

# SAAPS

## Satellite Anomaly Analysis and Prediction System

Technical Note 3

Satellite Anomaly Prediction Module

Version 0.3

ESA/ESTEC Contract No. 11974/96/NL/JG(SC)

P. Wintoft

6 February 2001

## Document status sheet

Technical Note 3, SAAPS-SAPM		
Version	Date	Comments
0.1	28 December 1999	
0.2	17 January 2000	
0.3	6 February 2001	Added Marecs Anomaly Predictions

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	Prediction of electron fluxes . . . . .	4
1.2	Prediction of satellite anomalies . . . . .	5
<b>2</b>	<b>Artificial neural networks</b>	<b>6</b>
2.1	Mathematical notation . . . . .	6
2.2	Multi-layer feed-forward NN . . . . .	6
2.2.1	Error measures . . . . .	7
2.2.2	Back-propagation learning . . . . .	8
2.2.3	Back-propagation with adaptive learning coefficient . . . . .	8
2.2.4	Marquardt-Levenberg algorithm . . . . .	9
2.3	Radial basis function neural networks . . . . .	9
2.4	Normalization . . . . .	10
2.5	Weight initialization . . . . .	10
2.6	Training set selection . . . . .	10
2.6.1	Classification problem . . . . .	10
2.6.2	Continuous valued problem . . . . .	11
2.7	Training strategies . . . . .	11
2.7.1	Training, validation and test . . . . .	11
2.7.2	Training with minimum set size . . . . .	12
<b>3</b>	<b>Fuzzy systems</b>	<b>14</b>
<b>4</b>	<b>Prediction of the electron flux from solar wind data</b>	<b>15</b>
4.1	GOES > 0.6 MeV . . . . .	17
4.2	GOES > 2 MeV . . . . .	17
4.3	LANL SOPA 50-75 keV . . . . .	17
4.4	LANL SOPA 75-105 keV . . . . .	17
4.5	LANL SOPA 105-150 keV . . . . .	17
4.6	LANL SOPA 150-225 keV . . . . .	17
4.7	LANL SOPA 225-315 keV . . . . .	17
4.8	LANL SOPA 315-500 keV . . . . .	17
4.9	LANL SOPA 500-750 keV . . . . .	17

---

<b>5 Prediction of satellite anomalies</b>	<b>18</b>
5.1 Daily predictions . . . . .	18
5.1.1 Meteosat-3 anomalies predicted using $\sum K_p$ . . . . .	18
5.1.2 Marecs-A anomalies predicted using $\sum K_p$ . . . . .	22
5.1.2.1 One day forecast . . . . .	22
5.1.2.2 Nowcast . . . . .	25
5.1.2.3 Summary . . . . .	30

# Chapter 1

## Introduction

This document will describe models and prediction techniques that could be useful for the SAAPS.

Chapters 1.1 and 1.2 examines past work done in the predictions of satellite anomalies and energetic electron flux in the magnetosphere. Chapters 2 and 3 describes the neural networks and the fuzzy systems. Chapters 4 and 5 explores the prediction of electron flux and satellite anomalies.

### 1.1 Prediction of electron fluxes

There are a large number of papers on the subject of energetic electron fluxes in the magnetosphere. Surprisingly, there are only few papers that examine the prediction of the electron flux. [Koons and Gorney, 1991] developed a neural network to predict the daily average electron flux for energies  $> 3$  MeV at geosynchronous orbit. The electron data was taken from the spectrometer for energetic electrons (SEE) on the (LANL?) 1982-019 satellite. The input to the network was a time delay line over the past 10 days of the daily sum  $K_p$ . The model was trained so that the day of the prediction was the same as the last day of the sum  $K_p$ , thus the model was trained to perform nowcasting. Then the model was tested to make one-day-ahead forecasts. The network model showed that the predictions were significantly more accurate than a linear filter. This model was extended to make predictions one day ahead and also included the electron flux itself at the input [Stringer and McPherron, 1993]. One-hour-ahead predictions of the GOES-7 electron fluxes has also been examined [Stringer et al., 1996]. The input to the network was the  $Dst$ ,  $K_p$ , the hourly average electron flux, and magnetic local time (MLT). [Freeman et al., 1998] developed a neural network to predict the slope and intercept of the electron power law using  $Dst$  and local time as inputs. The intercept ( $B$ ) and the slope ( $M$ ) relates the electron energy ( $E$ , in keV) to the differential flux ( $F$ , electrons/cm<sup>2</sup> s sr keV)

according to

$$\log F = B + M \log E.$$

The power law is valid for electron energies 100 keV to 1.5 MeV. The electron data was taken from the (LANL?) 1984-129 satellite.

## 1.2 Prediction of satellite anomalies

To predict satellite anomalies introduces another level of difficulty as compared to predicting the space environment. The occurrence of an anomaly depends also on the design and age of the satellite, and different anomaly types have different origin.

[*Andersson et al.*, 1999] ...

## Chapter 2

# Artificial neural networks

### 2.1 Mathematical notation

The notation and mathematics follow as closely as possible as that used by [Haykin, 1994].

### 2.2 Multi-layer feed-forward NN

The multi-layer feed-forward neural network (MLFFNN) consists of several layers of neurons. The activation, i.e. the weighted sum of the inputs, for neuron  $j$  at layer  $l$  is

$$v_j^{(l)}(n) = \sum_{i=0}^p w_{ji}^{(l)} y_i^{(l-1)}(n). \quad (2.1)$$

The input at neuron  $i$  at layer  $l-1$  for pattern  $n$  is  $y_i^{(l-1)}(n)$ . The activation is calculated by summing over all inputs ( $p$  neurons) multiplied with the weights  $w_{ji}^{(l)}$ . Then the activation is transformed to give the output at neuron  $j$  at layer  $l$

$$y_j^{(l)}(n) = \varphi(v_j^{(l)}(n)), \quad (2.2)$$

where the activation function  $\varphi$  can be chosen in several ways. It depends on the layer and the problem type.

For hidden layers, we will always chose the hyperbolic tangent for the activation function

$$\varphi(v) = a \tanh bv = a \frac{1 - e^{-bv}}{1 + e^{-bv}}. \quad (2.3)$$

The constants  $a$  and  $b$  determine the amplitude and slope of the activation function. When  $v \rightarrow \pm\infty$  then  $\varphi \rightarrow \pm a$ . The slope at  $\varphi(0)$  is  $ab$ . In practice,  $a$  and  $b$  are often set to 1. However, it is interesting to examine when  $b > 1$ , especially for classification problems.

For the output layer we will either choose the tanh function or a linear function

$$\varphi(v) = bv. \quad (2.4)$$

The choice depends on whether we want to solve a curve fitting problem or a classification problem. For a curve fitting, or function approximation, problem it is generally better to choose a linear output function in order to avoid saturation at large values. In the case of a classification problem the tanh output function should be used as we wish to push the output to either false or true ( $-a$  or  $+a$ ).

In most cases it is sufficient to have one hidden layer. This type of network will solve any continuous valued function approximation problem [Cybenko, 1989] or any convex classification problem. A network with two hidden layers will solve any (continuous or discontinuous) function approximation problem or any (convex or non-convex) classification problem.

### 2.2.1 Error measures

To train a MLFFNN known input-output pairs must be available. If we call the desired output for output neuron  $j$  to be  $d_j(n)$  then the error between the desired output and the network output becomes

$$e_j(n) = d_j(n) - o_j(n), \quad (2.5)$$

where  $o_j(n) = y_j^{(L)}(n)$  and  $L$  is the output layer.

The standard error measure used in back-propagation learning is the summed squared error. The instantaneous sum of squared errors is

$$E(n) = \frac{1}{2} \sum_j e_j^2(n). \quad (2.6)$$

The sum over  $j$  includes all the neurons in the output layer. Also summing over all pattern  $n$  we get the sum of squared errors

$$E = \sum_{n=1}^N E(n). \quad (2.7)$$

Finally, we may also use the average of squared errors

$$E_{\text{av}} = \frac{1}{N} E. \quad (2.8)$$

Other error measures can also be used.



### 2.2.2 Back-propagation learning

The goal of the error-back-propagation algorithm is to minimize the error  $E$  by adjusting the weights  $w_{ji}$ . This is achieved by changing the weights so that the error gradient  $\partial E/\partial w_{ji}$  is negative. As the error is only known at the output layer, the errors have to be back-propagated through the preceding layers to be able to calculate an error gradient. The weights are updated according to

$$w_{ji}^{(l)}(n+1) = w_{ji}^{(l)}(n) + \Delta w_{ji}^{(l)}(n), \quad (2.9)$$

where

$$\Delta w_{ji}^{(l)}(n) = \eta \delta_j^{(l)}(n) y_i^{(l-1)}(n) + \alpha \Delta w_{ji}^{(l)}(n-1). \quad (2.10)$$

Here  $\eta$  is the learning coefficient and  $\alpha$  is the momentum term. The local gradient  $\delta_j^{(l)}(n)$  is

$$\delta_j^{(L)}(n) = e_j^{(L)}(n) \varphi'(v_j^{(L)}(n)) \quad (2.11)$$

for a neuron at the output layer  $L$  and

$$\delta_j^{(l)}(n) = \varphi'(v_j^{(l)}(n)) \sum_k \delta_k^{(l+1)}(n) w_{kj}^{(l+1)}(n) \quad (2.12)$$

for a neuron at hidden layer  $l$ . The prime function for a tanh activation function is

$$\varphi'(v) = ab(1 - \varphi^2(v)). \quad (2.13)$$

### 2.2.3 Back-propagation with adaptive learning coefficient

The standard back-propagation algorithm is slow. To speed things up the learning coefficient  $\eta$  can be varied during training. The idea is that each weight should have its own learning coefficient. The coefficient should increase when the error gradient  $\partial E/\partial w_{ji}$  is negative, and decrease when the gradient is positive. This can be achieved with

$$\Delta \eta_{ji}(n) = \begin{cases} \kappa & \text{if } S_{ji}(n-1) D_{ji}(n) > 0 \\ -\beta \eta_{ji}(n) & \text{if } S_{ji}(n-1) D_{ji}(n) < 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.14)$$

where  $S_{ji}(n-1)$  and  $D_{ji}(n)$  are defined as, respectively

$$D_{ji}(n) = \frac{\partial E(n)}{\partial w_{ji}(n)} \quad (2.15)$$

and

$$S_{ji}(n-1) = (1 - \xi) D_{ji}(n-1) + \xi S_{ji}(n-1) \quad (2.16)$$

where  $\xi$  is a positive constant. The procedure is called the delta-bar-delta learning rule [Jacobs, 1988].

## 2.2.4 Marquardt-Levenberg algorithm

Both the standard backpropagation (BP) and the adaptive learning rate BP algorithms may demand long training times. A more efficient training procedure is the Marquardt-Levenberg backpropagation (MBP) algorithm [Hagan and Menhaj, 1994]. MBP is an approximation to the Newton method

$$\Delta \mathbf{w} = -[\nabla^2 E(\mathbf{w})]^{-1} \nabla E(\mathbf{w}), \quad (2.17)$$

where  $\mathbf{w}$  are the weights and the error function  $E$  is the sum of squared errors (Eq. 2.7). With this error function the terms in Eq. 2.17 can be rewritten using the Jacobian matrix. Neglecting higher order terms the Newton method becomes the Gauss-Newton method. MBP is a modification of the Gauss-Newton method and can be expressed as

$$\Delta \mathbf{w} = [J^T(\mathbf{w})J(\mathbf{w}) + \mu I]^{-1} J^T(\mathbf{w})\mathbf{e}(\mathbf{w}). \quad (2.18)$$

The parameter  $\mu$  is multiplied by a factor  $\beta$  whenever a step would result in an increase of the error function  $E$ , and divided by  $\beta$  whenever  $E$  decreases. When  $\mu$  is large the algorithm becomes the steepest descent with learning rate  $\eta = 1/\mu$ , and when  $\mu$  is small the algorithm becomes the Gauss-Newton.

The Jacobian matrix is

$$J(\mathbf{x}) = \begin{bmatrix} \frac{\partial e_1(\mathbf{w})}{\partial w_1} & \frac{\partial e_2(\mathbf{w})}{\partial w_1} & \cdots & \frac{\partial e_N(\mathbf{w})}{\partial w_1} \\ \frac{\partial e_1(\mathbf{w})}{\partial w_2} & \frac{\partial e_2(\mathbf{w})}{\partial w_2} & \cdots & \frac{\partial e_N(\mathbf{w})}{\partial w_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial e_1(\mathbf{w})}{\partial w_n} & \frac{\partial e_2(\mathbf{w})}{\partial w_n} & \cdots & \frac{\partial e_N(\mathbf{w})}{\partial w_n} \end{bmatrix}. \quad (2.19)$$

The number of columns  $N$  is equal to the product of the number of training examples  $Q$  and the number of layers  $L$ , thus  $N = Q \times L$ . The number of rows  $n$  is equal to the total number of weights in the network.

The advantage of the MBP is the much reduced training time. However, it also demands more memory to store the elements of the Jacobian matrix. We see that the size of the matrix depends both on the number of training examples and the total number of weights.

## 2.3 Radial basis function neural networks

The radial basis function neural network (RBFNN) is also a feed forward network, however, the hidden layer neurons use local transformation functions instead of the global functions used in the MLFFNN. The network has the following expression

$$F(\mathbf{x}) = \sum_{i=1}^N w_i \varphi(\|\mathbf{x} - \mathbf{x}_i\|), \quad (2.20)$$

where  $\mathbf{x}_i$ ,  $i = 1, 2, \dots, N$ , are the centres of the radial basis functions  $\varphi$ . The norm  $\|\cdot\|$  is taken as the Euclidean length. The activation function is the Gaussian function

$$\varphi(r) = \exp\left(-\frac{r^2}{2\sigma^2}\right), \quad (2.21)$$

where  $\sigma$  is the width of the function.

## 2.4 Normalization

## 2.5 Weight initialization

[*Nguyen and Widrow, 1990*]

## 2.6 Training set selection

### 2.6.1 Classification problem

To train the network we need to carefully select the examples that should go into the training set. Ideally we would like to find the decision boundary in the input space that separates the two classes from each other (e.g. anomaly and non-anomalies). According to [*Swingler, 1996*], it is not desirable to have the distribution of training examples equal to the the distribution of real world examples, as the class with more training examples will dominate the training. As an example we know that the Meteosat-3 satellite experienced an anomaly day in about 20% of the time. Thus by always saying that we will have no anomalies we would be correct in 80% of the time. A network trained on such a set will thus bias its performance to predict the no-anomaly events. This will be further examined in Chapter 5.

The goal of the training algorithm is to minimize the error between the desired output and the network output. The measure used is most often the summed squared error (Eq. 2.7). Assume we have two classes, class  $A$  and class  $B$ . The total number of class  $A$  examples is  $n_A$  and the total number of  $B$  examples is  $n_B$ . Then the total number of examples is

$$n = n_A + n_B \quad (2.22)$$

and the fraction of examples of each class is

$$\begin{aligned} f_A &= \frac{n_A}{n} \\ f_B &= \frac{n_B}{n} \end{aligned}, \quad (2.23)$$

where  $f_A + f_B = 1$ . Assume also that a fraction  $c_A$  of the examples in class  $A$  are correctly classified, and  $c_B$  of the examples in class  $B$ . It should be

noted that  $c_A \leq 1$  and  $c_B \leq 1$ , and that  $c_A + c_B \leq 2$ . The fraction of correct classification then becomes

$$c = f_{ACA} + f_{BCB} \quad (2.24)$$

$$= f_{ACA} + (1 - f_A)c_B. \quad (2.25)$$

This should now be compared with the summed squared error defined in Eq. 2.7. Rewriting Eq. 2.7 in terms of the individual contributions from the two classes we get

$$E = E_A + E_B, \quad (2.26)$$

where

$$\begin{aligned} E_A &= \frac{1}{2} \sum_{n \in n_A} (d(n) - y(n))^2 \\ E_B &= \frac{1}{2} \sum_{n \in n_B} (d(n) - y(n))^2. \end{aligned} \quad (2.27)$$

If we assume that we have a tanh transfer function at the output layer then the squared error  $(d(n) - y(n))^2$  is either 0 for a correct classification or 4 for an incorrect classification. Thus if all  $A$  classification are wrong then  $E_A = 1/2 n_A 4 = 2n_A$ . Thus Eq. 2.26 can be written as

$$E = E_A + E_B \quad (2.28)$$

$$= 2n_A(1 - c_A) + 2n_B(1 - c_B) \quad (2.29)$$

$$= 2n[f_A(1 - c_A) + f_B(1 - c_B)] \quad (2.30)$$

$$= 2n[f_A(1 - c_A) + (1 - f_A)(1 - c_B)] \quad (2.31)$$

$$= 2n[f_A + (1 - f_A) - f_{ACA} - (1 - f_A)c_B] \quad (2.32)$$

$$= 2n[1 - \{f_{ACA} + (1 - f_A)c_B\}]. \quad (2.33)$$

The goal of the training algorithm is to minimize  $E$  which, from Eq. 2.33, is identical to maximizing the total number of correct classifications  $c$ . Thus, depending on how we choose  $f_A$  it will effect the maximum value of  $c$ , and also the values of  $c_A$  and  $c_B$ .

## 2.6.2 Continuous valued problem

## 2.7 Training strategies

### 2.7.1 Training, validation and test

The data set that has been created from the development of the neural network should be divided into three subsets: a training set, a validation set, and a test set. The training set is used to train the network, i.e. the process when the weights are changed so that the error is minimized. The validation set is used to determine the optimal network with respect to the number of hidden neurons and the selection of input parameters. It can also be used to decide when to stop training. Finally, when the optimal network has been selected

it can then be tested on the test set to find the performance that can be expected in a real situation. The three sets should not have any overlapping examples. The sets should also have similar statistics.

The three sets should contain a sufficient number of examples so that the network can be properly trained, validated and tested. The number of examples needed is problem related; a small simple network needs fewer examples than a large complex network.

### 2.7.2 Training with minimum set size

Sometimes it is not possible to divide the data set into the three sets because the set is too small. Instead one can use only a training set and then estimate the size of the network that will generalize well given the size of the data set and the performance of the network. According to [Baum and Haussler, 1989] a binary classifying network (output 0 or 1) with one hidden layer will almost certainly provide good generalization if the number of examples in the training set  $N$  satisfies

$$N \geq \frac{32W}{\varepsilon} \ln \frac{32M}{\varepsilon} \quad (2.34)$$

where  $W$  is the number of weights in the network and  $M$  is the number of hidden neurons. If  $\varepsilon$  is small then the network performs very well on the training set. If the training set is small ( $N$  small) it is probable that the network simply has memorized the data and will thus generalize poorly. On the other hand if  $N$  is large it is likely that the network has found the underlying general relationship. Equation 2.34 represents the worst-case formula for estimating the training set size. In practice it seems that it is sufficient to keep the first term and dropping the constants [Haykin, 1994] so that we get

$$N > \frac{W}{\varepsilon}. \quad (2.35)$$

We can now illustrate equations 2.34 and 2.35 with the anomalies on Meteosat-3. During its operation the satellite experienced about 700 anomalies. If we assume a network with 10 input units (e.g.  $\sum K_p$  over 10 days) and 10 hidden units then  $M = 10$  and  $W = (10 + 1) \times 10 + 10 + 1$ . The minimum number of examples needed as a function of  $\varepsilon/2$  would vary according to Figure 2.1. Using the criterion from Equation 2.34 we see that  $N$  should be much larger than the 700 available examples, thus it can not be guaranteed that the network will generalize. Using the relaxed criterion (Eq. 2.35) we now see that the Meteosat-3 anomaly set is sufficiently large to be used for training when the fraction of errors ( $\varepsilon/2$ ) are larger than 0.1.

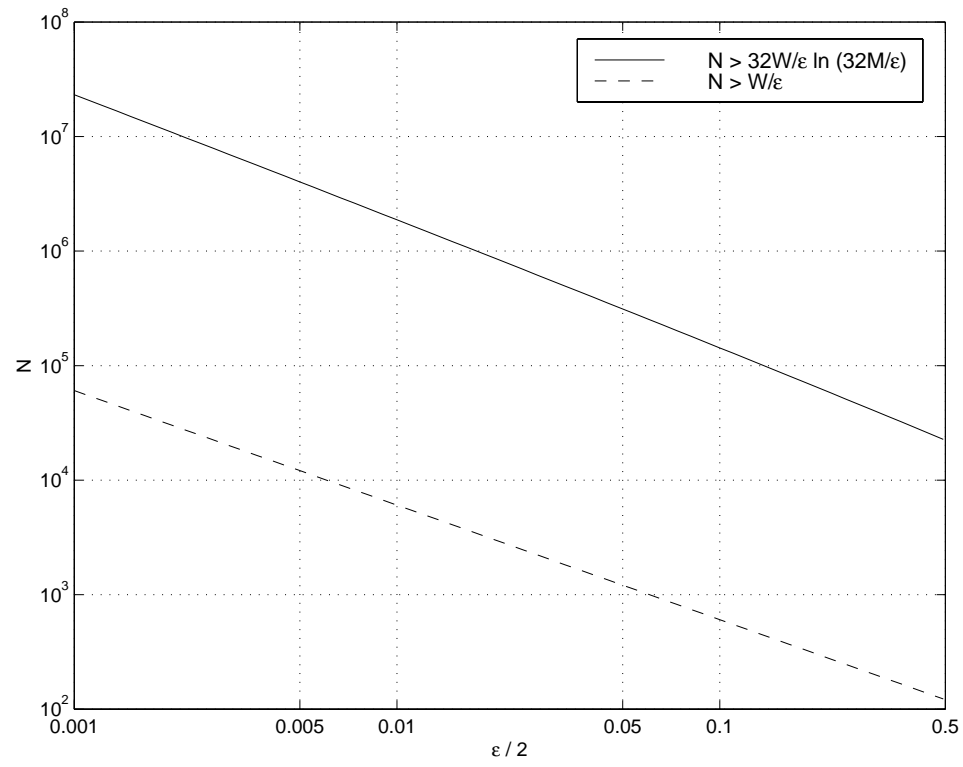


Figure 2.1: The minimum number of examples ( $N$ ) as a function of the fraction of errors ( $\epsilon/2$ ) for a network with 10 input units and 10 hidden units.

## Chapter 3

# Fuzzy systems

## Chapter 4

# Prediction of the electron flux from solar wind data

As described in section 1.1, past attempts to predict the electron flux at geosynchronous orbit has always used the local electron flux and/or geomagnetic indices ( $Kp$ ,  $Dst$ ). The emphasis has been on high energies (MeV) which are believed to give internal charging.

As it is the solar wind that drives magnetospheric activity it is natural to develop models to predict the electron flux using solar wind data. As the electron flux can vary by several orders of magnitude within 24 hours the time resolution of the predictions should be better than one day. If a time resolution of one hour is used the diurnal variation is captured. At the same time we avoid difficulties associated with substorm dynamics and the evolution of the solar wind from L1 to the Earth.

As the satellite measures the electron flux at a single point and at the same time moves in the 24-hour orbit around the Earth, it is not possible to distinguish between spatial and temporal variations. This problem can be solved in two ways: (1) The input is the solar wind data with 24 different networks predicting the flux in each local time sector; (2) The input is the solar wind data and the local time, and only one network predicting the flux in all local time sectors. It is difficult to say which method will work best. The first method produces simpler networks as each network only has to model one time sector. However, the training procedure is more complicated as 24 networks need to be trained. The number of available training examples will also be reduced by a factor of 24. The second method will produce a more complex network as it also has to model the local time variation. The training procedure is simpler, only one network, but the final prediction accuracy might be poorer as a collection of simpler specialized networks usually perform better than one general network. Both approaches shall be examined. The first approach is also similar to the hybrid neural network [Haykin, 1994].



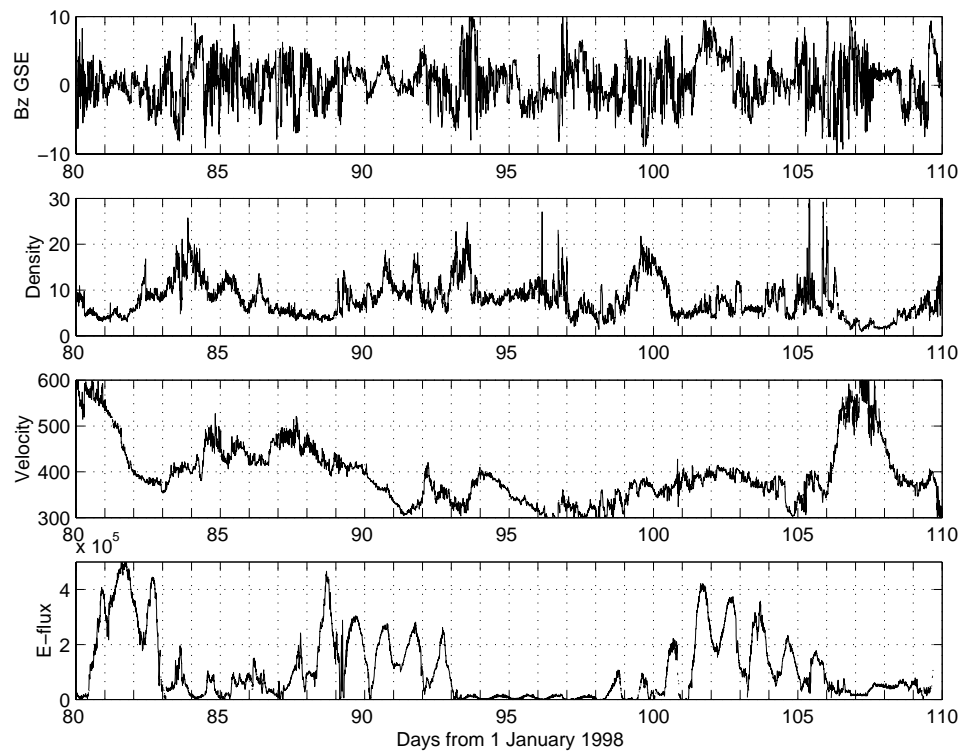


Figure 4.1: The solar wind magnetic field ( $B_z$ ), density, velocity, and the  $> 0.6$  MeV electron flux over 30 days in 1998.

- 4.1 GOES > 0.6 MeV
- 4.2 GOES > 2 MeV
- 4.3 LANL SOPA 50-75 keV
- 4.4 LANL SOPA 75-105 keV
- 4.5 LANL SOPA 105-150 keV
- 4.6 LANL SOPA 150-225 keV
- 4.7 LANL SOPA 225-315 keV
- 4.8 LANL SOPA 315-500 keV
- 4.9 LANL SOPA 500-750 keV

## Chapter 5

# Prediction of satellite anomalies

The prediction of satellite anomalies is a binary yes/no-problem, or a classification problem. This can be compared to the prediction of the electron flux which is a continuous valued problem. The prediction of an anomaly is either completely correct or completely wrong.

Essentially one would like to find a parameter space that can be separated into two classes that corresponds to anomaly and no-anomaly, respectively.

### 5.1 Daily predictions

#### 5.1.1 Meteosat-3 anomalies predicted using $\sum K_p$

As the daily average electron flux has been predicted using  $\sum K_p$  it is natural to develop an anomaly prediction model also using  $\sum K_p$ . This was examined by [Wu *et al.*, 1998]. The analysis will be repeated here, however, the definition of anomaly and non-anomaly will be slightly different. If one or several anomalies occur during a day, then that day is an anomaly day. If no anomalies occur during a day, that day is considered a non-anomaly day.

During the satellites 7 year lifetime (1988-1995) it experienced 724 anomalies. The anomalies were spread over 497 days, which should be compared to the total number of 2678 days for this period. The fraction of days with anomalies were thus 0.19. On average there were a day with anomalies every 5.4 days.

Using superposed epoch analysis one can discover trends in  $K_p$  that precede the anomalies. Figure 5.1 shows how the 3-hour  $K_p$  starts increase 6 days before the anomaly and reaches the maximum value one day ahead. However, one should also note the large variation of  $K_p$  for individual events as indicated by the standard deviation (dashed lines).

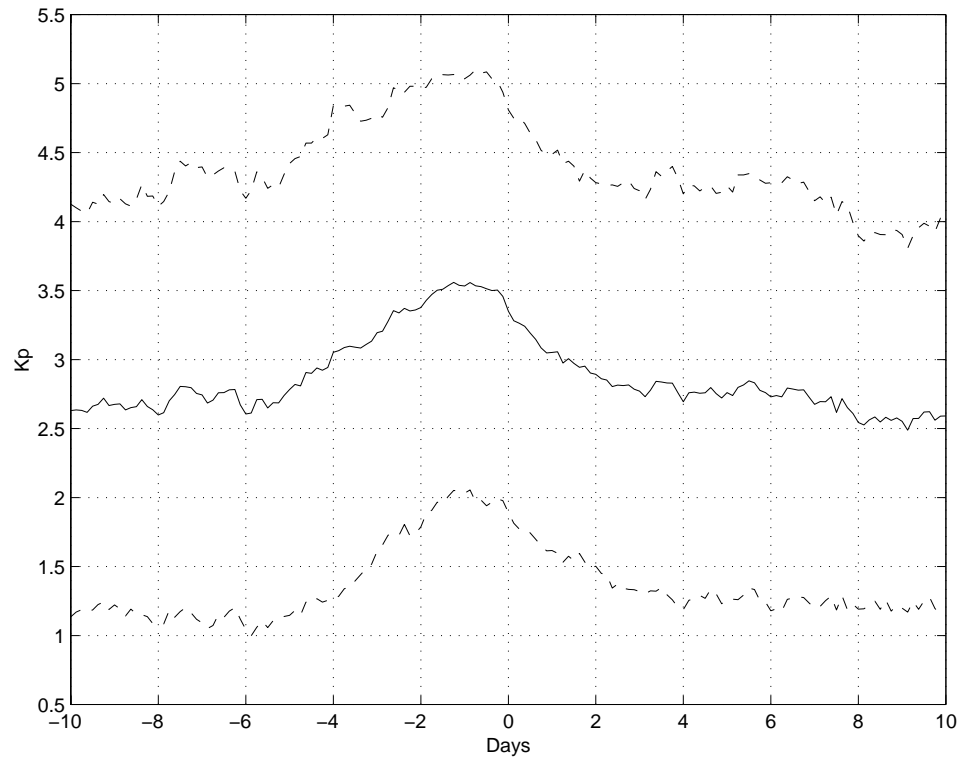


Figure 5.1: The 3-hour Kp superposed on anomaly events from Meteosat-3. The solid line shows the superposed values and the dashed lines are the standard deviation.

As was shown in Section 2.6 the minimization of the summed squared error is identical to maximizing the fraction of correct predictions  $c$ . The problem is to select a suitable training set as the number of days with anomalies is smaller than the number of days without anomalies.

To examine the effect of choosing the training and test sets with different distributions we construct a network with 8 input units, 3 hidden units and 1 output unit. The inputs are 8 days of  $\sum K_p$  and the output is whether there is an anomaly or not the next day. This network size produced similar errors on both the training set and the test set. A larger network always led to very small errors on the training set and larger errors on the test set. The total number of training examples has been fixed to  $n^{tr} = 497$  and then the fraction of anomalies in the training set  $f_A^{tr}$  has been varied over 0.1, 0.2, 0.3, 0.4, 0.5. This means that the number of anomalies in the training set  $n_A^{tr}$  has been 49, 99, 149, 198, 248. Then the networks have been tested on the test set where also the fraction of anomalies  $f_A^{ts}$  have been varied over 0.1, 0.2, 0.3, 0.4, 0.5. The result is shown in Figure 5.2. When we have a balanced training set ( $f_A^{tr} = 0.5$ ) then the fraction of correct classifications is close to constant ( $c^{ts} \approx 0.65$ ) independent of the test set distribution ( $f_A^{ts}$ ). When the training set is unbalanced ( $f_A^{tr} < 0.5$ ) makes  $c^{ts}$  dependent on the test set distribution, where  $c^{ts}$  increases for decreasing  $f_A^{ts}$ . The result is repeated in Figure 5.3 where each panel shows the variation of  $c^{ts}$  as a function of  $f_A^{ts}$  for constant  $f_A^{tr}$ . The solid-circle curves are the slices for constant  $f_A^{tr}$  from Figure 5.2. Also shown in the figure are the fractions of correct classifications for anomalies ( $c_A^{ts}$ , dash-square curve) and no-anomalies ( $c_B^{ts}$  dot-diamond curve), respectively. For the balanced training set  $c^{ts} \approx c_A^{ts} \approx c_B^{ts}$ . This we can understand from the training procedure which minimizes  $E$  (Eq. 2.7) or equivalently maximizes

$$c^{tr} = f_A^{tr} c_A^{tr} + f_B^{tr} c_B^{tr}. \quad (5.1)$$

As  $f_A^{tr} = f_B^{tr}$  both  $c_A$  and  $c_B$  will be maximized with equal weights. With an unbalanced training set  $c_A$  and  $c_B$  will be maximized unevenly, and if  $f_A^{tr} < f_B^{tr}$  then the training algorithm will bias toward maximizing  $c_B$ . This is seen in Figure 5.3 and is mostly pronounced for the case when  $f_A^{tr} = 0.1$ .

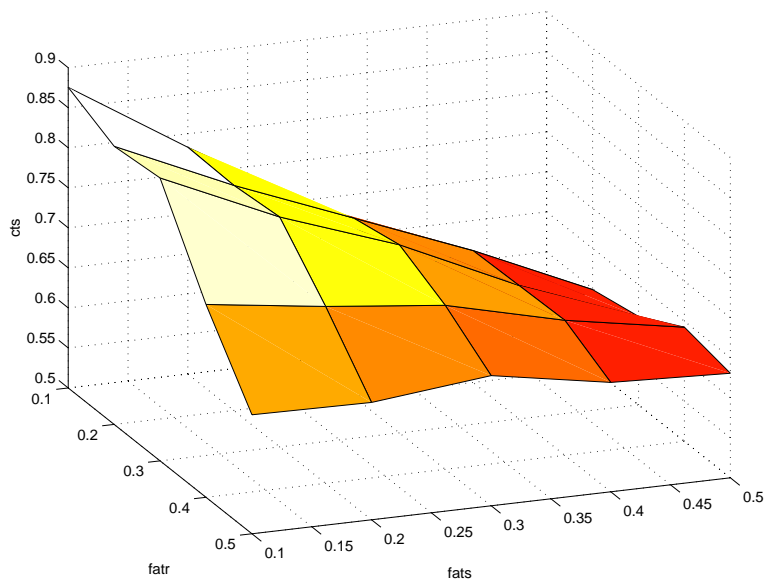


Figure 5.2: The fraction of correct predictions on the test set ( $cts$ ) as a function of the fraction of anomalies in the training set ( $fatr$ ) and the fraction of anomalies in the test set ( $fats$ ).

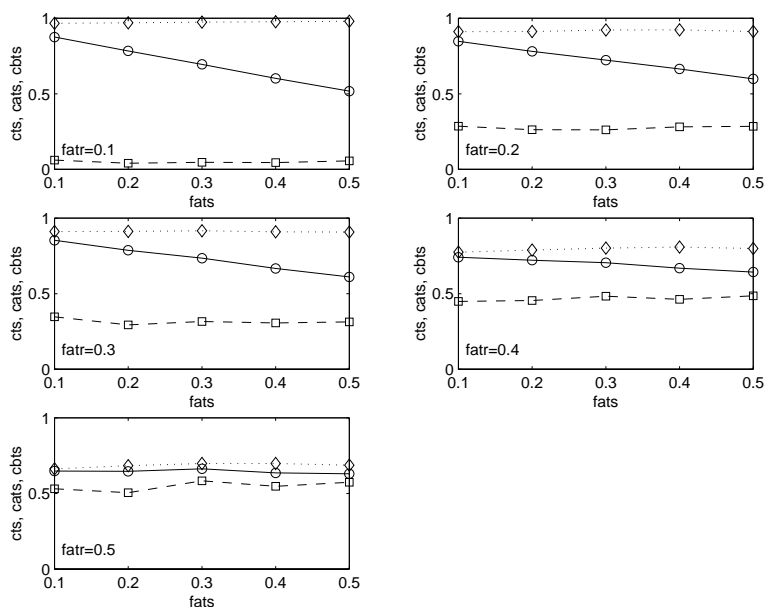


Figure 5.3: The fraction of correct predictions on the test set ( $cts$ ) as a function of the fraction of anomalies in the test set ( $fats$ ) for different values of the fraction of anomalies in the training set ( $fatr$ ) (solid curve). The fraction of correct predictions for the anomalies ( $cats$ , dashed curve) and the no-anomalies ( $cbts$ , dotted curve), respectively, are also shown.

### 5.1.2 Marecs-A anomalies predicted using $\sum K_p$

The Marecs-A anomaly data set consists of 2280 reported events over the 15 year period from 1982 to 1996. Most of the events are related to the power undervoltage indicator (A108) that tripped more than 1800 times. Daily anomaly data are created by defining an anomaly as whether there has been one or more reported events during a UT day. Subsequently, a no anomaly is defined if there are no reported events during a day. In Figures 5.1 and 5.2 in [Wintoft, 2000] the relative mutual information  $I$  between  $\sum K_p$  and daily Marecs anomalies were calculated. When all data was included a maximum of  $I = 0.125$  was achieved, while when selecting only the A108 anomalies the maximum was increased to  $I = 0.155$ , a 24 % increase. The A108 anomalies are primarily related to ESD effects and thus the stronger relationship to  $\sum K_p$ . Thus, the data set that is selected for the training of the neural network consists only of the A108 anomalies.

Using only the A108 anomalies the total data set consists of 724 anomalies. The total time period consists of 5152 days, on average the fraction of anomaly days is 0.14. As discussed in Section 5.1.1 it is important to have a balanced training set, and after balancing the data there are a total of 1448 examples with equal number of no anomalies and anomalies. This data set is divided into three sets of equal sizes: training set, validation set, and test set.

#### 5.1.2.1 One day forecast

Next, we construct a network for nowcasting and 1 day forecasts. We expect the nowcasting network to perform better than the 1 day forecast as the maximum  $I$  is reached with zero time shift between  $\sum K_p$  and the anomalies. The free parameters of the network are the time delay ( $T_1$ ) at the input and the number of hidden neurons ( $S_1$ ). The time delay is varied from -9 days up to the current time, and the number of hidden units are set to 3, 5, or 8. For each combination of  $T_1$  and  $S_1$  five networks are trained from different initial weights using the Marquardt-Levenberg algorithm. The training is stopped when the RMS error on the validation set starts increasing.

Figure 5.4 shows the result for the one day forecasts. Each point represents the RMS error for each network calculated on the three sets. The training set has the smallest error for a network with  $T_1 = -7$  days, the validation set at  $T_1 = -9$  days, and test set at  $T_1 = -3$  days. The spread of the data points in the vertical direction for each  $T_1$  represents networks with different  $S_1$  and different initial weights. The overall distribution is, however, quite flat and do not depend strongly on  $T_1$ . E.g. selecting the lowest validation error for each  $T_1$  gives only a maximum difference of less than 1 %. The choice of the value of  $T_1$  is thus not critical but we still select the network with the minimum validation error. That network have

$T_1 = -9$  days and  $S_1 = 3$  neurons. The optimal networks are summarised in Table 5.1.

Table 5.1: The optimal one day forecast networks determined from the three data set.

Data set	$T_1$	$S_1$	RMSE
Training	-7	8	0.726
Validation	-9	3	0.735
Test	-3	3	0.722

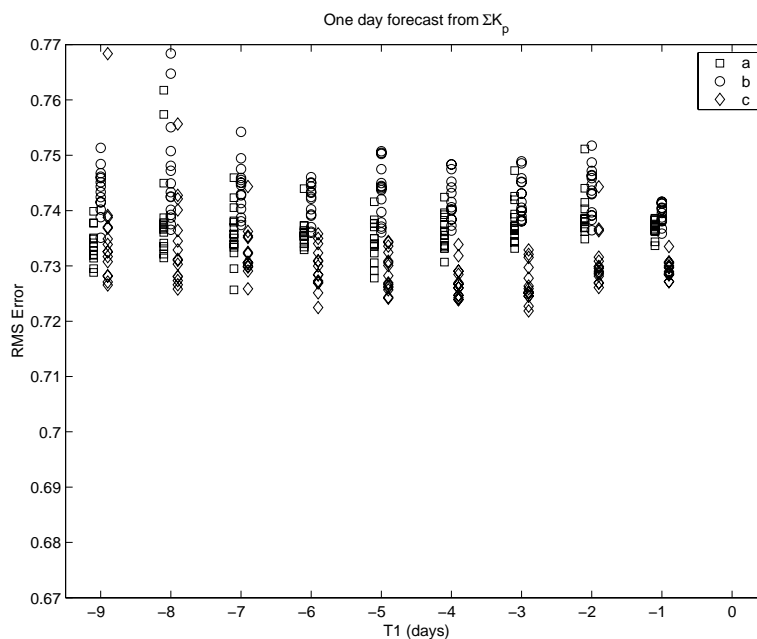


Figure 5.4: The figure shows the RMS error for the various networks that were trained for one day forecasts using  $\sum K_p$  as input. The labels are: a) training set, b) validation set, and c) test set.

We can now proceed to analyse the performance of the selected network ( $T_1 = -9$  and  $S_1 = 3$ ). We calculate the probability matrix (joint probability) for the three data sets and also for the set containing all three sets. The probabilities are shown in Table 5.2. We see that the probabilities for the different sets are quite similar with some small deviations. Approximately 36 % of the no anomalies are predicted and 30 % of the anomalies.

We can now also calculate the conditional probabilities from the joint probabilities. Let  $t$  be the observation and  $y$  the prediction from the network. Let the class  $A$  equal the “anomaly” and  $B$  the “no anomaly” classes, respectively. Then  $t$  can be either  $A$  or  $B$ . The conditional probability



Table 5.2: The probability matrices for the one day forecast network computed on the four sets: training set, validation set, test set, and combined training, validation, and test sets. Anomalies are the class A and no anomalies class B.

	Training		Validation		Test		All	
	B	A	B	A	B	A	B	A
B	0.365	0.135	0.363	0.137	0.357	0.143	0.362	0.138
A	0.212	0.288	0.181	0.320	0.184	0.316	0.192	0.308

$P(t = A|y = A)$  says “What is the probability the there is an anomaly given that the network predicts and anomaly?”. We can also turn it around to  $P(y = A|t = A)$  which states “What is the probability that the network predicts an anomaly given that there is an anomaly?”. The result is summarised in Table 5.3.

Table 5.3: The conditional probabilities for the one day forecast network computed on the four sets: training set, validation set, test set, and combined training, validation, and test sets. The observed output is  $t$  and the predicted output is  $y$ . The anomalies are class A and the no anomalies class B.

	Training	Validation	Test	All
$P(t = B y = B)$	0.633	0.668	0.660	0.653
$P(t = A y = A)$	0.681	0.700	0.689	0.690
$P(y = B t = B)$	0.730	0.726	0.715	0.724
$P(y = A t = A)$	0.577	0.639	0.632	0.616

One might suspect that some anomalies/no anomalies are easier to predict than others. Maybe some of the anomalies are more weakly related to  $\sum K_p$  than others. As the desired network output is -0.8 for no anomalies and +0.8 for anomalies we can examine how close to these desired values that the actual output is for specific forecasts. If the output is close to  $\pm 0.8$  it might be more accurate than a forecast that is close to zero. In Figure 5.5 the conditional probabilities are shown as a function of  $|y|$ . The predictions have been binned according to the limits in the figure and then the conditional probabilities have been calculated. Each bin contain approximately 300 examples, and the fraction of examples in each bin are indicated by the dashed line. When  $|y|$  is close to zero the probabilities are close to 0.5, i.e. close to pure guessing. All probabilities then increase with increasing  $|y|$  up to  $|y| < 0.5$ , and for larger  $|y|$  there is a big spread in the probabilities. The large spread can probably not be explained by the smaller number of examples in the last bin, there are still about 130 example (9 % of 1448). We see that almost all anomalies are predicted ( $P(y = A|t = A)$  is close to 1),

while very few of the no anomalies are predicted ( $P(y = B|t = B)$  is close to 0). The best result is achieved when  $0.375 < |y| < 0.5$  which happens for about 18 % of the examples. The results are summarized in Table 5.4.

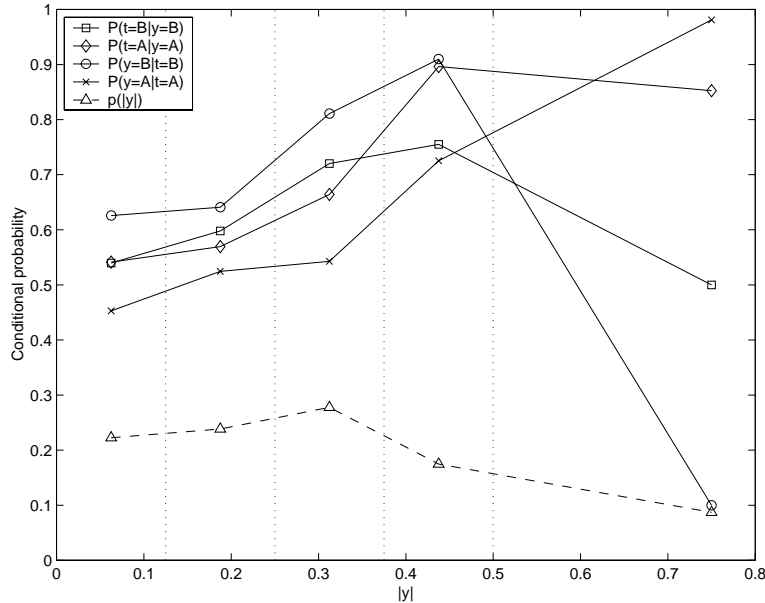


Figure 5.5: The conditional probabilities calculated for five different bins in  $|y|$ . Observations are denoted by  $t$  and predictions with  $y$ . The “no anomaly” class is called  $B$  and the “anomaly” class  $A$ . The dashed curve shows the fraction of examples in each bin.

### 5.1.2.2 Nowcast

We then repeat the same analysis for the nowcasting network and the result is shown in Figure 5.6. Comparing with Figure 5.4 it is seen that the errors are lower by about 7 %, a significant decrease. We expect that the errors should be lower from the analysis with the mutual information. The optimal networks based on the three data sets are summarised in Table 5.5. As for the one day forecast network, the RMS error for the nowcasting network is also quite insensitive on the time delay  $T_1$ .

Based on the performance from the validation set we choose the network with  $T_1 = -7$  days and  $S_1 = 3$  hidden neurons. Again we calculate the joint probabilities which are shown in Table 5.6.

The conditional probabilities are given in Table 5.7. We see that the probability that a predicted “no anomaly” is a “no anomaly” has increased by 11 % based on the “All” set. The probability that a predicted “anomaly” is an “anomaly” is increased by 3 %. However, the probability that an observed “no anomaly” is predicted as a “no anomaly” is decreased by 2 %.

Table 5.4: The conditional probabilities calculated for five different bins in  $|y|$ . Observations are denoted by  $t$  and predictions with  $y$ . The “no anomaly” class is called  $B$  and the “anomaly” class  $A$ . The last column shows the fraction of examples in each bin.

	$P(t y)$		$P(y t)$		$P( y )$
	$P(B B)$	$P(A A)$	$P(B B)$	$P(A A)$	
$0 \leq  y  < 0.125$	0.540	0.541	0.626	0.453	0.222
$0.125 \leq  y  < 0.250$	0.598	0.570	0.641	0.524	0.238
$0.250 \leq  y  < 0.375$	0.720	0.664	0.811	0.543	0.278
$0.375 \leq  y  < 0.500$	0.755	0.896	0.910	0.725	0.175
$0.500 \leq  y $	0.500	0.852	0.100	0.981	0.087

Table 5.5: The optimal nowcast networks determined from the three data set.

Data set	$T_1$	$S_1$	RMSE
Training	-9	8	0.675
Validation	-7	3	0.677
Test	-3	5	0.684

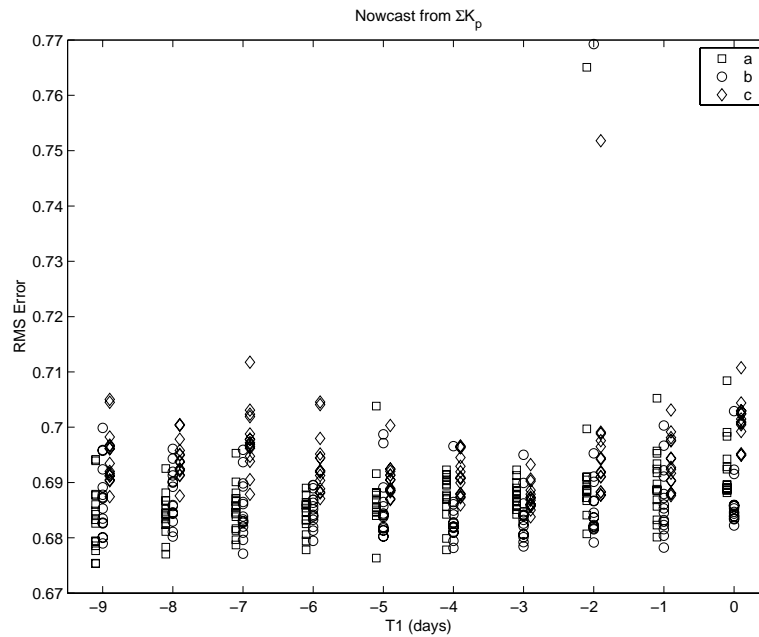


Figure 5.6: The figure shows the RMS error for the various networks that were trained for nowcasting using  $\sum K_p$  as input. The labels are: a) training set, b) validation set, and c) test set.

Table 5.6: The probability matrices for the nowcast network computed on the four sets: training set, validation set, test set, and combined training, validation, and test sets. Anomalies are the class A and no anomalies class B.

	Training		Validation		Test		All	
	B	A	B	A	B	A	B	A
B	0.348	0.152	0.357	0.143	0.353	0.147	0.353	0.147
A	0.141	0.359	0.139	0.361	0.116	0.384	0.132	0.368

The probability that an observed “anomaly” is predicted as an “anomaly” is increased by 19 %.

Table 5.7: The conditional probabilities for the nowcast network computed on the four sets: training set, validation set, test set, and combined training, validation, and test sets. The observed output is  $t$  and the predicted output is  $y$ . The anomalies are class A and the no anomalies class B.

	Training	Validation	Test	All
$P(t = B y = B)$	0.712	0.720	0.753	0.728
$P(t = A y = A)$	0.703	0.716	0.724	0.714
$P(y = B t = B)$	0.697	0.714	0.707	0.706
$P(y = A t = A)$	0.718	0.722	0.769	0.736

The variations of the conditional probabilities as a function of  $|y|$  for the nowcasting network are shown in Figure 5.7. The points in each bin now lies much closer to each other as compared to Figure 5.5. The probabilities also increase monotonically with increasing  $|y|$ . This suggests that the nowcast network works in a well behaved way, which the one day forecast network does not. For nowcasting with  $|y| > 0.375$  the predictions have probabilities of 0.8 or higher, and about 48 % of the examples are above that limit. The results are also summarized in Table 5.8.

Finally, we examine what the typical input patterns are for the situations when the predictions have outputs that are  $y < -0.5$ ,  $-0.2 < y < 0.2$ , and  $y > 0.5$ , respectively. By selecting all the input examples that give predictions in one of the three ranges we calculate the mean and standard deviation for each input. The result is shown in Figure 5.8. For inputs at time delays of about 7 days there is no significant difference for the three levels of output. But when we come closer to the present time the three curves starts deviating from each other. The clear “no anomaly” predictions have small values of  $\sum K_p$  at time 0, while the clear “anomaly” predictions have large values of  $\sum K_p$ . The hard-to-predicted cases have close to constant values of  $\sum K_p$  independent of the time delay and are very close to the average of all  $\sum K_p$  values of 22.

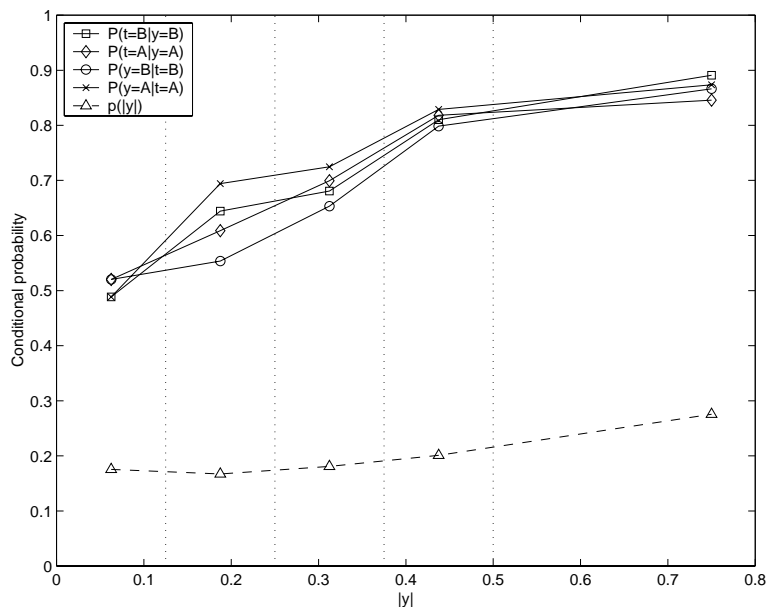


Figure 5.7: The conditional probabilities calculated for five different bins in  $|y|$ . Observations are denoted by  $t$  and predictions with  $y$ . The “no anomaly” class is called  $B$  and the “anomaly” class  $A$ . The dashed curve shows the fraction of examples in each bin.

Table 5.8: The conditional probabilities calculated for five different bins in  $|y|$ . Observations are denoted by  $t$  and predictions with  $y$ . The “no anomaly” class is called  $B$  and the “anomaly” class  $A$ . The last column shows the fraction of examples in each bin.

	$P(t y)$		$P(y t)$		$P( y )$
	$P(B B)$	$P(A A)$	$P(B B)$	$P(A A)$	
$0 \leq  y  < 0.125$	0.489	0.520	0.520	0.489	0.175
$0.125 \leq  y  < 0.250$	0.644	0.609	0.554	0.694	0.167
$0.250 \leq  y  < 0.375$	0.681	0.699	0.653	0.725	0.181
$0.375 \leq  y  < 0.500$	0.810	0.818	0.799	0.829	0.201
$0.500 \leq  y $	0.891	0.846	0.866	0.874	0.276

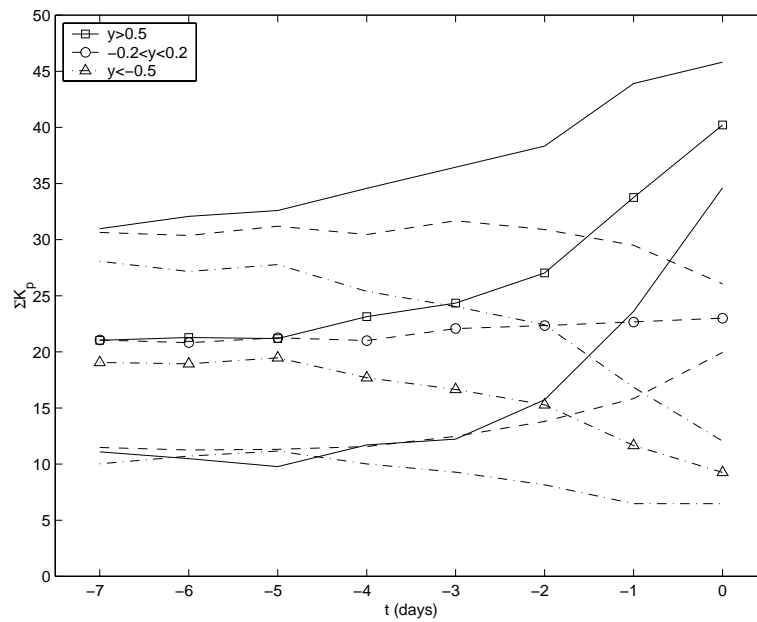


Figure 5.8: The average input values for different output values. The solid curves are for inputs which give outputs  $y > 0.5$ , the dashed curves for  $-0.2 < y < 0.2$ , and the dot-dashed curves for  $y < -0.5$ . The curves with symbols are the averages and the curves without symbols are the  $\pm 1$  standard deviation.

### 5.1.2.3 Summary

Two types of networks have been trained to predict Marecs-A A108 anomalies: one day forecast network and nowcast network. From a previous analysis it has been shown that a model for nowcasting should perform better than a forecasting network, this is also shown from the training and testing of the two networks. The nowcasting network gives, on average, a 72 % probability for correct predictions. However, for some predictions the results are poorer and for others the results are better. There exist a relation between the absolute value of the network output and the probability for correct prediction. When the output is close to zero the probability drops to 50 %, i.e. not better than guessing. When the output is below -0.375 or above +0.375 the probability increases to 84 %.

# Bibliography

- [*Andersson et al.*, 1999] Andersson, L., L. Eliasson, and P. Wintoft, Prediction of times with increased risk of internal charging on spacecraft, *Workshop on space weather*, ESTEC, Noordwijk, WPP-155, 427–430, 1999.
- [*Baum and Haussler*, 1989] Baum, E.B., and D. Haussler, What size net gives valid generalization?, *Neural Computation* 1, 151–160, 1989.
- [*Cybenko*, 1989] Cybenko, G., Approximation by superpositions of a sigmoidal function, *Math. Control Signals Systems* 2, 303–314, 1989.
- [*Hagan and Menhaj*, 1994] Hagan, M.T., and M. Menhaj, Training feed-forward networks with the Marquardt algorithm, *IEEE Trans. Neural Networks*, 5, 989–993, 1994.
- [*Haykin*, 1994] Haykin, S., *Neural Networks: A Comprehensive Foundation*, Macmillan College Publishing Company, Inc., New York, 1994.
- [*Jacobs*, 1988] Jacobs, R.A., Increased rates of convergence through learning rate adaption, *Neural Networks* 1, 295–307, 1988.
- [*Freeman et al.*, 1998] Freeman, J.W., T.P. O’Brien, A.A. Chan, and R.A. Wolf, Energetic electrons at geostationary orbit during the November 3-4, 1993 storm: Spatial/temporal morphology, characterization by a power law spectrum and, representation by an artificial neural network, *J. Geophys. Res.*, 103, 26,251–26,260, 1998.
- [*Koons and Gorney*, 1991] Koons, H.C., and D.J. Gorney, A neural network model of the relativistic electron flux at geosynchronous orbit, *J. Geophys. Res.*, 96, 5,549–5,556, 1991.
- [*Nguyen and Widrow*, 1990] Nguyen, D., and B. Widrow, Improving the learning speed of 2-layer neural networks by choosing initial values of the adaptive weights, in *Proc. IJCNN*, 3, 21–26, 1990.
- [*Stringer and McPherron*, 1993] Stringer, G.A., and R.L. McPherron, Neural networks and predictions of day-ahead relativistic electrons at



geosynchronous orbit, *Proceedings of the International Workshop on Artificial Intelligence Applications in Solar-Terrestrial Physics, Lund, Sweden, 22-24 September 1993*, Edited by J.A. Joselyn, H. Lundstedt, and J. Trolinger, Boulder, Colorado, 139–143, 1993.

[*Stringer et al., 1996*] Stringer, G.A., I. Heuten, C. Salazar, and B. Stokes, Artificial neural network (ANN) forecasting of energetic electrons at geosynchronous orbit, in *Radiation Belts: Models and Standards, Geophys. Monogr. Ser.*, vol 97, edited by J.F. Lemaire, D. Heynderickx, and D.N. Baker, AGU, Washington, D.C., 291–295, 1996.

[*Swingler, 1996*] Swingler, K., *Applying neural networks: a practical guide*, Academic Press Ltd, London, 1996.

[*Wu et al., 1998*] Wu, J.-G., H. Lundstedt, L. Andersson, L. Eliasson, and O. Norberg, Spacecraft anomaly forecasting using non-local environment data, *Study of plasma and energetic electron environment and effects (SPEE)*, ESTEC Technical Note, WP 220, 1998.

[*Wintoft, 2000*] Wintoft, P., SAAPS WP 120 – Technical Note 1, Satellite Anomaly Analysis Module, Version 0.2, 3 April 2000.