

## W4995 Applied Machine Learning

# Introduction

01/23/19

Andreas C. Müller

# Logistics

## Email

andreas.mueller@columbia.edu (NOT amueller who is someone else)

**CAs: Pranjal Bajaj, Ujjwal Peshin, Liyan Nie, Yao Fu, Luv Aggarwal, Sukriti Tiwari**

## Office Hours

- Andreas Müller Wednesdays 10am-11am, Interchurch 320 K
- CA office hours: TBA

# Logistics

- Course website <http://www.cs.columbia.edu/~amueller/comsw4995s19/>
- Six programming assignments
- Grade: 60% homeworks, 20% first exam, 20% second exam

# Slides and course materials

The screenshot shows a video player interface. At the top right is a timestamp **0:00:00**. Below it is a small video thumbnail showing a slide with the title "Slides and course materials". The main content area displays a slide with the title "What and Why of Machine Learning". To the right of the slide text is a note:  
As I told you, you can find the slides on the website or on github. I'll probably publish slides with and without notes. My goal is to write a note for each slide that says what I'm gonna say. So it's like a script for the lecture. Here's what this looks like. Unfortunately there'll be no video recording, but maybe these notes will be helpful. I might also publish some Jupyter notebooks at some point, if I feel they might help. In general, I'll mostly do standard slides, though.

Using markdown with remark. Press "p" for notes.

# Lecture Recordings

# Plagiarism and Code copying

- Homeworks are checked for plagiarism
- Copied code will result in 0 points for all involved
- Copying from my slides or online sources (Stack overflow, tutorials, etc. ) is fine.

# Scikit-learn Development



<http://scikit-learn.org/dev/developers/contributing.html>

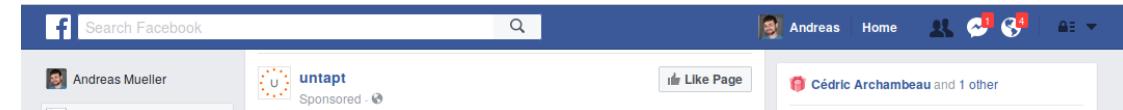
# Books

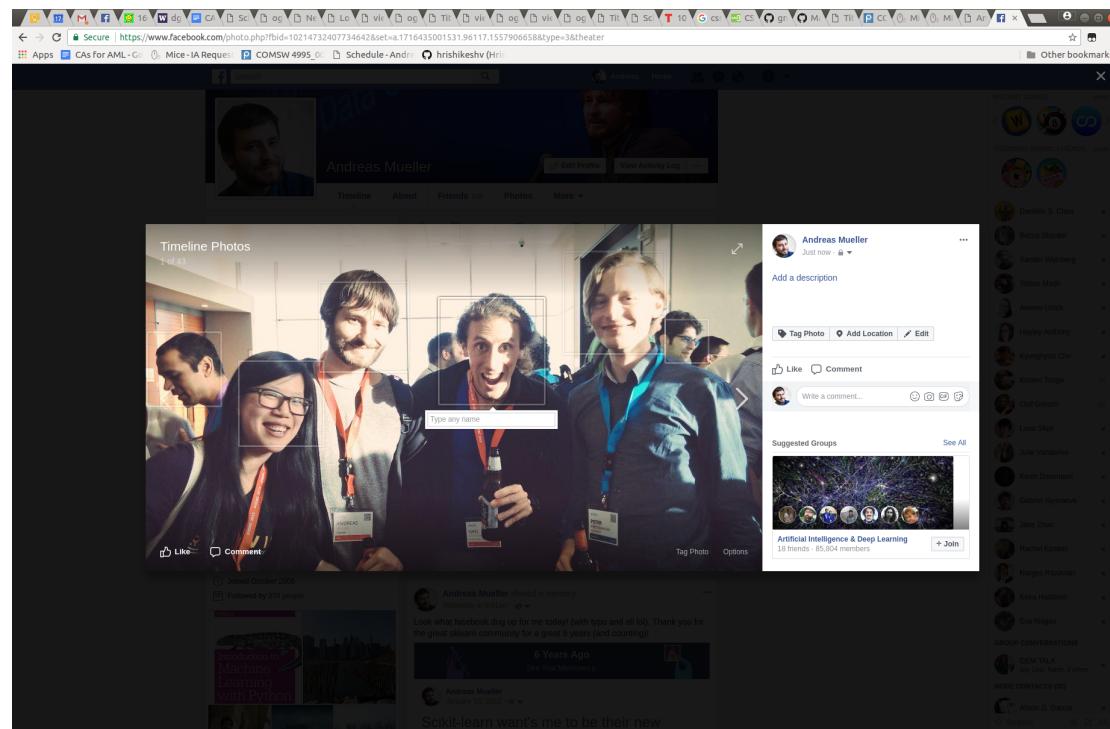
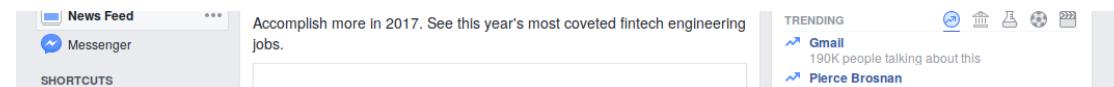




# What and Why of Machine Learning

# What is machine learning?





Search for people, places and things

Andreas Mueller added 8 new photos.

June 6 at 5:07pm · [View on Facebook](#)

Bye bye hong kong

Like · Comment · Share

1 Like · [Comment](#) · [Share](#)

Looks like such a vibrant city! I bet it was awesome 😊

June 6 at 7:39pm · [Like](#)

Write a comment...

SPONSORED

**Da ist Abwechslung drin!**  
franziskaner-weissbier.de

3 Sorten Franziskaner Weissbier in einem Pack. Hier klicken und 1 von 100 Packs gewinnen!

Singles auf Facebook

Schau dir Dating-Profile von Singles in deiner Nähe an.

Online Essen bestellen

Neu mit Lieferheld: PLZ eingeben, Lieferdienst finden und genießen!

Globaler Chauffeurservice

Fahren Sie eine Klasse besser zu günstigen Preisen – Blacklane ist weltweit verfügbar!

28,723 people like this

Sanssouci...

Geld vom Staat zurück!

Steuererklärung preiswert für Arbeitnehmer, Azubis Arbeitsuchende, Rentner Pensionäre etc.

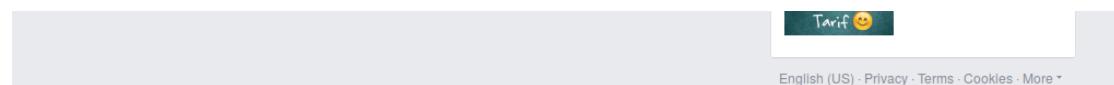
Like Page · 4 people like this page

Der neue WhatsApp Tarif!

eplus.de

WhatsApp SIM – der revolutionäre Prepaid Tarif

Der Prepaid



13 / 60

www.amazon.com/s/ref=nb\_sb\_noss\_1?url=search-alias%3Daps&field-keywords=machine learning&sprefix=machine+%2Caps&rh=

1-16 of 21,549 results for "machine learning"

Books > Machine Learning > Computer & Technology > Reference > Programming Algorithms > See more

Kindle Store > Computers & Technology > Computer Programming > See more

See All 30 Departments

Refine by Eligible for Free Shipping Free Shipping by Amazon Book Series Machine Learning Science and Statistics Adaptive Computation and Machine Learning series I Can Draw Use R! Book Language English Book Format Hardcover Paperback Kindle Edition HTML Avg. Customer Review 4.5 & Up 4.5 & Up 4.5 & Up 4.5 & Up

Machine Learning Resources Find tips and tricks with these featured titles on machine learning, algorithms, sensors, and more. Learn more

Related Searches: artificial intelligence, data mining, machine learning python

**Practical Machine Learning: Innovations in Recommendation** by Ted Dunning and Ellen Friedman (Apr 17, 2014)

\$0.00 Kindle Edition  
Autodelivered wirelessly

**Machine Learning: The Art and Science of Algorithms that Make Sense of Data** by Peter Flach (Nov 12, 2012)

\$44.96 Paperback Prime  
Only 18 left in stock - order soon.  
More Buying Choices - Paperback  
\$46.98 new (40 offers)  
\$35.00 used (14 offers)

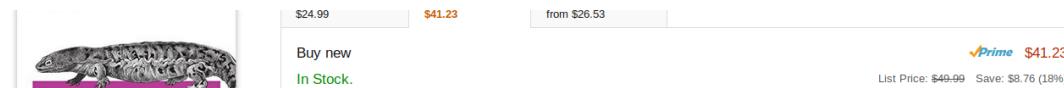
**Understanding Machine Learning: From Theory to Algorithms** by Shai Shalev-Shwartz and Shai Ben-David (May 19, 2014)

\$34.35 Kindle Edition  
Autodelivered wirelessly  
More Buying Choices - Hardcover  
\$50.92 new (10 offers)  
\$52.02 used (4 offers)

**Learning From Data** by Yaser S. Abu-Mostafa, Malik Magdon-Ismail and Hsuan-Tien Lin (Mar 27, 2012)

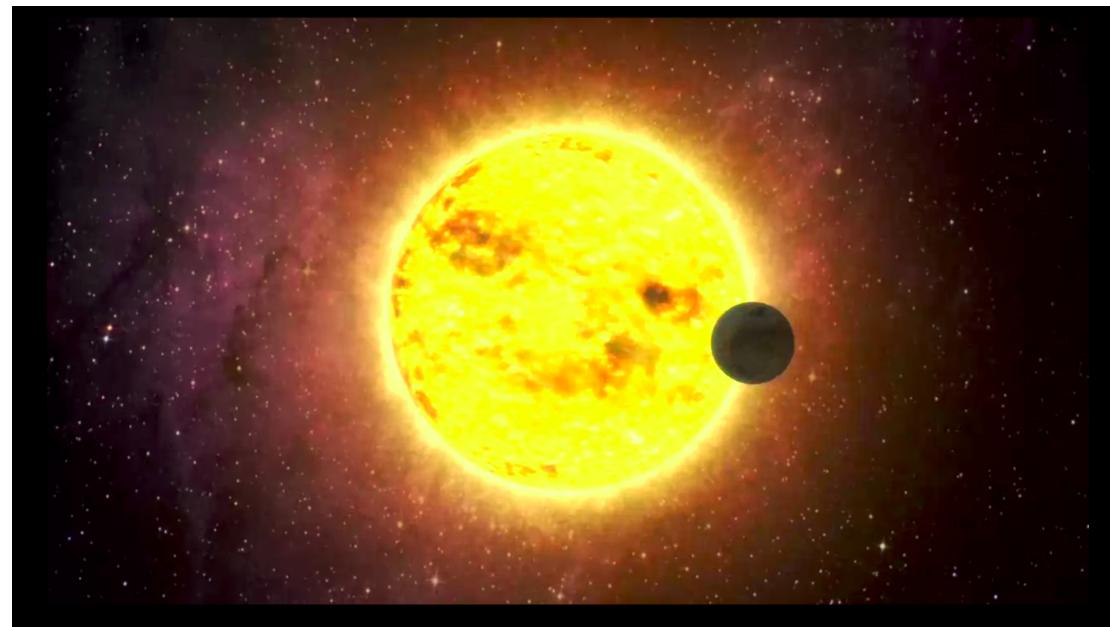
\$28.00 new (9 offers)  
\$40.00 used (13 offers)





15 / 60

# Science!



# Types of Machine Learning

# Types of Machine Learning

- Supervised
- Unsupervised
- Reinforcement

# Supervised Learning

$$(x_i, y_i) \propto p(x, y) \text{ i.i.d.}$$

$$x_i \in \mathbb{R}^p$$

$$y_i \in \mathbb{R}$$

$$f(x_i) \approx y_i$$

# Generalization

Not only  
also for new data:

$$\begin{aligned}f(x_i) &\approx y_i, \\f(x) &\approx y\end{aligned}$$

# Examples of Supervised Learning

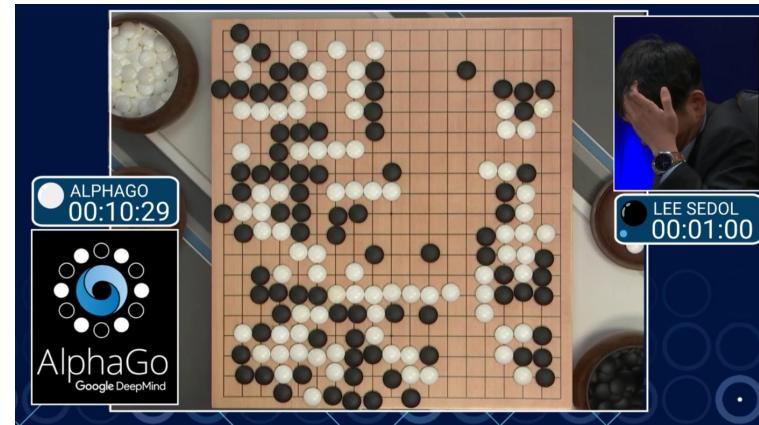
- spam detection
- medical diagnosis
- add click prediction

# Unsupervised Learning

$$x_i \propto p(x) \text{ i.i.d.}$$

Learn about  $p$ .

# Reinforcement Learning



# Explore & Learn

# Other kinds of learning

- Semi-supervised
- Active Learning
- Forecasting
- ...





# Classification and Regression

## Classification

- target  $y$  discrete
- Will you pass?

## Regression

- target  $y$  continuous
- How many points will you get in the exam?

# Relationship to Statistics

## Statistics

- model first
- inference emphasis

## Machine learning

- data first
- prediction emphasis

# Relationship to Statistics

## Statistics

- model first
- inference emphasis

## Machine learning

- data first
- prediction emphasis

# Guiding Principles in Machine Learning

# Goal considerations

# The Cost of Complex Systems

Data driven first? yes! (or maybe)

Machine Learning first: No!

**Thinking in Context!**  
**What is the baseline?**  
**What is the benefit?**

# Good and Bad Substitutes

# Communicating Results

# Explainable Results





36 / 60

# Sidebar: Ethical Considerations



<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

## Summary: Supervision

37 / 60

# Ethics: It's in the application!

# Data and Data Collection

# Free vs Expensive Data

Free

Expensive

# Free vs Expensive Data

## Free

Predict observable events

- Stock market
- Clicks
- House numbers

## Expensive

Automate complex process

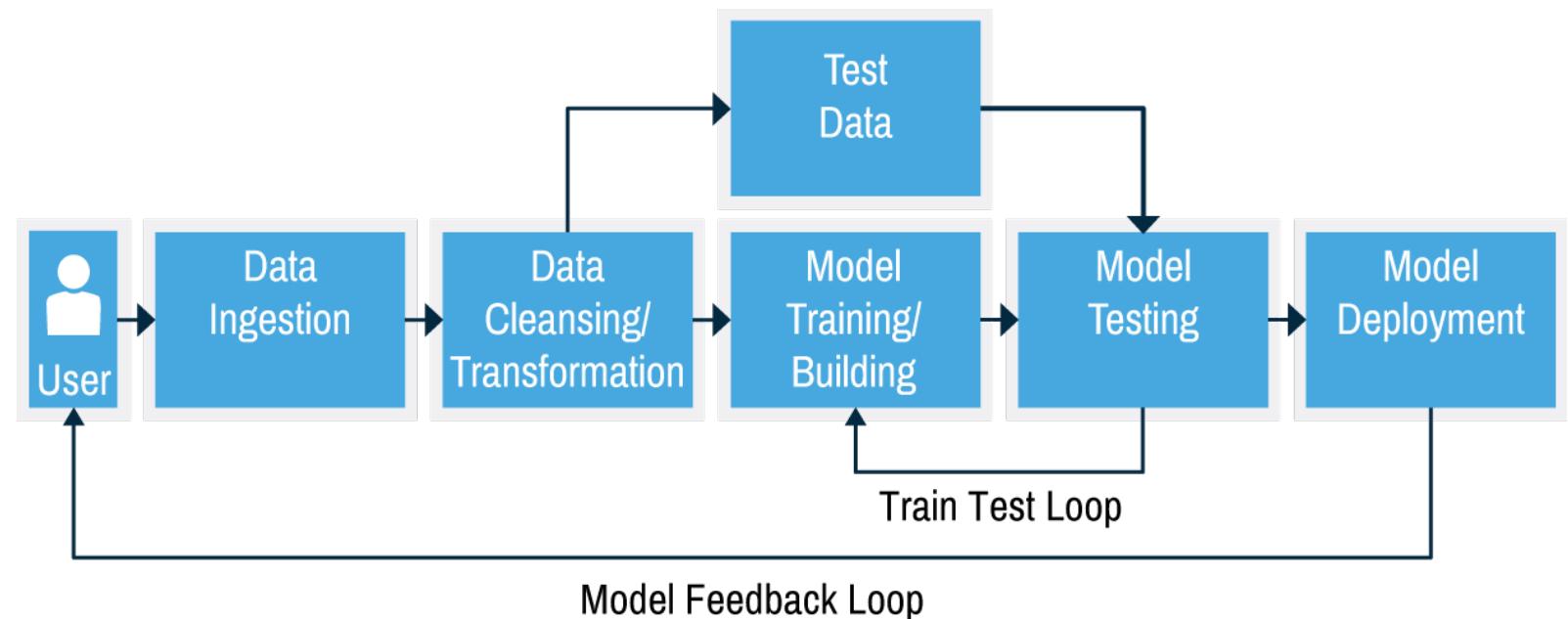
- Diagnosis
- Drug Trial
- Chip Design

# The cost (and benefit?) of BigData Subsample to RAM (which can be 512gb)

# Cornerstones of this Course

- Good software engineering practices
- Problem definition and success measures
- Feature engineering and data cleaning
- Strength and weaknesses of different algorithms
- Model selection best practices

# The Machine Learning Work-Flow



Taken from MAPR <https://www.mapr.com/ebooks/spark/08-recommendation-engine-spark.html>



# General coding guidelines

Programs must be written for people to read, and only incidentally for machines to execute.

Harold Abelson (wizard book)

Everyone knows that debugging is twice as hard as writing a program in the first place. So if you're as clever as you can be when you write it, how will you ever debug it?

Brian Kernighan

- Don't be clever!
- Make it readable!
- Future you is the most likely person to try to understand your code.

- Don't be clever!
- Make it readable!
- Future you is the most likely person to try to understand your code.
- Avoid writing code.

# Python basics

# Why Python?

- General purpose language
- Great libraries
- Easy to learn / use
- Contenders: R (Scala? Julia?)

# The two language problem

Python is sloooow...

- Numpy: C
- Scipy: C, fortran
- Pandas: Cython, Python
- Scikit-learn: Cython, Python

- CPvthon: C

52 / 60

# Python 2 vs Python 3

- “current” : (2.7), 3.6, 3.7
- Don't use Python 2

# Python ...

Package management:

- don't use system python!
- use Virtual environments
- understand pip (and wheels)
- probably use Conda (and anaconda or conda-forge)

# Python ...

Package management:

- don't use system python!
- use Virtual environments
- understand pip (and wheels)
- probably use Conda (and anaconda or conda-forge)

# Pip and conda and upgrades

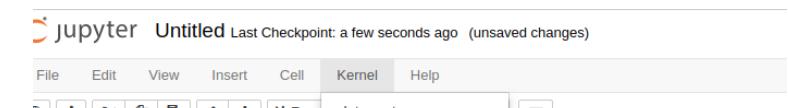
- Pip upgrade works on dependencies (unless you do -no-dep)
- pip has no dependency resolution!
- conda has dependency resolution
- Use conda environments!
- upgrading a conda package with pip (or vice versa) will break stuff!

# Environments and Jupyter Kernels

- Environment != kernels
- Use nb\_conda\_kernels or add environment kernels manually:

```
source activate myenv
python -m ipykernel install --user --name myenv --display-name "Python (myenv)"
source activate other-env
python -m ipykernel install --user --name other-env --display-name "Python (other-env)"
```

- <https://jakevdp.github.io/blog/2017/12/05/installing-python-packages-from-jupyter/>





# Dynamically typed, interpreted

- Invalid syntax lying around
- Code is less self-documenting

# Editors

- Flake8 / pyflake
- Scripted / weak typing: Have a syntax checker!
- write pep8 (according to the standard, not the tool)
- use autopep8 if you have code lying around

# Questions ?

