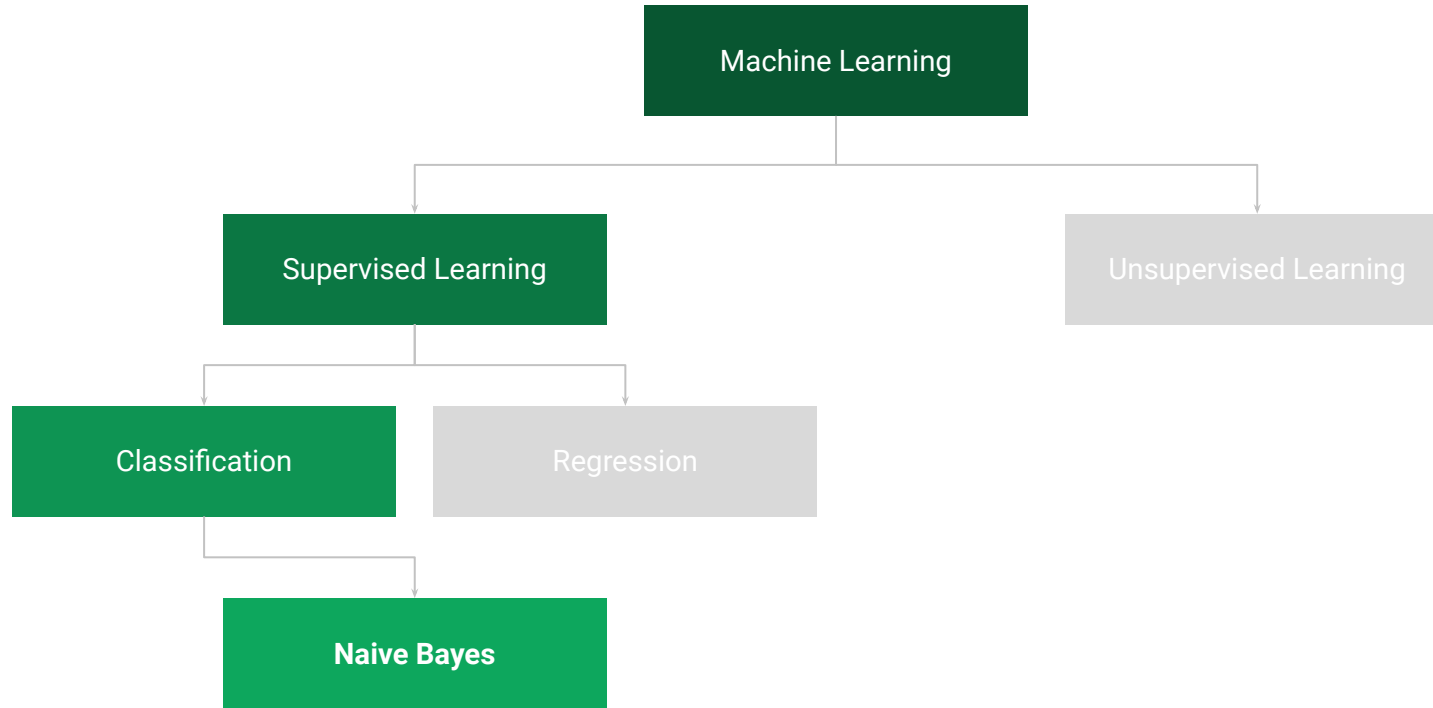


# Naive Bayes



# What is Naive Bayes?

A classification method that works on the principles of conditional probability as given by the Bayes' theorem:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$



Thomas Bayes  
1702 - 1761

# Why “Naive”?

Because it assumes that features are independent of each other.



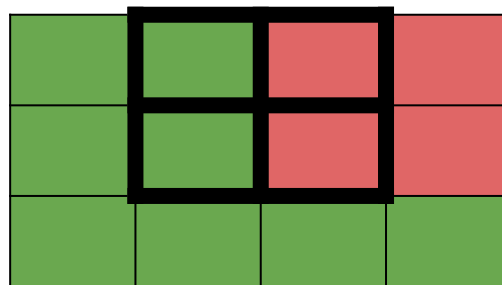
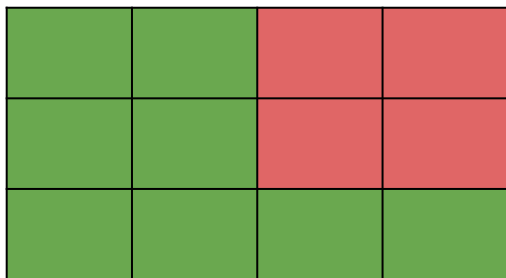
# Probabilities - Marginal Probability

- Tweets dataset (12) labelled either: positive (8 in green) or negative (4 in red)
- **Scenario A:** probability of one event happening,  $P(A)$  or  $P(B)$
- $P(\text{positive}) = N(\text{positive})/N(\text{total}) = 8/12 = 0.66$
- $P(\text{negative}) = 1 - P(\text{positive}) = 0.34$

**Dataset**

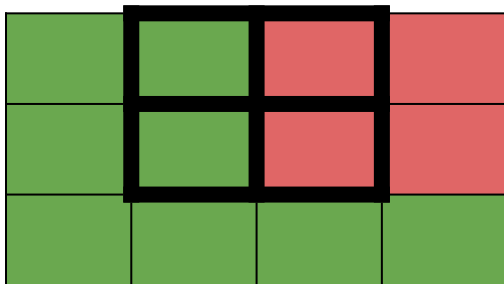

# Probabilities - Joint Probability

- **Scenario B:** probability of two (or more) events happening at the same time,  $P(A, B)$  or  $P(B, A)$
- I.e. Tweets that are positive and contain the word “happy”
- Overall we have 4 tweets containing the word “happy”, 2 of which are positive tweets
- $P(\text{“happy”}, \text{positive}) = P(\text{“happy”} \cap \text{positive}) = 2/20 = 0.1$



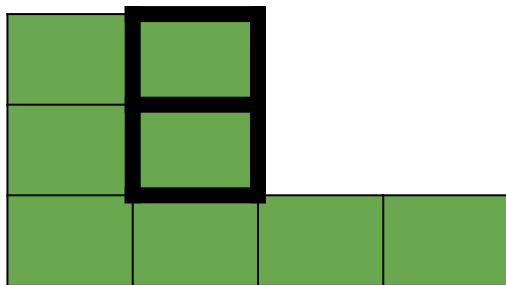
# Probabilities - Conditional Probability

- Wikipedia: “Conditional probability is a measure of the probability of an event occurring, given that another event has already occurred”
- **Scenario C:** probability of one (or more) events happening given an occurrence of another event,  $P(A \text{ given } B)$  or  $P(A | B)$
- I.e. Tweets that are positive given it contains the word “happy”
- Use joint probability
- $P(\text{positive} | \text{“happy”}) = P(\text{positive} \cap \text{“happy”}) / P(\text{“happy”}) = 2/4 = 0.5$



# Probabilities - Conditional Probability

- $P(\text{"happy"} \mid \text{positive}) = P(\text{"happy"} \cap \text{positive}) / P(\text{positive}) = 2/8 = 0.25$



- Note:  
 $P(A \mid B) \neq P(B \mid A)$ , but  $P(A \cap B) = P(B \cap A)$



# Introducing Bayes' Rule

- We saw that:
  - $P(\text{"happy"} \mid \text{positive}) = P(\text{"happy"} \cap \text{positive}) / P(\text{positive})$
  - $P(\text{positive} \mid \text{"happy"}) = P(\text{positive} \cap \text{"happy"}) / P(\text{"happy"})$
- Now we can derive Bayes' rule:
  - $P(\text{positive} \mid \text{"happy"}) = P(\text{"happy"} \mid \text{positive}) * P(\text{positive}) / P(\text{"happy"})$

# Bayes' Rule

In general:

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

- $P(X|Y) \Rightarrow$  **Posterior Probability**; Probability of X given Value of Y.
- $P(Y|X) \Rightarrow$  **Likelihood** of Y given X is True.
- $P(X) \Rightarrow$  **Prior Probability**; Probability of event X.
- $P(Y) \Rightarrow$  **Evidence**; Probability of event Y.

Posterior = Likelihood \* Prior / Evidence

# Bayes' Rule

What is the probability that there is fire given that there is smoke,  $P(\text{Fire}|\text{Smoke})$ ?

- Where:
  - $P(\text{Fire})$  is the Prior,
  - $P(\text{Smoke}|\text{Fire})$  is the Likelihood, and
  - $P(\text{Smoke})$  is the evidence:

$$P(\text{Fire}|\text{Smoke}) = P(\text{Smoke}|\text{Fire}) * P(\text{Fire}) / P(\text{Smoke})$$

# !!!QUESTION ALERT!!!

What is the probability that the sentence “**ML Bootcamp is awesome!**” is **positive**?

Given that:

40% of positive tweets contain the word “**awesome**”,

A total of 30% of the tweet contain the word “**awesome**”, and

60% of the total number of tweets are **positive**

1.  $P(\text{positive} \mid \text{“awesome”}) = 0.2$
2.  $P(\text{positive} \mid \text{“awesome”}) = 0.8$
3.  $P(\text{positive} \mid \text{“awesome”}) = 0.5$
4.  $P(\text{positive} \mid \text{“awesome”}) = 0.9$

# Naive Bayes

- Real-world:  $P(C \mid A_1, A_2 \dots A_n)$

$$P(C \mid A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n \mid C) P(C)}{P(A_1 A_2 \dots A_n)}$$

- Choose value of C that maximizes the top expression
- Basically:  $P(C \mid A_1 A_2 \dots A_n) = P(A_1, A_2, \dots, A_n \mid C) * P(C)$

# Naive Bayes

How do we estimate  $P(A_1, A_2, \dots, A_n | C)$  ?

- Assume independence for each A given C:

$$P(A_1, A_2, \dots, A_n | C) = P(A_1 | C) * P(A_2 | C) \dots P(A_n | C)$$

- We can estimate  $P(A_i | C)$

# Naive Bayes

- Under the previous independence assumptions

$$p(C_k \mid x_1, \dots, x_n) = \frac{1}{Z} p(C_k) \prod_{i=1}^n p(x_i \mid C_k)$$

- Where  $Z$  is the evidence,  $P(x)$ , and dependent only on  $x_1, \dots, x_n$  which is a constant if feature variables are known

# Naive Bayes Classifier

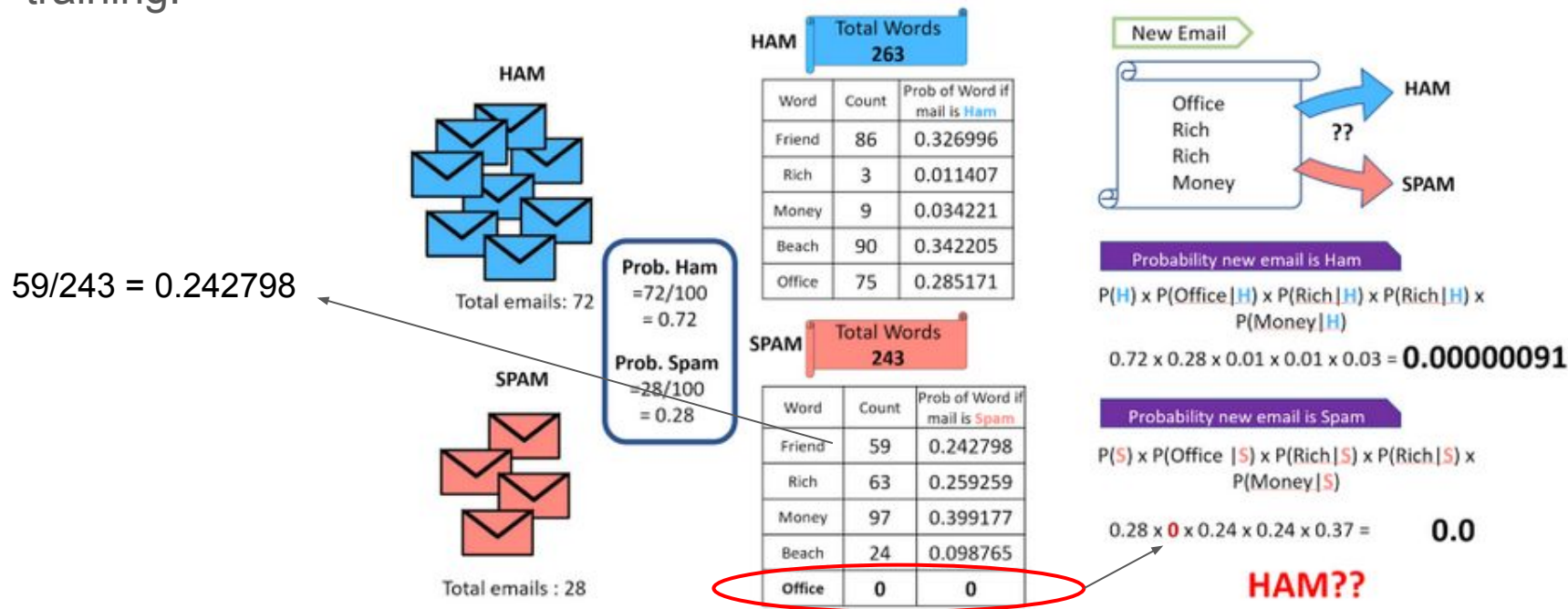
- With the previous model combined with a decision rule (*maximum a posteriori* or *MAP*) we get our classifier:

$$\hat{y} = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(C_k) \prod_{i=1}^n p(x_i \mid C_k)$$



# The Zero-Probability Problem!

- When an instance in test dataset has a category that was not present during training.



# The Zero-Probability Problem!

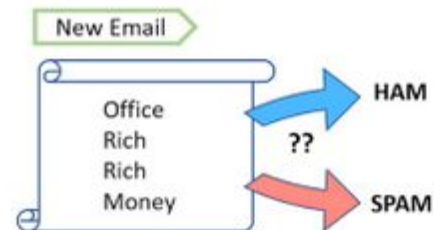
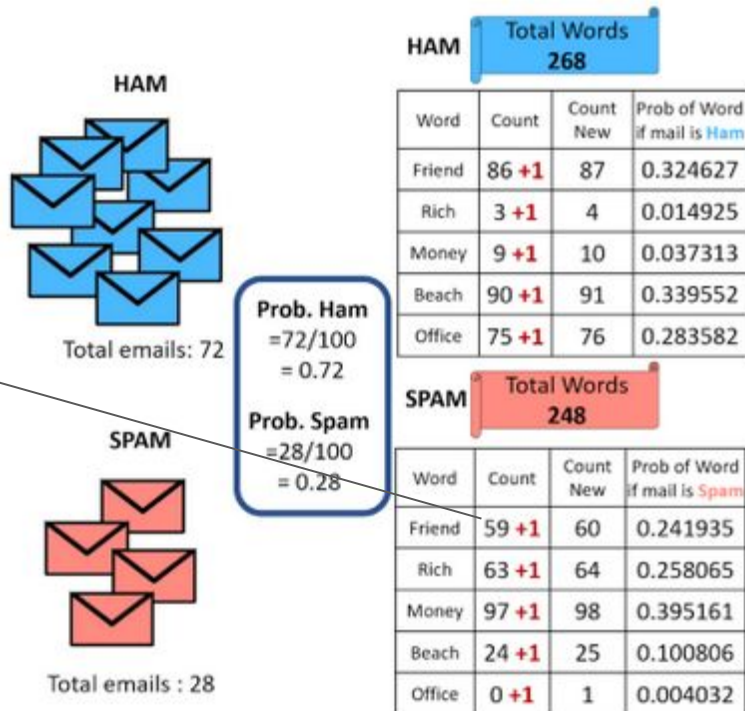
- Solution: Laplace smoothing
- Add one to the count for every attribute value-class combination
- Without smoothing:

$$P(A_i = a \mid C = c) = \frac{N_{ac}}{N_c}$$

- With smoothing:

$$P(A_i = a \mid C = c) = \frac{N_{ac} + 1}{N_c + N_i}$$

# The Zero-Probability Problem!



Probability new email is Ham

$$P(H) \times P(\text{Office}|H) \times P(\text{Rich}|H) \times P(\text{Rich}|H) \times P(\text{Money}|H)$$

$$0.72 \times 0.28 \times 0.01 \times 0.01 \times 0.04 = 0.0000017$$

Probability new email is Spam

$$P(S) \times P(\text{Office}|S) \times P(\text{Rich}|S) \times P(\text{Rich}|S) \times P(\text{Money}|S)$$

$$0.28 \times 0.004 \times 0.25 \times 0.25 \times 0.39 = 0.000029$$

**SPAM!!!**

# Vanishing Value Problem!

- Happens when multiplying many very small numbers
- Number get very small and close to zero => numerical underflow
- Can fix by computing the logarithm of the conditional probability

$$\begin{aligned}\log P(C|A) &\sim \log P(A|C) + \log P(A) \\ &= \sum_i \log(A_i|C) + \log P(A)\end{aligned}$$

# Handling Continuous Attributes

- One way is to discretize each continuous attribute and then replace the continuous attribute value with its corresponding discrete interval
- Another way is by assuming a certain form of probability distribution for the continuous variables and estimate the parameters of the distribution using the training data:

$$P(X|Y = c) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{\frac{-(x-\mu_c)^2}{2\sigma_c^2}}$$

# Naive Bayes Classifiers

- Multinomial Naive Bayes
  - Assumes discrete features
- Bernoulli Naive Bayes
  - Assumes binary features
- Gaussian Naive Bayes
  - Assumes continuous features

# Advantages

- Relatively simple to understand and implement
- Fast to train and classify compared to other models
- Can perform better than other models on a small dataset
- Not sensitive to irrelevant features
- Handles both continuous and discrete features
- Can be used for binary or multi-class problems
- Works well with high dimensions



# Disadvantages

- Naive! assumes all features are independent
- Treats all attributes equally
- Vanishing value
- Zero probability problem
- Smoothing is an over-head
- Continuous value attribute problem





# Trending Applications

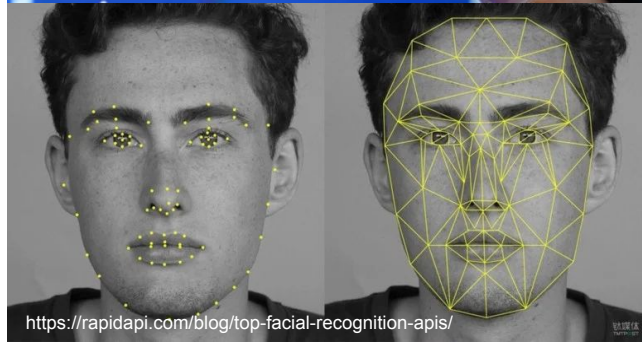
- Text classification, i.e.:
  - Spam filtering
  - Sentiment analysis
- Medicine, i.e.:
  - Disease detection based on genome wide data
- Image classification, i.e.:
  - Face recognition



<https://towardsdatascience.com/spam-filtering-using-naive-bayes-98a341224038>



<https://www.genengnews.com/insights/cytogenomic-analyses-for-genetic-disease-detection-and-diagnosis/>



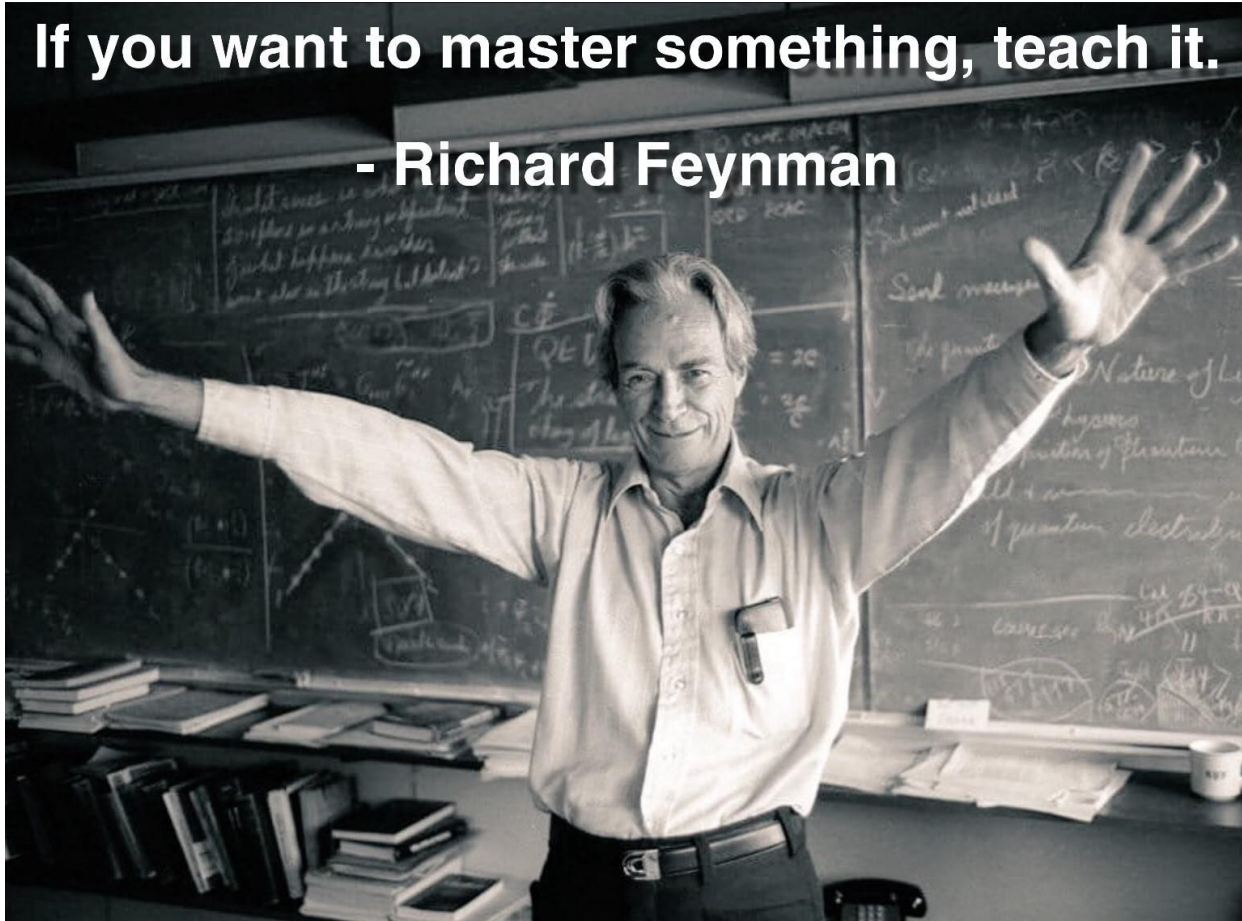
<https://rapidapi.com/blog/top-facial-recognition-apis/>

# Support Vector Machines

Hayden Fiege E.I.T.

**If you want to master something, teach it.**

**- Richard Feynman**



# What is SVM?

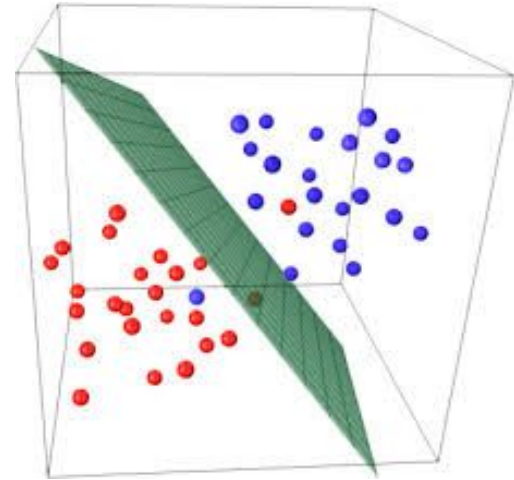
- A powerful and versatile ML model developed by Vladimir Vapnik and colleagues in the early 90's. SVM can be used for linear and non-linear classification and regression, for outlier detection, for clustering and more.

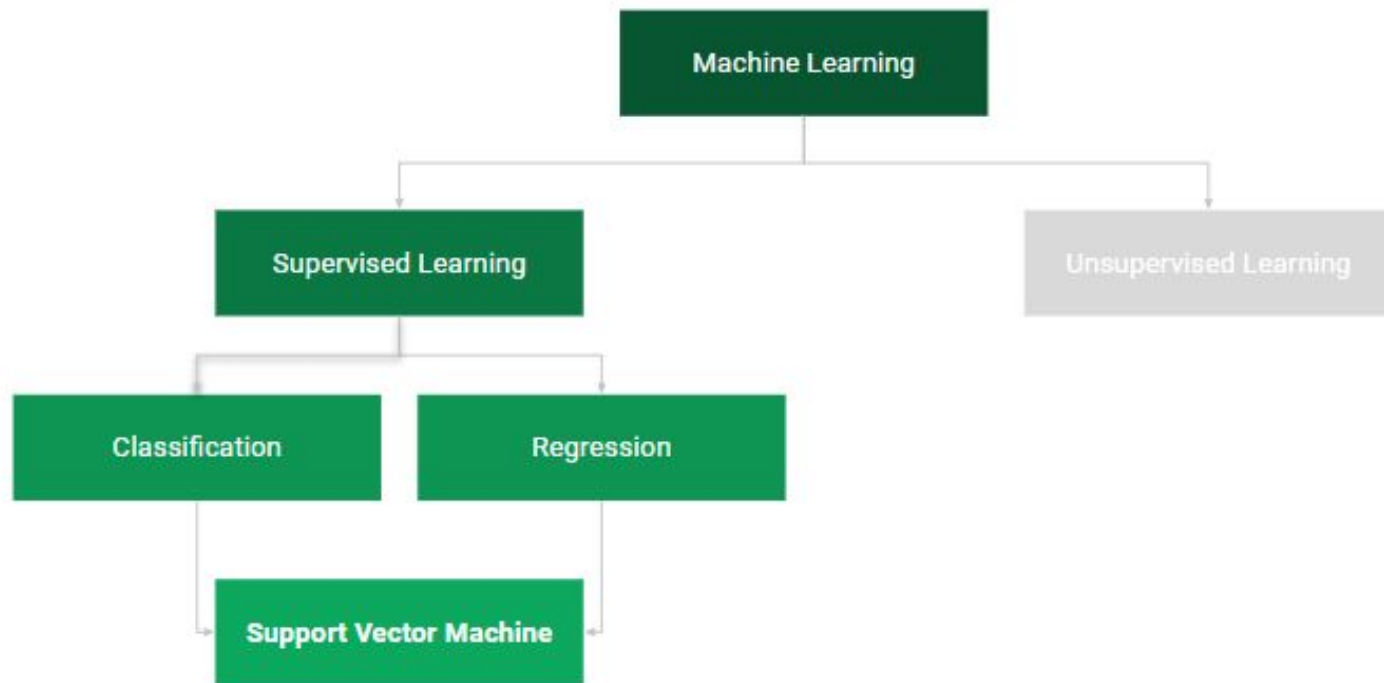
$$\min_{\beta^{\wedge}, \beta^{\vee}} \frac{1}{2} \sum_{i=1, j=1}^m (\beta_i^{\wedge} - \beta_i^{\vee}) (\beta_j^{\wedge} - \beta_j^{\vee}) K_{ij} + \sum_{i=1}^m (\varepsilon - y^{(i)}) \beta_i^{\wedge} + (\varepsilon + y^{(i)}) \beta_i^{\vee}$$

$$s.t : \sum_{i=1}^m (\beta_i^{\wedge} - \beta_i^{\vee}) = 0$$

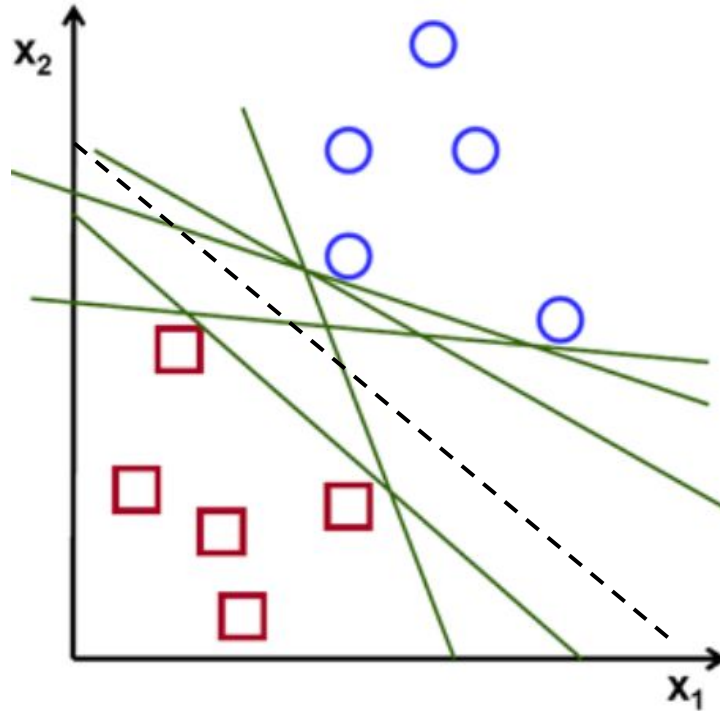
$$0 < \beta_i^{\wedge} < C, i = 1, 2, \dots, m$$

$$0 < \beta_i^{\vee} < C, i = 1, 2, \dots, m$$

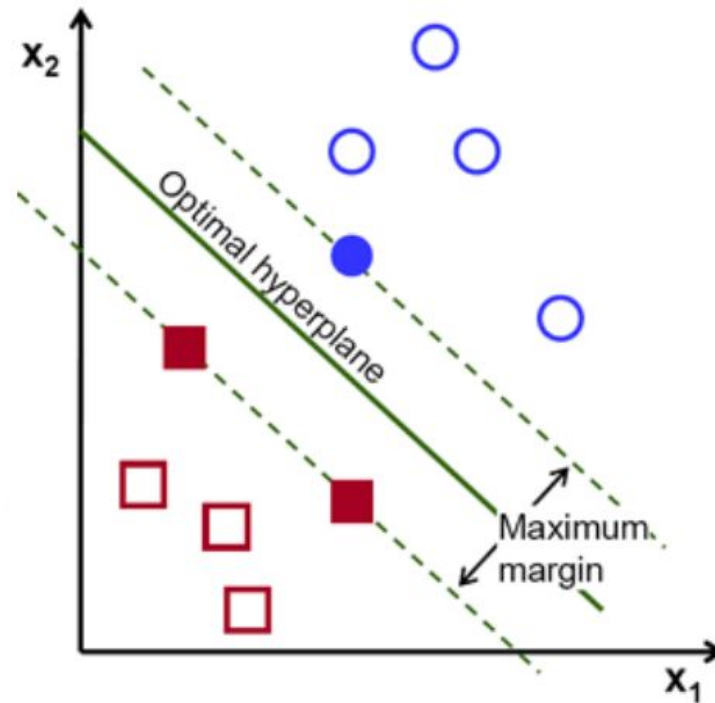




# Linear SVM – Multiple Lines

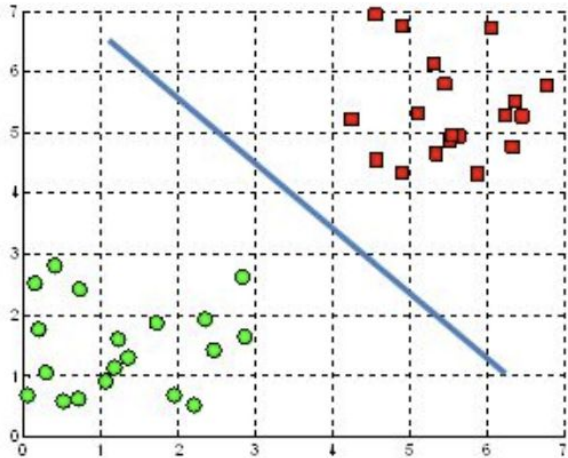


# Linear SVM – Optimal Hyperplane

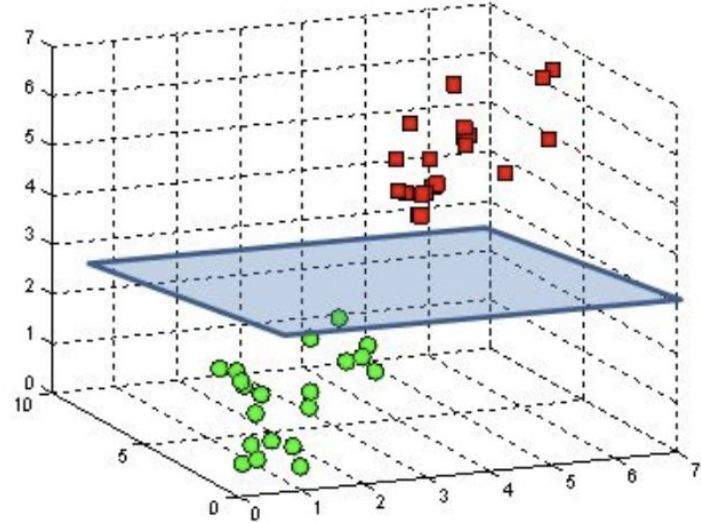


# What is a hyperplane?

A hyperplane in  $\mathbb{R}^2$  is a line



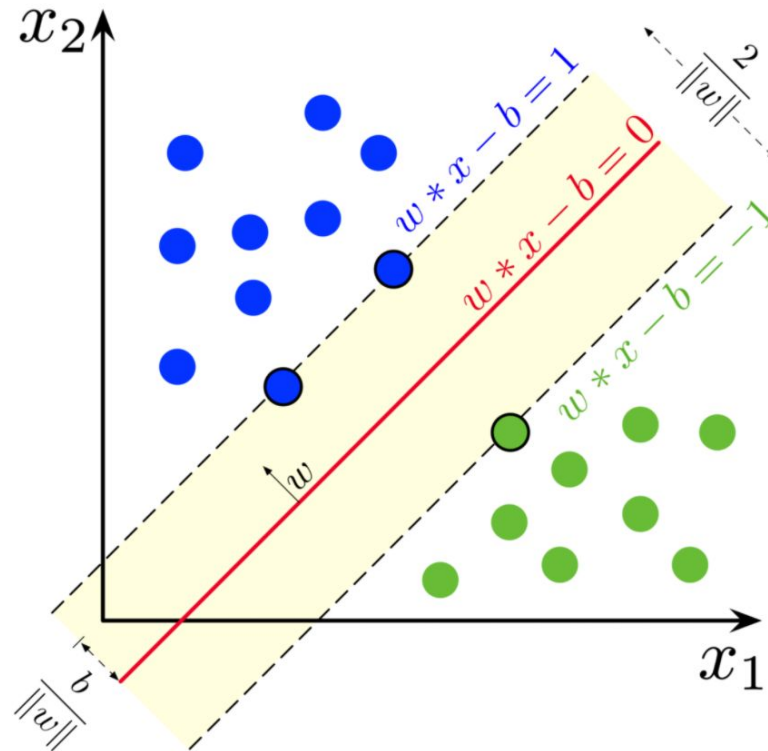
A hyperplane in  $\mathbb{R}^3$  is a plane



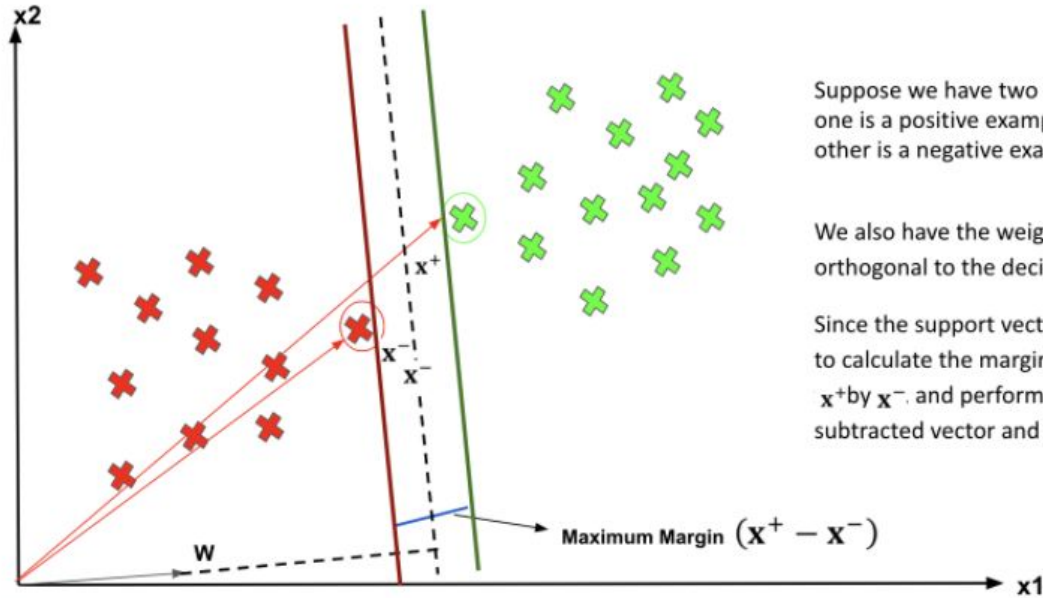
*Hyper plane equation*  $\rightarrow w \cdot x + b = 0$



# Linear SVM



# Max Margin



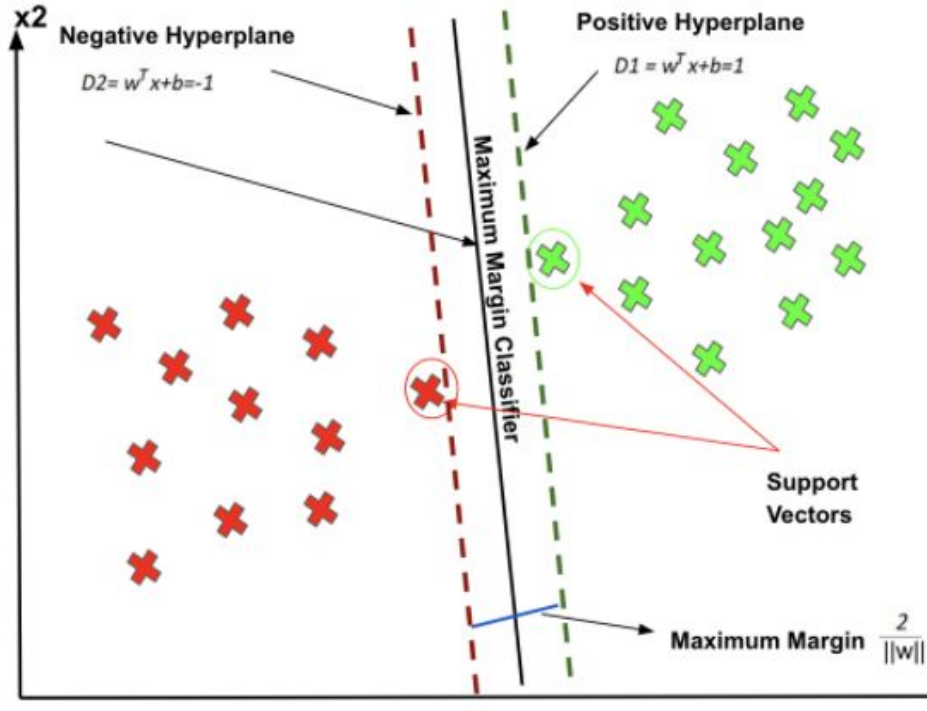
Suppose we have two support vectors where one is a positive example (e.g.,  $x^+$ ) and the other is a negative example (e.g.,  $x^-$ )

We also have the weight vector  $w$  which is orthogonal to the decision boundary

Since the support vectors define the margin, to calculate the margin, we can first subtract  $x^+$  by  $x^-$  and perform dot product on the subtracted vector and the unit vector of  $w$

$$(x^+ - x^-) \cdot \hat{w} = (x^+ - x^-) \cdot \frac{w}{\|w\|} = x^+ \cdot \frac{w}{\|w\|} - x^- \cdot \frac{w}{\|w\|}$$

# Max Margin



$$D1 = w^T x + b = 1 \quad w^T x + b - 1 = 0$$

$$D2 = w^T x + b = -1 \quad w^T x + b + 1 = 0$$

$$w^T x + b - 1 - (w^T x + b + 1)$$

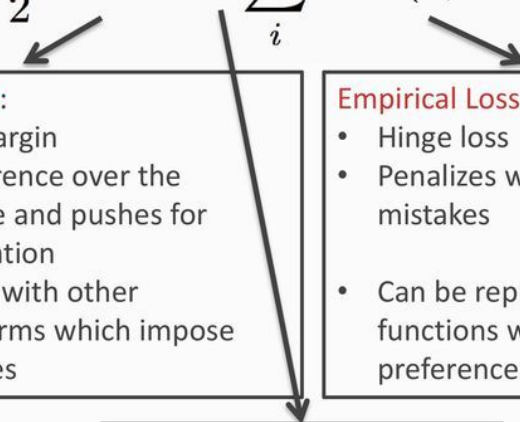


Solve algebraically

$$\frac{2}{\|w\|}$$

so to increase the *margin*,  
minimize the  $\frac{1}{2} \|w\|^2$

# SVM objective function

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$


## Regularization term:

- Maximize the margin
- Imposes a preference over the hypothesis space and pushes for better generalization
- Can be replaced with other regularization terms which impose other preferences

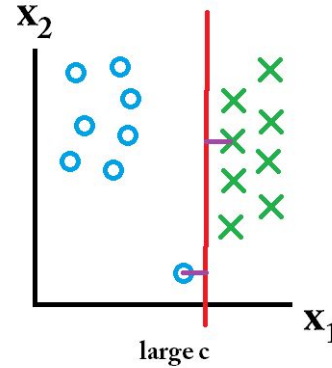
## Empirical Loss:

- Hinge loss
- Penalizes weight vectors that make mistakes
- Can be replaced with other loss functions which impose other preferences

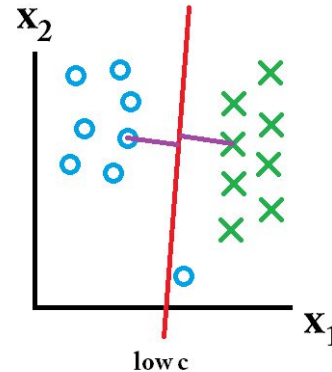
A **hyper-parameter** that controls the tradeoff between a large margin and a small hinge-loss

# C hyperparameter

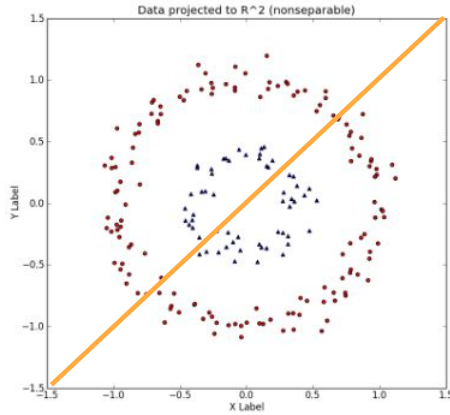
Large C: Lower bias, high  
variance  
i.e. More prone to over fitting



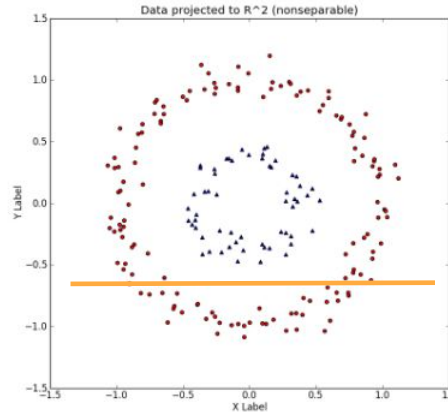
Small C: Higher bias, low  
variance  
i.e. More prone to under fitting



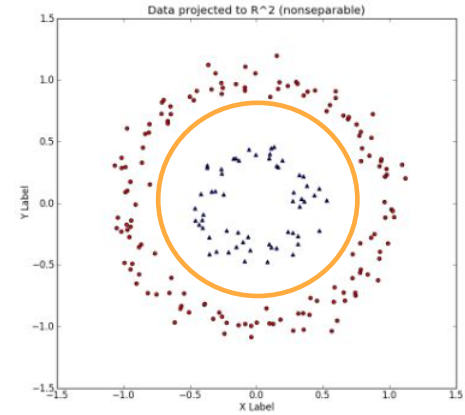
# How do we separate the data with SVM?



A

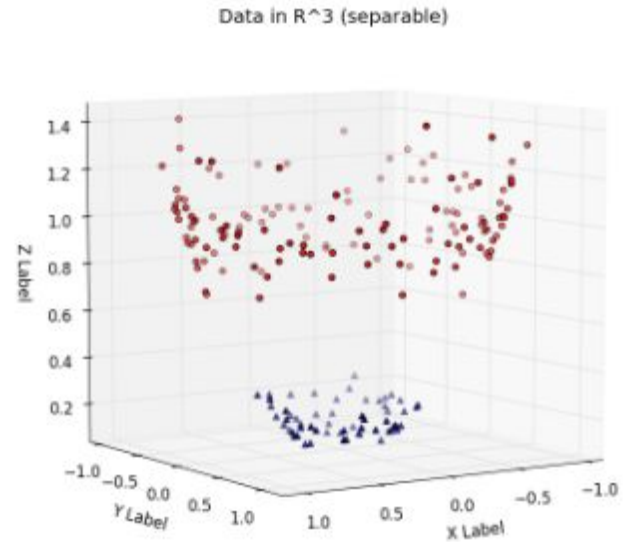
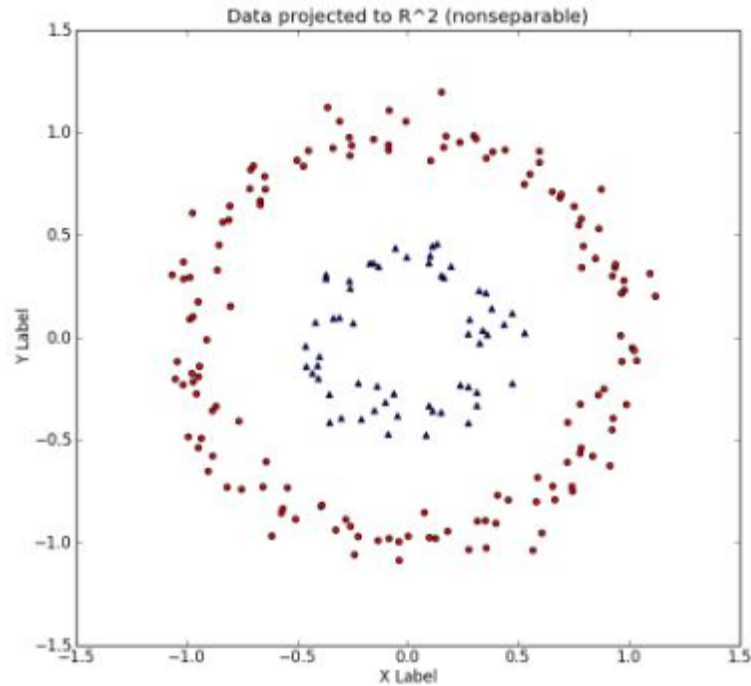


B

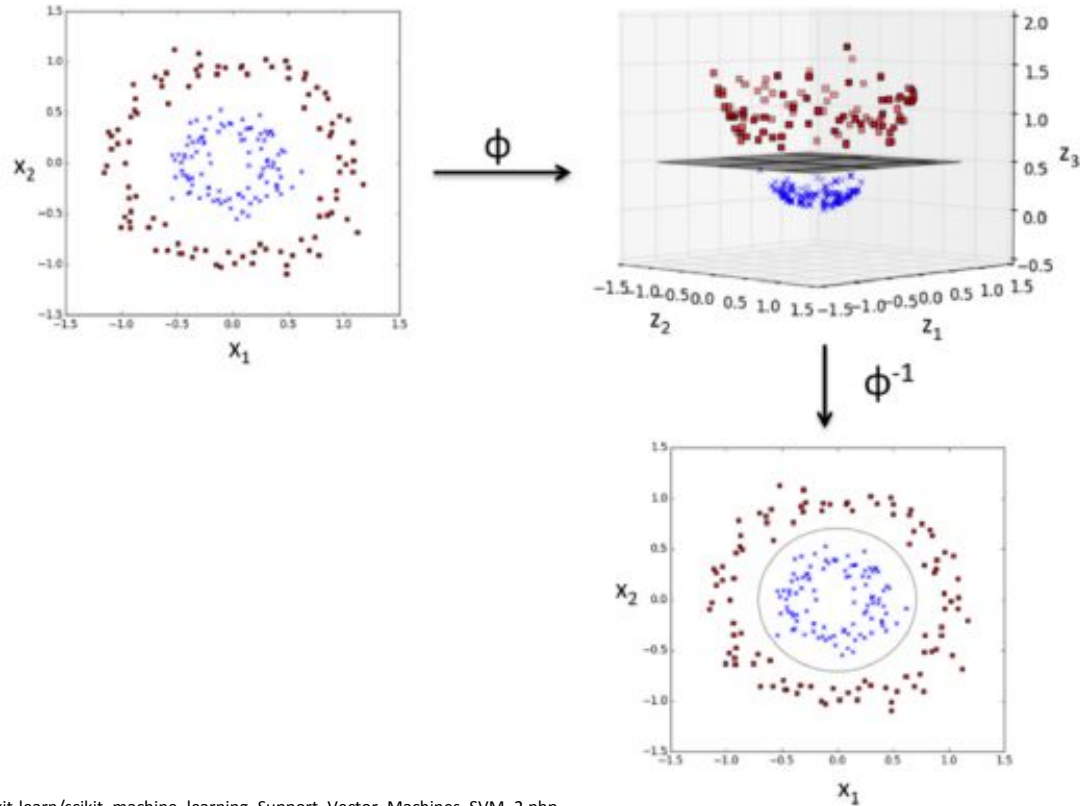


C

# Non-linear SVM



# Non-linear SVM





# Kernel Trick

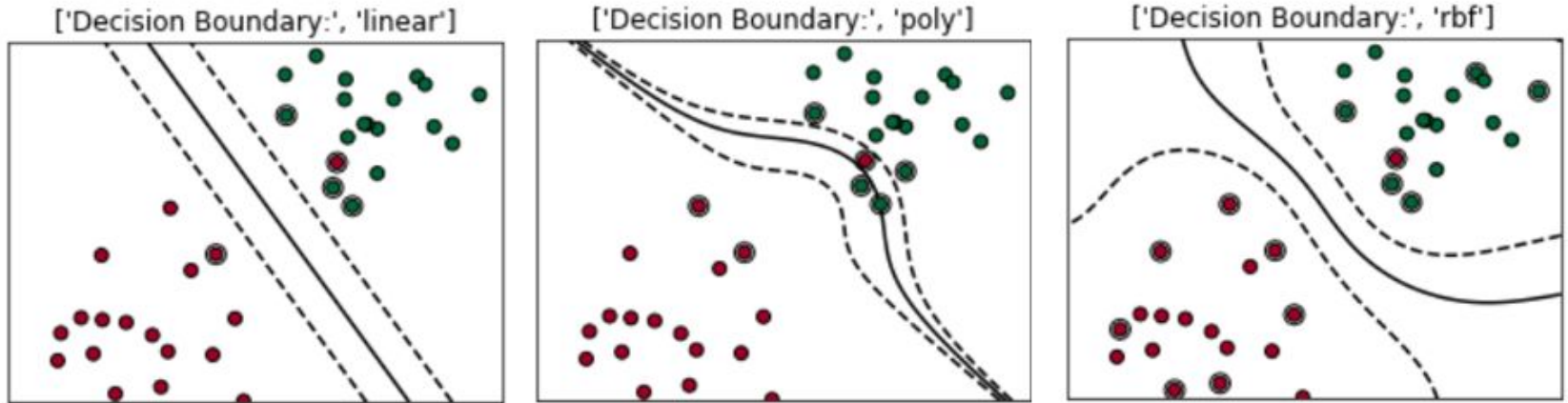
- It is a scientific term with the name trick in it (Exciting!)
- Provides mapping for linear learning algorithms to work on non-linear data
- The function must meet the definition provided in Mercers theorem to be a Kernel Function
- The functions employ mathematical tricks that allow the functions to operate in a higher-dimension without ever computing the coordinates of the data in that space

# Some of the Kernel Types

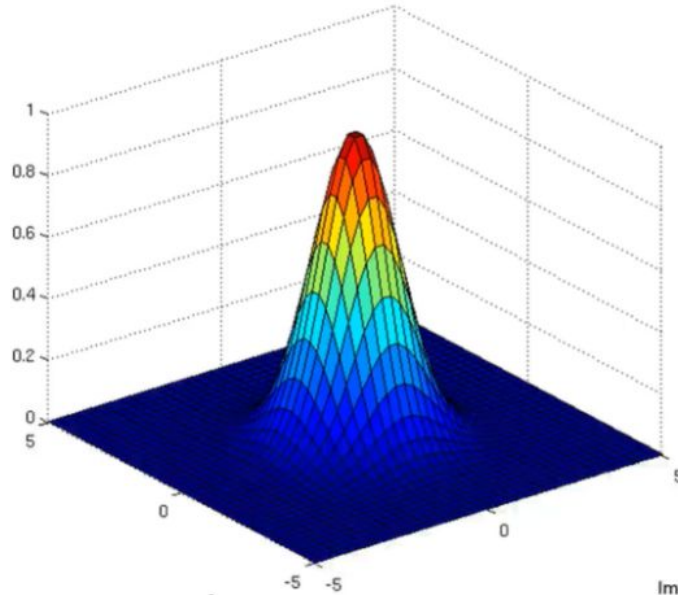
- polynomial:  $(\gamma \langle x, x' \rangle + r)^d$ , where  $d$  is specified by parameter `degree`,  $r$  by `coef0`.
- rbf:  $\exp(-\gamma \|x - x'\|^2)$ , where  $\gamma$  is specified by parameter `gamma`, must be greater than 0.
- sigmoid  $\tanh(\gamma \langle x, x' \rangle + r)$ , where  $r$  is specified by `coef0`.

Other kernels: String kernel, chi-square kernel, histogram intersection kernel, etc.

# Main Types of Kernels



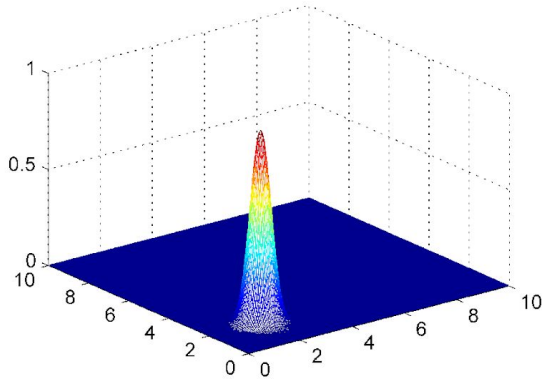
# Radial Basis Function (RBF) Kernel



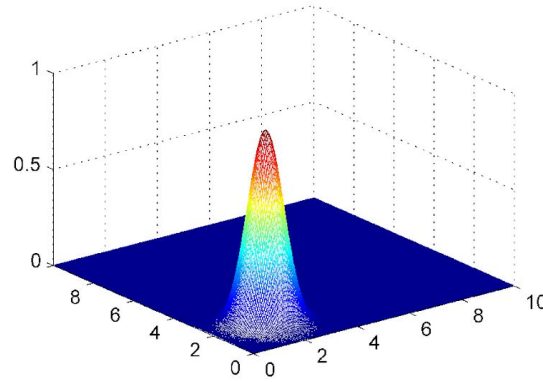
$$RBF: e^{-\gamma \|x-y\|^2}$$

Image source: <http://www.cs.toronto.edu/~duvenaud/cookbook/index.html>

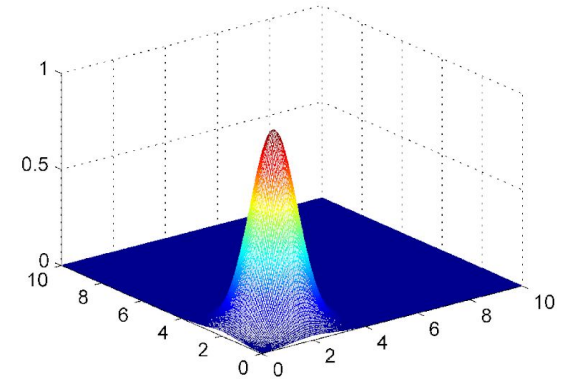
# Radial Basis Function (RBF) Kernel



$$\gamma = 12.5$$

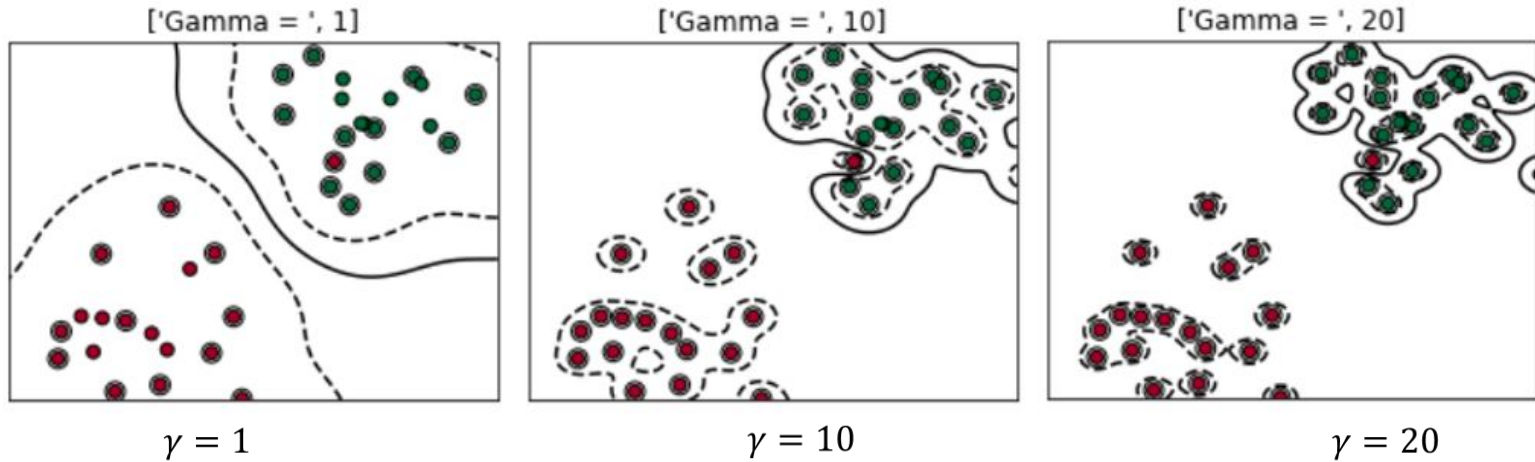


$$\gamma = 1.4$$

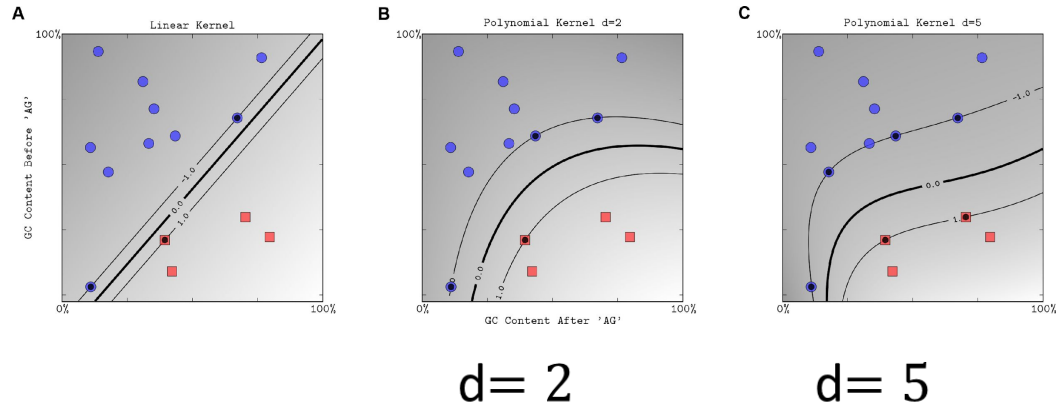


$$\gamma = 0.5$$

# Radial Basis Function (RBF) Kernel

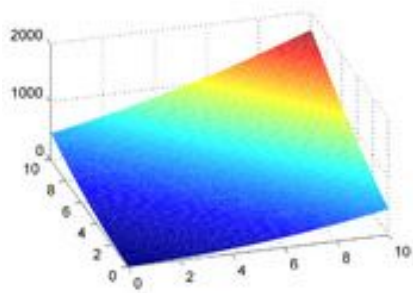


# Polynomial Kernel

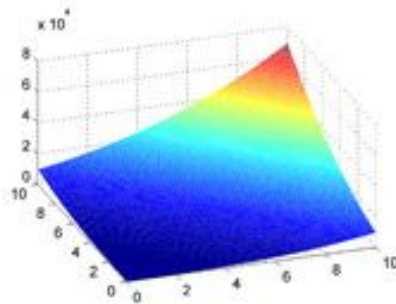


$$\text{Poly.} = (x \cdot y + r)^d$$

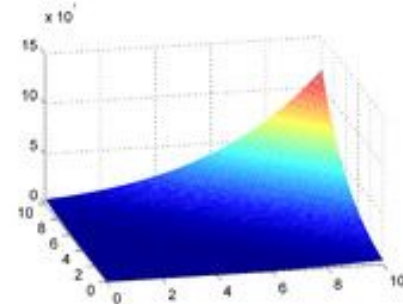
# Polynomial Kernel



$$d = 2$$



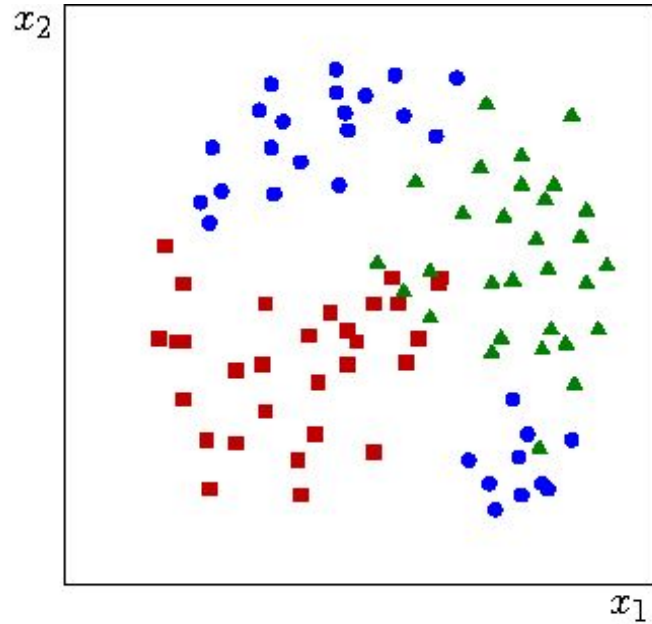
$$d = 3$$



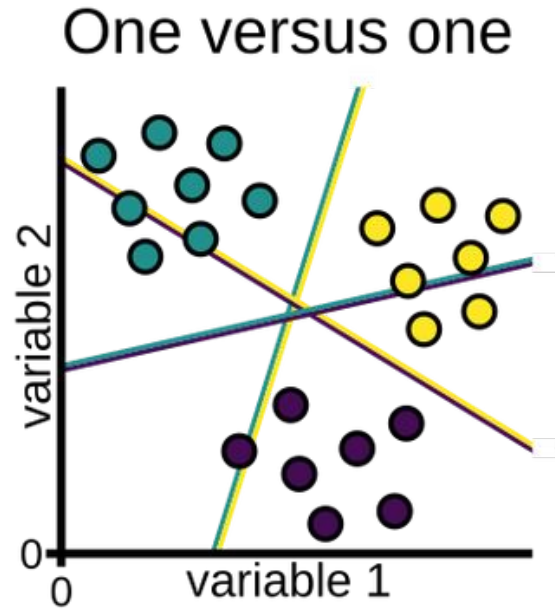
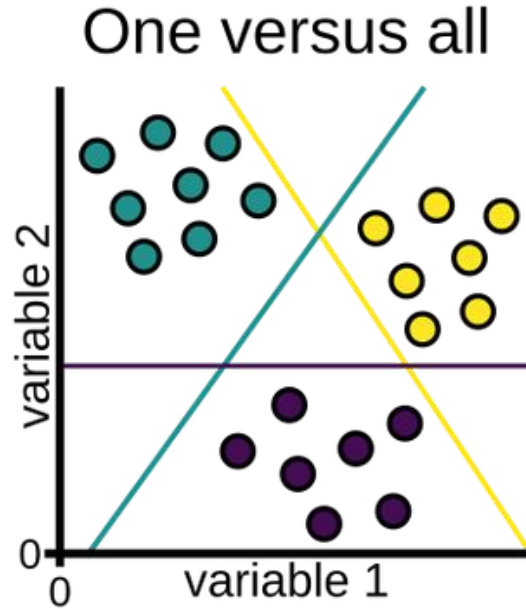
$$d = 5$$



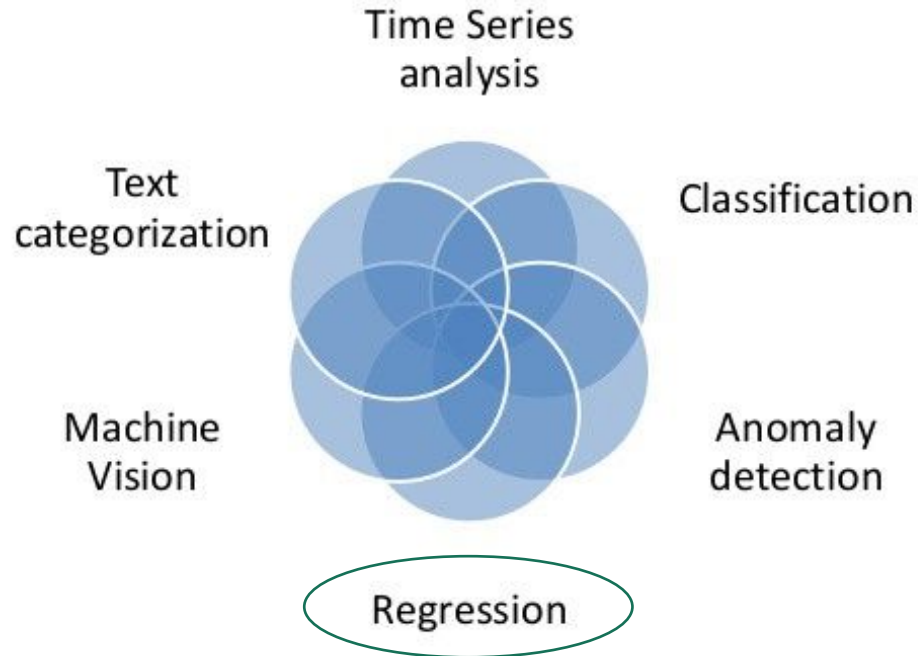
# Multi-class SVM



# Multi-class SVM



# Application of SVM



Thursday, August 7, 2014

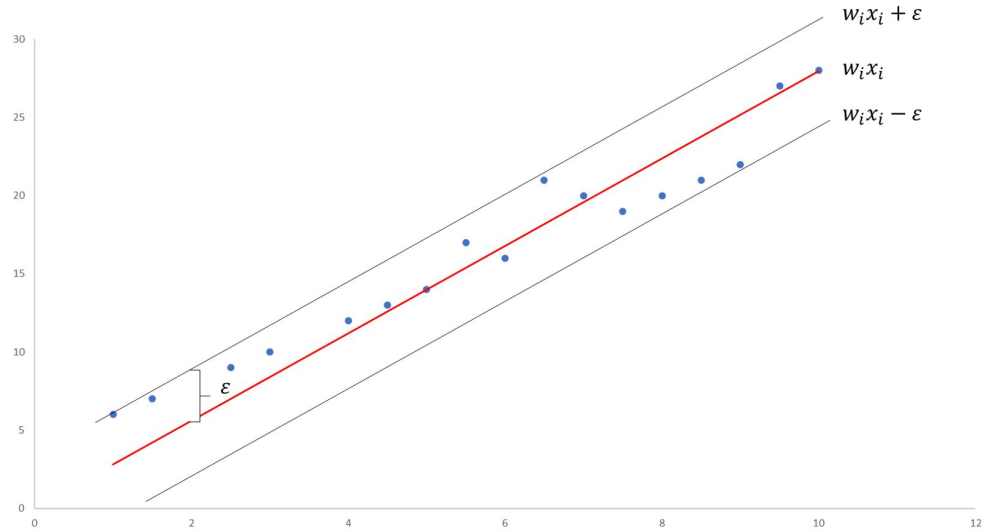
WITHOUT TEARS SERIES | [www.diggdata.in](http://www.diggdata.in)

4

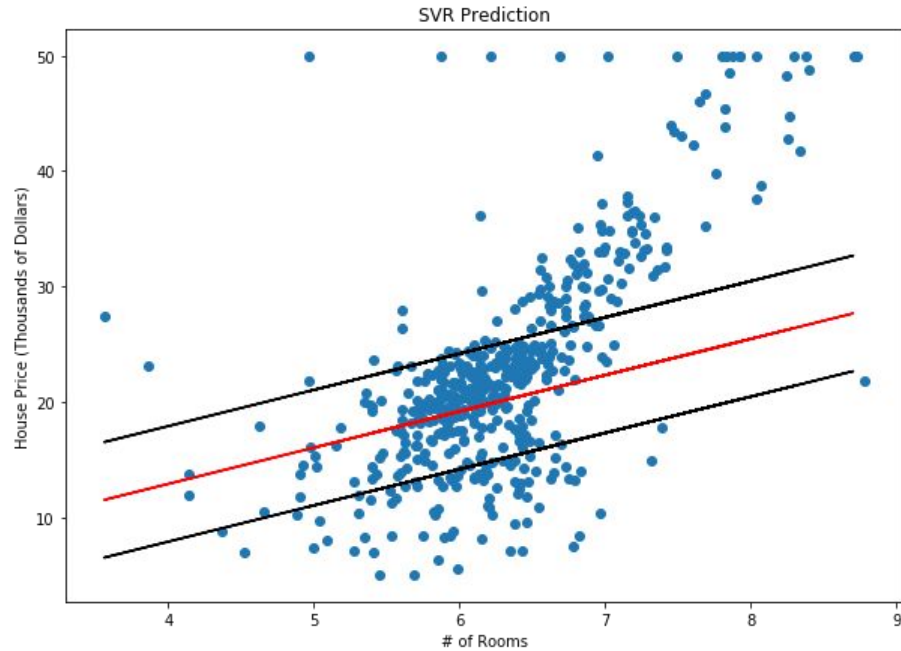
# Support Vector Regression (SVR)

- Works by changing the objective function
- We want to fit as many points on the street as possible
- Scikit-learn's `sklearn.svm.SVR` class is used for linear and non-linear SVM Regression
- Some of our inputs are kernel type, C, epsilon, and gamma for RBF, Polynomial, and Sigmoid kernels
- SVR does not scale well to large data sets

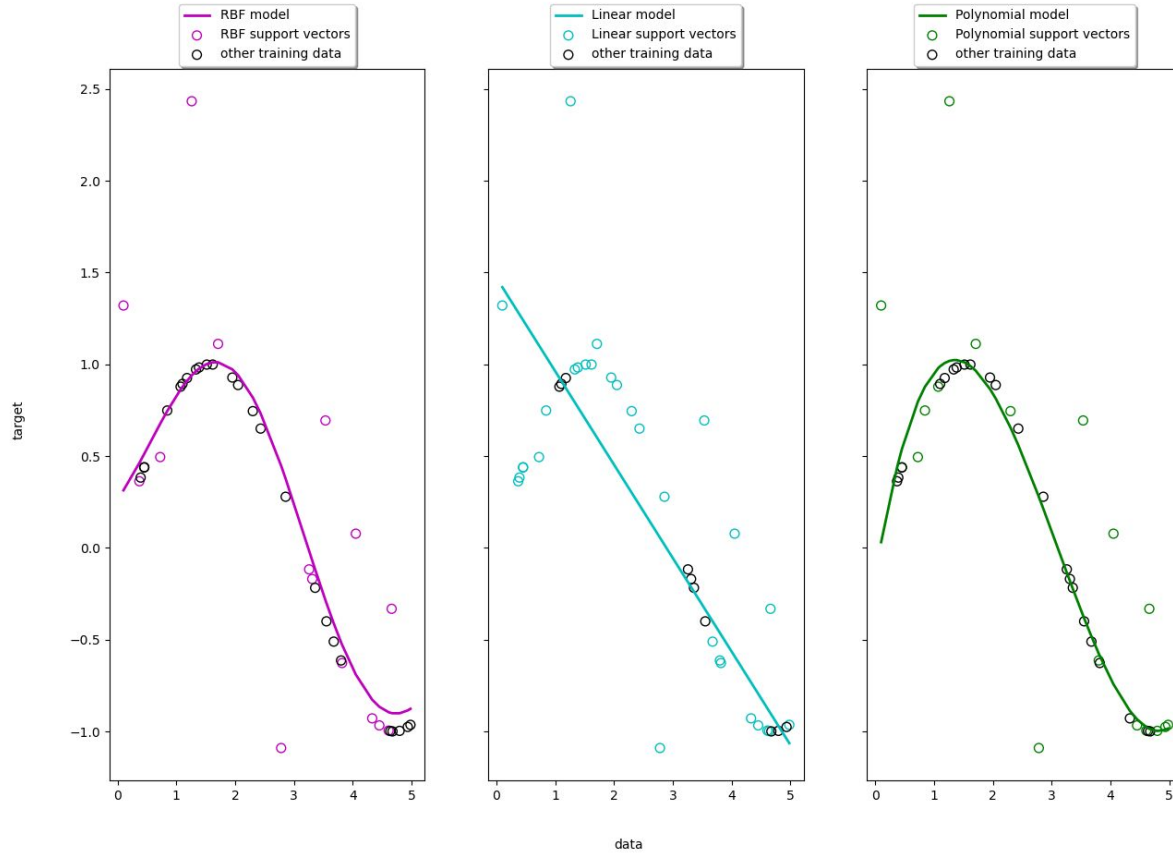
# SVR



# SVR



## Support Vector Regression



# Advantages

- Easy to implement and understand and provides powerful base-line predictions
- Few hyperparameters and therefore hyperparameter tuning
- Memory efficient
- Kernel Tricks
- Versatile in:
  - Types of data
  - Number of dimensions
  - Number of classes





# Disadvantages

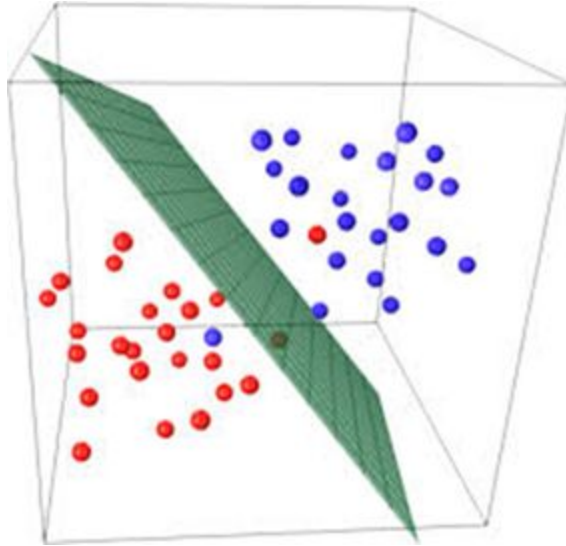
- Mapping the data to a higher dimension could lead to overfitting the model
- Results are dependent on our choice for the  $C$  parameter
- Choosing the correct kernel function for high dimensional data and data regularization are very crucial steps
- Performs poorly on very noisy or very large data sets
- Required balanced data



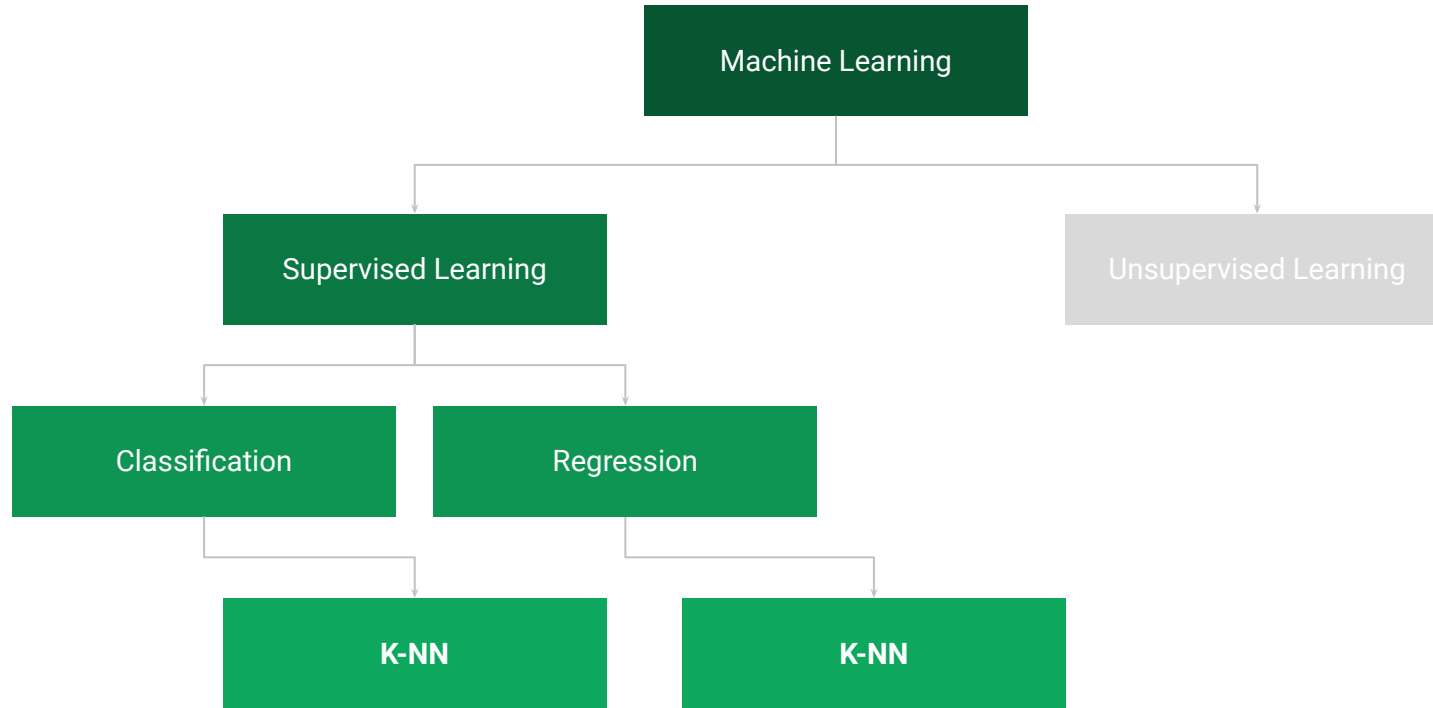
# Trending Applications

- General
  - Facial detection
  - Text and Hypertext Categorization
  - Handwriting classification
- Medical
  - Cancer Classification
  - Protein remote homology detection

# Questions?

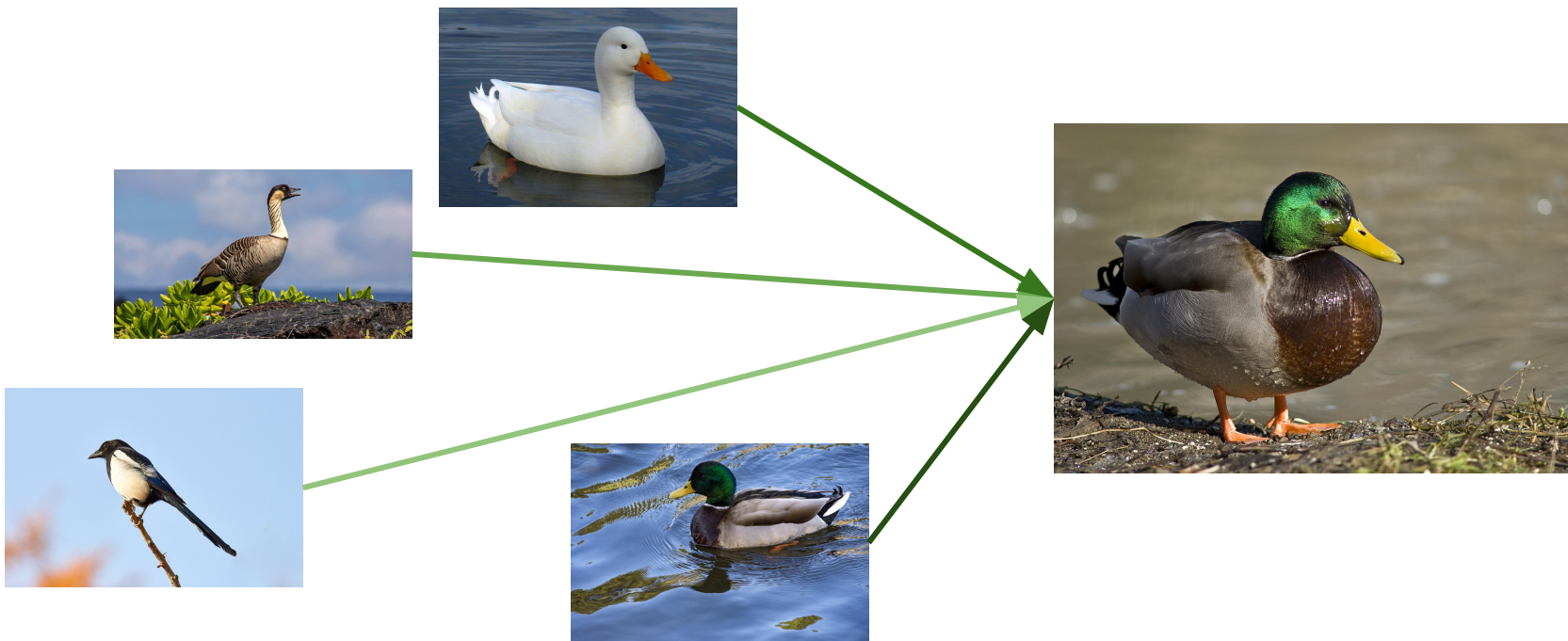


# K-Nearest Neighbors

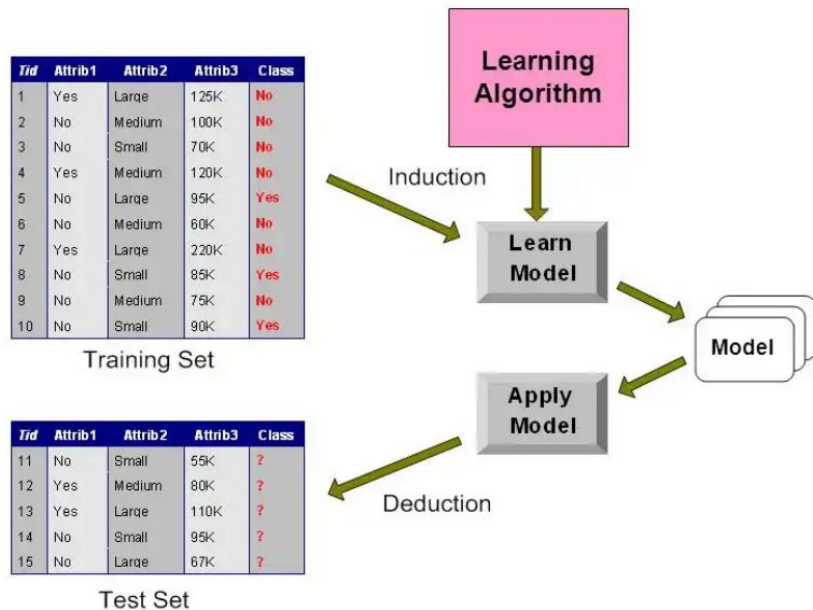


# What is K-Nearest Neighbors?

Basic idea: “If it walks like a duck, quacks like a duck, then it's probably a duck!”



# Laziness of K-NN!



## Set of Stored Cases

Atr1	.....	AtrN	Class
			A
			B
			B
			C
			A
			C
			B

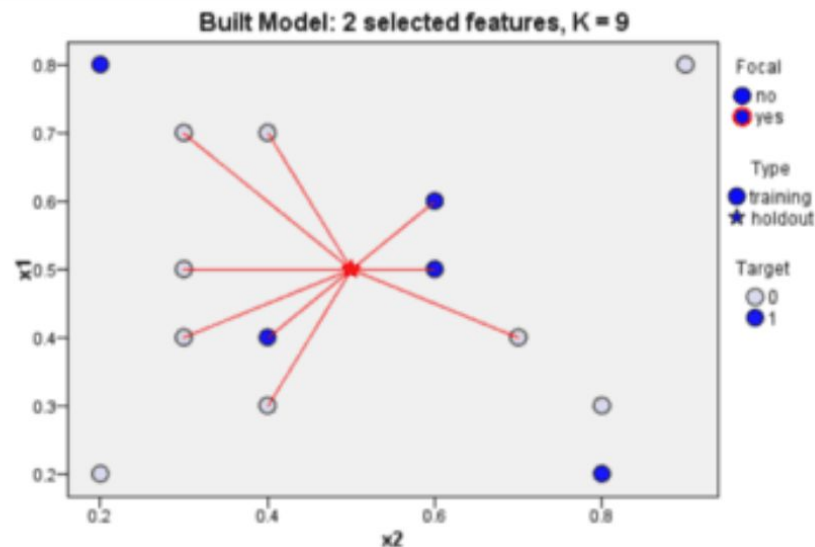
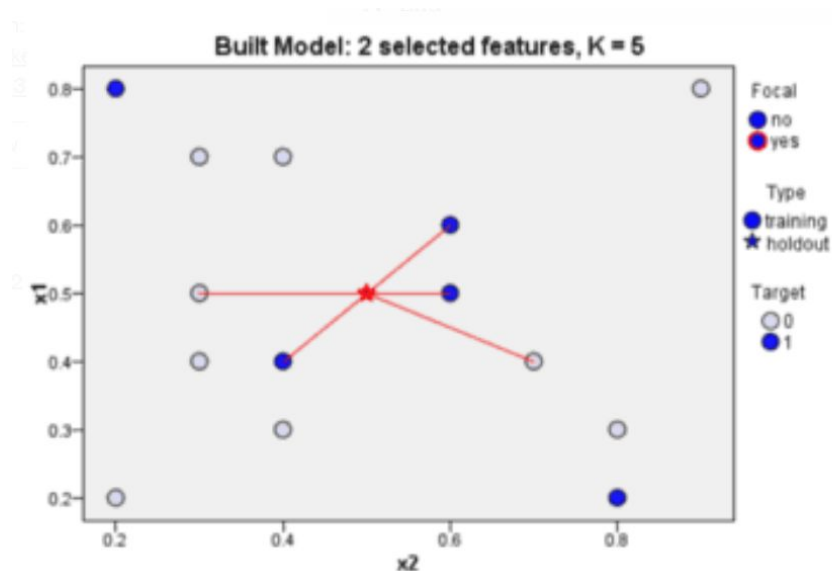
- Store the training records
- Use training records to predict the class label of unseen cases

## Unseen Case

Atr1	.....	AtrN

# K-Nearest Neighbors

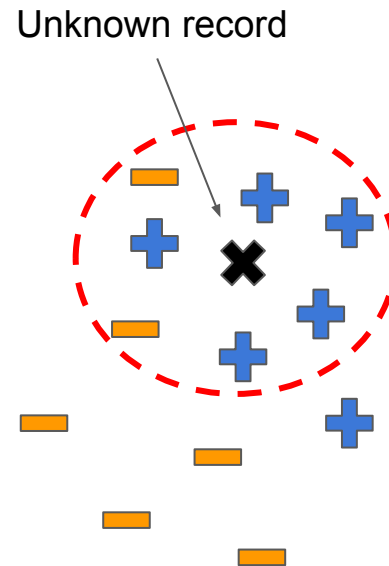
- Classification: k-nearest neighbors common class
- Regression: k-nearest neighbors average value





# K-Nearest Neighbors

- K-NN requires:
  - Sample records
  - Distance metric
  - Value of  $k$
  - Weights of neighbors
- To classify a new record:
  - Compute distances
  - Identify  $k$  nearest records
  - Use common class or average value of  $k$  nearest records



# Distance Functions

Various functions could be used, most commonly used in K-NN are:

- Minkowsky (for continuous variables)
- Euclidean (for continuous variables)
- Manhattan (for continuous variables)
- Hamming (for discrete variables)

# Minkowski, Manhattan and Euclidean Functions

- Used for continuous variables
- Computes distance between two points, A and B in a feature space with k dimensions:  $A = (x_1, x_2, \dots, x_k)$ ,  $B = (y_1, y_2, \dots, y_k)$

The diagram illustrates the relationship between the Minkowski distance function and its special cases, Manhattan and Euclidean distances. On the left, a green box labeled "Minkowski" is followed by the general formula:  $\left( \sum_{i=1}^k |x_i - y_i|^q \right)^{1/q}$ . Two arrows branch from this formula. The upper arrow, labeled "q=1", points to a green box labeled "Manhattan", which is followed by the formula  $\sum_{i=1}^k |x_i - y_i|$ . The lower arrow, labeled "q=2", points to a green box labeled "Euclidean", which is followed by the formula  $\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$ .

Minkowski

$$\left( \sum_{i=1}^k |x_i - y_i|^q \right)^{1/q}$$

q=1

Manhattan

$$\sum_{i=1}^k |x_i - y_i|$$

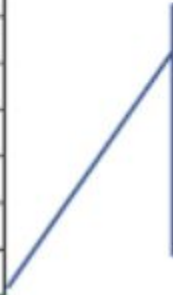
q=2

Euclidean

$$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

# Simple Example with Euclidean

Customer	Age	Loan	Default
John	25	40000	N
Smith	35	60000	N
Alex	45	80000	N
Jade	20	20000	N
Kate	35	120000	N
Mark	52	18000	N
Anil	23	95000	Y
Pat	40	62000	Y
George	60	100000	Y
Jim	48	220000	Y
Jack	33	150000	Y
Andrew	48	142000	?



We need to predict  
Andrew default status  
by using Euclidean  
distance

# Simple Example with Euclidean

Customer	Age	Loan	Default	Euclidean distance	Minimum Euclidean Distance
John	25	40000	N	1,02,000.00	
Smith	35	60000	N	82,000.00	
<b>Alex</b>	<b>45</b>	<b>80000</b>	<b>N</b>	<b>62,000.00</b>	<b>5</b>
Jade	20	20000	N	1,22,000.00	
<b>Kate</b>	<b>35</b>	<b>120000</b>	<b>N</b>	<b>22,000.00</b>	<b>2</b>
Mark	52	18000	N	1,24,000.00	
<b>Anil</b>	<b>23</b>	<b>95000</b>	<b>Y</b>	<b>47,000.01</b>	<b>4</b>
Pat	40	62000	Y	80,000.00	
<b>George</b>	<b>60</b>	<b>100000</b>	<b>Y</b>	<b>42,000.00</b>	<b>3</b>
Jim	48	220000	Y	78,000.00	
<b>Jack</b>	<b>33</b>	<b>150000</b>	<b>Y</b>	<b>8,000.01</b>	<b>1</b>
Andrew	48	142000	?		

Let assume  $K = 5$

Find minimum euclidean distance and rank in order (ascending)

In this case, 5 minimum euclidean distance.  
With  $k=5$ , there are two Default = N and three Default = Y out of five closest neighbors.

We can say Andrew default status is 'Y' (Yes)

# Hamming Distance Function

- Used for categorical variables
- Computes the number of instances in which corresponding symbols are different in two strings of equal length
- Number of attributes where  $x, x'$  differ

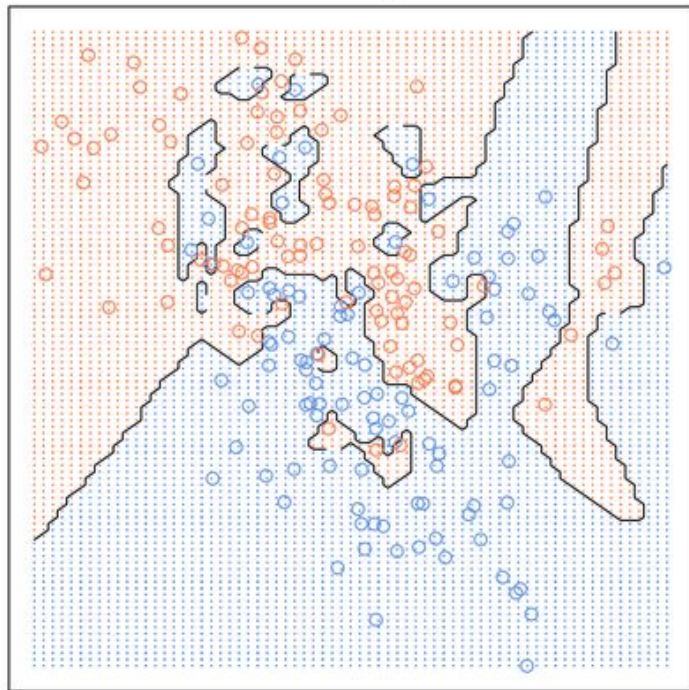
$$D(x, x') = \sum_d 1_{x_d \neq x'_d}$$

# Choosing Best K

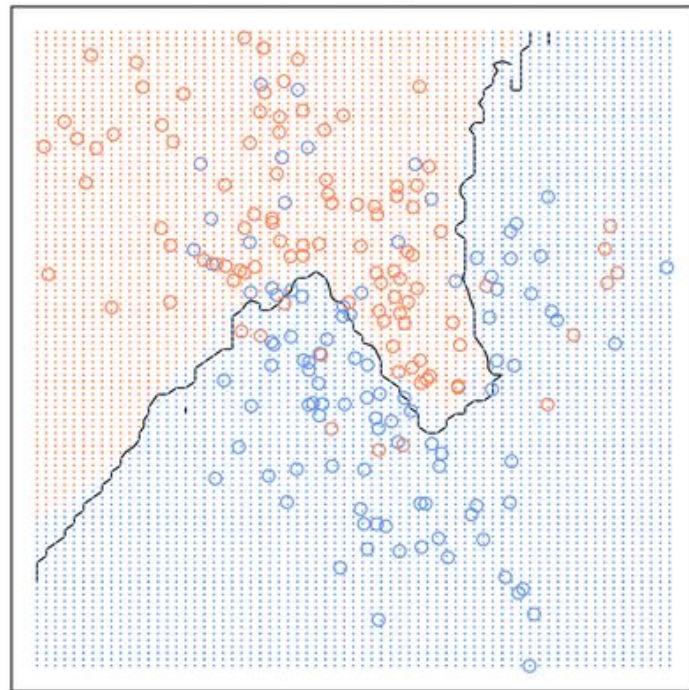
- No single value of K that will work with every dataset!
- Usually in binary classification, an odd value for K tends to work better than an even value
- Generally, increasing K makes it more precise; tends to smooth out decision boundaries, which will avoid overfitting at the cost of resolution
- We can inspect the data to choose an optimal K value, but usually we perform cross-validation to determine good K value

# Choosing Best K

1-nearest neighbours



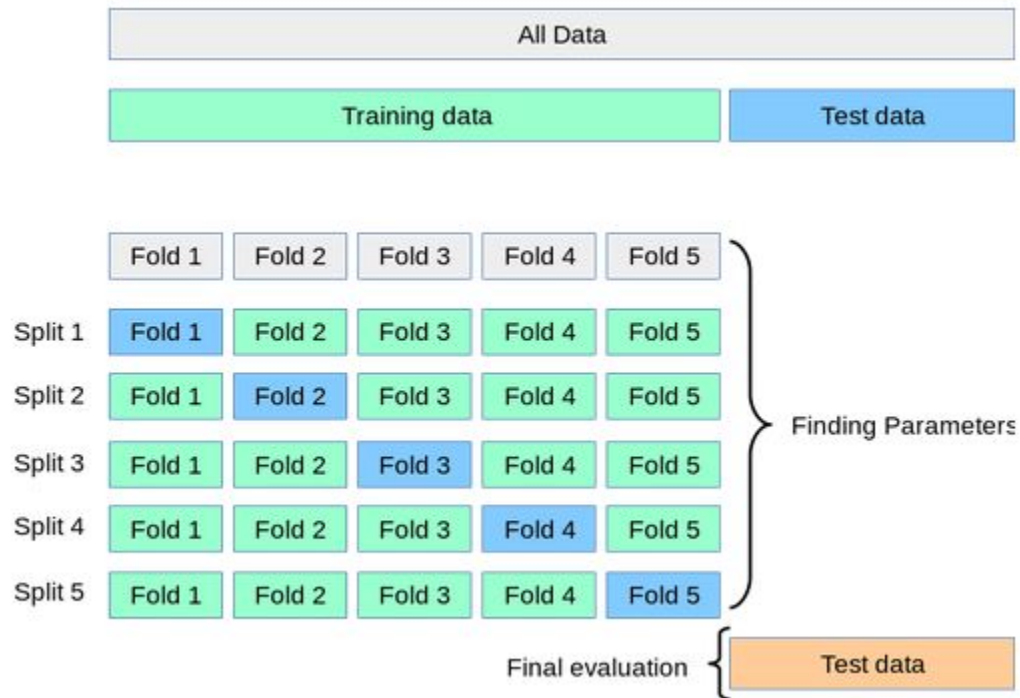
20-nearest neighbours





# Cross-Validation

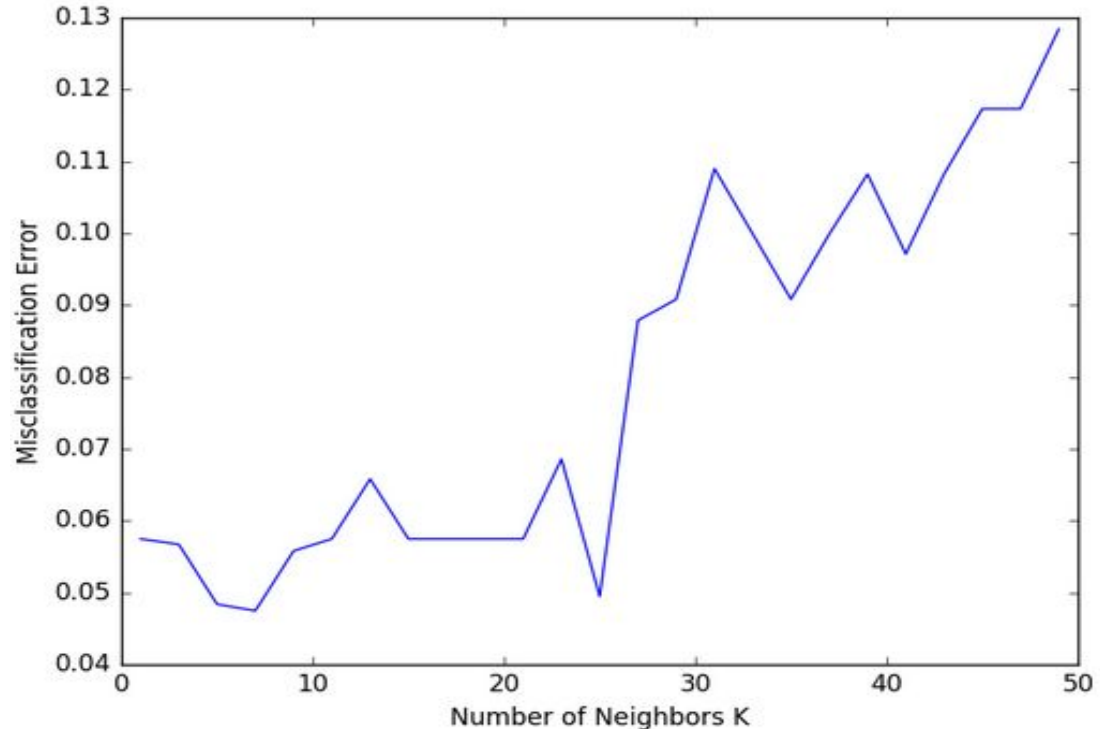
- Most commonly used is k-fold cross-validation



# !!!QUESTION ALERT!!!

In the graph below which is the best value of K to pick for our K-NN algorithm?

- 1- K = 7
- 2- K = 25
- 3- K = 48
- 4- K = 32

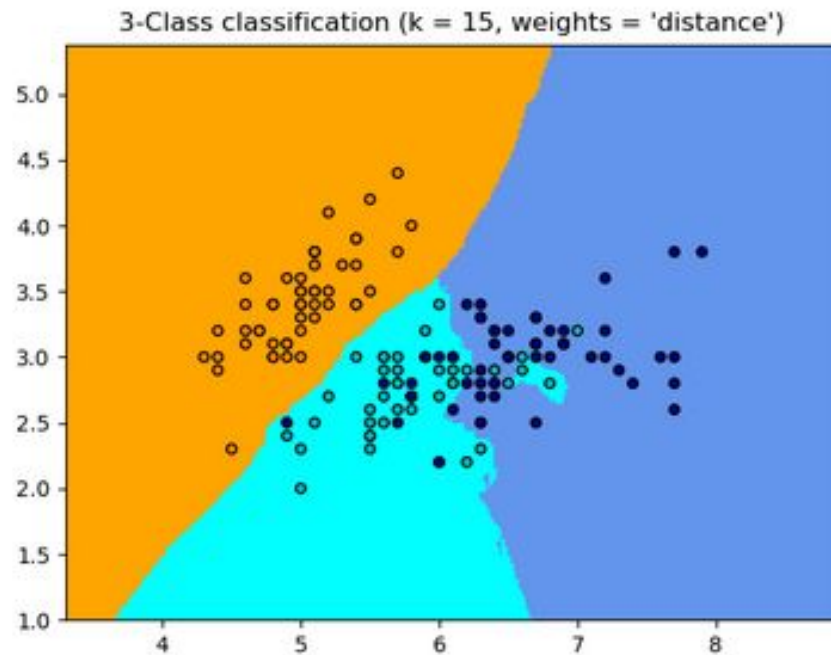
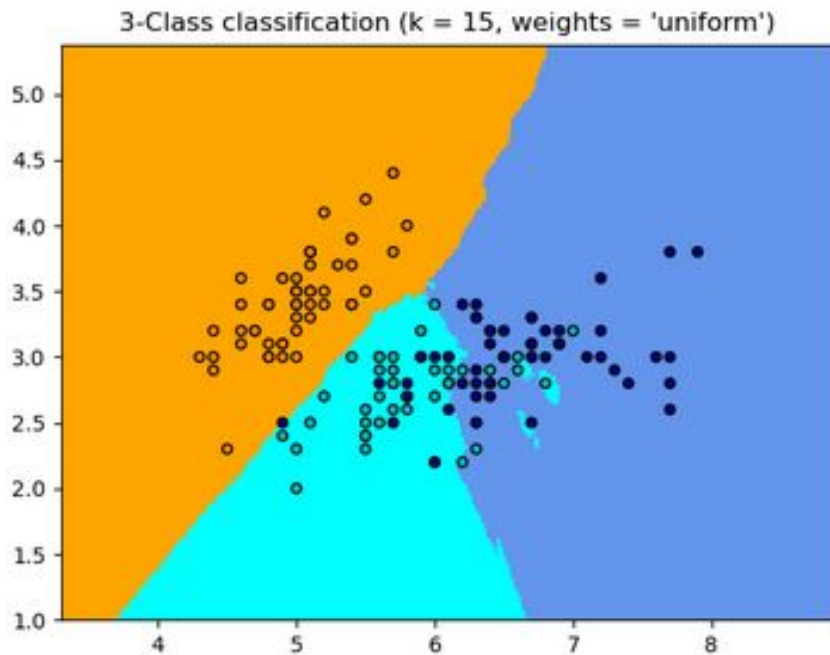


# Weights of Neighbors

This can be either:

- **Uniform:** where each point in the neighborhood is weighted equally towards the decision of the classification
- **Distance:** where points are weighted by the inverse of their distance. Closer points will have greater influence than points that are further away

# Weighted vs Uniform



# Advantages

- No training overhead therefore time efficient
- Very easy to implement as calculations only include the distance between different points
- As there is no training period thus new data can be added at any time since it won't affect the model
- Normalizing the data can dramatically improve the model accuracy
- Handles multi-class cases



# Disadvantages

- Very expensive to calculate distances when dataset is large
- Minimal training but expensive testing
- Need to determine the value of  $K$  (number of nearest neighbors)
- Doesn't work well in high dimensional datasets
- Sensitive to noisy and missing data
- Could produce different distance measures from the same training data computed before and after normalization/standardization
- Standardizing features is a drawback



# Trending Applications

- Usually in search applications when looking for similar items; i.e. document concept search
- Very popular in recommendation systems
- Used in medical sector, i.e. finding diabetics ratio