

MA8701 Advanced methods in statistical inference and learning

Part 5: Closing

Mette Langaas

4/13/23

Outline

- ▶ Learning outcomes and compulsory activity
- ▶ Take home messages from the Data analysis project
- ▶ Final evaluation
- ▶ Plan for the last three sessions!

Learning outcome

(The student should be able to)

1. Knowledge

- ▶ Understand and explain the central theoretical aspects in statistical inference and learning.
- ▶ Understand and explain how to use methods from statistical inference and learning to perform a sound data analysis.
- ▶ Be able to evaluate strengths and weaknesses for the methods and choose between different methods in a given data analysis situation.

2. Skills

Be able to analyse a dataset using methods from statistical inference and learning in practice (using R or Python), and give a good presentation and discussion of the choices done and the results found.

3. Competence

- ▶ The students will be able to participate in scientific discussions, read research presented in statistical journals.
- ▶ They will be able to participate in applied projects, and analyse data using methods from statistical inference and learning.

Compulsory activity 2023

- ▶ Data analysis project (analyse, write report, review)
- ▶ Article presentation (present and discuss)

Take home messages from the Data analysis project

Short overview of the five data analysis projects

Team SuperGreat

- ▶ Data set: Framingham coronary heart disease (CHD), $N=4238$, $p=15$ (categorical, binary, continuous), binary response (15.2% cases).
- ▶ Aim: Understanding effects of covariates for prediction of CHD (10 years follow-up) and compare complete case and single imputation results.
- ▶ Missing: 13.7% in total (highest for glucose with 9%)
- ▶ Methods used: Single imputation vs complete case, bootstrapping, lasso logistic and logistic regression, AIC.
- ▶ Result: important risk factors are age, male, systolic blood pressure, glucose (and for the imputed data also cigarettes per day smoked).

Team CDF

- ▶ Data set: wine quality, $N=6497$, $p=12$ (binary, continuous), binary response (from dividing approximately in two).
- ▶ Aim: investigate how different physiochemical variables affect wine quality.
- ▶ Missing: 22.5 % in total (but not all imputed)
- ▶ Methods used: Single imputation vs complete case, lasso logistic and logistic regression, polyheder inference, train-test split for ROC-AUC.
- ▶ Result: important variables for wine quality was volatile acidity, residual suger, free sulfur dioxide, total sulfr dioxide, sulphates and alcohol.

Team Balance

- ▶ Data set: robotic arm kinematic data, $N=17560$ but reduce to $N=176$ to avoid time series correlations, $p=28$.
- ▶ Aim: A theoretical model for the movement of the robot arm exists, involving trigonometrical functions - giving background to considering a sum of second order polynomials of the covariates. The aim is then to arrive at an interpretable simplified model.
- ▶ Missing: no missing data.
- ▶ Methods used: ACF/PACF, train-test split, OLS, elastic net, multi-sample splitting (median) on training data.
- ▶ Result: Only one covariate “left” after multi-splitting, and this covariate did not give a sensible physical interpretation.

Team JAA

- ▶ Data set: superconductor critical temperature, $N=21263$, $p=82$ (very multicollinear), response: continuous critical temperature.
- ▶ Aim: To construct a prediction method for the critical temperature of the superconductor so that the important factors influencing the critical temperature is understood.
“Critical temperature (of a substance) can be defined as the highest possible temperature value at which the substance can exist as a liquid.”
- ▶ Missing: no missing data (?)
- ▶ Methods: Forward selection with least squares, lasso, group lasso. Bootstrapping on single split data.
- ▶ Results: None of the models gave a small and interpretable model.

Team JKP

- ▶ Data set: Genome-wide association study, $N=1796$ $p=183155$ SNPs (for model selection) + 4 (23) clinical covariates, response=length of house sparrow wing at age 1 year.
- ▶ Aim: Which SNPs are associated with the response?
- ▶ Missing: present. Quality control defaults to removing SNPs and individuals with high missing rate. For the remaining missing data are imputed by single (mean) imputation (and often totally imputed SNPs are analysed).
- ▶ Method: Lasso regression (with snpnet and manual 10-fold CV for λ) and multi-sample splitting.
- ▶ Results: all SNPs had adjusted p -values of 1. No findings.

Group work

- 1) For all groups
 - ▶ What are you most proud of in your work?
 - ▶ What could have been done differently?
 - ▶ Choose one learning experience to share!
- 2) Specific questions for each group on paper hand-out
- 3) If you finish before we summarize: Discuss your study plan for the oral exam

Common themes

Negative (or no) results!

Badly behaved data

Regression model

A linear regression model (or the linear predictor in the GLM) is linear in the regression parameters, not necessarily in the covariates. In addition interaction term may be needed for a good model.

For tree-based methods any non-linearity in the covariates and interactions between covariates are *easily* picked up, but for methods like the lasso, we need to specify the linear predictor ourselves.

How to make sure the “right” linear predictor is used?

Methods for model selection

Comparing results across data sets and models/methods

Data leakage

Missing imputation

- ▶ Specification of the imputation model in missing imputation
- ▶ Should the analysis model response be a covariate in the imputation model?

Combining several challenges

There does (to my knowledge) not exist any off-shelf workflow to combine

- ▶ multiple imputation
- ▶ model selection
- ▶ and model evaluation.

How can that be done?

Evaluation: Oral exam

May 10, 15 and 22.

Pass/fail, with B as pass limit.

- ▶ On the last lecture (April 24) a list of five possible topics (questions) will be available at <https://wiki.math.ntnu.no/ma8701/2021v/exam>.
- ▶ If you want you may prepare a 5-10 minutes presentation of one of the topics (bring notes, but no slides, talk and write by hand) to be held in the start of the oral exam.
- ▶ The rest of the exam is general questions from the reading list (no notes)

Total duration < 30 minutes.

Plan ahead

- ▶ Next week: https://wiki.math.ntnu.no/ma8701/2023v/assignmentsap#presentation_schedule
- ▶ Monday April 24: Discussion on central topics for each part of the course, and present the five possible topics to prepare for the first part of the oral exam.