

# Introduction to Machine Learning

Live Class Notes

# Agenda

## Part 1

- Data pre-processing
- Objective functions and evaluation metrics
- Optimization techniques

## Part 2

- Supervised Learning Methods
  - Logistic Regression
  - Naive Bayes
  - k-Nearest Neighbours
  - SVM

## Part 3

- Anomaly Detection

# Data, data, data.....

According to the World Economic Forum, an astonishing 90% of the world's data has been generated in the last two years alone. It says that 2.5 quintillion ( $10^{18}$ ) bytes of data are produced by humans every day in 2020 and 463 exabytes (i.e. quintillion bytes) of data will be generated each day by humans by 2025!

But is all of this data really useful for ML algorithms and AI systems?  
How do we make data more meaningful?

“Garbage in, Garbage out”

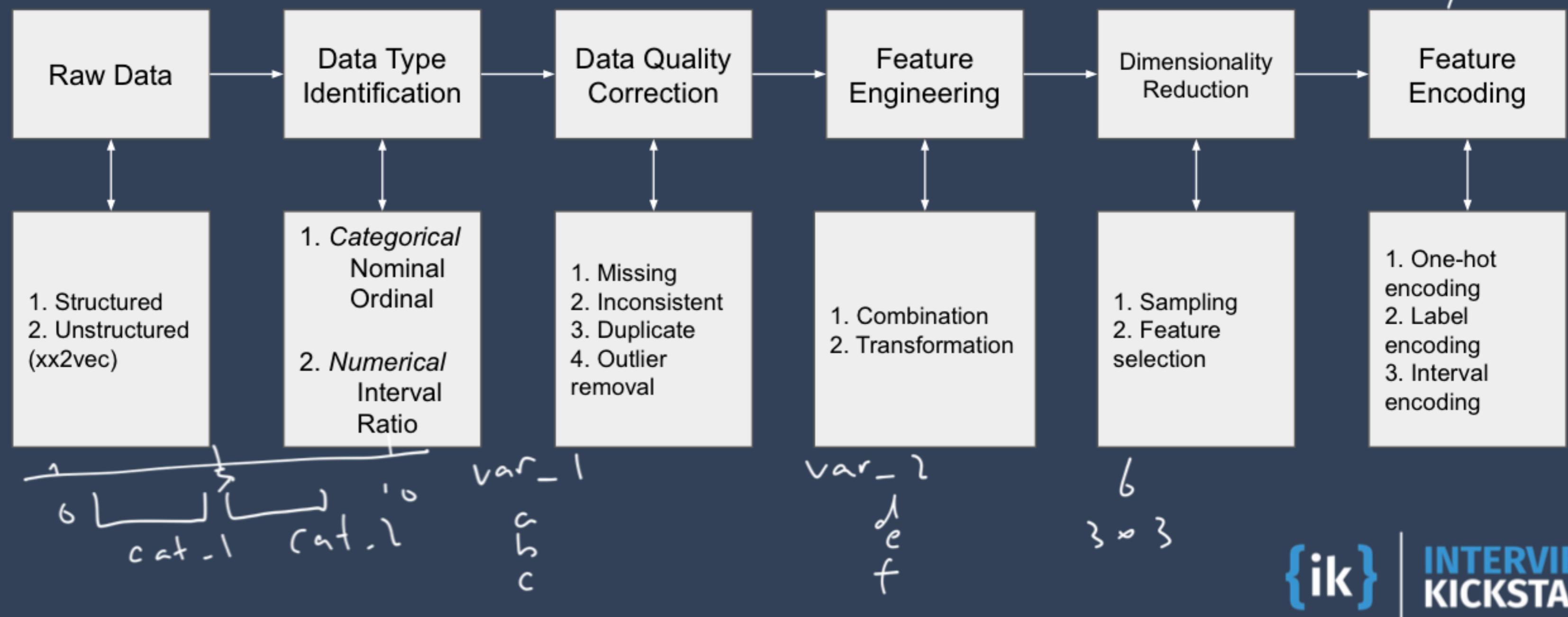
Before all the ML work, make sure your data is just relevant and logically correct.

# Data Preprocessing



- Data could take various shape, size and velocity!
  - Shape: Structured, Unstructured (Image, Audio, Video, Table, system logging, etc.)
  - Size: a dozen rows, 1M rows, ,....
  - Velocity: Streaming data, Transaction data, Daily summary, etc.
- But machines don't understand this data out-of-the-box due to difference in representations, encoding and data type
- *Data Preprocessing* transforms/encodes data to bring it to such a state that now the machine can easily parse it
  - In other words, the features of the data can now be easily interpreted by the ML algorithm

# Data Preprocessing



# Data Preprocessing

	Feature1	Feature2	Feature3
Categorical	Male	0.53	1
	Female	0.25	NULL
	Male	1.5	3
	Female	999.9	5

Annotations:

- An arrow points from the text "Categorical" to the first row of the table.
- An arrow points from the text "Numerical" to the value "999.9" in the Feature2 column.
- An arrow points from the text "Outlier" to the value "999.9" in the Feature2 column.
- An arrow points from the text "Nominal/Ordinal" to the value "5" in the Feature3 column.
- An arrow points from the text "Missing" to the value "NULL" in the Feature3 column.

- 1-hot encoding
  - Male - [0, 1]
  - Female - [1, 0]
- Label encoding
  - Male - 1
  - Female - 0

Question: Should we always remove outliers from our dataset during the preprocessing step?

- a. Yes
- b. No

depends

# Data Quality Correction

## 1. Imputation

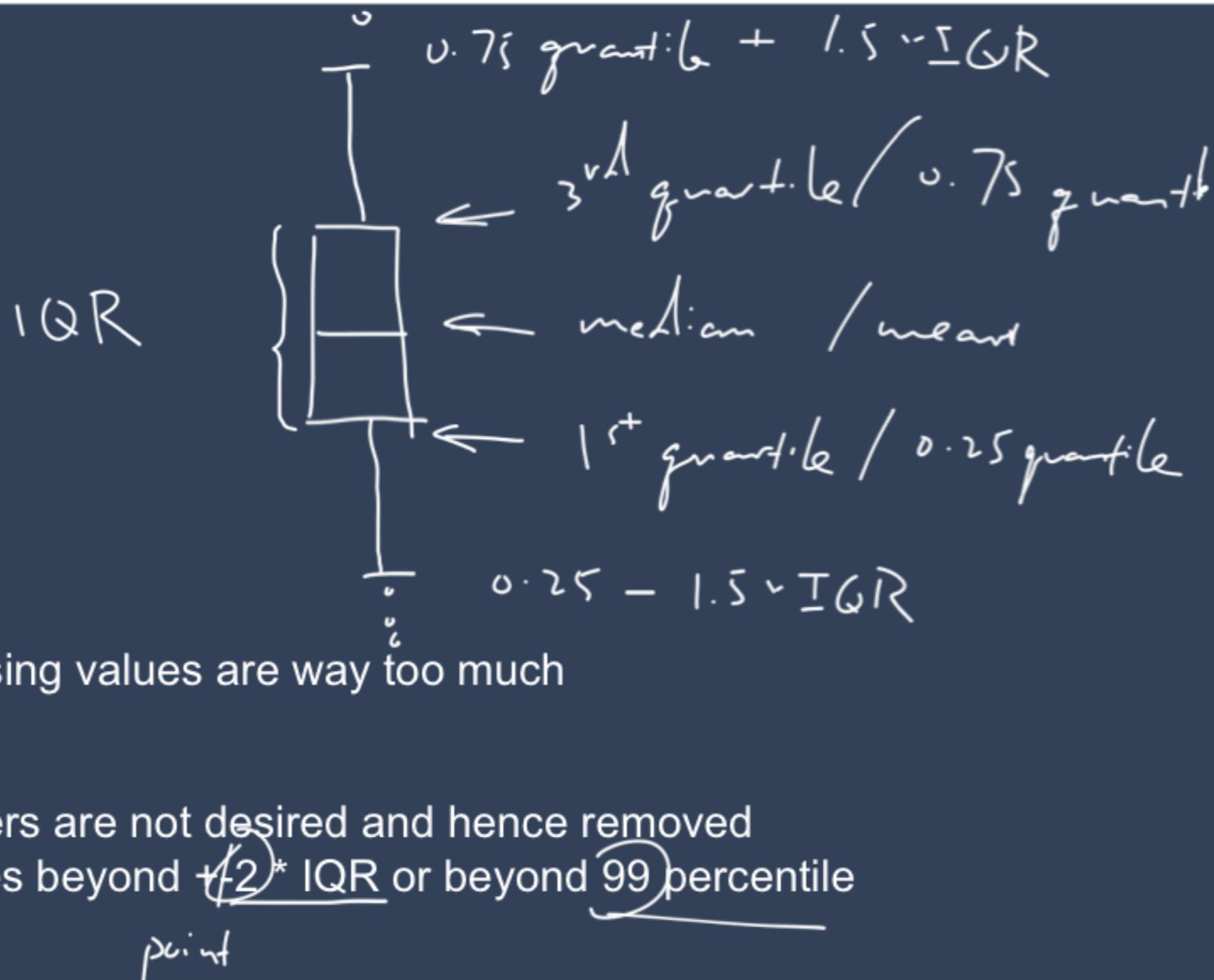
- a. Median imputation
- b. Mode imputation
- c. Gaussian imputation

## 2. Feature dropping

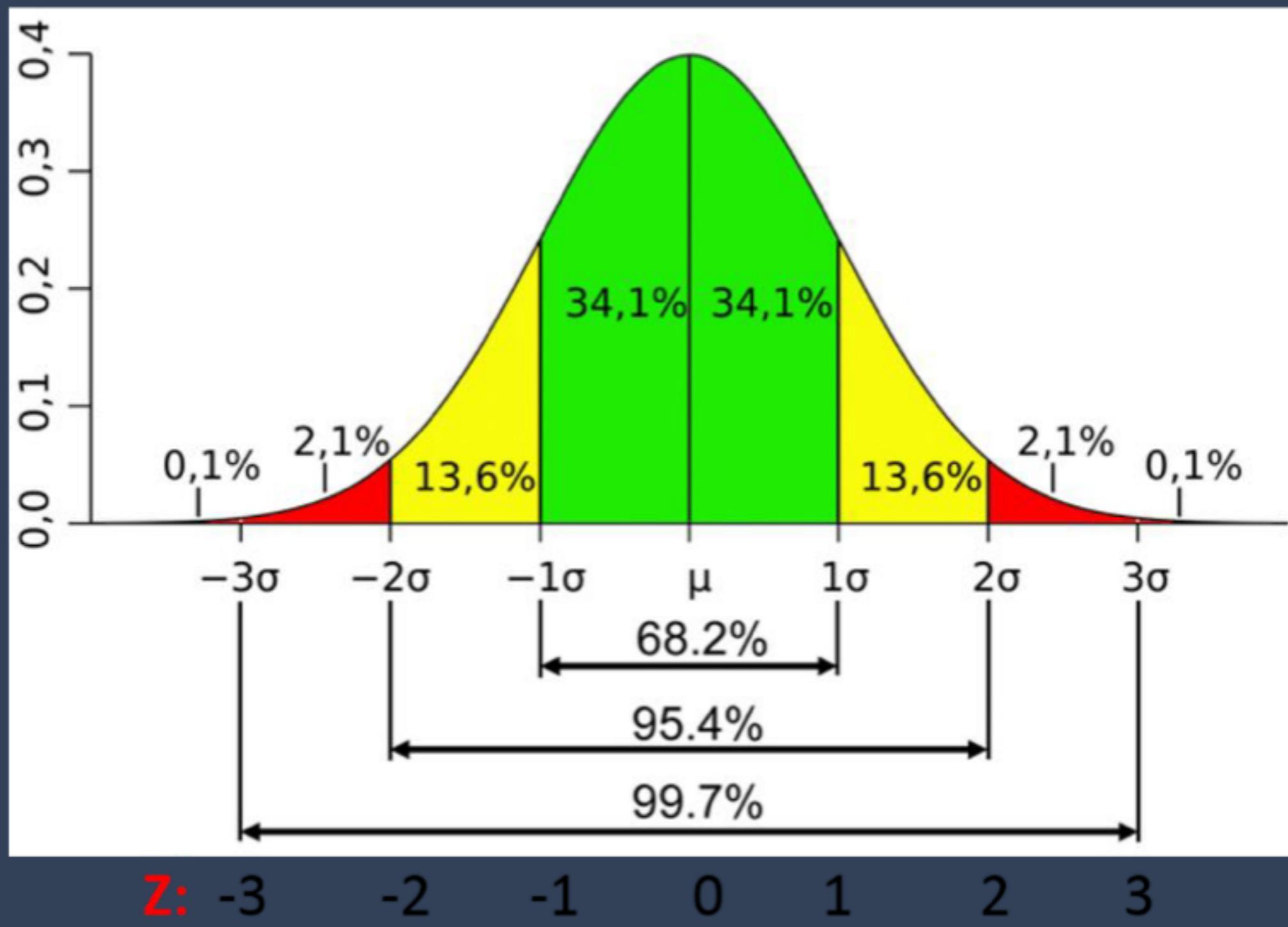
- a. Drop the feature if the number of missing values are way too much

## 3. Outlier removal

- a. There are many problems when outliers are not desired and hence removed
- b. Most common way is to remove values beyond ~~1.5 \* IQR~~ or beyond ~~99~~ percentile



## Imputation/Outlier Removal - Example



Now when we have made our data meaningful, what's next? How can ML algorithms consume it?

$$\min \sum (y_i - x_i \beta)^2 + \lambda \|\beta\|,$$

## Partition the data - Data classes

$\lambda$ : val error  
 $\lambda_1$   
 $\lambda_k$  . -

- Training Data
  - Goal is to find a model that optimizes the training objective metrics
  - This is the data on which machine learning algorithms are actually trained to build a model
  - Typically this is 70-80% of the overall data
- Validation Data — optional
  - This is the data on which hyperparameters of the model are tuned
  - Typically this is 5-10% of the overall data
- Test Data
  - This is the data on which model is evaluated before deployment and should not be used for parameter tuning
  - Typically this is 10-20% of the overall data

Question: Why do we keep aside the validation dataset?

- a. Parameter estimation
- b. Hyperparameter tuning ✓
- c. Both a & b
- d. None

Any good ML system should be able to quantify its goodness!

# Regression: Objective Functions

- Regression Loss

- Mean Square Error/Quadratic Loss/L2 Loss

- $MSE = \sum_{i=1}^n (y_i - \tilde{y}_i)^2 / n$

- Mean Absolute Error/L1 Loss

- $MAE = \sum_{i=1}^n |y_i - \tilde{y}_i| / n$

$y_i$  is the true target, while  $\tilde{y}_i$  is the model generated target,  $n$  is the number of *training* sample

$$R^2 = \frac{\text{var exp/true}}{\text{var in } \hat{y}}$$
$$= \frac{\sum (y_i - \tilde{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

0.5-quantile

if we keep outliers . and train with it



Question: Which is more robust to data with outliers?

- a. MSE
- b. MAE ✓
- c. RMSE  $\sqrt{\text{MSE}}$
- d. None



$$|y - \hat{y}|$$

# Regression: Evaluation Metrics

- All the metrics in previous slides
  - MSE
  - MAE
- MAPE (Mean Absolute Percentage Error)
  - $MAPE = \frac{1}{n} \sum |(y_i - \tilde{y}_i) / \underline{y_i}|$
  - Mean Absolute Percentage Error is also commonly used in evaluation time-series forecasting models



$y_i$  is the true target, while  $\tilde{y}_i$  is the model generated target, n is the number of *validation/testing* sample

# Classification: Objective Functions

- Classification Loss

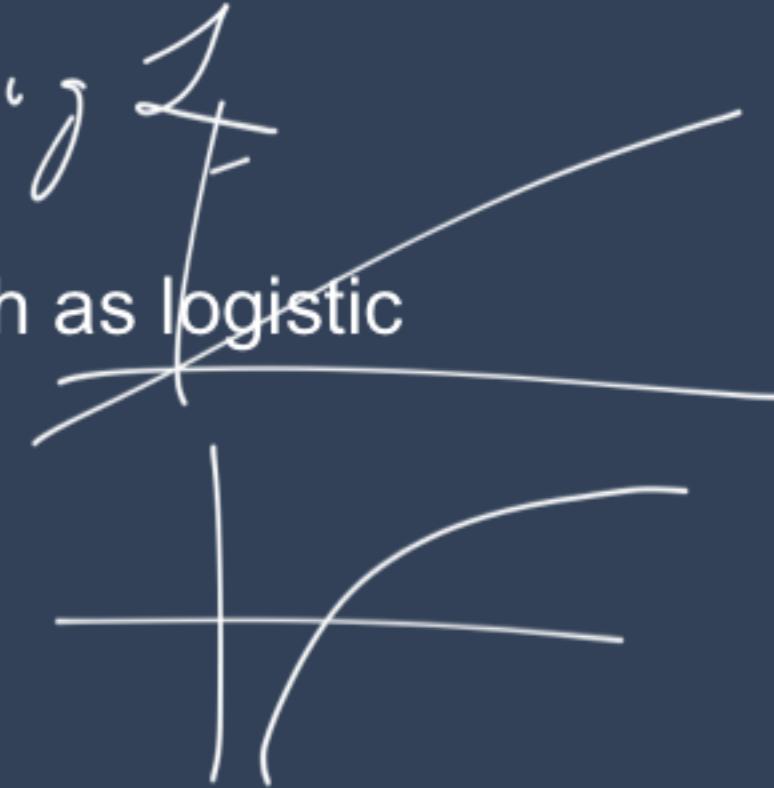
- Cross Entropy Loss

- $\text{Cross Entropy Loss} = -\sum(y_i \log(\tilde{y}_i) + (1 - y_i) \log(1 - \tilde{y}_i))$

- Used in binary and multi-class classification algorithms such as logistic regression, neural networks

maximum likelihood estimator (MLE)

$$\hat{\theta} = \arg \max_{\theta} \prod_{i=1}^n f(x_i; \theta)$$



# Classification: Evaluation Metrics

		Actual Values	
		Positive(1)	Negative(0)
Predicted Values	Positive(1)	TP	FP
	Negative(0)	FN	TN

- Confusion Matrix
  - It is a performance measurement for machine learning classification problem where output can be two or more classes.
  - Confusion table for 2 classes classification problem has 4 different combinations of predicted and actual values
  - It is extremely useful for measuring Recall, Precision, Specificity, Accuracy, and most importantly AUC-ROC curves

XMCC

## Classification: Evaluation Metrics

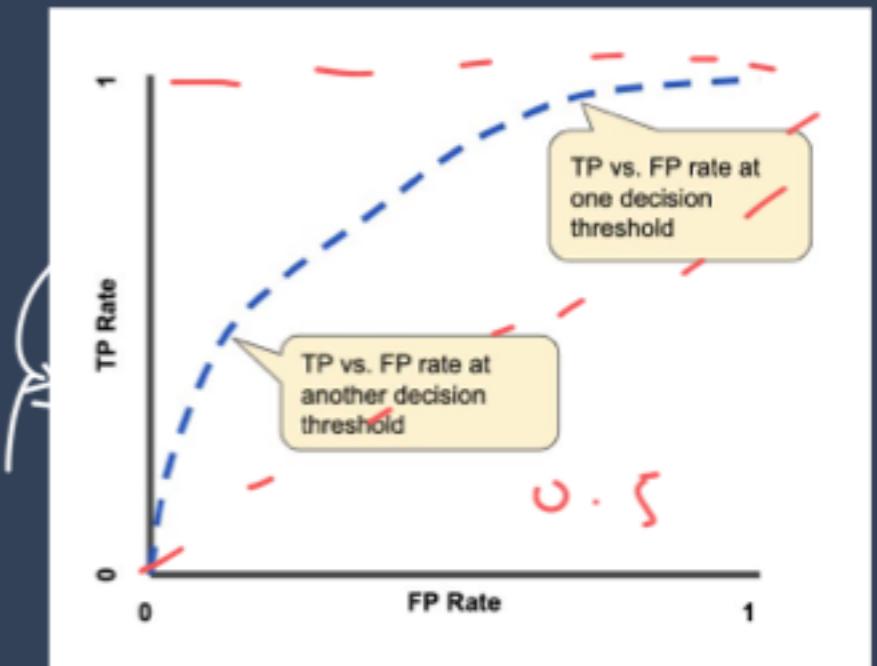
- Precision
  - $TP / (TP + FP)$
  - Represents % of predictions that are actually positive out of all we have predicted as positive
- Recall
  - $TP / (TP + FN)$
  - Represents coverage of the model in terms of how many we predicted correctly from all the positive classes
- F1-Score
  - $2 * Recall * Precision / (Recall + Precision)$
- Accuracy
  - $(TP + TN) / (TP + FP + TN + FN)$
  - Percentage of correct predictions

| U U U | P

# Classification: Evaluation Metrics

- AUC-ROC

- Area under curve - receiver operating characteristic curve
- An ROC curve is a graph showing the performance of a classification model at all classification **thresholds**.
- This curve plots two parameters:
  - Y axis -> True Positive Rate =  $TP / (TP + FN)$
  - X axis -> False Positive Rate =  $FP / (FP + TN)$
- Ranges between 0 and 1 with higher the better model performance
- For a model doing random prediction, AUC-ROC is 0.5
- For imbalance classes, AUC-ROC may not be the best metric. This is because a small number of correct/incorrect prediction can result in a large change in the ROC curve



$threshold = 0.01$   
 $predict > 0.01$   
 $predict < 0.01$

# Confusion Matrix

		Actual Values	
		Positive(1)	Negative(0)
Predicted Values	Positive(1)	100	25
	Negative(0)	10	100

- Precision =  $100 / (100 + 25) = 100 / 125 = 0.80$
- Recall =  $100 / (100 + 10) = 100 / 110 = 0.91$
- F-Score =  $(2 * 100 / 125 * 100 / 110) / (100 / 125 + 100 / 110) = 0.85$
- Accuracy =  $(100 + 100) / (100 + 25 + 10 + 100) = 0.85$

# Classification: Evaluation Metrics

- AUC-PR
  - Area under curve - Precision-Recall curves
  - It is a graph showing the performance of a classification model at all classification thresholds
  - This curve plots two parameters:
    - Y axis -> Precision =  $TP / (TP + FP)$
    - X axis -> Recall =  $TP / (TP + FN)$
  - Range between 0 and 1 with higher the better model performance
  - It is widely used for imbalance class problems

# ML Metrics

- MSE
  - $MSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 / n$
  - Mean squared error is widely used as a metric in regression and time-series forecasting models
- MAPE
  - $MAPE = 1/n \sum |(A_i - F_i) / A_i|$
  - Mean Absolute Percentage Error is commonly used in evaluation time-series forecasting models

True optimization is the revolutionary contribution of modern research to decision processes - George Dantzig

# Optimization

Optimization is the task of minimizing or maximizing an objective function described by a set of parameters  $\Theta$ .

- For example, in linear regression, we would like to find  $\Theta$  that minimize MSE i.e.  $L(\Theta) = \sum_{i=1}^n (y_i - \text{dot}(\Theta, x_i))^2 / n$
- Linear regression with one dimension feature  $L(\Theta) = \sum_{i=1}^n (y_i - (w_1 x_i + w_0))^2 / n$ . Here  $\Theta = [w_0, w_1]$

The following are the most common optimisation methods (and facilitators) used in various ML/DL algorithms:

- Gradient Descent
- Stochastic Gradient Descent(SGD)
- Adam Optimiser

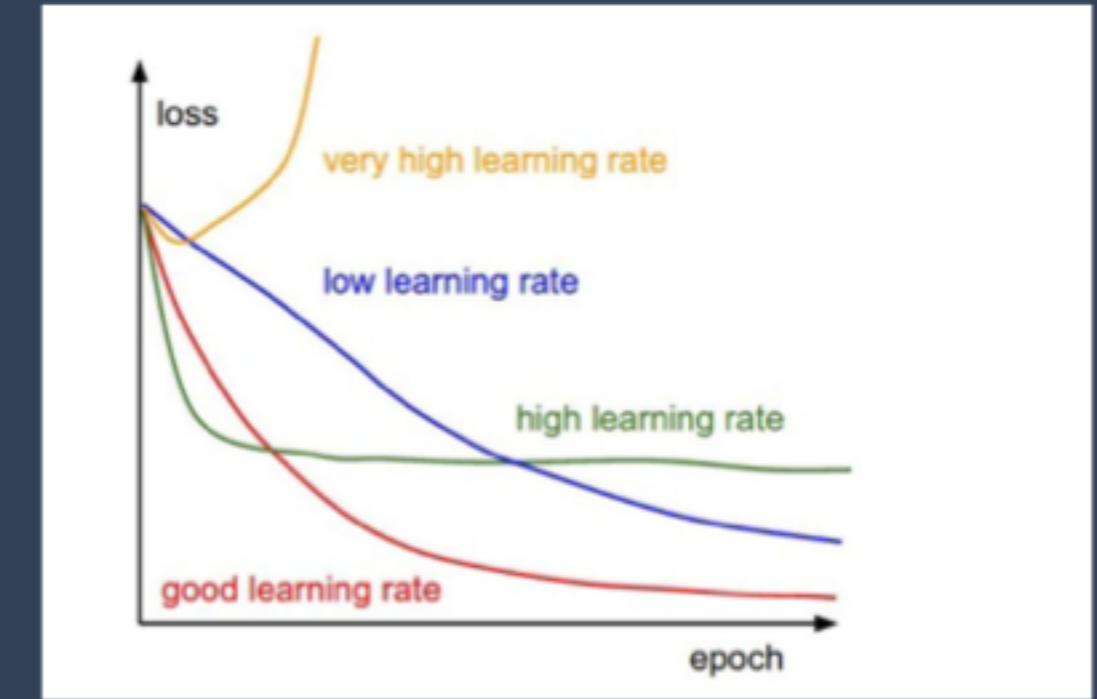
you + NN  
want to know

# Gradient Descent



- The goal of Gradient Descent is to minimize the objective convex function  $L(\cdot)$  via multiple iterations (epochs). It is implemented as follows:

- Randomly initialize values of all the parameters
- At step  $t$ , update values using-
  - $w_{t+1} = w_t - \text{learning rate} * (\delta J(\Theta)) / (\delta \text{weight})$
  - In our 1D linear regression example
    - $\delta L(\Theta) / (\delta w_{1t}) = -2/n * \sum_{i=1}^n x_i (y_i - (w_{1t}x_i + w_{0t}))$
    - $\delta L(\Theta) / (\delta w_{0t}) = -2/n * \sum_{i=1}^n (y_i - (w_{1t}x_i + w_{0t}))$
- Repeat until convergence condition is met



- Learning rate must be chosen wisely as:
  - if it is too small, then the model will take some time to converge.
  - if it is too large, model will keep on missing the convergence point by jumping abnormally

# Gradient Descent with Momentum

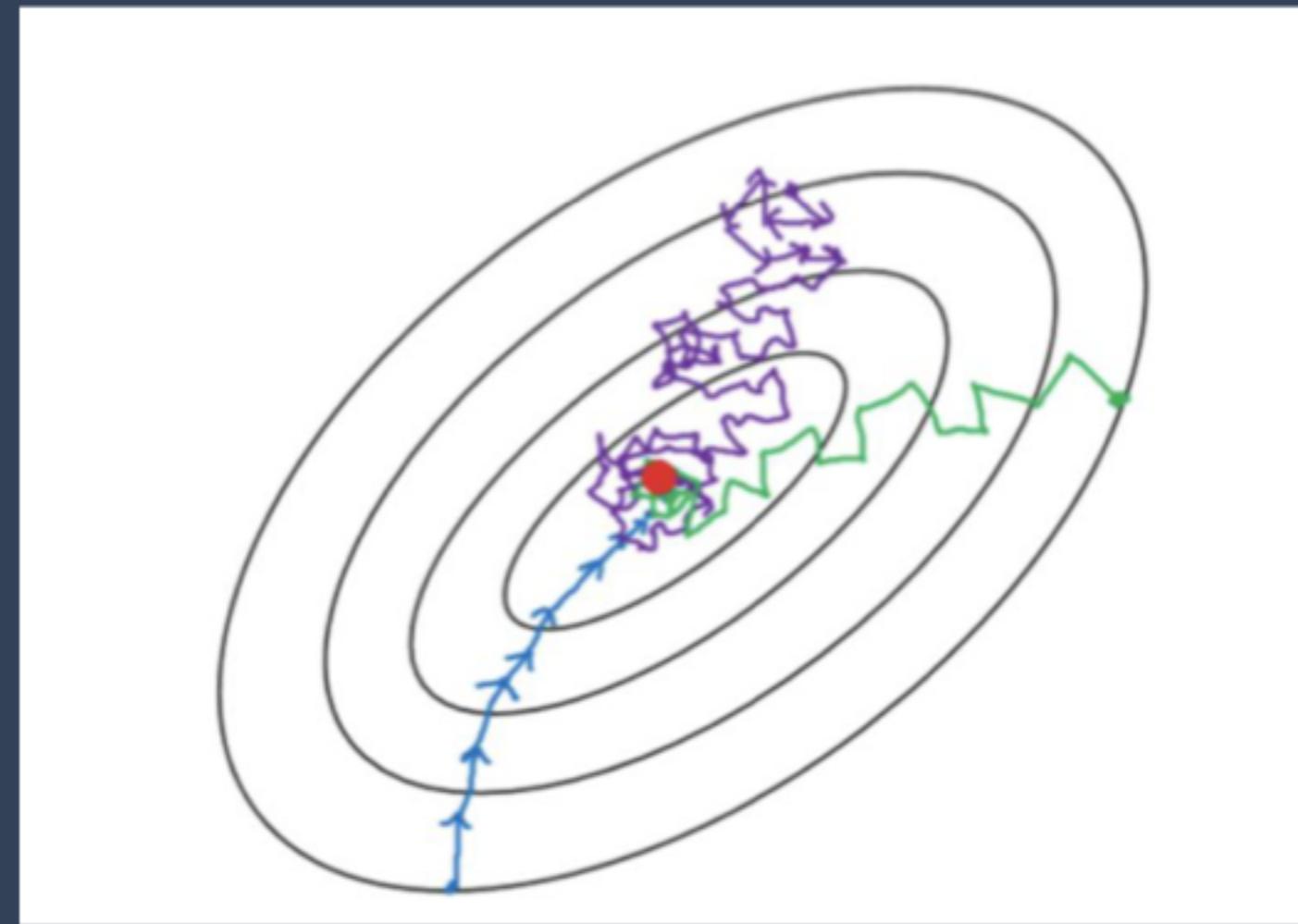
- A vanilla gradient descent model might converge slowly if the initialized values are far away from the minima
- To overcome this, we introduce the concept of momentum which accelerates the overall optimization process
- Momentum involves adding an additional hyperparameter that controls the amount of gradient history (momentum) to include in the update equation
  - $m_t = \beta m_{t-1} + (1 - \beta) \frac{\nabla L(\Theta)}{\text{learning\_rate}}$
  - $w_{t+1} = w_t - \text{learning\_rate} * m_t$
  - $m_t$  term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions. As a result, we gain faster convergence and reduced oscillation.

# Stochastic Gradient Descent(SGD)

- The vanilla gradient descent algorithm computes gradient for each weights using all data point which makes it very slow
- Stochastic Gradient Descent helps in this case by being “stochastic”
  - In plain terms, stochastic means random
  - In our 1D linear regression example
    - $\frac{\delta L(\Theta)}{\delta w_{1t}} = -\frac{2}{n} * \sum_{i=1}^n x_i(y_i - (w_{1t}x_i + w_{0t}))$
    - $\frac{\delta L(\Theta)}{\delta w_{2t}} = -\frac{2}{n} * \sum_{i=1}^n (y_i - (w_{1t}x_i + w_{0t}))$
- Stochastic Gradient Descent randomly picks one data point from the whole data set at each iteration to reduce the computations enormously
- Now the above step can also be done in batches, in which case, we call it *Batch Gradient Descent*

blue: all points  
smooth

green:



batch  
jitter  
smooth

min batch signal  
< jitter

Question: Which of the following statements is NOT TRUE?

- a. Purple: Stochastic Gradient Descent, Blue: Batch Gradient Descent
- b. Blue: Stochastic Gradient Descent, Purple: Batch Gradient Descent ? ✓
- c. Blue: Batch Gradient Descent, Green: Mini Batch Gradient Descent
- d. Blue: Batch Gradient Descent, Purple: Stochastic Gradient Descent

# Adam Optimizer

- Adam - Adaptive Moment Estimation
- It is an algorithm for solving optimization problems on top of gradient descent algorithm
- Adam optimizer has two components to it:
  - Momentum
  - Root Mean Square Propagation (RMSP)
- Adam optimizer control the rate of gradient descent in such a way that there is minimum oscillation when it reaches the global minimum while taking big enough steps so as to pass the local minima obstacles along the way

# Adam Optimizer

- Momentum
  - Finer control on *gradients*
  - $w_{t+1} = w_t - \alpha m_t$
  - $m_t = \beta m_{t-1} + (1 - \beta) \underline{(\delta L / \delta w_t)}$ 
    - $m_t$  = aggregate of gradients at time t
    - $m_{t-1}$  = aggregate of gradients at time t-1
    - $W_t$  = weights at time t
    - $W_{t+1}$  = weights at time t+1
    - $\alpha_t$  = learning rate at time t
    - $\partial L$  = derivative of Loss Function
    - $\partial W_t$  = derivative of weights at time t
    - $\beta$  = Moving average parameter

# Adam Optimizer

- Root Mean Square Propagation (RMSP)

- Finer control on *learning rates*
- $w_{t+1} = w_t - \frac{\alpha_t}{(v_t + \epsilon)^{1/2}} * (\delta L / \delta w_t)$
- $v_t = \beta v_{t-1} + (1 - \beta) (\delta L / \delta w_t)^2$ 
  - $W_t$  = weights at time t
  - $W_{t+1}$  = weights at time t+1
  - $\alpha_t$  = learning rate at time t
  - $\partial L$  = derivative of Loss Function
  - $\partial W_t$  = derivative of weights at time t
  - $V_t$  = sum of square of past gradients
  - $\beta$  = Moving average parameter
  - $\epsilon$  = A small positive constant

# Adam Optimizer

- $m_t = \beta_1 m_{t-1} + (1 - \beta_1) (\delta L / \delta w_t)$
- $v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\delta L / \delta w_t)^2$
- $w_{t+1} = w_t - \frac{\alpha_t}{(v_t + \epsilon)^{1/2}} * \underline{m_t}$

Commonly used values: 0.9 for  $\beta_1$ , 0.999 for  $\beta_2$ ,  $10^{-8}$  for  $\epsilon$ , and  $10^{-3}$  for  $\alpha_t$

Time for some action now!

$$y = \{0, 1\}$$

$$y \in \mathbb{R}$$

Question: Can we use Logistic Regression to solve regression problems?

a. Yes

b. No

$$y = \frac{1}{1 + \exp(-x\beta)}$$

$$\begin{aligned} z &= x\beta \\ \Rightarrow \beta &= -\ln\left(\frac{1}{y} - 1\right) = z \end{aligned}$$

Question: Is logistic regression a linear model given that the decision boundary is non-linear?

- a. Yes
- b. No

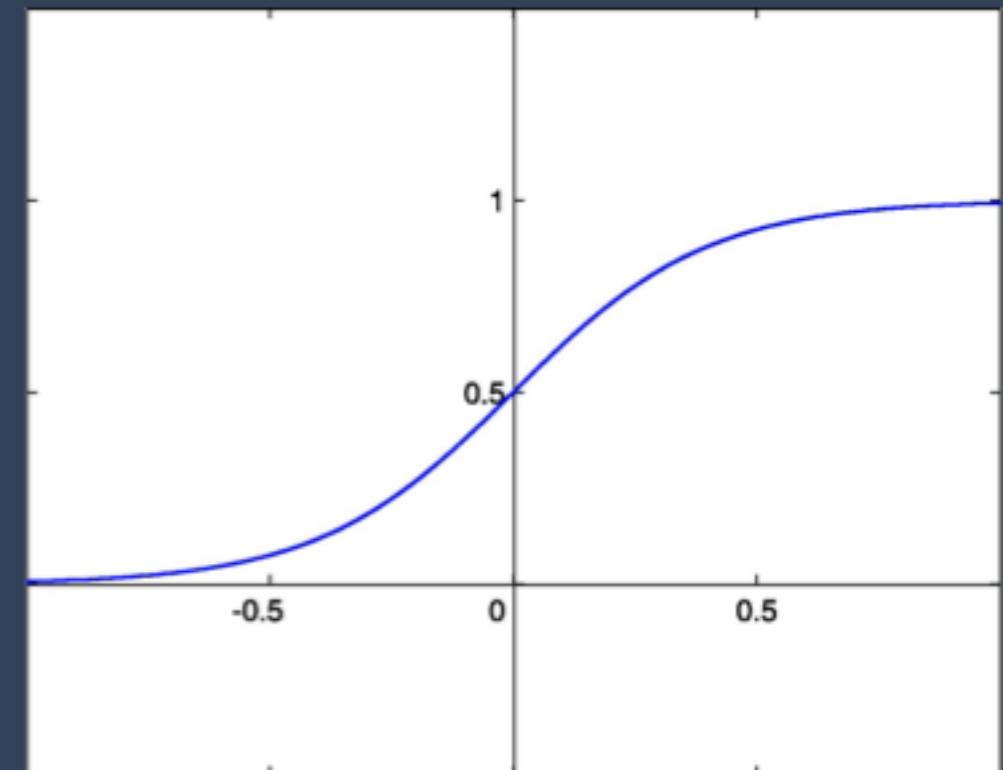
# Logistic Regression

$$\mathbb{R} \rightarrow (0, 1)$$

- Logistic Regression (LR) is primarily used as binary classification algorithm which classifies observations into binary class 0/1
- The core component of a LR algorithm is the “sigmoid” function which acts as a probabilistic component on top of the linear equation
- Consider a single input observation  $x$ , which we will represent by a vector of features  $[x_1, x_2, \dots, x_n]$
- The classifier output  $y$  as a probability, can be thresholded to either 1 or 0 represented by either  $P(Y=1 | x)$  or  $P(Y=0 | x)$
- Logistic Regression helps in accomplishing the above objective

# Logistic Regression

- Logistic regression solves the binary classification task by learning, from a training set, a vector of **weights** and a **bias** term
- Before jumping on the classification, let's recall our regression equation
  - $z = (\sum_{i=1}^n w_i * x_i) + b$ , or
  - $z = w \cdot x + b$ , represented by the dot product
- Now, recall the sigmoid function
  - $\sigma(z) = 1 / (1 + e^{-z})$
  - Sigmoid function squeezes the value between 0 and 1 and is ideal for representing the values as probability



# Logistic Regression

$$P_i \sim \text{Ber}(P_i)$$

each observation  $i$  ~  $\text{Ber}(P_i)$

- Hence, we easily move from regression -> classification

- $z = w \cdot x + b$
- $y_{\hat{}} = \sigma(z) = 1 / (1 + e^{-z})$
- $y_{\hat{}} = 1 / (1 + e^{-(w \cdot x + b)})$

- Thus,

- $y_{\hat{}} = 1 \quad \text{if } P(y=1 | x) > 0.5 \text{ (or some other thresholds)}$   
 $= 0 \quad \text{otherwise}$

$$P_i = y_{\hat{}} = \frac{1}{1 + e^{-x_i \beta}}$$

$$\begin{aligned} \frac{1}{P_i} &= 1 + e^{-x_i \beta} \\ x_i \beta &= -\log \frac{1 - P_i}{P_i} \\ &= \log \frac{P_i}{1 - P_i} \end{aligned}$$

# Logistic Regression

- How do we learn then?
- Well, Logistic Regression uses *cross-entropy* loss function to optimize for the parameter w
- We need a loss function that expresses, for an observation x, how close the classifier output is to the correct output (y, which is 0 or 1)
  - $L(y_{\hat{h}}, y) = \text{How much } y_{\hat{h}} \text{ differs from the true } y$

# Logistic Regression

- We'd like to learn weights that maximize the probability of the correct label  $p(y|x)$ . The probability of this can be expressed like a Bernoulli equation
  - $p(y | x) = y_{\text{hat}}^y(1 - y_{\text{hat}})^{1-y}$
- Now, we take log on both sides
  - $\log p(y | x) = \log (y_{\text{hat}}^y(1 - y_{\text{hat}})^{1-y})$
  - $\log p(y | x) = y \log y_{\text{hat}} + (1 - y) \log (1 - y_{\text{hat}})$
- In order to turn this into loss function, we flip the sign
  - $L = -\log p(y | x) = - (y \log y_{\text{hat}} + (1 - y) \log (1 - y_{\text{hat}}))$
- The above loss function can be optimized using gradient descent by taking the derivatives

negative  
loss = log loss

# Logistic Regression

- Logistic Regression could also suffer from overfitting and hence we need to regularize
- To avoid overfitting, a new regularization term  $R(\theta)$  is added to the objective function,  $\theta = [w_i, b]$ 
  - $\theta_{\hat{}} = \operatorname{argmax}_{\theta} \sum \log P(y^{(i)} | x^{(i)}) - R(\theta)$
- Here,  $R$  could be a L1 or L2 regularization function where former does feature selection while latter prevents large weights

# Logistic Regression - Step by step

X1	X2	Y
3	3	0
1	2	0
3	4	0
1	2	0
3	3	0
8	3	1
5	2	1
7	2	1
9	0	1
8	4	1

- Calculate a prediction using the current values of the coefficients

Let's start off by assigning 0.0 to each coefficient and calculating the probability of the first training instance that belongs to class 0.

$$\begin{aligned} b_0 &= 0.0 \\ b_1 &= 0.0 \\ b_2 &= 0.0 \end{aligned}$$

The first training instance is:  $x_1 = 3, x_2 = 3, Y = 0$

Using the above equation we can plug in all of these numbers and calculate a prediction:

$$\text{prediction} = 1 / (1 + e^{-(b_0 + b_1 * x_1 + b_2 * x_2)})$$

$$\text{prediction} = 1 / (1 + e^{-(0.0 + 0.0 * 3 + 0.0 * 3)})$$

$$\text{prediction} = 0.5$$

# Logistic Regression - Step by step

X1	X2	Y
3	3	0
1	2	0
3	4	0
1	2	0
3	3	0
8	3	1
5	2	1
7	2	1
9	0	1
8	4	1

- Calculate new coefficient values based on the error in the prediction

We can calculate the new coefficient values using a simple update equation.

$$b = b + \text{alpha} * (y - \text{prediction}) * \text{prediction} * (1 - \text{prediction}) * x$$

alpha is the learning rate and controls how much the coefficients (and therefore the model) changes or learns each time it is updated. We would use a value of 0.3 for this example.

Let's update the coefficients using the prediction (0.5) and coefficient values (0.0) from the previous section.

$$b_0 = b_0 + 0.3 * (0 - 0.5) * 0.5 * (1 - 0.5) * 1$$

$$b_1 = b_1 + 0.3 * (0 - 0.5) * 0.5 * (1 - 0.5) * 3$$

$$b_2 = b_2 + 0.3 * (0 - 0.5) * 0.5 * (1 - 0.5) * 3$$

or

$$b_0 = -0.0375, b_1 = -0.1125, b_2 = -0.1125$$

# Logistic Regression - Example

**Example:** Suppose that we are interested in the factors that influence whether a politician wins an election. The outcome variable is binary (0/1); win or lose.

The predictor variables of interest are:

- amount of money spent on the campaign
- amount of time spent campaigning negatively
- whether or not the candidate is an incumbent

# Naive Bayes

- Naive Bayes is a supervised probabilistic algorithm that utilizes probability theory and Bayes' Theorem to predict probability of occurrence of an event
- The algorithm assumes that each input variables is independent which is a naive assumption in reality and hence called as Naive Bayes
- For example,
  - if you use Naive Bayes for sentiment analysis, given the sentence ‘I like Game of Thrones’, the algorithm will look at the individual words and not the full sentence
  - In a sentence, words that stand next to each other influence the meaning of each other, and the position of words in a sentence is also important
  - However, for the algorithm, phrases like ‘I like Game of Thrones’, ‘Game of Thrones like I’, and ‘Thrones I like of Game’ are the same

# Naive Bayes

- To make classifications, we need to use X to predict Y. In other words, given a data point  $X = (x_1, x_2, \dots, x_n)$ , what the odd of this sample belongs to y. This can be rewritten as the following equation:

- $P(Y = y | X = (x_1, x_2, \dots, x_k))$

- Bayes Theorem

- $P(Y | X) = \frac{(P(X | Y) * P(Y))}{P(X)} \rightarrow \text{Posterior} = (\text{likelihood} * \text{prior}) / \text{evidence}$

- We can't get  $P(Y | X)$  directly but we can get  $P(X | Y)$  and  $P(Y)$  from the data

$$\begin{aligned} P(X | Y) &= P(X_1, X_2, \dots, X_n | Y) \\ &= P(X_1 | Y) P(X_2 | Y) \dots P(X_n | Y) \end{aligned}$$

# Naive Bayes

- As the number of feature grows, it becomes complex to calculate  $P(X | Y)$ . For example with 2 features( $X_1$  &  $X_2$ ), we would need to calculate following,
  - $P(X_1 = 1, X_2 = 1 | Y = 1)$ ,  $P(X_1 = 1, X_2 = 1 | Y = 0)$
  - $P(X_1 = 1, X_2 = 0 | Y = 1)$ ,  $P(X_1 = 1, X_2 = 0 | Y = 0)$
  - $P(X_1 = 0, X_2 = 1 | Y = 1)$ ,  $P(X_1 = 0, X_2 = 1 | Y = 0)$
  - $P(X_1 = 0, X_2 = 0 | Y = 1)$ ,
  - $P(X_1 = 0, X_2 = 0 | Y = 0) \rightarrow 1 - \text{others}$
- There are  $2^{k+1} - 1$  parameters where  $k$  is the number of features
- Solving for these number of parameters is time-consuming and hence we introduce the independence assumption!
  - If  $X_1$  and  $X_2$  are independent,  $P(X_1, X_2 | Y) = \underbrace{P(X_1 | Y) * P(X_2 | Y)}$

# Naive Bayes

- Naive Bayes Classifier:
  - $P(Y | X) = P(Y) \prod_i P(X_i | Y) / P(X)$
- For continuous features, there are two possible routes:
  - Discretization - Self explanatory
  - Continuous Naive Bayes
- In continuous Naive Bayes, we first assume the probability distribution from the following
  - Gaussian
  - Multinomial
  - Bernoulli
- And then, estimate the parameters of these distributions using the data to get the PDF

Question: Naive Bayes is not suitable for datasets with large numbers of numerical attributes.

- a. True ✓
- b. False

ind false

# Naive Bayes Example

Outlook	Temperature	Humidity	Windy	Play
sunny	hot	high	false	NO
sunny	hot	high	true	NO
overcast	hot	high	false	YES
rainy	mild	high	false	YES
rainy	cool	normal	false	YES
rainy	cool	normal	true	NO
overcast	cool	normal	true	YES
sunny	mild	high	false	NO
sunny	cool	normal	false	YES
rainy	mild	normal	false	YES
sunny	mild	normal	true	YES
overcast	mild	high	true	YES
overcast	hot	normal	false	YES
rainy	mild	high	true	NO

Weather dataset, from the [University of Edinburgh](#)

- In this case,  $X = (\text{Outlook}, \text{Temperature}, \text{Humidity}, \text{Windy})$ , and  $Y = \text{Play}$ .
  - Q: today = (sunny, hot, normal, false)
  - A:  $P(\text{play} | \text{today})$  and  $P(\text{Not play} | \text{today})$ ?
- $P(\text{play} | \text{today})$ 
$$= P(\text{today} | \text{play}) * P(\text{play}) / P(\text{today})$$
$$= \underline{P(\text{sunny} | \text{play})} * P(\text{hot} | \text{play}) * P(\text{normal} | \text{play}) * P(\text{false} | \text{play}) * P(\text{play}) / P(\text{today})$$
$$= \cancel{2/9} * 2/9 * 6/9 * 6/9 * 9/14 / P(\text{today})$$
$$= 0.0141 / \underline{P(\text{today})}$$
- $P(\text{Not play} | \text{today})$ 
$$= 3/5 * 2/5 * 1/5 * 2/5 * 5/14 / P(\text{today})$$
$$= 0.0068 / P(\text{today})$$

# Naive Bayes Example

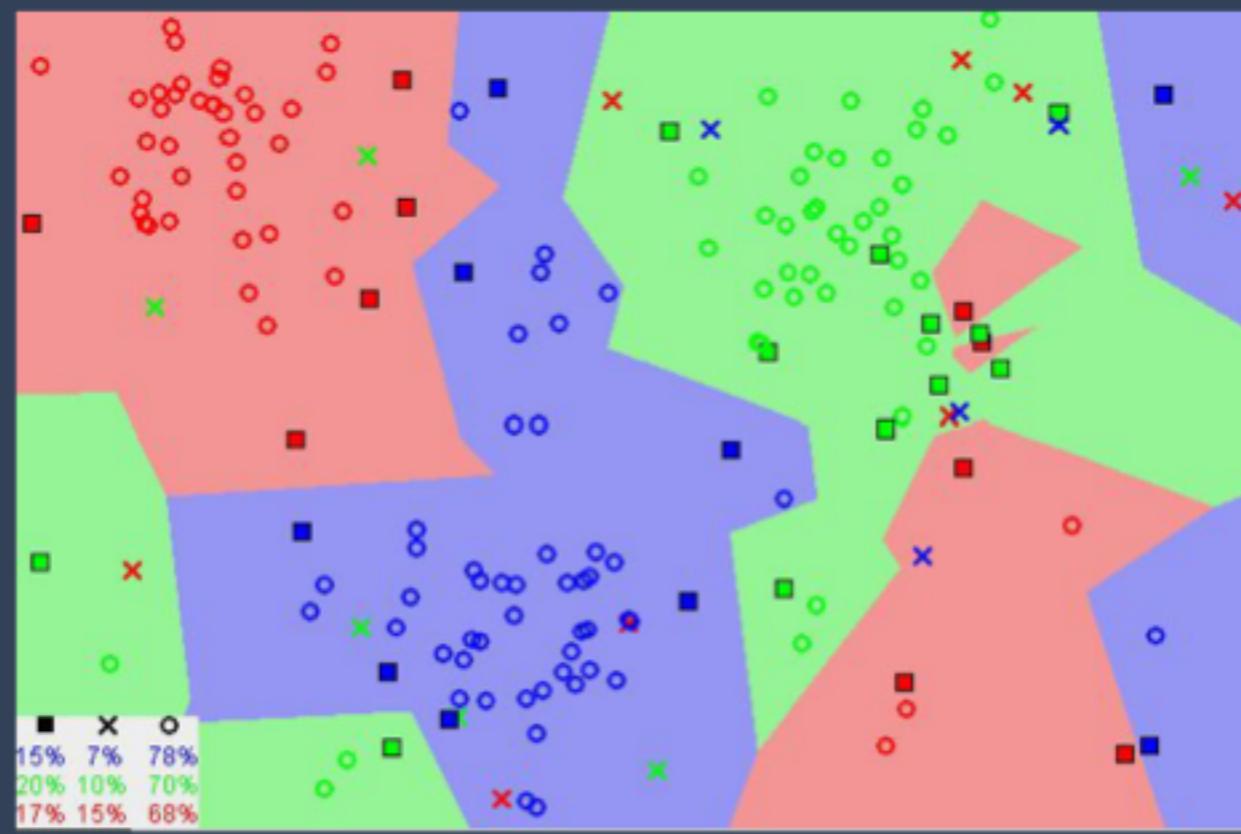
- $P(\text{play} | \text{today})$   
=  $0.0141 / P(\text{today})$
- $P(\text{Not play} | \text{today})$   
=  $0.0068 / P(\text{today})$

Conclusion:

- $P(\text{play} | \text{today}) > P(\text{Not play} | \text{today})$
- $P(\text{play} | \text{today}) = 0.67, P(\text{Not play} | \text{today}) = 0.33$  using the fact  $P(\text{play} | \text{today}) + P(\text{Not play} | \text{today}) = 1$

# k - Nearest Neighbors (kNN)

- kNN is a supervised machine learning algorithm used in both classification and regression problems
- The basic assumption here is that similar items appear in close **proximity** in a n-dimensional space



# k - Nearest Neighbors (kNN)

- One of the most important aspects while using a k-nn algorithm is to define the distance metric for calculating the neighbors such as
  - Euclidean
  - Cosine
- The choice of distance metric is a hyperparameter and depends on the type of dataset
- The number of neighbors  $k$  is also an hyperparameter which needs to be identified using the validation dataset

Question: What are the hyperparameters in kNN based modeling?

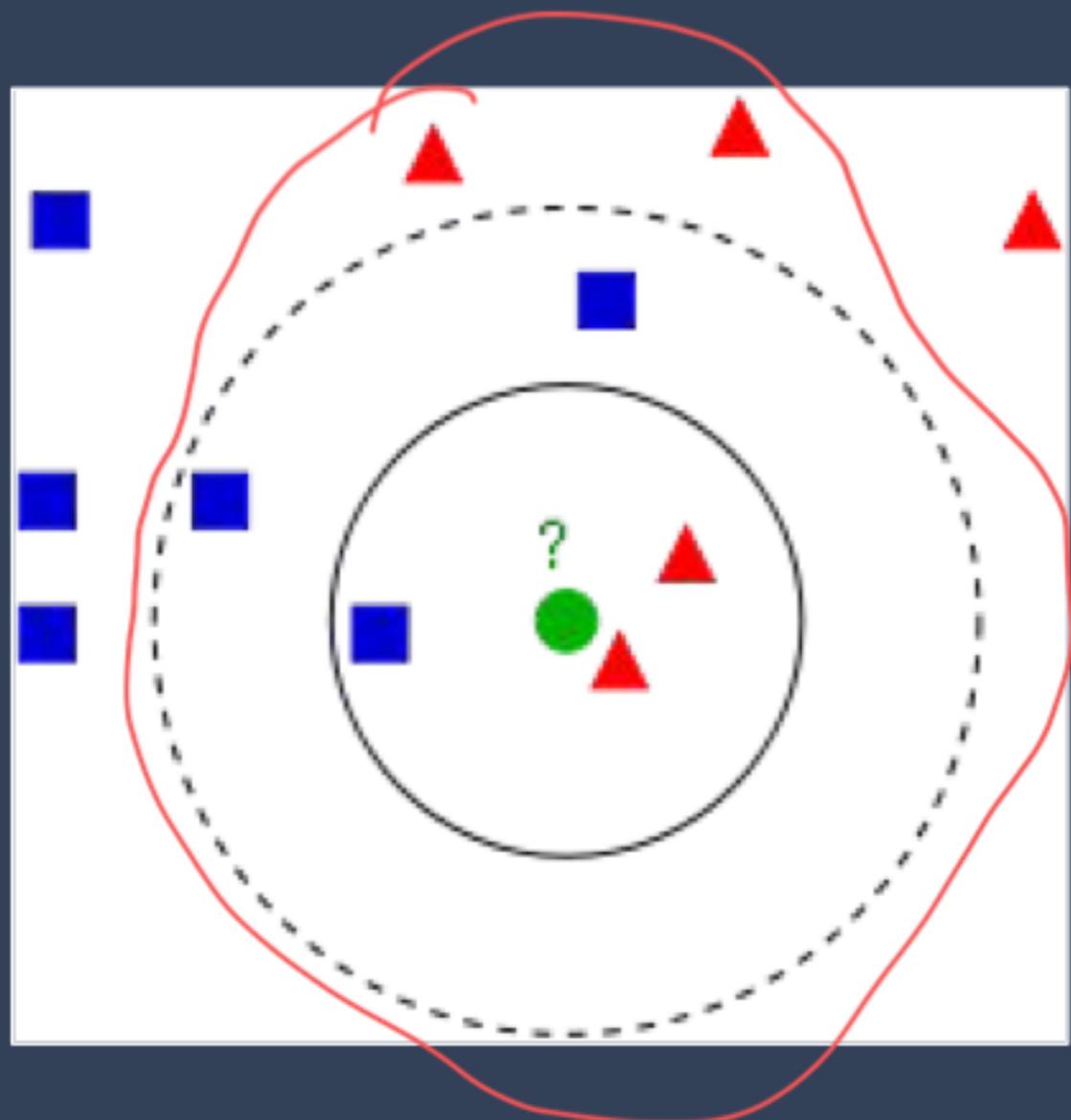
- a. Distance function
- b.  $k \rightarrow$  number of neighbors
- c. Both a & b ✓

# k - Nearest Neighbors (kNN)

## Algorithm

1. Set the value of k and choose the distance measure
2. For each data point
  - a. Calculate the distance between the input data point and the current data point from the data
  - b. Add the distance and the index of the example to an ordered collection
3. Sort the ordered collection of distances and indices in ascending order by the distances
4. Pick the first K entries from the sorted collection and the corresponding labels
5. In case of classification problem, could return the mode of the K labels
6. In case of regression problem, could return the mean of the K values

# k - Nearest Neighbors (kNN)



Which group (red, blue) the green dot belongs to?

$$k = 3$$

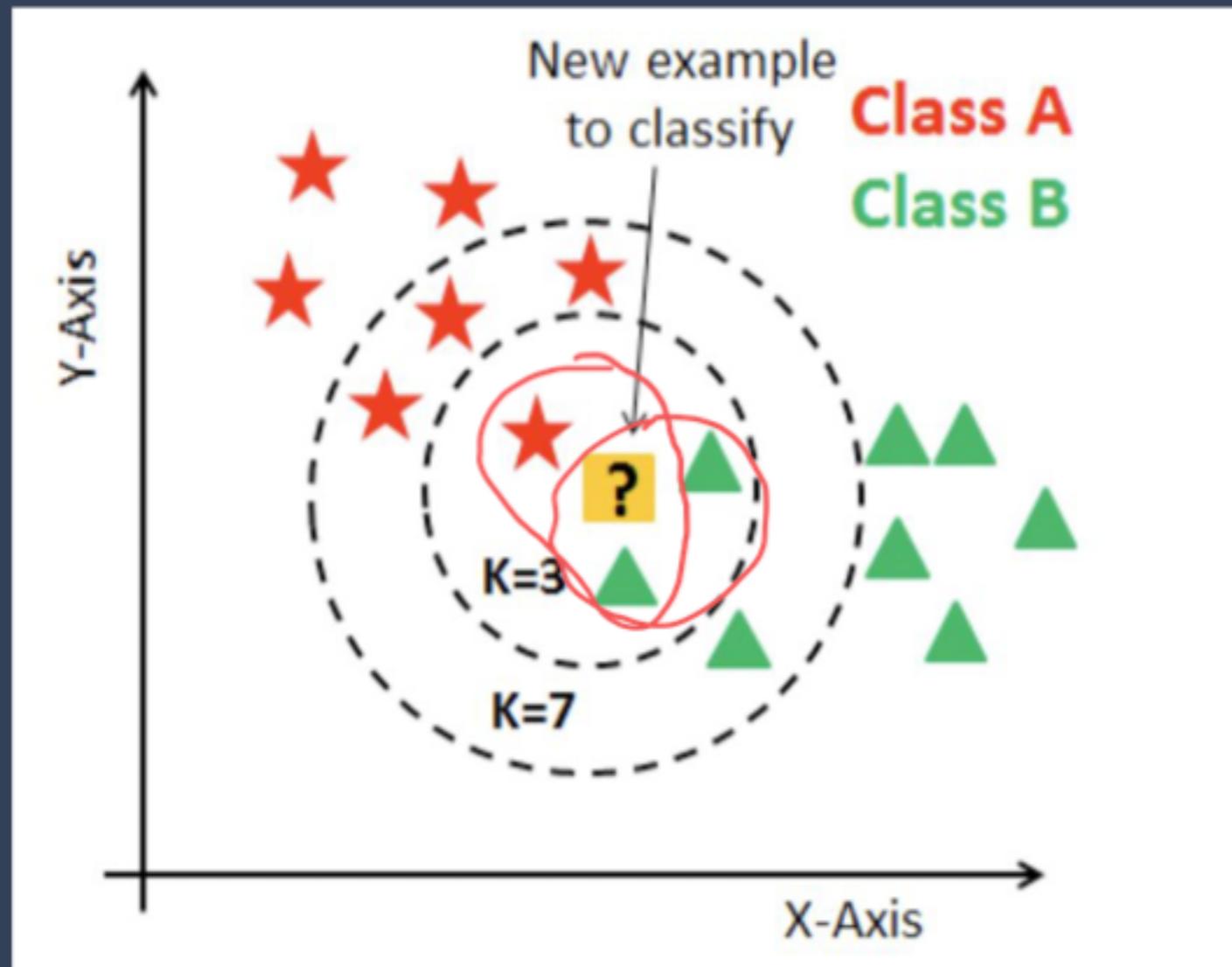
red

$$k = 5$$

blue

$$k = 7$$

$K = 0$

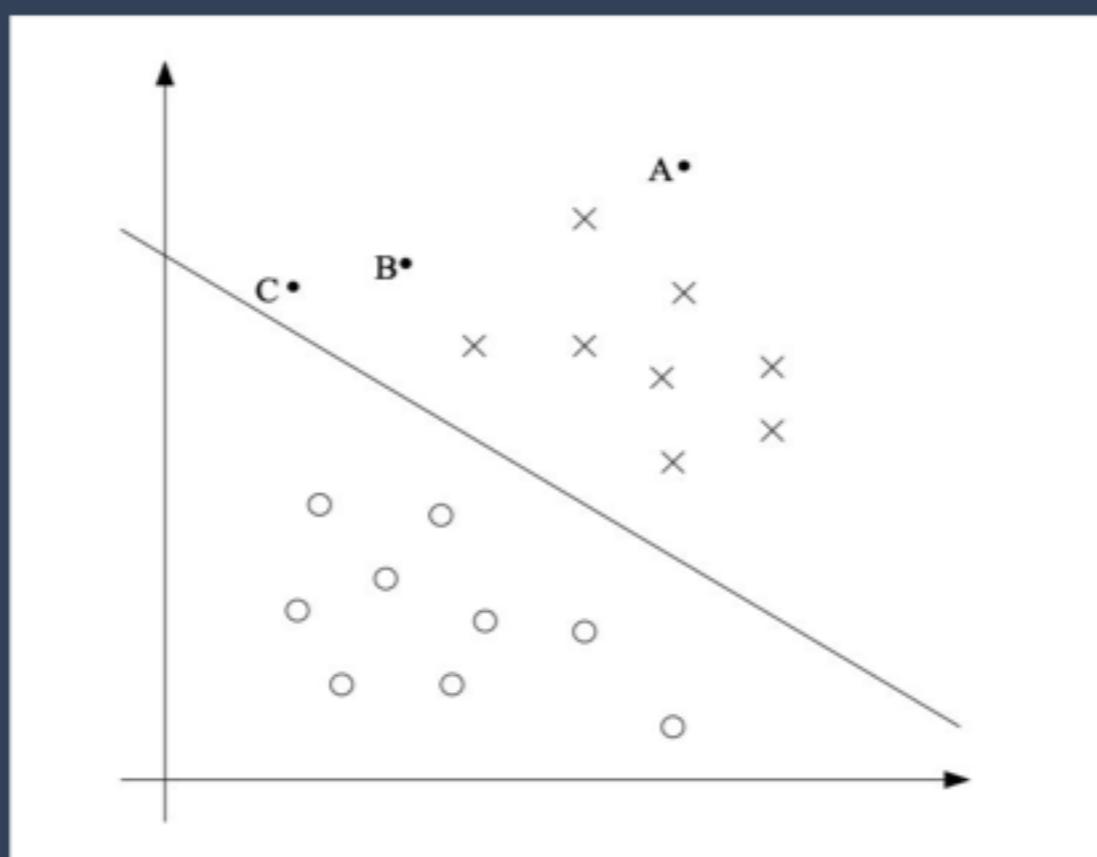


Question: In the above figure, which Class would Yellow(?) belong to when  $k=2$ ?

- a. Class A
- b. Class B

# SVM - Support Vector Machine

- SVM is a binary classification algorithm which falls in the category of maximum-margin classifier
- It is a discriminative model formally defined by a separating hyperplane. In 2-dimensional space, it is just called line.

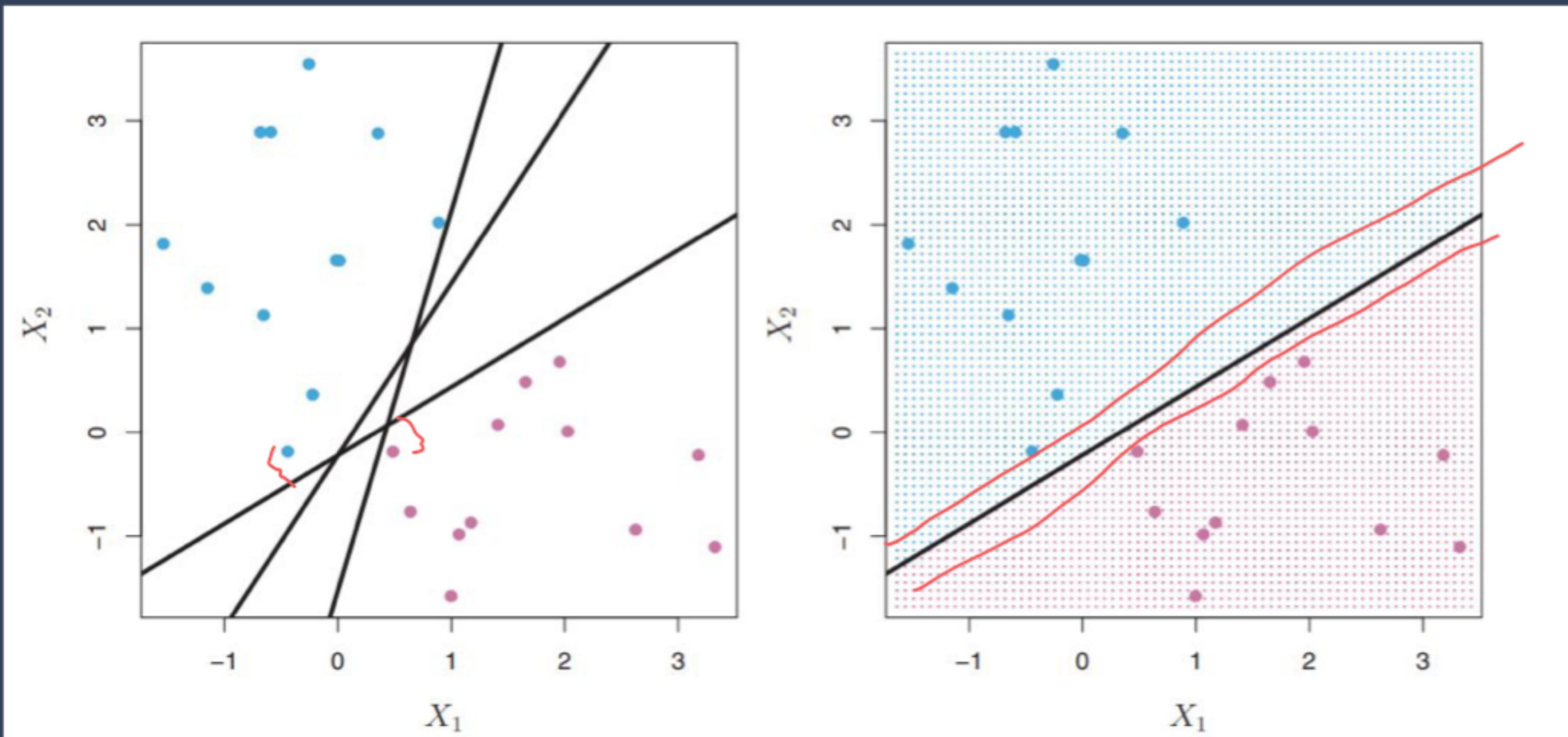


# SVM - Support Vector Machine

- The decision boundary can be represented using  $w^T \cdot x + b$
- Intuitively, the further a point is from the hyperplane on either side, the more confident we are about a prediction.
- The hyperplane  $w^T \cdot x + b$  then separates classes
- Hence, we try to create the hyperplane in such a manner that margins of the points from the separating hyperplane are as high as possible

svm maximizes the distances to the closest points  
on both sides

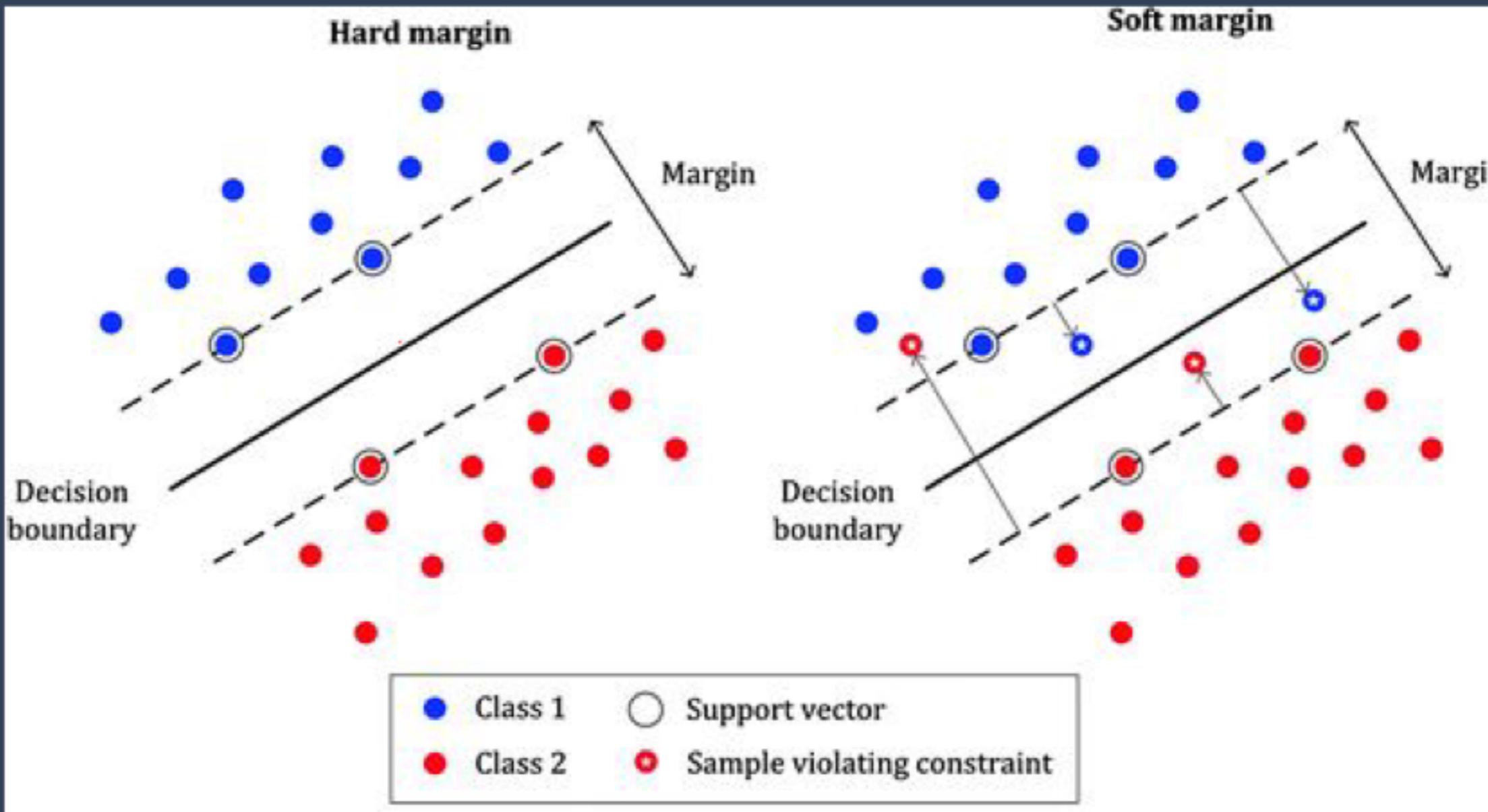
## SVM - Support Vector Machine



# SVM - Support Vector Machine

- Notation
  - Class label  $y \in \{-1, 1\}$
  - $h_{w,b}(x) = g(w^T \cdot x + b)$ , note here we separate out  $\theta$  as  $[w, b]$
  - Here,  $g(z) = 1$  if  $z \geq 0$ , and  $g(z) = -1$ , otherwise
- Functional margin
  - $\gamma^{(i)} = y^{(i)} (w^T \cdot x_i + b)$
  - if  $y^{(i)} = 1$ , then for the functional margin to be large we need  $w^T \cdot x^{(i)} + b$  to be a large positive number
  - if  $y^{(i)} = -1$ , then for the functional margin to be large we need  $w^T \cdot x^{(i)} + b$  to be a large negative number

# SVM - Support Vector Machine



# SVM - Support Vector Machine

- The magnitude of the functional margin depends on the magnitude of the parameters which is not very meaningful
  - Replacing  $(w, b)$  with  $(2w, 2b)$  results in multiplying our functional margin by a factor of 2 without changing the sign (prediction results)
  - We can add an arbitrary scaling constraints to  $(w, b)$  without changing the prediction sign
  - By exploiting our freedom to scale  $w$  and  $b$ , we can make the functional margin arbitrarily large without really changing anything meaningful
- Hence, it makes sense to impose some sort of normalization condition such as that  $\|w\|^2 = 1$

# SVM - Support Vector Machine



- Given a training set  $S = \{(x^{(i)}, y^{(i)}) ; i = 1, \dots, m\}$ , we also define the function margin of  $(w, b)$  with respect to  $S$  as the smallest of the functional margins of the individual training examples.

$$\circ \frac{\gamma^* = \min_{i=1, \dots, m} \gamma^{*(i)}}{w^T x_i + b}$$

$$\begin{pmatrix} w, b \\ y^{(i)}(w^T x_i + b) \end{pmatrix}$$

$$\hat{\gamma}(w, b)$$

- The optimization problem of SVM can now be defined as:

$$\circ \max_{w, b} \gamma^*(w, b)$$

$$\quad \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \gamma^*, i = 1, \dots, m$$

$$\quad \|w\| = 1$$

$$\hat{y} \in \{1, -1\}$$

$$\gamma(w, b)$$

# SVM - Support Vector Machine

- But the above formulation with  $\|w\|=1$  is non-convex and hence difficult to solve!
- What if we transform it as follows:
  - $\max_{\gamma, w, b} \gamma^{\wedge} / \|w\|$
  - s.t.  $y^{(i)} (w^T x^{(i)} + b) \geq \gamma^{\wedge}, i = 1, \dots, m$
- The rational is that
  - $\gamma^{\wedge(i)} = y^{(i)} ( (w / \|w\|)^T x^{(i)} + b / \|w\| )$  - geometric margin
  - If  $\|w\| = 1$ , then functional margin is same as geometric margin
- The geometric margin is invariant to rescaling of the parameters

# SVM -/ Support Vector Machine

- Recall: We can add an arbitrary scaling constraints to  $(w, b)$  without changing the prediction sign.
- Here, we introduce the scaling constraint that the functional margin of  $w, b$  with respect to the training set must be 1
  - $\gamma^* = 1$
- Then we can transform the following as
  - ~~$\max_{\gamma, w, b} \gamma^*/||w||$~~
  - ~~$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \gamma^*, i = 1, \dots, m$~~
  - $\left\{ \begin{array}{l} \max_{w, b} 1/||w|| \\ \text{s.t. } \underbrace{y^{(i)}(w^T x^{(i)} + b)}_{\geq 1}, i = 1, \dots, m \end{array} \right.$

# SVM - Support Vector Machine

- Also, maximizing  $1/\|w\|$  is same as minimizing  $\|w\|^2$ . Thus,
  - $\min_{w,b} 1/2 * \|w\|^2$ 
    - s.t.  $y^{(i)} (w^T x^{(i)} + b) \geq 1, i = 1, \dots, m$
- The above is an optimization problem with a convex quadratic objective and only linear constraints. Its solution gives us the *optimal margin classifier*
- The Lagrangian of the above equation would help us to solve it efficiently

# SVM - Support Vector Machine

## Lagrangian

- Consider a problem of the following form:
  - $\min_w f(w)$ 
    - s.t.  $h_i(w) = 0, i = 1, \dots, l$
- We define the *Lagrangian* to be
  - $L(w, \beta) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$
  - Here, the  $\beta_i$ 's are called the *Lagrange multipliers*
- We would then find and set  $L$ 's partial derivatives to zero and solve for  $w$  and  $\beta$

# SVM - Support Vector Machine

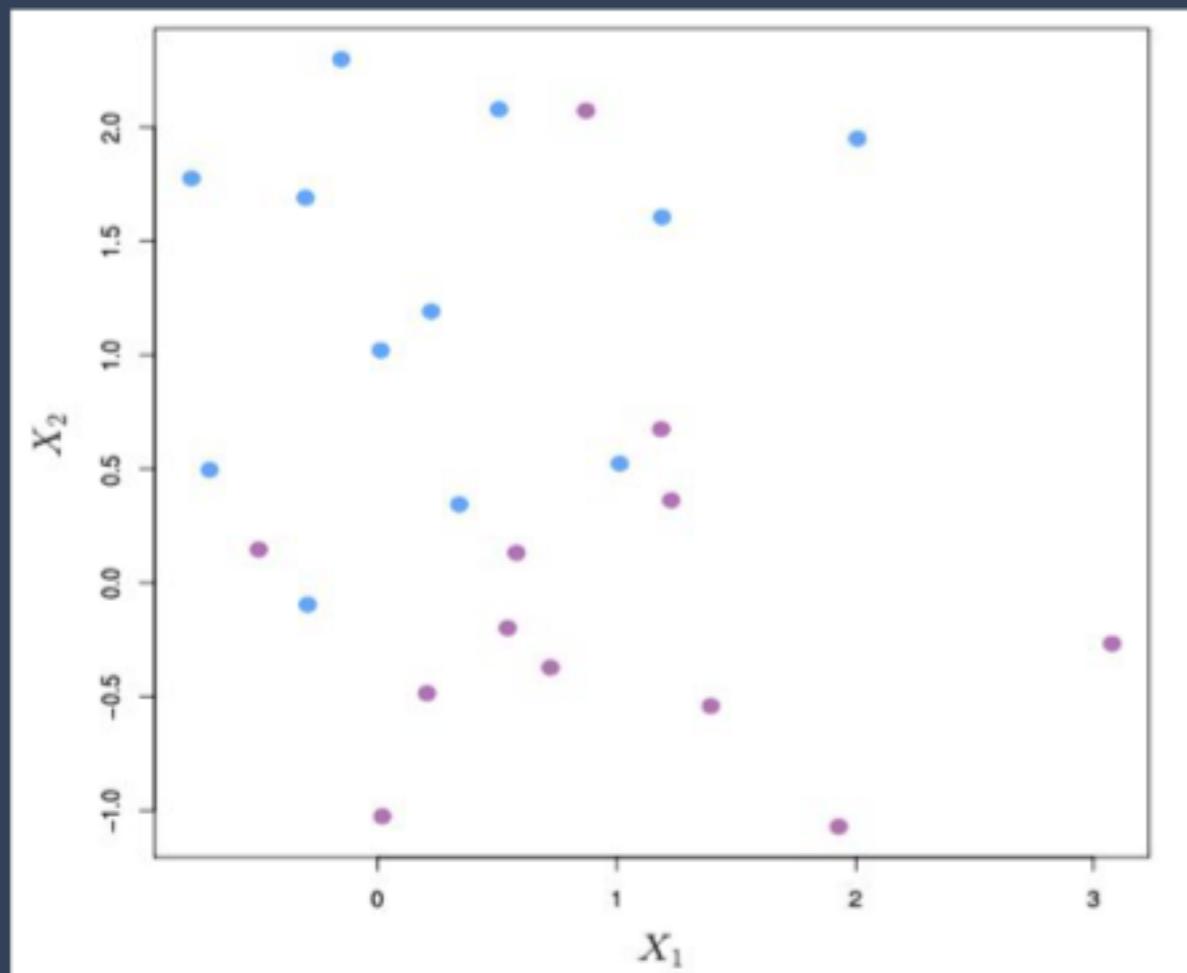
- Now remember our original optimization problem
  - $\min_{w,b} \frac{1}{2} * \|w\|^2$ 
    - s.t.,  $y^{(i)} (w^T x^{(i)} + b) \geq 1, i = 1, \dots, m$
- We can write the constraints as
  - $g_i(w) = -y^{(i)} (w^T x^{(i)} + b) + 1 \leq 0$
- When we construct the Lagrangian for our optimization problem we have:
  - $L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y^{(i)} (w^T x^{(i)} + b) - 1), \alpha_i \geq 0$
  - Taking the derivative of the above equation would solve for our optimal parameters

# SVM - Support Vector Machine

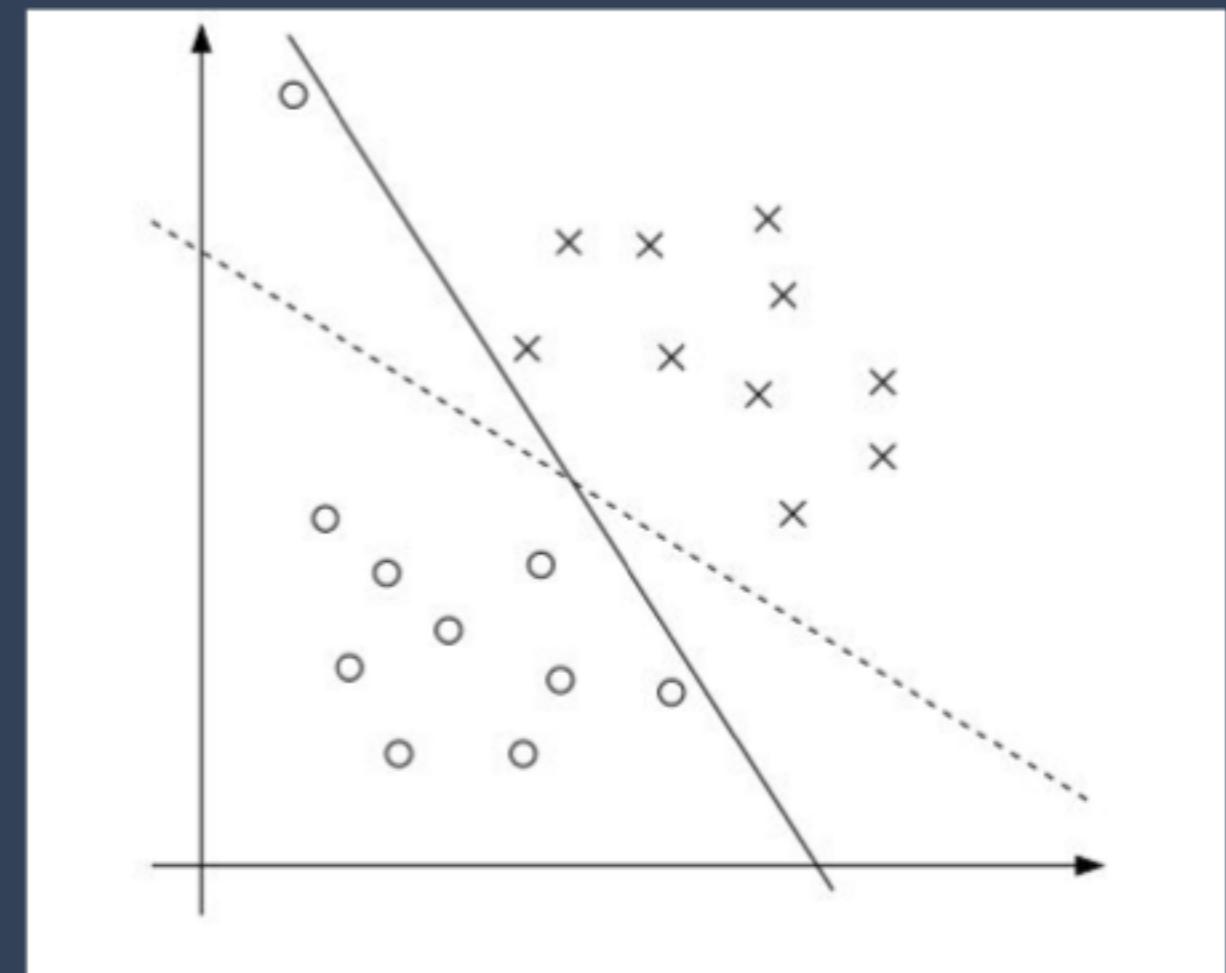
- Recall in SVM prediction, we use the sign of  $\theta^T \cdot x_{\text{new}}$  to decide the class for  $x_{\text{new}}$ 
  - $\theta^T \cdot x_{\text{new}} = w^T \cdot x_{\text{new}} + b$
  - $w^T = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$
  - $w^T \cdot x_{\text{new}} + b = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \cdot x_{\text{new}} + b$

# SVM - Support Vector Machine

Not separable by a hyperplane



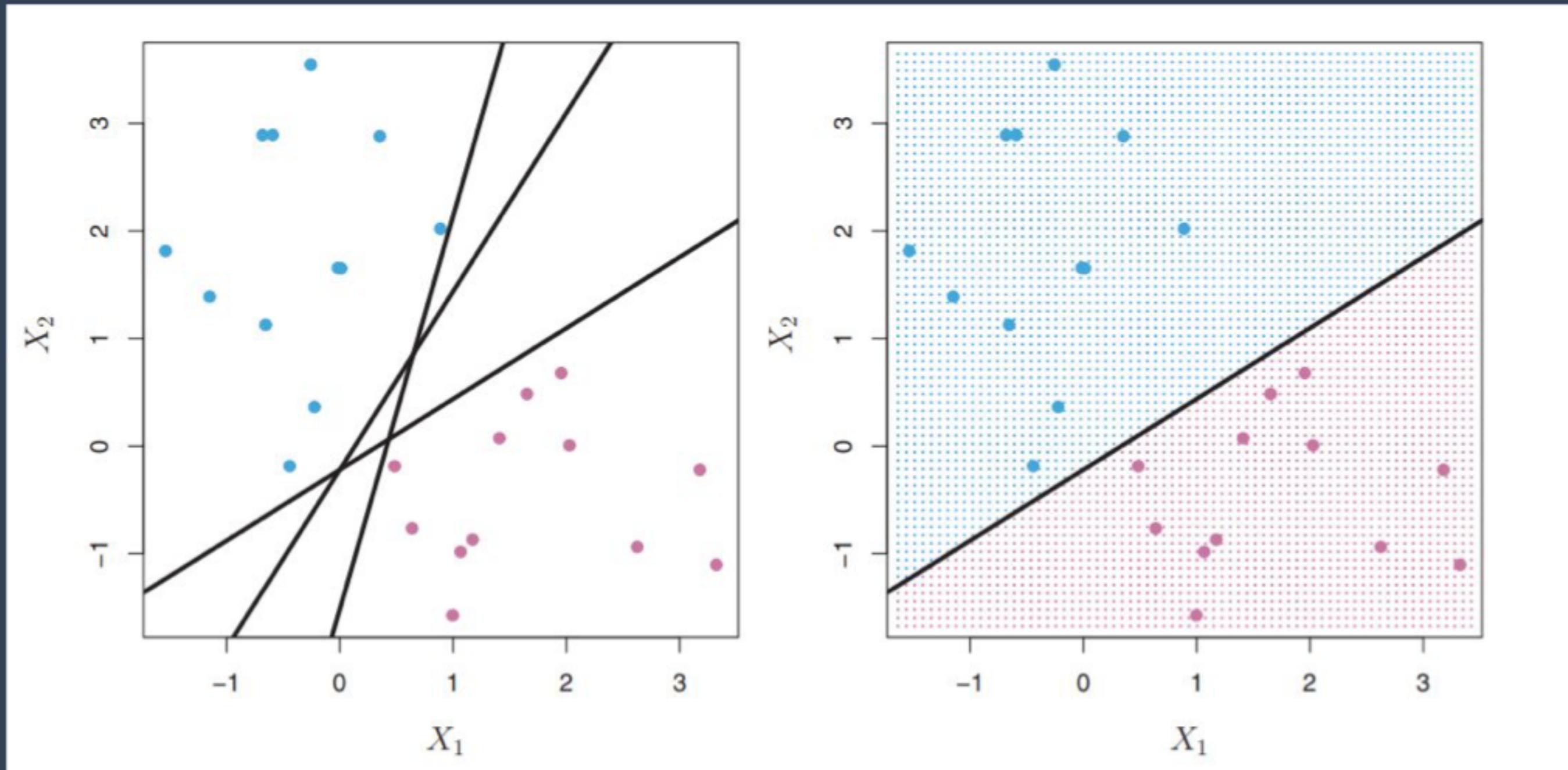
Which hyperplane is better?



# SVM - Soft Margin

- The derivation of the SVM as presented so far assumed that the data is linearly separable
- To make the algorithm work for non-linearly separable datasets as well as be less sensitive to outliers, we re-formulate our optimization (using L1 regularization) as follows:
  - $\min_{\gamma, w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \chi_i$ 
    - s.t.  $y^{(i)} (w^T x^{(i)} + b) \geq 1 - \chi_i, i = 1, \dots, m$
- The parameter  $C$  controls the relative weighting between the twin goals of making the  $\|w\|^2$  large and of ensuring that most examples have functional margin at least 1

# SVM - Support Vector Machine



# SVM - Kernel Trick

Fig.3

Linearly Separable

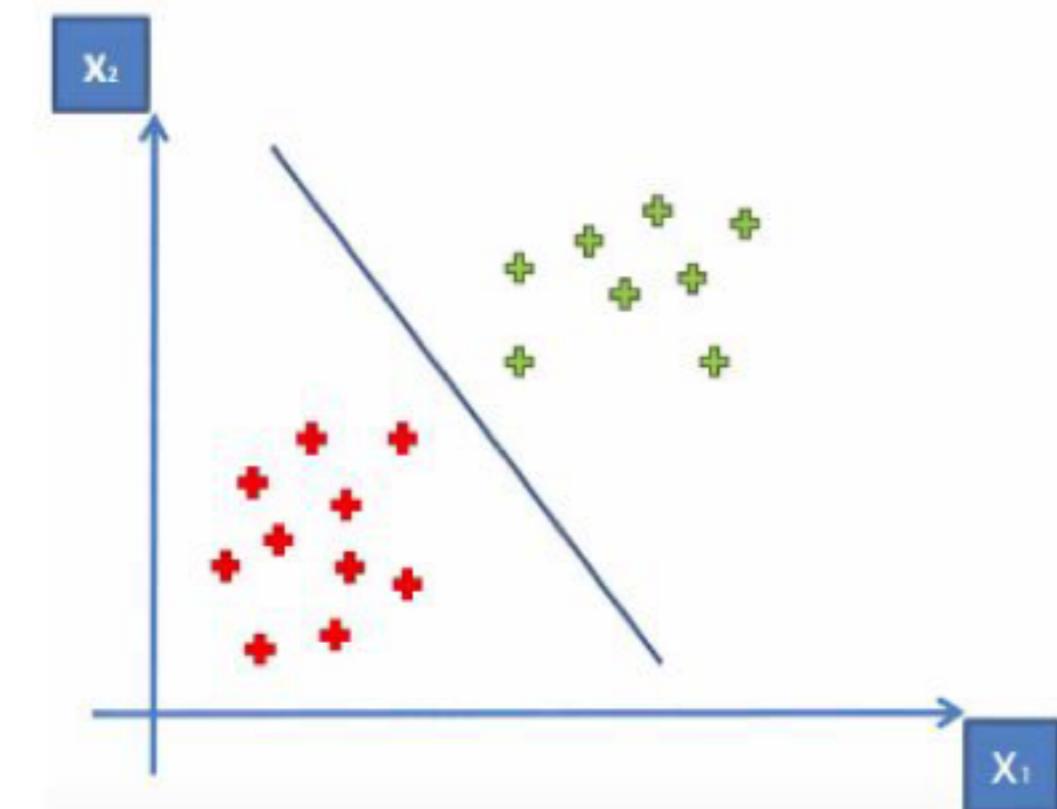
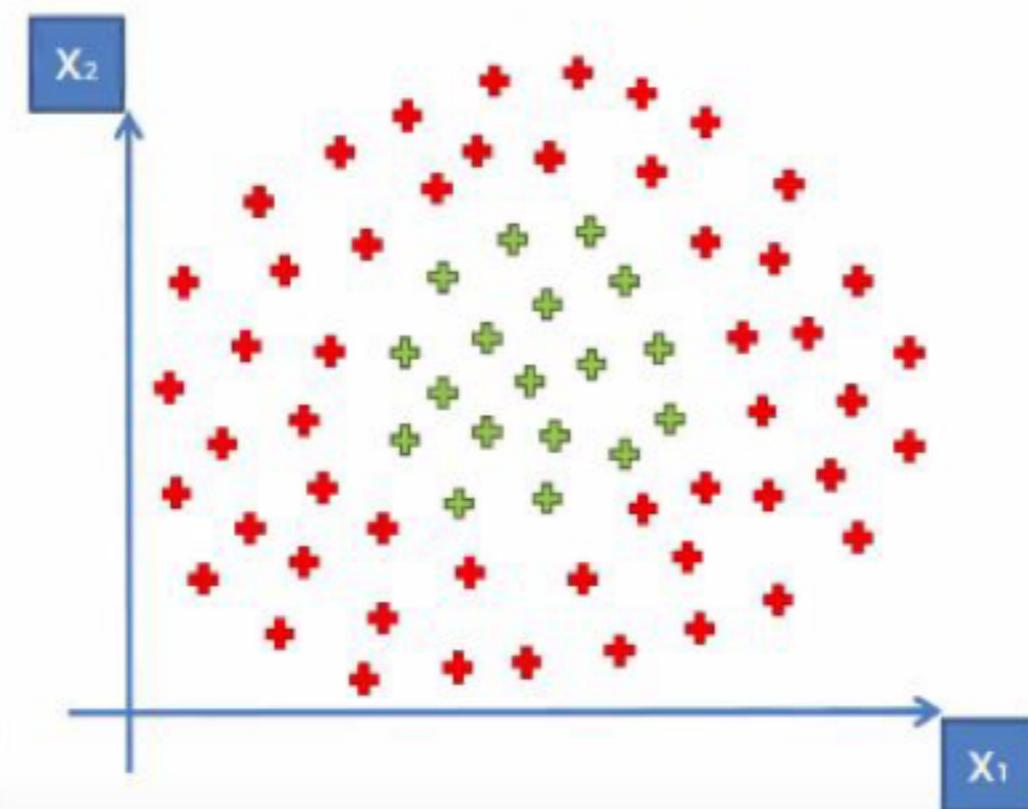


Fig.4

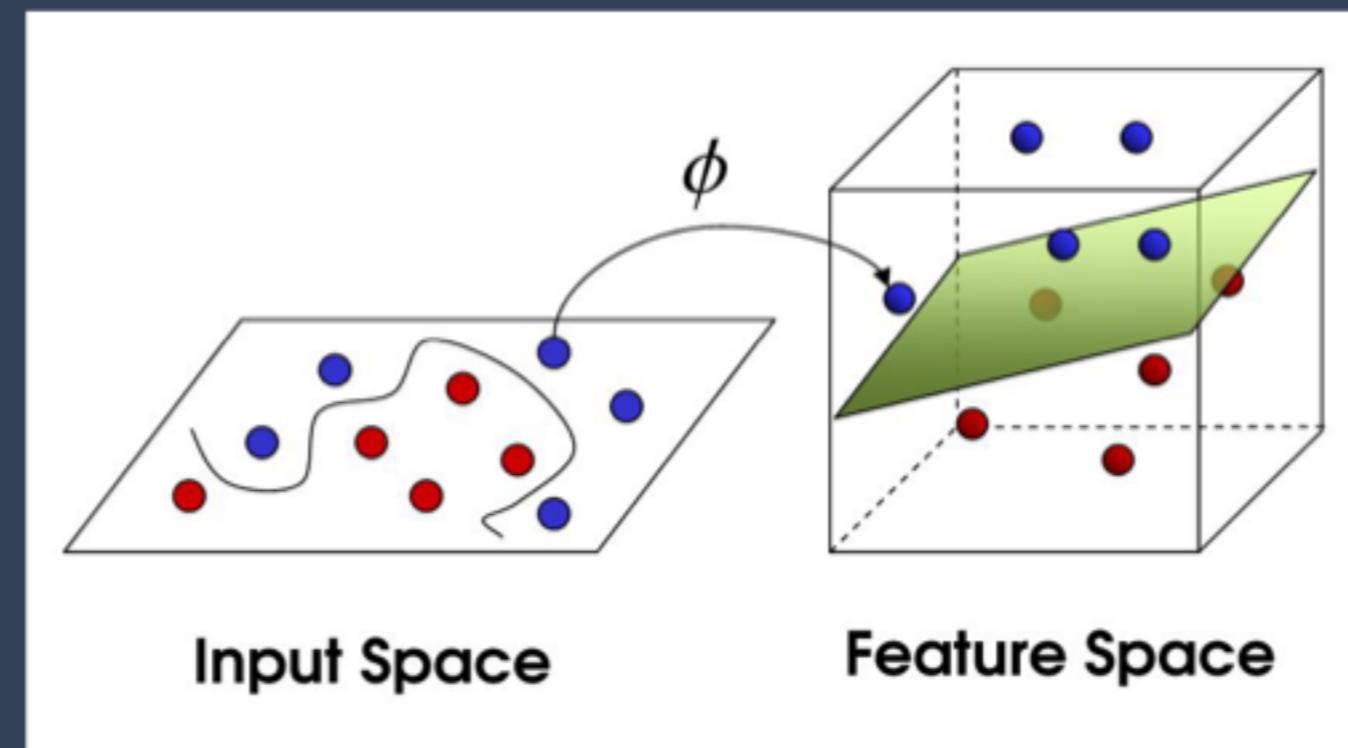
Not Linearly Separable





# SVM - Kernel Trick

- The function of kernel is to take data as input and transform it into the required form, similar to *feature mapping*, we transform  $x$  to  $[x \ x^2 \ x^3]^T$
- SVM algorithms use a set of mathematical functions that are defined as the kernel

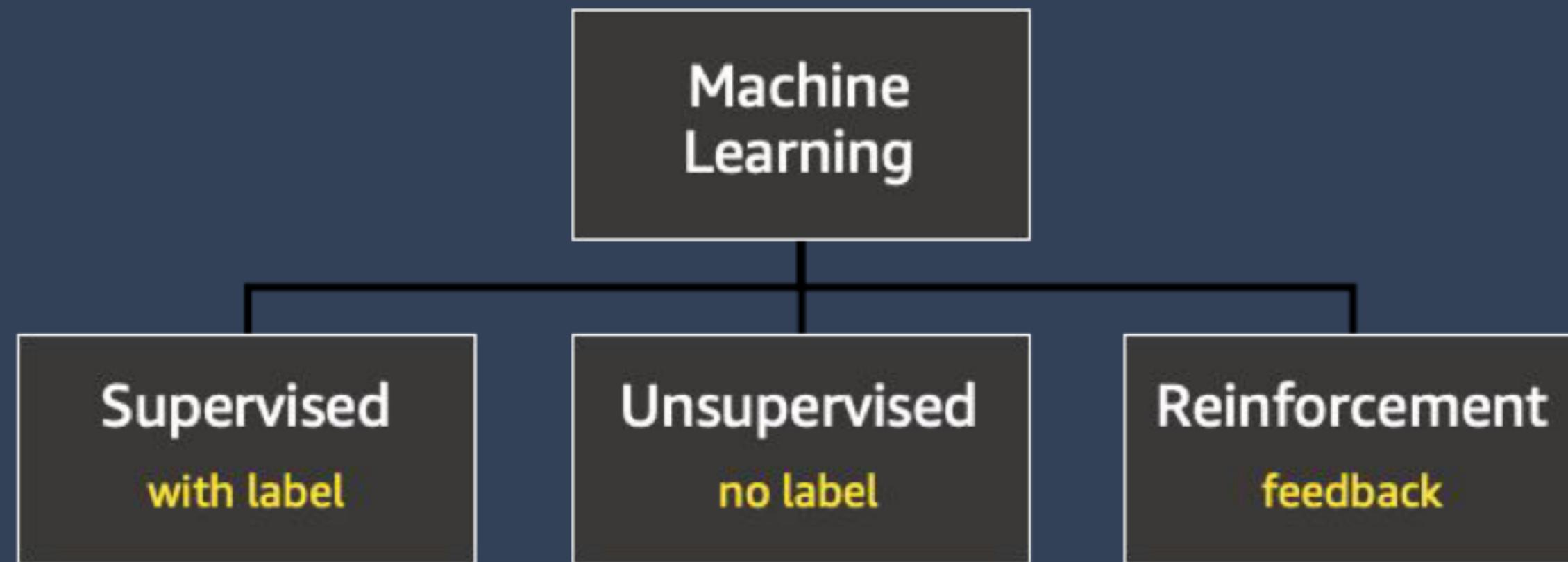


# SVM - In a nutshell

1. What is hyperplane?
  - a. If we have  $p$ -dimensional space, a hyperplane is a flat subspace with dimension  $p-1$  which separates point in the multi-dimensional space
2. If we divide data based on hyperplane, then there could be many possible hyperplanes to divide same set of numbers. How do we decide the finite number of hyperplanes to choose from?
  - a. Maximal margin classifier
3. What if the points are non-separable?
  - a. Soft-margin classifier
4. There are cases where we cannot draw a boundary that can separate two classes in linear fashion. What to do then?
  - a. Kernel transformation

# Anomaly Detection: A Flavor of Unsupervised Learning

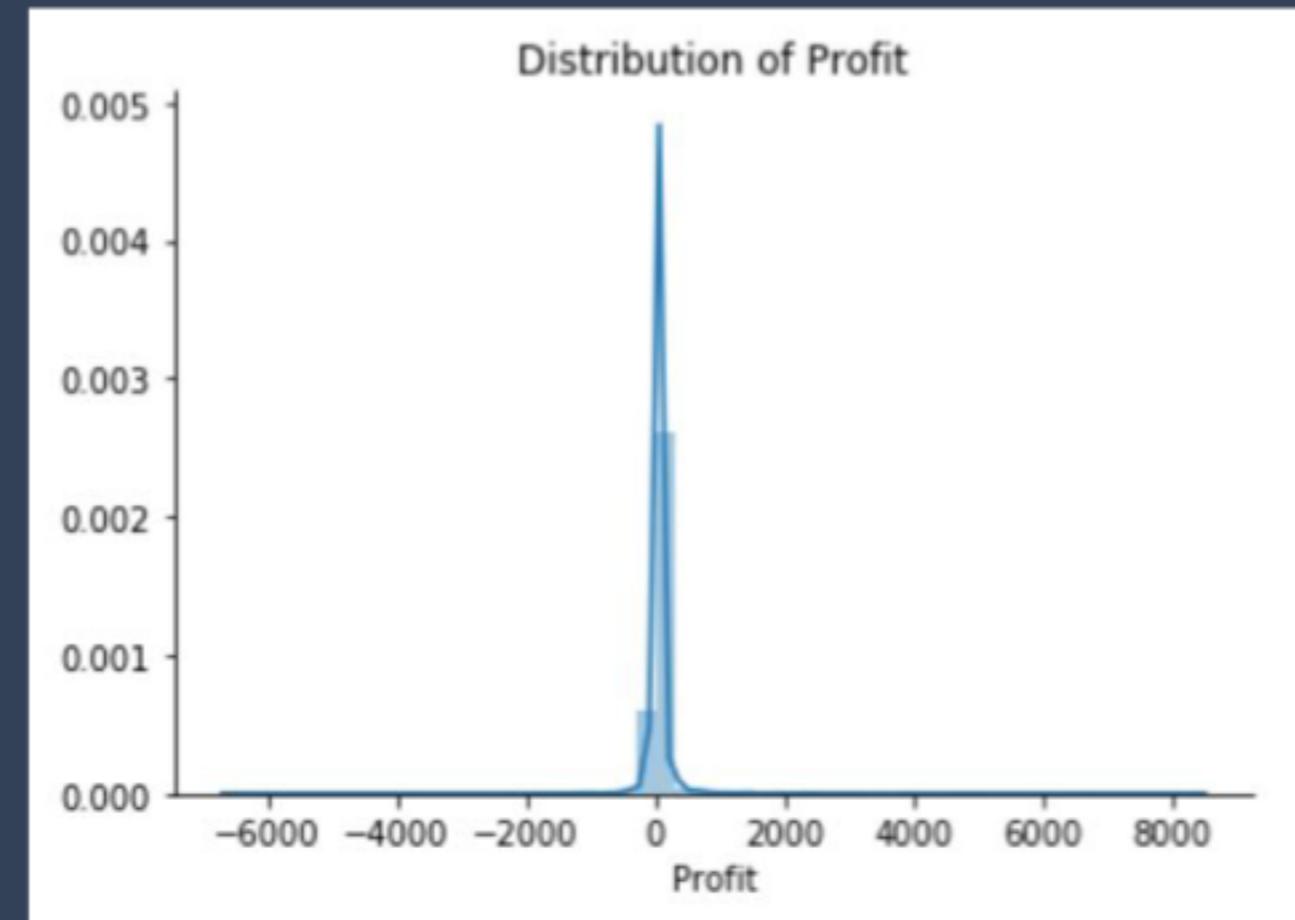
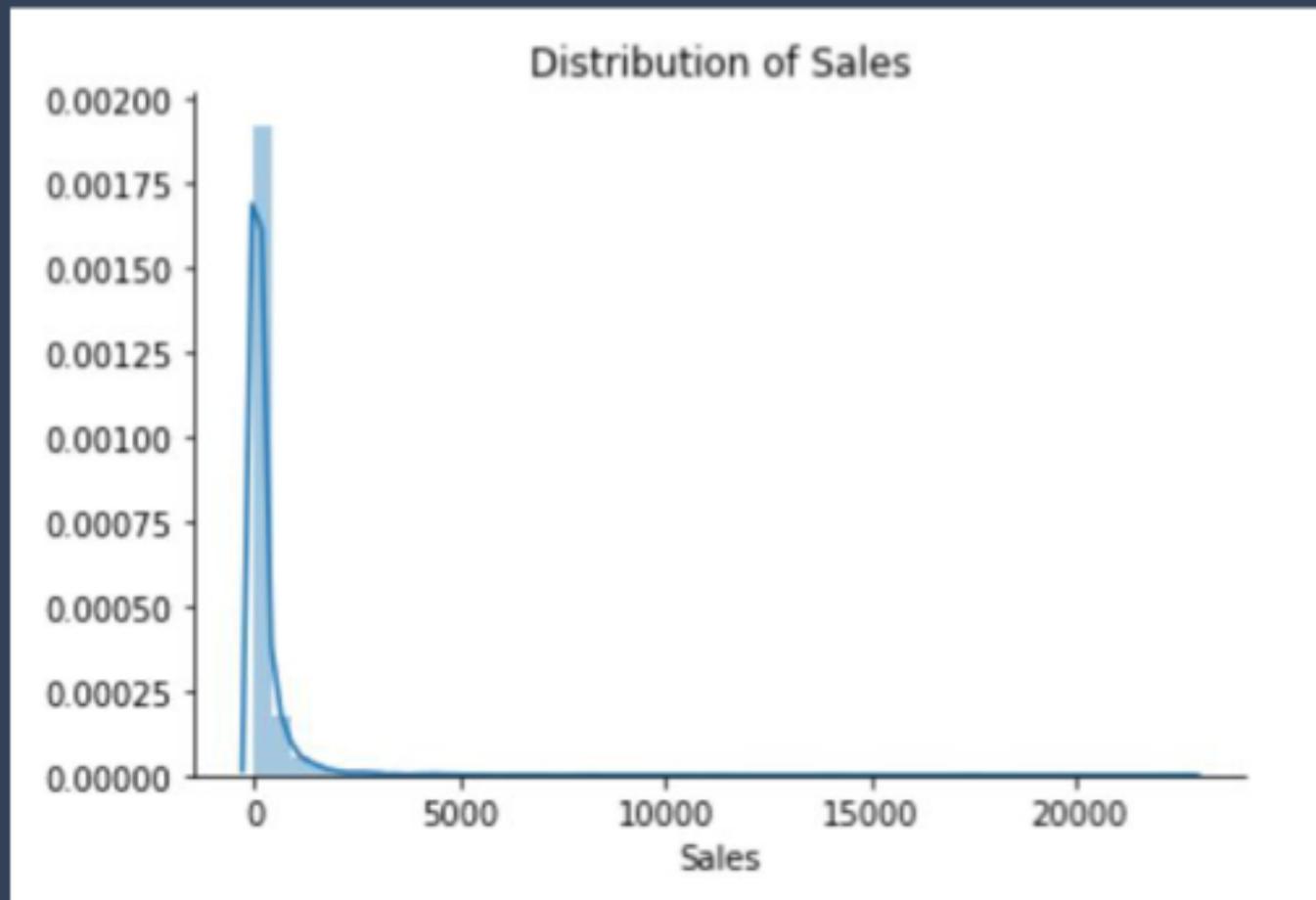
# Family of Machine Learning



# Anomaly Detection

- Anomaly detection is the process of identifying abnormal or unexpected events in datasets
- It is often applied on datasets where there are no-labels and hence it falls in the bucket of unsupervised learning
  - However, you can always train a supervised model on a set of correctly identified anomalies
- Anomaly detection can be accomplished by the following algorithms:
  - IQR
  - Isolation Forest
  - etc.

# Anomaly Detection



- In the figures above, there are values on left & right extremes of the mean density which are rare in occurrence and hence often referred to as **outliers**

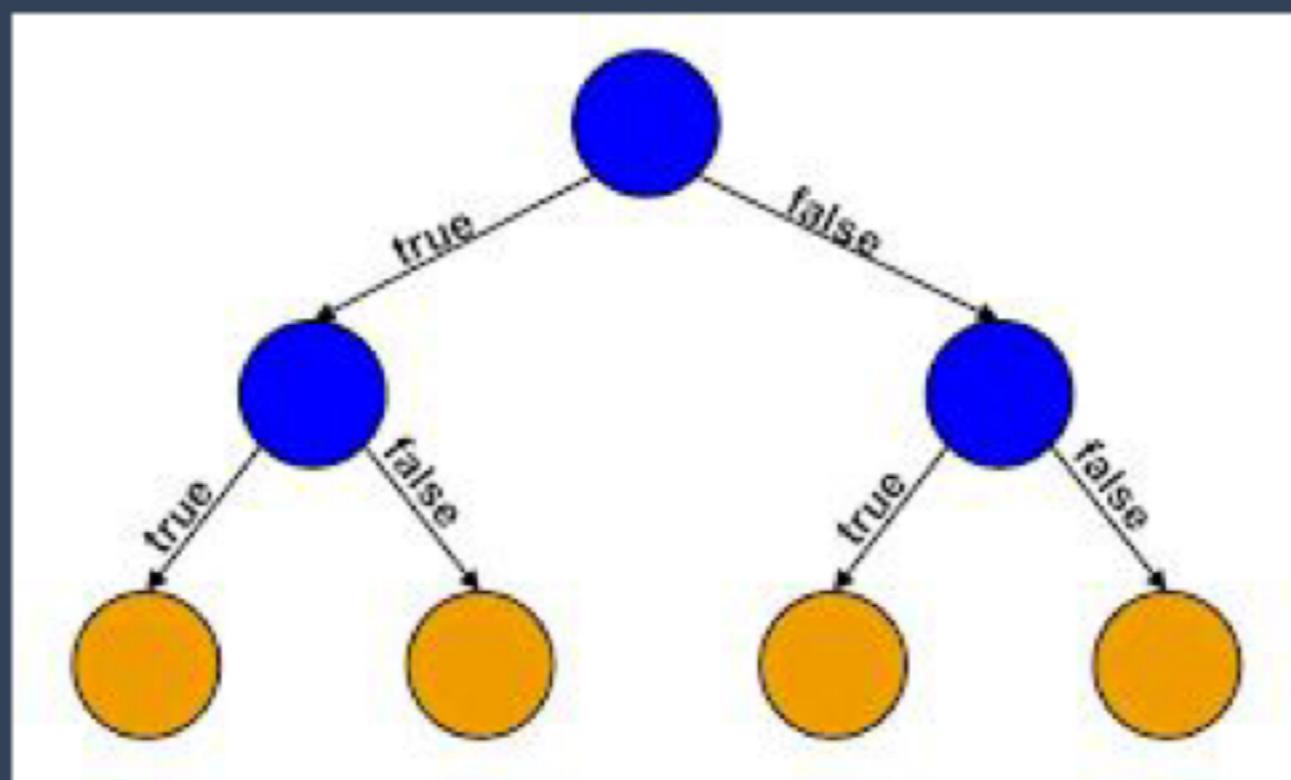
# Anomaly Detection

## IQR - Interquartile range

- The interquartile range defines the difference between the 3<sup>rd</sup>(Q3) and the 1<sup>st</sup>(Q1) quartile
  - $IQR = Q3 - Q1$
- Any value is referred to as anomaly if it satisfies any of the following 2 conditions
  - $> 1.5 * IQR$
  - $< 1.5 * IQR$
- IQR is used to identify univariate anomalies

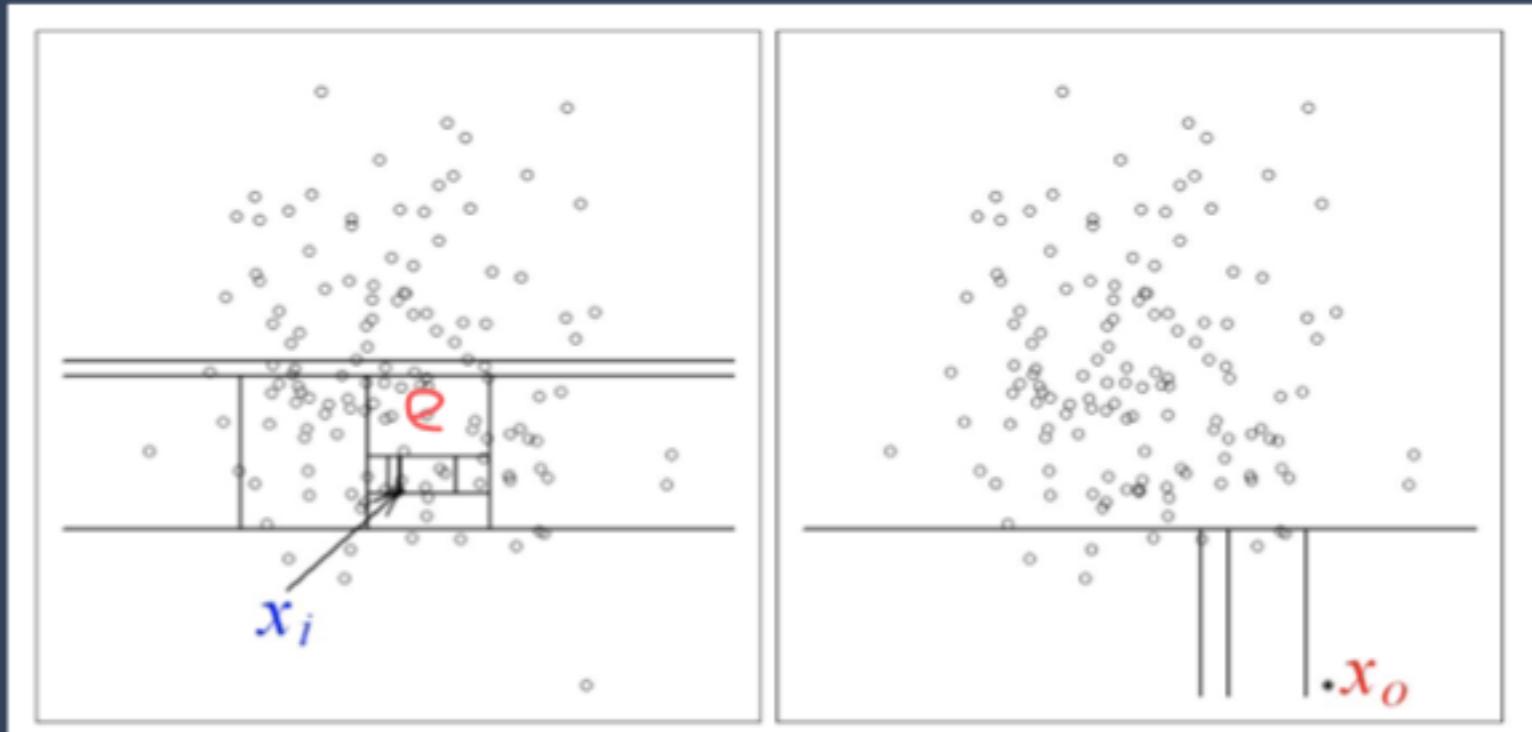
# Anomaly Detection - Isolation Forest

- Isolation Forest is a tree based *ensemble* algorithm to identify anomalies
- In these *trees*, partitions are created by first *randomly* selecting a feature and then selecting a *random* split value between the minimum and maximum value of the selected feature

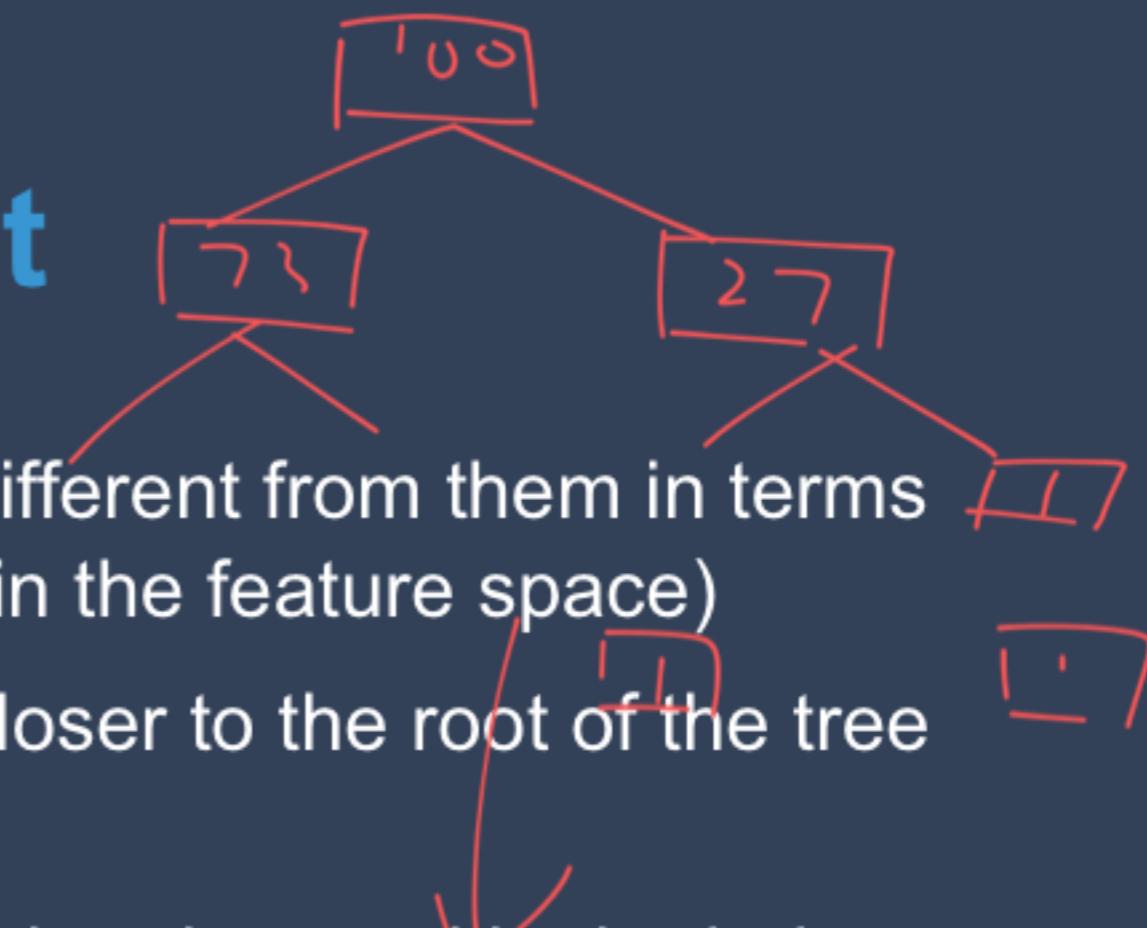


# Anomaly Detection - Isolation Forest

- Outliers are less frequent than regular observations and are different from them in terms of values (they lie further away from the regular observations in the feature space)
- By using such random partitioning, they should be identified closer to the root of the tree (shorter average path length), with fewer splits necessary
- The idea of identifying a normal vs. abnormal observation can be observed in the below figure [from the original paper].



A normal point (on the left) requires more partitions to be identified than an abnormal point (right).



# Anomaly Detection - Isolation Forest

- Anomaly score in case of Isolation Forest is defined as:
  - $s(x, n) = 2^{-E(h(x)) / c(n)}$ 
    - $h(x)$  is the path length of observation  $x$
    - $c(n)$  is a normalized constant for a data set of size  $n$ . Typically, we use the average path length of all samples. This has an approximation as  $\underline{2^*(\ln(n) + 0.577) - 2(n-1)/n}$ .
- Following decision can be made on the basis of anomaly score:
  - A score close to 1 indicates anomalies
    - from above formula, as  $E(h(x))$  goes to 0, anomaly scores go to 1
  - Score much smaller than 0.5 indicates normal observations

# Anomaly Detection

## Points to note

- Anomaly detection can be done for both Univariate and Multivariate datasets
- Anomaly detection is also used in following areas:
  - Data quality monitoring
  - KPI measurement
  - Systems with long lagging feedback period
- There is an inherent assumption in many algorithms that anomalies are rare and typically present in less than 5% of the population