

Applied Data Analysis and Machine Learning: Introduction to the course

Morten Hjorth-Jensen^{1,2}

¹Department of Physics, University of Oslo

²Department of Physics and Astronomy and National Superconducting Cyclotron Laboratory, Michigan State University

Dec 6, 2018

Suggested Literature

Possible textbooks

- HTF: Trevor Hastie, Robert Tibshirani, Jerome H. Friedman, The Elements of Statistical Learning, Springer
- AG: Aurelien Geron, Hands-On Machine Learning with Scikit-Learn and TensorFlow, O'Reilly

Teachers

Teachers :

1. Kristine B. Heine
2. Morten Hjorth-Jensen
3. Bendik Samseth

Lectures and ComputerLab

- Time: January 21 to February 1, 2019.
- Place: GANIL, Caen

Organization of the day.

1. Lectures 9am-12pm
2. Lunch 12pm-2pm
3. Computerlab 2pm-6pm

Learning outcomes

The course introduces a variety of central algorithms and methods essential for studies of data analysis and machine learning. The course is project based and through the various projects, normally three, you will be exposed to fundamental research problems in these fields, with the aim to reproduce state of the art scientific results. You will learn to develop and structure large codes for studying these systems, get acquainted with computing facilities and learn to handle large scientific projects. A good scientific and ethical conduct is emphasized throughout the course. More specifically, after this course you will

- Learn about basic data analysis, Bayesian statistics, Monte Carlo methods, data optimization and machine learning;
- Be capable of extending the acquired knowledge to other systems and cases;
- Have an understanding of central algorithms used in data analysis and machine learning;
- Gain knowledge of central aspects of Monte Carlo methods, Markov chains, Gibbs samplers and their possible applications;
- Understand linear methods for regression and classification;
- Learn about neural network, genetic algorithms and Boltzmann machines;
- Work on numerical projects to illustrate the theory. The projects play a central role and you are expected to know modern programming languages like Python or C++.

Topics covered in this course: Statistical analysis and optimization of data

The following topics will be covered

- Basic concepts, expectation values, variance, covariance, correlation functions and errors;
- Simpler models, binomial distribution, the Poisson distribution, simple and multivariate normal distributions;

- Central elements of Bayesian statistics and modeling;
- Central elements from linear algebra
- Cubic splines and gradient methods for data optimization
- Monte Carlo methods, Markov chains, Metropolis-Hastings algorithm, ergodicity;
- Linear methods for regression and classification;
- Estimation of errors using cross-validation, blocking, bootstrapping and jackknife methods;
- Principal Component Analysis and Clustering algorithms;

Topics covered in this course: Machine Learning

- Linear and non-linear regression
- Boltzmann machines;
- Neural networks;
- Decisions trees and random forests
- Support vector machines

Extremely useful tools, strongly recommended

Discussed at the lab sessions.

- GIT for version control
- ipython/jupyter notebook
- Anaconda and other Python environments

Course Content and detailed plan

Lectures are approximately 45 min each with a small break between each lecture. There is also a coffee break of 30 min in the morning sessions. It will most likely be scheduled around 1045am and is not marked in the program below. Lunch is from 12pm to 2pm. The lab sessions start at 2pm and end at 6pm. The acronyms are

- BS: Bendik Samseth
- KBH: Kristine B. Hein
- MHJ: Morten Hjorth-Jensen

Week 1, January 21-26.

Day	Lecture Topics and lecturer		P
Monday 21	9am-945am	Introduction and welcome (MHJ)	
	10am-1045am	Review of Python and Linear Algebra (MHJ)	
	1045am-1115am	Break	
	1115am-12pm	Getting started with Linear Regression (MHJ)	
	12pm-2pm	Lunch +own activities	
	2pm-6pm	Python installations and setups, anaconda and more	Exerci
Tuesday 22	9am-945am	Linear Regression (MHJ)	
	10am-1045am	Linear Regression, Lasso and Ridge (MHJ)	
	1045am-1115am	Break	
	1115am-12pm	Linear Regression, Lasso and Ridge (MHJ)	
	12pm-2pm	Lunch +own activities	
	2pm-6pm	Brief intro to scikit-learn	Exerci
Wednesday 23	9am-945am	Summary of linear regression (MHJ)	
	10am-1045am	Statistical analysis of data, bias and variance(MHJ)	
	1045am-1115am	Break	
	1115am-12pm	Statistical analysis of data, bias and variance (MHJ)	
	12pm-2pm	Lunch +own activities	
	2pm-6pm	More scikit-learn functionality	Exerci
Thursday 24	9am-945am	Statistical analysis, cross-validation and Bootstrap (MHJ)	
	10am-1045am	Statistical analysis, cross-validation and Bootstrap (MHJ)	
	1045am-1115am	Break	
	1115am-12pm	Optimization and gradient descent (MHJ)	
	12pm-2pm	Lunch +own activities	
	2pm-6pm	Statistics tools	Exerci
Friday 25	9am-945am	Optimization and Gradient descent (MHJ)	
	10am-1045am	Optimization and gradient descent(MHJ)	
	1045am-1115am	Break	
	1115am-12pm	Logistic regression and classification (MHJ)	
	12pm-2pm	Lunch +own activities	
	2pm-6pm	Gradient descent coding	Pro
Saturday 26	9am-945am	Logistic regression and classification (MHJ)	
	10am-1045am	Logistic Regression and classification (MHJ)	
	1045am-1115am	Break	
	1115am-12pm	Start with Neural networks (MHJ)	
	12pm-2pm	Lunch +own activities	
	2pm-6pm	Getting familiar with Tensorflow and Keras	Pro

Week 2, January 28-31.

Day	Lecture Topics and lecturer	
Monday 28	9am-945am	Neural networks (MHJ)
	10am-1045am	Neural networks and back propagation algorithm (MHJ)
	1045am-1115am	Break
	1115am-12pm	Neural networks and back propagation algorithm, setting up your c
	12pm-2pm	Lunch +own activities
	2pm-6pm	Tensorflow and Keras examples
Tuesday 29	9am-945am	Developing a neural network code (MHJ)
	10am-1045am	Convolutional Neural Network and Nuclear Physics experimen
	1045am-1115am	Break
	1115am-12pm	Convolutional Neural Networks (CNN) and NP experiments (
	12pm-2pm	Lunch +own activities
	2pm-6pm	Tensorflow/Keras and CNNs
Wednesday 30	9am-945am	Support Vector Machines (MHJ)
	10am-1045am	Support Vector Machines(MHJ)
	1045am-1115am	Break
	1115am-12pm	Support Vectort Machines and Decision trees (MHJ)
	12pm-2pm	Lunch +own activities
	2pm-6pm	Hints for project 2
Thursday 31	9am-945am	Decision Trees and Bagging (MHJ)
	10am-1045am	Bagging, Ensembles and Random Forests (MHJ)
	1045am-1115am	Break
	1115am-12pm	Summary of course and perspectives (MHJ)
	12pm-2pm	Lunch +own activities
	2pm-6pm	Finalize project 2

Teaching and projects

The course will be taught as an intensive course of duration of two weeks, with a total time of 30 h of lectures, 40 h of exercises and 2 projects that will be graded and form the final assessment. Each project counts 50% of the final grade. They have to be finalized three weeks after the course ended.

The final assignment will be graded with marks A, B, C, D, E and failed.

The organization of a typical course day is as follows:

Time	Activity
9am-12pm	Lectures, project relevant information and directed exercises
12pm-2pm	Lunch and own activities
2pm-6pm	Computational projects, exercises and hands-on sessions
6pm-7pm	Wrap-up of the day and eventual student presentations and/or further discussions