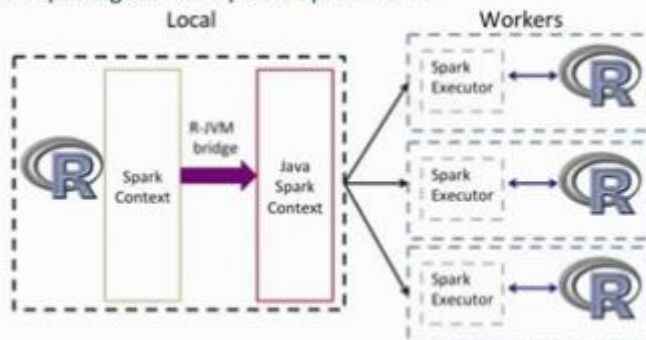## Lesson objectives

- Learn what SparkR is
- Learn why you would use SparkR
- List the features of SparkR
- Understand the interfaces into SparkR

## What is SparkR?

- An R package to use Apache Spark from R



- SparkR implements distributed dataframes
- SparkR supports operations like selection, filtering, aggregation etc.
- SparkR also supports MLIB

## Why use SparkR?

- R has usability issues with big data workflows
- R is very resource constrained so limited optimizations
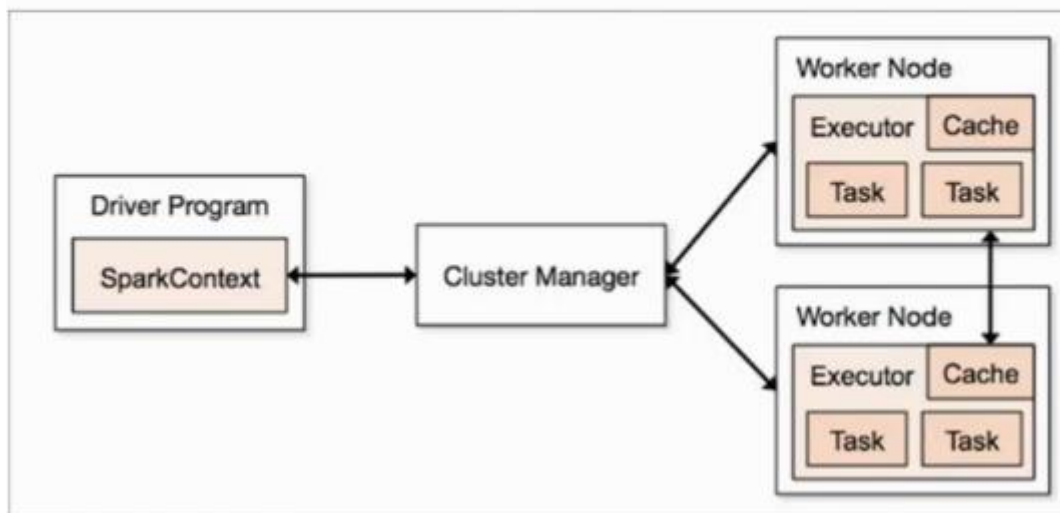- R has restricted machine learning

## Features of SparkR

- Scalability to many cores and machines
- Data Frame Optimizations
- Data Sources API
- RDDs as Distributed Lists
- Serializing closures
- Using existing R packages

# Interfaces

- Spark shell
- SparkR shell
- Rstudio
- Notebooks
- Data Scientist Workbench (an interactive service providing notebooks interface)
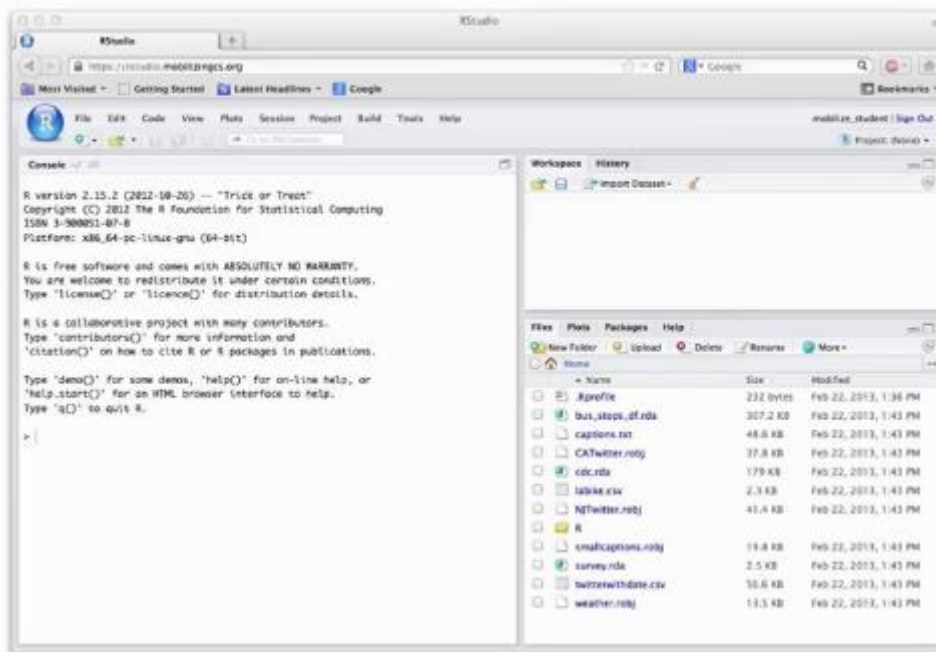
# SparkContext



# Spark SQL
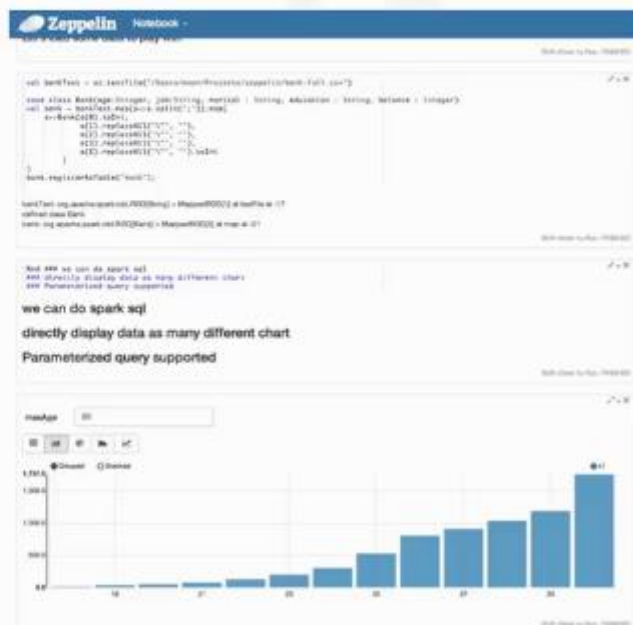
- Spark SQL is a Spark module for structured data processing
- The entry point into all relational functionality in Spark is the SQLContext
- There are several ways to interact with Spark SQL including SQL, the DataFrames API and the Datasets API:
  SQL
  JDBC/ODBC
  Datasets with a strongly-typed LINQ-like Query DSL

- SparkR works solely with DataFrames which requires Spark SQL

# Spark shell

- The following command is used to open Spark shell:
  - Spark-shell

```
$ ./bin/spark-shell
Spark context available as sc.
SQL context available as sqlContext.
Welcome to


   / __/__  ___ _____/ /__
  _\ \/ _ \/ _ `/ __/  '_/
 /___/ .__/\_,_/_/ /_/\_\   version 1.6.0-SNAPSHOT
    /_/

Using Scala version 2.11.7 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_66)
Type in expressions to have them evaluated.
Type :help for more information.

scala>
```

- create a SparkContext and an sqlContext

sc <- sparkR.init()
sqlContext <- sparkRSQL.init(sc)

# SparkR shell

- The following command is used to open SparkR shell:
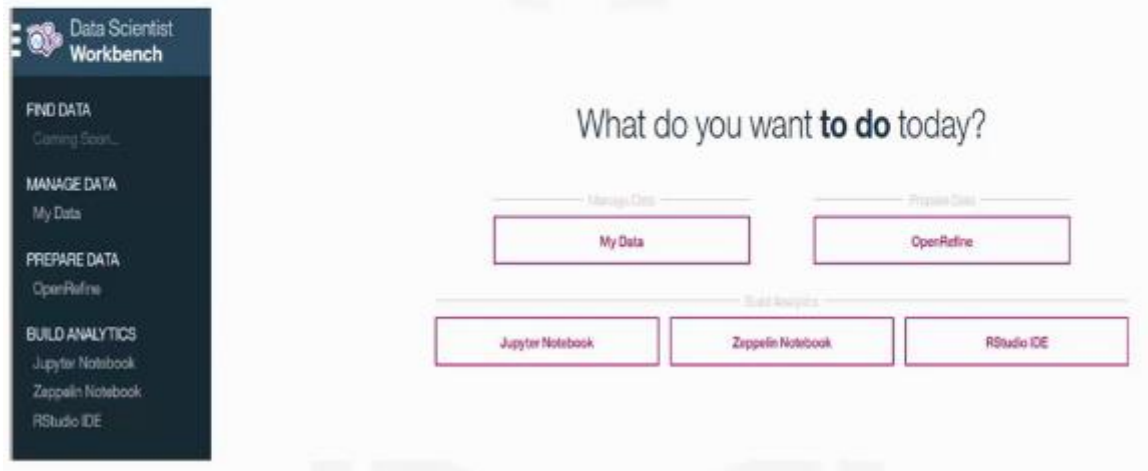  - sparkR shell

# RStudio



# Notebooks

- Zeppelin notebook

# Data Scientist Workbench (DSWB)



## Lesson objectives

- Understand how to use dataframes
- Learn to select data
- Learn to filter data
- Learn to aggregate data
- Learn to operate on columns
- Understand how to write SQL queries

## Dataframes

- Why dataframes
- Creating dataframes
    - Local dataframes
    - Dataframes from raw data
    - Dataframes from data sources

## Name Conflicts

- It is possible to have a name conflict between functions in SparkR and R
- The following have a conflict:
    - cov in package:stats
    - filter in package:stats
    - sample in package:base
    - table in package:base

    - It is also possible to have name conflicts between dplyr and SparkR

## Select columns

- select

```
SparkR::head(select(cars,cars$mpg))
```

```
Out[3]:     mpg
         1 21.0
         2 21.0
         3 22.8
         4 21.4
         5 18.7
         6 18.1
```

- selectExpr

```
SparkR:: head(selectExpr(df, "Day", "Day + 99", "Day * 3"))
```

```
  Day (Day + 99) (Day * 3)
1  1       100         3
2  2       101         6
3  3       102         9
4  4       103        12
5  5       104        15
6  6       105        18
```

# Operating on columns

- Operating on a specific column using $

```
cars$mpg <- cars$mpg/3.78541178
```

```
Out[4]:       mpg cyl disp  hp drat    wt  qsec vs am gear carb
         1 5.547613   6 160 110 3.90 2.620 16.46  0  1    4    4
         2 5.547613   6 160 110 3.90 2.875 17.02  0  1    4    4
         3 6.023123   4 108  93 3.85 2.320 18.61  1  1    4    1
         4 5.653282   6 258 110 3.08 3.215 19.44  1  0    3    1
         5 4.940017   8 360 175 3.15 3.440 17.02  0  0    3    2
         6 4.781514   6 225 105 2.76 3.460 20.22  1  0    3    1
```

- Operating on multiple columns using [ ]

```
[, c("you" , "me") ]
```

# Filter (select rows )

- filter by condition for one column

  filtered <- SparkR::head(SparkR::filter(cars, df$cyl == 6)
  collect(filtered)

  ```
      mpg cyl  disp  hp drat    wt  qsec vs am gear carb
  1 21.0   6 160.0 110 3.90 2.620 16.46  0  1    4    4
  2 21.0   6 160.0 110 3.90 2.875 17.02  0  1    4    4
  3 21.4   6 258.0 110 3.08 3.215 19.44  1  0    3    1
  4 18.1   6 225.0 105 2.76 3.460 20.22  1  0    3    1
  ```

- filter by condition for multiple columns
  - These two expression are same:

    SparkR::head(SparkR:: filter("cars, df$cyl == 6 and cars, cars$mpg < 20")
    SparkR::head(SparkR:: filter(cars, df$cyl == 6). SparkR::filter(cars, cars$mpg < 20)

# SparkSQL (row operations)

- Dataframes as SQL tables

  registerTempTable(cars,"cars")
  SparkR::head(sql(sqlContext, "SELECT * FROM cars WHERE cyl > 6"))

  ```
      mpg cyl  disp  hp drat   wt  qsec vs am gear carb
  1 18.7   8 360.0 175 3.15 3.44 17.02  0  0    3    2
  2 14.3   8 360.0 245 3.21 3.57 15.84  0  0    3    4
  3 16.4   8 275.8 180 3.07 4.07 17.40  0  0    3    3
  ```

# Aggregate data

- groupBy for one column (count, sum, average)
  - groupBy(cars, cars$mpg), count = n(cars$mpg))
  - groupBy(cars, cars$mpg), sum = sum(cars$mpg))
  - groupBy(cars, cars$mpg), average = avg(cars$hp))

- groupBy for multiple columns (count, sum, average)
  - groupBy(cars, "'class', 'year'"), avg = n('hwy'))

## Lesson objectives

- Understand machine learning
- Learn how to use GLM model

## Machine learning

- Example machine learning workflow



- Observation = rows in a dataframe
- Feature = columns in a dataframe

## MLlib and ML

- MLlib is Spark's machine learning (ML) library
- It divides into two packages:
  - spark.mllib contains the original API built on top of RDDs.
  - spark.ml provides higher-level API built on top of DataFrames
- Concepts in pipelines:
  - DataFrame
  - Transformer
  - Estimator
  - Pipeline
  - Parameter

# Pipelines components

- Transformers
  - A Transformer is an abstraction that includes feature transformers and learned models
- Estimators
  - An Estimator abstracts the concept of a learning algorithm that trains on data
- Pipeline
  - a sequence of algorithms to process and learn from data

# Implementing Linear Models in SparkR

- Prepare and load data
  - Load the data
  - Read the data in a Spark dataframe
  - Create factors
- Train the model
  - Formula
  - Dataset
  - Model
- Evaluate the model
  - Baseline reference
  - Predict()
- Implement model
  - Iterate through the training dataset
  - Find acceptable model

# GLM data load

- Download and save the data files using R
  - population_data_files_url <-
    'http://www2.census.gov/acs2013_1yr/pums/csv_pus.zip'
  - library(RCurl)
  - population_data_file <- getBinaryURL(population_data_files_url)
  - population_data_file_path <- '/nfs/data/2013-acs/csv_pus.zip'
  - population_data_file_local <- file(population_data_file_path, open = "wb")
- sparkContext
  - sc <- sparkR.init()
- sqlContext
  - sqlContext <- sparkRSQL.init(sc)

# GLM data prep

- Convert any categorical variable from a numeric variable into a factor for example
  - housing_df$ST <- cast(housing_df$ST, "string")
- Sanitize the data  for example remove nulls
  - housing_with_valp_df <- filter(housing_df,
    isNotNull(housing_df$VALP)
    & isNotNull(housing_df$TAXP)
    & isNotNull(housing_df$INSP)
    & isNotNull(housing_df$ACR))
- Create training and test data
  - housing_df_test <- sample(housing_with_valp_df,FALSE,0.1)

# GLM model training

- glm()
  - Fit a gaussian GLM model over the dataset.
    - model <- glm(Sepal_Length ~ Sepal_Width + Species, data = df, family = "gaussian")
- summary()
  - Model summary are returned in a similar format to R's native glm().
    - summary(model)
    - ##$devianceResiduals
    - ## Min     Max
    - ## -1.307112 1.412532
    - etc.......

# GLM model evaluation

- predict()
  - > preds <- predict(model, training)
  - > errors <- select(
    - preds, preds$label, preds$prediction, preds$aircraft_type,
    - alias(preds$label - preds$prediction, "error"))

# GLM model vizualization

- matplotlib
- Ggplot2
- Third party libraries such as prettyplotlib