

TH2

# Variable selection and neural networks for high-dimensional data analysis

## Application in infrared spectroscopy and chemometrics

NABIL BENOUDJIT

NOVEMBRE 2003

Thèse présentée en vue de l'obtention du grade  
de docteur en sciences appliquées

Faculté des sciences  
appliquées

Université catholique de Louvain



Université catholique de Louvain  
Faculté des Sciences Appliquées  
Département d'électricité  
Laboratoire de Microélectronique

Variable selection and neural networks for  
high-dimensional data analysis:  
application in infrared spectroscopy and  
chemometrics

Jury:

Prof. A. De Herde (président)	Université catholique de Louvain, Belgium
Prof. M. Verleysen (promoteur)	Université catholique de Louvain, Belgium
Prof. V. Wertz	Université catholique de Louvain, Belgium
Prof. M. Meurens	Université catholique de Louvain, Belgium
Dr. Ir. F. Rossi	Université Paris Dauphine, France
Prof. M. Saerens	Université catholique de Louvain, Belgium

Thèse présentée par  
**Nabil Benoudjit**  
en vue de l'obtention du grade de  
**Docteur en Sciences Appliquées**

Novembre 2003

# Contents

<b>1</b>	<b>Introduction</b>	<b>17</b>
1.1	Context and motivations . . . . .	17
1.2	Objectives of the thesis and our way to solutions . . . . .	20
1.3	Organization of the thesis . . . . .	21
<b>I</b>	<b>Artificial neural networks</b>	<b>25</b>
<b>2</b>	<b>Radial Basis Function Network (RBFN)</b>	<b>27</b>
2.1	Introduction . . . . .	27
2.2	Radial basis function network . . . . .	29
2.3	RBFN learning strategies . . . . .	30
2.3.1	Location of centroids . . . . .	31
2.3.2	Width factors . . . . .	32
2.3.3	Width Scaling Factor optimization (our approach) . . . . .	34
2.3.4	Optimal weights . . . . .	35
2.4	Comparison of RBF networks and multilayer perceptrons . . . . .	36
2.5	Normalized RBFN . . . . .	37
2.6	Summary of RBFN training procedure . . . . .	37
2.7	Conclusion . . . . .	38
<b>3</b>	<b>Model Selection</b>	<b>39</b>
3.1	Introduction . . . . .	39
3.2	Separate training, validation and test sets . . . . .	39
3.3	Model selection . . . . .	40
3.4	Cross-validation . . . . .	42
3.4.1	The hold-out method . . . . .	42

3.4.2	Cross-validation (Random subsampling) . . . . .	42
3.4.3	K-fold cross-validation . . . . .	43
3.4.4	Leave-one-out cross-validation . . . . .	44
3.5	Bootstrap . . . . .	46
3.6	Discussion . . . . .	47
<b>4</b>	<b>RBFN learning: width optimization</b>	<b>49</b>
4.1	Introduction . . . . .	49
4.2	Theoretical example . . . . .	49
4.2.1	Theoretical value of the optimal width of the Gaussian kernels . . . . .	51
4.2.2	Experimental value of the optimal width of the Gaussian kernels . . . . .	54
4.2.3	Results . . . . .	56
4.3	Analytical example . . . . .	61
4.3.1	Approximation results . . . . .	61
4.3.2	Comparison . . . . .	65
4.4	Conclusion . . . . .	68
<b>II</b>	<b>Variable selection</b>	<b>69</b>
<b>5</b>	<b>Spectrophotometric variable selection</b>	<b>71</b>
5.1	Introduction . . . . .	71
5.2	Beer-Lambert's law . . . . .	75
5.3	Spectrophotometric variable selection by linear models: state of the art . . . . .	76
5.3.1	Multiple linear regression (MLR) . . . . .	77
5.3.1.1	The coefficient of multiple determination	79
5.3.1.2	Collinearity . . . . .	79
5.3.2	Stepwise multiple linear regression (SMLR) . . . . .	80
5.3.2.1	Example . . . . .	82
5.3.3	Principal component regression (PCR) . . . . .	85
5.3.4	Partial least square regression (PLSR) . . . . .	87
5.4	Spectrophotometric variable selection by non-linear models	89
5.4.1	Forward-backward selection by RBF networks . . . . .	90
5.4.2	Error criterion . . . . .	92
5.5	Real-life examples . . . . .	93

5.5.1	Wine dataset . . . . .	93
5.5.2	Orange juice dataset . . . . .	99
5.5.3	Milk powder dataset . . . . .	106
5.5.4	Apple dataset . . . . .	112
5.6	Improvements . . . . .	118
5.6.1	Cross-validation . . . . .	118
5.6.2	Graphical detection of outliers . . . . .	120
5.6.3	Data distribution between training and validation sets . . . . .	125
5.6.4	Spectrophotometric variable selection by mutual information . . . . .	132
5.6.4.1	Mutual information . . . . .	132
5.6.4.2	Definition . . . . .	132
5.6.4.3	Computation of the mutual information . . . . .	134
5.6.4.4	Variable selection and validation by non-linear models . . . . .	134
5.6.4.5	Results . . . . .	135
5.7	Conclusion . . . . .	140
<b>6</b>	<b>Conclusion</b> . . . . .	<b>143</b>
6.1	Contributions and concluding remarks . . . . .	143
6.2	Perspectives and future work . . . . .	146
<b>A</b>		<b>147</b>

Université catholique de Louvain  
Faculté des sciences appliquées  
Département d'électricité  
Laboratoire de micro-électronique (Machine Learning Group)

Place du Levant 3, 1348 Louvain-la-Neuve, Belgique  
Tél. 32 (0) 10 47 25 40 – fax 32 (0) 10 47 25 98  
E-mail [benoudjit@dice.ucl.ac.be](mailto:benoudjit@dice.ucl.ac.be)  
<http://www.dice.ucl.ac.be/mlg/>



This thesis focuses particularly on the application of chemometrics in the field of analytical chemistry. Chemometrics (or multivariate analysis) consists in finding a relationship between two groups of variables, often called dependent and independent variables. In infrared spectroscopy for instance, chemometrics consists in the prediction of a quantitative variable (the obtention of which is delicate, requiring a chemical analysis and a qualified operator), such as the concentration of a component present in the studied product from spectral data measured on various wavelengths or wavenumbers (several hundreds, even several thousands). In this research we propose a methodology in the field of chemometrics to handle the spectrophotometric data which are often represented in high dimension. To handle these data, we first propose a new incremental method (step-by-step) for the selection of spectral data using linear and non-linear regression based on the combination of three principles: linear or non-linear regression, incremental procedure for the variable selection, and use of a validation set. This procedure allows on one hand to benefit from the advantages of non-linear methods to predict chemical data (there is often a non-linear relationship between dependent and independent variables), and on the other hand to avoid the overfitting phenomenon, one of the most crucial problems encountered with non-linear models. Secondly, we propose to improve the previous method by a judicious choice of the first selected variable, which has a very important influence on the final performances of the prediction. The idea is to use a measure of the mutual information between the independent and dependent variables to select the first one; then the previous incremental method (step-by-step) is used to select the next variables. The variable selected by mutual information can have a good interpretation from the spectrochemical point of view, and does not depend on the data distribution in the training and validation sets. On the contrary, the traditional chemometric linear methods such as PCR or PLSR produce new variables which do not have an obvious interpretation from the spectrochemical point of view. Four real-life datasets (wine, orange juice, milk powder and apples) are presented in order to show the efficiency and advantages of both proposed procedures compared to the traditional chemometric linear methods often used, such as MLR, PCR and PLSR.

