

NEURAL NETWORKS APPLIED TO AUTOMATIC FAULT DETECTION

Stefan Jakubek and Thomas Strasser

Institute for Machine and Process Automation
Vienna University of Technology
Gußhausstr. 27-29, 1040 Vienna, Austria

ABSTRACT

The objective of this paper is an automatic fault detection scheme for applications in the automotive industry. In order to increase data reliability and for the purpose of monitoring and control of test equipment a fault detection system based on multivariate data analysis is being developed. The detection scheme has to process up to several hundreds of different measurements at a time and check them for consistency. The main problem lies in the fact that besides the available data no further information is provided. The chosen approach models the distribution function of available fault-free data using ellipsoidal basis function networks. An important requirement for the fault detection scheme is that it should be able to automatically adapt itself to new data. The present paper is focused on this feature. It is demonstrated how a gradient optimization with algebraic constraints can be applied to adapt a pre-existing network to new data points. Numerical examples with actual data show that the proposed method produces good results.

1. INTRODUCTION

In order to achieve maximum productivity, modern measuring systems are now highly sophisticated, characterized by complex inter-relationships between sub-systems. The developers' growing need for reliable assessments gives rise to an online monitoring system that automatically checks measured data for consistency. The task is to analyze training data and to build a probability model based on these data using neural networks. However, in a high-tech environment there are often several hundreds of measurements to be processed simultaneously which entails the problem that any model has to be built in a high dimensional measurement space. The algorithm proposed in [1] uses principal component analysis to break down the high dimensionality of the problem. Since no information about functional relationships between data channels is known a-priori, a statistical approach was chosen. The distribution function of the available data-set is estimated using kernel regression techniques. Fault detection is accomplished by checking whether a new data record lies in a cluster of training data or not. Areas or clusters for valid data are determined through kernel regression. In order to reduce the computational effort the distribution functions are approximated by neural networks with ellipsoidal basis functions. In order to use as few basis functions as possible a new training algorithm for ellipsoidal basis function networks is presented. Even during normal operation of the engine test-bed new training data are provided for fault detection and it is required that the system is able to adapt itself autonomously to the new data. The present paper focuses on strategies and algorithms for recursive network training. The rest of the paper is organized as follows: Section 2 gives a brief introduction

into the fault detection algorithm in [1]. In section 3 strategies and decision layers for recursive training of an existing network are presented. Section 4 explains how single neurons or groups of neurons can be updated using gradient descent techniques. It is explained how constraints to the gradient can be used to guide the direction of convergence. Finally, the paper is concluded by some numerical results.

2. NETWORK TRAINING

After transforming the measurement channels using PCA one has to find a smooth approximation for the probability density function $f(\mathbf{x})$ where \mathbf{x} resembles an arbitrary point in the measurement space. A common way to accomplish this is placing a *kernel function* $K(\mathbf{x})$ about each data point \mathbf{x}_i [2]. In the present application Gaussian kernels were used. With d as the dimension of the data record one gets for the approximation $\hat{f}(\mathbf{x})$ of the probability density function

$$\hat{f}(\mathbf{x}) = \frac{1}{nh_1 \dots h_d} \sum_{i=1}^n \left[\prod_{j=1}^d K\left(\frac{x_j - t_{ij}}{h_j}\right) \right]. \quad (1)$$

Here n denotes the number of available records for training, x_j is the j -th component of \mathbf{x} , t_{ij} is the j -th component of the i -th data point and $h_1 \dots h_d$ are so-called *bandwidth* parameters. Taking into consideration the large number of data both the sum and the product in (1) make this formula time-consuming to evaluate. Moreover, it is difficult or even impossible to discover structures, like isolated clusters in the data. Both problems give rise to develop a powerful method to reduce the number of parameters in $\hat{f}(\mathbf{x})$ whilst keeping the degree of information almost constant. A neural network with ellipsoidal basis functions can be used to approximate the map $\mathbf{x} \rightarrow \hat{f}(\mathbf{x})$.

2.1. Neuron and network structure

In our design a neuron $n(\mathbf{x})$ has an ellipsoidal Gaussian basis function:

$$n(\mathbf{x}) = \gamma \exp\left(-\frac{1}{2} \tilde{\mathbf{x}}^T \mathbf{A} \tilde{\mathbf{x}}\right) \quad (2)$$

The weight is denoted as γ and $\tilde{\mathbf{x}}$ denotes the vector from the neuron's center $\boldsymbol{\mu}$ to \mathbf{x} : $\tilde{\mathbf{x}} = \mathbf{x} - \boldsymbol{\mu}$. The matrix \mathbf{A} is symmetric and contains "spread parameters" which are determined using a Taylor approximation of the target function about the neuron's center. With the use of this neuron structure the network output is calculated in

the following way (where n_{neu} is the number of neurons):

$$\tilde{f}(\mathbf{x}) = \sum_{i=1}^{n_{neu}} \gamma_i \exp\left(-\frac{1}{2} \tilde{\mathbf{x}}^T \mathbf{A}_i \tilde{\mathbf{x}}\right) \quad (3)$$

Now it is possible to obtain clusters in the training data for the plausibility check. Details about this algorithm can be found in [1].

3. RECURSIVE NETWORK TRAINING STRATEGIES

If a record of new data points arrives, it is subdivided into two classes. For this purpose it is investigated, whether a data point fits into one of the existing models. Every model uses one or more principal axes, that span the *model space*. If the column vectors of the matrix \mathbf{S}_i span the model space of the i -th model, and \mathbf{x}_k is a new data point, the error vector $\mathbf{e}_{i,k}$ describes the orthogonal distance from \mathbf{x}_k to the i -th model space [3]:

$$\mathbf{e}_{i,k} = \underbrace{[\mathbf{I} - \mathbf{S}_i(\mathbf{S}_i^T \mathbf{S}_i)^{-1} \mathbf{S}_i^T]}_{\mathbf{\Gamma}_i} \mathbf{x}_k \quad (4)$$

An error measure can thus be defined as

$$J_{i,k} = \|\mathbf{e}_{i,k}\|_2^2 = \mathbf{x}_k^T \mathbf{\Gamma}_i^T \mathbf{\Gamma}_i \mathbf{x}_k \quad (5)$$

The first decision layer selects those data points that can be assigned to a pre-existing model i (i.e. for these data points for a certain i $J_{i,k}$ falls below a threshold). These points can be used to update existing models. All the remaining data points lie in new areas of the measurement space and are used to train new models. They usually emerge when new domains of the operating envelope are being tested. The next decision layer looks at pre-existing models and the selected data that fit into these models. Again, the data set is separated into those points that lie in the vicinity of existing neurons and those that lie in unpopulated areas. The former are used to update the neuron-parameters, the latter are taken into account by placing new neurons.

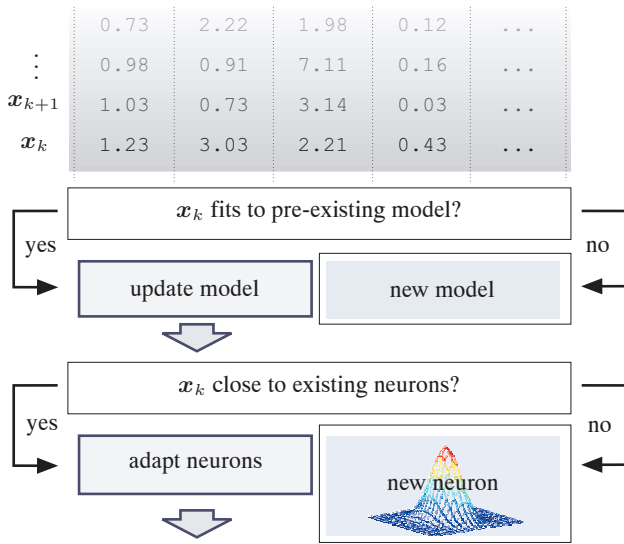


Figure 1: Decision layers for network update

4. UPDATING NEURONS

As mentioned above after new training data are available an update of the existing network is necessary. For that case we develop $n(\mathbf{x})$ into a Taylor series with respect to the parameters γ and \mathbf{A} :

$$n(\mathbf{x}) \doteq n_0(\mathbf{x}) + \left. \frac{\partial n(\mathbf{x})}{\partial \gamma} \right|_0 (\gamma - \gamma_0) + \sum_{i,j} \left. \frac{\partial n(\mathbf{x})}{\partial a_{i,j}} \right|_0 (a_{i,j} - a_{i,j,0}) \quad (6)$$

Explicit formulas for the derivatives are:

$$\left. \frac{\partial n(\mathbf{x})}{\partial \gamma} \right|_0 = \frac{n_0(\mathbf{x})}{\gamma_0}, \quad (7)$$

$$\left. \frac{\partial n(\mathbf{x})}{\partial a_{i,i}} \right|_0 = n_0(\mathbf{x})(-0.5 \tilde{x}_i^2), \quad \left. \frac{\partial n(\mathbf{x})}{\partial a_{i,j}} \right|_0 = n_0(\mathbf{x})(-\tilde{x}_i \tilde{x}_j) \quad (8)$$

With the desired output p_k of the network for data point # k the network error becomes

$$e_k = \sum_{i=1}^{n_{neu}} n_i(\mathbf{x}_k) - p_k \quad (9)$$

where $n_i(\mathbf{x}_k)$ stands for the output of neuron i . Next, the network is subdivided into those neurons to be adapted ($^a n(\mathbf{x})$) and those to remain unchanged ($^u n(\mathbf{x})$):

$$e_k = \sum_{i=n_{a1}}^{n_{al}} {}^a n_i(\mathbf{x}_k) + \underbrace{\sum_{i=n_{u1}}^{n_{um}} {}^u n_i(\mathbf{x}_k)}_{\nu_k} - p_k = \sum {}^a n(\mathbf{x}_k) - \nu_k. \quad (10)$$

Here, n_{al} and n_{um} denote the number of neurons to be adapted or unchanged respectively. After a variation of neuron parameters γ and $a_{i,j}$ the new predicted network error e_{k+} at data point k becomes

$$e_{k+} = e_k + [\nabla^a n(\mathbf{x}_k)]^T \cdot \boldsymbol{\vartheta} \quad (11)$$

where $\boldsymbol{\vartheta}$ stands for a vector of parameter variations for one neuron, $\nabla^a n(\mathbf{x}_k)$ is the column vector containing all parameter derivatives of one neuron as given in (7) to (8). Forming a column vector ${}_a \nabla(\mathbf{x}_k)$ of parameter derivatives for all neurons ${}^a n_i(\mathbf{x}_k)$ we get

$${}_a \nabla(\mathbf{x}_k) = \begin{bmatrix} \nabla^a n_{a1}(\mathbf{x}_k) \\ \nabla^a n_{a2}(\mathbf{x}_k) \\ \vdots \\ \nabla^a n_{al}(\mathbf{x}_k) \end{bmatrix} \quad (12)$$

and with a column vector $\boldsymbol{\theta}$ of parameter variations for all neurons ${}^a n(\mathbf{x}_k)$ one gets

$$e_{k+} = e_k + {}_a \nabla^T \cdot \boldsymbol{\theta} \quad (13)$$

For a gradient descent one has to choose

$$\boldsymbol{\theta} = {}_a \nabla \rho. \quad (14)$$

with $\rho \in \mathbb{R}$. Then, e_{k+} becomes

$$e_{k+} = e_k + ({}_a \nabla^T \cdot {}_a \nabla) \cdot \rho \quad (15)$$

The goal is to reduce the network error at one specific data point k by a certain percentage:

$$e_{k+} = \lambda e_k \quad (16)$$

For $\lambda = 0$ the network is completely adapted to the data point, for $\lambda = 1$ no adaptation to the current data point is carried out. With this approach we get for ρ :

$$\rho = (\lambda - 1)e_k / (\mathop{\Delta}\limits_a^T \mathop{\Delta}\limits_a) \quad (17)$$

4.1. Neuron weighting

It turned out, that the derivatives with respect to the spread parameters $\left. \frac{\partial n(\mathbf{x})}{\partial a_{i,j}} \right|_0 = n_0(\mathbf{x})(-\tilde{x}_i \tilde{x}_j)$ are almost always significantly higher than $\left. \frac{\partial n(\mathbf{x})}{\partial \gamma} \right|_0 = \frac{n_0(\mathbf{x})}{\gamma_0}$. In the vicinity of a neuron's (say Neuron n_l) center, however, this leads to the situation that the gradient with respect to the $a_{i,j}$ of a neighbouring neurons can have more influence on the network output, than parameter changes of n_l itself. Adapting only parameters of n_l also has the advantage, that the influence of the data point is more bounded. Therefore, if a data point \mathbf{x}_k lies in the vicinity of a neuron, instead of the actual gradient, an alternative gradient is used, that emphasizes this neuron. A data point \mathbf{x}_k is said to lie in the vicinity of a neuron j if

$$\exp\left(-\frac{1}{2}\tilde{\mathbf{x}}^T \mathbf{A}_j \tilde{\mathbf{x}}\right) > \xi. \quad (18)$$

Here, ξ is a "proximity threshold", in practice $\xi = 0.8$ turned out to be a good value. If the above condition (18) is fulfilled, say for neuron 2, then instead of θ as given in (14) one has to use

$$\theta = \underbrace{\begin{bmatrix} \nabla^a n_{a1}(\mathbf{x}_k) \\ 100 \nabla^a n_{a2}(\mathbf{x}_k) \\ \vdots \\ \nabla^a n_{al}(\mathbf{x}_k) \end{bmatrix}}_{\mathop{\Delta}\limits_a^{mod}} \rho. \quad (19)$$

In (19) the weight of neuron 2 raised by a factor of 100 against the remaining neurons which has the effect that mainly the parameters of this neuron are used to adapt the network to data point \mathbf{x}_k . Using (19) the scaling factor ρ now becomes

$$\rho = (\lambda - 1)e_k / (\mathop{\Delta}\limits_a^T \mathop{\Delta}\limits_a^{mod}). \quad (20)$$

4.2. Fixed points

Changes both of the neuron weight γ and of the spread parameters $a_{i,j}$ entail a more or less global change of the neuron output. In many cases this is not desirable, since we only want to adjust the network output locally. Therefore it is desirable to adapt one or more neurons to a data point as effectively as possible (e.g. using gradient methods) and at the same time ensure that the overall change of the network output remains within some limits. This can be achieved by introducing **fixed points** \mathbf{p}_f . The concept is to ensure that the output of a neuron (say neuron 2) remains unchanged at a fixed point, and, under this constraint find an "optimal" vector θ of parameter variations. Computing the gradient $\mathop{\Delta}\limits_a(\mathbf{p}_f)$ at the fixed point \mathbf{p}_f after adaptation to the k -th data point the neuron output at \mathbf{p}_f changes as follows:

$$n_{k+}(\mathbf{p}_f) - n_k(\mathbf{p}_f) \doteq \mathop{\Delta}\limits_a(\mathbf{p}_f) \cdot \theta \quad (21)$$

Since we only want to affect the output of neuron 2 $\mathop{\Delta}\limits_a(\mathbf{p}_f)$ becomes

$$\mathop{\Delta}\limits_a(\mathbf{p}_f) = \begin{bmatrix} \mathbf{0} \\ \nabla^a n_{a2}(\mathbf{p}_f) \\ \mathbf{0} \\ \vdots \end{bmatrix} \quad (22)$$

In order to leave the neuron output unchanged one simply has to ensure that $\mathop{\Delta}\limits_a(\mathbf{p}_f) \cdot \theta = 0$ i.e. θ must lie in the null space of $\mathop{\Delta}\limits_a(\mathbf{p}_f)$ as was suggested in [4, 5]:

$$\theta \in \mathcal{N}(\mathop{\Delta}\limits_a(\mathbf{p}_f)) \quad (23)$$

Usually, there is more than one fixed point to be observed so that a matrix Φ is constructed by stacking the $\mathop{\Delta}\limits_a(\mathbf{p}_{f,i})$ of all fixed points $\mathbf{p}_{f,i}$ as row vectors:

$$\Phi = \begin{bmatrix} \mathop{\Delta}\limits_a^T(\mathbf{p}_{f,1}) \\ \mathop{\Delta}\limits_a^T(\mathbf{p}_{f,2}) \\ \mathop{\Delta}\limits_a^T(\mathbf{p}_{f,3}) \\ \vdots \end{bmatrix} \quad (24)$$

Again, a suitable choice of the variation vector θ must lie in the right null -space of Φ . Let \mathbf{N} denote a matrix whose columns span $\mathcal{N}(\Phi)$, e.g. $\Phi \cdot \mathbf{N} = \mathbf{0}$. The optimal projection $\mathop{\Delta}\limits_a^{opt}$ of $\mathop{\Delta}\limits_a^{mod}$ on $\mathcal{N}(\Phi)$ thus yields (see [5])

$$\mathop{\Delta}\limits_a^{opt} = \mathbf{N}(\mathbf{N}^T \mathbf{N})^{-1} \mathbf{N}^T \cdot \mathop{\Delta}\limits_a^{mod} \quad (25)$$

Using $\theta = \rho \mathop{\Delta}\limits_a^{opt}$ one gets for ρ

$$\rho = (\lambda - 1)e_k / (\mathop{\Delta}\limits_a^T \mathop{\Delta}\limits_a^{opt}) \quad (26)$$

The way $\mathop{\Delta}\limits_a^{opt}$ is determined in (25) assures that the variation vector θ lies as closely to $\mathop{\Delta}\limits_a^{mod}$ as possible while satisfying all constraints. It remains to determine a reasonable amount of fixed points for each neuron. The following idea produced good results: If d denotes the dimensionality of the network then d fixed points are chosen on the contour ellipse defined trough $-\frac{1}{2}\tilde{\mathbf{x}}^T \mathbf{A} \tilde{\mathbf{x}} = \alpha$ with $\log \alpha \in]0; 1[$. Ideally, these fixed points are placed on the vertices of the contour ellipse (Fig. 2). Thus, the neuron is constrained along the whole contour line (indicated red in Fig. 2) during the parameter update. The logarithmic constant α determines the contour line to be fixed. The smaller α is chosen ($\sim \log(0.2)$), the less global are the effects on the neuron output due to parameter variations. On the other hand, larger values of α ($\alpha > \log(0.7)$) can lead to significant global changes of the neuron shape as can easily be imagined. There

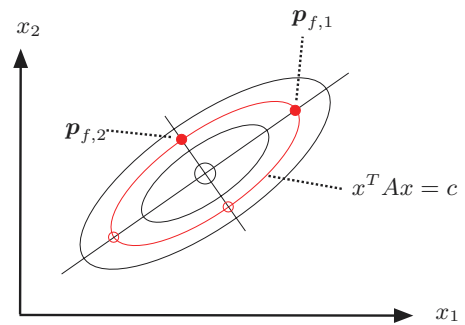


Figure 2: Choosing fixed points on contour lines

is one critical situation when a data point lies in the vicinity of a fixed contour line. Then, the denominator in (26) tends towards zero ($\mathop{\Delta}\limits_a^T \mathop{\Delta}\limits_a^{opt} \rightarrow 0$) which leads to numerical difficulties and parameter divergence. In this case one either has to do without fixed points when adapting to this particular data point (which bears the risk of unreasonable parameter changes) or the network is not updated to the data point at all. For the sake of stability of the algorithm we used the latter alternative.

5. RESULTS

Figure 3 shows an example for a one-dimensional network consisting of two neurons. The data that the network should adapt to are indicated by black dots. It can be seen that the fixed points keep the effects of adaptation local to the data points. Figures 4 and 5 depict a

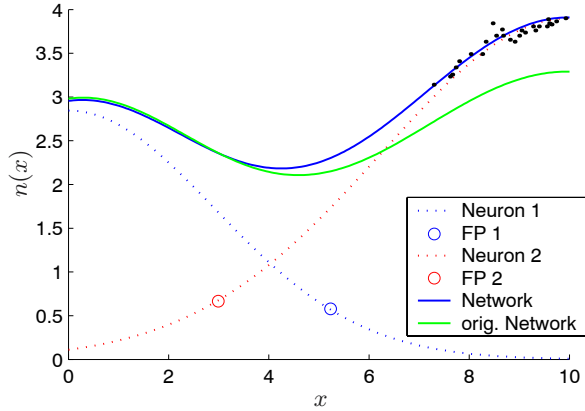


Figure 3: One-dimensional example

network trained with actual data before and after adaptation to a new data record. The model was built on PCA-channels one and two. Black dots represent the desired new network output p_k . It can be seen that the adaptation produces satisfactory results.

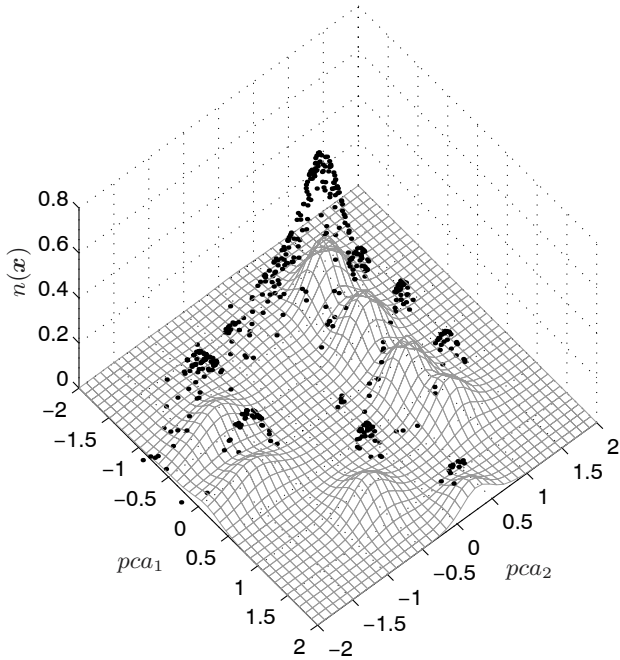


Figure 4: Two-dimensional example, prior to adaption

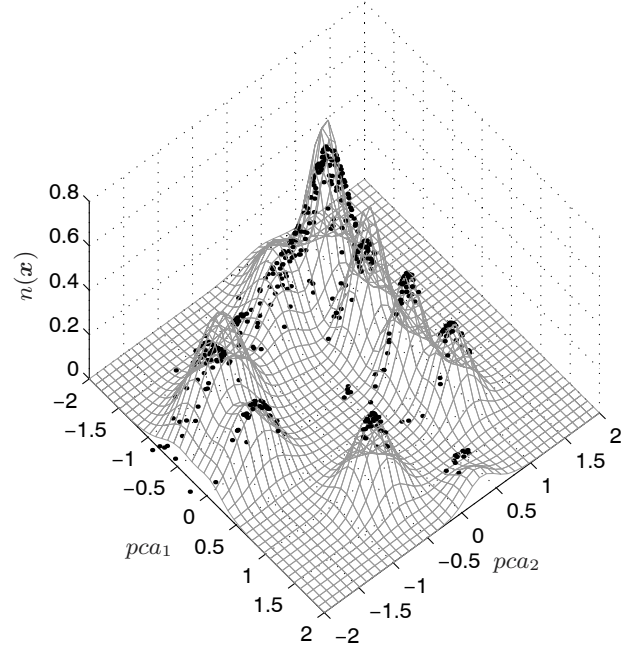


Figure 5: Two-dimensional example, after adaption

6. CONCLUSION

In this paper an approach for fault detection in automotive measuring systems was presented that automatically adapts itself if new data are available. The main problem was to develop an algorithm for adapting neurons to a new target function. For this purpose we introduced an approach for parameter adaption using gradient descent, neuron weighting and fixed points. Application to both a simulated example and to real data show that the proposed method yields good results. The next development stage features an extended gradient method that moves the centers of existing neurons in addition to the parameter changes already demonstrated in this paper. Moreover it is planned to replace the PCA in the preprocessing step by another suitable method.

7. REFERENCES

- [1] S. Jakubek and T. Strasser. Fault-diagnosis using neural networks with ellipsoidal basis functions. In *Proceedings of the 20th American Control Conference (ACC 2002)*, Anchorage, AK, pages 3846–3851, May 2002.
- [2] Kyle S. Cranmer. Kernel estimation for parametrization of discriminant variable distributions. <http://www.wisconsin.cern.ch/~cranmer/keys.html>, 1999.
- [3] J. Chen and R.J. Patton. *Robust Model-Based Fault Diagnosis for Dynamic Systems*. Kluwer Academic Publishers, Boston/Dordrecht/London, 1999.
- [4] S. Schweid and T.K. Sarkar. Projection minimization techniques for orthogonal QMF filters with vanishing moments. *IEEE Transactions on Circuits and Systems II*, 42:694–701, 1995.
- [5] T.J. Hebert and Keming Lu. Expectation-maximization algorithms, null spaces, and MAP image restoration. *IEEE Transactions on Image Processing*, 4:1084–1095, 1995.