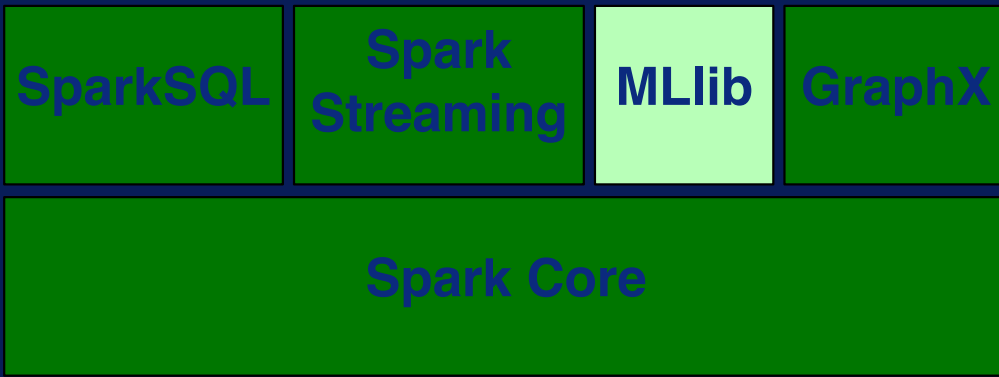# Spark MLlib

# After this video you will be able to..

- Describe what MLlib is

- List main categories of techniques available in MLlib.

- Explain code segments containing MLlib algorithms.

## Spark MLlib

- Scalable machine learning library
- Provides distributed implementations of common machine learning algorithms and utilities
- Has APIs for Scala, Java, Python, and R

# MLlib Algorithms & Techniques

- Machine Learning
  - Classification, regression, clustering, etc.
  - Evaluation metrics
- Statistics
  - Summary statistics, sampling, etc.

  *correlations to sample dataset*

  *mean, stdev, etc ;*
- Utilities
  - Dimensionality reduction, transformation, etc.

  *methods for preprocessing the data*

# MILib Example – Summary Statistics

- Compute column summary statistics

```
from pyspark.mllib.stat import Statistics
```
**1** → import

```
# Data as RDD of Vectors
dataMatrix = sc.parallelize([ [1, 2, 3], [4, 5, 6], [7, 8, 9], [10, 11, 12] ])
```
**2**

create RDD of vectors with data
→ each vector → a column in a data matrix

```
# Compute column summary statistics.
summary = Statistics.colStats(dataMatrix)
print(summary.mean())
print(summary.variance())
print(summary.numNonzeros())
```

**3** → invokes column stats function to compute summary statistics for each column

**4** → print mean, variance and number of non-zero entries for each column

# MLlib Example – Classification

- Build decision tree model for classification

*(6 steps to:)*

```
from pyspark.mllib.tree import DecisionTree, DecisionTreeModel
from pyspark.mllib.util import MLUtils

# Read and parse data
data = sc.textFile("data.txt")



# Decision tree for classification
model = DecisionTree.trainClassifier
            (parsedData, numClasses=2)
print(model.toDebugString())
model.save(sc, "decisionTreeModel")
```

1 → import decision tree module

2 → import MLUtils module

3

4 → build decision tree to classify in two classes

5 → print model

6 → save model in a file

# MLlib Example – Clustering

- Build k-means model for <u>clustering</u>

```python
from pyspark.mllib.clustering import KMeans, KMeansModel
from numpy import array

# Read and parse data
data = sc.textFile("data.txt")
parsedData = data.map(lambda line:
        array([float(x) for x in line.split(' ')]))

# k-means model for clustering
clusters = Kmeans.train (parsedData, k=3)

print(clusters.centers)
```

**1** — imports

**2**

**3** — using space as limiter

**4** — kmeans built by dividing the parsedData into 3 clusters

**5** — cluster centers printed out

# Main Take-Aways

- MLlib is Spark's machine learning library.
  - Distributed implementations
- Main categories of algorithms and techniques:
  - Machine learning
  - Statistics
  - Utility for ML pipeline