

## Appendices

# A

---

## R Packages Used in this Book

The table below contains an overview of R packages mentioned (but not necessarily treated in the same detail) in the book.

Table A.1. R packages used in this book

---

ada	elasticnet	mclust	rda
ALS	fastICA	meboot	relaxo
AMORE	fingerprint	msProstate	robustbase
boost	FRB	neuralnet	rpart
boot	glmnet	nnet	rrcov
BootPR	gpls	paltran	sfsmisc
bootstrap	gtools	VPdtw	som
caMassClass	ipred	pls	spls
ChemometricsWithR	kohonen	plsgenomics	stats
class	lars	plspm	subselect
cluster	lasso2	ppls	TIMP
DAIM	leaps	PROcess	tree
dtw	lpc	ptw	wccsom
e1071	lspls	randomForest	xcms
EffectiveDose	MASS		

---

---

## References

1. R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2010. ISBN 3-900051-07-0.
2. W.N. Venables, D.M. Smith, and the R development Core Team. An introduction to R, December 2009. Version 2.10.1.
3. T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer, New York, 2001.
4. K. Varmuza and P. Filzmoser. *Introduction to Multivariate Statistical Analysis in Chemometrics*. Taylor & Francis - CRC Press, Boca Raton, FL, USA, 2009.
5. J. Kalivas. Two data sets of near infrared spectra. *Chemom. Intell. Lab. Syst.*, 37:255–259, 1997.
6. M. Forina, C. Armanino, M. Castino, and M. Ubigli. Multivariate data analysis as a discriminating method of the origin of wines. *Vitis*, 25:189–201, 1986.
7. B.-L. Adam, Y. Qu, J.W. Davis, M.D. Ward, M.A. Clements, L.H. Cazares, O.J. Semmes, P.F. Schellhammer, Y. Yasui, Z. Feng, and G.L. Wright. Serum protein fingerprinting coupled with a pattern-matching algorithm distinguishes prostate cancer from benign prostate hyperplasia and healthy men. *Cancer Res.*, 62(13):3609–3614, July 2002.
8. Y. Qu, B.-L. Adam, Y. Yasui, M.D. Ward, L.H. Cazares, P.F. Schellhammer, Z. Feng, O.J. Semmes, and G.L. Wright. Boosted decision tree analysis of surface-enhanced laser desorption/ionization mass spectral serum profiles discriminates prostate cancer from noncancer patients. *Clin Chem*, 48(10):1835–43, October 2002.
9. T.G. Bloemberg, J. Gerretzen, H.J.P. Wouters, J. Gloerich, M. van Dael, H.J.C.T. Wessels, L.P. van den Heuvel, P.H.C. Eilers, L.M.C. Buydens, and R. Wehrens. Improved parametric time warping for proteomics. *Chemom. Intell. Lab. Systems*, 2010.
10. A. Savitsky and M.J.E. Golay. Smoothing and differentiation of data by simplified least squares procedures. *Anal. Chem.*, 36:1627–1639, 1964.
11. W.S. Cleveland. Robust locally weighted regression and smoothing scatterplots. *J. Am. Stat. Assoc.*, 74:829–836, 1979.
12. G.P. Nason. *Wavelet methods in statistics with R*. Springer, New York, 2008.
13. P. Geladi, D. MacDougall, and H. Martens. Linearization and scatter-correction for NIR reflectance spectra of meat. *Appl. Spectr.*, 39:491–500, 1985.

14. T. Næs, T. Isaksson, and B.R. Kowalski. Locally weighted regression and scatter correction for near-infrared reflectance data. *Anal. Chem.*, 62:664–673, 1990.
15. H. Sakoe and S. Chiba. Dynamic programming algorithm optimization for spoken word recognition. *IEEE Trans. Acoust., Speech, Signal Process.*, 26:43–49, 1978.
16. L.R. Rabiner, A.E. Rosenberg, and S.E. Levinson. Considerations in dynamic time warping algorithms for discrete word recognition. *IEEE Trans. Acoust., Speech, Signal Process.*, 26:575–582, 1978.
17. C.P. Wang and T.L. Isenhour. Time-warping algorithm applied to chromatographic peak matching gas chromatography / Fourier Transform infrared / Mass Spectrometry. *Anal. Chem.*, 59:649–654, 1987.
18. N.P.V. Nielsen, J.M. Carstensen, and J. Smedsgaard. Aligning of single and multiple wavelength chromatographic profiles for chemometric data analysis using correlation optimized warping. *J. Chrom. A*, 805:17–35, 1998.
19. P.H.C. Eilers. Parametric time warping. *Anal. Chem.*, 76:404–411, 2004.
20. R. de Gelder, R. Wehrens, and J.A. Hageman. A generalized expression for the similarity spectra: application to powder diffraction pattern classification. *J. Comput. Chem.*, 22(3):273–289, 2001.
21. D. Clifford, G. Stone, I. Montoliu, S. Rezzi, F.-P. Martin, P. Guy, S. Bruce, and S. Kochhar. Alignment using variable penalty dynamic time warping. *Anal. Chem.*, 81:1000–1007, 2009.
22. W. Windig, J. Phalp, and A. Payna. A noise and background reduction method for component detection in liquid chromatography/mass spectrometry. *Anal. Chem.*, 68:3602–3606, 1996.
23. T. Giorgino. Computing and visualizing dynamic time warping alignments in R: the dtw package. *J. Stat. Softw.*, 31(7), 2009.
24. J.E. Jackson. *A User's Guide to Principal Components*. Wiley, Chichester, 1991.
25. I.T. Jolliffe. *Principal Component Analysis*. Springer, New York, 1986.
26. K. Mardia, J. Kent, and J. Bibby. *Multivariate Analysis*. Academic Press, 1979.
27. W. Härdle and L. Simar. *Applied Multivariate Statistical Analysis*. Springer, Berlin, 2nd edition, 2007.
28. K.R. Gabriel. The biplot graphic display of matrices with application to principal component analysis. *Biometrika*, 58:453–467, 1971.
29. J.C. Gower and D.J. Hand. *Biplots*. Number 54 in Monographs on Statistics and Applied Probability. Chapman and Hall, London, UK, 1996.
30. K. Baumann. Uniform-length molecular descriptors for quantitative structure-property relationships (QSPR) and quantitative structure-activity relationships (QSAR): classification studies and similarity searching. *Trends Anal. Chem.*, 18(1):36–46, 1999.
31. T.F. Cox and M.A.A. Cox. *Multidimensional Scaling*. Chapman and Hall, 2001.
32. I. Borg and P.J.F. Groenen. *Modern Multidimensional Scaling*. Springer, 2nd edition, 2005.
33. J.C. Gower. Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika*, 53:325–328, 1966.
34. B.D. Ripley. *Pattern recognition and neural networks*. Cambridge University Press, 1996.

35. J.H. Friedman and J.W. Tukey. A projection pursuit algorithm for exploratory data analysis. *IEEE Trans. Comput.*, C23:881–889, 1974.
36. P.J. Huber. Projection pursuit. *The Annals of Statistics*, 13:435–475, 1985.
37. J.H. Friedman. Exploratory projection pursuit. *Journal of the American Statistical Association*, 82:249–266, 1987.
38. A. Hyvärinen, J. Karhunen, and E. Oja. *Independent Component Analysis*. Wiley, Chichester, 2001.
39. T.M. Cover and J.A. Thomas. *Elements of Information Theory*. Wiley, Chichester, 1991.
40. A. Hyvärinen and E. Oja. Independent component analysis: algorithms and applications. *Neural Networks*, 13:411–430, 2000.
41. C. Spearman. “General intelligence”, objectively determined and measured. *Am. J. Psychol.*, 15:201–293, 1904.
42. T. Kohonen. *Self-Organizing Maps*. Number 30 in Springer Series in Information Sciences. Springer, Berlin, 3 edition, 2001.
43. R. Wehrens and L.M.C. Buydens. Self- and super-organising maps in R: the kohonen package. *Journal of Statistical Software*, 21(5), 9 2007.
44. A. Ultsch. Self-organizing neural networks for visualization and classification. In O. Opitz, B. Lausen, and R. Klar, editors, *Information and Classification – Concepts, Methods and Applications*, pages 307–313. Springer Verlag, 1993.
45. W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer, New York, fourth edition, 2002. ISBN 0-387-95457-0.
46. R. Wehrens and E. Willighagen. Mapping databases of x-ray powder patterns. *R News*, 6(3):24–28, August 2006.
47. M. Eisen, P.T. Spellman, P.O. Brown, and D. Botstein. Cluster analysis and display of genome-wide expression patterns. *Proc. Natl. Acad. Sci. USA*, 95:14863–14868, December 1998.
48. L. Kaufman and P.J. Rousseeuw. *Finding Groups in Data – An Introduction to Cluster Analysis*. John Wiley & Sons, New York, 1990.
49. G.J. McLachlan and D. Peel. *Finite Mixture Models*. John Wiley & Sons, New York, 2000.
50. C. Fraley and A.E. Raftery. Model-based clustering, discriminant analysis, and density estimation. *J. Am. Stat. Assoc.*, 97:611–631, 2002.
51. A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *J. R. Statist. Soc. B*, 39(1):1–38, 1977.
52. G.J. McLachlan and T. Krishnan. *The EM Algorithm and Extensions*. John Wiley & Sons, 1997.
53. H. Akaike. A new look at the statistical model identification. *IEEE Trans. Automatic Control*, 19:716–723, 1974.
54. G. Schwarz. Estimating the dimension of a model. *Ann. Statist.*, 6:461–464, 1978.
55. C. Fraley and A.E. Raftery. Enhanced software for model-based clustering, discriminant analysis, and density estimation: MCLUST. *J. Classif.*, 20:263–286, 2003.
56. C. Fraley. Algorithms for model-based gaussian hierarchical clustering. *SIAM J. Scient. Comput.*, 20:270–281, 1998.
57. J.D. Banfield and A.E. Raftery. Model-based Gaussian and non-Gaussian clustering. *Biometrics*, 49:803–821, 1993.
58. L. Hubert. Comparing partitions. *J. Classif.*, 2:193–218, 1985.

59. W.M. Rand. Objective criteria for the evaluation of clustering methods. *J. Am. Stat. Assoc.*, 66:846–850, 1971.
60. J. Vesanto and E. Alhoniemi. Clustering of the self-organising map. *IEEE Trans. Neural Netw.*, 11:586–600, 2000.
61. E.B. Fowlkes and C.L. Mallows. A method for comparing two hierarchical clusterings. *J. Am. Stat. Assoc.*, 78:553–584, 1983. Including discussion.
62. L.A. Goodman and W.H. Kruskal. Measures of association for cross classifications. *J. Am. Statist. Assoc.*, 49:732–764, 1954.
63. M. Meila. Comparing clusterings – an information-based distance. *J. Multivar. Anal.*, 98(5):873–895, 2007.
64. G.J. McLachlan. *Discriminant Analysis and Statistical Pattern Recognition*. Wiley-Interscience, 2004.
65. M. Stone. Cross-validated choice and assessment of statistical predictions. *J. R. Statist. Soc. B*, 36:111–147, 1974. Including discussion.
66. B. Efron and R.J. Tibshirani. *An Introduction to the Bootstrap*. Chapman and Hall, New York, 1993.
67. R.A. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7:179–188, 1936.
68. J. Friedman. Regularized discriminant analysis. *J. Am. Stat. Assoc.*, 84:165–175, 1989.
69. S. Dudoit, J. Fridlyand, and T.P. Speed. Comparison of discrimination methods for the classification of tumors using gene expression data. *J. Am. Stat. Assoc.*, 97:77–87, 2002.
70. D. Hand and K. Yu. Idiot’s Bayes – not so stupid after all? *Int. Statist. Rev.*, 69:385–398, 2001.
71. R. Tibshirani, T. Hastie, B. Narashimhan, and G. Chu. Class prediction by nearest shrunken centroids with applications to dna microarrays. *Statistical Science*, 18:104–117, 2003.
72. Y. Guo, T. Hastie, and R. Tibshirani. Regularized linear discriminant analysis and its application in microarrays. *Biostatistics*, 8(1):86–100, 2007.
73. L. Breiman, J.H. Friedman, R.A. Olshen, and C.J. Stone. *Classification and Regression Trees*. Wadsworth, 1984.
74. J.R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
75. J.R. Quinlan. The induction of decision trees. *Mach. Learn.*, 1(1):81–106, 1986.
76. T.M. Therneau and E.J. Atkinson. An introduction to recursive partitioning using the RPART routines. Technical Report 61, Mayo Foundation, September 1997.
77. V. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, 1995.
78. N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines and other kernel-based Learning Methods*. Cambridge University Press, 2000.
79. B. Schölkopf and A.J. Smola. *Learning with kernels*. MIT Press, Cambridge, MA, 2002.
80. F. Rosenblatt. *Principles of neurodynamics*. Spartan Books, Washington DC, 1962.
81. D.E. Rumelhard and J.L. McClelland, editors. *Parallel distributed processing: explorations in the microstructure of cognition. Volume 1: Foundations*. MIT Press, Cambridge MA, 1986.
82. Frauke Günther and Stefan Fritsch. neuralnet: Training of neural networks. *The R Journal*, 2(1):30–38, June 2010.

83. B.-H. Mevik and R. Wehrens. The pls package: principal component and partial least squares regression in R. *J. Stat. Soft.*, 18(2), 2007.
84. B.-H. Mevik and H.R. Cederkvist. Mean squared error of prediction (MSEP) estimates for principal component regression (PCR) and partial least squares regression (PLSR). *J. Chemom.*, 18:422–429, 2004.
85. A.S. Barros and D.N. Rutledge. Genetic algorithms applied to the selection of principal components. *Chemom. Intell. Lab. Syst.*, 40:65–81, 1998.
86. B.S. Dayal and J.F. MacGregor. Improved PLS algorithms. *J. Chemom.*, 11:73–85, 1997.
87. H. Martens and T. Næs. *Multivariate Calibration*. Wiley, Chichester, 1989.
88. S. Rännar, F. Lindgren, P. Geladi, and S. Wold. The PLS Kernel algorithm for data sets with many variables and fewer objects. Part 1: Theory and algorithm. *J. Chemom.*, 8:111, 1994.
89. S. de Jong. SIMPLS: an alternative approach to partial least squares regression. *Chemom. Intell. Lab. Syst.*, 18:251–263, 1993.
90. I.E. Frank and J.H. Friedman. A statistical view of some chemometrics regression tools. *Technometrics*, 35:109–135, 1993.
91. S. Wold, N. Kettaneh-Wold, and B. Skagerberg. Nonlinear PLS modeling. *Chemom. Intell. Lab. Syst.*, 7:53–65, 1989.
92. K. Hasegawa, T. Kimura, Y. Miyashita, and K. Funatsu. Nonlinear partial least squares modeling of phenyl alkylamines with the monoamine oxidase inhibitory activities. *J. Chem. Inf. Comput. Sci.*, 36:1025–1029, 1996.
93. K. Jorgensen, V. H. Segtnan, K. Thyholt, and T. Næs. A comparison of methods for analysing regression models with both spectral and designed variables. *J. Chemometr.*, 18:451–464, 2004.
94. B. Ding and R. Gentleman. Classification using penalized partial least squares. *J. Comput. Graph. Stat.*, 14:280–298, 2005.
95. B.D. Marx. Iteratively reweighted partial least squares estimation for generalized linear regression. *Technometrics*, 38:374–381, 1996.
96. C.J.F. ter Braak and S. Juggins. Weighted averaging partial least squares regression WAPLS: an improved method for reconstructing environmental variables from species assemblages. *Hydrobiologia*, 269:485–502, 1993.
97. M. Tenenhaus, V. Esposito Vinzi, Y.M. Chatelin, and C. Lauro. PLS path modelling. *Comput. Stat. Data Anal.*, 48:159–2005, 2005.
98. N. Krämer, A.-L. Boulesteix, and G. Tutz. Penalized partial least squares with applications to B-spline transformations and functional data. *Chemom. Intell. Lab. Syst.*, 94:60–69, 2008.
99. H. Chun and S. Keles. Sparse partial least squares for simultaneous dimension reduction and variable selection. *J. Royal Stat. Soc. – Series B*, 72:3–25, 2010.
100. A.E. Hoerl. Application of ridge analysis to regression problems. *Chemical Engineering Progress*, 58:54–59, 1962.
101. A.E. Hoerl and R.W. Kennard. Ridge regression: biased estimation for non-orthogonal problems. *Technometrics*, 8:27–51, 1970.
102. A.E. Hoerl, R.W. Kennard, and K.F. Baldwin. Ridge regression: some simulations. *Commun. Stat. – Simul. Comput.*, 4:105–123, 1975.
103. J.F. Lawless and P. Wang. A simulation study of ridge and other regression estimators. *Commun. Stat. – Theory and Methods*, 5:303–323, 1976.
104. M. Stone and R.J. Brooks. Continuum regression: cross-validated sequentially constructed prediction embracing ordinary least squares, partial least squares

- and principal components regression (with discussion). *J. R. Statist. Soc.*, 52:237–269, 1990.
105. S. de Jong and H.A.L. Kiers. Principal covariates regression: Part I. Theory. *Chemom. Intell. Lab. Syst.*, 14:155–164, 1992.
  106. R.W. Kennard and L. Stone. Computer aided design of experiments. *Technometrics*, 11:137–148, 1969.
  107. C.D. Brown and H.T. Davis. Receiver operating characteristic curves and related decision measures: a tutorial. *Chemom. Intell. Lab. Syst.*, 80:24–38, 2006.
  108. C.L. Mallows. Some comments on Cp. *Technometrics*, 15:661–675, 1973.
  109. P. Craven and G. Wahba. Smoothing noisy data with spline functions. *Numer. Math.*, 31:377–403, 1979.
  110. S. Smit, M.J. van Breemen, H.C.J. Hoefsloot, A.K. Smilde, J.M.F.G. Aerts, and C.G. de Koster. Assessing the statistical validity of proteomics based biomarkers. *Anal. Chim. Acta*, 592:210–217, 2007.
  111. A.C. Davison and D.V. Hinkley. *Bootstrap Methods and their Applications*. Cambridge University Press, Cambridge, 1997.
  112. B. Efron and R. Tibshirani. Improvements on cross-validation: the .632+ bootstrap method. *J. Am. Stat. Assoc.*, 92:548–560, 1997.
  113. B. Efron. Bootstrap methods: another look at the jackknife. *Ann. Stat.*, 7:1–26, 1979.
  114. L. Breiman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.
  115. L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
  116. Y. Freund and R.E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1):119–139, 1997.
  117. V. Svetnik, A. Liaw, C. Tong, J.C. Culberson, R.P. Sheridan, and B.P. Feuston. Random forest: a classification and regression tool for compound classification and qsar modeling. *J. Chem. Inf. Comput. Sci.*, 43(6):1947–58, 2003.
  118. G. Michailides, K. Johnson, and M. Culp. ada: an R package for stochastic boosting. *J. Stat. Softw.*, 17(2), 2006.
  119. J.H. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. *Ann. Stat.*, 28:337–374, 2000.
  120. R.E. Schapire, Y. Freund, P. Bartlett, and W.S. Lee. Boosting the margin: a new explanation for the effectiveness of voting methods. *Ann. Stat.*, 26:1651–1686, 1998.
  121. R. Tibshirani. Regression shrinkage and selection via the lasso. *J. Royal. Statist. Soc B*, 58:267–288, 1996.
  122. R. Wehrens and W.E. van der Linden. Bootstrapping principal-component regression models. *J. Chemom.*, 11(2):157–171, 1997.
  123. A.H. Land and A.G. Doig. An automatic method for solving discrete programming problems. *Econometrica*, 28:497–520, 1960.
  124. G.M. Furnival and G.M. Wilson. Regression by leaps and bounds. *Technometrics*, 16:499–511, 1974.
  125. B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32:407–499, 2004.
  126. H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *J. Royal. Stat. Soc. B*, 67:301–320, 2005.
  127. S. Kirkpatrick, C.D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220:671–680, 1983.



128. V. Cerny. A thermodynamical approach to the travelling salesman problem: an efficient simulation algorithm. *Journal of Optimization Theory and Applications*, 45:41–51, 1985.
129. N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, and E. Teller. Equations of state calculations by fast computing machines. *J. Chem. Phys.*, 21:1087–1092, 1953.
130. V. Granville, M. Krivanek, and J.-P. Rasson. Simulated annealing: a proof of convergence. *IEEE Trans. Patt. Anal. Machine Intell.*, 16:652–656, 1994.
131. A.P. Duarte Silva. Efficient variable screening for multivariate analysis. *J. Mult. Anal.*, 76:35–62, 2001.
132. D.E. Goldberg. *Genetic Algorithms in Search, Optimization and Machine Learning*. Kluwer Academic Publishers, Boston, MA., 1989.
133. R. Leardi. Genetic algorithms in chemometrics and chemistry: a review. *J. Chemom.*, 15:559–569, 2001.
134. J. Shao. Linear model selection by cross-validation. *J. Am. Statist. Assoc.*, 88:486–494, 2003.
135. K. Baumann, H. Albert, and M. von Korff. A systematic evaluation of the benefits and hazards of variable selection in latent variable regression. Part I. Search algorithm, theory and simulations. *J. Chemometr.*, 16:339–350, 2002.
136. K. Baumann, H. Albert, and M. von Korff. A systematic evaluation of the benefits and hazards of variable selection in latent variable regression. Part II. Practical applications. *J. Chemometr.*, 16:351–360, 2002.
137. M. Hubert. Robust calibration. In S.D. Brown, R. Tauler, and B. Walczak, editors, *Comprehensive Chemometrics – Chemical and Biochemical Data Analysis*, chapter 3.07, pages 315–343. Elsevier, 2009.
138. C. Croux and G. Haesbroeck. Principal components analysis based on robust estimators of the covariance or correlation matrix. *Biometrika*, 87:603–618, 2000.
139. P. Rousseeuw. Least median of squares regression. *J. Am. Stat. Assoc.*, 79:871–880, 1984.
140. P.J. Rousseeuw and K. van Driessen. A fast algorithm for the minimum covariance determinant estimator. *Technometrics*, 41:212–223, 1999.
141. M. Hubert, P.J. Rousseeuw, and K. Vanden Branden. ROBPCA: a new approach to robust principal component analysis. *Technometrics*, 47:64–79, 2005.
142. V. Todorov and P. Filzmoser. An object oriented framework for robust multivariate analysis. *J. Stat. Softw.*, 32(3):1–47, 2009.
143. M. Hubert and K. Vanden Branden. Robust methods for partial least squares regression. *J. Chemom.*, 17:537–549, 2003.
144. B. Liebmann, P. Filzmoser, and K. Varmuza. Robust and classical PLS regression compared. *J. Chemom.*, 24:111–120, 2009.
145. S. Wold, H. Antti, F. Lindgren, and J. Ohman. Orthogonal signal correction of near-infrared spectra. *Chemom. Intell. Lab. Syst.*, 44:175–185, 1998.
146. O. Svensson, T. Kourti, and J.F. MacGregor. A comparison of orthogonal signal correction algorithms and characteristics. *J. Chemom.*, 16:176–188, 2002.
147. J. Trygg and S. Wold. Orthogonal projections to latent structures (O-PLS). *J. Chemom.*, 16:119–128, 2002.
148. M. Barker and W. Rayens. Partial least squares for discrimination. *J. Chemom.*, 17:166–173, 2003.
149. D.V. Nguyen and D. Rocke. Tumor classification by partial least squares using microarray gene expression data. *Bioinformatics*, 18:39–50, 2002.

150. A.L. Boulesteix. PLS dimension reduction for classification with high-dimensional microarray data. *Stat. Appl. Genet. Mol. Biol.*, 3, 2004. Article 33.
151. O.E. de Noord. Multivariate calibration standardization. *Chemom. Intell. Lab. Syst.*, 25:85–97, 1994.
152. Y. Wang, D.J. Veltkamp, and B.R. Kowalski. Multivariate instrument standardization. *Anal. Chem.*, 63:2750–2756, 1994.
153. Z.Y. Wang, T. Dean, and B.R. Kowalski. Additive background correction in multivariate instrument standardization. *Anal. Chem.*, 67:2379–2385, 1995.
154. E. Bouveresse, D.L. Massart, and P. Dardenne. Modified algorithm for standardization of near-infrared spectrometric instruments. *Anal. Chem.*, 67:1381–1389, 1995.
155. Y. Wang, M.J. Lysaght, and B.R. Kowalski. Improvement of multivariate calibration through instrument standardization. *Anal. Chem.*, 64:764–771, 1992.
156. R. Tauler, S. Lacorte, and D. Barceló. Application of multivariate self-modeling curve resolution to the quantitation of trace levels of organophosphorus pesticides in natural waters from interlaboratory studies. *J. Chromatogr. A*, 730:177–183, 1996.
157. W.H. Lawton and E.A. Sylvestre. Self-modeling curve resolution. *Technometrics*, 13:617–633, 1971.
158. R. Tauler. Multivariate curve resolution applied to second-order data. *Chemom. Intell. Lab. Syst.*, 30:133–146, 1995.
159. A. de Juan, M. Maeder, M. Martinez, and R. Tauler. Combining hard- and soft-modelling techniques to solve kinetic problems. *Chemom. Intell. Lab. Syst.*, 54:49–67, 2000.
160. M. Maeder. Evolving factor analysis for the resolution of overlapping chromatographic peaks. *Anal. Chem.*, 59:527–530, 1987.
161. H.R. Keller and D.L. Massart. Peak purity control in liquid chromatography with photodiode-array detection by a fixed size moving window evolving factor analysis. *Anal. Chim. Acta*, 246:379–390, 1991.
162. F. Questa Sanchez, M.S. Khots, D.L. Massart, and J.O. de Beer. Algorithm for the assessment of peak purity in liquid chromatography with photodiode-array detection. *Anal. Chim. Acta*, 285:181–192, 1994.
163. W. Windig and J. Guilment. Interactive self-modeling mixture analysis. *Anal. Chem.*, 63:1425–1432, 1991.
164. F. Cuesta Sanchez, B. van de Bogaert, S.C. Rutan, and D.L. Massart. Multivariate peak purity approaches. *Chemom. Intell. Lab. Syst.*, 36:153–164, 1996.
165. P. Stilbs. Molecular self-diffusion coefficients in Fourier transform nuclear magnetic resonance spectrometric analysis of complex mixtures. *Anal. Chem.*, 53:2135–2137, 1981.
166. P. Stilbs. Fourier-transform pulsed-gradient spin-echo studies of molecular diffusion. *Progr. NMR Spectrosc.*, 19:1–45, 1987.
167. R. Huo, R. Wehrens, J. van Duynhoven, and L.M.C. Buydens. Assessment of techniques for DOSY NMR data processing. *Anal. Chim. Acta*, 490:231–251, 2003.
168. R. Huo, R. Wehrens, and L.M.C. Buydens. Improved DOSY NMR data processing by data enhancement and combination of multivariate curve resolution with non-linear least squares fitting. *J. Magn. Reson.*, 169:257–269, 2004.

169. K.M. Mullen, I.H.M. van Stokkum, and V.V. Mihaleva. Global analysis of multiple gas chromatography-mass spectrometry (GC/MS) data sets: a method for resolution of co-eluting components with comparison to MCR-ALS. *Chemom. Intell. Lab. Syst.*, 95:150–163, 2009.
170. G. Munoz and A. de Juan. pH- and time-dependent hemoglobin transitions: a case study for process modelling. *Anal. Chim. Acta*, 595:198–208, 2007.

---

# Index

- Akaike's information criterion (AIC),
  - 91, 180
- Artificial neural networks (ANNs),
  - 141–144
- Backpropagation networks, *see* Artificial neural networks (ANNs)
- Bagging, 196–197
- Baseline removal, 18–20
- Bayesian information criterion (BIC),
  - 91, 180
- Bias, 149, 177, 183, 184
- Binning, 11, 16
- Biomarkers, 38, 103
- Boosting, 202–204
- Bootstrap, 177, 186–195
  - .632 estimate, 187
  - BC $\alpha$  confidence intervals, 193
  - nonparametric, 186
  - parametric, 186
  - percentile confidence intervals, 191
  - studentized confidence intervals, 192
- Breakdown point, 236
- Bucketing, *see* Binning
- Chromatography, 21
- Classification and regression trees (CART), 126–135
- Clustering, 79–99
  - average linkage, 81
  - comparing clusterings, 95–97
  - complete linkage, 80
  - hierarchical, 80–84
  - k-means, 85–87
  - k-means clustering, 68
  - k-medoids, 87–90
  - single linkage, 80
  - Ward's method, 81
- Common factors, 63
- Component Detection Algorithm (CODA), 25, 31
- Crossvalidation, 109–111, 177, 181–184, 245
  - double, 183
  - generalized, 182
  - leave-multiple-out, 110, 183, 232
  - leave-one-out (LOO), 110, 181
  - ten-fold, 110, 183
- Data sets
  - gas chromatography, 8, 19
  - gasoline, 7, 18–19, 35, 53, 168–169, 177, 184, 196, 223, 230
  - LC-MS, 11–12, 21–26, 29–30
  - prostate, 9–11, 14, 33–35, 48, 119, 138, 196, 201, 244–250
  - shootout (NIR), 252–254
  - UV, 255, 258–267
  - wine, 9, 36–37, 46, 49, 51, 58, 63, 81, 85, 87, 91, 96, 106, 115, 117, 122, 138, 221, 229, 236
- Discriminant analysis, 104–118
  - canonical, 114
  - diagonal, 119–120
  - Fisher LDA, 111–114
  - linear, 105–108
  - model-based, 116–118
  - PCDA, 244–248

- PLSDA, 248–250
  - quadratic, 114–116
  - regularized, 118–121
  - shrunk centroid, 120–121
- Dual representation, 137
- Elastic net, 216
- Entropy, 130
- Error estimates, 178–179
- Expectation-maximization (EM)
  - algorithm, 90, 92
- Factor analysis, 63–65
- False positive rate, *see* Specificity
- Feed-forward networks, *see* Artificial neural networks (ANNs)
- Finite mixture modelling, *see* Mixture modelling
- Gas chromatography, 8
- Generalized inverse, 148, 262
- Gini index, 130
- Independent component analysis (ICA), 60–62
- Jackknife, 184
- k-nearest-neighbours (KNN), 122–126
- Kennard-Stone algorithm, 176
- Kernel functions, 137
- LDA, *see* Discriminant analysis
- LOO, *see* crossvalidation
- Loss function, 135, 163, 178
- Mahalanobis distance, 105, 107, 110, 115, 123, 212
- Mallows'  $C_p$ , 180
- Mass spectrometry
  - coupled to liquid chromatography (LC-MS), 11
- Metropolis criterion, 218
- Minimum covariance determinant (MCD), 236
- Mixture modelling, 90–94
- Model selection, 179
- Model-based clustering, *see* Mixture modelling
- Moore-Penrose inverse, *see* Generalized inverse
- Multi-layer perceptrons, *see* Artificial neural networks (ANNs)
- Multidimensional scaling (MDS), 57–60, 77
  - classical, 58
  - non-metric MDS, 58
  - Sammon mapping, 58
- Multiplicative scatter correction (MSC), 18, 177
- Near-infrared (NIR) spectroscopy, 7
- Neural networks, *see* Artificial neural networks (ANNs)
  - hidden layer, 141
  - transfer functions, 142
- NP-complete, 4, 130
- Nuclear Magnetic Resonance (NMR), 7, 13, 20, 22, 31, 33
- OPLS, 240–243
- Orthogonal signal correction (OSC), 240
- Outliers, 87, 204, 235–240
- Overfitting, 132, 143, 144, 153, 167, 168, 176, 203, 250
- Peak distortion, 15, 16, 22, 27, 29
- Peak picking, 31–33
- Penalization, 149, 163
- Principal component analysis (PCA), 43–56
  - biplot, 54
  - explained variance, 45
  - loading plot, 47
  - loadings, 43, 45
  - robust, 235–240
  - score plot, 46
  - scores, 43, 45
- Principal coordinate analysis, *see* Multidimensional scaling (MDS)
- Projection pursuit, 60, 237
- Pruning, 132
- QDA, *see* Discriminant analysis
- Rand index (adjusted), 95
- Random Forests, 197–201

- Recall rate, *see* Sensitivity
- Receiver operating characteristic (ROC), 179
- Regression
  - logistic, 170
  - multiple, 145–149
  - PCR, 149–155
  - PLS, 155–163
  - Ridge, 163–164
- Root-mean-square error (RMS), 152, 153, 178, 180
- Savitsky-Golay filter, 16
- Scaling, 33–38
  - autoscaling, 35, 104
  - double centering, 58
  - length scaling, 34
  - mean-centering, 35
  - Pareto scaling, 38
  - range scaling, 34
  - standard normal variate scaling, 37
  - standardization, 35
  - variance scaling, 34, 35
- Self-organising maps
  - initialization, 96
- Self-organizing maps (SOMs), 67–78
  - codebook vectors, 67
  - initialization, 69
  - learning rate, 68
  - topology, 70
  - U-matrix, 72
- Sensitivity, 178
- Shrinkage, *see* Penalization
- Simulated annealing, 218–225
- Singular value decomposition (SVD), 45
- Smoothing, 13–18
  - running mean, 15
  - running median, 16
- Sparseness, 136
- Specific factors, 63
- Specificity, 178
- Support Vector Machines (SVMs), 136–141
- Tanimoto distance, 77
- True positive rate, *see* Sensitivity
- Uniquenesses, 63
- Validation, 103
  - test and training sets, 104, 176
- Variable selection
  - stepwise, 210
- Varimax rotation, 65
- Wavelets, 16