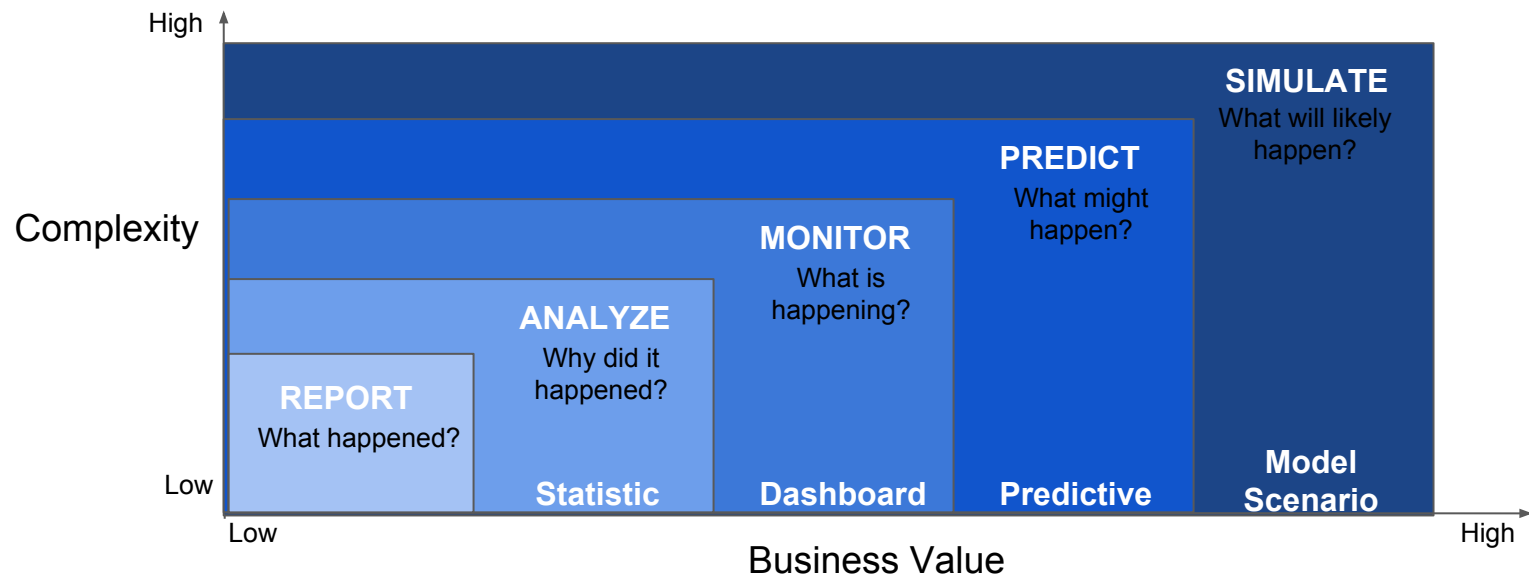# Melek Machine Learning

Data Science Indonesia

# Contents

- Introducing Machine Learning
- Supervised Learning
- Unsupervised Learning
- Evaluation
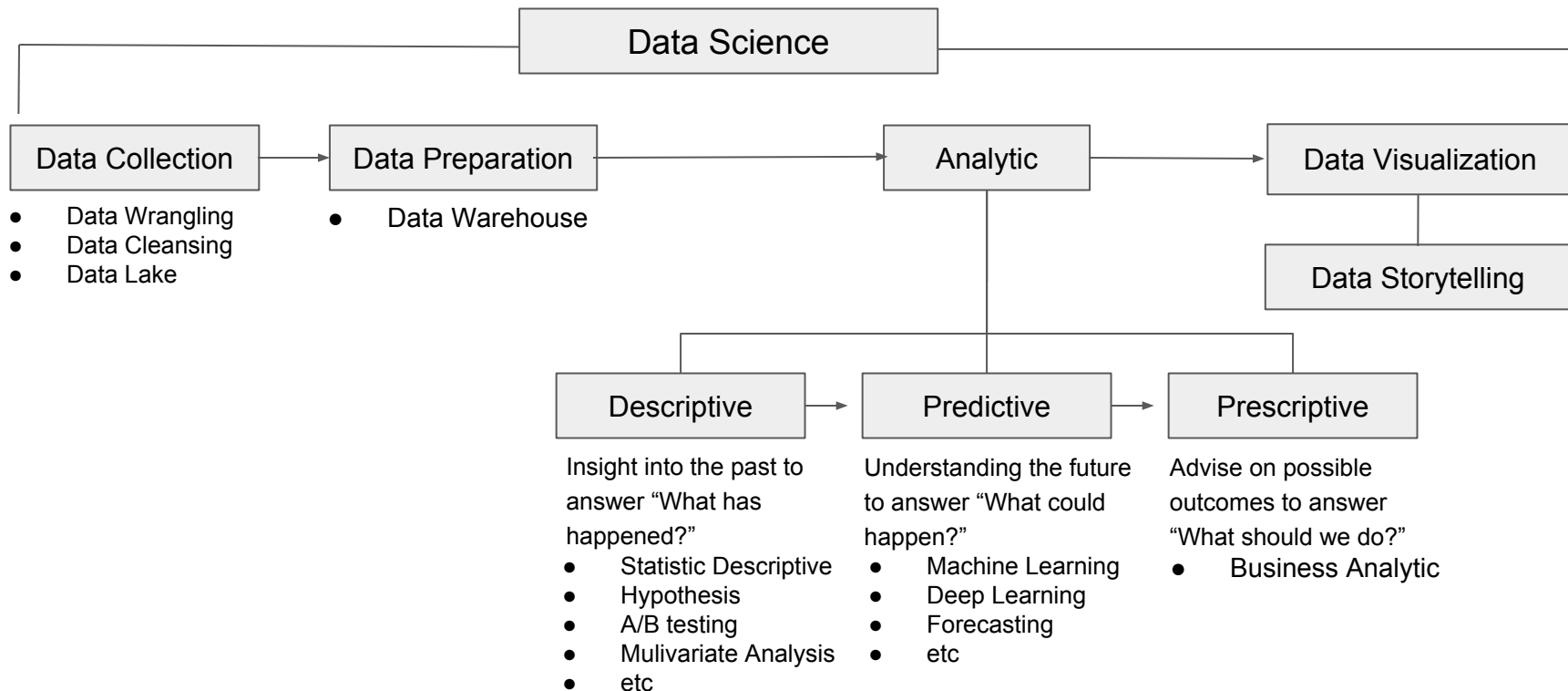- Tuning Parameter
- Analytics

# Where is the location of Machine learning?

From impact to business side:

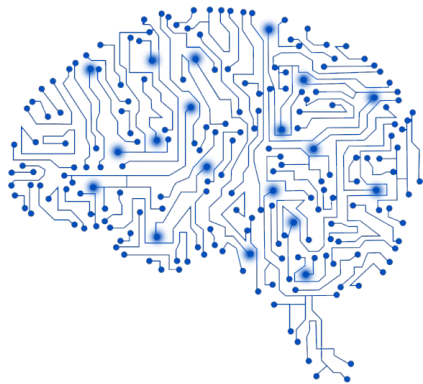# Where is the location of Machine learning?

From Data Science Side:



Data Science

Data Collection
- Data Wrangling
- Data Cleansing
- Data Lake

Data Preparation
- Data Warehouse

Analytic

Data Visualization

Data Storytelling

Descriptive

Insight into the past to answer "What has happened?"
- Statistic Descriptive
- Hypothesis
- A/B testing
- Mulivariate Analysis
- etc

Predictive

Understanding the future to answer "What could happen?"
- Machine Learning
- Deep Learning
- Forecasting
- etc

Prescriptive

Advise on possible outcomes to answer "What should we do?"
- Business Analytic

# So, What is Machine Learning?

**Wikipedia:**

Machine learning is a subset of artificial intelligence in the field of computer science that often uses statistical techniques to give computers the ability to "learn" (i.e., progressively improve performance on a specific task) with data, without being explicitly programmed.



Falalala~

# So, What is Machine Learning?

**Wikipedia:**

Machine learning is a subset of artificial intelligence in the field of computer science that often uses statistical techniques to give computers the ability to "learn" (i.e., progressively improve performance on a specific task) with data, without being explicitly programmed.

**Me:**

Machine Learning is a thing labeler!

# So, What is Machine Learning?

**Thing labeler,** talking your description of something and telling you what label it should get

Example:

We have data:

- *Kaki empat*
- *Memiliki Ekor*
- *Berbulu*
- *Suka colek kaki orang yang makan di warteg*

What should label that appropriate with this object?

→ Machine Learning Model →

- 98% Probability is "Kucing"
- 65% Probability is "Anjing"
- 6% Probability is "Kuda nil"

# So, What is Machine Learning?

why you should be excited about thing labeler



What is this animal?

# So, What is Machine Learning?

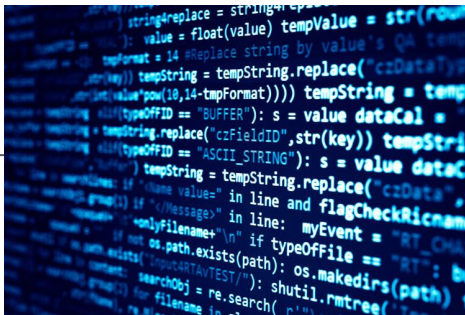Why you should be excited about thing labeler

                    ⟶                    Kucing!

# So, What is Machine Learning?
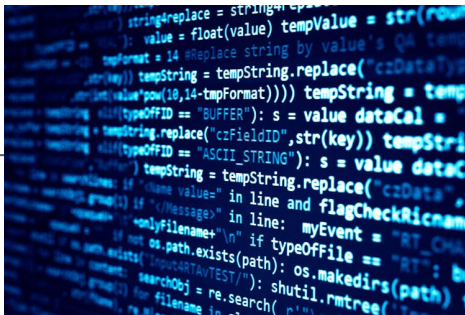
Why you should be excited about thing labeler



Kucing!

Machine learning is a new programming paradigm, a new way of communicating your wishes to a computer.

# So, What is Machine Learning?

Why you should be excited about thing labeler



Kucing!

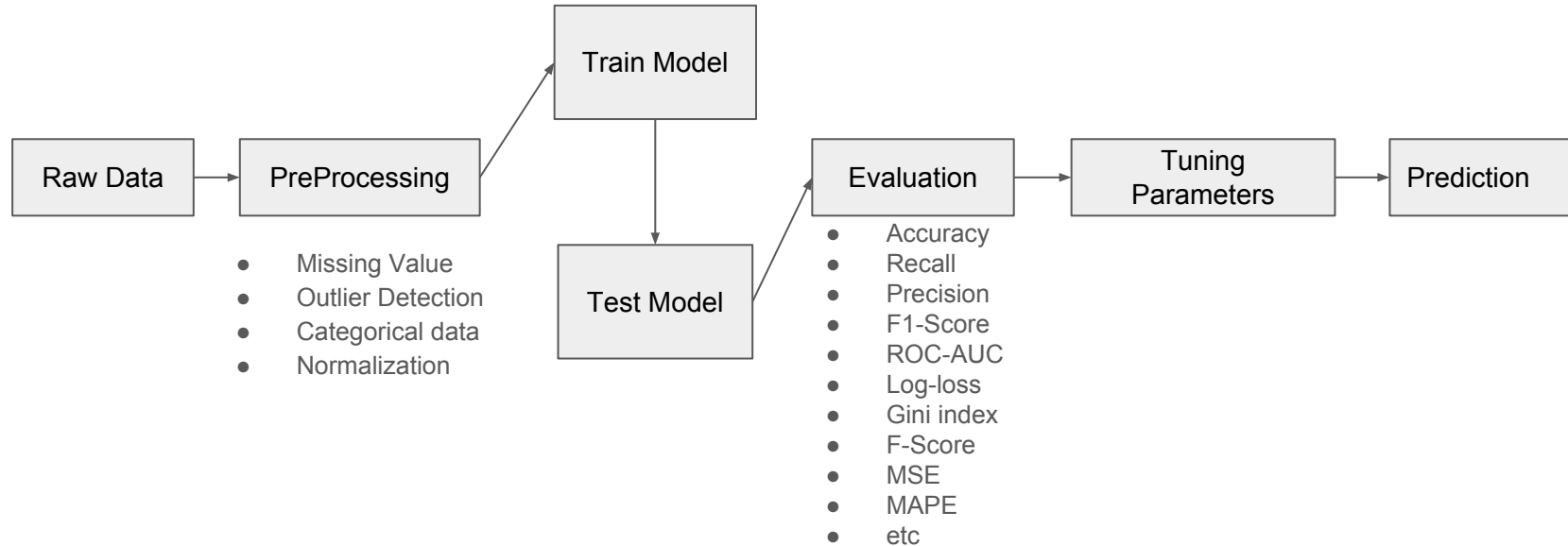Explain with examples (data), not instructions

# Why Machine Learning?
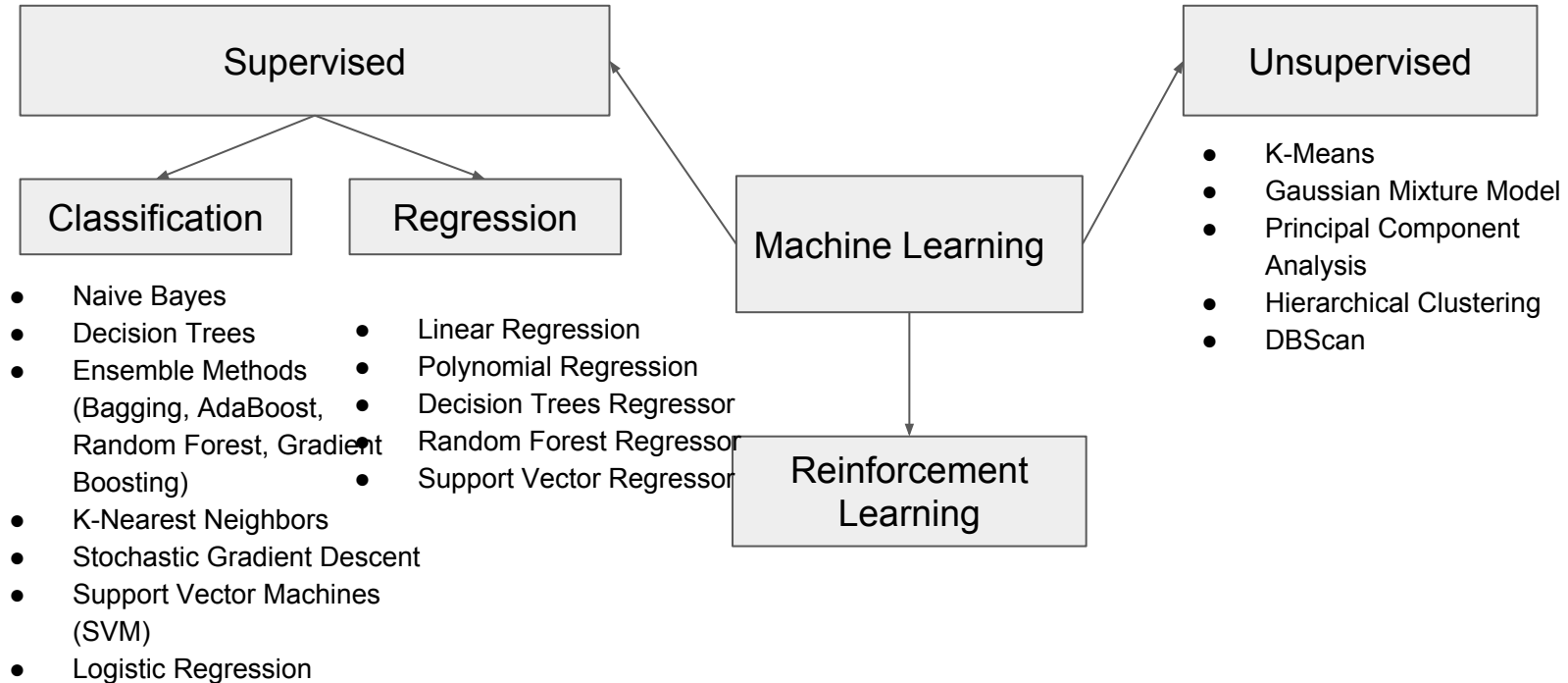
Machine Learning is used to solve business problem

My boss asked me that he wants to:

- wants to know about our customer behavior ---> *ML can help you!*
- wants to know which customer dare to pay expensive ---> *ML can help you!*
- wants to know customer who will stop using our product ---> *ML can help you!*
- wants to sell the same product with different price on each customer at the same time ---> *ML can help you!*
- wants to know the customer who will pay the credit until the end ---> *ML can help you!*
- wants to recommend products according to the needs of every customer ---> *ML can help you!*
- wants to know if there is fraud to the customer in using our product ---> *ML can help you!*
- wants to know the market potential for new innovation products ---> *ML can help you!*
- And many more

# How to build Machine Learning Models?

# What are kinds of Machine Learning model?

**Supervised**

**Classification**

- Naive Bayes
- Decision Trees
- Ensemble Methods (Bagging, AdaBoost, Random Forest, Gradient Boosting)
- K-Nearest Neighbors
- Stochastic Gradient Descent
- Support Vector Machines (SVM)
- Logistic Regression

**Regression**

- Linear Regression
- Polynomial Regression
- Decision Trees Regressor
- Random Forest Regressor
- Support Vector Regressor

**Machine Learning**

**Reinforcement Learning**

**Unsupervised**

- K-Means
- Gaussian Mixture Model
- Principal Component Analysis
- Hierarchical Clustering
- DBScan

# Supervised Learning

# Supervised Learning

## What is Supervised Learning?

Supervised learning is where you have input variables (x) and an output variable (Y) and you use an algorithm to learn the mapping function from the input to the output.

## Algorithms:

- K-Nearest Neighbor
- Decision Tree
- Support Vector Machine
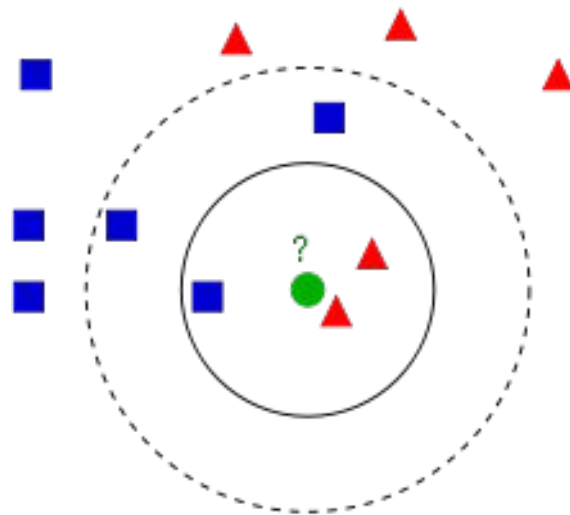
# Supervised Learning

- Regression
- Classification

Supervised learning in classification handles estimator for discrete labels rather than continuous labels.

Classification > binary classification, multiclass classification, multilabel classification

# K-Nearest Neighbor

- K-NN is an instance-based learning or lazy learning
- K-NN using distance to measure the likelihood of the class
- Number of K take as a comparison of the likelihood

# Pros and Cons

Pros

- Simple Algorithm
- Versatile (Classification and Regression)
- Does not assume any probability distribution on the input data

Cons

- Requires high memory
- Computationally Expensive
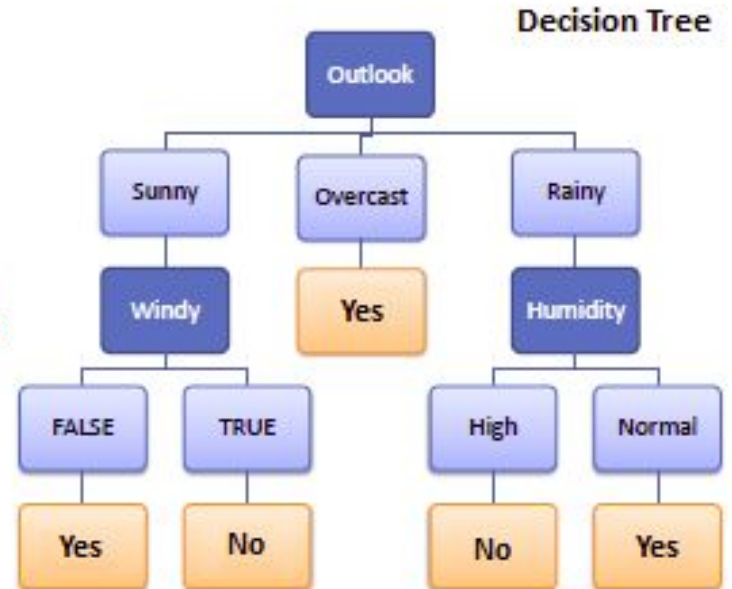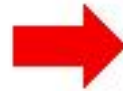- Lazy Learning

# Decision Tree

Metrics that DT consider:

- Information Gain
- Entropy

Algorithms:

- ID3
- C4.5
- C5.0
- CART

# Decision Tree

Predictors

Target

| Outlook | Temp. | Humidity | Windy | Play Golf |
|---|---|---|---|---|
| Rainy | Hot | High | False | No |
| Rainy | Hot | High | True | No |
| Overcast | Hot | High | False | Yes |
| Sunny | Mild | High | False | Yes |
| Sunny | Cool | Normal | False | Yes |
| Sunny | Cool | Normal | True | No |
| Overcast | Cool | Normal | True | Yes |
| Rainy | Mild | High | False | No |
| Rainy | Cool | Normal | False | Yes |
| Sunny | Mild | Normal | False | Yes |
| Rainy | Mild | Normal | True | Yes |
| Overcast | Mild | High | True | Yes |
| Overcast | Hot | Normal | False | Yes |
| Sunny | Mild | High | True | No |

**Decision Tree**

Outlook

Sunny — Overcast — Rainy

Windy — Yes — Humidity

FALSE — TRUE

Yes — No

High — Normal

No — Yes

# Step 1: Calculate Entropy

$$E(S) = \sum_{i=1}^{c} - p_i \log_2 p_i$$

$$E(T,X) = \sum_{c \in X} P(c)E(c)$$

| Play Golf | |
|-----------|-----|
| Yes | No |
| 9 | 5 |

| | | Play Golf | | |
|---------|----------|-----|-----|-----|
| | | Yes | No | |
| Outlook | Sunny | 3 | 2 | 5 |
| | Overcast | 4 | 0 | 4 |
| | Rainy | 2 | 3 | 5 |
| | | | | 14 |

Entropy(PlayGolf) = Entropy (5,9)

    = Entropy (0.36, 0.64)

    = - (0.36 log$_2$ 0.36) - (0.64 log$_2$ 0.64)

    = 0.94

**E**(PlayGolf, Outlook) = **P**(Sunny)***E**(3,2) + **P**(Overcast)***E**(4,0) + **P**(Rainy)***E**(2,3)

    = (5/14)*0.971 + (4/14)*0.0 + (5/14)*0.971

    = 0.693

# Step 2: Calculate Information Gain

$$Gain(T,X) = Entropy(T) - Entropy(T,X)$$

G(PlayGolf, Outlook) = E(PlayGolf) – E(PlayGolf, Outlook)

= 0.940 – 0.693 = 0.247

| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Outlook | Sunny | 3 | 2 |
| | Overcast | 4 | 0 |
| | Rainy | 2 | 3 |
| Gain = 0.247 | | | |

| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Temp. | Hot | 2 | 2 |
| | Mild | 4 | 2 |
| | Cool | 3 | 1 |
| Gain = 0.029 | | | |

| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Humidity | High | 3 | 4 |
| | Normal | 6 | 1 |
| Gain = 0.152 | | | |

| | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Windy | False | 6 | 2 |
| | True | 3 | 3 |
| Gain = 0.048 | | | |

# Step 3: Choose Root Node

| | | Play Golf | |
|---|---|---|---|
| | ★ | Yes | No |
| **Outlook** | Sunny | 3 | 2 |
| | Overcast | 4 | 0 |
| | Rainy | 2 | 3 |
| | Gain = 0.247 | | |

| Outlook | Temp | Humidity | Windy | Play Golf |
|---|---|---|---|---|
| Sunny | Mild | High | FALSE | Yes |
| Sunny | Cool | Normal | FALSE | Yes |
| Sunny | Cool | Normal | TRUE | No |
| Sunny | Mild | Normal | FALSE | Yes |
| Sunny | Mild | High | TRUE | No |
| Overcast | Hot | High | FALSE | Yes |
| Overcast | Cool | Normal | TRUE | Yes |
| Overcast | Mild | High | TRUE | Yes |
| Overcast | Hot | Normal | FALSE | Yes |
| Rainy | Hot | High | FALSE | No |
| Rainy | Hot | High | TRUE | No |
| Rainy | Mild | High | FALSE | No |
| Rainy | Cool | Normal | FALSE | Yes |
| Rainy | Mild | Normal | TRUE | Yes |

# Step 3: Node Leaf

| Temp | Humidity | Windy | Play Golf |
|------|----------|-------|-----------|
| Hot | High | FALSE | Yes |
| Cool | Normal | TRUE | Yes |
| Mild | High | TRUE | Yes |
| Hot | Normal | FALSE | Yes |

Outlook
- Sunny
- Overcast
  - Play=Yes
- Rainy

# Step 4: Calculate Entropy & Information Gain

| Temp | Humidity | Windy | Play Golf |
|------|----------|-------|-----------|
| Mild | High | FALSE | Yes |
| Cool | Normal | FALSE | Yes |
| Mild | Normal | FALSE | Yes |
| Cool | Normal | TRUE | No |
| Mild | High | TRUE | No |

Outlook

Sunny    Overcast    Rainy

Windy    Play=Yes

FALSE    TRUE

Play=Yes    Play=No

# Step 5: Calculate Information Gain

**R₁: IF** (Outlook=Sunny) AND (Windy=FALSE) **THEN** Play=Yes

**R₂: IF** (Outlook=Sunny) AND (Windy=TRUE) **THEN** Play=No

**R₃: IF** (Outlook=Overcast) **THEN** Play=Yes

**R₄: IF** (Outlook=Rainy) AND (Humidity=High) **THEN** Play=No

**R₅: IF** (Outlook=Rain) AND (Humidity=Normal) **THEN** Play=Yes

# Pros and Cons

Pros:

- Easy to Explain
- Follows the same approach  with human
- Interpretation of a complex decision tree can be simplified by visualization

Cons:

- High probability of overfitting
- Calculations can be complex when there are many class labels

# Support Vector Machine

# Pros and Cons

Pros:

- Works well with clear margin
- Effective in High Dimensional Spaces
- Effective in cases where num. Of dimension is greater than num. Of samples

Cons:

- High computation
- Bad at a lot of noise

# Let's Code

# Unsupervised Learning

# Unsupervised Learning

What is Unsupervised Learning?

Algorithms:

- K-Means
- Hierarchy clustering
- DBSCAN

# K-Means

K-Means is the 'go-to' clustering algorithm for many simply because it is fast, easy to understand, and available everywhere (there's an implementation in almost any statistical or machine learning tool you care to use) <- Really Suitable with large data



© 2012 Ted Goff

"Here's a list of 100,000 warehouses full of data. I'd like you to condense them down to one meaningful warehouse."

# Pros and Cons

- Applicable only when mean is defined, then what about categorical data?
- Need to specify k, the number of clusters, in advance. If you have a good intuition for how many clusters the dataset your exploring has then great, otherwise you might have a problem
- Unable to handle noisy data and outliers
- Not suitable to discover clusters with non-convex shapes
- Doesn't consider the proportion of different cluster
- Doesn't consider the variance of different cluster

# Hierarchy Clustering

Hierarchy algorithms: Create a hierarchical

decomposition of the set of data (or objects) using

some criterion

In short, decompose data objects into a several
levels of nested partitioning (tree of clusters), called
a dendrogram.

A clustering of the data objects is obtained by
cutting the dendrogram at the desired level, then
each connected component forms a cluster.

# Pros and Cons

- **do not scale well**: time complexity of at least O(n^2), where n is the number of total objects
- **can never undo** what was done previously
- Similar to K-Means we are stuck **choosing the number of clusters** (not easy in EDA), or trying to discern some natural parameter value from a plot that may or may not have any obvious natural choices.



"What we've done is make it dramatically easier to navigate the corporate hierarchy."

# DBSCAN

DBSCAN is a density based algorithm – it assumes clusters for dense regions. It is also the first actual clustering algorithm we've looked at: it doesn't require that every point be assigned to a cluster and hence doesn't partition the data, but instead extracts the 'dense' clusters and leaves sparse background classified as 'noise'.

# Pros and Cons

- Epsilon is a distance value, so you need to survey the distribution of distances in your dataset to attempt to get an idea of where it should lie. In practice, however, this isn't an especially intuitive parameter, nor is it easy to get right.
- Sensitive to the parameter

# Let's Code

# Evaluation

# Evaluation - Classification

- Cross Validation
  - Stratified Cross Validation
  - One Leave Out Cross Validation
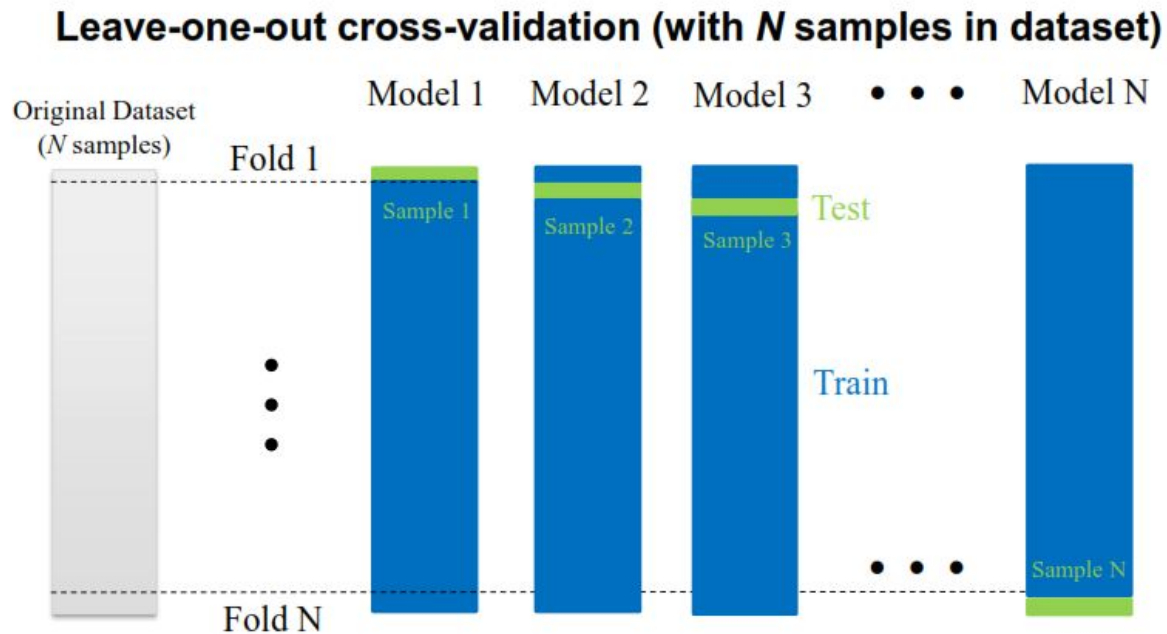- Metrices
  - Accuracy
  - Precision
  - Recall
  - F1-Score

# Cross Validation

# Stratified Cross Validation



## Stratified Cross-validation

| fruit_label | fruit_name |
|---|---|
| 1 | Apple |
| 1 | Apple |
| 1 | Apple |
| 1 | Apple |
| 1 | Apple |
| 2 | Mandarin |
| ... | ... |
| 3 | Orange |
| ... | ... |
| 4 | Lemon |
| 4 | Lemon |
| 4 | Lemon |
| 4 | Lemon |
| 4 | Lemon |

Fold 1

Test 20%

Train 80%

Stratified folds each contain a proportion of classes that matches the overall dataset. Now, all classes will be fairly represented in the test set.

# Leave One Out Cross Validation



Leave-one-out cross-validation (with $N$ samples in dataset)

# Confusion Matrix



## Binary prediction outcomes

|  | **Predicted** negative | **Predicted** positive |
|---|---|---|
| **True** negative | TN | FP |
| **True** positive | FN | TP |

Label 1 = positive class (class of interest)

Label 0 = negative class (everything else)

TP = true positive
FP = false positive (Type I error)
TN = true negative
FN = false negative (Type II error)

# Confusion Matrix Visualization



digits dataset: positive class (black) is digit 1, negative class (white) all others

# Accuracy

| | Predicted negative | Predicted positive | |
|---|---|---|---|
| True negative | TN = 400 | FP = 7 | |
| True positive | FN = 17 | TP = 26 | |
| | | | $N = 450$ |

$$\text{Accuracy} = \frac{TN+TP}{TN+TP+FN+FP}$$

$$= \frac{400+26}{400+26+17+7}$$

$$= 0.95$$

# Recall

| | Predicted negative | Predicted positive | |
|---|---|---|---|
| True negative | TN = 400 | FP = 7 | |
| True positive | FN = 17 | TP = 26 | |
| | | | $N = 450$ |

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

$$= \frac{26}{26+17}$$

$$= 0.60$$

Recall is also known as:
- True Positive Rate (TPR)
- Sensitivity
- Probability of detection

# Precision

|  | Predicted negative | Predicted positive |
|---|---|---|
| True negative | TN = 400 | FP = 7 |
| True positive | FN = 17 | TP = 26 |

N = 450

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

$$= \frac{26}{26+7}$$

$$= 0.79$$

# F1-Score

$$F_1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall} = \frac{2 \cdot TP}{2 \cdot TP + FN + FP}$$

# Evaluation - Clustering

# Tuning Parameter (Decision Tree)

**max_depth**

The first parameter to tune is max_depth. This indicates how deep the tree can be. The deeper the tree, the more splits it has and it captures more information about the data.

**min_samples_split**

*min_samples_split* represents the minimum number of samples required to split an internal node. This can vary between considering at least one sample at each node to considering all of the samples at each node. When we increase this parameter, the tree becomes more constrained as it has to consider more samples at each node.

# Tuning Parameter (Decision Tree)

**min_samples_leaf**

*min_samples_leaf* is The minimum number of samples required to be at a leaf node. This parameter is similar to *min_samples_splits*, however, this describe the minimum number of samples of samples at the leafs, the base of the tree.

**max_features**

max_features represents the number of features to consider when looking for the best split.

# Analytics

Machine Learning Use Cases

- Supervised Learning
  a. Use case 1
  b. Use case 2
  c. Use case 3
- Unsupervised Learning
  a. Use case 1
  b. Use case 2
  c. Use case 3

# Use Case Machine Learning

How machine learning implemented in real case

# Rekomendasi

# Google Search



| | |
|---|---|
| When to use | ✕ 🔍 |
| 🔍 when to use **spearman correlation** | ↖ |
| 🔍 when to use **chi square test** | ↖ |
| 🔍 when to use **pearson or spearman** | ↖ |
| 🔍 when to use **whom** | ↖ |
| **Google bahkan tahu pertanyaanmu :))** | |
| 🔍 when to use **in on at** | ↖ |

| | |
|---|---|
| review f | ✕ 🔍 |
| 🕐 review f**emale daily dr g brightening peeling gel** | ✕ |
| 🕐 review f**emale daily nature republic lemon peeling gel** | ✕ |
| 🕐 review f**emale daily mizon apple smoothie peeling gel** | ✕ |
| 🔍 review f**oundation wardah** | ↖ |
| 🔍 review f**ilm target** | ↖ |

# Credit Scoring



Poor    723    Excellent

300     850

| Payment History | Age & Type of Credit | % of Credit Limit Used | Total Balances/Debt | Recent Credit Behavior | Available Credit |
|---|---|---|---|---|---|
| Extremely Influential | Highly Influential | Highly Influential | Moderately Influential | Less Influential | Least Influential |

# Bot

# Study Case

# Data Science Indonesia

**Data Science Indonesia (DSI)**, adalah sebuah Komunitas yang didirikan pada bulan Mei 2015 yang terdiri dari sekumpulan ilmuwan, seniman dan pembelajar yang ingin membangun budaya Data Driven di Indonesia dengan menginspirasi, mengajarkan serta menawarkan nilai dari sebuah data melalui pendekatan Data Science

Visi Kami:
•Bersama masyarakat menciptakan ekosistem inovasi berbasis data untuk meningkatkan kesejahteraan masyarakat

Misi :
•Menjadi mitra bagi sektor publik maupun swasta untuk mengeskplorasi pendekatan data science dalam mencari solusi atas tantangan yang ada
•Meningkatkan pengetahuan dan kesadaran masyarakat terhadap data science
•Menjadi wadah bagi masyarakat untuk berjejaring dalam konteks pemanfaatan data science

DATA SCIENCE INDONESIA

http://datascience.or.id/daftar-anggota-dsi/

Contact:

contact@datascience.or.id