

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

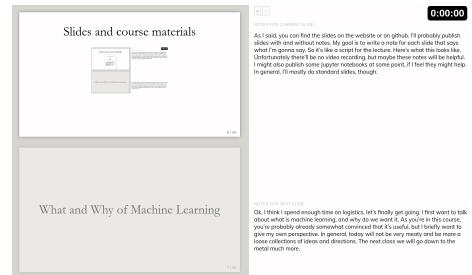
Office Hours

- 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



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

Lecture Recordings

Plagiarism and Code copying

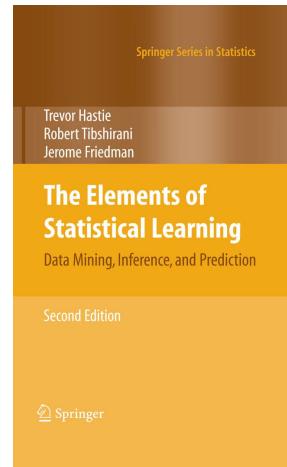
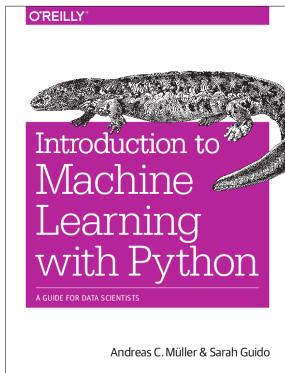
- Homeworks are checked for plagiarism
- Copied code will result in 0 points for all involved
- Copying from my slides or online sources (Stack overflow, tutorials, etc.) is fine.

Scikit-learn Development



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

Books

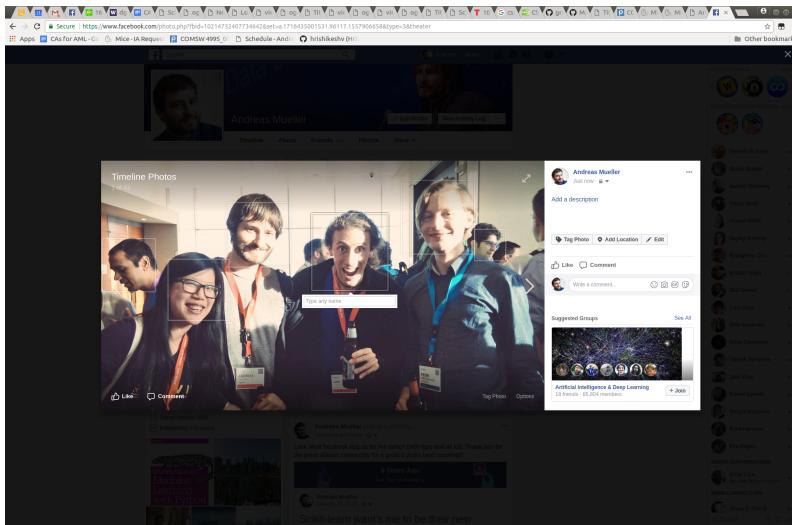


<http://ciml.info/>

What and Why of Machine Learning

What is machine learning?

The screenshot shows a Facebook news feed for user Andreas Mueller. On the left, there's a sidebar with links like News Feed, Messenger, and various Explore options. The main feed displays a sponsored post from 'untapt' with the headline: "Accomplish more in 2017. See this year's most coveted fintech engineering jobs." Below it, another post from Juliana Vanderlee says: "Juliana Vanderlee is 😊 watching Tom Price's confirmation hearing." There are also sections for trending topics and people you may know.



Search for people, places and things

Andreas Mueller added 8 new photos.

June 6 at 5:07pm

Bye bye hong kong

Like · Comment · Share

George Pineda, Cesar Hernandez, Michael Cheung and 6 others like this.

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

June 6 at 7:39pm · Like

Write a comment...

SPONSORED See All

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

Globaler Chauffeurservice
blacklane.com
Fahren Sie eine Klasse besser zu günstigen Preisen – Blacklane ist weltweit verfügbar!

28,723 people like this

Sansouci...
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 Tarif 😊

English (US) · Privacy · Terms · Cookies · More

The screenshot shows the Amazon search results for the query "machine learning". The top navigation bar includes links for Account & Lists, Today's Deals, Gift Cards, Sell, Help, Go, and a search bar. A Father's Day banner is visible on the right. The main search results display five books:

- Practical Machine Learning: Innovations in Recommendation** by Ted Dunning and Ellen Friedman (Apr 17, 2014) - \$0.00 Kindle Edition.
- Machine Learning: The Art and Science of Algorithms that Make Sense of Data** by Peter Flach (Nov 12, 2012) - \$35.00 Kindle Edition.
- Understanding Machine Learning: From Theory to Algorithms** by Shai Shalev-Shwartz and Shai Ben-David (May 19, 2014) - \$34.35 Kindle Edition.
- Learning From Data** by Yaser S. Abu-Mostafa, Malik Magdon-Ismail, and Hsuan-Tien Lin (Mar 27, 2012) - \$30.00 new (9 offers), \$40.00 used (13 offers).

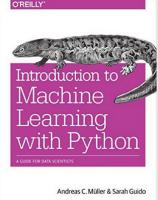
On the left sidebar, there are filters for Department (e.g., Books, Kindle Store), Refine by (e.g., Kindle Store, Book Series, Book Language, Book Format), and Avg Customer Review (e.g., 4.5+ stars). A "See all 21,049 results for 'machine learning'" link is also present.

Introduction to Machine Learning with Python: A Guide for Data Scientists 1st Edition

by Andreas C. Müller (Author), Sarah Guido (Author)

★☆☆☆☆ 9 customer reviews

[Look inside](#)



Kindle \$24.99 Paperback \$41.23 Other Sellers from \$26.53

Buy new In Stock. Ships from and sold by Amazon.com. Gift-wrap available.

Want it tomorrow, Jan. 19? Order within 5 hrs 29 mins and choose One-Day Shipping at checkout. Details

Prime \$41.23 List Price: \$49.99 Save: \$8.76 (18%) 36 New from \$26.53

Qty: 1 Add to Cart Turn on 1-Click ordering

Ship to: Andreas C Mueller- New York - 10009

More Buying Choices 51 used & new from \$26.53 See All Buying Options

ISBN-13: 978-1449369415 ISBN-10: 1449369413 Why is ISBN important?

Have one to sell? Sell on Amazon

Add to List Share Email Facebook Twitter

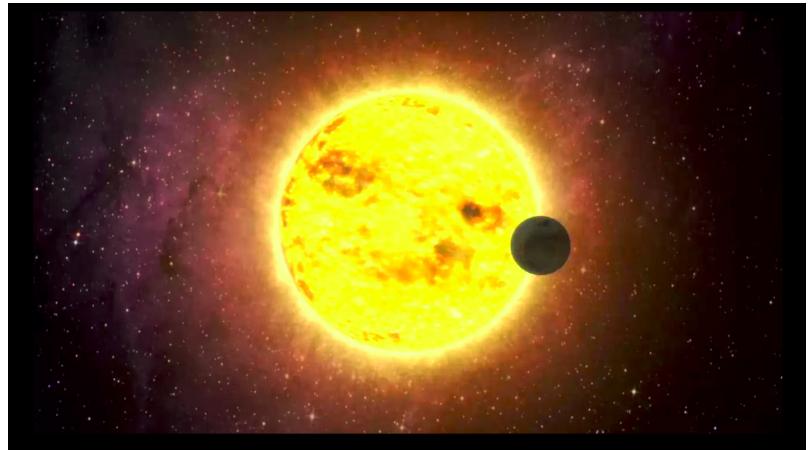
Save up to 90% on textbooks Shop now

Frequently Bought Together

 Total price: \$108.29

Add all three to Cart Add all three to List

Science!



Types of Machine Learning

Types of Machine Learning

- Supervised
- Unsupervised
- Reinforcement

Supervised Learning

$(x_i, y_i) \propto p(x, y)$ 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:

$$f(x_i) \approx y_i,$$

$$f(x) \approx y$$

Examples of Supervised Learning

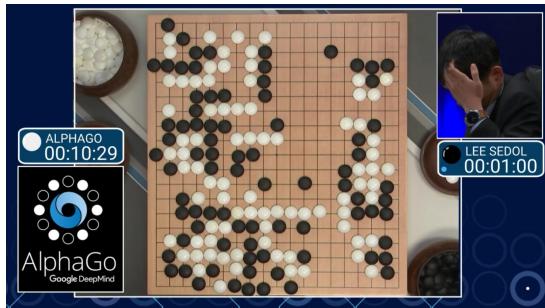
- spam detection
- medical diagnosis
- add click prediction

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

These recommendations are based on items you own and more.

view All | New Releases | Coming Soon

1.  **R for Data Science: Import, Tidy, Transform, Visualize, and Model Data**
by Hadley Wickham (January 5, 2017)
Average Customer Review: ★★★★☆ (4)
In Stock
List Price: \$39.99
Price: \$32.91
28 used & new from \$28.91
 I own it Not interested Rate this item [Add to Cart](#) [Add to Wish List](#)

Recommended because you purchased Data Science from Scratch: First Principles with Python and more ([See the](#))

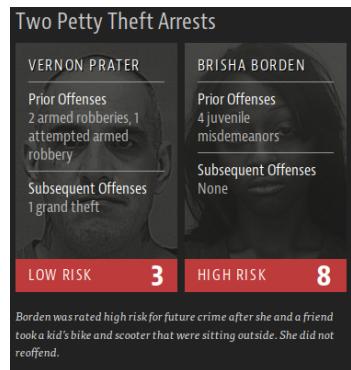
2.  **Storytelling with Data: A Data Visualization Guide for Business Professionals**
by Cole Nussbaumer Knaflic (November 2, 2015)
Average Customer Review: ★★★★☆ (137)
In Stock
List Price: \$39.95
Price: \$29.71
94 used & new from \$20.00
 I own it Not interested Rate this item [Add to Cart](#) [Add to Wish List](#)

Recommended because you purchased Mindset: The New Psychology of Success and more ([See the](#))

3.  **Deep Learning (Adaptive Computation and Machine Learning series)**
by Ian Goodfellow (November 18, 2016)
Average Customer Review: ★★★★☆ (20)
In Stock
List Price: \$69.00
Price: \$68.34
18 used & new from \$68.34
 I own it Not interested Rate this item [Add to Cart](#) [Add to Wish List](#)

Recommended because you purchased Data Science from Scratch: First Principles with Python and more ([See the](#))

Sidebar: Ethical Considerations



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

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

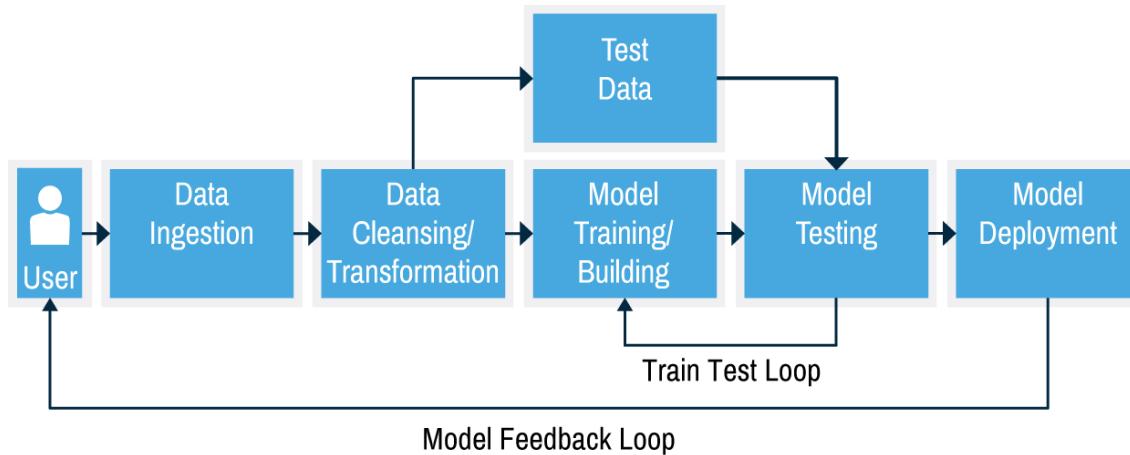
- Diagnosis
- Drug Trial
- Chip Design

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
- CPython: C

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

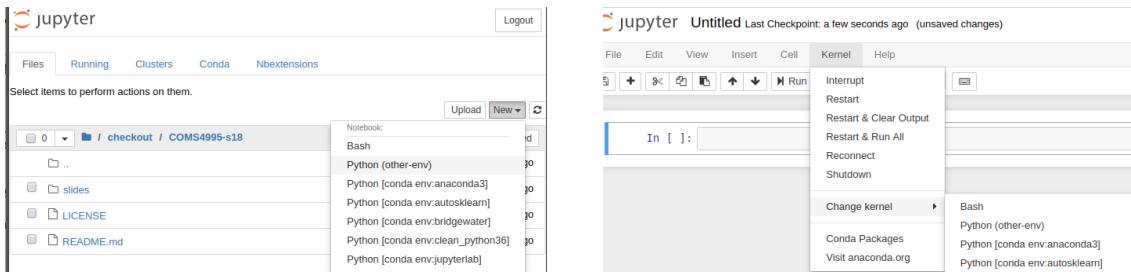
- Pip upgrade works on dependencies (unless you do -no-dep)
- pip has no dependency resolution!
- conda has dependency resolution
- Use conda environments!
- upgrading a conda package with pip (or vice versa) will break stuff!

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 ?