

MACHINE LEARNING OVERVIEW

GETTING THE MOST OUT OF YOUR DATA

Purushottam Darshankar

AGENDA

- ✓ Introduction to Machine learning and its application
- ✓ Overview of Specific methods
 - Supervised
 - Regression
 - Classification
 - Unsupervised
 - Clustering / Segmentation
 - Anomaly Detection
 - Dimensionality Reduction
 - Reinforced
- ✓ Best Practices of Machine Learning

Wikipedia:

“Machine learning is a scientific discipline that deals with the construction and study of algorithms that can learn from data. Such algorithms operate by building a model based on inputs and using that to make predictions or decisions, rather than following only explicitly programmed instructions.”

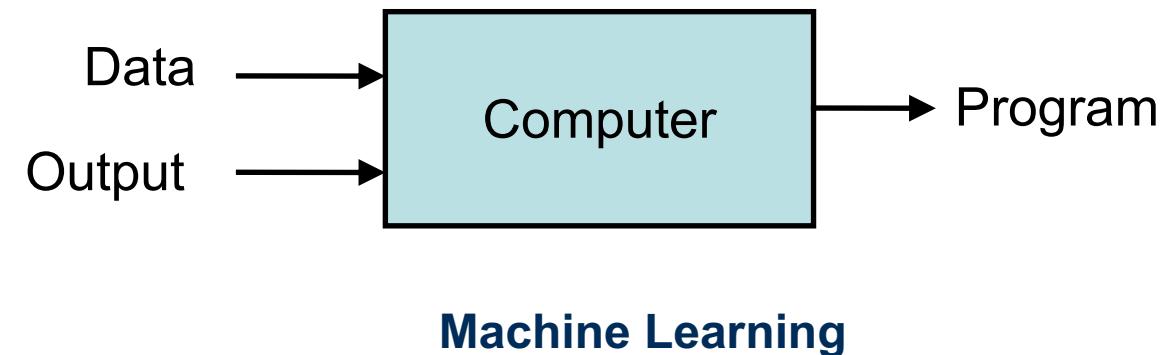
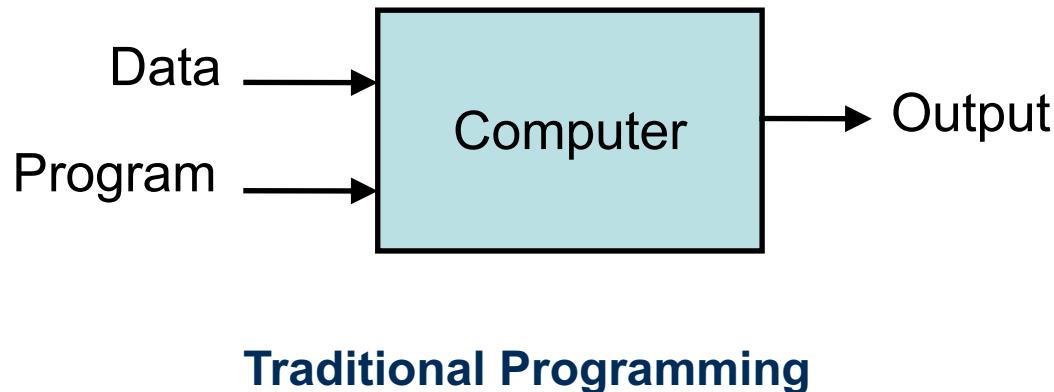
MACHINE LEARNING

WHAT IS IT?

"MACHINE LEARNING REFERS TO A SYSTEM CAPABLE OF THE AUTONOMOUS ACQUISITION AND INTEGRATION OF KNOWLEDGE."

"FIELD OF STUDY THAT GIVES COMPUTERS THE ABILITY TO LEARN WITHOUT BEING EXPLICITLY PROGRAMMED."

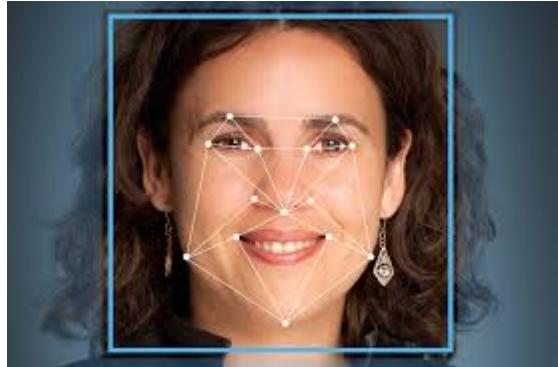
"GETTING COMPUTERS TO PROGRAM THEMSELVES."



MACHINE LEARNING WHEN WOULD WE USE MACHINE LEARNING?

- ✓ When patterns exists in our data
 - Even if we don't know what they are
- ✓ We can not pin down the functional relationships mathematically
 - Else we would just code up the algorithm
- ✓ When we have lots of (unlabeled) data
 - Labeled training sets harder to come by
 - Data is of high-dimension - for example, sensor data
 - Want to “discover” lower-dimension representations -Dimension reduction
- ✓ Rapidly changing phenomena
 - credit scoring, financial modeling
 - diagnosis, fraud detection

MACHINE LEARNING FEW EXAMPLES



Text / Image / Voice Recognition



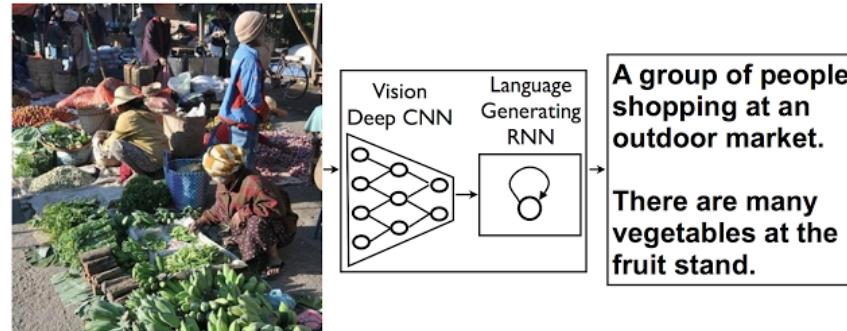
Spam Mail



Self Driving Cars



Recommendations



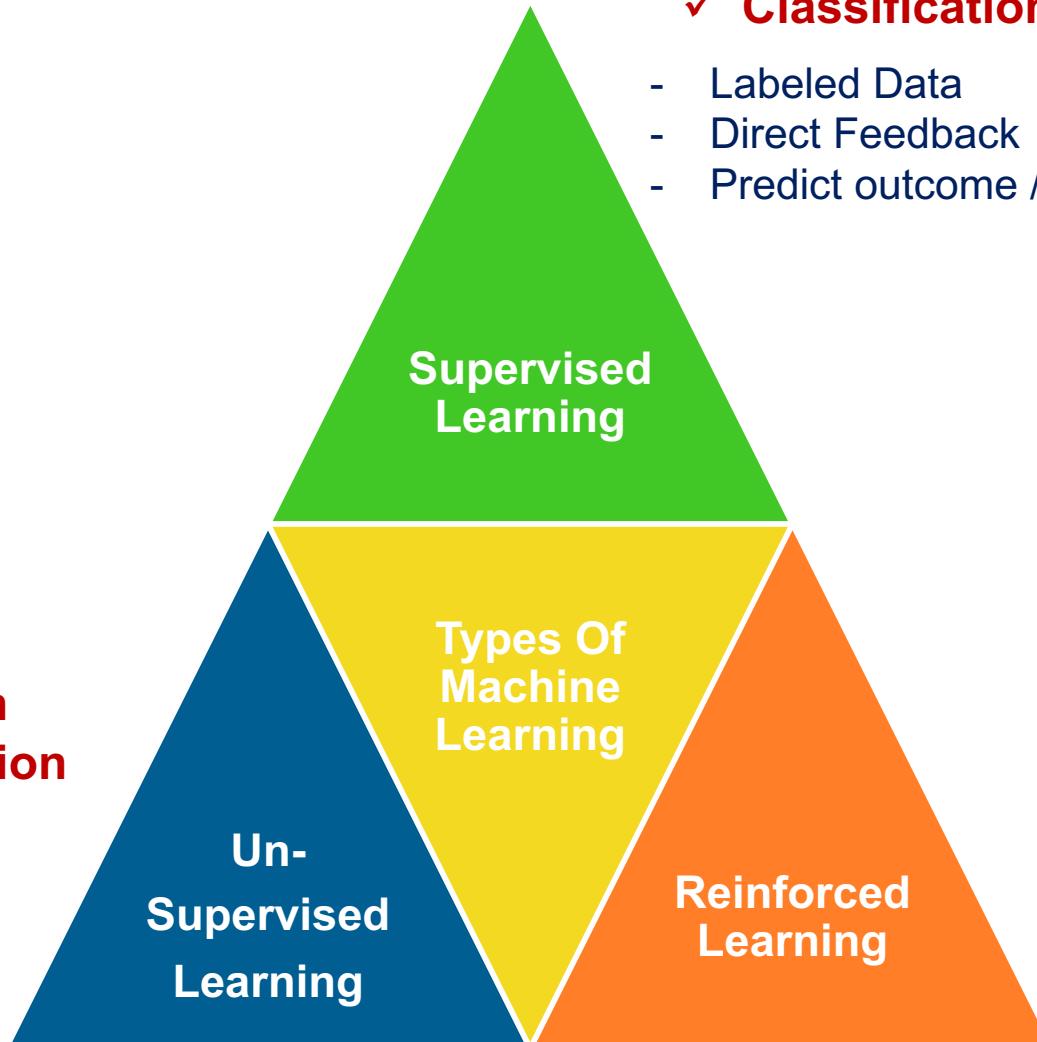
Picture Annotations

... and many more

MACHINE LEARNING

TYPES OF LEARNING

- ✓ **Clustering**
- ✓ **Anomaly Detection**
- ✓ **Dimension Reduction**
 - No Labels
 - No Feedback
 - “Find Hidden Structure”



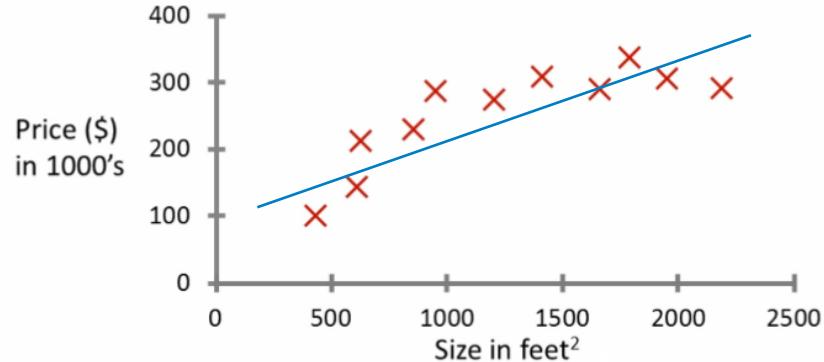
- ✓ **Regression**
- ✓ **Classification**
 - Labeled Data
 - Direct Feedback
 - Predict outcome / Future

- Decision Process
- Reward System
- Learn Series of actions

MACHINE LEARNING

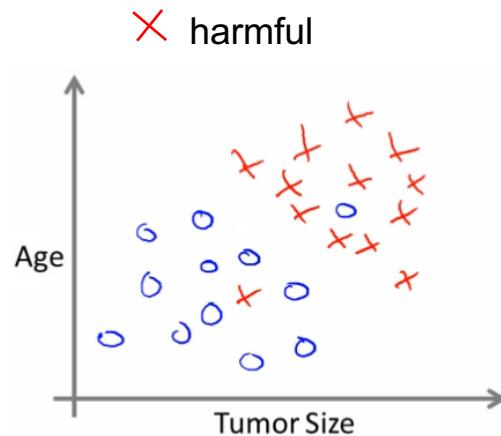
TYPES OF LEARNING

Housing price prediction.



Regression

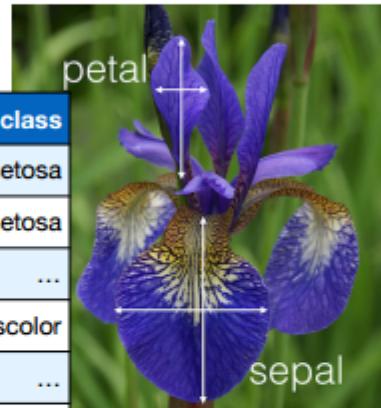
Linear & Polynomial Regression,
Decision Tree, Random Forest



MACHINE LEARNING NOMENCLATURES

Instances (samples, observations)

	sepal_length	sepal_width	petal_length	petal_width	class
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
...
50	6.4	3.2	4.5	1.5	veriscolor
...
150	5.9	3.0	5.1	1.8	virginica



Features (attributes, dimensions)

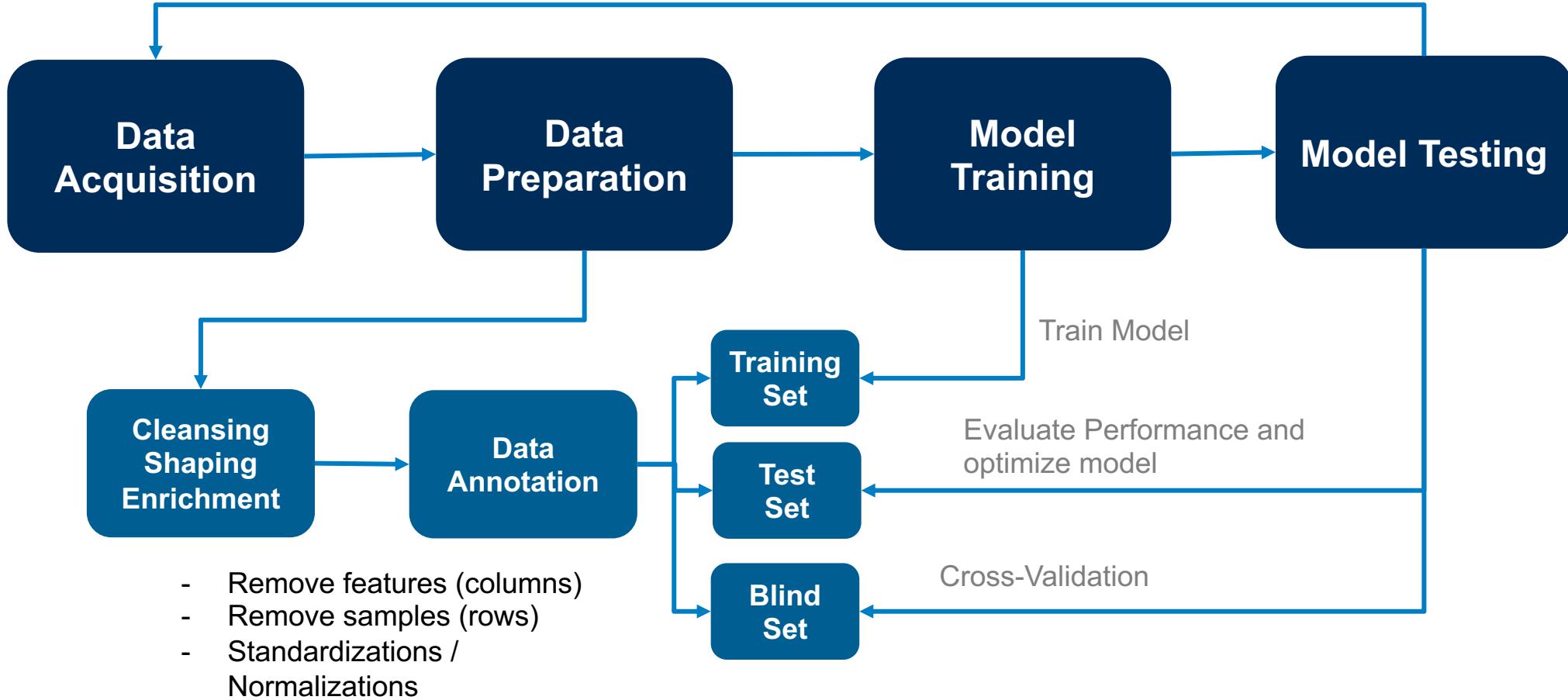
Classes (targets)

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1D} \\ x_{21} & x_{22} & \cdots & x_{2D} \\ x_{31} & x_{32} & \cdots & x_{3D} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{ND} \end{bmatrix}$$

$$\mathbf{y} = [y_1, y_2, y_3, \dots, y_N]$$

MACHINE LEARNING TYPICAL FLOW DIAGRAM

Iterative





OVERVIEW OF SPECIFIC MACHINE LEARNING METHODS

- ✓ Linear Regression
- ✓ Polynomial Regression
- ✓ Logistic Regression
- ✓ Decision Tree
- ✓ Naive Bayes
- ✓ Support Vector Machine
- ✓ K-Nearest Neighbour
- ✓ SVD / PCA

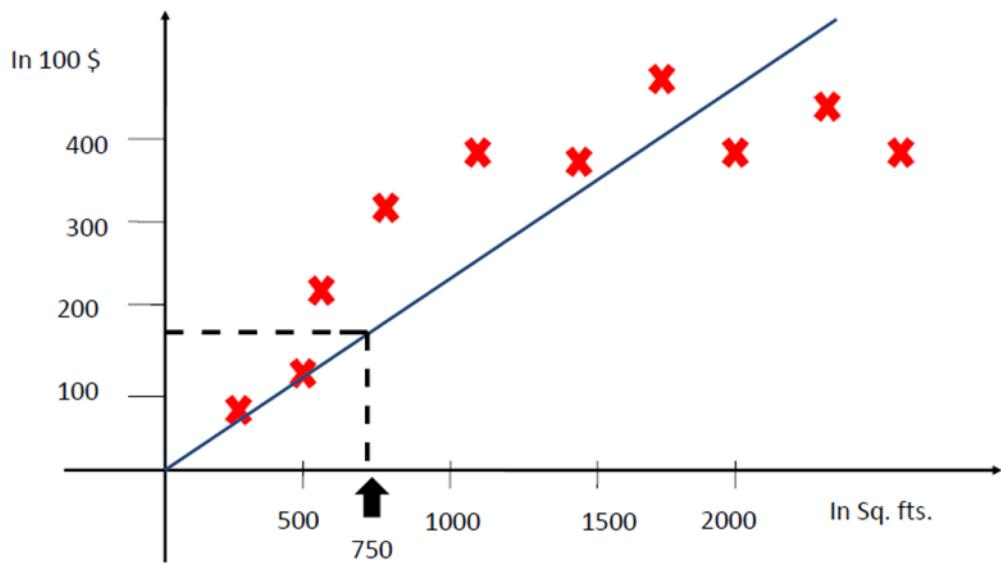
“CLASSICAL” REGRESSION

$$y_{\text{pred}} = \Theta_0 + \Theta_1 X$$

```
model= LinearRegression()  
model.fit(X, y)  
  
y_predicted = model.predict(X)
```

PREDICT LINEAR REGRESSION

Predict housing price given the sq. ft.



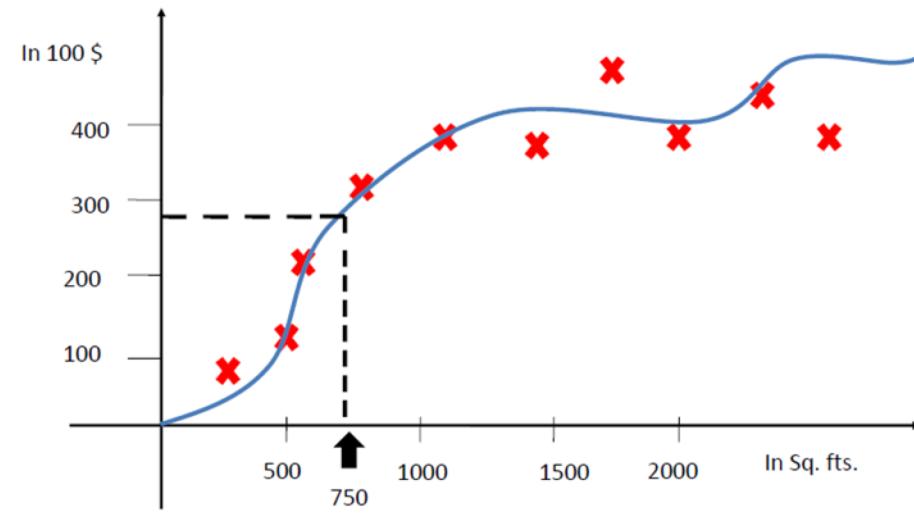
n – number of training sample

p – number of features/attributes

X – input variable

Y – output variable

(X(i), Y(i)) – ith training sample



In general –

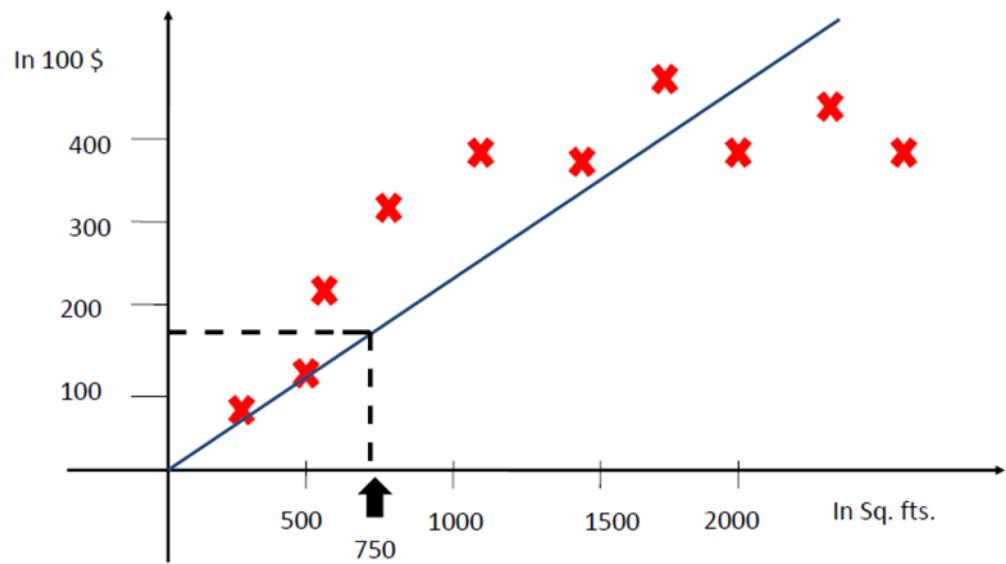
$$Y = (y_1, \dots, y_n)^T, \quad n \times 1$$

$$X = (X_1, \dots, X_p), \quad n \times p$$

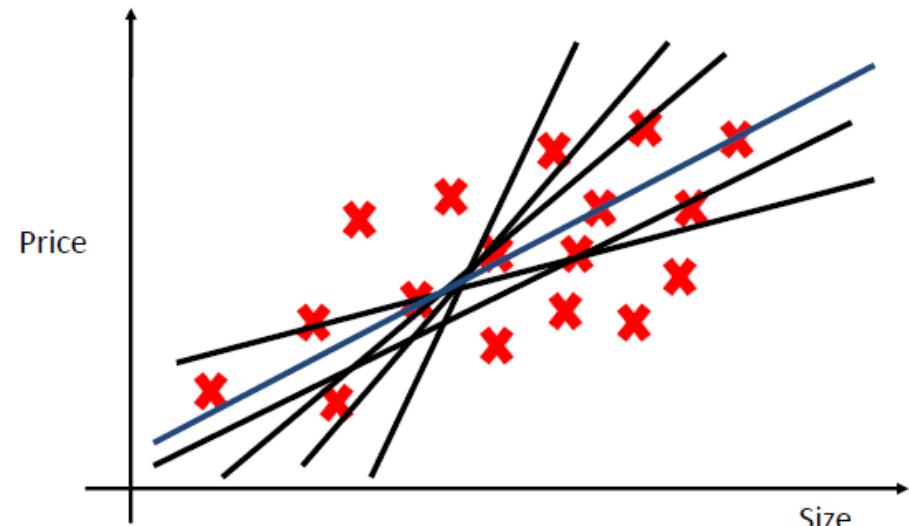
find hypothesis h that is linear combination and
use $\hat{y} = x^T \cdot h$ to predict y

PREDICT LINEAR REGRESSION

Predict housing price given the sq. ft.



$$h(X) = \Theta_0 + \Theta_1X$$



What is the best value of Θ_0 and Θ_1

n – number of training sample

p – number of features/attributes

X – input variable

Y – output variable

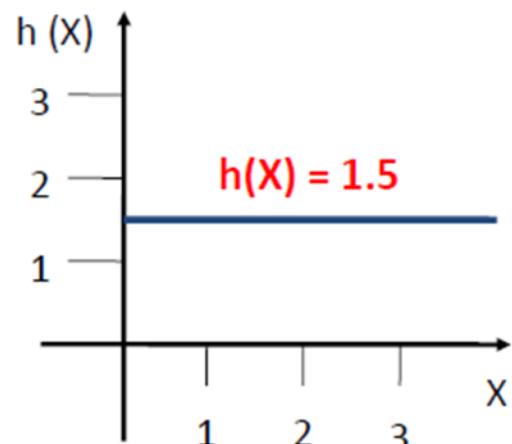
($X(i)$, $Y(i)$) – ith training sample



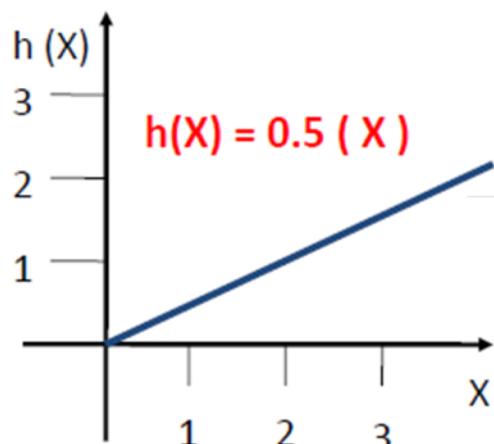
PREDICT LINEAR REGRESSION

What is the best value of Θ_0 and Θ_1

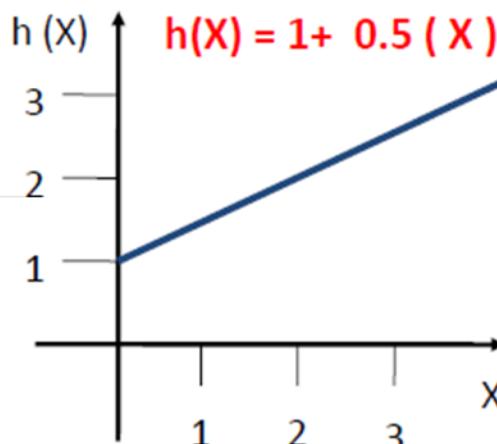
$$h(X) = \Theta_0 + \Theta_1 X$$



$$\begin{aligned}\Theta_0 &= 1.5 \\ \Theta_1 &= 0\end{aligned}$$



$$\begin{aligned}\Theta_0 &= 0 \\ \Theta_1 &= 0.5\end{aligned}$$



$$\begin{aligned}\Theta_0 &= 1 \\ \Theta_1 &= 0.5\end{aligned}$$

Choose the value of Θ_0 and Θ_1 such that $h_\Theta(X)$ is close to Y , for all training set

$$\text{minimize } (h_\Theta(X) - Y)^2$$

Predicted Given

Root Mean Error Square

Cost Function

$$\text{minimize } J(\Theta_0, \Theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\Theta(X^i) - Y^i)^2$$

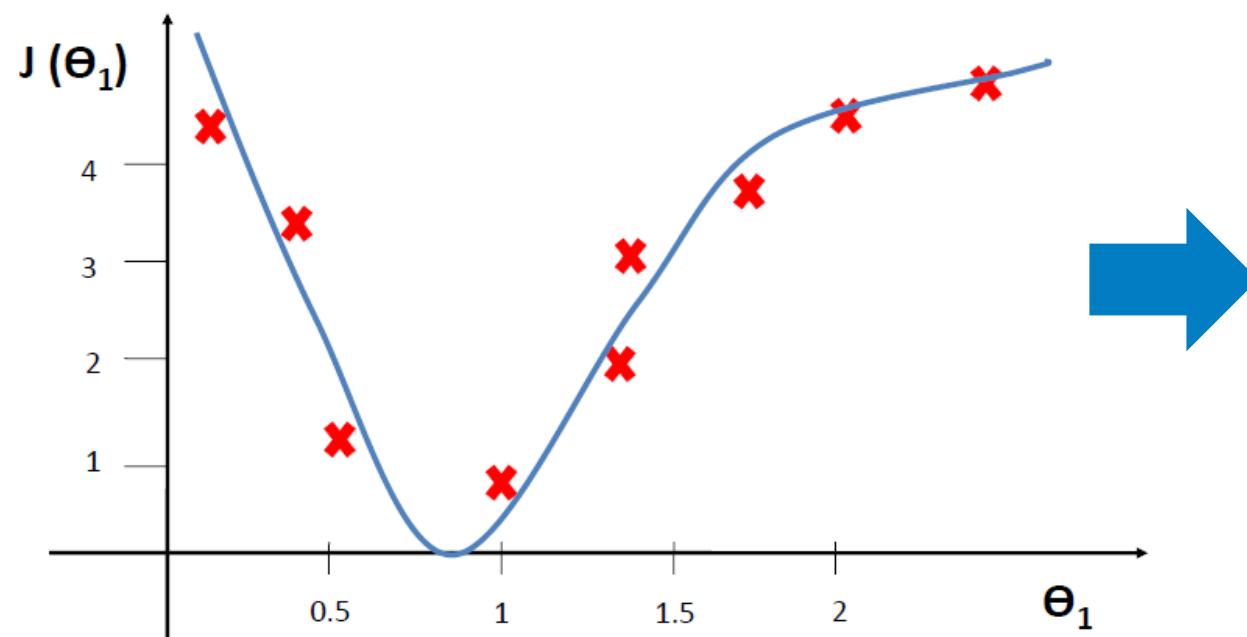
PREDICT LINEAR REGRESSION

What is the best value of Θ_0 and Θ_1

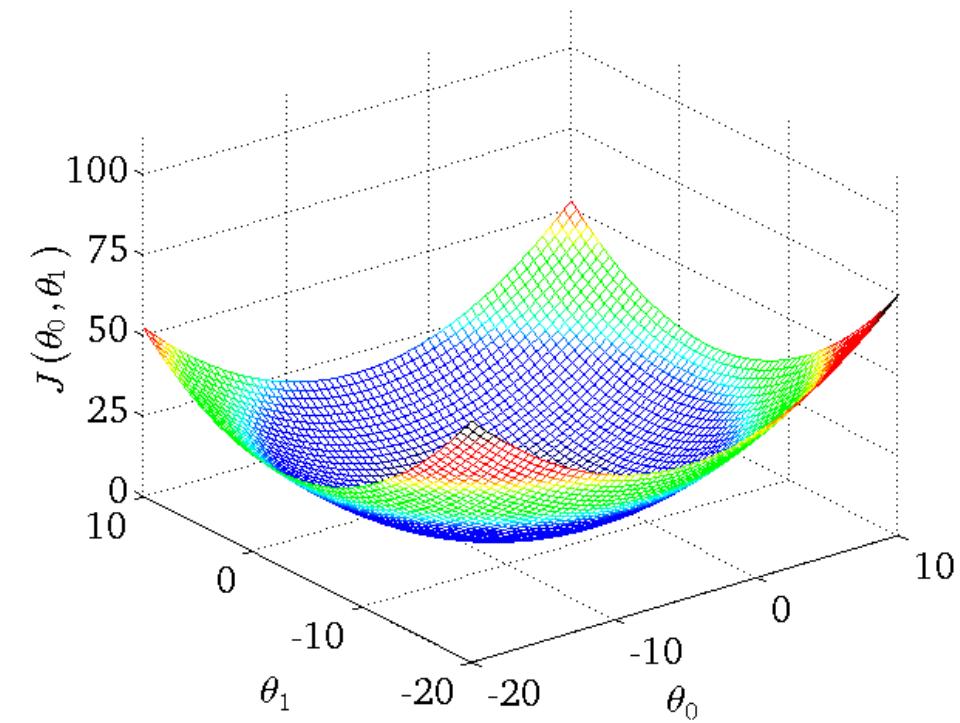
For $\Theta_0 = 0$

$$h(X) = \Theta_0 + \Theta_1 X$$

Cost function Vs Θ_1



Cost function Vs Θ_0 & Θ_1



Large number of iterations to find local minimum

Cost Function

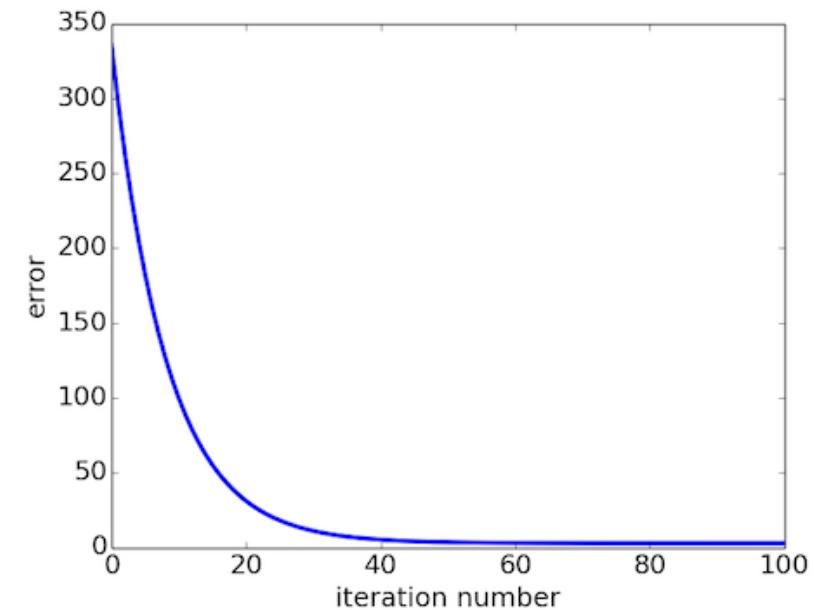
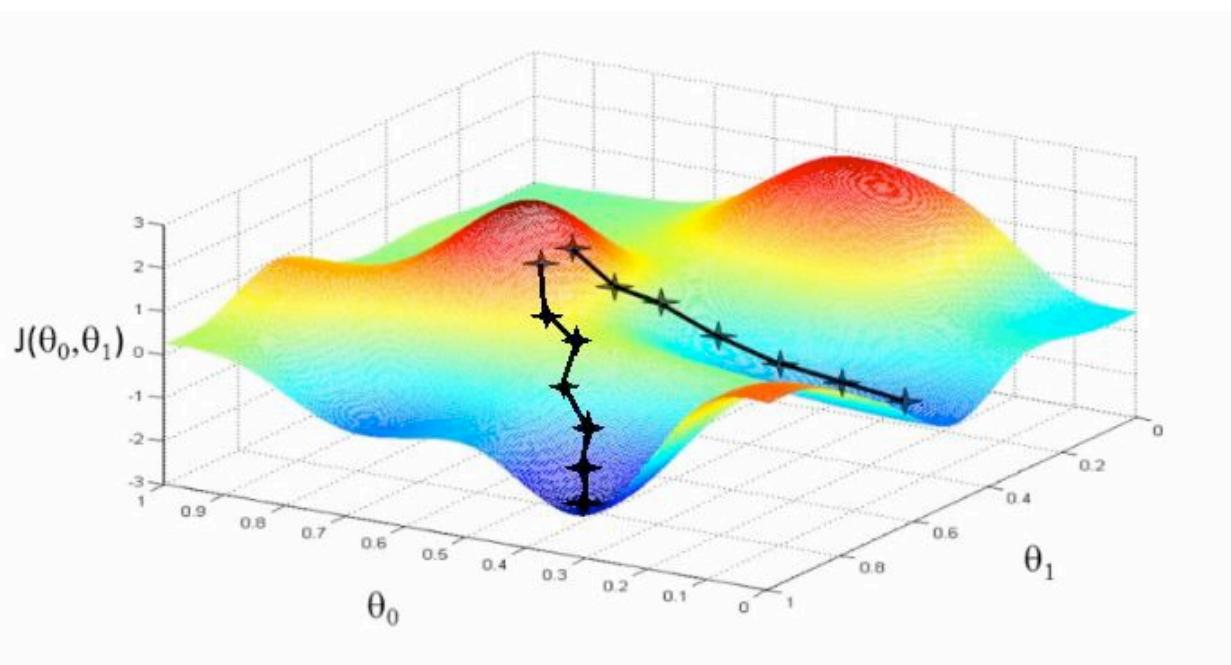
$$\text{minimize } J(\Theta_0, \Theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\Theta(X^i) - Y^i)^2$$

PREDICT LINEAR REGRESSION

What is the best value of Θ_0 and Θ_1

$$h(X) = \Theta_0 + \Theta_1 X$$

Gradient Descent Algorithm

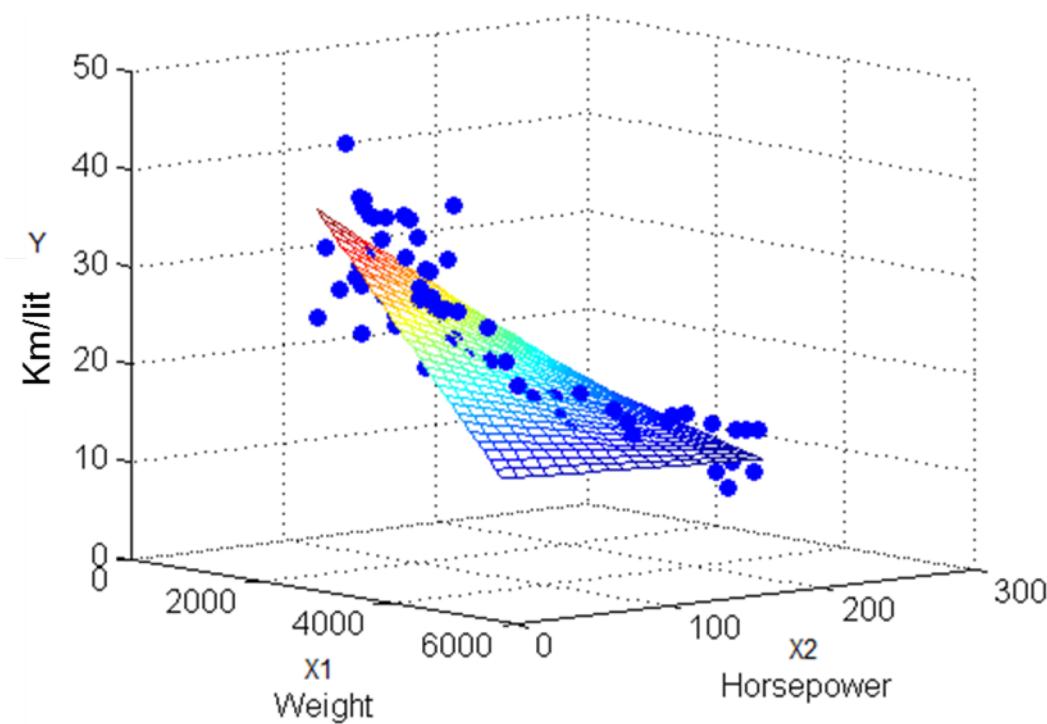


Error decreases with every iterations as it takes steps to negative gradient and hence much faster

PREDICT MULTIPLE LINEAR REGRESSION

What is the best value of $\Theta_0, \Theta_1, \Theta_2$

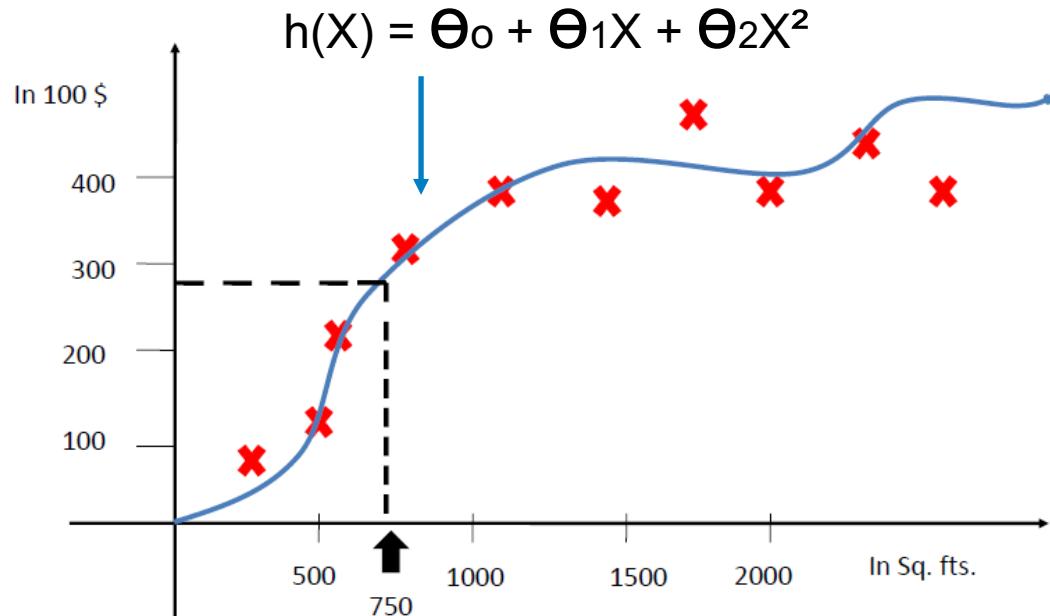
$$h(X) = \Theta_0 + \Theta_1 X_1 + \Theta_2 X_2$$



PREDICT POLYNOMIAL REGRESSION

[Demo](#)

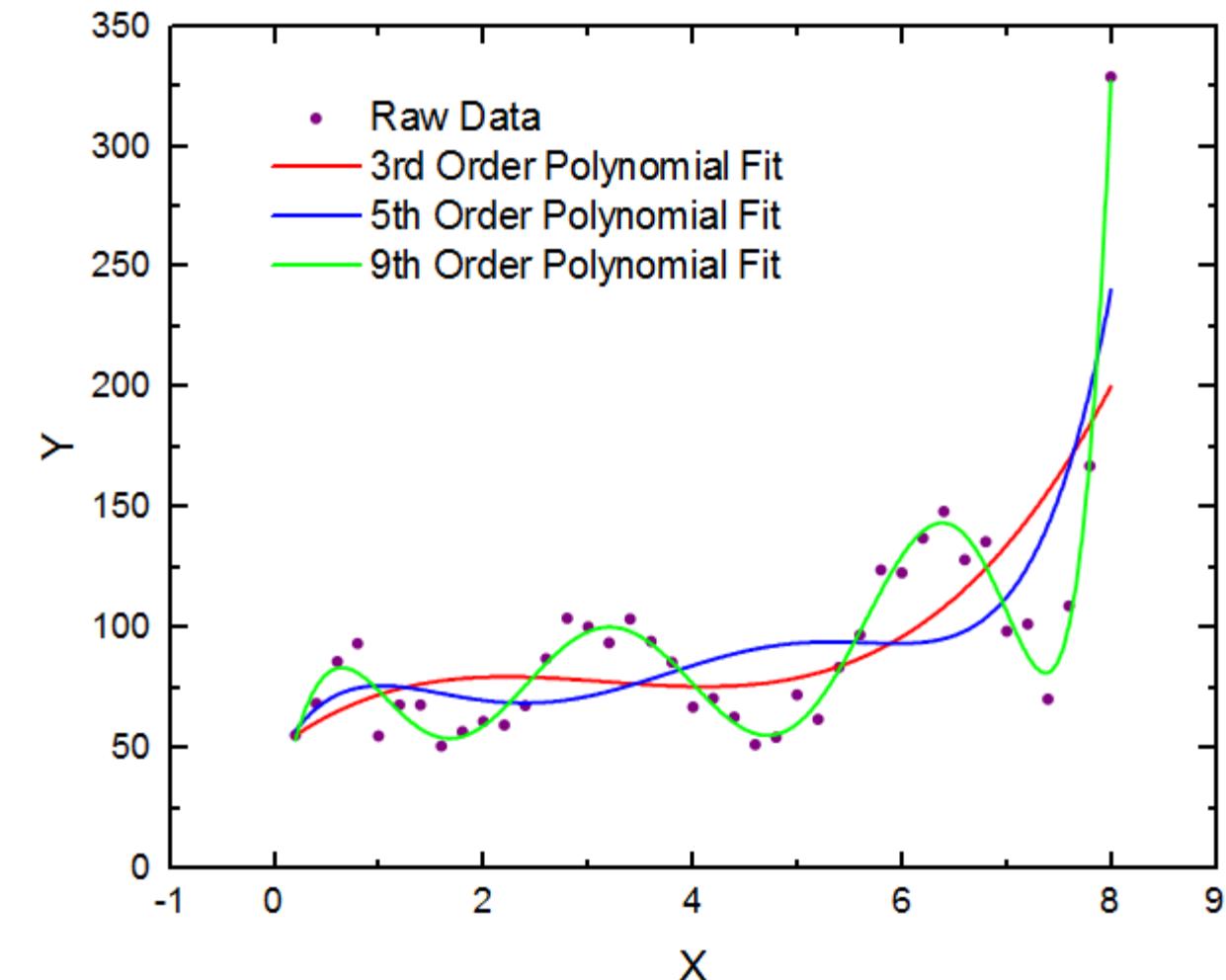
What is the best value of $\Theta_0, \Theta_1, \Theta_2$



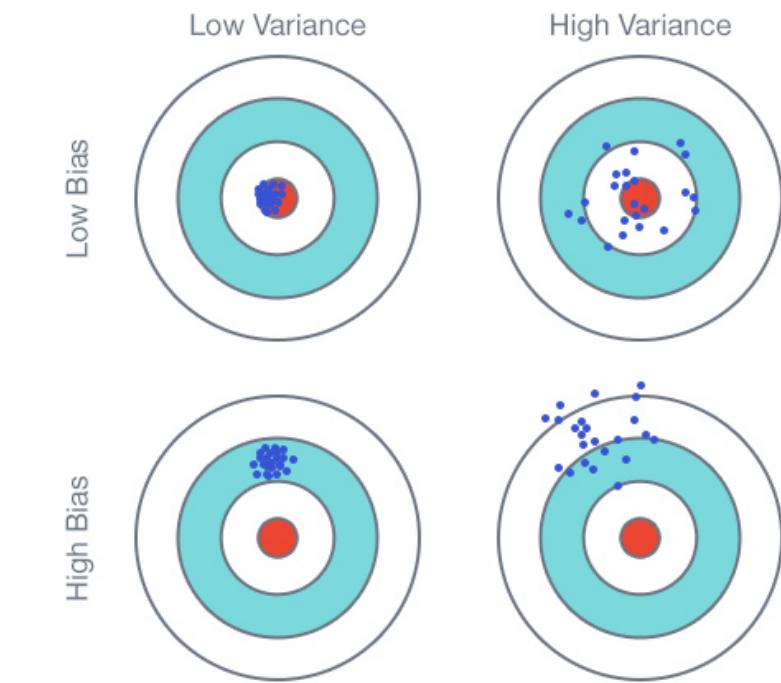
```
quadratic_featurizer = PolynomialFeatures(degree)
X_quadratic = quadratic_featurizer.fit_transform(X)

regressor_quadratic = LinearRegression()
regressor_quadratic.fit(X_quadratic, y)
```

PREDICT BIAS , VARIANCE AND OVERFITTING



With high variance the data is more spread out



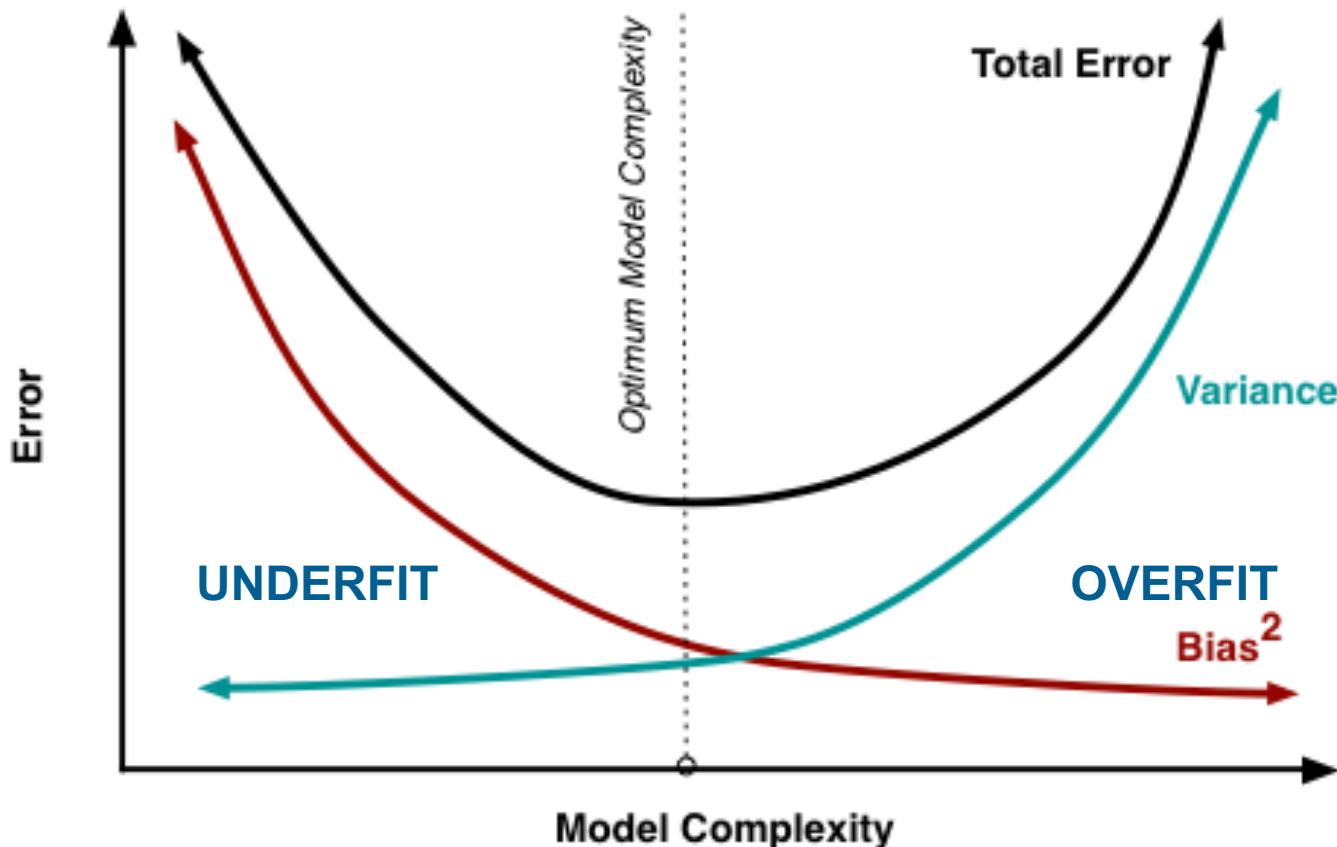
With high Bias the data is “biased” away from center

Bulls eye reflect the perfect Model

$$f(x) = a_0 + a_1x + a_2x^2 + \dots + a_{d-1}x^{d-1} + a_dx^d$$

PREDICT BIAS , VARIANCE AND OVERFITTING

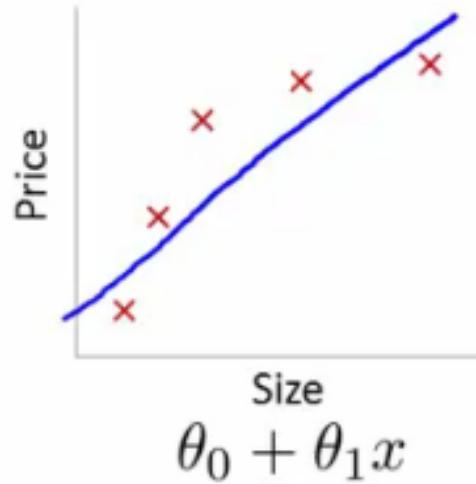
Bias and Variance Trade OFF



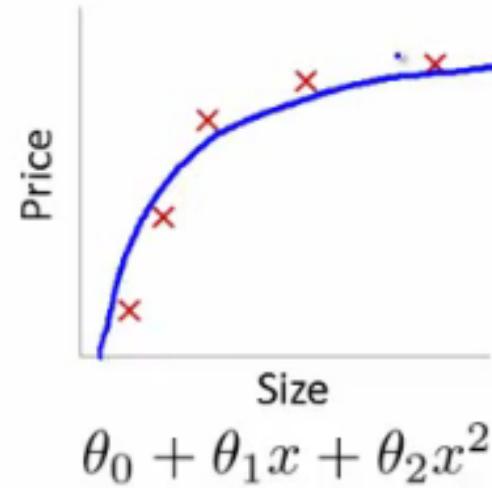
- As the degree gets higher the function becomes more complex
- Low dimensional polynomial has high bias – it under fits the data
- High dimensional data may fit for train data but gives more variance for future prediction – it over fits the data

Mean square error:

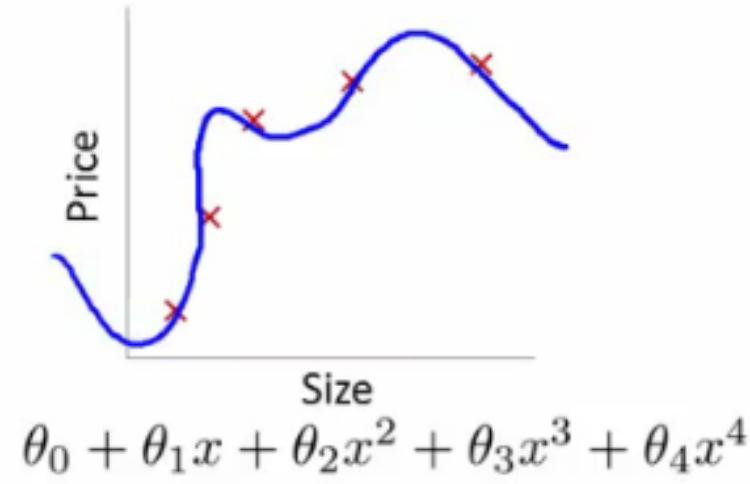
$$\begin{aligned} r(d, \theta) &= E[(d - \theta)^2] = (E[d] - \theta)^2 + E[(d - E[d])^2] \\ &= \text{Bias}^2 + \text{Variance} \end{aligned}$$

PREDICT RIGHT FIT SOLUTION

High bias
(underfit)



"Just right"



High variance
(overfit)

REGULARISATION RIDGE AND LASSO

- When large number of features are present with many of them correlated , we use regularization techniques to solve over-fitting problem.
- It does so by penalizing the bent of the regression line that tries to closely match the noisy data points. (penalizing the magnitude of coefficients of polynomials)
- It leads to smoothening of the regression line and thus prevents over-fitting
- You penalize your loss function by adding a multiple of an L1 (LASSO) or an L2 (Ridge) norm of your weights vector w

RIDGE Regression (L2)

Performs L2 regularization, i.e. adds penalty equivalent to **square of the magnitude** of coefficients

Ridge Regression

$$\min_w \||Xw - y\|_2^2 + \alpha \||w\||_2^2$$

```
from sklearn.linear_model import Ridge
ridgereg = Ridge(alpha=0, normalize=True)
ridgereg.fit(X_train, y_train)
y_pred = ridgereg.predict(X_test)
```

LASSO Regression (L1)

Performs L1 regularization, i.e. adds penalty equivalent to **absolute value of the magnitude** of coefficients

Lasso Regression

$$\min_w \frac{1}{2n_{samples}} \|Xw - y\|_2^2 + \alpha \rho \||w\||_1 + \frac{\alpha(1-\rho)}{2} \||w\||_2^2$$

```
from sklearn.linear_model import Lasso
model = Lasso(alpha=0.001, normalize=True)
model.fit(X_train, y_train)
```

PREDICT LOGISTIC REGRESSION

Often we have to resolve questions with binary or yes/no outcomes. For example:

Does a patient have cancer?

Will the customer buy my product?

Will a team win the next game?

Will I get the loan?

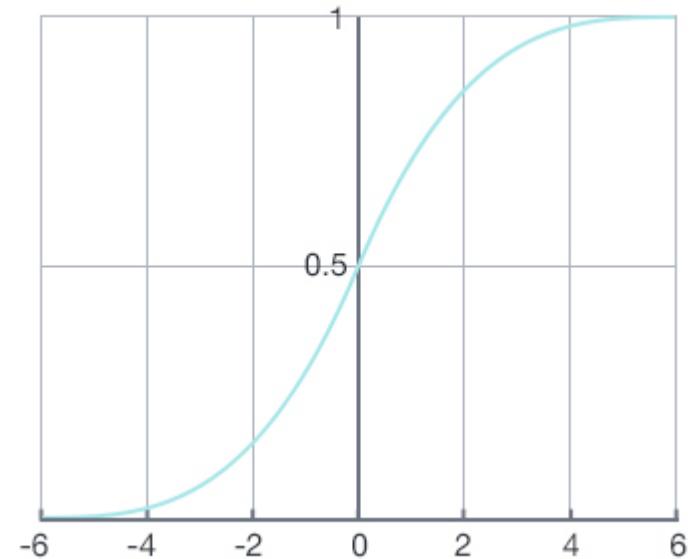
It is used to predict value that leads to classify it into two classes. This is classification but since it makes use of regression with logistic function it is called logistic regression.

Linear regression : $h(X) = \Theta_0 + \Theta_1 X$

Logistic regression : $h(X) = g(\Theta_0 + \Theta_1 X)$ so that $0 < h(X) < 1$

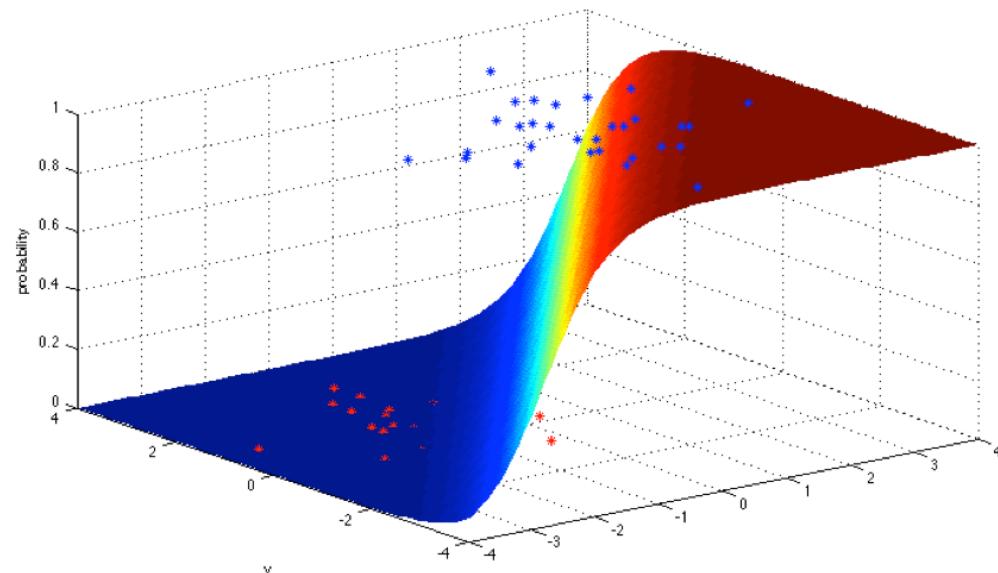
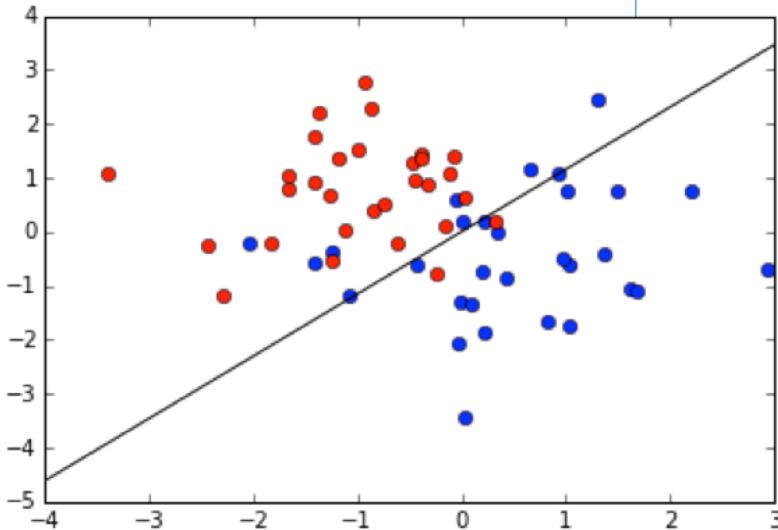
Where g is **sigmoid function** (logistic function)

$$g(z) = \frac{1}{1+e^{-z}}$$



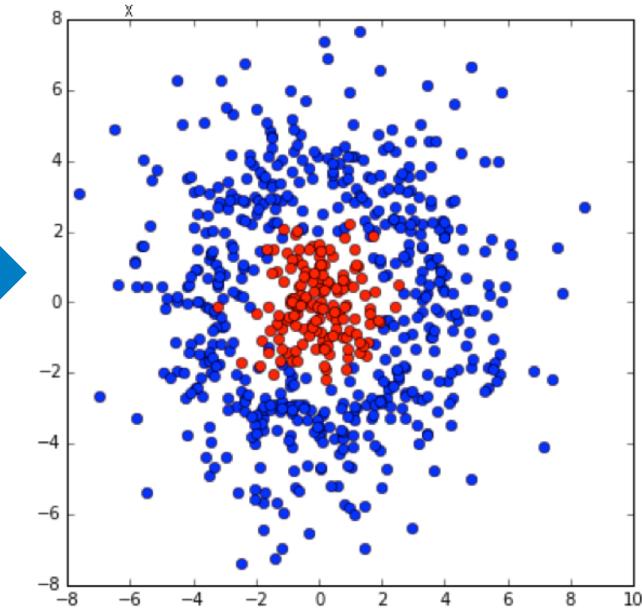
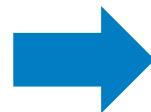
Applying linear regression model for classification is not great idea

PREDICT LOGISTIC REGRESSION



A logistic regression model makes a linear boundary between the two classes

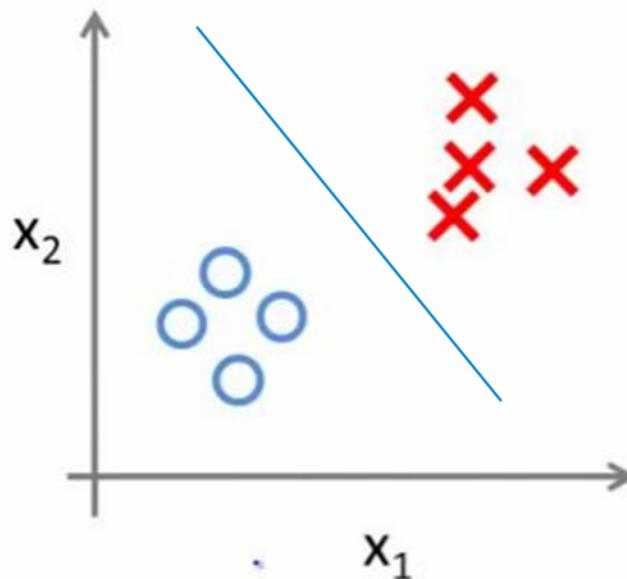
Logistic regression cannot draw a more complicated boundary than a line and hence not suitable for data set with many features



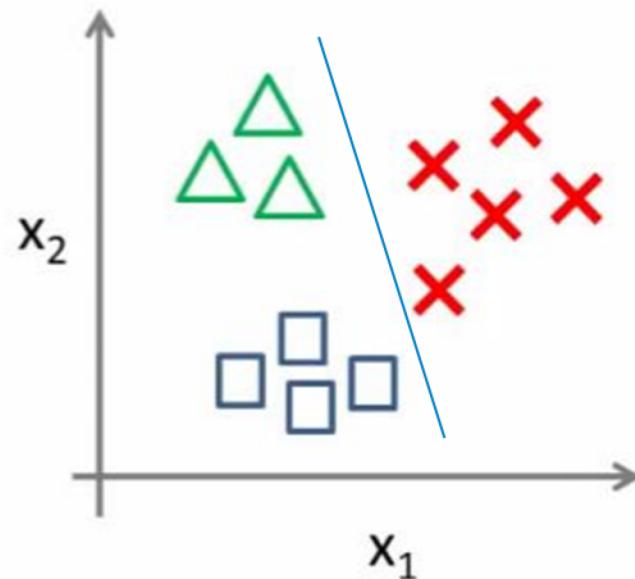
```
from sklearn.linear_model import LogisticRegression  
model = LogisticRegression()  
model.fit(X,y)
```

PREDICT MULTICLASS LOGISTIC REGRESSION

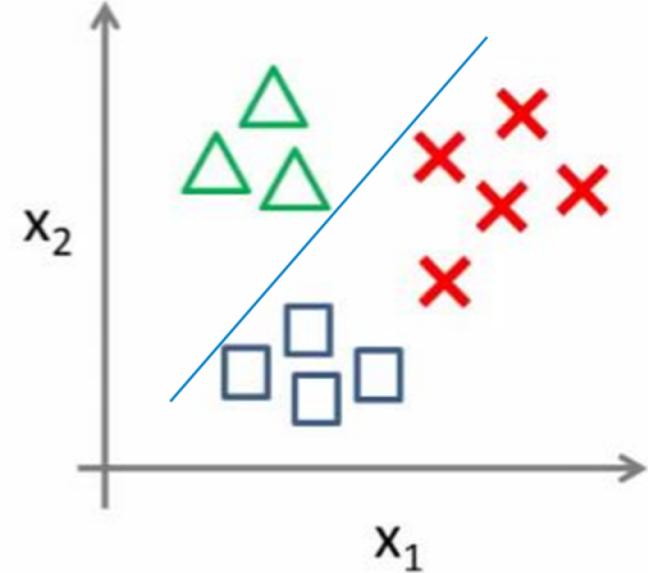
Binary classification:



Multi-class classification:



Multi-class classification:



Precision

- How often does our algorithm have false positives

= true positives / # predicted positive

= true positives / (true positive + false positive)

- High precision is good (i.e. closer to 1)

- You want a big number, because you want false positive to be as close to 0 as possible

$$P = \frac{TP}{TP + FP}$$

Accuracy:

$$acc = \frac{TP + TN}{TP + TN + FP + FN}$$

Recall

- How sensitive is our algorithm?

Of all patients in set that actually have cancer, what fraction did we correctly detect

= true positives / # actual positives

= true positive / (true positive + false negative)

- High recall is good (i.e. closer to 1)

- You want a big number, because you want false negative to be as close to 0 as possible

$$R = \frac{TP}{TP + FN}$$

F1Score (fscore)

Fscore is like taking the average of precision and recall giving a higher weight to the lower value

$$F1 = 2 \frac{PR}{P + R}$$

DECISION TREES

```
from sklearn.tree import DecisionTreeClassifier  
model = DecisionTreeClassifier()  
model.fit(X, y)
```

Tuning Parameters :

criterion = 'gini', 'entropy'
max_depth = 1, 2, 3..

DECISION TREES

How does it work? A simple example

Suppose we have the given data set

- ✓ 9 - Yes
- ✓ 5 - No

We want to know if Sachin will play golf or not on Day-15?



Training Samples
9 Yes / 5 No

Day	Outlook	Humidity	Windy	Play
D1	Sunny	High	Weak	No
D2	Sunny	High	Strong	No
D3	Overcast	High	Weak	Yes
D4	Rain	High	Weak	Yes
D5	Rain	Normal	Weak	Yes
D6	Rain	Normal	Strong	No
D7	Overcast	Normal	Strong	Yes
D8	Sunny	High	Weak	No
D9	Sunny	Normal	Weak	Yes
D10	Rain	Normal	Weak	Yes
D11	Sunny	Normal	Strong	Yes
D12	Overcast	High	Strong	Yes
D13	Overcast	Normal	Weak	Yes
D14	Rain	High	Strong	No
D15	Rain	High	Weak	?

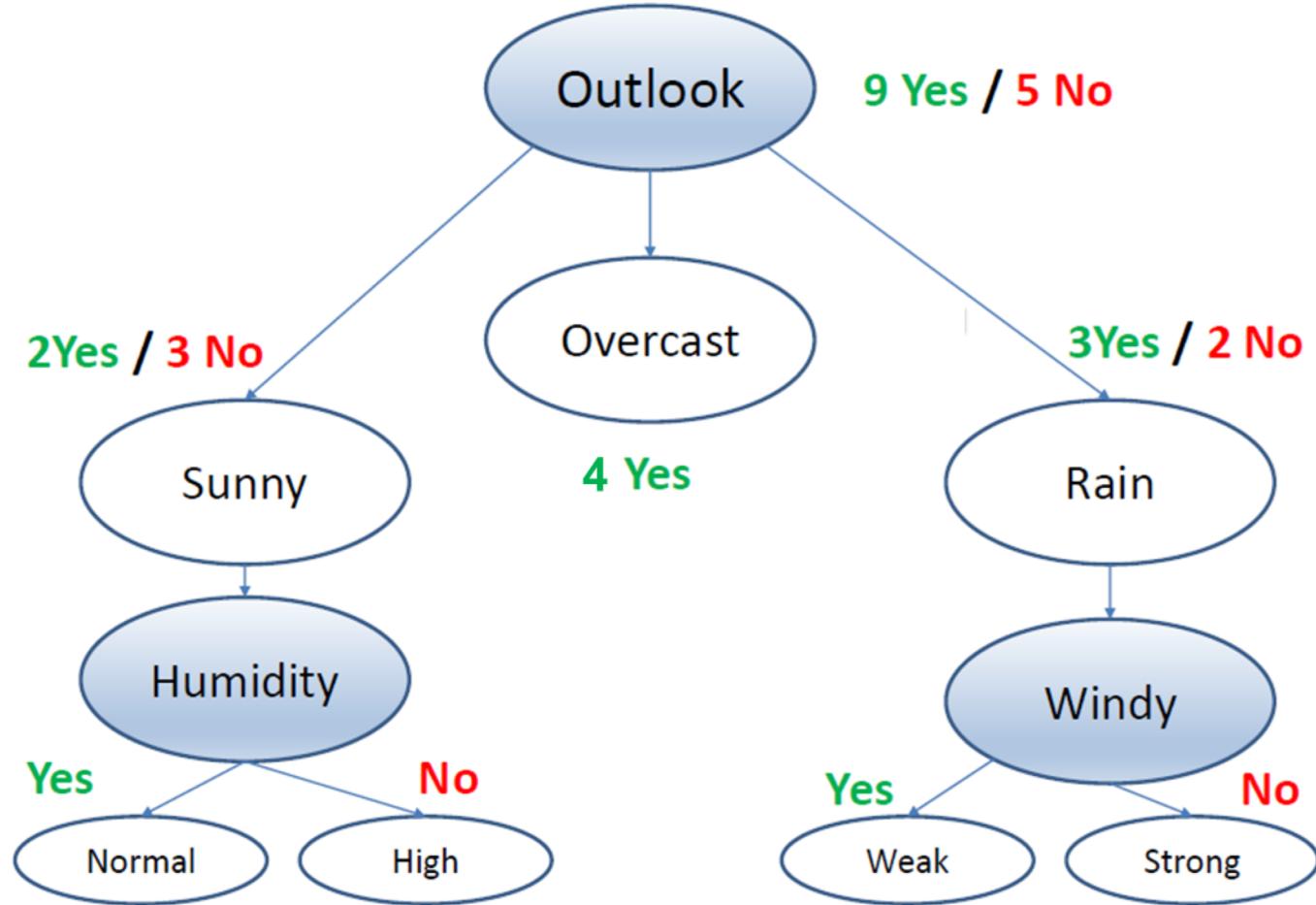
DECISION TREES

- If outlook is overcast , Sachin always plays a golf
- If outlook is sunny and humidity is high , he will not play
- If outlook is Rainy and there is a strong wind , he will not play



Logical Formulation:

(Outlook = Sunny \wedge Humidity=Normal)
 \vee (Outlook = Overcast)
 \vee (Outlook = Rain \wedge Wind = Weak)



How to split the attributes ???

ENTROPY CALCULATION

WHICH ATTRIBUTE IS BEST?

If,

- S is a sample of training examples
- p_{\oplus} is the proportion of positive examples in S
- p_{\ominus} is the proportion of negative examples in S

Entropy measures the impurity of S

$$\text{Entropy}(S) \equiv H(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

Pure Set (4 Yes / 0 No)

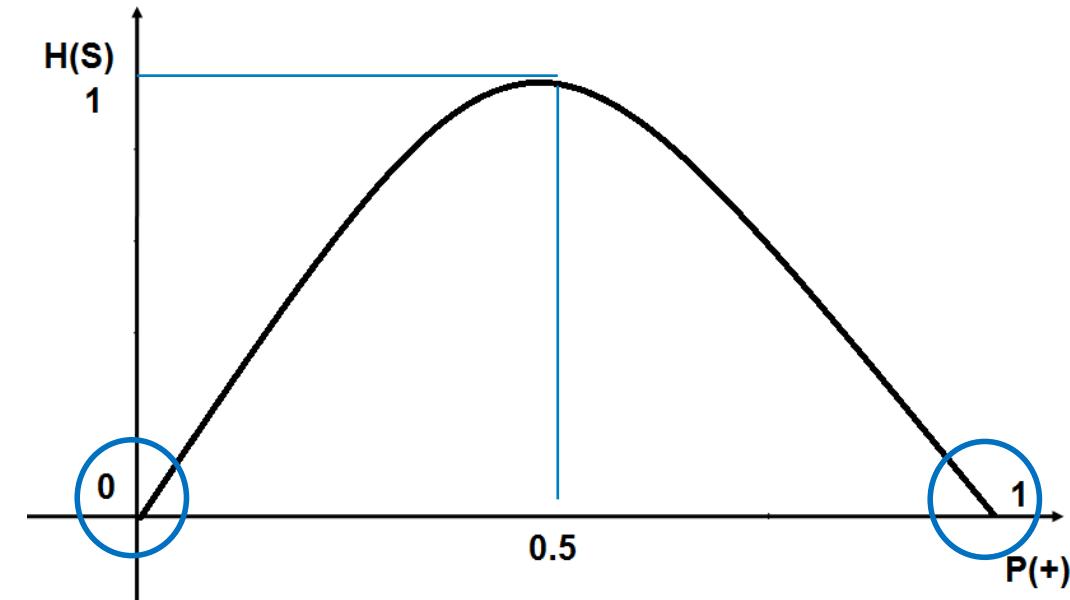
$H(S) = 0$ if sample is pure (all + or all -)

$$H(S) = -4/4\log_2(4/4) - 0/4\log_2(4/4) = 0$$

Impure Set (3 Yes / 3 No)

$$H(S) = -3/6\log_2(3/6) - 3/6\log_2(3/6) = 1$$

$$H(S) = 1 \text{ if } p_{\oplus} = p_{\ominus} = 0.5$$



INFORMATION GAIN

WHICH ATTRIBUTE IS BEST?

Step-1) Calculate entropy of the target ...required in step-2

$$E(S) = \sum_{i=1}^c - p_i \log_2 p_i$$

Play Golf	
Yes	No
9	5



$$\begin{aligned}\text{Entropy(PlayGolf)} &= \text{Entropy}(5,9) \\ &= \text{Entropy}(0.36, 0.64) \\ &= -(0.36 \log_2 0.36) - (0.64 \log_2 0.64) \\ &= 0.94\end{aligned}$$

INFORMATION GAIN

WHICH ATTRIBUTE IS BEST?

Step-2) The dataset is then split on the different attributes. The entropy for each branch is calculated.

		Play Golf	
		Yes	No
Outlook	Sunny	3	2
	Overcast	4	0
	Rainy	2	3
		Gain = 0.247	

		Play Golf	
		Yes	No
Windy	False	6	2
	True	3	3
		Gain = 0.048	

		Play Golf	
		Yes	No
Humidity	High	3	4
	Normal	6	1
		Gain = 0.152	

$$Gain(T, X) = Entropy(T) - Entropy(T, X)$$

$$E(T, X) = \sum_{c \in X} P(c)E(c)$$

		Play Golf		
		Yes	No	
Outlook	Sunny	3	2	5
	Overcast	4	0	4
	Rainy	2	3	5
				14



$$\begin{aligned}
 E(\text{PlayGolf}, \text{Outlook}) &= P(\text{Sunny}) * E(3,2) + P(\text{Overcast}) * E(4,0) + P(\text{Rainy}) * E(2,3) \\
 &= (5/14) * 0.971 + (4/14) * 0.0 + (5/14) * 0.971 \\
 &= 0.693
 \end{aligned}$$



$$\begin{aligned}
 G(\text{PlayGolf}, \text{Outlook}) &= E(\text{PlayGolf}) - E(\text{PlayGolf}, \text{Outlook}) \\
 &= 0.940 - 0.693 = 0.247
 \end{aligned}$$

INFORMATION GAIN

WHICH ATTRIBUTE IS BEST?

Step-3) Choose attribute with the largest information gain as the decision node.

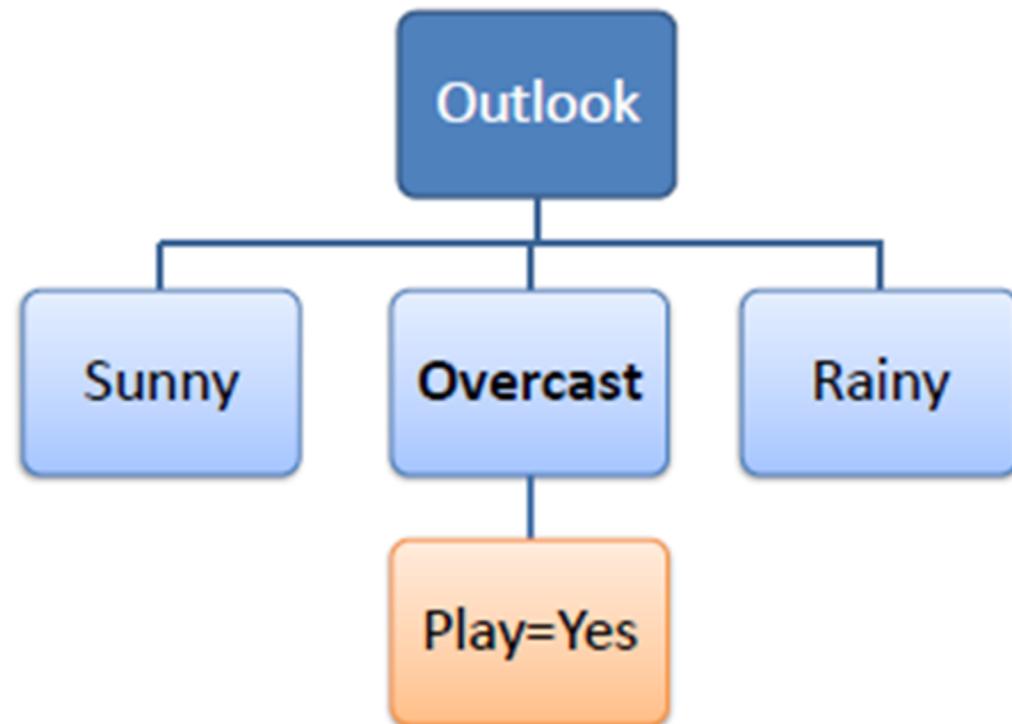
Outlook		Play Golf	
		Yes	No
Sunny	Yes	3	2
	Overcast	4	0
	Rainy	2	3
Gain = 0.247			

INFORMATION GAIN

WHICH ATTRIBUTE IS BEST?

Step-4a) A branch with entropy of 0 is a leaf node.

Humidity	Windy	Play Golf
High	FALSE	Yes
Normal	TRUE	Yes
High	TRUE	Yes
Normal	FALSE	Yes
High	FALSE	Yes

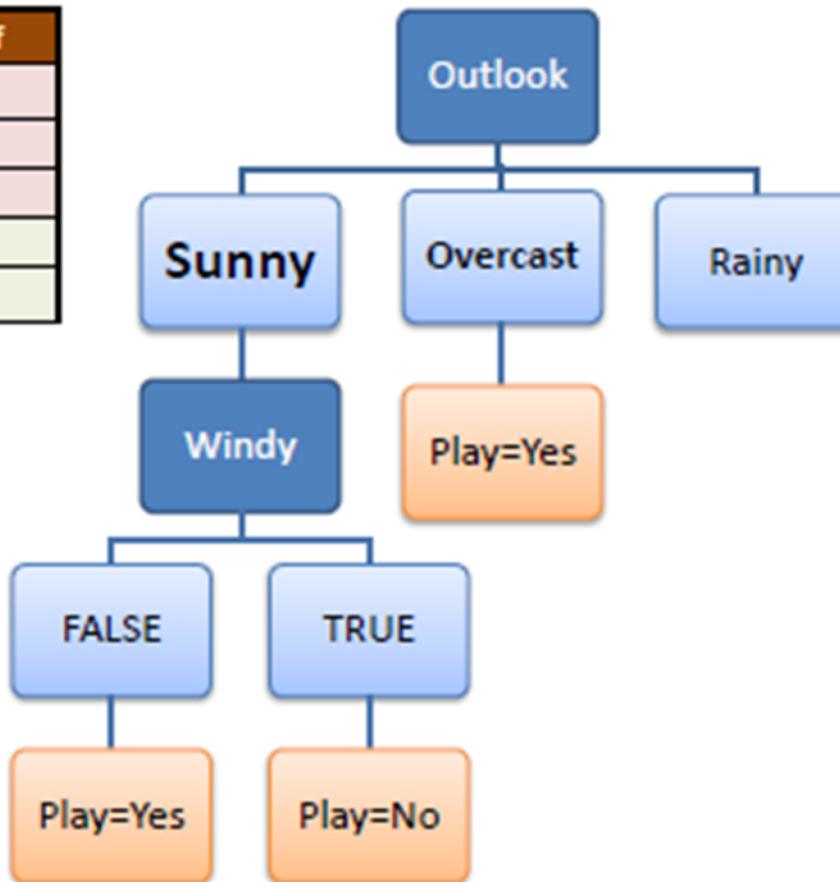


INFORMATION GAIN

WHICH ATTRIBUTE IS BEST?

Step-4a) A branch with entropy more than 0 needs further splitting.

Humidity	Windy	Play Golf
High	FALSE	Yes
Normal	FALSE	Yes
Normal	FALSE	Yes
Normal	TRUE	No
High	TRUE	No



INFORMATION GAIN

WHICH ATTRIBUTE IS BEST?

Step-5 – The algorithm is run recursively on the non-leaf branches, until all data is classified.

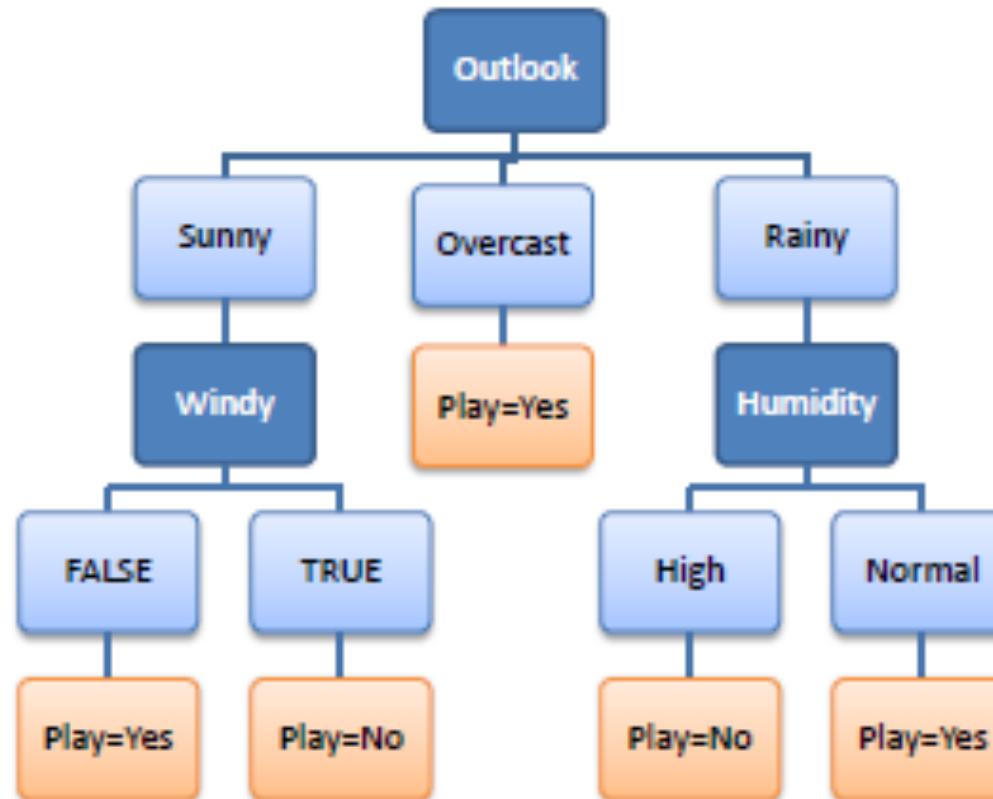
R₁: IF (Outlook=Sunny) AND
(Windy=FALSE) THEN Play=Yes

R₂: IF (Outlook=Sunny) AND
(Windy=TRUE) THEN Play=No

R₃: IF (Outlook=Overcast) THEN
Play=Yes

R₄: IF (Outlook=Rainy) AND
(Humidity=High) THEN Play=No

R₅: IF (Outlook=Rain) AND
(Humidity=Normal) THEN
Play=Yes

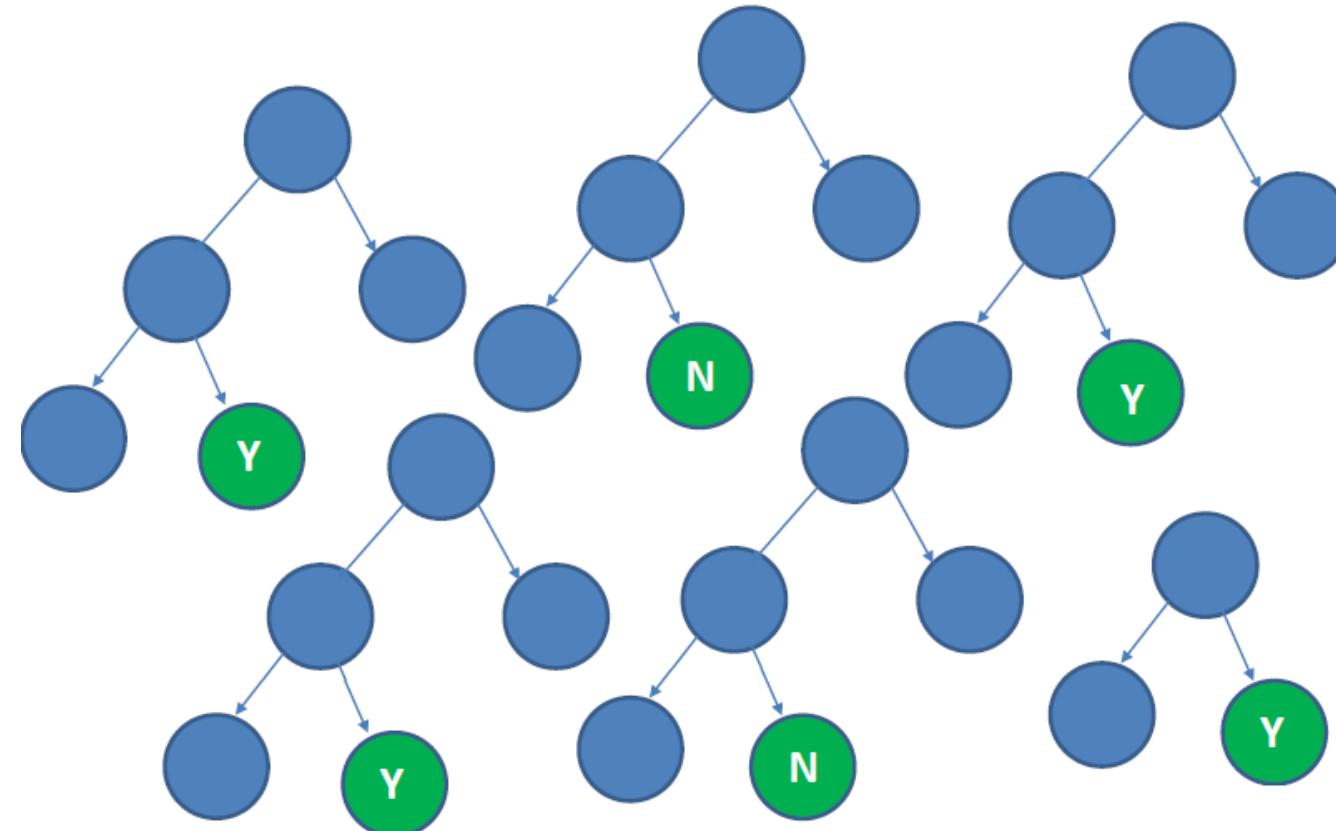


RANDOM FOREST

HOW IT WORKS?

- An ensemble classifier using many decision tree models.
 - In a dataset, where M is the total number of input attributes in the dataset, only m attributes are chosen at random for each tree where $m < M$

Not suitable
for prediction
of continuous
attribute.



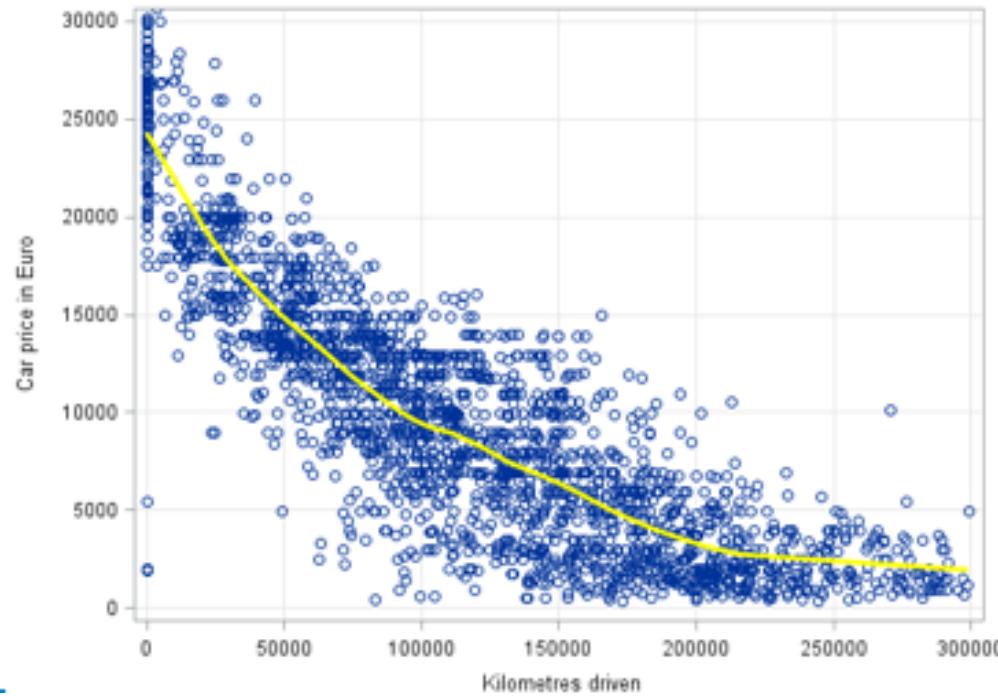
Decision trees pros and cons

pros

- Interaction between variables
- Interpretable rules
- Missing values easy to incorporate.

cons

- May become Unstable if data set is improper
- “Lack-of-Smoothnes”
- prone to errors in classification, owing to differences in perceptions



NAIVE BAYES CLASSIFIER

BAYES THEOREM

Likelihood of evidence 'X' if Hypothesis 'C' is true

$$P(c | X) = \frac{P(X | C) P(C)}{P(X)}$$

Posterior Probability of 'C' Given the evidence X

Prior Probability

Prior probability that the evidence itself is true

Frequency Table

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

		Play Golf	
		Yes	No
Humidity	High	3	4
	Normal	6	1

		Play Golf	
		Yes	No
Windy	False	6	2
	True	3	3

NAIVE BAYES WILL SACHIN PLAY ?

$$P(c|X) = P(X|C) * P(C) / P(X)$$

$$\begin{aligned}P(\text{Yes}|\text{Sunny}) &= P(\text{Sunny}|\text{Yes}) * P(\text{Yes}) / P(\text{Sunny}) \\&= (3/9) * (9/14) / (5/14)\end{aligned}$$

$$P(x | C) = P(\text{Sunny} | \text{yes}) = 3/9$$

		Playing 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	5/9	

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

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

NAIVE BAYES | WILL SACHIN PLAY ?

Outlook : Rainy

Humidity : Normal

Windy : Strong

Will Sachin play ?

Likelihood of Yes =

$$P(\text{Outlook: Rainy} \mid \text{Yes}) * P(\text{Humidity: Normal} \mid \text{Yes}) * P(\text{Windy: Strong} \mid \text{Yes}) * P(\text{Yes})$$

Likelihood of No =

$$P(\text{Outlook: Rainy} \mid \text{No}) * P(\text{Humidity: Normal} \mid \text{No}) * P(\text{Windy: Strong} \mid \text{No}) * P(\text{Yes})$$

$$P(\text{yes}) = \frac{\text{Likelihood of Yes}}{(\text{Likelihood of Yes} + \text{Likelihood of No})}$$

$$P(\text{No}) = \frac{\text{Likelihood of No}}{(\text{Likelihood of Yes} + \text{Likelihood of No})}$$

NAIVE BAYES BASED ON DISTRIBUTION OF FEATURES

- **GaussianNB**: The features follow a normal distribution. Instead of discrete counts, we have continuous features
- **MultinomialNB**: Text classification data are typically represented as word vector counts
- **BernoulliNB**: The binomial model is useful if your feature vectors are binary

SUPPORT VECTOR MACHINE (SVM)

```
from sklearn import svm  
model = svm.SVC()  
model.fit(X, y)
```

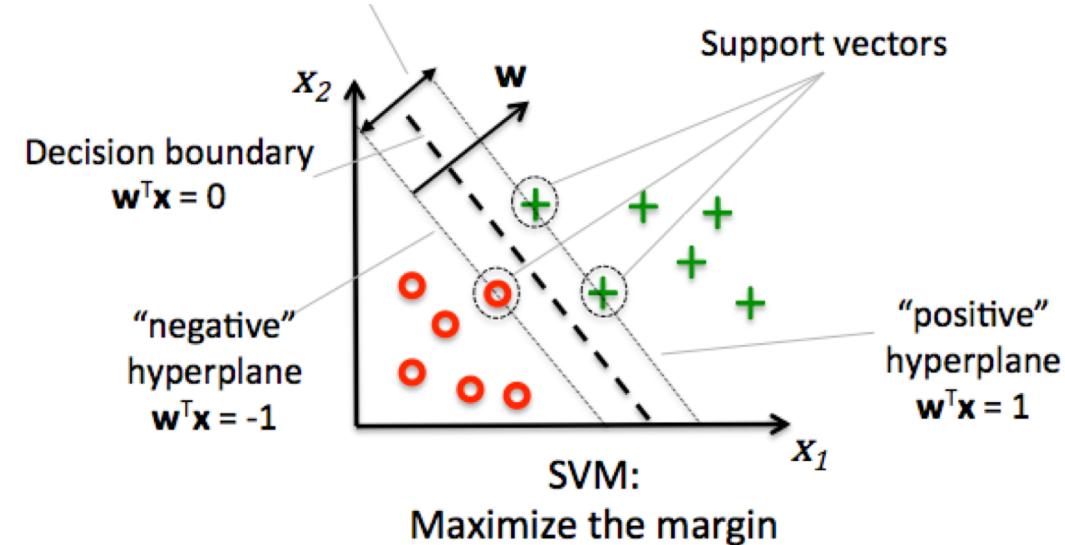
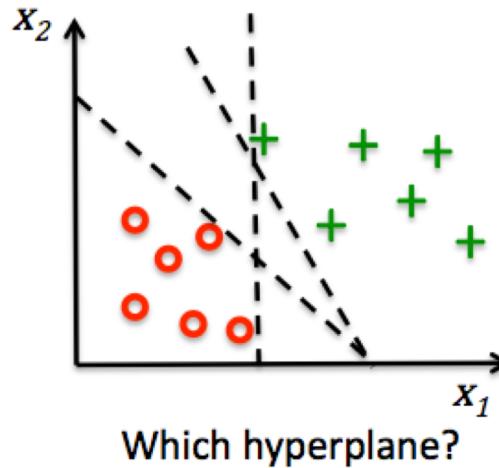
Tuning Parameters :

kernel = ' linear ', C = 1 or 0.1 etc

kernel = ' rbf ', C , gamma=0, 0.1 or 10

SVM HOW IT WORKS

A Support Vector Machine (SVM) performs classification by constructing an N -dimensional hyper plane that optimally separates the data into two categories.



Generally SVM deals with

- more than two attribute variables
- separating the points with non-linear curves
- handling the cases where clusters cannot be completely separated
- handling classifications with more than two categories.

Linear SVM

$$x_i \cdot x_j$$

Non-linear SVM

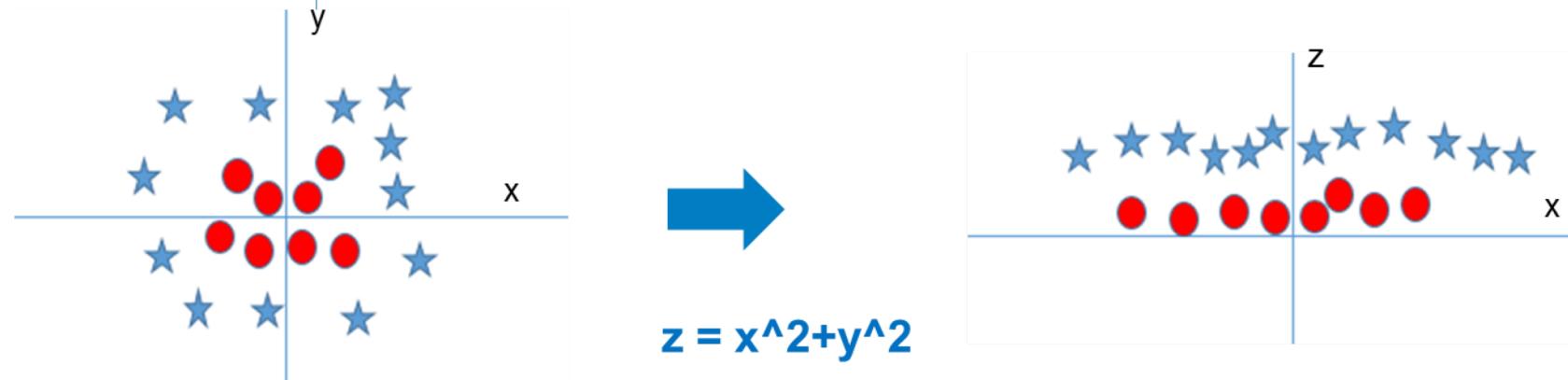
$$\phi(x_i) \cdot \phi(x_j)$$

Kernel function

$$k(x_i \cdot x_j)$$

SVM KERNEL TRICK

Demo

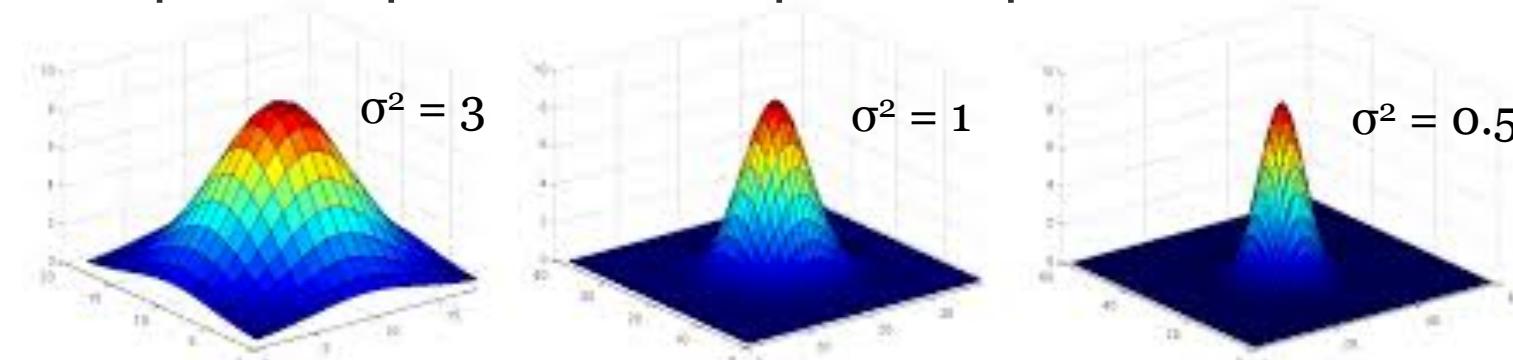


Kernels are functions which takes low dimensional input space and transform it to a higher dimensional space.

It is mostly useful in non-linear separation problem.

It converts not separable problem to separable problem.

σ^2 is a Gaussian parameter



* You can develop your own basis function

K-NEAREST NEIGHBORS (KNN) CLASSIFIER

```
from sklearn.neighbors import KNeighborsClassifier  
model = KNeighborsClassifier(n_neighbors=3)  
model.fit(X,y)
```

Tuning Parameters :

metric = 'minkowski', Euclidean, Manhattan
n_neighbors = 3, 5, 7, 9..

KNN | HOW IT WORKS ?

K nearest neighbors is a simple algorithm that stores all available instances and classifies new instance based on a similarity measure (e.g., distance functions).

- Decide K (generally odd number)
- Calculate the distance between any two points
- Find the nearest neighbors based on these pairwise distances
- Majority vote on a class labels based on the nearest neighbor list of K distances

Demo

KNN | HOW IT WORKS ?

Distance functions

Euclidean

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

Manhattan

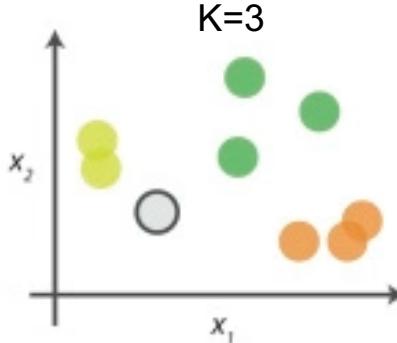
$$\sum_{i=1}^k |x_i - y_i|$$

Minkowski

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

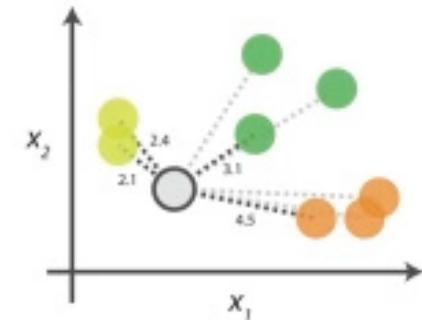
Also called as Lazy learning

0. Look at the data



Say you want to classify the grey point into a class. Here, there are three potential classes - lime green, green and orange.

1. Calculate distances



Start by calculating the distances between the grey point and all other points.

2. Find neighbours

Point	Distance	Nearest Neighbour
...	2.1	1st NN
...	2.4	2nd NN
...	3.1	3rd NN
...	4.5	4th NN

Next, find the nearest neighbours by ranking points by increasing distance. The nearest neighbours (NNs) of the grey point are the ones closest in dataspace.

3. Vote on labels

Class	# of votes
lime green	2
green	1
orange	1

Class lime green wins the vote!
Point grey is therefore predicted to be of class lime green.

Vote on the predicted class labels based on the classes of the k nearest neighbours. Here, the labels were predicted based on the k=3 nearest neighbours.

UNSUPERVISED LEARNING

CLUSTERING K MEANS

K-MEANS CLUSTERING

HOW IT WORKS ?

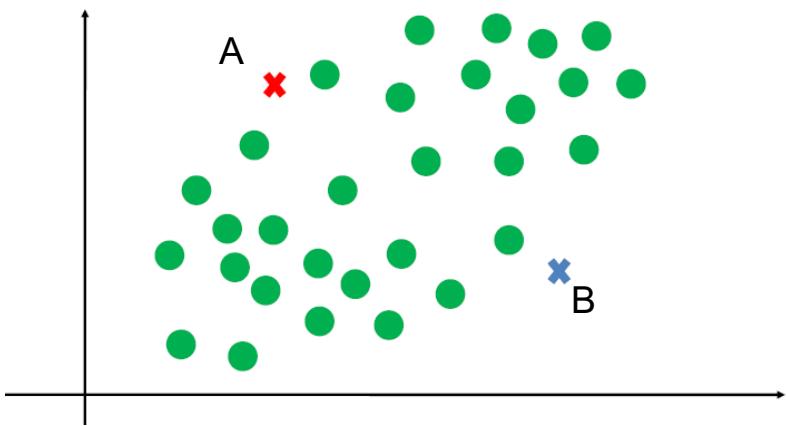
Simple and easy way to classify a given data set through a certain number of clusters (k clusters)

- Define k centers randomly , one for each cluster
- Take each point belonging to a given data set and associate it to the nearest center
- Re-calculate k new centroids as barycenter of the clusters resulting from the previous step
- Repeat this procedure. As a result of this, we may notice that the k centers change their location step by step until no more changes are done or in other words centers do not move any more.

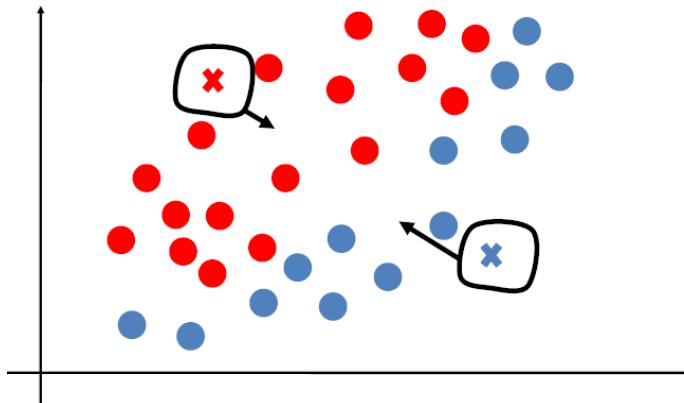
K-MEANS CLUSTERING

HOW IT WORKS ?

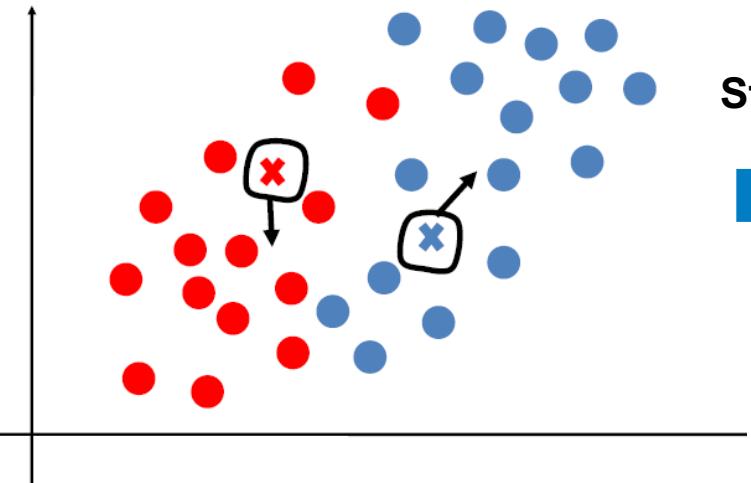
K = 2 (2 cluster)



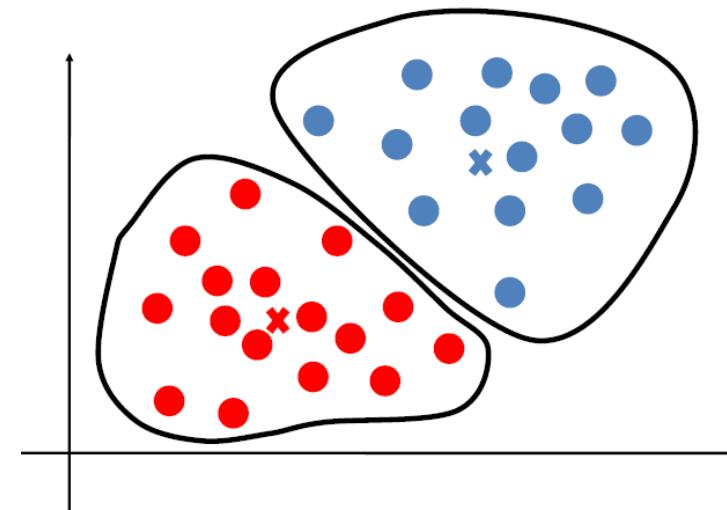
Step-1



Step-2

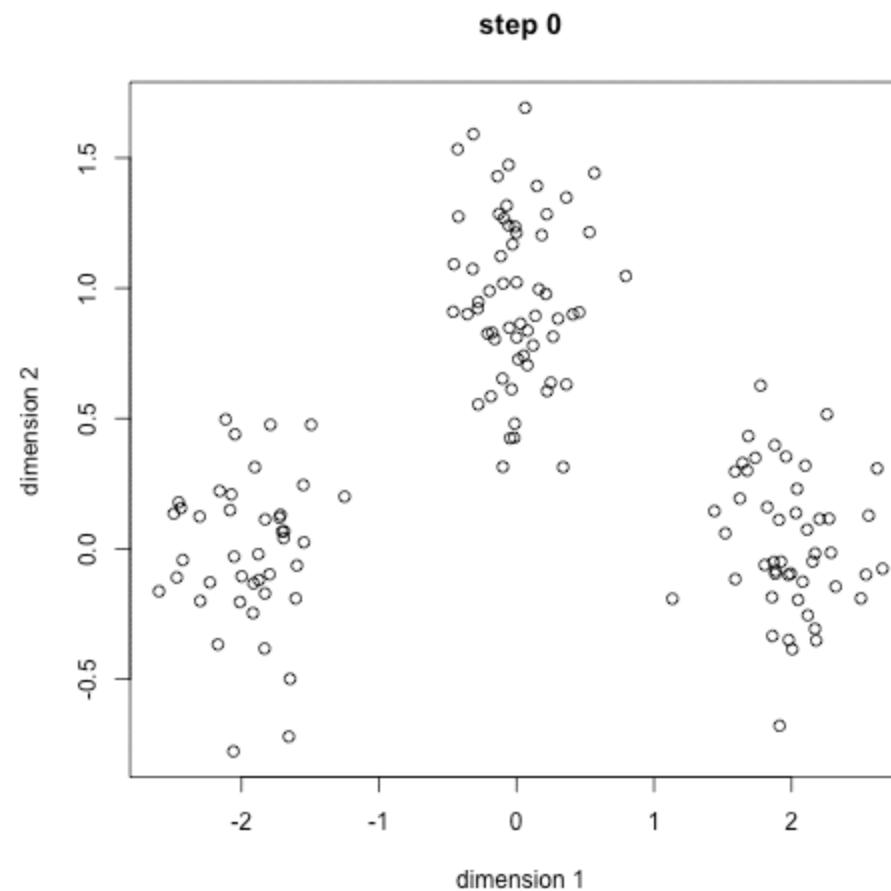


Step-n



K-MEANS CLUSTERING

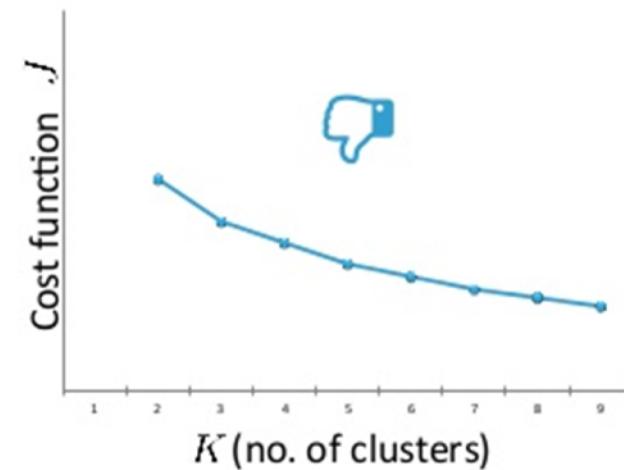
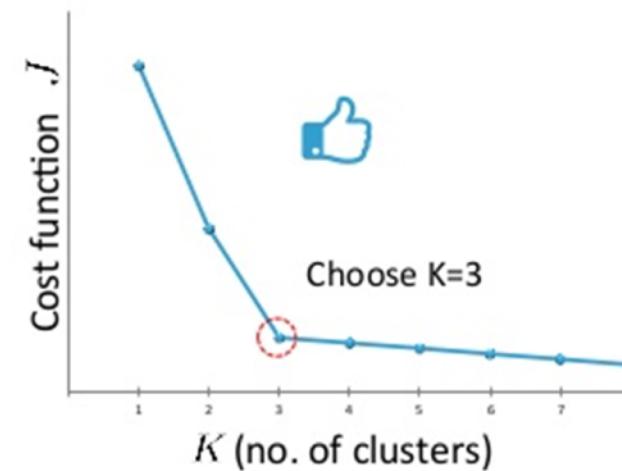
HOW IT WORKS ?



Elbow Criterion

- Draw graph of cost function (sum of square of variance within cluster) against the number of clusters
- The first clusters will not add much information (explain a lot of variance), but at some point the marginal gain will drop, giving an angle in the graph.
- The number of clusters are chosen at this point, hence the “elbow criterion”.

Elbow method:

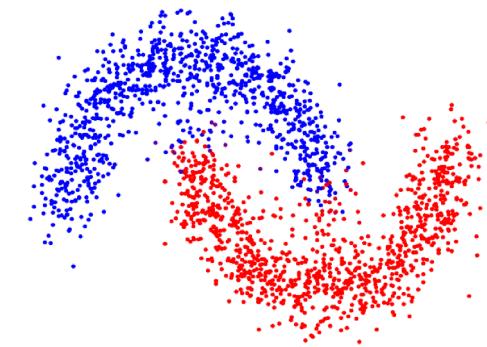
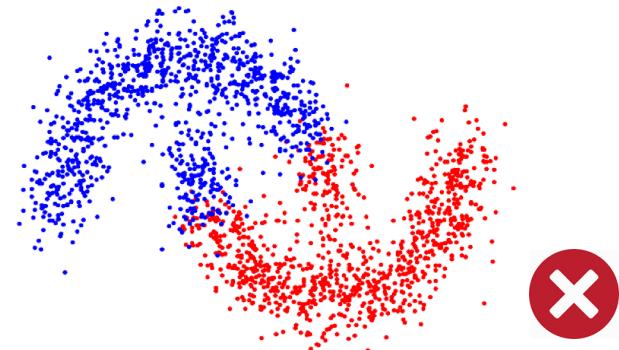


PROS

- Fast, robust and easier to understand.
- Gives best result when data set are distinct or well separated from each other.

CONS

- If there are two highly overlapping data then k-means will not be able to resolve that there are two clusters.
- Algorithm fails for non-linear data set.
- Applicable only when mean is defined i.e. fails for categorical data.
- Unable to handle noisy data and outliers.



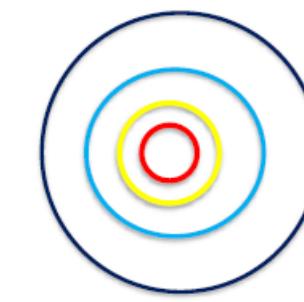
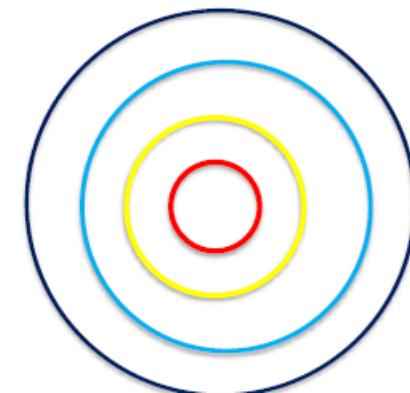
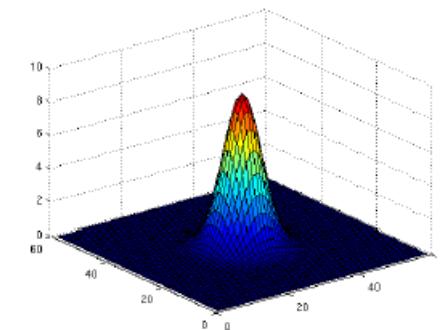
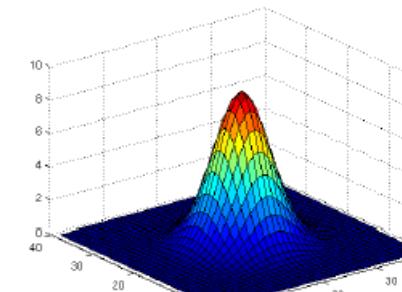
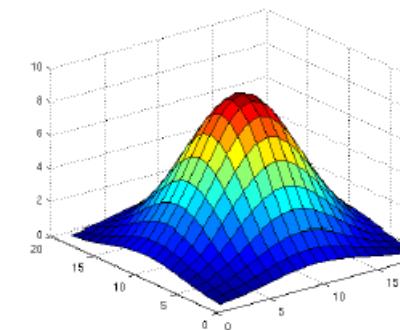
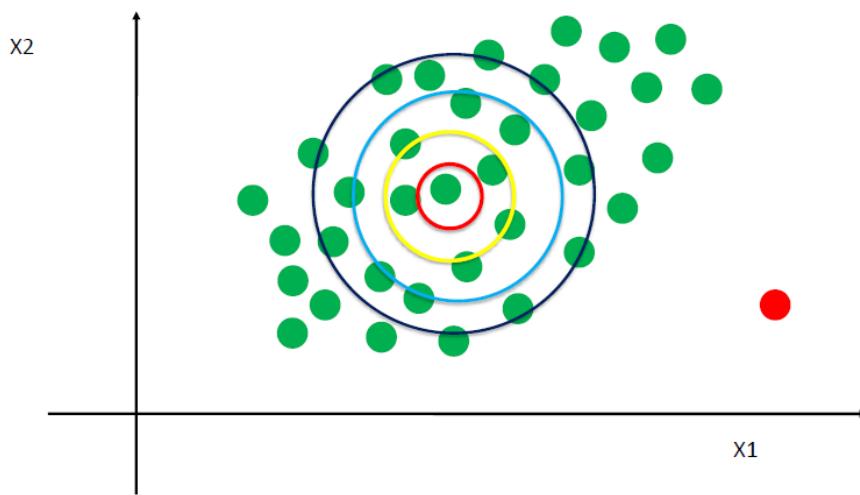
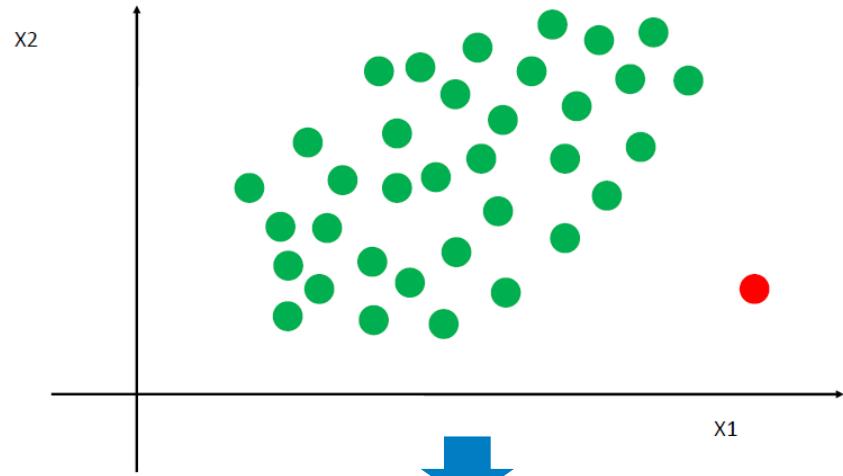
DBSCAN

Density-based spatial
clustering of
applications with noise

ANAMOLY DETECTION

ANAMOLY DETECTION

HOW IT WORKS

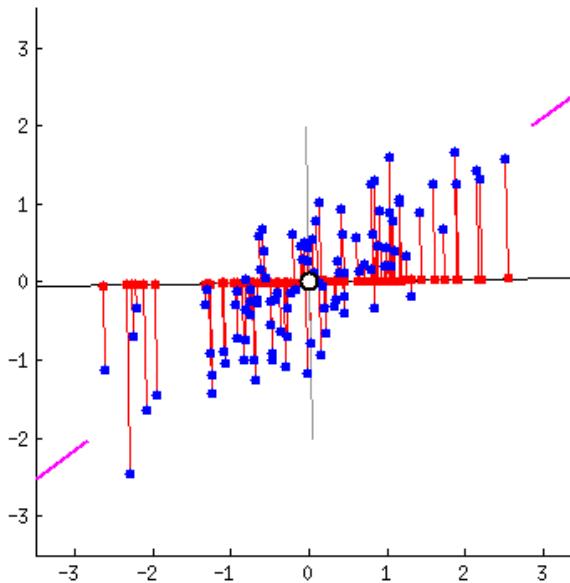
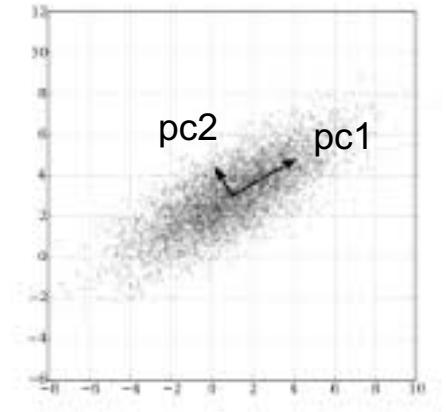


DIMENSIONALITY REDUCTION – PCA / SVD

PRINCIPAL COMPONENT

ANALYSIS

- PCA is transformation onto a new feature subspace
- Find directions of maximum variance
- Retain most of the information



The Math behind

$$\begin{bmatrix} p_{11} & p_{21} \\ \vdots & \vdots \\ \vdots & \vdots \\ p_{1n} & p_{2n} \end{bmatrix} = \begin{bmatrix} x_{11} & x_{21} \\ \vdots & \vdots \\ \vdots & \vdots \\ x_{1n} & x_{2n} \end{bmatrix} \begin{bmatrix} w_{11} & w_{21} \\ w_{12} & w_{22} \end{bmatrix}$$

With two dimensions

$$\mathbf{P} = \mathbf{X} \mathbf{W}$$

In general

w_{11} and w_{12} are the *loadings* corresponding to the first principle component.

w_{21} and w_{22} are the *loadings* corresponding to the second principle component.

It turns out that the columns of \mathbf{W}
are the eigenvalue vectors of the matrix $\mathbf{X}^T \mathbf{X}$

PRINCIPAL COMPONENT

DIMENSION REDUCTION

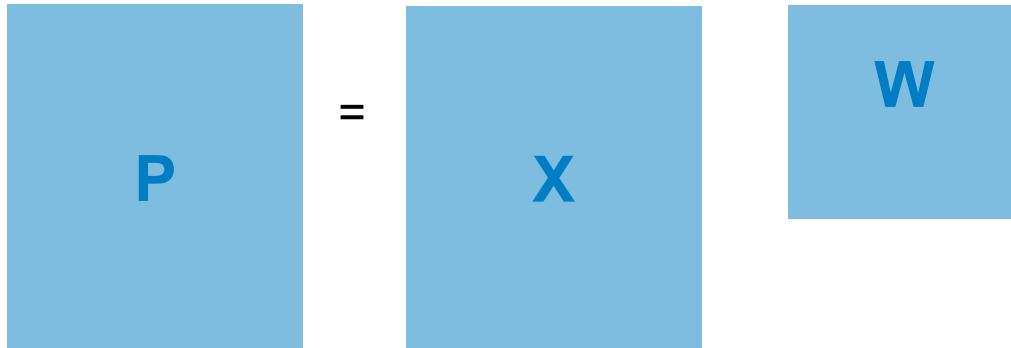
$$P = XW$$

Now only take the first L columns of W

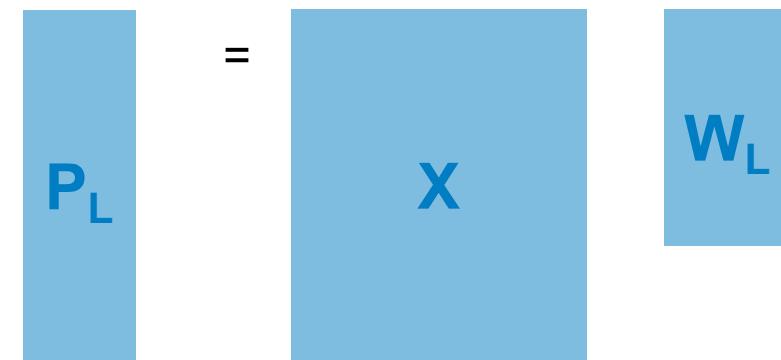
$$P_L = X W_L$$

For example for visualization only use the first 2 or 3 columns so that P_L only has 2 or 3 columns that can be visualized in scatter or contour plots

(10000 by 100) (10000 by 100) (100 by 100)



(10000 by 2) (10000 by100) (100 by2)



Applications of PCA

- ✓ Dimension reduction
- ✓ Data compression (communication system)
- ✓ Feature extraction (computer vision)
- ✓ Visualisation (machine learning)
 - ✓ → Outlier / anomalie detectie
- ✓ PCA regression
- ✓ Use PC instead of the original inputs

SINGULAR VALUE DECOMPOSITION

Matrix SVD decomposition: $A = U\Sigma V^T$

$$\begin{matrix} m \times n \\ A \end{matrix} = \begin{matrix} m \times n \\ U \end{matrix} \begin{matrix} n \times n \\ \Sigma \end{matrix} \begin{matrix} n \times n \\ V^T \end{matrix}$$

Σ Diagonal with r singular values
[could be a large number]

$$A = \begin{bmatrix} .96 & 1.72 \\ 2.28 & .96 \end{bmatrix} = U\Sigma V^T = \begin{bmatrix} .6 & -.8 \\ .8 & .6 \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} .8 & -.6 \\ .6 & -.8 \end{bmatrix}^T$$

- Any real $m \times n$ matrix A can be decomposed uniquely $A = U\Sigma V^T$.
- U is $m \times n$ and column orthonormal ($U^T \cdot U = I$) (the columns of U are called **Eigen vectors**)
- Σ is $n \times n$ and diagonal (**Covariance Matrix**)
 - σ_i are called singular values of A
 - It is assumed that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$
- V is $n \times n$ and orthonormal ($V \cdot V^T = V^T \cdot V = I$)

SINGULAR VALUE DECOMPOSITION

Matrix SVD decomposition:

$$A = U \Sigma V^T$$

$$\begin{matrix} A \\ A_k \end{matrix} = \begin{matrix} u_k & U \end{matrix}$$

$$\begin{matrix} \Sigma \\ \Sigma_k \end{matrix} \quad \begin{matrix} v^T_k \\ V^T \end{matrix}$$

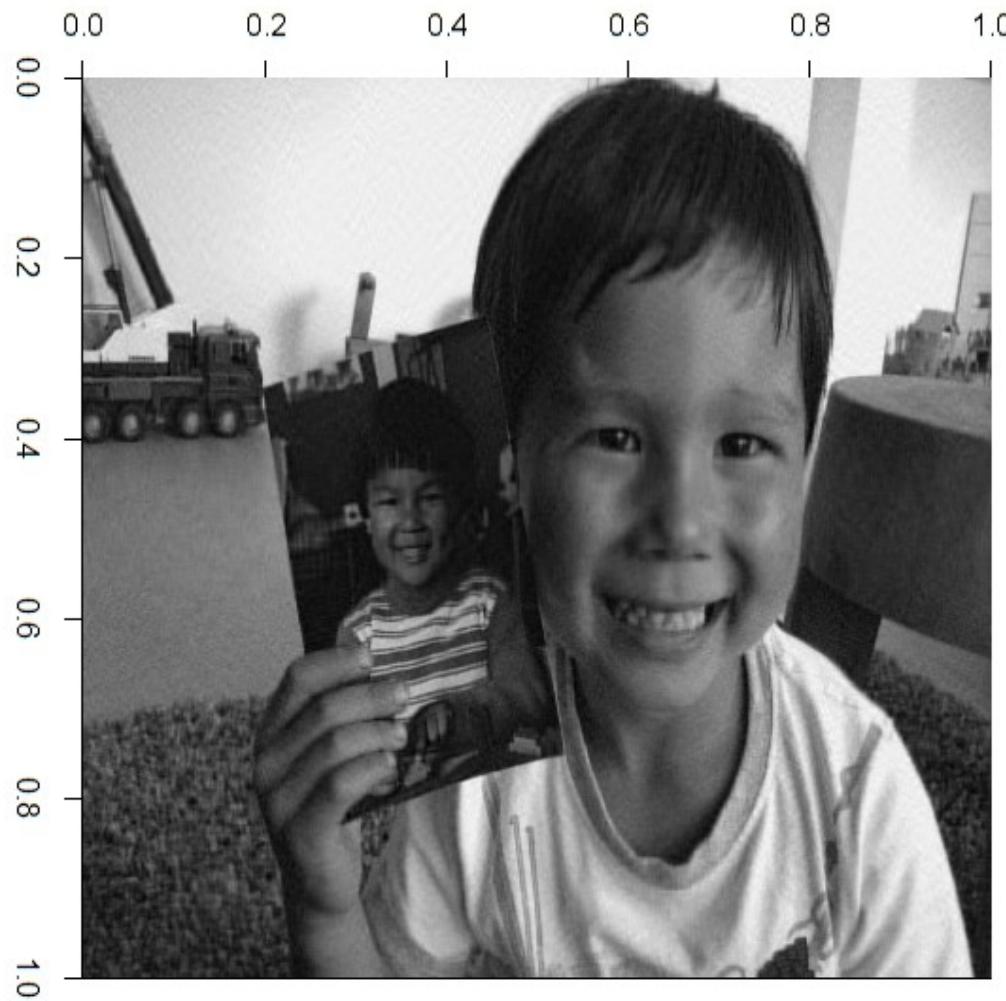
Σ Diagonal with r singular values
[could be a large number]

Σ_k Take only $k << r$ singular values

A datapoint d can now be represented by k dimensional point

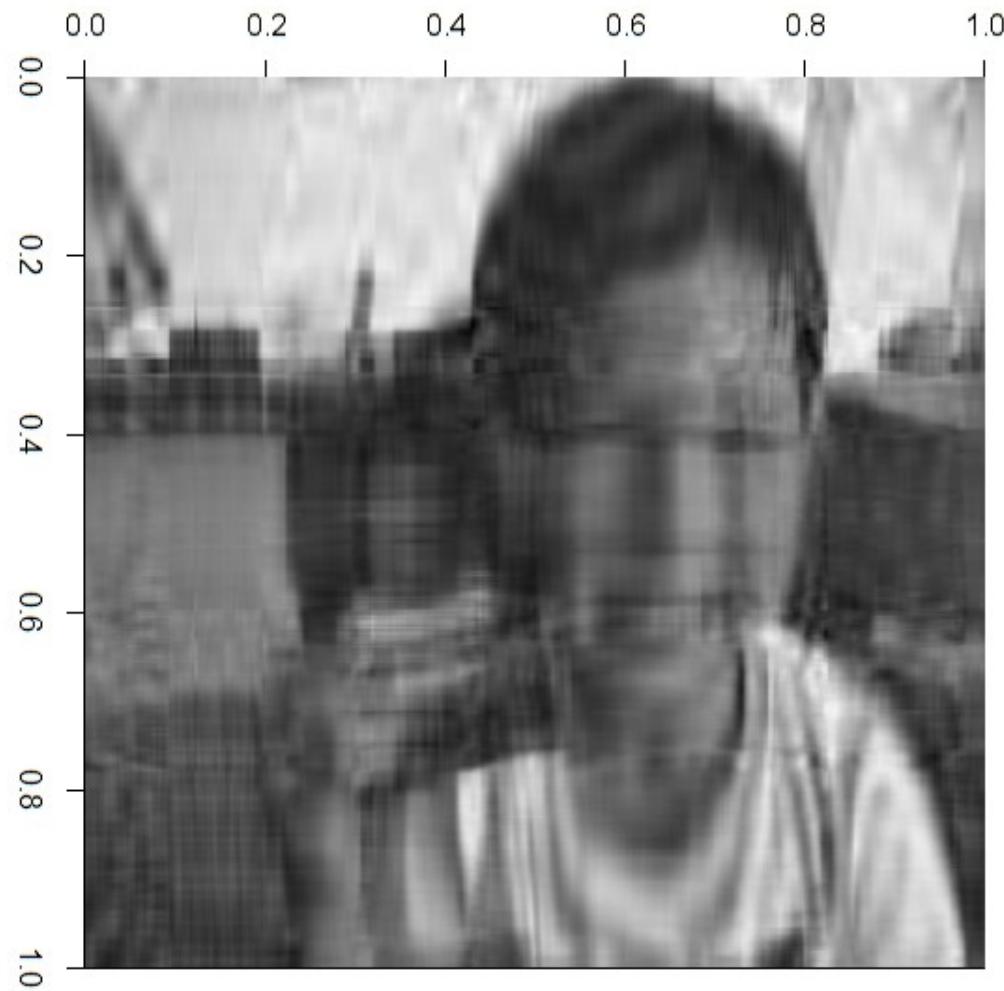
$$\hat{d} = U_k^T d.$$

SVD EXAMPLE | AS AN EXPERIMENT



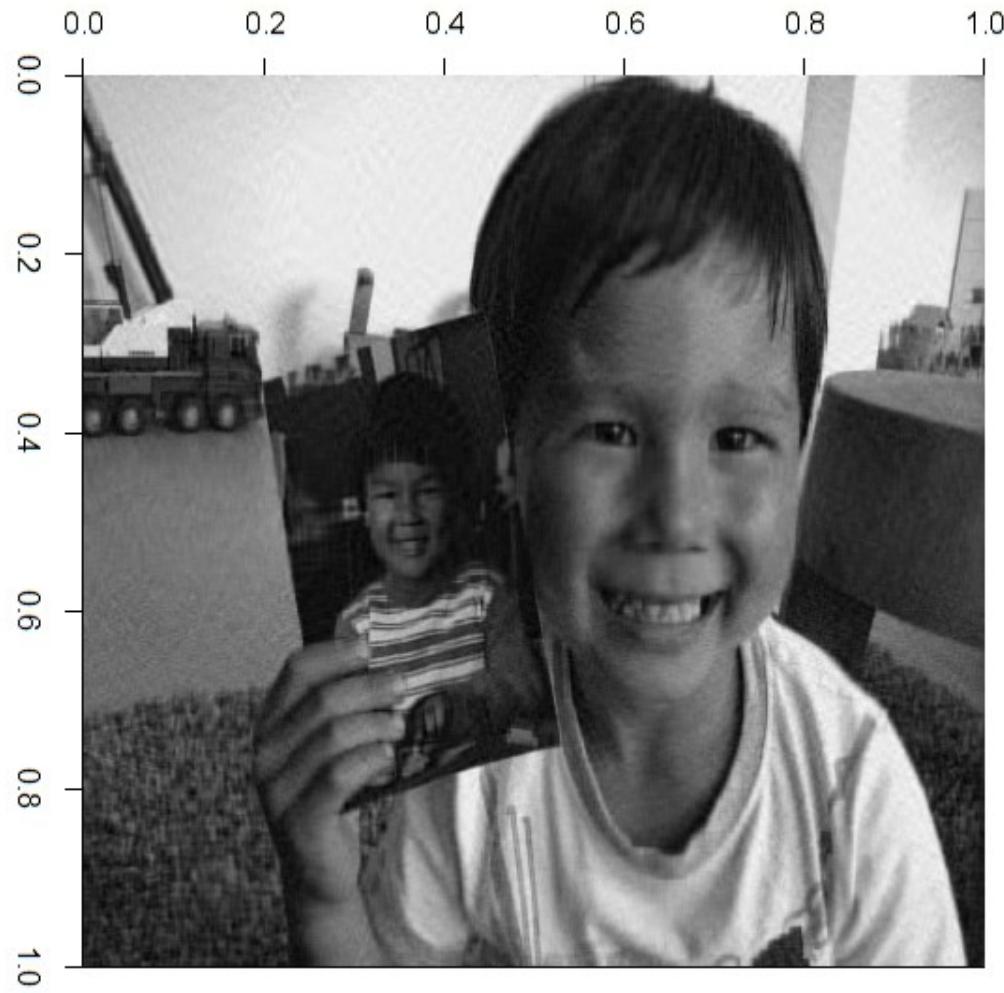
Original
2448 X 3264 ~ 8 mln numbers

SVD EXAMPLE AS AN EXPERIMENT

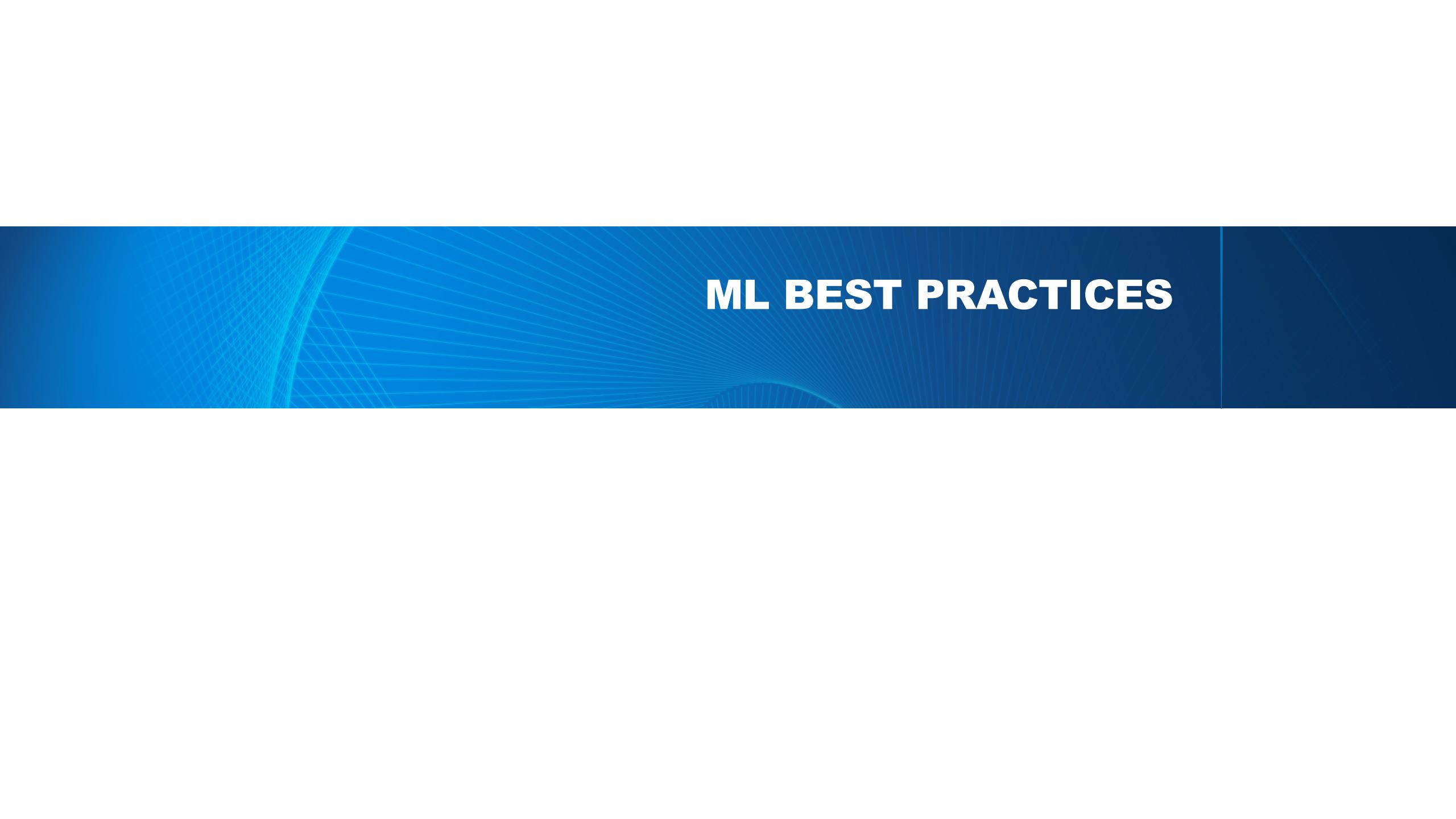


**SVD: 15 largest SV's
1% of the data**

SVD EXAMPLE AS AN EXPERIMENT



**SVD: 75 largest V's
5% of the data**

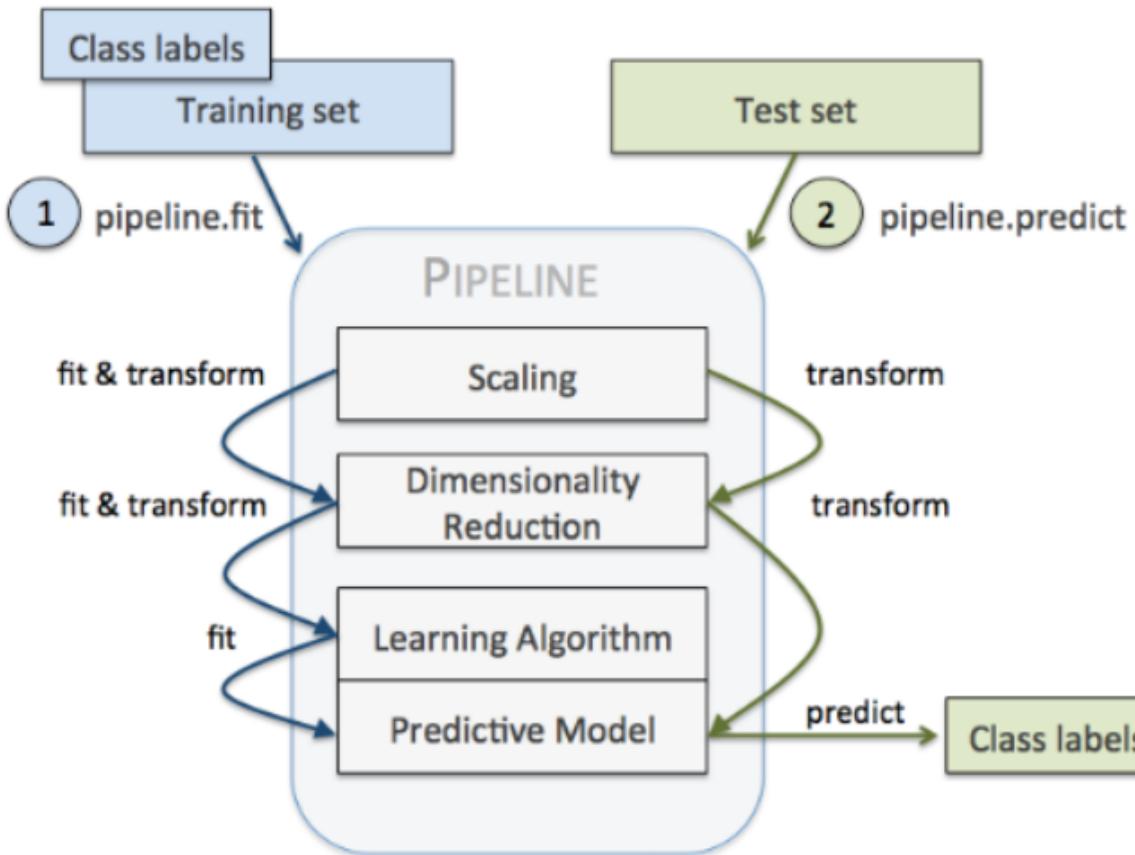


ML BEST PRACTICES

BEST PRACTICES STANDARDISATION AND OUTLIERS

- Data set has features whose values are significantly different in magnitude and range.
 - This disparity can degrade the performance
 - To mitigate this transform your interval feature values to be on a similar scale, commonly standardizing to have mean 0 and variance 1
-
- You will often find outliers in your data—observations that are very distinct from the others in one or more of the feature values.
 - Although outliers can certainly be very informative and can identify anomalies that deserve special attention, they can be quite detrimental to training an effective generalizable model.
 - In general, outliers drag clusters artificially toward the outside of your feature space
 - Determine whether an outlying value is simply an invalid or erroneous entry that can be disregarded

BEST PRACTICES PIPELINES



```
import numpy as np
from sklearn.svm import SVR
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler

x,y = np.load('data.npz')
x_test = np.linspace(0, 200)

model = Pipeline([
    ('standardize', StandardScaler()),
    ('svr', SVR(kernel='rbf', verbose=0, C=5e6,
                epsilon=20))])
model.fit(x[:, np.newaxis], y)
y_test = model.predict(x_test[:, np.newaxis])
```

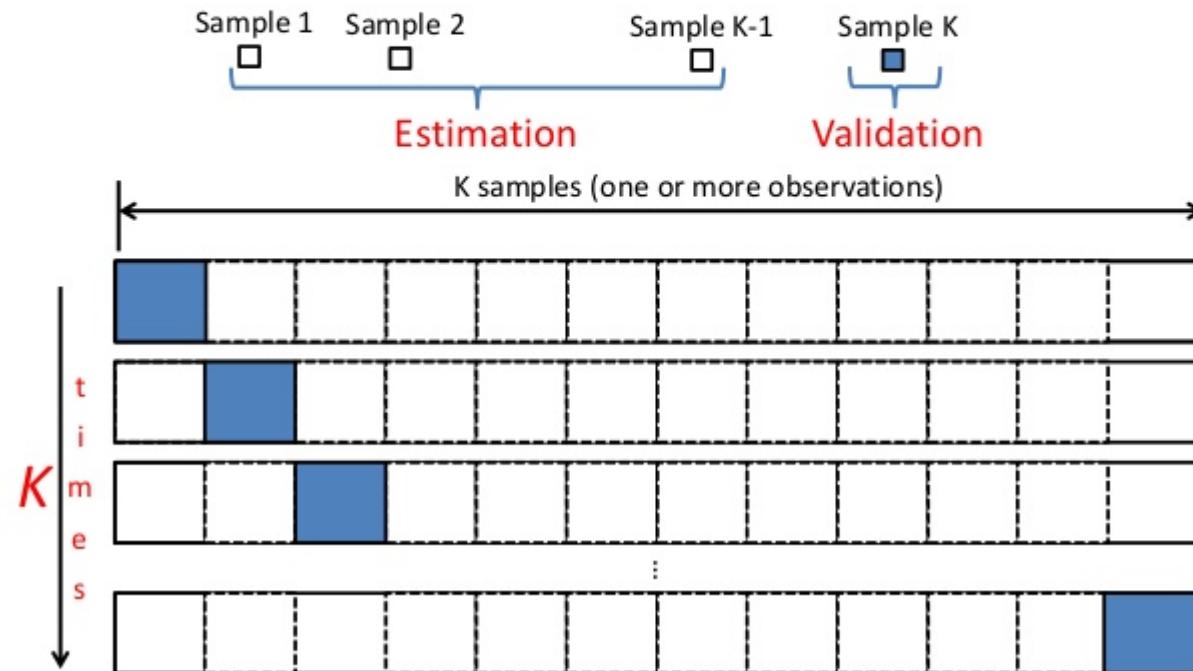
regularization term

Pipelines - Sequence of data transforms to be chained together culminating in a modeling process that can be evaluated.

BEST PRACTICES K- FOLD CROSS VALIDATION

- Cross-validation avoids overlapping test sets
 - First step: split data into k subsets of equal size
 - Second step: use each subset in turn for testing, the remainder for training
- Called k-fold cross-validation ($K = 10$)
- The process is repeated K times (folds)
- Often the subsets are stratified (ranked) before the cross-validation is performed
- The error estimates are averaged to yield an overall error estimate

- K-fold cross-validation:



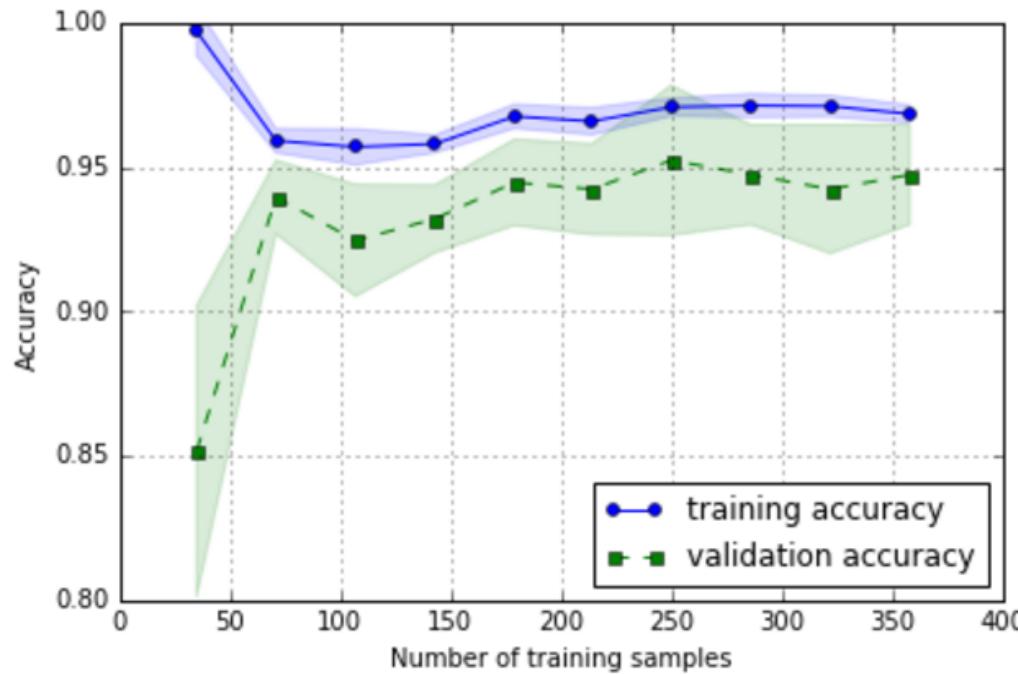
BEST PRACTICES

PERFORMANCE EVALUATION MATRIX – GRID SEARCH

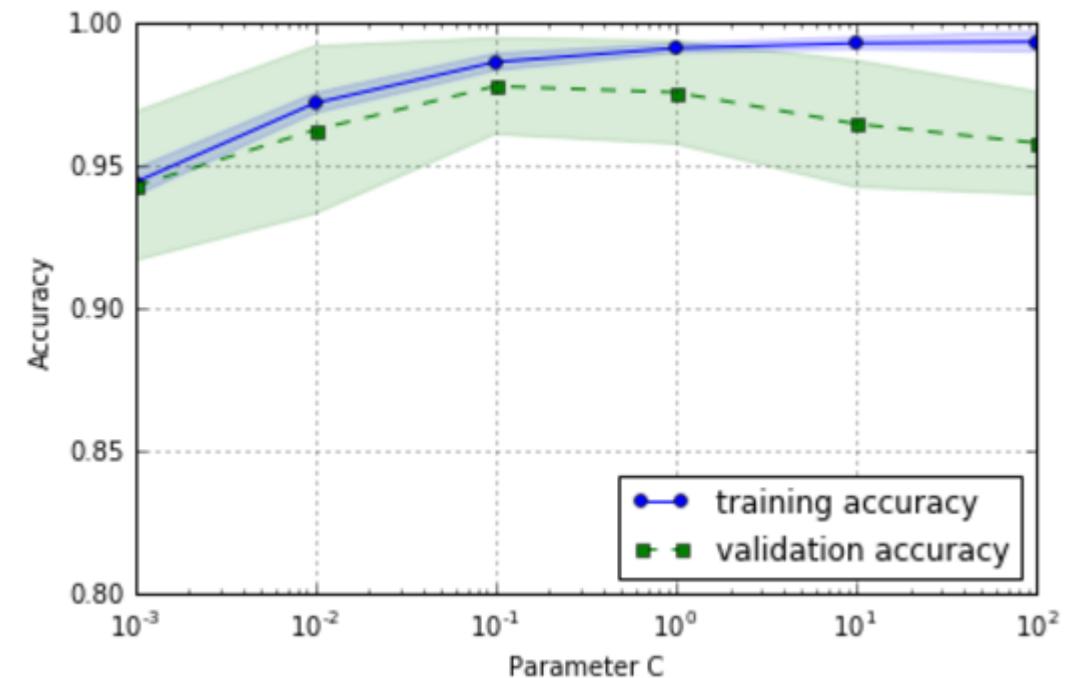
The general belief is that more samples are required for better accuracy.

Demo

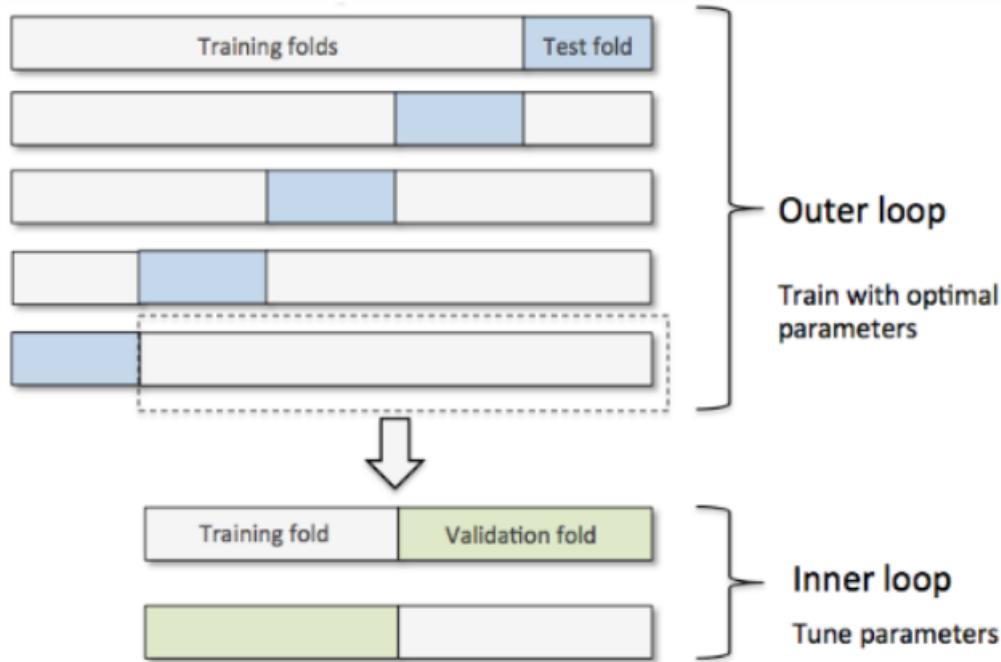
Learning Curve



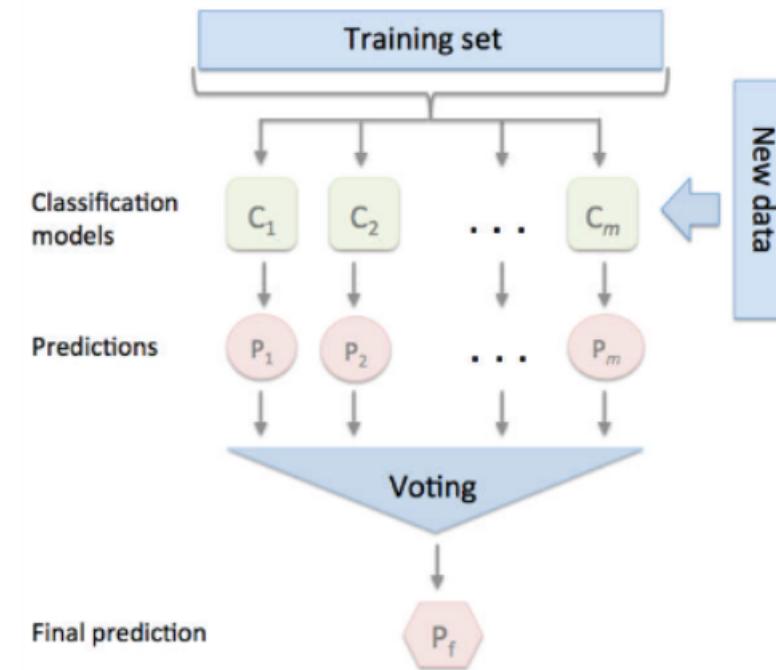
Validation Curve



Grid Search + Cross Validation

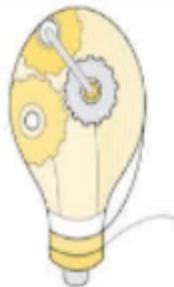


Majority voting Classifier



BEST PRACTICES

MACHINE LEARNING PLATFORMS



Google Prediction API



THANK YOU

PURUSHOTTAM DARSHANKAR