

Classification Models

Ajanda

01 Classification Methods

02 Logistic Regression

03 Decision Trees

04 KNN

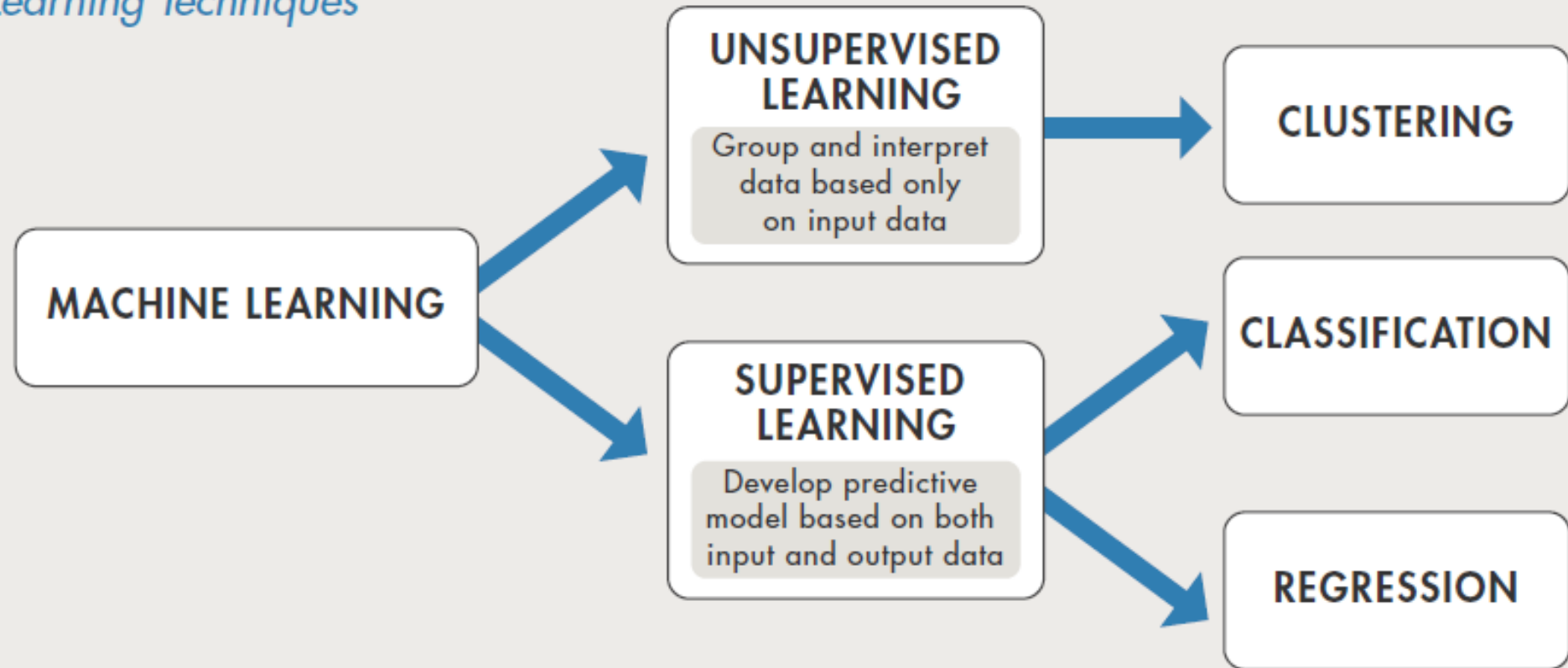
05 Performance Metrics

06 Practice on Mock Data

01

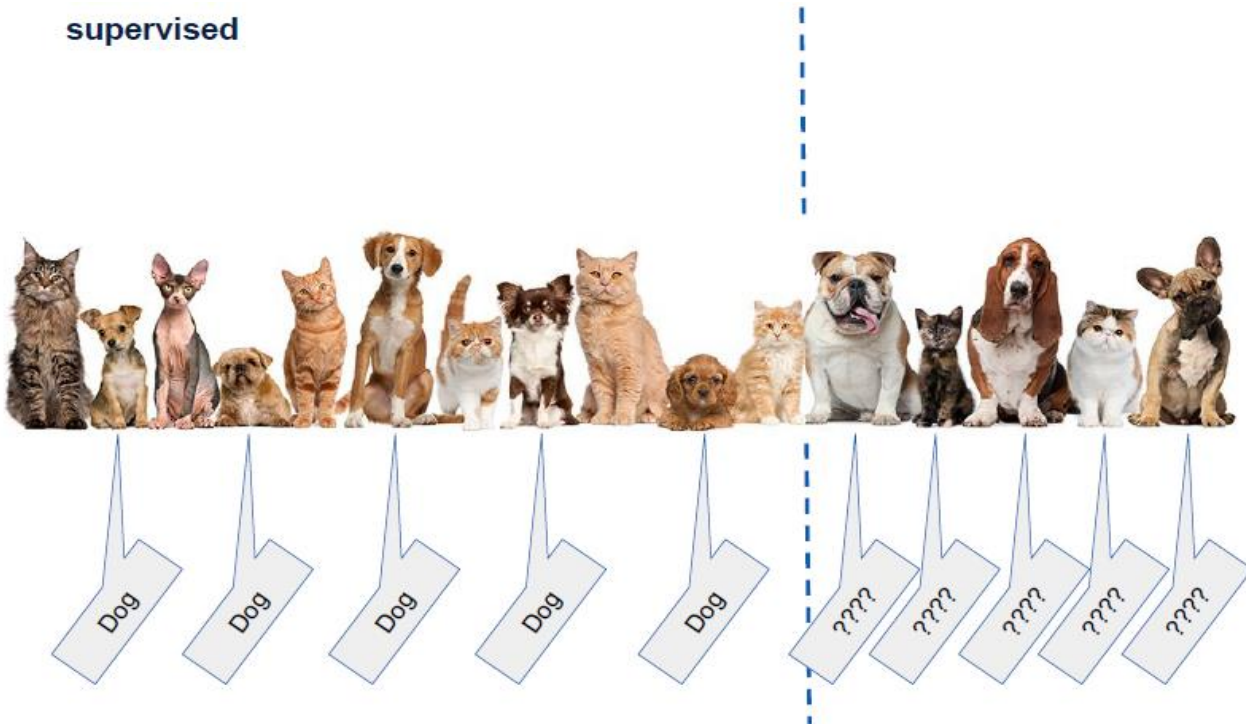
Classification Methods

Machine Learning Techniques



Supervised Learning

supervised



Training Dataset

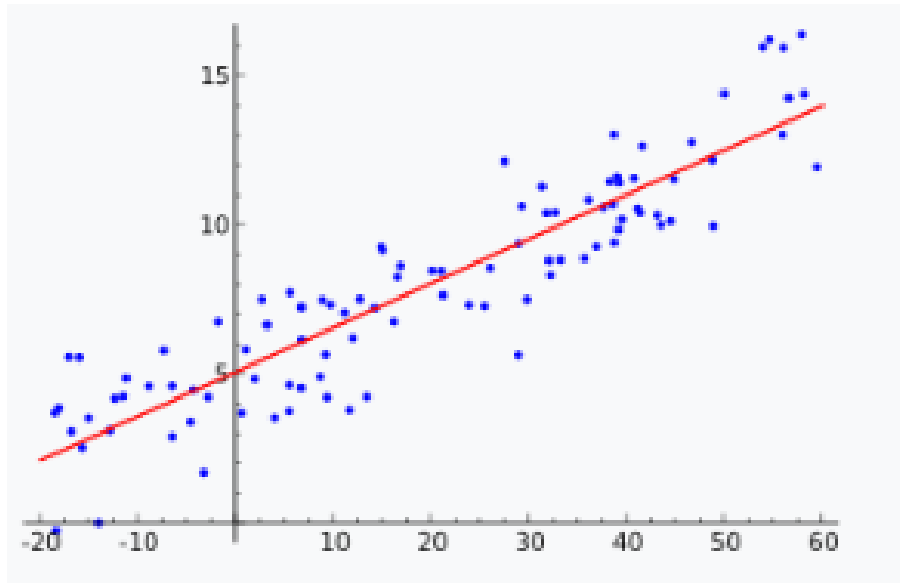
n training samples
or instances

		Predictors (Features)				Labels (output Y)
		x_1	x_2	x_p	Y
Samples	1					
	n					
		p Predictors				

02

Logistic Regression

What We Learned? -- Linear Regression



$$y = \beta_0 + \beta_1 x$$

where

y : model estimation

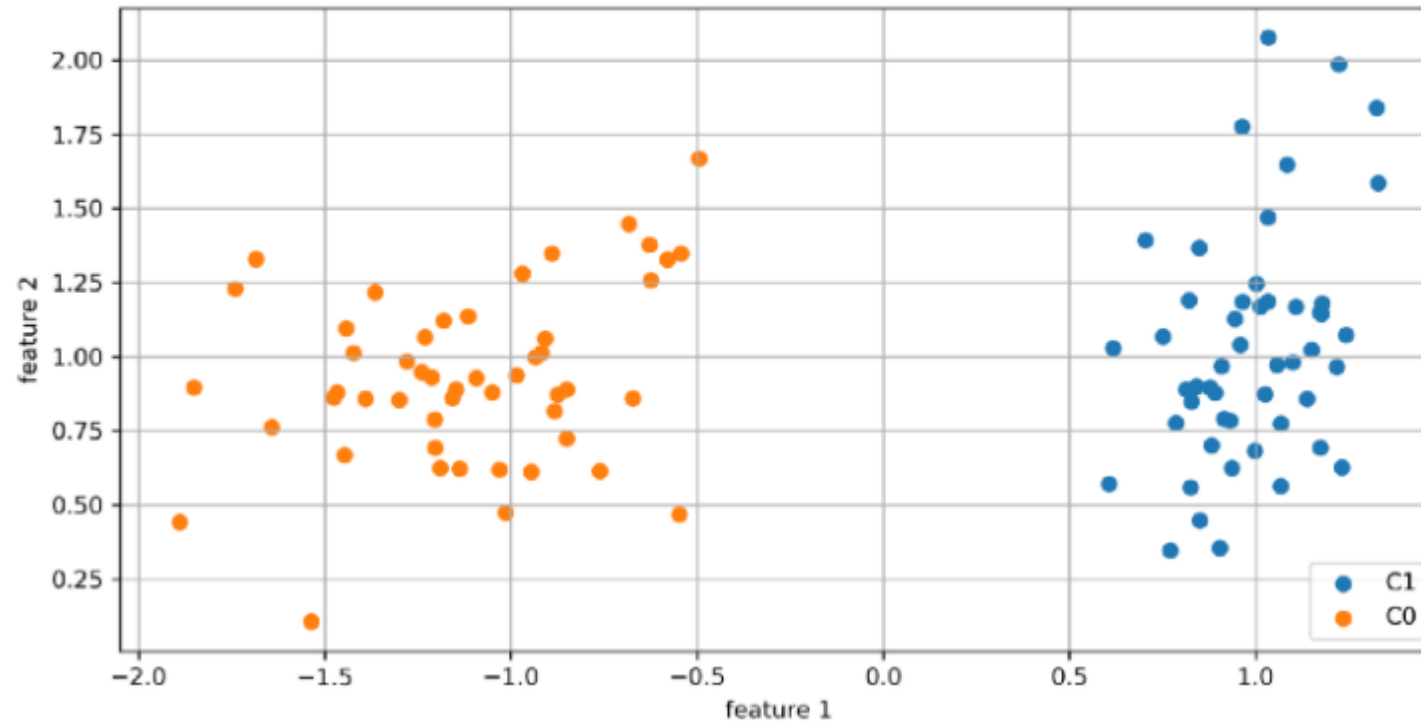
x : the input variable or predictor

B_0 : intercept ($x=0$)

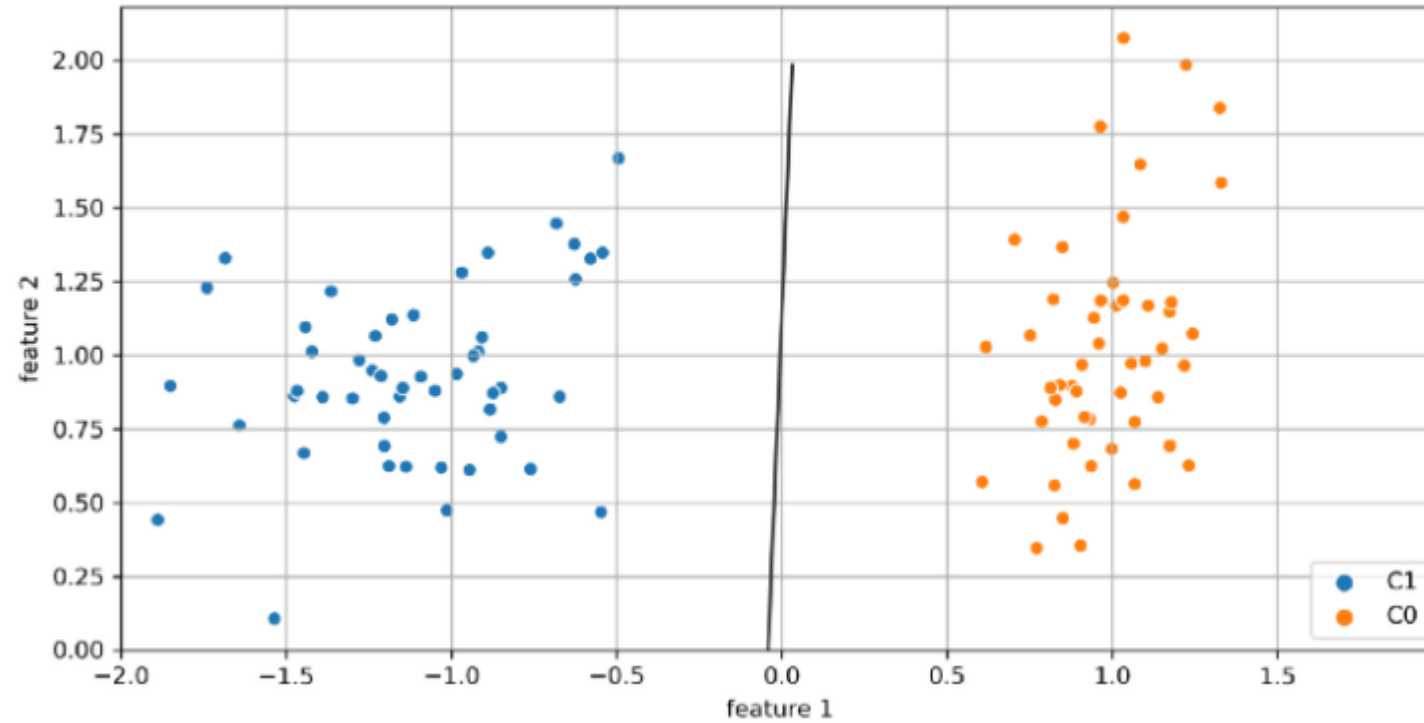
B_1 : slope (coef of feature)

$$J(\beta_0, \beta_1) = \frac{1}{2m} \sum_{i=1}^m [y_i - (\beta_0 + \beta_1 x_i)]^2$$

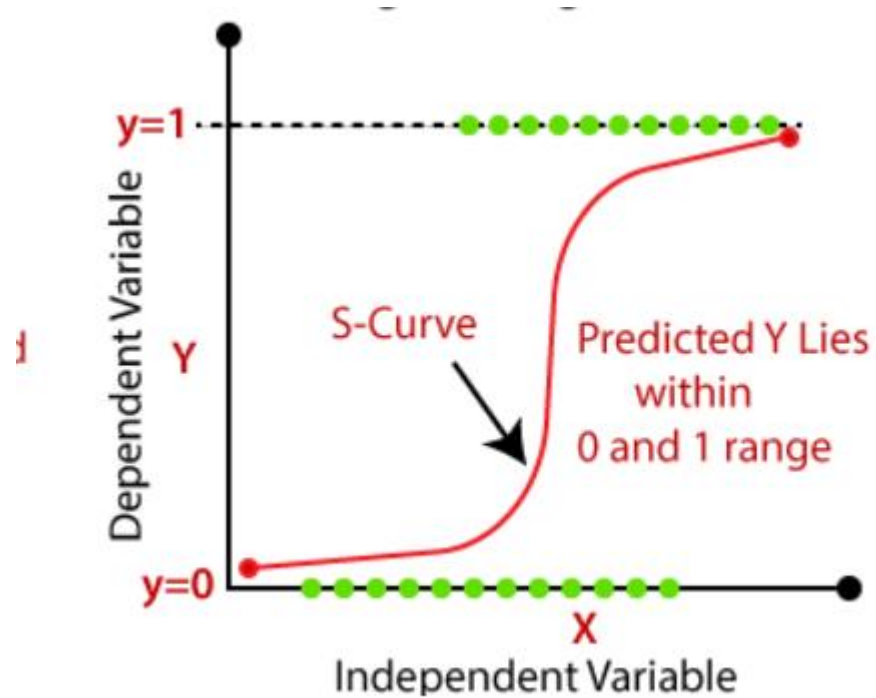
Logistic Regression with 2 Feature



Decision Boundary



Sigmoid Function



- Binary & Categorical target
- A form of linear regression
- Sigmoid Function

$$y = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

Assumptions & Advantages

Assumptions

- Linear relationship between features and target
- Minimum or no relationship between features
- Relationship between one feature and target is independent of other feature and target
- No extreme Outliers
- Sample size is large enough

Advantages

- Results interpretability
- Ease of use
- Low computational cost

Overfitting & Regularization

- Lasso (L1) Regularization

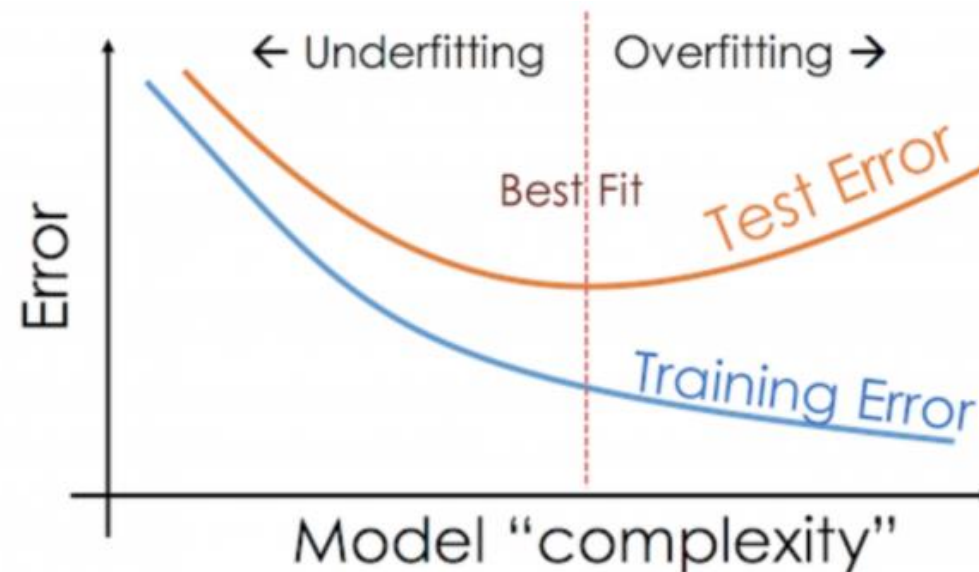
$$\sum_{i=1}^n (Y_i - \sum_{j=1}^p X_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

Cost function

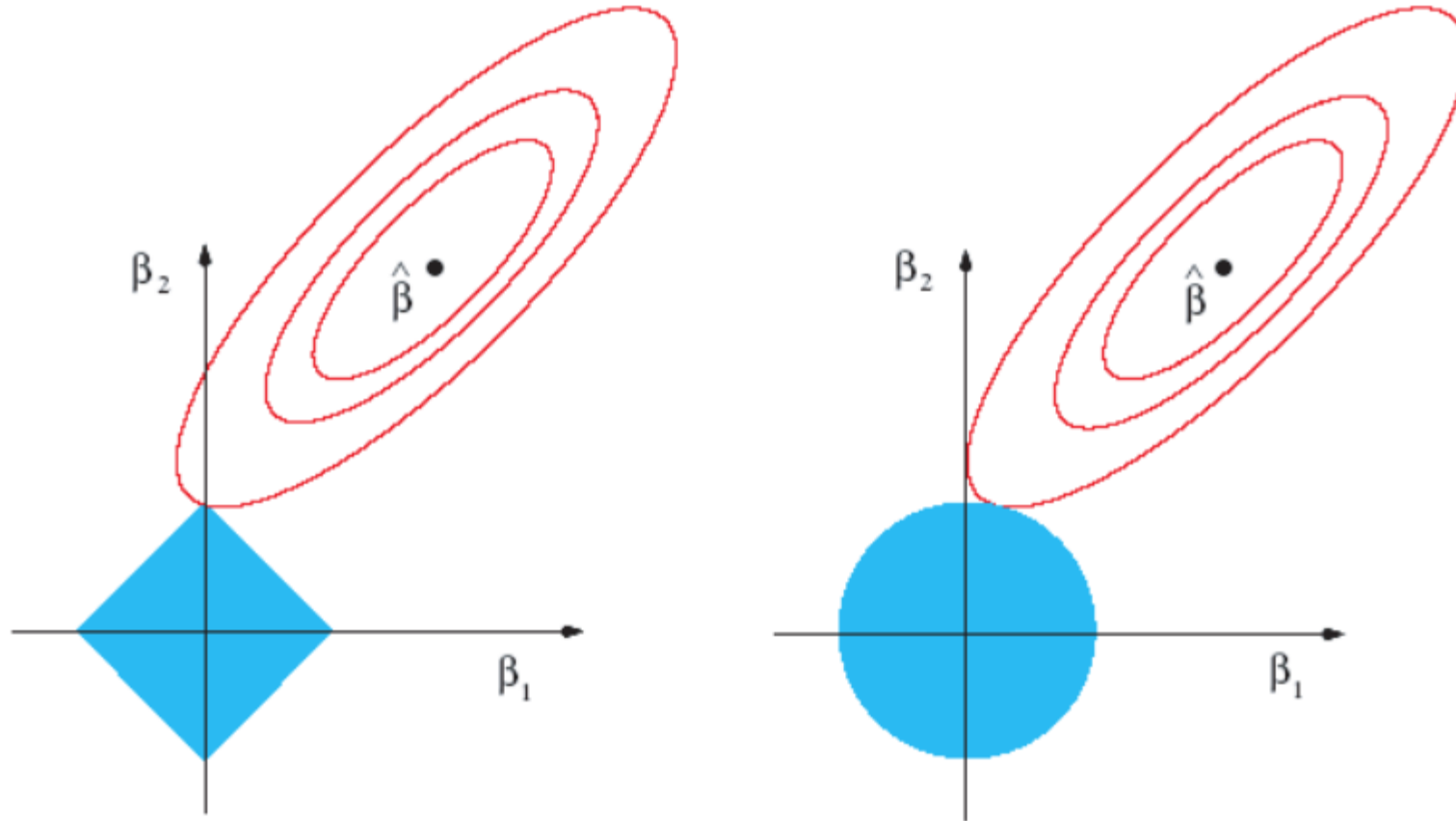
- L2: Ridge (L2) Regularization

$$\sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

Cost function

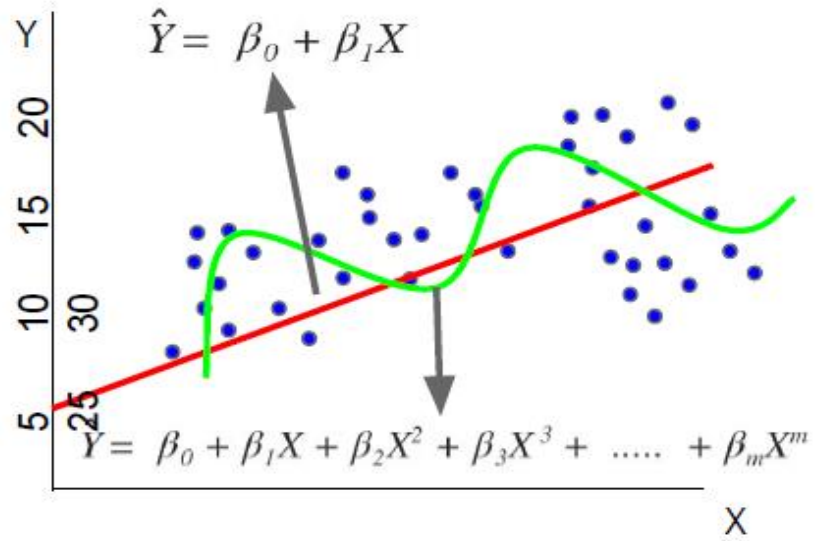


Lasso vs Ridge Regularization

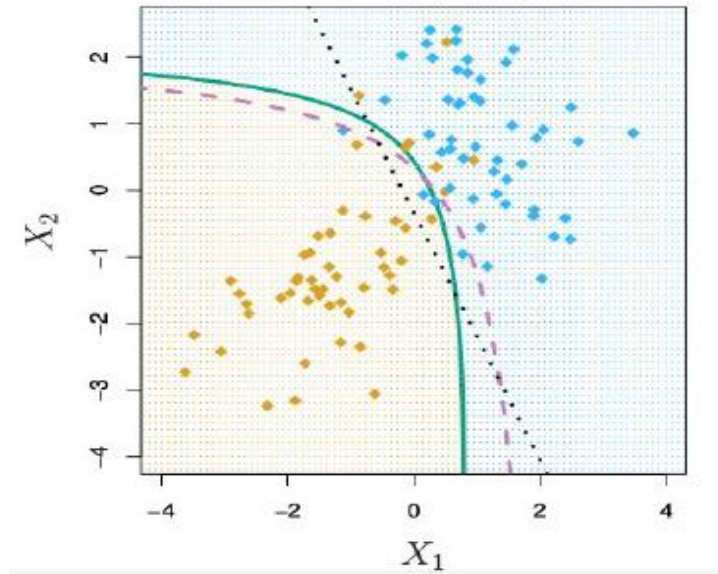


Summary

Regression



Classification

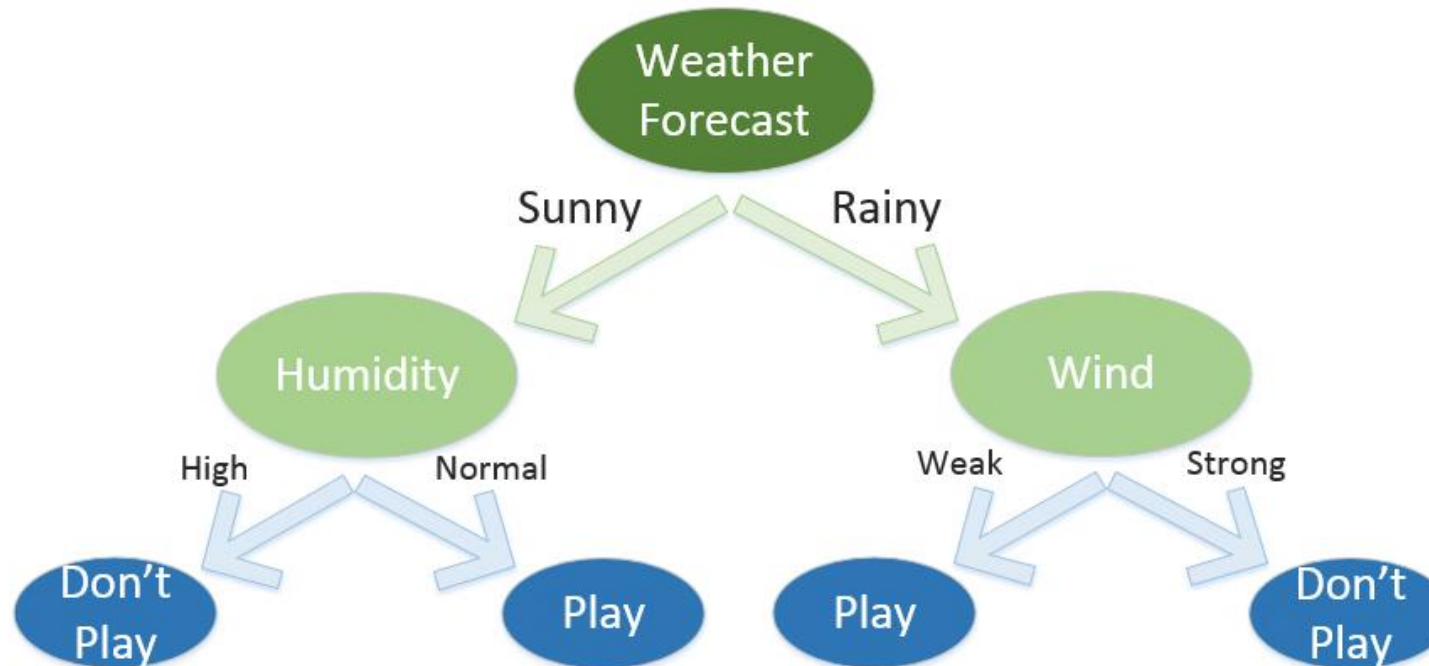


03

Decision Trees

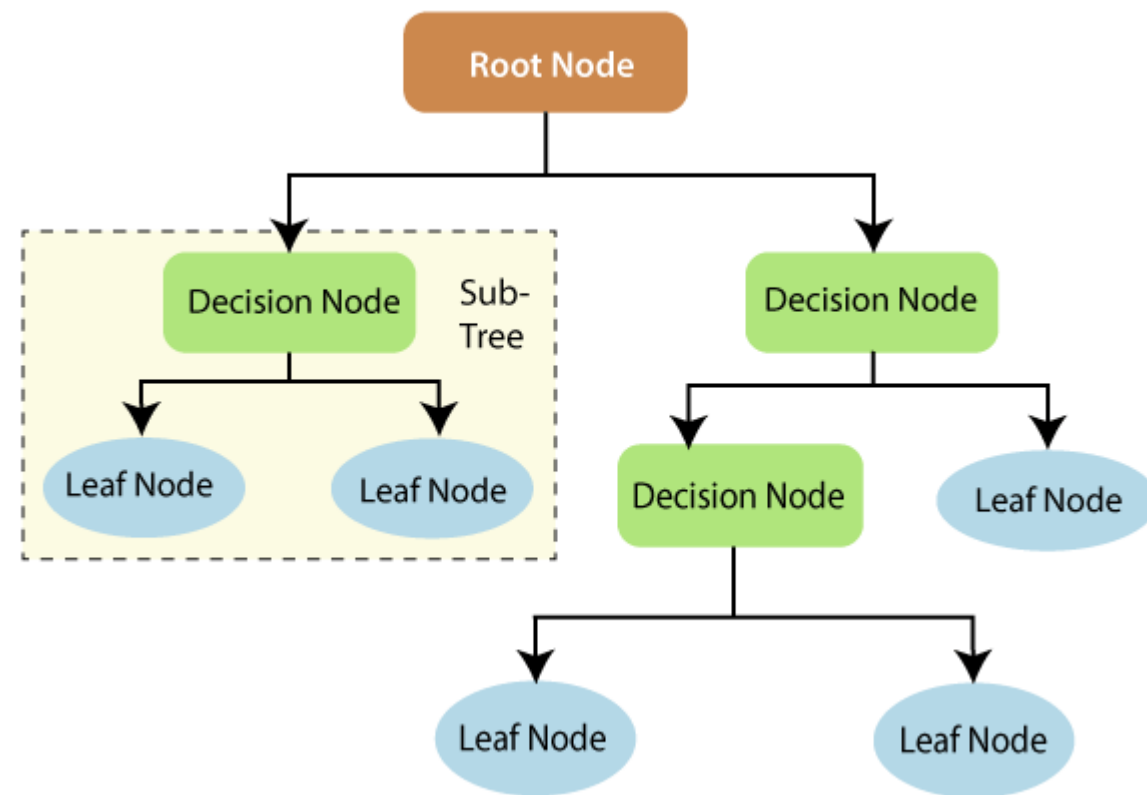
Decision Trees

- One of the most widely used and practical methods for inductive inference
- **CARTs** can be used for classification or regression problems (mostly classification)
- Weaknesses are relatively minor and can be largely avoided.



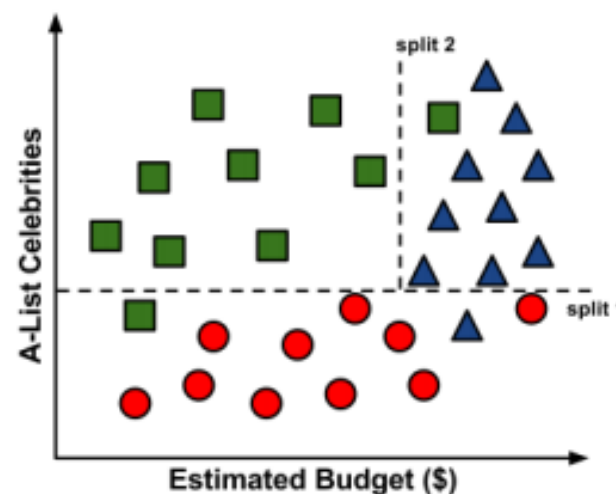
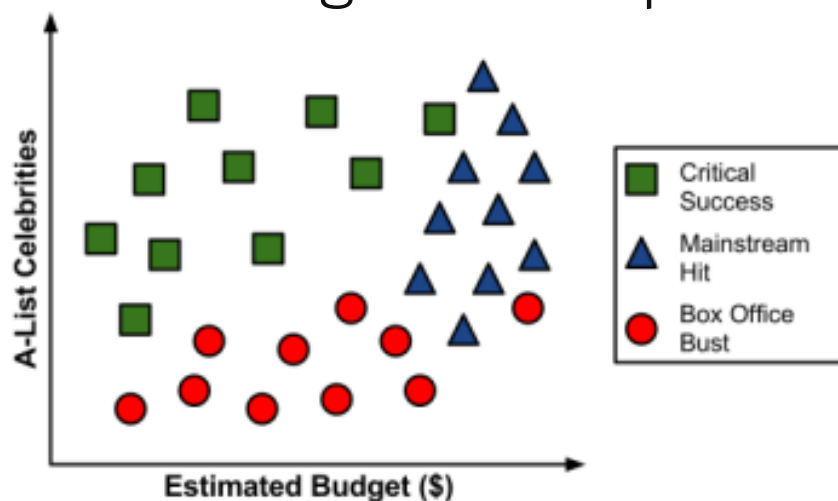
Structure of a Decision Tree

- A flow-chart-like tree structure
- **Decision node** denotes a test on an attribute
- **Branch** represents an outcome of the test
- **Leaf nodes** represent class labels or class distribution



Divide and Conquer

- Decision trees are built using a heuristic called recursive partitioning
- It splits the data into smaller and smaller subsets of similar classes
- Divide and conquer the nodes until a stopping criterion is reached
 - All (or nearly all) of the examples at the node have the same class
 - There are no remaining features to distinguish among examples
 - The tree has grown to a predefined size limit



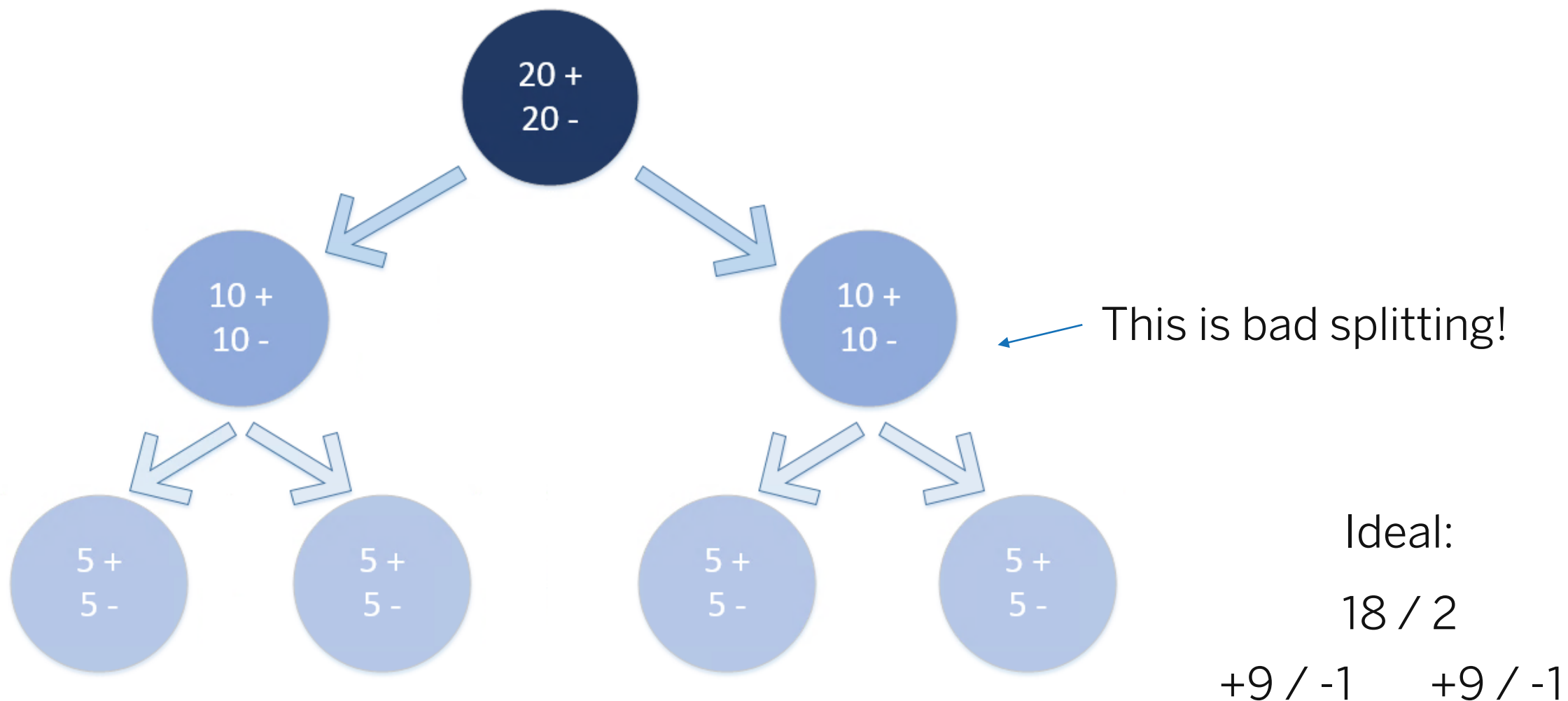
Learning the Tree

■ Many algorithms: ID3, J48, C4.5, C5.0

■ Algorithmic issues:

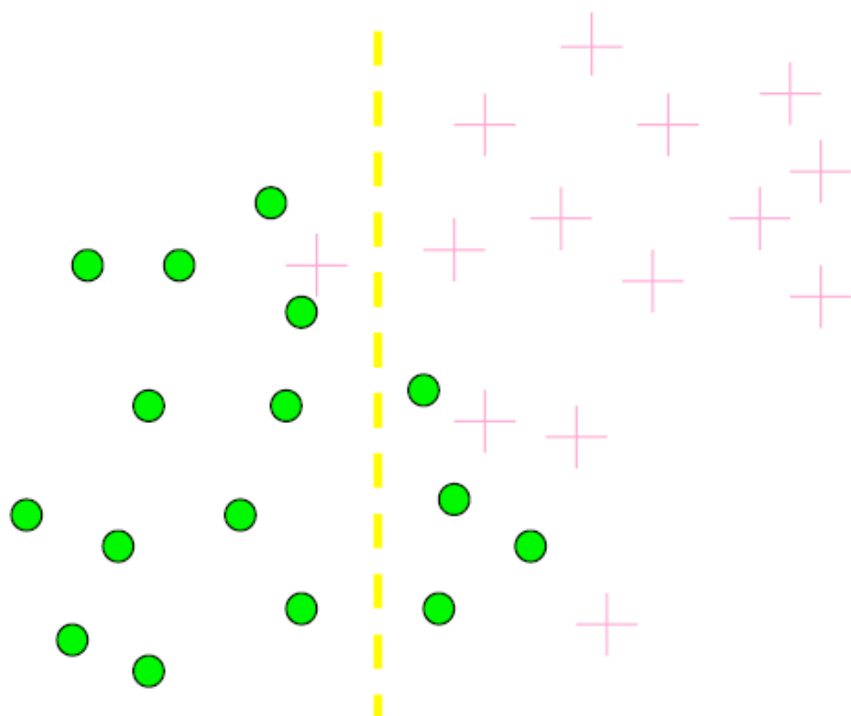
- **Determine how to split the records**
 - How to specify the attribute test condition?
 - How to determine the best split?
- **Determine when to stop splitting**
- **How to classify a leaf node**
 - Assign the majority class
 - If leaf is empty, assign the class with the most popular label

Which Feature is Best for Splitting?



Determining the Best Split

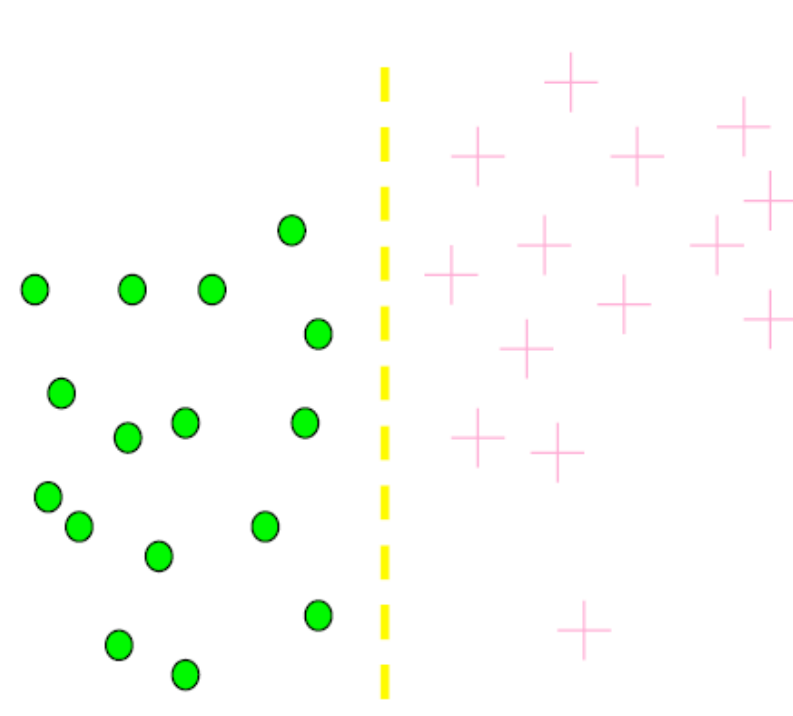
**Split over whether
Balance exceeds 50K**



Less or equal 50K

Over 50K

**Split over whether
applicant is employed**



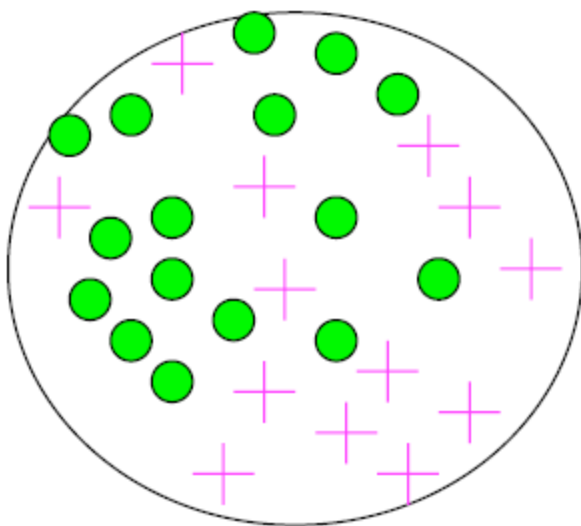
Unemployed

Employed

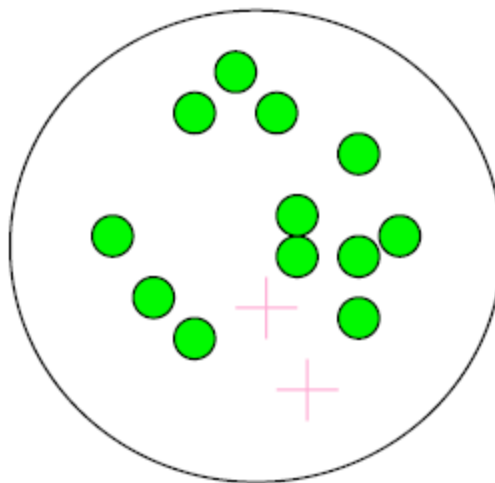
Impurity

We need a measure of impurity for class label distribution on a node

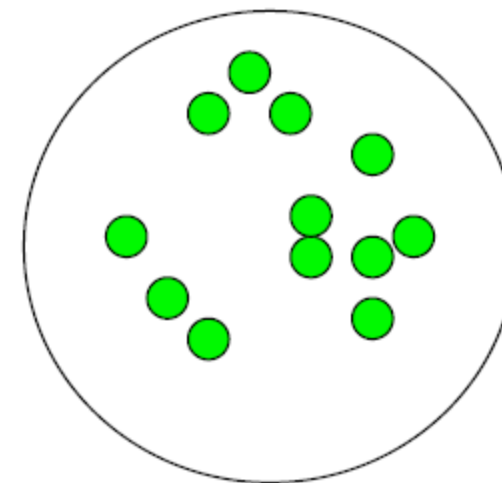
Very impure group



Less impure

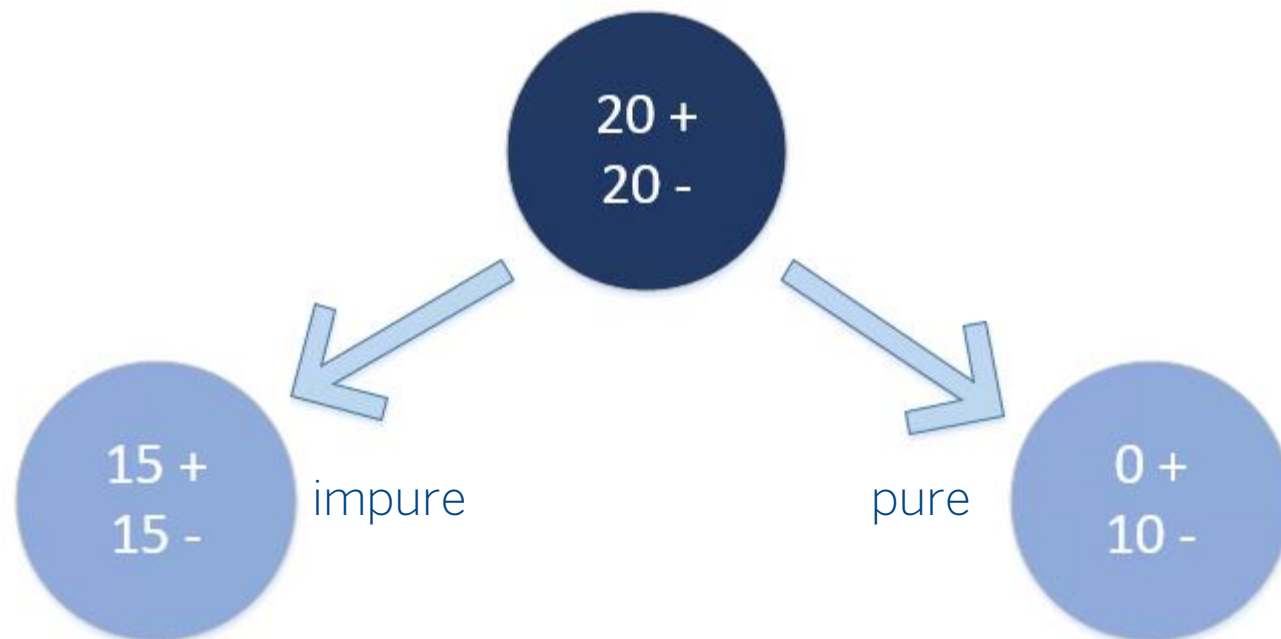


Minimum impurity



Common Measures of Impurity

- Entropy
- Gini Index
- Misclassification Error



If the distribution is **less uniform, the node is purer.**

Information & Entropy

Information : Amount of uncertainty (amount of surprise) in the outcome

- Flipping coin
- Rolling a dice

$$I(X) = \log_2 \frac{1}{p(x)} = -\log_2 p(x)$$

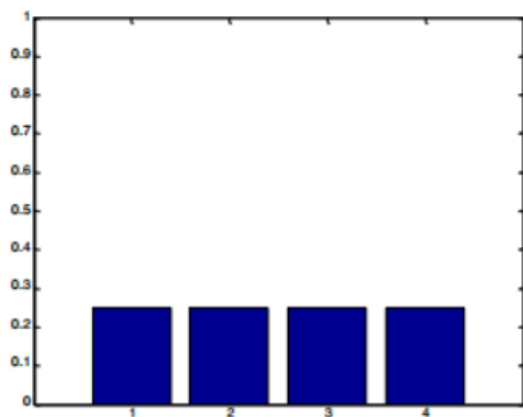
■ If the probability of an event is small and it happens, the information is large

- Flipping «head» on a coin $> I = -\log_2 1/2 = 1$
- Rolling «6» on a dice $> I = -\log_2 1/6 = 2.58$

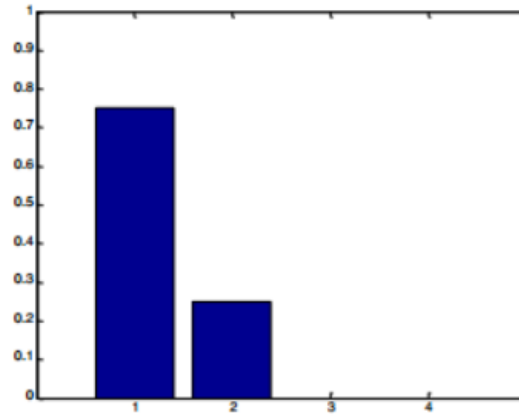
Information & Entropy

Entropy : The expected amount of information when observing the output of a random variable X :

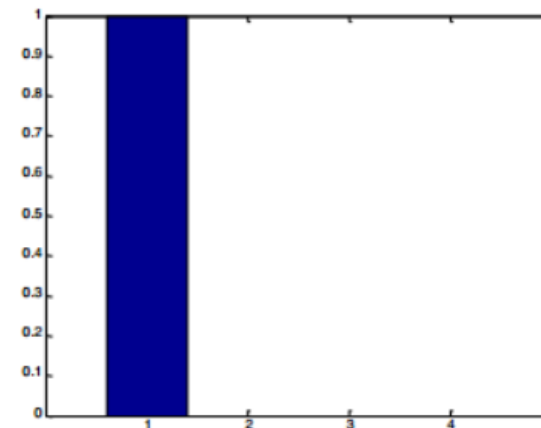
$$H(X) = E(I(X)) = \sum_i p(x_i) I(x_i) = - \sum_i p(x_i) \log_2 p(x_i)$$



$$\begin{aligned} H(x) &= .25 \log 4 + .25 \log 4 + \\ &\quad .25 \log 4 + .25 \log 4 \\ &= \log 4 = 2 \text{ bits} \end{aligned}$$



$$\begin{aligned} H(x) &= .75 \log 4/3 + .25 \log 4 \\ &\approx .8133 \text{ bits} \end{aligned}$$



$$\begin{aligned} H(x) &= 1 \log 1 \\ &= 0 \text{ bits} \end{aligned}$$

Information Gain

Information Gain : Reduction in uncertainty by knowing Y

$$IG(X,Y)=H(Y) - H(Y|X)$$

- Information gain tells us the importance of the given feature
- Determines the most useful one for discriminating between the classes to be learned.

Calculation of Information Gain



Weighted Average Entropy of Children

$$= \left(\frac{10}{18} \times 0.881\right) + \left(\frac{8}{18} \times 0.544\right) = 0.731$$

Information Gain :
0.991 - 0.731 = 0.26

Other Impurity Measures

Gini :

$$GINI(t) = 1 - \sum_i p(i | t)^2 \quad p(i | t) \text{ is the relative frequency of class } i \text{ at node } t$$

Misclassification Error :

$$\text{Classification error}(t) = 1 - \max_i p(i | t)$$

+	2
-	4

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Gini} = 1 - (2/6)^2 - (4/6)^2 = 0.444$$

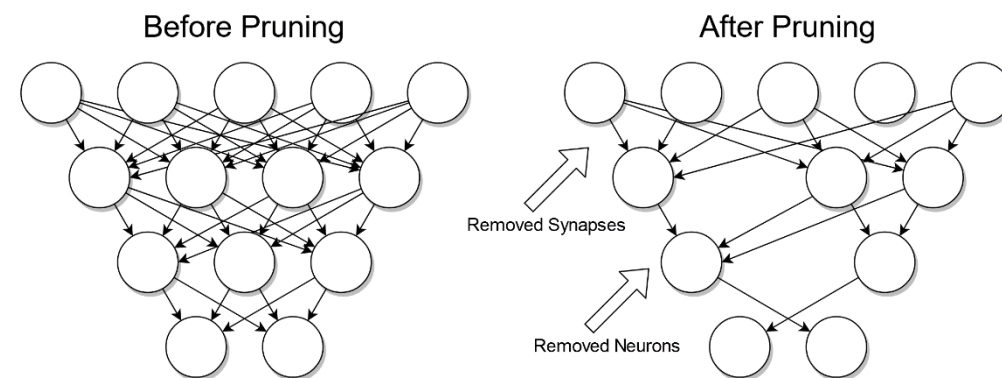
$$\text{Entropy} = - (2/6) \log_2 (2/6) - (4/6) \log_2 (4/6) = 0.92$$

$$\text{Error} = 1 - \max (2/6, 4/6) = 1 - 4/6 = 1/3$$

Generalization of Decision Trees

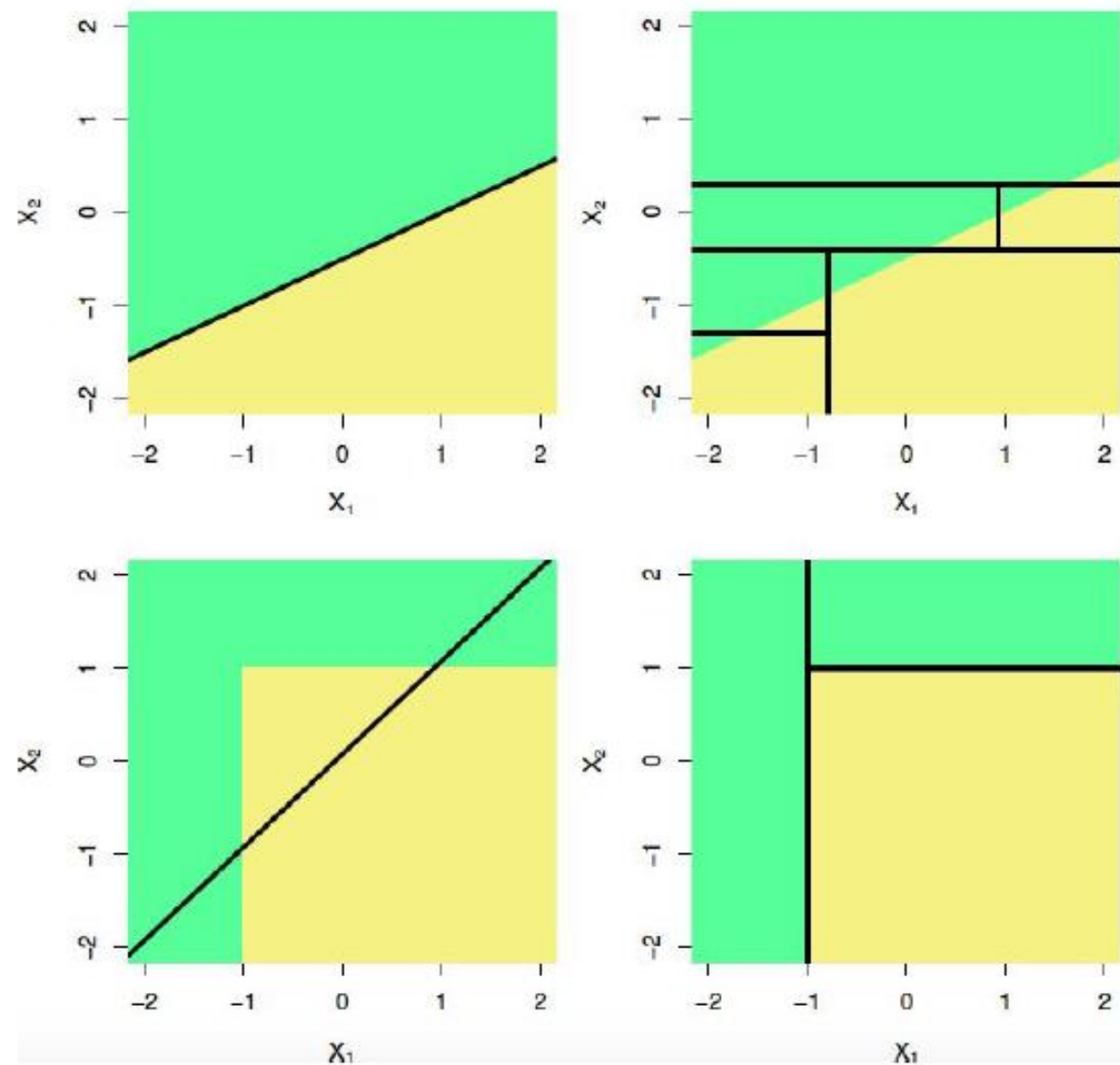
Pruning: Replacing a subtree with a leaf with the most common classification in the subtree.

- Pre Pruning (Early Stopping Rule) : Stop growing the tree before it perfectly classifies the training data.
- Post Pruning : Grow full tree, then prune.



Decision Boundaries (Linear vs Trees)

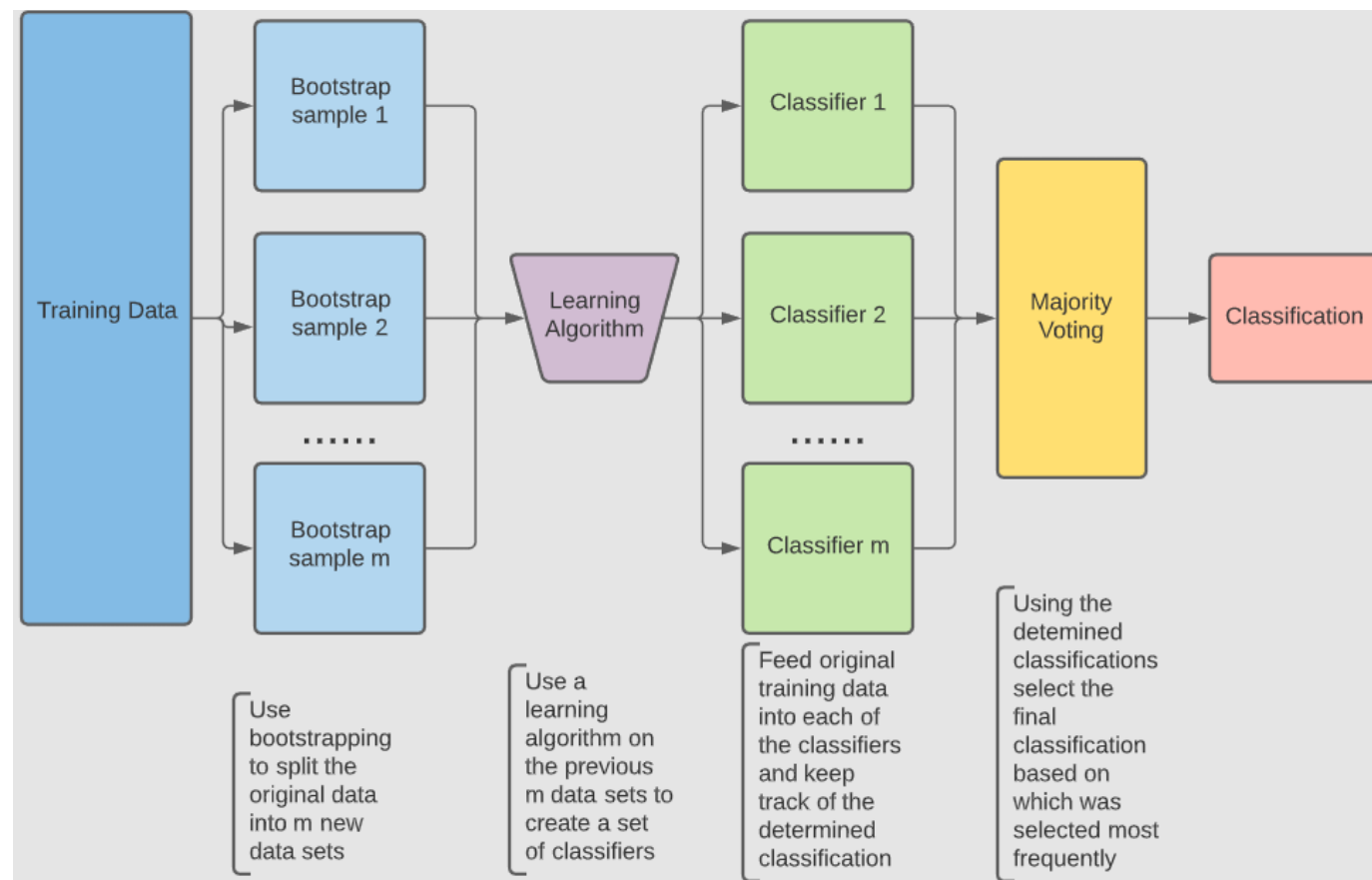
- Decision boundary is parallel to axes because test condition involves a single feature at a time
- Tree models better reflect nonlinear relationships between feature and target.
- if the relationship is expected to be linear, linear models might be better.



Bagging (Bootstrap Aggregation)

■ Random Forest

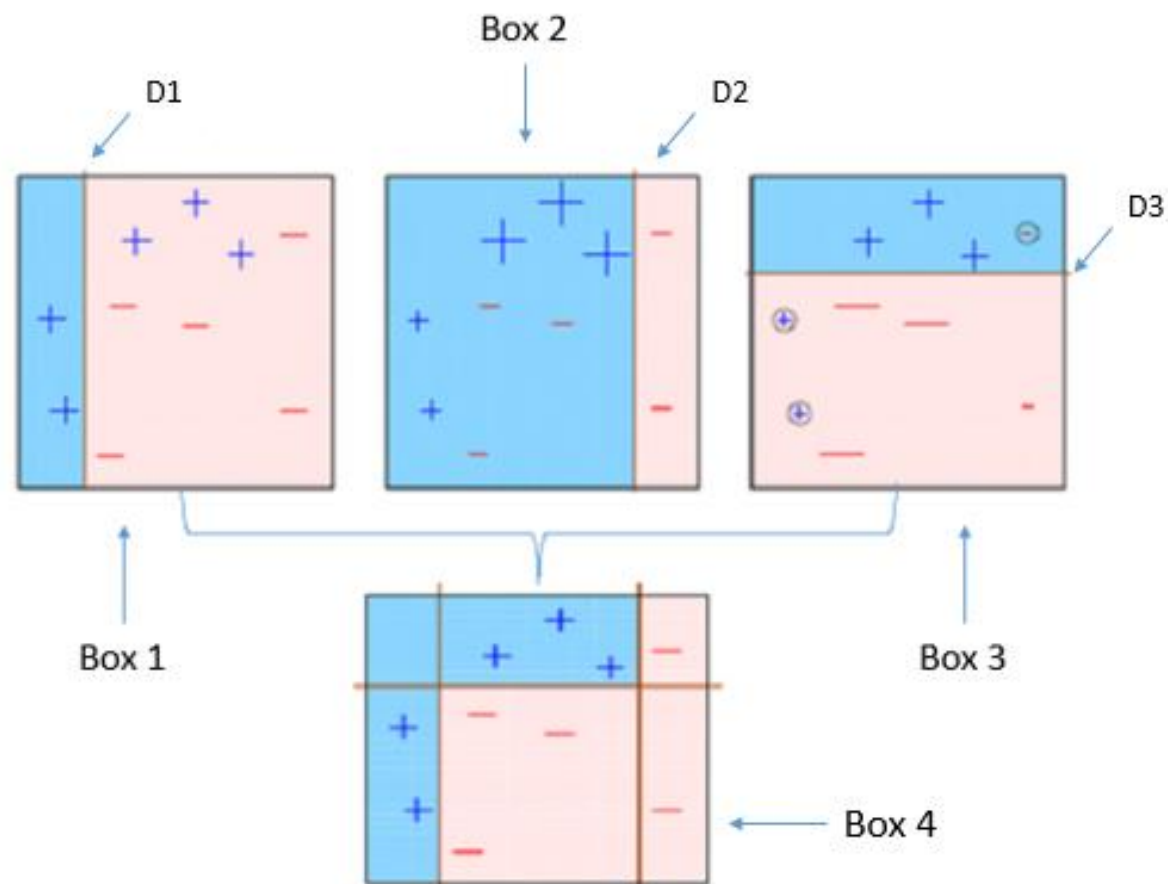
General procedure that can be used to **reduce the variance** for the decision tree algorithm that have high variance.



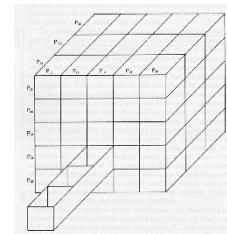
Boosting

■ Adaptive Boosting (AdaBoost)

General ensemble method that creates a strong classifier from a number of **weak classifiers**.



Category Theory



Morphological Analysis

Homologic Analysis

Boosting

Reinforcement Learning

Deep Learning

Bagging

Computability Theory

Topology



Clustering Topology

Multi Class Classification Topology

Binary Classification Topology

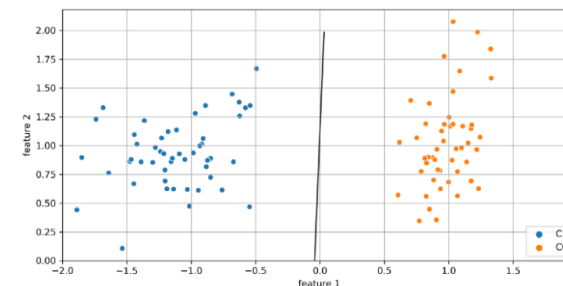
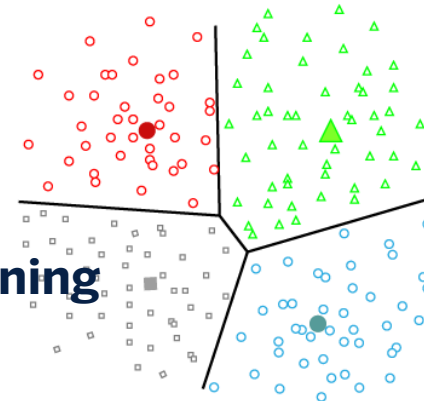
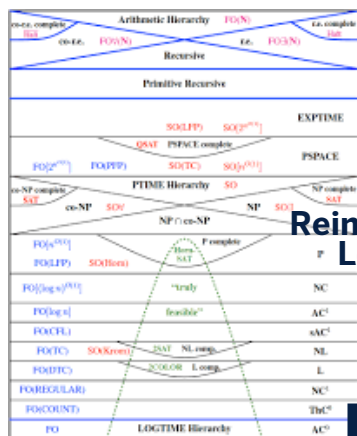
Entropy D

Entropy C

Entropy B

Entropy A

Problem

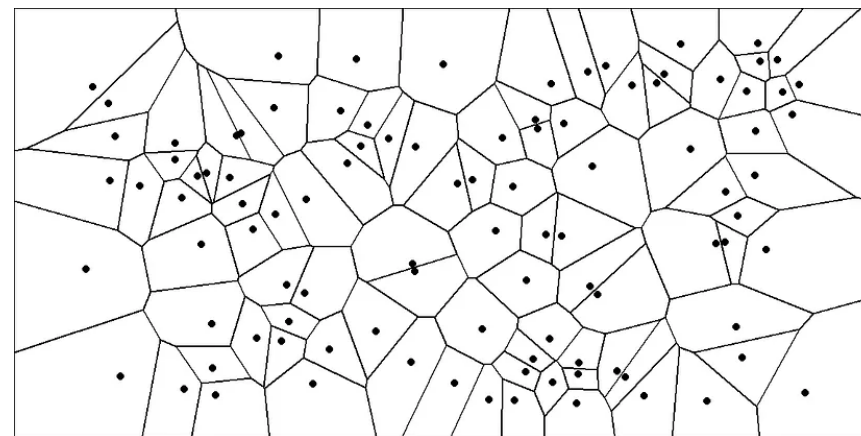
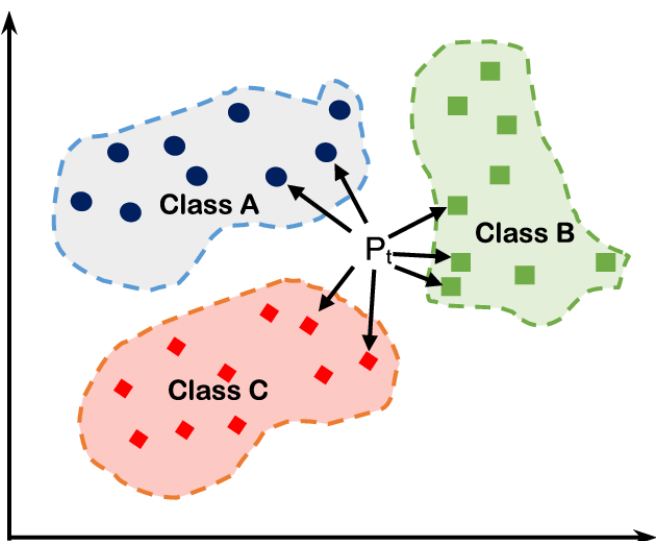


04

K Nearest Neighbors

k Nearest Neighbours

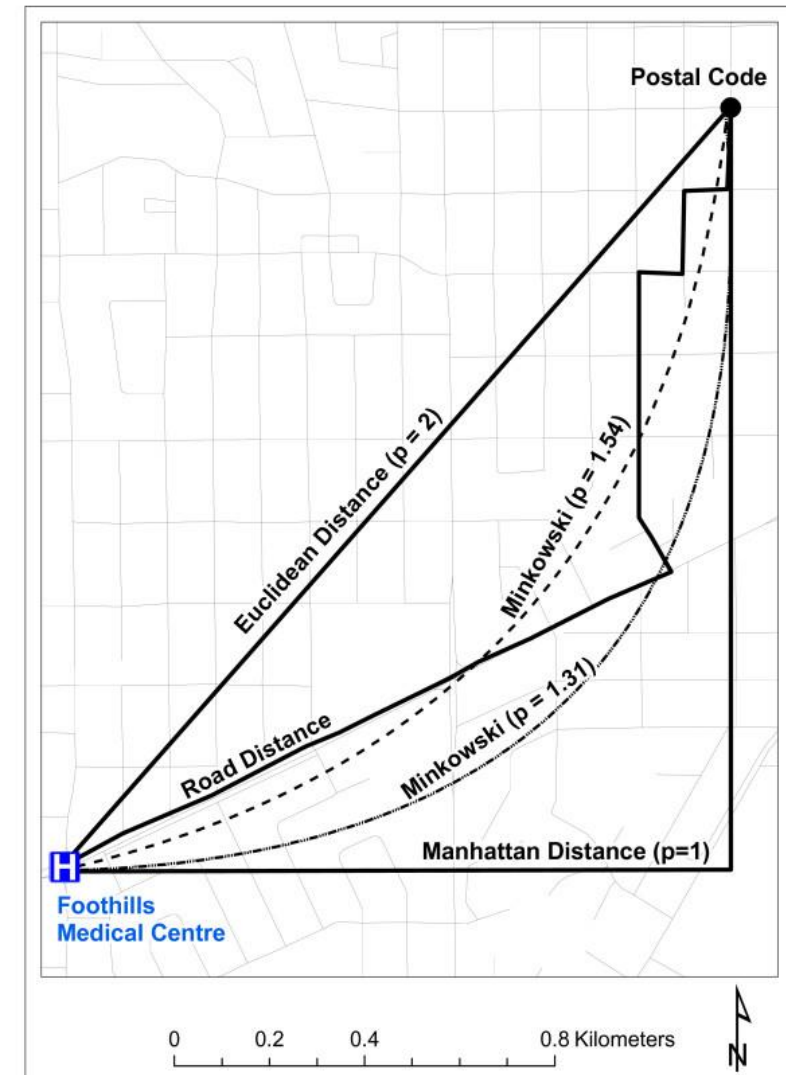
- One of the easy to implement non-probabilistic supervised learning algorithms.
- **KNN** can be used for classification or regression problems (mostly classification)
- Algorithm assumes that similar things exist in close proximity. It categorizes objects based on the classes of their nearest neighbours in the dataset.



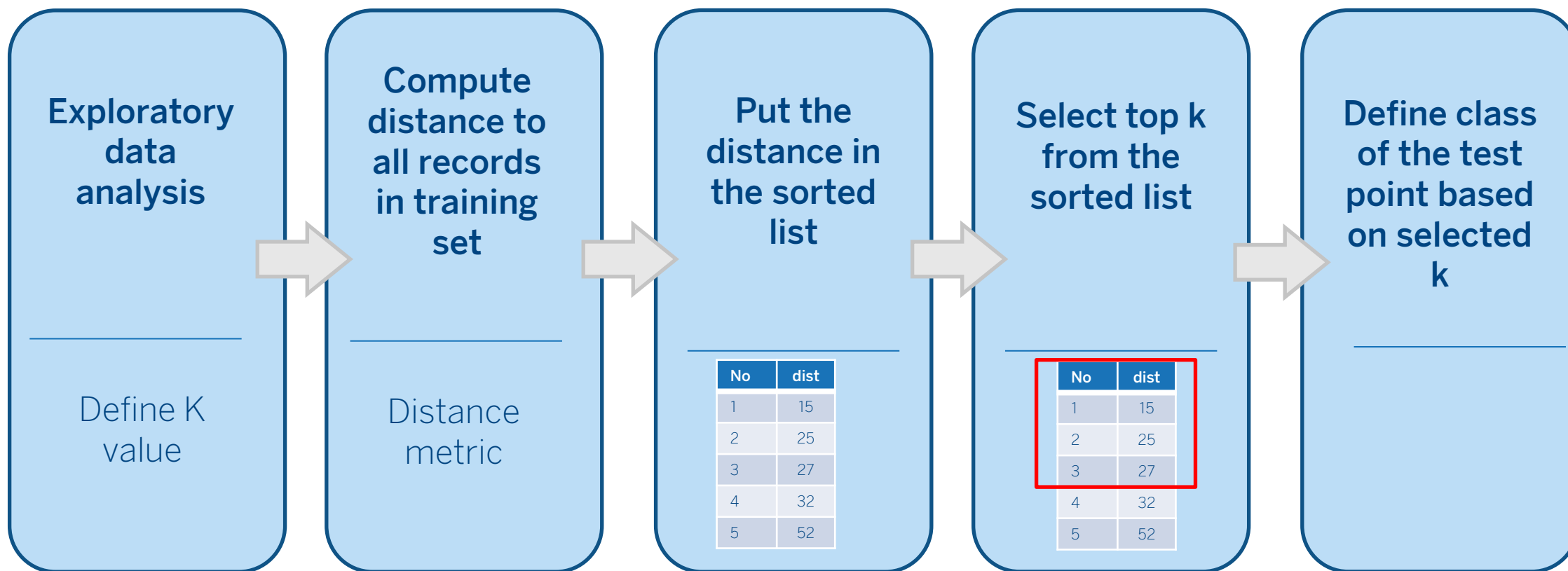
1-NN Pattern : Voronoi Diagrams

Measure of Distance in Data Mining

- Euclidean $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$
- Manhattan $|x_1 - x_2| + |y_1 - y_2|$
- Minkowski $d(p, q) = \left(\sum_{i=1}^n |p_i - q_i|^c \right)^{\frac{1}{c}}$

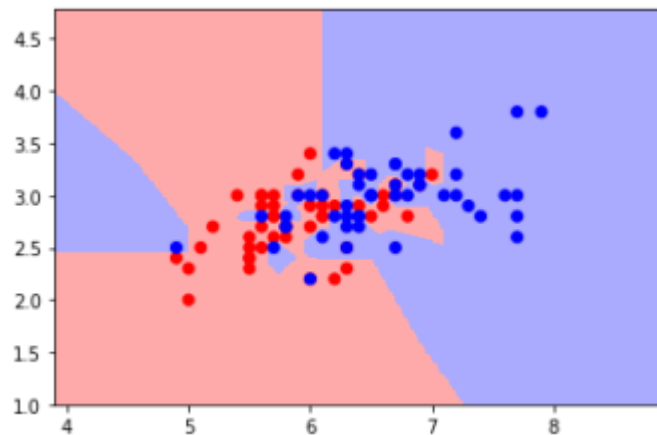


k Nearest Neighbours

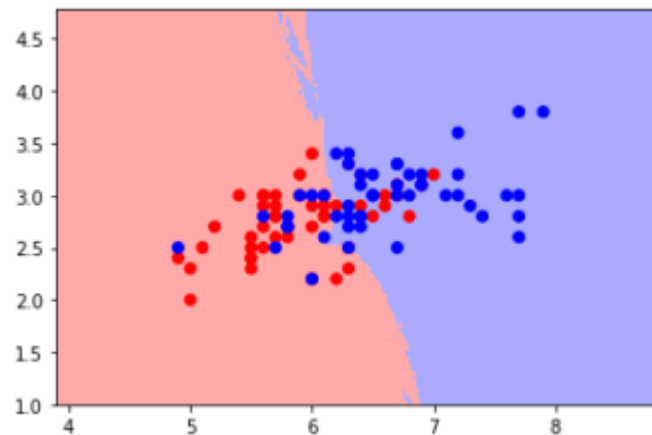


Choosing the k Value : Variance & Bias

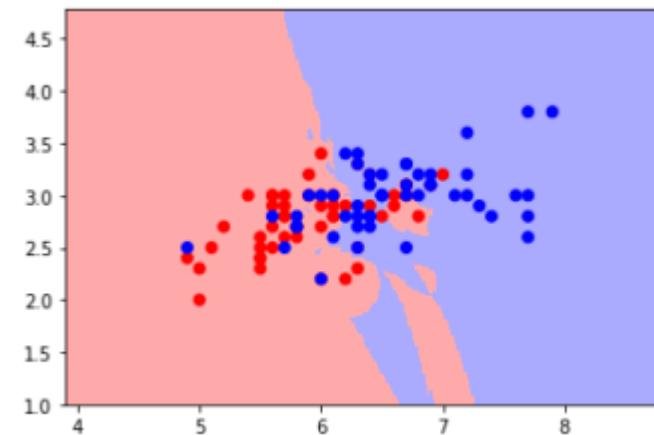
1



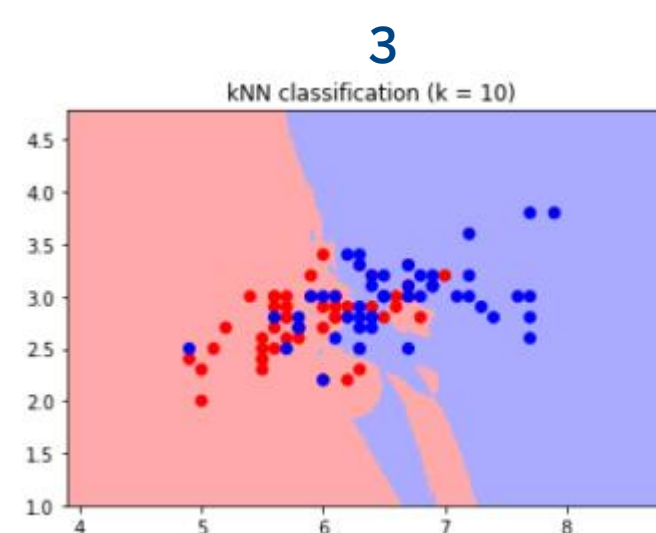
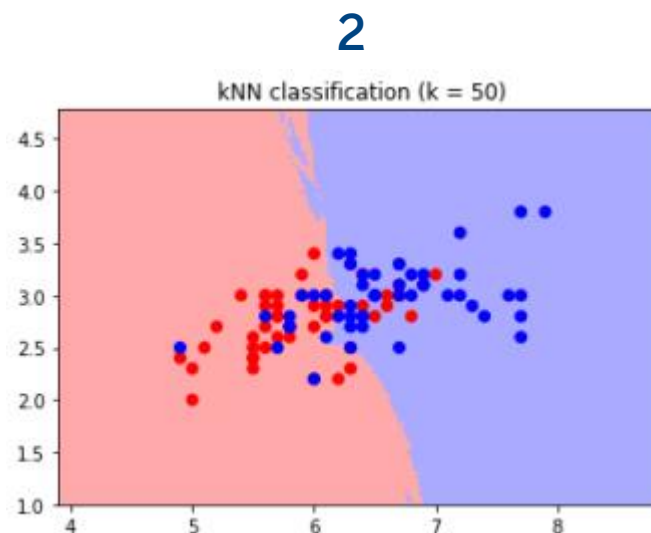
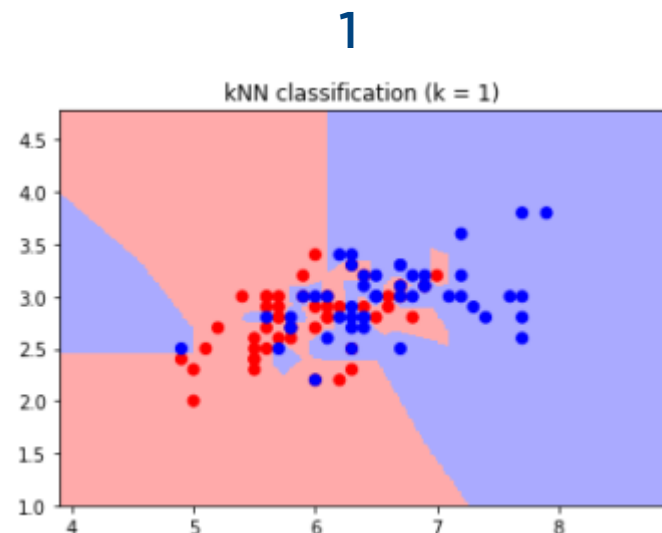
2



3



Choosing the k Value : Variance & Bias

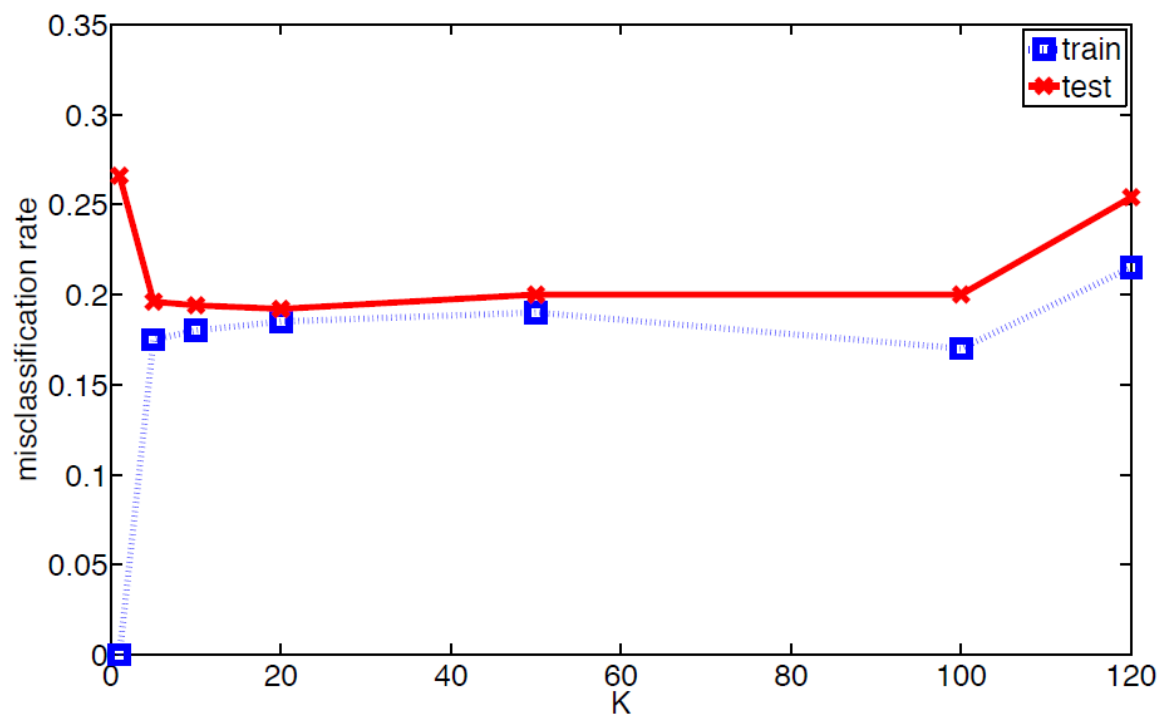


$$2 > 3 > 1$$

- **Small k values** cause **high variance** and an unstable output
- **Large k values** lead to **high bias** by classifying everything as the more probable class

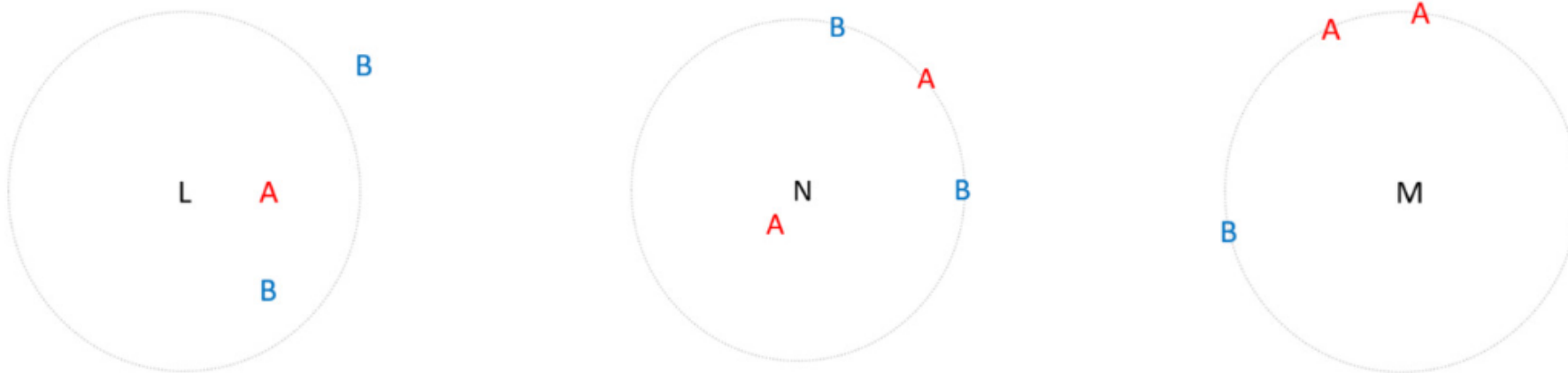
Choosing the k Value : Error Rate

- Rule of thumb is using **sqrt(N)** as k
- Or, iterate with different k and interpret the error rate



$$\text{err}(f, \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(f(\mathbf{x}_i) \neq y_i)$$

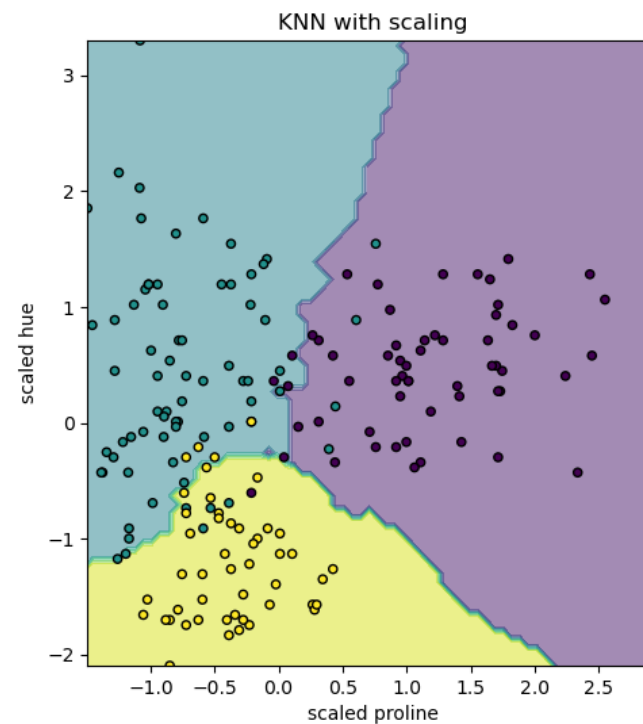
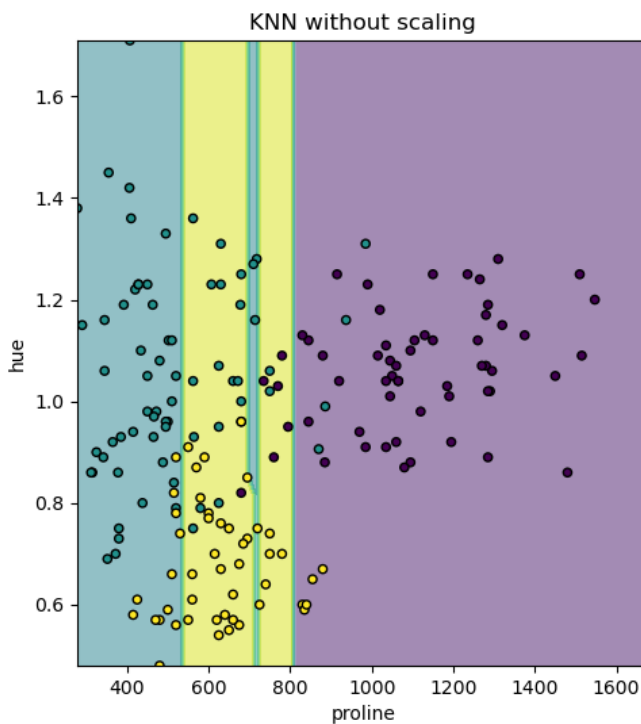
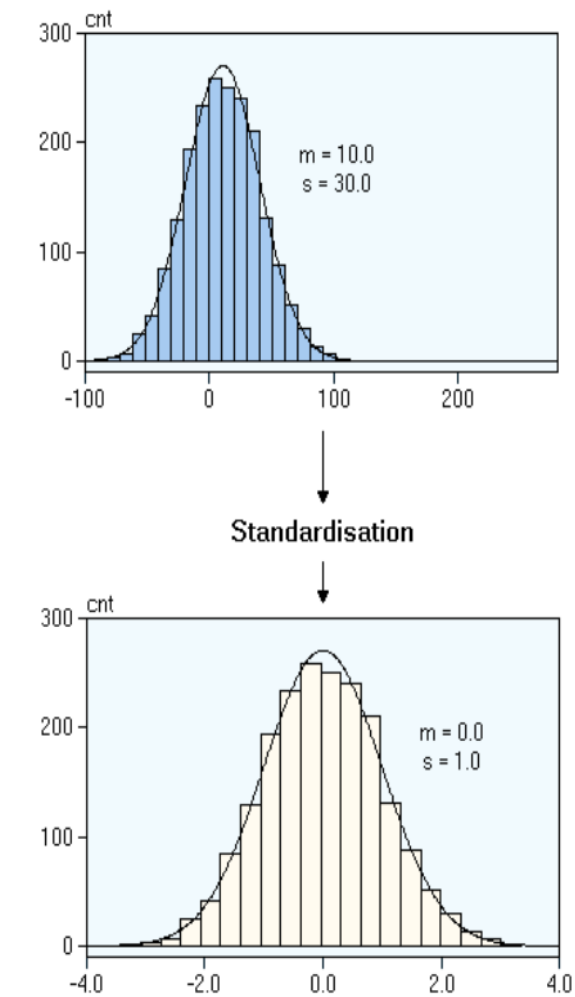
Choosing the k Value : Ties



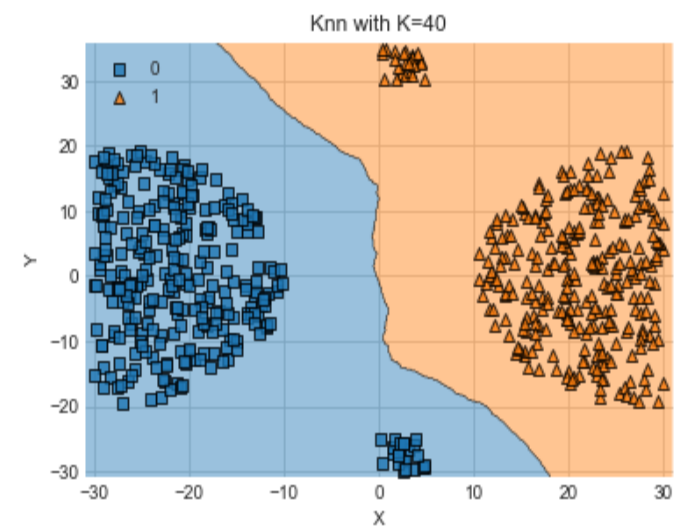
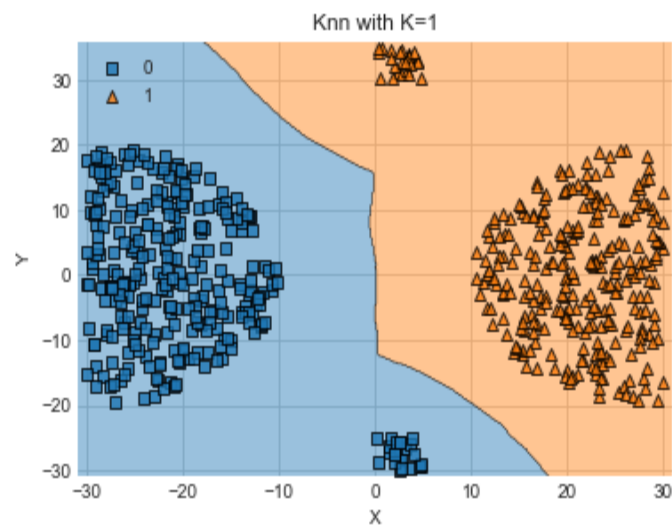
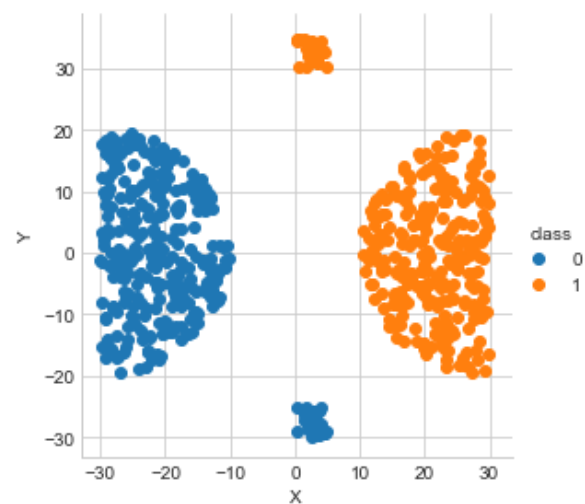
A tie can occur when two or more points are equidistant from an unclassified observation

- Choose different k. For binary classification generally odd numbers are used
- Randomly choose
- Allow observations in until natural stop point

Important Notes : Scaling



Important Notes : Outliers



Best Used, Advantages, Disadvantages

Best Used

When you need a simple algorithm to establish benchmark learning rules.

When memory usage of the trained model is a lesser concern.

When prediction speed of the trained model is a lesser concern.

When data set is small.

Advantages

The algorithm is simple and easy to implement.

The algorithm is versatile. It can be used for classification, regression, and search.

There's no need to tune several parameters or make additional assumptions.

Disadvantages

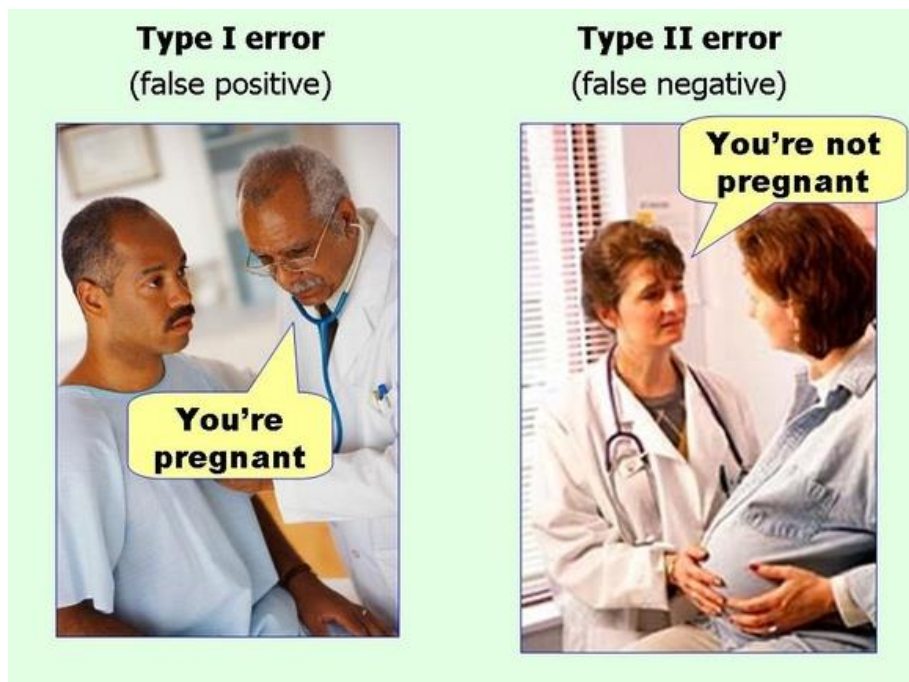
The algorithm gets significantly slower as the number of examples and/or predictors/independent variables increase.

When data set is big.

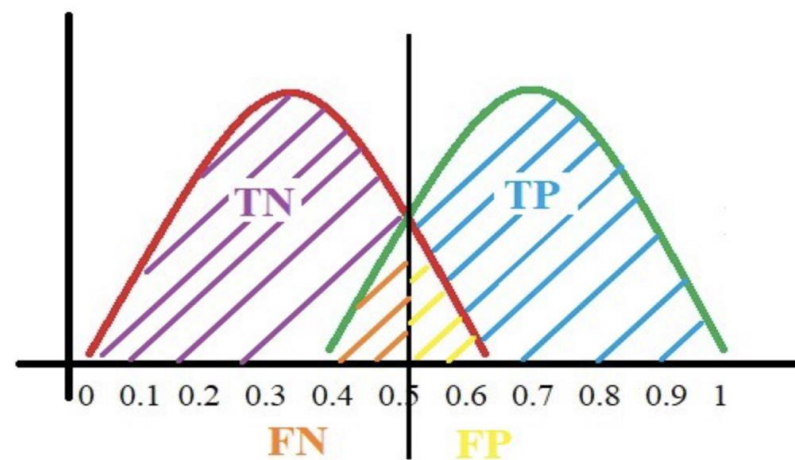
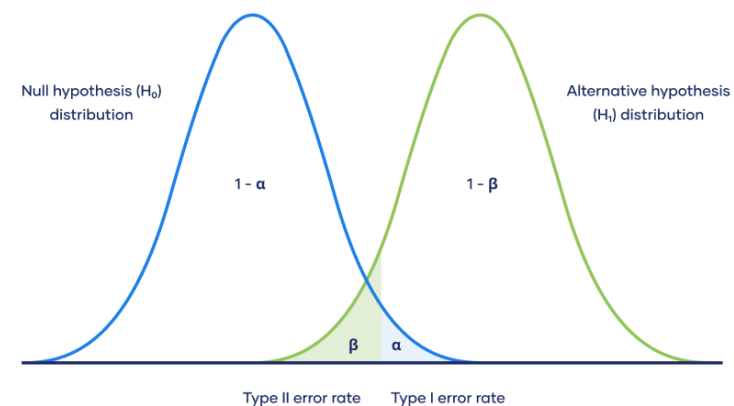
05

Performance Metrics

Performance Metrics : Error Types Revisited



Probability of making Type I and Type II errors



Performance Metrics : Confusion Matrix

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

Why accuracy is not enough?

Performance Metrics : Precision, Recall

Which type of error — FPs or FNs — is more undesirable?

	Predicted Positive	Predicted Negative	
Actual Positive	TP <i>True Positive</i>	FN <i>False Negative</i>	Sensitivity $\frac{TP}{(TP + \text{FN})}$
Actual Negative	FP <i>False Positive</i>	TN <i>True Negative</i>	Specificity $\frac{TN}{(TN + FP)}$
	Precision $\frac{TP}{(TP + FP)}$	Negative Predictive Value $\frac{TN}{(TN + FN)}$	Accuracy $\frac{TP + TN}{(TP + TN + FP + FN)}$

	Cancer	No Cancer
Cancer	80	80
No Cancer	20	820

Accuracy : 90%, Precision: 20%, Recall: %50

- **Precision** measures the extent of **error caused by False Positives (FPs)**
- **Recall** measures the extent of **error caused by False Negatives (FNs)**
- Usually increasing precision will decrease recall, and vice versa.

Performance Metrics : F1 Score

What if, the FN and FP errors equally undesirable?

Set 1: (0.5, 0.5)

Set 2: (0.05, 0.95)

Mean Type	Set 1 Result	Set 2 Result
Arithmetic	0.5	0.5
Geometric	0.5	0.21
Harmonic	0.5	0.09

$$\text{F1 Score} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Performance Metrics : ROC

		Predicted	
		Spam	Not
Actual	Spam	800 (TP)	100 (FN)
	Not	500 (FP)	8600 (TN)

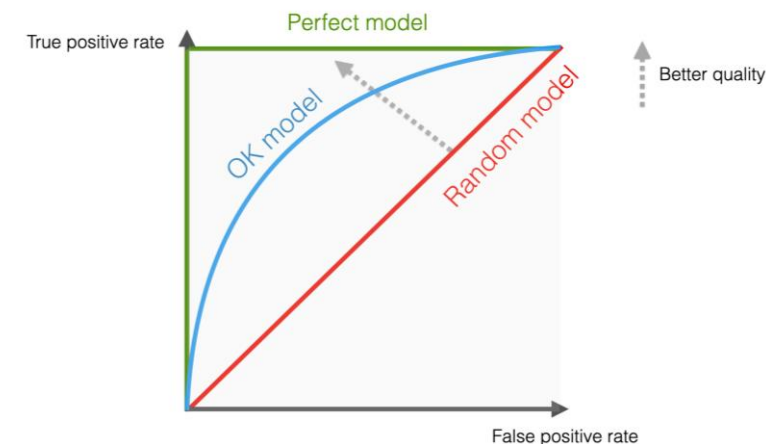
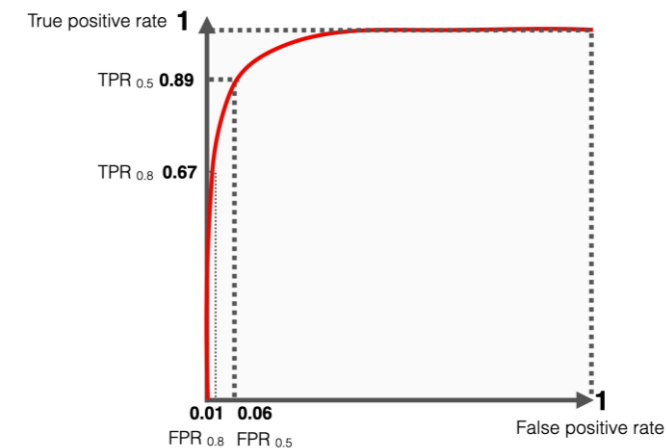
		Predicted	
		Spam	Not
Actual	Spam	600 (TP)	300 (FN)
	Not	100 (FP)	9000 (TN)

		Predicted	
		Spam	Not
Actual	Spam	200 (TP)	700 (FN)
	Not	10 (FP)	9090 (TN)

Decision threshold	0.5	0.8	0.95
Recall	0.89	0.67	0.22
False positive rate	0.06	0.01	0.001

$$\text{TPR / Recall / Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\begin{aligned} \text{FPR} &= 1 - \text{Specificity} \\ &= \frac{\text{FP}}{\text{TN} + \text{FP}} \end{aligned}$$



Performance Metrics : AUC

