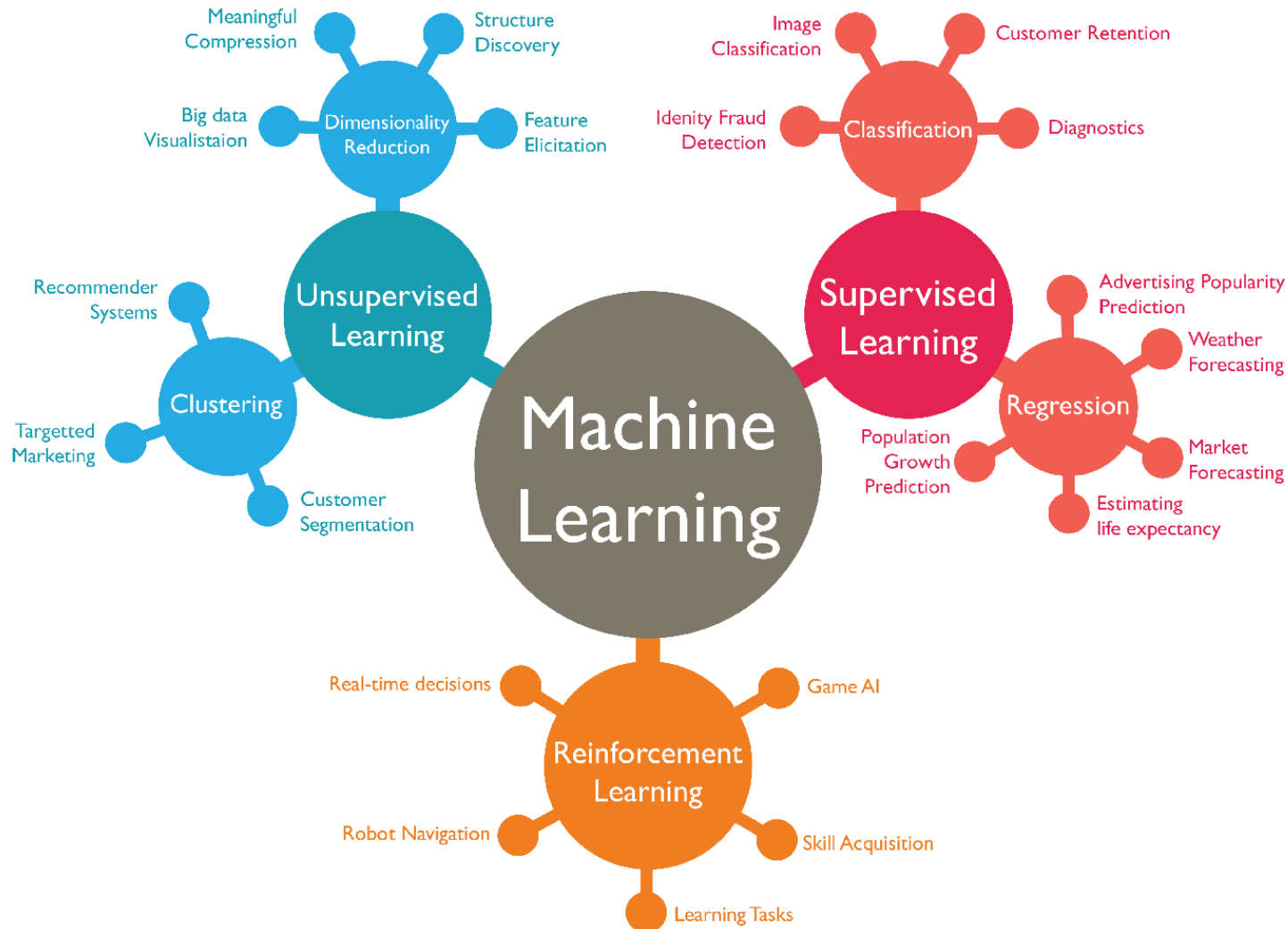


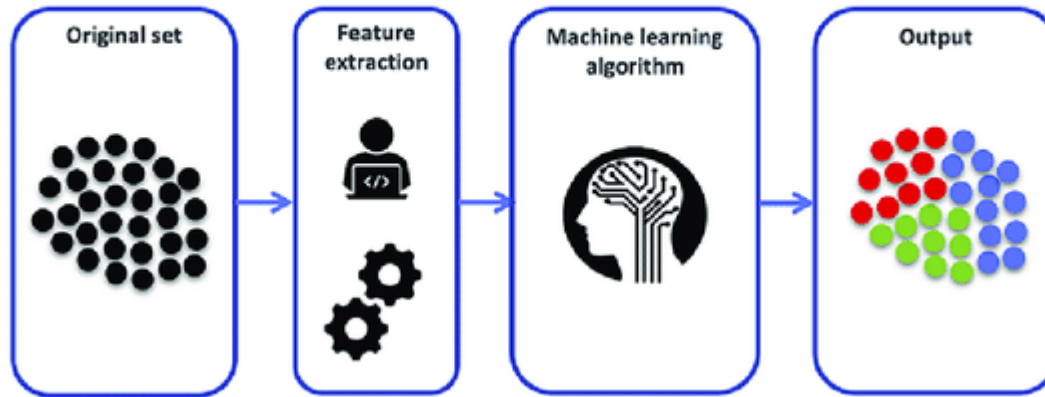
Today Focus



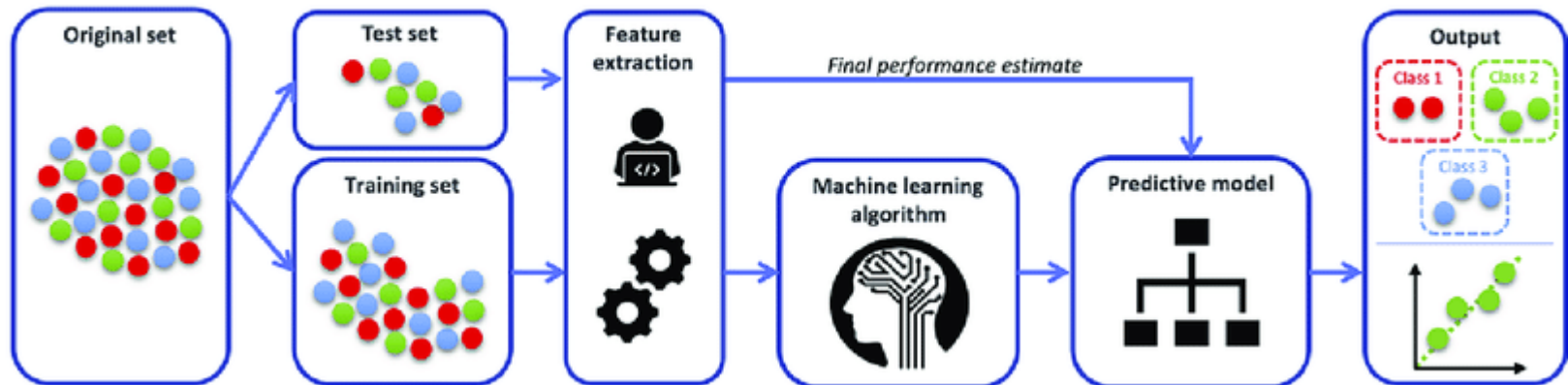


ML Process

UNSUPERVISED LEARNING

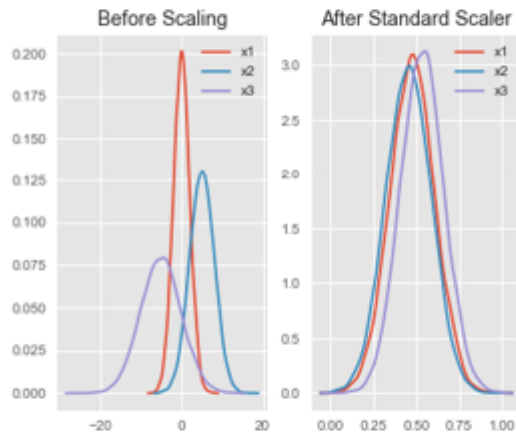


SUPERVISED LEARNING

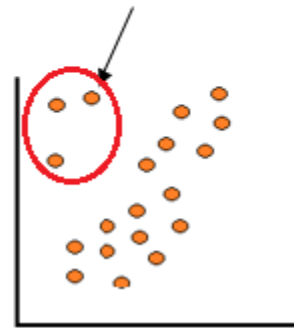


Today Focus

Normalization



Outliers



Imputation

	First	Second	Third
0	100.0	30.0	NaN
1	90.0	45.0	40.0
2	NaN	55.0	80.0
3	95.0	NaN	98.0

Encoding

Food Name	Apple	Chicken	Broccoli
Apple	1	0	0
Chicken	0	1	0
Broccoli	0	0	1

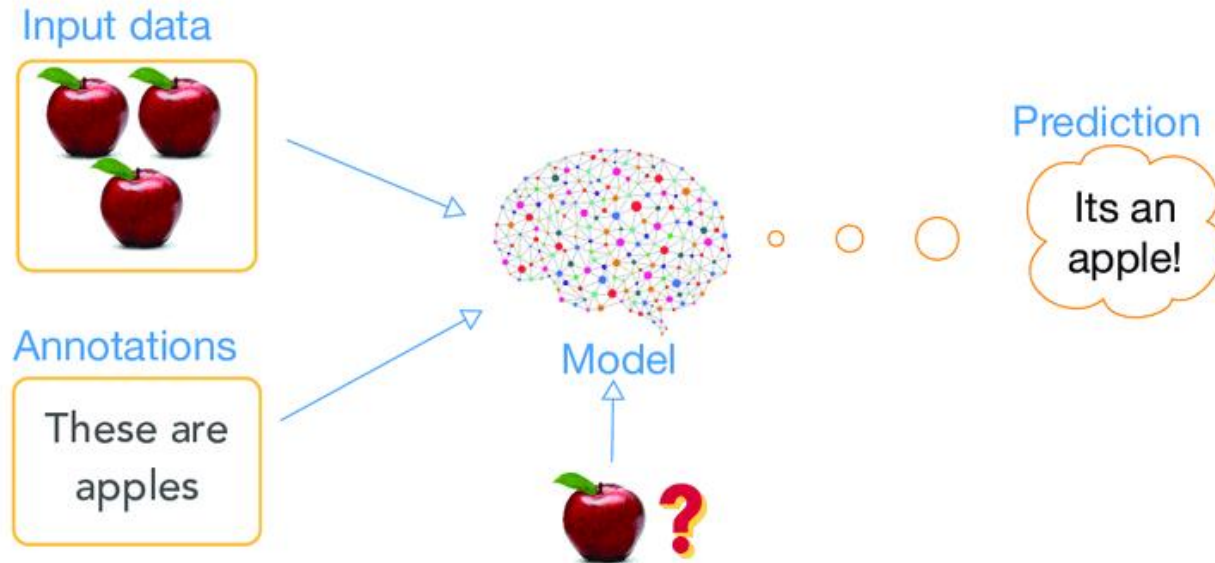
Data Preprocessing



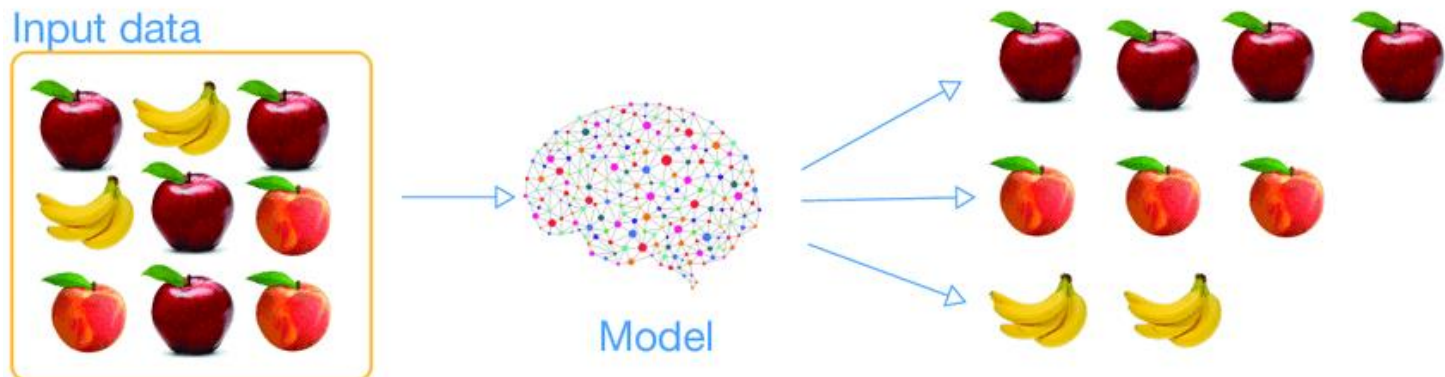
INVESTIC

Classified by Type

supervised learning



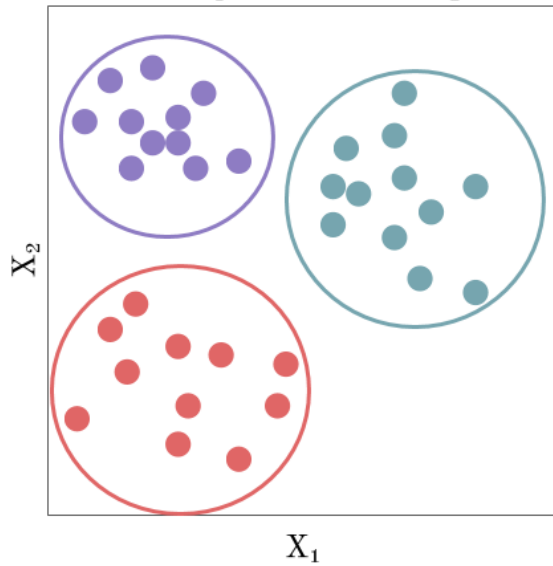
unsupervised learning



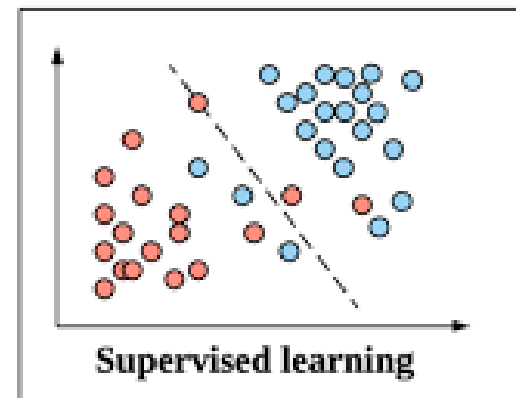
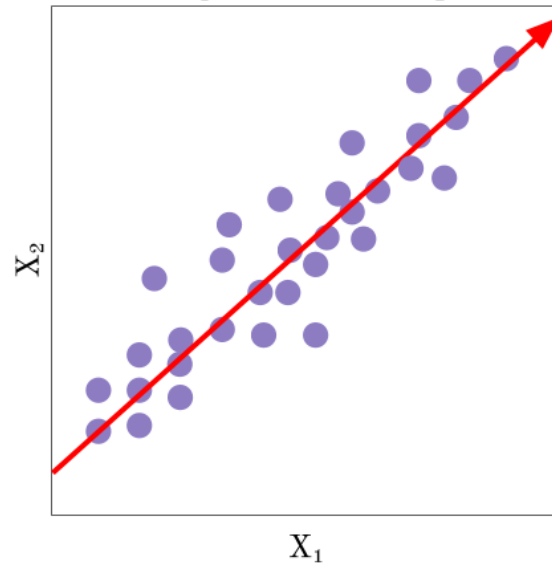
INVESTIC

Which is which

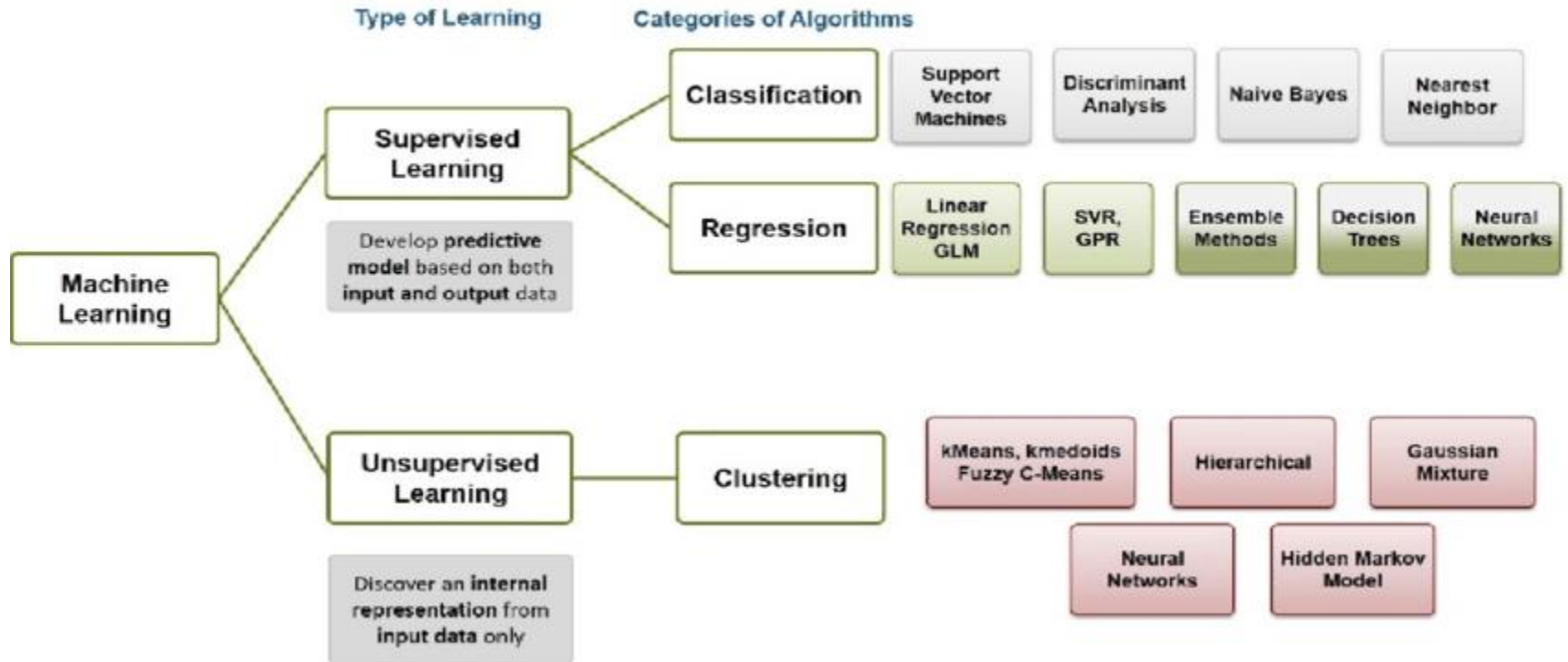
Unsupervised Learning



Supervised Learning



Regression & Classification

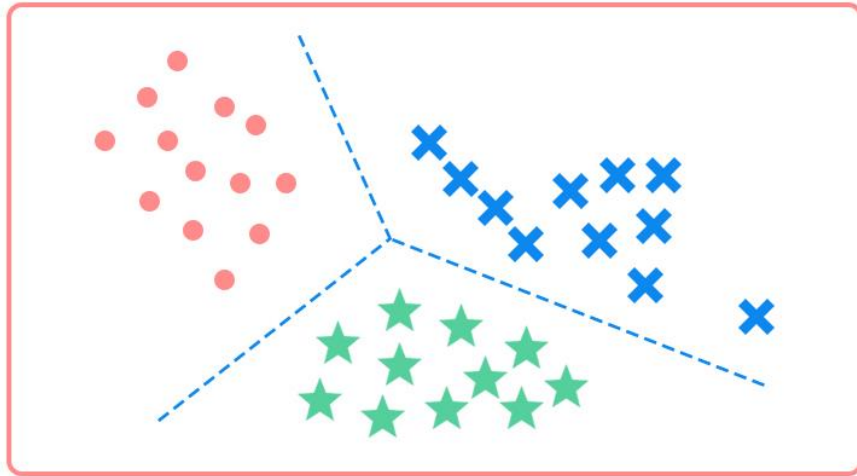


Classification vs Clustering



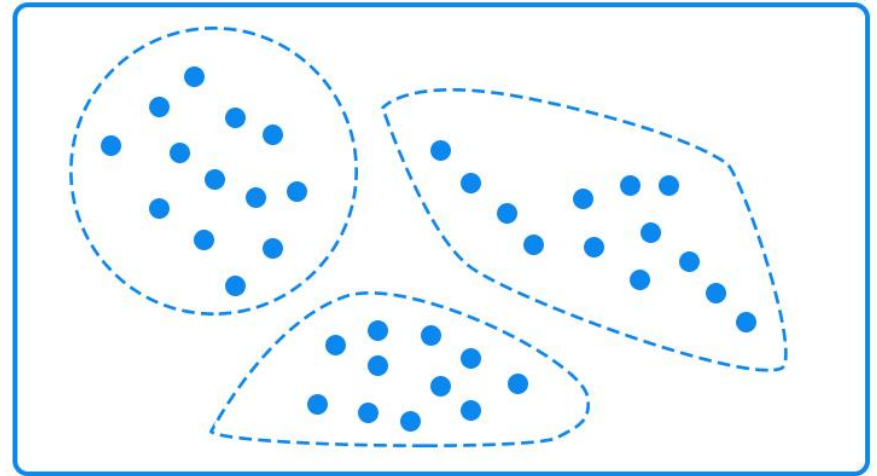
Supervised vs. Unsupervised Learning

Classification



Supervised learning

Clustering

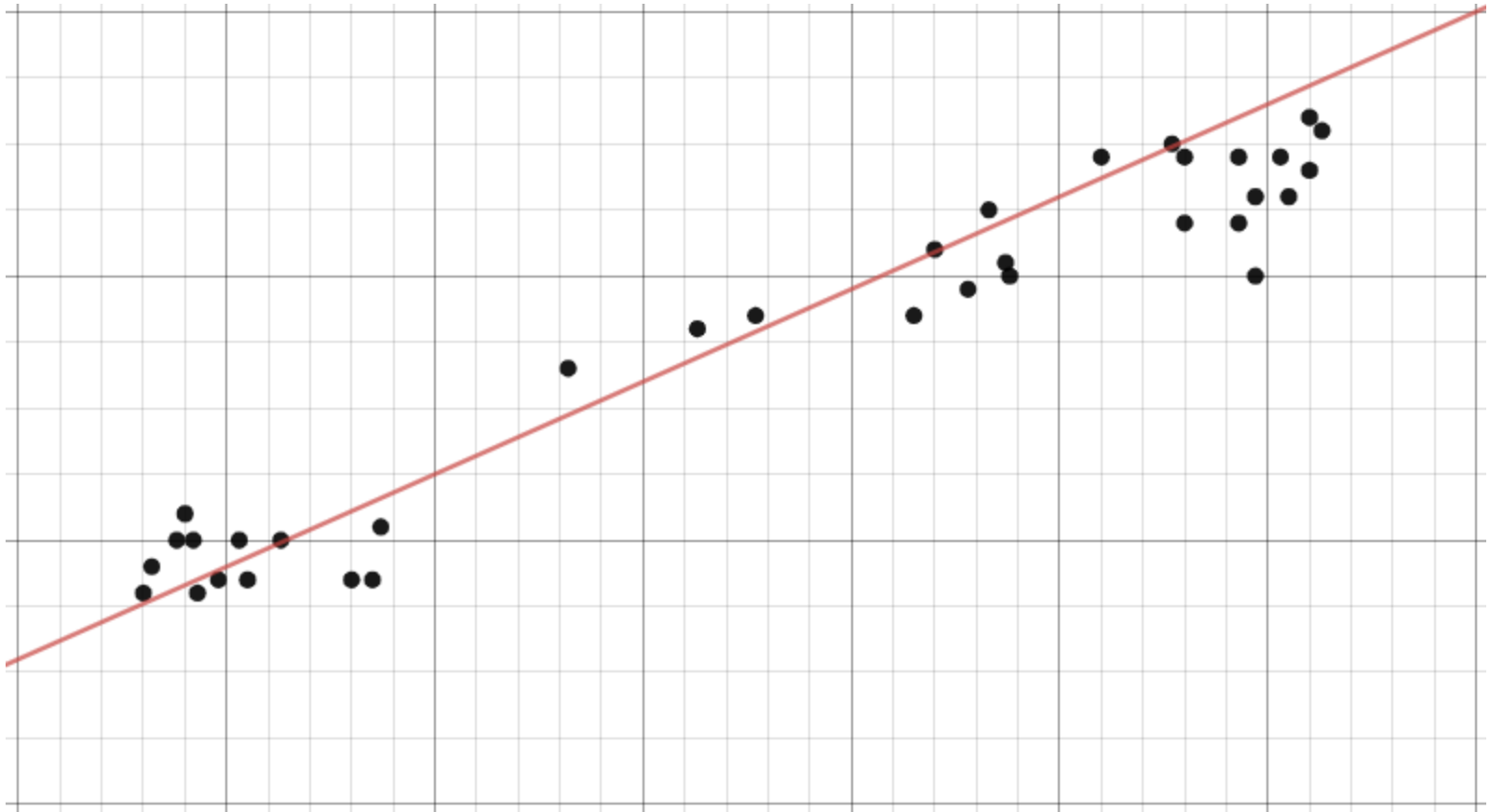


Unsupervised learning



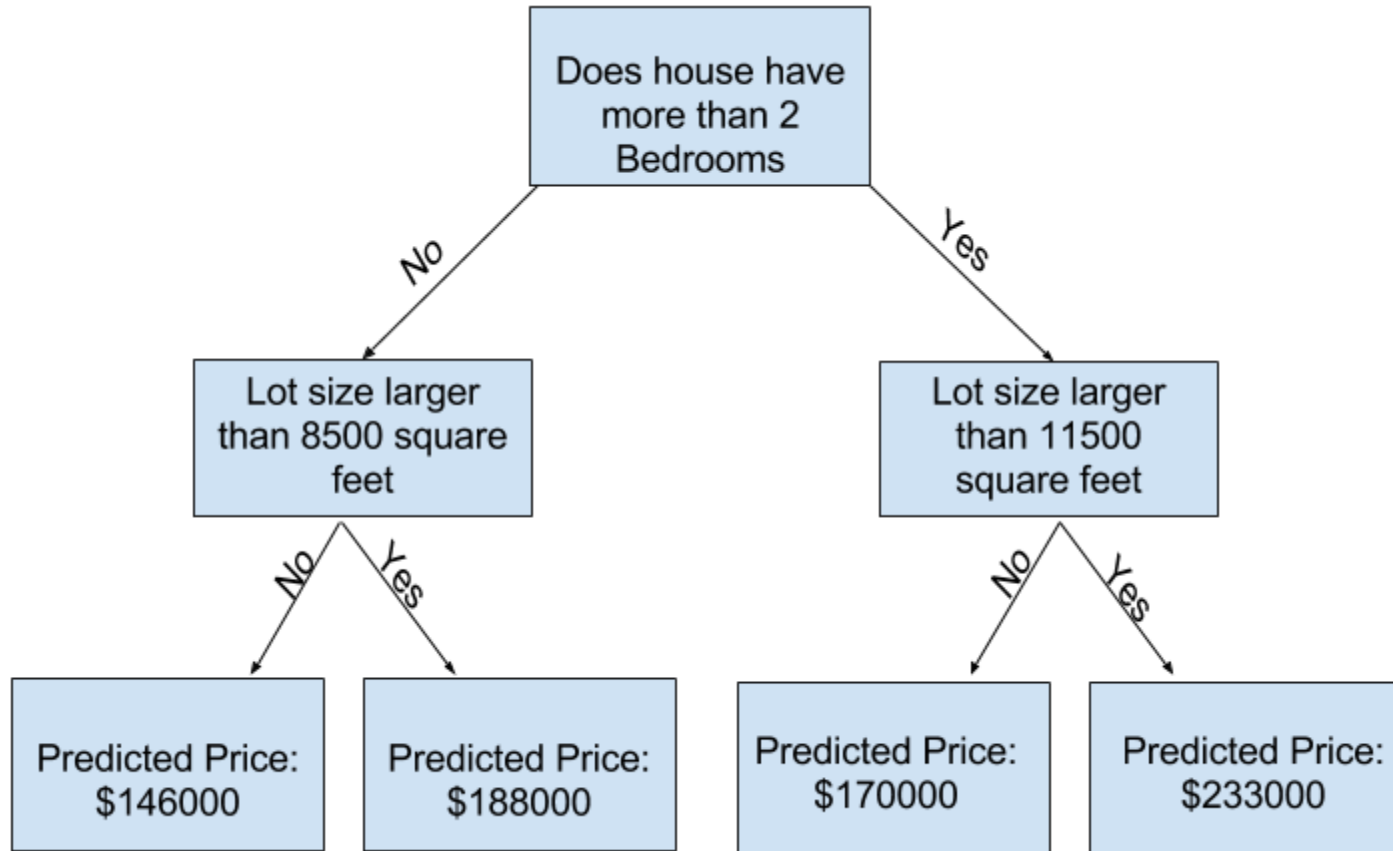


Regression Learning Model



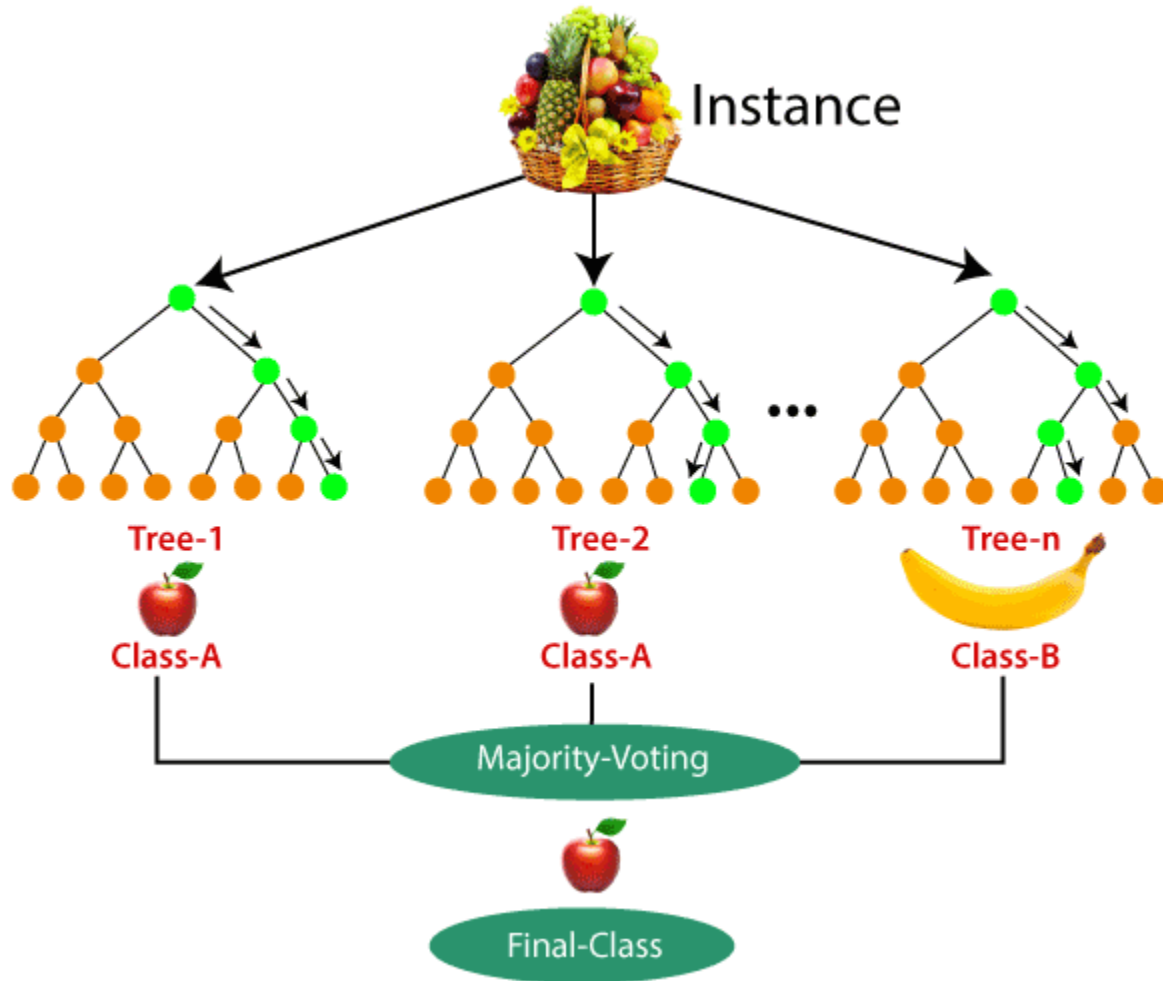


Regression Learning Model



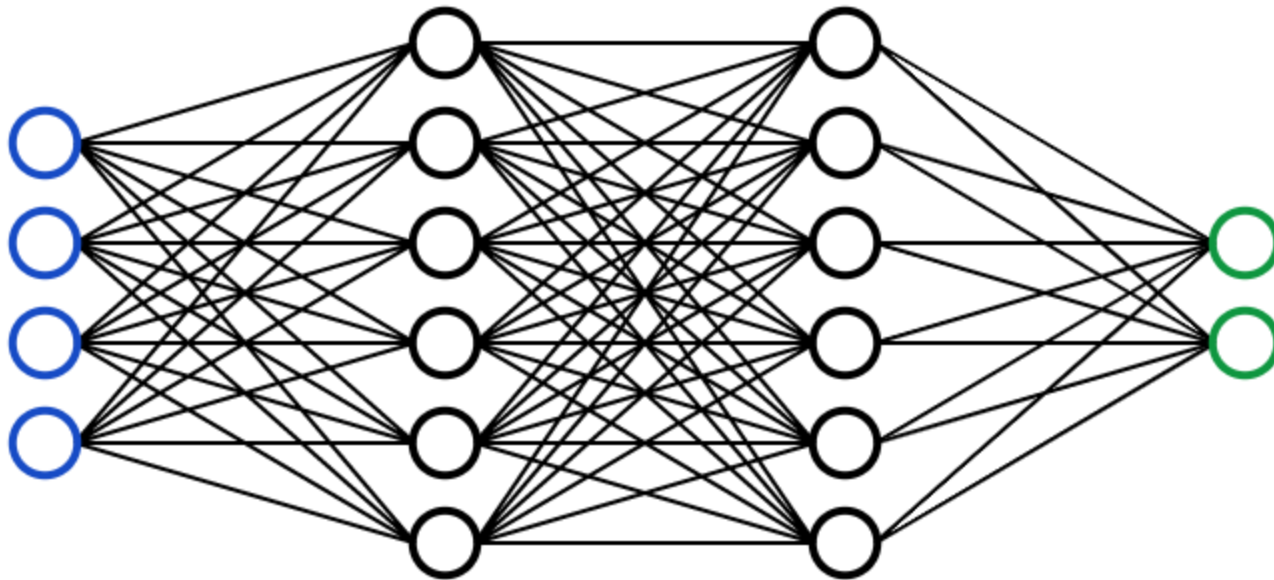


Regression Learning Model



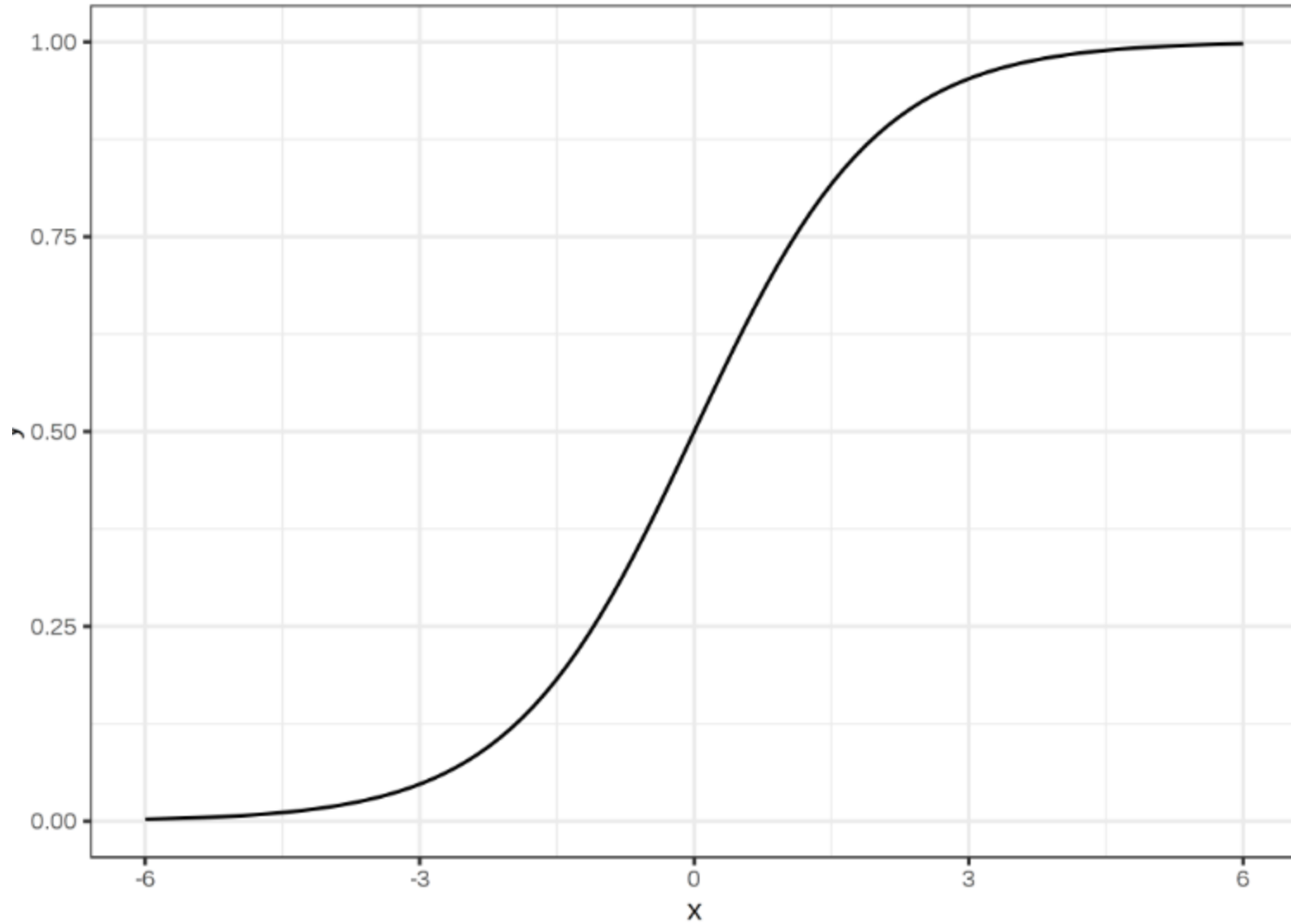


Regression Learning Model



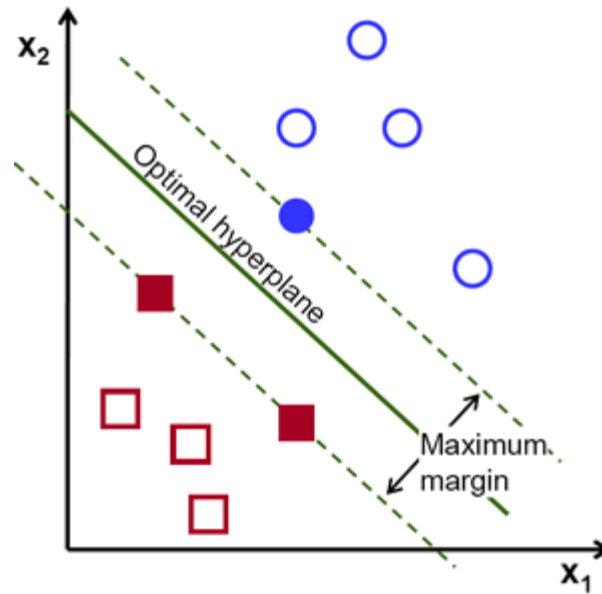


Classification Learning Model



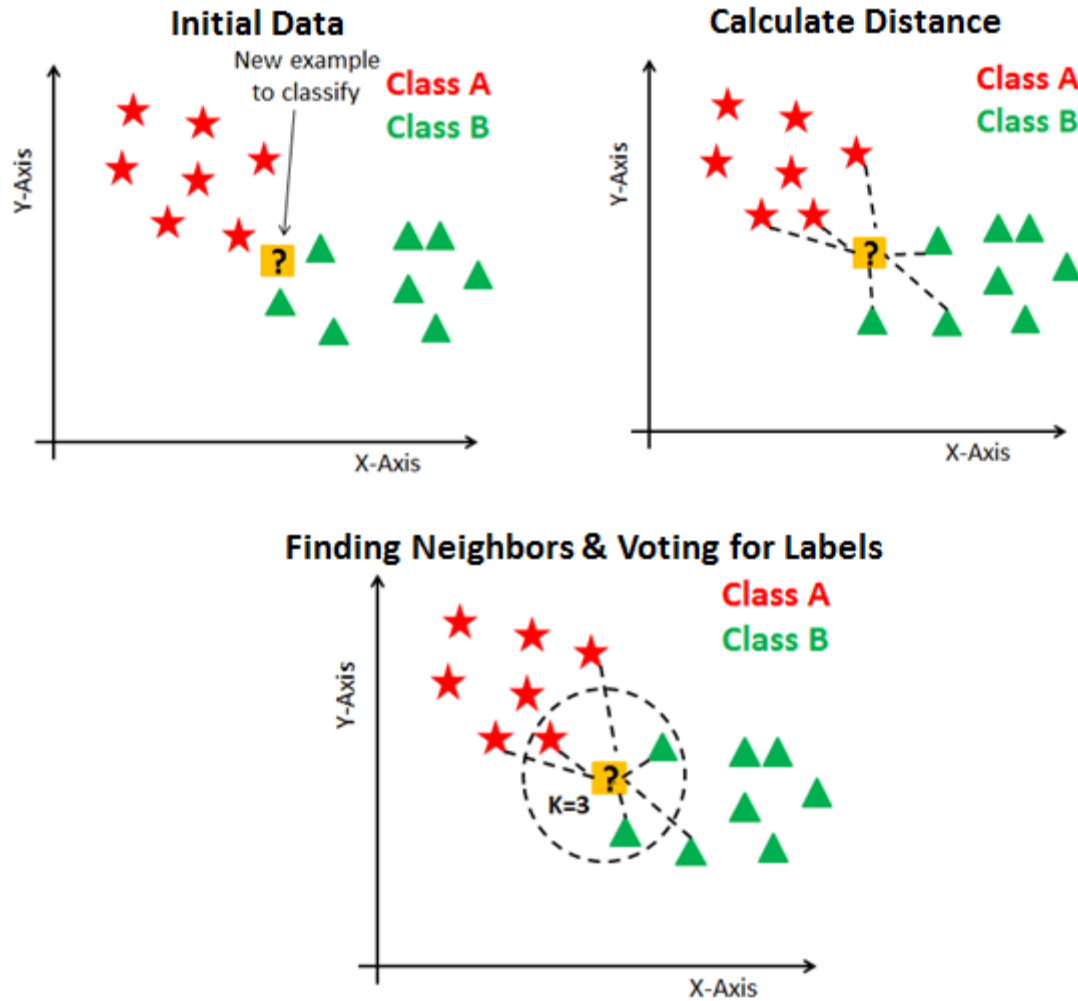


Classification Learning Model





Classification Learning Model





Classification Learning Model

Naive Bayes

$$P(y|X) = \frac{P(X|y) * P(y)}{P(X)}$$

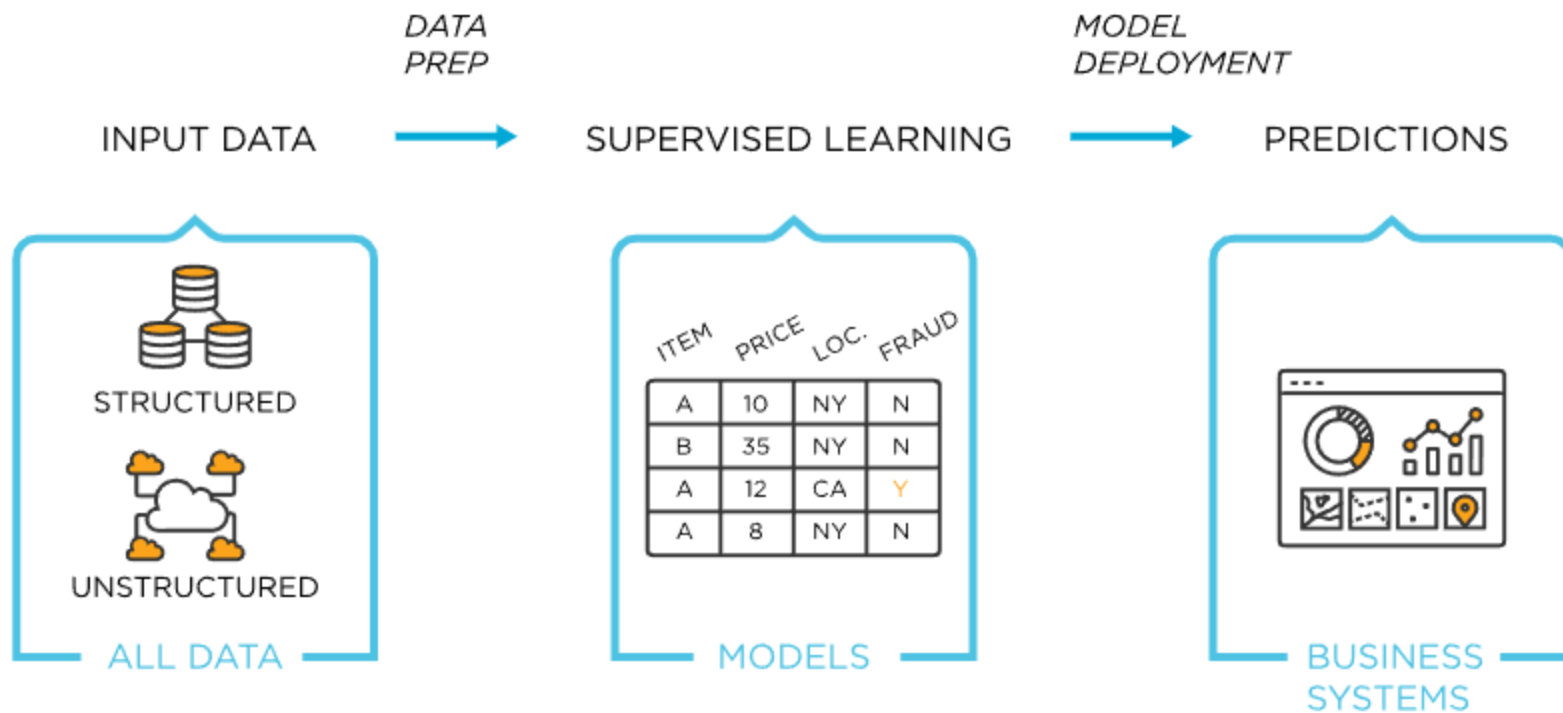
“What is the probability of y given X?”

$$P(y|X) \propto P(X|y) * P(y)$$

The goal is to find the class y with the maximum proportional probability.

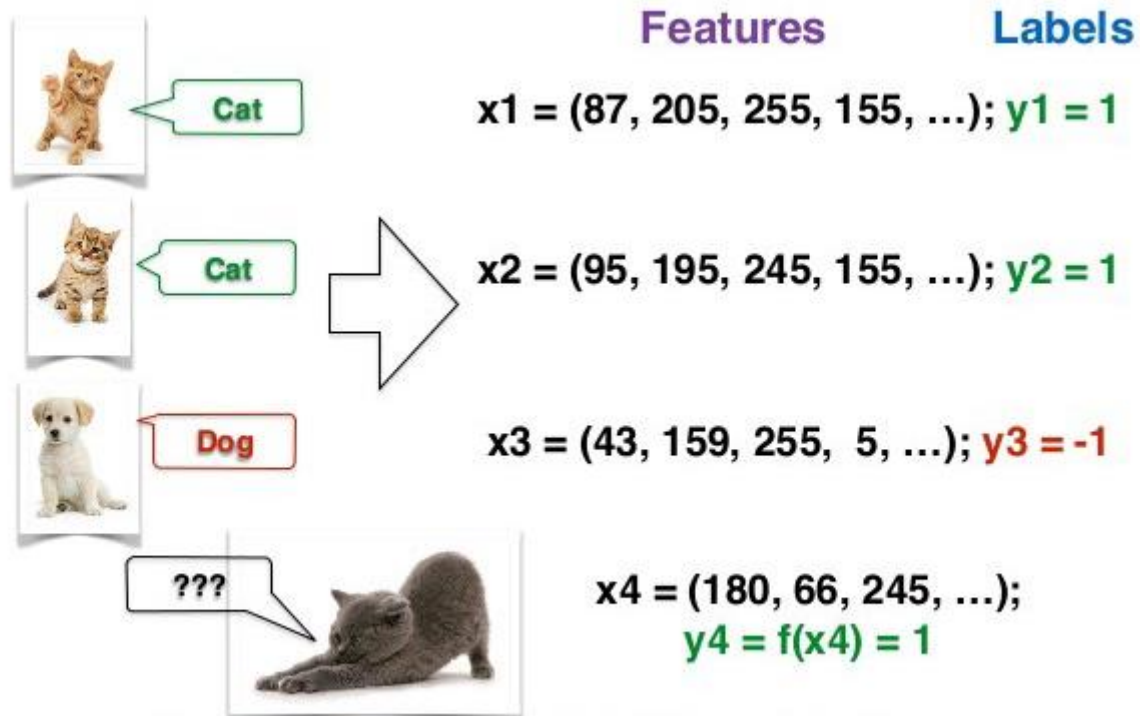


Supervised Learning





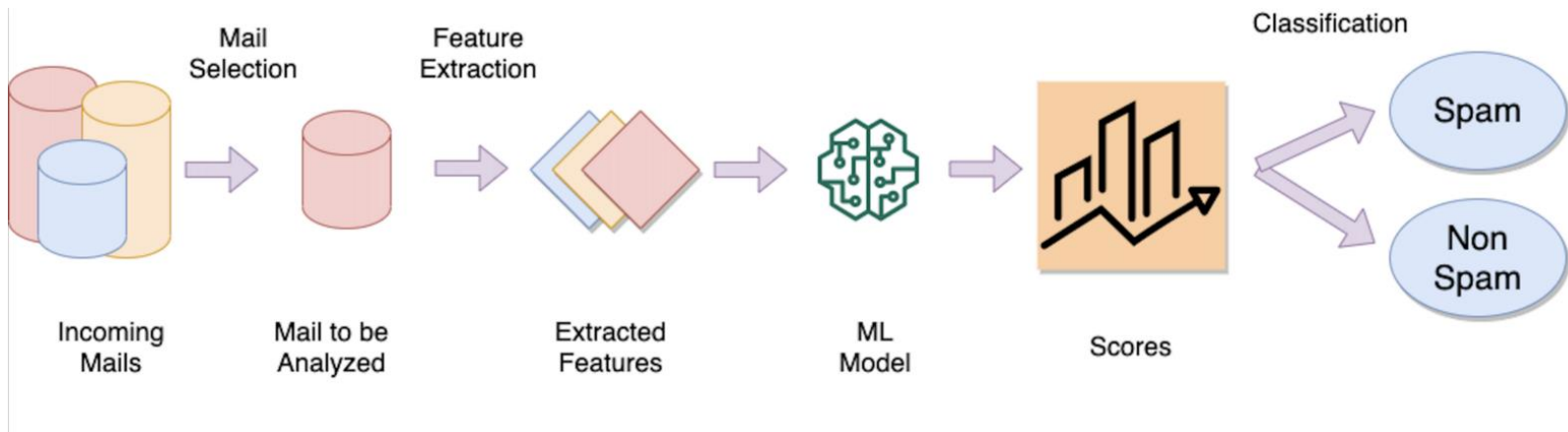
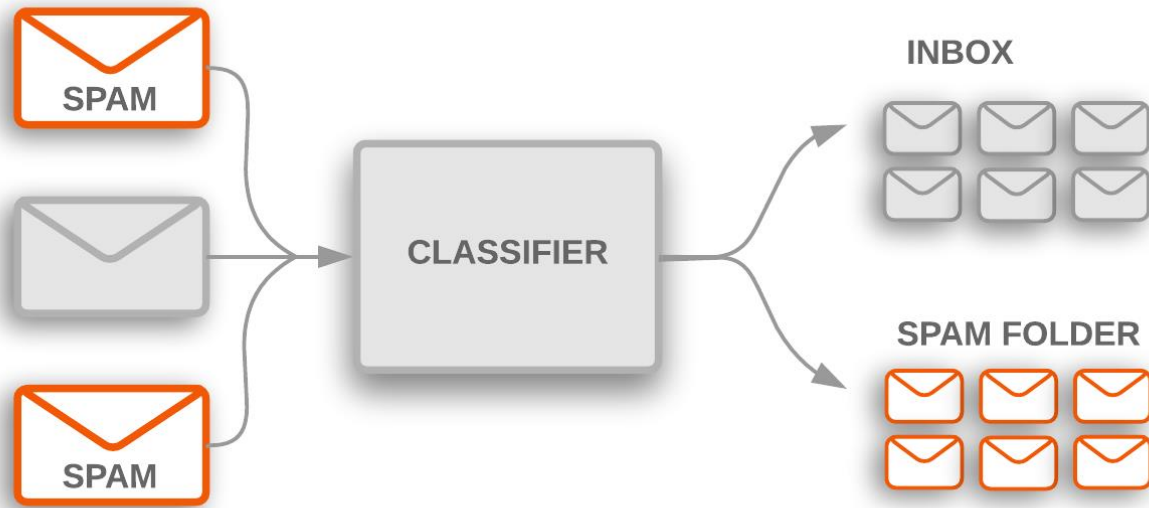
Supervised Learning



Adopted from P.Vincent http://videlectures.net/deeplearning2015_vincent_machine_learning/



Complete picture





ML Process

—

1. Prepare & Transform Data
 1. Normalizing/standardizing
 2. Label Encoding
2. Modeling
3. Evaluating





ML Process

–

1. Prepare & Transform Data

1. Outlier

2. Feature Scaling

1. Normalization

2. Standardization

3. Label encoding/ One hot encoding





ML Process Simple

–

1. Prepare & Transform Data
2. Model
3. Rumble !!
4. Backtest



Example



ML Process

–

1. Prepare & Transform Data
 1. Outlier
 2. Feature Scaling
 1. Normalization
 2. Standardization
 3. Label encoding/ One hot encoding
2. Split Train/ Test Data

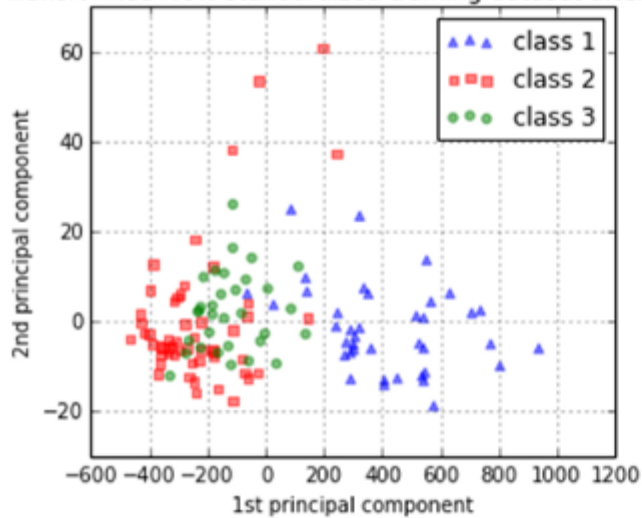




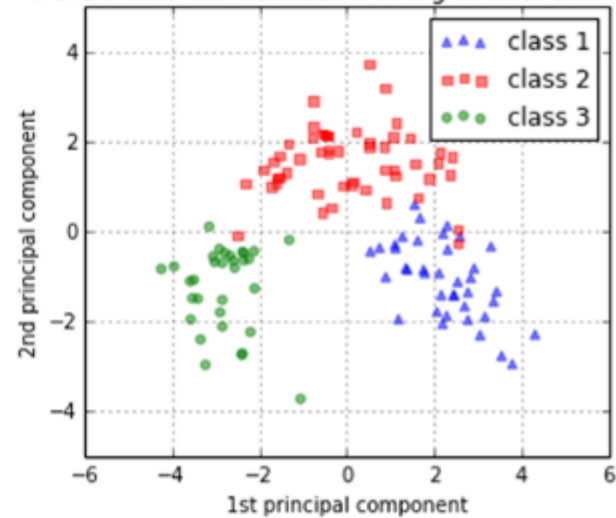
Is feature scaling that matter??

—

Transformed NON-standardized training dataset after PCA

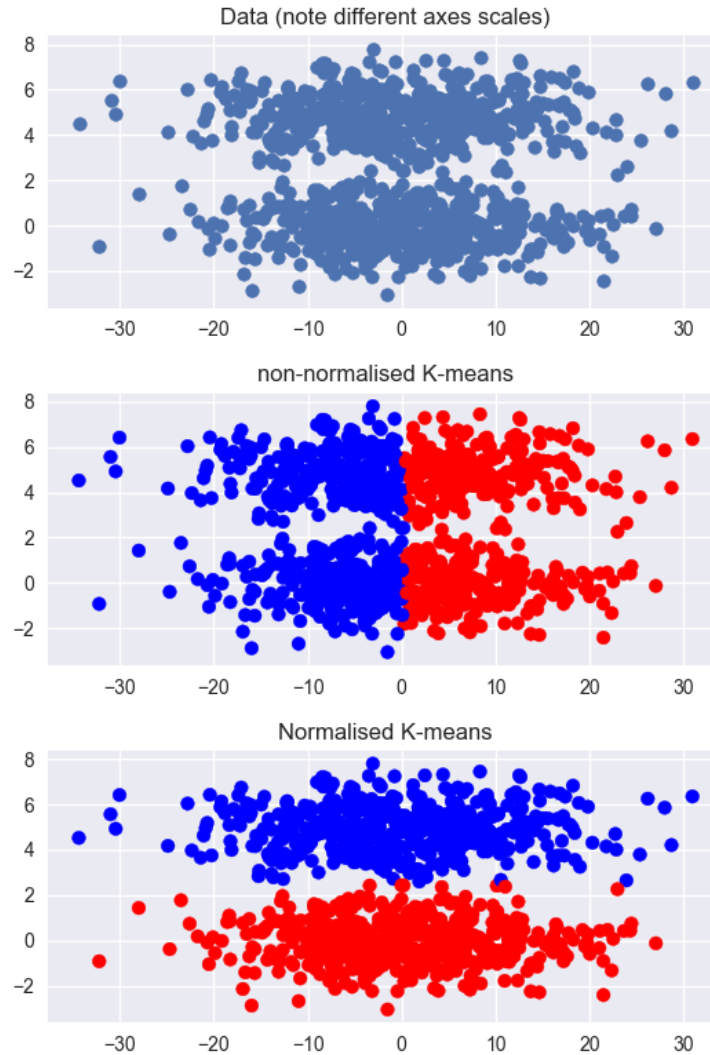


Transformed standardized training dataset after PCA





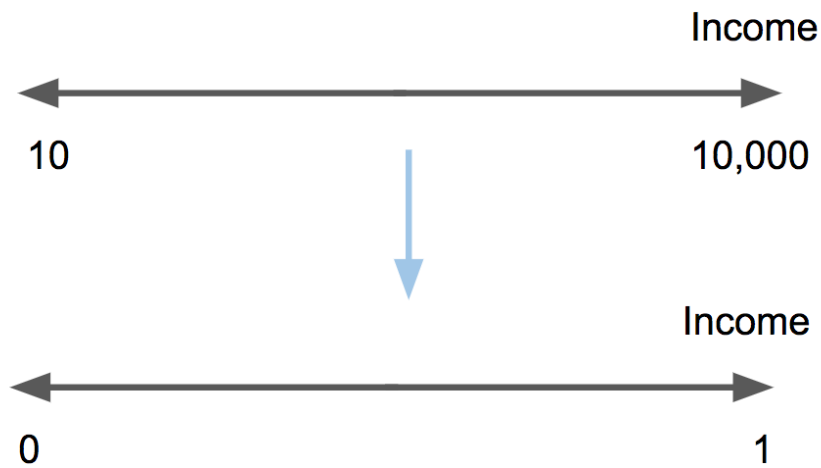
Is feature scaling that matter??





Feature Scaling

Normalize



Standardize

$$z = \frac{x - \mu}{\sigma}$$

μ = Mean

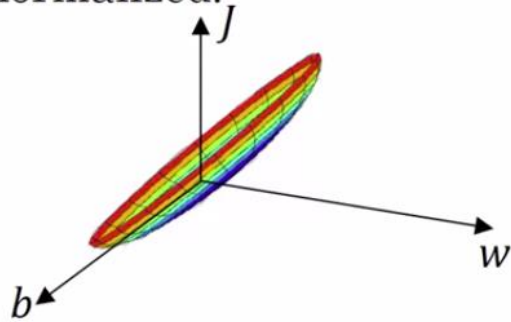
σ = Standard Deviation



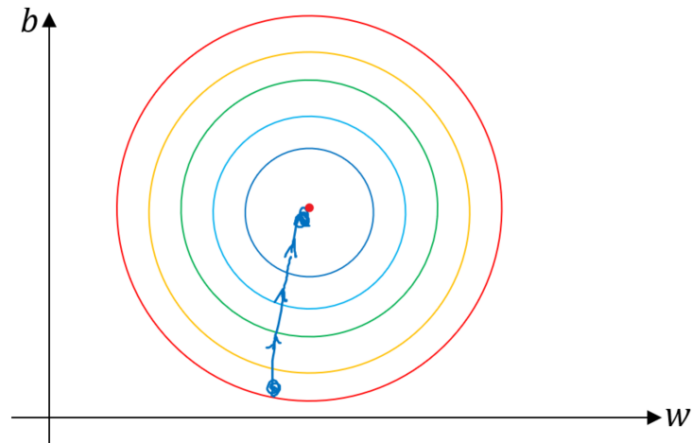
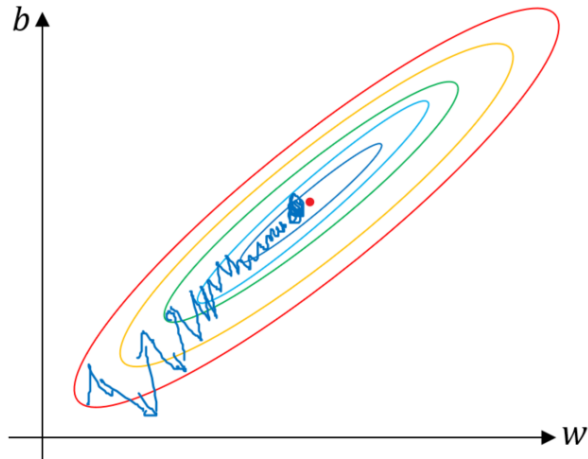
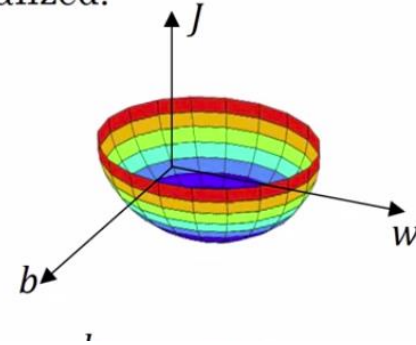


Why Normalized

Unnormalized:

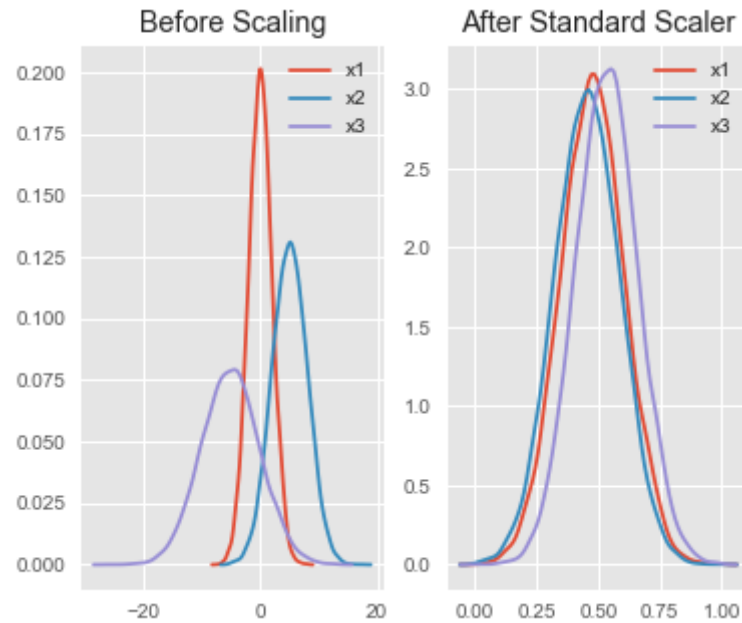


Normalized:



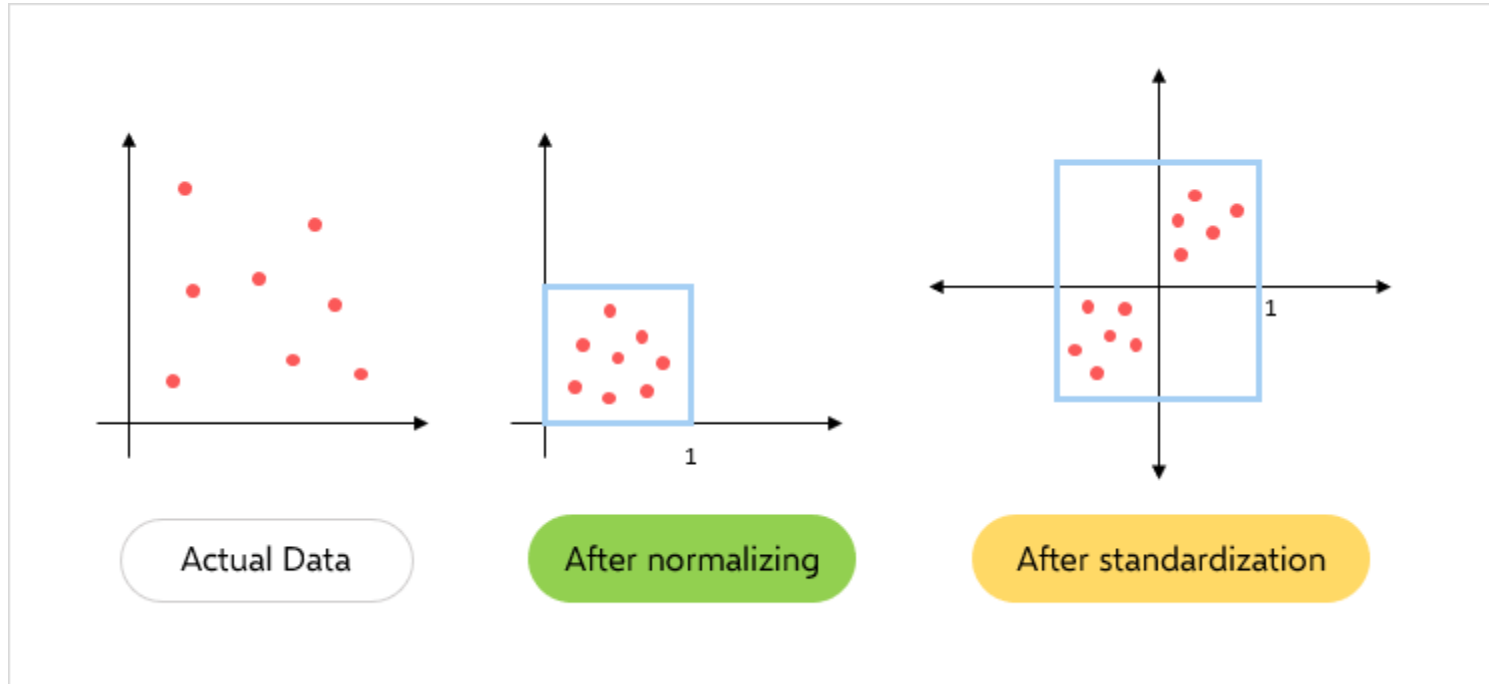


Why Standardized





Normalized vs Standardized





Normalized vs Standardized

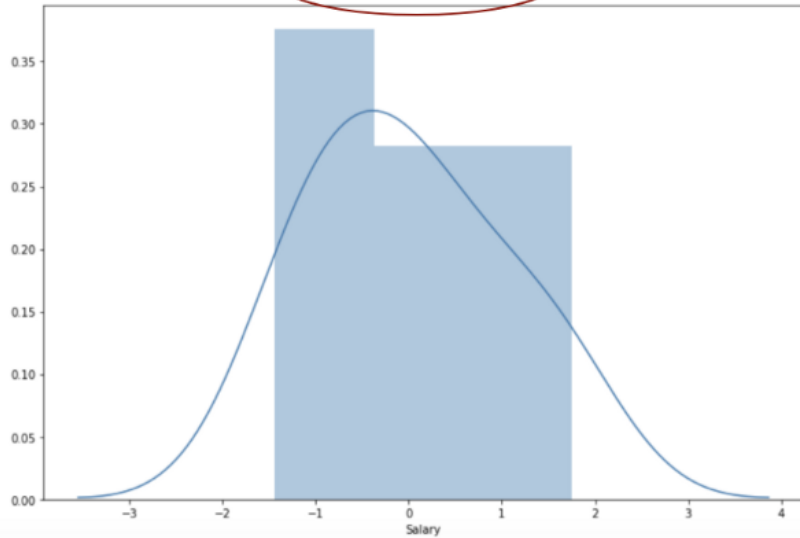
—

Column: Salary

Standard Deviation (Salary):
Max-Min Normalization (0.33) < Standardisation (1.05)

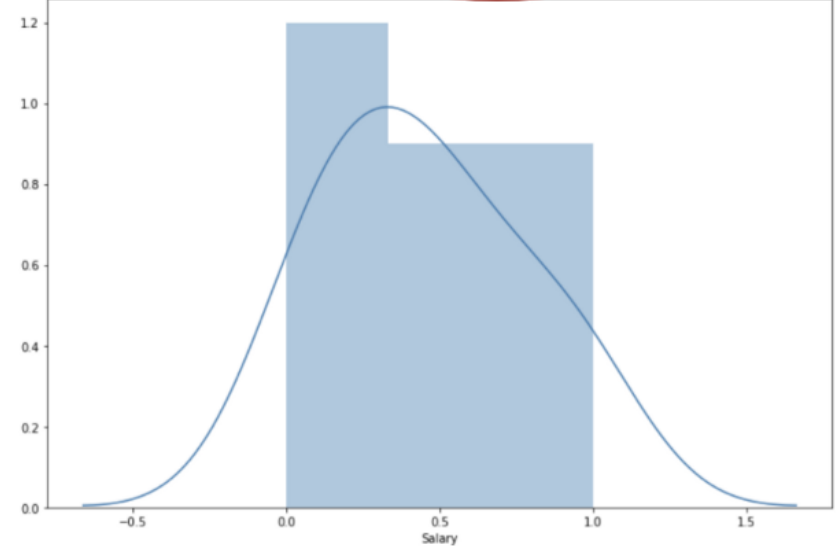
Standardisation

Standard Deviation of sc_Salary is 1.0540925533894598



Max-Min Normalisation

Standard Deviation of df_MinMax_Salary is 0.33040284015892535



INVESTIC



Sklearn Class

Name	Sklearn_Class
StandardScaler	StandardScaler
MinMaxScaler	MinMaxScaler
MaxAbsScaler	MaxAbsScaler
RobustScaler	RobustScaler
QuantileTransformer-Normal	QuantileTransformer(output_distribution='normal')
QuantileTransformer-Uniform	QuantileTransformer(output_distribution='uniform')
PowerTransformer-Yeo-Johnson	PowerTransformer(method='yeo-johnson')
Normalizer	Normalizer





Normalized vs Standardized

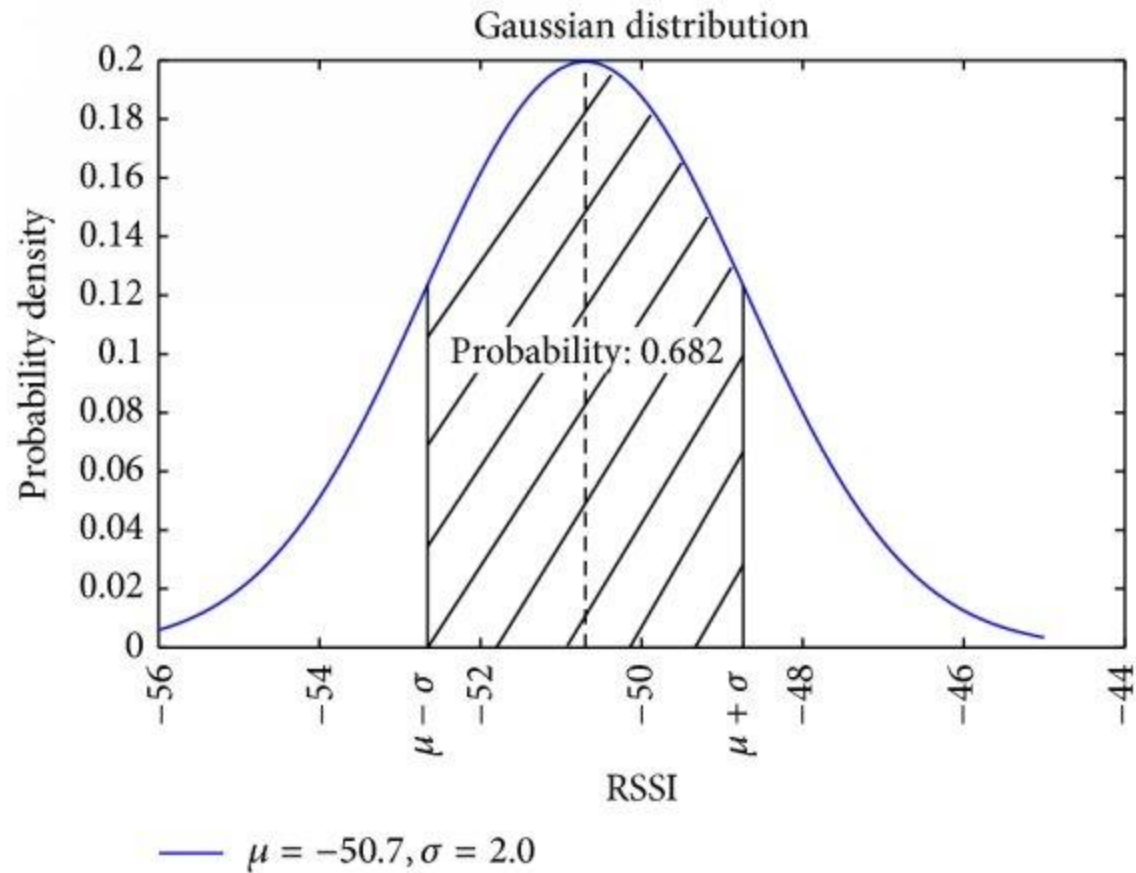
–
Normalization is good to use when you know that the distribution of your data does not follow a Gaussian distribution

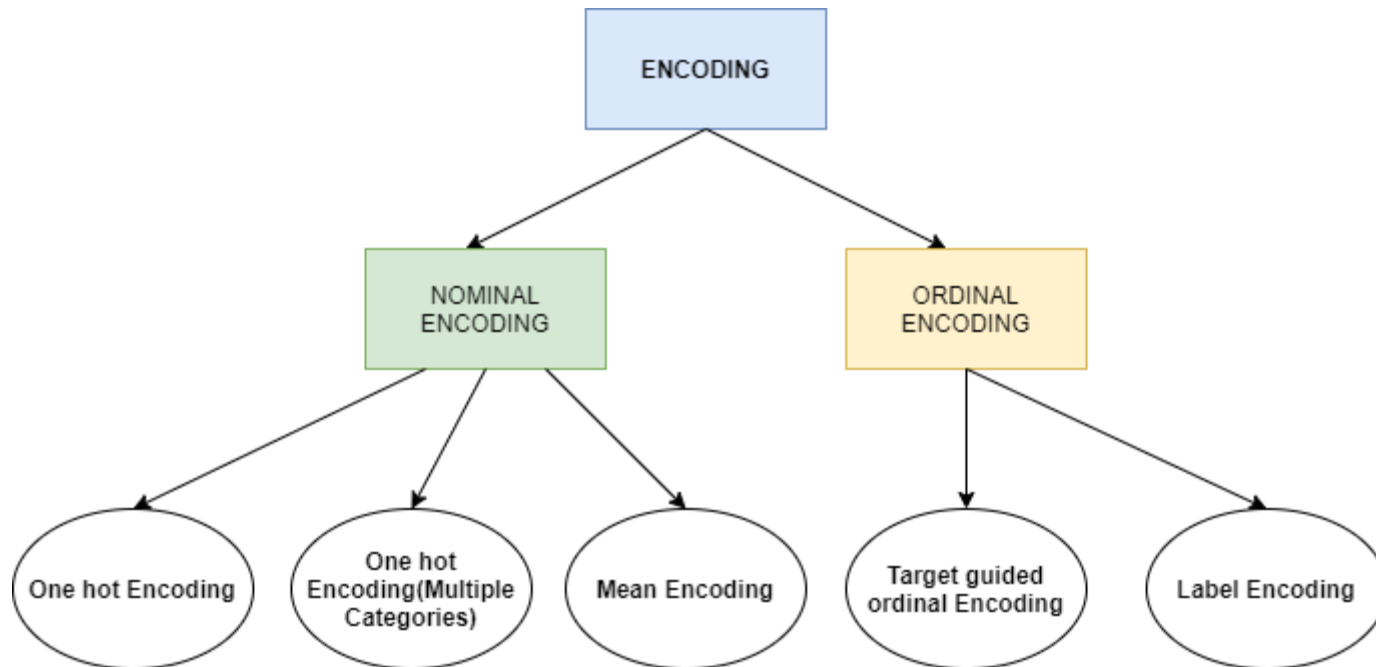
This can be useful in algorithms that do not assume any distribution of the data like K-Nearest Neighbors and Neural Networks.





Gaussian = Normal Distribution







Encoding

Nominal categorical variables are those for which we do not have to worry about the arrangement of the categories.

Male and Female.

Different states like NY, FL, NV, TX

Ordinal categories are those in which we have to worry about the rank. These categories can be rearranged based on ranks.

education level (PHD-1, masters-2, bachelors-3).



ML Process

–

1. Prepare & Transform Data

1. Outlier

2. Feature Scaling

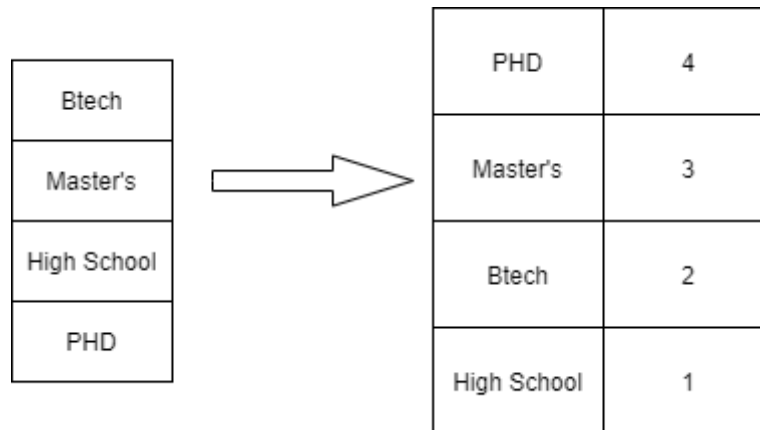
1. Normalization

2. Standardization

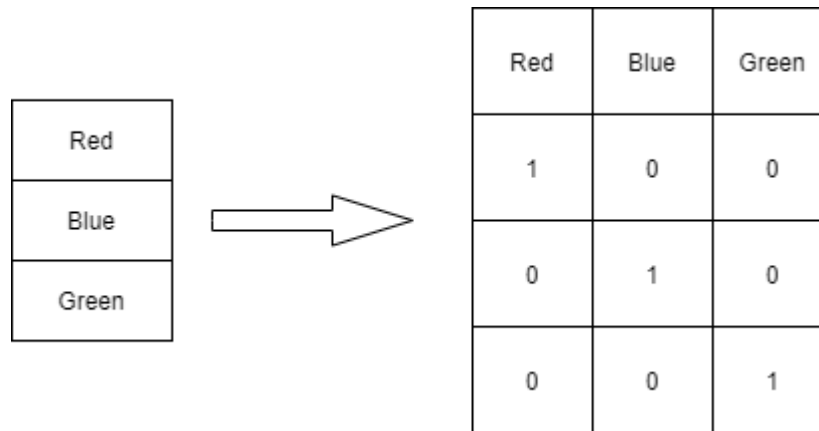
3. Label encoding/ One hot encoding



Label Encoding

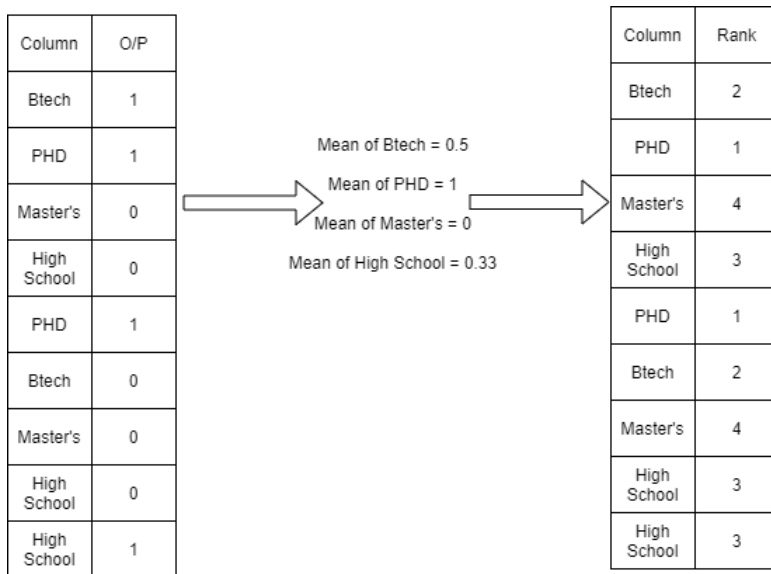


One Hot Encoding

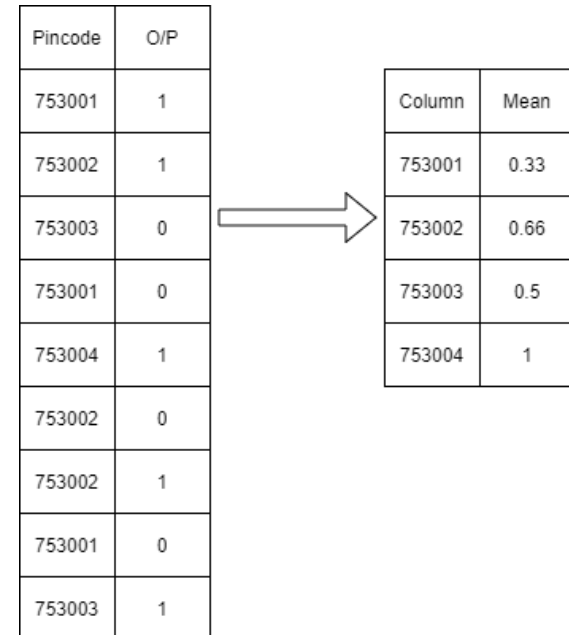


Other Encoding

Target guided ordinal categories

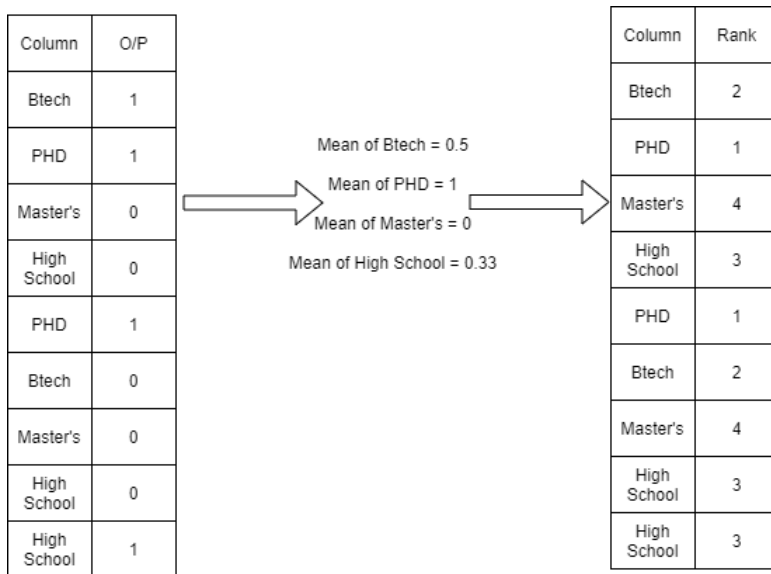


Mean Encoding

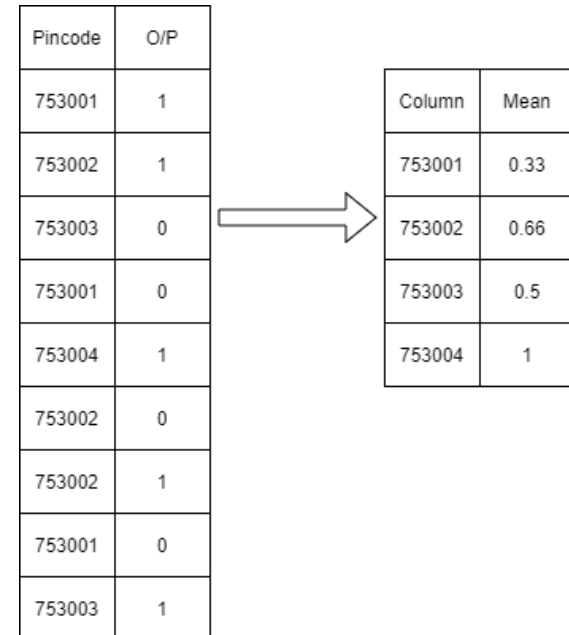


Other Encoding

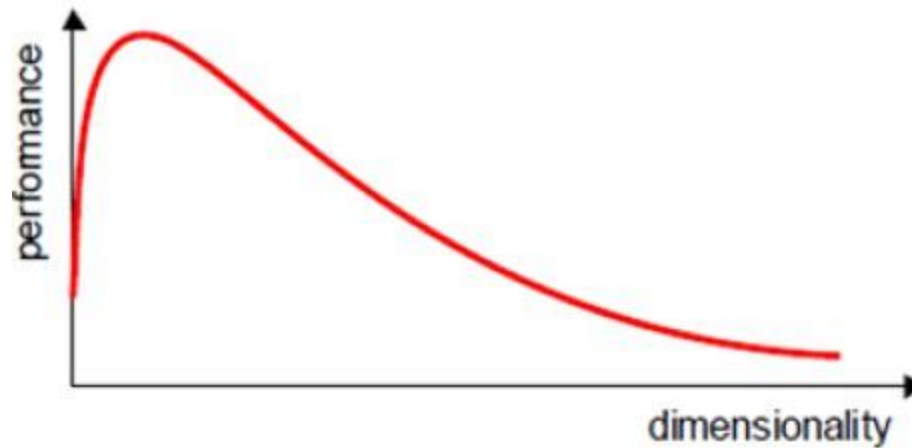
Target guided ordinal categories



Mean Encoding

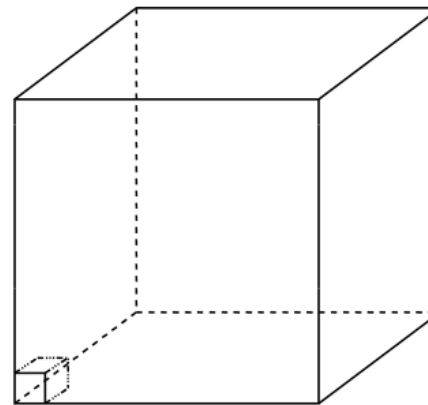
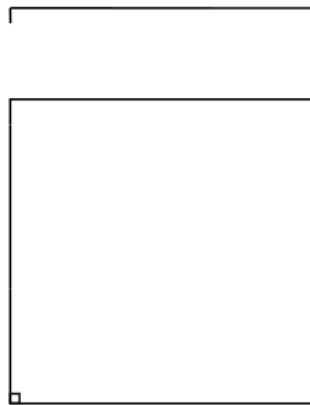


Curse of dimensionality



Curse of dimensionality

- k NN breaks down in high-dimensional space
 - “Neighborhood” becomes very large.
- Assume 5000 points uniformly distributed in the unit hypercube and we want to apply 5-nn. Suppose our query point is at the origin.
 - In 1-dimension, we must go a distance of $5/5000 = 0.001$ on the average to capture 5 nearest neighbors
 - In 2 dimensions, we must go $\sqrt{0.001}$ to get a square that contains 0.001 of the volume.
 - In d dimensions, we must go $(0.001)^{1/d}$



Source: mathematics stack exchange



INVESTIC

Curse of Dimensionality

- Low dimension \rightarrow good performance for nearest neighbor.
- As dataset grows, the nearest neighbors are near and carry similar labels.
- Curse of dimensionality: in high dimensions, almost all points are far away from each other.

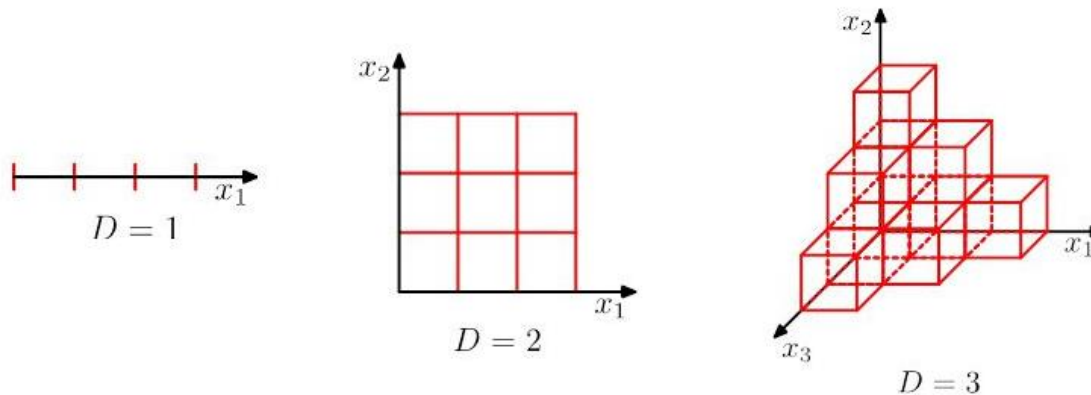


Figure Bishop 1.21