



Introduction to Machine Learning

Chakrit Phain

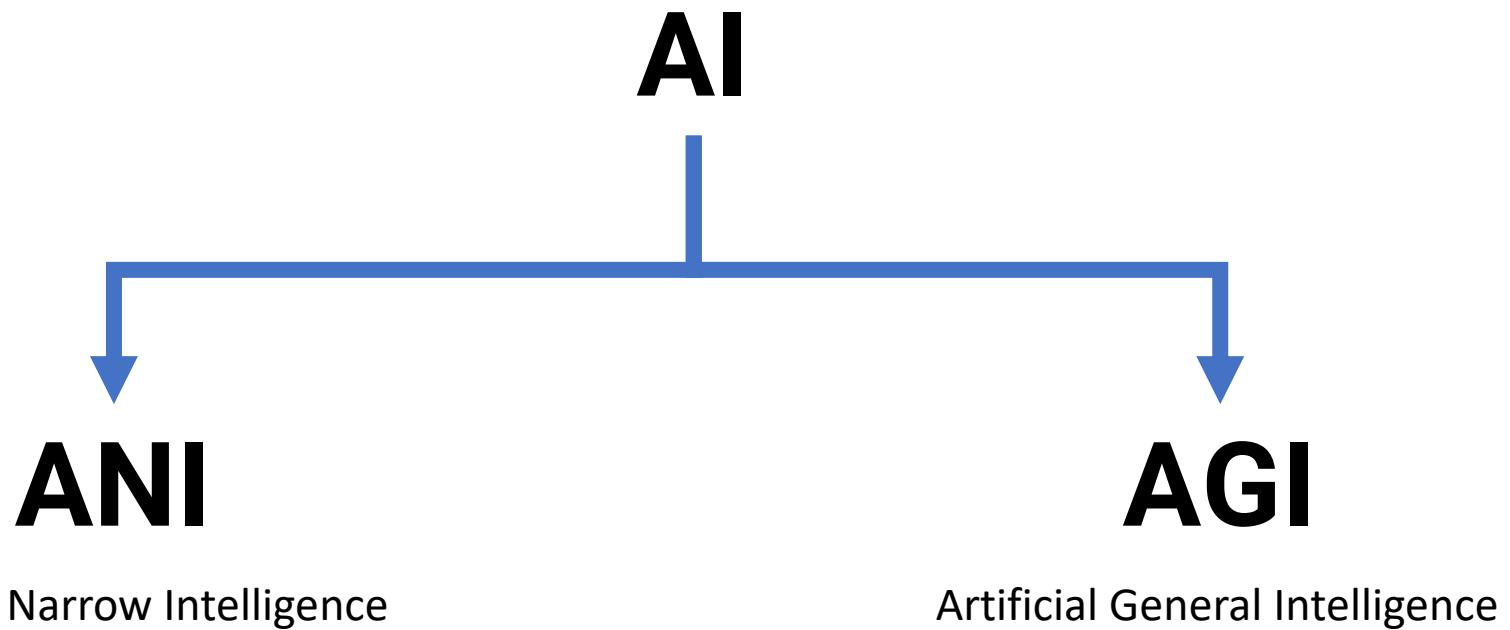
Runchana Laohwawat





What is AI?





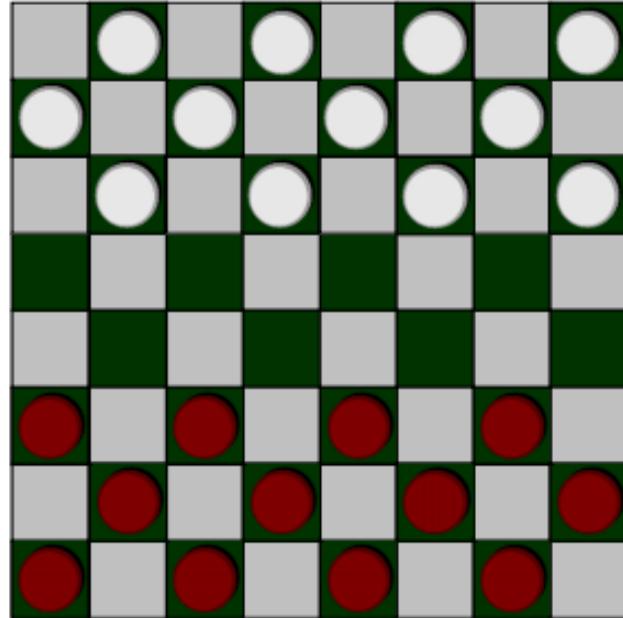
Artificial Narrow Intelligence

Artificial General Intelligence

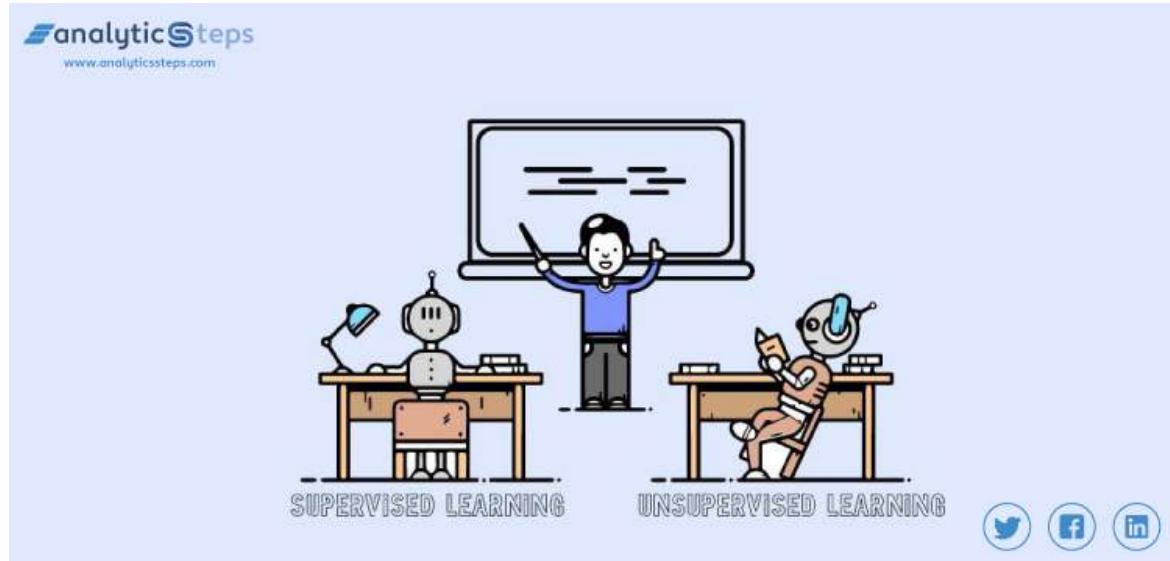


What is Machine Learning?

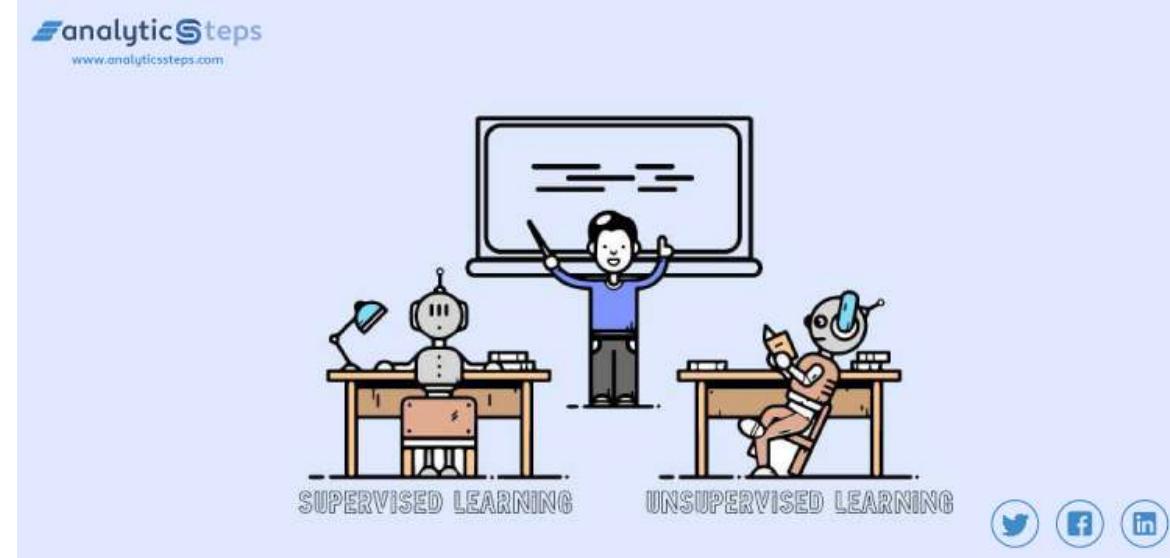




Machine Learning is the field of study that gives the computer the ability to learn without being explicitly programmed.
- Arthur Samuel (1959)



Machine Learning is any process by which a system improves performance from experience.
- Herbert Simon (Turing Award 1975)



Machine learning (ML) is the study of computer algorithms that improve automatically through experience.

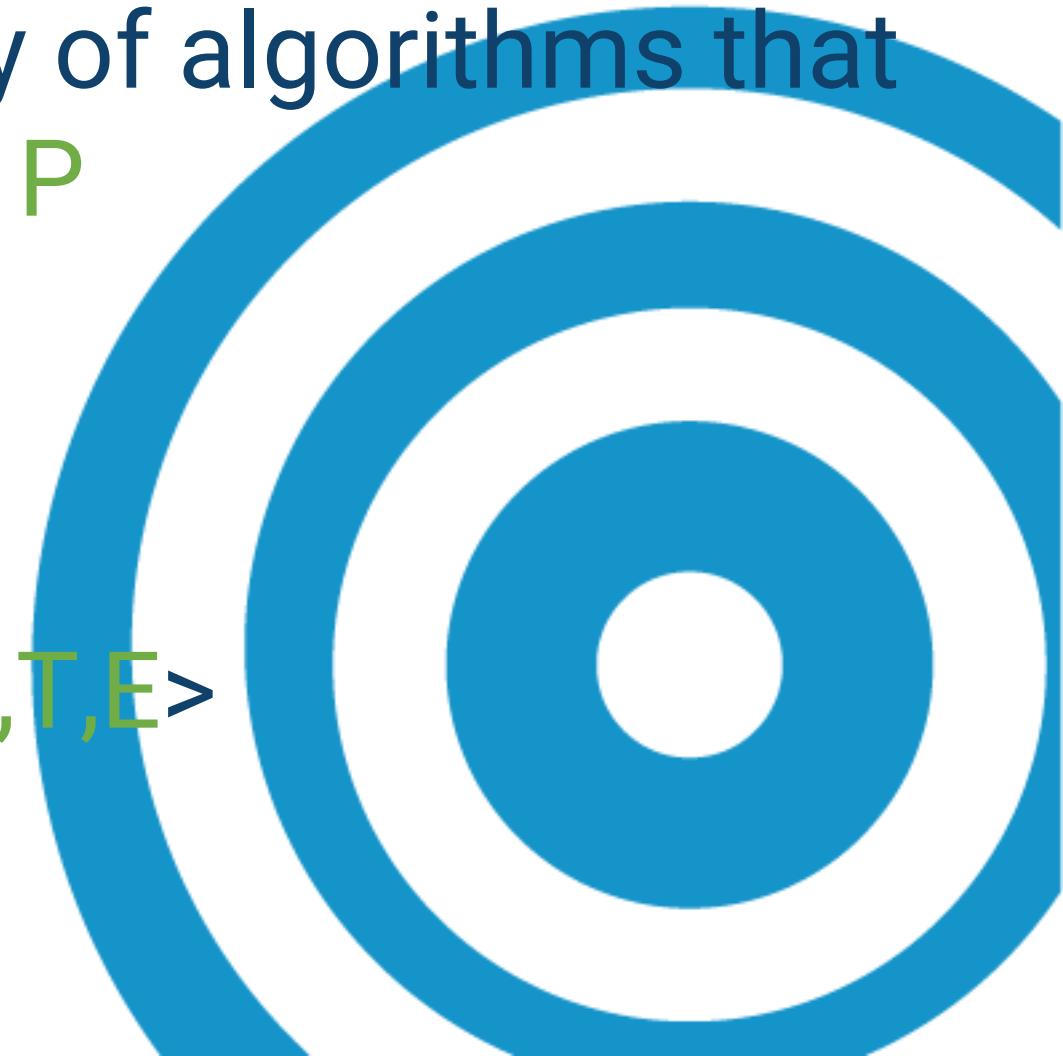
- Tom Mitchell (1998)

Definition by Tom Mitchell (1998)

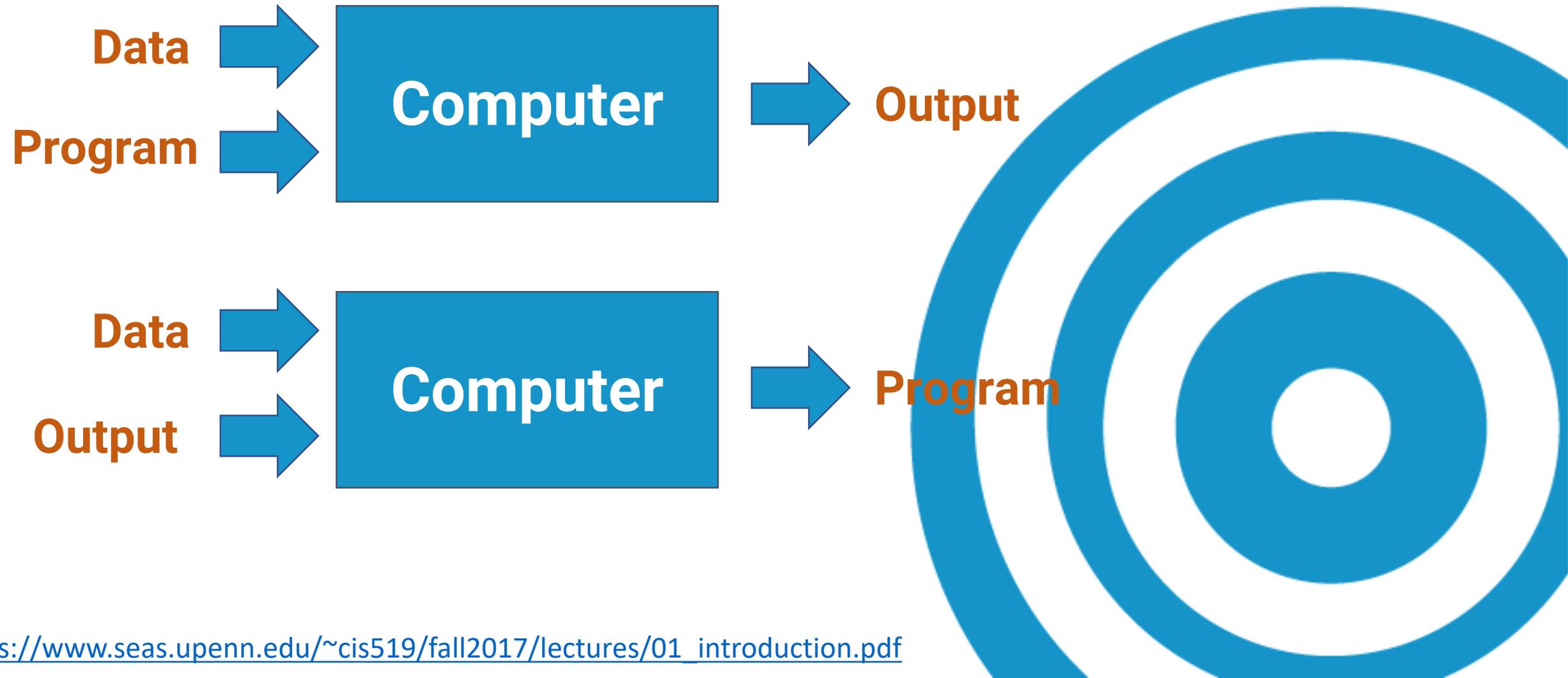
Machine Learning is the study of algorithms that

- improve their performance P
- at some task T
- with experience E

well-defined learning task: <P,T,E>

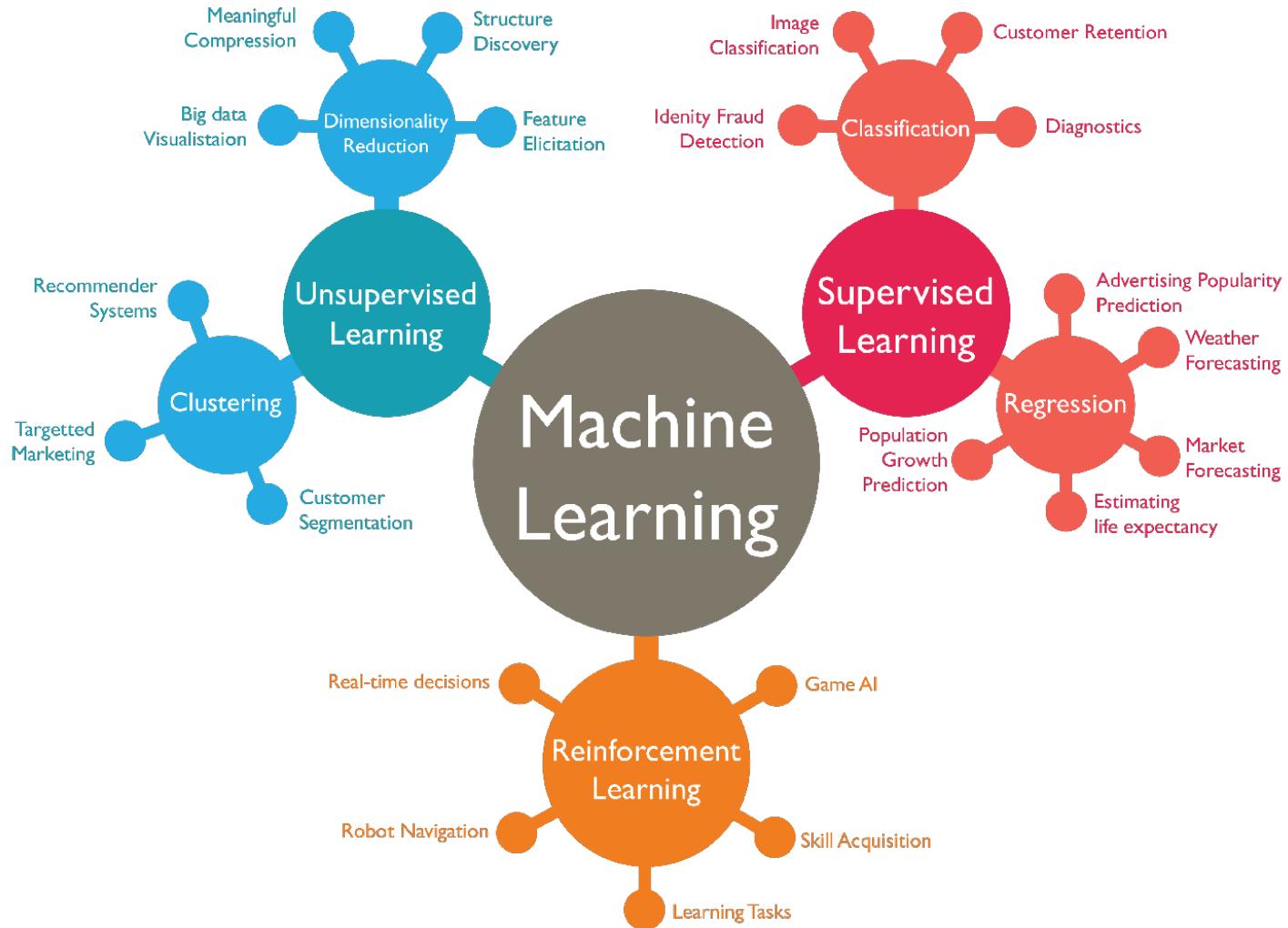


Machine Learning Definition



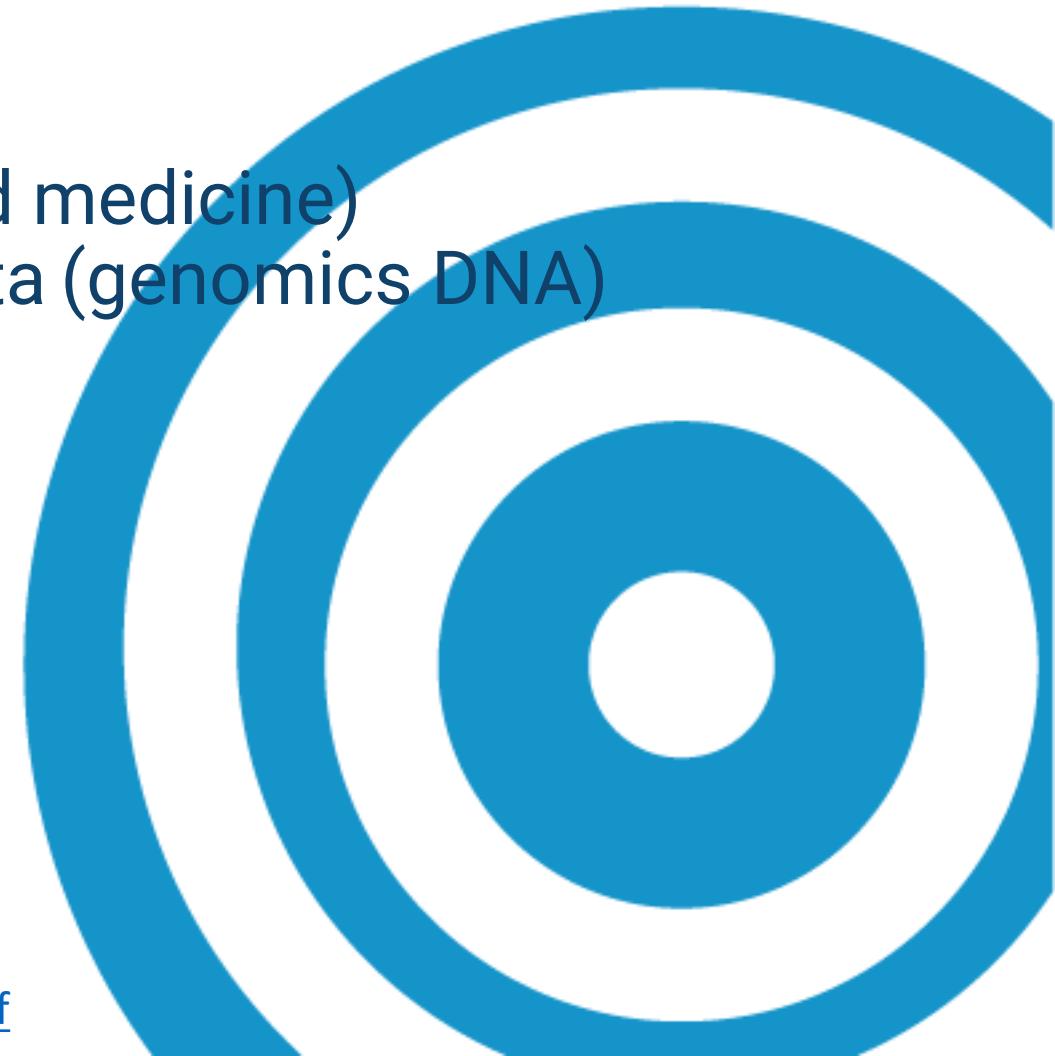


Machine Learning



When Do We Use Machine Learning?

- Human expertise does not exist
- Humans can't explain their expertise
- Model must be customized (personalized medicine)
- Models are based on huge amounts of data (genomics DNA)



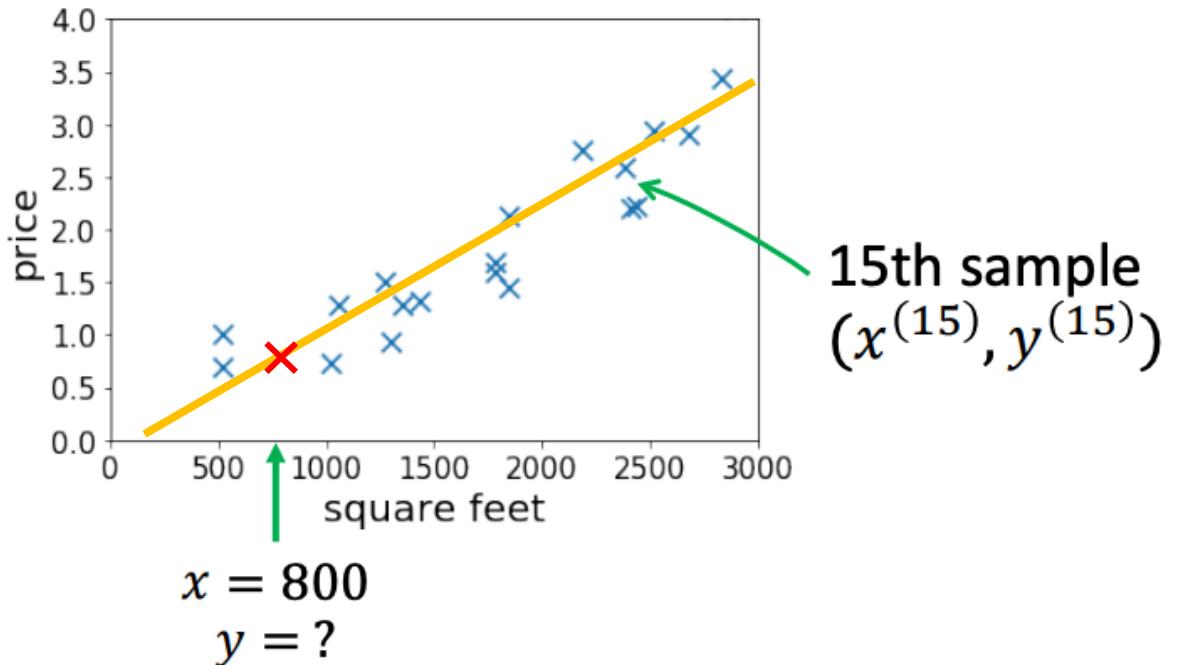


Supervised Learning



Housing Price Prediction

- Given: a dataset that contains n samples
 $(x^{(1)}, y^{(1)}), \dots (x^{(n)}, y^{(n)})$
- Task: if a residence has x square feet, predict its price?

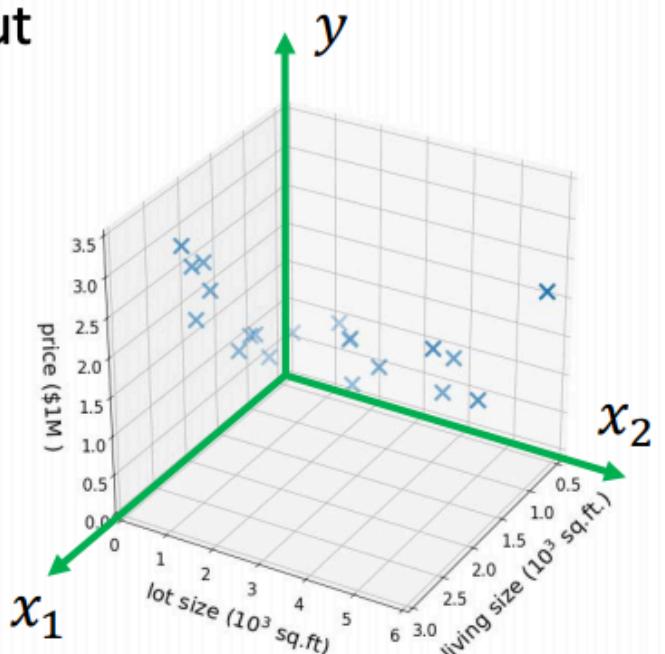


Housing Price Prediction

- Suppose we also know the lot size
- Task: find a function that maps

$$\begin{array}{ccc} (\text{size, lot size}) & \rightarrow & \text{price} \\ \underbrace{\phantom{\text{size, lot size}}}_{\text{features/input}} & & \underbrace{\phantom{\text{price}}}_{\text{label/output}} \\ x \in \mathbb{R}^2 & & y \in \mathbb{R} \end{array}$$

- Dataset: $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$
- where $x^{(i)} = (x_1^{(i)}, x_2^{(i)})$
- “Supervision” refers to $y^{(1)}, \dots, y^{(n)}$



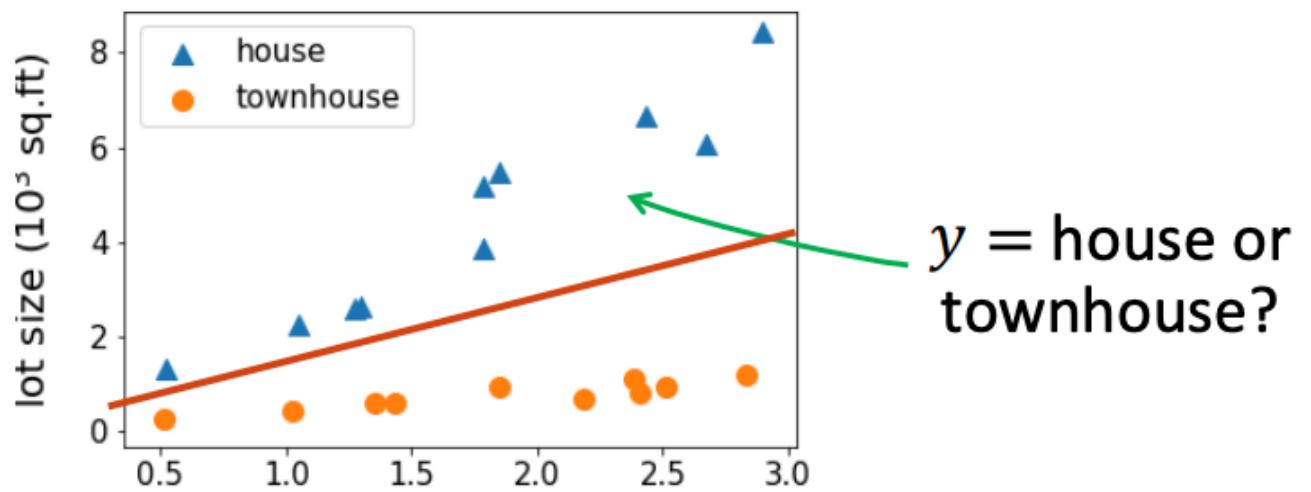
Housing Price Prediction

- $x \in \mathbb{R}^d$ for large d
- E.g.,

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ \vdots \\ x_d \end{bmatrix} \begin{array}{l} \text{--- living size} \\ \text{--- lot size} \\ \text{--- \# floors} \\ \text{--- condition} \\ \text{--- zip code} \\ \vdots \end{array} \longrightarrow y \text{ --- price}$$

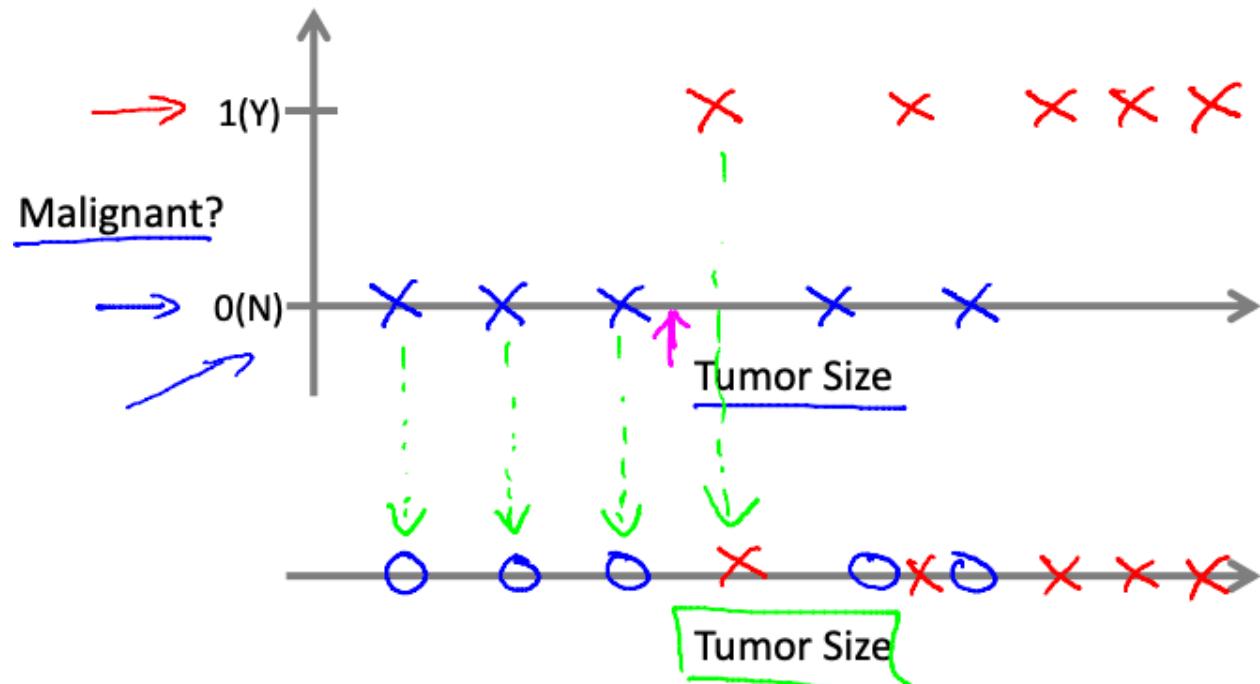
Regression vs Classification

- regression: if $y \in \mathbb{R}$ is a continuous variable
 - e.g., price prediction
- classification: the label is a discrete variable
 - e.g., the task of predicting the types of residence
(size, lot size) → house or townhouse?



Lecture 3&4:
classification

Breast cancer

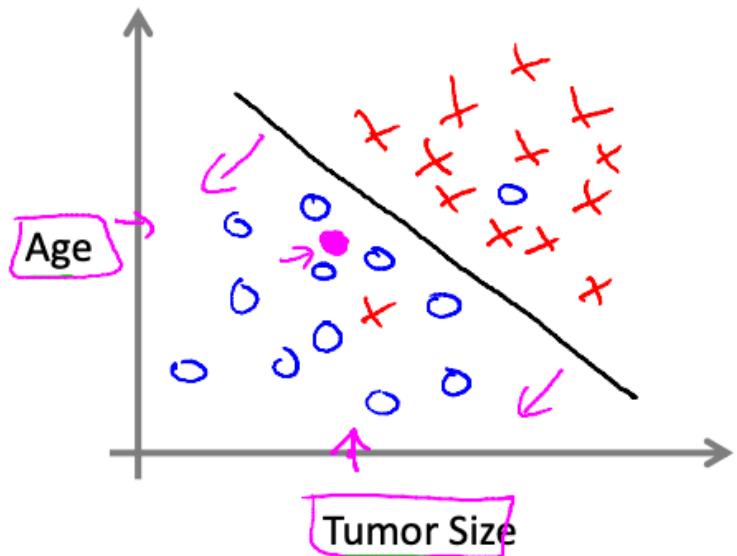


Classification

Discrete valued output (0 or 1)

0, 1, 2, 3
benign type 1 cancer

Breast cancer



- Clump Thickness
- Uniformity of Cell Size
- Uniformity of Cell Shape
- ...

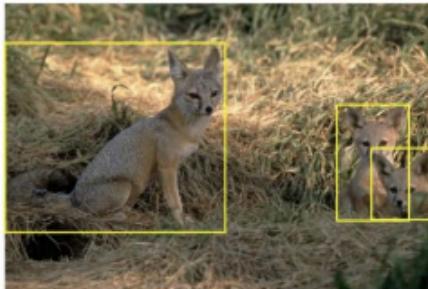
Supervised Learning in Computer Vision

- Image Classification
 - x = raw pixels of the image, y = the main object

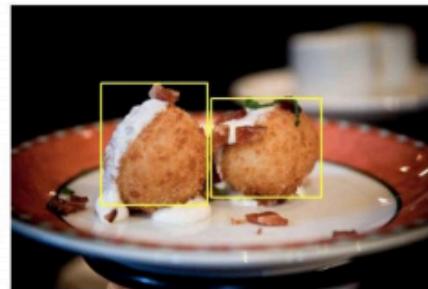


Supervised Learning in Computer Vision

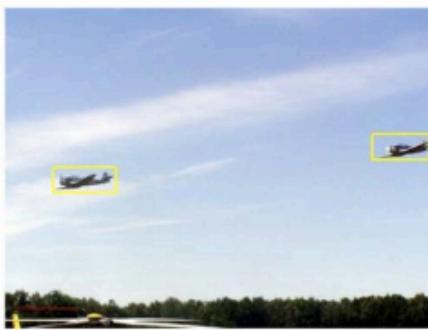
- Object localization and detection
 - x = raw pixels of the image, y = the bounding boxes



kit fox



croquette



airplane



frog

ImageNet Large Scale Visual Recognition Challenge. Russakovsky et al.'2015



Unsupervised Learning





Supervised Learning in Natural Language Processing

Machine translation

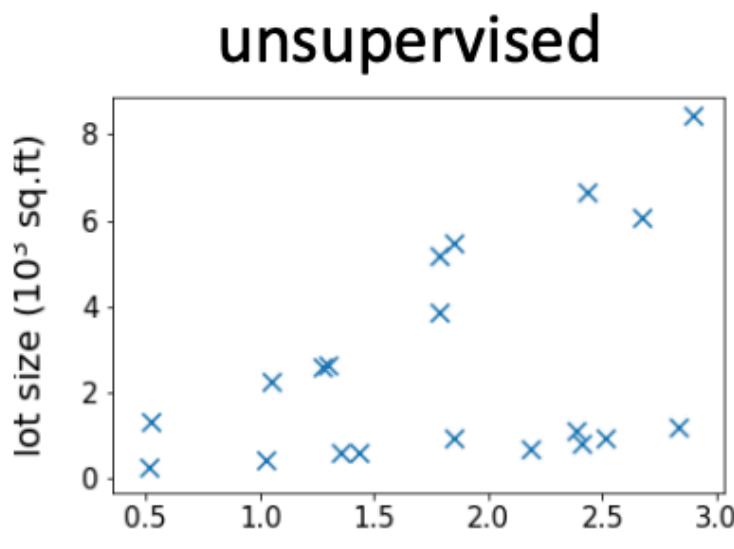
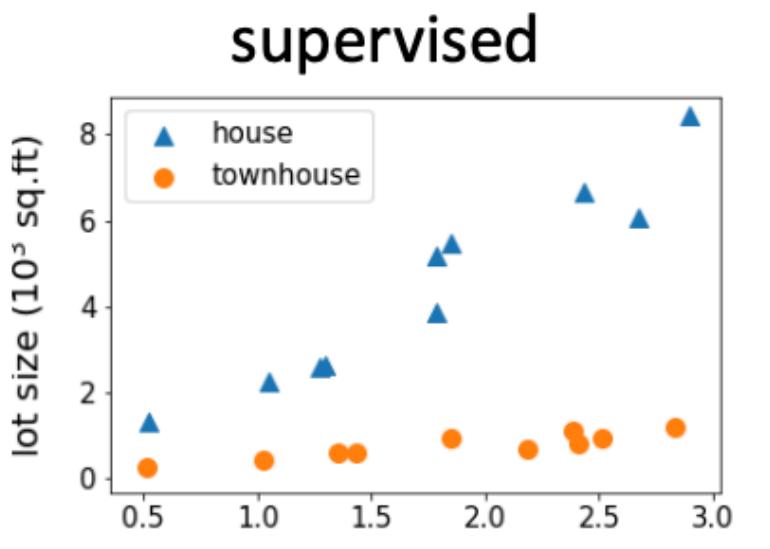
Google แปลภาษา

The screenshot shows the Google Translate interface. The source text 'Hello world' is in English, and the target language is set to Thai ('ไทย'). The translated text is 'สวัสดีชาวโลก' (Swāsdi chāw lok) with a phonetic transcription below it. The interface includes standard translation controls like microphone, speaker, and share icons.

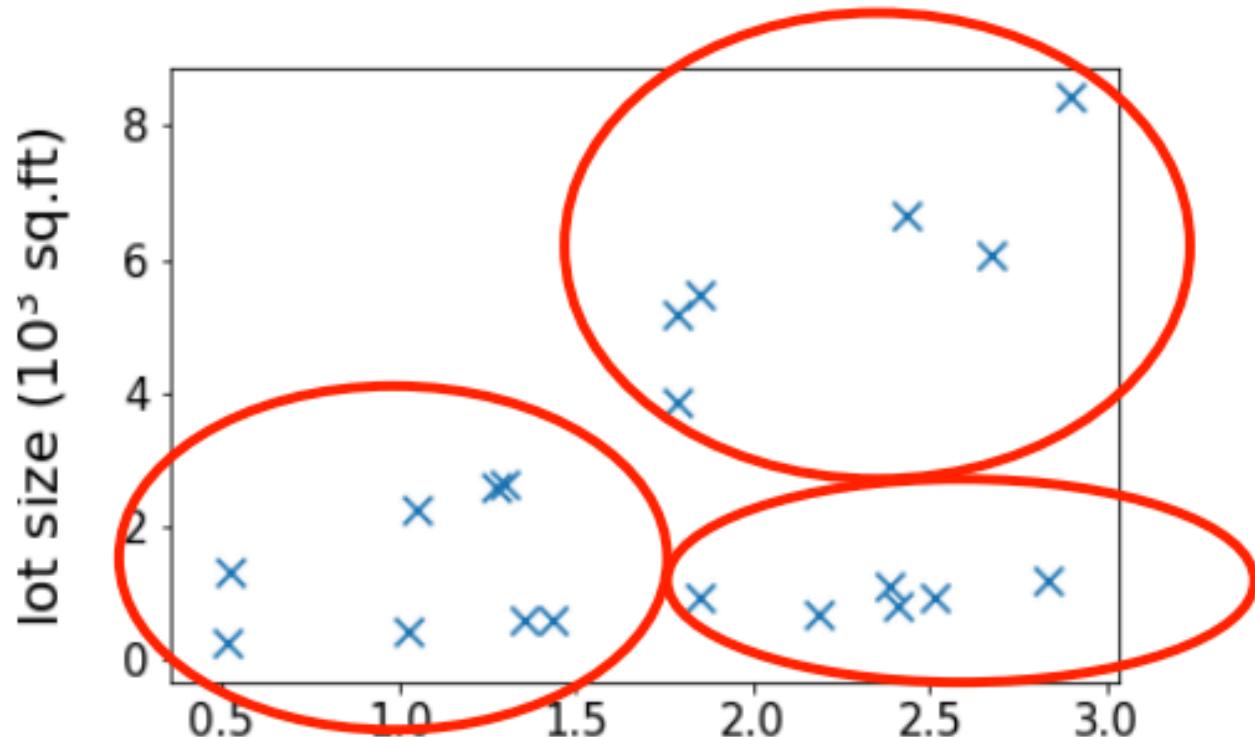
ผู้รายงานติดต่อที่นี่

Unsupervised Learning

- Dataset contains **no labels**: $x^{(1)}, \dots x^{(n)}$
- **Goal** (vaguely-posed): to find interesting structures in the data



Clustering



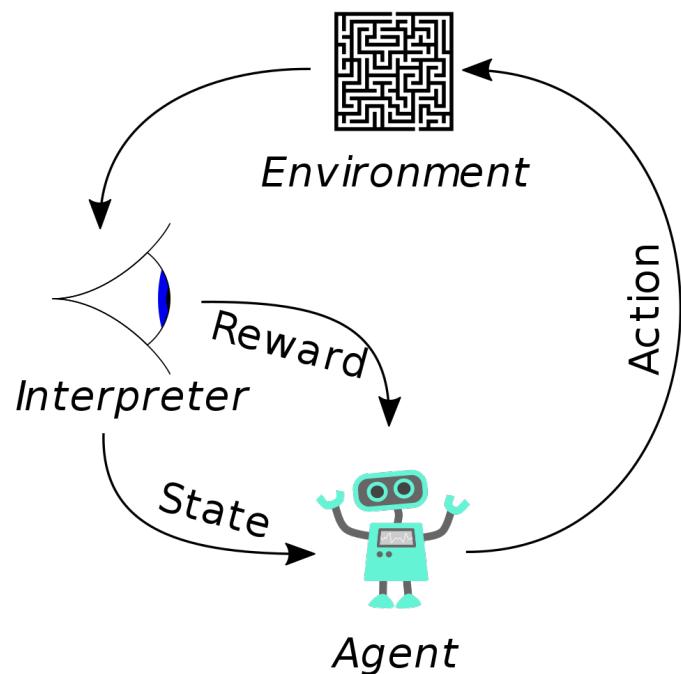


Reinforcement Learning



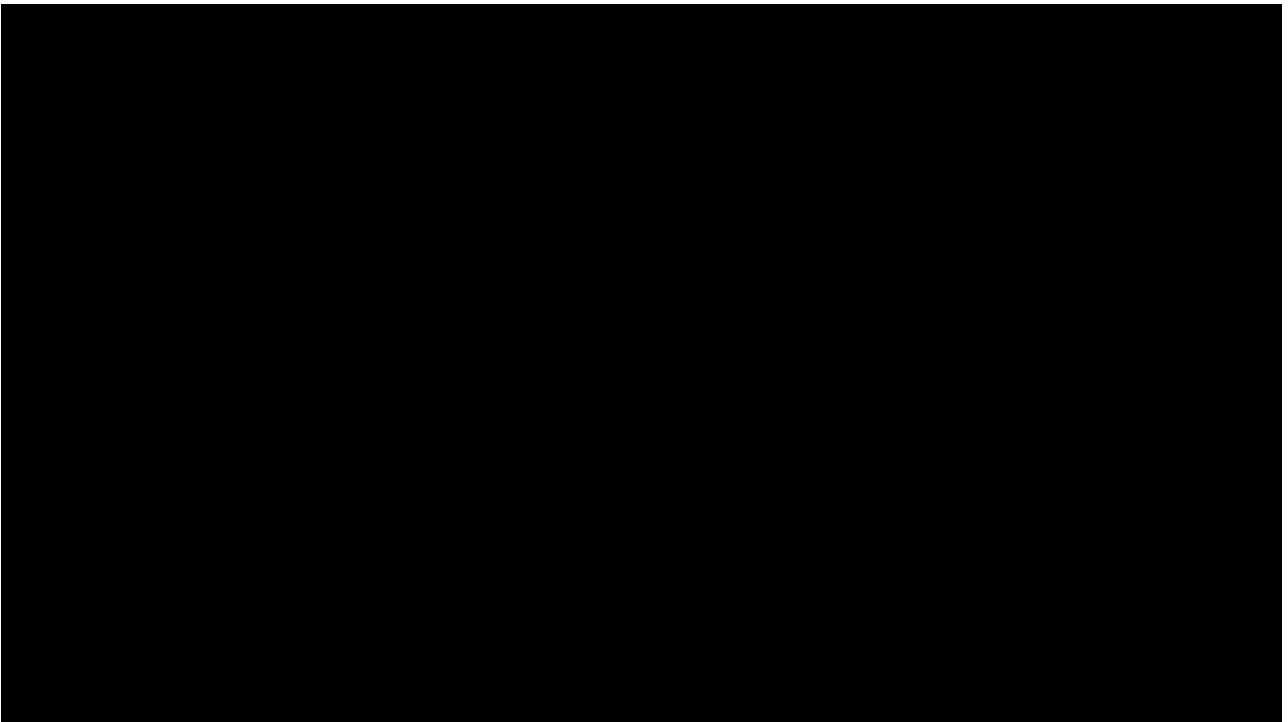
Reinforcement Learning

Reinforcement learning (RL) is an area of machine learning concerned with how software agents ought to take actions in an environment in order to maximize the notion of cumulative reward





Reinforcement Learning





Machine Learning Process

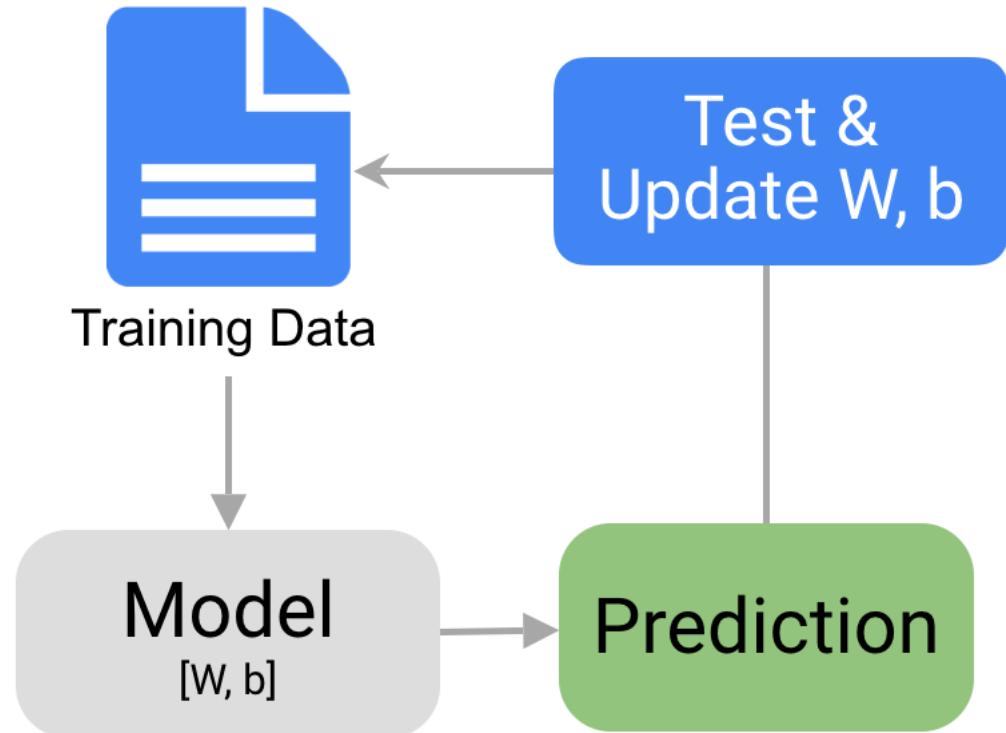




Machine Learning Process



Machine Learning Process



Step 2 Data preparation



Explore Data

Mean	Sum of all values Total number of values
Median	Middle value(when data are arranged in order)
Mode	Most common value

Central tendency
of a distribution

Variance	how far a set of numbers are spread out from mean
Interquartile range	divides a data set into quartiles.
Standard deviation	dispersion of a set of data from mean

Measure of Variation

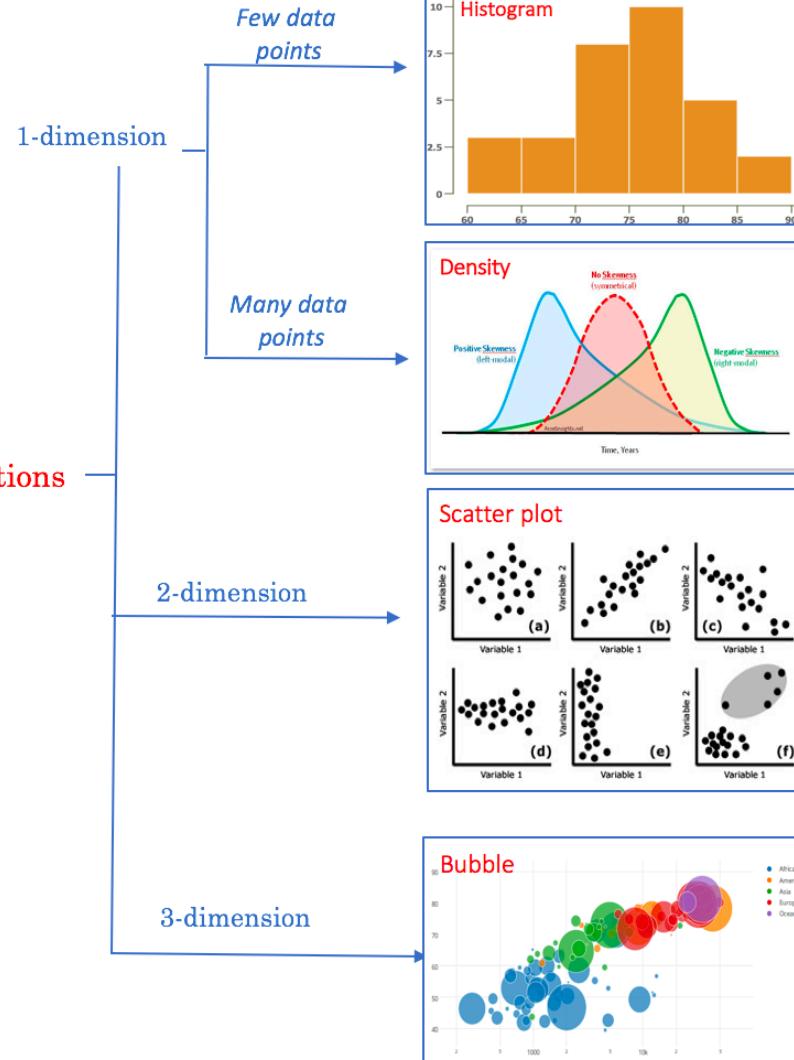
Descriptive statistics

EDA Methods

Visualizations

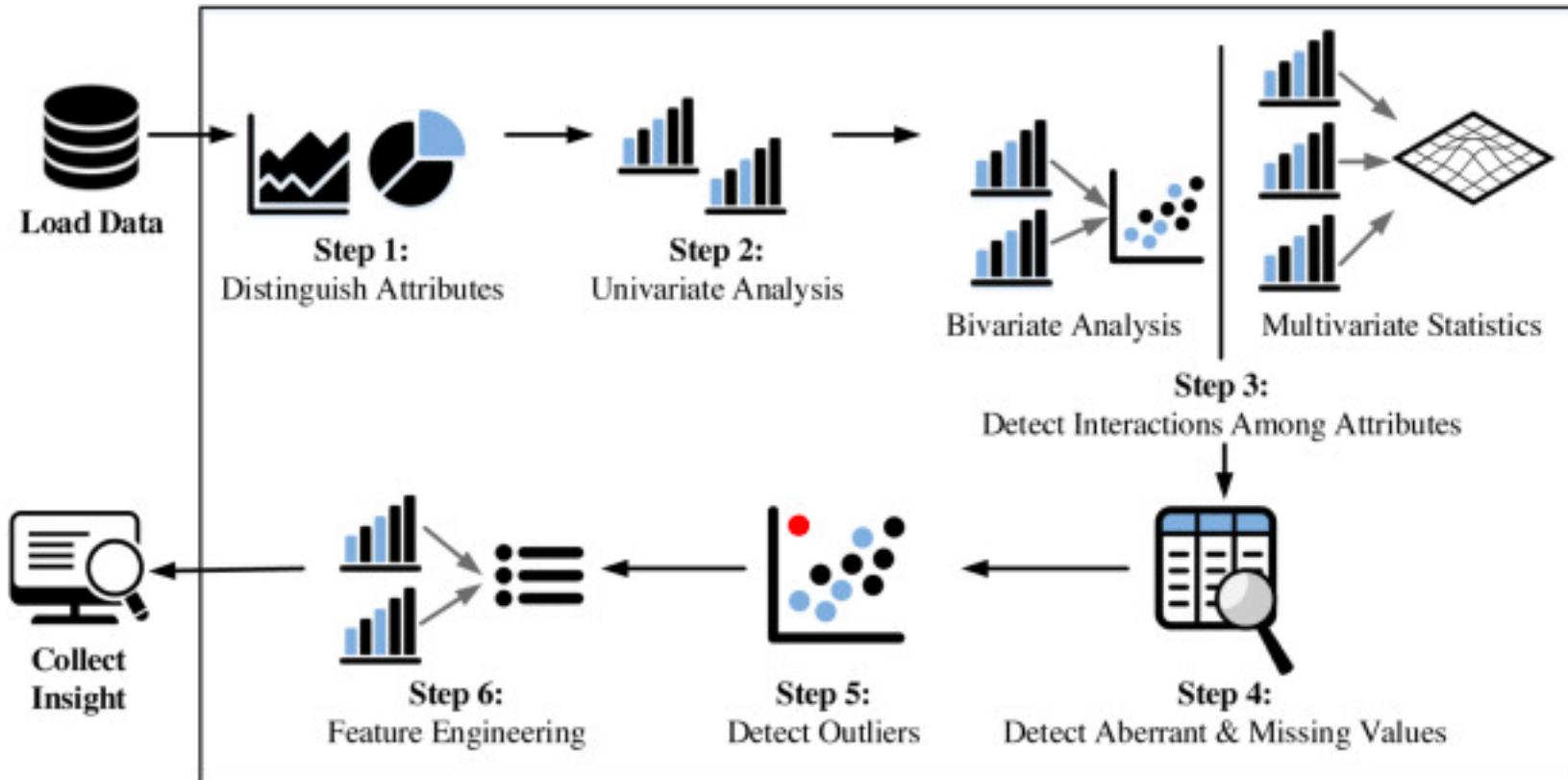
Skewness	Measure of symmetry
Kurtosis	Kurtosis is a measure of “peakedness” relative to a Gaussian shape

Skewness & Kurtosis



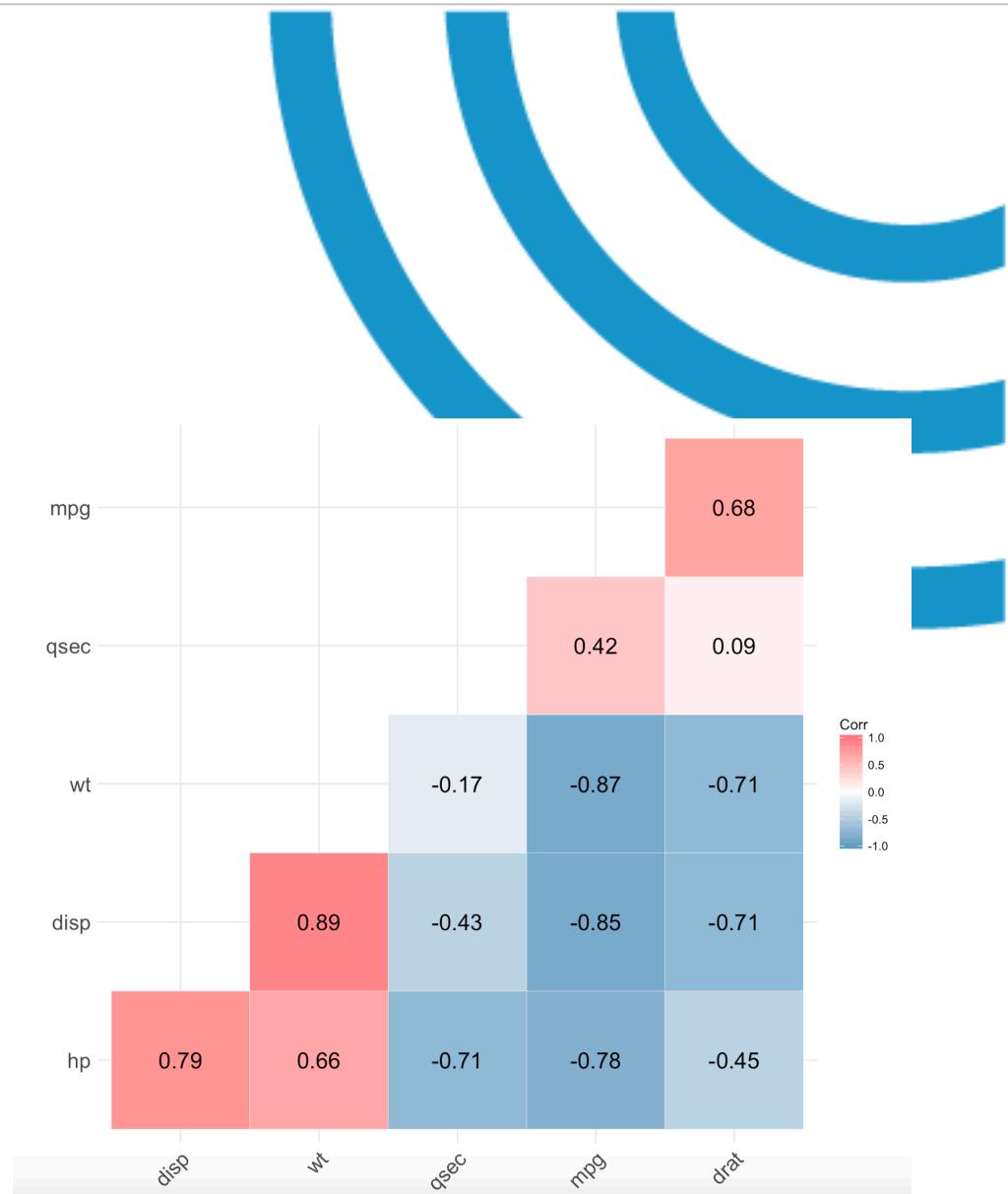
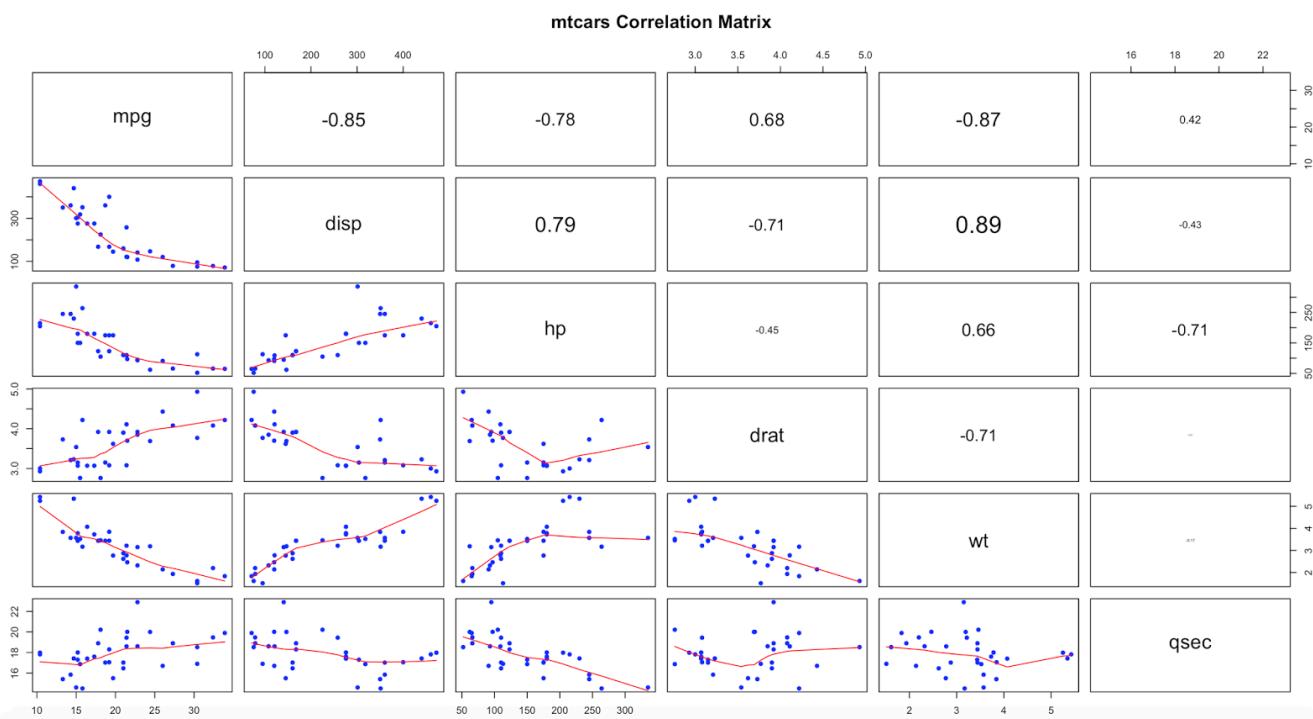


Explore Data



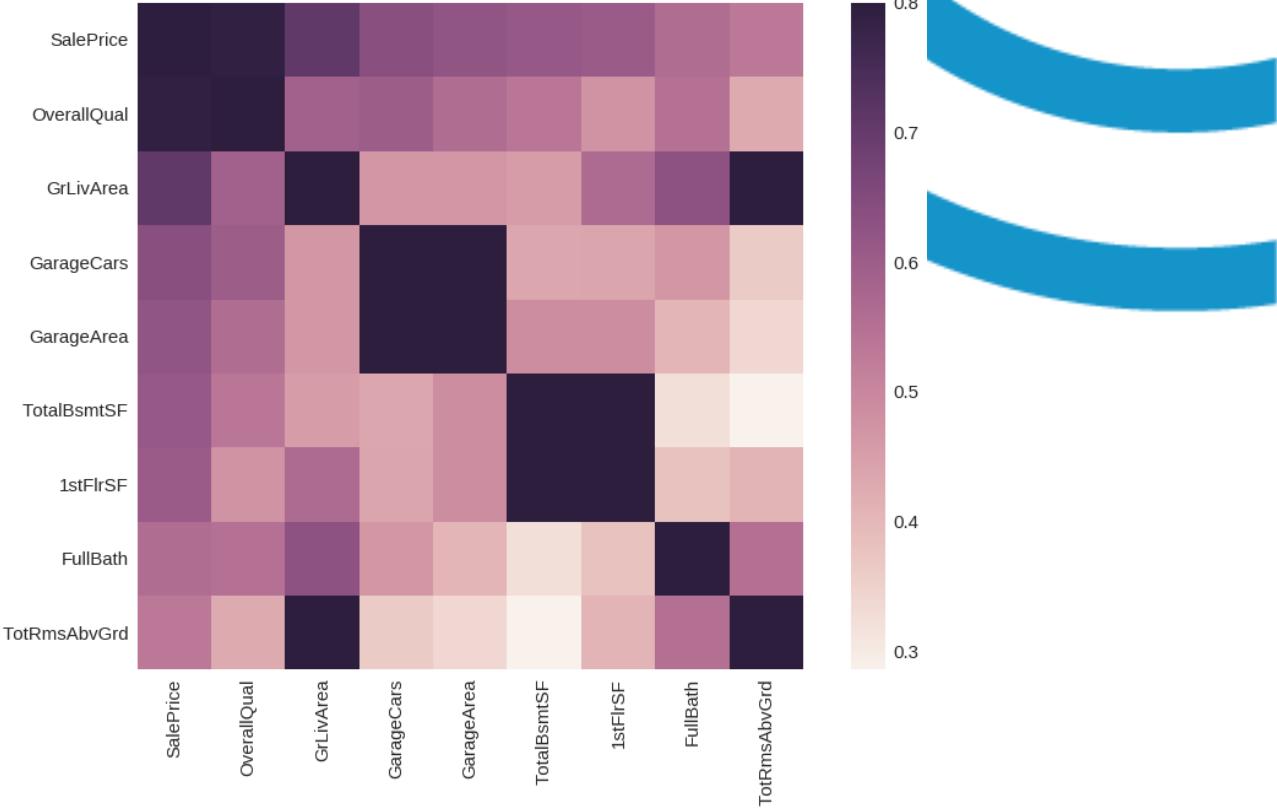
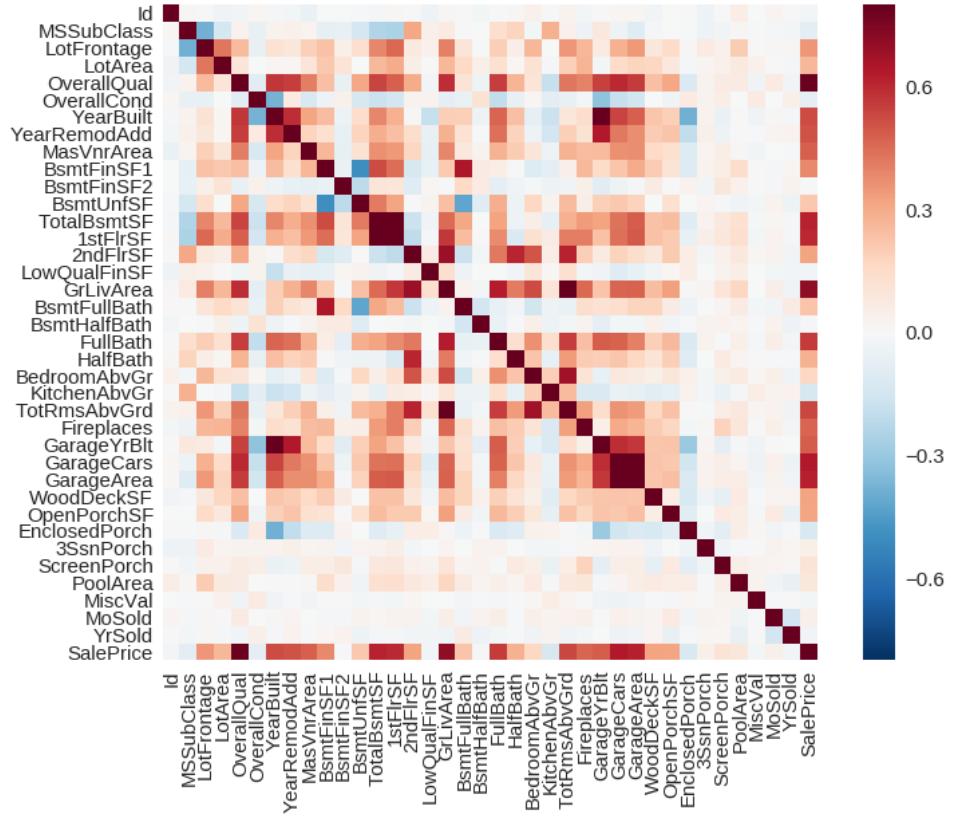


Feature Engineering





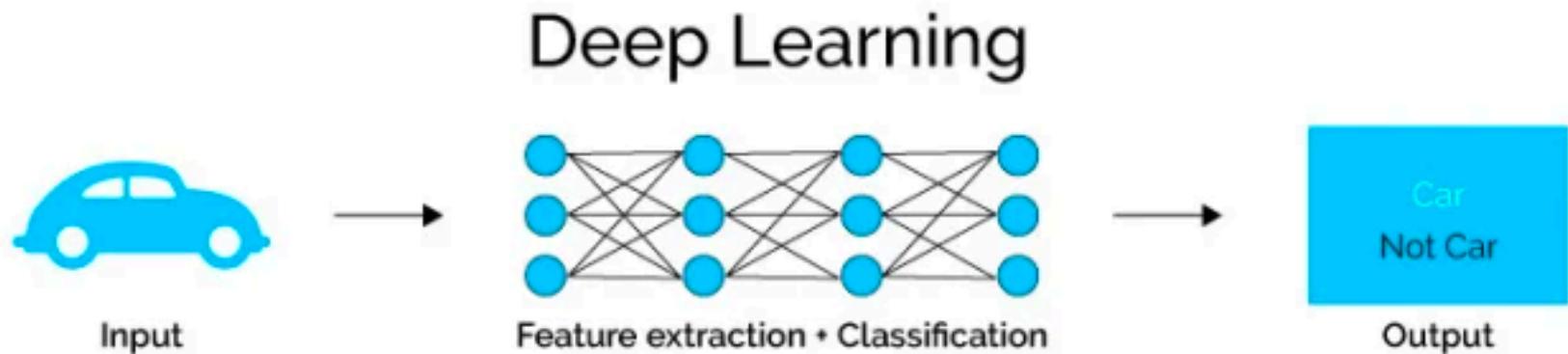
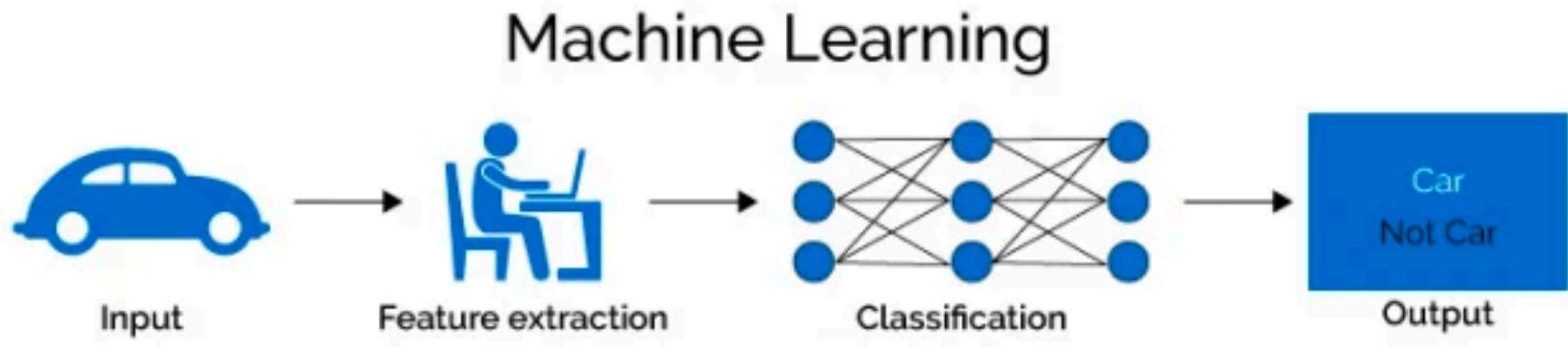
Exploratory Data Analysis Example



<https://towardsdatascience.com/predicting-house-prices-with-linear-regression-machine-learning-from-scratch-part-ii-47a0238aeac1>

<https://colab.research.google.com/drive/1YEgrxbMQCONVeKDfXg49YUyD3lpwlGxu#scrollTo=DWQVJtT2uBSG>

Feature Extraction



Step 3 Choosing Model



MachineLearning Overview

MACHINE LEARNING IN EMOJI

BecomingHuman.AI

SUPERVISED

human builds model based on input / output

UNSUPERVISED

human input, machine output
human utilizes if satisfactory

REINFORCEMENT

human input, machine output
human reward/punish, cycle continues

BASIC REGRESSION

LINEAR

`Linear_model.LinearRegression()`
Lots of numerical data



LOGISTIC

`Linear_model.LogisticRegression()`
Target variable is categorical



CLUSTER ANALYSIS

K-MEANS

`cluster.KMeans()`
Similar datum into groups based on centroids



ANOMALY DETECTION

`covariance.EllipticEnvelope()`
Finding outliers through grouping



CLASSIFICATION



NEURAL NET

`neural_network.MLPClassifier()`

Complex relationships. Prone to overfitting
Basically magic.



K-NN

`neighbors.KNeighborsClassifier()`

Group membership based on proximity



DECISION TREE

`tree.DecisionTreeClassifier()`

If/then/else. Non-contiguous data.
Can also be regression.



RANDOM FOREST

`ensemble.RandomForestClassifier()`

Find best split randomly
Can also be regression



SVM

`svm.SVC()` `svm.LinearSVC()`

Maximum margin classifier. Fundamental Data Science algorithm



NAIVE BAYES

`GaussianNB()` `MultinomialNB()` `BernoulliNB()`

Updating knowledge step by step with new info



FEATURE REDUCTION

T-DISTRIB STOCHASTIC NEIB EMBEDDING



`manifold.TSNE()`
Visual high dimensional data. Convert similarity to joint probabilities

PRINCIPLE COMPONENT ANALYSIS



`decomposition.PCA()`
Distill feature space into components that describe greatest variance

CANONICAL CORRELATION ANALYSIS



`decomposition.CCA()`
Making sense of cross-correlation matrices

LINEAR DISCRIMINANT ANALYSIS



`lda.LDA()`
Linear combination of features that separates classes

OTHER IMPORTANT CONCEPTS

BIAS VARIANCE TRADEOFF

UNDERFITTING / OVERFITTING

INERTIA

ACCURACY FUNCTION

$\frac{TP+TN}{TP+TN+FP+FN}$

Precision Function

$\frac{TP}{TP+FP}$

Specificity Function

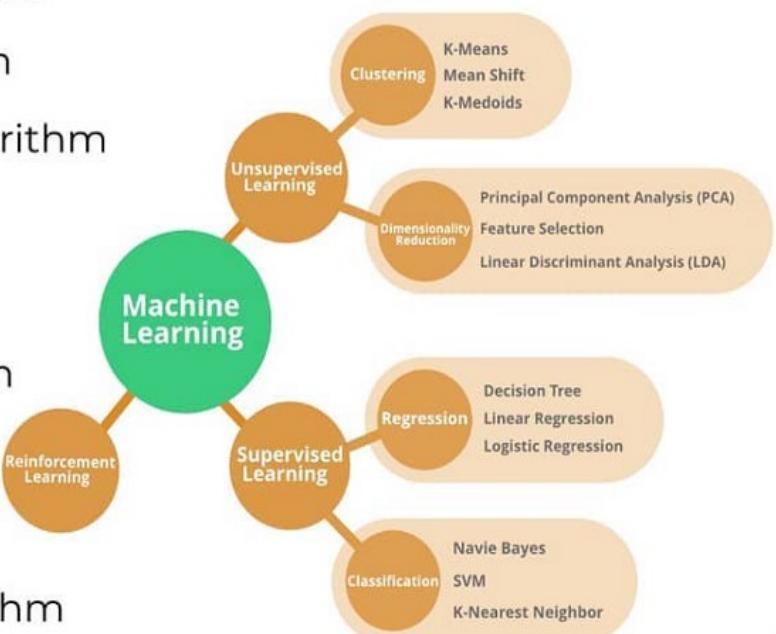
$\frac{TN}{TN+FP}$

Sensitivity Function

$\frac{TP}{TP+FN}$

Top 10 Algorithms every Machine Learning Engineer should know

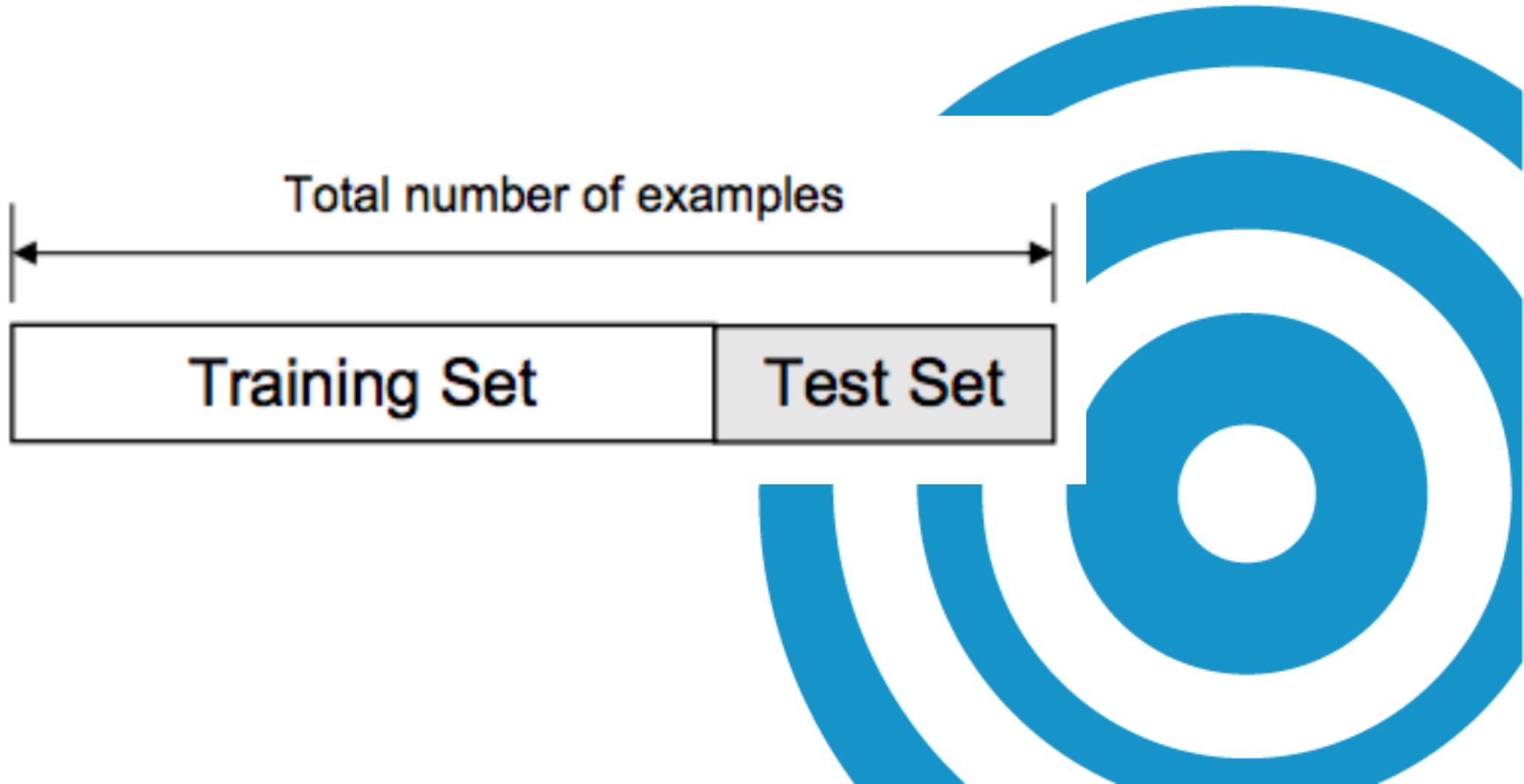
1. Naïve Bayes Classifier Algorithm
2. K Means Clustering Algorithm
3. Support Vector Machine Algorithm
4. Apriori Algorithm
5. Linear Regression Algorithm
6. Logistic Regression Algorithm
7. Decision Trees Algorithm
8. Random Forests Algorithm
9. K Nearest Neighbours Algorithm
10. Artificial Neural Networks Algorithm



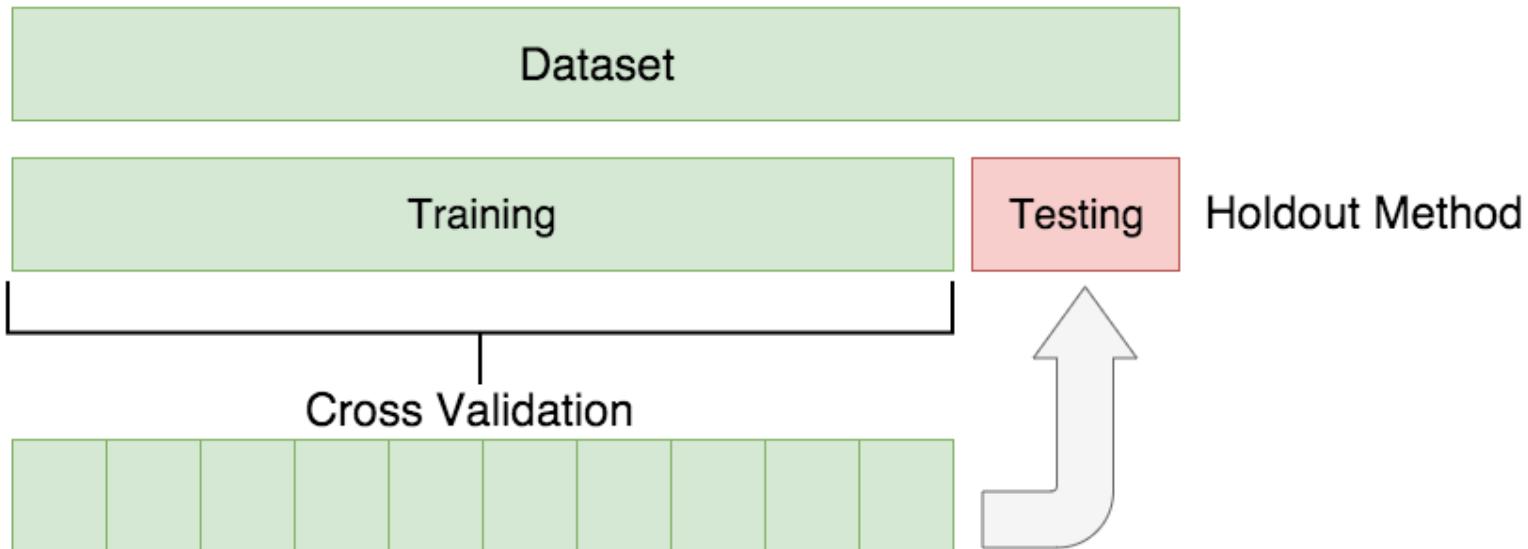
Step 4 Training



Train/Test Split



Cross Validation



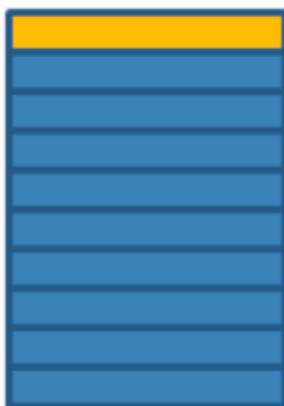
Data Permitting:



K-Folds Cross Validation

Validation Set
Training Set

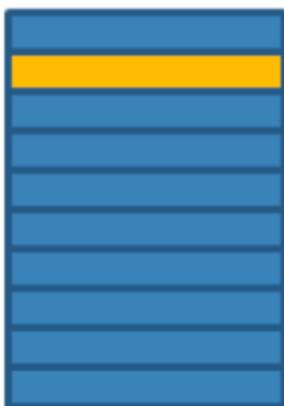
Round 1



Validation
Accuracy:

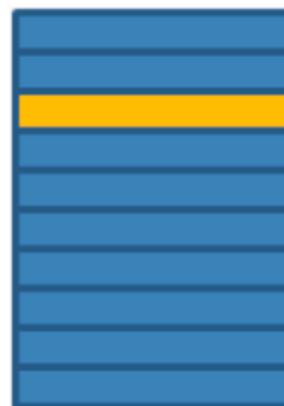
93%

Round 2



90%

Round 3



91%

Round 10

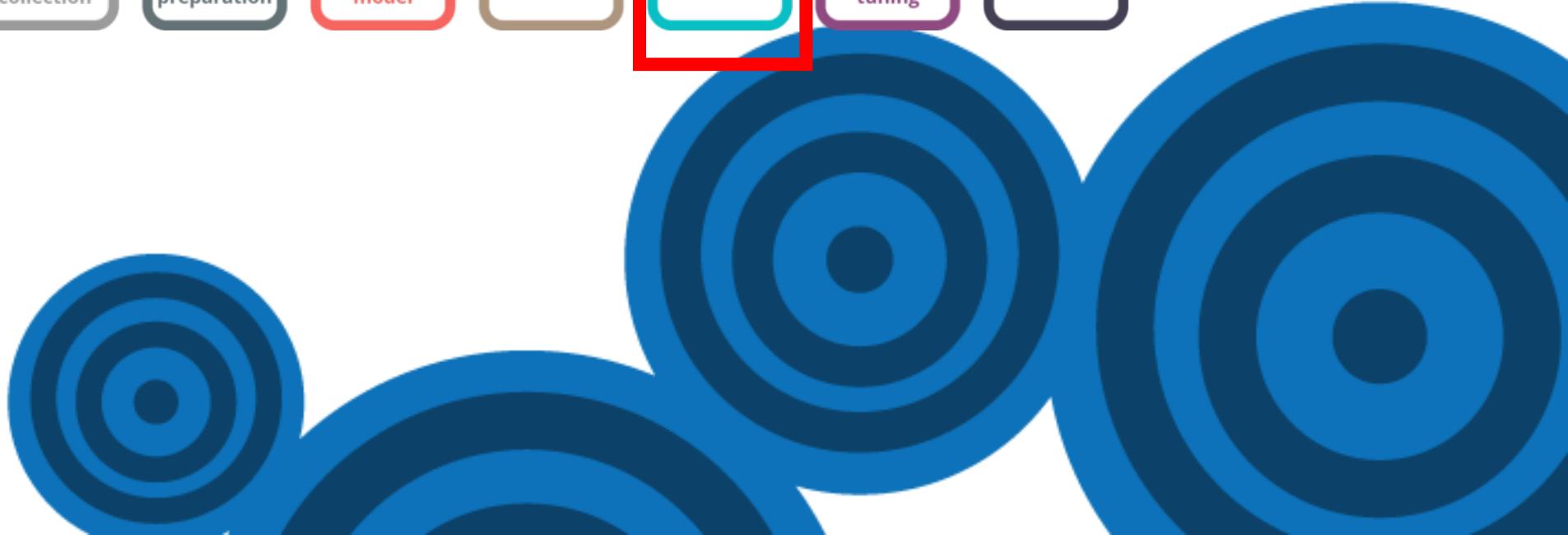


95%

...

Final Accuracy = Average(Round 1, Round 2, ...)

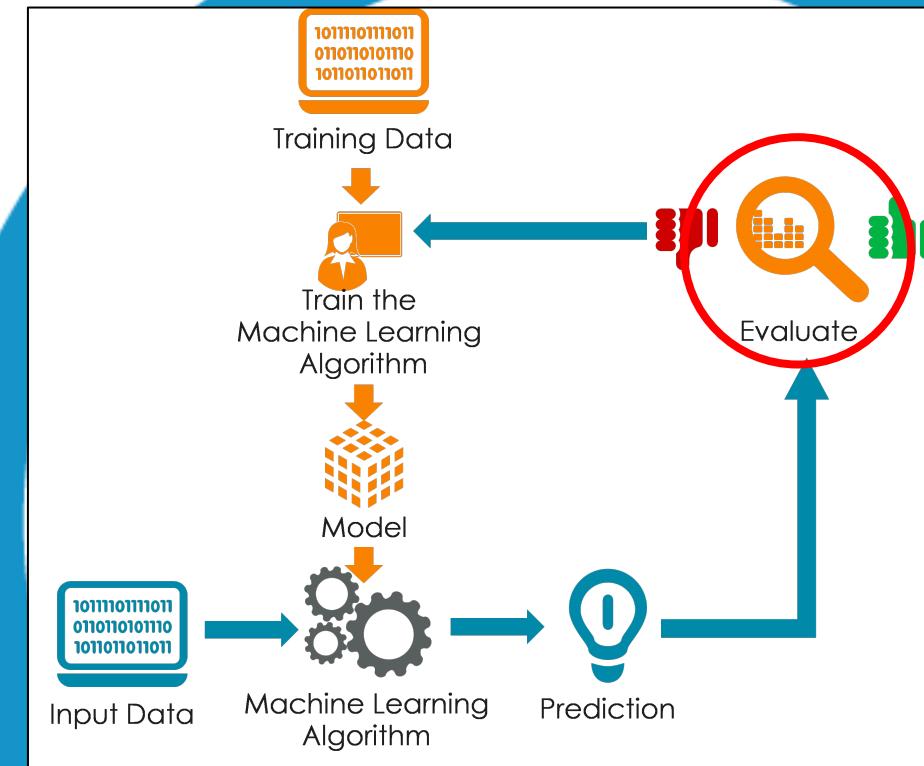
Step 5 Evaluation



Accuracy

តួនាទីនៃការឲ្យរាយ	ពូល	ទ្វាន់	តុក
អើយកើយ	Yes	Yes	True
វីងយ៉ាង	Yes	Yes	True
មិញគេល់បៀង	Yes	Yes	True
ជុំបីបោះពុំ	Yes	Yes	True
កំរើយទៅខ្លួន	Yes	Yes	True
និងទីនេះ	Yes	No	False

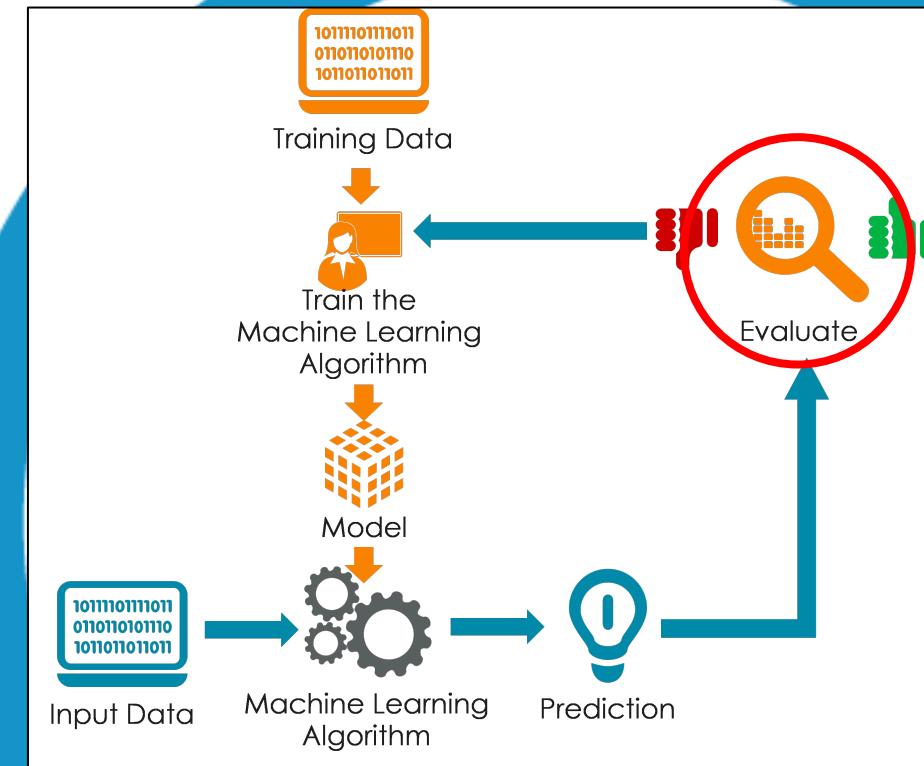
$$\text{Accuracy} = 5/6 = 0.833$$



Accuracy

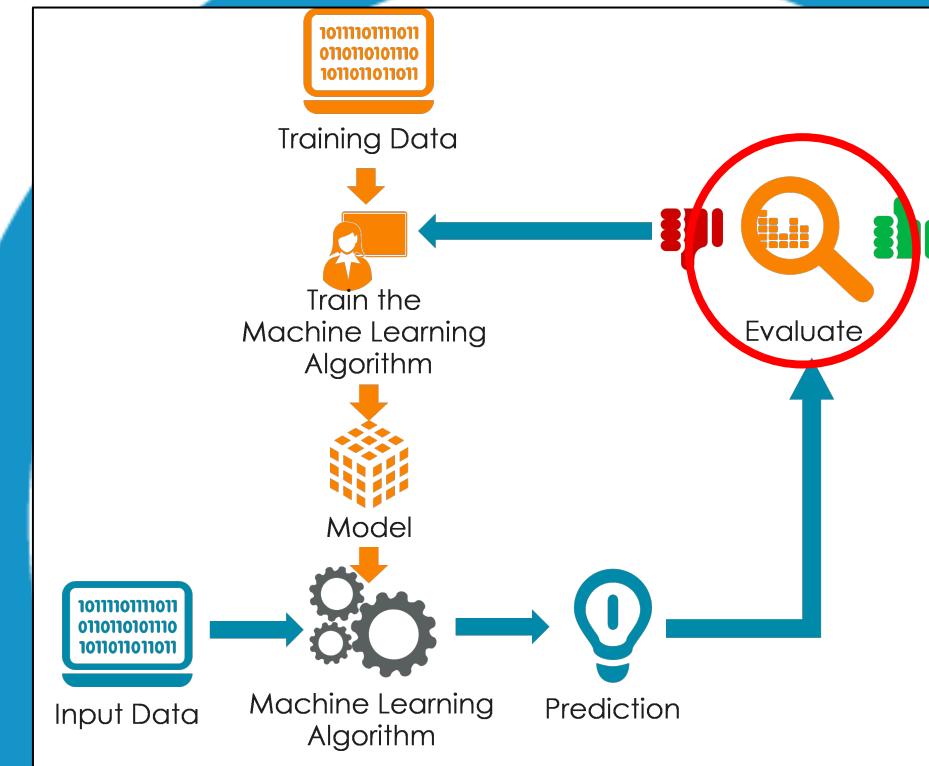
იკრთიდ კოვიდ-19	თობ	ჯრებ	გუკ
ეიქავა	მარტივ	მარტივ	True
ვარდა	მარტივ	მარტივ	True
შეივლებულება	მარტივ	მარტივ	True
ჯირის განვითარება	მარტივ	მარტივ	True
გარეული განვითარება	მარტივ	მარტივ	True
ვარდა მოვარდი	მარტივ	მარტივ	False

$$\text{Accuracy} = 5/6 = 0.833$$



Confusion Matrix

		Predicted Class		Sensitivity $\frac{TP}{(TP + FN)}$
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	Specificity $\frac{TN}{(TN + FP)}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	
	Precision $\frac{TP}{(TP + FP)}$	Negative Predictive Value $\frac{TN}{(TN + FN)}$	Accuracy $\frac{TP + TN}{(TP + TN + FP + FN)}$	

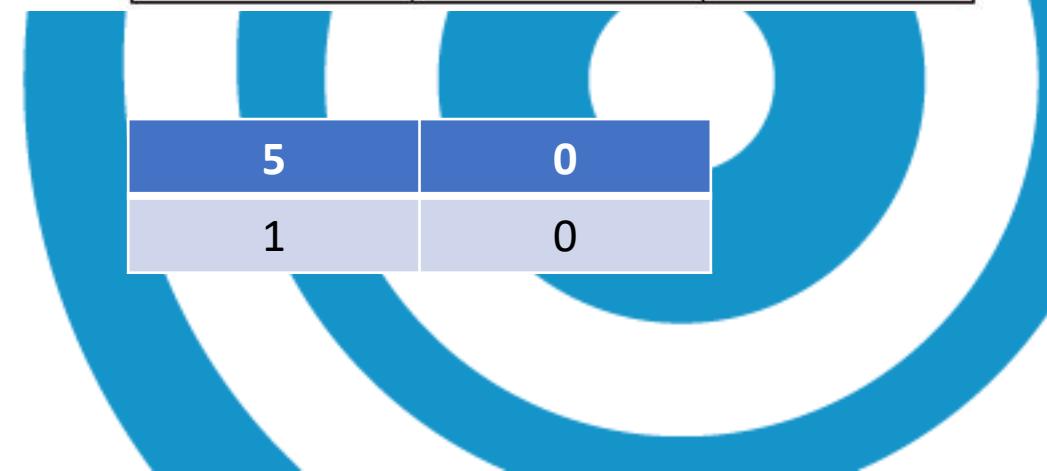


Accuracy

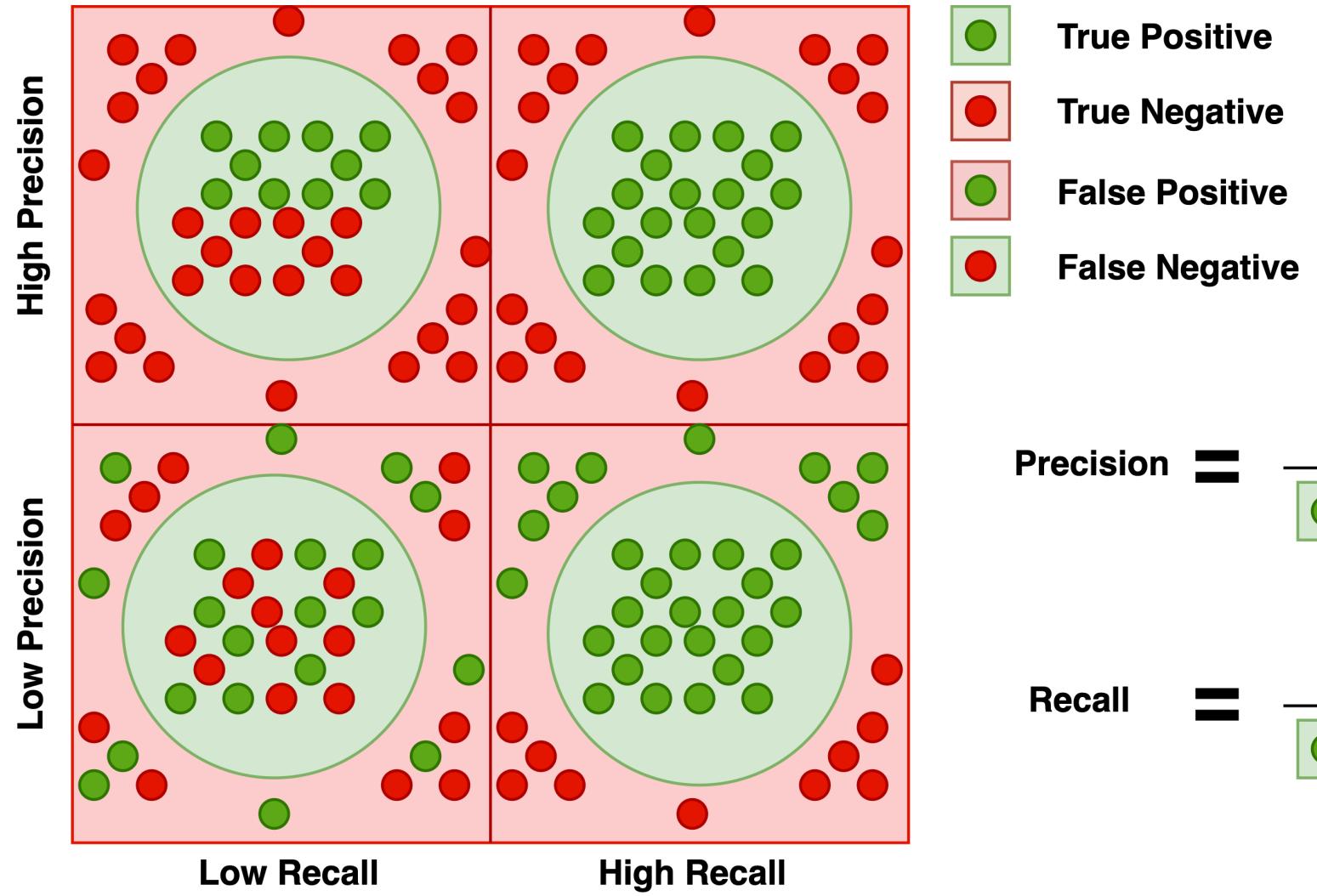
โคตรดิโวติดโควิด-19	ตอบ	จริง	ถูก
เอี้ยกัวย	ไม่ติด	ไม่ติด	True
อึ้งยัง	ไม่ติด	ไม่ติด	True
เชียวนเล่นนีง	ไม่ติด	ไม่ติด	True
จิวแปะกง	ไม่ติด	ไม่ติด	True
กໍວຍເຈັ້ງ	ไม่ติด	ไม่ติด	True
ແອຣ໌ ພອຕາເຕອර໌	ไม่ติด	ຕັດ	False

		Predicted Class		Sensitivity $\frac{TP}{(TP + FN)}$
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	Sensitivity $\frac{TP}{(TP + FN)}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	
		Precision $\frac{TP}{(TP + FP)}$	Negative Predictive Value $\frac{TN}{(TN + FN)}$	Accuracy $\frac{TP + TN}{(TP + TN + FP + FN)}$

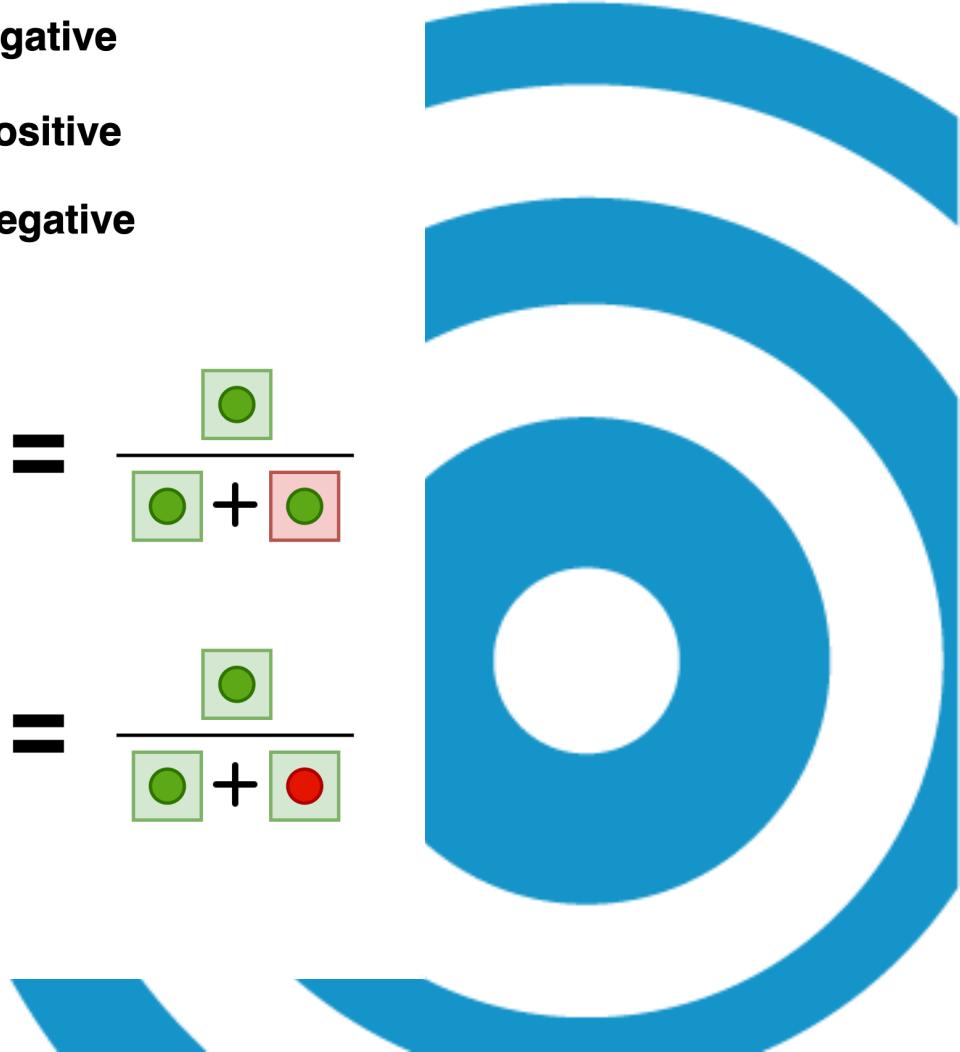
Accuracy = $5/6 = 0.833$



Precision & Recall



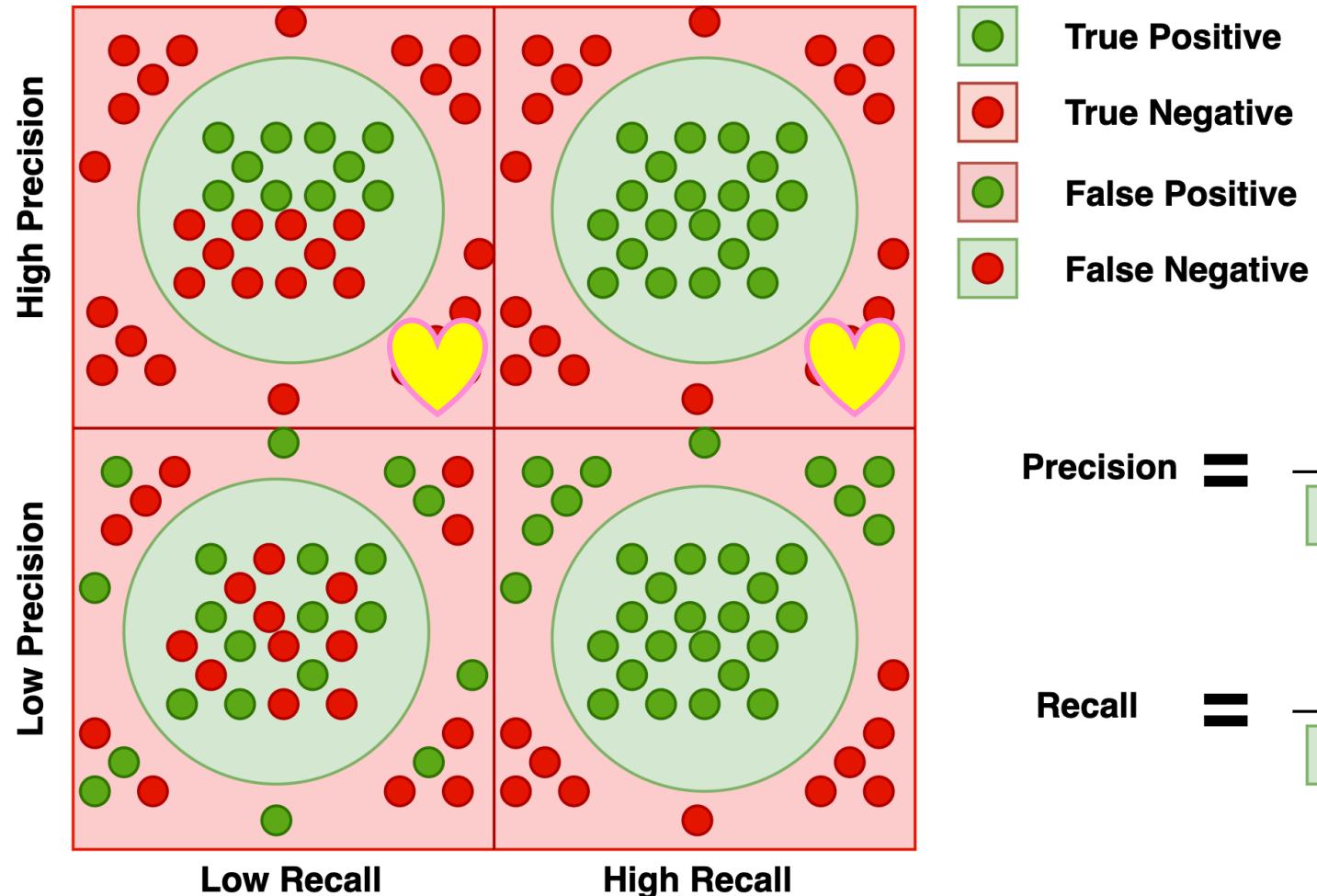
$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$
$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$



Precision & Recall

Which email is a spam

- Spam
- Not Spam



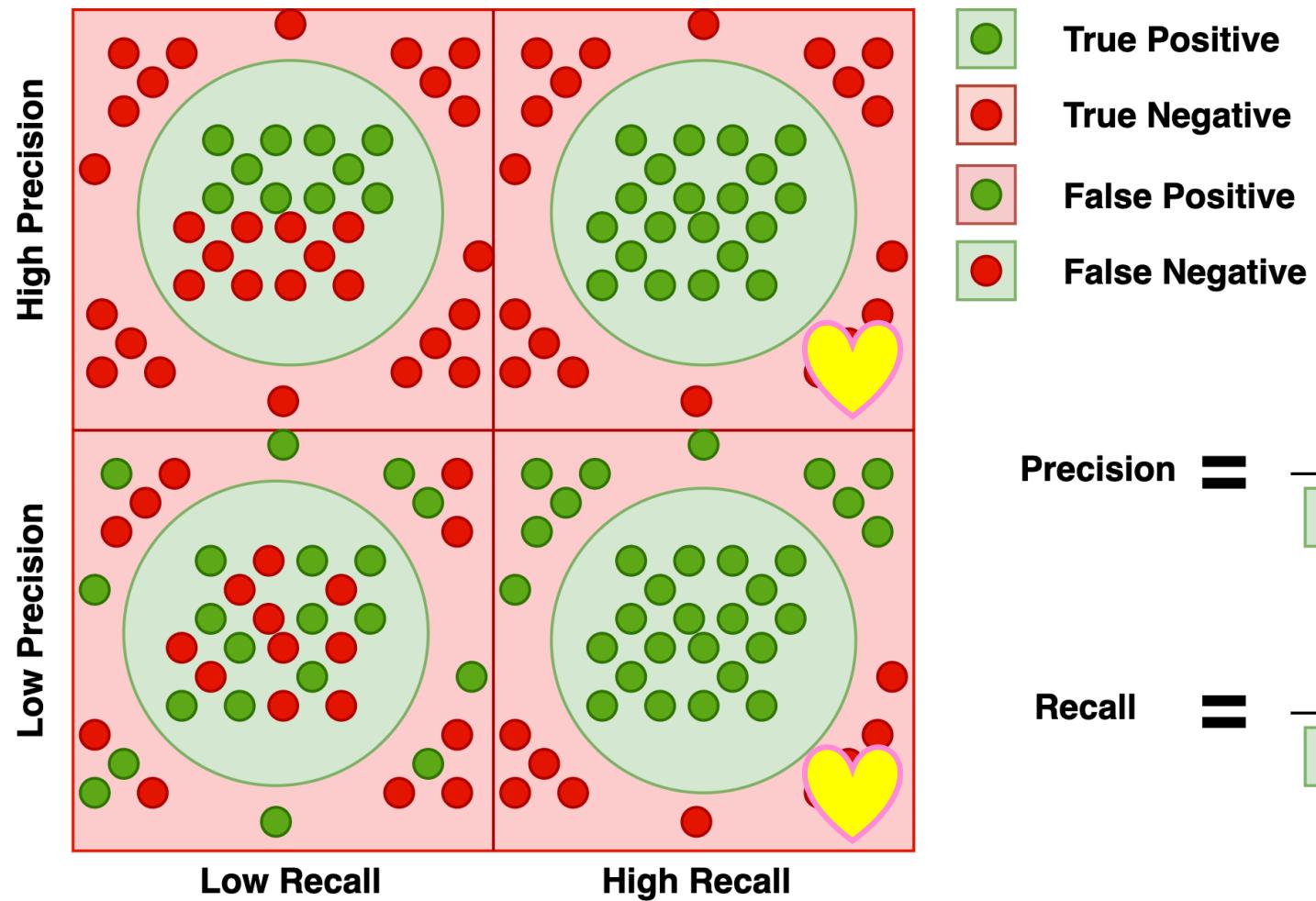
$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Precision & Recall

Who is getting infected covid-19

- Infect
- Not Infect

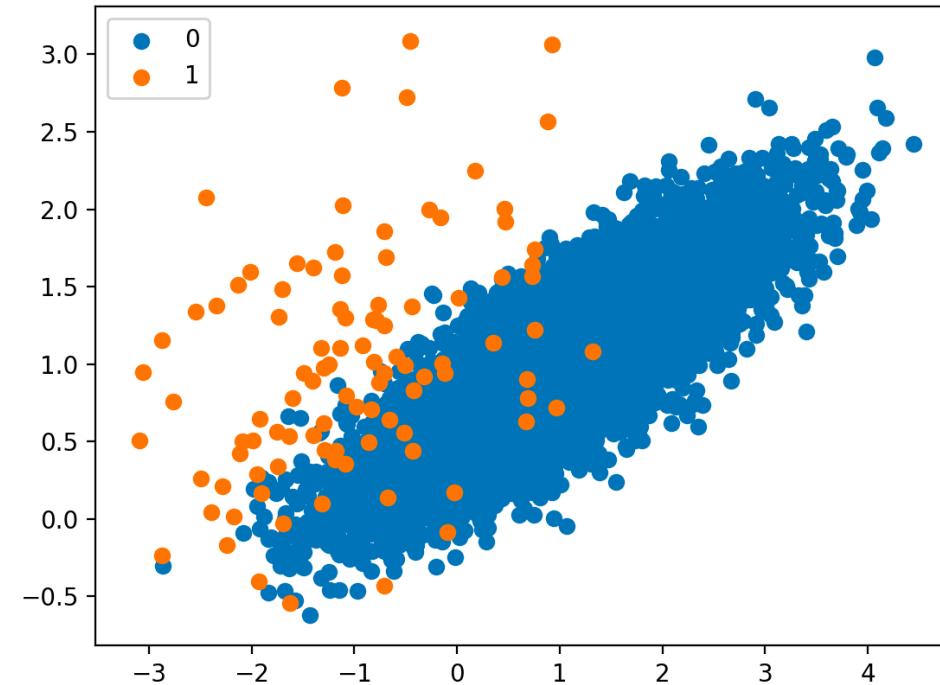


$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

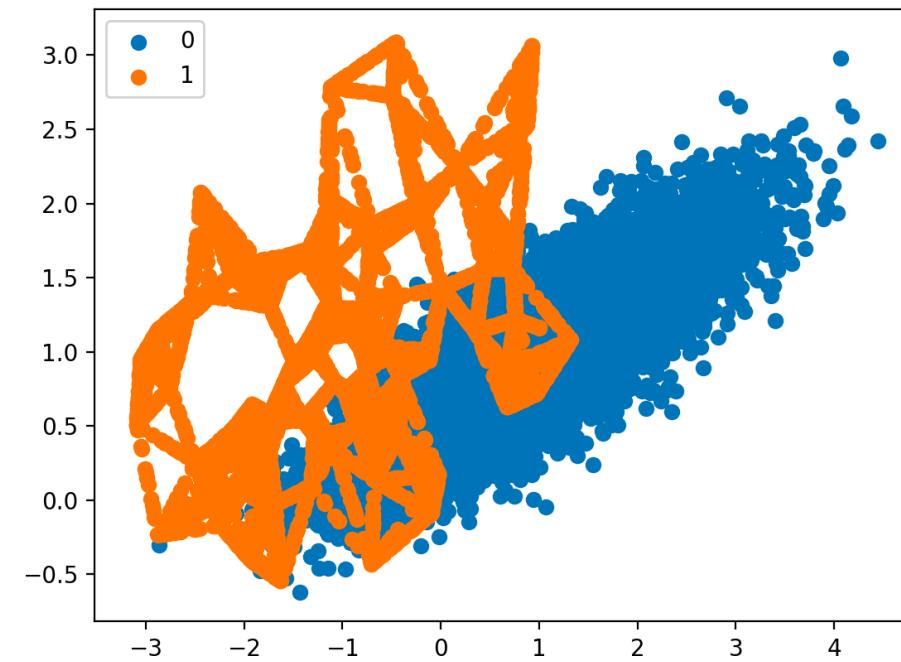
$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

Supervised Learning and Imbalance data

Oversampling



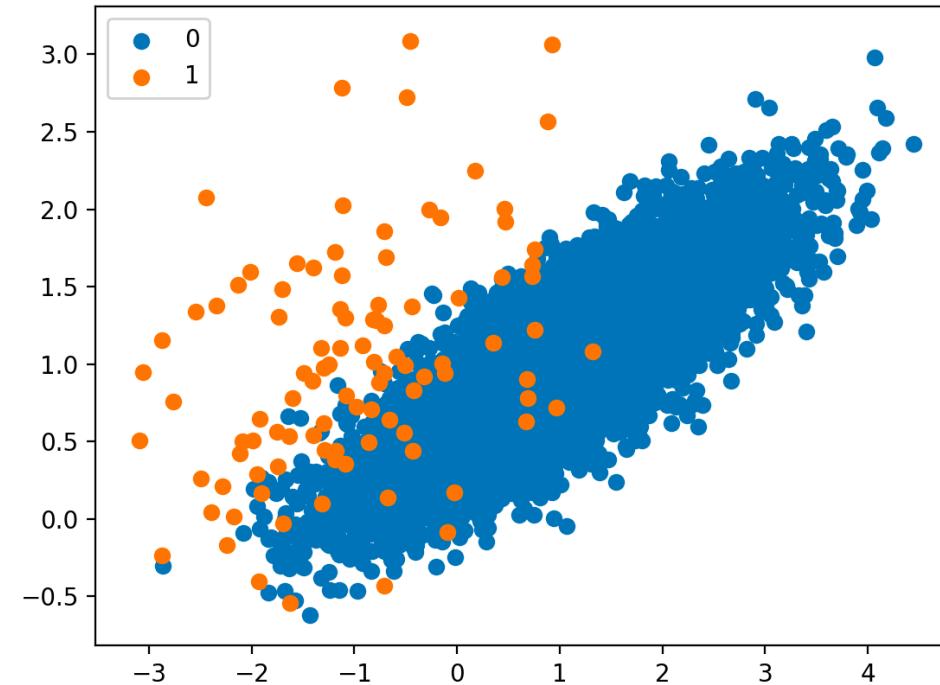
Original dataset



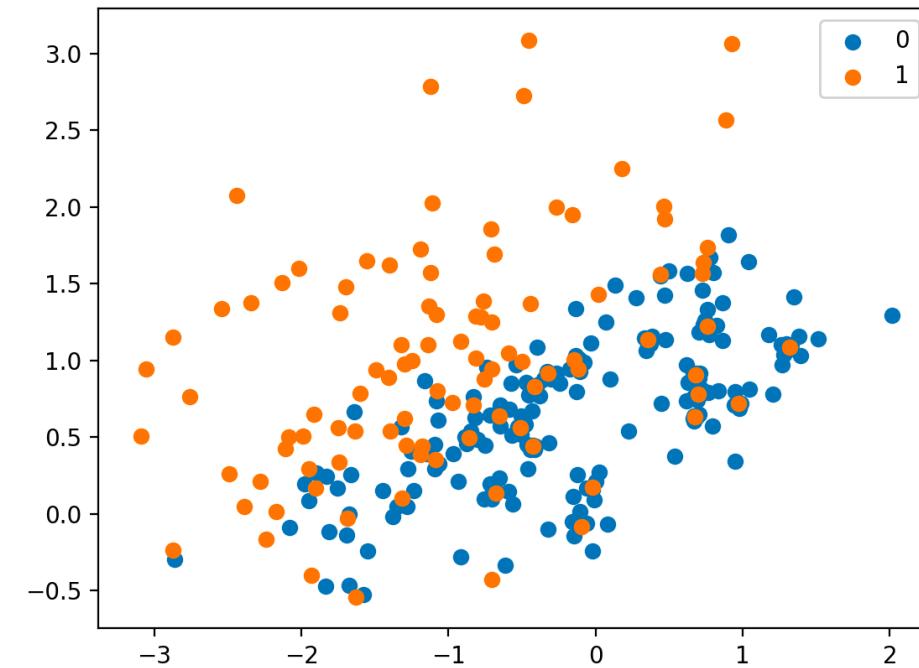
Over sampling with smote

Supervised Learning and Imbalance data

Undersampling



Original dataset

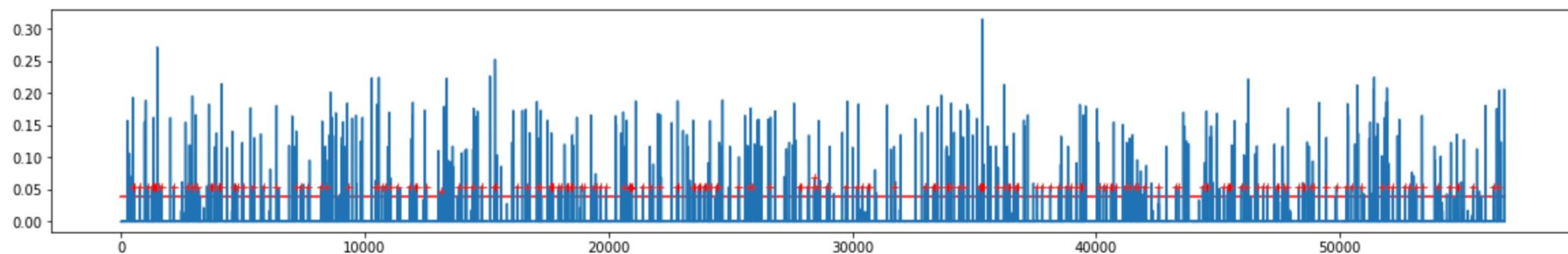


Undersampling

Precision & Recall in Real World

precision recall f1-score support

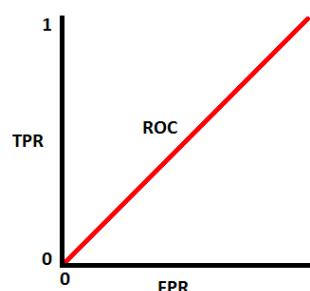
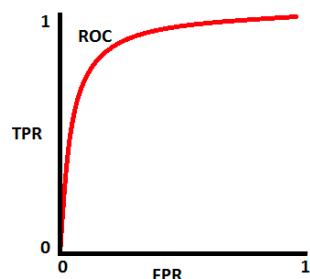
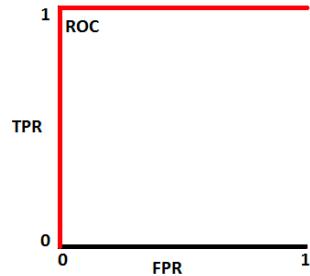
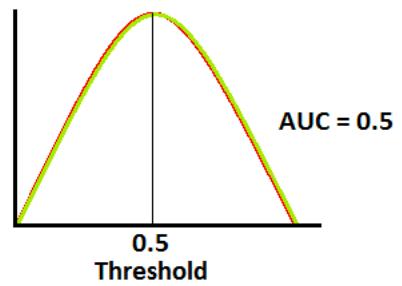
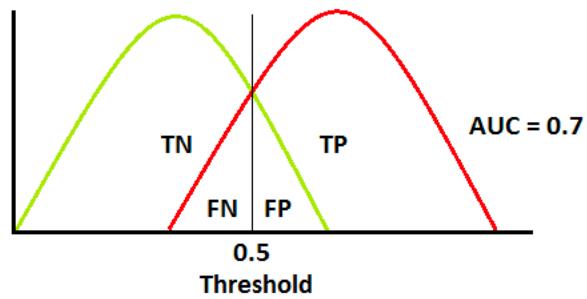
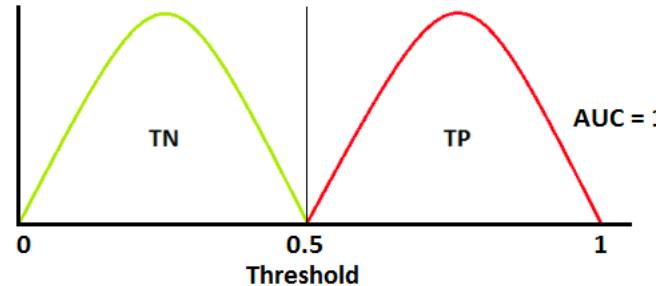
	0	1.00	1.00	56225
	1	0.40	0.57	525
accuracy			0.99	56750
macro avg	1.00	0.70	0.78	56750
weighted avg	0.99	0.99	0.99	56750



เมื่อกำหนด threshold ที่ 0.04 พนว่า

- recall ของ anomaly เท่ากับ 0.4 หมายความว่าเรา Pred เป็น Anomaly ใน 100% เรา Pred ถูก 40%
- precision ของ anomaly เท่ากับ 1 หมายความว่า ใน 40% จาก recall เรา Pred ถูกหมด 100%

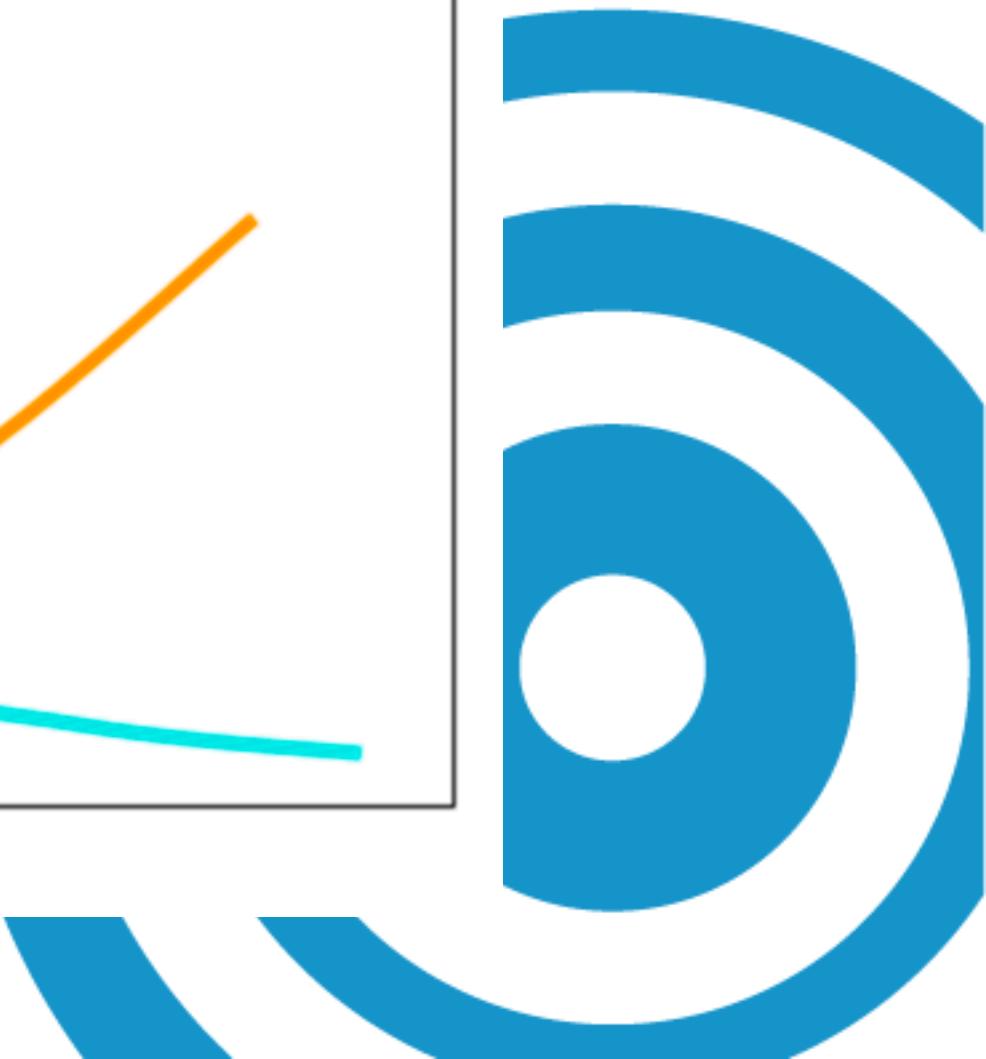
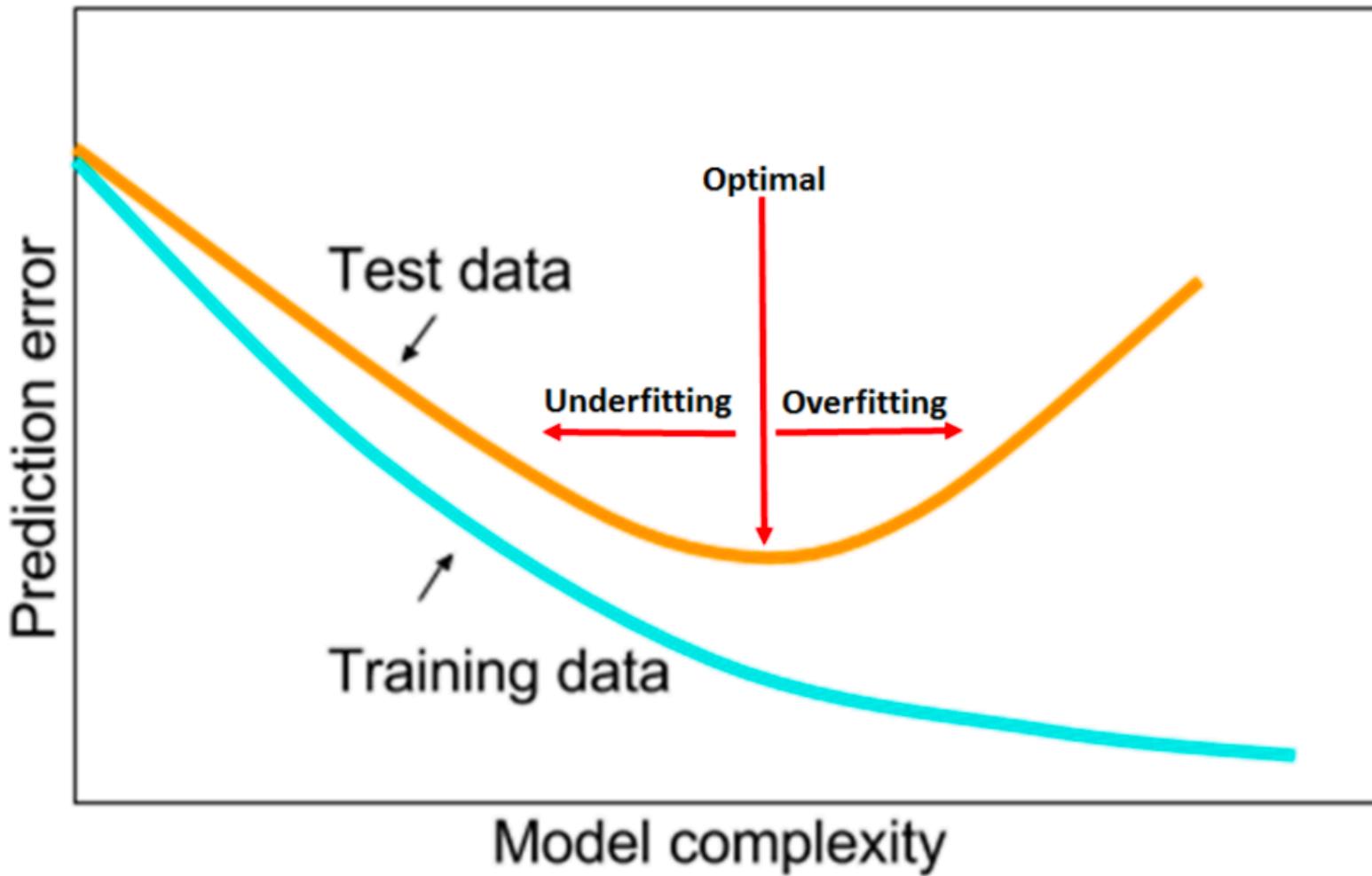
ROC and AUC



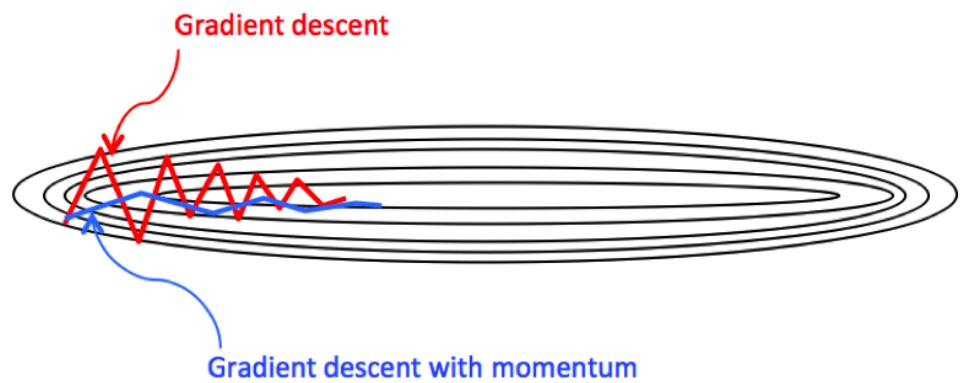
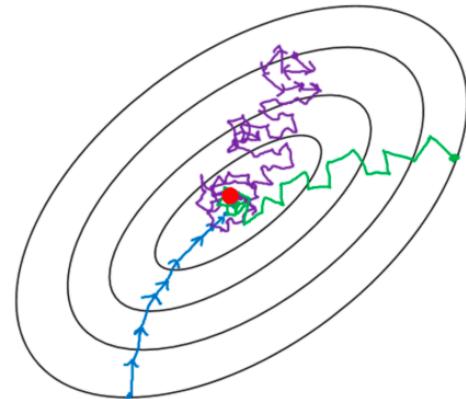
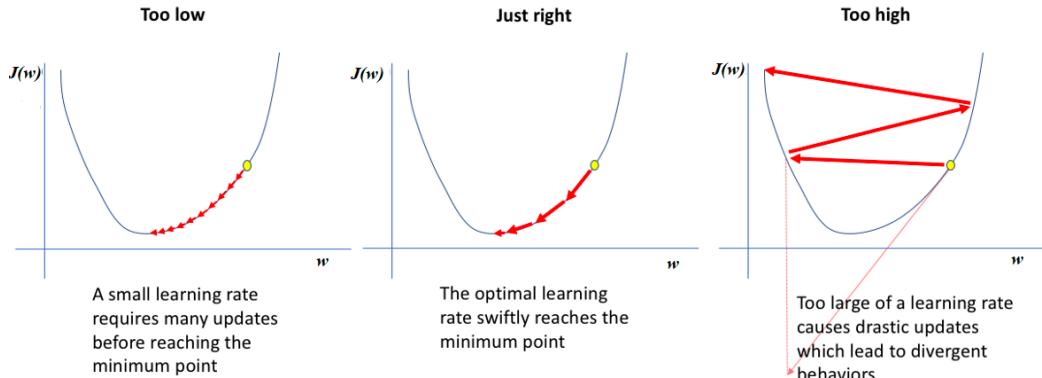
Step 6-7 Parameter tuning & Prediction



Overfit

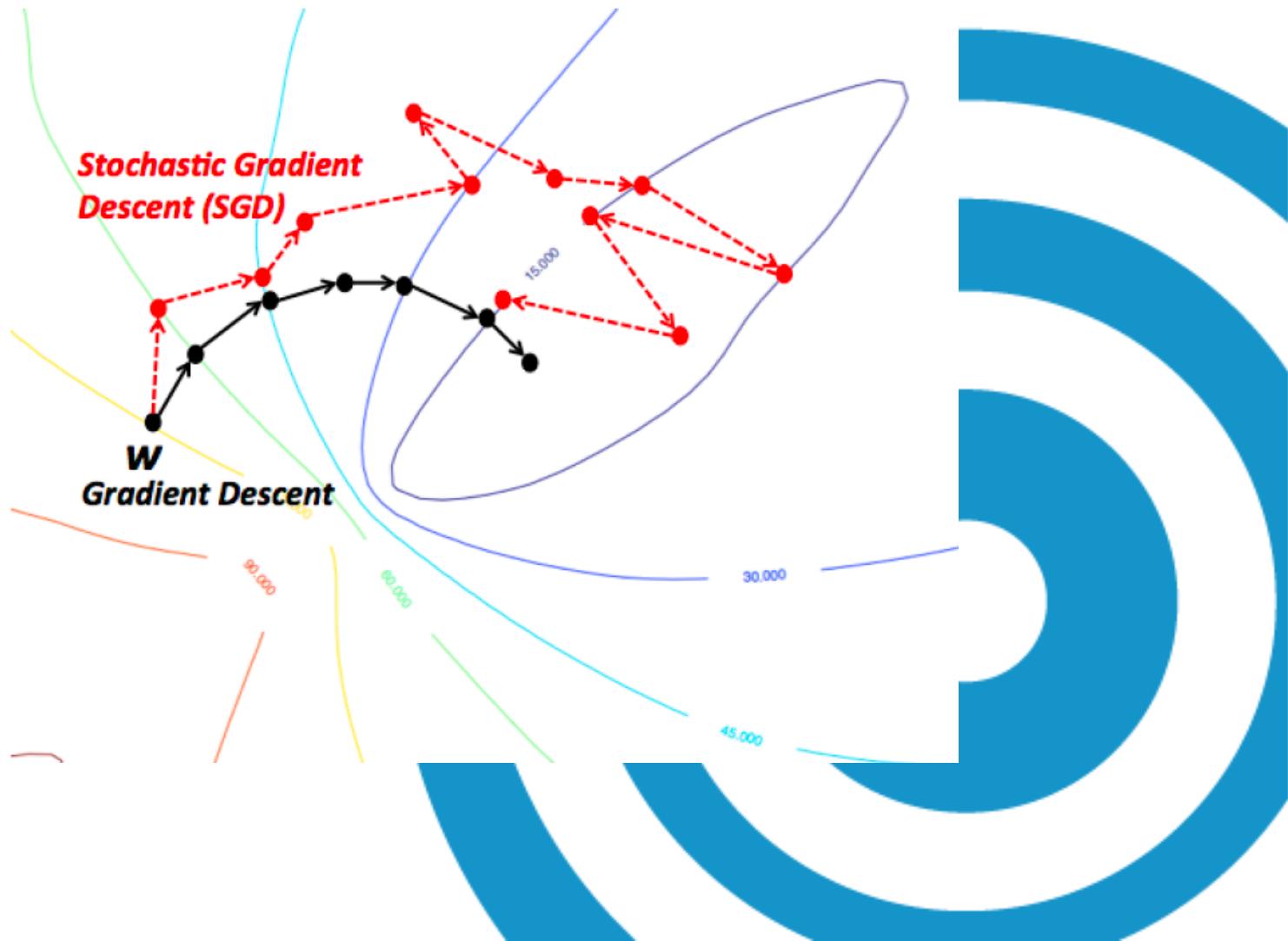


Hyper-parameter Tuning



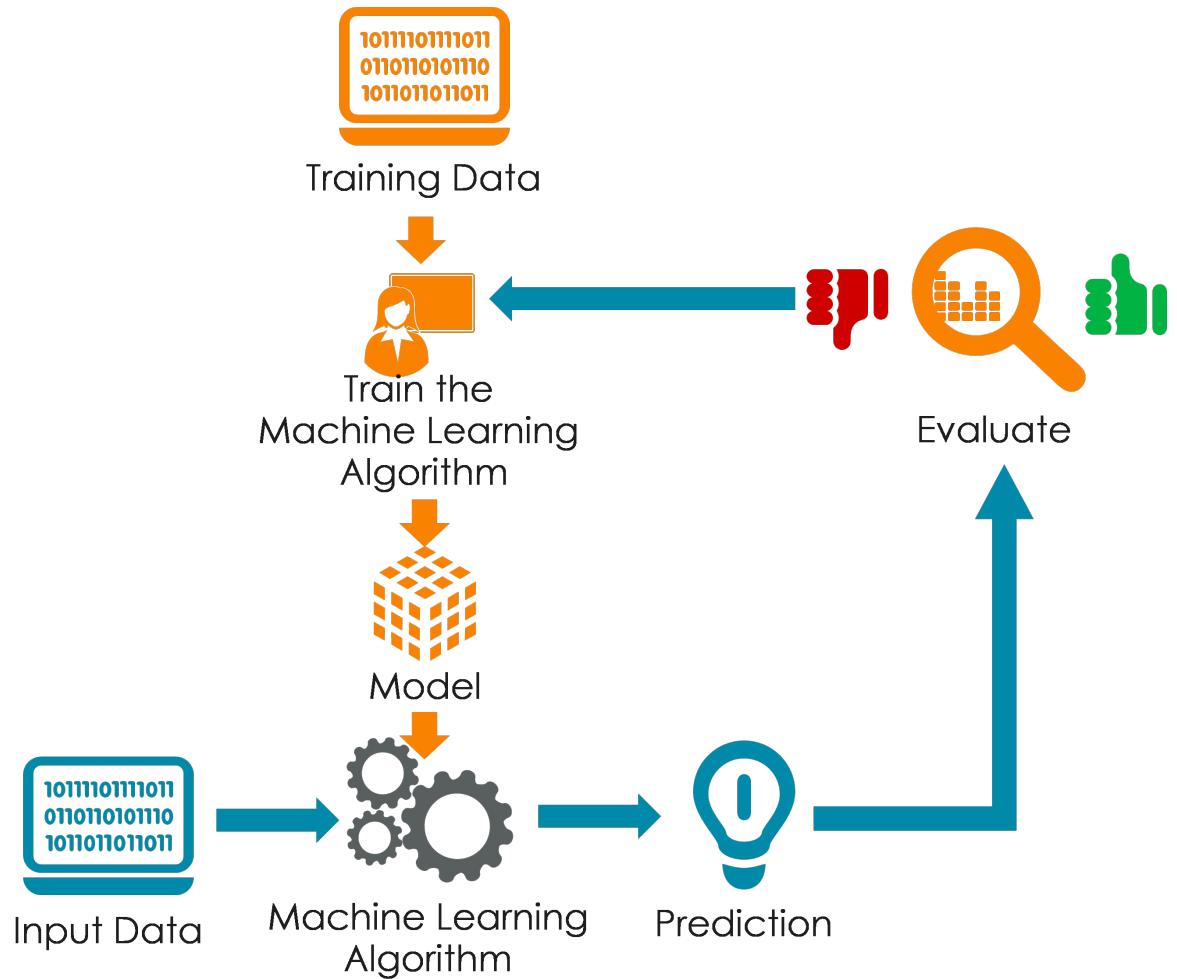
Optimizers

- [SGD](#)
- [RMSprop](#)
- [Adam](#)
- [Adadelta](#)
- [Adagrad](#)
- [Adamax](#)
- [Nadam](#)
- [Ftrl](#)





Prediction



The Machine Learning Life Cycle



1. Define Project Objectives

- Specify business problem
- Acquire subject matter expertise
- Define unit of analysis and prediction target
- Prioritize modeling criteria
- Consider risks and success criteria
- Decide whether to continue

2. Acquire & Explore Data

- Find appropriate data
- Merge data into single table
- Conduct exploratory data analysis
- Find and remove any target leakage
- Feature engineering

3. Model Data

- Variable selection
- Build candidate models
- Model validation and selection

4. Interpret & Communicate

- Interpret model
- Communicate model insights

5. Implement, Document & Maintain

- Set up batch or API prediction system
- Document modeling process for reproducibility
- Create model monitoring and maintenance plan



Thankyou

