8 | 279 | 268 | 58 | 6 | 1368 | 7.9 | 118 | 299 | 82 | 78 | 0.2 | 1529 | 7.7 | 27 | 136 | 29 | 76 | 0 | 1523 | 58 | 69 | 98 | 77 | 55 | 88 | 71 | 87 | 100 |
| D-19/6/91 | 35389 | 0.4 | 7.6 | 157 | 392 | 232 | 53 | 5 | 1565 | 7.7 | 163 | 274 | 54 | 6 | 1517 | 7.8 | 103 | 250 | 76 | 82 | 0.2 | 1630 | 7.6 | 18 | 65 | 22 | 82 | 0 | 1610 | 37 | 72 | 98 | 83 | 74 | 89 | 83 | 91 | 100 |
| D-20/6/91 | 36708 | 3.9 | 7.8 | 196 | 403 | 200 | 66 | 5 | 1592 | 7.8 | 301 | 236 | 70 | 4.5 | 1478 | 7.7 | 130 | 288 | 74 | 84 | 0.2 | 1469 | 7.6 | 12 | 54 | 14 | 86 | 0 | 1483 | 57 | 69 | 96 | 91 | 81 | 94 | 87 | 93 | 100 |
| D-21/6/91 | 33064 | 1.2 | 7.7 | 266 | 438 | 200 | 62 | 4.5 | 1396 | 7.9 | 217 | 214 | 63 | 4.5 | 1423 | 7.9 | 137 | 269 | 74 | 76 | 0.2 | 1471 | 7.8 | 13 | 104 | 16 | 78 | 0 | 1441 | 37 | 65 | 96 | 91 | 61 | 95 | 76 | 92 | 100 |
| D-24/6/91 | 31949 | 3.3 | 7.3 | 133 | 310 | 208 | 57 | 3.5 | 1134 | 7.5 | 137 | 168 | 64 | 2.5 | 1133 | 7.6 | 104 | 212 | 64 | 81 | 0.2 | 1143 | 7.6 | 6 | 90 | 12 | 83 | 0 | 1226 | 24 | 62 | 92 | 94 | 58 | 96 | 71 | 94 | 100 |
| D-25/6/91 | 35195 | 3.2 | 7.6 | 134 | 404 | 190 | 58 | 5.5 | 1305 | 7.7 | 207 | 212 | 59 | 5.5 | 1341 | 7.8 | 109 | 325 | 64 | 81 | 0.2 | 1286 | 7.6 | 17 | 95 | 26 | 73 | 0 | 1190 | 47 | 70 | 97 | 84 | 71 | 87 | 77 | 86 | 100 |
| D-26/6/91 | 34886 | 0.9 | 7.6 | 178 | 310 | 198 | 67 | 5.5 | 1560 | 7.7 | 156 | 178 | 66 | 3 | 1549 | 7.7 | 100 | 248 | 74 | 89 | 0.1 | 1590 | 7.7 | 13 | 62 | 22 | 82 | 0 | 1635 | 36 | 58 | 97 | 87 | 75 | 93 | 80 | 89 | 100 |
| D-27/6/91 | 33708 | 0.5 | 7.8 | 171 | 380 | 164 | 83 | 4.5 | 1565 | 7.8 | 216 | 190 | 80 | 4.5 | 1589 | 7.8 | 119 | 279 | 52 | 100 | 0.1 | 1571 | 7.7 | 12 | 74 | 13 | 83 | 0 | 1560 | 45 | 73 | 98 | 90 | 74 | 93 | 81 | 92 | 100 |
| D-28/6/91 | 32253 | 1.9 | 7.7 | 140 | 380 | 178 | 65 | 4.5 | 1668 | 7.7 | 147 | 182 | 68 | 3.5 | 1789 | 7.8 | 110 | 272 | 66 | 85 | 0.2 | 1898 | 7.8 | 15 | 92 | 22 | 77 | 0 | 1883 | 25 | 64 | 96 | 86 | 66 | 89 | 76 | 88 | 100 |
| D-30/6/91 | 29793 | 2.4 | 7.6 | 137 | 372 | 170 | 72 | 4 | 1704 | 7.7 | 143 | 172 | 63 | 3.3 | 1616 | 7.7 | 100 | 290 | 76 | 82 | 0.5 | 1705 | 7.8 | 10 | 78 | 10 | 100 | 0 | 1817 | 30 | 56 | 85 | 90 | 73 | 93 | 79 | 94 | 100 |
| D-1/10/91 | 32208 | 0.5 | 7.6 | 145 | 258 | 194 | 59 | 2.5 | 985 | 7.6 | 161 | 204 | 56 | 2.8 | 991 | 7.7 | 94 | 196 | 82 | 71 | 0.2 | 1050 | 7.8 | 19 | 55 | 20 | 88 | 0 | 1124 | 42 | 60 | 93 | 80 | 72 | 87 | 79 | 90 | 100 |
| D-2/10/91 | 33649 | 1.4 | 7.5 | 147 | 351 | 230 | 47 | 3.5 | 1395 | 7.7 | 169 | 264 | 48 | 4.2 | 1390 | 7.8 | 99 | 253 | 86 | 77 | 0.3 | 1160 | 7.8 | 24 | 89 | 17 | 85 | 0 | 1068 | 41 | 67 | 93 | 76 | 65 | 84 | 79 | 93 | 100 |
| D-3/10/91 | 34536 | 0.5 | 7.6 | 133 | 316 | 188 | 62 | 2.5 | 1273 | 7.5 | 219 | 218 | 60 | 3.5 | 1322 | 7.6 | 69 | 264 | 92 | 72 | 0.2 | 1373 | 7.7 | 20 | 104 | 22 | 76 | 0 | 1371 | 69 | 58 | 94 | 71 | 61 | 85 | 67 | 88 | 100 |
| D-4/10/91 | 33178 | 1.4 | 7.6 | 155 | 392 | 182 | 65 | 4.5 | 1198 | 7.7 | 197 | 180 | 67 | 4 | 1264 | 7.8 | 34 | 356 | 78 | 80 | 0.4 | 1245 | 7.8 | 12 | 68 | 15 | 72 | 0 | 1248 | 40 | 57 | 90 | 65 | 81 | 92 | 83 | 92 | 100 |
| D-5/10/91 | 33695 | 1.1 | 7.3 | 133 | 520 | 196 | 85 | 3.3 | 1258 | 7.6 | 155 | 128 | 80 | 3 | 1251 | 7.7 | 87 | 194 | 90 | 75 | 0.2 | 1130 | 7.8 | 16 | 44 | 19 | 81 | 0 | 1165 | 44 | 60 | 93 | 82 | 77 | 88 | 92 | 91 | 100 |
| D-6/10/91 | 30442 | 4.5 | 7.8 | 152 | 318 | 204 | 56 | 3.5 | 1615 | 7.9 | 135 | 236 | 57 | 4.5 | 1595 | 7.9 | 103 | 198 | 74 | 76 | 0.3 | 1630 | 7.9 | 13 | 62 | 13 | 86 | 0 | 1534 | 24 | 69 | 93 | 87 | 69 | 91 | 81 | 94 | 100 |
| D-8/10/91 | 29448 | 3 | 7.6 | 115 | 272 | 266 | 38 | 2.5 | 1418 | 7.7 | 117 | 290 | 40 | 3.5 | 1385 | 7.7 | 75 | 248 | 84 | 69 | 0.2 | 1446 | 7.6 | 18 | 52 | 18 | 69 | 0 | 1498 | 36 | 71 | 96 | 76 | 79 | 84 | 81 | 93 | 99 |
| D-9/10/91 | 33623 | 0.5 | 7.8 | 135 | 366 | 240 | 49 | 4.5 | 1344 | 7.7 | 149 | 354 | 46 | 7 | 1274 | 7.7 | 112 | 220 | 110 | 60 | 0.2 | 1182 | 7.8 | 29 | 84 | 23 | 82 | 0 | 1144 | 25 | 69 | 97 | 74 | 62 | 79 | 77 | 90 | 100 |
| D-11/10/91 | 30927 | 0.4 | 7.8 | 184 | 424 | 366 | 52 | 6.5 | 1365 | 7.7 | 167 | 304 | 51 | 6 | 1384 | 7.8 | 75 | 212 | 88 | 73 | 0.2 | 1353 | 7.7 | 13 | 64 | 15 | 80 | 0 | 1374 | 55 | 71 | 97 | 83 | 70 | 93 | 85 | 96 | 100 |
| D-12/10/91 | 34823 | 0.3 | 7.9 | 170 | 332 | 226 | 65 | 5.4 | 1219 | 7.9 | 163 | 272 | 59 | 5.5 | 1223 | 7.9 | 94 | 192 | 78 | 77 | 0.2 | 1134 | 7.8 | 21 | 24 | 18 | 83 | 0 | 1105 | 42 | 71 | 96 | 78 | 88 | 88 | 93 | 92 | 100 |
| D-13/10/91 | 34018 | 0.9 | 7.6 | 153 | 464 | 274 | 56 | 5.5 | 1149 | 7.7 | 142 | 284 | 54 | 6 | 1165 | 7.7 | 82 | 228 | 74 | 76 | 0.2 | 1157 | 7.7 | 21 | 92 | 18 | 82 | 0.1 | 1180 | 42 | 74 | 97 | 74 | 60 | 86 | 80 | 93 | 98 |
| D-15/10/91 | 42876 | 0.2 | 7.8 | 133 | 349 | 310 | 48 | 6.5 | 875 | 7.8 | 128 | 356 | 41 | 6 | 870 | 7.8 | 67 | 172 | 64 | 78 | 0.2 | 887 | 7.8 | 16 | 55 | 19 | 78 | 0 | 976 | 48 | 82 | 98 | 76 | 68 | 88 | 84 | 94 | 100 |
| D-16/10/91 | 34820 | 0.3 | 8.1 | 185 | 439 | 256 | 56 | 7.5 | 2210 | 7.9 | 180 | 316 | 53 | 7 | 2070 | 7.9 | 121 | 247 | 102 | 73 | 0.3 | 1770 | 7.8 | 33 | 47 | 21 | 80 | 0 | 1539 | 33 | 68 | 96 | 73 | 81 | 82 | 89 | 92 | 100 |
| D-17/10/91 | 31780 | 0.1 | 7.7 | 175 | 457 | 262 | 60 | 4.5 | 1700 | 7.7 | 197 | 356 | 60 | 7.5 | 1698 | 7.8 | 119 | 251 | 98 | 82 | 0.3 | 1659 | 7.8 | 25 | 64 | 24 | 80 | 0.1 | 1631 | 40 | 73 | 96 | 85 | 75 | 86 | 86 | 91 | 98 |
| D-18/10/91 | 33370 | 0.1 | 7.8 | 223 | 511 | 202 | 63 | 4.7 | 1473 | 7.8 | 228 | 276 | 59 | 5.5 | 1584 | 7.7 | 122 | 273 | 86 | 72 | 0.2 | 1510 | 7.8 | 22 | 81 | 22 | 82 | 0 | 1504 | 47 | 69 | 96 | 82 | 70 | 90 | 84 | 89 | 100 |
| D-19/10/91 | 34408 | 0.3 | 8 | 174 | 442 | 268 | 58 | 5.7 | 1306 | 7.9 | 180 | 302 | 64 | 7.2 | 1316 | 7.9 | 121 | 275 | 96 | 71 | 0.4 | 1219 | 7.9 | 20 | 54 | 18 | 78 | 0 | 1241 | 33 | 68 | 94 | 84 | 80 | 89 | 88 | 93 | 100 |
| D-20/10/91 | 32720 | 1.6 | 7.8 | 235 | 489 | 252 | 65 | 4.5 | 2110 | 7.8 | 244 | 268 | 64 | 5 | 2100 | 7.9 | 135 | 303 | 88 | 68 | 0.2 | 2120 | 7.6 | 17 | 140 | 24 | 71 | 0.2 | 2080 | 45 | 67 | 96 | 80 | 54 | 89 | 71 | 91 | 96 |
| D-22/10/91 | 28707 | 0.4 | 7.8 | 117 | 296 | 142 | 72 | 3 | 1494 | 7.8 | 108 | 158 | 67 | 3 | 1516 | 7.8 | 96 | 250 | 68 | 91 | 0.2 | 1564 | 7.8 | 23 | 77 | 19 | 80 | 0.1 | 1680 | 11 | 57 | 93 | 76 | 69 | 80 | 74 | 87 | 97 |
| D-23/10/91 | 36182 | 1.5 | 7.7 | 195 | 380 | 216 | 67 | 4.5 | 1542 | 7.8 | 170 | 240 | 67 | 3.5 | 1311 | 7.8 | 119 | 238 | 82 | 81 | 0.2 | 1443 | 7.9 | 37 | 100 | 18 | 89 | 0.1 | 1365 | 30 | 66 | 94 | 69 | 58 | 81 | 74 | 92 | 99 |
| D-24/10/91 | 34364 | 1.2 | 7.9 | 191 | 400 | 184 | 74 | 6.5 | 1384 | 7.8 | 192 | 164 | 67 | 5 | 1461 | 7.8 | 114 | 274 | 66 | 97 | 0.4 | 1480 | 7.8 | 12 | 84 | 15 | 93 | 0.1 | 1463 | 41 | 60 | 92 | 90 | 70 | 94 | 79 | 92 | 99 |
| D-25/10/91 | 35400 | 0.7 | 7.6 | 156 | 364 | 194 | 64 | 5.5 | 1680 | 7.6 | 169 | 222 | 56 | 5 | 1637 | 7.6 | 89 | 274 | 86 | 61 | 0.3 | 1537 | 7.7 | 21 | 64 | 18 | 78 | 0 | 1840 | 47 | 61 | 94 | 76 | 70 | 87 | 82 | 91 | 100 |
| D-26/10/91 | 30964 | 3.3 | 7.7 | 220 | 540 | 184 | 62 | 3.5 | 1445 | 7.7 | 197 | 140 | 59 | 3 | 1414 | 7.6 | 119 | 207 | 86 | 81 | 0.2 | 1428 | 7.9 | 16 | 27 | 15 | 71 | 0.2 | 1337 | 40 | 39 | 93 | 85 | 87 | 93 | 95 | 92 | 96 |
| D-27/10/91 | 35573 | 7.3 | 7.6 | 176 | 333 | 178 | 64 | 3.5 | 1627 | 7.7 | 170 | 178 | 67 | 3 | 1684 | 7.6 | 119 | 247 | 106 | 77 | 0.2 | 1720 | 7.6 | 16 | 67 | 23 | 80 | 0.1 | 1799 | 40 | 40 | 95 | 85 | 73 | 91 | 80 | 91 | 99 |
| D-29/10/91 | 29801 | 1.6 | 7.7 | 172 | 400 | 136 | 70 | 1.5 | 1402 | 7.7 | 182 | 178 | 65 | 2.5 | 1417 | 7.8 | 123 | 263 | 106 | 70 | 0.2 | 1421 | 7.7 | 15 | 59 | 22 | 82 | 0.1 | 1468 | 32 | 40 | 88 | 88 | 78 | 91 | 85 | 84 | 97 |
| D-30/10/91 | 31524 | 1.6 | 7.9 | 183 | 478 | 204 | 65 | 6 | 1798 | 7.9 | 197 | 214 | 66 | 7.2 | 1814 | 7.9 | 119 | 365 | 120 | 70 | 2.5 | 1713 | 7.8 | 18 | 90 | 21 | 86 | 0 | 1568 | 40 | 44 | 65 | 85 | 75 | 90 | 81 | 90 | 98 |
| D-1/8/91 | 29834 | 3 | 7.4 | 160 | 348 | 194 | 62 | 3 | 1720 | 7.5 | 148 | 172 | 65 | 2.5 | 1729 | 7.6 | 105 | 265 | 80 | 83 | 0.2 | 1780 | 7.7 | 15 | 95 | 28 | 77 | 0 | 1772 | 29 | 54 | 92 | 86 | 64 | 91 | 73 | 86 | 99 |
| D-2/8/91 | 28492 | 2.6 | 7.5 | 124 | 281 | 172 | 66 | 3 | 1520 | 7.5 | 117 | 172 | 65 | 2 | 1479 | 7.6 | 77 | 221 | 90 | 69 | 0 | 1535 | 7.8 | 12 | 75 | 20 | 86 | 0 | 1549 | 34 | 48 | 99 | 84 | 66 | 90 | 73 | 88 | 99 |
| D-4/8/91 | 24978 | 0.5 | 7.3 | 146 | 288 | 124 | 68 | 2 | 1210 | 7.4 | 145 | 109 | 69 | 1.3 | 1164 | 7.4 | 81 | 200 | 52 | 67 | 0.1 | 1202 | 7.5 | 20 | 84 | 29 | 83 | 0 | 1259 | 40 | 52 | 92 | 75 | 70 | 86 | 79 | 77 | 99 |
| D-5/8/91 | 29719 | 0.2 | 7.6 | 133 | 284 | 186 | 71 | 5 | 1114 | 7.6 | 136 | 194 | 68 | 2.5 | 1095 | 7.6 | 61 | 160 | 52 | 94 | 0.1 | 1076 | 7.8 | 11 | 60 | 21 | 91 | 0 | 1100 | 55 | 73 | 96 | 82 | 63 | 92 | 79 | 89 | 100 |
| D-6/8/91 | 29741 | 0.5 | 7.9 | 151 | 316 | 196 | 64 | 2.5 | 948 | 7.8 | 163 | 200 | 69 | 3 | 904 | 7.9 | 90 | 220 | 71 | 82 | 0.2 | 929 | 8 | 16 | 44 | 22 | 86 | 0 | 951 | 45 | 65 | 95 | 82 | 80 | 89 | 86 | 89 | 100 |
| D-7/8/91 | 29027 | 0.4 | 7.6 | 136 | 328 | 186 | 68 | 3 | 899 | 7.6 | 132 | 170 | 72 | 2.5 | 921 | 7.6 | 73 | 192 | 61 | 84 | 0.2 | 872 | 7.8 | 11 | 76 | 23 | 77 | 0 | 898 | 40 | 64 | 94 | 85 | 60 | 92 | 77 | 88 | 99 |
| D-8/8/91 | 30211 | 0.5 | 7.6 | 114 | 521 | 506 | 44 | 7.5 | 866 | 7.5 | 113 | 498 | 44 | 8 | 882 | 7.6 | 58 | 137 | 65 | 79 | 0.1 | 880 | 7.9 | 11 | 44 | 20 | 78 | 0 | 884 | 49 | 87 | 99 | 81 | 68 | 90 | 92 | 96 | 100 |
| D-9/8/91 | 30848 | 0.2 | 7.7 | 142 | 376 | 144 | 71 | 3 | 940 | 7.6 | 129 | 164 | 70 | 7.5 | 918 | 7.6 | 119 | 255 | 86 | 79 | 0.2 | 933 | 7.8 | 9 | 57 | 14 | 83 | 0 | 947 | 40 | 48 | 97 | 85 | 78 | 94 | 85 | 90 | 99 |
| D-11/8/91 | 17527 | 0.6 | 7.5 | 150 | 171 | 172 | 37 | 1.4 | 732 | 7.5 | 113 | 220 | 38 | 1 | 731 | 7.5 | 42 | 113 | 49 | 67 | 0.1 | 691 | 7.5 | | 39 | 16 | 85 | 0 | 728 | 63 | 78 | 90 | 74 | 66 | 93 | 77 | 91 | 99 |
| D-12/8/91 | 33331 | 0.2 | 7.6 | 92 | 233 | 234 | 38 | 1.4 | 829 | 7.6 | 103 | 172 | 57 | 1.5 | 852 | 7.6 | 65 | 167 | 97 | 59 | 0.2 | 879 | 7.7 | 8 | 47 | 18 | 78 | 0 | 929 | 37 | 44 | 87 | 88 | 72 | 91 | 80 | 92 | 99 |
| D-13/8/91 | 27998 | 0.6 | 7.5 | 138 | 268 | 154 | 66 | 1.7 | 890 | 7.5 | 105 | 166 | 64 | 1.5 | 880 | 7.7 | 65 | 157 | 97 | 65 | 0.2 | 827 | 7.8 | 8 | 33 | 13 | 85 | 0 | 858 | 38 | 42 | 87 | 88 | 79 | 94 | 88 | 92 | 99 |
| D-14/8/91 | 32845 | 0.2 | 7.6 | 84 | 251 | 98 | 71 | 2 | 866 | 7.6 | 110 | 104 | 67 | 1.5 | 877 | 7.6 | 54 | 161 | 66 | 70 | 0.3 | 840 | 7.6 | 7 | 49 | 17 | 87 | 0 | 879 | 51 | 37 | 80 | 87 | 70 | 92 | 81 | 83 | 99 |
| D-16/8/91 | 27933 | 0.2 | 7.6 | 158 | 375 | 178 | 61 | 3.5 | 1049 | 7.7 | 153 | 168 | 60 | 3 | 992 | 7.7 | 49 | 177 | 56 | 71 | 0.1 | 910 | 7.9 | 9 | 103 | 30 | 64 | 0.1 | 828 | 68 | 67 | 97 | 82 | 42 | 94 | 73 | 83 | 99 |
| D-18/8/91 | 27527 | 0.2 | 7.3 | 191 | 240 | 166 | 74 | 3 | 1072 | 7.4 | 130 | 156 | 76 | 2.5 | 1023 | 7.4 | 80 | 274 | 71 | 78 | 0.3 | 990 | 7.5 | 8 | 44 | 11 | 100 | 0 | 999 | 39 | 55 | 88 | 90 | 70 | 96 | 82 | 93 | 100 |
| D-19/8/91 | 32363 | 0.1 | 7.6 | 159 | 310 | 146 | 69 | 1.6 | 1096 | 7.6 | 131 | 166 | 71 | 1.7 | 1083 | 7.7 | 98 | 169 | 64 | 84 | 0.2 | 1112 | 7.9 | 21 | 59 | 16 | 70 | 0 | 1083 | 25 | 61 | 91 | 79 | 65 | 87 | 81 | 89 | 99 |
| D-20/8/91 | 31437 | 0.5 | 7.6 | 132 | 304 | 148 | 65 | 2 | 939 | 7.7 | 147 | 156 | 62 | 1.8 | 974 | 7.8 | 80 | 155 | 62 | 77 | 0.1 | 1008 | 7.9 | 14 | 42 | 13 | 83 | 0 | 1012 | 46 | 60 | 94 | 83 | 73 | 89 | 86 | 91 | 100 |
| D-21/8/91 | 31914 | 2 | 7.7 | 127 | 274 | 144 | 72 | 2 | 1031 | 7.6 | 124 | 162 | 69 | 3 | 1048 | 7.6 | 80 | 157 | 69 | 83 | 0.2 | 1020 | 7.8 | 9 | 35 | 16 | 83 | 0 | 1053 | 40 | 60 | 93 | 85 | 78 | 90 | 79 | 91 | 100 |
| D-22/8/91 | 28088 | 0.2 | 7.5 | 153 | 307 | 124 | 82 | 2.5 | 1044 | 7.6 | 163 | 136 | 71 | 2.5 | 1039 | 7.7 | 97 | 188 | 62 | 92 | 0.2 | 1045 | 7.9 | 10 | 46 | 12 | 90 | 0 | 1038 | 41 | 54 | 94 | 90 | 76 | 94 | 85 | 90 | 100 |
| D-23/8/91 | 27838 | 0.1 | 7.6 | 179 | 265 | 128 | 72 | 1.8 | 992 | 7.6 | 102 | 120 | 85 | 2 | 1012 | 7.7 | 88 | 188 | 66 | 85 | 0.1 | 1036 | 7.9 | 11 | 54 | 16 | 83 | 0 | 1044 | 14 | 45 | 95 | 88 | 71 | 94 | 80 | 89 | 100 |
| D-25/8/91 | 29271 | 0.4 | 7.5 | 99 | 585 | 140 | 71 | 4.5 | 962 | 7.6 | 103 | 194 | 62 | 4.5 | 966 | 7.6 | 61 | 129 | 55 | 84 | 0.3 | 993 | 7.7 | 25 | 95 | 26 | 77 | 0 | 968 | 41 | 72 | 93 | 59 | 26 | 75 | 84 | 81 | 100 |
| D-26/8/91 | 32723 | 0.2 | 7.7 | 93 | 252 | 176 | 57 | 2.3 | 894 | 7.7 | 108 | 146 | 66 | 3 | 873 | 7.7 | 63 | 224 | 55 | 78 | 0.2 | 915 | 7.9 | 19 | 54 | 6 | 100 | 0 | 942 | 40 | 62 | 93 | 70 | 76 | 80 | 79 | 97 | 100 |
| D-27/8/91 | 33535 | 0.3 | 7.8 | 192 | 346 | 172 | 69 | 4 | 988 | 7.8 | 210 | 192 | 69 | 4.5 | 991 | 7.7 | 100 | 215 | 80 | 74 | 0.1 | 966 | 7.9 | 17 | 88 | 16 | 90 | 0 | 950 | 40 | 58 | 98 | 83 | 59 | 91 | 75 | 91 | 100 |
| D-28/8/91 | 32922 | 0.3 | 7.4 | 139 | 367 | 180 | 64 | 3 | 1060 | 7.5 | 163 | 200 | 63 | 3.5 | 1040 | 7.6 | 105 | 250 | 70 | 86 | 0.1 | 1152 | 7.7 | 25 | 84 | 20 | 84 | 0 | 1136 | 40 | 65 | 97 | 76 | 66 | 82 | 77 | 89 | 99 |

| D-29/8/91 | 32190 | 0.3 | 7.3 | 200 | 545 | 258 | 65 | 4 | 1260 | 7.4 | 191 | 226 | 67 | 3.5 | 1198 | 7.5 | 115 | 244 | 77 | 77 | 0.1 | 1351 | 7.7 | 21 | 71 | 27 | 71 | 0 | 1326 | 40 | 66 | 97 | 82 | 71 | 90 | 87 | 90 | 100 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D-30/8/91 | 30488 | 0.2 | 7.5 | 152 | 300 | 132 | 70 | 4.5 | 1073 | 7.4 | 150 | 210 | 60 | 4.5 | 1081 | 7.4 | 93 | 233 | 64 | 84 | 0.3 | 1188 | 7.3 | 17 | 55 | 18 | 80 | 0 | 1224 | 40 | 70 | 93 | 82 | 76 | 90 | 82 | 86 | 100 |