



INTRODUCTION TO DATA ANALYTICS

Classification Analysis

Dr. Rathachai Chawuthai

Department of Computer Engineering

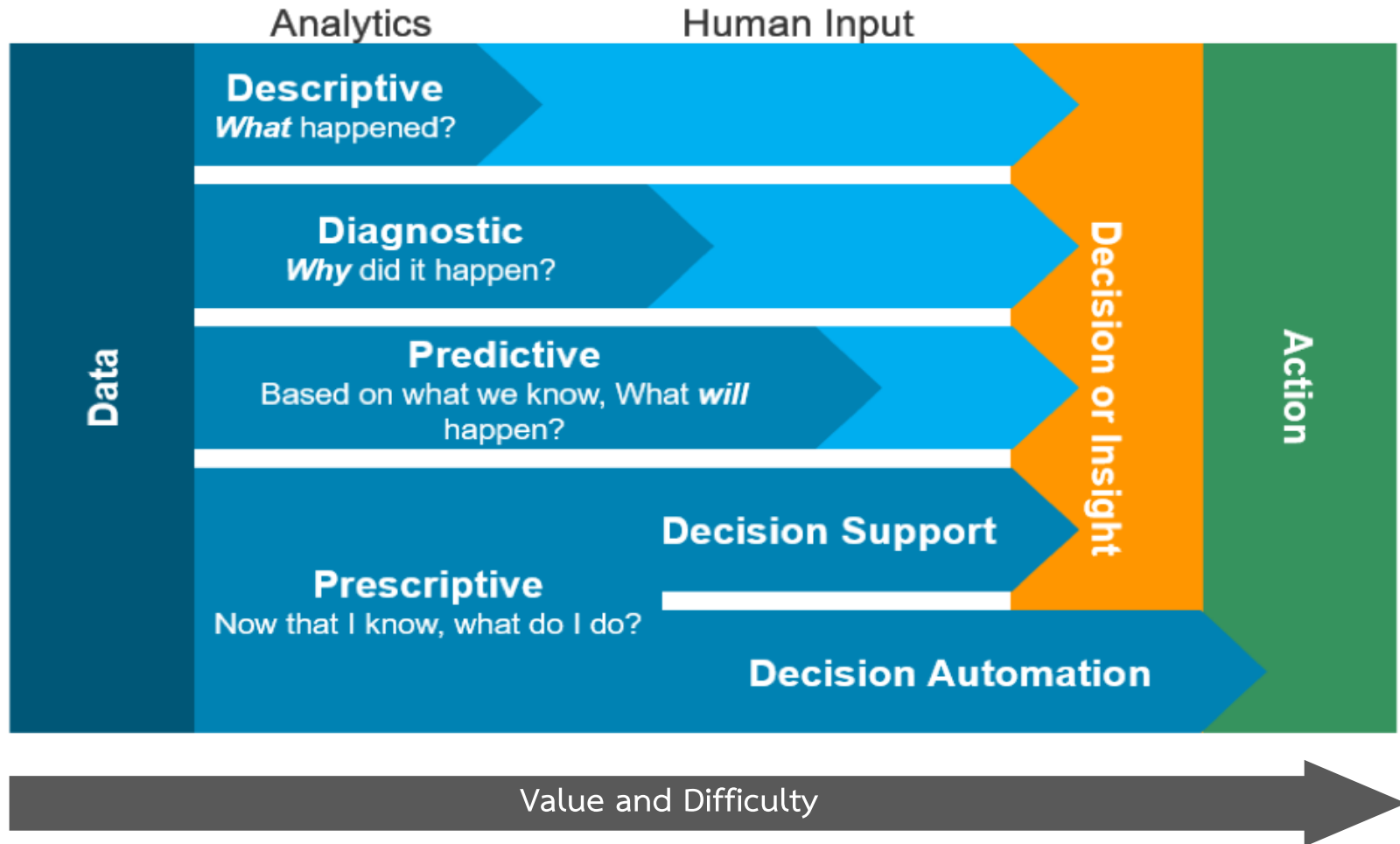
Faculty of Engineering

King Mongkut's Institute of Technology Ladkrabang

Agenda

- Classification
- Classifier I
- Evaluation Methods
- Classifiers II
- Regularization

Data Analytics



Machine Learning

```
graph TD; ML[Machine Learning] --> SL[Supervised Learning]; ML --> UL[Unsupervised Learning]; SL --> R[Regression]; SL --> C[Classification]; UL --> Cl[Clustering]; R --- RL["• Linear Regression<br>• Polynomial Regression"]; C --- CL["• Decision Tree<br>• Logistic Regression<br>• Neural Network<br>• etc."]; Cl --- CL2["• K-Means<br>• DB-SCAN<br>• etc."];
```

Supervised Learning

Develop predictive model based on both input and output data

Regression

- Linear Regression
- Polynomial Regression

Classification

- Decision Tree
- Logistic Regression
- Neural Network
- etc.

Unsupervised Learning

Develop predictive model based on both input and output data

Clustering

- K-Means
- DB-SCAN
- etc.

Classification



Classes



Classes



Classes



Classes

Table



- 4 leg
- sit (yes)
- sleep (no)

Bed



- 4 legs
- mattress (yes)
- sleep (yes)

Sofa



- 4 legs
- mattress (yes)
- sit (yes)

Classes

Table



Class/
Label

- 4 leg
- sit (yes)
- sleep (no)

Attributes/
Features

- 4 legs
- mattress (yes)
- sleep (yes)

Bed



Sofa



- 4 legs
- mattress (yes)
- sit (yes)

Classes

Table



- 4 leg
- sit (yes)
- sleep (no)

- 4 legs
- mattress (yes)
- sleep (yes)

Bed



Sofa



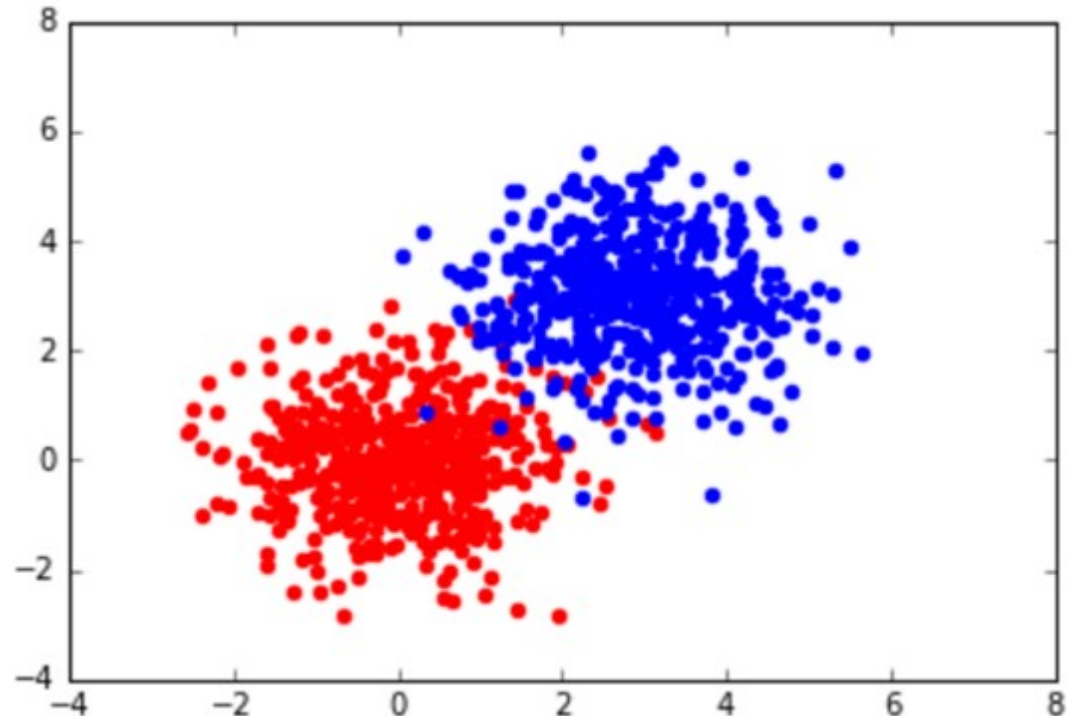
- 4 legs
- mattress (yes)
- sit (yes)

Classification

- Class named by Label
- Attributes or Features

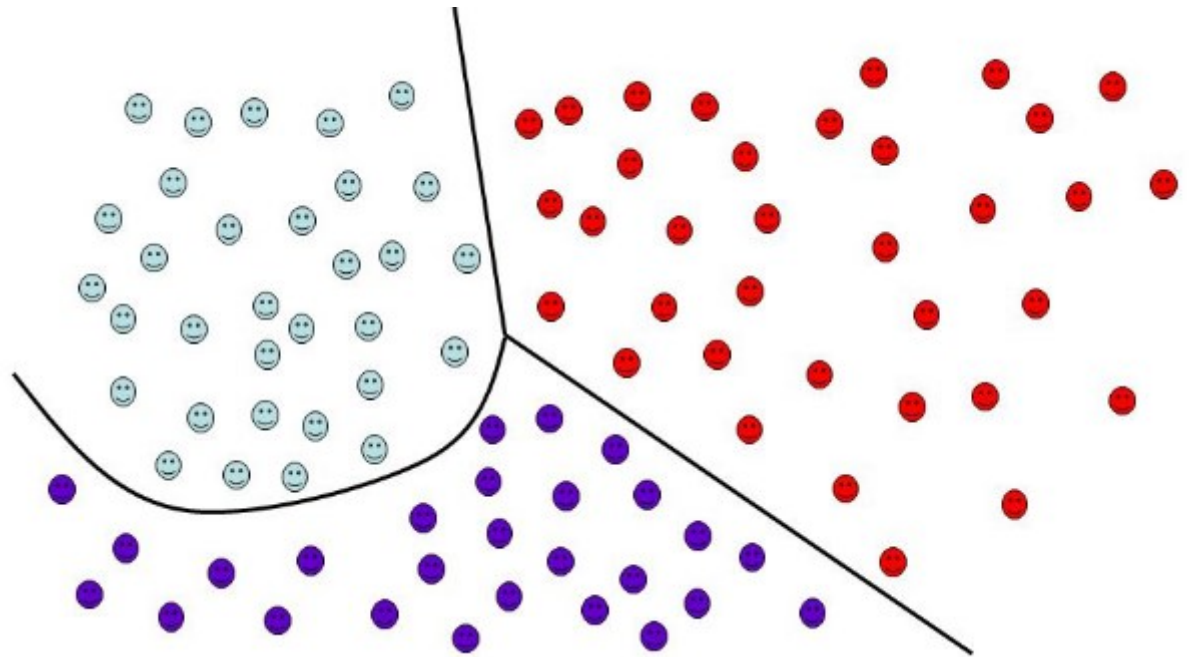
Binary Classification

- Has 2 Classes
- We often have to deal with the simple task of Binary Classification. Some examples are:
 - Sentiment Analysis (positive/negative),
 - Spam Detection (spam/not-spam),
 - Fraud Detection (fraud/not-fraud).

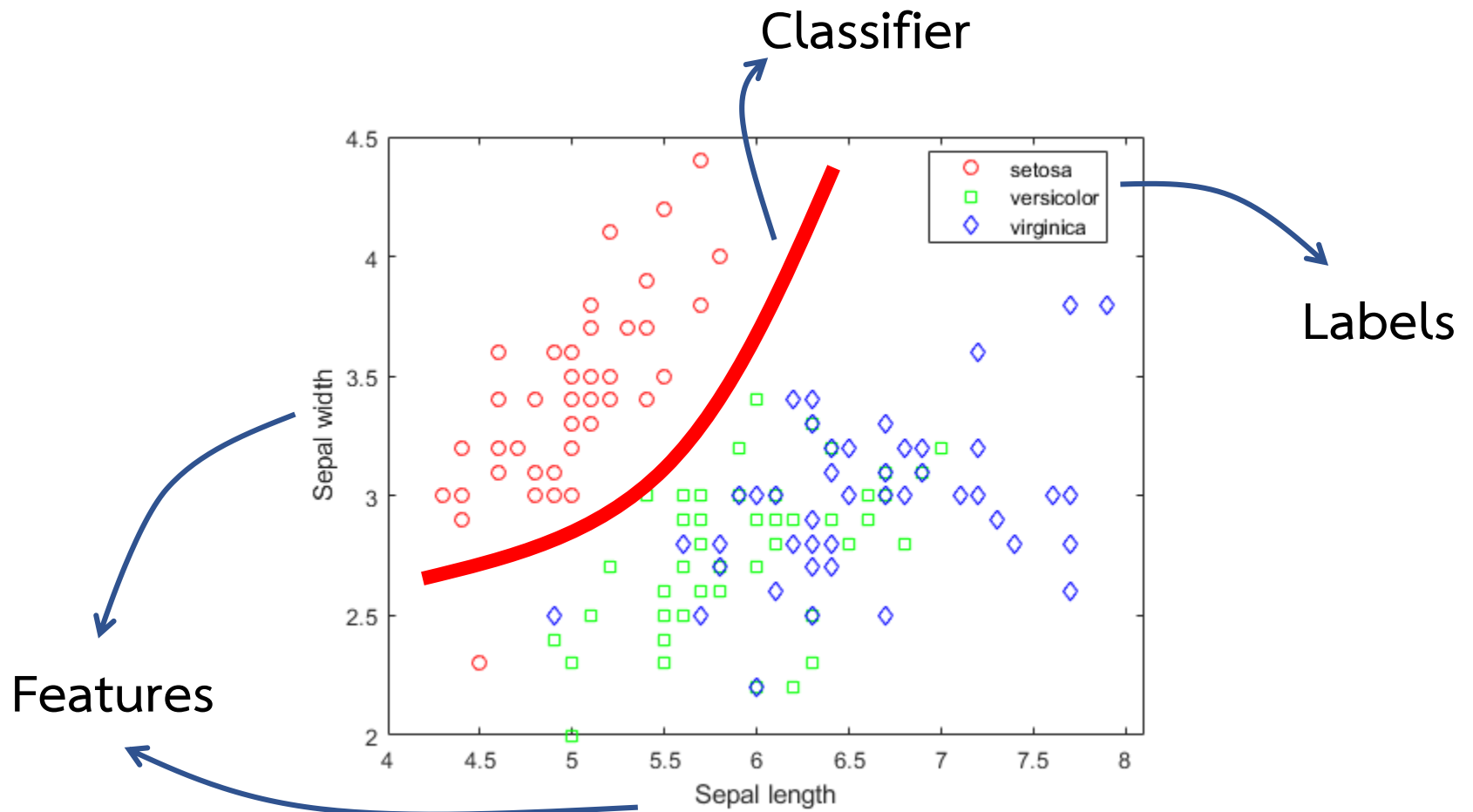


Multi-Class Classification

- Has more than 2 classes
- For example,
 - Genres of Movies
 - Grades
 - Sentiment Analysis (Positive, Neutral, Negative)
 - etc.



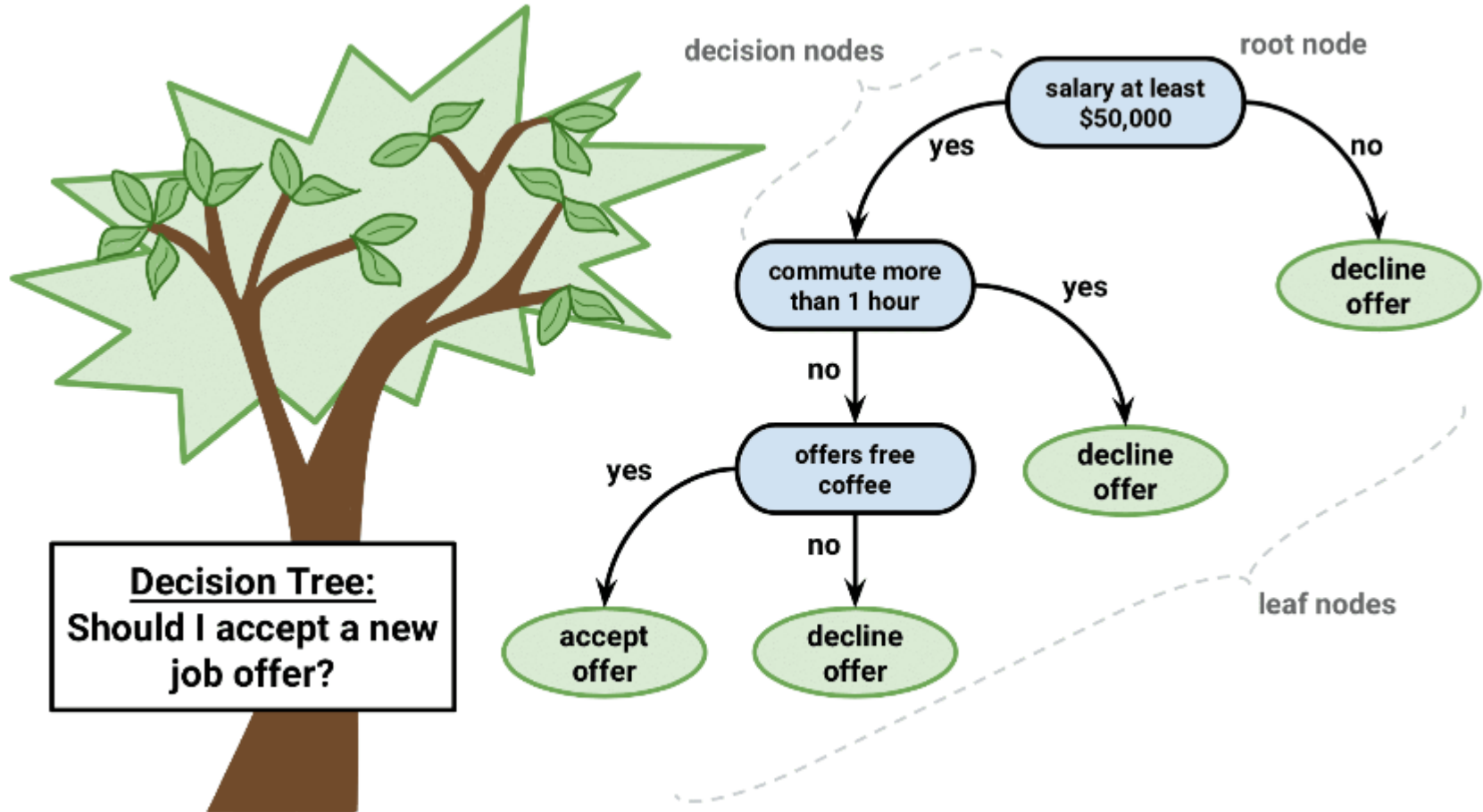
Classification



Classifier I



Decision Tree



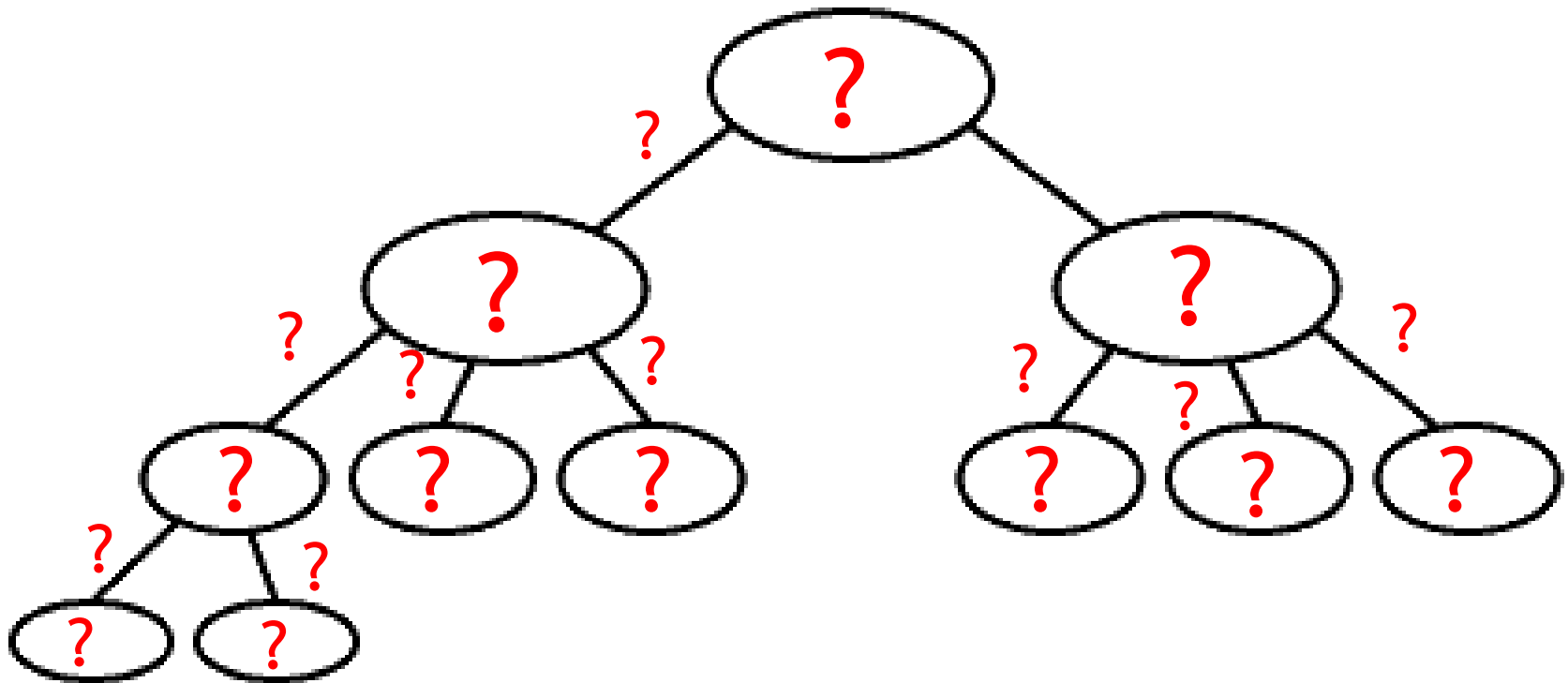
Decision Tree

- A decision tree is a decision support tool that uses a tree-like graph or model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility. It is one way to display an algorithm that only contains conditional control statements.
- Decision trees are commonly used in operations research, specifically in decision analysis, to help identify a strategy most likely to reach a goal, but are also a popular tool in machine learning.
- A decision tree consists of three types of nodes:
 - Decision nodes – typically represented by squares
 - Chance nodes – typically represented by circles
 - End nodes – typically represented by triangles

Case

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

Tree



Case

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

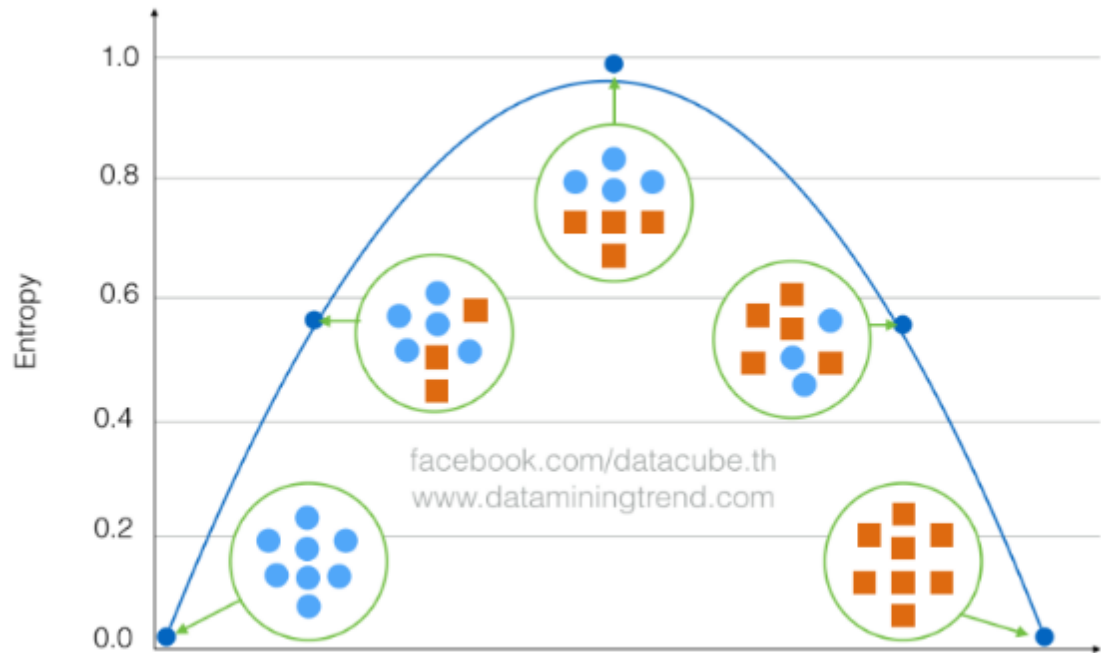
Case

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

Information Gain

- Same things --> Low
- Different things --> High

$$Entropy = \sum_{i=1}^C -p_i * \log_2(p_i)$$



Decision Tree: Construction

Finding entropy H of every attribute

1) Headache

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

Headache = YES

Headache = no

$H(\text{Flu} \mid \text{Headache}) =$

-

-

Decision Tree: Construction

Finding entropy H of
every attribute

1) Headache

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

Headache = YES

$H(\text{Flu} \mid \text{Headache}) =$

-

$\frac{4}{7}$

Headache = no

$\frac{3}{7}$

Decision Tree: Construction

Finding entropy H of every attribute

1) Headache

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

Headache = YES

$H(\text{Flu} \mid \text{Headache}) =$

$$-\frac{4}{7} \left[\boxed{} + \boxed{} \right]$$

Flu=YES

Flu=no

Headache = no

$$-\frac{3}{7} \left[\boxed{} + \boxed{} \right]$$

Flu=YES

Flu=no

Decision Tree: Construction

Finding entropy H of every attribute

1) Headache

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

Headache = YES

Headache = no

$H(\text{Flu} \mid \text{Headache}) =$

$$- \frac{4}{7} \left[\frac{3}{7} \log_2 \left(\frac{3}{4} \right) + \frac{1}{7} \log_2 \left(\frac{1}{4} \right) \right] - \frac{3}{7} \left[\quad + \quad \right]$$

↑
↑

Flu=YES
Flu=no

Decision Tree: Construction

Finding entropy H of every attribute

1) Headache

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

Headache = YES

$H(\text{Flu} \mid \text{Headache}) =$

$$- \frac{4}{7} \left[\frac{3}{4} \log_2 \left(\frac{3}{4} \right) + \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right]$$

Headache = no

$$- \frac{3}{7} \left[\frac{0}{3} \log_2 \left(\frac{0}{4} \right) + \frac{3}{3} \log_2 \left(\frac{3}{3} \right) \right]$$

Flu=YES

Flu=no

Decision Tree: Construction

Finding entropy H of every attribute

1) Headache

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

$$\begin{aligned} H(\text{Flu} \mid \text{Headache}) = & -\frac{4}{7} \left[\frac{3}{4} \log_2 \left(\frac{3}{4} \right) + \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right] - \frac{3}{7} \left[\frac{0}{3} \log_2 \left(\frac{0}{4} \right) + \frac{3}{3} \log_2 \left(\frac{3}{3} \right) \right] \\ & = 0.464 \end{aligned}$$

Note: In the original image, red arrows point from the terms $\frac{0}{3} \log_2 \left(\frac{0}{4} \right)$ and $\frac{3}{3} \log_2 \left(\frac{3}{3} \right)$ to a red '0' above the second bracket, indicating that these terms are zero.

Decision Tree: Construction

Finding entropy H
of every attribute

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no

$$H(\text{Flu} \mid \text{Headache}) = 0.464$$

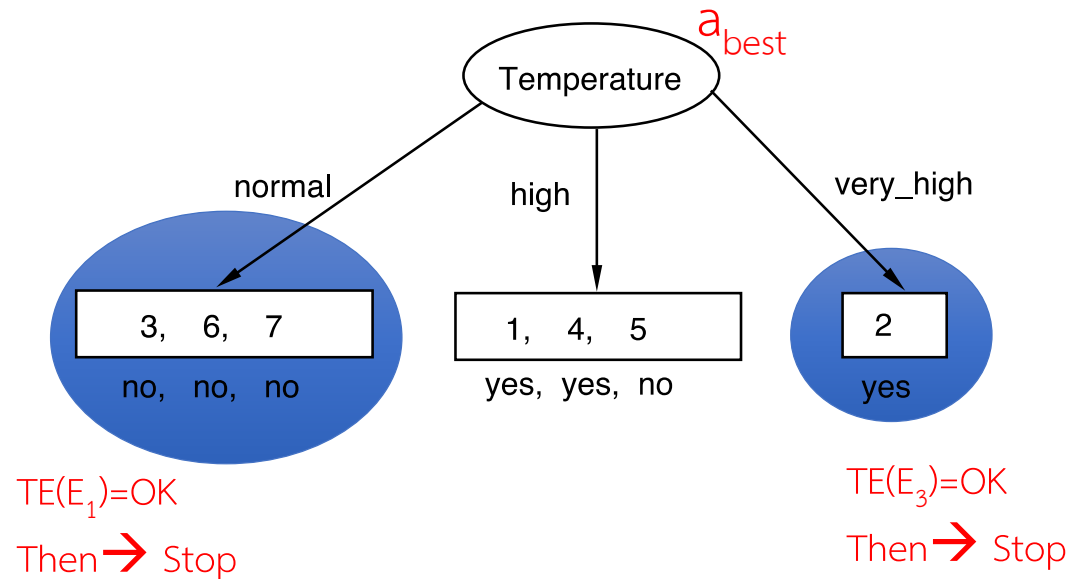
$$H(\text{Flu} \mid \text{Temperature}) = 0.394$$

$$H(\text{Flu} \mid \text{Nausea}) = 0.620$$

Lowest

Decision Tree: Construction

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no



$$H(\text{Flu} \mid \text{Headache}) = 0.464$$

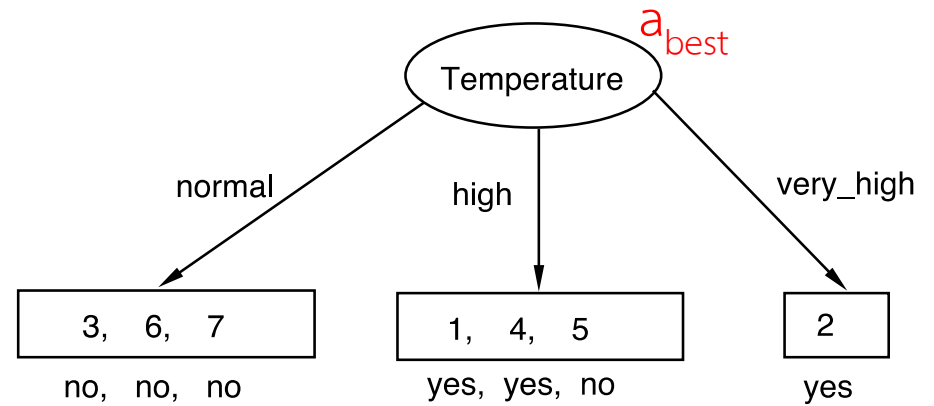
$$H(\text{Flu} \mid \text{Temperature}) = 0.394$$

$$H(\text{Flu} \mid \text{Nausea}) = 0.620$$

Lowest

Decision Tree: Construction

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no



Finding entropy H of attributes from the remainings

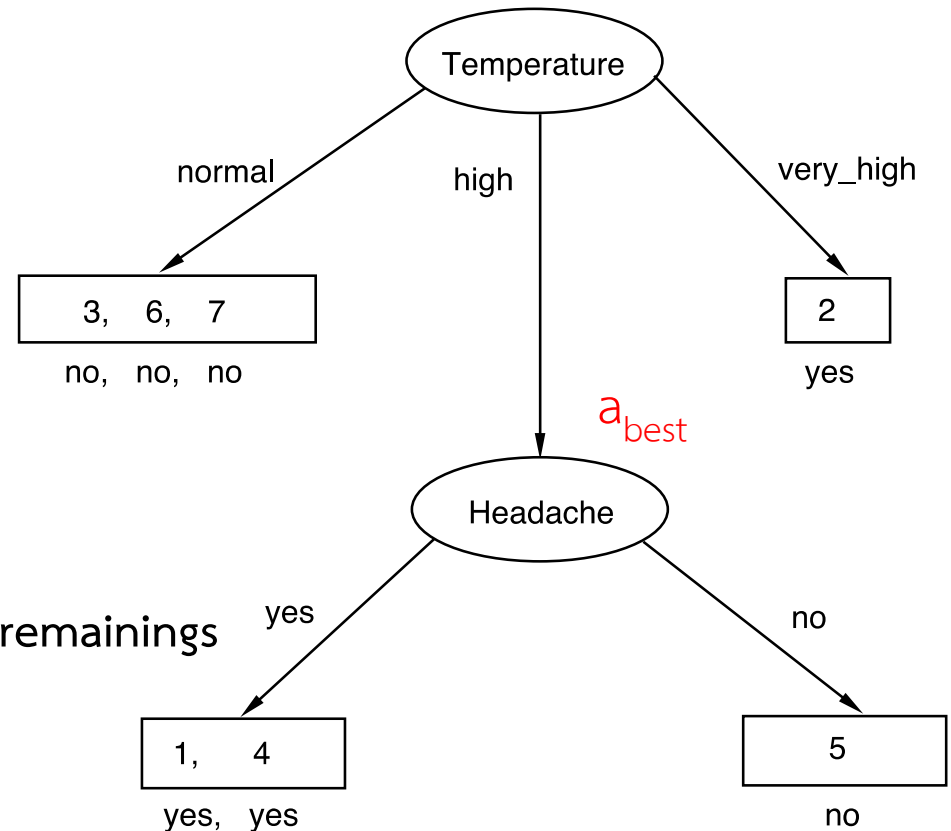
$$H(\text{Flu}_{\text{Temperature} = \text{high}} \mid \text{Headache}) = 0$$

Lowest

$$H(\text{Flu}_{\text{Temperature} = \text{high}} \mid \text{Nausea}) = 0.667$$

Decision Tree: Construction

CASE	TEMPERATURE	HEADACHE	NAUSEA	Flu
1	high	YES	no	YES
2	very_high	YES	YES	YES
3	normal	no	no	no
4	high	YES	YES	YES
5	high	no	YES	no
6	normal	YES	no	no
7	normal	no	YES	no



Finding entropy H of attributes from the remainings

$$H(\text{Flu}_{\text{Temperature} = \text{high}} \mid \text{Headache}) = 0$$

$$H(\text{Flu}_{\text{Temperature} = \text{high}} \mid \text{Nausea}) = 0.667$$

Evaluation Methods



Accuracy

- Accuracy is the most popular performance measure used and for good reason. It's extremely helpful, simple to compute and to understand. It is the proportion of the correctly classified samples and all the samples.

```
1 from sklearn.metrics import accuracy_score
2
3 print accuracy_score(y_test, model.predict(X_test))  # 0.98
4
```

TP | TN | FP | FN

- There are other ways to measure different aspects of performance. In classic machine learning nomenclature, when we're dealing with binary classification, the classes are: **positive** and **negative**. Think of these classes in the context of disease detection:
 - **positive** – we predict the disease is present
 - **negative** – we predict the disease is not present.
- Let's now define some notations:
 - **TP** – True Positives (Samples the classifier has correctly classified as positives)
 - **TN** – True Negatives (Samples the classifier has correctly classified as negatives)
 - **FP** – False Positives (Samples the classifier has incorrectly classified as positives)
 - **FN** – False Negatives (Samples the classifier has incorrectly classified as negatives)

TP | TN | FP | FN

	predicted: YES	predicted: NO
actual: YES	4211 ✓ (true positive)	181 ✗ (false negative)
actual: NO	94 ✗ (false positive)	1207 ✓ (true negative)

Evaluation Methods

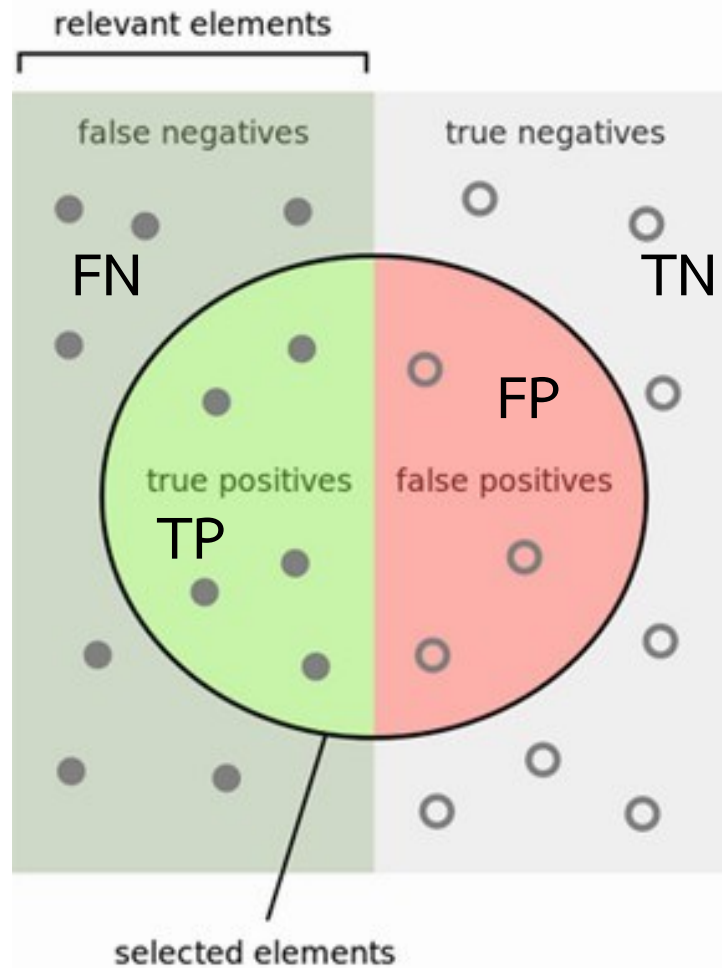
$$\textit{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Evaluation Methods

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

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

Precision & Recall



How many selected items are relevant?

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

How many relevant items are selected?

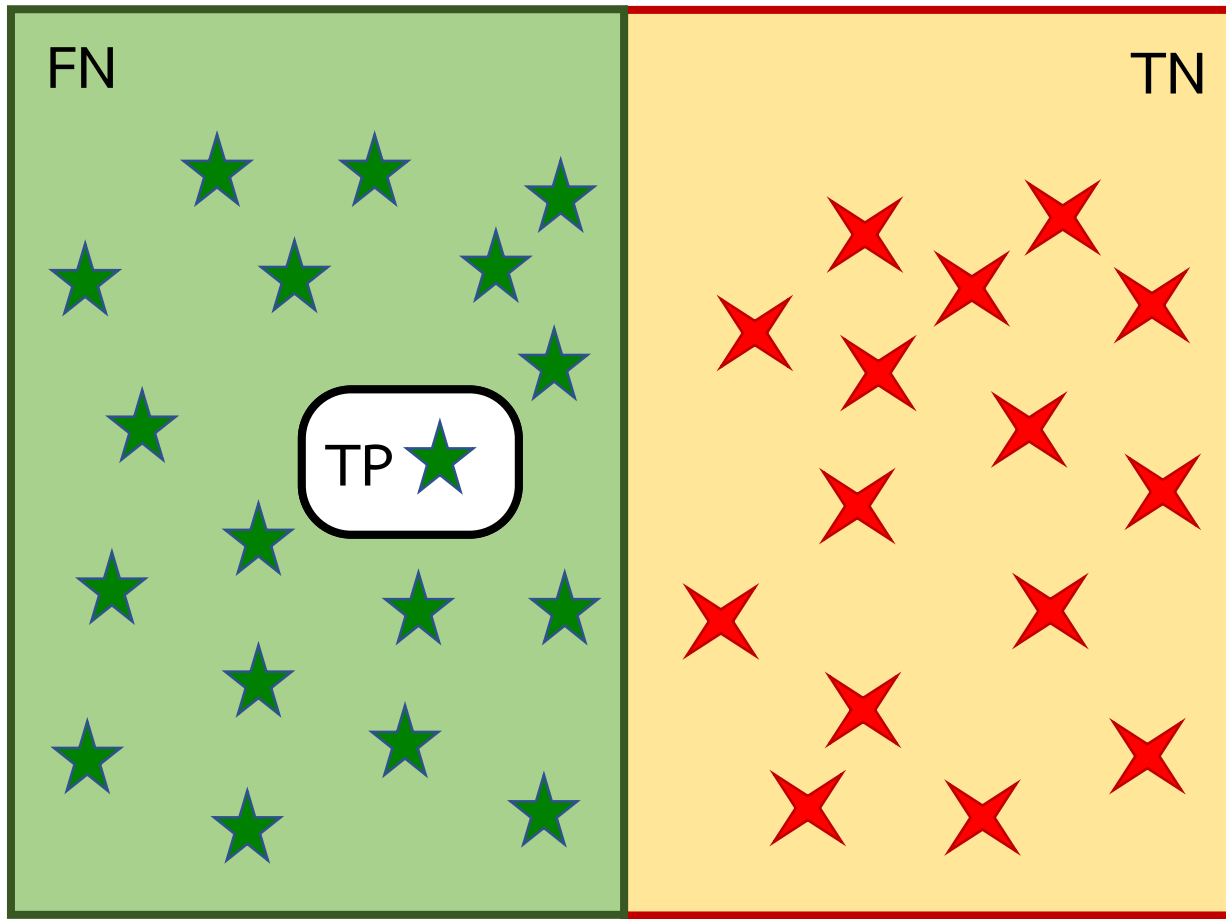
$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Question

Recall = 0.99

is it good?

Perfect



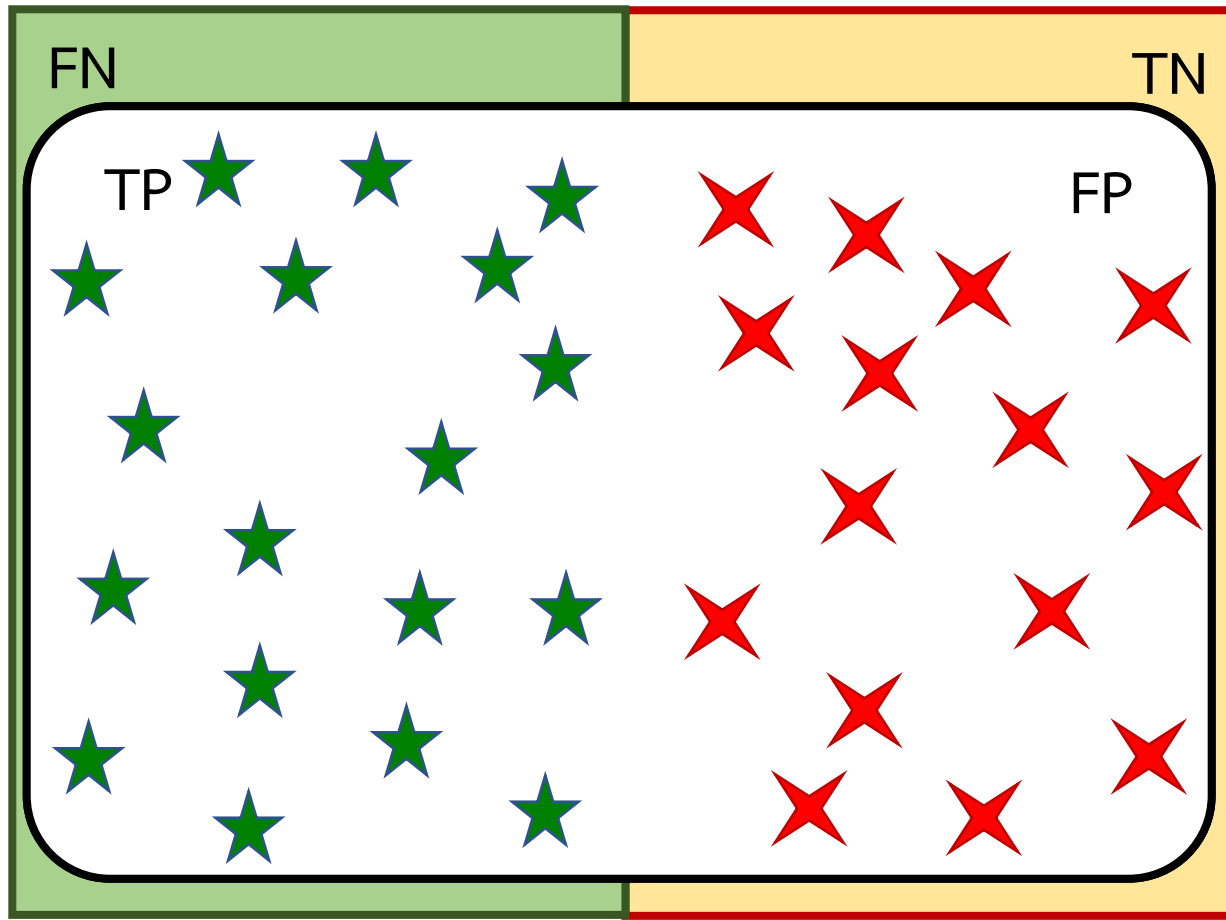
Precision = 1

Question

Recall = 0.99

is it good?

Perfect



Recall = 1

Evaluation Methods

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

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

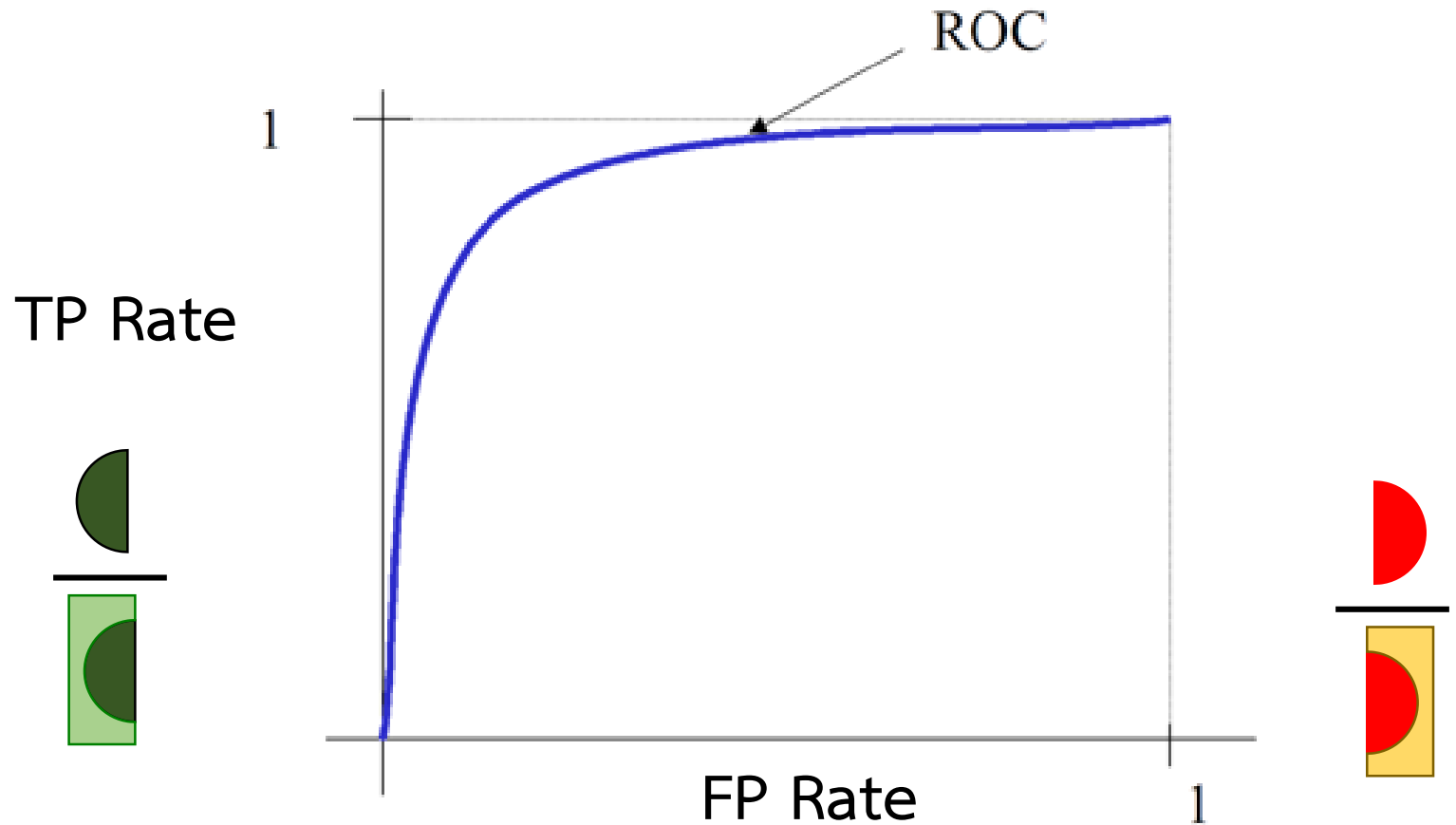
$$F1 = 2 \times \frac{\textit{Precision} \times \textit{Recall}}{\textit{Precision} + \textit{Recall}}$$

Perfect



ROC Curve

- Receive Operating Characteristic

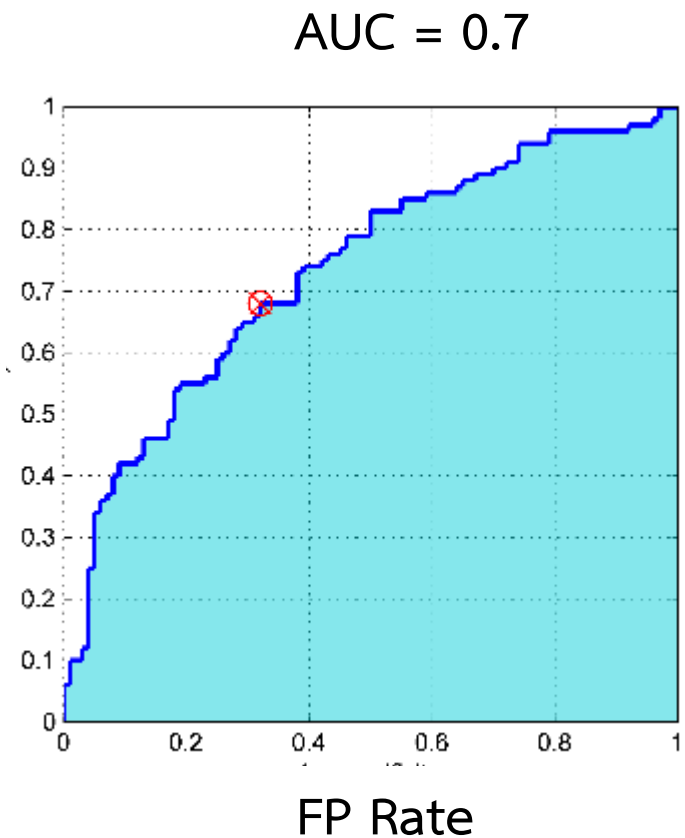


AUC : Area Under the ROC Curve

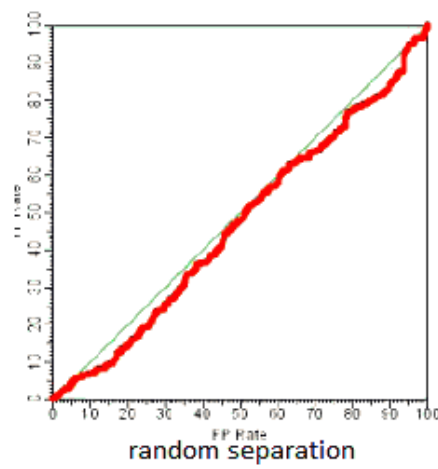
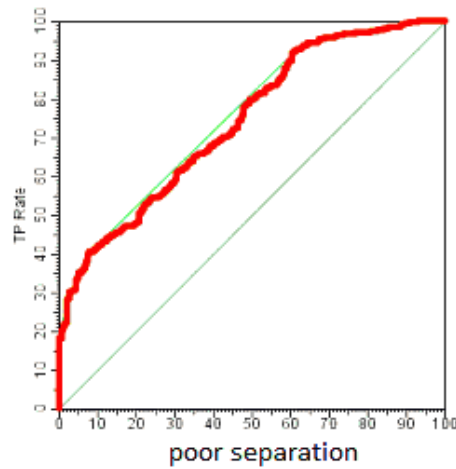
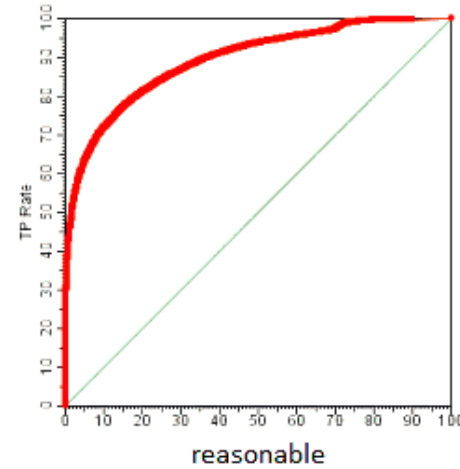
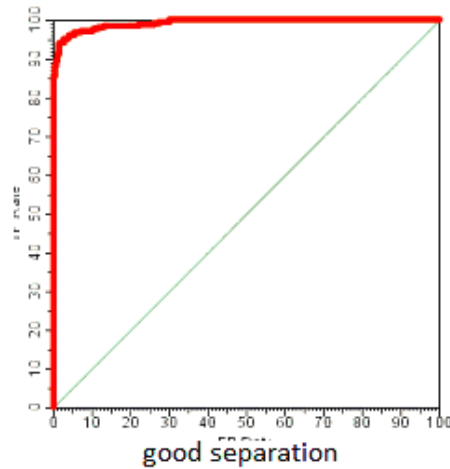
The graph at right shows three ROC curves representing excellent, good, and worthless tests plotted on the same graph. An area of 1 represents a perfect test; an area of 0.5 represents a worthless test. A rough guide for classifying the accuracy of a diagnostic test is the traditional academic point system:

- .90-1 = excellent (A)
- .80-.90 = good (B)
- .70-.80 = fair (C)
- .60-.70 = poor (D)
- .50-.60 = fail (F)

TP Rate



AUC

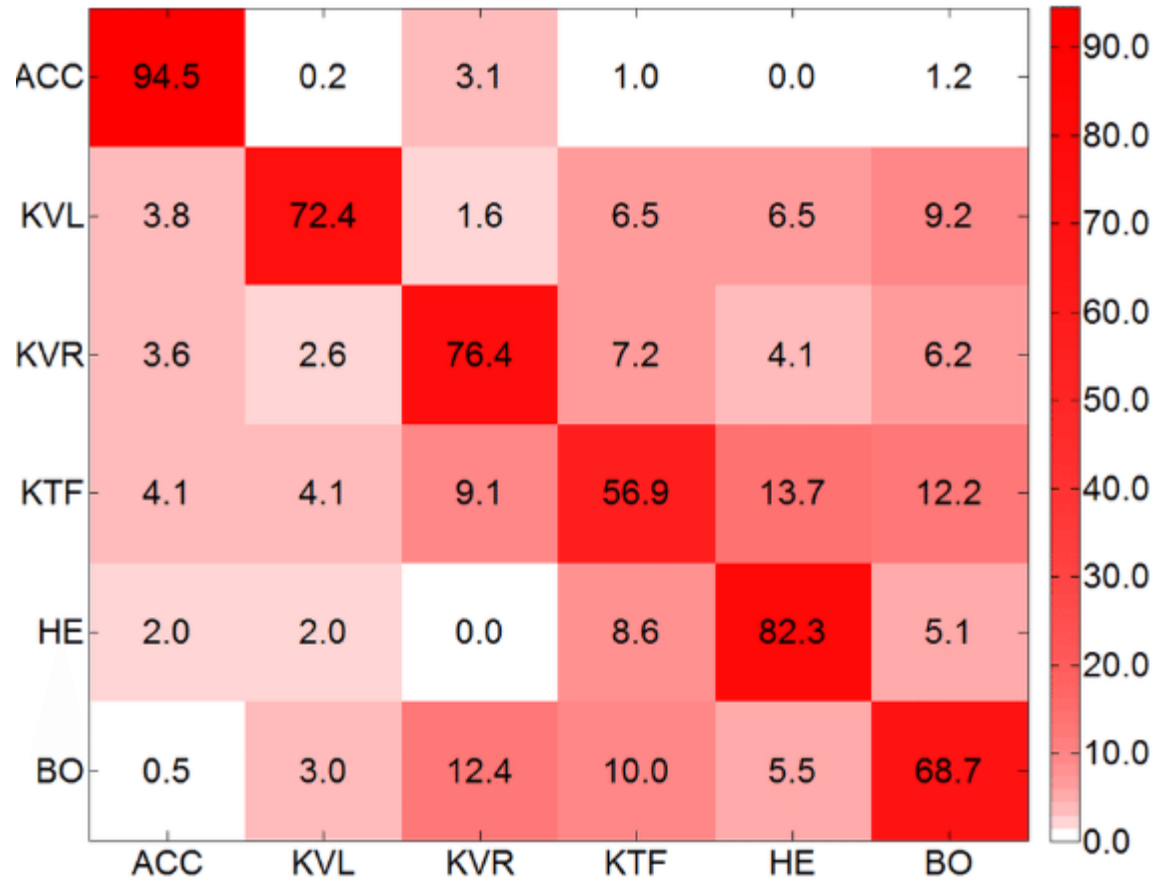


Confusion Matrix

		Predicted Class		
		Cat	Dog	Rabbit
Actual Class	Cat	5	2	0
	Dog	3	3	2
	Rabbit	0	1	11

Confusion Matrix : Heatmap

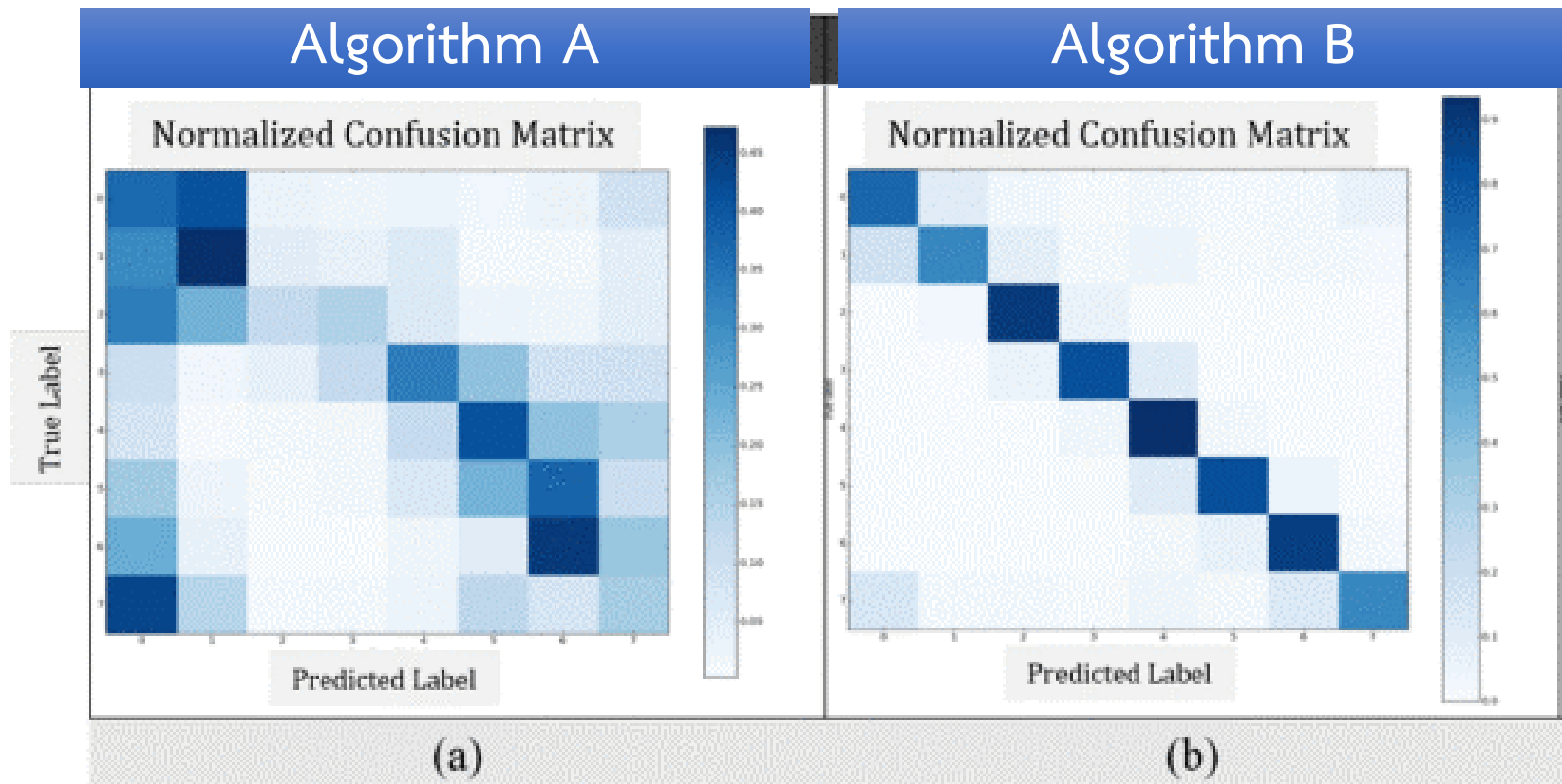
Actual
Class



Predicted Class

Confusion Matrix : Comparison

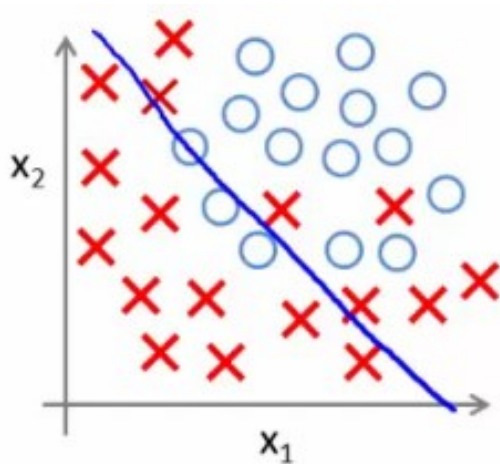
Which one is better?



Classification Methods II



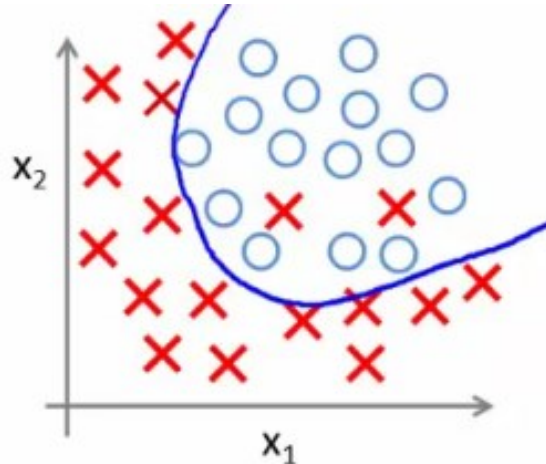
Underfitting vs.. Overfitting



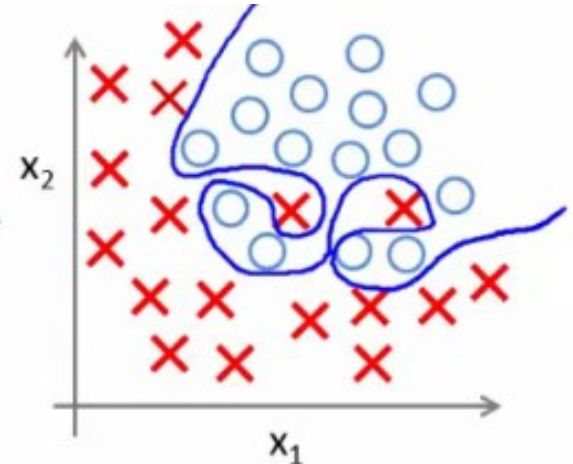
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

(g = sigmoid function)

UNDERFITTING
(high bias)



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



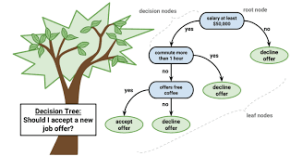
$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

OVERFITTING
(high variance)

Classification Techniques

- Linear classifiers
 - Fisher's linear discriminant
 - Logistic regression
 - Naive Bayes classifier
 - Perceptron
- Support vector machines
- Quadratic classifiers
- Kernel estimation
 - k-nearest neighbor
- Boosting (meta-algorithm)
- Decision trees
- Neural networks
- Learning vector quantization

Decision Tree



`DecisionTreeClassifier` is a class capable of performing multi-class classification on a dataset.

As with other classifiers, `DecisionTreeClassifier` takes as input two arrays: an array `X`, sparse or dense, of size `[n_samples, n_features]` holding the training samples, and an array `Y` of integer values, size `[n_samples]`, holding the class labels for the training samples:

```
>>> from sklearn import tree
>>> X = [[0, 0], [1, 1]]
>>> Y = [0, 1]
>>> clf = tree.DecisionTreeClassifier()
>>> clf = clf.fit(X, Y)
```

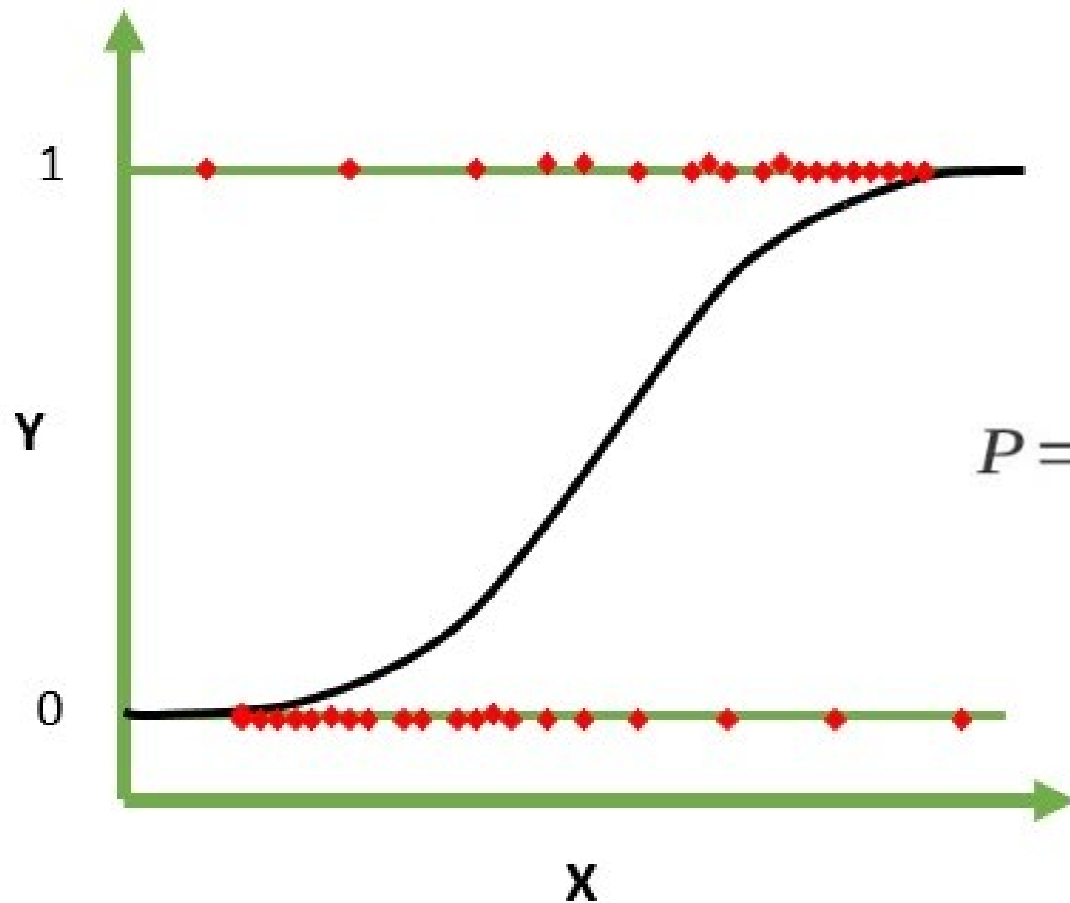
After being fitted, the model can then be used to predict the class of samples:

```
>>> clf.predict([[2., 2.]])
array([1])
```

Alternatively, the probability of each class can be predicted, which is the fraction of training samples of the same class in a leaf:

```
>>> clf.predict_proba([[2., 2.]])
array([[ 0.,  1.]])
```


Logistic Regression



$$P = \frac{e^{(\beta + \alpha_1 X_1 + \dots + \alpha_n X_n)}}{1 + e^{(\beta + \alpha_1 X_1 + \dots + \alpha_n X_n)}}$$

Logistic Regression

```
from sklearn.linear_model import LogisticRegression
from sklearn import metrics

clf = LogisticRegression()
clf.fit(X_train, y_train)
```

Logistic Regression

```
from sklearn import model_selection
from sklearn.model_selection import cross_val_score
kfold = model_selection.KFold(n_splits=10, random_state=7)

clf = LogisticRegression()

results = model_selection.cross_val_score
         (clf, X, y, cv=kfold, scoring= 'accuracy')

print("10-fold accuracy: %.3f" % (results.mean()))
```

Logistic Regression

```
from sklearn import model_selection
from sklearn.model_selection import cross_val_score
kfold = model_selection.KFold(n_splits=10, random_state=7)

clf = LogisticRegression()

results = model_selection.cross_val_score
        (clf, X, y, cv=kfold, scoring= 'accuracy')

print("10-fold accuracy: %.3f" % (results.mean()))
```

Logistic Regression

```
from sklearn.metrics import confusion_matrix  
  
confusion_matrix = confusion_matrix(y_test, y_pred)  
  
print(confusion_matrix)
```

```
[[10872, 109]  
 [ 1122, 254]]
```

The result is telling us that we have 10872+254 correct predictions and 1122+109 incorrect predictions.

Logistic Regression

```
from sklearn.metrics import classification_report  
  
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score
0	0.91	0.99	0.95
1	0.70	0.18	0.29
avg / total	0.88	0.90	0.87

Naive Bayes

- $P(c|x)$ is the posterior probability of class (*target*) given predictor (*attribute*)
- $P(c)$ is the prior probability of class.
- $P(x|c)$ is the likelihood which is the probability of predictor given class.
- $P(x)$ is the prior probability of predictor.

The diagram shows the Naive Bayes formula with arrows pointing from descriptive labels to the terms in the equation:

$$P(c | x) = \frac{P(x | c) P(c)}{P(x)}$$

Labels and their corresponding terms:

- Likelihood** points to $P(x | c)$
- Class Prior Probability** points to $P(c)$
- Posterior Probability** points to $P(c | x)$
- Predictor Prior Probability** points to $P(x)$

$$P(c | \mathbf{X}) = P(x_1 | c) \times P(x_2 | c) \times \cdots \times P(x_n | c) \times P(c)$$

Naive Bayes

$$P(c | x) = \frac{P(x | c)P(c)}{P(x)}$$

Frequency Table		Play Golf	
		Yes	No
Outlook	Sunny	3	2
	Overcast	4	0
	Rainy	2	3



Likelihood Table		Play Golf		
		Yes	No	
Outlook	Sunny	3/9	2/5	5/14
	Overcast	4/9	0/5	4/14
	Rainy	2/9	3/5	5/14
		9/14	5/14	

$$P(x | c) = P(\text{Sunny} | \text{Yes}) = 3 / 9 = 0.33$$

$$P(x) = P(\text{Sunny}) = 5 / 14 = 0.36$$

$$P(c) = P(\text{Yes}) = 9 / 14 = 0.64$$

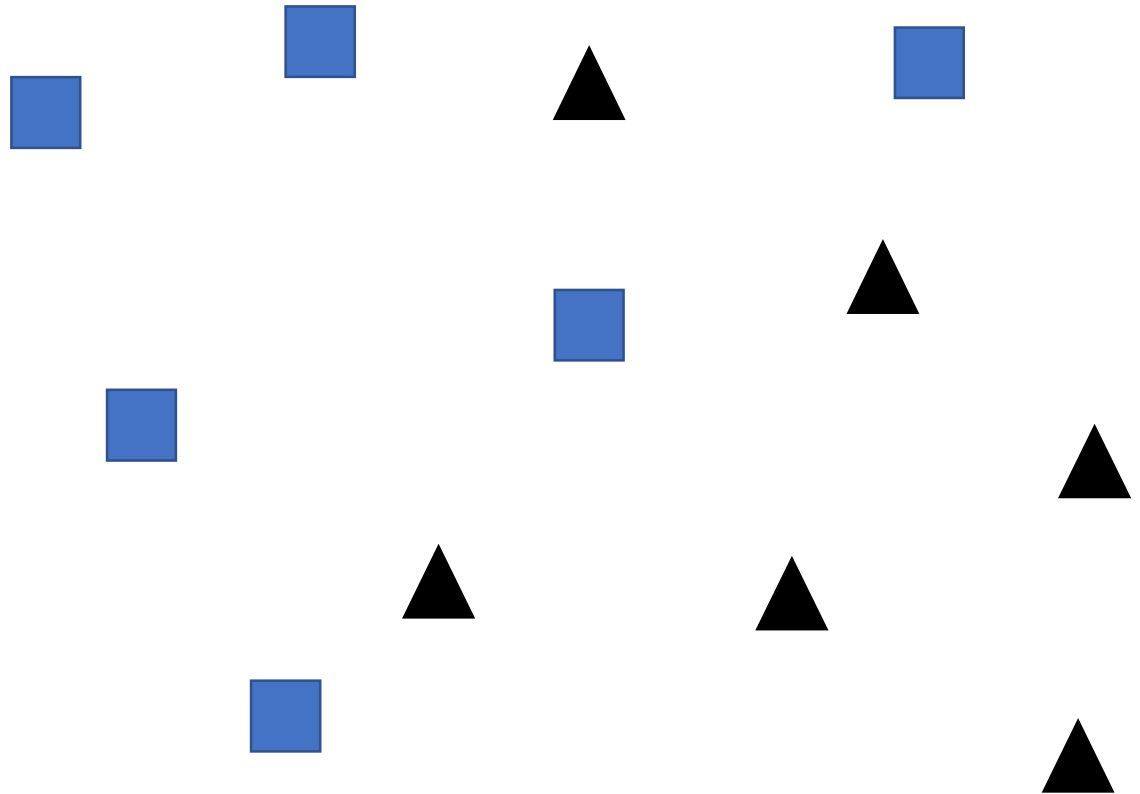
Posterior Probability:

$$P(c | x) = P(\text{Yes} | \text{Sunny}) = 0.33 \times 0.64 \div 0.36 = 0.60$$

Naive Bayes

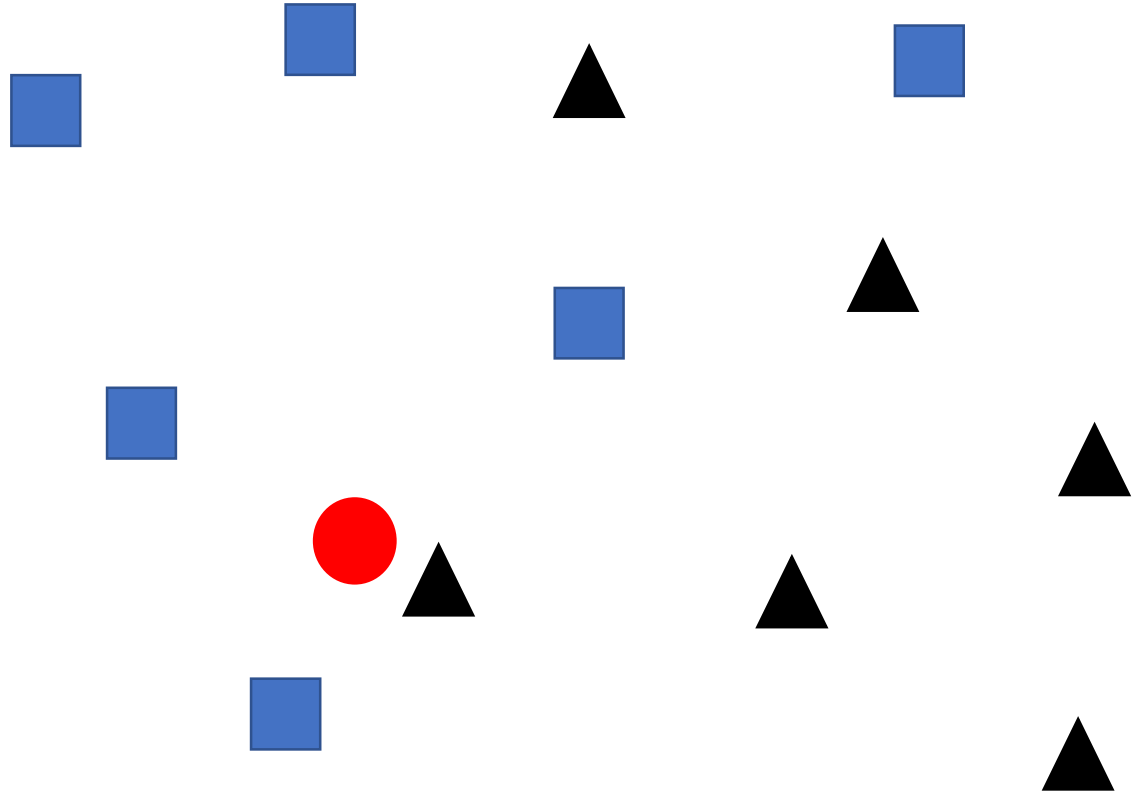
```
>>> import numpy as np
>>> X = np.array([[-1, -1], [-2, -1], [-3, -2], [1, 1], [2, 1], [3, 2]])
>>> Y = np.array([1, 1, 1, 2, 2, 2])
>>> from sklearn.naive_bayes import GaussianNB
>>> clf = GaussianNB()
>>> clf.fit(X, Y)
GaussianNB(priors=None)
>>> print(clf.predict([[-0.8, -1]]))
[1]
```

K-Nearest Neighbor



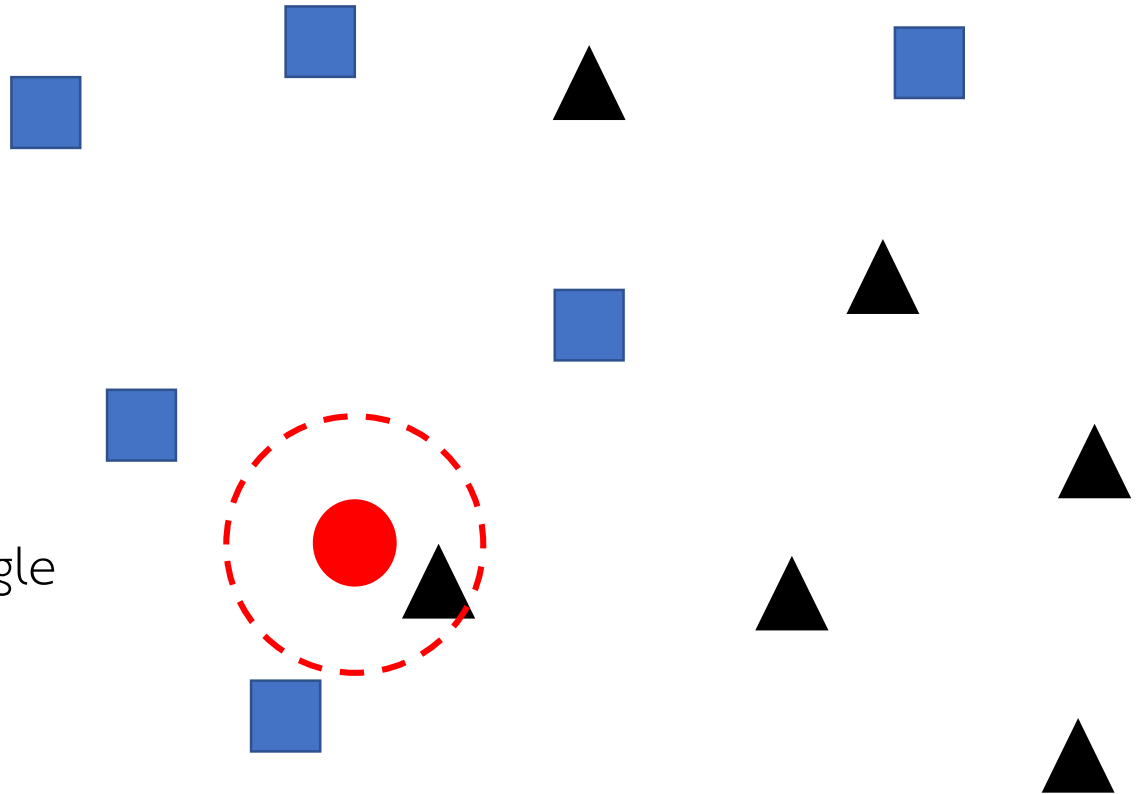
K-Nearest Neighbor

- Which class?



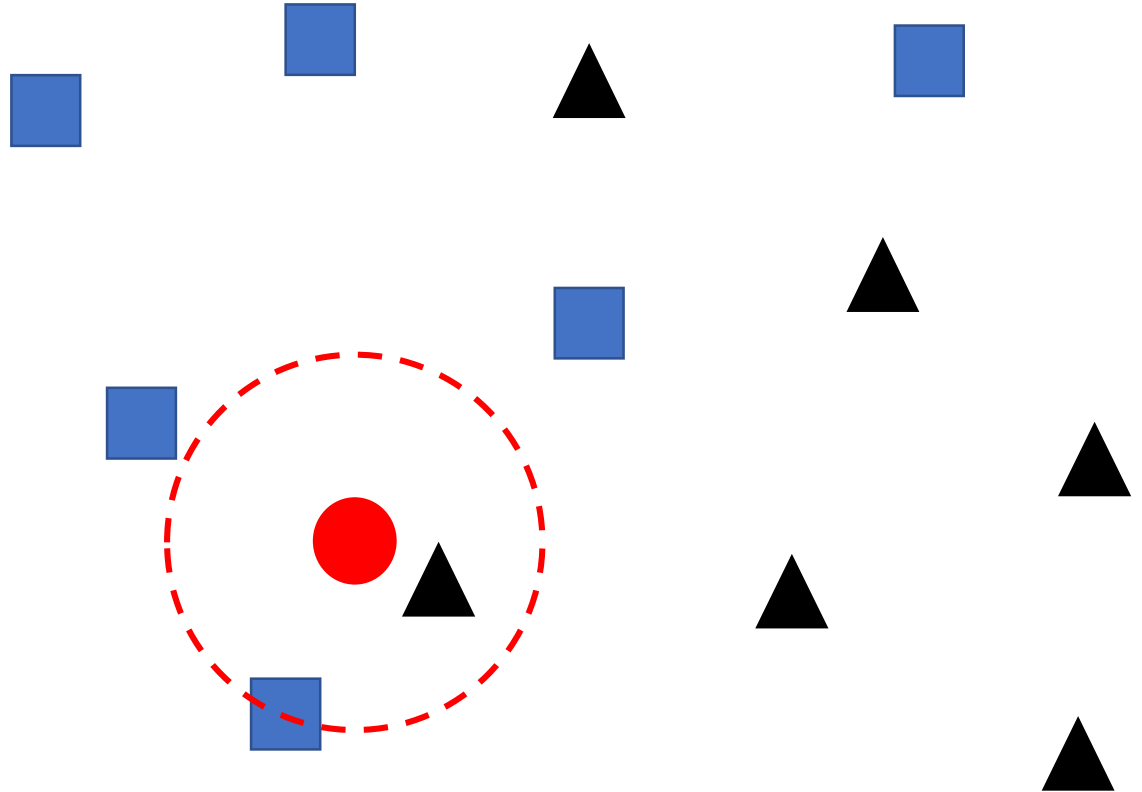
K-Nearest Neighbor

- $k=1$
- It is classified as a triangle



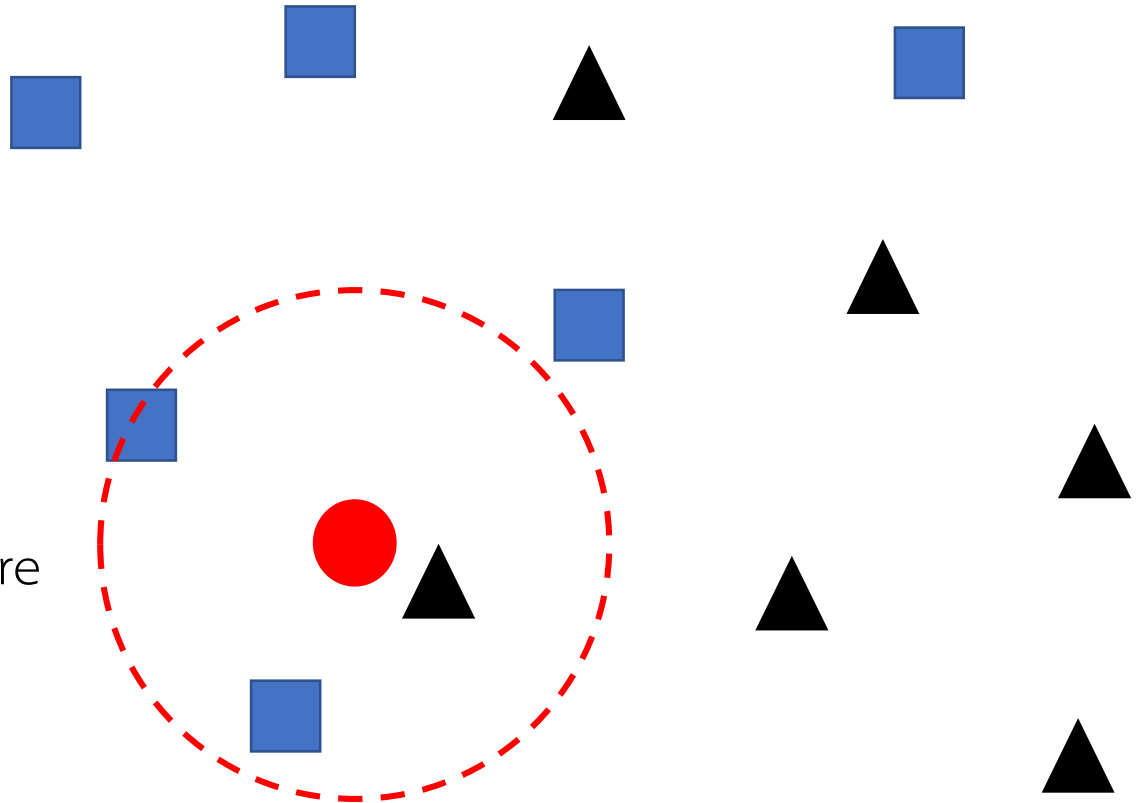
K-Nearest Neighbor

- $k=2$
- It is classified as ???



K-Nearest Neighbor

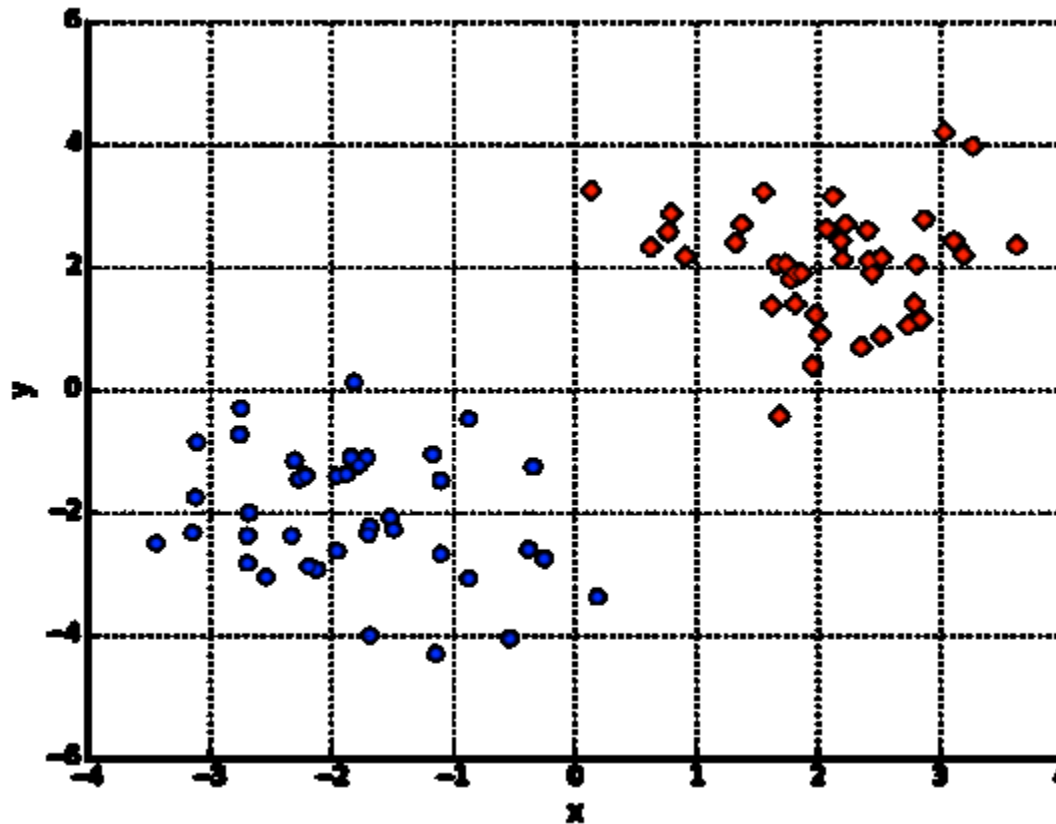
- $k=3$
- It is classified as a square



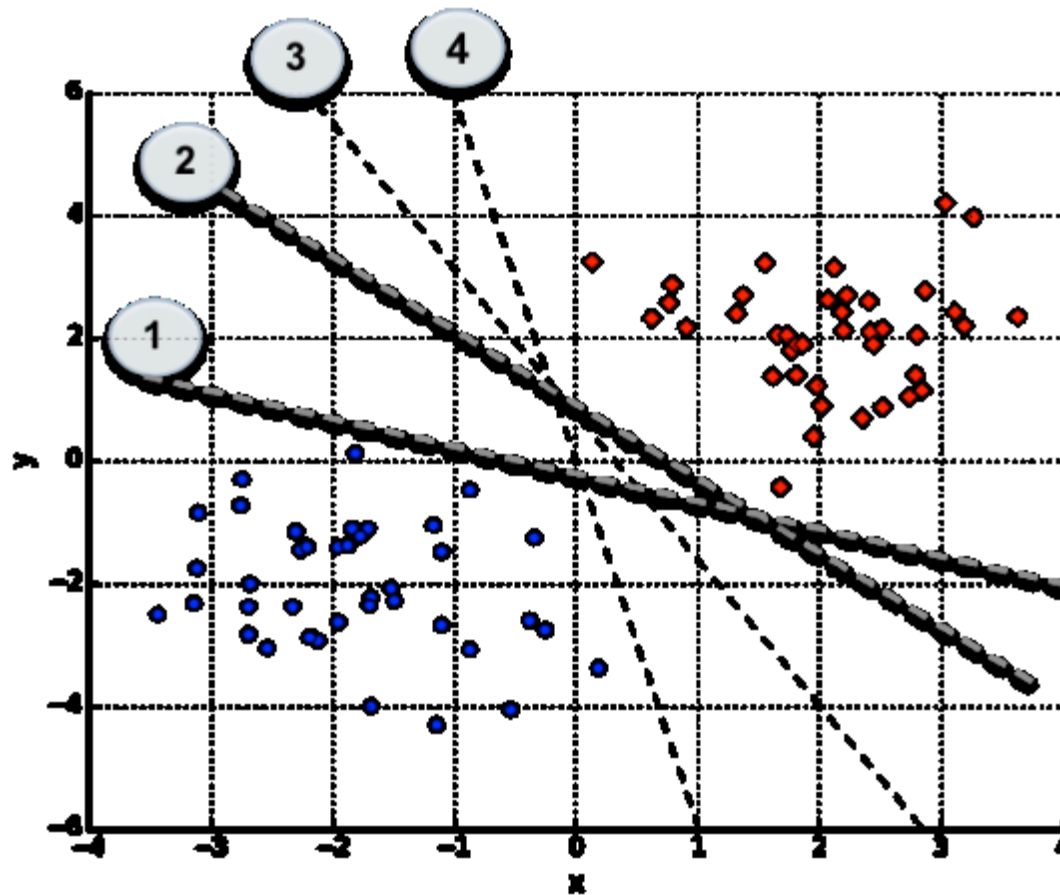
K-Nearest Neighbor

```
>>> X = [[0], [1], [2], [3]]
>>> y = [0, 0, 1, 1]
>>> from sklearn.neighbors import KNeighborsRegressor
>>> neigh = KNeighborsRegressor(n_neighbors=2)
>>> neigh.fit(X, y)
KNeighborsRegressor(...)
>>> print(neigh.predict([[1.5]]))
[ 0.5]
```

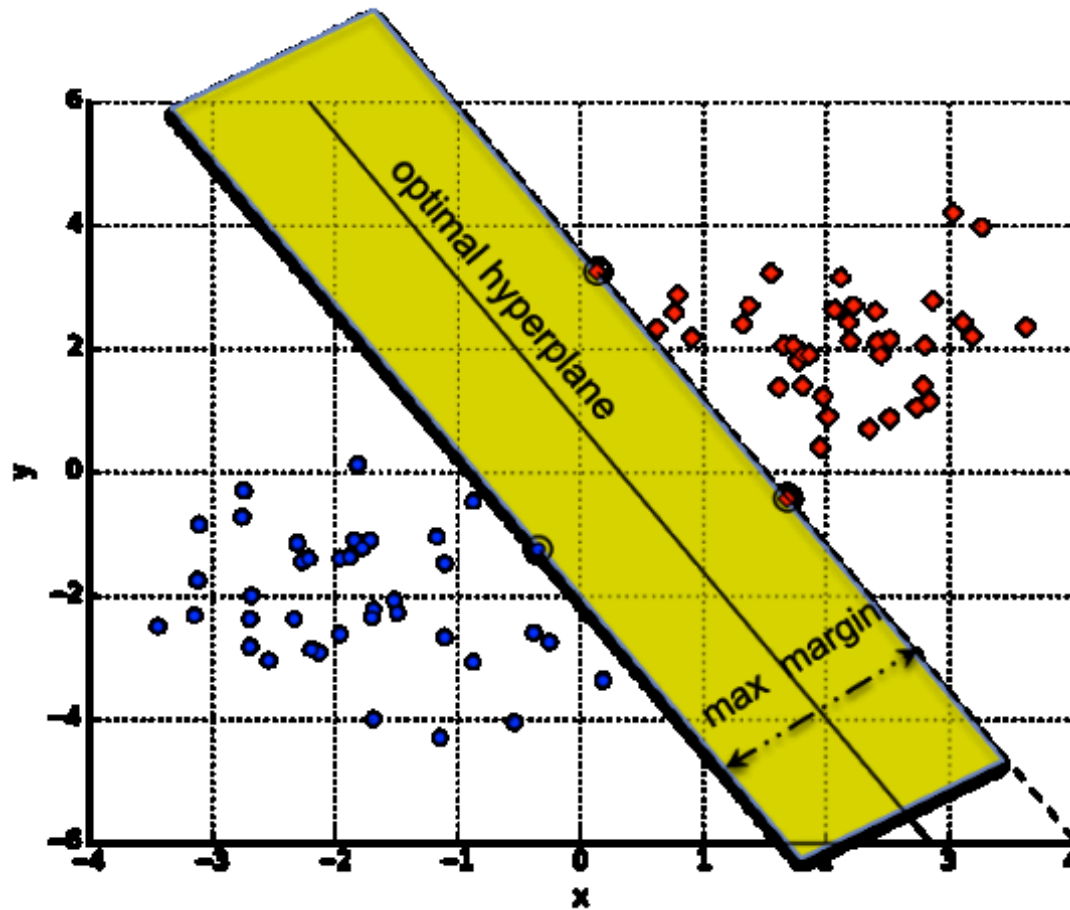
Support Vector Machines (SVM)



Support Vector Machines (SVM)



Support Vector Machines (SVM)



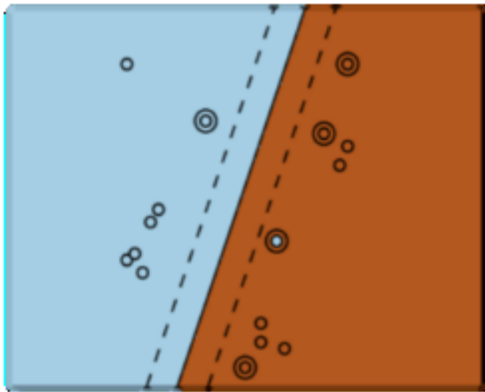
Support Vector Machines (SVM)

- Kernels

- **Gaussian Radial Basis Function (RBF):** $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$
- **Polynomial:** $K(x_i, x_j) = (\gamma x_i^T x_j + c)^d$
- **Sigmoid:** $K(x_i, x_j) = \tanh(\gamma x_i^T x_j + c)$

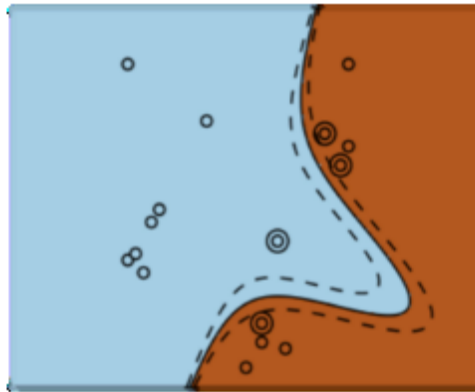
Support Vector Machines (SVM)

Linear Kernel



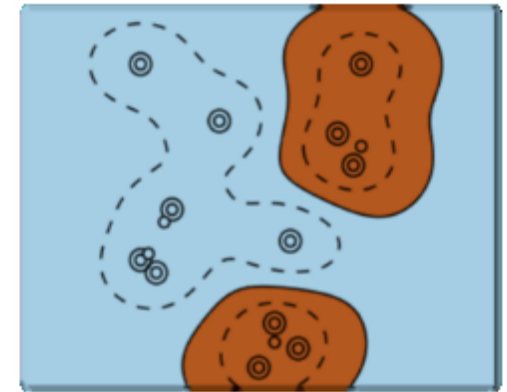
C hyperparameter

Polynomial Kernel



*C plus gamma, degree and
coefficient hyperparameters*

RBF Kernel



*C plus gamma
hyperparameter*

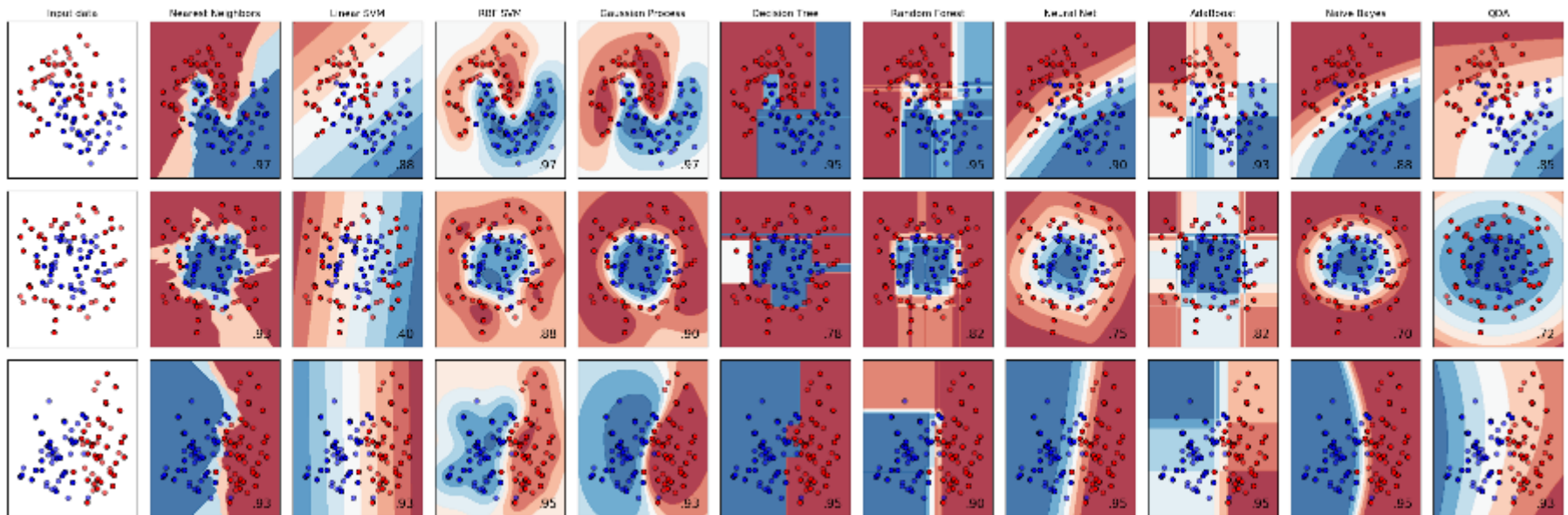
Support Vector Machines (SVM)

```
>>> from sklearn import svm
>>> X = [[0, 0], [1, 1]]
>>> y = [0, 1]
>>> clf = svm.SVC()
>>> clf.fit(X, y)
SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```

Python Code

- Go to

http://scikit-learn.org/stable/auto_examples/classification/plot_classifier_comparison.html



Regularization

Regularization

- The minimization

$$\min_f |Y_i - f(X_i)|^2$$

may be attained with zero errors.

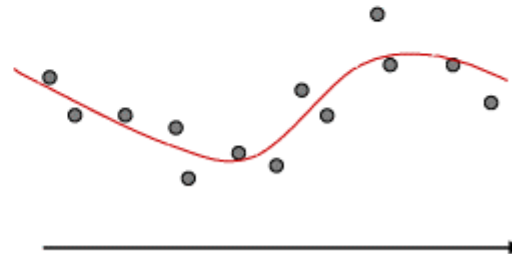
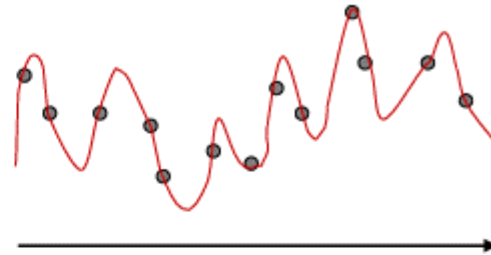
But the function may not be unique.



- Regularization

$$\min_{f \in H} \sum_{i=1}^n |Y_i - f(X_i)|^2 + \lambda \|f\|_H^2$$

- Regularization with smoothness penalty is preferred for uniqueness and smoothness.
- Link with some RKHS norm and smoothness



L1 vs. L2

L1 regularization on least squares:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 + \lambda \sum_{i=1}^k |w_i|$$

L2 regularization on least squares:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 + \lambda \sum_{i=1}^k w_i^2$$

L2 regularization	L1 regularization
Computational efficient due to having analytical solutions	Computational inefficient on non-sparse cases
Non-sparse outputs	Sparse outputs
No feature selection	Built-in feature selection

“

In the spirit of science, there really is
no such thing as a 'failed experiment.'
Any test that yields valid data
is a valid test.

”

Adam Savage