

## Preface

In psychology, *regression* means degeneration, or return to a prior, lower level of development. And, indeed, speaking of statistical methods sounds somewhat outdated today — all those neural networks and fuzzy systems being now available.

However, it has been recognized that the new soft computing methods are no panacea. These methods cannot find models any better than the more conventional ones if there is not enough information available in the data to start with. On the other hand, many celebrated soft computing applications could have been solved with age-old methods — assuming that somebody were familiar with them. After all, neural networks can only operate on the statistical properties visible in the data. Why not concentrate on these statistical properties directly?

The starting point here is that *when analysing complex systems, simple methods are needed*. When the system under study is very complicated, one needs data analysis methods that are reliable and fast and that give possibility for closer analysis.

Unfortunately, the statistical literature seems to be mathematically rather unpenetrable for normal engineers (for those that are not too ashamed to admit that). Researchers seem to represent the results in such sophisticated forms that the very simple and elegant ideas remain hidden — the problem is that the powerful methods may not be applied in practice. What is more, different methods have been developed in different research communities: It is difficult to see the connections between the methods when the approaches and notations differ. New methods are constantly being developed; there exists no one-volume book that would cover the whole field in sufficient detail.

This text tries to show how simple the basic ideas are and how closely related different methods turn out to be. The approach is rather pragmatic, many proofs being omitted for readability. All the methods are presented in a homogeneous framework. It is crucial that a student recognizes that all formulas and algorithms are based on simple underlying principles; everything can be questioned and nothing needs to be believed as some kind of unpenetrable “secret wisdom”. One objective is to show that there is still room for new ideas and innovations. As the reader can see, there exist plenty of ideas waiting to be explored and exploited — indeed, one is very near to frontier science, maybe able to cross the boundary towards new discoveries (examples of such exploratory experiments are indicated by stars “\*”).

The theoretical methods are supported by Matlab routines. The implemented “vanilla” algorithms are by no means optimized; their main purpose is also to show that the methods are by no means unpenetrable.

In addition to the printed version, this report is available in public domain in PostScript format. The Matlab Toolbox and the textual material can be accessed through the HTML page at the Internet address

[http://saato014.hut.fi/hyotyniemi/publications/01\\_reportXXX.htm](http://saato014.hut.fi/hyotyniemi/publications/01_reportXXX.htm).

The earlier version of this text, with the name “Multivariate Regression — Techniques and Tools” was published in 2001. Despite its various shortcomings, it received a positive acceptance. The printed version was “sold out” a long time ago. This interest was the motivation to try and fix some of the holes that were left open.

In the previous version, applications were not discussed: Now there is the last chapter that tries to present examples of not so evident but potential ways of applying the new tools. Second, the fields of physical first-principles modeling and data-oriented modeling are not so distinct that no connections could be found; the new appendices concentrates on this kind of discussions, fitting the structural considerations into the domain of numeric manipulations in a more or less seamless way.

National Technology Agency of Finland (TEKES) has provided funding during the research under several project frames, and this support is gratefully acknowledged.

A handwritten signature in black ink, appearing to read "Heikki Hyötyniemi". The signature is fluid and cursive, with a long horizontal stroke at the end.

Heikki Hyötyniemi

## List of symbols

The same variable names are used consistently in the theoretical discussion and in the accompanying **Regression Toolbox** for **Matlab**, if possible.

- $A, B, C, D$ : Matrices determining a state-space system
- $c$ : Arbitrary constant, scalar or vector
- $i, j$ : Matrix and vector indices
- $e, E$ : Measurement error vector and matrix, dimensions  $m \times 1$  and  $k \times m$ , respectively
- $\epsilon$ : State error, dimension  $n \times 1$
- $f, F$ : Vector and matrix defining a linear mapping
- $g(\cdot)$ : Any function (scalar or vector-valued)
- $\gamma, \Gamma$ : Constraint vector and matrix, respectively
- $J(\cdot)$ : Cost criterion
- $I, I_n$ : Identity matrix (of dimension  $n \times n$ )
- $k$ : Time index, sample number (given in parentheses)
- $K$ : Kalman gain
- $m$ : Dimension of the output space
- $M$ : Arbitrary matrix (or vector)
- $n$ : Dimension of the input space / Non-compressed feature space
- $N$ : Dimension of the latent space / Number of levels or substructures
- $P$ : Error covariance matrix, dimension  $d \times d$
- $p(\cdot)$ : Probability density
- $R$ : Covariance matrix (or, more generally, *association matrix*)
- $u, U$ : Input vector and matrix, dimensions  $\nu \times 1$  and  $k \times \nu$ , respectively
- $v, V$ : Combined data vector and matrix ( $x$  and  $y$  together), dimensions  $m + n \times 1$  and  $k \times m + n$ , respectively
- $\nu$ : Unprocessed measurement vector / Dimension of the vector  $u$  in dynamic systems
- $w, W$ : Weight vector and matrix, respectively
- $x, X$ : Data vector and matrix, dimensions  $n \times 1$  and  $k \times n$ , respectively. In the case of a dynamic system, state vector and matrix, dimensions  $d \times 1$  and  $k \times d$ , respectively.

- $y, Y$ : Output data vector and matrix, dimensions  $m \times 1$  and  $k \times m$ , respectively
- $z, Z$ : Latent data vector and matrix, dimensions  $N \times 1$  and  $k \times N$ , respectively
- $\xi, \zeta$ : Arbitrary scalars
- $\lambda, \Lambda$ : Vector of eigenvalues and eigenvalue matrix, respectively
- $\mu, \eta$ : Lagrange multipliers
- $\theta, \phi$ : Reduced base, dimensions  $n \times N$  and  $m \times N$ , respectively
- $\Theta, \Phi$ : Matrices of data basis vectors, dimensions  $n \times n$  and  $m \times m$ , respectively

## Notations

- $\mathbf{M}$ : Unprocessed data
- $\bar{M}$ : Mean value of  $M$  (columnwise mean matrix if  $M$  is matrix)
- $\hat{M}$ : Estimate of  $M$
- $\tilde{M}$ : Error in  $M$  / Erroneous  $M$
- $M'$ :  $M$  modified (somehow)
- $M^i$ : Power of  $M$  / Level  $i$  data structure
- $M_i$ : Column  $i$  for matrix  $M$  / Element  $i$  for vector-form  $M$
- $M^T$ : Transpose of  $M$
- $M^\dagger$ : Pseudoinverse of  $M$  (for definition, see page 21)
- $E\{M\}$ : Expectation value of  $M$
- $M_{\text{test}}, M_{\text{est}}$ : Independent testing data or run-time data
- $( M_1 \mid M_2 )$ : Partitioning of a matrix
- $M_{\xi \times \zeta}$ : Matrix dimensions ( $\xi$  rows,  $\zeta$  columns)

## Abbreviations

- CCA/CCR: Canonical Correlation Analysis/Regression (page 102)
- CR: Continuum Regression (page 98)
- CA: Cluster Analysis (page 205)
- CLR: Constrained Linear Regression (page ??)

- DA or FDA: (Fisher) Discriminant Analysis (page 208)
- EIV: Error In Variables model (page 68)
- FOBI: Fourth-Order Blind Identification (page 115)
- GHA: Generalized Hebbian Algorithm (page 138)
- ICA/ICR: Independent Component Analysis/Regression (page 109)
- MLR: Multi-Linear Regression (page 63)
- NNR: Neural Networks based Regression (page 129)
- OLS: Orthogonal Least Squares (page 72)
- PCA/PCR: Principal Component Analysis/Regression (page 77)
- PLS: Partial Least Squares regression (page 95)
- RR: Ridge Regression (page 73)
- SOM: Self-Organizing Map (page 126)
- SPC: Statistical Process Control (page 90)
- SSI: SubSpace Identification (page 143)
- TLS: Total Least Squares (page 185)



# Contents

<b>I</b>	<b>Theoretical Toolbox</b>	<b>9</b>
<b>1</b>	<b>Introduction to Multivariate Modeling</b>	<b>11</b>
1.1	About systems and models . . . . .	11
1.2	About mathematical tools . . . . .	14
1.2.1	Challenge of high dimensions . . . . .	14
1.2.2	About matrices . . . . .	15
1.2.3	Optimization . . . . .	19
1.2.4	Lagrange multipliers . . . . .	20
<b>2</b>	<b>About Distributions</b>	<b>23</b>
2.1	Data mining . . . . .	23
2.2	Normal distribution . . . . .	24
2.2.1	About distribution parameters . . . . .	26
2.2.2	Association matrices . . . . .	27
2.2.3	$\chi^2$ distribution . . . . .	29
2.3	Motivation of modeling approaches . . . . .	30
2.3.1	Why linear models? . . . . .	30
2.3.2	Why sum-of-error-squared criteria? . . . . .	32
2.4	Tackling with real-world data . . . . .	33
2.4.1	Gaussian mixture models . . . . .	33
2.4.2	Example: Types of “Natural Data” . . . . .	34
2.4.3	Outliers . . . . .	35
2.5	Excursion: Networks and <i>power law</i> . . . . .	36
<b>3</b>	<b>Understanding Data</b>	<b>39</b>
3.1	From intuition to information . . . . .	39
3.1.1	Some philosophy . . . . .	40
3.1.2	Implementing structure on the data . . . . .	41

3.1.3	Experiment design . . . . .	43
3.2	Selection of variables . . . . .	44
3.2.1	Feature extraction . . . . .	44
3.2.2	Special challenge: Dynamic systems . . . . .	45
3.3	Data preprocessing . . . . .	47
3.3.1	Reaching “well-behavedness” of data . . . . .	47
3.3.2	“Operating point” . . . . .	48
3.3.3	Data scaling . . . . .	50
3.4	Model construction and beyond . . . . .	51
3.4.1	Analysis and synthesis . . . . .	51
3.4.2	Validating the model . . . . .	52
3.4.3	Cross-validation . . . . .	53
3.5	Summary: Modeling procedures . . . . .	53
3.6	Case studies . . . . .	55
3.6.1	Analysis of the paper machine dry line . . . . .	55
3.6.2	Modeling of flotation froth . . . . .	58
<b>4</b>	<b>“Quick and Dirty”</b>	<b>63</b>
4.1	Linear regression model . . . . .	63
4.1.1	Least-squares solution . . . . .	63
4.1.2	Piece of analysis . . . . .	65
4.1.3	Multivariate case . . . . .	67
4.2	“Colored noise” . . . . .	68
4.2.1	Error in variables . . . . .	68
4.2.2	Instrumental variables . . . . .	69
4.3	Collinearity . . . . .	70
4.3.1	Example: When variables are redundant . . . . .	70
4.3.2	Patch fixes . . . . .	72
<b>5</b>	<b>Tackling with Redundancy</b>	<b>77</b>
5.1	Some linear algebra . . . . .	77
5.1.1	On spaces and bases . . . . .	77
5.1.2	About linear mappings . . . . .	78
5.1.3	Data model revisited . . . . .	79
5.2	Principal components . . . . .	81
5.2.1	Eigenproblem properties . . . . .	83
5.2.2	Analysis of the PCA model . . . . .	84



5.2.3	Another view of “information” . . . . .	85
5.2.4	Selection of basis vectors . . . . .	86
5.3	Practical aspects . . . . .	88
5.3.1	Regression based on PCA . . . . .	88
5.3.2	Other applications . . . . .	89
5.3.3	Analysis tools . . . . .	89
5.3.4	Calculating eigenvectors in practice . . . . .	91
5.4	New problems . . . . .	92
5.4.1	Experiment: “Associative regression”* . . . . .	92
<b>6</b>	<b>Bridging Input and Output</b>	<b>95</b>
6.1	Partial least squares . . . . .	95
6.1.1	Maximizing correlation . . . . .	95
6.2	Continuum regression . . . . .	98
6.2.1	On the correlation structure . . . . .	98
6.2.2	Filling the gaps . . . . .	99
6.2.3	Further explorations* . . . . .	100
6.3	Canonical correlations . . . . .	102
6.3.1	Problem formulation . . . . .	102
6.3.2	Analysis of CCA . . . . .	104
6.3.3	Regression based on PLS and CCA . . . . .	105
6.3.4	Further ideas* . . . . .	105
<b>7</b>	<b>Towards the Structure</b>	<b>109</b>
7.1	Factor analysis . . . . .	110
7.2	Independent components . . . . .	111
7.2.1	Why independence? . . . . .	111
7.2.2	Measures for independence . . . . .	111
7.2.3	ICA vs. PCA . . . . .	112
7.3	Eigenproblem-oriented ICA algorithms . . . . .	112
7.3.1	Data whitening . . . . .	114
7.3.2	Deformation of the distribution . . . . .	115
7.3.3	Further explorations* . . . . .	117
7.4	Beyond independence . . . . .	120
7.4.1	Sparse coding . . . . .	121
<b>8</b>	<b>Regression vs. Progression</b>	<b>125</b>
8.1	Neural clustering . . . . .	125

8.1.1	Self-organizing maps . . . . .	126
8.1.2	“Expectation Maximizing SOM” . . . . .	127
8.1.3	Radial basis function regression . . . . .	128
8.2	Feedforward networks . . . . .	129
8.2.1	Perceptron networks . . . . .	130
8.2.2	Back-propagation of errors . . . . .	131
8.2.3	Relations to subspace methods . . . . .	134
8.3	“Anthropomorphic models” . . . . .	136
8.3.1	Hebbian algorithms . . . . .	136
8.3.2	Generalized Hebbian algorithms . . . . .	138
8.3.3	Further extensions . . . . .	138
8.4	Cybernetic neurons* . . . . .	139
<b>9</b>	<b>Application to Dynamic Models</b>	<b>143</b>
9.1	Representing dynamics . . . . .	143
9.1.1	Capturing history . . . . .	143
9.1.2	State-space models . . . . .	144
9.2	Subspace identification . . . . .	145
9.2.1	Stochastic models . . . . .	145
9.2.2	Stochastic-deterministic models . . . . .	149
9.3	Practical aspects . . . . .	150
9.3.1	Comparisons . . . . .	151
9.3.2	Emulating the process . . . . .	151
9.3.3	Preprocessing and postprocessing . . . . .	152
9.4	Case study: Towards “smart devices” . . . . .	153
9.4.1	Data based data reconciliation . . . . .	154
9.4.2	Connections to AI* . . . . .	157
<b>10</b>	<b>Relations to Systems Engineering</b>	<b>159</b>
10.1	MIMO vs. SISO systems . . . . .	159
10.2	Dimension reduction . . . . .	160
10.2.1	About state-space systems . . . . .	160
10.2.2	Preliminary experiments . . . . .	162
10.2.3	Balanced realizations . . . . .	163
10.2.4	Eliminating states . . . . .	166
10.3	State estimation . . . . .	166
10.3.1	Kalman filter . . . . .	167

10.3.2	Optimality vs. reality . . . . .	168
10.3.3	Reducing the number of measurements . . . . .	169
10.4	SISO identification . . . . .	170
10.4.1	Black-box model . . . . .	170
10.4.2	Recursive least-squares algorithm . . . . .	171
10.4.3	Structure of dynamic data . . . . .	173
10.4.4	Further analysis: System identifiability* . . . . .	174
<b>11</b>	<b>Towards “Emergent Models”</b>	<b>179</b>
11.1	Capturing semantics in data . . . . .	179
11.1.1	What is “semantics”? . . . . .	180
11.1.2	Neocybernetic starting points . . . . .	181
11.1.3	Modeling chemical systems . . . . .	181
11.2	From constraints to degrees of freedom . . . . .	184
11.2.1	Constraint-based models . . . . .	184
11.2.2	Total Least Squares . . . . .	185
11.2.3	Emergent models . . . . .	187
11.2.4	Examples . . . . .	189
11.3	Case studies . . . . .	191
11.3.1	Characterizing the state in practical processes . . . . .	191
11.3.2	Case 1: Modeling an industrial nickel plating process <sup>1</sup> . . . . .	192
11.3.3	Case 2: Modeling genetic networks and metabolic systems <sup>2</sup> . . . . .	194
11.4	Towards “artificial cells” . . . . .	197
<b>A</b>	<b>Structure <i>from</i> Data</b>	<b>205</b>
A.1	Cluster analysis . . . . .	205
A.1.1	K-means algorithm . . . . .	206
A.1.2	EM algorithm . . . . .	206
A.2	Classification . . . . .	208
A.2.1	Fisher discriminant analysis . . . . .	208
A.2.2	Support Vector Machines (SVM) . . . . .	210
<b>B</b>	<b>Structure <i>into</i> Data</b>	<b>213</b>
B.1	Data reconciliation . . . . .	213
B.1.1	Explicit constraints on parameters . . . . .	216

---

<sup>1</sup>The simulations were carried out by Mr. Hans-Christian Pfisterer

<sup>2</sup>The simulations were carried out by Mr. Olli Haavisto, M.Sc.

B.2 Observing functional hierarchy . . . . .	218
B.3 Dimensional analysis . . . . .	220
B.3.1 Fixing missing data . . . . .	223

<b>II Practical Toolbox</b>	<b>227</b>
-----------------------------	------------

## A Good Book Must Have Obscure Quotes

A theory has only the alternative of being wrong. A model has a third possibility — it might be right but irrelevant.

— M. Eigen

Statistics in the hands of an engineer are like a lamppost to a drunk — they're used more for support than illumination.

— B. Sangster

Like other occult techniques of divination, the statistical method has a private jargon deliberately contrived to obscure its methods from non-practitioners.

— G. O. Ashley

Make the model as simple as possible — but not simpler!

— A. Einstein

There are three kinds of lies: lies, damned lies, and statistics.

— B. Disraeli

In earlier times, they had no statistics, and so they had to fall back on lies.

— Stephen Leacock

Torture the data long enough and they will confess to anything.

— unknown

If at first it doesn't fit, fit, fit again.

— J. McPhee

Data! Data! Data! I can't make bricks without clay!

— S. Holmes

...

Why's the bell-shaped curve called normal?

Is it normal to be so formal?

There's nothing mean about the mean.

Its just average, as is clearly seen.

And what's so standard about that deviation?

It's a really malicious creation.

Confusing students is its only function.

It frustrates and mystifies, in conjunction.

...

— *"On statistical terminology"* by C. Lation