Applied Machine Learning Introduction to the Course



CHALMERS

Richard Johansson

richajo@chalmers.se

Welcome to the course!

- Machine learning is increasingly popular among students
 - our courses take increasing volumes
 - many thesis projects develop or apply ML models
- ...and in industry, public sector
 - many companies come to us looking for students
 - joint research projects
- Why the fuss and why now?

Success stories: image recognition



Success stories: machine translation



Decoder



Data

Under the bonnet

How a self-driving car works







50X BOOST IN DEEP LEARNING IN 3 YEARS



AlexiVet training throughput based on 20 iterations, CPU: Ix ES-2680v3 12 Core 2.5GHz. 128G8 System Memory, Ubuntu 14.04







Applications...











Topics covered in the course

- The usual "zoo": a selection of machine learning models
 - what's the idea behind them?
 - · how are they implemented? (at least on a high level)
 - what are the use cases?
 - how can we apply them practically?
- But hopefully also the "real-world context":
 - extended "messy" practical assignments requiring that you think of what you're doing
 - invited talks from industry and/or the healthcare sector
 - annotation of data, evaluation
 - ethical and legal issues, interpretability



Practical issues about the course

Fundamental concepts in machine learning

Machine learning libraries in Python



Course webpage

• The official course webpage is the Canvas page https://chalmers.instructure.com/courses/33104



People involved in the course

- Richard: examiner, responsible for the course
- Anton, Jack, Newton, Selma, Philipp, Laleh, Styrbjörn: helping you with the assignments



Structure of teaching

- Lecture discussions Tuesdays and Fridays 13–15
 - we will use a flipped classroom format with pre-recorded lectures you are expected to have watched before the session
 - summary and discussion of the content of the recorded lectures
 - interactive coding
 - solving a few exercises when we have time
 - feel free to ask questions before the session!
- Assistance sessions Thursdays 13-17
 - our TAs help you work on your assignments
 - please let me know if it's too crowded
 - in a computer lab room (with possibly additional remote sessions)

Assignments

- Five compulsory assignments:
- $\operatorname{PA}\ 1$ intro to the ML workflow, decision trees
- PA 2 random forests
- PA 3 text classification
- PA 4 neural network software
- PA 5 medical image classification
- We will use the **Python** programming language
- Please refer to the course PM for details about grading
- Assignments are done in groups

Programming assignment 1

- Warmup lab exercise: quick tour of the scikit-learn library
- Introduction to decision trees
- For a high grade: implement decision tree regression
- Assistance sessions this Thursday
- Submission deadline: January 27

Literature

- We won't follow a book closely, but we'll give pointers to reading material in this book:
 - Machine Learning: A course for engineers and scientists by Lindholm et al: http://smlbook.org/
- And additional papers to read for some topics
- Some notes to complement the lectures
- Example code will be posted on the course page

Additional material along the way

- Exercise sheets, old exams
- Online quizzes



Exam, mid-March

- This is a take-home exam: a written assignment
- Will be available during the whole exam period
- Two-part structure:
 - 1. a first compulsory part about basic concepts: you need to answer most of these questions correctly to pass
 - 2. a second optional part that requires more insight: answer these questions for a higher grade

Student representatives

- If you're interested in being a student representative, please send me an email!
- The workload is light and there will be a small reward...





Practical issues about the course

Fundamental concepts in machine learning

Machine learning libraries in Python



Predictive models

- Given some object, make a prediction
 - is this patient diabetic?
 - what animal does this image show?
 - what is the market value of this apartment?
 - what are the phonemes contained in this speech signal?

Predictive models

- Given some object, make a prediction
 - is this patient diabetic?
 - what animal does this image show?
 - what is the market value of this apartment?
 - what are the phonemes contained in this speech signal?
- The goal of machine learning is to build the predictive models by observing data

Predictive models

- Given some object, make a prediction
 - is this patient diabetic?
 - what animal does this image show?
 - what is the market value of this apartment?
 - what are the phonemes contained in this speech signal?
- The goal of machine learning is to build the predictive models by observing data
- Contrast: expert-defined or data-driven







Why machine learning?

Why would we want to "learn" the function from data instead of just implementing it?

- Usually because we don't really know how to write down the function by hand
 - speech recognition
 - image classification
 - machine translation
 - ...
- Might not be necessary for limited tasks where we know
- What is more expensive in your case? knowledge or data?

Don't forget your domain expertise!

ML makes some tasks automatic, but we still need our brains:

- defining the tasks, terminology, evaluation metrics
- annotating (hand-labeling) training and testing data
- designing features
- error analysis

Example: is the patient diabetic?





Example: is the patient diabetic?



• In order to predict, we make some measurements of properties we believe will be useful: these are called the **features**

More terminology: what is the output?

- Classification: learning to output a category label
 - spam/non-spam; positive/negative; ...
- Regression: learning to guess a number
 - value of a share; number of stars in a review; ...



How is the training signal provided?

- In supervised learning, the training set consists of input-output pairs
- our goal is to learn to produce the outputs

age	workclass	fniwgt	education	education-num	marital- status	occupation	relationship	race	sex	capital-gain	capital-loss	hours- per-week	native- country	arge
39	State-gov	77516	Bachelors	13	Never- married	Adm-clerical	Not-In-family	White	Male	2174	0	40	United- States	<=50K
50	Self-emp- not-inc	83311	Bachelors	13	Married- chr-spouse	Exec-managerial	Husband	White	Male	0	0	13	United- States	<=50K
38	Private	215646	HS-grad	9	Divorced	Handlers- cleaners	Not-in-family	White	Male	0	0	40	United- States	<=50K
53	Private	234721	11th	7	Married- clv-spouse	Handlers- cleaners	Husband	Black	Male	0	0	40	United- States	<-50K
28	Private	338409	Bachelors	13	Married- chr-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	a=50H
		F1	F2	F3		F4	F5		F6	F	7 F8	F	9 F	NSD
0	13558	.30	4305.35	0.31754	162.17	30 1.8727	791e+06	215	.3590	4287.8	7 102	27.030	2 17	.284
1	6191	.96	1623.16	0.26213	53.38	94 8.0344	467e+05	87	.2024	3328.9	39	38.546	86	.021
2	7725	.98	1726.28	0.22343	67.28	87 1.0756	648e+06	81	.7913	2981.04	4 29	38.811	9 9	.275
3	8424	.58	2368.25	0.28111	67.83	25 1.2104	172e+06	109	.4390	3248.2	2 70	39.065	1 15	.851
4	7460	.84	1736.94	0.23280	52.41	23 1.0210)20e+06	94	.5234	2814.4	2 41	39.914	7 7	.962



Types of supervision: alternatives

- Unsupervised learning: we are given "unorganized" data
 - our goal is to discover some structure



- Reinforcement learning: our problem is formalized as a game
 - an agent carries out actions and receives rewards





Example: Fisher's iris data





	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa



Approach 1: linear separator

if $0.85 \cdot \text{petal_length} + 2.42 \cdot \text{petal_width} \ge 8.34$: return virginica else

return versicolor



Approach 2: if/then/else tree



CHALMERS | 🕘 UNIVERSITY OF GOTHENBURG

Basic supervised machine learning workflow





Basic ML methodology: evaluation

- Select an evaluation procedure (a "metric") such as
 - classification accuracy: proportion correct classifications?
 - mean squared error often used in regression
 - or some domain-specific metric
- Compare to one or more baselines
 - trivial solution
 - rule-based solution
 - existing solution
- Apply your model to a held-out test set and evaluate
 - the test set must be different from the training set
 - also: don't optimize on the test set; use a development set or cross-validation!

Managing your data





Managing your data





Managing your data





Managing your data for evaluation and cross-validation







Practical issues about the course

Fundamental concepts in machine learning

Machine learning libraries in Python



Use cases for machine learning

• **Standard use cases**: standard solutions are available

• Special cases: we may need to tailor our own solutions







The Python machine learning ecosystem (selection)





Machine learning software: a small sample

- General-purpose software, large collections of algorithms:
 - scikit-learn: http://scikit-learn.org
 - Python library will be used in this course
 - Weka: http://www.cs.waikato.ac.nz/ml/weka
 - Java library with nice user interface
- Special-purpose software, small collections of algorithms:
 - Keras, PyTorch, TensorFlow, JAX for neural networks
 - LibSVM/LibLinear for support vector machines
 - XGboost, lightgbm for tree ensembles
 - ...
- large-scale learning in distributed architectures:
 - Spark MLLib
 - H2O

Scikit-learn toy example

See also

https://scikit-learn.org/stable/getting_started.html



Up next

- Thursday: lab sessions for programming assignment 1
- Topic of Friday's discussion:
 - decision trees
 - ensembles and random forests
 - generalization, under/overfitting
- Please prepare for assignment 1 by reading my code and the extra reading on decision trees