**Spark**

User-friendly API for data transformation

Large and active community

Platform components - SQL

Multitenancy

**H$_2$O**

Memory efficient

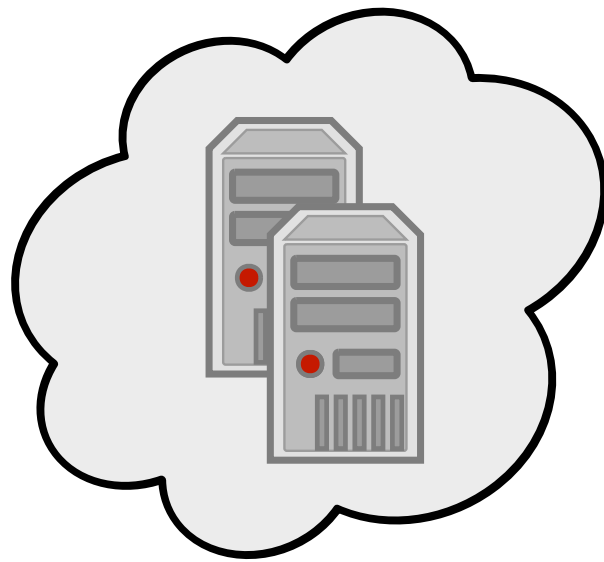Performance of computation

Machine learning algorithms

Parser, GUI, R-interface
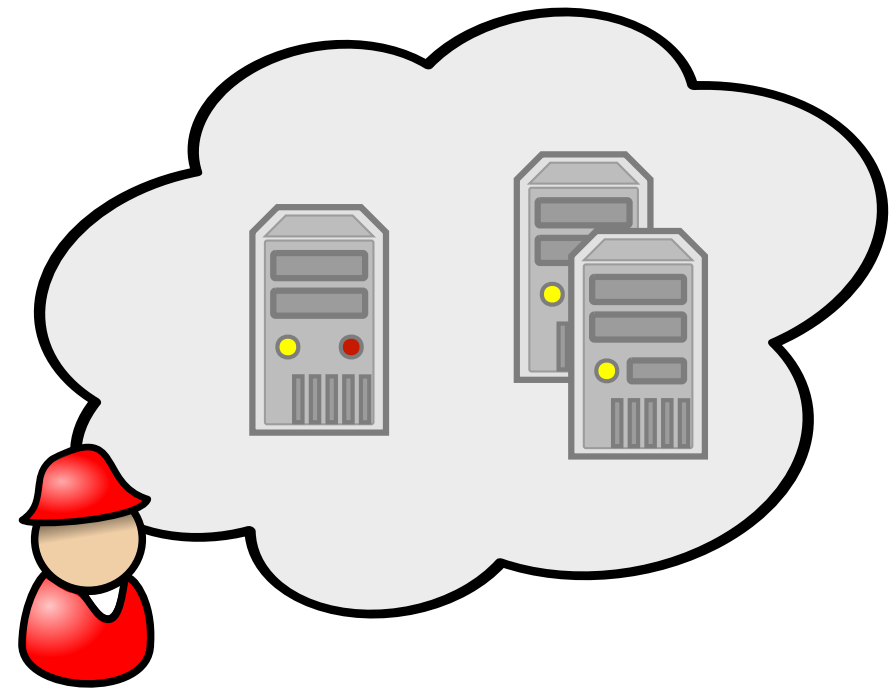
0xdata

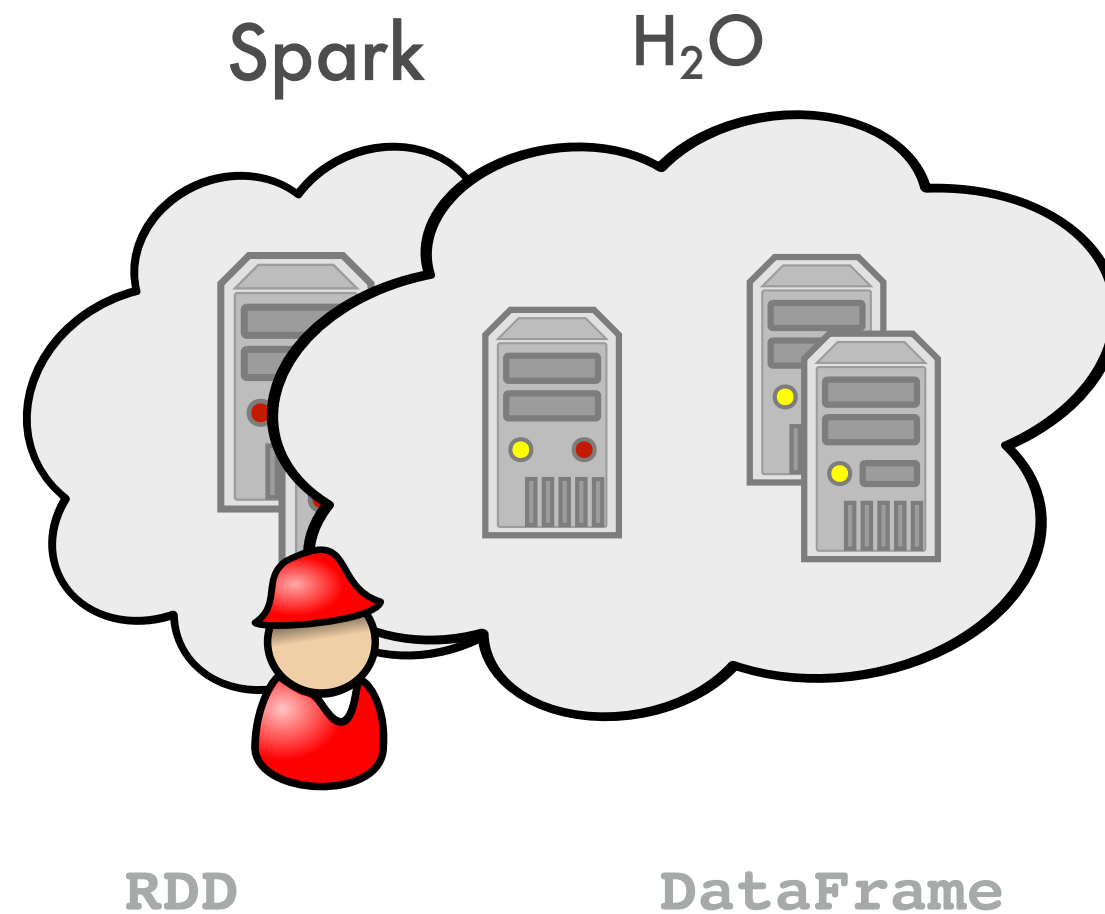# Sparkling Water

Spark

$H_2O$



+

RDD
immutable
world

DataFrame
mutable
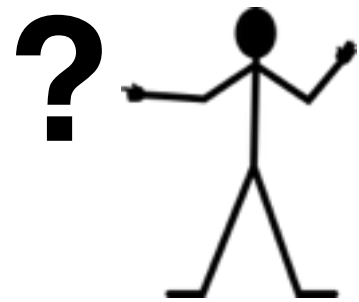world

# Sparkling Water

# Sparkling Water

**Provides**

Transparent integration into Spark ecosystem

Pure H2ORDD encapsulating $H_2O$ DataFrame

Transparent use of $H_2O$ data structures and algorithms with Spark API

**Excels in Spark workflows requiring advanced Machine Learning algorithms**

# Sparkling Water Design

# Data Distribution



Data Source (e.g. HDFS)

Sparkling Water Cluster

Spark Executor JVM

H2O

Spark RDD

H2O RDD

Spark Executor JVM

H2O

Spark Executor JVM

H2O

RDDs and DataFrames share same memory space

# Hands-On #1
# Sparkling Shell

# Sparkling Water Requirements

**Linux or Mac OS X**

**Oracle Java 1.7+**

**Spark 1.1.0**

**Provided on USB stick
or download from
http://meetups.h2o.ai/sw.zip**

# Sparkling Water Download

**http://h2o.ai/download/**

# Where is the code?

**https://github.com/h2oai/sparkling-water/blob/master/examples/scripts/**

# Flight delays prediction

**"Build a model using weather and flight data to predict delays of flights arriving to Chicago O'Hare International Airport"**

# Example Outline

Load & Parse CSV data from 2 data sources

Use Spark API to filter data, do SQL query for join

Create regression models

Use models to predict delays

Graph residual plot from R

# Install and Launch

**Unpack zip file**

and

**Point SPARK_HOME to your Spark 1.2.0 installation**

It is on USB stick

and

**Launch bin/sparkling-shell**

# What is Sparkling Shell?

**Standard spark-shell**

**With additional Sparkling Water classes**

Spark Master address

```
export MASTER="local-cluster[3,2,1024]"

spark-shell \
    --jars sparkling-water.jar
```

JAR containing Sparkling Water

# Lets play with Sparkling shell...

# Create H$_2$O Client

Contains implicit utility functions

Demo specific classes

Size of demanded H$_2$O cloud

```scala
import org.apache.spark.h2o._
import org.apache.spark.examples.h2o._

val h2oContext = new H2OContext(sc).start()
import h2oContext._
```

Regular Spark context provided by Spark shell

# Is Spark Running?

**Go to http://localhost:4040**

# Is H₂O running?

## http://localhost:54321/flow/index.html

# Load Data #1

**Load weather data into RDD**

```
val weatherDataFile =
  "examples/smalldata/
Chicago_Ohare_International_Airport.csv"

val wrawdata = sc.textFile(weatherDataFile,3)
                    .cache()

val weatherTable = wrawdata
      .map(_.split(","))
      .map(row => WeatherParse(row))
      .filter(!_.isWrongRow())
```

Regular Spark API

Ad-hoc Parser

# Weather Data

```scala
case class Weather( val Year   : Option[Int],
                    val Month  : Option[Int],
                    val Day    : Option[Int],
                    val TmaxF  : Option[Int],   // Max temperatur in F
                    val TminF  : Option[Int],   // Min temperatur in F
                    val TmeanF : Option[Float], // Mean temperatur in F
                    val PrcpIn : Option[Float], // Precipitation (inches)
                    val SnowIn : Option[Float], // Snow (inches)
                    val CDD    : Option[Float], // Cooling Degree Day
                    val HDD    : Option[Float], // Heating Degree Day
                    val GDD    : Option[Float]) // Growing Degree Day
```

Simple POJO to hold one row of weather data

# Load Data #2

**Load flights data into H2O frame**

```scala
import java.io.File

val dataFile =
  "examples/smalldata/year2005.csv.gz"

val airlinesData = new DataFrame(new File(dataFile))
```

Shortcut for data load and parse

# Where is the data?

**Go to http://localhost:54321/flow/
index.html**

# Use Spark API for Data Filtering

Create a cheap wrapper around $H_2O$ DataFrame

```
// Create RDD wrapper around DataFrame
val airlinesTable : RDD[Airlines]
                  = asRDD[Airlines](airlinesData)

// And use Spark RDD API directly
val flightsToORD = airlinesTable
    .filter( f => f.Dest == Some("ORD") )
```

Regular Spark RDD call

# Use Spark SQL to Data Join

```scala
import org.apache.spark.sql.SQLContext
// We need to create SQL context
implicit val sqlContext = new SQLContext(sc)
import sqlContext._

flightsToORD.registerTempTable("FlightsToORD")
weatherTable.registerTempTable("WeatherORD")
```

Make context implicit to share it with h2oContext

# Join Data based on Flight Date

```scala
val joinedTable = sql(
  """SELECT
    | f.Year,f.Month,f.DayofMonth,
    | f.CRSDepTime,f.CRSArrTime,f.CRSElapsedTime,
    | f.UniqueCarrier,f.FlightNum,f.TailNum,
    | f.Origin,f.Distance,
    | w.TmaxF,w.TminF,w.TmeanF,
    | w.PrcpIn,w.SnowIn,w.CDD,w.HDD,w.GDD,
    | f.ArrDelay
    | FROM FlightsToORD f
    | JOIN WeatherORD w
    | ON f.Year=w.Year AND f.Month=w.Month
    |   AND f.DayofMonth=w.Day""".stripMargin)
```

# Split data

```scala
import hex.splitframe.SplitFrame
import hex.splitframe.SplitFrameModel.SplitFrameParameters

val sfParams = new SplitFrameParameters()
sfParams._train = joinedTable
sfParams._ratios = Array(0.7, 0.2)
val sf = new SplitFrame(sfParams)

val splits = sf.trainModel().get._output._splits
val trainTable = splits(0)
val validTable = splits(1)
val testTable  = splits(2)
```

Result of SQL query is implicitly converted into H2O DataFrame

# Launch H₂O Algorithms

```scala
import hex.deeplearning._
import hex.deeplearning.DeepLearningModel
        .DeepLearningParameters

// Setup deep learning parameters
val dlParams = new DeepLearningParameters()
dlParams._train = trainTable
dlParams._response_column = 'ArrDelay
dlParams._valid = validTable
dlParams._epochs = 100
dlParams._reproducible = true
dlParams._force_load_balance = false


// Create a new model builder
val dl = new DeepLearning(dlParams)

val dlModel = dl.trainModel.get
```

Blocking call

# Make a prediction

```scala
// Use model to score data
val dlPredictTable = dlModel.score(testTable)('predict)

// Collect predicted values via RDD API
val predictionValues = asSchemaRDD(dlPredictTable)
                            .collect
                            .map (row =>

     if (row.isNullAt(0))
        Double.NaN
     else
        row(0)
```

# Hands-On #2

# Can I access results from R?

# YES!

# Requirements

**R 3.1.2+**

**RStudio**

**H2O R package**

# Install R package

**You can find R package on USB stick**

1. Open RStudio

2. Click on "Install Packages"

3. Select h2o_0.1.22.99999.tar.gz file from USB



install.packages("sparkling-water-meetup/R/h2o_0.1.22.99999.tar.gz", repos = NULL, type = "source")

# Generate R code

**In Sparkling Shell:**

```
import org.apache.spark.examples.h2o.DemoUtils.residualPlotRCode

residualPlotRCode(
    predictionH2OFrame, 'predict,
    testFrame, 'ArrDelay)
```

**Utility generating R code to show residuals plot for predicted and actual values**

# Residuals Plot in R

```r
# Import H2O library and initialize H2O client
library(h2o)

h = h2o.init()

# Fetch prediction and actual data, use remembered keys
pred = h2o.getFrame(h, "dframe_b5f449d0c04ee75fda1b9bc865b14a69")
act = h2o.getFrame (h, "frame_rdd_14_b429e8b43d2d8c02899ccb61b72c4e57")

# Select right columns
predDelay = pred$predict
actDelay = act$ArrDelay

# Make sure that number of rows is same
nrow(actDelay) == nrow(predDelay)

# Compute residuals
residuals = predDelay - actDelay

# Plot residuals
compare = cbind(
    as.data.frame(actDelay$ArrDelay),
    as.data.frame(residuals$predict))

plot( compare[,1:2] )
```
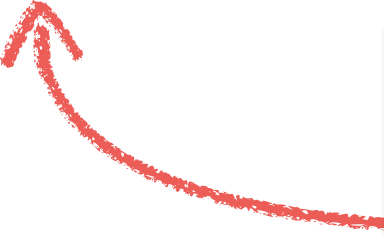
References of data

# Warning!

**If you are running R v3.1.0 you will see different plot:**



**Why? Float number handling was changed in that version. Our recommendation is to upgrade your R to the newest version.**

# Try GBM Algo

```scala
import hex.tree.gbm.GBM
import hex.tree.gbm.GBMModel.GBMParameters

val gbmParams = new GBMParameters()
gbmParams._train = trainTable
gbmParams._response_column = 'ArrDelay
gbmParams._valid = validTable
gbmParams._ntrees = 100

val gbm = new GBM(gbmParams)
val gbmModel = gbm.trainModel.get

// Print R code for residual plot
val gbmPredictTable = gbmModel.score(testTable)('predict)
printf( residualPlotRCode(gbmPredictTable, 'predict, testTable,
'ArrDelay) )
```
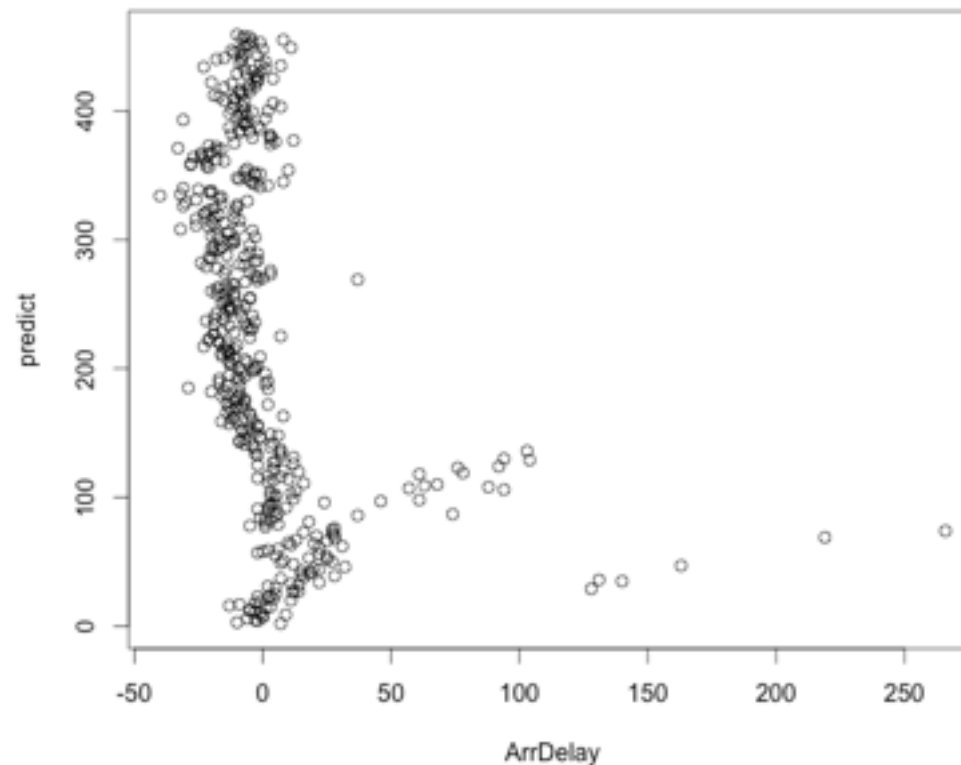
# Residuals plot for GBM prediction

# Hands-On #3

# How Can I Develop and Run Standalone App?

# Requirements

**Idea or Eclipse**

**Git**

# Use Sparkling Water Droplet

**Clone H2O Droplets repository**

```
git clone https://github.com/h2oai/h2o-droplets.git

cd h2o-droplets/sparkling-water-droplet/
```

# Generate IDE project

**For Idea**

`./gradlew idea`

**For Eclipse**

`./gradlew eclipse`

**... add import project into your IDE**

# Create An Application

```scala
object AirlinesWeatherAnalysis {

  /** Entry point */
  def main(args: Array[String]) {
    // Configure this application
    val conf: SparkConf = new SparkConf().setAppName("Flights Water")
    conf.setIfMissing("spark.master", sys.env.getOrElse("spark.master", "local"))

    // Create SparkContext to execute application on Spark cluster
    val sc = new SparkContext(conf)
    // Start H2O cluster only
    new H2OContext(sc).start()

    // User code
    // . . .
  }
}
```

Create
Spark Context
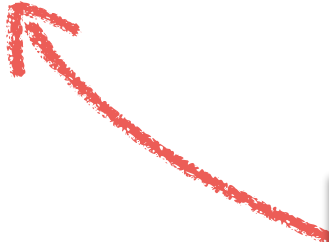
Create H2O context
and start H2O on top
of Spark

# Build the Application

**Build and test**

```
./gradlew build shadowJar
```

**Create an assembly
which can be submitted
to Spark cluster**

# Run code on Spark

```bash
#!/usr/bin/env bash

APP_CLASS=water.droplets.AirlineWeatherAnalysis
FAT_JAR_FILE="build/libs/sparkling-water-droplet-app.jar"
MASTER=${MASTER:-"local-cluster[3,2,1024]"}
DRIVER_MEMORY=2g

$SPARK_HOME/bin/spark-submit "$@" \
  --driver-memory $DRIVER_MEMORY \
  --master $MASTER \
  --class "$APP_CLASS" $FAT_JAR_FILE
```

# It is Open Source!

**You can participate in**

**H2O Scala API**

**Sparkling Water testing**

**Mesos, Yarn, workflows (PUBDEV-23,26,27,31-33)**

**Spark integration**

**MLLib Pipelines**

**Check out our JIRA
at http://jira.h2o.ai**

# Come to Meetup

**http://www.meetup.com/Silicon-Valley-Big-Data-Science/**

# More info

**Checkout H2O.ai Training Books**

**http://learn.h2o.ai/**

**Checkout H2O.ai Blog for Sparkling Water tutorials**

**http://h2o.ai/blog/**

**Checkout H2O.ai Youtube Channel**

**https://www.youtube.com/user/0xdata**

**Checkout GitHub**

**https://github.com/h2oai/sparkling-water**

# And the winner is ...