

Multivariate Statistics and Machine Learning

Unit code:	MATH38161
Credit Rating:	10
Unit level:	Level 3
Teaching period(s):	Semester 1
Offered by	School of Mathematics
Available as a free choice unit?:	N

Requisites

Prerequisite

- [MATH20701 - Probability 2](#) (Compulsory)
- [MATH20802 - Statistical Methods](#) (Compulsory)

Desirable

- Good working knowledge in the R statistical programming language
- Good grades in MATH20802

Aims

To familiarise students with the fundamental concepts and ideas underlying multivariate statistical data analysis methods and related supervised and unsupervised machine learning approaches for pattern recognition and classification, as well as with their practical implementation and application using the R statistical programming language.

Overview

Multivariate statistical models and methods are essential for analysing complex-structured and possibly high-dimensional data from any areas of science and industry, ranging from biology and medicine, and genetics to finance and sociology. Multivariate statistics also provides the foundation of many machine learning algorithms.

In the first part of this module covers the foundations of multivariate data analysis, e.g., multivariate random variables, covariance and correlation, and multivariate regression. In addition, related approaches such dimension reduction and latent variable models are discussed.

The second part of the course is concerned with multivariate approaches for statistical learning in supervised and unsupervised settings, including techniques from machine learning, and their application in pattern recognition, classification, and high-dimensional data analysis.

Learning outcomes

On successful completion of the course students will be able to:

- use the programming language R for multivariate data analysis and graphical presentation

- apply dimension reduction techniques such as PCA and CCA
- perform clustering and classification using tools from both statistics and machine learning
- make good choices among available parametric and nonparametric approaches
- analyse high-dimensional data sets with suitable regularisation techniques

Assessment methods

- Other - 20%
- Written exam - 80%

Assessment Further Information

- Coursework (1 written project): 10 hours weighting 20%
- End of semester examination: 2 hours weighting 80%

Syllabus

- Multivariate normal model (4 lectures): distributional properties, estimation of covariance and correlation matrix both in large and small sample settings (using likelihood and regularised/shrinkage estimation), connection with multivariate regression.
- Dimension reduction and latent variable models (4 lectures): whitening transformations, Principle Components Analysis (PCA), Canonical Correlation Analysis (CCA), Factor Analysis (FA)
- Unsupervised learning / clustering (4 lectures): model-based clustering (Gaussian mixture models, EM algorithm), algorithmic approaches (e.g. K-means, hierarchical clustering)
- Supervised learning / classification (4 lectures): Diagonal, Linear, and Quadratic Discriminant Analysis (DDA, LDA, QDA) and regularised versions for high-dimensional data analysis.
- Nonlinear and Nonparametric models (4 lectures): Nonlinear regression, decision trees, random forest, neural networks, Gaussian processes.

Recommended reading

- Härdle, W.K., and L. Simar. 2015. *Applied Multivariate Statistical Analysis*. Fourth edition.
- Hastie, T., R. Tibshirani, and J. Friedman. 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
- James, G., D. Witten, T. Hastie and R. Tibshirani. 2013. *An Introduction to Statistical Learning with Applications in R*. Springer.
- Marden, J.I. 2015. *Multivariate Statistics: Old School*. <http://stat.istics.net/Multivariate>
- Rogers, S. and M. Girolami. 2017. *A first course in machine learning (2nd Edition)*. Chapman and Hall / CRC.

Feedback methods

Computer labs will provide an opportunity for students to try out the methods on real data and to get feedback from the instructor. Coursework and tutorials also provide an opportunity for students to receive feedback. Students can also get feedback on their understanding directly from the lecturer, for example during the lecturer's office hour or after class.

Study hours

- Lectures - 22 hours (11 x 2 hours)
- Tutorials - 10 hours (10 x 1 hours) - 5 tutorials will be computer based sessions
- Independent study hours - 60 hours

Teaching staff

Korbinian Strimmer - Unit coordinator