

# What is Machine Learning (ML) ?

- prediction



# What is Machine Learning (ML) ?

- informed prediction
  - prediction is better than random guess
  - method: learn from existing data
  - goal: Generalization. Predicting on new, unseen data



# Where is ML used ?

**Everywhere !!**

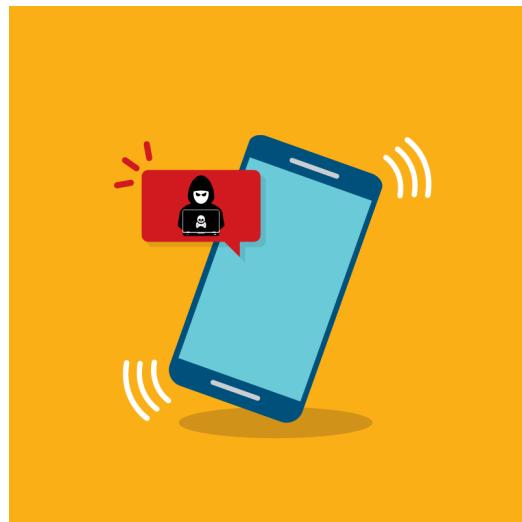
- targeted advertising
  - Why does Facebook seem to know what I'm thinking ?



# Where is ML used ?

Spam detection

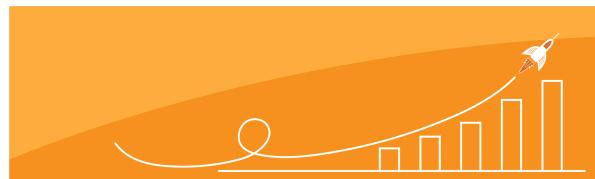
- You are a winner !</ul>



# Where is ML used ?

Forecasting

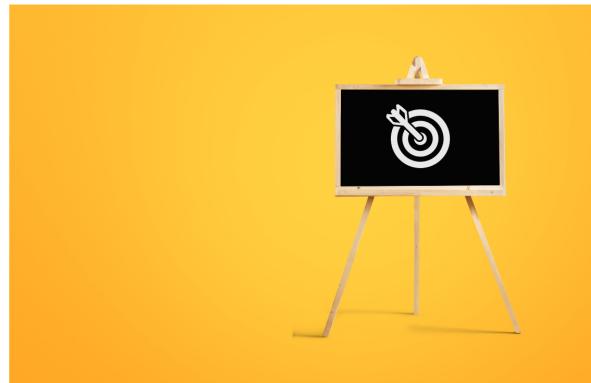
- Sales
- Logistics
- Where's my Uber ?



# Where is ML used ?

Anomaly detection

- Credit card fraud



# Uses in Finance

- Model prices, risk
  - hedging
- Trading signals
- forecast sales
- Predict defaults, prepayments



# Not just numeric data !

- Images
  - Satellite:
    - Counting cars in a parking lot to forecast sales
    - How full is that oil tank ?



- Did the CFO really mean what he said ? facial signals for confidence/evasiveness
- Text
  - Twitter sentiment as a signal ?
  - SEC filings
  - Derive industry groups by clustering press releases

# What you need to succeed: An inquiring mind



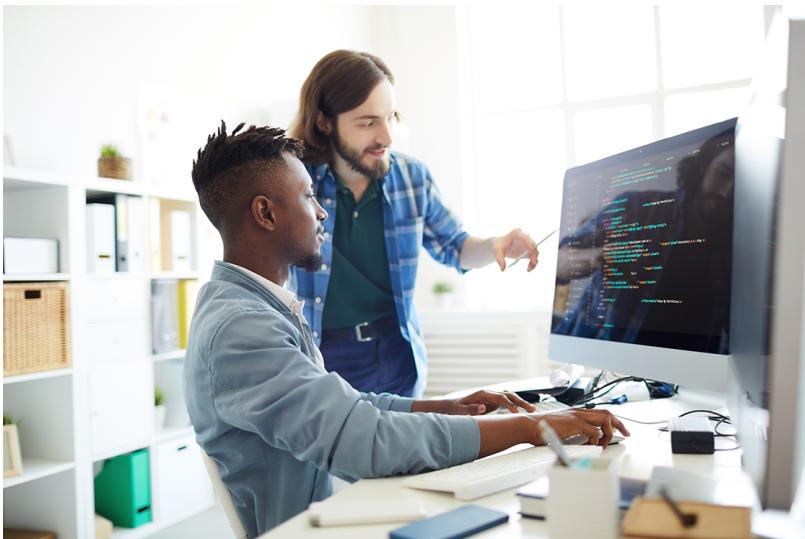
- Approach this topic like a Scientist
  - Find a problem, gather data, formulate a hypothesis, test.
  - Repeat.

# What you need to succeed: Technical skills

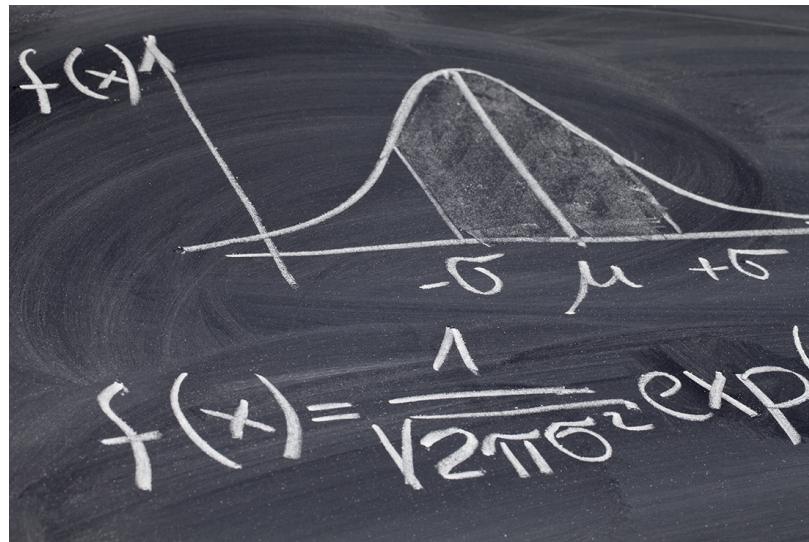


You are engineers !

# What you need to succeed: Solid programming skills



# What you need to succeed: Some math/statistics



To be a successful data scientist, you need to understand the machinery.

It is not enough to know an API.

# What you need to succeed Self-motivation and energy



- Willingness to pick up tools/skills outside of lectures
- You are engineers, nothing is too hard !

# Technical prerequisites

- Python
  - Object oriented (OO) Python
  - Numpy
  - Pandas
  - Matplotlib
- Some statistics (e.g., regression)
- Some math
  - comfort with Matrix/Vector notation

# Textbooks

## Python Data Science Handbook

VanderPlas (<https://jakevdp.github.io/PythonDataScienceHandbook/>)

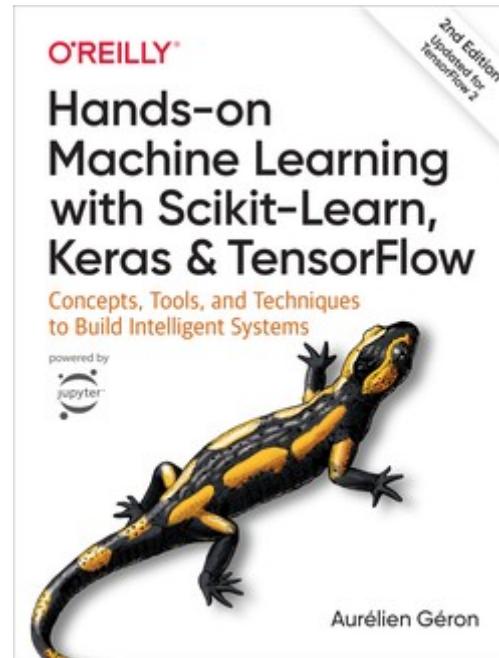
- Online !
- Solid foundation to acquire pre-requisites
  - Jupyter
  - Numpy
  - Pandas
  - Matplotlib
  - Quick view of models



# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2nd Edition

Geron (<https://www.oreilly.com/library/view/hands-on-machine-learning/9781492032632/>).

- Assumes you know the pre-requisites
  - More detailed chapters on various models



# Notebooks

The real learning comes from active "doing" (playing with notebooks) rather than passive "reading".

The course material and both textbooks have code repositories (including notebooks) on GitHub

- The VanderPlas "book" is actually a notebook !

All of the lectures and assignments are presented as notebooks too !

Detailed instructions on obtaining course materials may be found at [Setting up your ML Environment \(Setup\\_ML\\_Environment\\_NYU.ipynb\)](#).

# Machine Learning *using* Scikit-Learn (sklearn), TensorFlow (TF)

sklearn is a popular library for Machine Learning. We will be using for the Classical ML part of the series.

TensorFlow is another popular library that we will use in the second part of the course.

We are learning **Machine Learning**, not sklearn/TensorFlow !

Tools are a means, not an ends

- Goal is to understand ML independent of the toolset
- You can be an expert in sklearn/TensorFlow and still not understand ML

# Teaching method

**Iterative:** visit the problem many times, at increasing levels of focus

- Top-down vs bottom-up
- Motivate: very high level view
  - Know WHAT we are trying to achieve

- Understand: medium level view



- Deep understanding
  - math, statistics

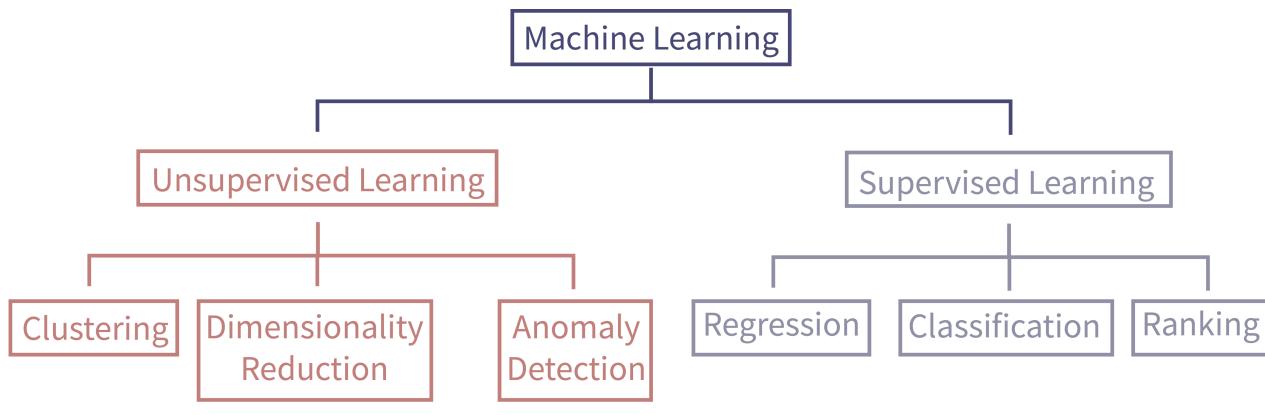


Bonus: advice from a practitioner

# **Dimensions of ML**

## Types of learning

- Supervised
- Unsupervised



</td></tr></table>

## Types of targets

- Continuous
- Discrete/Categorical

## Types of features

- numeric
- image
- text

# Challenges of ML



- You need data to train, often a lot of it
  - Not always easy to get
    - supervised: needs to be labeled
  - Quality issues
  - Is the training data representative of "the real world" for which you are designing ?

- Overfitting and Underfitting
  - Overfit: good training accuracy, poor generalization
  - Underfit: lost opportunity

- Engineering meaningful features is key
  - Data transformations
    - Create features that aid prediction
    - Art and science
  - Deep Learning may view feature engineering as part of the problem, not the solution !

- Testing and validation
  - An honest test uses held-out data
  - Training data is a precious resource; painful to hold some out

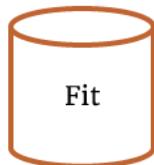
# ML in one slide



Observations:  
features  
features/labels

$$x,^1 x,^2 \dots, x^N$$

$$y,^1 y,^2 \dots, y^N$$



Create a predictor (estimator, hypothesis)  $y^* = h(x^*)$

- Predicts an outcome when presented with features  $x^*$
- Prediction based on the subset of  $x,^1 x,^2 \dots, x^N$
- That looks like  $x^*$



Predict

Scikit Learn

```
from sklearn import mod_class  
model = mod_class.Model()
```

```
model.fit(x,y)
```

```
y_star =  
model.predict(x_star)
```

In [3]: `print("Done")`

Done