

# Course and Examination Fact Sheet: Spring Semester 2024

# 10,365: Computational Statistics

ECTS credits: 4

# Overview examination/s

(binding regulations see below) decentral - Written work, Digital, Group work group grade (100%) Examination time: Term time

# Attached courses

Timetable -- Language -- Lecturer 10,365,1.00 (GSERM) Computational Statistics -- English -- <u>Audrino Francesco</u>

# **Course information**

# Course prerequisites

Advanced knowledge in statistics and econometrics.

## Learning objectives

- Students will gain an advanced knowledge on the statistical aspects related to the use of machine learning techniques needed to analyze large or high-dimensional datasets.
- Students will learn how to apply machine learning tools in a responsible way and will properly apply the methods on a concrete dataset of their choice and prepare a research paper summarizing their results.

# Course content

Computational Statistics is the area of specialization within statistics that includes statistical visualization and other computationally-intensive methods of statistics for mining large, nonhomogeneous, multi-dimensional datasets so as to discover knowledge in the data. As in all areas of statistics, probability models are important, and results are qualified by statements of confidence or of probability. An important activity in computational statistics is model building and evaluation.

First, the basic multiple linear regression is reviewed. Then, some nonparametric procedures for regression and classification are introduced and explained. In particular, Kernel estimators, smoothing splines, classification and regression trees, additive models, projection pursuit and eventually neural nets will be considered, where some of them have a straightforward interpretation, other are useful for obtaining good predictions.

The main problems arising in computational statistics like the curse of dimensionality will be discussed. Moreover, the goodness of a given (complex) model for estimation and prediction is analyzed using resampling, bootstrap and cross-validation techniques.

# Course structure and indications of the learning and teaching design

### Outline:

### 1. Overview of supervised learning

Introductory examples, two simple approaches to prediction, statistical decision theory, local methods in high dimensions, structured regression models, bias-variance tradeoff, multiple testing and use of p-values.

### 2. Linear methods for regression



Multiple regression, analysis of residuals, subset selection and coefficient shrinkage.

#### 3. Methods for classification

Bayes classifier, linear regression of an indicator matrix, discriminant analysis, logistic regression.

#### 4. Nonparametric density estimation and regression

Histogram, kernel density estimation, kernel regression estimator, local polynomial nonparametric regression estimator, smoothing splines and penalized regression.

#### 5. Model assessment and selection

Bias, variance and model complexity, bias-variance decomposition, optimism of the training error rate, AIC and BIC, cross-validation, boostrap methods.

### 6. Flexible regression and classification methods

Additive models; multivariate adaptive regression splines (MARS); neural networks; projection pursuit regression; classification and regression trees (CART).

#### 7. Bagging and Boosting

The bagging algorithm, bagging for trees, subagging, the AdaBoost procedure, steepest descent and gradient boosting.

#### 8. Introduction to the idea of a Superlearner

## **Course literature**

### Main references:

- F. Audrino, Lecture Notes.
- Hastie T., Tibshirani, R. and Friedman, J. (2001). *The elements of statistical learning: data mining, inference and prediction,* Springer Series in Statistics, Springer, Canada.
- Bühlmann, P. and van de Geer, S. (2011). <u>Statistics for High-Dimensional Data: Methods, Theory and Applications.</u> Springer.
- van der Laan, M.J. and Rose, S. (2011). Targeted Learning: Causal Inference for Observational and Experimental Data. Springer.

References to related published papers / chapters of books will be given during the course.

### Additional course information

--

# Examination information

## Examination sub part/s

# 1. Examination sub part (1/1)

#### Examination modalities

Examination type	Written work
Responsible for organisation	decentral
Examination form	Written work
Examination mode	Digital
Time of examination	Term time
Examination execution	Asynchronous
Examination location	Off Campus

Fact sheet version: 1.0 as of 04/12/2023, valid for Spring Semester 2024



Grading type	Group work group grade
Weighting	100%
Duration	

Examination languages Question language: English

Answer language: English

Remark

#### Examination-aid rule Free aids provision

Basically, students are free to choose aids. Any restrictions are defined by the faculty members in charge of the examination under supplementary aids.

#### Supplementary aids

- -

# Examination content

### Outline:

### 1. Overview of supervised learning

Introductory examples, two simple approaches to prediction, statistical decision theory, local methods in high dimensions, structured regression models, bias-variance tradeoff, multiple testing and use of p-values.

### 2. Linear methods for regression

Multiple regression, analysis of residuals, subset selection and coefficient shrinkage.

### 3. Methods for classification

Bayes classifier, linear regression of an indicator matrix, discriminant analysis, logistic regression.

### 4. Nonparametric density estimation and regression

Histogram, kernel density estimation, kernel regression estimator, local polynomial nonparametric regression estimator, smoothing splines and penalized regression.

#### 5. Model assessment and selection

Bias, variance and model complexity, bias-variance decomposition, optimism of the training error rate, AIC and BIC, cross-validation, boostrap methods.

#### 6. Flexible regression and classification methods

Additive models; multivariate adaptive regression splines (MARS); neural networks; projection pursuit regression; classification and regression trees (CART).

### 7. Bagging and Boosting

The bagging algorithm, bagging for trees, subagging, the AdaBoost procedure, steepest descent and gradient boosting.

### 8. Introduction to the idea of a superlearner

# Examination relevant literature

F. Audrino, Lecture Notes, available on Canvas before the beginning of the course.

Fact sheet version: 1.0 as of 04/12/2023, valid for Spring Semester 2024



## Please note

Please note that only this fact sheet and the examination schedule published at the time of bidding are binding and takes precedence over other information, such as information on StudyNet (Canvas), on lecturers' websites and information in lectures etc.

Any references and links to third-party content within the fact sheet are only of a supplementary, informative nature and lie outside the area of responsibility of the University of St.Gallen.

Documents and materials are only relevant for central examinations if they are available by the end of the lecture period (CW21) at the latest. In the case of centrally organised mid-term examinations, the documents and materials up to CW 13 are relevant for testing.

Binding nature of the fact sheets:

- Course information as well as examination date (organised centrally/decentrally) and form of examination: from bidding start in CW 04 (Thursday, 25. Januar 2024);
- Examination information (supplementary aids, examination contents, examination literature) for decentralised examinations: in CW 12 (Monday, 18 March 2024);
- Examination information (supplementary aids, examination contents, examination literature) for centrally organised mid-term examinations: in CW 13 (Monday, 25 March 2024);
- Examination information (regulations on aids, examination contents, examination literature) for centrally organised examinations: Starting with de-registration period in CW 15 (Monday, 08. April 2024).