

BUKALAPAK

Natural Language Processing

Basic Class

Afif A. Iskandar



About Me



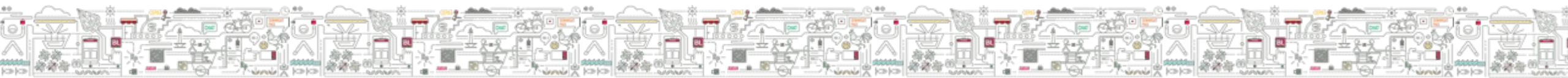
Name : Afif Akbar Iskandar

Role : AI Scientist

Company : Bukalapak

Specialization :

- Computer Vision
- Machine Learning
- Deep Learning
- Natural Language Processing



About Me



Educational Background :

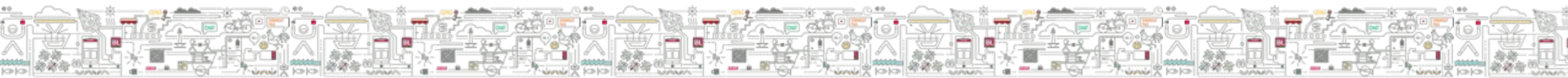
- Bachelor of Mathematics at Universitas Indonesia (2011)
- Master of Computer Science at Universitas Indonesia (2015)

Working Experience :

- Data Scientist (2015-Now)

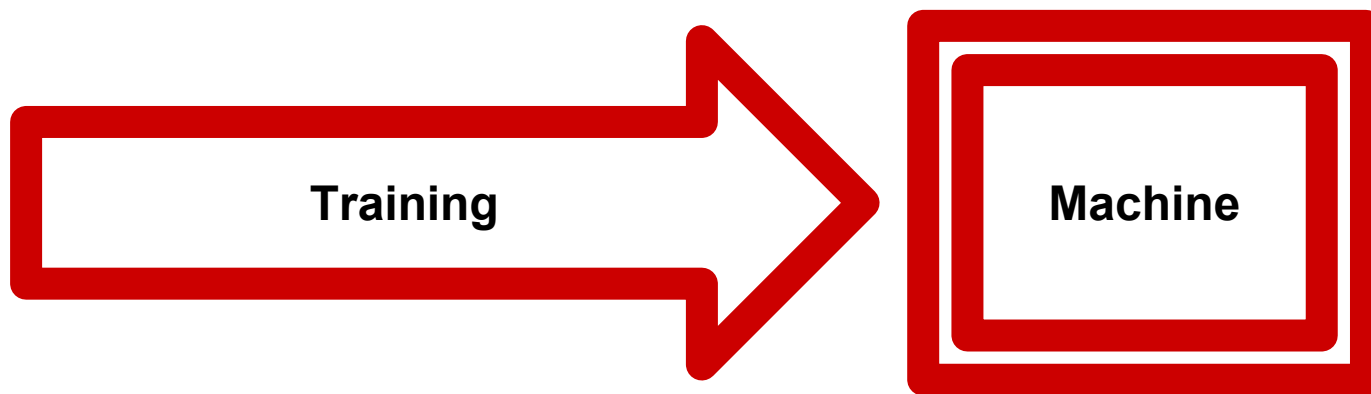
OUTLINE

- Machine Learning Review
- Text Preprocessing
- Text Classification (Spam Detector)

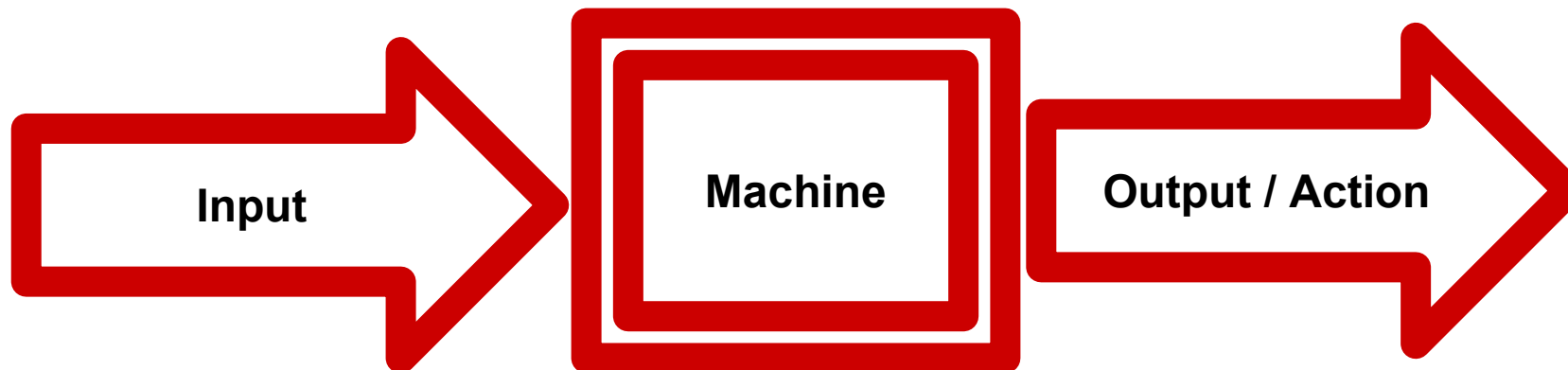


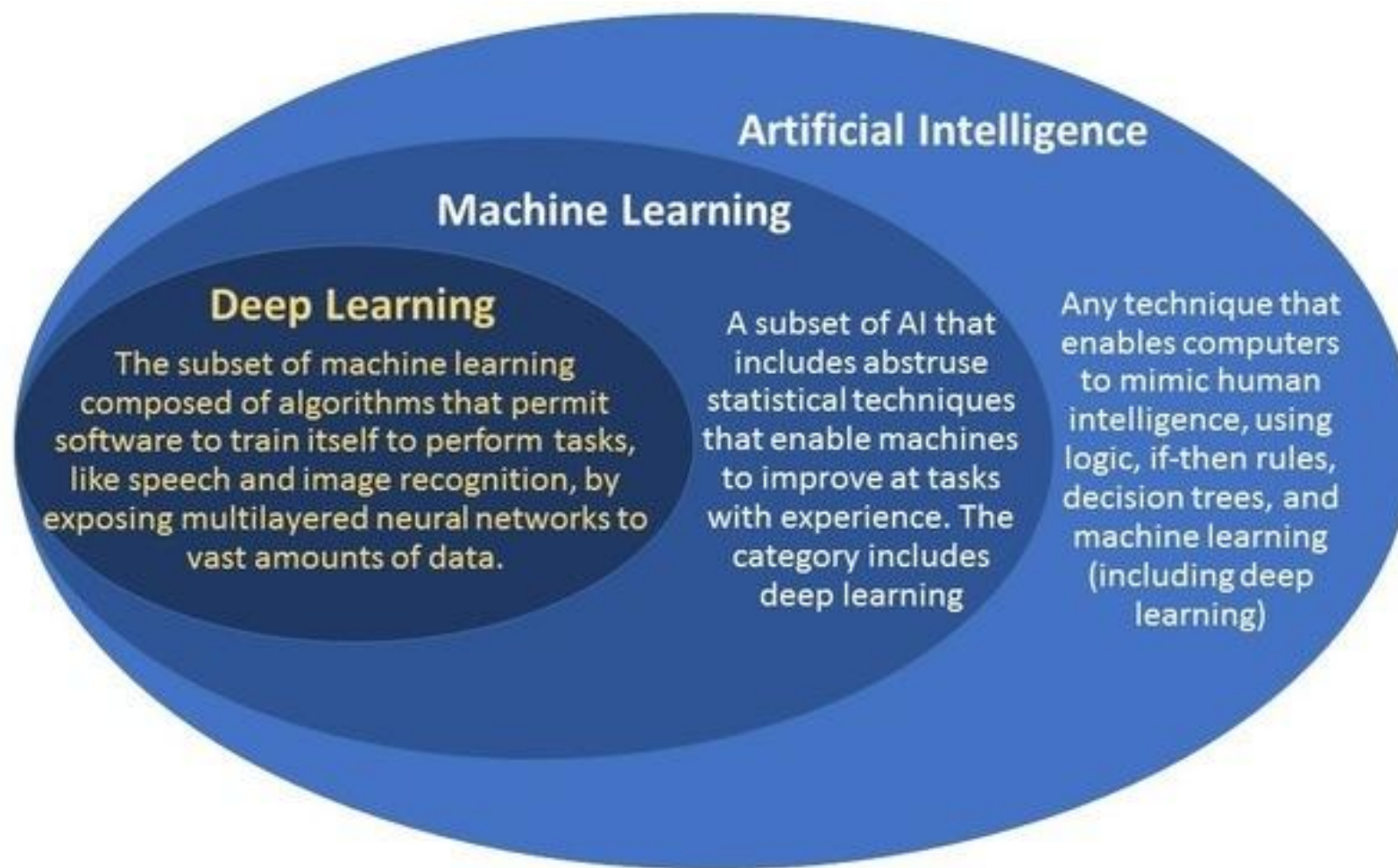
A very simplistic view of Machine Learning

Let's get the machine to learn stuff, by training it thousands times, million times, billion times



...after it finished learning from the intensive training...





Type of Machine Learning



Supervised Learning



Laptop

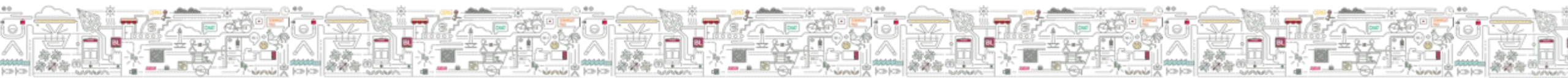


Tablet

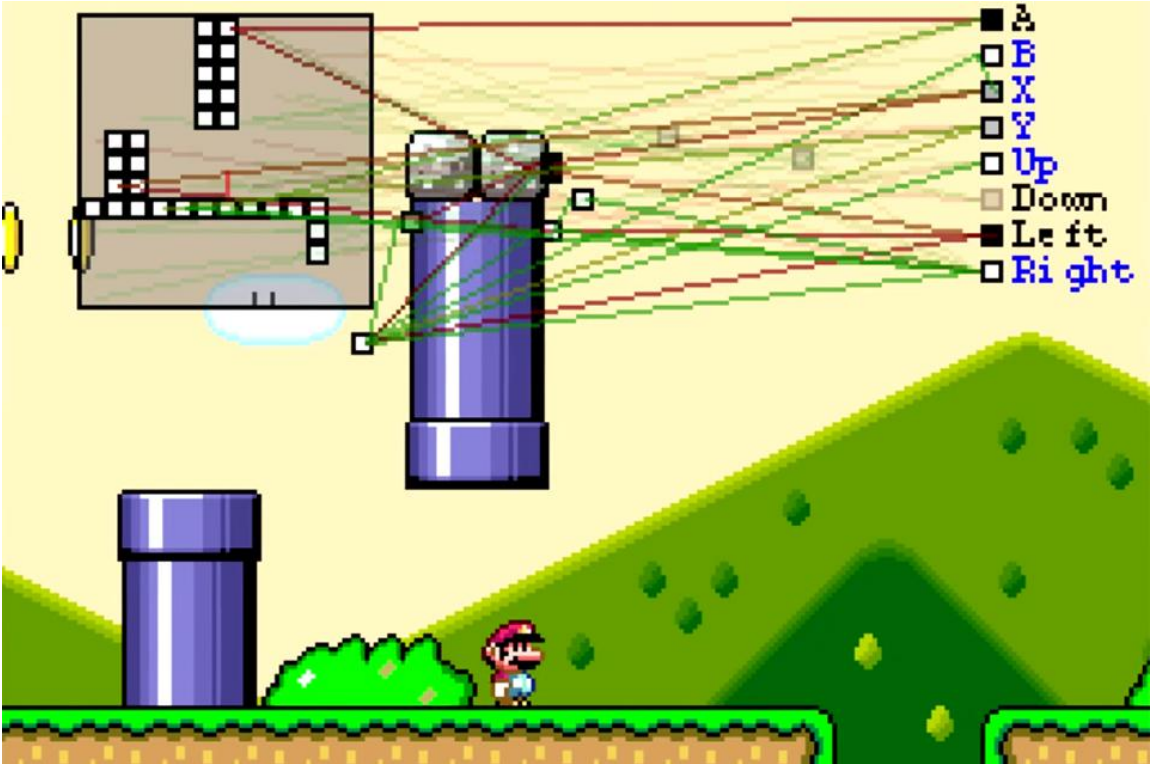
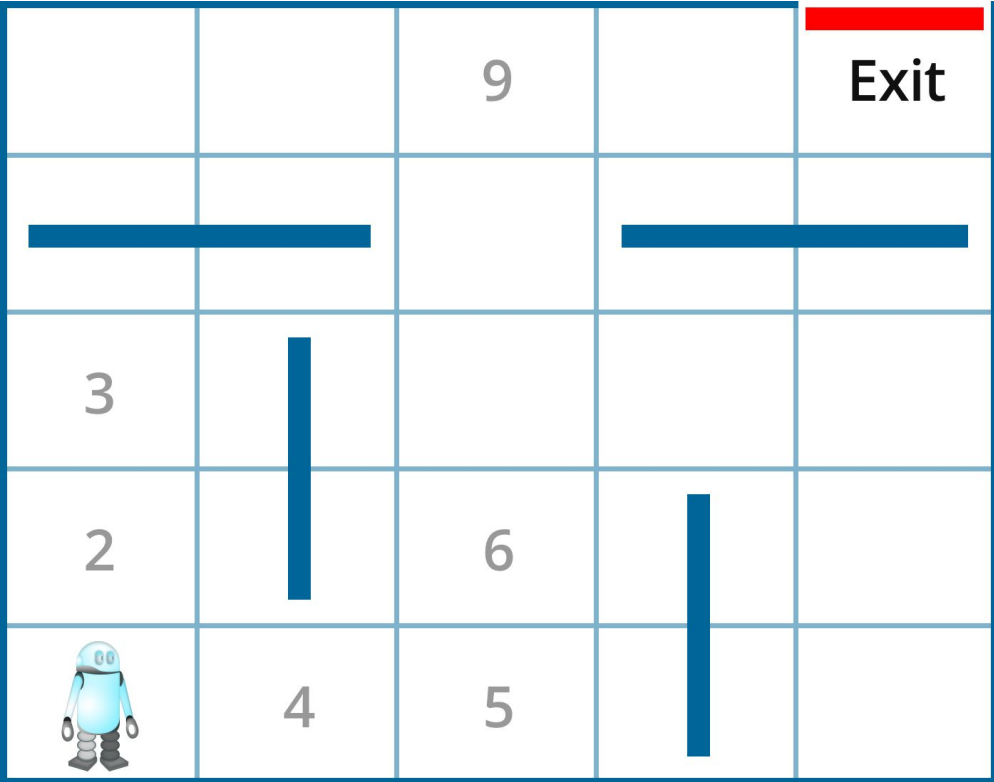


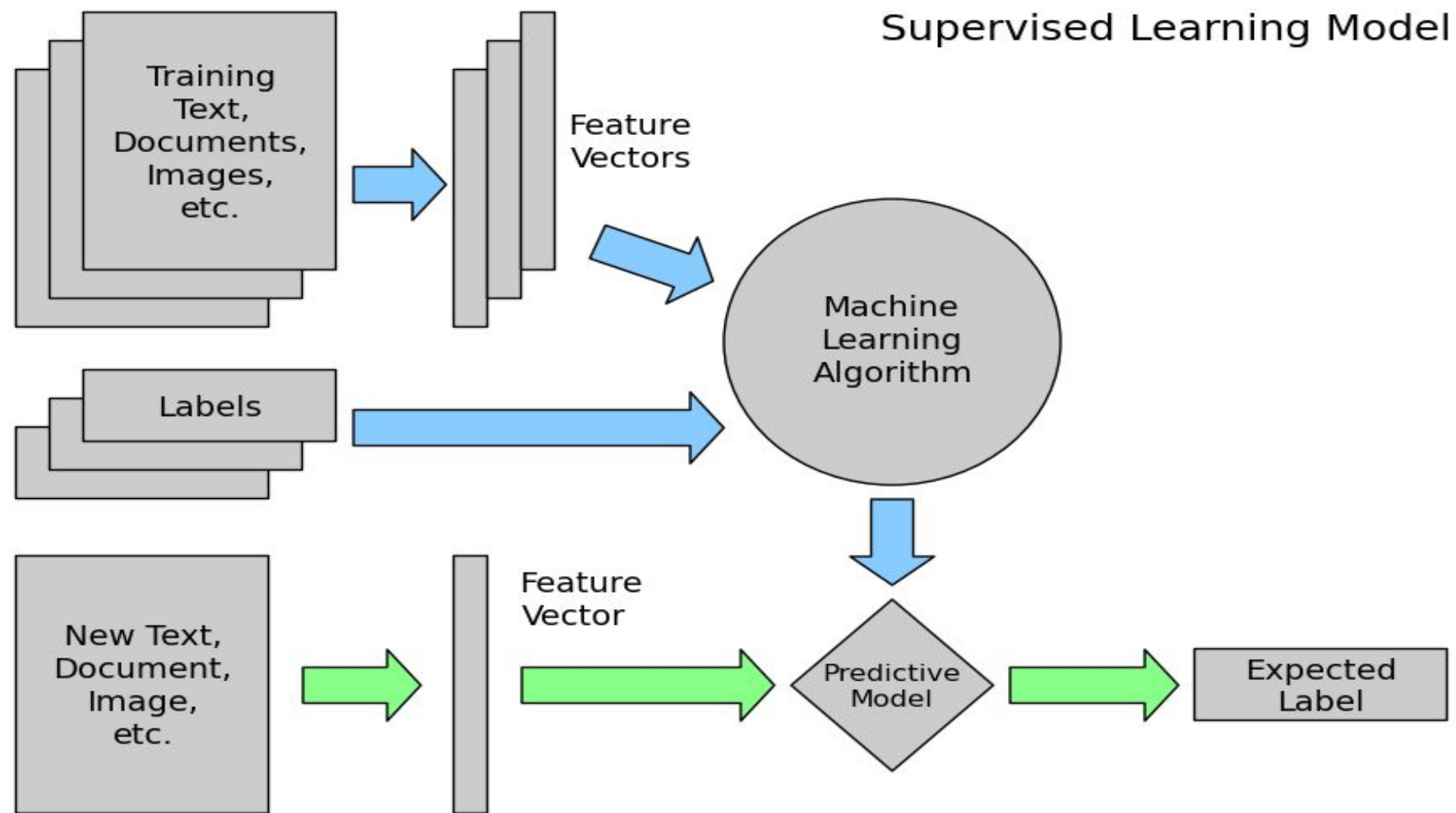
???

Unsupervised Learning



Reinforcement Learning



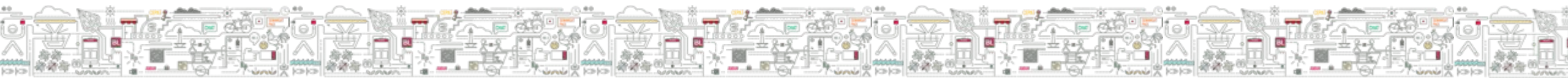


Extracting Features from Text



Bag of Words Model

- Count Vectorizer
- Term Frequency–Inverse Document Frequency



Tf-Idf

- TF: Term Frequency, which measures how frequently a term occurs in a document

$TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document})$

- IDF: Inverse Document Frequency, which measures how important a term is

$IDF(t) = \log_e(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it})$



Spam Detector using Supervised Learning

Go to Jupyter Notebook



Thank You

Bukalapak

