# DATA SCIENCE

## SYD DAT 6

## Week 6 – Cloud computing, Big Data & Spark
## Wednesday 16th November

1. Cloud Computing
2. Data Stores and Computation
3. Big Data
4. Spark
5. Lab
6. Real World Problem
7. Review

# CLOUD

**https://aws.amazon.com/ec2/pricing/on-demand/**
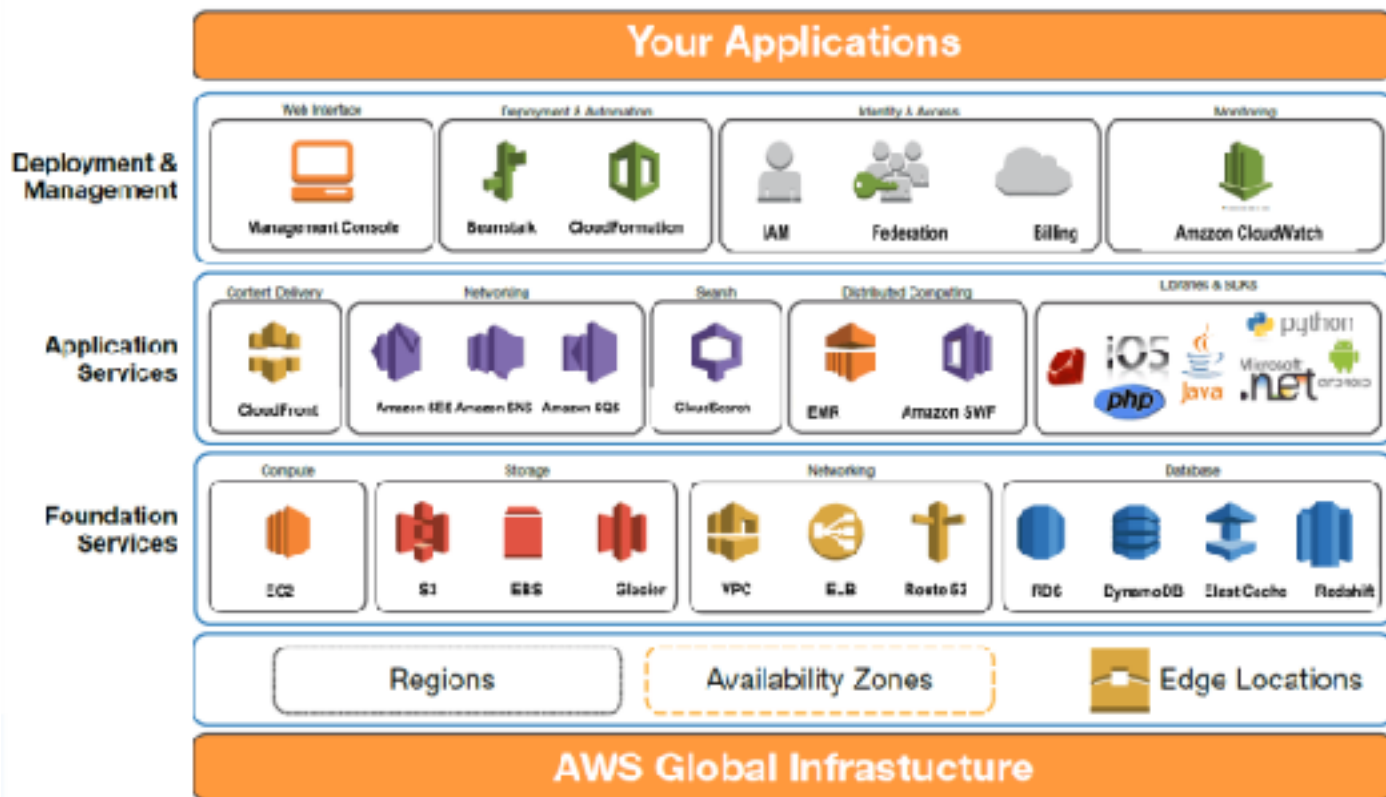
| Linux | RHEL | SLES | Windows | Windows with SQL Standard | Windows with SQL Web |
| --- | --- | --- | --- | --- | --- |

Windows with SQL Enterprise

Region: US East (N. Virginia)

| | vCPU | ECU | Memory (GiB) | Instance Storage (GB) | Linux/UNIX Usage |
| --- | --- | --- | --- | --- | --- |
| **General Purpose - Current Generation** | | | | | |
| t2.nano | 1 | Variable | 0.5 | EBS Only | $0.0065 per Hour |
| t2.micro | 1 | Variable | 1 | EBS Only | $0.013 per Hour |
| t2.small | 1 | Variable | 2 | EBS Only | $0.026 per Hour |
| t2.medium | 2 | Variable | 4 | EBS Only | $0.052 per Hour |
| t2.large | 2 | Variable | 8 | EBS Only | $0.104 per Hour |

Docker containers wrap up a piece of software in a complete filesystem that contains everything it needs to run: code, runtime, system tools, system libraries – anything you can install on a server. This guarantees that it will always run the same, regardless of the environment it is running in.



‣ Lightweight

‣ Open

‣ Secure

# DATA STORES AND COMPUTATION

## Sources

## Structured, relational data

## HOW DO DATA LAKES WORK?

The concept can be compared to a water body, a lake, where water flows in, filling up a reservoir and flows out.

**1** The incoming flow represents multiple raw data archives ranging from emails, spreadsheets, social media content, etc.

**STRUCTURED DATA**
1. Information in rows and columns
2. Easily ordered and processed with data mining tools

**UNSTRUCTURED DATA**
1. Raw, unorganized data
2. Emails
3. PDF files
4. Images, video and audio
5. Social media tools

**2** The reservoir of water is a dataset, where you run analytics on all the data.

**3** The outflow of water is the analyzed data.

**4** Through this process, you are able to "sift" through all the data quickly to gain key business insights.

# NO-SQL

| SQL | NoSQL |
|---|---|
| ‣ Traditional rows and columns data | ‣ No well defined data structure |
| ‣ Strict structure / Primary Keys | ‣ Works better for unstructured data |
| ‣ Entire column for each feature | ‣ Cheaper hardware |
| ‣ Industry standard | ‣ Popular among Startups |

| SQL | NoSQL |
|:---:|:---:|
| ‣ MySQL | ‣ MongoDB |
| ‣ Oracle | ‣ CouchDB |
| ‣ Postgres | ‣ Redis |
| ‣ SQLite | ‣ Cassandra |
| ‣ SQLServer | ‣ Neo4j |
| ‣ Redshift | ‣ HBase |

## JSON - JavaScript Object Notation

‣ Human readable data with attribute-value pairs.

‣ What is inside the curly brackets is an object

‣ In the object we declare variables with 'attribute' : 'value' pairs

```
1  var json = {
2    "firstName": "John",
3    "lastName": "Smith",
4    "age": 25,
5    "address": {
6      "streetAddress": "34 York St",
7      "city": "Sydney",
8      "state": "NSW",
9      "postalCode": "2000"
10   },
11   "phoneNumbers": [
12     {
13       "type": "home",
14       "number": "02 95999999"
15     },
16     {
17       "type": "office",
18       "number": "0431 111 111"
19     }
20   ],
21   "children": [],
22   "spouse": null
23 }
```

‣ Webservices provide application programming interfaces (APIs) are now usually transferring data via JSON

‣ Underlying document databases like MongoDB

‣ Increasingly common data format

# BIG DATA

Data comes from people, technology systems and sensors in the environment

‣ Organisations web applications

‣ External web applications (APIs)

‣ Devices with sensors (Infrastructure, Internet of Things)

# Hortonworks Data Platform

Hortonworks

**GOVERNANCE & INTEGRATION**

Data Workflow, Lifecycle & Governance

Falcon
Sqoop
Flume
NFS
WebHDFS

**DATA ACCESS**

| Batch | Script | SQL | NoSQL | Stream | Others |
|-------|--------|-----|-------|--------|--------|
| Map Reduce | Pig | Hive/Tez HCatalog | HBase Accumulo | Storm | In-Memory Analytics ISV Engines |

**YARN : Data Operating System**

**HDFS**
(Hadoop Distributed File System)

**DATA MANAGEMENT**

**SECURITY**

Authentication
Authorization
Accounting
Data Protection

Storage: HDFS
Resources: YARN
Access: Hive, ...
Pipeline: Falcon
Cluster: Knox

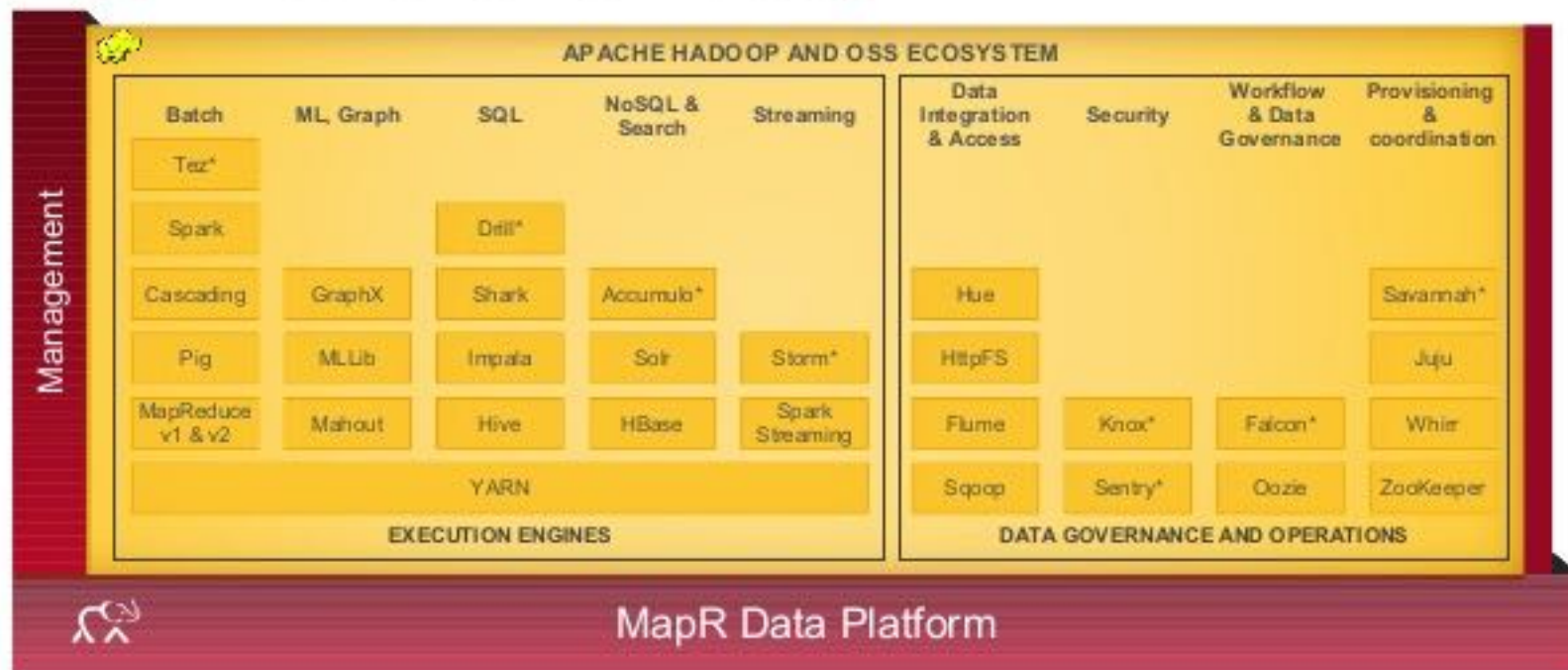**OPERATIONS**

Provision, Manage & Monitor

Ambari
Zookeeper

Scheduling

Oozie

# cloudera



CDH
100% OPEN SOURCE

| CLOUD | USER INTERFACE | WORKFLOW MGMT | METADATA |
| --- | --- | --- | --- |
| WH<br>WHIRR | HU<br>HUE | OO<br>OOZIE | |

INTEGRATION

SQ
SQOOP

FL
FLUME

FILE
FUSE DFS

REST
WEBHDFS
HTTPFS

SQL
ODBC
JDBC

BATCH PROCESSING

| HI<br>HIVE | PI<br>PIG | MA<br>MAHOUT | DF<br>DATAFU |
| --- | --- | --- | --- |

BATCH COMPUTE

| MR<br>MAPREDUCE | MR2<br>MAPREDUCE2 |
| --- | --- |

RESOURCE MGMT
& COORDINATION

| YA<br>YARN | ZO<br>ZOOKEEPER |
| --- | --- |

STORAGE

| HDFS<br>HADOOP DFS | HB<br>HBASE |
| --- | --- |

REAL-TIME ACCESS
& COMPUTE

IM
IMPALA

AC
ACCESS

MS
META
STORE

## MapR Distribution for Hadoop

**APACHE HADOOP AND OSS ECOSYSTEM**

| Management | **EXECUTION ENGINES** | | | | | **DATA GOVERNANCE AND OPERATIONS** | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Batch | ML, Graph | SQL | NoSQL & Search | Streaming | Data Integration & Access | Security | Workflow & Data Governance | Provisioning & coordination |
| | Tez* | | | | | | | | |
| | Spark | | Drill* | | | | | | |
| | Cascading | GraphX | Shark | Accumulo* | | Hue | | | Savannah* |
| | Pig | MLLib | Impala | Solr | Storm* | HttpFS | | | Juju |
| | MapReduce v1 & v2 | Mahout | Hive | HBase | Spark Streaming | Flume | Knox* | Falcon* | Whirr |
| | YARN | | | | | Sqoop | Sentry* | Oozie | ZooKeeper |

**MapR Data Platform**

# BIG DATA FILE FORMATS

Avro stores the data definition in JSON format making it easy to read and interpret, the data itself is stored in binary format making it compact and efficient. Avro files include markers that can be used to splitting large data sets into subsets suitable for MapReduce processing.

a new columnar storage format for Hadoop. Parquet can handle complex nested data structures.

Optimized Row Columnar (ORC). A self-describing type-aware columnar file format designed for Hadoop ecosystem workloads. The columnar format lets the reader read, decompress, and process only the columns that are required for the current query.

# SPARK

Spark is a fast and general processing engine compatible with Hadoop data. It can process data in HDFS, HBase, Cassandra, Hive, and any Hadoop InputFormat. It is designed to perform both batch processing (similar to MapReduce) and new workloads like streaming, interactive queries, and machine learning.
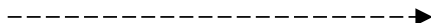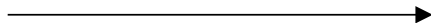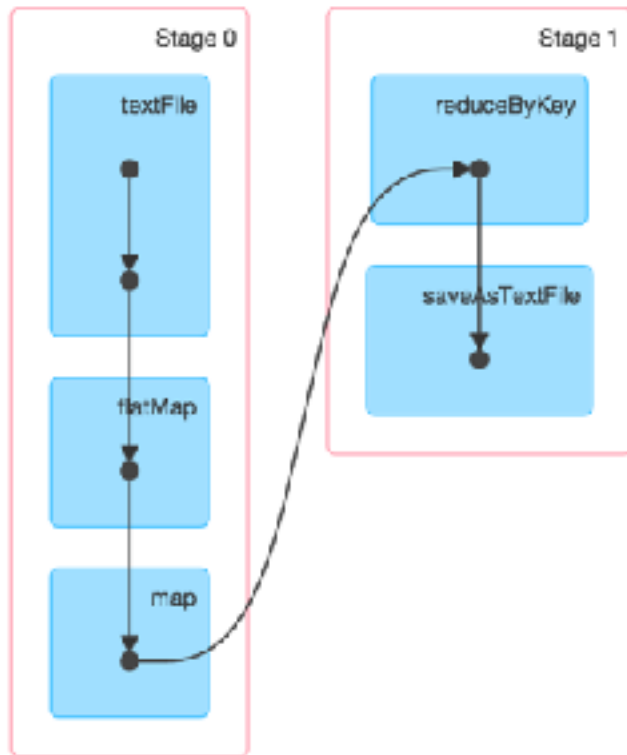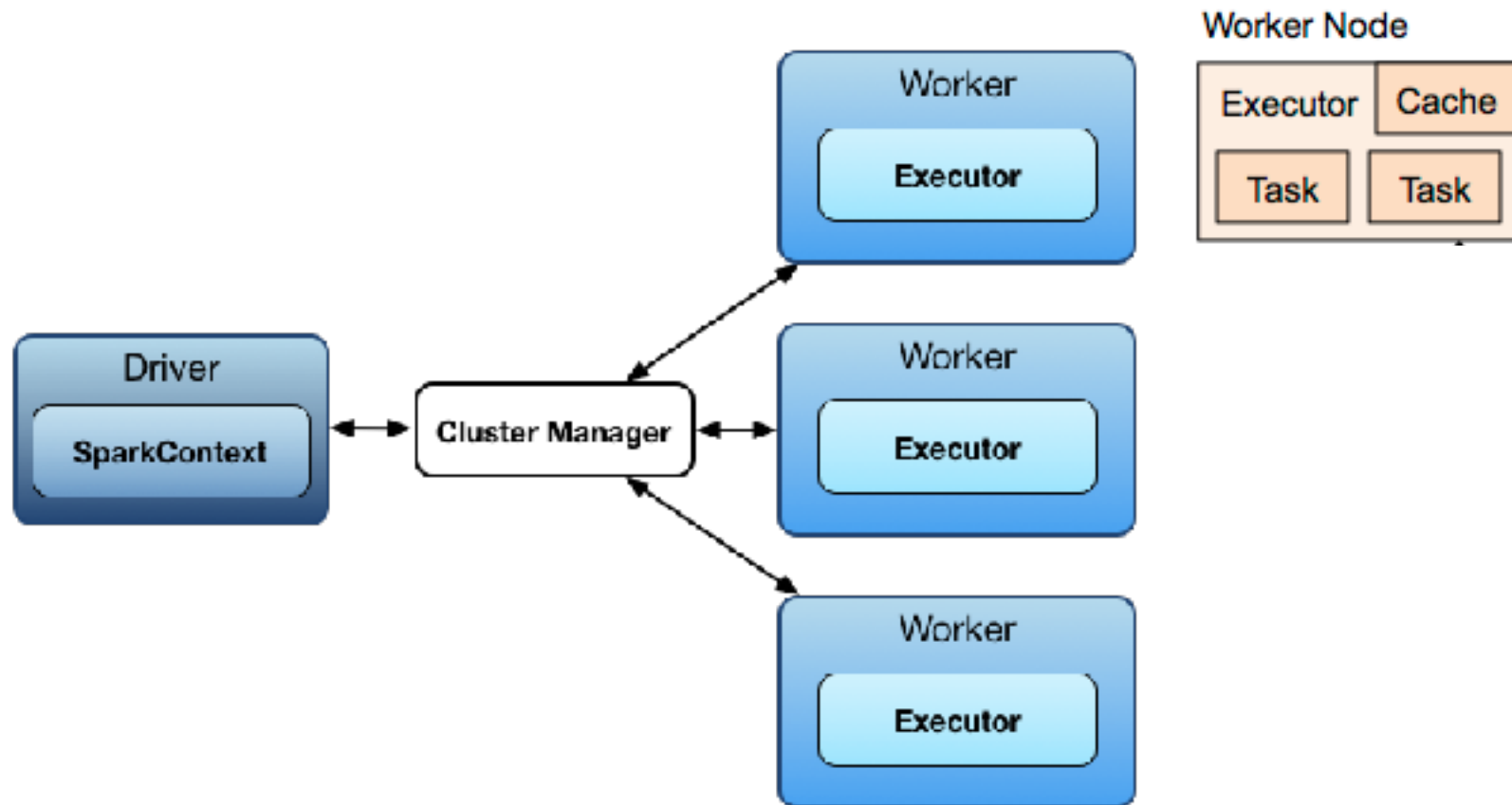
Spark revolves around the concept of a resilient distributed dataset (RDD), which is a fault-tolerant collection of elements that can be operated on in parallel.

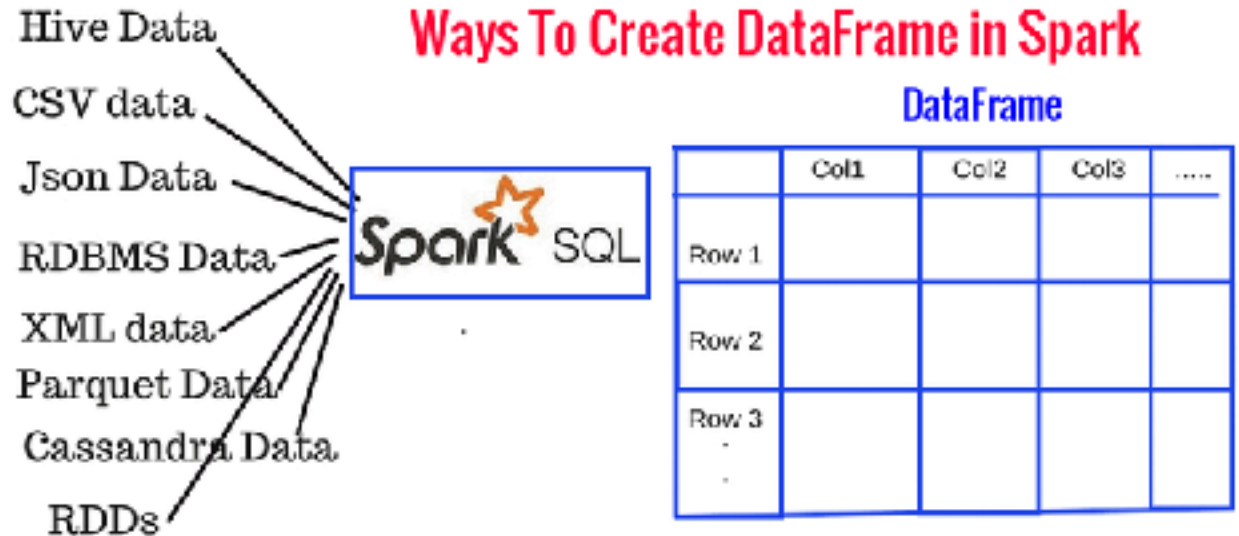There are two ways to create RDDs:

1. Parallelizing an existing collection in your driver program

2. Referencing a dataset in an external storage system, such as a shared filesystem, HDFS, HBase, or any data source offering a Hadoop InputFormat

▾ DAG Visualization

Worker Node

| Executor | Cache |
| --- | --- |
| Task | Task |

Worker

Executor

Worker

Executor

Worker

Executor

Driver

SparkContext

Cluster Manager

DataFrames API was inspired by data frames in R and Pandas in Python. DataFrames integrates with Python, Java, Scala and R.



Ways To Create DataFrame in Spark

Spark Operations = TRANSFORMATIONS + ACTIONS

Create RDD

Transformation

Lineage

RDD

Action

Result

= easy     = medium

## Essential Core & Intermediate Spark Operations

**TRANSFORMATIONS**

**General**
- map
- filter
- flatMap
- mapPartitions
- mapPartitionsWithIndex
- groupBy
- sortBy

**Math / Statistical**
- sample
- randomSplit

**Set Theory / Relational**
- union
- intersection
- subtract
- distinct
- cartesian
- zip

**Data Structure / I/O**
- keyBy
- zipWithIndex
- zipWithUniqueID
- zipPartitions
- coalesce
- repartition
- repartitionAndSortWithinPartitions
- pipe

**ACTIONS**

- reduce
- collect
- aggregate
- fold
- first
- take
- forEach
- top
- treeAggregate
- treeReduce
- forEachPartition
- collectAsMap

- count
- takeSample
- max
- min
- sum
- histogram
- mean
- variance
- stdev
- sampleVariance
- countApprox
- countApproxDistinct

- takeOrdered

- saveAsTextFile
- saveAsSequenceFile
- saveAsObjectFile
- saveAsHadoopDataset
- saveAsHadoopFile
- saveAsNewAPIHadoopDataset
- saveAsNewAPIHadoopFile

Use the programming guides on the Spark Website:

http://spark.apache.org/docs/latest/programming-guide.html



You can select which programming language API you want to write Spark code with :
- Scala
- Jave
- Python (PySpark)
- R (SparkR)

**Data types**

**Basic statistics**

‣ summary statistics

‣ correlations

‣ stratified sampling

‣ hypothesis testing

‣ streaming significance testing

‣ random data generation

**Classification and regression**

‣ linear models (SVMs, logistic regression, linear regression)

‣ naive Bayes

‣ decision trees

‣ ensembles of trees (Random Forests and Gradient-Boosted Trees)

‣ isotonic regression

**Collaborative filtering**

‣ alternating least squares (ALS)

**Clustering**

‣ k-means

‣ Gaussian mixture

‣ power iteration clustering (PIC)

‣ latent Dirichlet allocation (LDA)

‣ bisecting k-means

‣ streaming k-means

**Dimensionality reduction**

‣ singular value decomposition (SVD)

‣ principal component analysis (PCA)

**Feature extraction and transformation**

**Frequent pattern mining**

‣ FP-growth

‣ association rules

‣ PrefixSpan

**Evaluation metrics**

**PMML model export**

**Optimization (developer)**

# COURSE FEEDBACK

# LAB

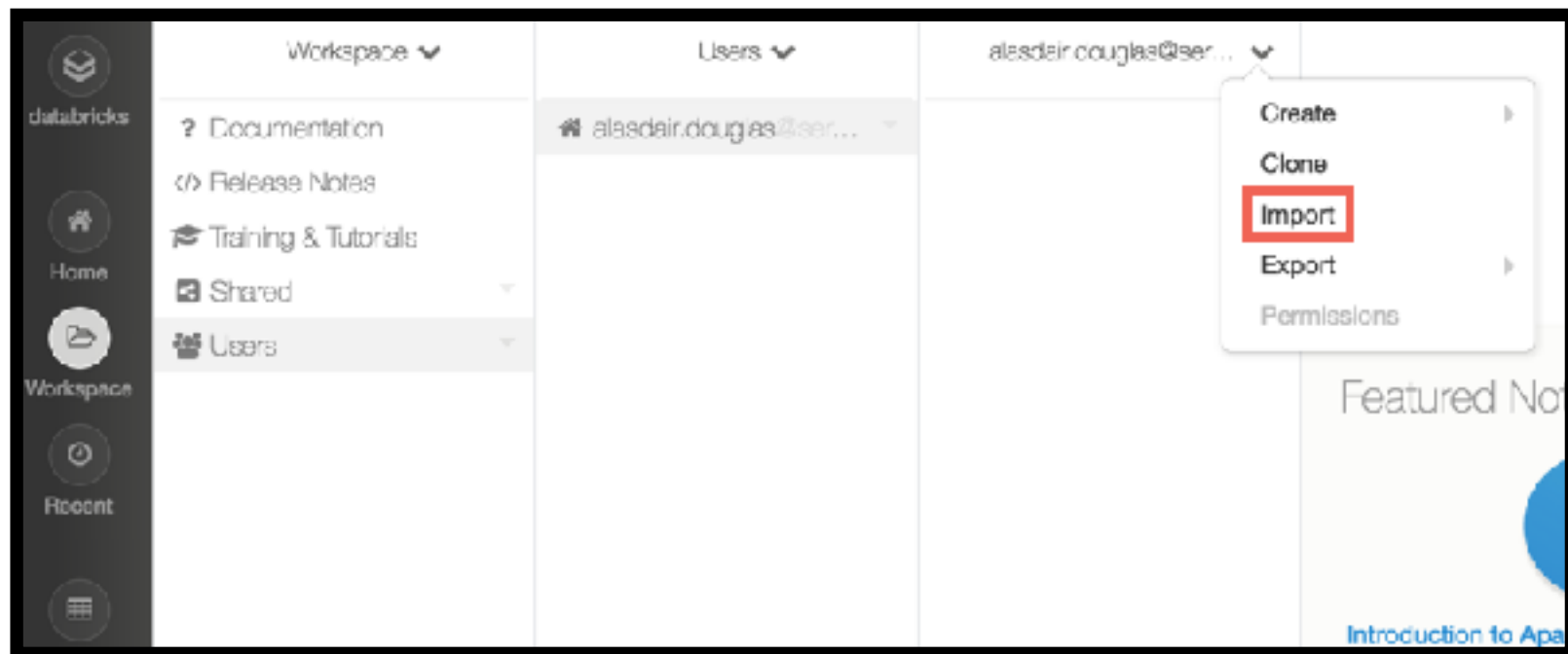‣ **Log into Databricks community edition** <u>**https://databricks.com/try-databricks**</u>

‣ <u>**https://docs.databricks.com/_static/notebooks/gentle-introduction-to-apache-spark.html**</u>

**https://docs.databricks.com/_static/notebooks/gentle-introduction-to-apache-spark.html**

1.  re-name your labs with lab_name.<yourname>.ipynb  (to prevent a conflict)

2.  cd <path to the root of your SYD_DAT_6 local repo>

3.  commit your changes ahead of sync

    - git status

    - git add .

    - git commit -m "descriptive label for the commit"

    - git status

4.  download new material from official course repo (upstream) and merge it

    - git checkout master  (ensures you are in the master branch)

    - git fetch upstream

    - git merge upstream/master

# HOMEWORK

**Homework**

- Read Natural Language Processing website – http://www.nltk.org/ (5 mins)
- Read and be able to explain one use case of the Alchemy API (10 mins)
- Download and install NLTK for Python (10mins)

**Reading**

- For those that want a good foundation in data stores and architecture, the redbook 5th edition is a good reference http://www.redbook.io/