

Module 1 :

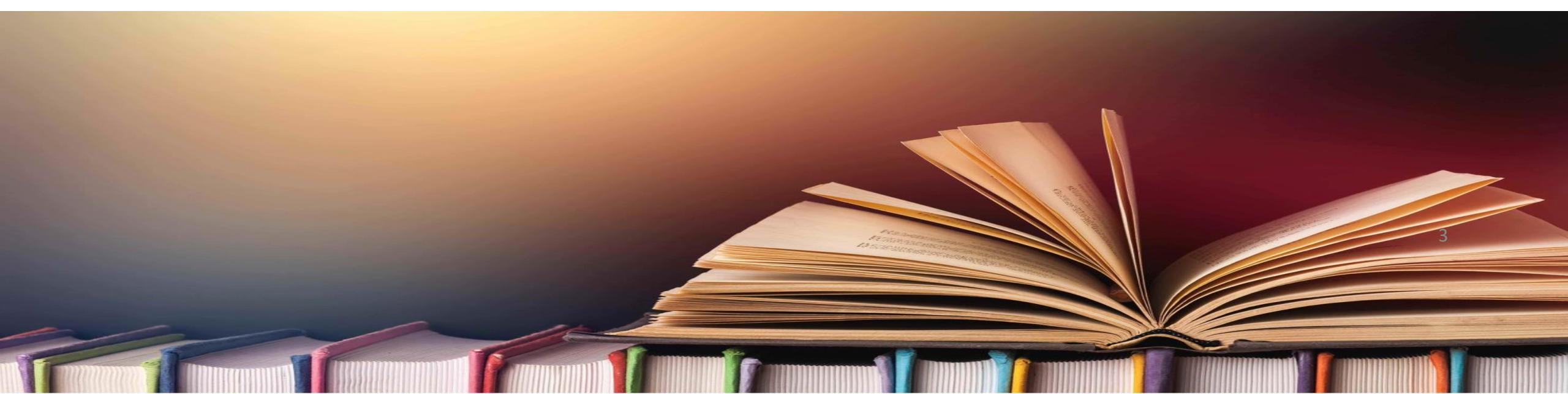
Machine Learning Review

Machine Learning
Introduction



Bibliography

- Deep Learning book (Goodfellow, Bengio, Courville)
- Machine Learning @ Stanford (Prof Andrew Ng)
- Hands-On Machine Learning with Scikit-Learn & Tensorflow (Aurélien Géron)

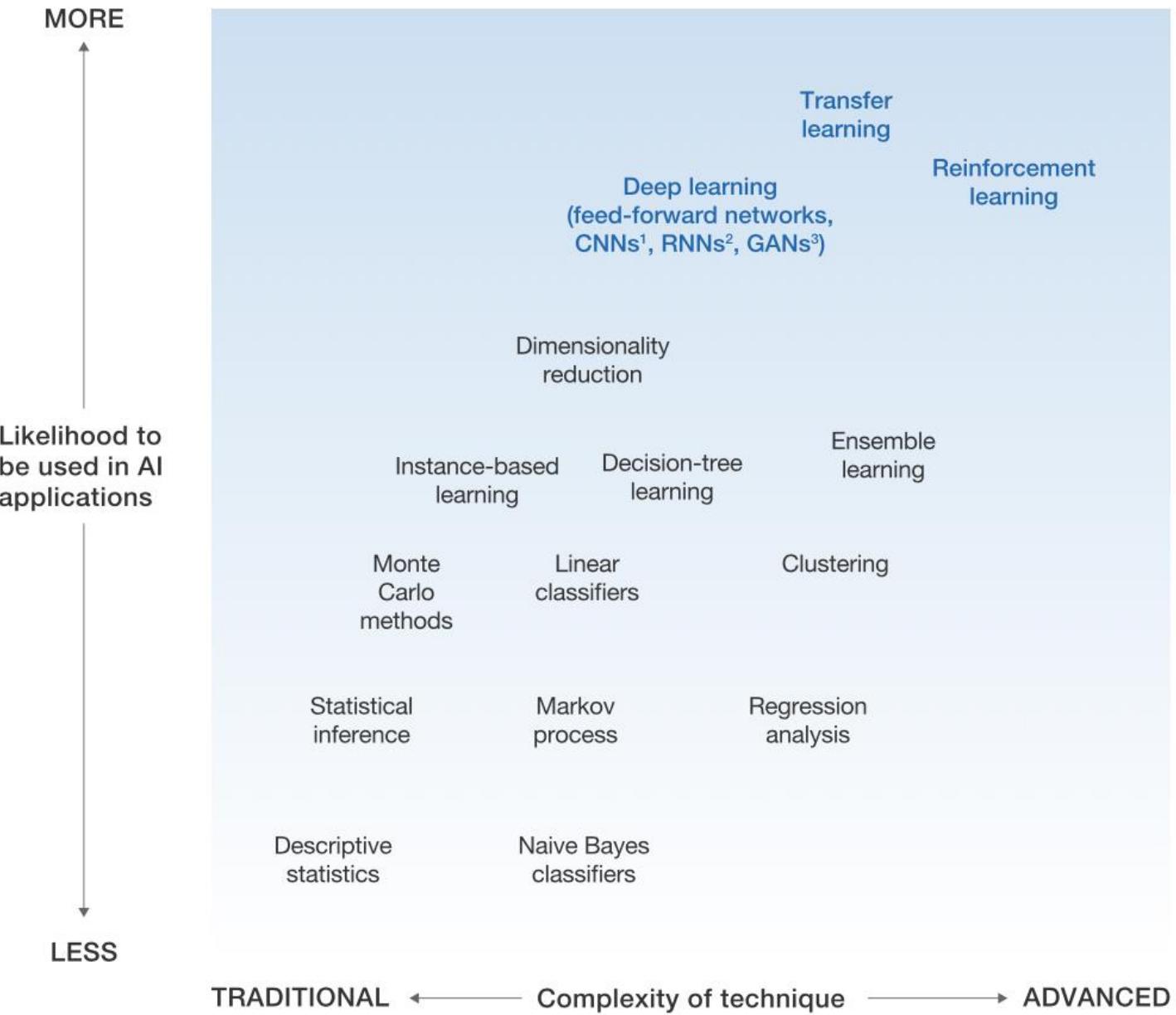




Learning Objectives

- What is Machine Learning ?
- Learning pillars
- Main Challenges
- Data Preparation
- Optimization
- Technical Details
- How to choose an algorithm ?

Differences between ML and Statistical Modeling



THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

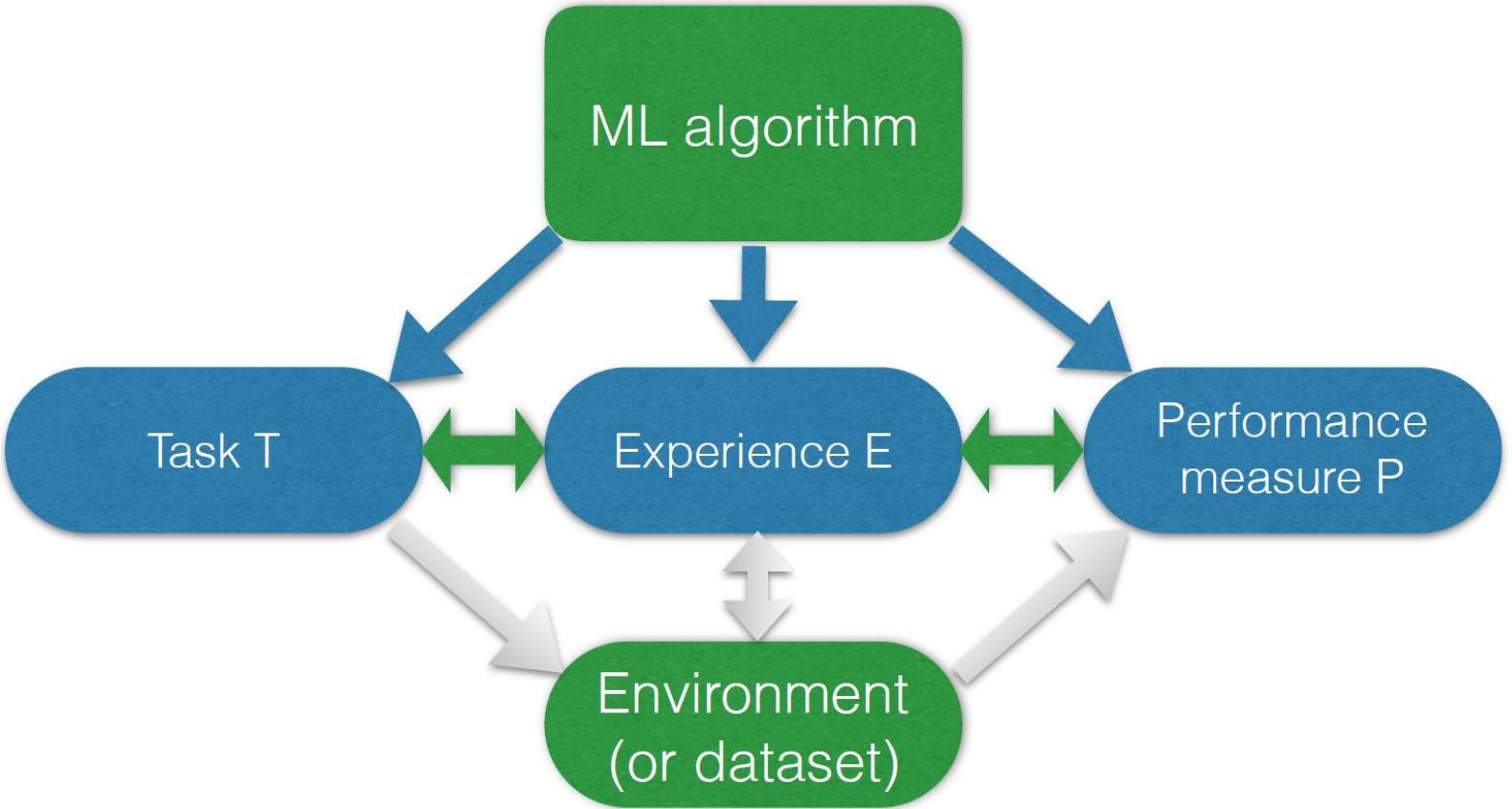
JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.



What is Machine Learning ?

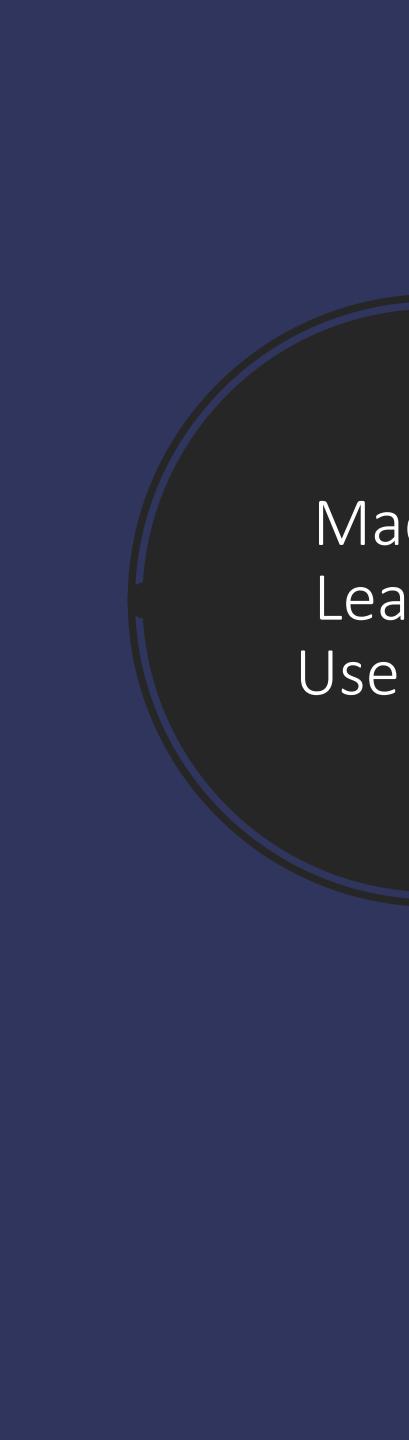
"Can machines do what we (as thinking entities) can do?" A. Turing

Definition



« A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E. »

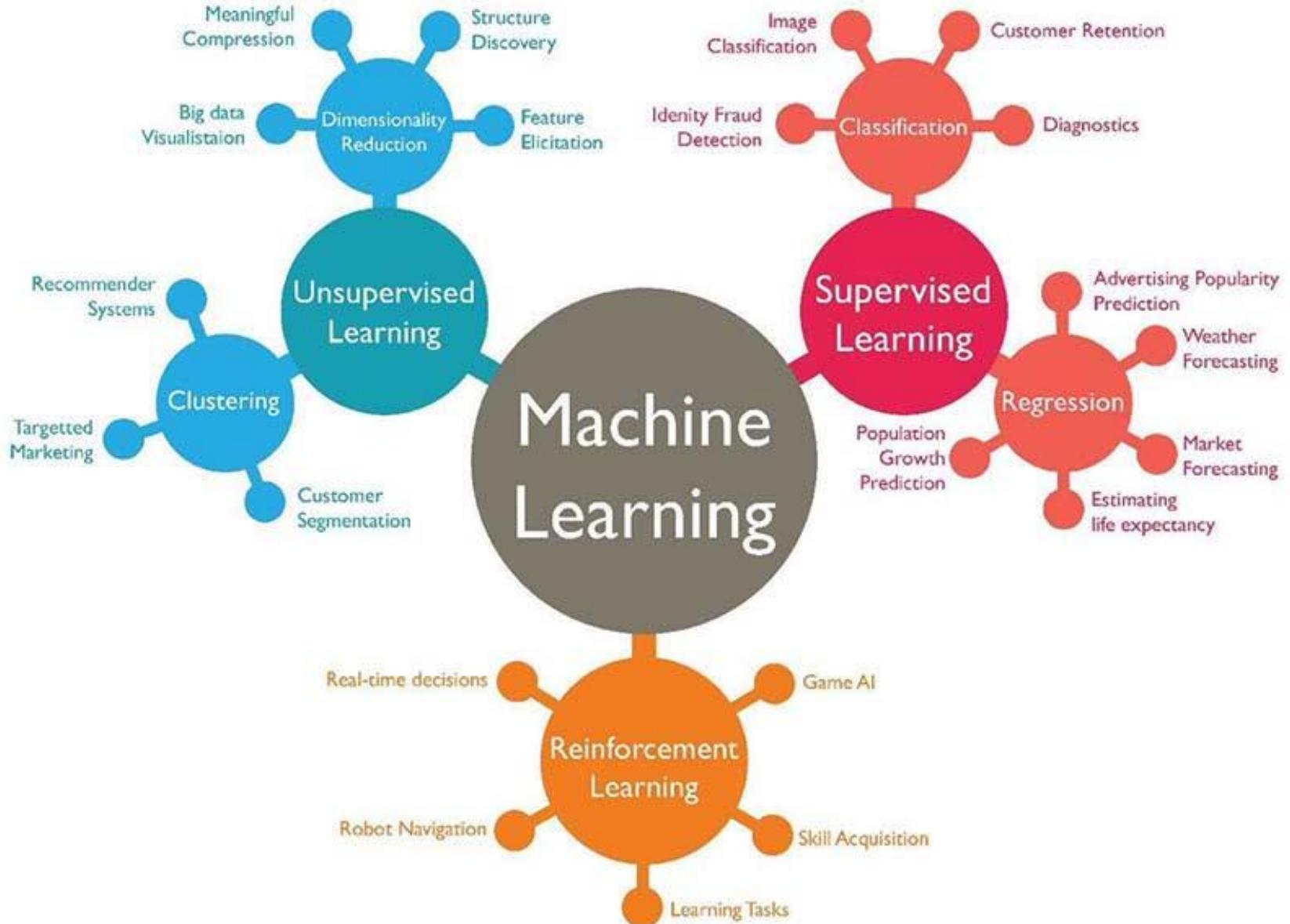
Tom M. Mitchell (1997)



Machine Learning Use Cases

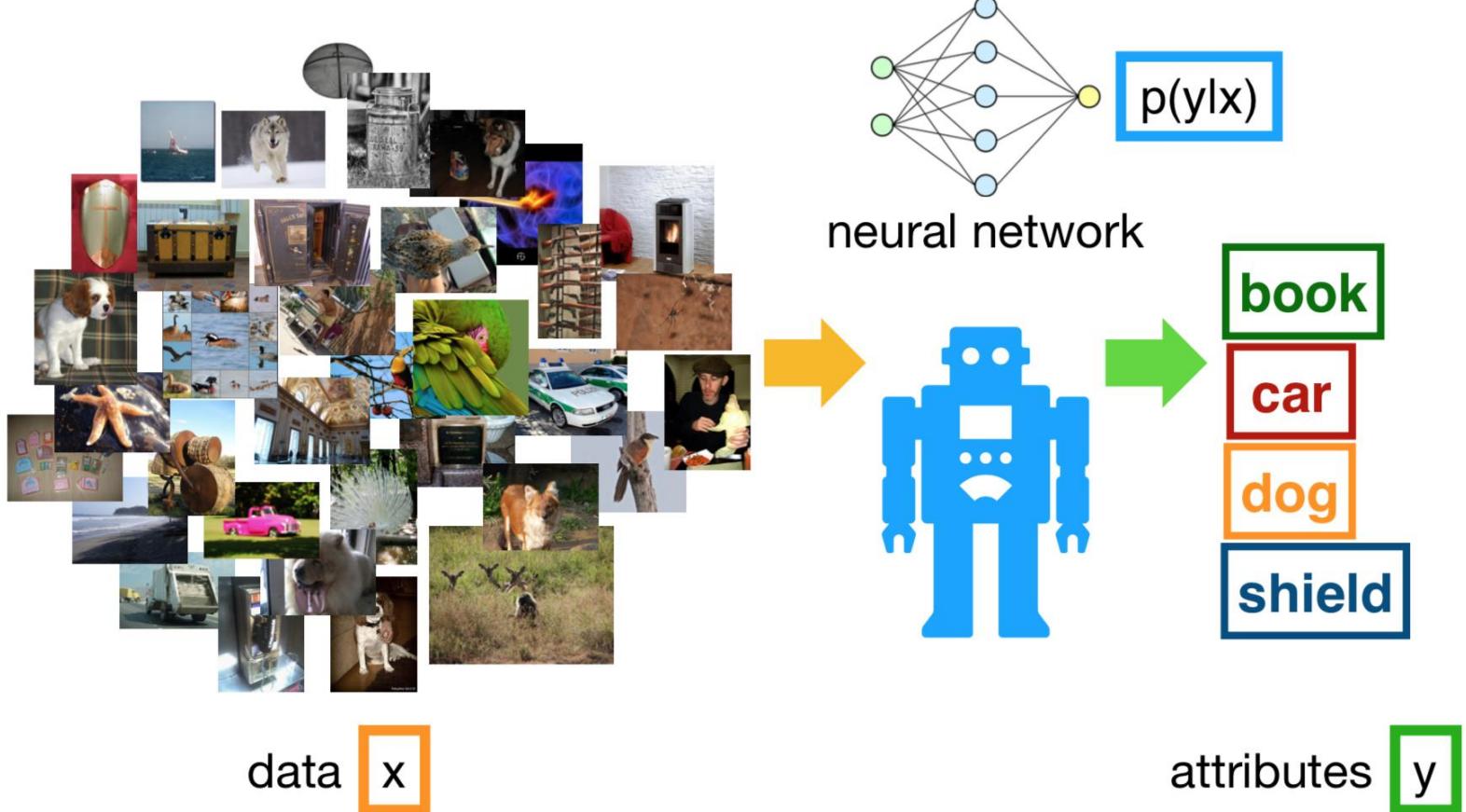
- Problems for which existing solutions require a lot of hand-tuning of long lists of rules
 - ML *simplifies* code and *performs better*
- Complex problems for which there is no good solution at all using a traditional approach
 - ML can *find a solution*
- Fluctuating environments
 - ML system can *adapt* to new data
- Getting insights about complex problems and large amounts of data

Learning Pillars



Supervised Learning

- Prediction of an output y given an input x



Regression



\$82000



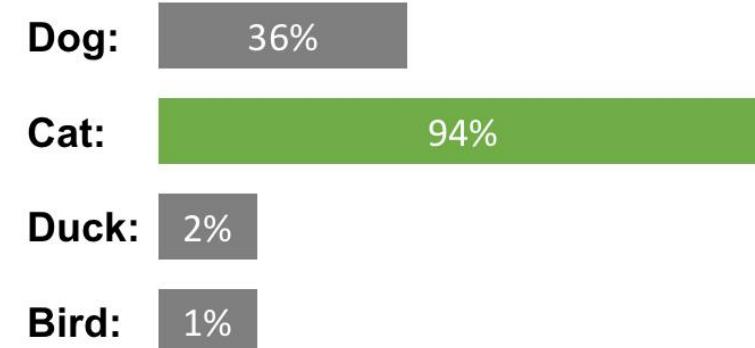
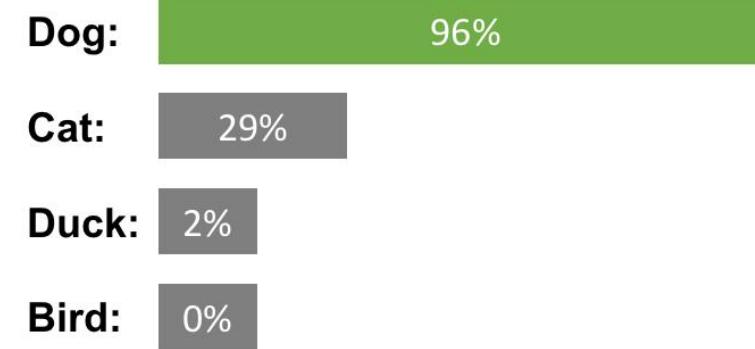
\$55500



???

Predict results within *a continuous output*

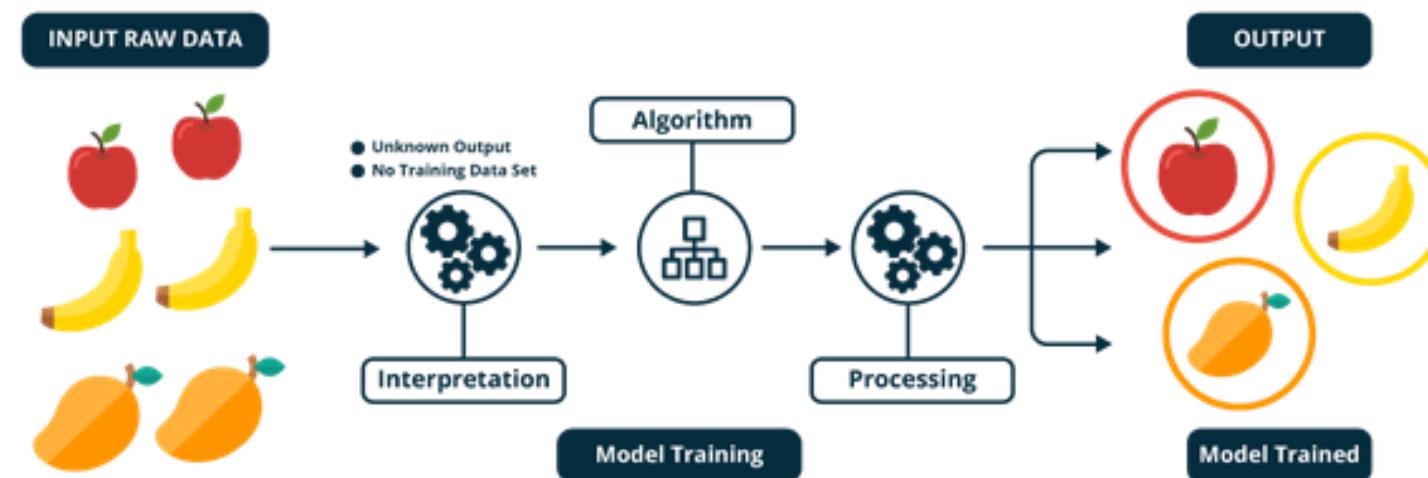
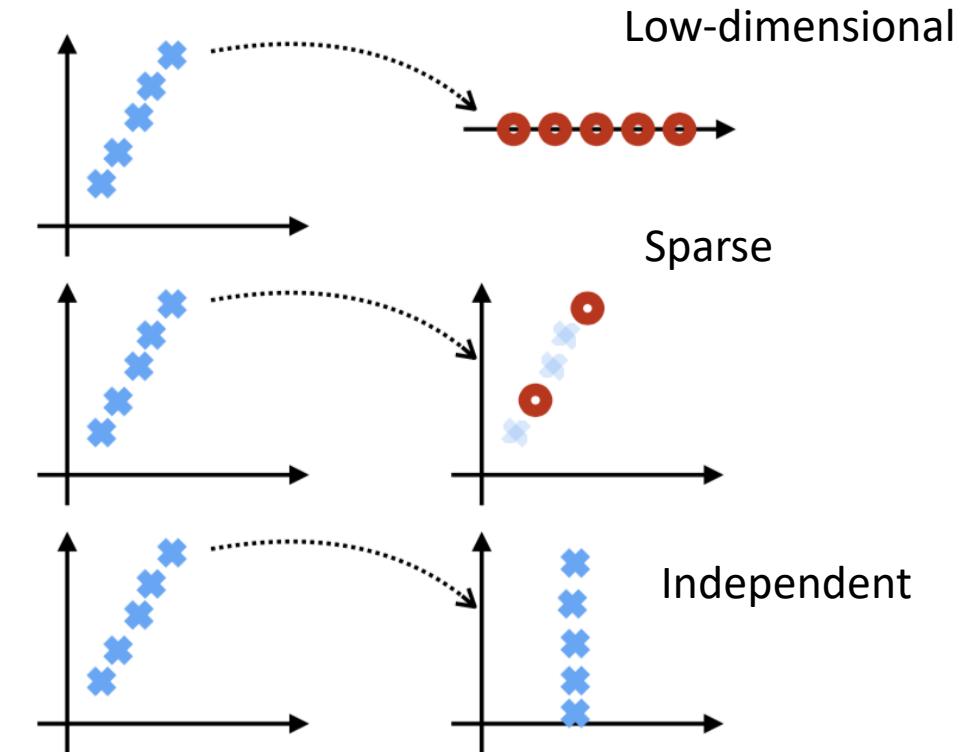
Classification



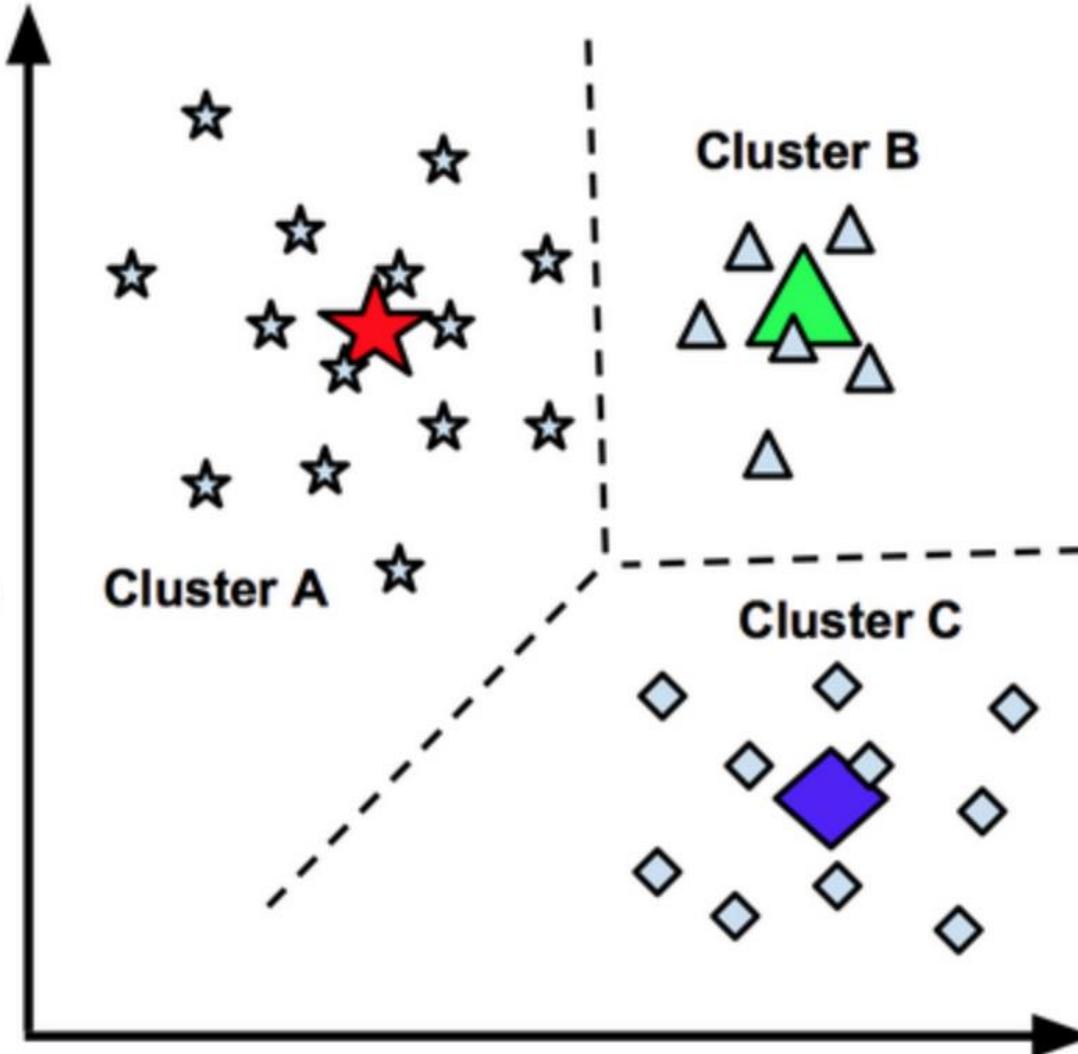
categorize new inputs as belonging to one of a set of categories
→ Predict results within a *discrete output (categories)*

Unsupervised Learning

- Find a *suitable data representation*
 - Preserving all task-relevant information
 - Simpler than the original data and easier to use



Clustering



Create a **set of categories**, for which individual data instances have a set of **common or similar characteristics**.

Representation Learning

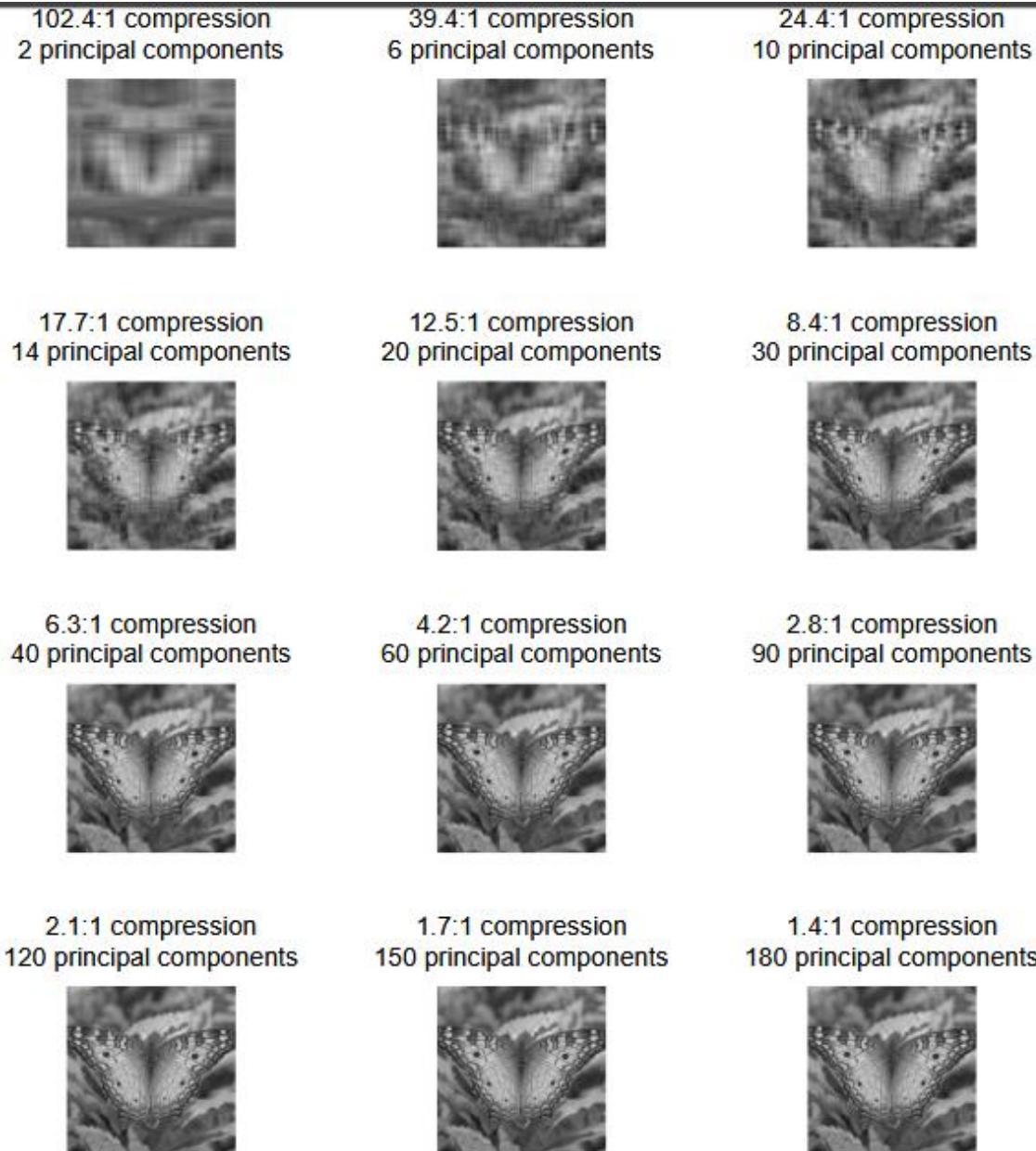


Figure 10: The visual effect of retaining principal components

More
nuances...



Supervised Classification



Semi-supervised Learning



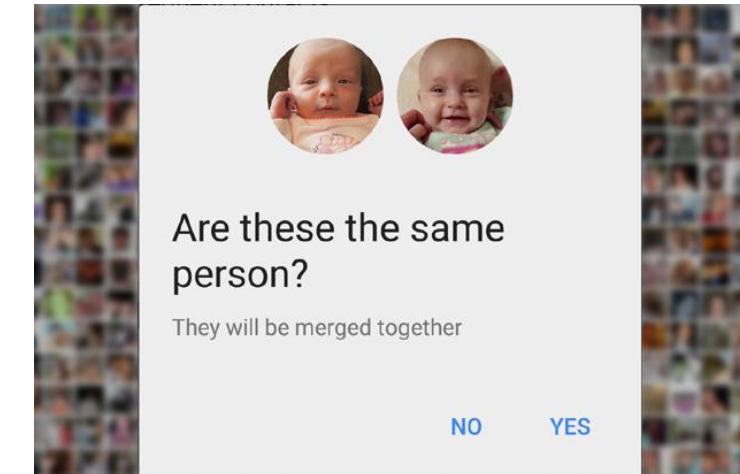
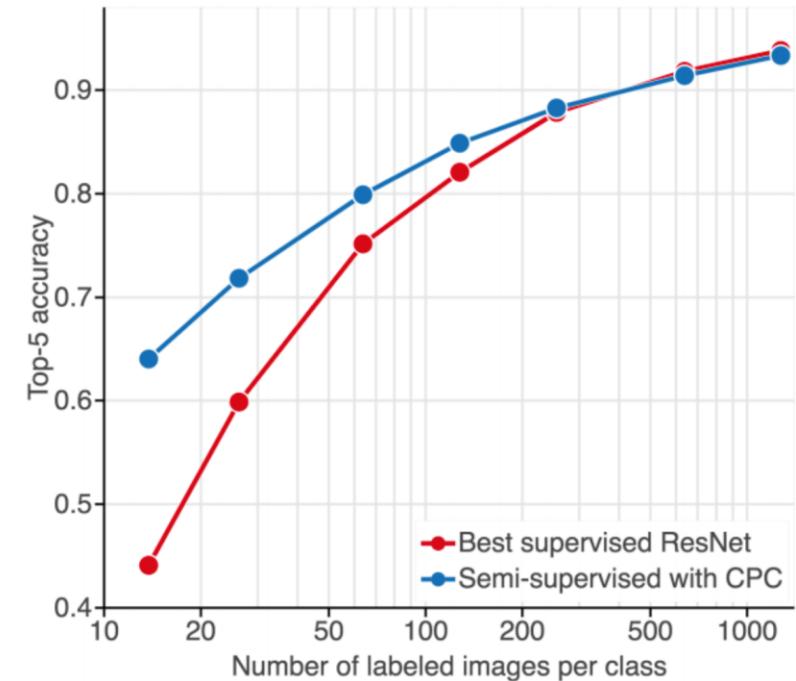
Transfer Learning



Self-taught Learning

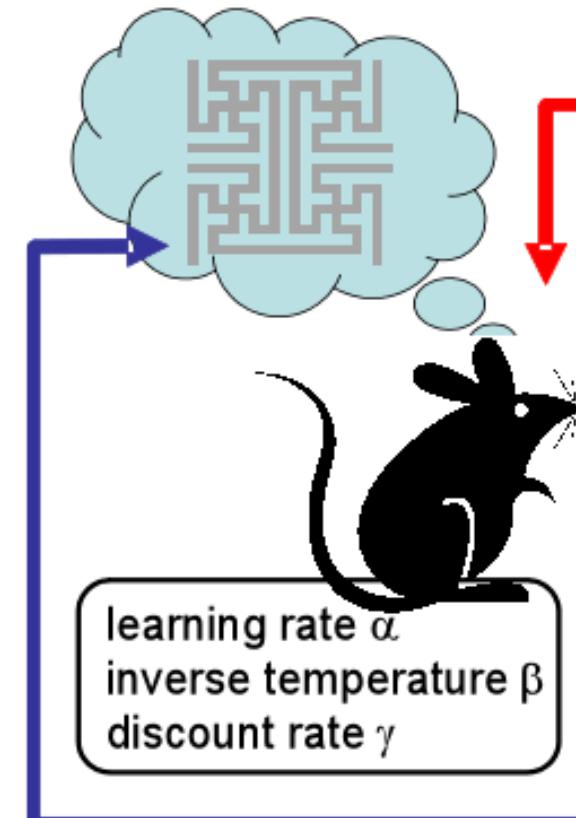
Semi-supervised Learning

- Partially labelled training data
 - Lots of unlabeled data and a little bit of labelled data
- Why bother ?
- Most semi-supervised learning algorithms are combinations of unsupervised and supervised algorithms
 - Example : Google Photos



Reinforcement Learning

internal state

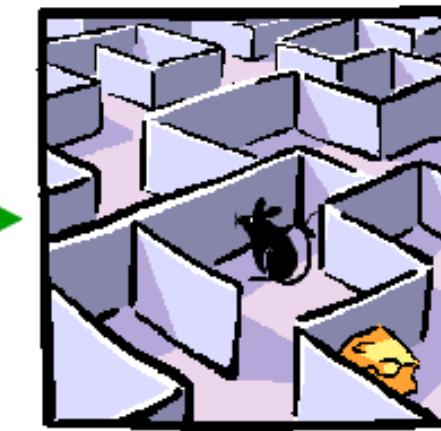


observation



reward

environment



action



<https://www.youtube.com/watch?v=Bui3DWs02h4>

Two-Minute Papers



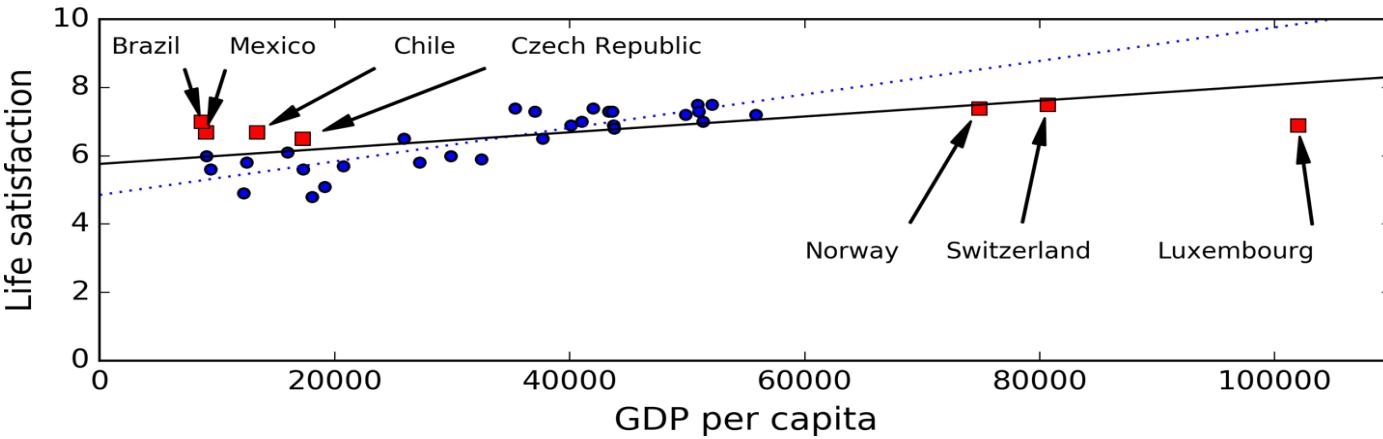
Main Challenges

- Bad data
- Bad model



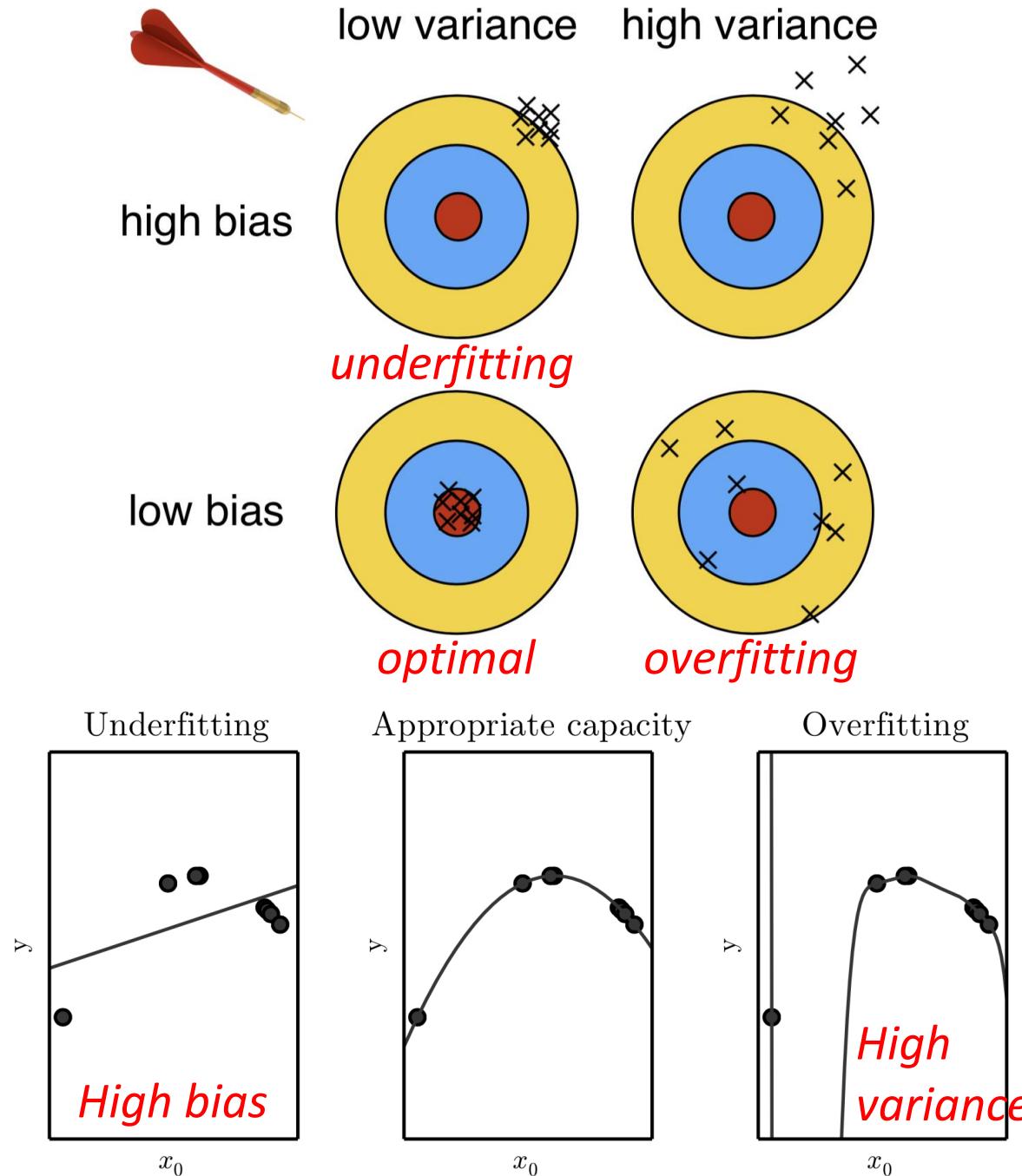
- Insufficient training data
 - Try to get more data, data augmentation,...

- Non-representative training data

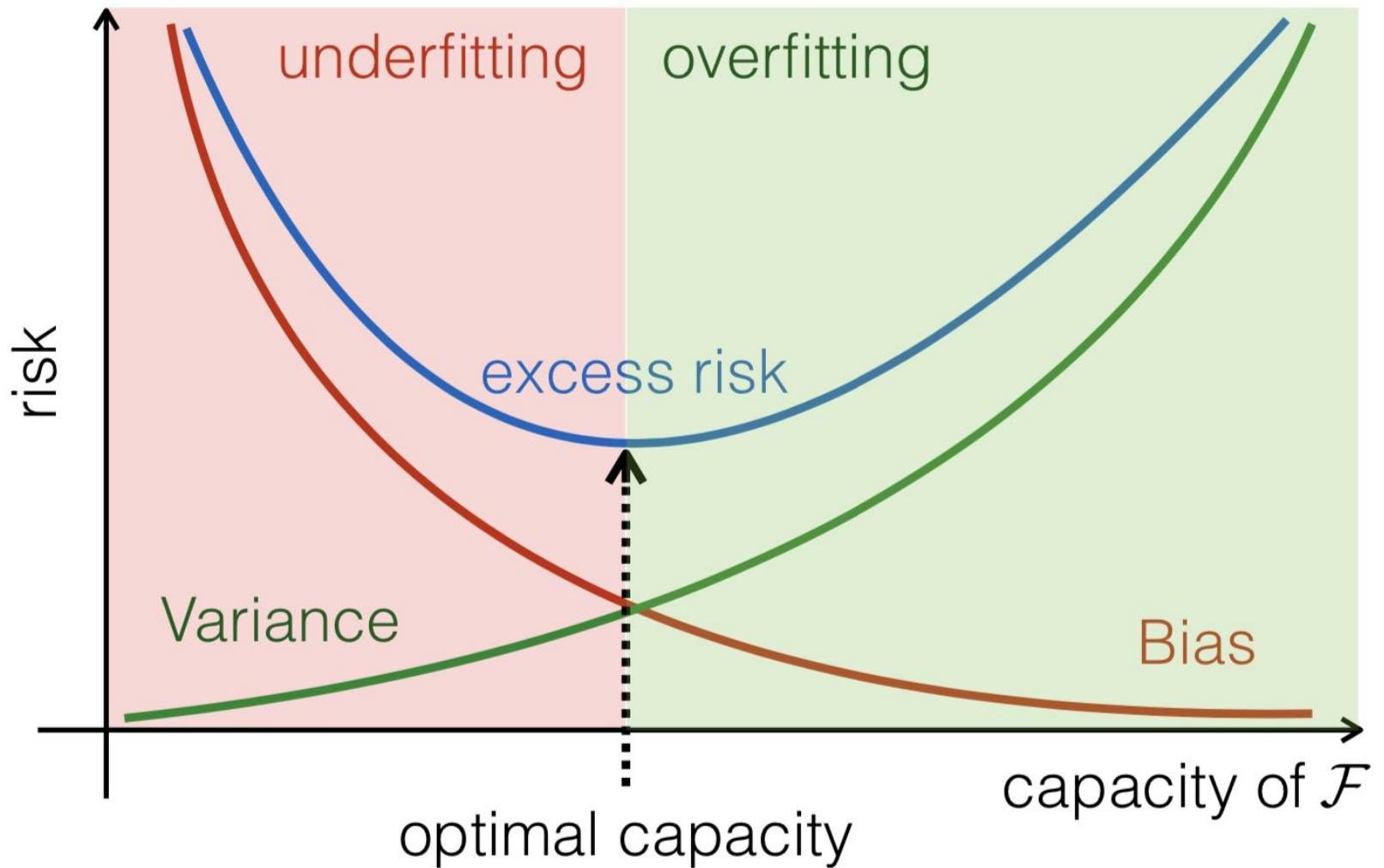


- Poor quality data
 - Data cleaning
- Irrelevant Features
 - Feature selection (select most useful features)
 - Feature extraction (combine features to produce a more useful one)

Bad Model : Overfitting and Underfitting



Model :
Bias,
variance and
capacity





Cross-validation

- Validation technique that assesses how the forecasts of a model generalize on unseen data
- A score can be computed on the forecasts across multiple test sets

Cross-validation : Holdout

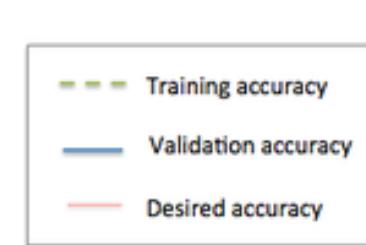
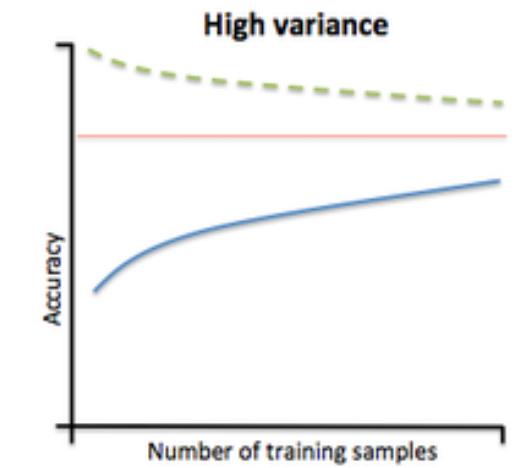
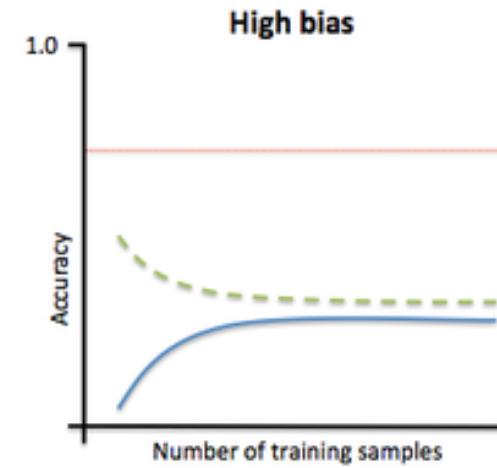
- Separate the data into 3 sets
 - Training set for training
 - Validation/dev set to fine-tune the model hyperparameters
 - Test set to estimate the FINAL model performance
- Separation depends on size of the dataset



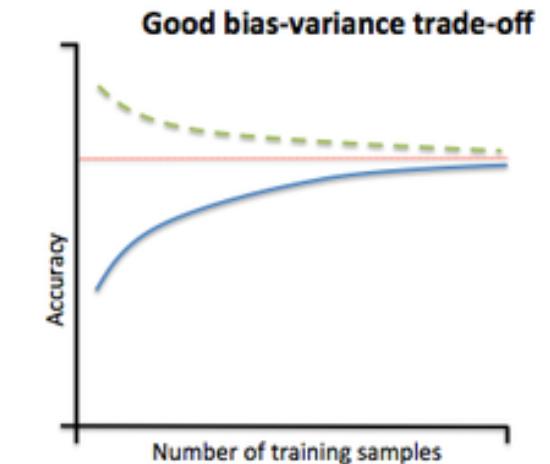
Validation set: use cases

- **Hyperparameter tuning** : train the model on the training set and evaluate its performance on the validation set
 - adjust the hyperparameters to find the best configuration, ie guide decisions and monitor training process

- **Early Stopping**: stop training early to prevent overfitting

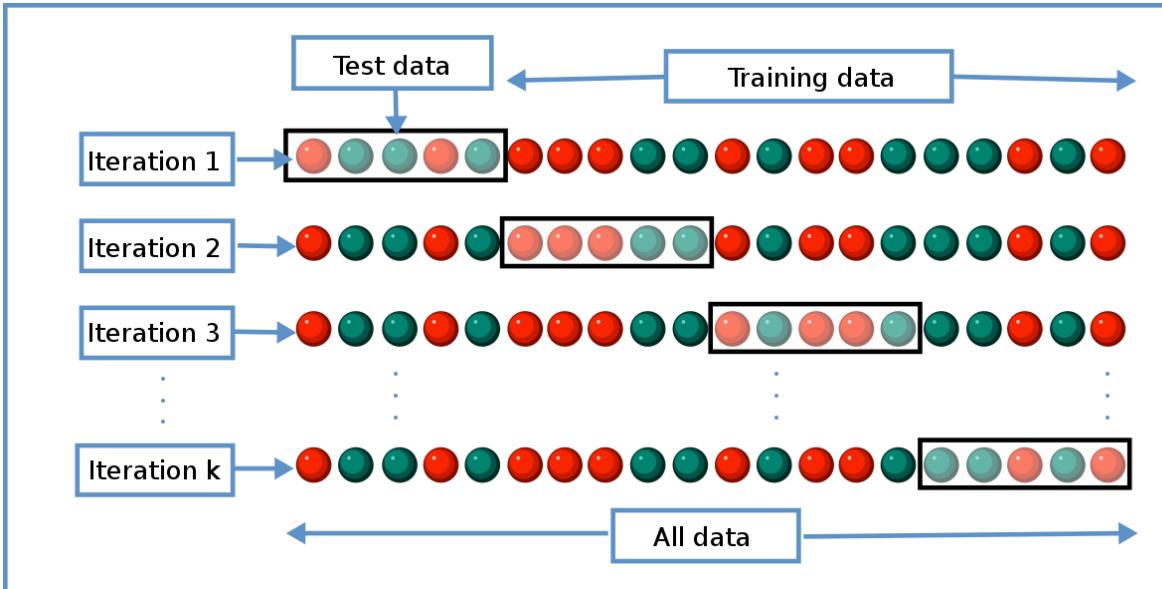


- **Model Selection**: If you're experimenting with different algorithms or architectures, compare their performance on the validation set to choose the best one **BEFORE** proceeding to the final evaluation on the test set.



Cross-validation : k-fold

- To avoid “wasting” too much training data
- test on **multiple splits** so we can get a better idea on how the model will perform on unseen data



- In principle better than using the holdout method because the holdout method score is dependent on how the data is split into train and test sets

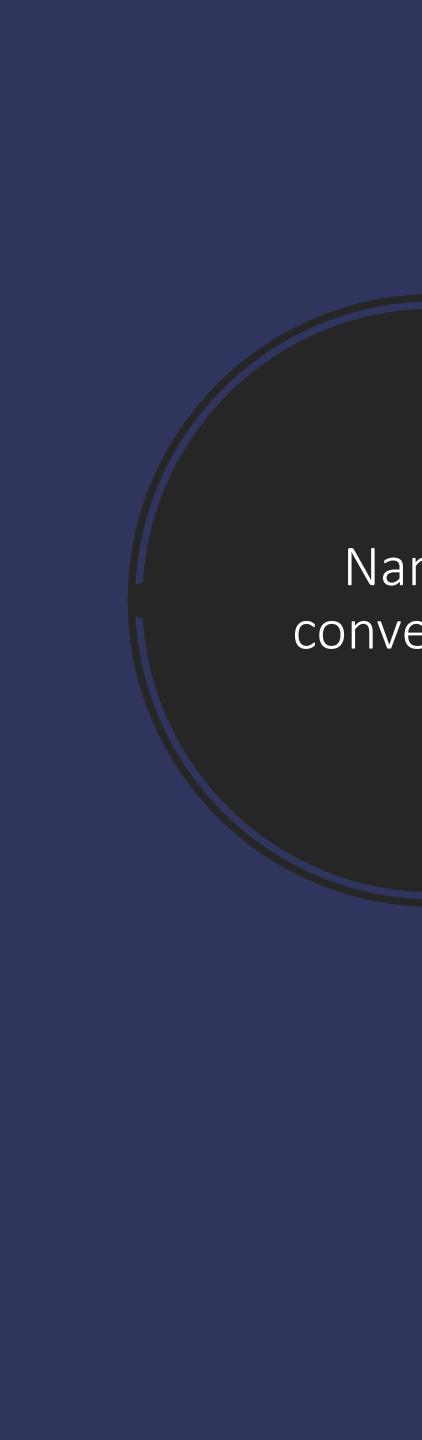
It should not be used with time series data



Data Preparation

Raw Data

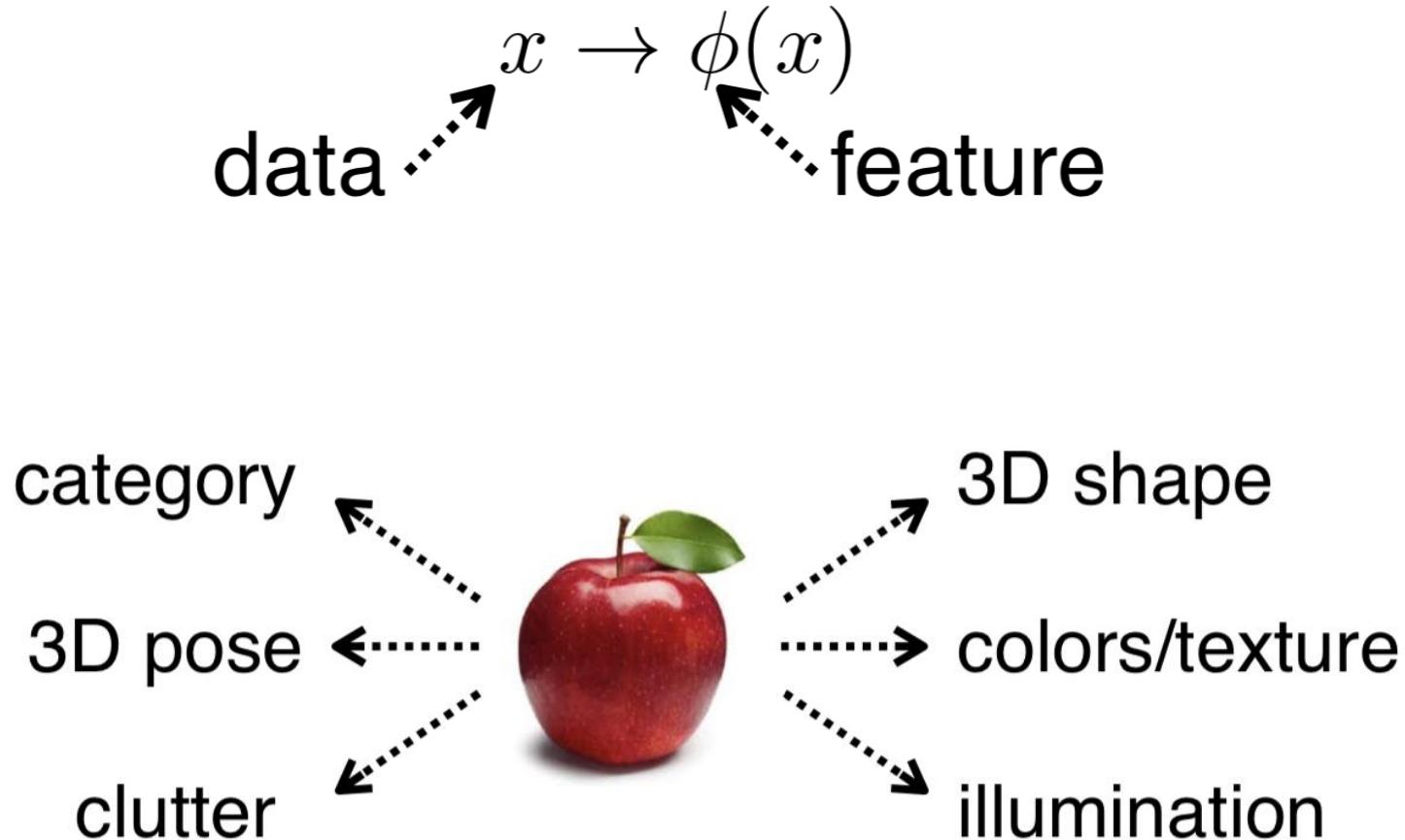
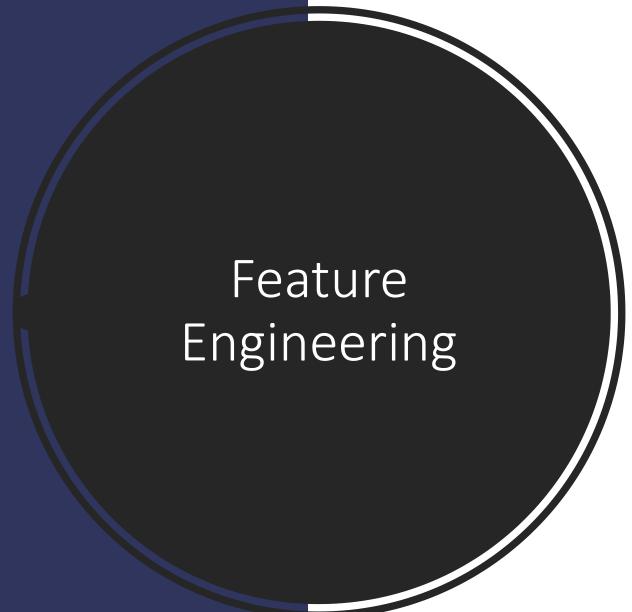
Type	Examples
Nominal	Name, Phone, Gender,...
Ordinal	1. Satisfied, 2. Indifferent, 3. Dissatisfied
Interval	Temperature, time, IQ test values,...
Ratio	<2 hours, 3-4 hours, 4-5 hours,...
Dichotomous	Sex (F/M), completed (Y/N),...



Naming conventions

- **Data Cleaning** : identify and fix errors in the data prior to modeling (outliers, missing values, ...)
- **Feature Engineering** : create new input variables from raw data (*next slide*)
- **Data Wrangling** : more general or colloquial term for data preparation that might include some data cleaning and feature engineering.

Data often encoded into **more focused relevant information** (features or internal representation) to simplify the decision





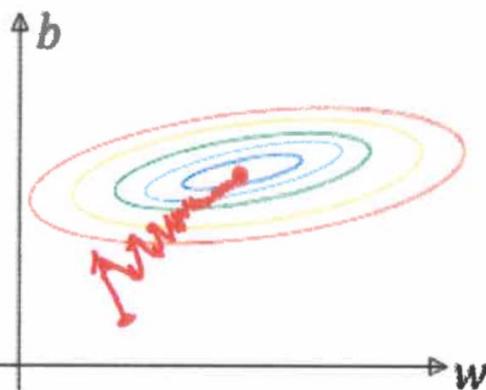
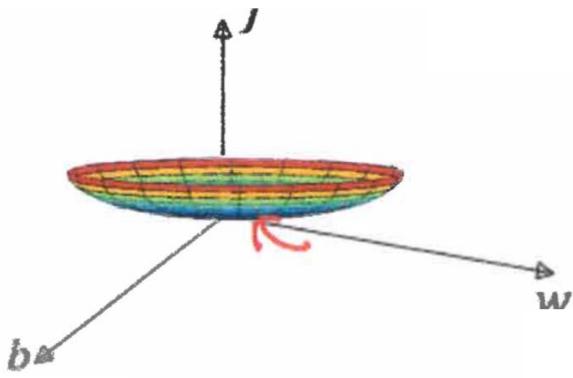
1. Data Pre-Processing

- Convert labels into one-hot encoded vectors to avoid biases
 - [1 2 3 0 2] → [0100 0010 0001 1000 0010]
- Shuffling to ensure independent samples and an unbiased estimate
- Remove under-represented categories

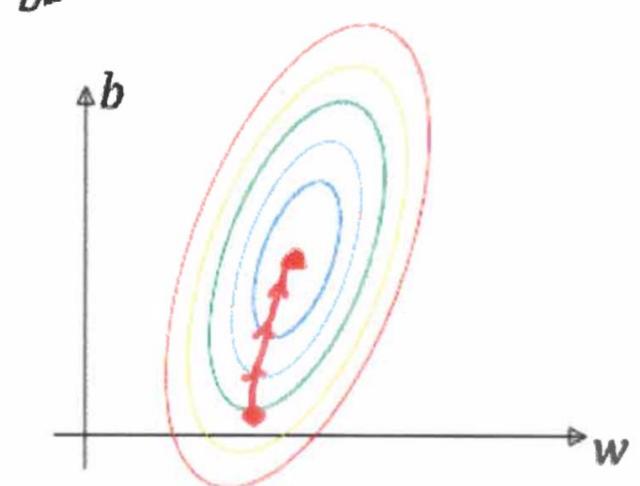
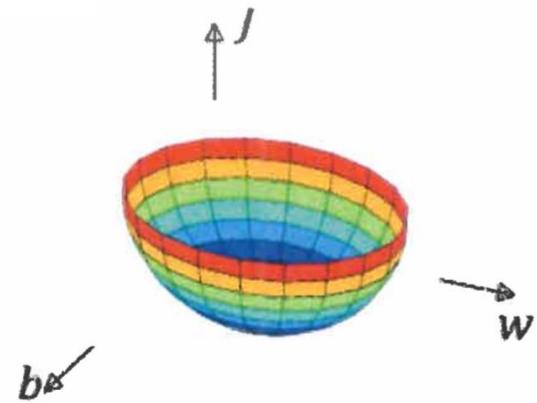
2. Standardization (Normalization)

- for faster convergence and more reliable predictions
- Applied to training data (`scaler.fit_transform`)
 - Learn mean and standard deviation, which are then applied to the validation and test samples (`scaler.transform`)

Unnormalized :



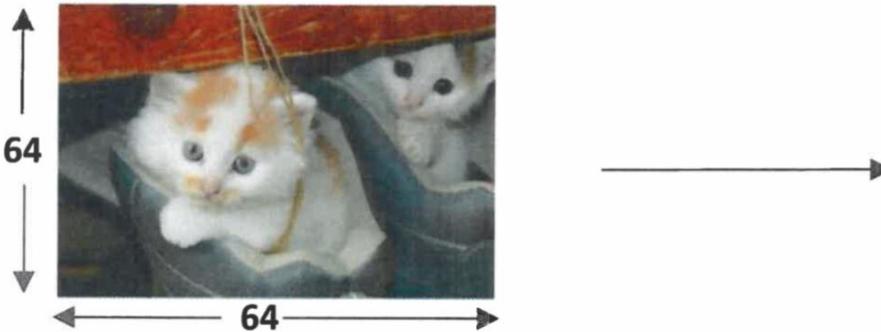
Normalized :



Data Preparation Example : Binary Classification

1. Data Preprocessing :

- Image stored in 3 matrices (red, green, blue color channels)
 - Similar size as the image



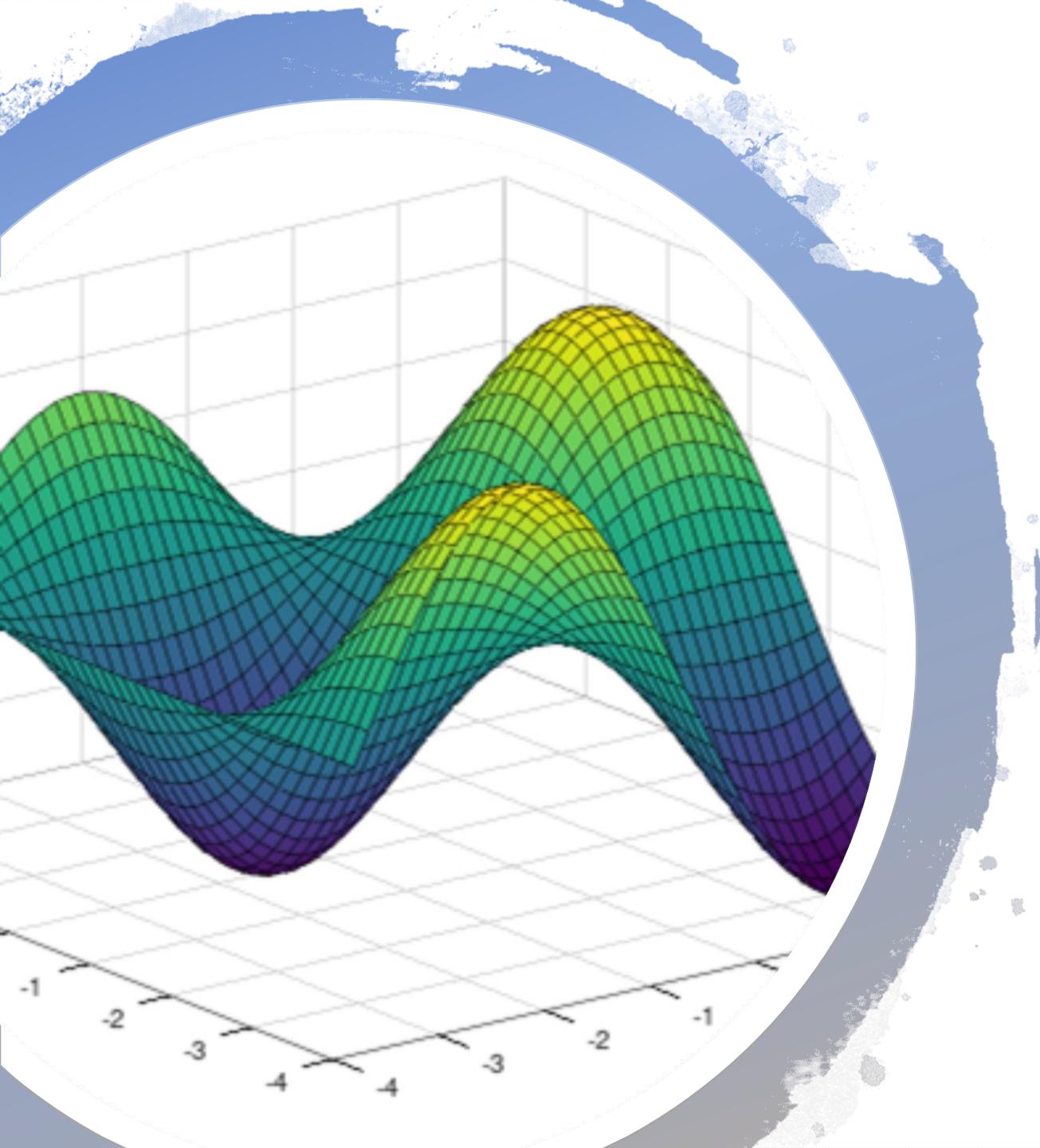
	Blue			
Green	255	134	93	22
Red	255	134	202	22
255	231	42	22	4
123	94	83	2	192
34	44	187	92	34
34	76	232	124	94
67	83	194	202	

$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix}$$

red
green
blue

- Reshape the values into a **feature vector x**

2. Normalization : divide by 255 (more convenient for data)



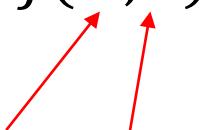
Optimization

Loss and Cost Functions

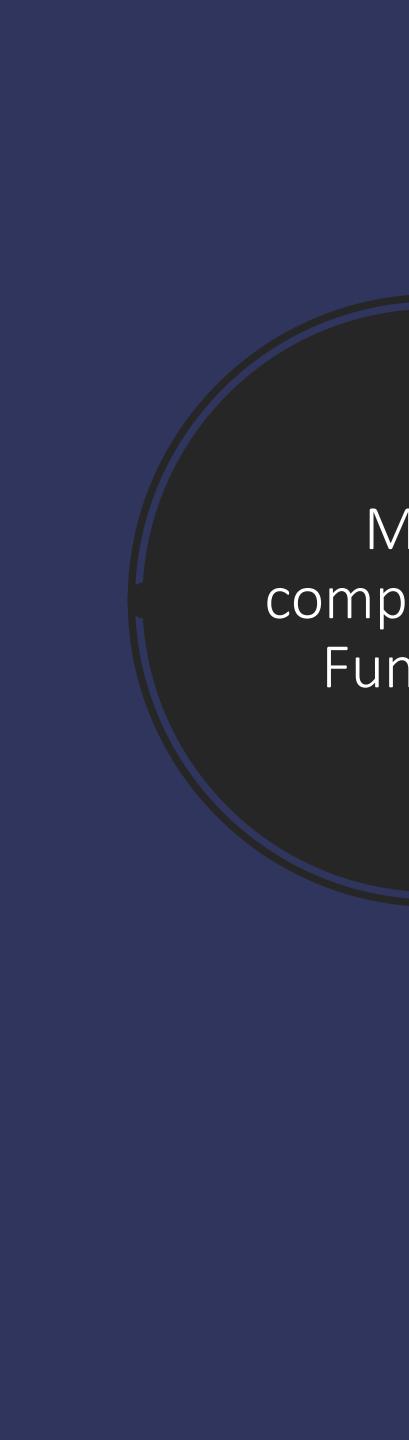
- **Loss function** $L(\hat{y}^{(i)}, y^{(i)})$, also called error function, measures **how different** the prediction $\hat{y} = f(x)$ and the desired output y are

- **Cost function** $J(w, b)$ is the average of the loss function on the **entire training set**

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$


weights, biases

- Goal of the optimization is to find the **parameters** $\theta = (w, b)$ that minimize the cost function



More complex Loss Function

- Can include more features :
 - **Regularization terms**: incorporate additional information to prevent overfitting, like L1 (lasso) and L2 (ridge)
 - **Feature Importance Weights**: assign different degrees of importance to different features (more or less relevant features)
 - **Attention Mechanisms**: give different levels of importance to different parts of the input sequence (transformers).
 - **Multi-Task Learning**: the loss function combines the losses from each task, possibly with different weights to reflect the importance of each task.

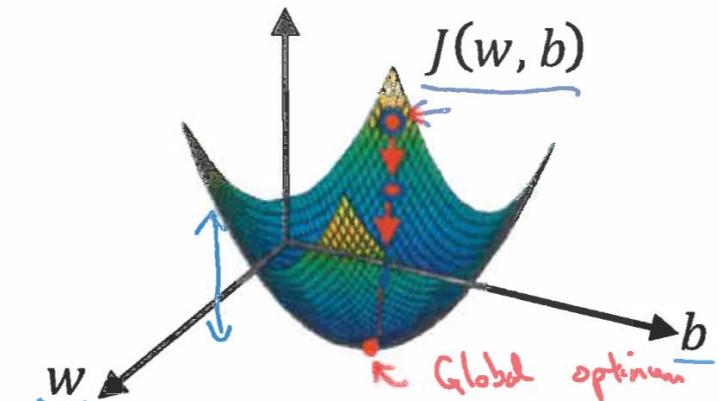
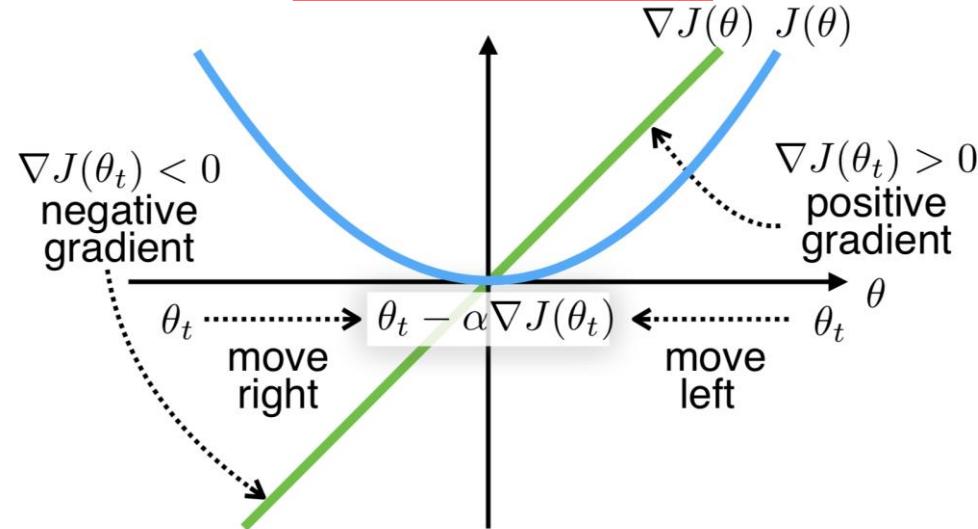
Gradient Descent

- iterative method to find the parameters $\theta = (w, b)$ that minimize $J(\theta)$

- Learning rate α

gradient descent

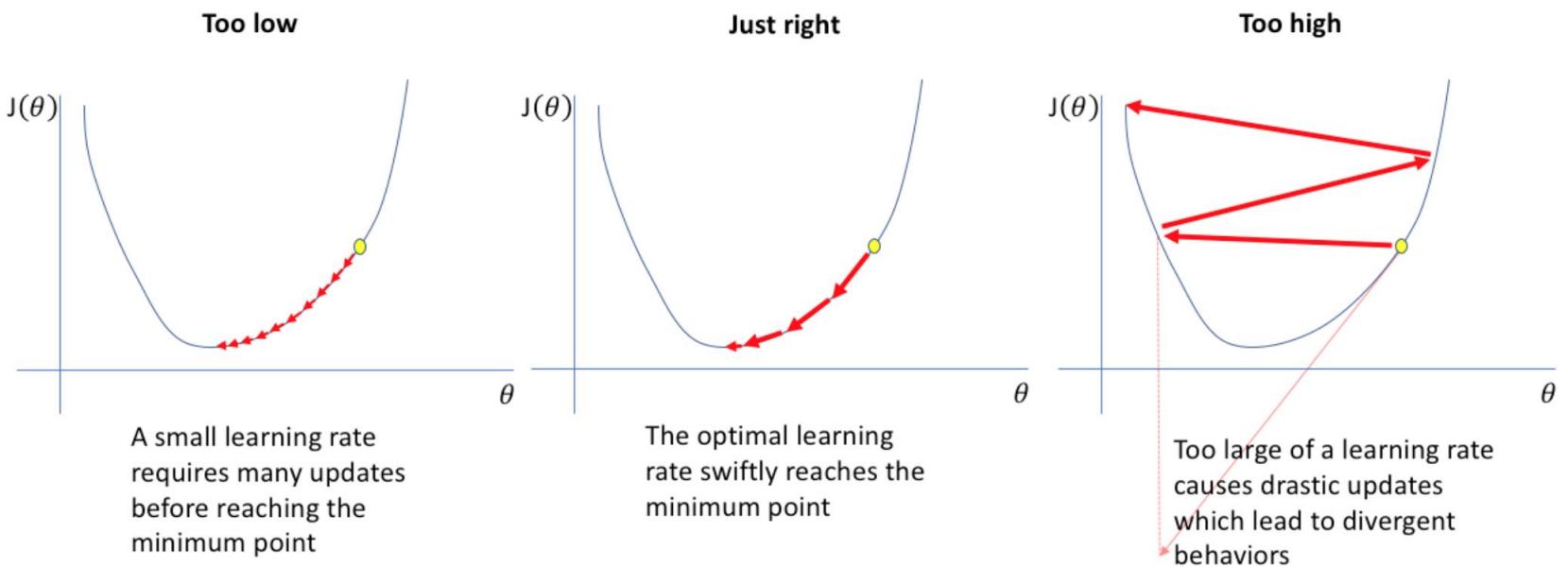
$$\theta_{t+1} = \theta_t - \alpha \nabla J(\theta_t)$$



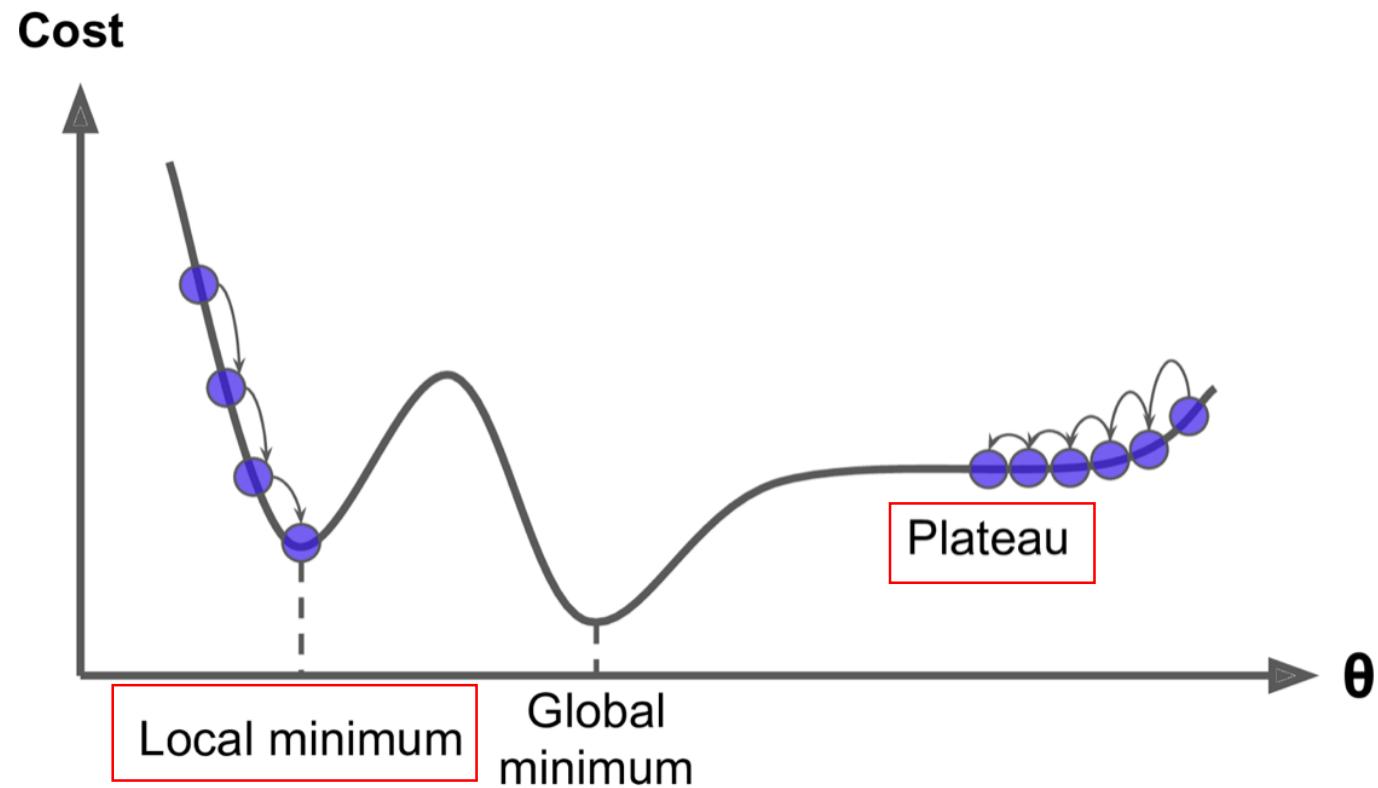
$$\nabla J(w) = \frac{dJ(w, b)}{dw}$$

$$\nabla J(b) = \frac{dJ(w, b)}{db}$$

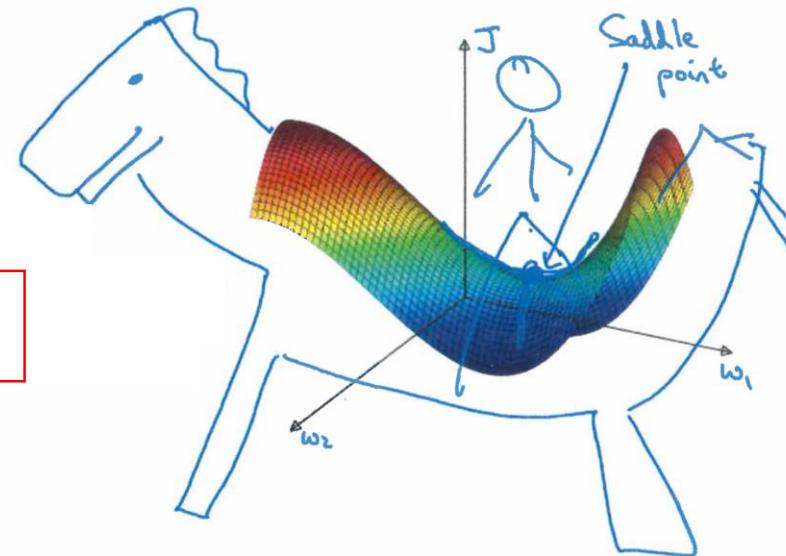
- Has a significant impact on the model performance, while being **one of the most difficult parameters to set**



Optimization Pitfalls



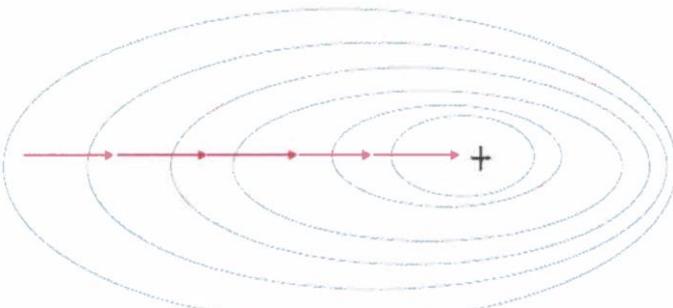
Saddle point



Mini-batch Gradient Descent and Stochastic Gradient Descent

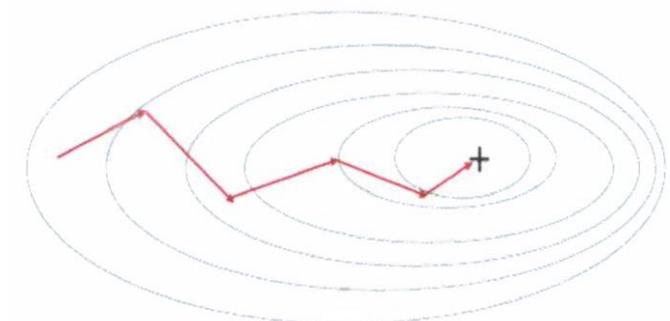
- At each iteration :

Gradient descent (GD)



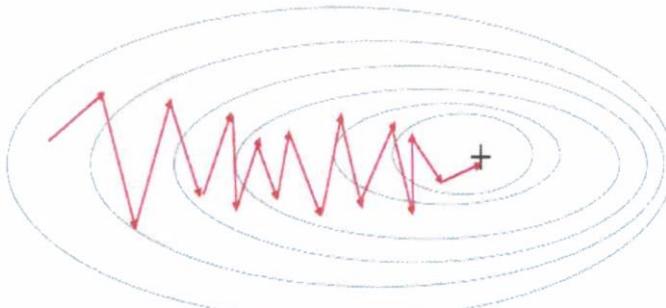
the *whole* training set

Mini-batch gradient
descent



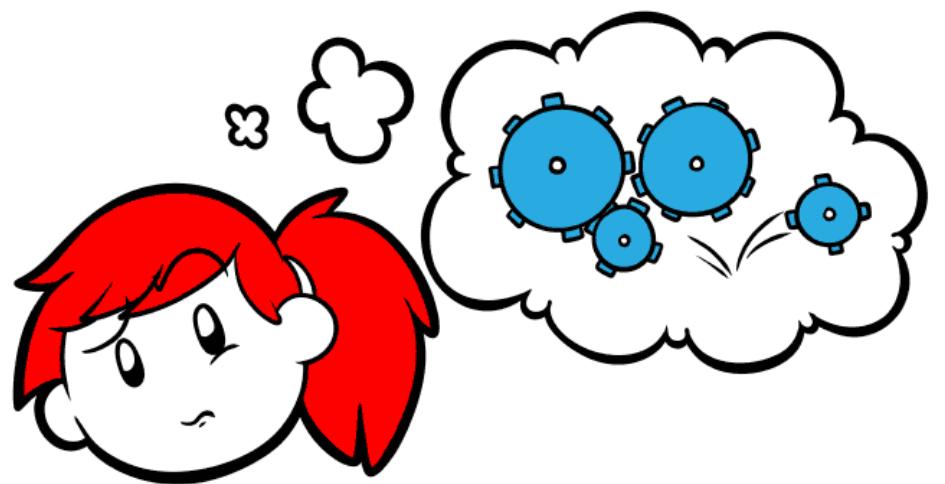
a *batch* of samples

Stochastic gradient
descent (SGD)



1 sample

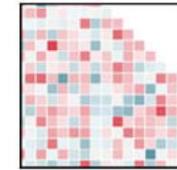
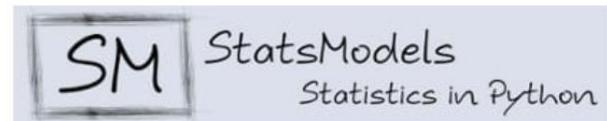
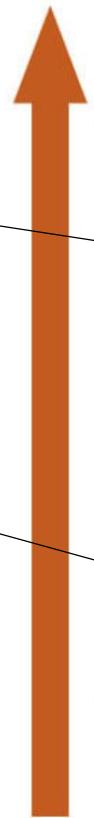
Batch size choice *typically*
32,64,128,256,512



Technical Details

Tools

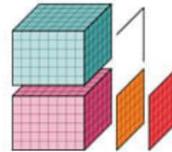
data structure & analysis



Seaborn

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



xarray



scikit
learn

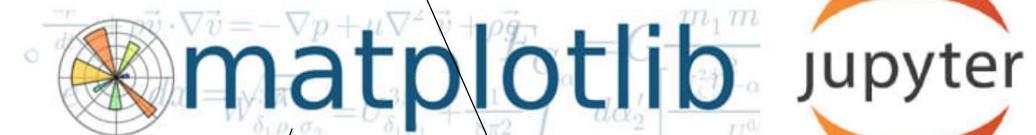


scikit-image
image processing in python

main package for scientific computing in Python



NumPy



matplotlib



jupyter

interactive coding environments embedded in a webpage

Python-based ecosystem of open-source software for mathematics, science, and engineering.



famous library to plot graphs in Python

IP[y]:
IPython

provides simple and efficient tools for data mining and data analysis

h5py : common package to interact with a dataset that is stored on an H5 file

Notebooks (1)

Use Anaconda Navigator

New environment with tensorflow :

conda create -n tf tensorflow

conda activate tf

Notebooks (2)

Add libraries :

conda install matplotlib

conda install scikit-learn

conda install pandas

conda install seaborn

conda install tensorflow

conda install git

conda install pydot

conda install graphviz

conda install anaconda-client

Anaconda cloud : file sharing

- You can run on the Anaconda cloud to get faster results
 - Log into <https://anaconda.org/>
 - conda install anaconda-client
- Command : anaconda upload notebook_M1D4_BON.ipynb
- Results : https://anaconda.org/username/notebook_m1d4_bon

Cloud computing

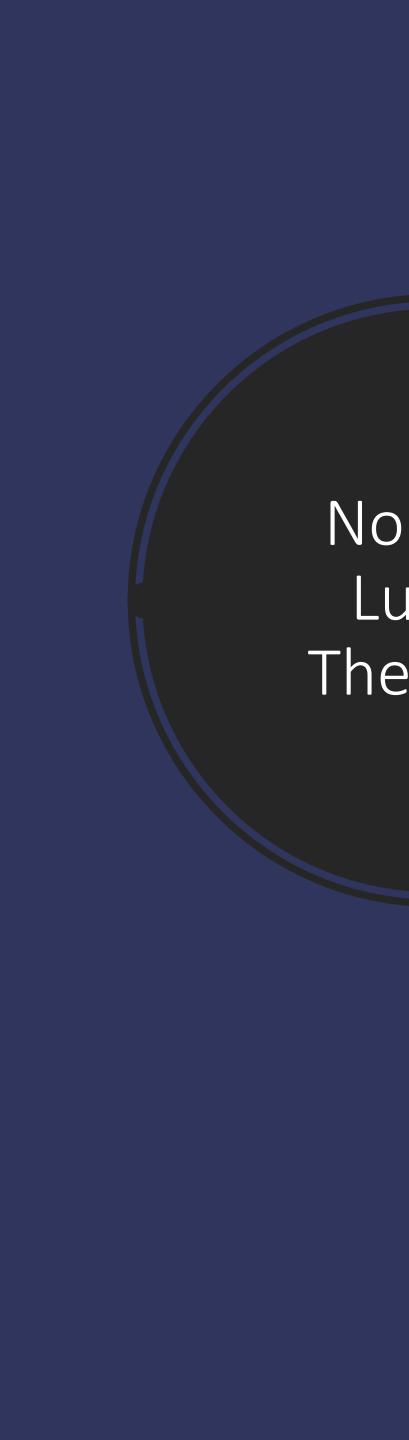
	Format	Pros	Cons
Google Colab	.ipynb	No setup required Access to GPU/TPU Integration with Google Drive	Limited runtime and resources for FREE version
Kaggle Kernels	.ipynb	Access to Kaggle datasets Competitive challenges	Limited resource availability Kernel runtime limits
Microsoft Azure	.ipynb	Integration with Microsoft services Access to Azure services	Limites resources and usage
IBM Watson studio	.ipynb	Integration with IBM services Easy data access	Limited resources for FREE version
AWS Educate	.ipynb .py	Real-world cloud experience Access to various AWS services	Limited credits, requires setup
Free Tier Cloud VM	.py	More flexibility Full control over environment	Setup required, limited resources, potential charges if usage exceeds FREE tier limits

Google Colab

- Here : <https://colab.research.google.com/>
- Import the Jupyter notebook and execute



How to find its way in
the algorithm jungle ?

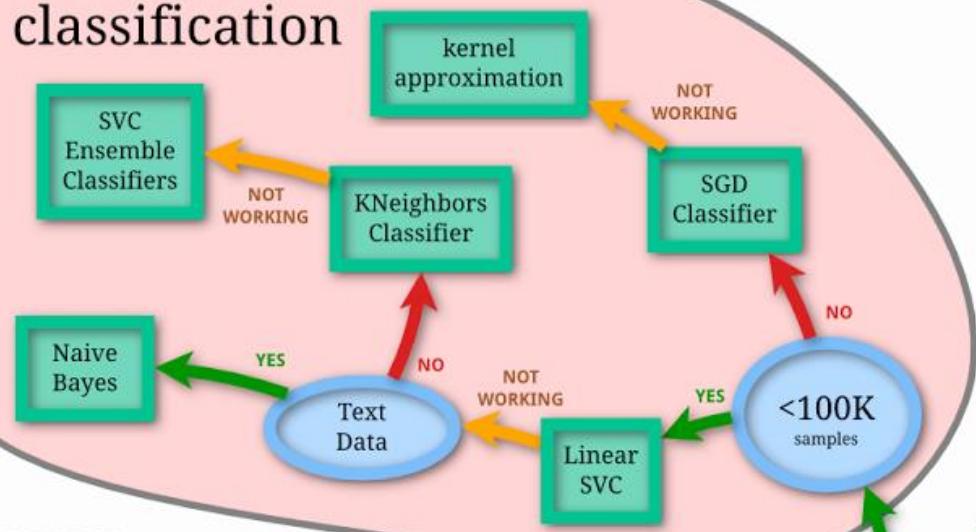


No Free Lunch Theorem

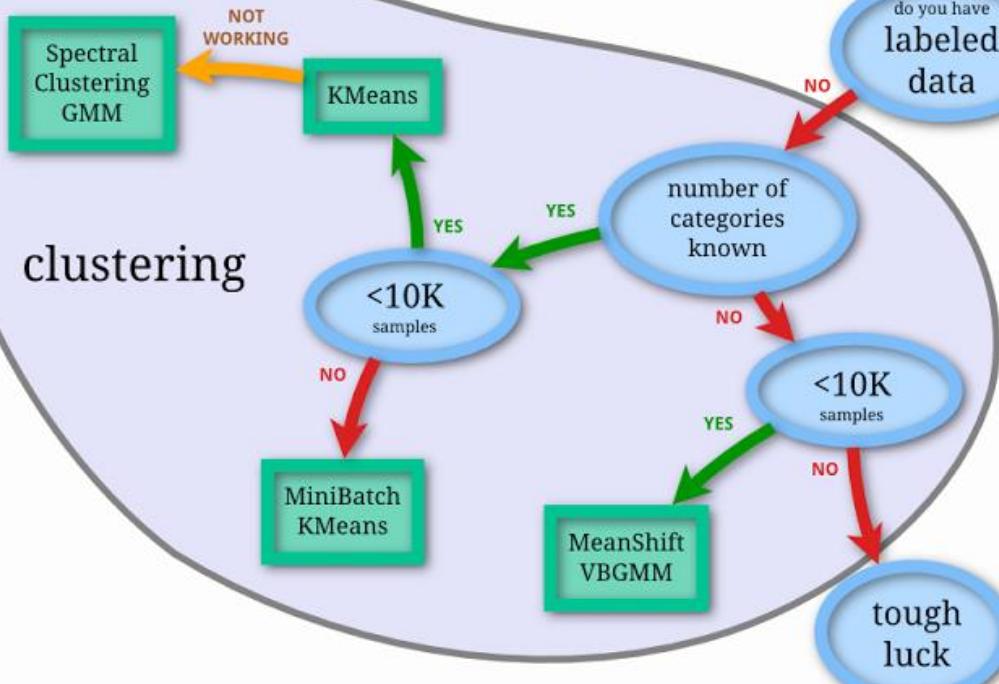
- If you make **no assumption about the data**, there is no reason to prefer one model over any other
- No model that is **a priori guaranteed to work better**
- In practice, **impossible to test all the models**, so make some **reasonable assumptions** about the data and evaluate only a few models

scikit-learn algorithm cheat-sheet

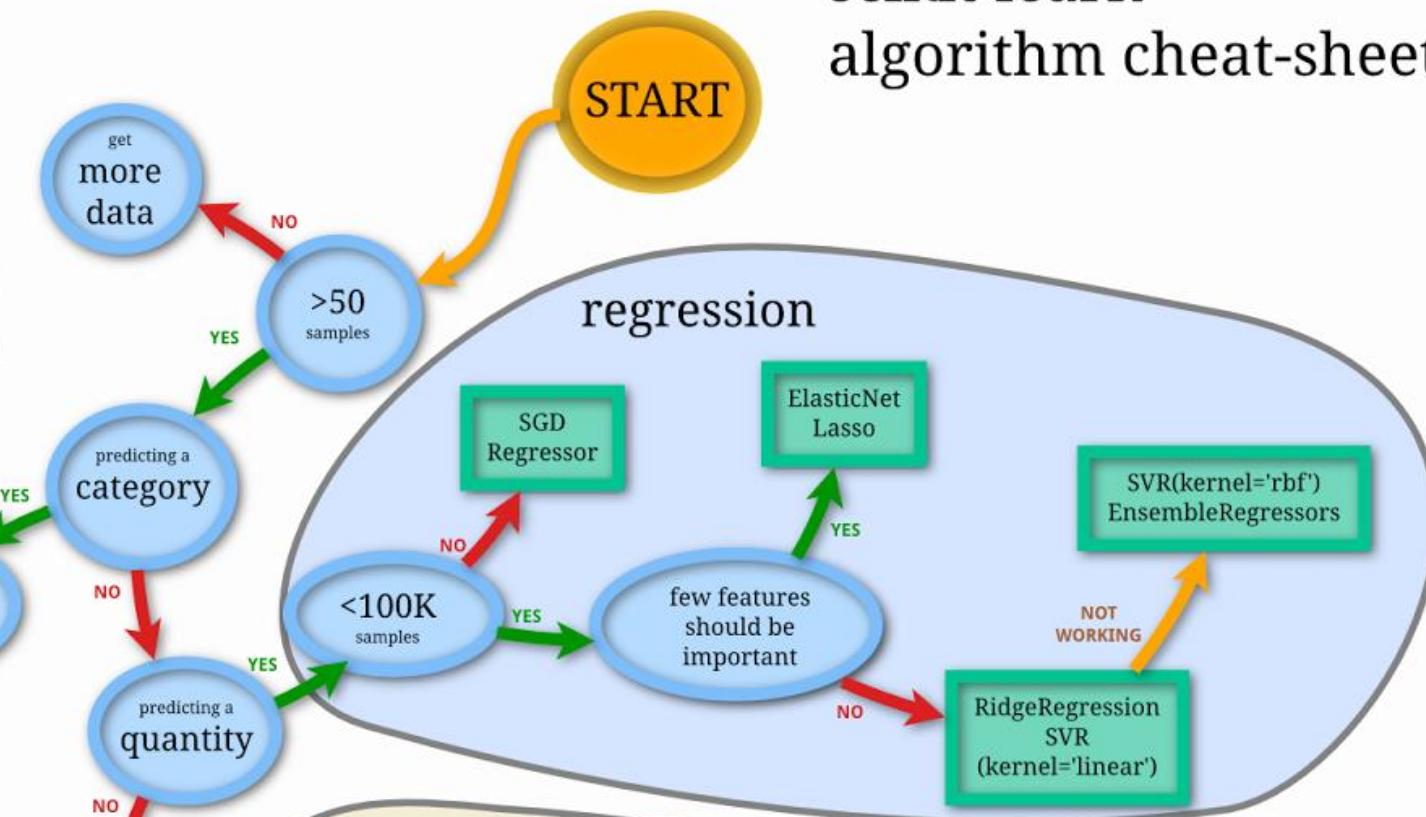
classification



clustering



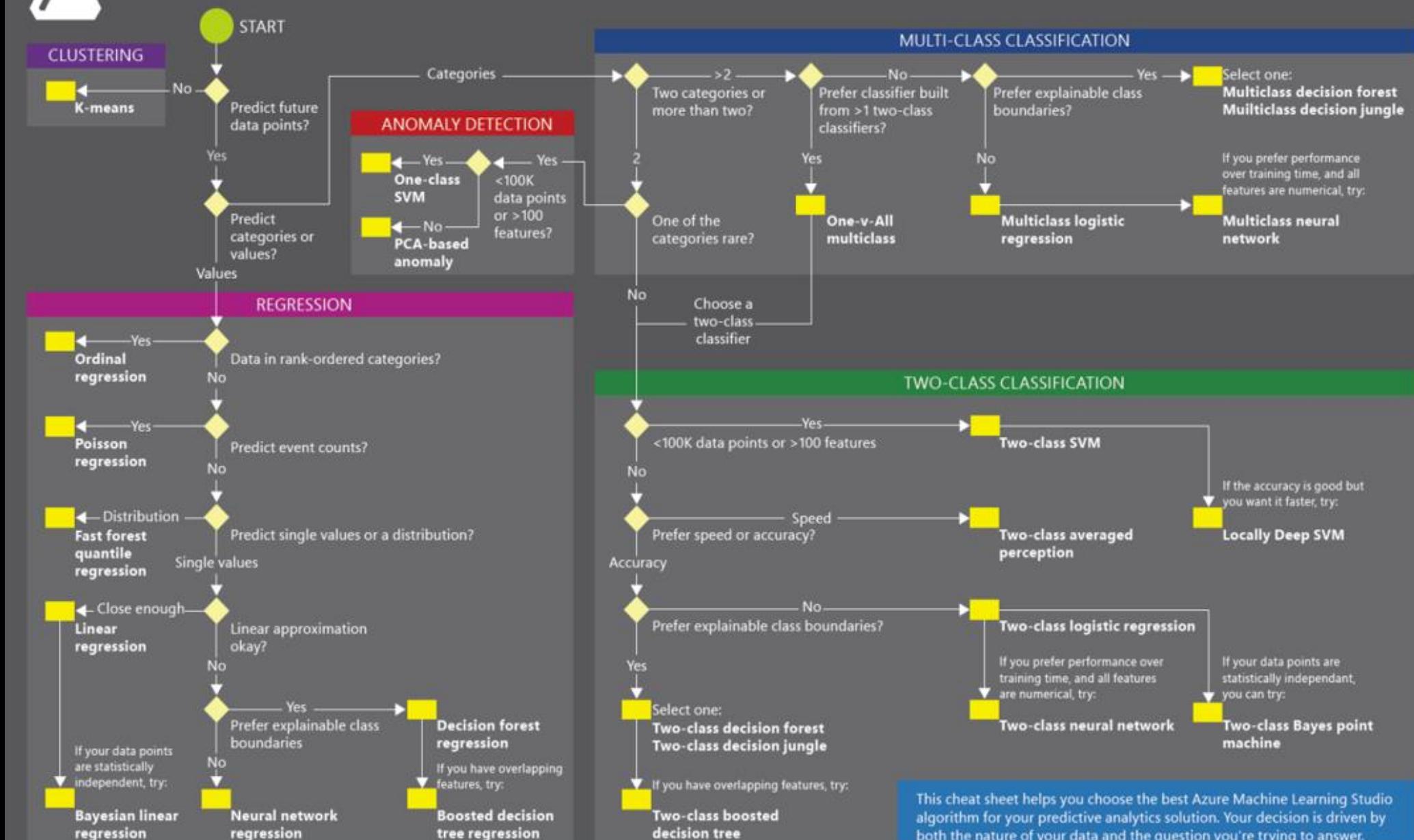
regression



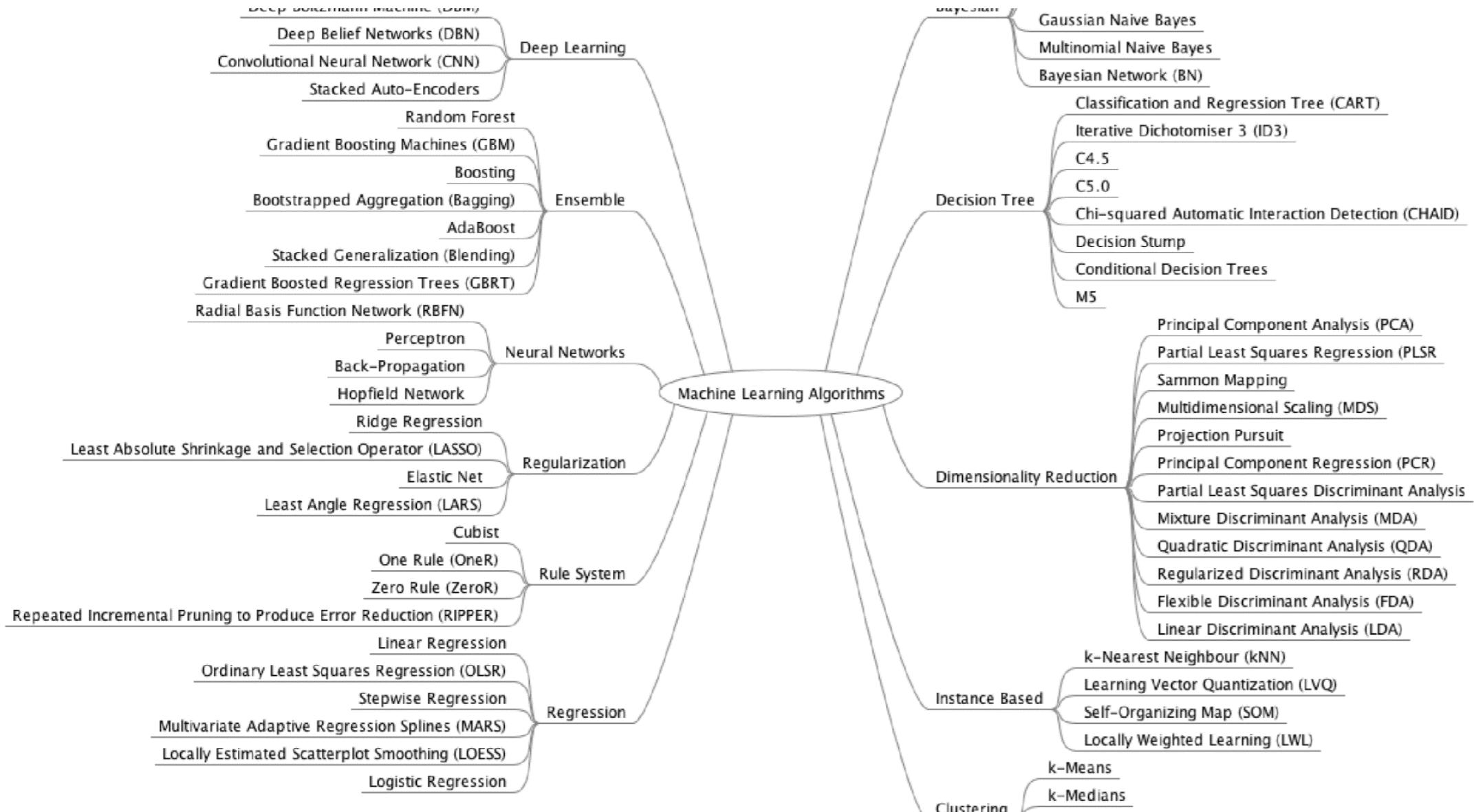
dimensionality reduction

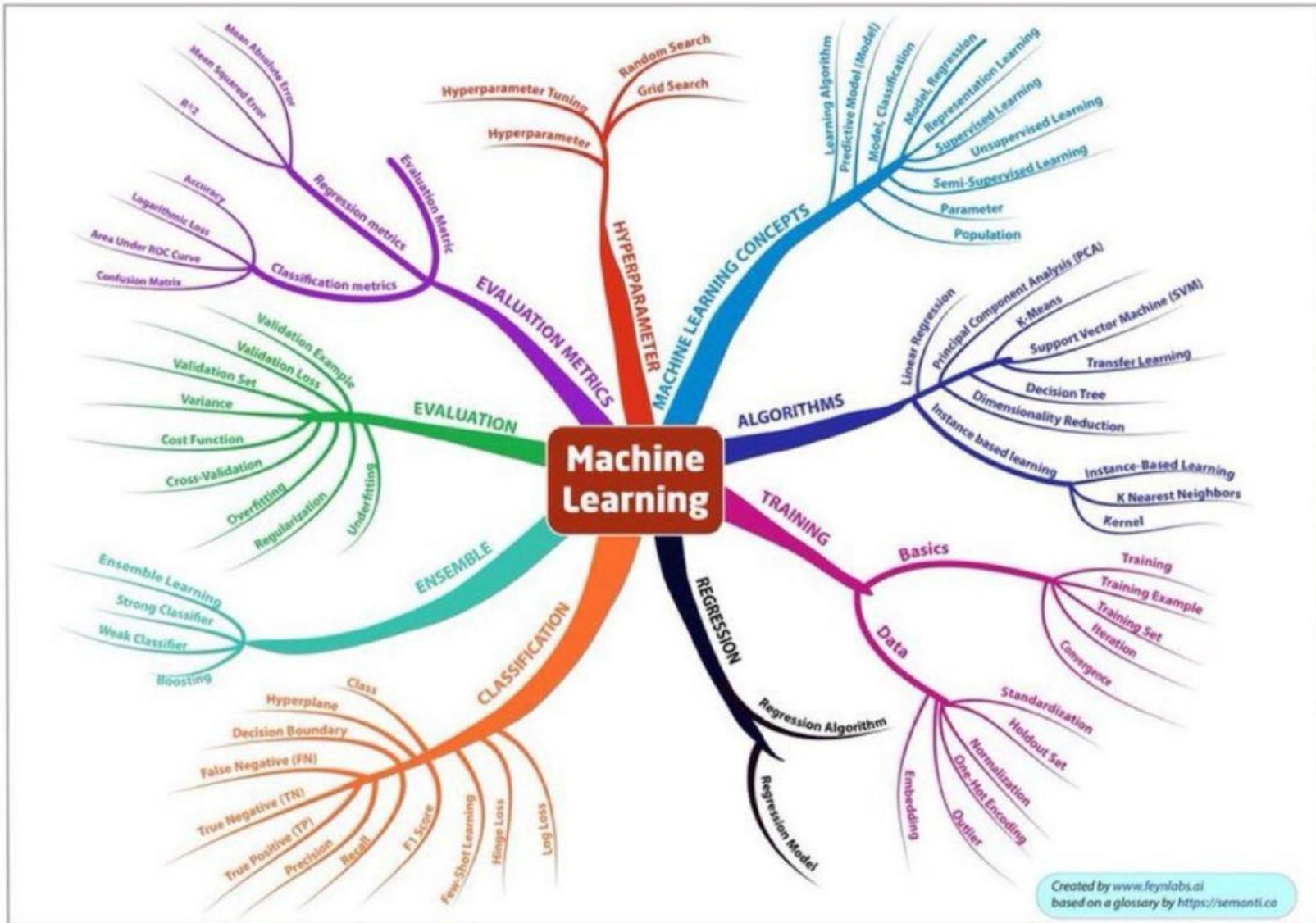


Microsoft Azure Machine Learning: Algorithm Cheat Sheet



This cheat sheet helps you choose the best Azure Machine Learning Studio algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the question you're trying to answer.





General Guidelines

- **Data Amount (number of instances):**
 - More data helps
 - Complex models require more data
 - **Rule of Thumb:** you should have at least 10-20 times as many data points as you have model parameters. This is not a strict rule, but it can help guide you in selecting an appropriate model complexity
- **Feature amount (Number of Dimensions/Features):**
 - **Curse of dimensionality:** As the number of features increases, the amount of data required to generalize well increases exponentially
 - **Feature Selection and Dimensionality Reduction :** reduce the number of dimensions and focus on the most informative ones
 - **Regularization:** help constrain the model's parameters and prevent it from overfitting
 - **Complex Models and Features:** start with simpler models. Complex models with many features can easily overfit the training data and struggle to generalize

The Lesson

The message is always to fit a simple model first, and then only adopt a more complicated ML model if the extra predictive accuracy (value) it provides is worth it.





Exercise

- Create your own summary (graph) for :
 - Regression
 - Classification

A large black circle with a white border, centered on a dark blue background. The word "Quiz" is written in white inside the circle.

Quiz

<https://b.socrative.com/login/student/>

Room : CONTI6128