



RIGA**CODING**SCHOOL

# Machine Learning with Python Anaconda Data Science Distribution



# VALDIS SAULESPURĒNS



- Izglītība: Maģistra grāds datorzinātnēs
- Pieredze programmēšanā: 20+ gadi
- Pythons + Datu Analīze: 10+ gadus
- Specialitāte: grafu teorija sociālo tīklu analizēšanā
- Hobiji: prāta spēles, riteņbraukšana, šahs

Contact: [valdis.s.coding@gmail.com](mailto:valdis.s.coding@gmail.com)



# Data Lake or Data Swamp?



# Buzzword bingo



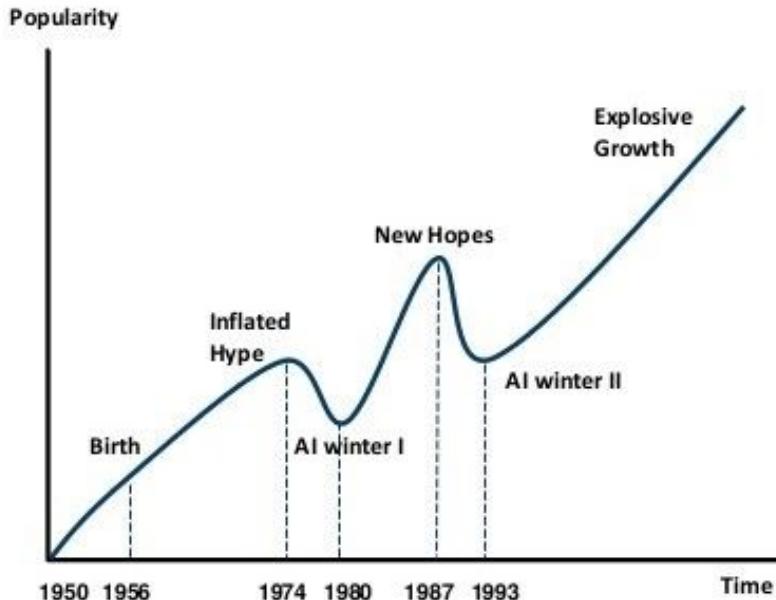
- Big Data
- Data Mining (datizrace)
- Machine Learning – subset of AI
- Data Science – statistics
- Big Data or Pokemon
- <https://pixelastic.github.io/pokemonorbigdata/>

Random Forests	Neural Network	Reinforcement Learning	Supervised Learning	Cognitive Computing
Caffe	Support Vector Machine	Artificial Intelligence	Python	Cloud
Unstructured Data	Bot	DATA SCIENCE BUZZWORD BINGO (free square)	K-means	GPU
Spark	Data Wrangling	Deep Learning	Ensemble	Machine Learning
Keras	Tensorflow	Big Data	Algorithm	Feature Engineering



# Brief History of AI

AI HAS A LONG HISTORY OF BEING “THE NEXT BIG THING”...



## Timeline of AI Development

- **1950s-1960s:** First AI boom - the age of reasoning, prototype AI developed
- **1970s:** AI winter I
- **1980s-1990s:** Second AI boom: the age of Knowledge representation (appearance of expert systems capable of reproducing human decision-making)
- **1990s:** AI winter II
- **1997:** Deep Blue beats Gary Kasparov
- **2006:** University of Toronto develops Deep Learning
- **2011:** IBM's Watson won Jeopardy
- **2016:** Go software based on Deep Learning beats world's champions

# What is machine learning?

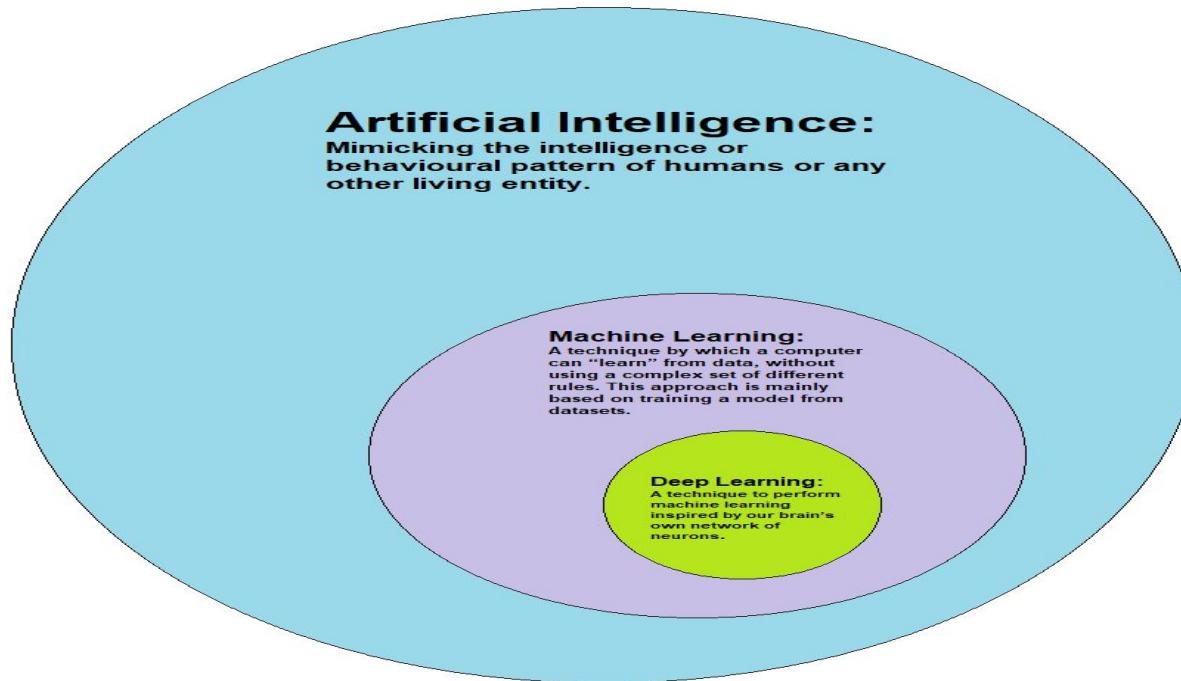


“Machine learning is the science of getting computers to act without being explicitly programmed.” – Stanford ML course

“The field of Machine Learning seeks to answer the question “How can we build computer systems that automatically improve with experience, and what are the fundamental laws that govern all learning processes?” – [Carnegie Mellon University](#)

“Machine learning can’t get something from nothing...what it does is get more from less.” – Dr. Pedro Domingo, University of Washington

# AI, ML, Deep Learning



# Brief History of Data Analysis



- What is this ? :)



# Brief History of Data Analysis



- ~ 18,000BC – Uganda, Ishango Bone
- ~ 2400BC – Babylon abacus, libraries
- 300BC – 48AD – Library of Alexandria



# Brief History of Data Analysis



- How about this modern recreation of a 2000 years old device?



# Brief History of Data Analysis



- ~ 100-200AD  
**Antikythera Mechanism**

Predicting:

- Astrology
- Astronomy
- Olympics
- Calendar



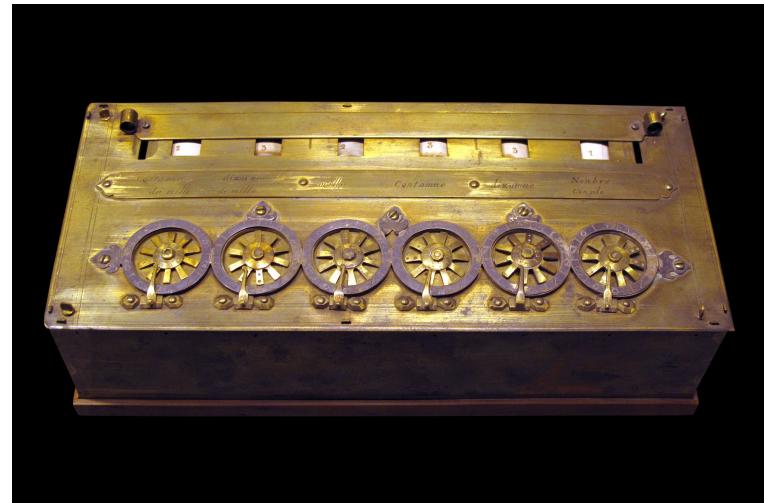
# Brief History of Data Analysis



- **1642 Blaise Pascal's  
Pascaline**

Performs:

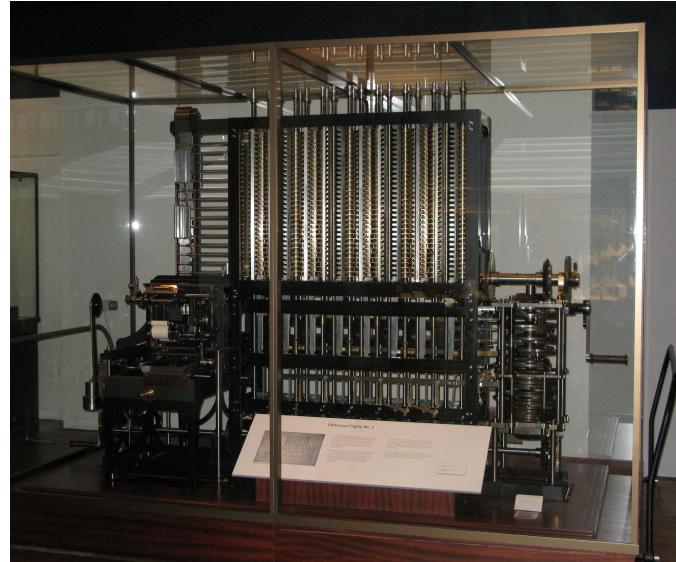
- Addition
- Subtraction
- Multiply/Divide using Add/Sub
- 1649 Royal Patent by Louis XIV



# Brief History of Data Analysis



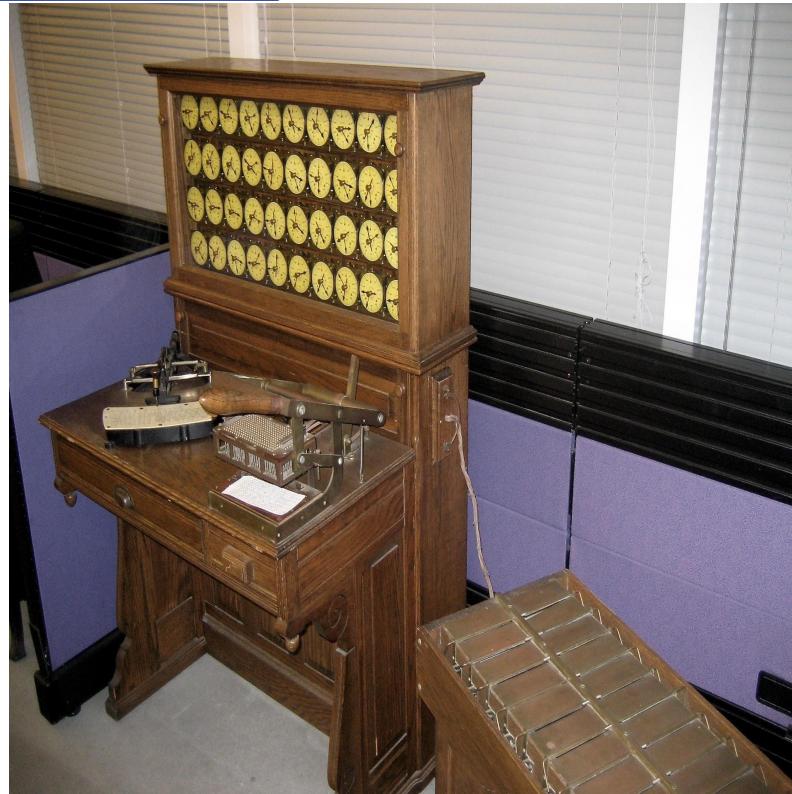
- **Charles Babbage Difference Engine**
  - Ada Lovelace - first programmer?
- Performs:
  - Arithmetic
  - Derivation
  - Power Series
  - Curve Fitting



# Brief History of Data Analysis II



- 1663 – London, J.Graunt mortality analysis
- 1865 – banker H. Furnese business intelligence
- 1880-90 US Census Hollerith Machine -> IBM
- 1928 – F. Pfleumer magnetic tape invention



# Brief History of Data Analysis



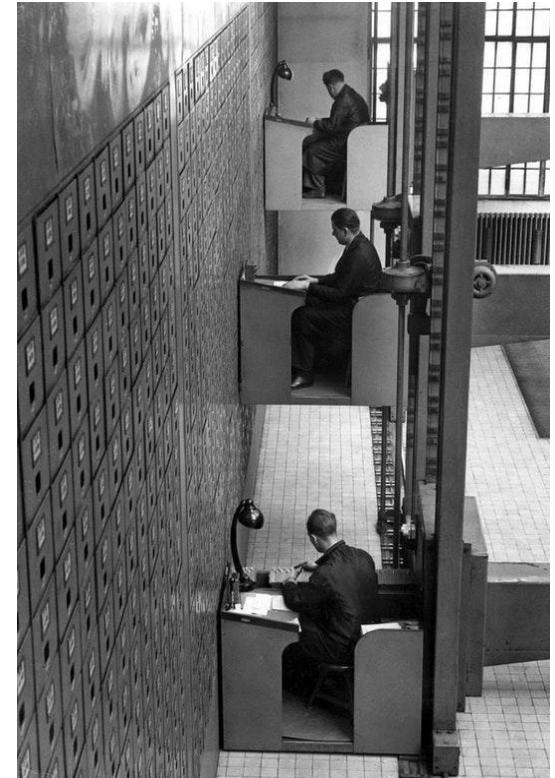
- 1940s - First General Purpose Electronic Computer ENIAC (Zuse mechanical)
- Turing complete
- von Neumann architecture
- most computers work the same today



# Brief History of Data Analysis



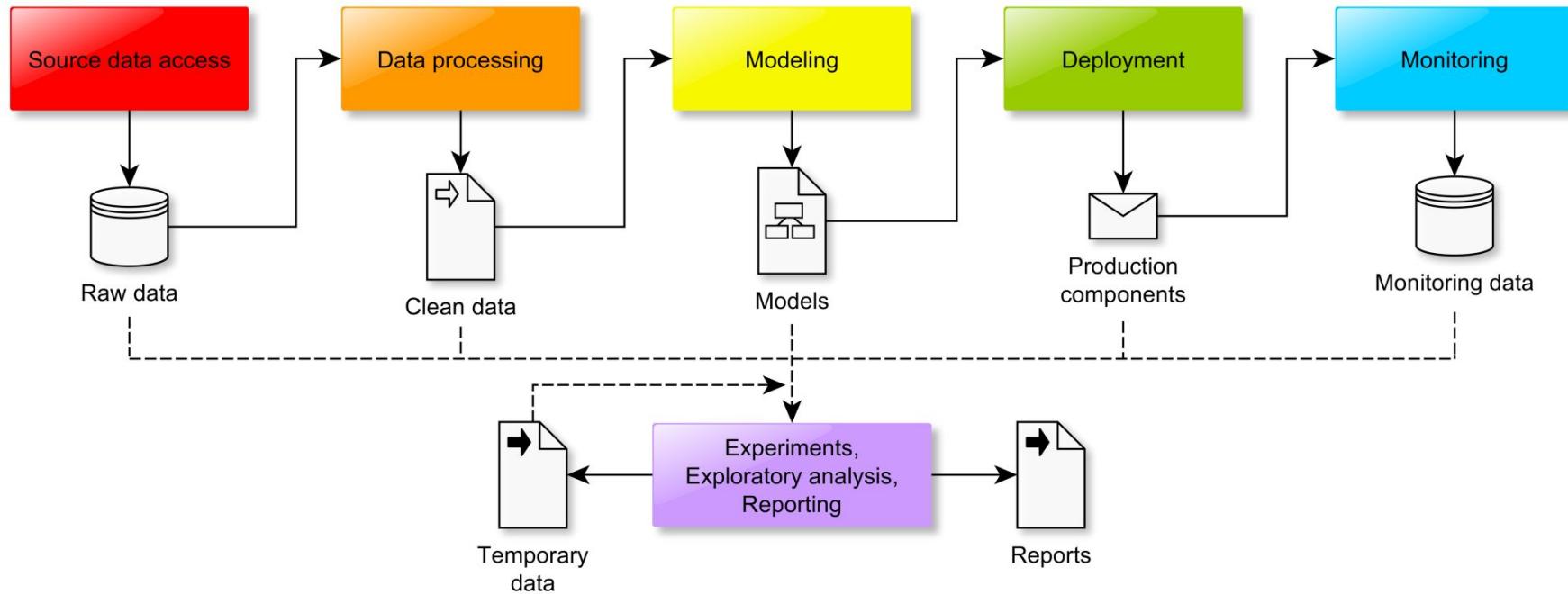
- 1950s - Flat Files
- 1958 – IBM's Luhn defines Business Intelligence
- 1960s - CODASYL
- 1970s – Codd's relational DBs -> SQL
- 1980s – Data Warehouses / Marts
- 2000s – Big Data / noSQL DBs



# BIG DATA LANDSCAPE 2017



# Full Analysis Framework



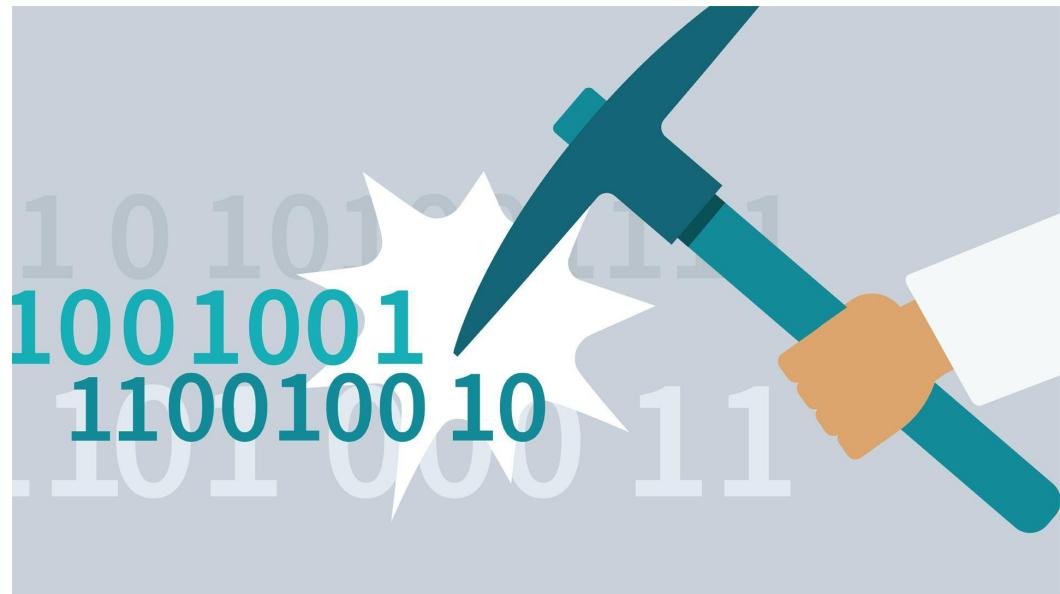
# Data Analysis Sandbox



# Data Mining



- Anomalies
- Classification
- Clusters
- Dimension Reduction
- Regression
- Relationship finding
- Summarization / Visualization



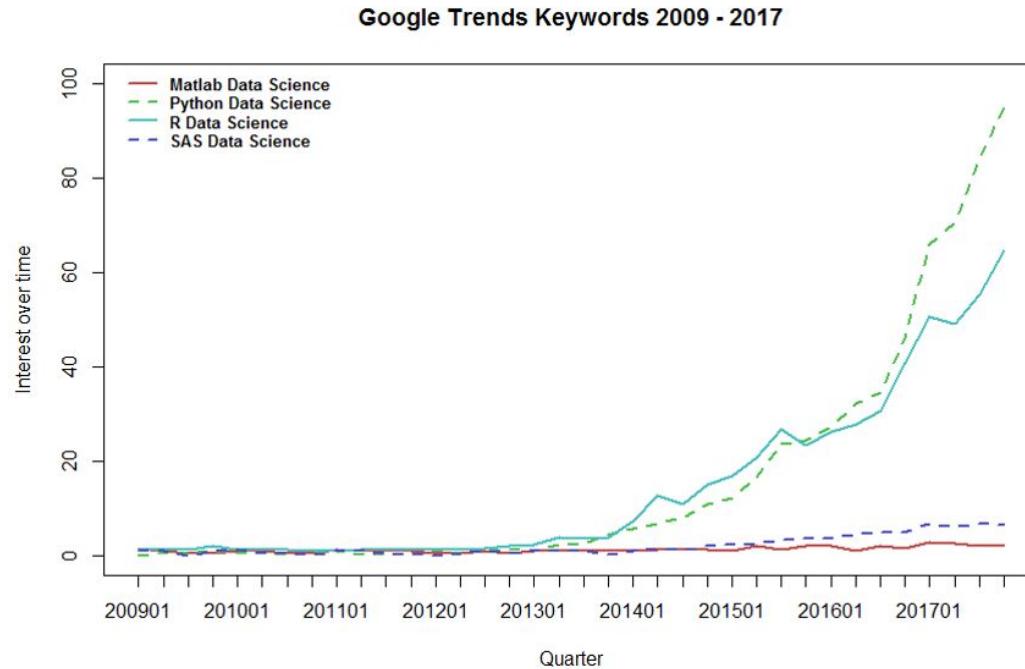
# Building a Pipeline



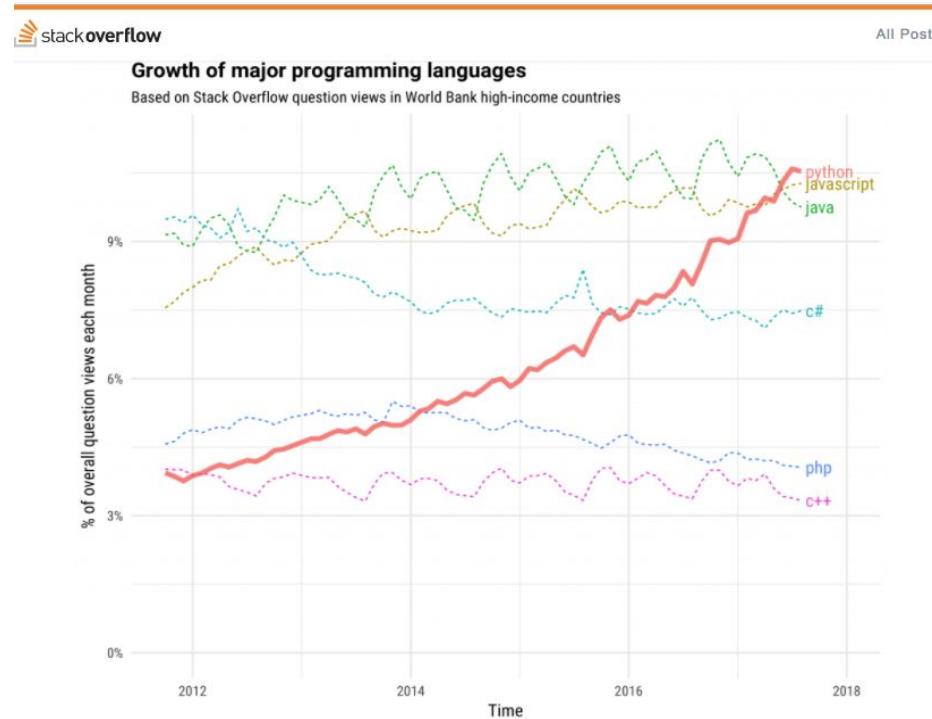
- Data Security
- Cleanup
- Organising Database
- Analysis
- Visualization – Dashboard
- Emphasis on Analysis less on Infrastructure



# Why Python?



# Why Python? Part 2



# Why Python?



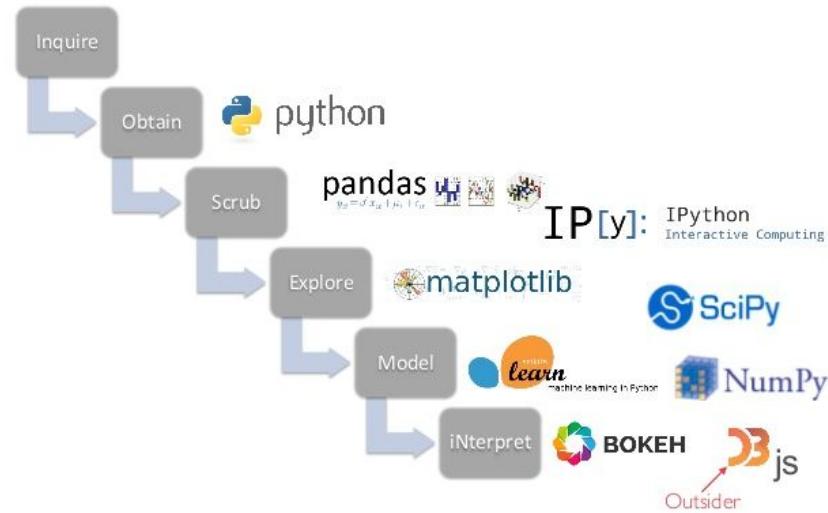
- Readable
- Interactive
- Interpreted
- Dynamic
- Object-oriented
- Portable
- High level
- Extensible in C++ & C



# Python Ecosystem



PYTHON IS IOSEMN



# Anaconda Distribution



## ANACONDA POWERS THE MODERN ANALYTICS STACK

APP	Notebooks	Embeddable Dashboards	Data Services	Visual Apps
VIZ	Plots	Interactive Viz	Big Data	Maps & GIS
STORYBOARD	Notebooks	Interactive Exploration	Visual Programming	Data IDEs
ANALYTICS	Data Prep	Stats	ML & Ensembles	Deep Learning
	Text & NLP	Video/Image/Audio Mining	Graph & Network	Simulation & Optimization
DATA	Hadoop & Hive	Spark	NoSQL	DW & SOL
HARDWARE	Servers	Clusters	GPUs & High End Workstations	Files & Web Services

# scikit-learn machine learning library



- Tools for data mining and data analysis
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license
- 
- Source: <https://github.com/scikit-learn/scikit-learn>



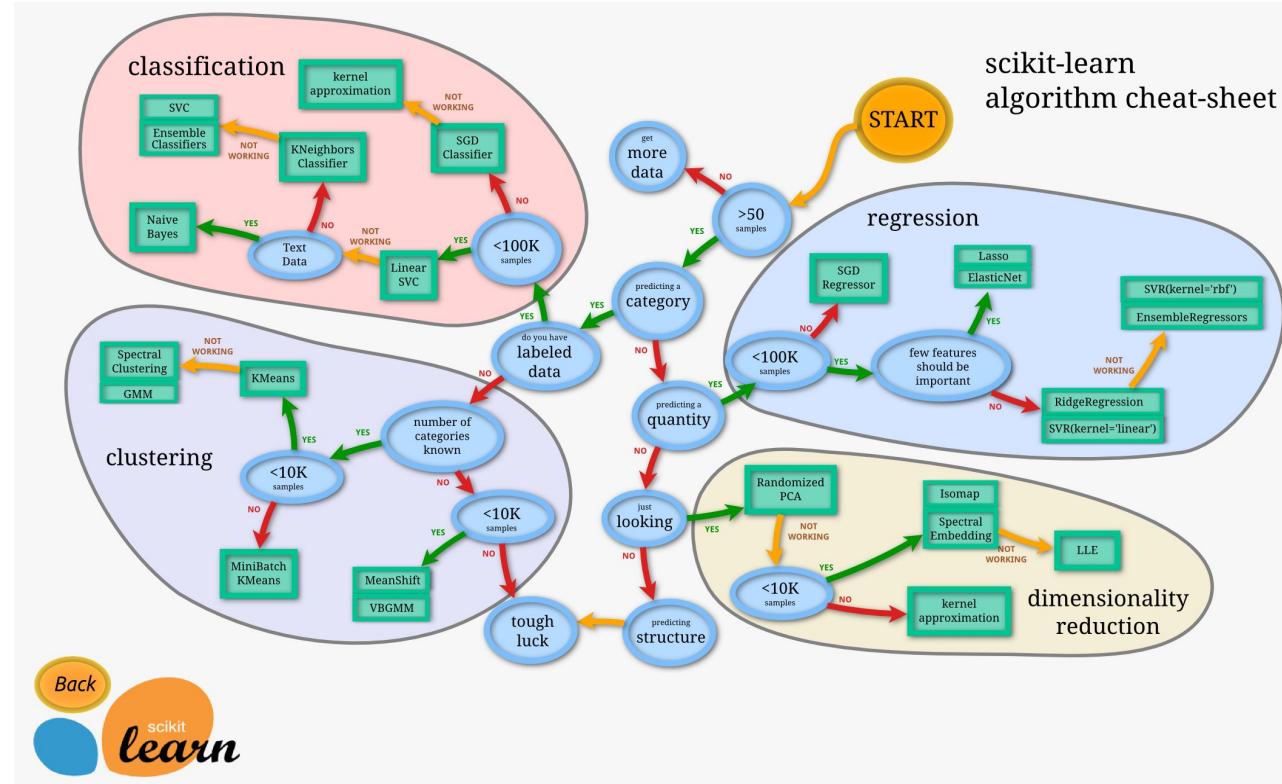
# scikit-learn machine learning library



- Classification
- Regression
- Clustering
- Dimensionality Reduction
- Model Selection
- Preprocessing



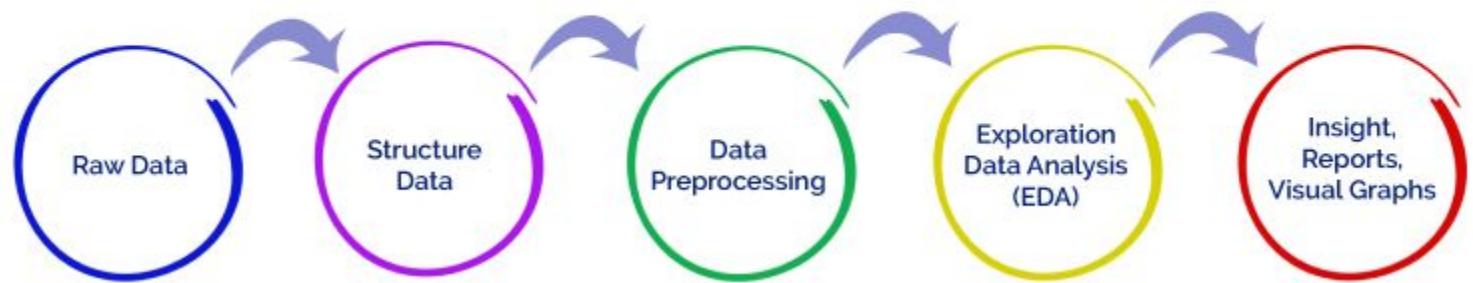
# scikit-learn machine learning library



# Data Cleaning / Wrangling



## Data Preparation



# Tidy Data



country	year	cases	population
Afghanistan	1990	745	1637071
Afghanistan	2000	1666	2095360
Brazil	1999	37737	17206362
Brazil	2000	80488	17404898
China	1999	21258	127215272
China	2000	21666	128042583

variables

country	year	cases	population
Afghanistan	1999	745	1637071
Afghanistan	2000	1666	2095360
Brazil	1999	37737	17206362
Brazil	2000	80488	17404898
China	1999	21258	127215272
China	2000	21666	128042583

observations

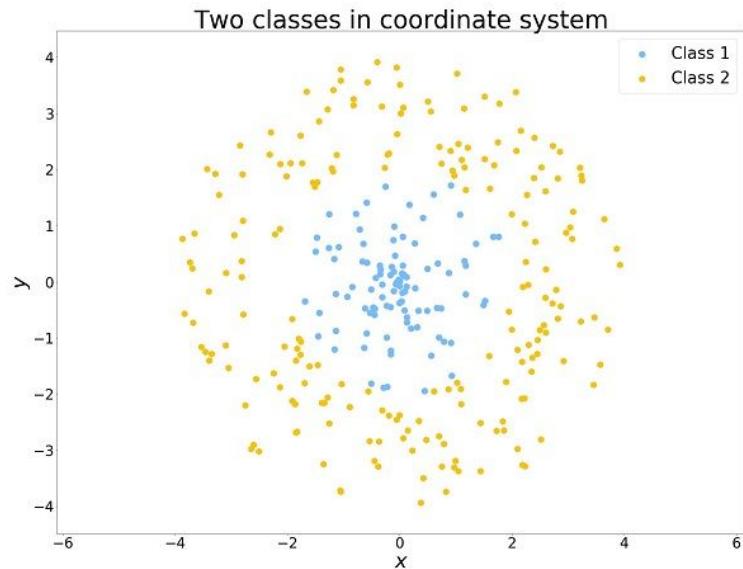
country	year	cases	population
Afghanistan	1999	745	1637071
Afghanistan	2000	1666	2095360
Brazil	1999	37737	17206362
Brazil	2000	80488	17404898
China	1999	21258	127215272
China	2000	21666	128042583

values

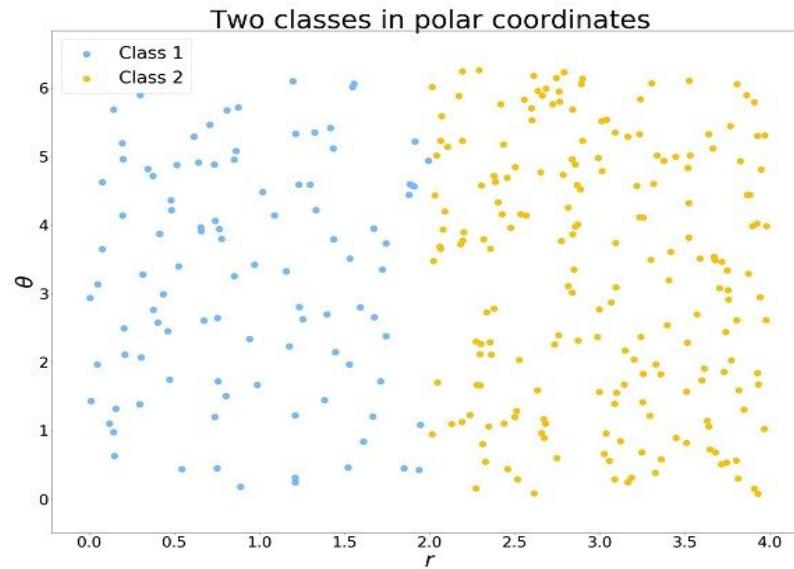
1. Each variable -> column.
2. Each observation -> row.
3. Each type of obseration -> table.

<https://vita.had.co.nz/papers/tidy-data.pdf>

# Feature Engineering



# Feature Engineering



$$r = \sqrt{x^2 + y^2} \quad \theta = \arctan \frac{y}{x}$$

# Feature Engineering

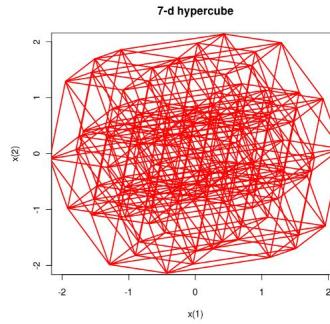
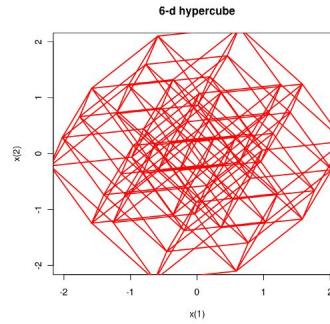
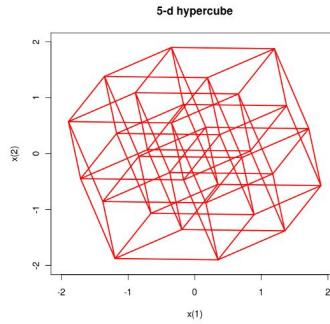
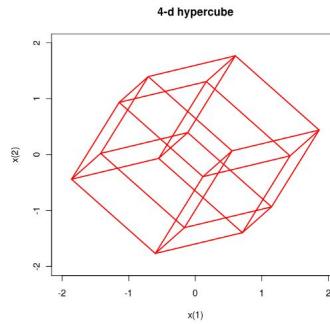
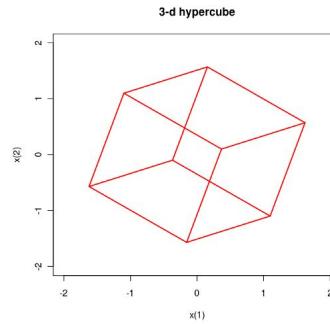
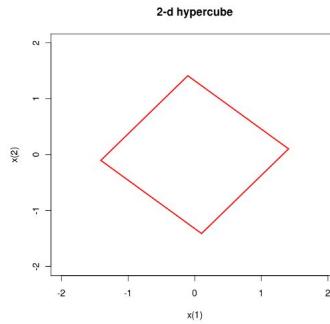


- Categorical Data -> Category Codes
- Categorical Data -> One-Hot Encoding (Only one True value)

The diagram illustrates the process of feature engineering, specifically one-hot encoding. On the left, a table shows a column of categorical data: Color, with entries Red, Red, Yellow, Green, and Yellow. A large yellow arrow points from this table to the right, indicating the transformation. On the right, a second table shows the resulting one-hot encoded matrix. The columns are labeled Red, Yellow, and Green. The rows correspond to the categories in the first table. The values in the matrix are binary (0 or 1), where a '1' indicates the presence of the category and '0' indicates its absence. For example, the first two rows (Red) have a '1' in the 'Red' column and '0's in the other two columns. The third row (Yellow) has a '1' in the 'Yellow' column and '0's in the others. The fourth row (Green) has a '0' in the 'Red' column and a '1' in the 'Yellow' column, with '0's in the others. The fifth row (Yellow) has a '0' in the 'Red' column and a '0' in the 'Yellow' column, with a '1' in the 'Green' column.

Color	Red	Yellow	Green
Red	1	0	0
Red	1	0	0
Yellow	0	1	0
Green	0	0	1
Yellow	0	0	1

# PCA and Curse of Dimensionality

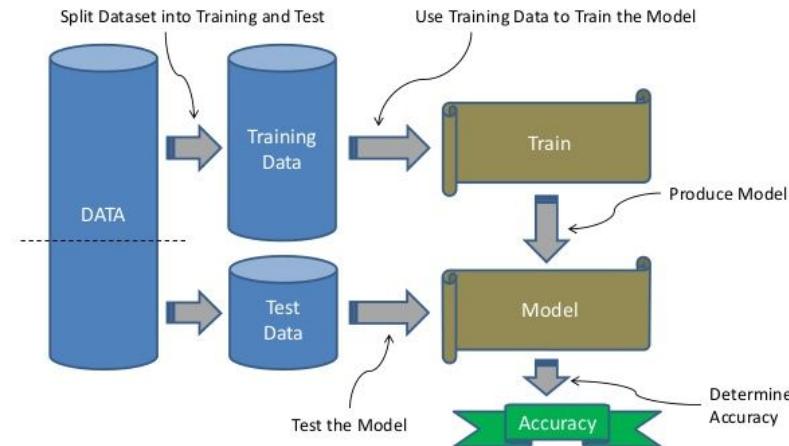


# Training Data vs Testing Data

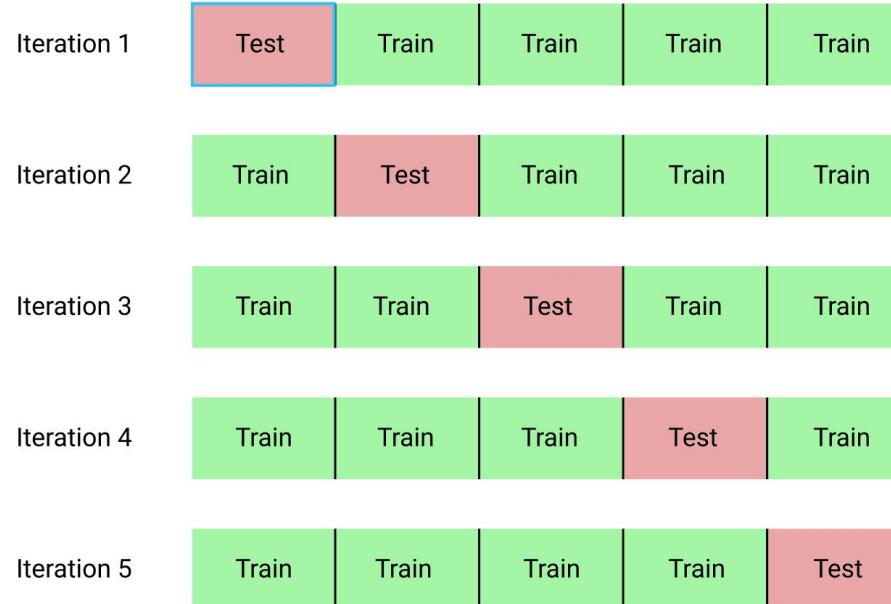


## It's About Training

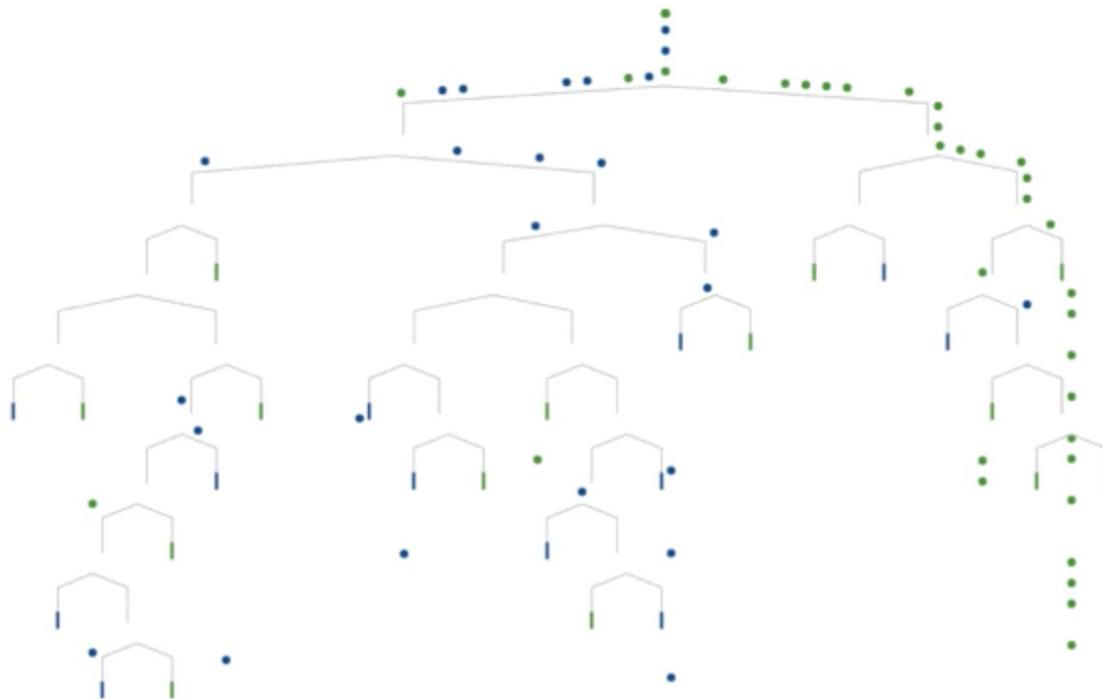
Machine Learning is about using data to train a model



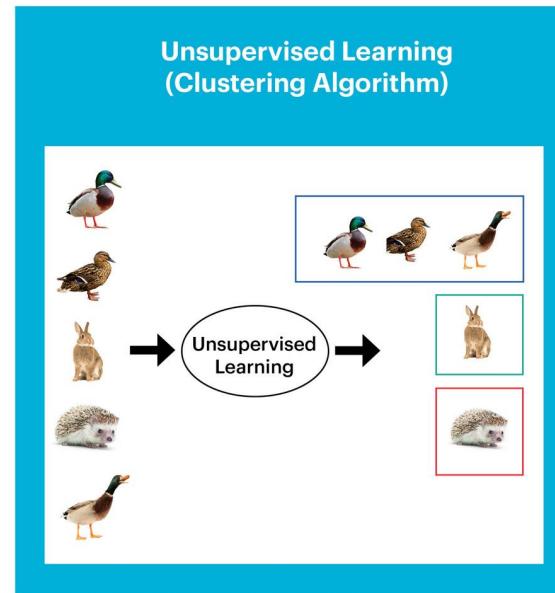
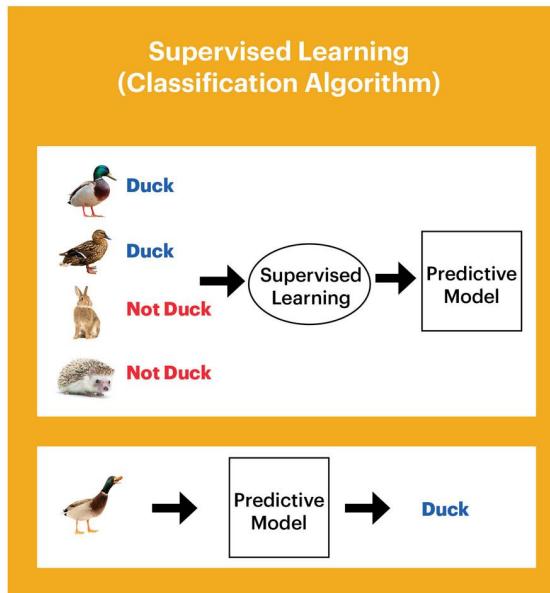
# Cross Validation



# Decision Trees Classifiers

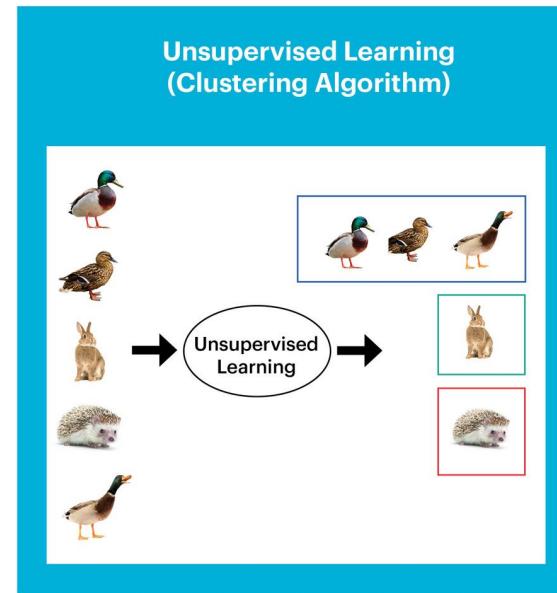
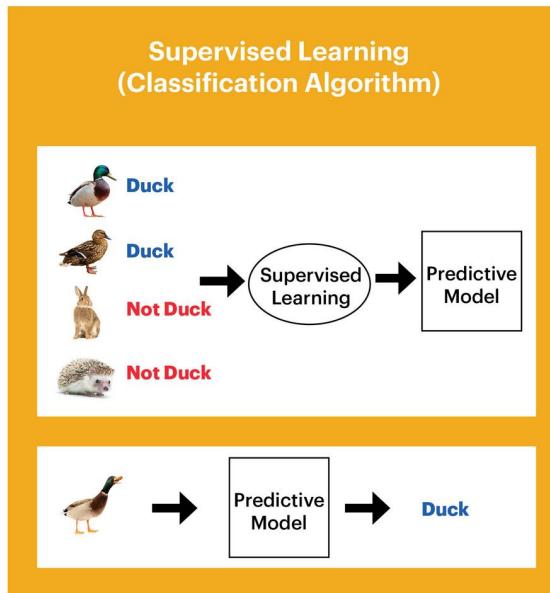


# Supervised vs Unsupervised



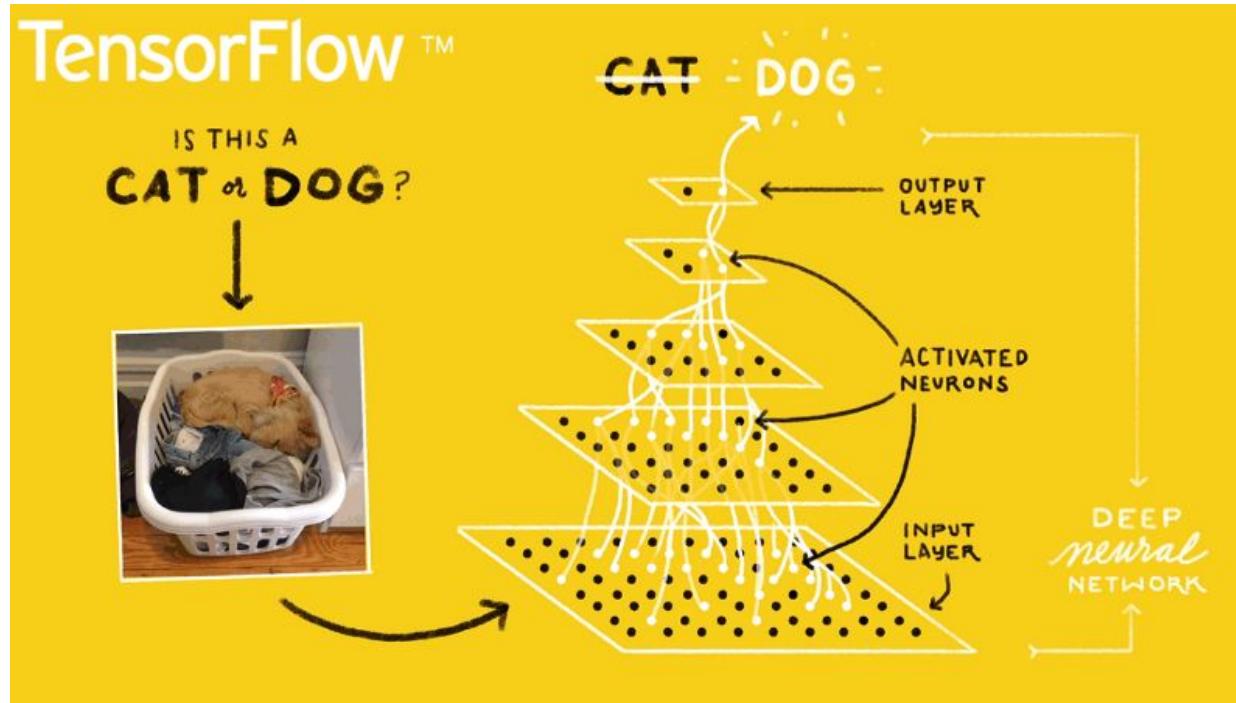
Western Digital.

# Supervised vs Unsupervised



Western Digital.

# Deep Learning with Keras, TensorFlow



# Deep Learning with Keras, TensorFlow



TensorFlow - The core open source library to help you develop and train ML models

<https://www.tensorflow.org/>

Changes Often! :(



Keras - Wrapper Library for TensorFlow(and some others)

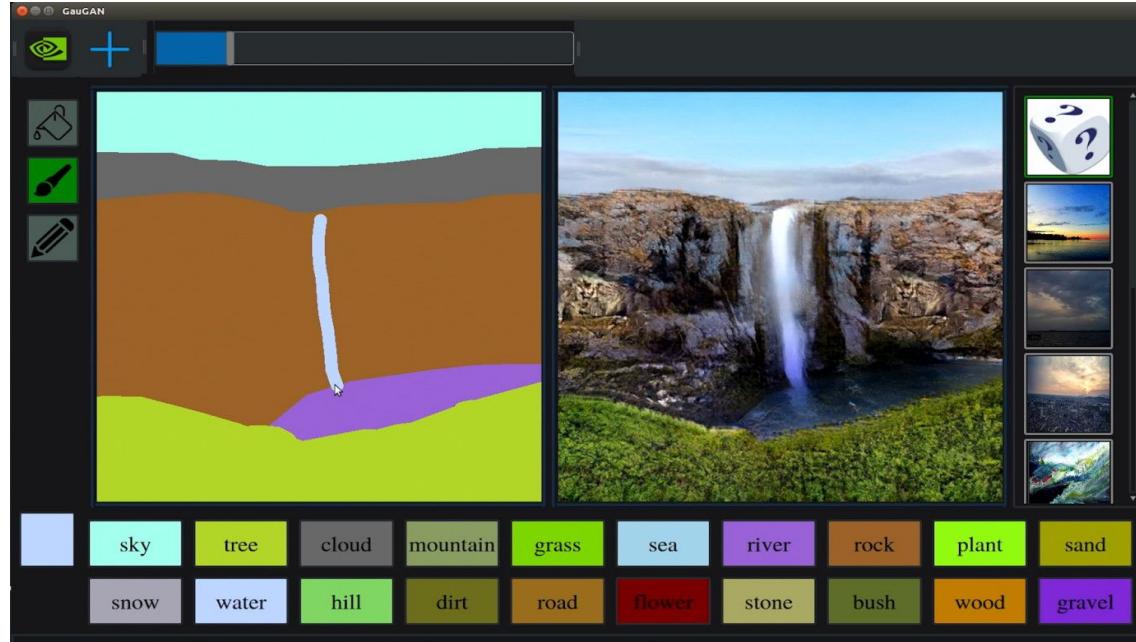
<https://keras.io/>



Example Image Classification Lab:

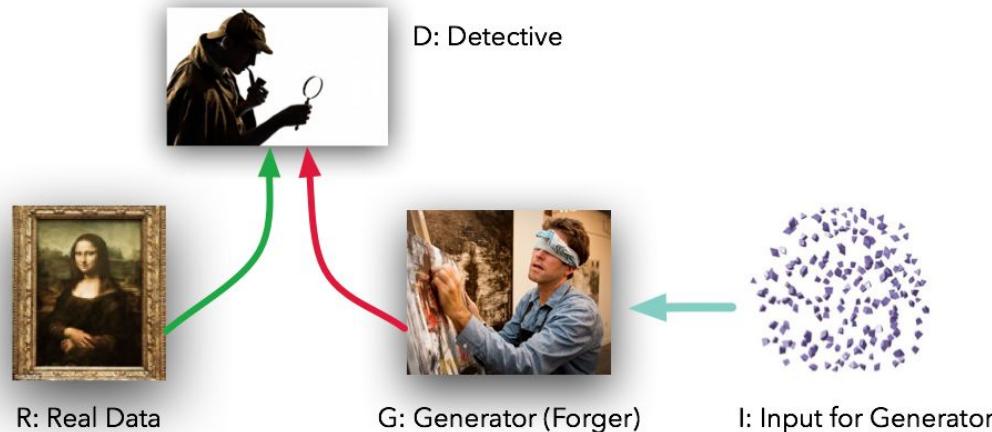
[http://colab.research.google.com/github/ValRCS/RCS\\_Data\\_Analysis\\_Python\\_2019/  
blob/master/Keras\\_TensorFlow\\_Image\\_Recognition/keras\\_image\\_recognition\\_classifier\\_in\\_class.ipynb](http://colab.research.google.com/github/ValRCS/RCS_Data_Analysis_Python_2019/blob/master/Keras_TensorFlow_Image_Recognition/keras_image_recognition_classifier_in_class.ipynb)

# New! GAN (Generative Adversarial Networks)



<https://techcrunch.com/2019/03/18/nvidia-ai-turns-sketches-into-photorealistic-landscapes-in-seconds/>

# GAN - in 50 lines of Python



<https://pytorch.org/>

<https://medium.com/@devnag/generative-adversarial-networks-gans-in-50-lines-of-code-pytorch-e81b79659e3f>

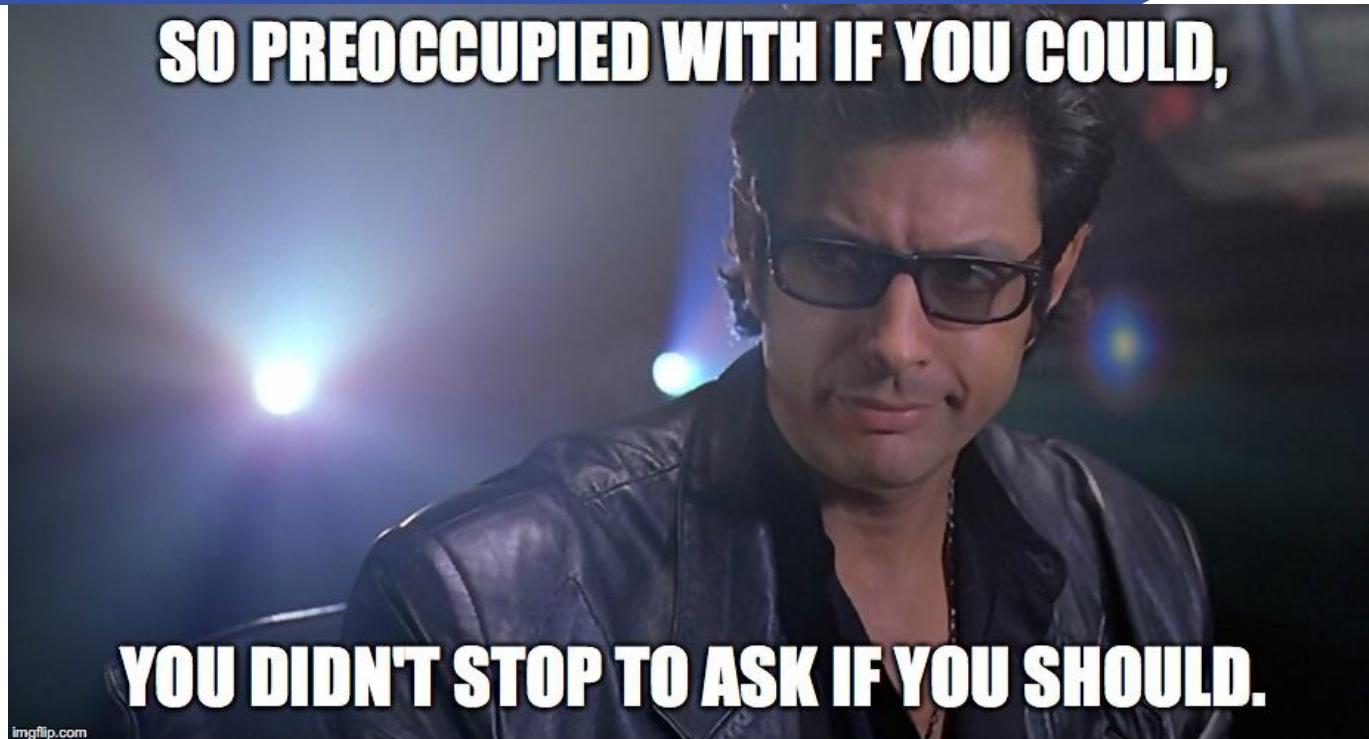
# Avoiding Pitfalls



- Overfitting
- Data Dredging / p-hacking
- <https://xkcd.com/882/>



# GDPR, ethics

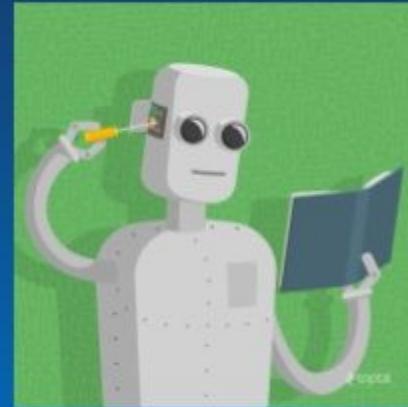


# Practical ML Applications



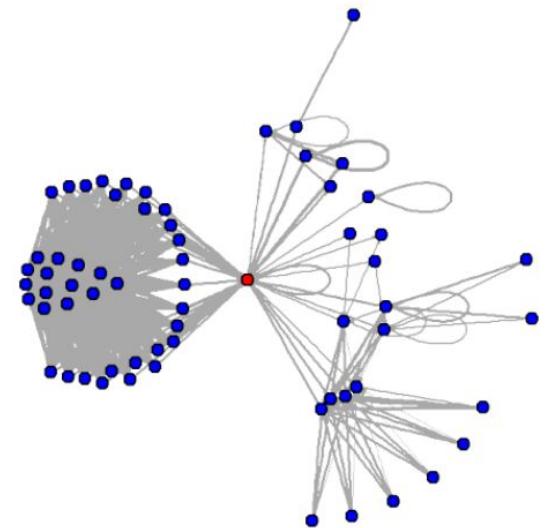
## Machine Learning: Practical applications

- Spam filtering
- Fraud detection
- Optical character recognition
- Medical diagnoses
- Anticipate customer needs
- Route requests efficiently
- Anticipate deals



# Projects

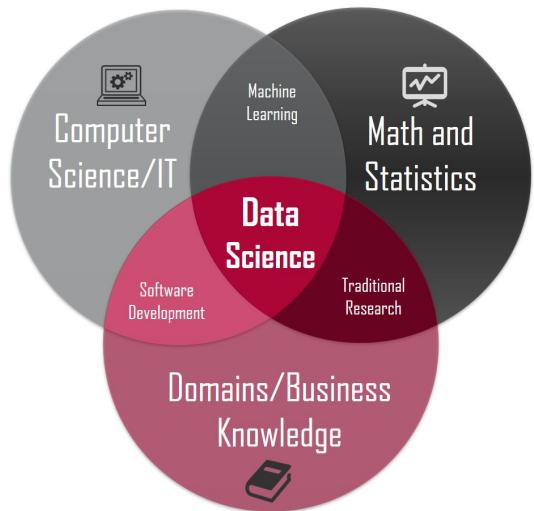
- Recommendation System / Churn Prediction
- Web Comments Sentiment Analysis
- Network Analysis (including blockchain)
- Visualizations with PowerBI( or Tableau)
- Dashboards with Dash/plotly



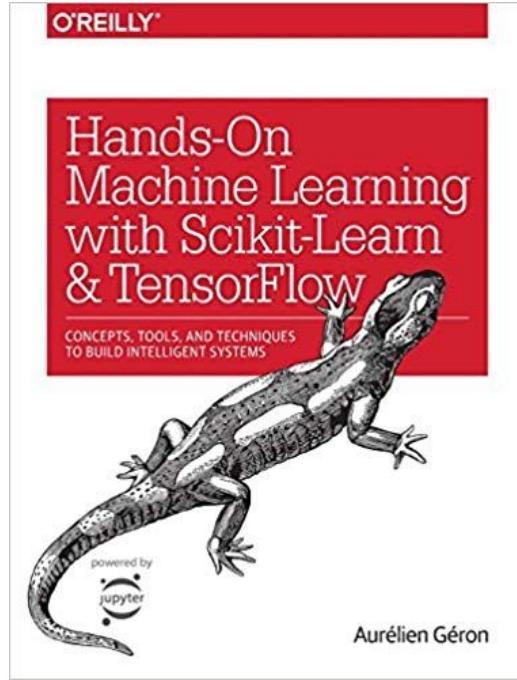
# Goals



- Access structured/unstructured data
- Clean data
- Apply correct methods for analysis
- Visualize Results



# Books



- General Python resources
- <https://automatetheboringstuff.com/>
- <https://www.py4e.com/lessons>

The BEST book on Python and Machine Learning

<https://www.amazon.com/Hands-Machine-Learning-Skikit-TensorFlow/dp/1491962291>

# Requirements



- Analytical / Logical Mind
- Helpful but NOT required knowledge:
- Comfortable in command line
- SQL
- Statistics
- Helpful: a computer with a minimum of 8GB RAM
- <https://www.anaconda.com/download/> (3.6+)





# PALDIES!

