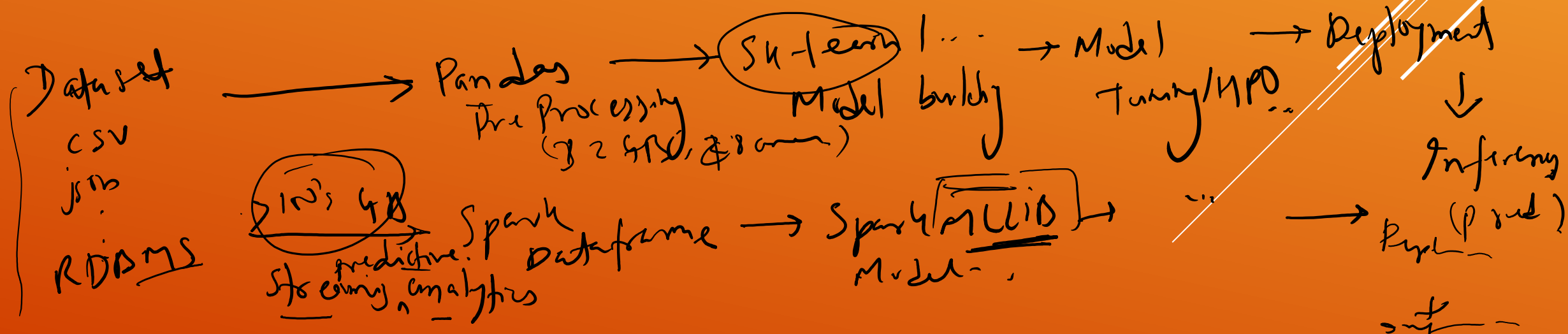


Spark MLlib machine learning

- ~~Connect~~ to batch/ real time streaming sources
- Data to be cleansed and transformed into a stream, stored in memory
- Models can or cannot be built in real time depending on the requirement and availability of libraries
- ~~Do~~ predictions in batch/real time
- ~~Receive~~, store and communicate with external system



Spark ML Lib

- **Packages** – spark.mllib, spark.ml
- **Data types** – Local vector, Labelled point
- **Local vector** – vector of double values – Dense and Sparse
- e.g, a vector (1.0, 0.0, 3.0) can be represented in dense format as [1.0, 0.0, 3.0] or in sparse format as (3, [0, 2], [1.0, 3.0])
- **Labeled point** – contains the target variable(outcome) and list of features (predictors) ()
- **Process** - Load data into RDD -> Transform RDD – Filter, Type conversion, centering, Scaling, etc. -> Convert to labeled point -> Split training and testing -> Create model -> Tune -> Perform predictions

Representation
of data from vectors

Chosen

Not

Chosen

Target variable

Location	Gender
0	0
1	1
2	1
3	0

(0, 1, 2, 3)

(0, 1, 1, 0)

100 rows
Chosen
Not Chosen

70:30
Train
Test

Test
20
10
20 -> 15

Inform

(1.0, 2.0, 0.0, 3.0)

1100

Sparse -> (4, [0, 1, 3], [1.0, 2.0, 3.0])

ML Pipelines

Spam / Not spam
feature engineering → TF-IDF → Logistic regression → prediction

Shortcut of completely ML pipeline in a single line of code.

- Standard APIs for machine learning algorithms to make it easier to combine multiple algorithms into a single pipeline, or workflow.
- Inspired from SciKit-learn where multiple libraries of pre-processing, feature engineering and modelling are combined to give an output in using single line of code.

- Three major parts:

1. **DataFrame** – to hold the dataset

2. **Transformer** – Apply Functions/Models to datasets (may add columns),
.transform() — function applied

3. **Estimator** – Converts a dataset into a Transformer, e.g. training the datasets using .fit() — Resource intensive fit() — action trans starts.

- Sometimes pipelines may also contain params which are the args passed to a transformer.

TF — TOS → TensorFlow on Spark
→ Horovod on Spark.

Feature engineering in Spark

admiral Term 2
Span/Not Span
Span
NS
NS

- **Extraction:** Extracting features from “raw” data.
- **TF-IDF** - Feature vectorization method widely used in text mining to reflect the importance of a term to a document in the corpus.
- **TF:** Both HashingTF and CountVectorizer can be used to generate the term frequency vectors.
- HashingTF is a Transformer which takes sets of terms and converts those sets into fixed-length feature vectors.
- **IDF** is an Estimator which is fit on a dataset and produces an IDFModel. The IDFModel takes feature vectors (generally created from HashingTF or CountVectorizer) and scales each column.
- **Word2Vec** - Estimator which takes sequences of words representing documents and trains a Word2VecModel. The model maps each word to a unique fixed-size vector. The Word2VecModel transforms each document into a vector using the average of all words in the document
- **CountVectorizer** - CountVectorizer and CountVectorizerModel aim to help convert a collection of text documents to vectors of token counts.

Term
1 Word
1 Document
= 1 sentence

1, 2
am, 2,

TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

Term frequency

Inverse document frequency

Number of times term t appears in a doc, d

$$\log \frac{1 + \overset{\text{\# of documents}}{n}}{1 + \underset{\text{Document frequency of the term } t}{df(d, t)}} + 1$$

$$= 2 \times 1 = \boxed{2}$$

$$\log \frac{1 + 2 + 1}{1 + 2}$$

$$= \log \frac{3 + 1}{3 + 1}$$

$$= \log 1 + 1$$

$$= 0 + 1$$

$$= 1$$

Selection: Selecting a subset from a larger set of features

- VectorSlicer - Transformer that takes a feature vector and outputs a new feature vector with a sub-array of the original features. It is useful for extracting features from a vector column. *→ v2 2-map*
- Rformula - selects columns specified by an R model formula. Currently supports a limited subset of the R operators, including '~', '.', ':', '+', and '-'.
- ChiSqSelector - It operates on labeled data with categorical features. ChiSqSelector uses the Chi-Squared test of independence to decide which features to choose.

Name, id, tpr, ...

> 5

(select a part
of where
'id > 5')
id > 5)

R developers
worrying

Spurious ML models

Spam

not spam

Not so

may be spam

may be not

Spam

VS.

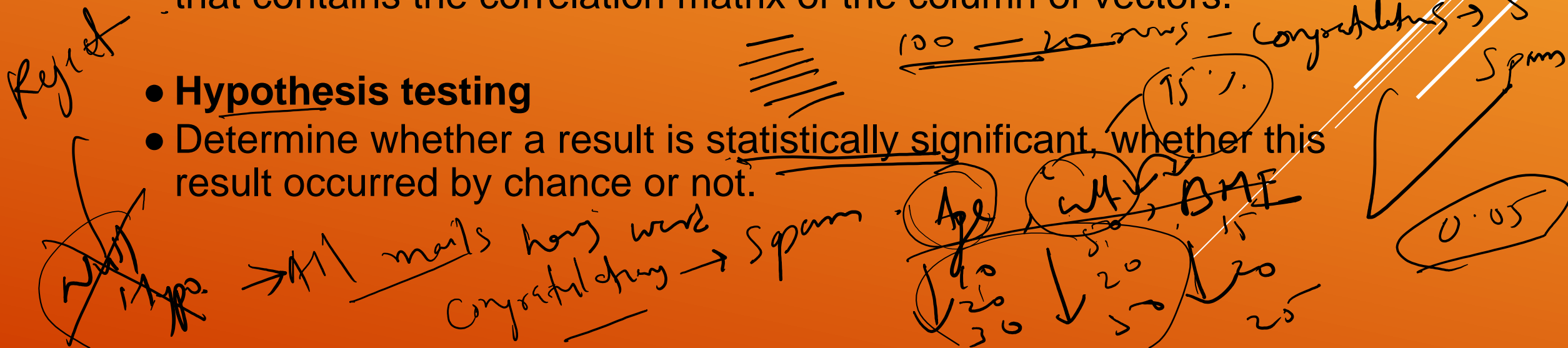
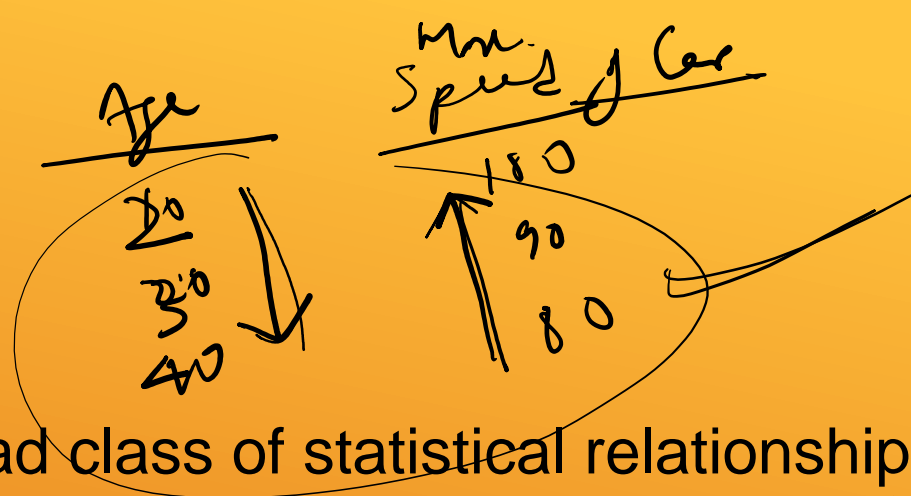
More..

• Correlation

- Correlation is any of a broad class of statistical relationships involving dependence, though in common usage it most often refers to how close two variables are to having a linear relationship with each other.
- Correlation computes the correlation matrix for the input Dataset of Vectors using the specified method. The output will be a DataFrame that contains the correlation matrix of the column of vectors.

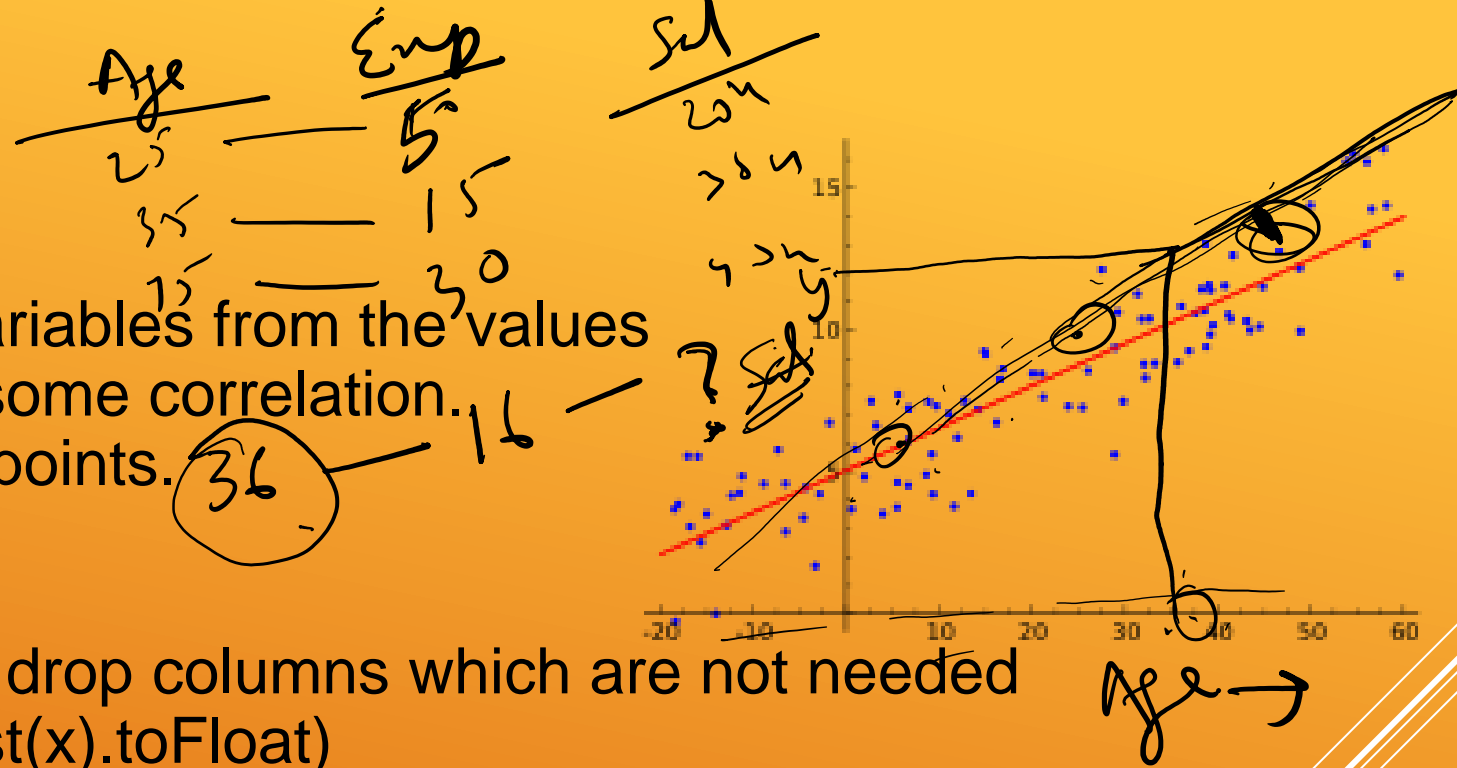
• Hypothesis testing

- Determine whether a result is statistically significant, whether this result occurred by chance or not.



Linear regression

- Estimate value of dependent variables from the values of independent variables with some correlation.
- Draw the best line to fit plotted points.



Example -

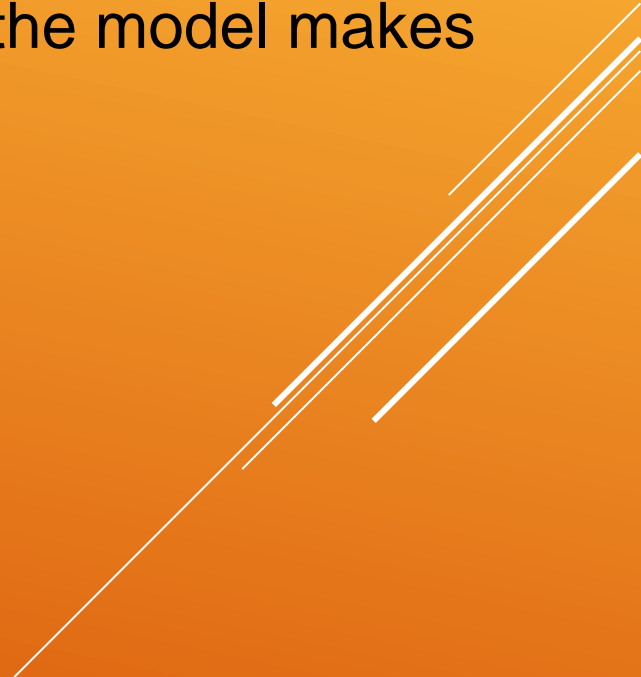
- Convert input string to vector & drop columns which are not needed
- Create rDD vectors.dense(atList(x).toFloat)
- Create labelled vectors & drop low correlation vectors (should be continuous)
- Run LR model on training data Print coefficients, intercepts, summary, lr.fit(trainingdata)
- Predict
- Evaluate using MSE (average((predicted value - actual value) ** 2, should be low), R square should be high as close to 1 as possible

Mean Squared Error <<<

K-Means Clustering

- Attempts to split data into K- groups that are closest to K centroids
- Un Supervised learning – uses only position of each data point
- Randomly pick K centroids -> Assign each data point to the centroid it's closest to
- Re-compute centroids based on average position of each centroid's points
- Iterate until points stop changing assignment to centroids
- If you want to predict the cluster for new points, just find the centroid they are closest to

Classification

- **Binary Classification** is the task of predicting a binary label. E.g., is an email spam or not spam?
 - Supervised learning – uses only position of each data point
 - For binary classification problems, the algorithm outputs a binary logistic regression model. Given a new data point, denoted by x , the model makes predictions by applying the logistic function
- 
- A series of three parallel white diagonal lines in the bottom right corner of the slide.

Model selection and hyper parameter tuning

- An important task in ML is *model selection*, or using data to find the best model or parameters for a given task.
- **Estimator**: algorithm or Pipeline to tune
- Set of **ParamMaps** : parameters to choose from, sometimes called a “parameter grid” to search over
 - Split the input data into separate training and test datasets.
 - For each (training, test) pair, they iterate through the set of ParamMaps
 - Identify the best ParamMap, CrossValidator finally re-fits the Estimator using the best ParamMap and the entire dataset.
- **Evaluator**: metric to measure how well a fitted Model does on held-out test data i.e evaluate the Model’s performance using the Evaluator.
- RegressionEvaluator , a BinaryClassificationEvaluator for binary data, or a MulticlassClassificationEvaluator



Thank you!

