

# Day 2

# Welcome to Day 2 Spark Machine Learning

Please perform PRE-WORK

1. Access the virtual Lab using link <https://html.inspiredvlabs.com> Use the username TEKMD190-XX (replace XX with your number) and password

**TEKmd189!23**

<https://tinyurl.com/bdtSparkML>

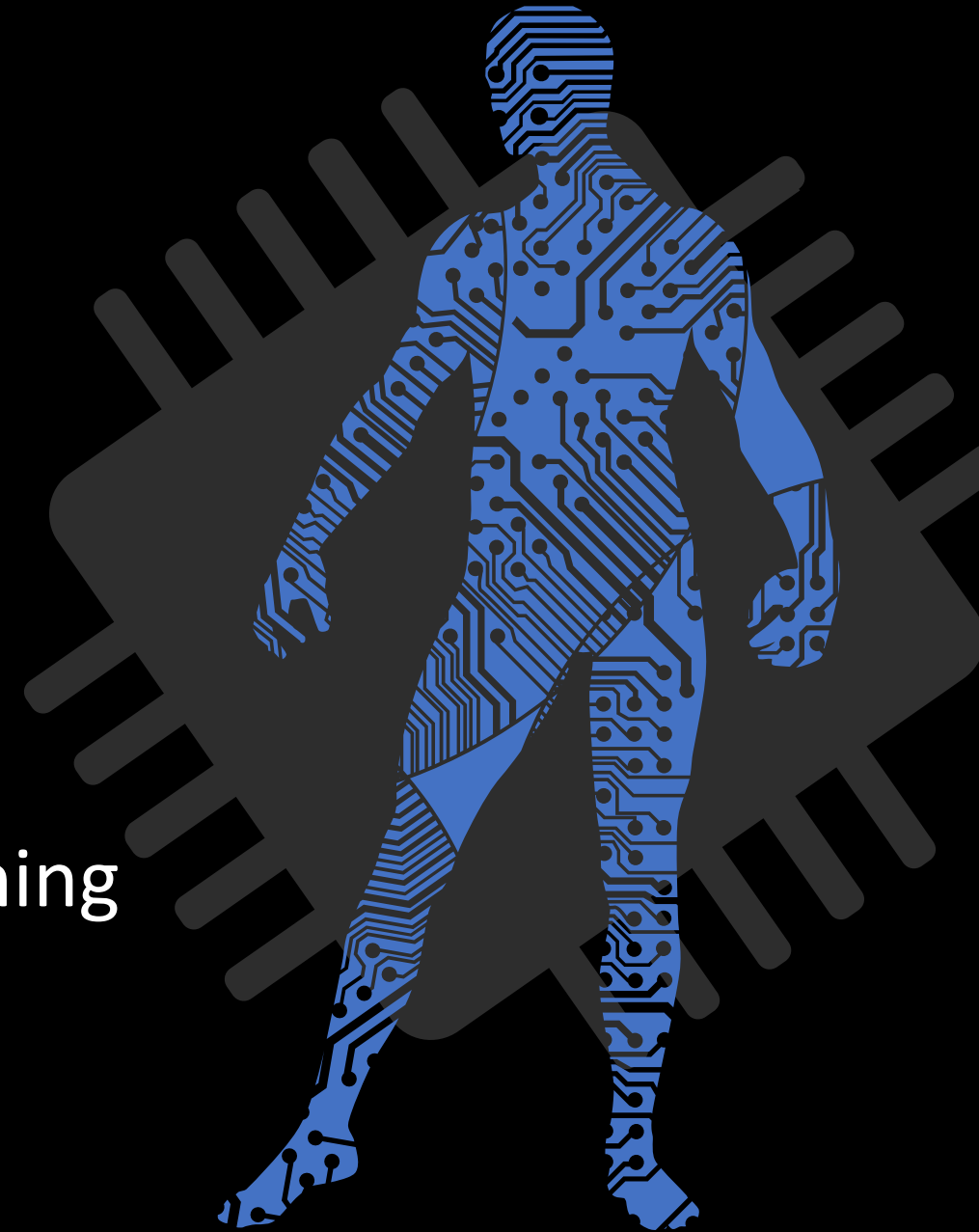
<https://tinyurl.com/bdtSparkMLSlides>

Last Name	First Name	User-Name
ALAMELU KRISHNAN	SURESH	TEKMD190-01
BRAGG	JAMES	TEKMD190-02
GARIMELLA	HIMABINDU	TEKMD190-03
JAJOO	SANDESH	TEKMD190-04
JOSHI	CECIL	TEKMD190-05
LANNERS	QUINN	TEKMD190-06
MADAS	ANANT	TEKMD190-07
MAXSON	CRAIG	TEKMD190-08
MISRA	SMARAN	TEKMD190-09
PASUMARTHY	VENKATA	TEKMD190-10
POUDYAL	BASHUDEV	TEKMD190-11
RUTHERFORD IV	ROB ROY	TEKMD190-12
SAMALLA	APARNA	TEKMD190-13
SCHLEY	DANIEL	TEKMD190-14
VONGSOUVANH	BOUAPHAYCHIT	TEKMD190-15
ZHENG	LEI	TEKMD190-16

We will be starting soon

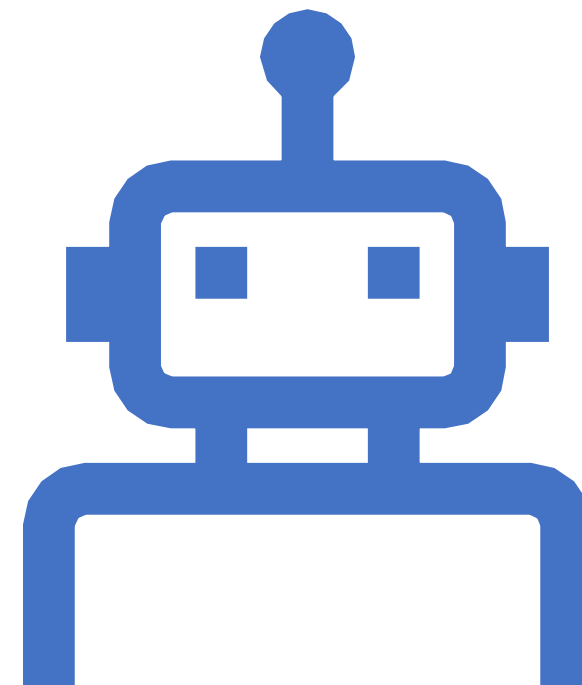
# Agenda – Day 2

1. Recap of Day 1
2. Use Case Real Estate (Homework)
3. Project Initiation
4. Spark ML Processing
5. Feature Engineering and Data Cleaning
6. More Algorithms
7. Model Evaluations and Metrics



# Recap – Day 1

1. What is Machine Learning?
2. Spark Overview (Ecosystems before and after)
3. Spark ML Development
4. Machine Learning Techniques (CCRA)
  - a. Supervised
  - b. Unsupervised
5. Machine Learning Development (DIAPERS)
6. Data (Clean, Coverage, Complete)
7. Statistics Brush up
8. Popular Algorithms
9. Linear Regression
10. Multiple hands-on exercises



# Use case: Boston Real Estate Data

- Dataset – from scikit-learn datasets
- Linear Regression example with following features:

**CRI:** Per capita crime rate by town

**ZN:** Proportion of residential land zoned for lots over 25,000 sq. ft

**INDUS:** Proportion of non-retail business acres per town

**CHAS:** Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)

**NOX:** Nitric oxide concentration (parts per 10 million)

**RM:** Average number of rooms per dwelling

**AGE:** Proportion of owner-occupied units built prior to 1940

**DIS:** Weighted distances to five Boston employment centers

**RAD:** Index of accessibility to radial highways

**TAX:** Full-value property tax rate per \$10,000

**PTRATIO:** Pupil-teacher ratio by town

**B:**  $1000(B_k - 0.63)^2$ , where  $B_k$  is the proportion of [people of African American descent] by town

**LSTAT:** Percentage of lower status of the population

**PRICE:** Median value of owner-occupied homes in \$1000s



**Spark Lab/Spark Linear Regression Real Estate**

# Project Initiation Income Prediction

# Spark Income Prediction

- Objective is to predict if a person's income will be  $\geq 50K$  or  $< 50K$  using number of features
- Code examples for Transformer, Estimator, Evaluator and Pipeline
- Will implement different algorithms with this dataset

# Data Scaling



# Data Preparation for Machine Learning

## Summary

- In this example, we will look at different ways of scaling the data
- ***MinMaxScaler*** – set the range 0 to 1 to scale values
- ***StandardScaler*** – standardize so that the mean is 0 and standard deviation is 1
- ***Normalizer*** – Rescale each sample
- ***Train and Test split*** – split the data into training set and testing set

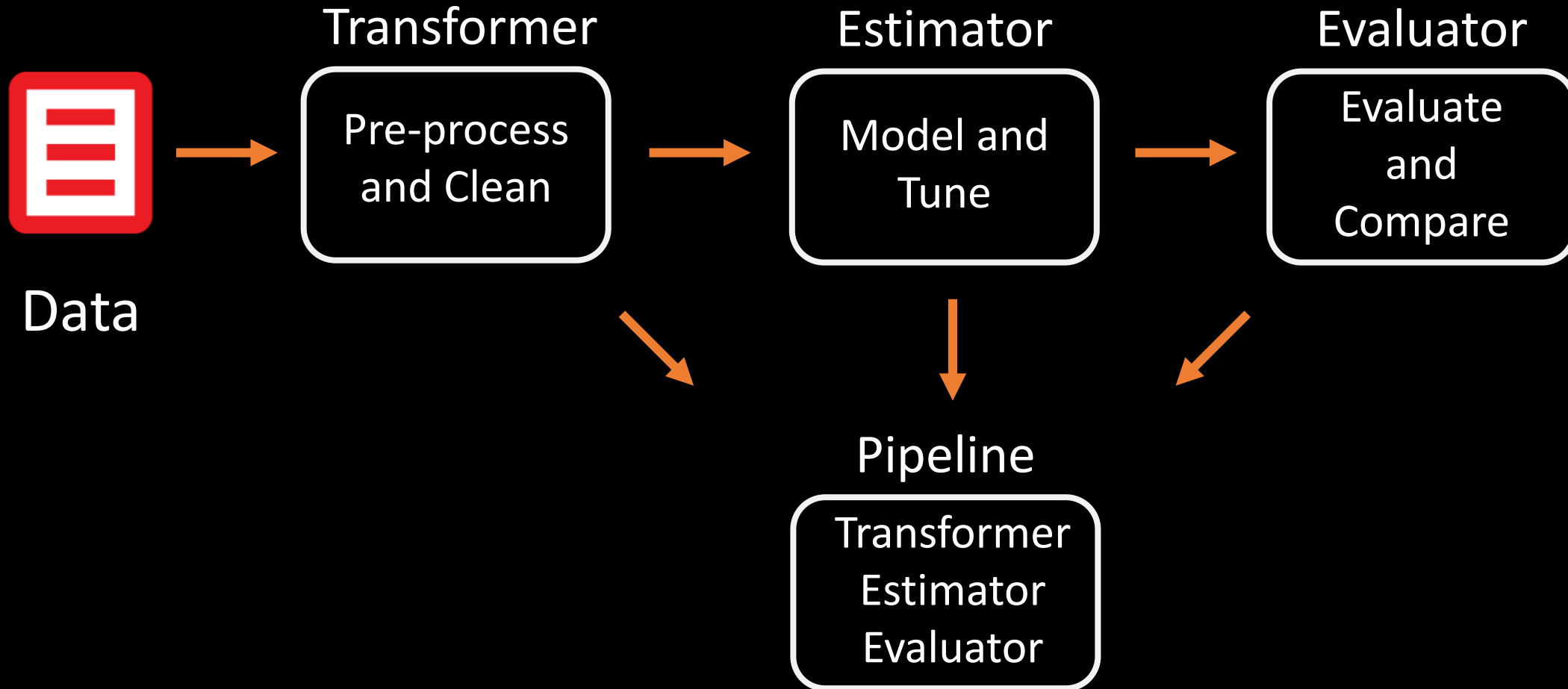
# Spark ML Processing

# Spark ML



Data

# Spark ML



# Evaluator Types

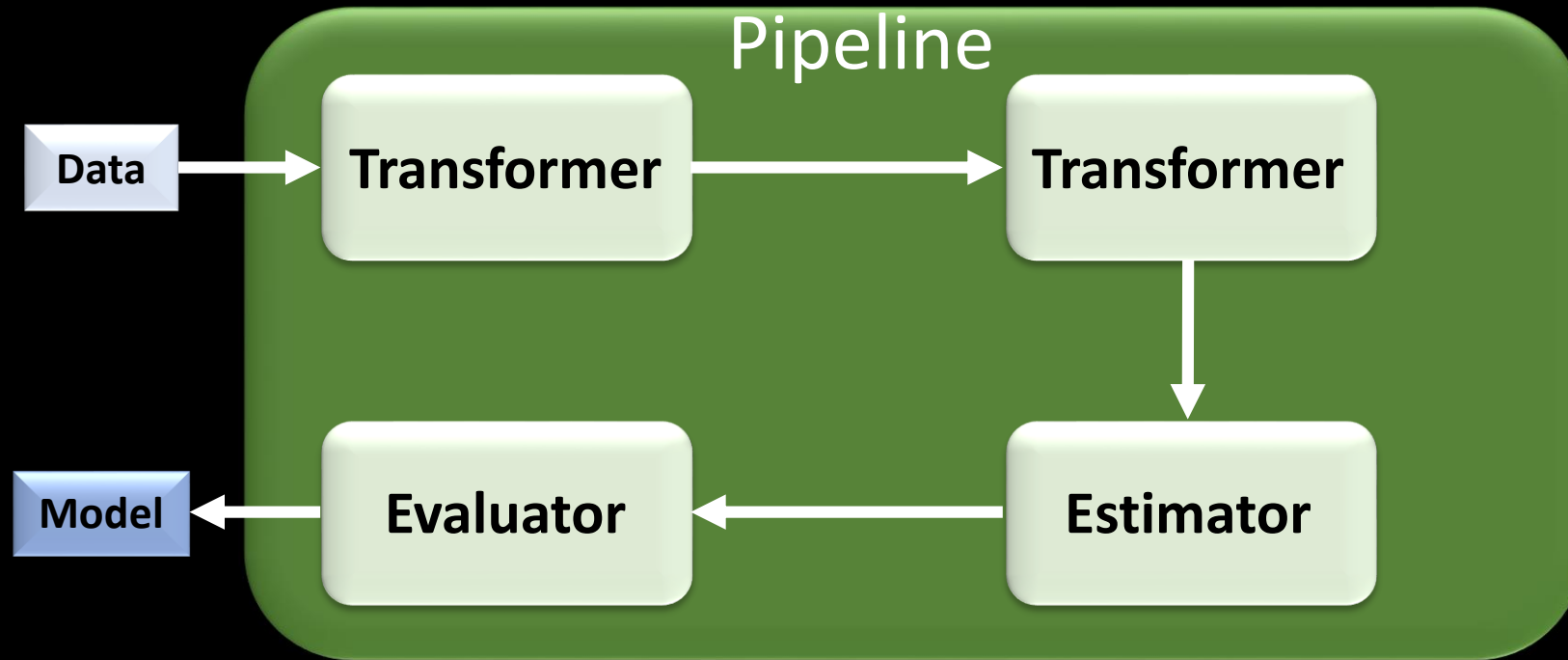
Evaluator	Description
BinaryClassifierEvaluator	Binary Classification model evaluator
MultiClassClassificationEvaluator	Multiple Class Classification model evaluator
RegressionEvaluator	Regression model evaluator
ClusteringEvaluator	Clustering model evaluator

```
from pyspark.ml.evaluation import BinaryClassificationEvaluator
```

```
# Evaluate model
```

```
evaluator = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction")  
evaluator.evaluate(predictions)
```

# Pipeline

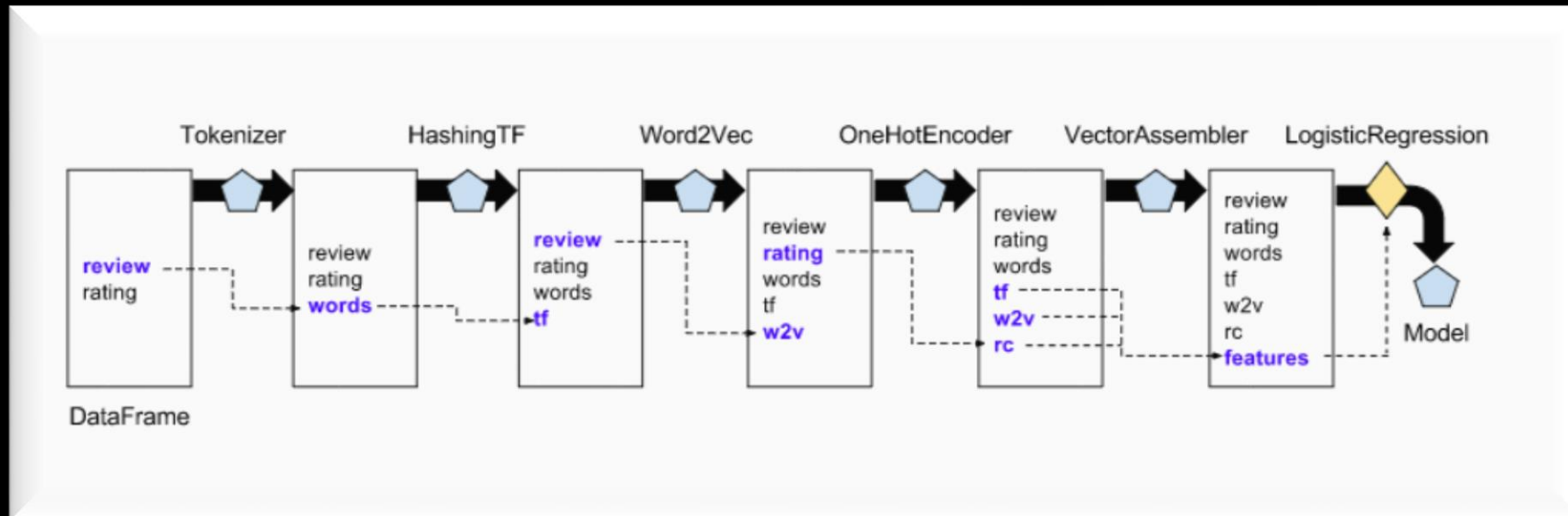


- A machine learning work-flow
- Made up of number of stages
- Can be persisted

# Pipeline Example

# Configure pipeline stages

```
tok      = Tokenizer(inputCol="review", outputCol="words")
htf      = HashingTF(inputCol="words", outputCol="tf", numFeatures=200)
w2v      = Word2Vec(inputCol="review", outputCol="w2v")
ohe      = OneHotEncoder(inputCol="rating", outputCol="rc")
va       = VectorAssembler(inputCols=["tf", "w2v", "rc"], outputCol="features")
lr       = LogisticRegression(maxIter=10, regParam=0.01) # Build the pipeline
pipeline = Pipeline(stages=[tok, htf, w2v, ohe, va, lr]) # Fit the pipeline
model    = pipeline.fit(train_df)
```



# Project - Dataframe

Age	Sex	Race	Income
39	Male	White	<=50K
50	Female	Black	>50K
38	Female	Asian	<=50K
53	Male	Other	>50K



# Project - StringIndexer

Age	Sex	Race	Income
39	0	0	0
50	1	1	1
38	1	2	0
53	0	3	1

# Project – One-Hot Encoding Estimator

Age	Sex	Race	Race_ White	Race_ Black	Race_ Asian	Race_ Other	Income
39	0	0	1	0	0	0	0
50	1	1	0	1	0	0	1
38	1	2	0	0	1	0	0
53	0	3	0	0	0	1	1

# Project – Sparse Vector

Age	Sex	Race	Race_White	Race_Black	Race_Asian	Race_Other	Income
39	0	0	1	0	0	0	0
50	1	1	0	1	0	0	1
38	1	2	0	0	1	0	0
53	0	3	0	0	0	1	1

# Sparse Vector Encode: [0, 8, [0,3], [39, 1]]

Age	Sex	Race	Race_White	Race_Black	Race_Asian	Race_Other	Income
39	0	0	1	0	0	0	0
50	1	1	0	1	0	0	1
38	1	2	0	0	1	0	0
53	0	3	0	0	0	1	1

# Sparse Vector Encode: [0, 8, [0,1,4,7], [50,1,1,1] ]

Age	Sex	Race	Race_White	Race_Black	Race_Asian	Race_Other	Income
39	0	0	1	0	0	0	0
50	1	1	0	1	0	0	1
38	1	2	0	0	1	0	0
53	0	3	0	0	0	1	1

# Sparse Vector Assembler:

Age	Sex	Race	Race_White	Race_Black	Race_Asian	Race_Other	Income
39	0	0	1	0	0	0	0
50	1	1	0	1	0	0	1
38	1	2	0	0	1	0	0
53	0	3	0	0	0	1	1

**[0, 8, [0,3], [39,1 ] .....[0,8 ], [0,6,7], [53, 1, 1]]**

# Spark ML

## Training

**model.fit()**

Learn from data

## Score the Model

**model.evaluate()**

Compare predictions with actual labels

## Evaluate with Test Data

**model.transform()**

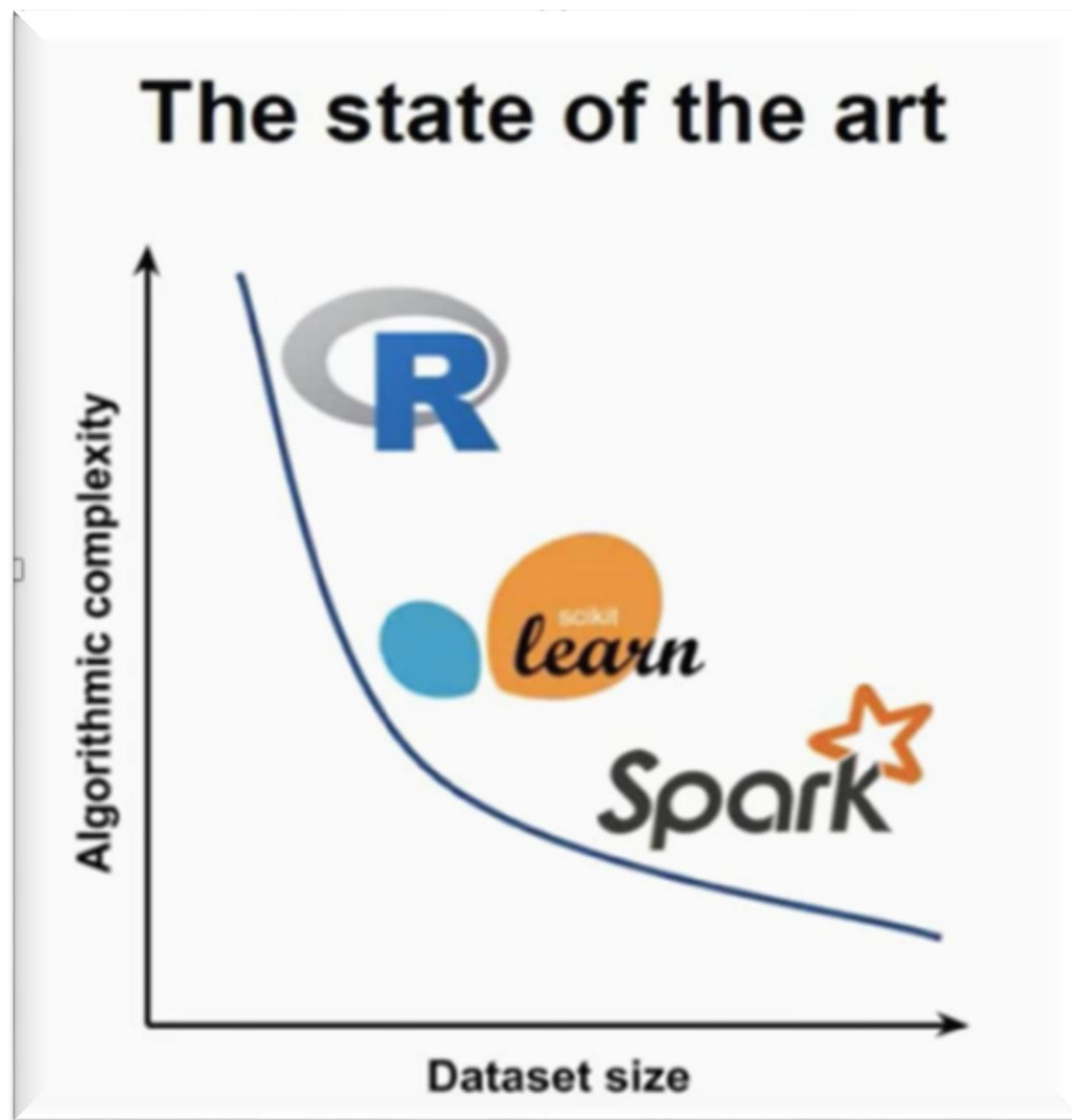
Make predictions

# Spark and Sci-kit Learn

Function	Spark ML	Scikit-Learn
Small and Medium datasets (in megabytes)	Spark will definitely perform better (depend on your deployment environment)	For larger datasets scikit-learn will require lot of RAM
Distributed Deployment	Spark is designed for it	Need distributed support
Model building	Less flexible – does not have lot of APIs. You will have to write more code	More flexible and efficient
Visualization	Cumbersome to implement	Hands-down more elegant support



# Spark v/s Other ML Implementations



# Feature Engineering and Data Cleaning



# Titanic Dataset

- Will use titanic data set to predict whether a passenger survived or not
- Data consists of following:
  - Pclass - Passenger Class (1 = 1st; 2 = 2nd; 3 = 3rd)
  - sex – Gender
  - sibsp - Number of siblings/Spouses
  - parch - Number of parents/children
  - fare - travel fare
  - Embarked - (C = Cherbourg; Q = Queenstown; S = Southampton)
  - boat - Life boat
  - body - Body identification number
  - home.dest – Destination
  - ticket - Ticket number
  - cabin - Cabin number
  - name - Passenger name
  - survived - (0: No, 1: Yes)

# Data Cleaning & Feature Engineering

## Summary

- In this example, we are going perform cleaning of data
- Look for missing values
- Binary categorical data replacement e.g. Male, Female
- N-nary categorical data replacement e.g. Embarked
- Using Imputer to replace missing age values
- Drop an entire sample due to many missing features

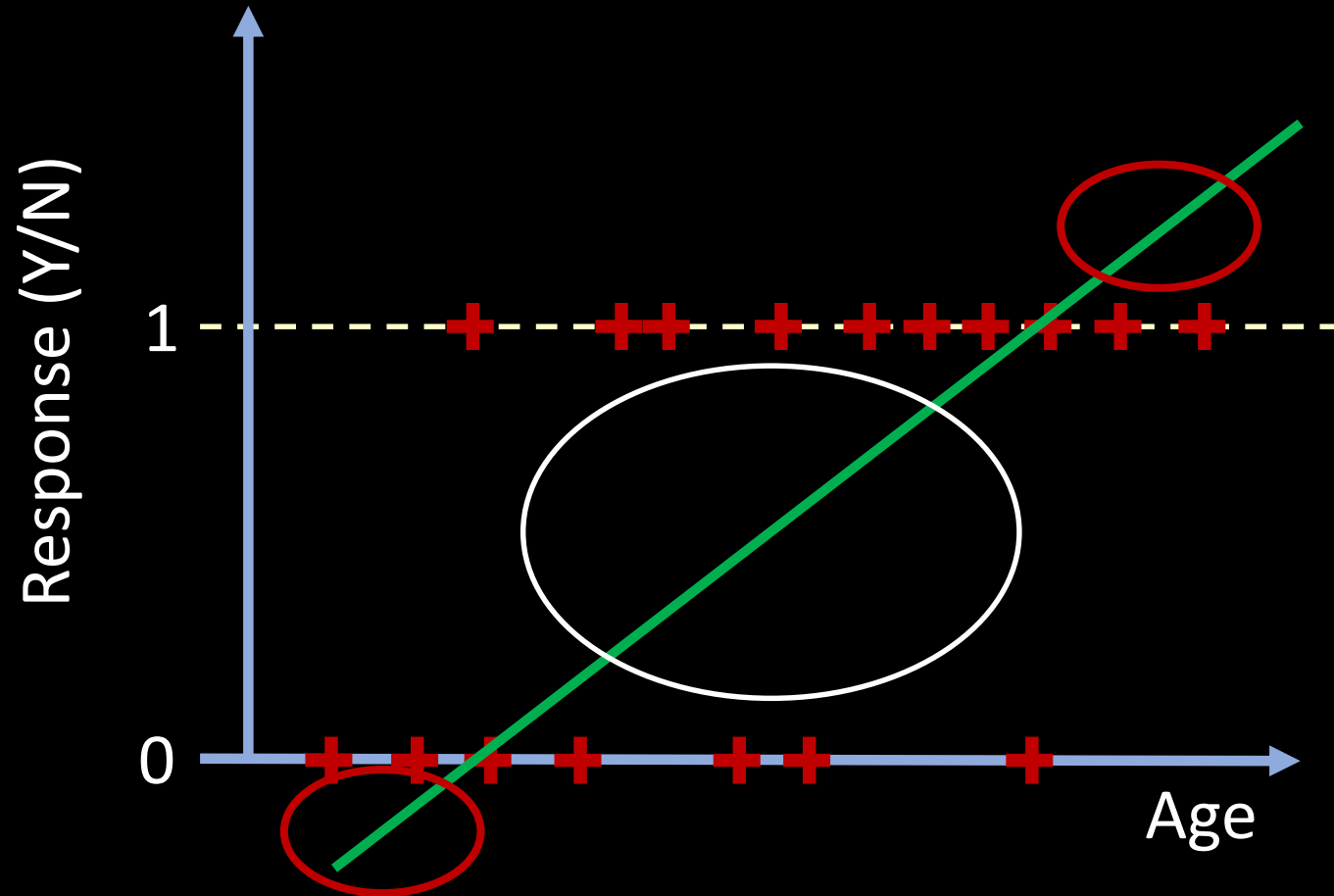
# Logistic Regression

# Example: Logistic Regression

## Social Network Ads

User ID	Gender	Age	Estimated	
			Salary	Purchased
15624510	Male	19	19000	Yes
15810944	Male	35	20000	No
15668575	Female	26	43000	No
15603246	Female	27	57000	Yes
15804002	Male	19	76000	No
15728773	Male	27	58000	No
15598044	Female	27	84000	No
15694829	Female	32	150000	Yes
15600575	Male	25	33000	No
15727311	Female	35	65000	No
15570769	Female	26	80000	No
15606274	Female	26	52000	No

# Linear Regression Challenges



# Logistic Regression

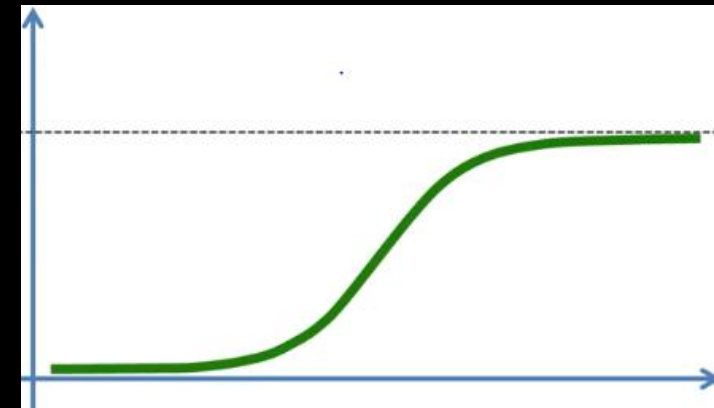
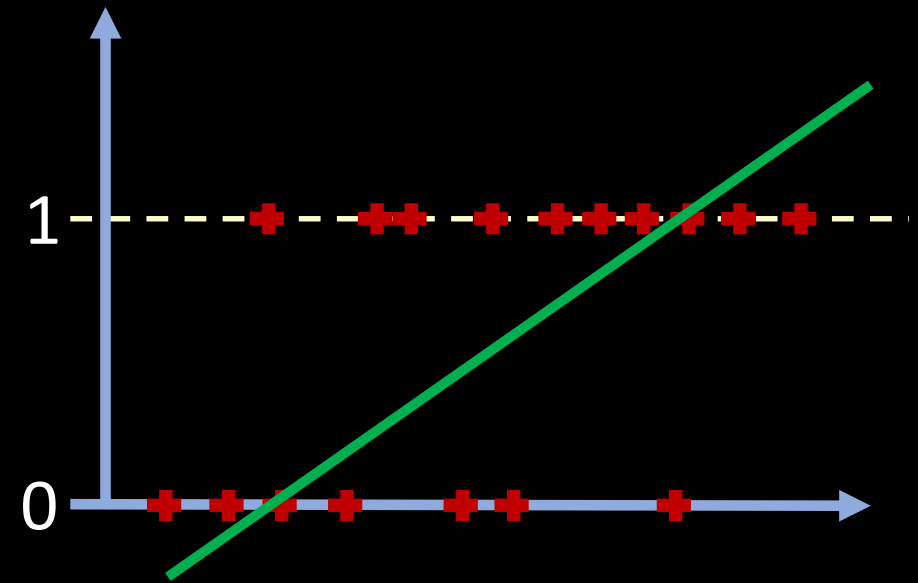
$$y = b_0 + b_1 * x_1$$

Sigmoid Function

$$p = \frac{1}{1 + e^{-y}}$$

$$\ln \left( \frac{p}{1 - p} \right) = b_0 + b_1 * x$$

This is the formula for logistic regression

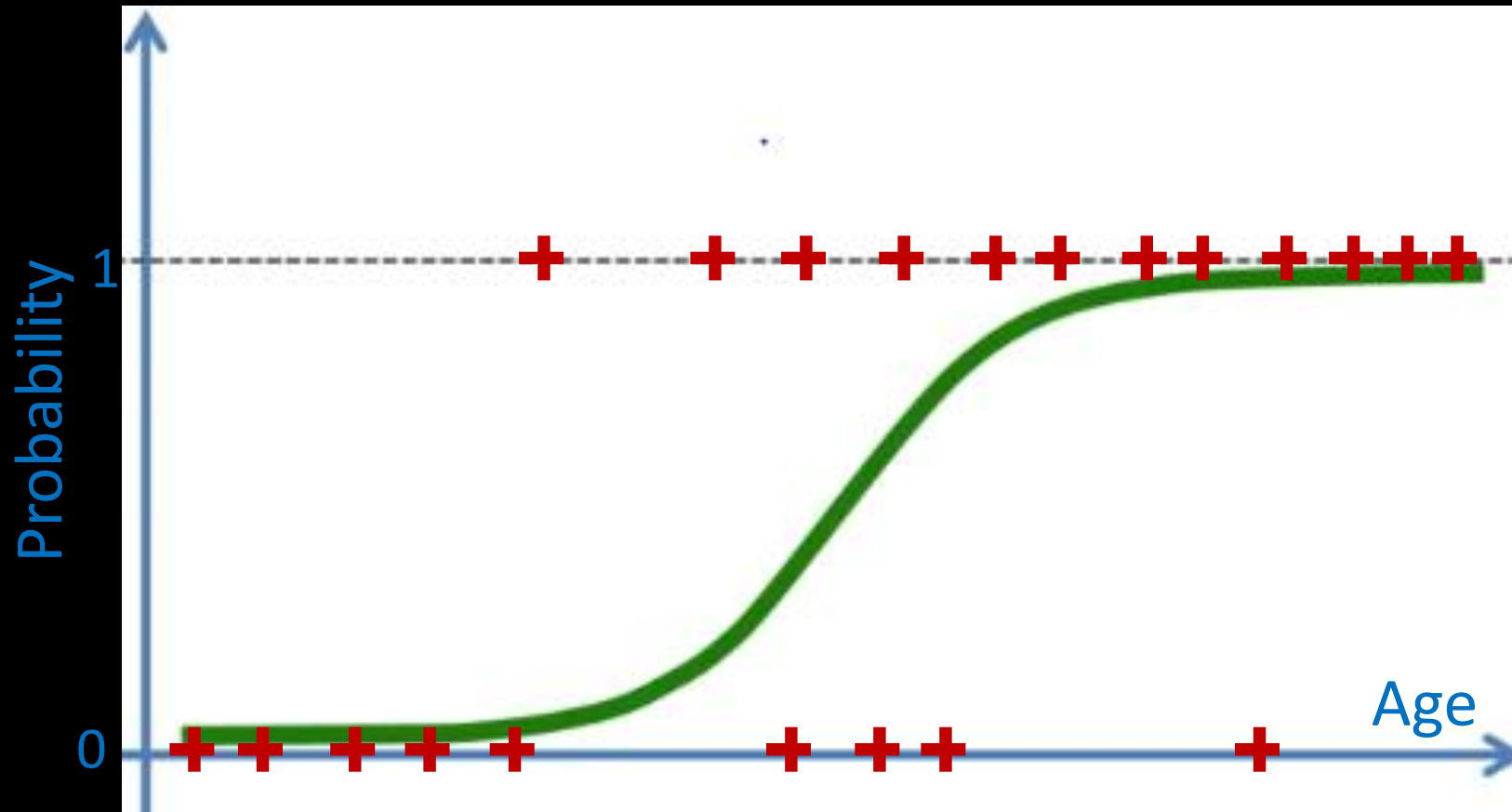




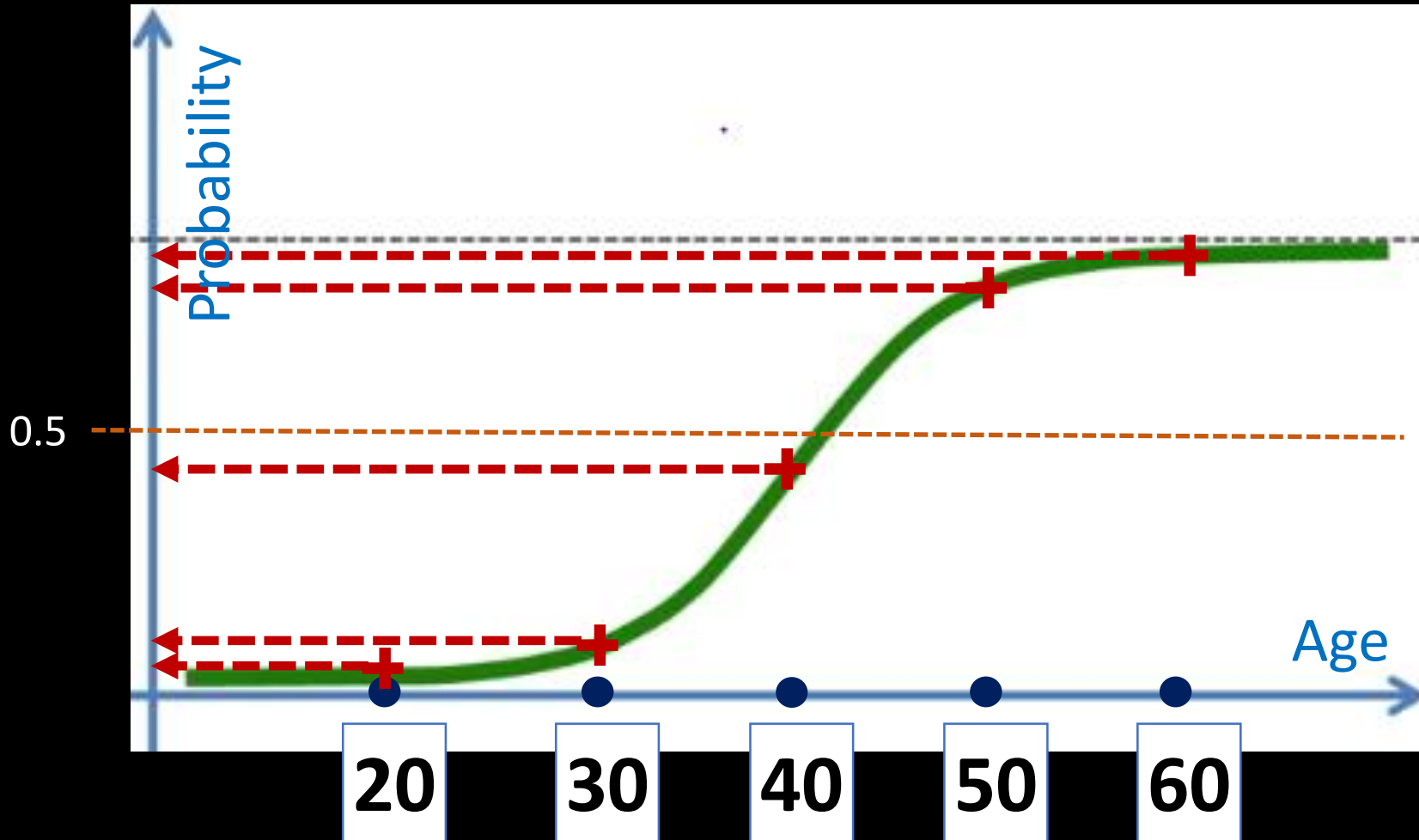
# Logistic Regression

$$\ln\left(\frac{p}{1-p}\right) = b_0 + b_1 * x$$

Age = Independent Variable (X)  
Probability = Dependent Variable (y)



# Logistic Regression



positive class = 1

negative class = 0

If probability  $\geq 0.5$ : predict 1

If probability  $< 0.5$ : predict 0

# Linear v/s Logistic Regression

Criteria	Linear Regression	Logistic Regression
Basic Definition	Data is modelled as a straight line	Is used to model certain probability of event/class happening e.g. pass/fail, yes/no, etc.
Linear relation between independent variable and dependent variable	Required	Not Required
Independent variables	Can be correlated to with each other (mostly in multiple regression)	Should not be correlated with each other
Predictions	Values can be any number	Values between 0 and 1

# Model Metrics

# Model Performance Metrics

- **False Positive**

- Model predicted a positive outcome, but it was negative.  
We predicted an event that did not occur
- It is like a warning sign e.g. the prediction was that an earthquake will occur, but it did not occur

- **False Negative**

- Model predicted that there won't be an event, but the event occurred
- Also have True Positive and True Negative

# Confusion Matrix

		Predictions	
		0	1
Actual	0	True Negative	False Positive
	1	False Negative	True Positive

# Confusion Matrix - Accuracy

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

# Classification Report

	precision	recall	f1-score	support
0	0.85	0.94	0.89	64
1	0.86	0.69	0.77	36

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

When it predicts a positive result, how often is it correct?  
Goal is to limit number of false positive (FP)

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

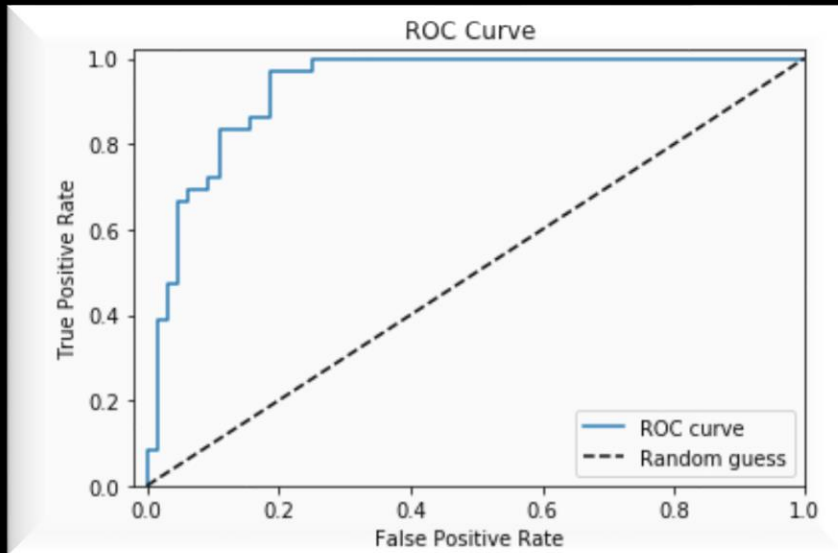
When it actually predicts the positive result, how often does it predict correctly?  
Goal is to limit number of false negatives (FN)

$$f1score = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

It is harmonic mean of precision and recall. Useful when we need to take both precision and recall into account



# ROC Curve



ROC Curve - Receiver Operating Characteristic (ROC) curve - it is a probability curve

- A visual way to measure performance of **binary classifiers**
- When it is actually the negative result, how often does the model predict incorrectly?
- Want the curve to be as far away from the random guess line
- Area under the curve (AUC) represents the degree of separation
- Higher the AUC, the model is better at predicting 0 as 0s and 1 as 1s

# Logistic Regression Demo



## Use Case: Credit Card Application

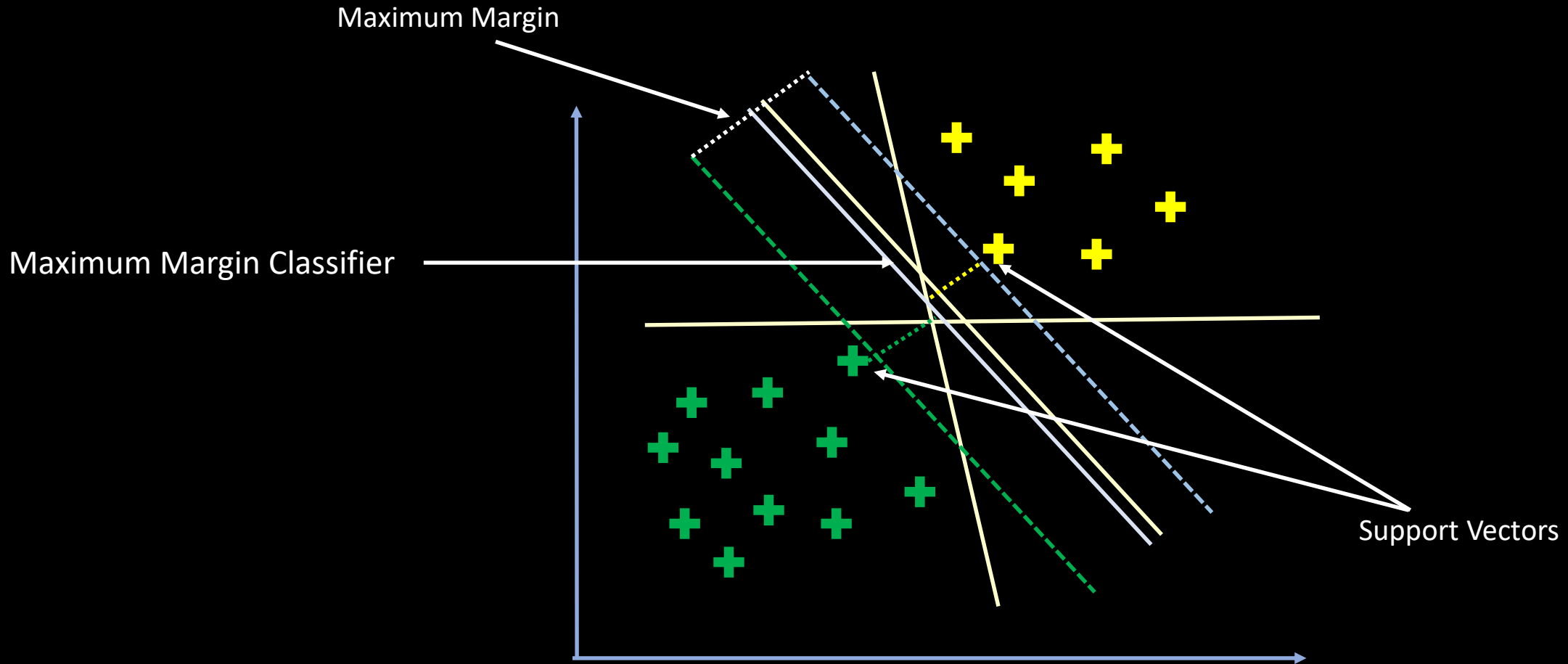
- Based on 15 features (not named), predict whether a new customer's credit card application should be approved
- Predicting a “class” – 1 or 0
- Dataset - Credit\_Card\_Application\_s.csv

# Classification: Support Vector Machine

# Support Vector Machine

- Support Vector Machine (SVM) is a supervised machine learning algorithm
- Supports both Classification and Regression, however it is primarily used for Classification

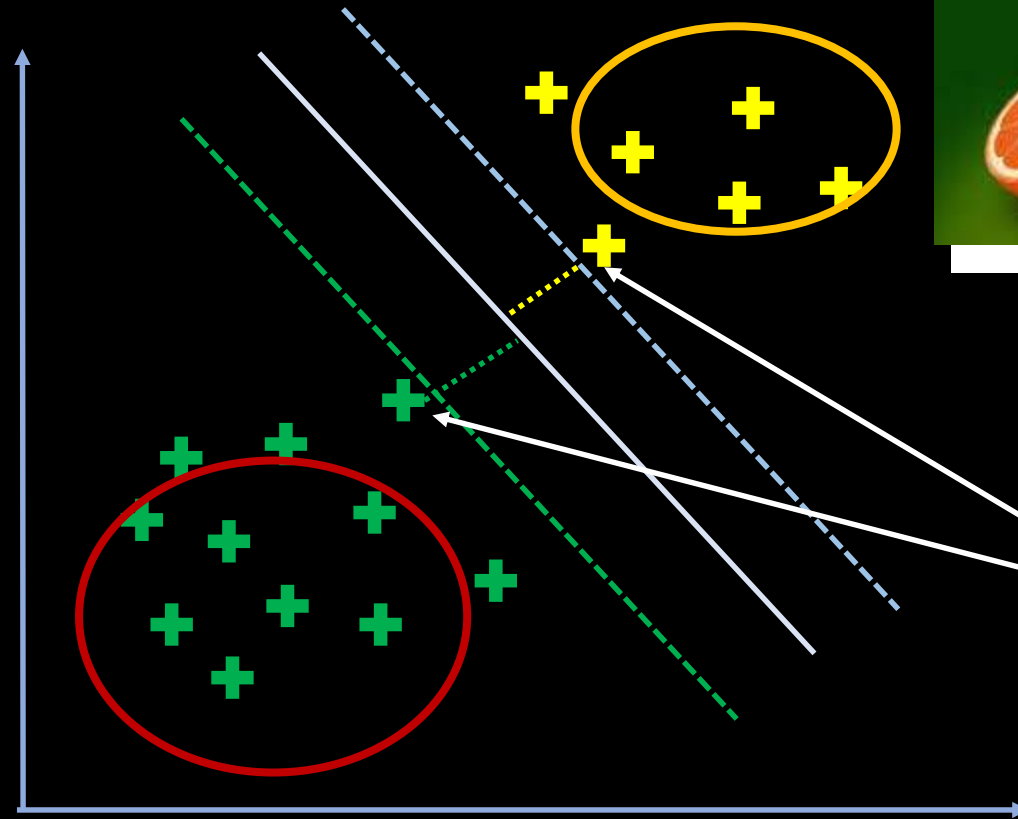
# SVM – How does it work?



# Support Vector Machine (SVM)



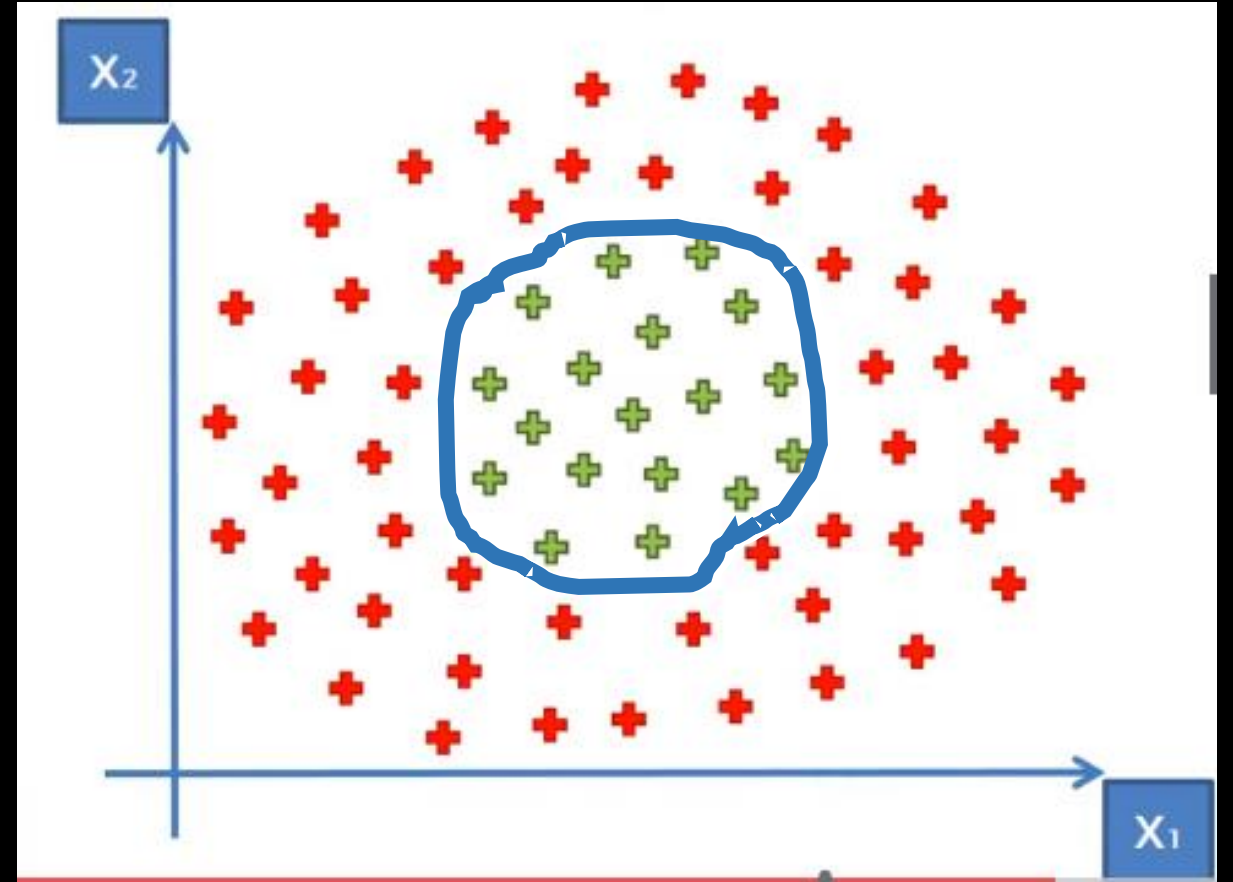
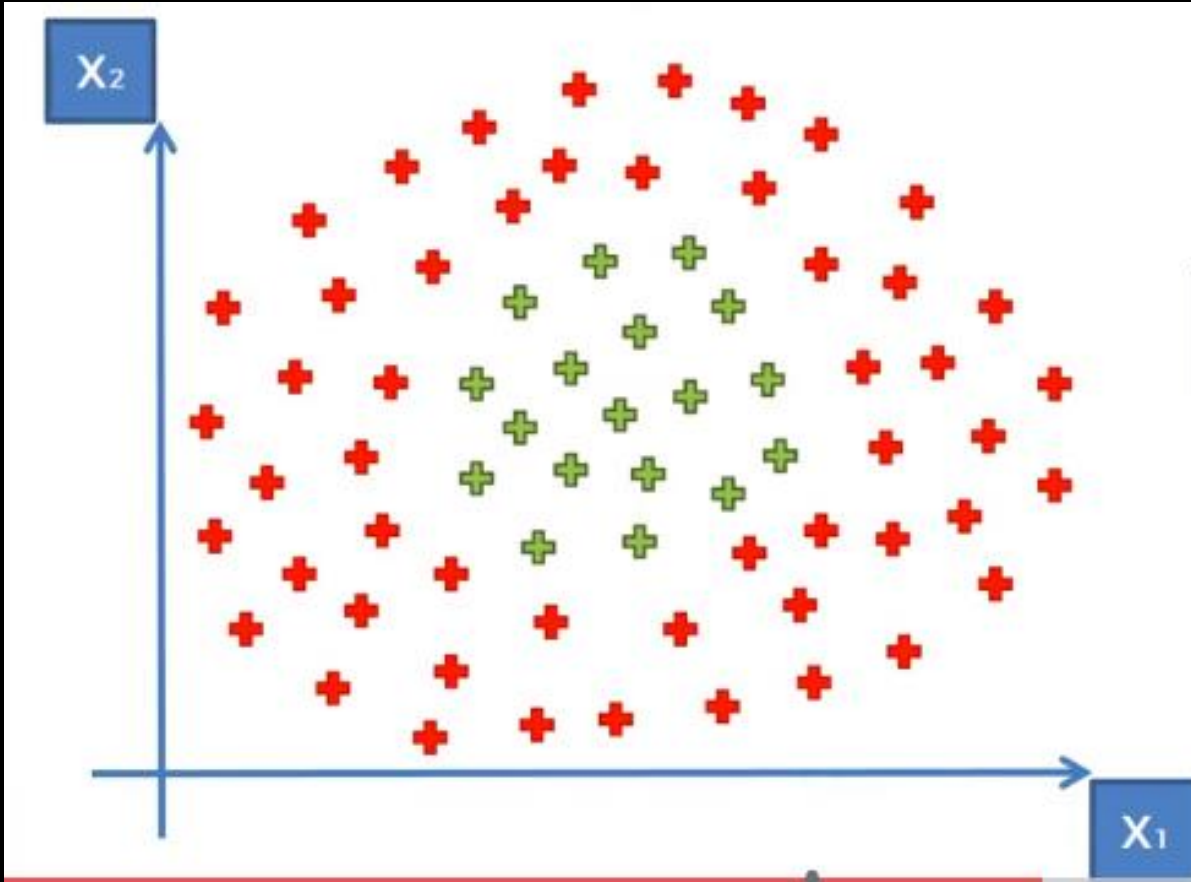
# Support Vector Machines (SVM)



Support Vectors



# Kernel – SVM: Non-Linear Data



# SVM-Classification Demo

# Real Life Use case : Customer Churn



# Use case: Financial Customer Churn Data

- Dataset – Customer information with a financial institution
- If doing classification with SVM then it can be applied to most of the datasets where logistic regression is used
- Dataset has following features:

**RowNumber:** Dataset row number

**CustomerId:** Customer Id

**Surname:** Last name of the person

**CreditScore:** Credit Score of the person

**Geography:** Country of residence

**Gender:** Person's Gender

**AGE:** Age of the person

**Tenure:** How long has the person owned the card

**Balance:** Outstanding balance

**NumOfProducts:** Number of products owned by the person with company

**HasCrCard:** Person has credit card

**IsActiveMember:** Is the person active member of the company

**EstimatedSalary:** Estimated salary of the person

**Exited:** Did the person stay or leave

- Dataset - Churn\_Modelling.csv

# Project Code Walkthrough

# Spark Income Prediction

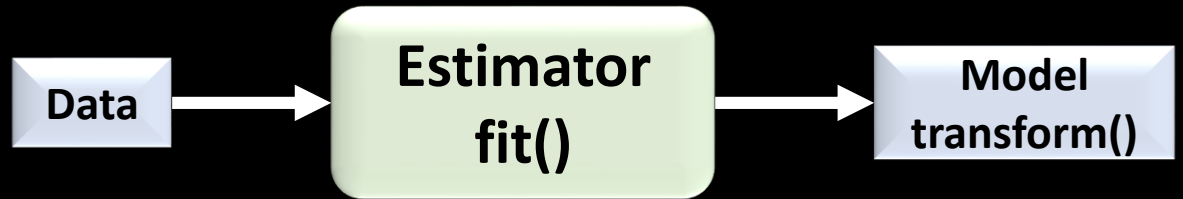
- Objective is to predict if a person's income will be  $\geq 50K$  or  $< 50K$  using number of features
- Code examples for Transformer, Estimator, Evaluator and Pipeline
- Will implement different algorithms with this dataset

# Transformer



- A Transformer is an algorithm that is used to transform one data frame into another data frame
- Feature Extraction (initial processing)
- Transform data into format required by ML Algorithms
- Take input column and transform it into output column
- For example:
  - Sentences into words – Tokenizer
  - Convert categorical data into numbers e.g. Male = 1, Female = 0
  - Normalize the data

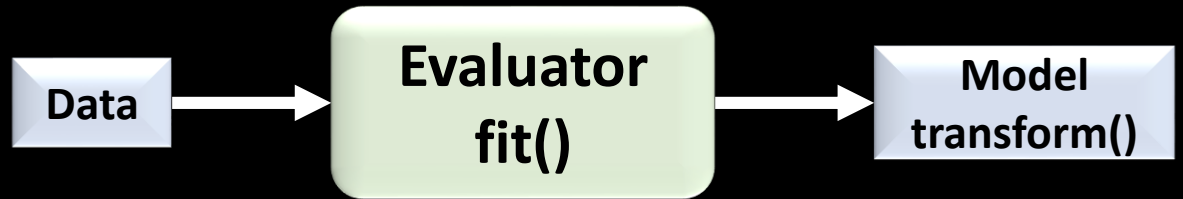
# Estimator



- It is a kind of transformer
- Algorithm that trains (does fit) on the data
- Resulting model is a type of transformer
- For example:
  - `LogisticRegression.fit()` → LogisticRegression Model
  - The `StringIndexer` is a type of estimator



# Evaluator



- Evaluate model performance
- Assists in automating model tuning process:
  - Select best model for making predictions
  - Compare model performance

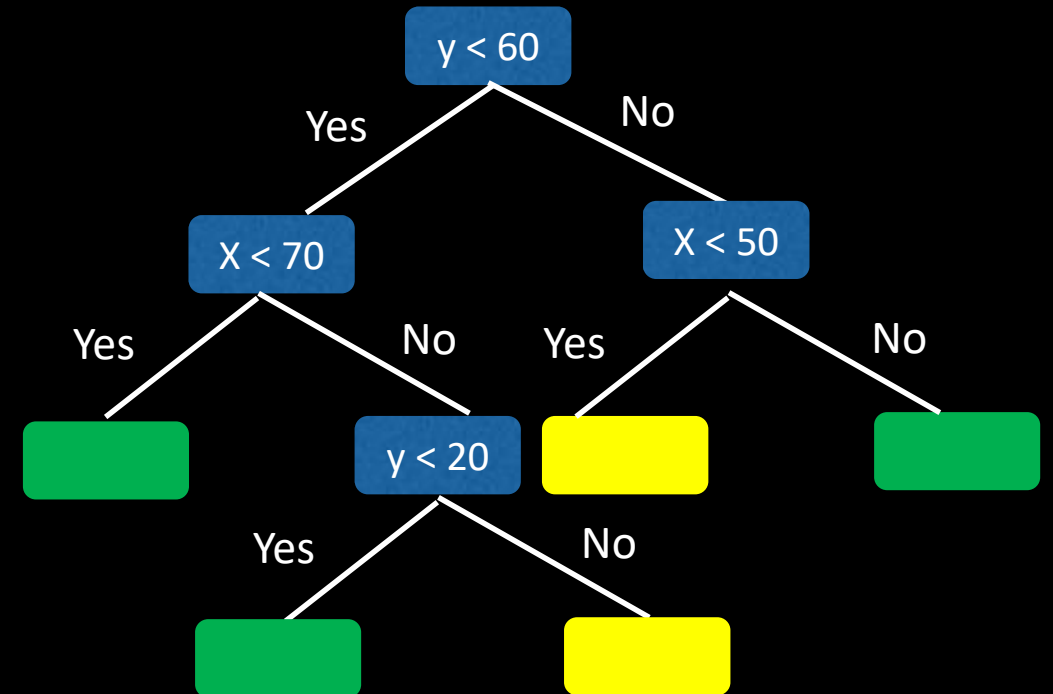
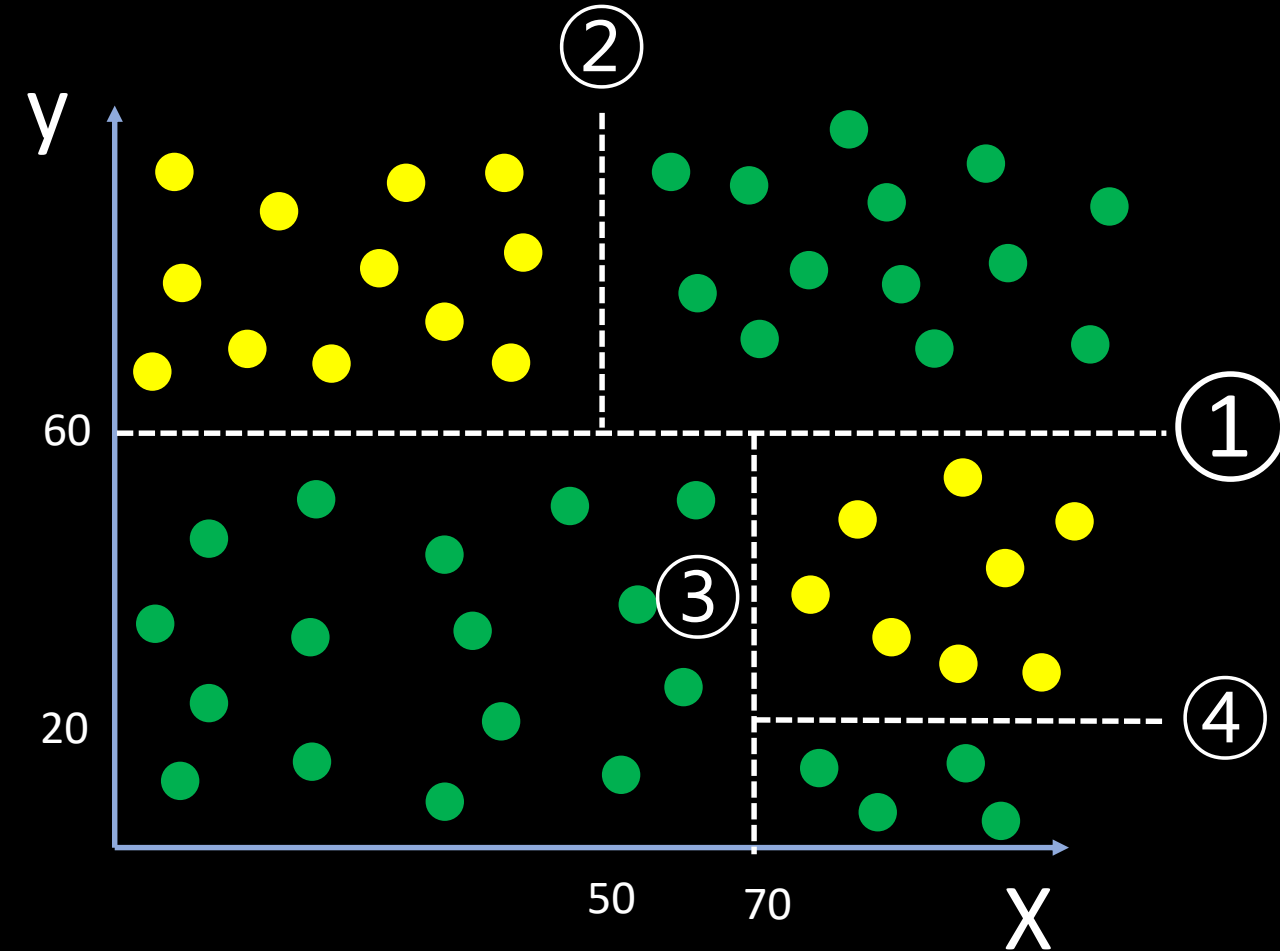
# Classification: Decision Trees

# Decision Trees

## Classification and Regression Tree (CART)

- CART model is a binary tree
- Classification Trees
  - Classify data into categorical variables e.g. buy not buy
- Regression Trees
  - Predict outcome that can be real numbers e.g. temperature, person's salary, etc.
- Examine Classification Trees

# Decision Trees



# Decision Trees - Splits

- Cost function is used to determine the splits
- Regression Trees uses mean squared error (MSE) function
- Classification Trees use Gini cost function – Gini Index
- Measure of how good the split is
- A perfect separation will have value of 0, whereas 50-50 split will be the worst (value 0.5)

# Decision Tree Regressor (Example)

# Decision Tree Classification - Project

# Classification: Random Forest



# Random Forest

- Random Forest is an ensemble classifier that uses multiple decision trees
- Ensemble Models
  - Results from multiple models are combined
  - This result is better than results from individual models
  - Each individual tree predicts a best class, final result is based on taking majority votes for a class

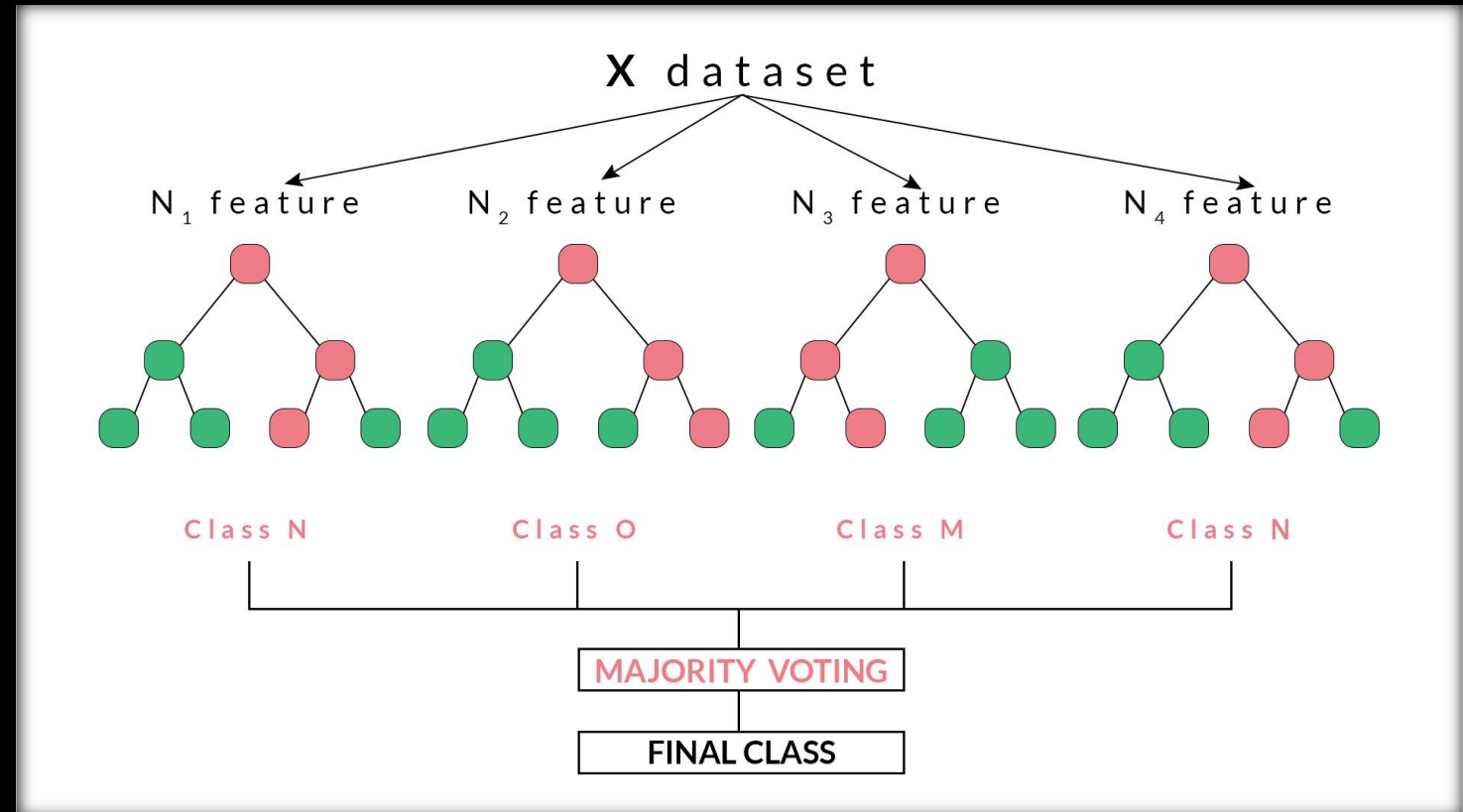
# Random Forest

Pick  $n$  data points from the training set

Build a decision tree using these  $n$  points

Choose  $N$  number of trees & repeat above 2 steps

Take new data point and apply it to all  $N$  trees and get prediction. Class with majority WINS



# Random Forest



- Microsoft used it in Kinect
- <https://www.microsoft.com/en-us/research/wpcontent/uploads/2016/02/BodyPartRecognition.pdf>

# Random Forest Regressor - Demo

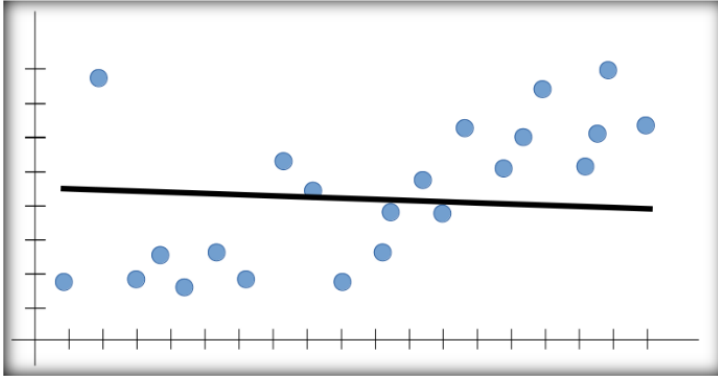
# Random Forest Classifier - Project



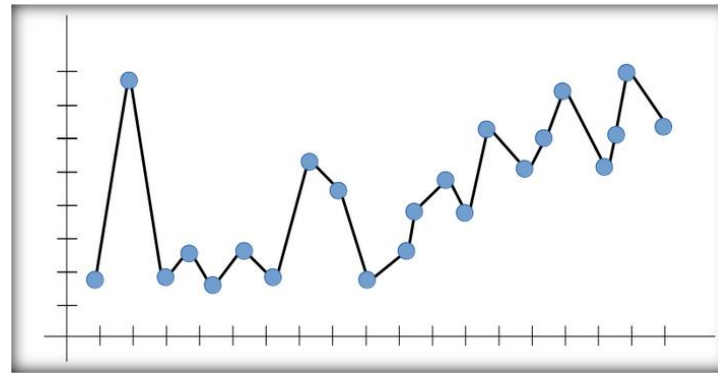
## Use Case: Heart Disease

- Based on 13 features predict if a patient will have heart disease
- Features include age, gender, chest pain, resting blood pressure, cholesterol, fasting blood sugar, resting ECG, max heart rate, exercise induced angina, etc.
- Predicting a “target” – 1 or 0
- Dataset - heart.csv

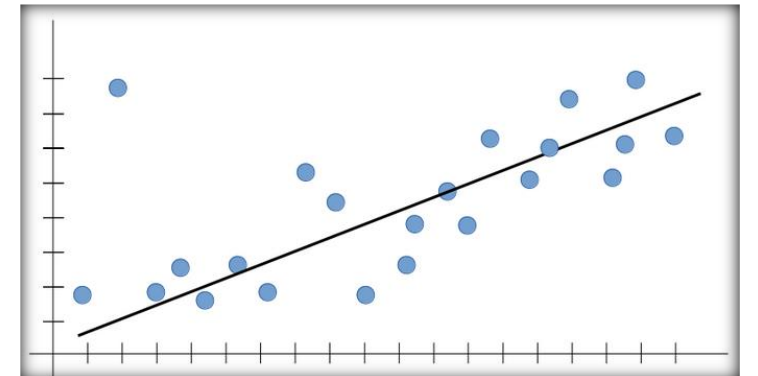
# Model Performance Metrics



Under fitting (High Bias)



Over fitting (High Variance)



Bias-Variance Trade off

# Under and Over Fitting



# Evaluation Model

- Evaluation is an estimate of how well our model is performing
- However, it is not a guarantee of performance
- Want to see techniques to create useful estimates of performance
- Have already seen train-test split
- *K*-fold Cross validation is another one

# K-Fold Cross Validation

Training Set										
1	2	3	4	5	6	7	8	9	10	.69
1	2	3	4	5	6	7	8	9	10	.64
1	2	3	4	5	6	7	8	9	10	.73
1	2	3	4	5	6	7	8	9	10	.82
1	2	3	4	5	6	7	8	9	10	.64
1	2	3	4	5	6	7	8	9	10	.70
1	2	3	4	5	6	7	8	9	10	.68
1	2	3	4	5	6	7	8	9	10	.71
1	2	3	4	5	6	7	8	9	10	.70
1	2	3	4	5	6	7	8	9	10	.69

# K-Fold Cross Validation

Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Fold 6	Fold 7	Fold 8	Fold 9	Fold 10	Mean
.69	.64	.73	.82	.64	.70	.68	.71	.70	.69	.70

Best Result

Worst Result

# Improving Performance

---

- Select multiple algorithms on the dataset and measure performance of each
- **Parameters Tuning**
  - **Model Parameter** – these are learnt by the model and are internal
  - **Hyper Parameter** – these are like “knobs” that we use to tune the algorithm
    - Selecting number of clusters, tree depth, how many iterations, number of trees in random forest, etc.



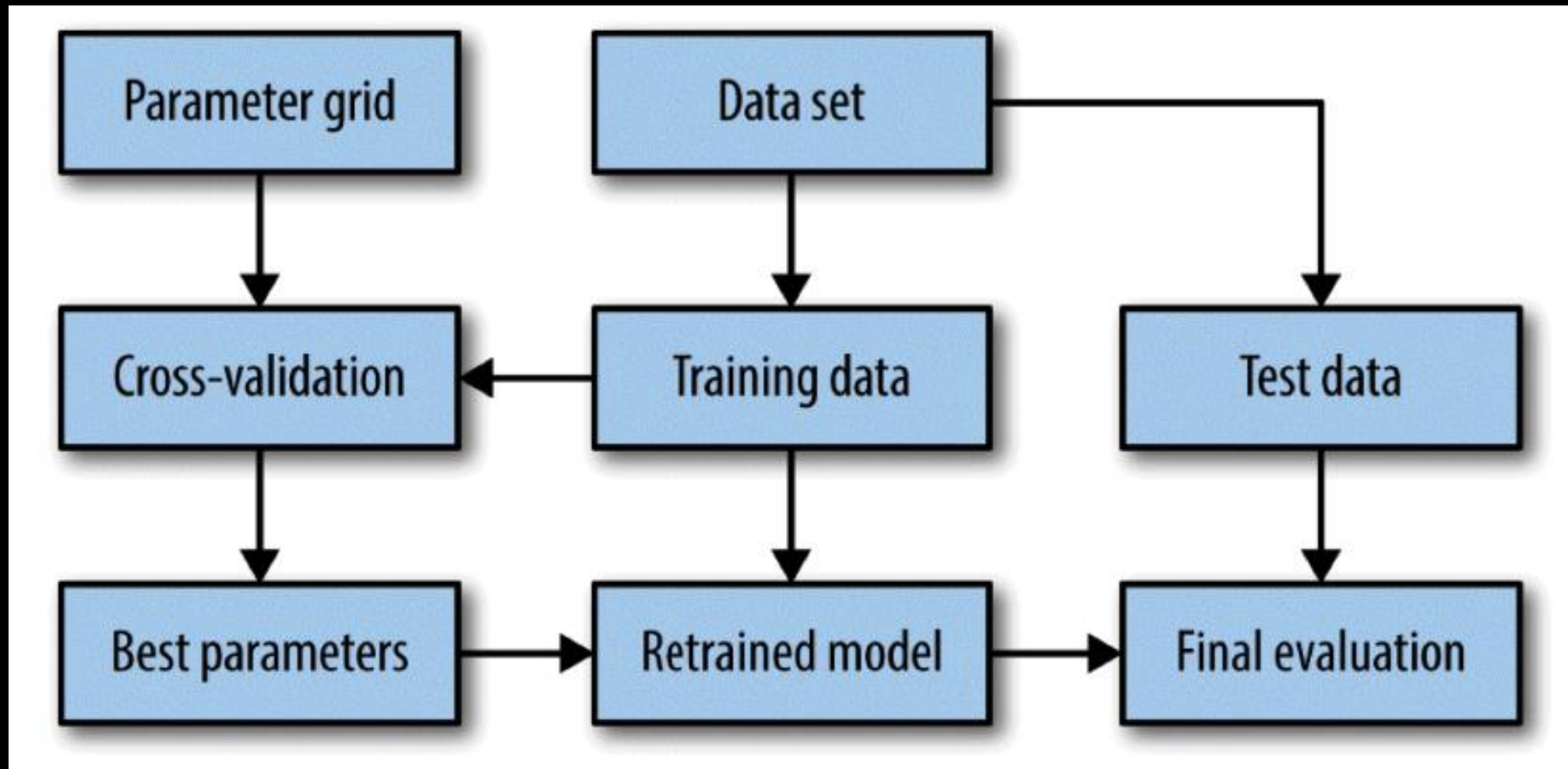
# Model Plot



# Grid Search – Hyper Parameter tuning

- Grid search will build and evaluate a model for each combination of algorithm parameters specified
- Example in the “Spark Project – Binary Classification”
- Implementation of “ParamGridBuilder” and “CrossValidator”

# Model Evaluation – Summary Takeaway



# Extra Credits





# Extra Credits

- There are 3 datasets:
  1. Churn Modelling
  2. Credit Card Application
  3. Heart Disease
- Your Work
  - Pick anyone of these datasets
  - Apply the techniques that you have learnt related to Feature Cleaning, Feature selection, etc.
  - Choose 1 or more algorithms and make predictions
  - Tomorrow afternoon we will review it with the class
- Description of datasets in following slides

# Use case: Financial Customer Churn Data

- Dataset – Customer information with a financial institution
- If doing classification with SVM then it can be applied to most of the datasets where logistic regression is used
- Dataset has following features:

**RowNumber:** Dataset row number

**CustomerId:** Customer Id

**Surname:** Last name of the person

**CreditScore:** Credit Score of the person

**Geography:** Country of residence

**Gender:** Person's Gender

**AGE:** Age of the person

**Tenure:** How long has the person owned the card

**Balance:** Outstanding balance

**NumOfProducts:** Number of products owned by the person with company

**HasCrCard:** Person has credit card

**IsActiveMember:** Is the person active member of the company

**EstimatedSalary:** Estimated salary of the person

**Exited:** Did the person stay or leave

- Dataset - Churn\_Modelling.csv



## Use Case: Credit Card Application

- Based on 15 features (not named), predict whether a new customer's credit card application should be approved
- Predicting a “class” – 1 or 0
- Dataset - Credit\_Card\_Application\_s.csv



## Use Case: Heart Disease

- Based on 13 features predict if a patient will have heart disease
- Features include age, gender, chest pain, resting blood pressure, cholesterol, fasting blood sugar, resting ECG, max heart rate, exercise induced angina, etc.
- Predicting a “target” – 1 or 0
- Dataset - heart.csv