



Association for
Computing Machinery

Advancing Computing as a Science & Profession

Austin SIGKDD Spark talk

Random Forest and Decision Trees in Spark MLlib

Tuhin Mahmud
Sep 20th, 2017

About Me

Tuhin Mahmud

Advisory Software Engineer , IBM

tuhinm@hotmail.com

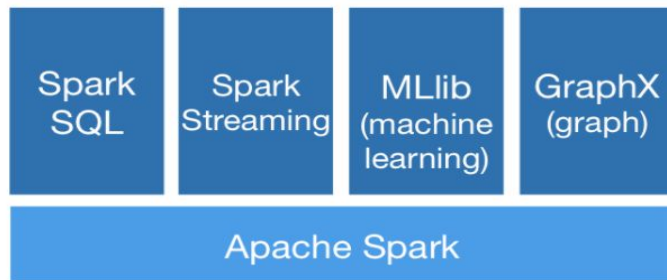
<https://www.linkedin.com/in/tuhinmahmud/>

Agenda

- Spark and MLLib Overview
- Decision Trees in Spark MLib
- Random Forest in Spark MLib
- Demo

What is MLlib?

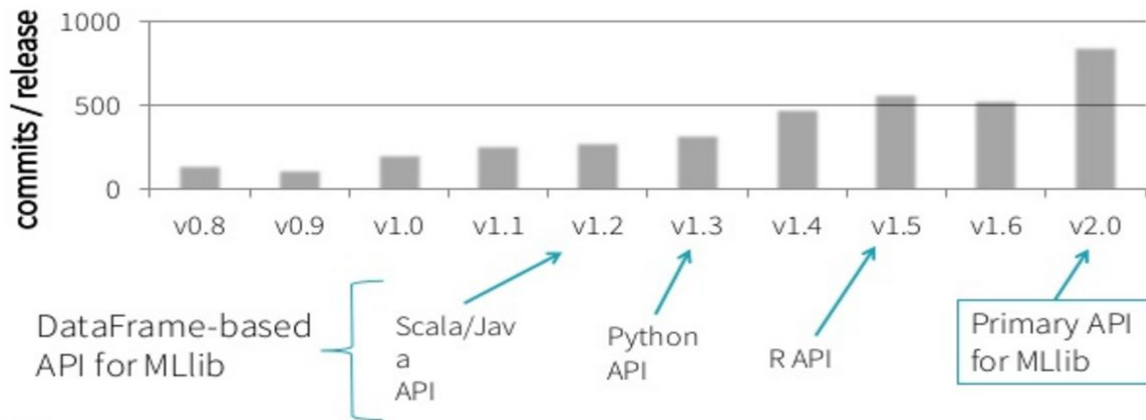
- MLlib is Spark's machine learning (ML) library.
- Its goal is to make practical machine learning **scalable** and **easy to use**.
- ML algorithms include Classification , Regression , Decision Trees and random Forests, Recommendation, Clustering, Topic Modeling (LDA), Distributed linear Algebra(SVD,PCA) and many more.



Spark MLlib Trajectory

 Apache Spark MLlib 2.0 Preview: Data Science and Production - Databricks

MLlib trajectory



 databricks

Spark MLlib

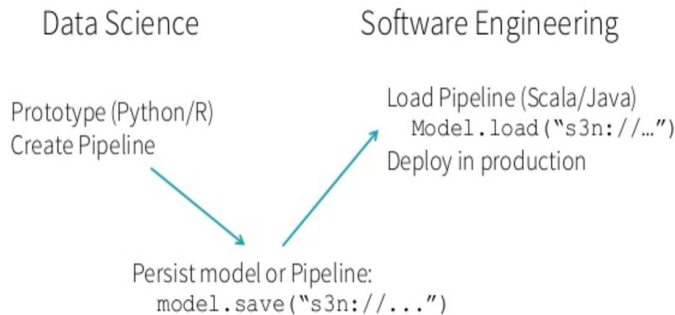
- Initially developed by MLbase team in AmpLab (UC Berkeley) - supported only Scala and Java
- Shipped to Spark v0.8 (Sep 2013)
- Current release [Spark 2.2.0 released](#) (Jul 11, 2017)
- MLlib 2.2:
 - **DataFrame-based API** becomes the primary API for MLlib
 - `org.apache.spark.ml` and `pyspark.ml`
 - Ease of use
 - Java, Scala, **Python**, **R**
 - interoperates with [NumPy](#) in Python
 - Any Hadoop data source (e.g. HDFS, HBase, or local files), making it easy to plug into **Hadoop workflows**.

Spark Mllib in production

- ML persistence
 - Saving and loading
- **Pipeline** MLib standardizes APIs for machine learning algorithms to make it easier to combine multiple algorithms into a single pipeline, or workflow. This section covers the key concepts introduced by the Pipelines API, where the pipeline concept is mostly inspired by the [scikit-learn](#) project.

Some concepts[5]

- **DataFrame**
- **Transformer** - one dataframe to another
- **Estimator** - fits on a dataframe to produce a transformer
- **Pipeline** A pipeline chains multiple transformers and Estimators to specify
- **Parameter**



MLlib 2.2.0 API

1. MLlib: RDD-based API (maintenance mode)

2. MLlib: DataFrame-based API

Why is MLlib switching to the DataFrame-based API?

- **DataFrames** provide a more **user-friendly API** than RDDs.
 - a. The many benefits of DataFrames **include Spark Datasources**,
 - b. **SQL/DataFrame queries**, Tungsten and Catalyst optimizations,
 - c. uniform APIs across languages.
- The DataFrame-based API for MLlib provides a uniform API across ML algorithms and across multiple languages.
- DataFrames facilitate practical **ML Pipelines**, particularly feature transformations. See the [Pipelines guide](#) for details.

What is “Spark ML”?

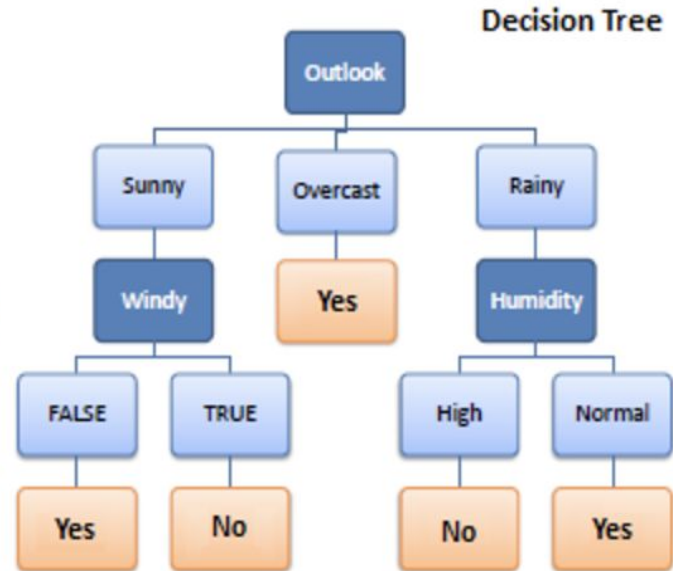
- “Spark ML” is **not an official name** but occasionally used to refer to the MLlib DataFrame-based API.

Agenda

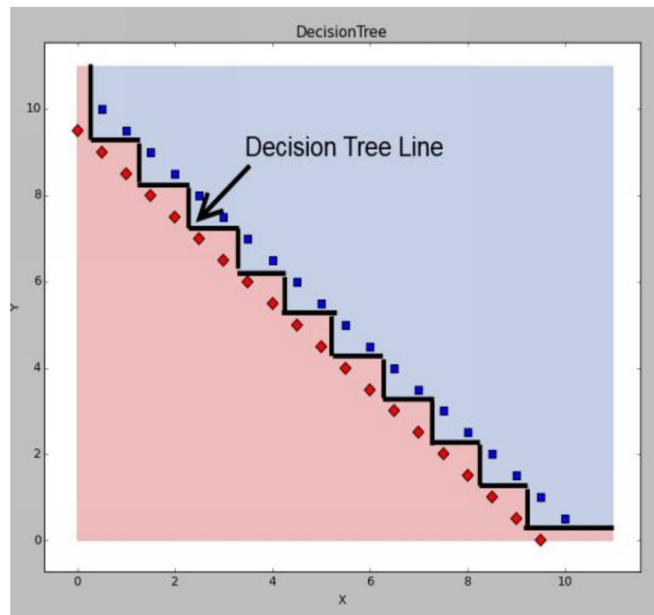
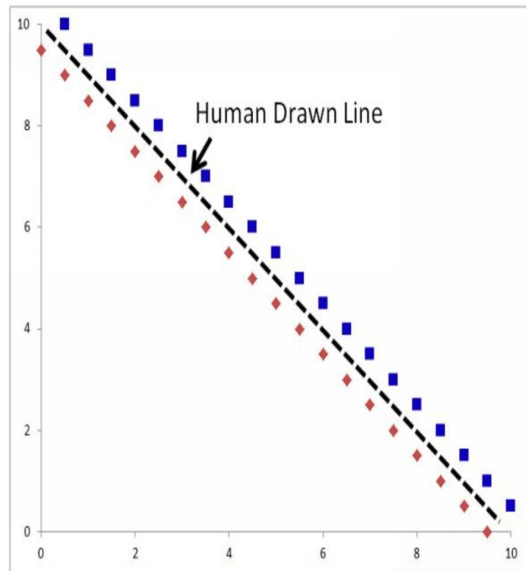
- Spark and MLlib Overview
- Decision Tree in Spark MLlib
- Random Forest in Spark MLlib
- Demo

Decision Tree

Predictors				Target
Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	False	No
Rainy	Hot	High	True	No
Overcast	Hot	High	False	Yes
Sunny	Mild	High	False	Yes
Sunny	Cool	Normal	False	Yes
Sunny	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Rainy	Mild	High	False	No
Rainy	Cool	Normal	False	Yes
Sunny	Mild	Normal	False	Yes
Rainy	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Sunny	Mild	High	True	No



Decision Tree Applied



DecisionTrees packages in MLlib

1. `from pyspark.mllib.tree import DecisionTree, DecisionTreeModel` (RDD based)
2. `from pyspark.ml.classification import DecisionTreeClassifier` (Dataframe based)

DecisionTree Classifier (MLlib Dataframe based API)

```
>>> from pyspark.ml.linalg import Vectors
>>> from pyspark.ml.feature import StringIndexer
>>> df = spark.createDataFrame([
...     (1.0, Vectors.dense(1.0)),
...     (0.0, Vectors.sparse(1, [], []))], ["label", "features"])
>>> stringIndexer = StringIndexer(inputCol="label", outputCol="indexed")
>>> si_model = stringIndexer.fit(df)
>>> td = si_model.transform(df)
>>> dt = DecisionTreeClassifier(maxDepth=2, labelCol="indexed")
>>> model = dt.fit(td)
>>> model.numNodes
```

Decision Tree

hyper-parameters for decision tree in MLlib

- **numClasses:** How many classes are we trying to classify?
- **categoricalFeaturesInfo:** A specification whereby we declare what features are **categorical features** and should not be treated as numbers
- **impurity:** A measure of the homogeneity of the labels at the node. Currently in Spark, there are two measures of impurity with respect to classification: **Gini** and **Entropy**
- **maxDepth:** A **stopping criterion** which limits the depth of constructed trees. Generally, deeper trees lead to more accurate results but run the risk of overfitting.
- **maxBins:** Generally, increasing the number of bins **allows the tree to consider more values** but also increases computation time.

$$Entropy = \sum_j -p_j \log_2 p_j$$

Where p_j is the frequency of label j at a node.

The **Gini Index**: is a measure of how often a randomly chosen element would be misclassified if it were randomly given a label according to the distribution of labels at a given node.

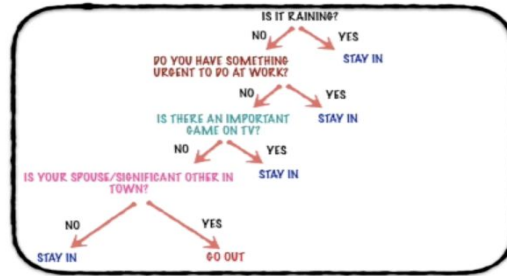
$$Gini Index = 1 - \sum_j p_j^2$$

Decision Trees are prone to Overfitting

A DECISION TREE CAN BE USED TO SOLVE MACHINE LEARNING PROBLEMS

A DECISION TREE PREDICTS THE OUTCOME GIVEN THE VALUES OF INPUT VARIABLES

INPUT VARIABLES/
PREDICTORS



OUTCOME/OUTPUT
VARIABLES

DECISION TREES ARE VERY PRONE
TO THE RISK OF OVERFITTING

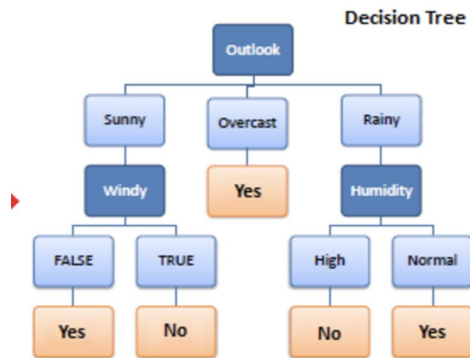
ENSEMBLE LEARNING CAN MITIGATE
THE RISK OF OVERFITTING

Spark MLlib -Decision Tree Example Notebook

1. Notebook using small dataset golf play

http://nbviewer.jupyter.org/github/tuhinmahmud/sigkdd_austin/blob/master/SparkMLlibPyspark.golf.ipynb

	Play	Outlook	NumericalTemp	NumericalHumidity	Windy
0	0.0	sunny	85	85	False
1	0.0	sunny	80	90	True
2	1.0	overcast	83	86	False
3	1.0	rainy	70	96	False
4	1.0	rainy	68	80	False
5	0.0	rainy	65	70	True
6	1.0	overcast	64	65	True
7	0.0	sunny	72	95	False
8	1.0	sunny	69	70	False
9	1.0	sunny	75	80	False



Agenda

- Spark and MLlib Overview
- Decision Tree in Spark MLlib
- Random Forest in Spark MLlib
- Demo

Random Forest as Ensemble

A MACHINE LEARNING ENSEMBLE IS A
COLLECTION OF MODELS

THE MODELS IN THE ENSEMBLE CAN BE
BASED ON DIFFERENT TECHNIQUES
TRAINED ON DIFFERENT TRAINING SETS
USING DIFFERENT FEATURES
USING DIFFERENT VALUES OF PARAMETERS

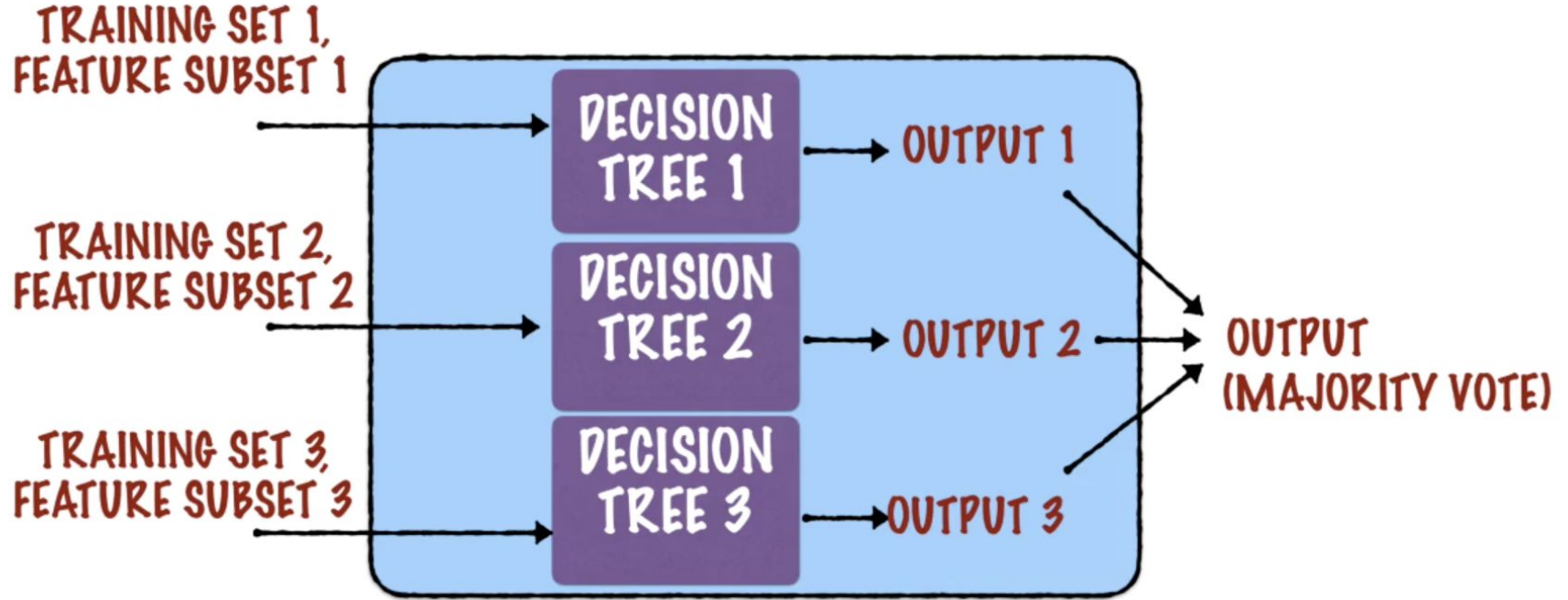
A RANDOM FOREST IS AN
ENSEMBLE OF DECISION TREES

EACH DECISION TREE IN THE ENSEMBLE IS

✓ (BAGGING)
TRAINED ON DIFFERENT TRAINING SETS

USING DIFFERENT FEATURES
(A RANDOMLY SELECTED SUBSET OF FEATURES)

RANDOM FOREST



Random Forest

The name “Random Forest” comes from combining the **randomness** that is used to pick the **subset of data** with having a bunch of Decision trees.

A random forest (RF) is a collection of **tree predictors**

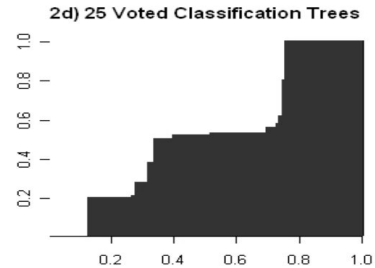
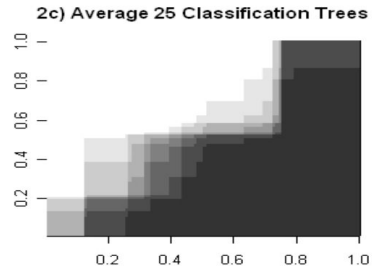
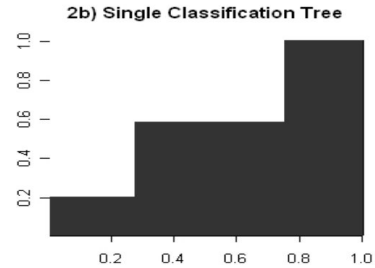
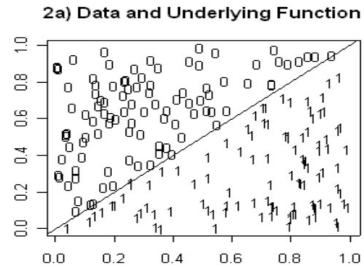
$$f(x, T, \Theta_k), k = 1, \dots, K$$

where the Θ_k are i.i.d random vectors.

The trees are combined by

- voting (for classification)
- averaging (for regression).

Random Forest

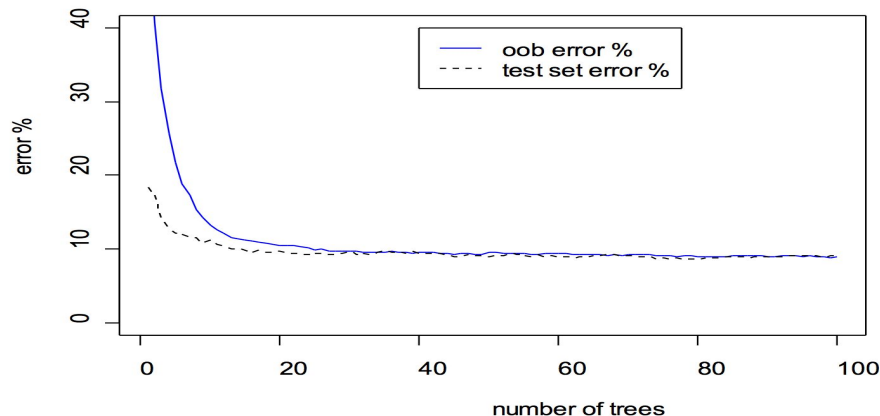


Random Forest - Out of bag Error

Out-of-bag Error Estimate

- Average over the cases within each class to get a classwise out-of-bag error rate.
- Average over all cases to get an overall out-of-bag error rate.

In random forests, there is no need for cross-validation to get an unbiased estimate of the test set error. It is estimated internally,



RandomForest in MLlib

DataFrame Based API

```
class pyspark.ml.classification.RandomForestClassifier(self, featuresCol="features", labelCol="label", predictionCol="prediction",  
probabilityCol="probability", rawPredictionCol="rawPrediction", maxDepth=5, maxBins=32, minInstancesPerNode=1, minInfoGain=0.0,  
maxMemoryInMB=256, cacheNodeIds=False, checkpointInterval=10, impurity="gini", numTrees=20, featureSubsetStrategy="auto", seed=None,  
subsamplingRate=1.0)[source]¶
```

```
>>> import numpy>>> from numpy import allclose  
>>> from pyspark.ml.linalg import Vectors  
>>> from pyspark.ml.feature import StringIndexer  
>>> df = spark.createDataFrame([  
...   (1.0, Vectors.dense(1.0)),  
...   (0.0, Vectors.sparse(1, [], [])), ["label", "features"])  
>>> stringIndexer = StringIndexer(inputCol="label", outputCol="indexed")  
>>> si_model = stringIndexer.fit(df)  
>>> td = si_model.transform(df)  
>>> rf = RandomForestClassifier(numTrees=3, maxDepth=2, labelCol="indexed", seed=42)  
>>> model = rf.fit(td)
```

<https://spark.apache.org/docs/2.2.0/ml-classification-regression.html#random-forest-classifier>

Random Forest - parameters

RandomForest - some parameters (older RDD based but some apply to dataframe based)

- **numTrees**: Number of trees in the resulting forest. **Increasing the number of trees decreases model variance.**
- **featureSubsetStrategy**: Specifies a method which produces a number of how many features are selected for training a single tree.
- **seed**: Seed for random generator initialization, since RandomForest depends on random selection of features and rows

The parameters **numTrees** and **maxDepth** are often referenced as **stopping criteria**.

Random Forest -parameters

Spark also provides additional parameters to **stop tree growing** and produce fine-grained trees:

- **minInstancesPerNode:** A node is not split anymore, if it would provide left or right nodes which would contain smaller number of observations than the value specified by this parameter. Default value is 1, but typically for regression problems or large trees, the value should be higher.
- **minInfoGain:** Minimum information gain a split must get. Default value is 0.0.

MLlib - Labeled point vector (RDD based)

Labeled point vector

- Prior to running any supervised machine learning algorithm using Spark MLlib, we must **convert our dataset into a labeled point vector**.
 - ```
val higgs = response.zip(features).map {
 case (response, features) =>
 LabeledPoint(response, features) }
higgs.setName("higgs").cache()
```
- An example of a labeled point vector follows:  

```
(1.0, [0.123, 0.456, 0.567, 0.678, ..., 0.789])
```

# MLlib - data caching

## Data caching

Many machine learning algorithms are iterative in nature and thus require multiple passes over the data.

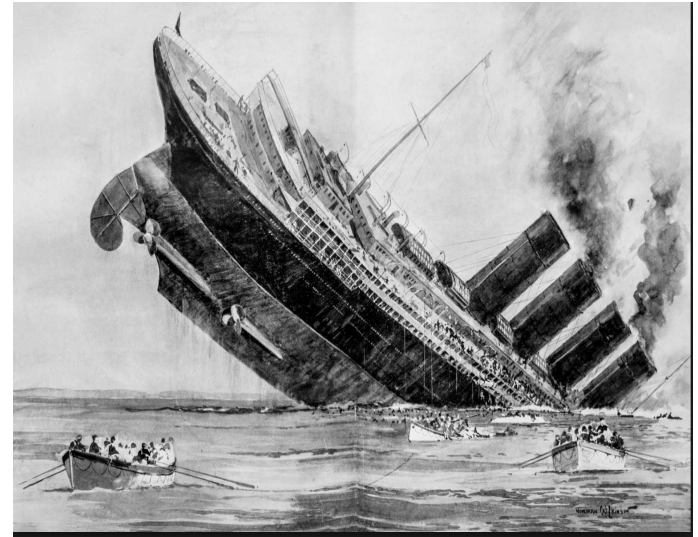
Spark provides a way to persist the data in case we need to iterate over it. Spark also publishes several `StorageLevels` to allow storing data with various options:

- `NONE`: No caching at all
- `MEMORY_ONLY`: Caches RDD data only in memory
- `DISK_ONLY`: Write cached RDD data to a disk and releases from memory
- `MEMORY_AND_DISK`: Caches RDD in memory, if it's not possible to offload data to a disk
- `OFF_HEAP`: Use external memory storage which is not part of JVM heap

## Making sense of a tragedy - titanic dataset

- Notebook using titanic dataset and MLlib spark dataframe based apis for Decision Tree and Random Forest

[http://nbviewer.jupyter.org/github/tuhinmahmud/sigkdd\\_austin/blob/master/SparkMLLibTitanicNewDFbasedAPI.ipynb](http://nbviewer.jupyter.org/github/tuhinmahmud/sigkdd_austin/blob/master/SparkMLLibTitanicNewDFbasedAPI.ipynb)

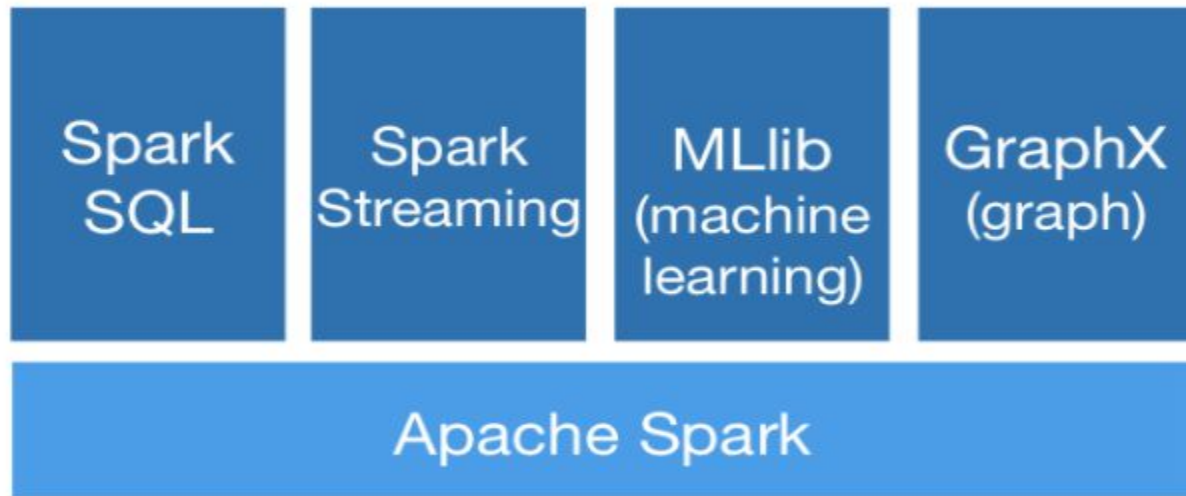


# Notebooks

1. Notebook using small dataset golf play  
[http://nbviewer.jupyter.org/github/tuhinmahmud/sigkdd\\_austin/blob/master/SparkMllibPyspark.golf.ipynb](http://nbviewer.jupyter.org/github/tuhinmahmud/sigkdd_austin/blob/master/SparkMllibPyspark.golf.ipynb)
2. Notebook using titanic dataset and MLib spark dataframe based apis for Decision Tree and Random Forest  
[http://nbviewer.jupyter.org/github/tuhinmahmud/sigkdd\\_austin/blob/master/SparkMLLibTitanicNewDFbasedAPI.ipynb](http://nbviewer.jupyter.org/github/tuhinmahmud/sigkdd_austin/blob/master/SparkMLLibTitanicNewDFbasedAPI.ipynb)

**THANK YOU!**

## Back slide : Spark Stack



# Reference

1. <https://www.slideshare.net/databricks/apache-spark-mllib-20-preview-data-science-and-production>
2. <https://spark.apache.org/mllib/>
3. <https://spark.apache.org/docs/latest/ml-guide.html>
4. [http://www.saedsayad.com/decision\\_tree.htm](http://www.saedsayad.com/decision_tree.htm)
5. <http://www.math.usu.edu/adele/RandomForests/ENAR.pdf>
6. <https://spark.apache.org/docs>
7. <https://spark.apache.org/docs/2.2.0/ml-classification-regression.html#random-forest-classifier>
8. <https://spark.apache.org/docs/2.2.0/api/python/pyspark.ml.html#pyspark.ml.classification.DecisionTreeClassifier>
9. Books: Machine learning with Random Forests And Decision Trees- A visual Guide for Beginner - by Scott Hartshorn
10. **Machine Learning - Decision Trees and Random Forests - by** Loonycorn