



OBJECT  
COMPUTING



# Data Science with Groovy

**Presented by  
Dr Paul King**

# Dr Paul King



*OCI Groovy Lead*

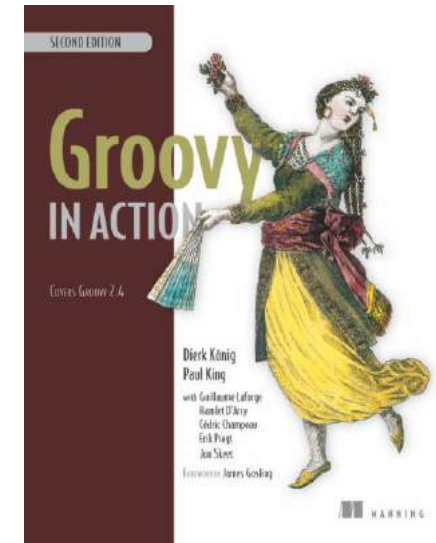
*V.P. and PMC Chair Apache Groovy*

*Author:* <https://www.manning.com/books/groovy-in-action-second-edition>

<https://speakerdeck.com/paulk/groovy-data-science>

<https://github.com/paulk-asert/groovy-data-science>

@paulk\_asert



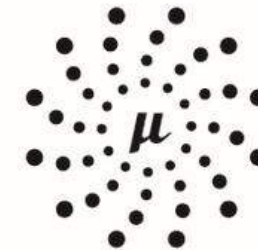
# WE ARE SOFTWARE ENGINEERS.



We deliver mission-critical software solutions that accelerate innovation within your organization and stand up to the evolving demands of your business.



Grails



MICRONAUT

# WE ARE SOFTWARE ENGINEERS.

We deliver mission-critical software solutions that accelerate innovation within your organization and stand up to the evolving demands of your business.

## AREAS OF EXPERTISE



### INDUSTRIAL IoT

We equip industrial environments with seamless connectivity and real-time analytics that reduce operating costs and deliver on customer demands.



### MACHINE LEARNING

We can modernize your legacy applications and enable scalable, integrated AI capabilities tailored to your unique sets of data.



### BLOCKCHAIN CONSULTING

We are at the forefront of blockchain technology, and we have practical, real-world experience with its implementation.



### CLOUD SOLUTIONS

We combine software engineering expertise with cloud-native architecture to accelerate innovation within your organization.

# What is Groovy?

*It's like a super version of Java:*

- *Supports most Java syntax but allows simpler syntax for many constructs*
- *Supports all Java libraries but provides many extensions and its own productivity libraries*
- *Has both a static and dynamic nature*
- *Extensible language and tooling*



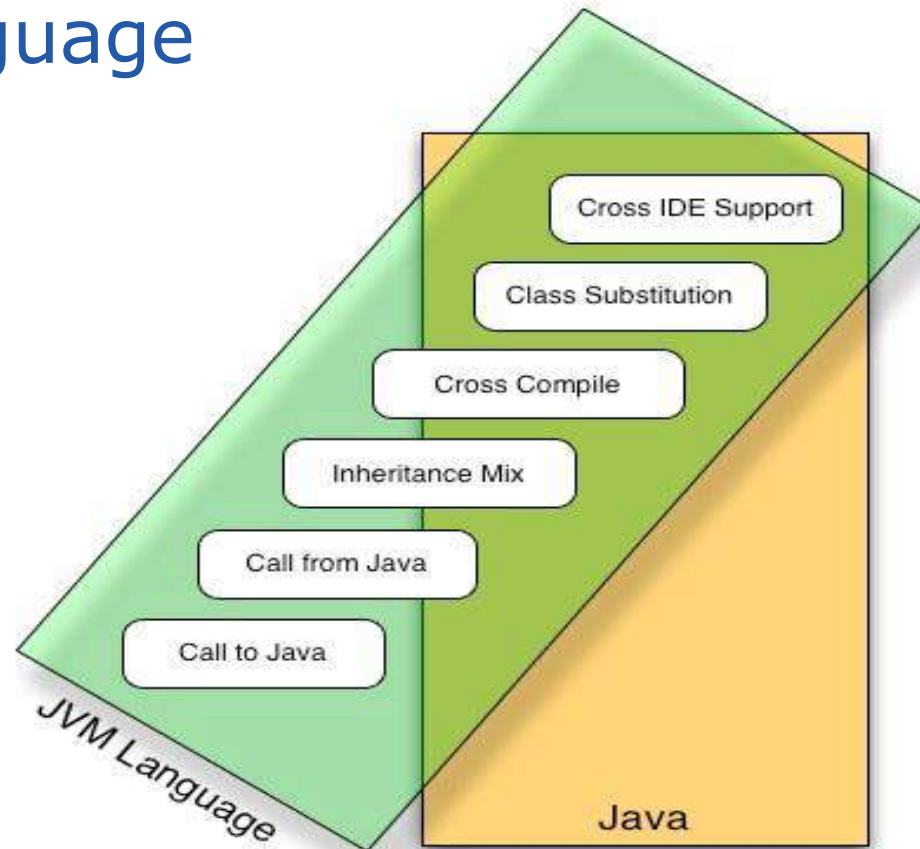
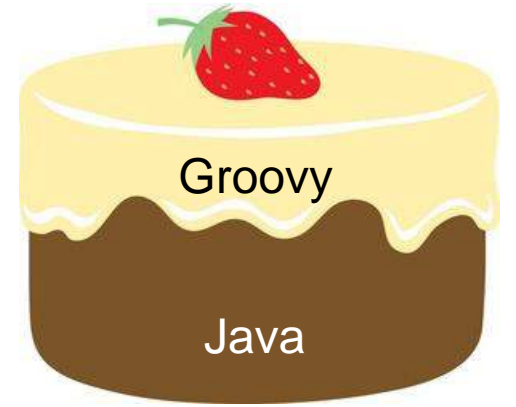


Groovy has a low learning curve, particularly for Java developers

# Groovy is Java's friend

## Seamless integration

- IDEs provide cross-language compile, navigation, and refactoring
- Arbitrarily mix source language
- Drop-in replace any class
- Overloaded methods
- Syntax alignment
- Shared data types



# Java-like or concise shortcuts

```
import java.util.List;
import java.util.ArrayList;

class Main {
    private List keepShorterThan(List strings, int length) {
        List result = new ArrayList();
        for (int i = 0; i < strings.size(); i++) {
            String s = (String) strings.get(i);
            if (s.length() < length) {
                result.add(s);
            }
        }
        return result;
    }

    public static void main(String[] args) {
        List names = new ArrayList();
        names.add("Ted"); names.add("Fred");
        names.add("Jed"); names.add("Ned");
        System.out.println(names);
        Main m = new Main();
        List shortNames = m.keepShorterThan(names, 4);
        System.out.println(shortNames.size());
        for (int i = 0; i < shortNames.size(); i++) {
            String s = (String) shortNames.get(i);
            System.out.println(s);
        }
    }
}
```

```
names = ["Ted", "Fred", "Jed", "Ned"]
println names
shortNames = names.findAll{ it.size() < 4 }
println shortNames.size()
shortNames.each{ println it }
```





# Concise syntax including DSL friendly

```
import java.util.List;
import java.util.ArrayList;

class Main {
    private List keepShorterThan(List strings, int length) {
        List result = new ArrayList();
        for (int i = 0; i < strings.size(); i++) {
            String s = (String) strings.get(i);
            if (s.length() < length) {
                result.add(s);
            }
        }
        return result;
    }

    public static void main(String[] args) {
        List names = new ArrayList();
        names.add("Ted"); names.add("Fred");
        names.add("Jed"); names.add("Ned");
        System.out.println(names);
        Main m = new Main();
        List shortNames = m.keepShorterThan(names, 4);
        System.out.println(shortNames.size());
        for (int i = 0; i < shortNames.size(); i++) {
            String s = (String) shortNames.get(i);
            System.out.println(s);
        }
    }
}
```

```
names = ["Ted", "Fred", "Jed", "Ned"]
println names
shortNames = names.findAll{ it.size() < 4 }
println shortNames.size()
shortNames.each{ println it }
```



given the names "Ted", "Fred", "Jed" and "Ned"  
display all the names  
display the number of names having size less than 4  
display the names having size less than 4



# Metaprogramming

```
public final class Person {
    private final String first;
    private final String last;

    public String getFirst() {
        return first;
    }

    public String getLast() {
        return last;
    }

    @Override
    public int hashCode() {
        final int prime = 31;
        int result = 1;
        result = prime * result + ((first == null)
            ? 0 : first.hashCode());
        result = prime * result + ((last == null)
            ? 0 : last.hashCode());
        return result;
    }

    public Person(String first, String last) {
        this.first = first;
        this.last = last;
    }
}
// ...
```



```
// ...
@Override
public boolean equals(Object obj) {
    if (this == obj)
        return true;
    if (obj == null)
        return false;
    if (getClass() != obj.getClass())
        return false;
    Person other = (Person) obj;
    if (first == null) {
        if (other.first != null)
            return false;
    } else if (!first.equals(other.first))
        return false;
    if (last == null) {
        if (other.last != null)
            return false;
    } else if (!last.equals(other.last))
        return false;
    return true;
}

@Override
public String toString() {
    return "Person(first:" + first
        + ", last:" + last + ")";
}
boilerplate
```

Groovy's metaprogramming capabilities allow complex design patterns to be encapsulated within simple constructs.

```
@Immutable class Person {
    String first, last
}
```



# Matrix manipulation example

```
import org.apache.commons.math3.linear.*;

public class MatrixMain {
    public static void main(String[] args) {
        double[][] matrixData = { {1d,2d,3d}, {2d,5d,3d}};
        RealMatrix m = MatrixUtils.createRealMatrix(matrixData);

        double[][] matrixData2 = { {1d,2d}, {2d,5d}, {1d, 7d}};
        RealMatrix n = new Array2DRowRealMatrix(matrixData2);

        RealMatrix o = m.multiply(n);

        // Invert p, using LU decomposition
        RealMatrix oInverse = new LUdecomposition(o).getSolver().getInverse();

        RealMatrix p = oInverse.scalarAdd(1d).scalarMultiply(2d);

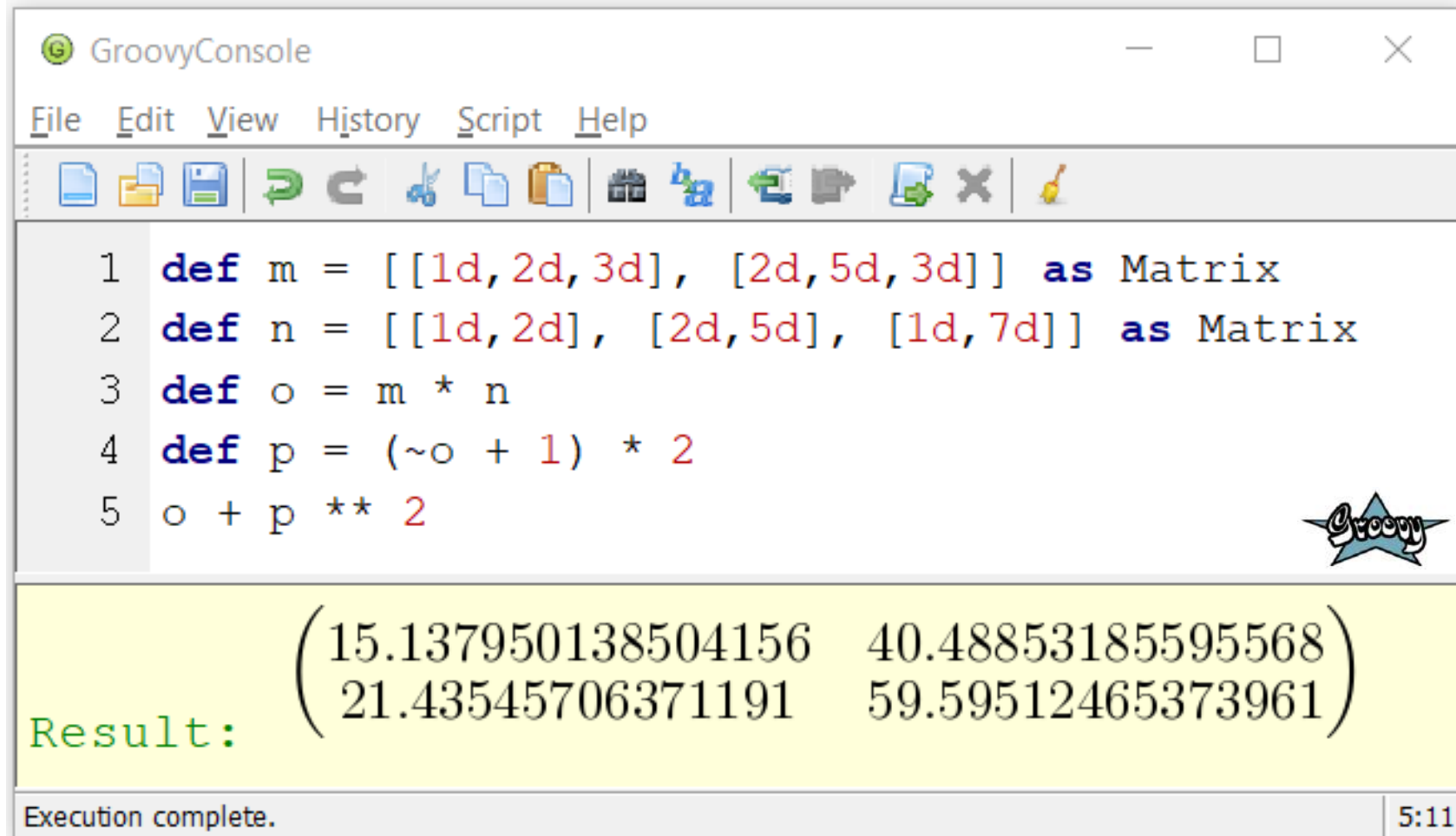
        RealMatrix q = o.add(p.power(2));

        System.out.println(q);
    }
}
```



```
Array2DRowRealMatrix{{15.1379501385,40.488531856},{21.4354570637,59.5951246537}}
```

# Operator overloading and extensible tools



The screenshot shows a GroovyConsole window with a menu bar (File, Edit, View, History, Script, Help) and a toolbar with various icons. The main area contains five lines of Groovy code:

```
1 def m = [[1d, 2d, 3d], [2d, 5d, 3d]] as Matrix
2 def n = [[1d, 2d], [2d, 5d], [1d, 7d]] as Matrix
3 def o = m * n
4 def p = (~o + 1) * 2
5 o + p ** 2
```

Below the code, the result is displayed in a yellow box:

Result:  $\begin{pmatrix} 15.137950138504156 & 40.48853185595568 \\ 21.43545706371191 & 59.59512465373961 \end{pmatrix}$

At the bottom left, it says "Execution complete." and at the bottom right, the time "5:11" is shown.

```
GroovyConsole
File Edit View History Script Help
[Icons]
1 @Grab('org.apache.opennlp:opennlp-tools:1.9.0')
2 import opennlp.tools.langdetect.*
3
4 u = 'http://apache.forsale.plus/opennlp/models/langdetect/1.8.3/langdetect-183.bin'
5 d = new LanguageDetectorME(new LanguageDetectorModel(new URL(u)))
6 a = 'Bienvenue à Paris'
7 b = 'Velkommen til København'
8
9 assert d.predictLanguage(a).lang == d.predictLanguage(b).lang

Exception thrown

Assertion failed:

assert d.predictLanguage(a).lang == d.predictLanguage(b).lang
| |                | | | | |                | |
| |                | 'fra'| | |                | 'dan'
| |                |     | | |                'Velkommen til København'
| |                |     | | dan (0.024727160814654276)
| |                |     | opennlp.tools.langdetect.LanguageDetectorME@232a5780
| |                |     false
| |                'Bienvenue à Paris'
| fra (0.018630393459062138)
opennlp.tools.langdetect.LanguageDetectorME@232a5780

at ConsoleScript12.run(ConsoleScript12:9)

Execution terminated with exception. 9:62
```

showing @Grab, power asserts and natural language processing

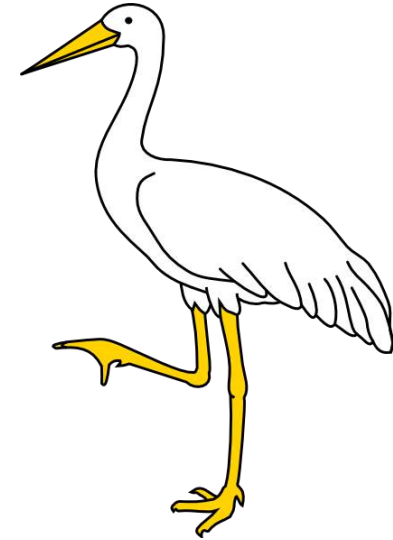
# REPL for Groovy (groovysh)

```
paulk@pop-os: ~  
paulk@pop-os:~$ groovysh  
Groovy Shell (3.0.0-beta-1, JVM: 12.0.1)  
Type ':help' or ':h' for help.  
-----  
groovy:000> :grab com.mitchtalmadge:asciidata:1.4.0  
groovy:000> import com.mitchtalmadge.asciidata.graph.ASCIIGraph as Graph  
==> com.mitchtalmadge.asciidata.graph.ASCIIGraph as Graph  
groovy:000> double[] sinWave() { (0..<60).collect{ Math.sin(it * ((Math.PI * 4) / 60)) } }  
==> true  
groovy:000> println Graph.fromSeries(sinWave()).withNumRows(15).plot()  
0.99  
0.85  
0.71  
0.57  
0.43  
0.28  
0.14  
0.00  
-0.14  
-0.28  
-0.43  
-0.57  
-0.71  
-0.85  
-0.99  
==> null  
groovy:000> █
```

The groovysh REPL is where you can program interactively.

# Command chains and extensible type system

```
cranes have 2 legs  
tortoises have 4 legs  
  
there are 7 animals  
there are 20 legs  
  
display solution
```



```
Cranes 4  
Tortoises 3
```

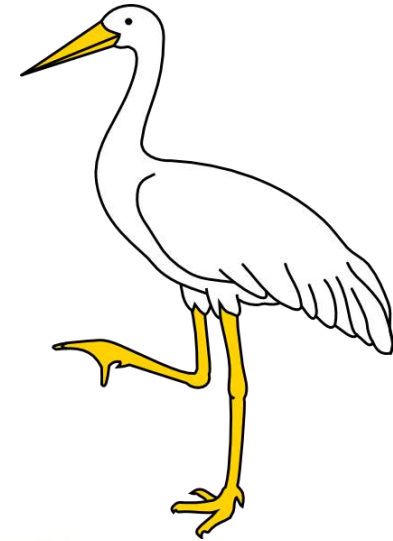
Text file? No it's a statically typed Groovy program which won't compile if the animal names aren't valid.  
A constraint programming solver is used to find the solution.

# Command chains and extensible type system

```
cranes have 2 legs  
tortoises have 4 legs  
millipedes have 1000 legs  
there are 8 animals  
there are 1020 legs
```

```
display solution
```

```
Cranes 4  
Tortoises 3  
Millipedes 1
```





# Groovy frameworks plus entire Java ecosystem



Rich ecosystem.



# Groovy frameworks plus entire Java ecosystem



```
class MathSpec extends Specification {  
    def "maximum of two numbers"() {  
        expect:  
            Math.max(a, b) == c  
  
        where:  
            a | b | c  
            1 | 3 | 3  
            7 | 4 | 7  
            0 | 0 | 0  
    }  
}
```



MICRONAUT



Jenkins

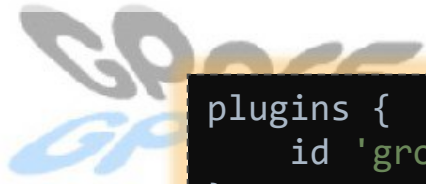
Rich ecosystem.



X

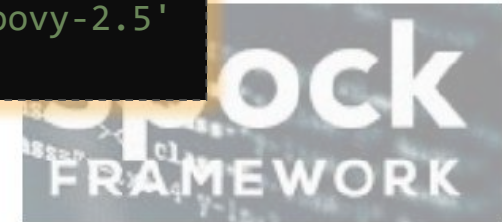


# Groovy frameworks plus entire Java ecosystem



MICRONAUT

```
plugins {  
    id 'groovy'  
}  
  
repositories {  
    jcenter()  
}  
  
dependencies {  
    implementation 'org.codehaus.groovy:groovy-all:2.5.7'  
  
    testImplementation 'org.spockframework:spock-core:1.2-groovy-2.5'  
}
```



Rich ecosystem.

# Groovy frameworks plus entire Java ecosystem



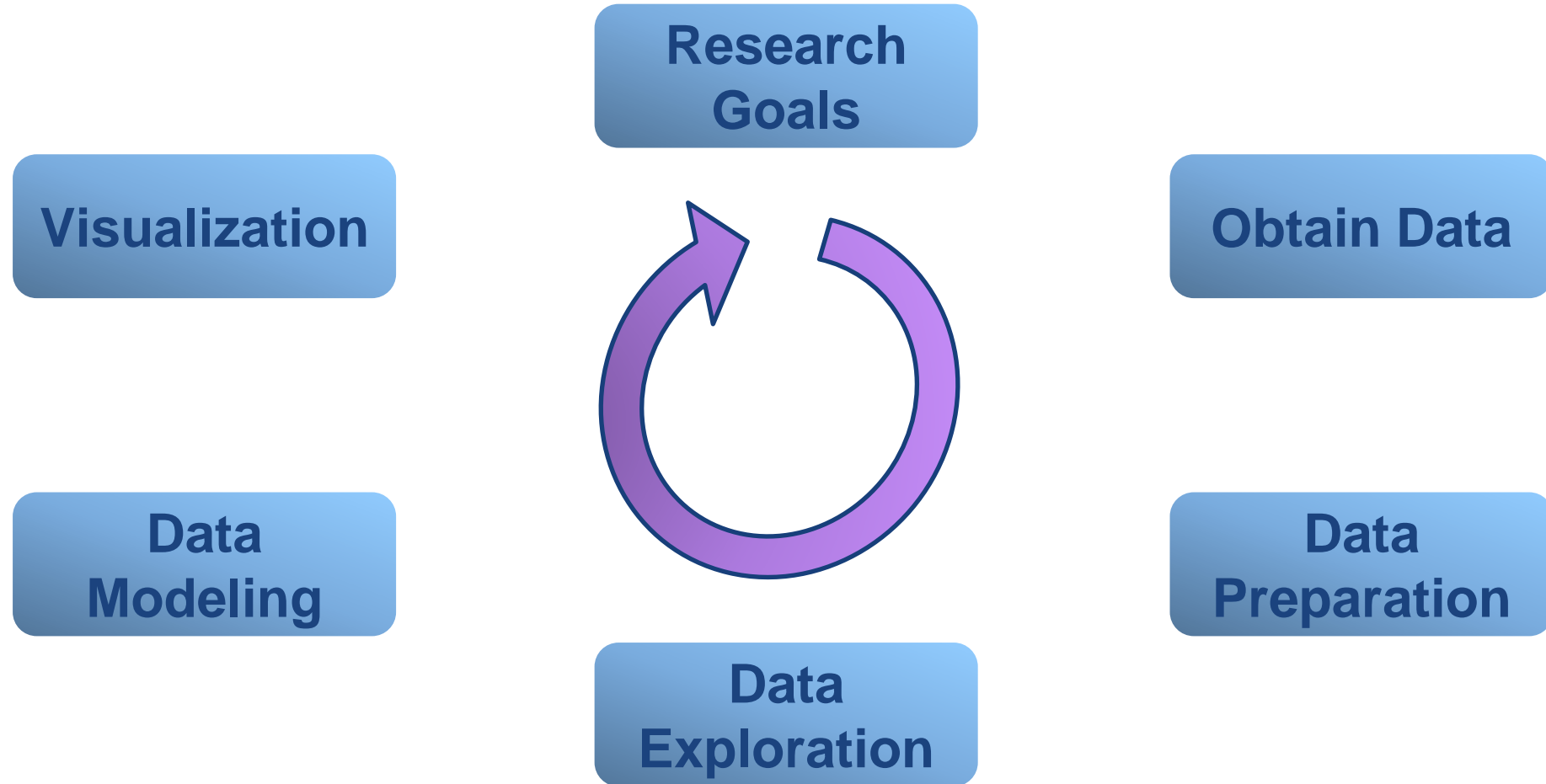
MICRONAUT

<i>Data Parallelism: Fork/Join Map/Reduce</i>	<i>Fixed coordination (for collections)</i>
<i>Actors</i>	<i>Explicit coordination</i>
<i>Safe Agents</i>	<i>Delegated coordination</i>
<i>Dataflow</i>	<i>Implicit coordination</i>

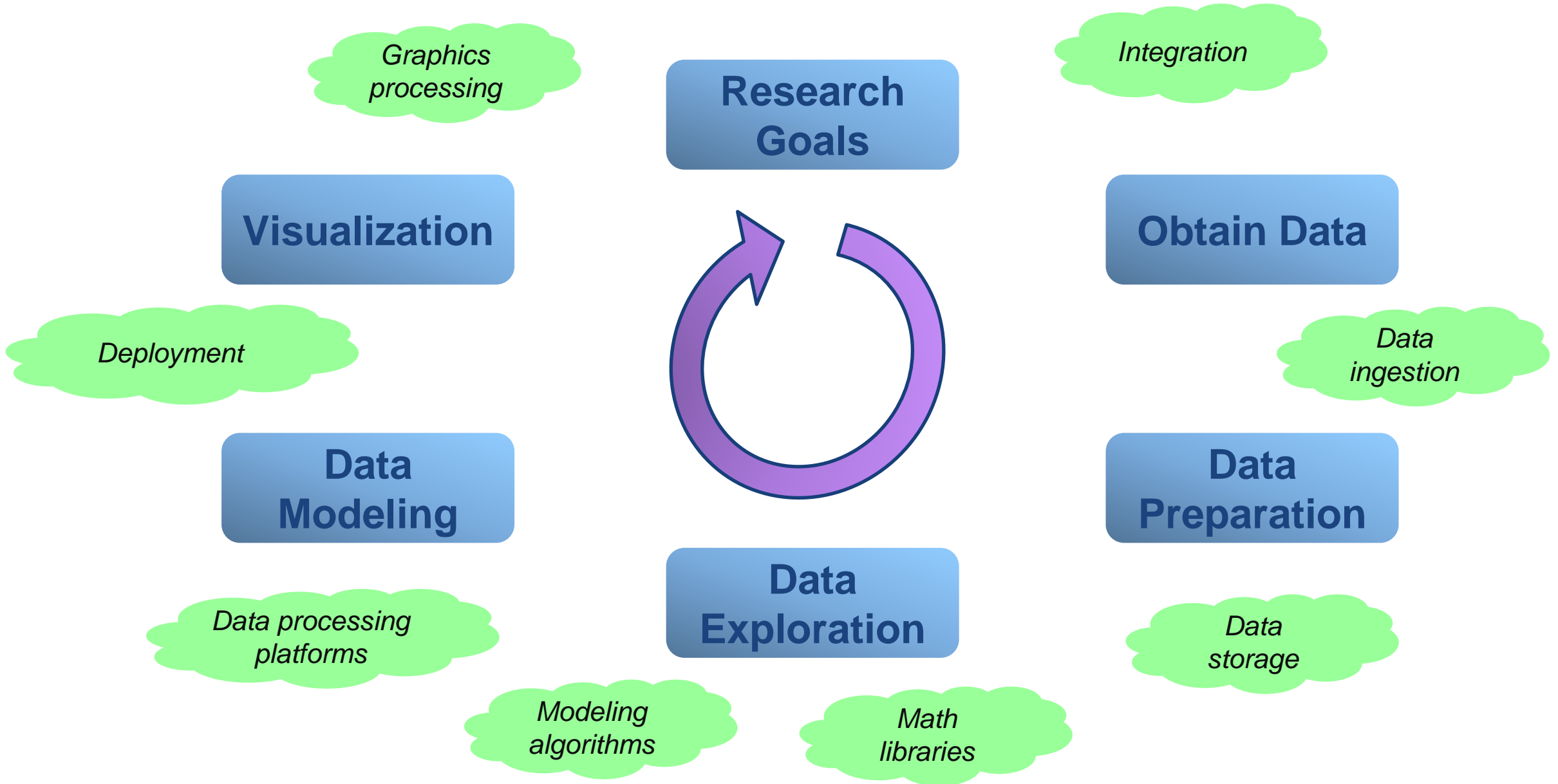
Rich ecosystem.







# Data Science Process





# Data Science Process



# Libraries/Tools/Frameworks for Data Science - R

	COMMITTS	CONTRIBUTORS	FEATURES	
DATA MANIPULATION	 dplyr	4 354	136	<ul style="list-style-type: none"> <li>powerful library for data wrangling</li> <li>works with local data frames and remote database tables</li> <li>precise and simple command syntax</li> </ul>
	data.table	3 211	43	<ul style="list-style-type: none"> <li>quick aggregation of large data</li> <li>laconic flexible syntax and a wide suite of useful functions</li> <li>friendly file reader and parallel file writer</li> </ul>
	lubridate	1 427	45	<ul style="list-style-type: none"> <li>a set of functions to work with date and time format</li> <li>easy and fast parsing of date-time data</li> <li>expanded mathematical operations on time data</li> </ul>
	jsonlite	908	11	<ul style="list-style-type: none"> <li>robust and quick parsing JSON objects in R</li> <li>great tool for interacting with web APIs and building pipelines</li> <li>functions to stream, validate, and prettify JSON data</li> </ul>
GRAPHIC DISPLAYS	 ggplot2	3 903	133	<ul style="list-style-type: none"> <li>powerful implementation of the grammar of graphics visualization</li> <li>developed static graphics system</li> <li>takes care of plot specifications</li> </ul>
	Corrplot	299	8	<ul style="list-style-type: none"> <li>abilities to visualize correlation matrices and confidence intervals</li> <li>contains algorithms to do matrix reordering</li> <li>flexible appearance details settings</li> </ul>
	lattice	132	0	<ul style="list-style-type: none"> <li>high-level visualization system</li> <li>emphasis on multivariate data</li> <li>efficiently copes with nonstandard requirements</li> </ul>
HTML WIDGETS	 plotly	2 989	26	<ul style="list-style-type: none"> <li>rich features and plenty of available charts</li> <li>web-based toolbox for building visualizations</li> <li>abilities to make ggplot2 graphics interactive</li> </ul>
	ggvis	2 159	21	<ul style="list-style-type: none"> <li>implementation of an interactive grammar of graphic</li> <li>incorporates shiny reactive programming model and dplyr grammar of data transformation</li> </ul>
	DT DataTables	1 919	21	<ul style="list-style-type: none"> <li>displays R matrices and data frames as interactive HTML tables</li> <li>creates sortable tables with a minimum of code</li> <li>many useful features and styling options for tables</li> </ul>
	 rCharts	638	11	<ul style="list-style-type: none"> <li>interactive JS charts from R</li> <li>tools for creation, customization, and sharing</li> </ul>

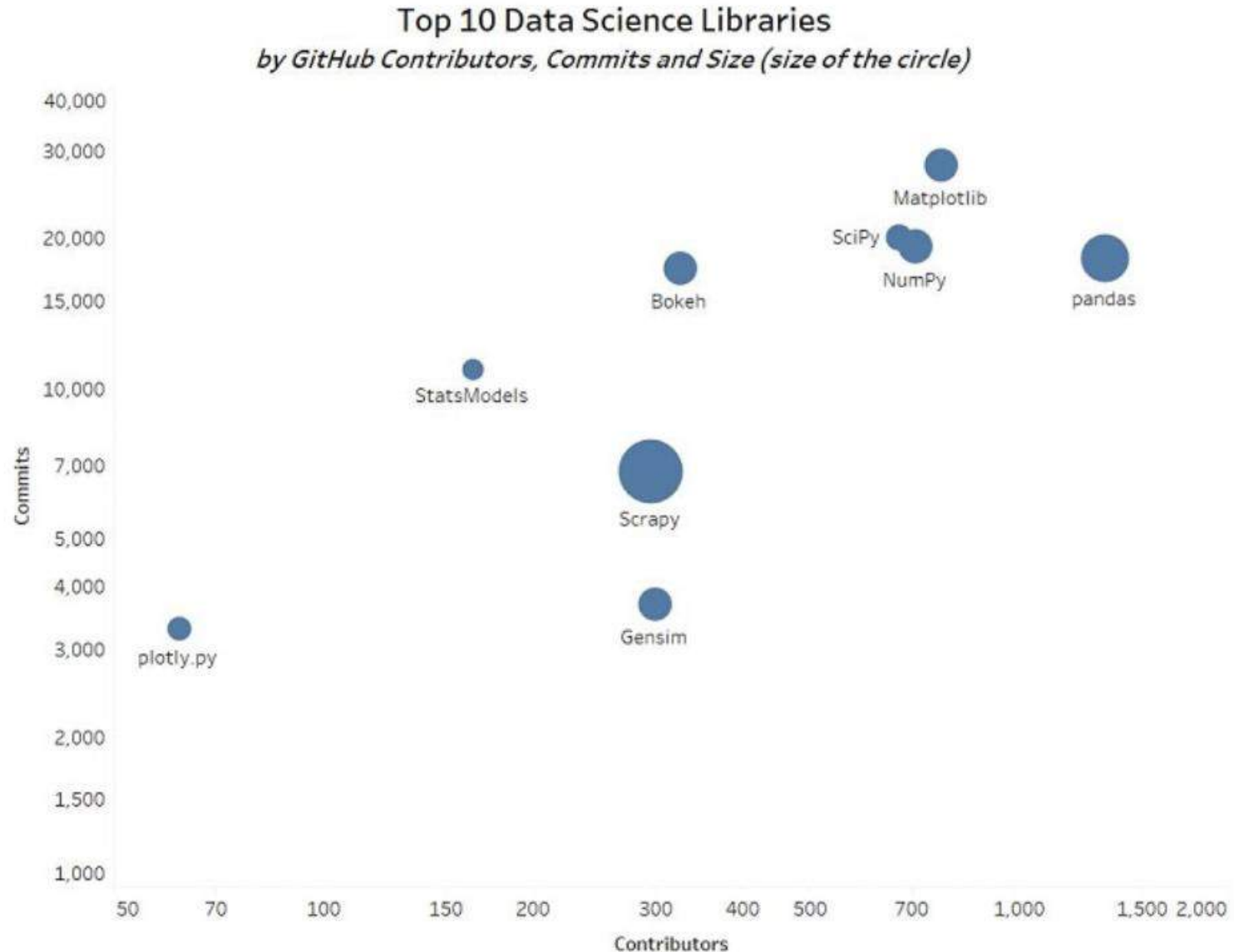
	COMMITTS	CONTRIBUTORS	FEATURES	
REPRODUCIBLE RESEARCH	 knitr	5 467	96	<ul style="list-style-type: none"> <li>transparent tool for easy dynamic report generation in R</li> <li>enables integration of R code into LaTeX, LyX, HTML, Markdown, AsciiDoc, and reStructuredText documents</li> </ul>
	 rmarkdown	2 297	56	<ul style="list-style-type: none"> <li>next generation implementation of R Markdown based on pandoc</li> <li>many static and dynamic output formats</li> <li>abilities to define new formats for custom publishing requirements</li> </ul>
	slidify	302	7	<ul style="list-style-type: none"> <li>generates reproducible html5 slides from r markdown</li> <li>allows embedded code chunks and mathematical formulas</li> <li>rich sharing and customizing opportunities</li> </ul>
MACHINE LEARNING	mlr	3 915	55	<ul style="list-style-type: none"> <li>extensible framework for classification, regression, survival analysis, and clustering</li> <li>easy extension mechanism through S3 inheritance</li> </ul>
	<i>dmlc</i> XGBoost	3 188	259	<ul style="list-style-type: none"> <li>implementation of the Gradient Boosted Decision Trees algorithm</li> <li>reach tools for regression, classification, and ranking problems</li> <li>high speed and performance</li> </ul>
	caret	1 659	59	<ul style="list-style-type: none"> <li>many models for classification and regression</li> <li>powerful tools and algorithms for creating predictive models</li> </ul>
	gbm	731	26	<ul style="list-style-type: none"> <li>represents Generalized Boosted Regression Models</li> <li>includes plenty of regression methods</li> <li>tools variable selection and final stage precision modeling</li> </ul>
	Prophet	190	20	<ul style="list-style-type: none"> <li>high-quality forecasts for time series data</li> <li>manages data that has multiple seasonality with linear or non-linear growth</li> <li>robust to missing data, shifts in the trend, and large outliers</li> </ul>
	randomForest	56	0	<ul style="list-style-type: none"> <li>implements Breiman's random forest algorithm for classification and regression</li> <li>builds multiple decision trees and gives back the mean prediction of the individual trees</li> </ul>

Updated: December 2017

Created by  ActiveWizards

Source: <https://www.kdnuggets.com/2018/05/top-20-r-libraries-data-science-2018.html>

# Libraries/Tools/Frameworks for Data Science - Python



Source: <https://www.kdnuggets.com/2018/11/top-python-data-science-libraries.html>



# Obtain data

## Identify sources

- Internal/External (data providers)
- Databases, files, events, IoT
- Unstructured/structured
- Text, XML, JSON, YAML, CSV/spreadsheets, audio, images, video, web services
- Apache Tika
- JExcel
- Apache POI

## Clarify ownership

## Check quality

```
def jsonSlurper = new JsonSlurper()
def object = jsonSlurper.parseText('{ "myList": [4, 8, 15, 16, 23, 42] }')
assert object instanceof Map
assert object.myList instanceof List
assert object.myList == [4, 8, 15, 16, 23, 42]
```

```
def response = new XmlSlurper().parseText(books)
def authorResult = response.value.books.book[0].author
assert authorResult.text() == 'Miguel de Cervantes'
```

```
def qry = 'SELECT * FROM Author'
assert sql.rows(qry, 1, 3)*.firstname == ['Dierk', 'Paul', 'Guillaume']
assert sql.rows(qry, 4, 3)*.firstname == ['Hamlet', 'Cedric', 'Erik']
assert sql.rows(qry, 7, 3)*.firstname == ['Jon']
```

# Data preparation

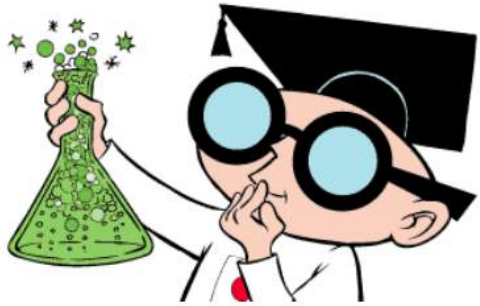
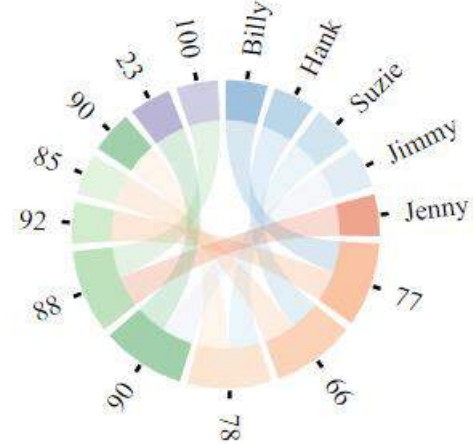
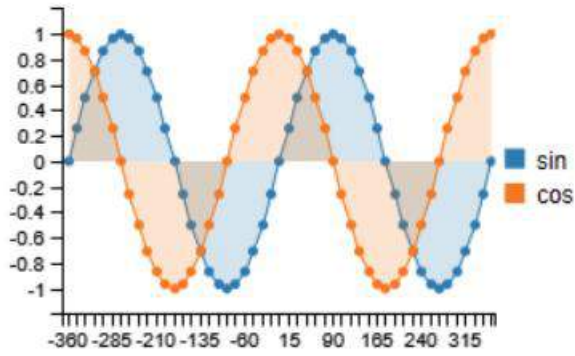
## Collections:

- Java: lists, maps, sets
- Groovy: literal notation, GDK methods, GPath
- Libraries
  - Google Guava <https://github.com/google/guava>
  - Apache Common Collections <https://commons.apache.org/collections/>

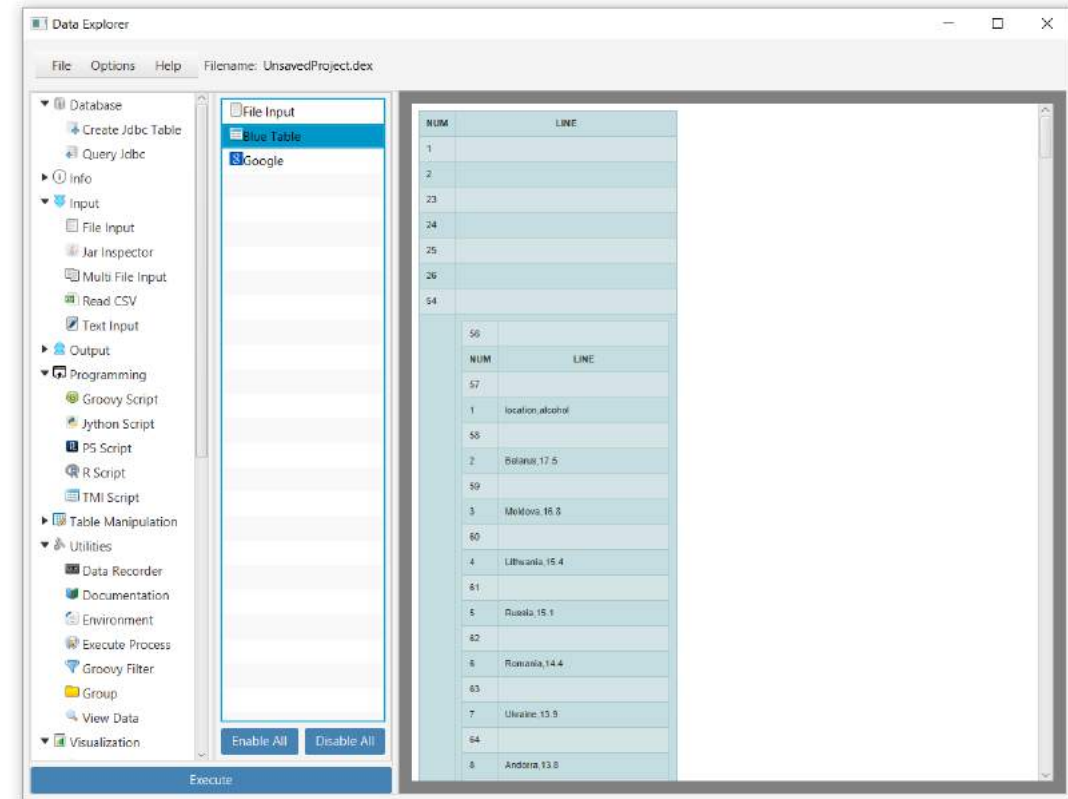
## DataFrames:

- Joinery <https://cardillo.github.io/joinery/>  
A data frame implementation in the spirit of Pandas or R data frames with show/plot
- Tablesaw <https://jtablesaw.github.io/tablesaw/>  
Java dataframe and visualization library
- Apache Spark DataFrames <https://spark.apache.org/docs/latest/sql-programming-guide.html>  
Equivalent to a table in a relational database or a data frame in R/Python
- Paleo <https://github.com/netzwerg/paleo>  
Immutable Java 8 data frames with typed columns (including primitives)

# Data exploration - DEX



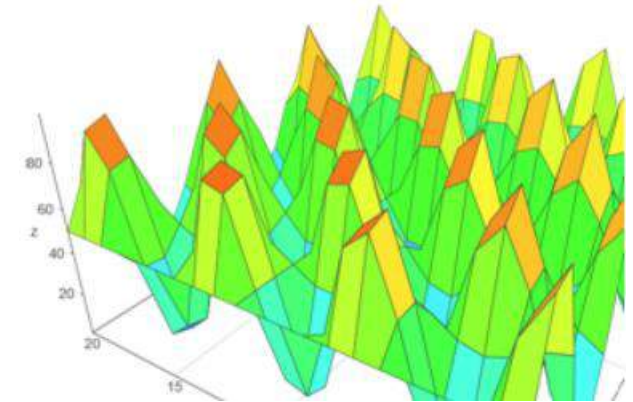
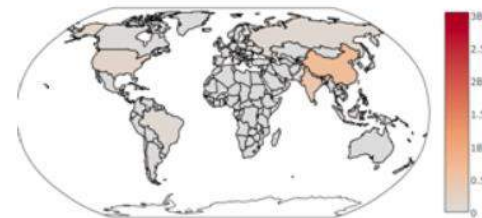
## DEX : The Data Explorer



## Dex

Dex : The data explorer is a powerful tool for data science. It is written in Groovy and Java on top of JavaFX and offers the ability to:

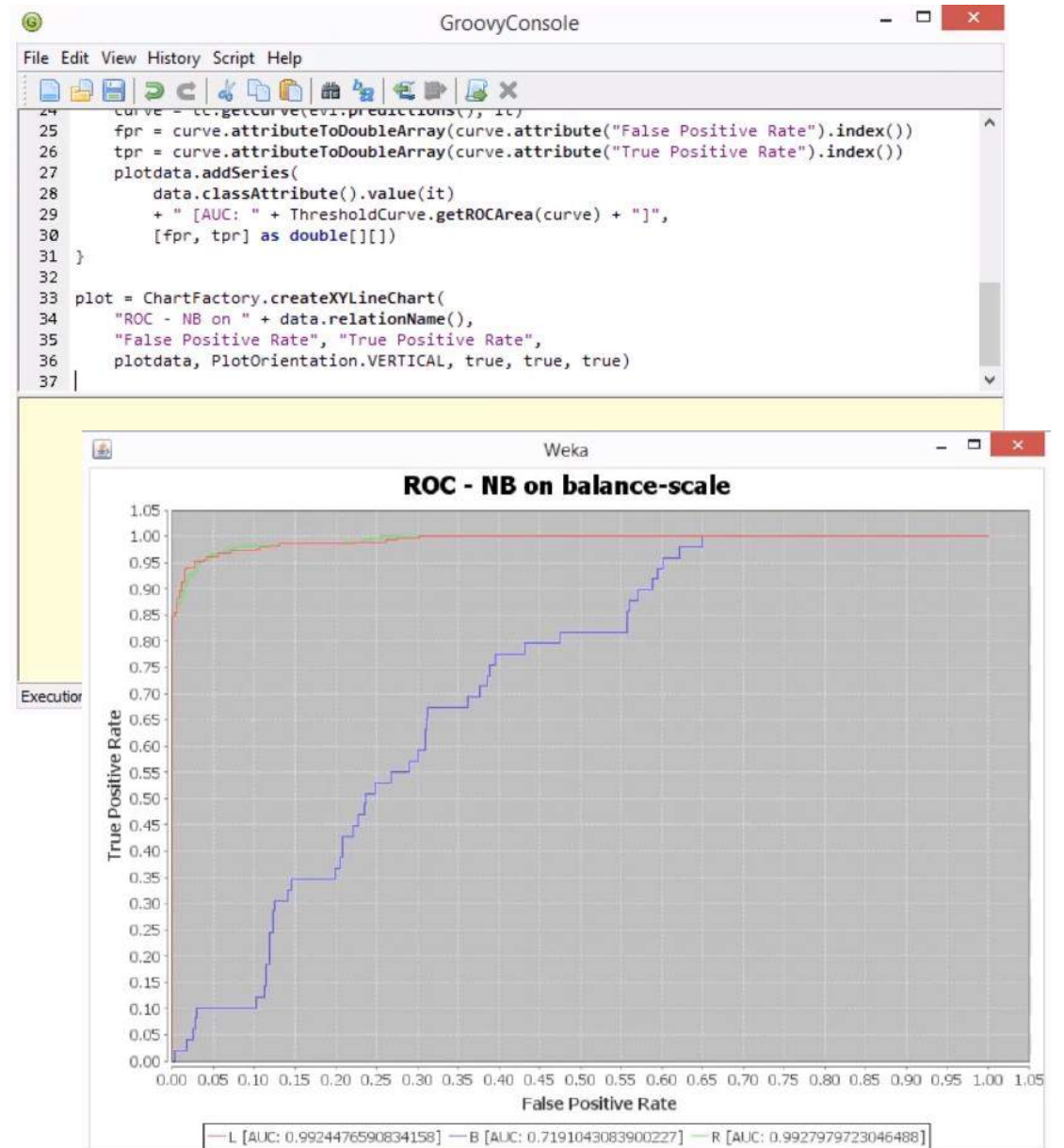
- Read in data from a variety of sources such as files, programs and a variety of databases.
- Transform the data in a powerful ways.
- Apply powerful machine learning to the data via SMILE and R integration.
- Visualize the data in over 50 distinct ways.
- Output the data to a variety of databases and file formats.
- Extend Dex from within via templates and internal scripting. Save the extensions to reuse later.



# Data exploration - Weka

The screenshot shows the Weka Explorer GUI. The 'Current relation' is 'autos' with 205 instances. The 'Selected attribute' is 'normalized-losses', which is a numeric attribute with 51 distinct values and a sum of weights of 205. A histogram for this attribute is displayed at the bottom right, showing a distribution of values from 65 to 256. The 'Attributes' list on the left includes 'normalized-losses', 'make', 'fuel-type', 'aspiration', 'num-of-doors', 'body-style', 'drive-wheels', 'engine-location', and 'wheel-base'.

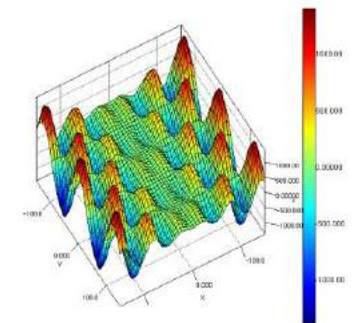
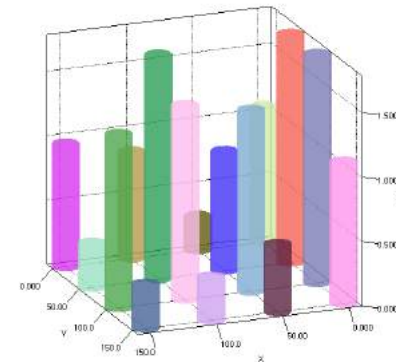
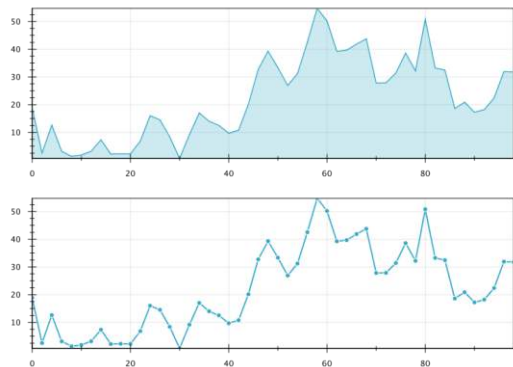
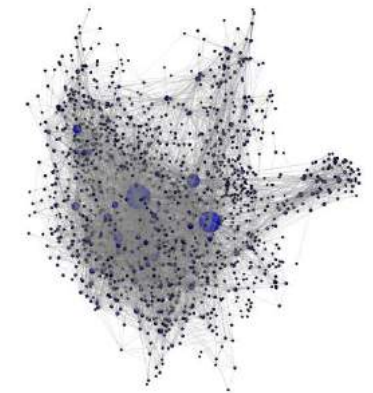
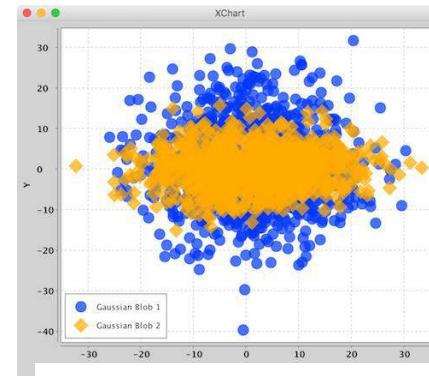
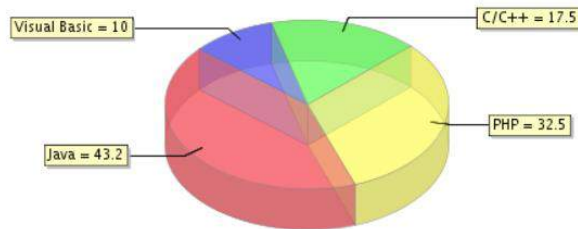
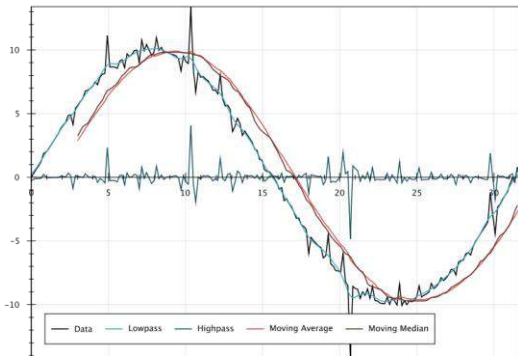
The screenshot shows the Weka GUI Chooser window. The 'Tools' menu is open, and 'Groovy console' is selected, which is associated with the keyboard shortcut Ctrl+G. Other options in the menu include Package manager (Ctrl+U), ArffViewer (Ctrl+A), SqlViewer (Ctrl+S), Bayes net editor (Ctrl+N), and Jython console (Ctrl+J). The 'Applications' section on the right includes Explorer, Experimenter, KnowledgeFlow, Workbench, and Simple CLI.



More info: <https://www.youtube.com/watch?v=7quZv6WCTQc>

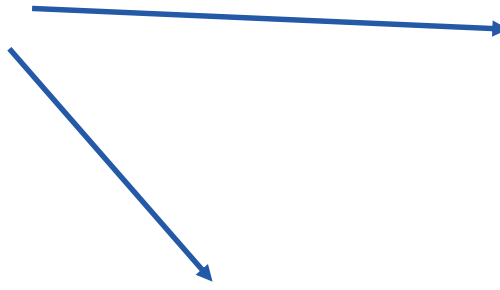
# Visualization

- Open-Source plotting libraries: GRAL, JFreeChart, Xchart, JMathPlot, Jzy3d, JavaFX, GroovyFX
- Types: Scatter plot, Line plot, Area plot, Pie plot, Donut plot, Horizontal bars, Vertical bars, Bubble plot, Radar plot, Box plot, Raster plot, Contour plot, Gauge plot, Step plot, Gantt plot, Histogram
- Features: Multiple axes, Logarithmic axes, Line smoothing, Data filtering/aggregation



# Notebooks supporting Groovy

- Jupyter/beakerx
- Apache Zeppelin
- GroovyLab
- Seco



The screenshot shows a Jupyter notebook interface with the following content:

```
In [1]: // generate a random walk
steps = 10000
random = new Random()
def walk(ssize) {
    def cur = 0.0
    def x = []
    def y = []
    for (i in 0..steps) {
        y[i] = cur
        x[i] = 1
        cur += random.nextGaussian() * ssize
    }
    return [x:x, y:y]
}
// now x is time for the x-axis, and y is the random variable
beer = walk(10)
whiskey = walk(100)
OutputCell.HIDDEN

In [2]: p = new TimePlot(title: "Drunken Sailor Walks", showLegend: true, lodThreshold: 2000)
p << new Line(x:beer.x, y:beer.y, displayName:"Beer Walk")
p << new Line(x:whiskey.x, y:whiskey.y, displayName:"Whiskey Walk")
```

The notebook displays two plots:

- Bar Chart Demo:** A bar chart titled "Bar Chart Demo" showing "Temperature ° Celcius" on the y-axis (0 to 30) and "Alkali" on the x-axis (Lithium, Sodium, Potassium, Rubidium). The chart displays two series: "Solid" (orange bars) and "Liquid" (yellow bars). A legend on the right shows "All", "Solid", and "Liquid" with checkboxes.
- Drunken Sailor Walks:** A line plot titled "Drunken Sailor Walks" showing two random walks over time (00:00.000 to 00:10.000). The y-axis ranges from -10,000 to 5,000. The "Beer Walk" is a blue line that stays near zero, while the "Whiskey Walk" is an orange line that fluctuates significantly between -10,000 and 5,000. A legend on the right shows "All", "Beer Walk", and "Whiskey Walk" with checkboxes.

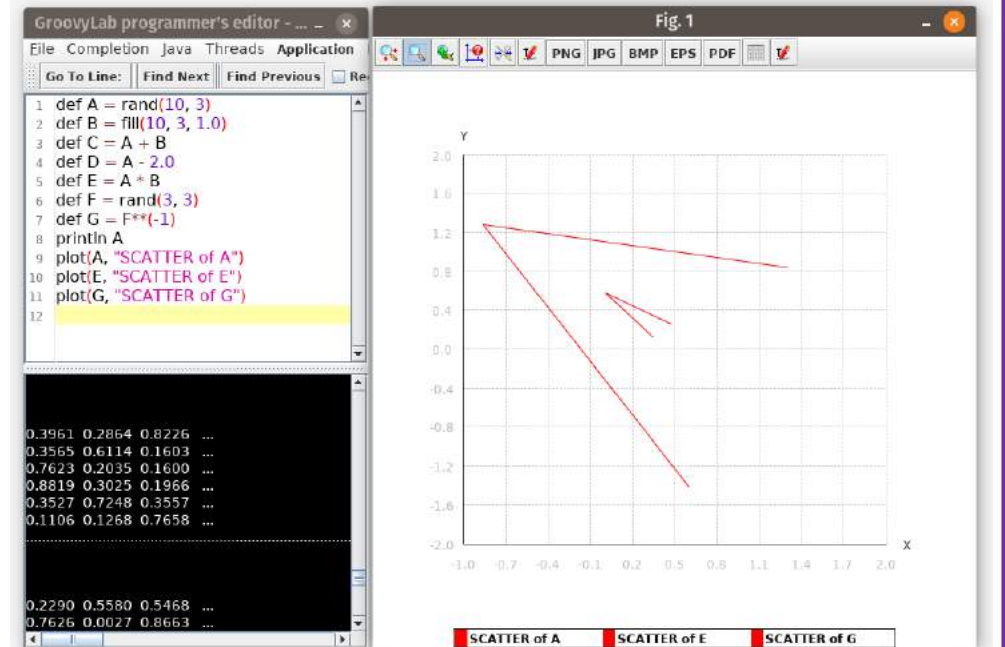
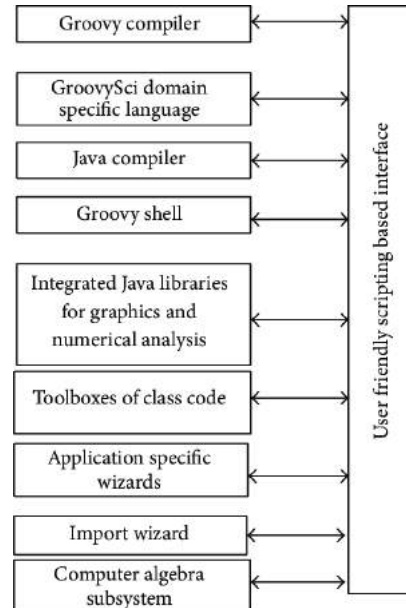
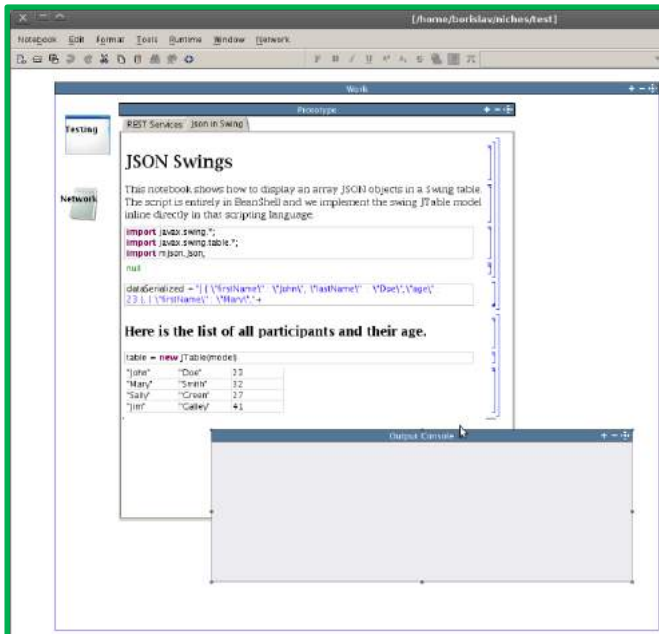
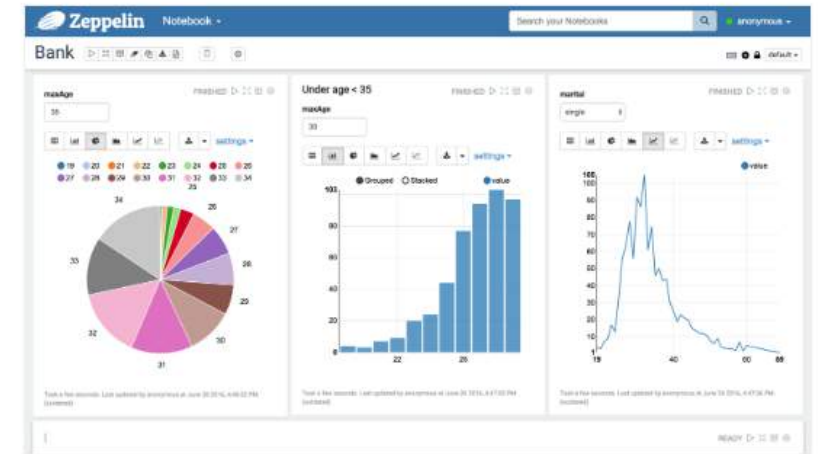
# Notebooks supporting Groovy

- Jupyter/beakerx
- Apache Zeppelin
- GroovyLab
- Seco

## What is Apache Zeppelin?

Multi-purpose notebook which supports 20+ language backends

- 🏭 Data Ingestion
- 👁 Data Discovery
- 🔧 Data Analytics
- 👤 Data Visualization & Collaboration



# Data preparation

## IO:

- Java in-built
- Groovy enhancements
- Libraries
  - Apache Commons IO <https://commons.apache.org/io/>
  - AOL Cyclops-React <https://github.com/aol/cyclops-react>

## Web scraping:

- GEB

## Batch processing/ETL:

- Spring Batch <http://projects.spring.io/spring-batch/>
- Apache Camel
- Apache Gobblin



# Math libraries

Math and stats libraries

- Java: log, exp, other basic methods.

Libraries:

- Apache Commons Math <http://commons.apache.org/math/> statistics, optimization, and linear algebra
- Apache Mahout <http://mahout.apache.org/> linear algebra, plus distributed linear algebra and machine learning
- JBlas <http://jblas.org/> optimized and very fast linear algebra package that uses the BLAS library

Also, many of the machine learning libraries come with some math functionality, often linear algebra, stats, and optimization.

# Data storage

## Database:

- Java's JDBC
- Groovy JDBC enhancements and Datasets
- Other NOSQL
  - Apache Cassandra
  - MongoDB
  - Apache CouchDB
  - Neo4j

```
import com.gmongo.GMongo
```

```
def db = new GMongo().getDB('athletes')
db.athletes.drop()
db.athletes << [
    first: 'Paul',
    last: 'Tergat', dob: '1969-06-17',
    runs: [[distance: 42195, time: 2 * 60 * 60 + 4 * 60 + 55,
            venue: 'Berlin', when: '2003-09-28']]
]
```

```
println "World records following $marathon3.venue $marathon3.when:"
def t = new Traversal()
for (Path p in t.description().breadthFirst().
    relationships(MyRelationshipTypes.supercedes).
    evaluator(Evaluators.fromDepth(1)).
    uniqueness(Uniqueness.NONE).
    traverse(marathon3)) {
    def newRecord = p.endNode()
    println "$newRecord.venue $newRecord.when"
}
```

```
Graph g = new Neo4jGraph(graphDb)
def pretty = { it.collect { "$it.venue $it.when" }.join(', ' ) }
def results = []
g.V('venue', 'London').fill(results)
println 'London world records: ' + pretty(results)
```

```
results = []
g.V('venue', 'London').in('supercedes').fill(results)
println 'World records after London: ' + pretty(results)
```

# Machine learning and data mining

Machine learning and data mining libraries:

- Weka <http://www.cs.waikato.ac.nz/ml/weka/>
- JavaML <http://java-ml.sourceforge.net/>  
old and reliable ML library but not actively updated
- Smile <http://haifengl.github.io/smile/>  
ML library under active development
- JSAT <https://github.com/EdwardRaff/JSAT>  
contains many machine learning algorithms
- H2O <http://www.h2o.ai/>  
distributed Java ML framework also supports Scala, R and Python
- Apache Mahout <http://mahout.apache.org/>  
distributed linear algebra framework
- Apache Spark ml and mllib  
<https://github.com/apache/spark/tree/master/examples/src/main/java/org/apache/spark/examples>

# Neural networks and text processing

## Neural networks:

- Encog <http://www.heatonresearch.com/encog/>
- DeepLearning4j <http://deeplearning4j.org/>  
natively supports Keras models

## Text processing

- Java: StringTokenizer, the java.text package, regular expressions
- Groovy: Regex enhancements, templates, String enhancements
- Libraries:
  - Apache Lucene <https://lucene.apache.org/>
  - Stanford CoreNLP <http://stanfordnlp.github.io/CoreNLP/>
  - Apache OpenNLP <https://opennlp.apache.org/>
  - LingPipe <http://alias-i.com/lingpipe/>
  - GATE <https://gate.ac.uk/>
  - MALLET <http://mallet.cs.umass.edu/>
  - Smile <http://haifengl.github.io/smile/>

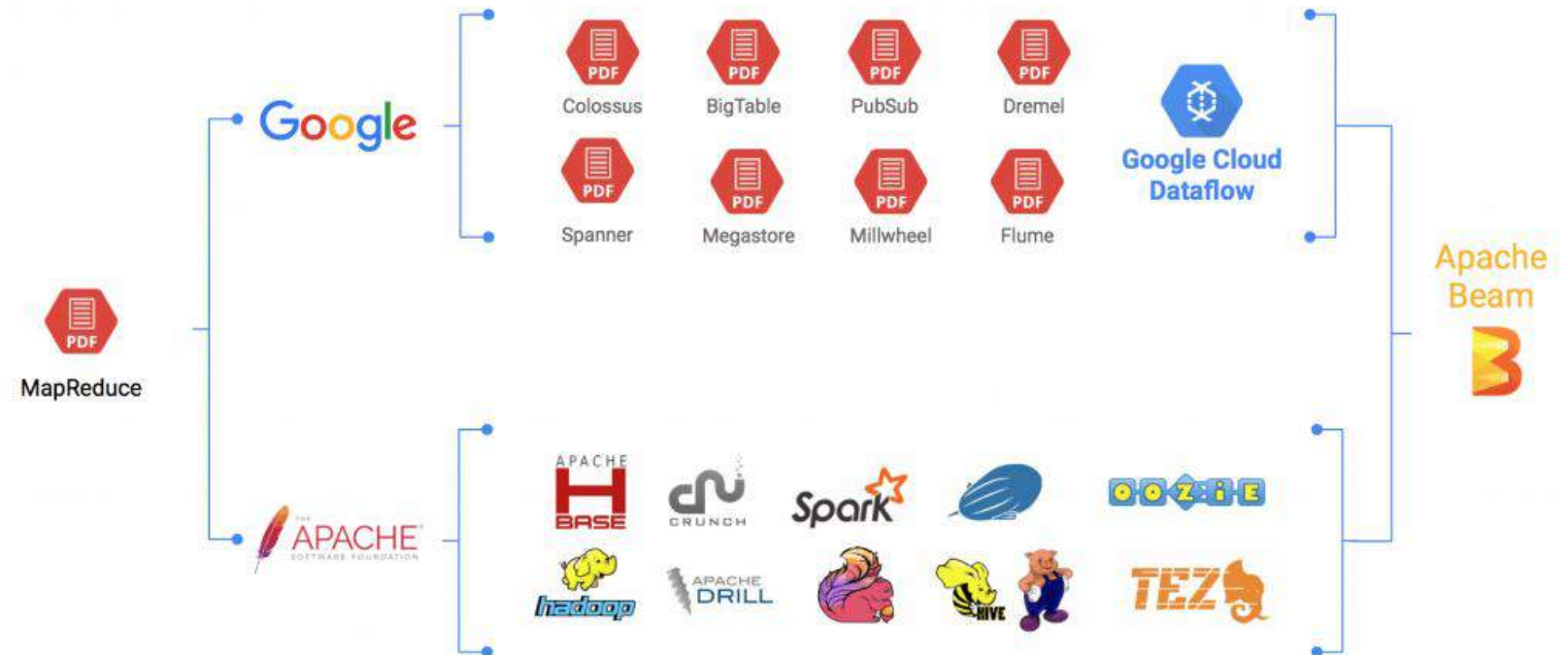
# Scaling up

Microservice frameworks:

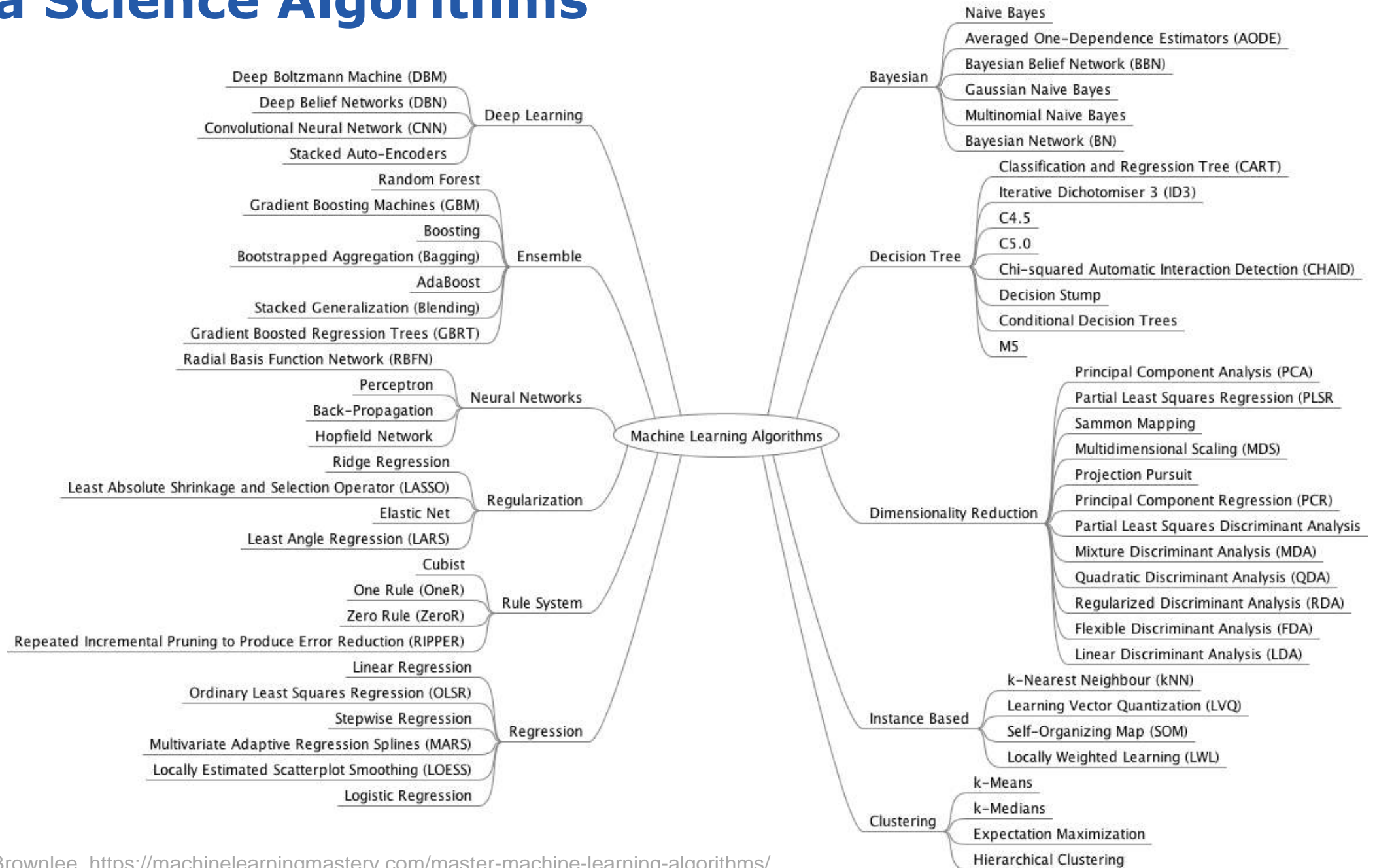
- Micronaut

Distributed data processing:

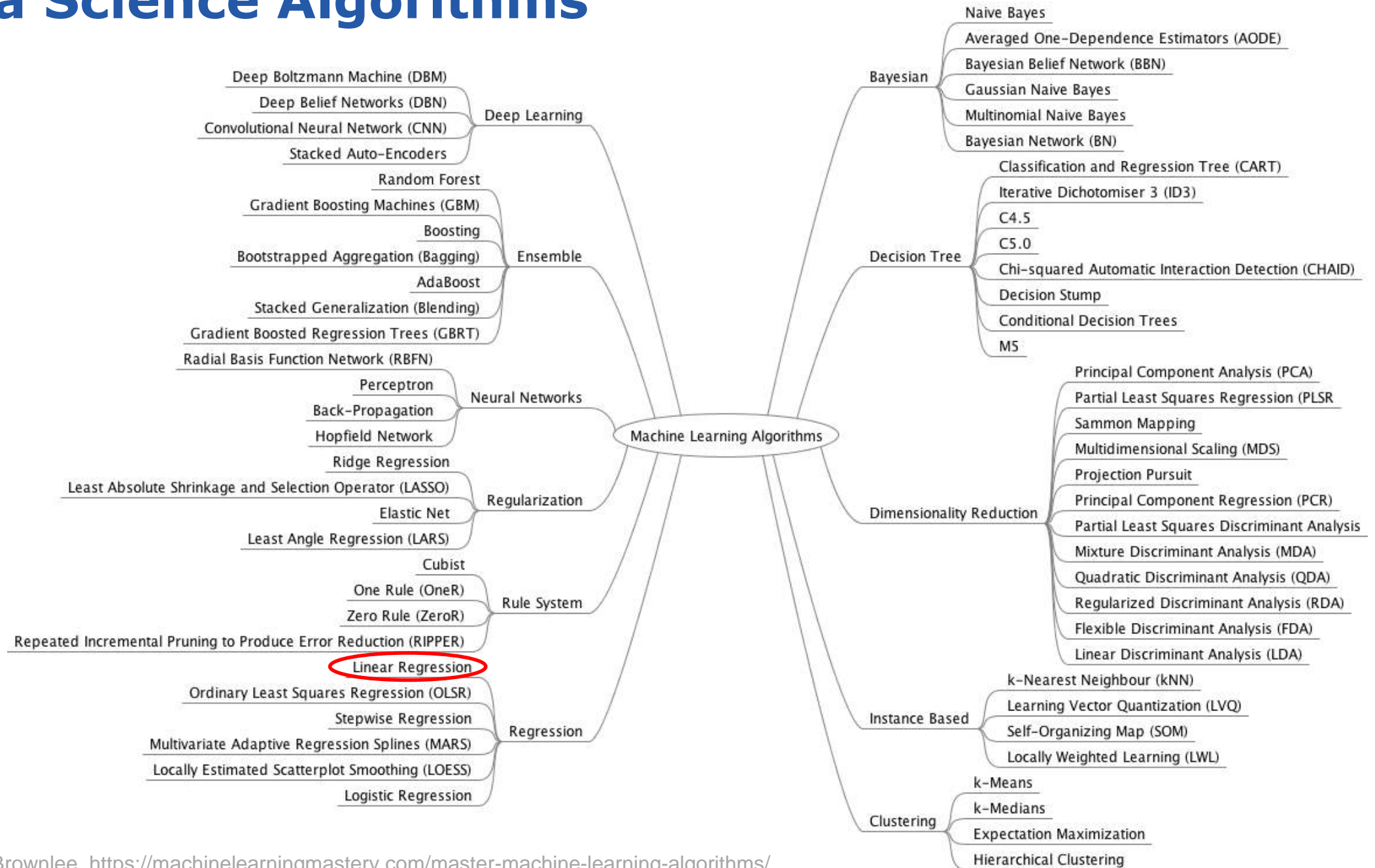
- Apache Hadoop
- Apache Spark
- Apache Flink
- Apache Samza
- Apache Kafka
- Apache Storm
- Apache Apex
- Apache Beam
- Apache Ignite
- Apache Nifi



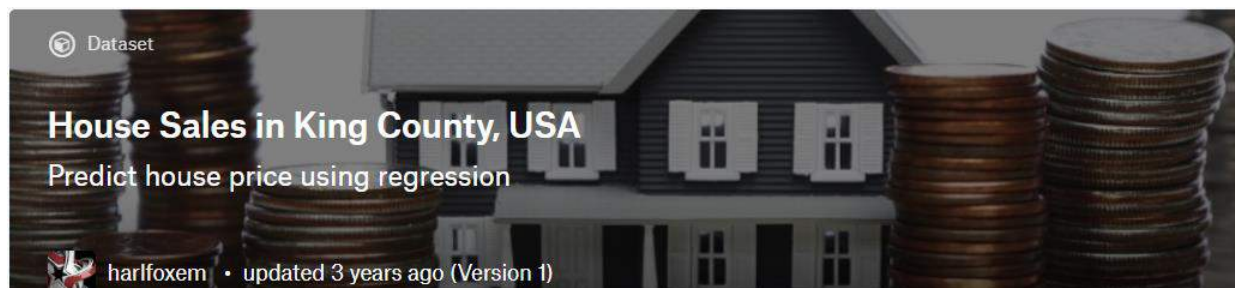
# Data Science Algorithms



# Data Science Algorithms




# House price predictions



 **Usability** 8.8

 **License** CC0: Public Domain

 **Tags** finance, home

## Description


This dataset contains house sale prices for King County, which includes Seattle. It includes homes sold between May 2014 and May 2015.

It's a great dataset for evaluating simple regression models.

Data (778 KB)



## Data Sources

 kc\_hous... 21.6k x 21

## About this file

19 house features plus the price and the id columns, along with 21613 observations.

Dataset: <https://www.kaggle.com/harlfoxem/housesalesprediction/>

See also: <https://nicolai92.github.io/posts/pentaho-weka-prediction>

## Columns

# id a notation for a house

A date Date house was sold

# price Price is prediction target

# bedrooms Number of Bedrooms/House

# bathrooms Number of bathrooms/House

# sqft\_living square footage of the home

# sqft\_lot square footage of the lot

# floors Total floors (levels) in house

A waterfront House which has a view to a waterfront

A view Has been viewed

A condition How good the condition is ( Overall )

A grade overall grade given to the housing unit, based on King County grading system

# sqft\_above square footage of house apart from basement

# sqft\_basement square footage of the basement

# yr\_built Built Year

# yr\_renovated Year when house was renovated

# zipcode zip

# lat Latitude coordinate

# long Longitude coordinate

# sqft\_living15 Living room area in 2015(implies-- some renovations)  
This might or might not have affected the lotsize area

# sqft\_lot15 lotSize area in 2015(implies-- some renovations)



# House price predictions – view data graphically



Commons CSV™

commons™

[Math]

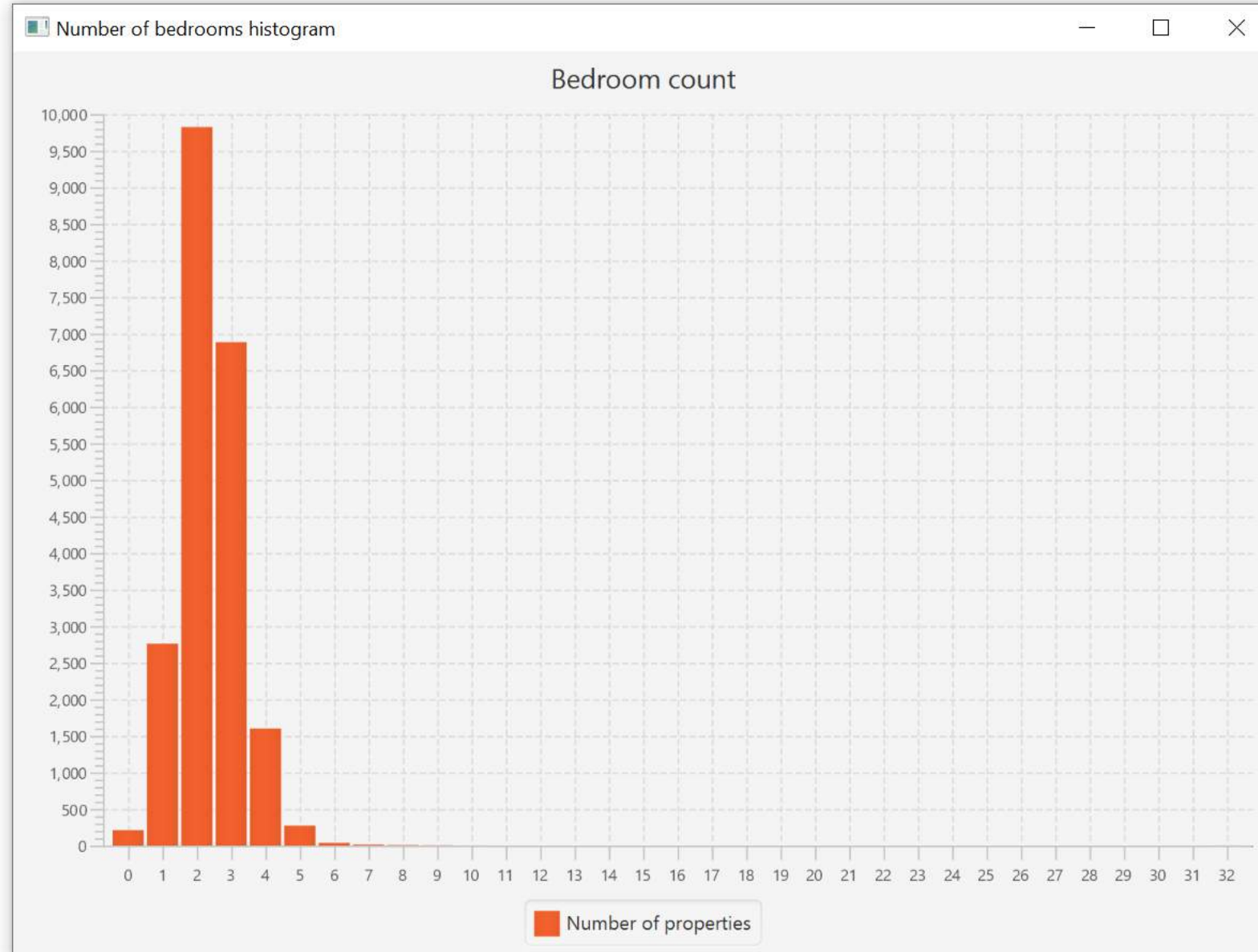


```
import org.apache.commons.math3.random.EmpiricalDistribution
import static groovyx.javafx.GroovyFX.start
import static org.apache.commons.csv.CSVFormat.RFC4180 as CSV

def file = 'kc_house_data.csv' as File
def csv = CSV.withFirstRecordAsHeader().parse(new FileReader(file))
def all = csv.collect { it.bedrooms.toInteger() }
def dist = new EmpiricalDistribution(all.max()).tap{ load(all as double[]) }
def bins = dist.binStats.withIndex().collectMany { v, i -> [i.toString(), v.n] }

start {
    stage(title: 'Number of bedrooms histogram', show: true, width: 800, height: 600) {
        scene {
            barChart(title: 'Bedroom count', barGap: 0, categoryGap: 2) {
                series(name: 'Number of properties', data: bins)
            }
        }
    }
}
```

# House price predictions – view data graphically



# House price predictions – view stats

```
import org.apache.commons.math3.stat.descriptive.SummaryStatistics
import static org.apache.commons.csv.CSVFormat.RFC4180 as CSV

def file = 'kc_house_data.csv' as File
def csv = CSV.withFirstRecordAsHeader().parse(new FileReader(file))
def all = csv.collect { it.bedrooms.toInteger() }

def stats = new SummaryStatistics()
all.each { stats.addValue(it as double) }
println stats.summary
```

```
n: 21613
min: 0.0
max: 33.0
mean: 3.370841623097218
std dev: 0.930061831147451
variance: 0.8650150097573497
sum: 72854.0
```

# House price predictions – investigate outliers

```
import org.apache.commons.csv.CSVFormat.RFC4180 as CSV

def file = 'kc_house_data.csv' as File
def csv = CSV.withFirstRecordAsHeader().parse(new FileReader(file))
csv.findAll{ it.bedrooms.toInteger() > 10 }.each{ println it.toMap() as TreeMap }
```

```
[bathrooms:3, bedrooms:11, condition:3, date:20140821T000000, floors:2, grade:7, id:1773100755, lat:47.556,
long:-122.363, price:520000, sqft_above:2400, sqft_basement:600, sqft_living:3000, sqft_living15:1420,
sqft_lot:4960, sqft_lot15:4960, view:0, waterfront:0, yr_built:1918, yr_renovated:1999, zipcode:98106]
[bathrooms:1.75, bedrooms:33, condition:5, date:20140625T000000, floors:1, grade:7, id:2402100895, lat:47.6878,
long:-122.331, price:640000, sqft_above:1040, sqft_basement:580, sqft_living:1620, sqft_living15:1330,
sqft_lot:6000, sqft_lot15:4700, view:0, waterfront:0, yr_built:1947, yr_renovated:0, zipcode:98103]
```

# House price predictions – investigate outliers

opencsv

```
import com.opencsv.bean.*
import groovy.transform.ToString

@ToString(includeNames = true)
class House {
    @CsvBindByName
    Integer bedrooms

    @CsvBindByName
    String bathrooms

    @CsvBindByName(column = 'sqft_lot')
    Integer area_lot
}

def file      = 'kc_house_data.csv' as File
def builder  = new CsvToBeanBuilder(new FileReader(file))
def records  = builder.withType(House).build().parse()

records.findAll{ it.bedrooms > 10 }.each{ println it }
```

```
House(bedrooms:11, bathrooms:3, area_lot:4960)
House(bedrooms:33, bathrooms:1.75, area_lot:6000)
```

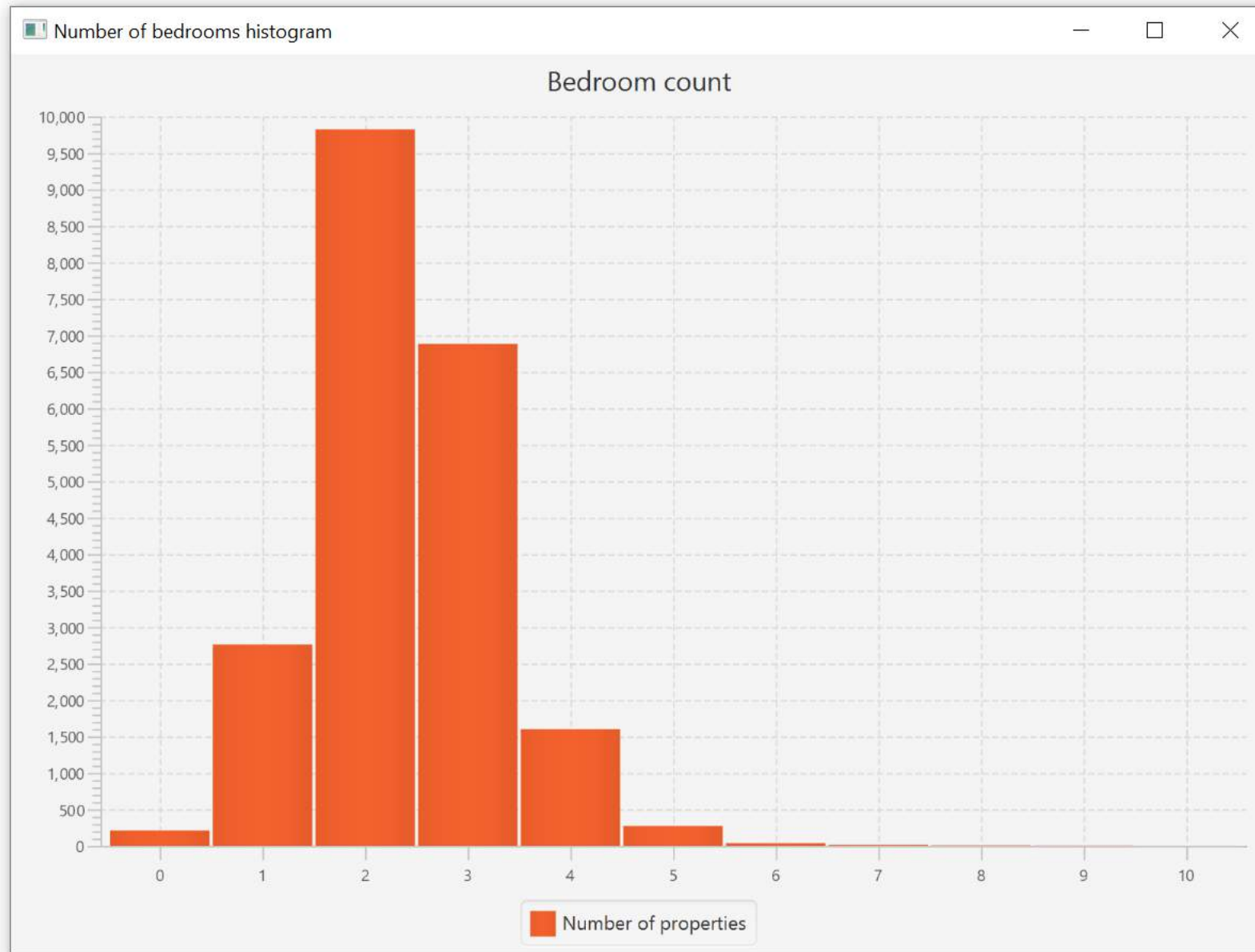
# House price predictions – remove outlier and view

```
import org.apache.commons.math3.random.EmpiricalDistribution
import static groovyx.javafx.GroovyFX.start
import static org.apache.commons.csv.CSVFormat.RFC4180 as CSV

def file = 'kc_house_data.csv' as File
def csv = CSV.withFirstRecordAsHeader().parse(new FileReader(file))
def all = csv.collect { it.bedrooms.toInteger() }.findAll { it < 30 }
def dist = new EmpiricalDistribution(all.max()).tap { load(all as double[]) }
def bins = dist.binStats.withIndex().collectMany { v, i -> [i.toString(), v.n] }

start {
    stage(title: 'Number of bedrooms histogram', show: true, width: 800, height: 600) {
        scene {
            barChart(title: 'Bedroom count', barGap: 0, categoryGap: 2) {
                series(name: 'Number of properties', data: bins)
            }
        }
    }
}
```

# House price predictions – remove outlier and view



# House price predictions – remove outlier and view

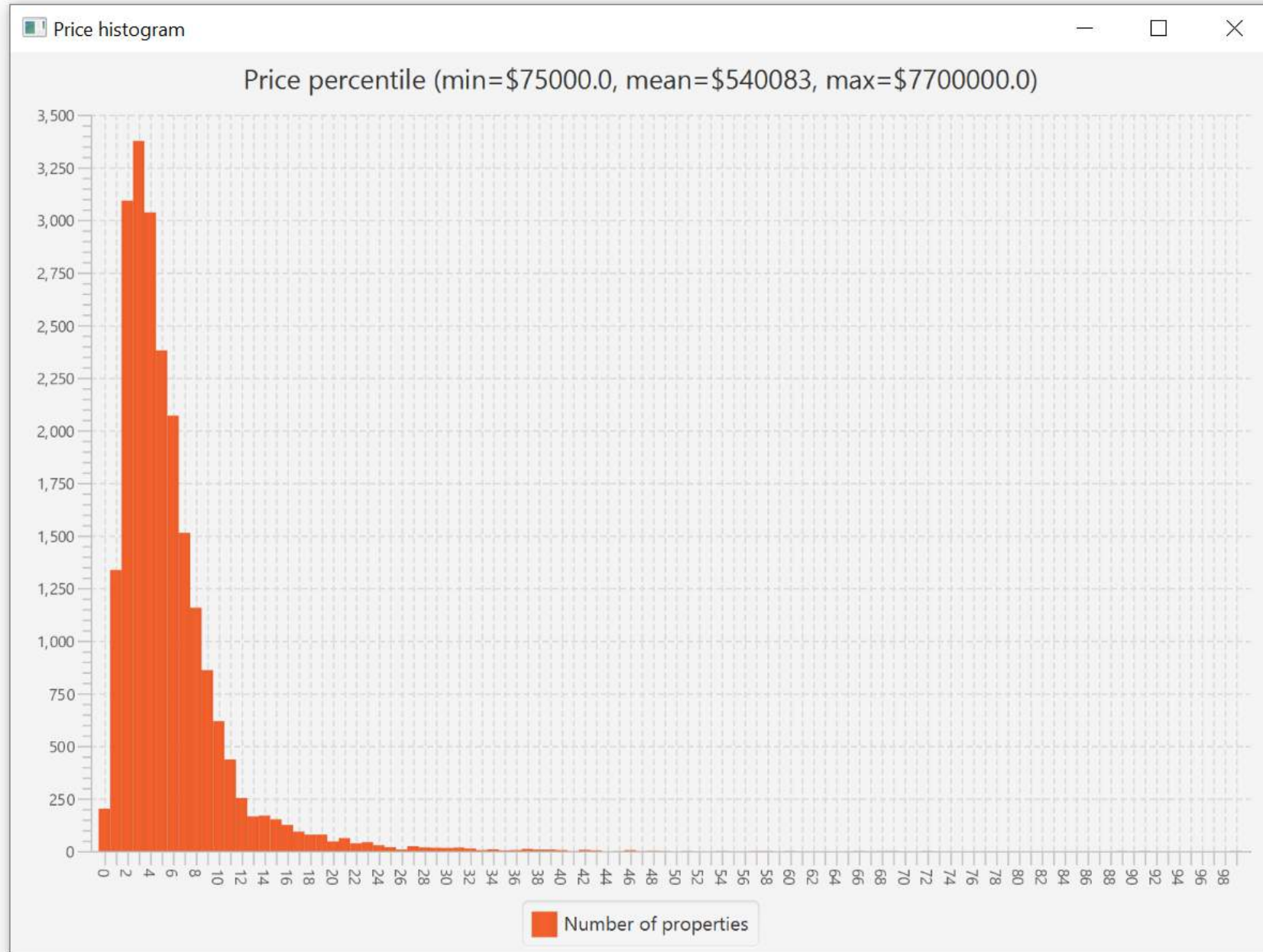
```
import org.apache.commons.math3.random.EmpiricalDistribution
import static groovyx.javafx.GroovyFX.start
import static org.apache.commons.csv.CSVFormat.RFC4180 as CSV

def file = 'kc_house_data.csv' as File
def csv = CSV.withFirstRecordAsHeader().parse(new FileReader(file))
def all = csv.findAll { it.bedrooms.toInteger() < 30 }.collect { it.price.toDouble() }
def info = new SummaryStatistics(); all.each(info::addValue)
def head = "Price percentile (min=\$$info.min, mean=\$${info.mean as int}, max=\$$info.max)"
def dist = new EmpiricalDistribution(100).tap{ load(all as double[]) }
def bins = dist.binStats.withIndex().collectMany { v, i -> [i.toString(), v.n] }

start {
  stage(title: 'Price histogram', show: true, width: 800, height: 600) {
    scene {
      barChart(title: head, barGap: 0, categoryGap: 0) {
        series(name: 'Number of properties', data: bins)
      }
    }
  }
}
```



# House price predictions – remove outlier and view



# House price predictions – linear regression

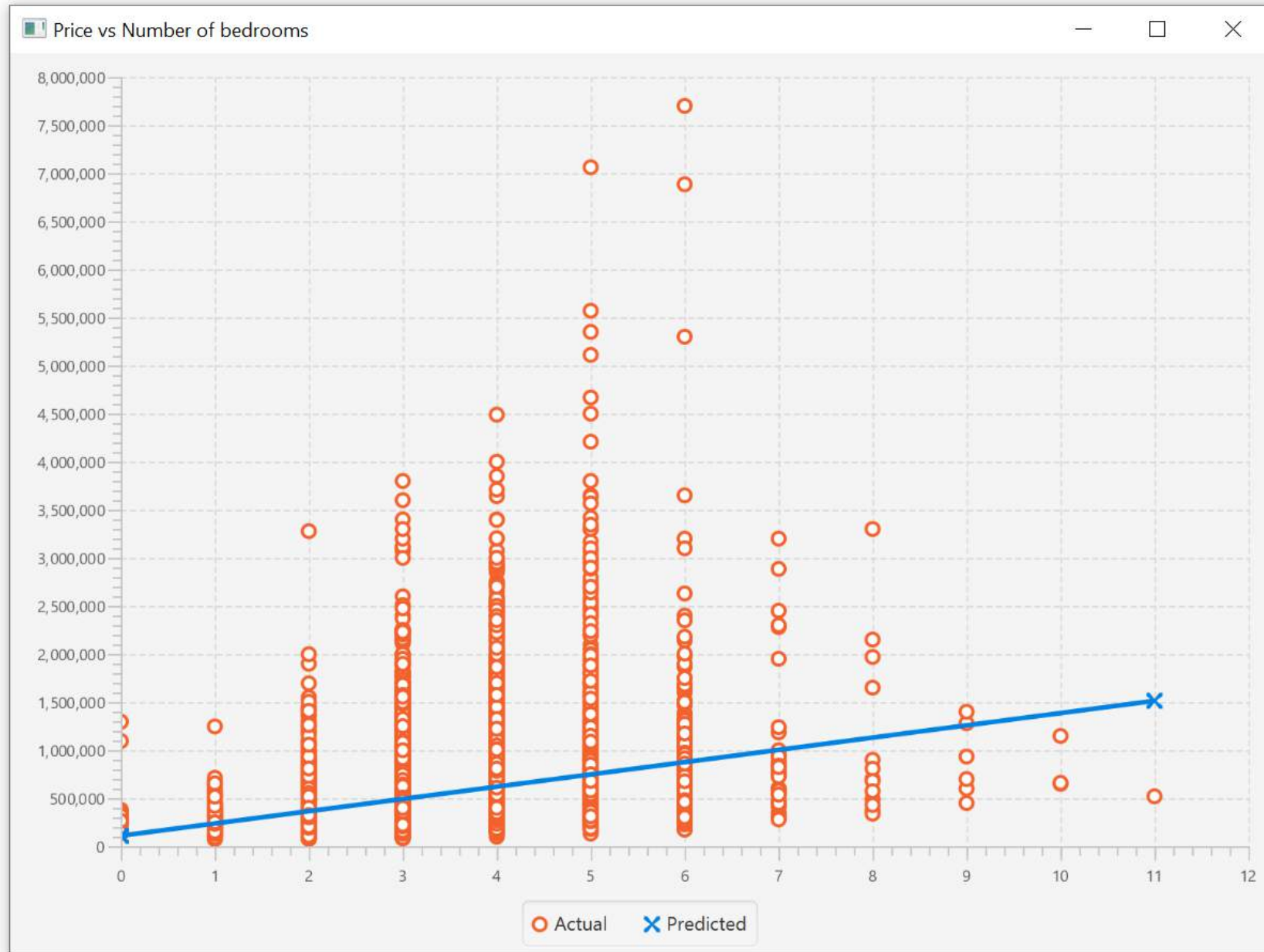
```
import org.apache.commons.math3.stat.regression.SimpleRegression
import static groovyx.javafx.GroovyFX.start
import static org.apache.commons.csv.CSVFormat.RFC4180 as CSV

def file = 'kc_house_data.csv' as File
def csv = CSV.withFirstRecordAsHeader().parse(new FileReader(file))
def all = csv.collect { [it.bedrooms.toDouble(), it.price.toDouble()] }.findAll{ it[0] < 30 }

def reg = new SimpleRegression().tap{ addData(all as double[][] ) }
def (minB, maxB) = [all.min{ it[0] }[0], all.max{ it[0] }[0]]
def predicted = [[minB, reg.predict(minB)], [maxB, reg.predict(maxB)]]

start {
    stage(title: 'Price vs Number of bedrooms', show: true, width: 800, height: 600) {
        scene {
            scatterChart {
                series(name: 'Actual', data: all)
            }
            lineChart {
                series(name: 'Predicted', data: predicted)
            }
        }
    }
}
```

# House price predictions – linear regression



# House price predictions - Weka

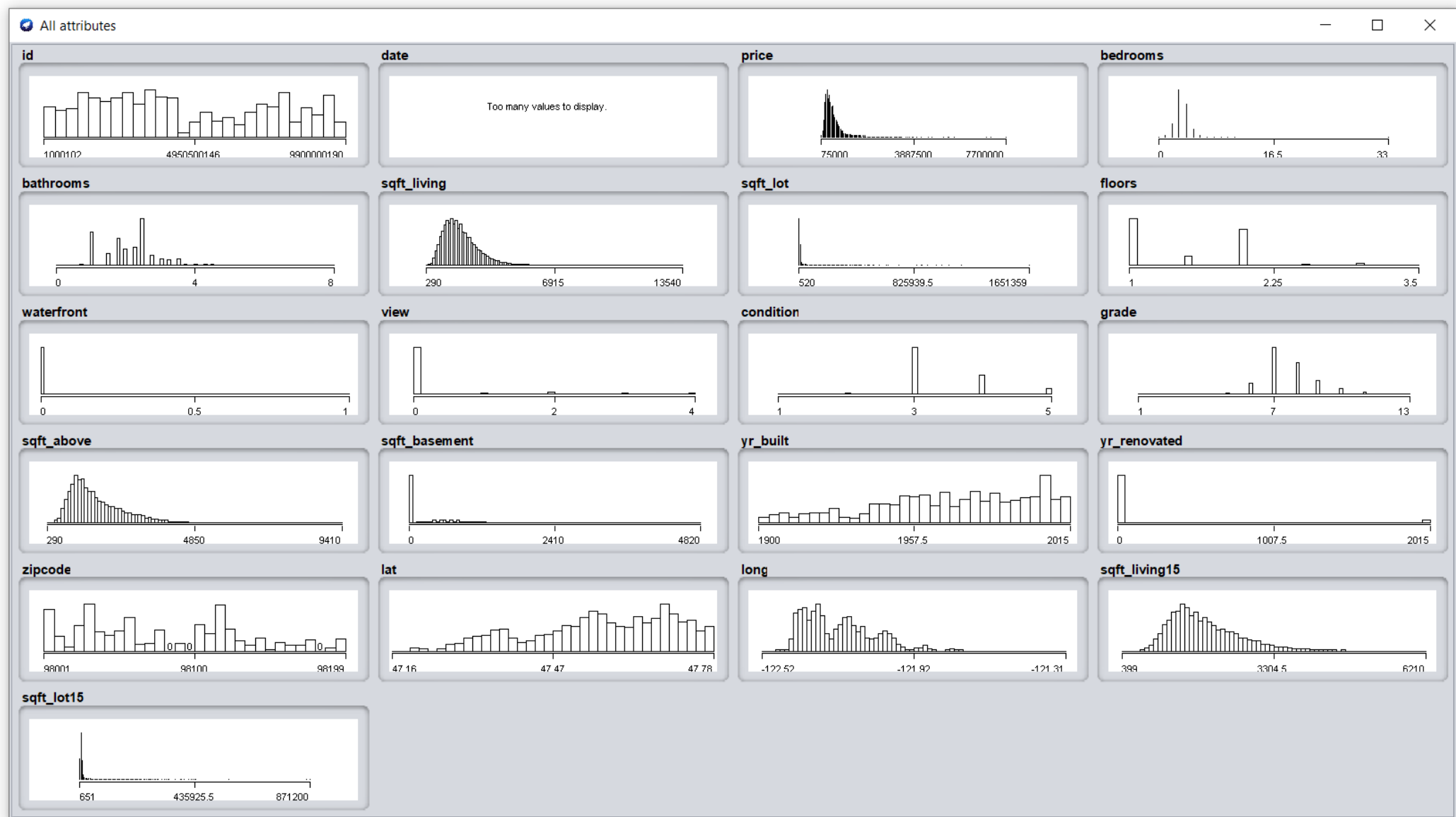


The screenshot shows the 'Weka Explorer' window. At the top, there are tabs for 'Preprocess', 'Classify', 'Cluster', 'Associate', 'Select attributes', and 'Visualize'. Below these are buttons for 'Open file...', 'Open URL...', 'Open DB...', 'Generate...', 'Undo', 'Edit...', and 'Save...'. A 'Filter' section has a 'Choose' button and a dropdown menu set to 'None', with 'Apply' and 'Stop' buttons. The 'Current relation' section shows 'Relation: kc\_house\_data', 'Attributes: 21', 'Instances: 21613', and 'Sum of weights: 21613'. The 'Attributes' section has buttons for 'All', 'None', 'Invert', and 'Pattern', and a list of attributes with checkboxes. The 'price' attribute is selected. The 'Selected attribute' section shows 'Name: price', 'Missing: 0 (0%)', 'Distinct: 4028', 'Type: Numeric', and 'Unique: 2395 (11%)'. Below this is a table of statistics for the 'price' attribute.

Statistic	Value
Minimum	75000
Maximum	7700000
Mean	540088.142
StdDev	367127.196

Below the statistics is a dropdown menu set to 'Class: price (Num)' and a 'Visualize All' button. A histogram shows the distribution of house prices, with the x-axis ranging from 75000 to 7700000. The status bar at the bottom shows 'OK', a 'Log' button, and a kiwi bird icon with 'x 0'.

# House price predictions



# House price predictions

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Attribute Evaluator: Choose **CorrelationAttributeEval**

Search Method: Choose **Ranker -T -1.7876931348823157E308 -N -1**

Attribute Selection Mode:  Use full training set  
 Cross-validation Folds: 10 Seed: 1

(Num) price

Start Stop

Result list (right-click for options)

21:53:12 - Ranker + CorrelationAttributeE

Attribute selection output

Search Method: Attribute ranking.

Attribute Evaluator (supervised, Class (numeric): 3 price): Correlation Ranking Filter

Ranked attributes:

0.70204	6	sqft_living
0.66743	12	grade
0.60557	13	sqft_above
0.58538	20	sqft_living15
0.52514	5	bathrooms
0.39729	10	view
0.32382	14	sqft_basement
0.30835	4	bedrooms
0.307	18	lat
0.26637	9	waterfront
0.25679	8	floors
0.12643	16	yr_renovated
0.08966	7	sqft_lot
0.08245	21	sqft_lot15
0.05401	15	yr_built
0.03636	11	condition
0.02163	19	long
0.00587	2	date
-0.01676	1	id
-0.0532	17	zipcode

Selected attributes: 6,12,13,20,5,10,14,4,18,5,8,16,7,21,15,11,19,2,1,17 : 20

Status: OK Log x 0

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose **None** Apply Stop

Current relation: Relation: house\_train-weka.filters.unsupervise... Instances: 17290 Attributes: 6 Sum of weights: 17290

Selected attribute: Name: price Type: Numeric Missing: 0 (0%) Distinct: 3361 Unique: 1919 (11%)

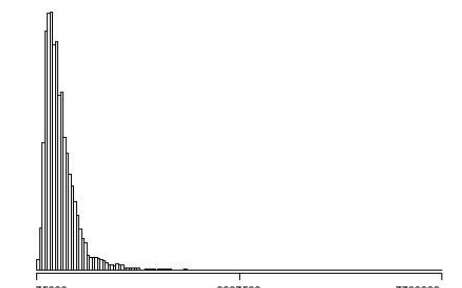
Statistic	Value
Minimum	75000
Maximum	7700000
Mean	533247.493
StdDev	365935.931

Attributes: All None Invert Pattern

No.	Name
<input checked="" type="checkbox"/>	1 price
<input type="checkbox"/>	2 bathrooms
<input type="checkbox"/>	3 sqft_living
<input type="checkbox"/>	4 grade
<input type="checkbox"/>	5 sqft_above
<input type="checkbox"/>	6 sqft_living15

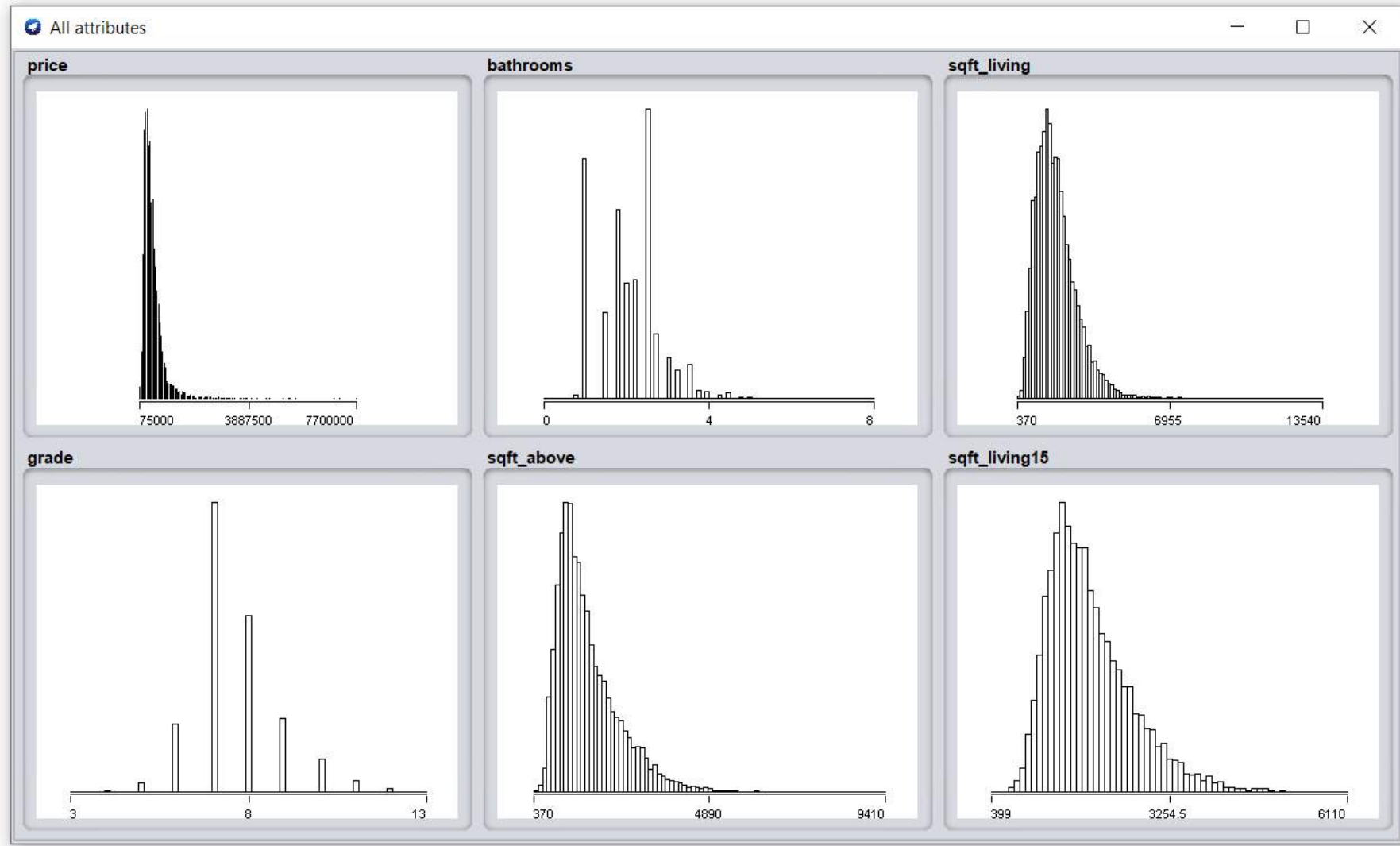
Remove

Class: sqft\_living15 (Num) Visualize All



Status: OK Log x 0

# House price predictions



# House price predictions

The screenshot shows the Weka Explorer interface. The 'Classifier' tab is active, and a RandomForest classifier is selected with the following configuration: `RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1`. The 'Test options' section shows 'Cross-validation' selected with 'Folds' set to 10. The 'Classifier output' pane displays the following text:

```
Relation: house_train-weka.filters.unsupervised.attribute.Remove-R1-2-weka.f
Instances: 17290
Attributes: 6
price
bathrooms
sqft_living
grade
sqft_above
sqft_living15
Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

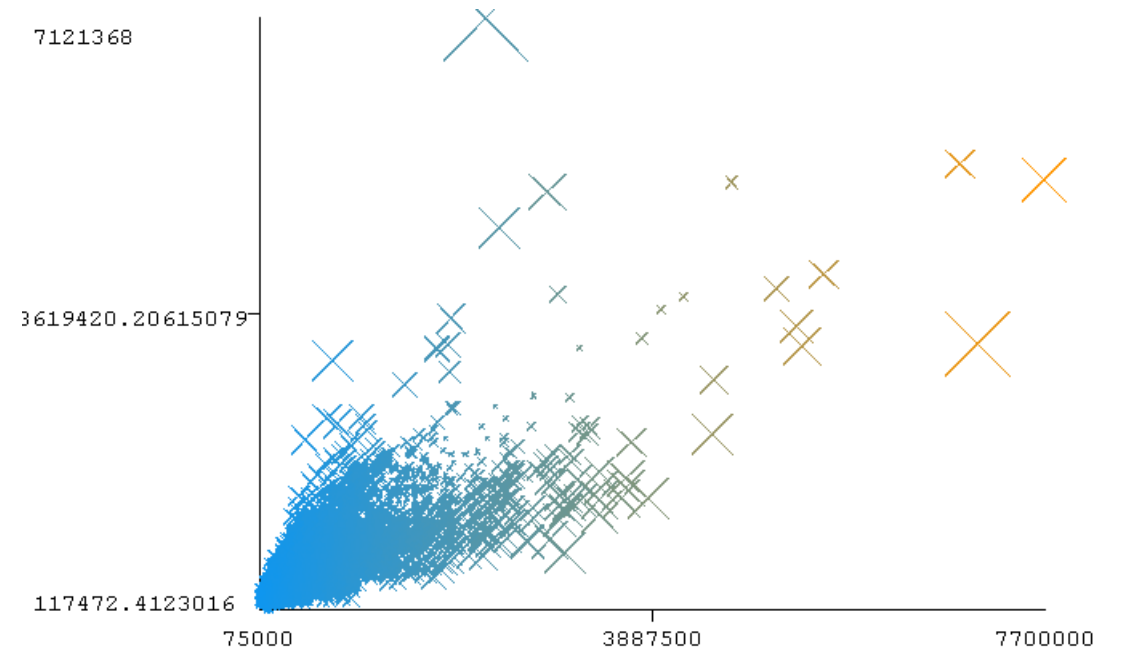
RandomForest

Bagging with 100 iterations and base learner

...
taken to build model: 4.44 seconds

Cross-validation ===
Summary ===
relation coefficient          0.7716
absolute error              151470.7871
mean squared error         233856.8653
relative absolute error     65.2858 %
relative squared error     63.9063 %
Total Number of Instances  17290
```

A context menu is open over the 'Result list' section, with 'Save model' selected. The status bar at the bottom shows 'OK' and a 'Log' button.





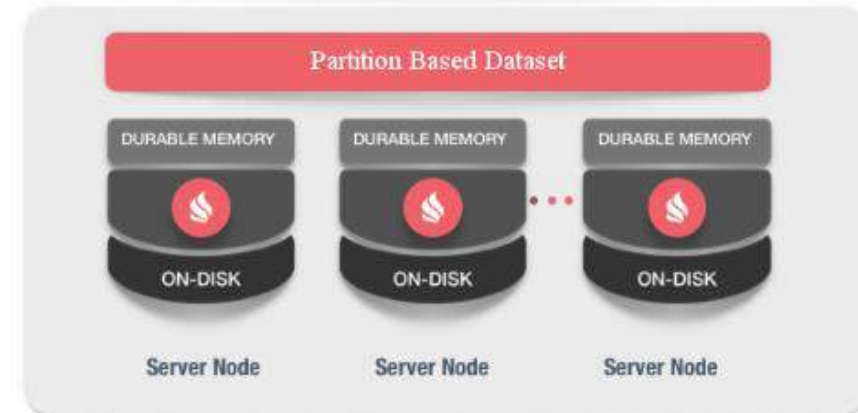
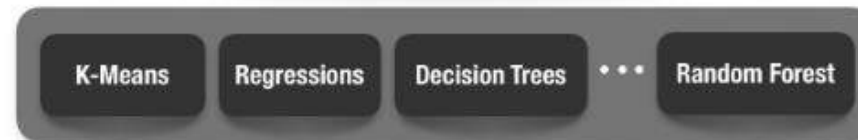
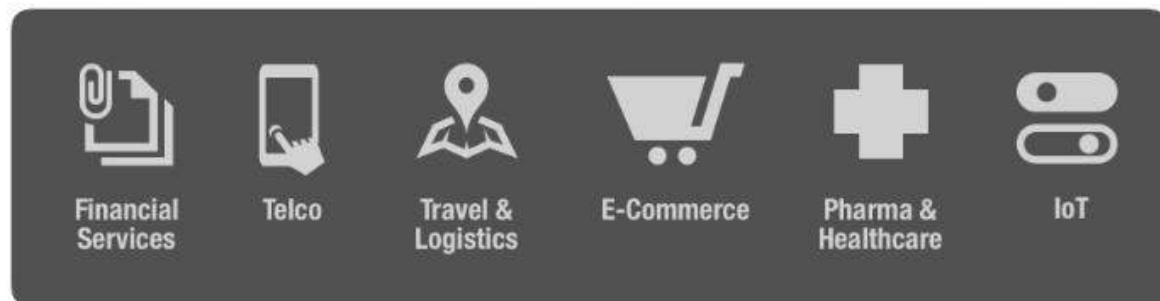
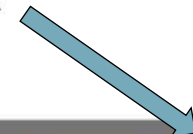
# House price predictions – split dataset

```
def full = getClass().classLoader.getResource('kc_house_data.csv').file as File
def parent = full.parentFile
def lines = full.readlines()
def (trainLines, testLines) = lines.chop(lines.size() * 0.8 as int, -1)

def train = new File(parent, 'house_train.csv')
train.text = trainLines.join('\n')

def test = new File(parent, 'house_test.csv')
test.delete()
test << lines[0] << '\n' << testLines.join('\n')
```

# Scaling up



# House price predictions – linear regression

```
Ignite ignite = Ignition.start("/path/to/example-ignite.xml")
ignite.withCloseable {
    println ">>> Ignite grid started."
    def dataCache = new SandboxMLCache(ignite).fillCacheWith(BOSTON_HOUSE_PRICES)
    def trainer = new LinearRegressionLSQRTrainer()
        .withEnvironmentBuilder(defaultBuilder().withRNGSeed(0))
    def vectorizer = new DummyVectorizer().labeled(Vectorizer.LabelCoordinate.FIRST)
    def split = new TrainTestDatasetSplitter().split(0.8)
    def mdl = trainer.fit(ignite, dataCache, split.trainFilter, vectorizer)
    def metric = new RegressionMetrics().withMetric { it.r2() }
    def score = Evaluator.evaluate(dataCache, split.testFilter, mdl, vectorizer, metric)
```

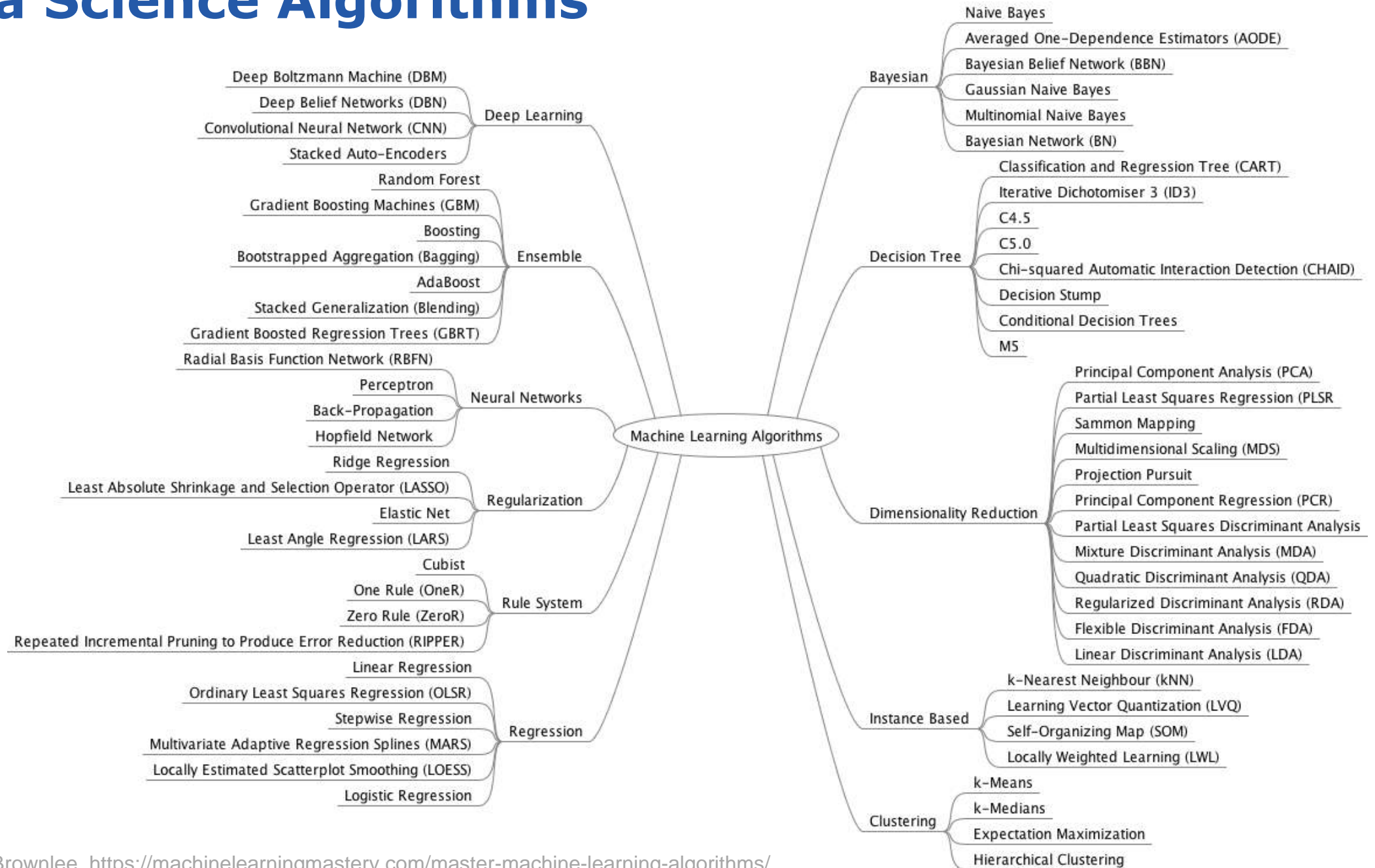


```
    println ">>> Model: " + toString(mdl)
    println ">>> R^2 score: " + score
    dataCache.destroy()
}
```

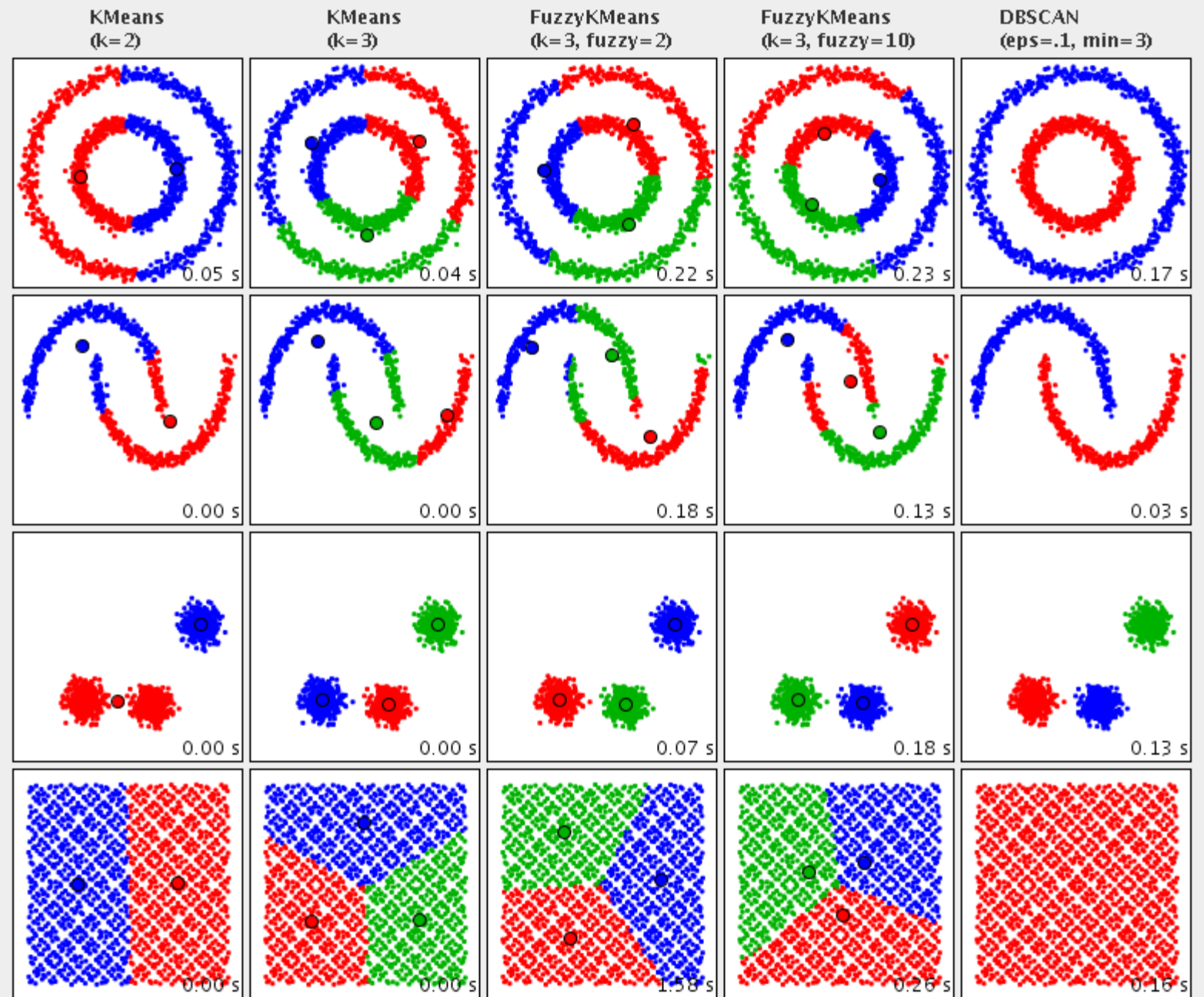
```
static toString(mdl) {
    def sign = { val -> val < 0 ? '-' : '+' }
    def valIdx = { idx, val -> sprintf '%.2f*f%d', val, idx }
    def valIdxSign = { idx, val -> sign(val) + valIdx(idx, Math.abs(val)) }
    def valSign = { val -> sign(val) + sprintf '%.2f', Math.abs(val) }
    def w = mdl.weights
    def i = mdl.intercept
    def result = [valIdx(0, w.get(0)), *(1..<w.size()).collect{ valIdxSign(it, w.get(it)) }, valSign(i)]
    result.join(' ')
}
```

```
[09:17:59] Ignite node started OK (id=0b9592e1)
[09:17:59] Topology snapshot [ver=1, locNode=0b9592e1, servers=1, clients=0,
state=ACTIVE, CPUs=12, offheap=6.4GB, heap=7.1GB]
>>> Ignite grid started.
>>> Model: 0.04*f0 - 0.05*f1 - 0.79*f2 - 10.01*f3 + 0.28*f4 + 0.01*f5 -
0.89*f6 + 0.54*f7 - 0.00*f8 - 0.18*f9 - 0.01*f10 + 0.09*f11 - 0.15*f12 + 14.10
>>> R^2 score: 0.412878309612252
[09:17:59] Ignite node stopped OK [uptime=00:00:00.626]
```

# Data Science Algorithms

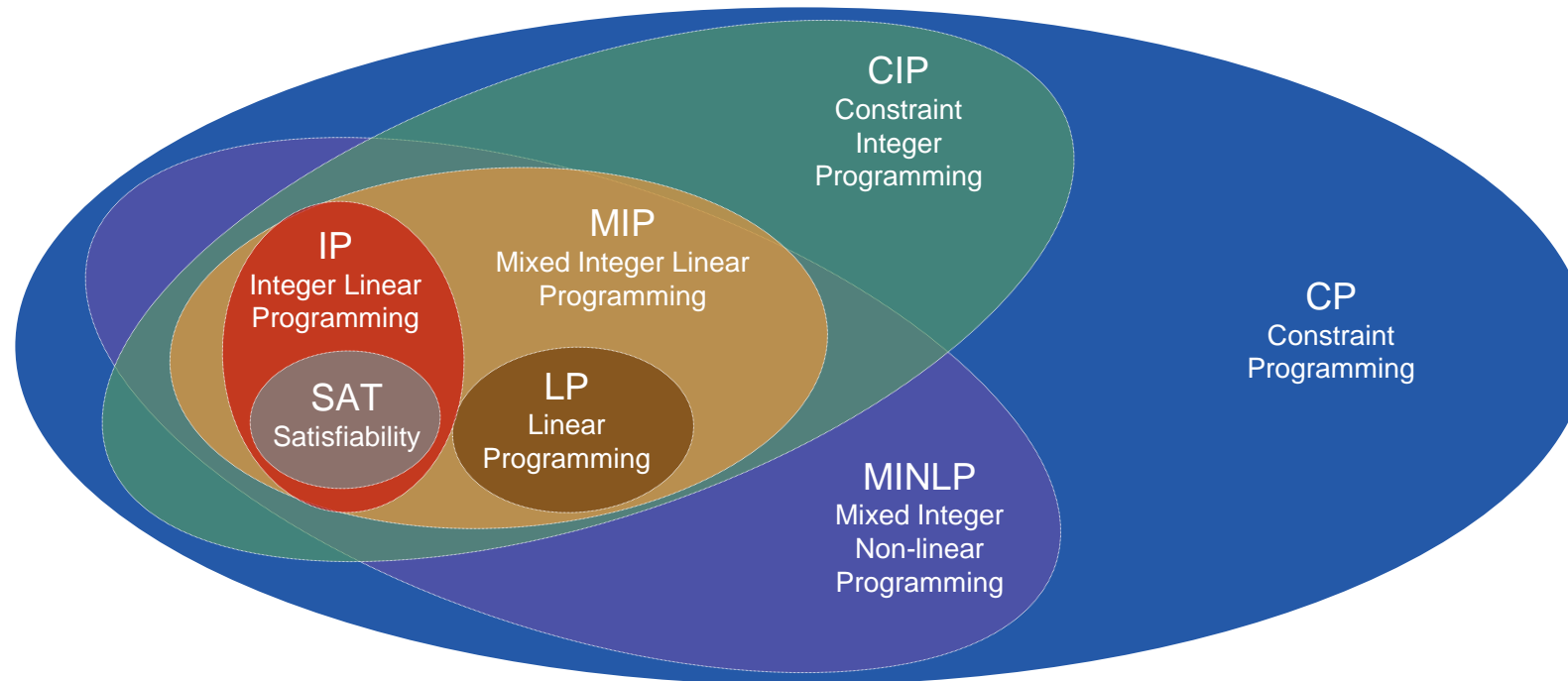


# Data Science Algorithms: Classification



# Computational problems in Computer Science

- Kinds: decision problems, search problems, counting problems, optimization problems
- Paradigms: brute force, divide & conquer, search & enumerate, randomized, complexity reduction, recursive, back tracking, graph
- Optimization approaches: linear programming, dynamic programming, greedy, heuristics



# Constraint Programming

## Where is it used?

- Factory/assembly line scheduling
- Workforce management/scheduling (call centres, aircraft crew, etc)
- Vehicle routing/traffic planning/traffic management
- Packing problems
- Timetabling (exams, lectures, trains, gate arrival at airports)
- Configuration and design
- Supertree construction (bioinformatics)
- Network design (telecommunications)
- Military logistics (air cover for naval fleet)
- Aircraft maintenance schedules
- Recruitment (personnel selection)
- Financial trading systems
- Language detection
- Operations research

# Constraint Programming

Typical domains:

- boolean domains, where only true/false constraints apply (SAT problem)
- integer domains
- rational domains
- interval domains, in particular for scheduling problems
- linear domains, where only linear functions are described and analyzed (although approaches to non-linear problems do exist)
- finite domains, where constraints are defined over finite sets
- mixed domains, involving two or more of the above



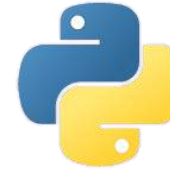
# Cryptarithmic

$$\begin{array}{r} \phantom{+} \phantom{=} \phantom{M} \phantom{O} \phantom{N} \phantom{E} \phantom{D} \\ \phantom{+} \phantom{=} \phantom{M} \phantom{O} \phantom{N} \phantom{E} \phantom{D} \\ + \phantom{=} \phantom{M} \phantom{O} \phantom{N} \phantom{E} \phantom{D} \\ = \phantom{M} \phantom{O} \phantom{N} \phantom{E} \phantom{D} \phantom{Y} \\ \phantom{+} \phantom{=} \phantom{M} \phantom{O} \phantom{N} \phantom{E} \phantom{D} \\ \phantom{+} \phantom{=} \phantom{M} \phantom{O} \phantom{N} \phantom{E} \phantom{D} \end{array}$$

Replace the letters with decimal digits to make a valid arithmetic sum

# Cryptarithmic: Brute force

```
def solutions():  
    # Letters = ('s', 'e', 'n', 'd', 'm', 'o', 'r', 'y')  
    all_solutions = list()  
    for s in range(1, 10):  
        for e in range(0, 10):  
            for n in range(0, 10):  
                for d in range(0, 10):  
                    for m in range(1, 10):  
                        for o in range(0, 10):  
                            for r in range(0, 10):  
                                for y in range(0, 10):  
                                    if len({s, e, n, d, m, o, r, y}) == 8:  
                                        send = 1000 * s + 100 * e + 10 * n + d  
                                        more = 1000 * m + 100 * o + 10 * r + e  
                                        money = 10000 * m + 1000 * o + 100 * n + 10 * e + y  
  
                                        if send + more == money:  
                                            all_solutions.append((send, more, money))  
  
    return all_solutions  
  
print(solutions())
```



	S	E	N	D	
+	M	O	R	E	
=	M	O	N	E	Y

```
[(9567, 1085, 10652)]
```

# Cryptarithmic: Brute force

```
def solutions() {
  // letters = ['s', 'e', 'n', 'd', 'm', 'o', 'r', 'y']
  def all_solutions = []
  for (s in 1..<10)
    for (e in 0..9)
      for (n in 0..9)
        for (d in 0..9)
          for (m in 1..9)
            for (o in 0..9)
              for (r in 0..9)
                for (y in 0..9)
                  if ([s, e, n, d, m, o, r, y].toSet().size() == 8) {
                    def send = 1000 * s + 100 * e + 10 * n + d
                    def more = 1000 * m + 100 * o + 10 * r + e
                    def money = 10000 * m + 1000 * o + 100 * n + 10 * e + y
                    if (send + more == money)
                      all_solutions.add([send, more, money])
                  }
                }
            }
          }
        }
      }
    }
  }

  return all_solutions
}

print(solutions())
```



	S	E	N	D	
+	M	O	R	E	
=	M	O	N	E	Y

```
[[9567, 1085, 10652]]
```

# Cryptarithmic: Brute force

```
from itertools import permutations
```

```
def solution2():
```

```
    letters = ('s', 'e', 'n', 'd', 'm', 'o', 'r', 'y')
```

```
    digits = range(10)
```

```
    for perm in permutations(digits, len(letters)):
```

```
        sol = dict(zip(letters, perm))
```

```
        if sol['s'] == 0 or sol['m'] == 0:
```

```
            continue
```

```
        send = 1000 * sol['s'] + 100 * sol['e'] + 10 * sol['n'] + sol['d']
```

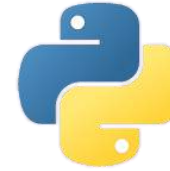
```
        more = 1000 * sol['m'] + 100 * sol['o'] + 10 * sol['r'] + sol['e']
```

```
        money = 10000 * sol['m'] + 1000 * sol['o'] + 100 * sol['n'] + 10 * sol['e'] + sol['y']
```

```
        if send + more == money:
```

```
            return send, more, money
```

```
print(solution2())
```



```
      S E N D
+     M O R E
=     M O N E Y
```

```
(9567, 1085, 10652)
```

# Cryptarithmic: Brute force

```
def solution2() {  
  def digits = 0..9  
  for (p in digits.permutations()) {  
    if (p[-1] < p[-2]) continue  
    def (s, e, n, d, m, o, r, y) = p  
    if (s == 0 || m == 0)  
      continue  
    def send = 1000 * s + 100 * e + 10 * n + d  
    def more = 1000 * m + 100 * o + 10 * r + e  
    def money = 10000 * m + 1000 * o + 100 * n + 10 * e + y  
    if (send + more == money)  
      return [send, more, money]  
  }  
}  
  
print(solution2())
```



	S	E	N	D	
+	M	O	R	E	
=	M	O	N	E	Y

[9567, 1085, 10652]

# Cryptarithmic: Constraint programming

```
from csp import Constraint, CSP
from typing import Dict, List, Optional

class SendMoreMoneyConstraint(Constraint[str, int]):
    def __init__(self, letters: List[str]) -> None:
        super().__init__(letters)
        self.letters: List[str] = letters

    def satisfied(self, assignment: Dict[str, int]) -> bool:
        # not a solution if duplicate values
        if len(set(assignment.values())) < len(assignment):
            return False

        # if all vars assigned, check if correct
        if len(assignment) == len(self.letters):
            s: int = assignment["S"]
            e: int = assignment["E"]
            n: int = assignment["N"]
            d: int = assignment["D"]
            m: int = assignment["M"]
            o: int = assignment["O"]
            r: int = assignment["R"]
            y: int = assignment["Y"]
            send: int = s * 1000 + e * 100 + n * 10 + d
            more: int = m * 1000 + o * 100 + r * 10 + e
            money: int = m * 10000 + o * 1000 + n * 100 + e * 10 + y
            return send + more == money

        return True
```



```
  S E N D
+ M O R E
= M O N E Y
```

```
...
letters = ["S", "E", "N", "D", "M", "O", "R", "Y"]
possible_digits = {}
for letter in letters:
    possible_digits[letter] = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
possible_digits["M"] = [1] # can't start with 0 constraint
csp: CSP[str, int] = CSP(letters, possible_digits)
csp.add_constraint(SendMoreMoneyConstraint(letters))
solution: Optional[Dict[str, int]] = csp.backtracking_search()
if solution is None:
    print("No solution found!")
else:
    print(solution)
```

```
{'S': 9, 'E': 5, 'N': 6, 'D': 7,
'M': 1, 'O': 0, 'R': 8, 'Y': 2}
```

# Cryptarithmic: Constraint programming

```
@Grab('org.choco-solver:choco-solver:4.10.0')
import org.chocosolver.solver.Model
import org.chocosolver.solver.variables.IntVar

def model = new Model("SEND+MORE=MONEY")
def S = model.intVar("S", 1, 9)
def E = model.intVar("E", 0, 9)
def N = model.intVar("N", 0, 9)
def D = model.intVar("D", 0, 9)
def M = model.intVar("M", 1, 9)
def O = model.intVar("0", 0, 9)
def R = model.intVar("R", 0, 9)
def Y = model.intVar("Y", 0, 9)

model.allDifferent(S, E, N, D, M, O, R, Y).post()
...
```

Solution: S=9, E=5, N=6,  
D=7, M=1, 0=0, R=8, Y=2,



```
  S E N D
+  M O R E
=  M O N E Y
```

```
...
IntVar[] ALL = [
    S, E, N, D,
    M, O, R, E,
    M, O, N, E, Y]
int[] COEFFS = [
    1000, 100, 10, 1,
    1000, 100, 10, 1,
    -10000, -1000, -100, -10, -1]
model.scalar(ALL, COEFFS, "=", 0).post()

model.solver.findSolution()
```

# Dietary restrictions

	Bread	Milk	Cheese	Potato	Fish	Yogurt
<i>Cost</i>	2.0	3.5	8.0	1.5	11.0	1.0
<i>Protein (g)</i>	4.0	8.0	7.0	1.3	8.0	9.2
<i>Fat (g)</i>	1.0	5.0	9.0	0.1	7.0	1.0
<i>Carbohydrates (g)</i>	15.0	11.7	0.4	22.6	0.0	17.0
<i>Calories</i>	90	120	106	97	130	180

Minimise cost of diet given:

- Must be at least 300 calories
- Not more than 10 grams of protein
- Not less than 10 grams of carbohydrates
- Not less than 8 grams of fat
- At least 0.5 units of fish
- No more than 1 unit of milk





# Dietary restrictions (Excel)

Solver Parameters

Set Objective:

To:  Max  Min  Value Of:

By Changing Variable Cells:

Subject to the Constraints:

- 
- 
- 
- 
- 
- 

Make Unconstrained Variables Non-Negative

Select a Solving Method:

Solving Method

Select the GRG Nonlinear engine for Solver Problems that are smooth nonlinear. Select the LP Simplex engine for linear Solver Problems, and select the Evolutionary engine for Solver problems that are non-smooth.

	<i>Bread</i>	<i>Milk</i>	<i>Cheese</i>	<i>Potato</i>	<i>Fish</i>	<i>Yogurt</i>	<i>Total</i>
<i>Cost</i>	2	3.5	8	1.5	11	1	0
<i>Protein</i>	4	8	7	1.3	8	9.2	0
<i>Fat</i>	1	5	9	0.1	7	1	0
<i>Carbs</i>	15	11.7	0.4	22.6	0	17	0
<i>Calories</i>	90	120	106	97	130	180	0



# Dietary restrictions (Excel)

Solver Parameters

Set Objective:

To:  Max  Min  Value Of:

By Changing Variable Cells:

Subject to the Constraints:

- 
- 
- 
- 
- 
- 

Make Unconstrained Variables Non-Negative

Select a Solving Method:

Solving Method  
Select the GRG Nonlinear engine for Solver Problems that are smooth nonlinear. Select the LP Simplex engine for linear Solver Problems, and select the Evolutionary engine for Solver problems that are non-smooth.

	<i>Bread</i>	<i>Milk</i>	<i>Cheese</i>	<i>Potato</i>	<i>Fish</i>	<i>Yogurt</i>	<i>Total</i>
	0	0.053599	0.449499	1.865168	0.5	0	
<i>Cost</i>	2	3.5	8	1.5	11	1	12.08134
<i>Protein</i>	4	8	7	1.3	8	9.2	10
<i>Fat</i>	1	5	9	0.1	7	1	8
<i>Carbs</i>	15	11.7	0.4	22.6	0	17	42.9597
<i>Calories</i>	90	120	106	97	130	180	300



# Dietary restrictions (Google sheets)

	<i>Bread</i>	<i>Milk</i>	<i>Cheese</i>	<i>Potato</i>	<i>Fish</i>	<i>Yogurt</i>	<i>Total</i>	Min	Max
	0.0000	0.0536	0.4495	1.8652	0.5000	0.0000			
<i>Cost</i>	2	3.5	8	1.5	11	1	12.08133788		
<i>Protein</i>	4	8	7	1.3	8	9.2	10		1
<i>Fat</i>	1	5	9	0.1	7	1	8	8	
<i>Carbs</i>	15	11.7	0.4	22.6	0	17	42.95969651	10	
<i>Calories</i>	90	120	106	97	130	180	300	300	
Min					0.5				
Max		1							

Objective Cell:

I4 Update Clear

Objective Sense:

minimise
  target value: 0  
 maximise

Variable Cells:

C3:H3 Add Update Delete

Unconstrained variables non-negative

Constraints:

<Add new constraint>  
 I5 <= K5  
 I6 >= J6  
 I7 >= J7  
 I8 >= J8  
 G3 >= G10  
 D3 <= D11



# Diet problem (ojalgo)

```
def model = new ExpressionsBasedModel()

def bread = model.addVariable("Bread").lower(0)
def milk = model.addVariable("Milk").lower(0).upper(1)
def cheese = model.addVariable("Cheese").lower(0)
def potato = model.addVariable("Potato").lower(0)
def fish = model.addVariable("Fish").lower(0.5)
def yogurt = model.addVariable("Yogurt").lower(0)

def cost = model.addExpression("Cost")
cost.set(bread, 2.0).set(milk, 3.5).set(cheese, 8.0).set(potato, 1.5).set(fish,
11.0).set(yogurt, 1.0)

def protein = model.addExpression("Protein").upper(10)
protein.set(bread, 4.0).set(milk, 8.0).set(cheese, 7.0).set(potato,
1.3).set(fish, 8.0).set(yogurt, 9.2)

def fat = model.addExpression("Fat").lower(8)
fat.set(bread, 1.0).set(milk, 5.0).set(cheese, 9.0).set(potato, 0.1).set(fish,
7.0).set(yogurt, 1.0)

def carbs = model.addExpression("Carbohydrates").lower(10)
carbs.set(bread, 15.0).set(milk, 11.7).set(cheese, 0.4).set(potato,
22.6).set(fish, 0.0).set(yogurt, 17.0)

def calories = model.addExpression("Calories").lower(300)
calories.set(bread, 90).set(milk, 120).set(cheese, 106).set(potato, 97).set(fish,
130).set(yogurt, 180)

def result = model.minimise()

// for a variation, see:
// https://www.ojalgo.org/2019/05/the-diet-problem/
```



```
...
OPTIMAL 0.0 @ { 0.0, 0.0, 0.40813898143741034,
1.8538791051880057, 0.5916230366492151, 0.0 }
#####
0 <= Bread: 0
0 <= Milk: 0 <= 1
0 <= Cheese: 0.408139
0 <= Potato: 1.853879
0.5 <= Fish: 0.591623
0 <= Yogurt: 0
10 <= Carbohydrates: 42.060923
8 <= Fat: 8.0
Cost: 12.553784
Protein: 10.0 <= 10
300 <= Calories: 300.0
#####
```

# Diet problem (Choco)

```
def model = new Model("Diet problem")
def unbounded = 100000

// scale quantities by 10, coefficients by 10, products by 100
def bread = model.intVar("Bread", 0, unbounded)
def milk = model.intVar("Milk", 0, 10)
def cheese = model.intVar("Cheese", 0, unbounded)
def potato = model.intVar("Potato", 0, unbounded)
def fish = model.intVar("Fish", 5, unbounded)
def yogurt = model.intVar("Yogurt", 0, unbounded)
IntVar[] all = [bread, milk, cheese, potato, fish, yogurt]

def cost = model.intVar("Cost", 0, unbounded)
model.scalar(all, [20, 35, 80, 15, 110, 10] as int[], "=", cost).post()

def protein = model.intVar("Protein", 0, 1000)
model.scalar(all, [40, 80, 70, 13, 80, 92] as int[], "=", protein).post()

def fat = model.intVar("Fat", 800, unbounded)
model.scalar(all, [10, 50, 90, 1, 70, 10] as int[], "=", fat).post()

def carbs = model.intVar("Carbohydrates", 1000, unbounded)
model.scalar(all, [150, 117, 4, 226, 0, 170] as int[], "=", carbs).post()

def calories = model.intVar("Calories", 30000, unbounded)
model.scalar(all, [900, 1200, 1060, 970, 1300, 1800] as int[], "=", calories).post()

model.setObjective(Model.MINIMIZE, cost)

def found = model.solver.findSolution()
if (found) {
  all.each { println "$it.name: ${it.value / 10}" }
  [carbs, fat, protein, calories, cost].each { println "$it.name: ${it.value / 100}" }
} else {
  println "No solution"
}
```



```
Bread: 0
Milk: 0
Cheese: 0.5
Potato: 1.9
Fish: 0.5
Yogurt: 0
Carbohydrates: 43.14
Fat: 8.19
Protein: 9.97
Calories: 302.3
Cost: 12.35
```

*Choco supports RealVar but doesn't have as rich a set of possible constraints for such vars.*

# Diet problem (Choco)

```
def model = new Model("Diet problem")
def unbounded = 1000.0d
def precision = 0.00001d

// scale quantities by 10, coefficients by 10, products by 100
def bread = model.realVar("Bread", 0.0, unbounded, precision)
def milk = model.realVar("Milk", 0.0, 1.0, precision)
def cheese = model.realVar("Cheese", 0.0, unbounded, precision)
def potato = model.realVar("Potato", 0.0, unbounded, precision)
def fish = model.realVar("Fish", 0.5, unbounded, precision)
def yogurt = model.realVar("Yogurt", 0.0, unbounded, precision)
RealVar[] all = [bread, milk, cheese, potato, fish, yogurt]

def scalarIbex = { coeffs, var ->
  def (a, b, c, d, e, f) = coeffs
  model.realIbexGenericConstraint("$a*{0}+$b*{1}+$c*{2}+$d*{3}+$e*{4}+$f*{5}={6}",
    [*all, var] as RealVar[]).post();
}

def cost = model.realVar("Cost", 0.0, unbounded, precision)
scalarIbex([2.0, 3.5, 8.0, 1.5, 11.0, 1.0], cost)
def protein = model.realVar("Protein", 0.0, 10.0, precision)
scalarIbex([4.0, 8.0, 7.0, 1.3, 8.0, 9.2], protein)
def fat = model.realVar("Fat", 8.0, unbounded, precision)
scalarIbex([1.0, 5.0, 9.0, 0.1, 7.0, 1.0], fat)
def carbs = model.realVar("Carbohydrates", 10.0, unbounded, precision)
scalarIbex([15.0, 11.7, 0.4, 22.6, 0.0, 17.0], carbs)
def calories = model.realVar("Calories", 300, unbounded, precision)
scalarIbex([90, 120, 106, 97, 130, 180], calories)

model.setObjective(Model.MINIMIZE, cost)

def found = model.solver.findSolution()
```

```
def pretty = { var ->
  def bounds = found.getRealBounds(var)
  printf "%s: %.6f .. %.6f%n", var.name, *bounds
}

if (found) {
  all.each { pretty(it) }
  [carbs, fat, protein, calories, cost].each { pretty(it) }
} else {
  println "No solution"
}
```

```
Bread: 0.025131 .. 0.025137
Milk: 0.000009 .. 0.000010
Cheese: 0.428571 .. 0.428571
Potato: 1.848118 .. 1.848124
Fish: 0.561836 .. 0.561836
Yogurt: 0.000007 .. 0.000010
Carbohydrates: 42.316203 .. 42.316211
Fat: 8.000000 .. 8.000005
Protein: 9.997920 .. 9.997926
Calories: 300.000000 .. 300.000008
Cost: 12.431241 .. 12.431245
```

*Choco does have a plugin (via JNI) for the Ibex C++ constraint processing library which does handle real numbers.*

# Diet problem (Apache Commons Math)

```
import org.apache.commons.math3.optim.linear.*
import org.apache.commons.math3.optim.nonlinear.scalar.GoalType
import static org.apache.commons.math3.optim.linear.Relationship.*

def cost = new LinearObjectiveFunction([2.0, 3.5, 8.0, 1.5, 11.0, 1.0] as double[], 0)

static scalar(coeffs, rel, val) { new LinearConstraint(coeffs as double[], rel, val) }

def bread_min = scalar([1, 0, 0, 0, 0, 0], GEQ, 0)
def milk_min = scalar([0, 1, 0, 0, 0, 0], GEQ, 0)
def milk_max = scalar([0, 1, 0, 0, 0, 0], LEQ, 1)
def cheese_min = scalar([0, 0, 1, 0, 0, 0], GEQ, 0)
def potato_min = scalar([0, 0, 0, 1, 0, 0], GEQ, 0)
def fish_min = scalar([0, 0, 0, 0, 1, 0], GEQ, 0.5)
def yogurt_min = scalar([0, 0, 0, 0, 0, 1], GEQ, 0)
def protein = scalar([4.0, 8.0, 7.0, 1.3, 8.0, 9.2], LEQ, 10)
def fat = scalar([1.0, 5.0, 9.0, 0.1, 7.0, 1.0], GEQ, 8)
def carbs = scalar([15.0, 11.7, 0.4, 22.6, 0.0, 17.0], GEQ, 10)
def calories = scalar([90, 120, 106, 97, 130, 180], GEQ, 300)

LinearConstraintSet constraints = [bread_min, milk_min, milk_max, fish_min, cheese_min,
                                  potato_min, yogurt_min, protein, fat, carbs, calories]

def solution = new SimplexSolver().optimize(cost, constraints, GoalType.MAXIMIZE)

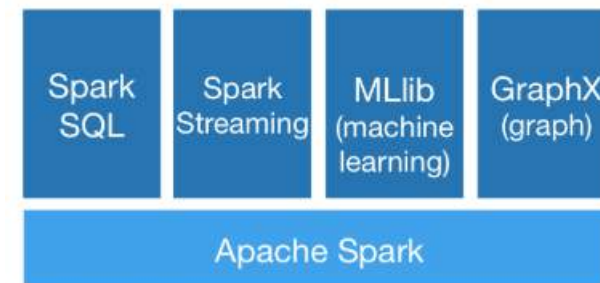
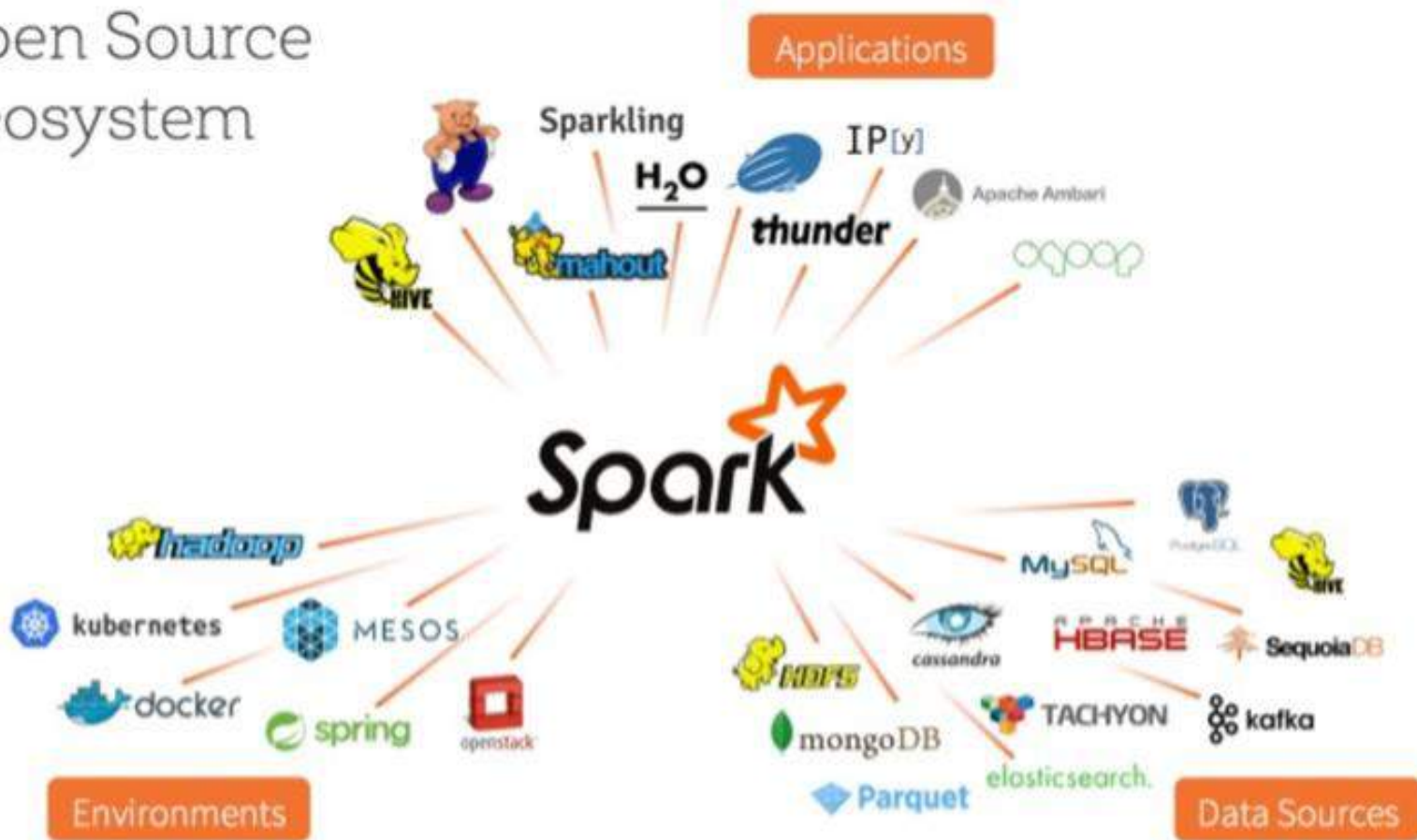
if (solution != null) {
    println "Opt: $solution.value"
    println solution.point.collect{ sprintf '%.2f', it }.join(', ')
}
```



```
Opt: 12.553783912422674
-0.00, -0.00, 0.41, 1.85, 0.59, 0.00
```

# Apache Spark

Open Source  
Ecosystem

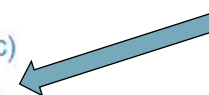


## Programming Guides:

- [Quick Start](#): a quick introduction to the Spark API; start here!
- [RDD Programming Guide](#): overview of Spark basics - RDDs
- [Spark SQL, Datasets, and DataFrames](#): processing structure
- [Structured Streaming](#): processing structured data streams with
- [Spark Streaming](#): processing data streams using DStreams
- [MLlib](#): applying machine learning algorithms
- [GraphX](#): processing graphs

## API Docs:

- [Spark Scala API \(Scaladoc\)](#)
- [Spark Java API \(Javadoc\)](#)
- [Spark Python API \(Sphinx\)](#)
- [Spark R API \(Roxygen2\)](#)
- [Spark SQL, Built-in Functions \(MkDocs\)](#)





# Apache Spark

```
import scala.Tuple2
import org.apache.spark.sql.SparkSession

def SPACE = ~" "
def spark = SparkSession.builder()
    .config("spark.master", "local")
    .appName("JavaWordCount")
    .getOrCreate()

def lines = spark.read().textFile('/path/to/peppers.txt').javaRDD()
def words = lines.flatMap(s -> SPACE.split(s).iterator())
def ones = words.mapToPair(s -> new Tuple2<>(s, 1))
def counts = ones.reduceByKey{ i1, i2 -> i1 + i2 }
def output = counts.collect()
for (tuple in output) {
    println tuple._1() + ": " + tuple._2()
}
spark.stop()
```



```
peppers: 4
where's: 1
if: 1
a: 3
picked: 4
Piper: 4
of: 4
Peter: 4
peck: 4
the: 1
pickled: 4
```






# Apache Beam Overview

Apache Beam is an open source, unified model for defining both batch and streaming data-parallel processing pipelines.

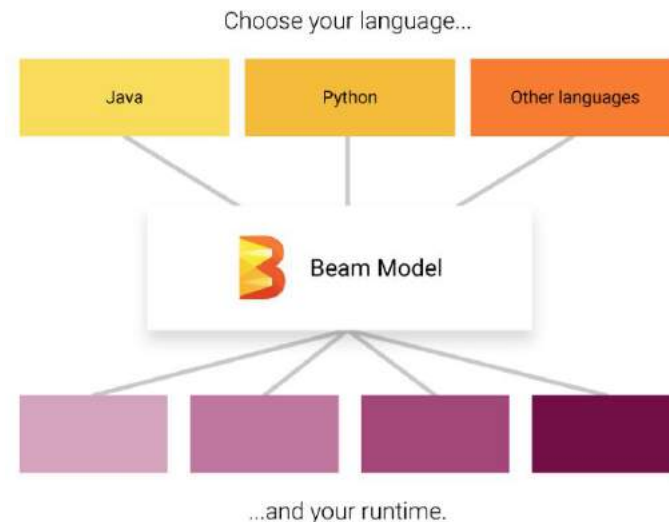
## Apache Beam Pipeline Runners

The Beam Pipeline Runners translate the data processing pipeline you define with your Beam program into the API compatible with the distributed processing back-end of your choice. When you run your Beam program, you'll need to specify an [appropriate runner](#) for the back-end where you want to execute your pipeline.

Beam currently supports Runners that work with the following distributed processing back-ends:

- Apache Apex 
- Apache Flink 
- Apache Gearpump (incubating) 
- Apache Samza 
- Apache Spark 
- Google Cloud Dataflow 

**Note:** You can always execute your pipeline locally for testing and debugging purposes.




## Apache Beam SDKs

The Beam SDKs provide a unified programming model that can represent and transform data sets of any size, whether the input is a finite data set from a batch data source, or an infinite data set from a streaming data source. The Beam SDKs use the same classes to represent both bounded and unbounded data, and the same transforms to operate on that data. You use the Beam SDK of your choice to build a program that defines your data processing pipeline.

Beam currently supports the following language-specific SDKs:

- Java 
- Python 
- Go 

A Scala  interface is also available as [Scio](#).



# Apache Beam Overview

Apache Beam is an open source, unified model for defining both batch and streaming data-parallel processing pipelines.

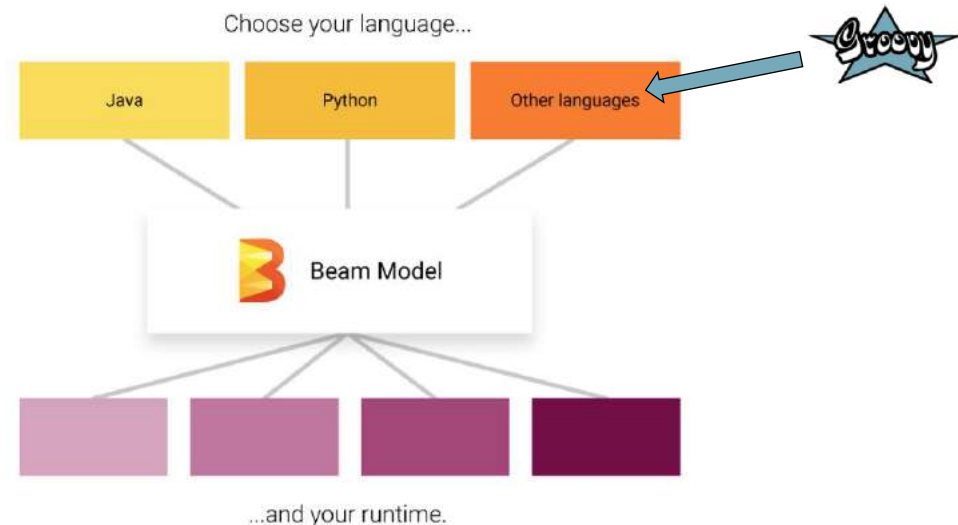
## Apache Beam Pipeline Runners

The Beam Pipeline Runners translate the data processing pipeline you define with your Beam program into the API compatible with the distributed processing back-end of your choice. When you run your Beam program, you'll need to specify an [appropriate runner](#) for the back-end where you want to execute your pipeline.

Beam currently supports Runners that work with the following distributed processing back-ends:

- Apache Apex 
- Apache Flink 
- Apache Gearpump (incubating) 
- Apache Samza 
- Apache Spark 
- Google Cloud Dataflow 




**Note:** You can always execute your pipeline locally for testing and debugging purposes.




## Apache Beam SDKs

The Beam SDKs provide a unified programming model that can represent and transform data sets of any size, whether the input is a finite data set from a batch data source, or an infinite data set from a streaming data source. The Beam SDKs use the same classes to represent both bounded and unbounded data, and the same transforms to operate on that data. You use the Beam SDK of your choice to build a program that defines your data processing pipeline.

Beam currently supports the following language-specific SDKs:

- Java 
  - Python 
  - Go 
- An arrow points from a star-shaped logo to the Java SDK.

A Scala  interface is also available as [Scio](#).

# Apache Beam Word Count – Quickstart guide

## Apache Beam Java SDK Quickstart

Source: <https://beam.apache.org/get-started/quickstart-java/>

This Quickstart will walk you through executing your first Beam pipeline to run `WordCount`, written using Beam's `Java SDK`, on a `runner` of your choice.

## Get the WordCount Code

The easiest way to get a copy of the `WordCount` pipeline is to use the following command to generate a simple Maven project that contains Beam's `WordCount` examples and builds against the most recent Beam release:

Unix

PowerShell

```
PS> mvn archetype:generate `
-D archetypeGroupId=org.apache.beam `
-D archetypeArtifactId=beam-sdks-java-maven-archetypes-examples `
-D archetypeVersion=2.13.0 `
-D groupId=org.example `
-D artifactId=word-count-beam `
-D version="0.1" `
-D package=org.apache.beam.examples `
-D interactiveMode=false
```

## Run WordCount

For Unix shells:

Direct

Apex

Flink-local

Flink-cluster

Spark

Dataflow

Samza-local

Nemo

Jet

```
$ mvn compile exec:java -Dexec.mainClass=org.apache.beam.examples.WordCount \
-Dexec.args="--inputFile=pom.xml --output=counts" -Pdirect-runner
```

# Apache Beam Word Count – Kata

See: <https://beam.apache.org/blog/2019/05/30/beam-kata-release.html>

The screenshot displays an IDE window for the Apache Beam Word Count Kata. The main editor shows the `Task.java` file with the following code:

```
public static void main(String[] args) {
    String[] lines = {
        "apple orange grape banana apple banana",
        "banana orange banana papaya"
    };

    PipelineOptions options = PipelineOptionsFactory.fromArgs(args).create();
    Pipeline pipeline = Pipeline.create(options);

    PCollection<String> wordCounts =
        pipeline.apply(Create.of(Arrays.asList(lines)));

    PCollection<String> output = applyTransform(wordCounts);

    output.apply(Log.ofElements());

    pipeline.run();
}

static PCollection<String> applyTransform(PCollection<String> input) {
    return input
        .apply(FlatMapElements.into(TypeDescriptors.strings())
            .via(line -> Arrays.asList(line.split(" "))))
        .apply(Count.perElement())
        .apply(ParDo.of(processElement(element, out) -> {
            out.output(element.getKey() + ":" + element.getValue());
        })));
}
```

The right-hand pane shows the **Word Count Pipeline** description:

**Kata:** Create a pipeline that counts the number of words.

Please output the count of each word in the following format:

```
word:count
ball:5
book:3
```

**Hint** ▼  
Refer to your katas above.

The bottom pane shows the test results for `TaskTest`:

```
Tests passed: 1 of 1 test - 1 s 87 ms
Task :examples-word_count-word_count:classes UP-TO-DATE
Task :util:compileTestJava UP-TO-DATE
Task :util:compileTestGroovy NO-SOURCE
Task :util:processTestResources NO-SOURCE
Task :util:testClasses UP-TO-DATE
Task :examples-word_count-word_count:compileTestJava UP-TO-DATE
Task :examples-word_count-word_count:compileTestGroovy NO-SOURCE
Task :examples-word_count-word_count:processTestResources NO-SOURCE
Task :examples-word_count-word_count:testClasses UP-TO-DATE
Task :examples-word_count-word_count:test
Deprecated Gradle features were used in this build, making it incompatible with Gradle 5.0.
See https://docs.gradle.org/4.8/userguide/command_line_interface.html#sec:command_line_warnings
BUILD SUCCESSFUL in 2s
6 actionable tasks: 2 executed, 4 up-to-date
9:20:46 AM: Tasks execution finished 'examples-word_count-word_count:cleanTest :examples-word_count-word_count:test --tests "org.apache.beam.learning.katas.examples.wordcount.TaskTestJ"'
```

The **Select Course** dialog box shows a list of available courses. The selected course is **Beam Katas - Java**. The dialog also displays the instructor's name, **Henry Suryawirawan**, and the Apache Beam website URL, <https://beam.apache.org/>. The dialog includes a search bar, a list of courses, and buttons for **Join** and **Cancel**.

# Apache Beam Word Count – GroovyConsole single file

```
import ...

static PCollection applyTransform(PCollection input) {
    ProcessFunction asWords = line -> line.split(" ").toList()

    def kv2out = new DoFn<KV, String>() {
        @ProcessElement
        void processElement(@Element KV element, OutputReceiver<String> out) {
            out.output(element.key + ":" + element.value)
        }
    }

    return input
        .apply(into(strings()).via(asWords))
        .apply(Count.perElement())
        .apply(ParDo.of(kv2out))
}

def lines = ["apple orange grape banana apple banana",
            "banana orange banana papaya"]

def pipeline = Pipeline.create()
def counts   = pipeline.apply(Create.of(lines))
def output   = applyTransform(counts)

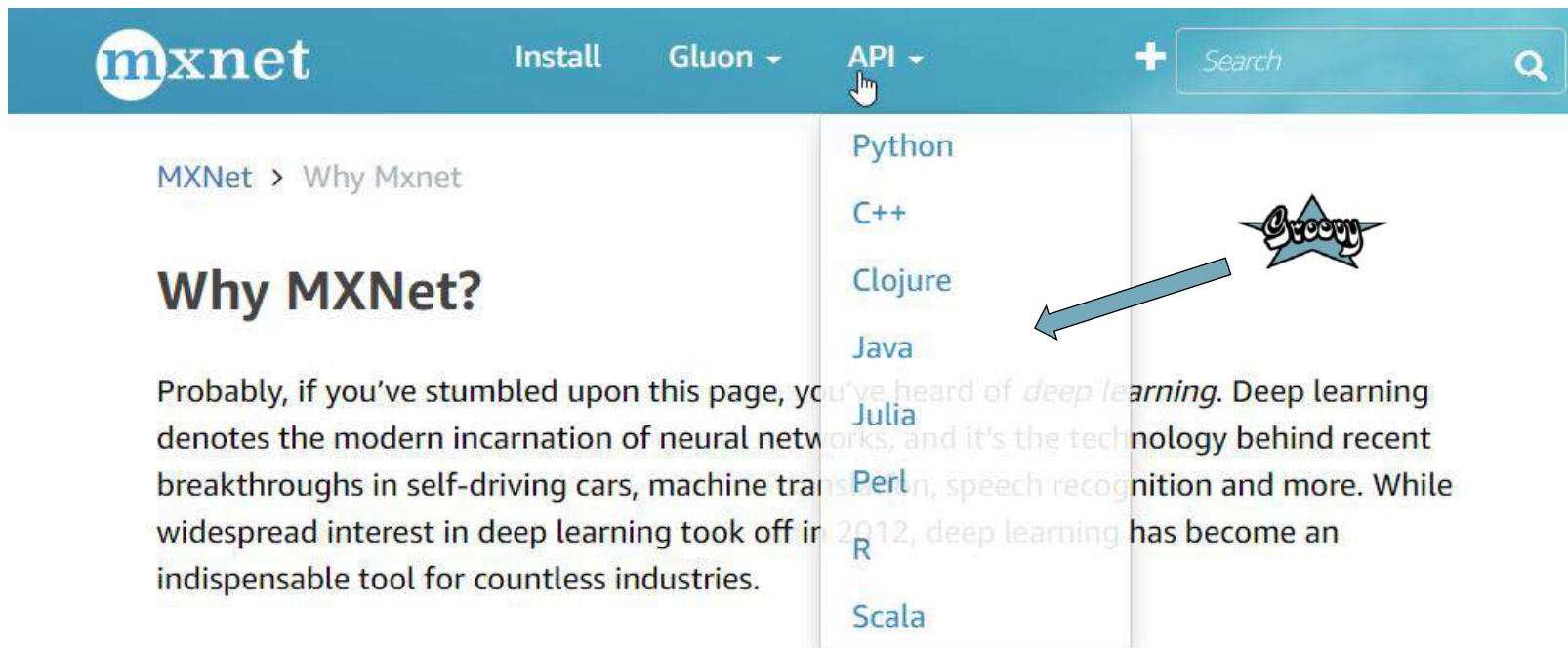
output.apply(Log.ofElements())

pipeline.run()
```

```
Jun 17, 2019 11:48:16 AM org.slf4j.Logger$info$0 call
INFO: banana:4
Jun 17, 2019 11:48:16 AM org.slf4j.Logger$info$2 call
INFO: orange:2
Jun 17, 2019 11:48:16 AM org.slf4j.Logger$info call
INFO: grape:1
Jun 17, 2019 11:48:16 AM org.slf4j.Logger$info$1 call
INFO: papaya:1
Jun 17, 2019 11:48:16 AM org.slf4j.Logger$info$3 call
INFO: apple:2
```



# Apache MXNET



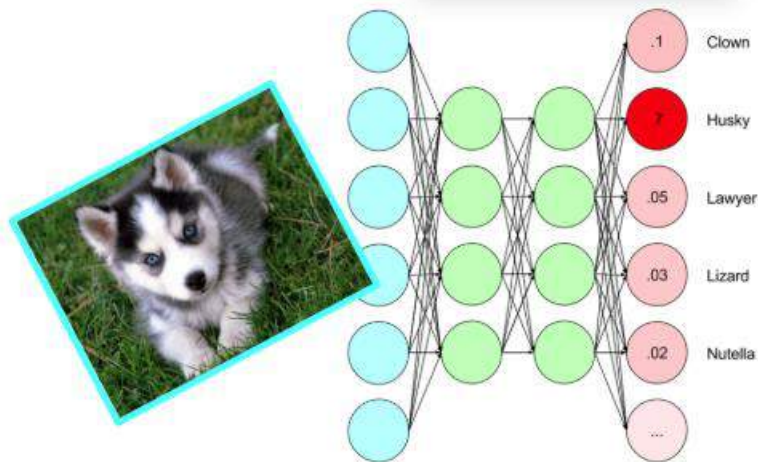
The screenshot shows the Apache MXNet website's navigation bar. The 'API' dropdown menu is open, listing programming languages: Python, C++, Clojure, Java, Julia, Perl, R, and Scala. A blue star with the text '@Droopy' and an arrow points to the 'Java' option. The main content area shows the breadcrumb 'MXNet > Why Mxnet' and the heading 'Why MXNet?'. Below the heading is a paragraph of text about deep learning.

mxnet Install Gluon API Python C++ Clojure Java Julia Perl R Scala

MXNet > Why Mxnet

## Why MXNet?

Probably, if you've stumbled upon this page, you've heard of *deep learning*. Deep learning denotes the modern incarnation of neural networks, and it's the technology behind recent breakthroughs in self-driving cars, machine translation, speech recognition and more. While widespread interest in deep learning took off in 2012, deep learning has become an indispensable tool for countless industries.



# Apache MXNET Object detection

```
import org.apache.mxnet.infer.javaapi.ObjectDetector
import org.apache.mxnet.javaapi.*

static void downloadUrl(String url, String filePath) {
    def destination = filePath as File
    if (!destination.exists()) {
        destination.bytes = new URL(url).bytes
    }
}

static downloadModelImage() {
    String tempDirPath = System.getProperty('java.io.tmpdir')
    println "tempDirPath: $tempDirPath"
    def imagePath = tempDirPath + "/inputImages/resnetssd/dog-ssd.jpg"
    String imgURL = "https://s3.amazonaws.com/model-server/inputs/dog-ssd.jpg"
    downloadUrl(imgURL, imagePath)
    def modelPath = tempDirPath + "/resnetssd/resnet50_ssd_model"
    println "Download model files, this can take a while..."
    String modelURL = "https://s3.amazonaws.com/model-server/models/resnet50_ssd/"
    downloadUrl(modelURL + "resnet50_ssd_model-symbol.json",
        tempDirPath + "/resnetssd/resnet50_ssd_model-symbol.json")
    downloadUrl(modelURL + "resnet50_ssd_model-0000.params",
        tempDirPath + "/resnetssd/resnet50_ssd_model-0000.params")
    downloadUrl(modelURL + "synset.txt",
        tempDirPath + "/resnetssd/synset.txt")
    [imagePath, modelPath]
}

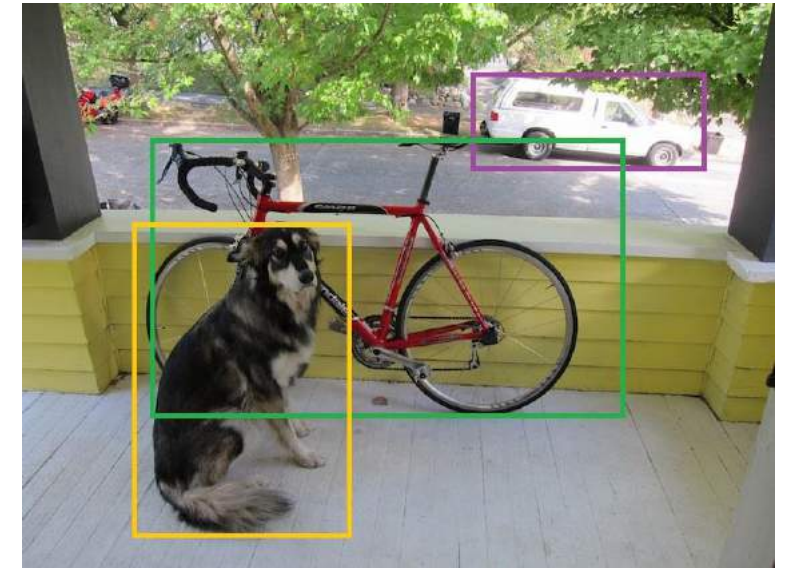
...
```





# Apache MXNET Object detection

```
...  
  
static detectObjects(String modelPath, String imagePath, inputShape) {  
    def context = [Context.cpu()]  
    def inputDescriptor = new DataDesc("data", inputShape, DType.Float32(), "NCHW")  
    def objDet = new ObjectDetector(modelPath, [inputDescriptor], context, 0)  
    return objDet.imageObjectDetect(ObjectDetector.loadImageFromFile(imagePath), 3)  
}  
  
def (imagePath, modelPath) = downloadModelImage()  
def (width, height) = [512, 512]  
Shape inputShape = new Shape([1, 3, width, height])  
def results = detectObjects(modelPath, imagePath, inputShape).sum()  
  
for (r in results) {  
    print "Class: $r.className"  
    printf " with probability: %.3f%n", r.probability  
    def coord = [r.XMin * width, r.XMax * height, r.YMin * width, r.YMax * height]  
    println "Coord: ${coord.collect { sprintf '%.2f', it }.join(', ')}\n"  
}
```



Class: car with probability: 0.998  
Coord: 312.21, 456.01, 72.03, 150.66

Class: bicycle with probability: 0.905  
Coord: 155.96, 383.84, 149.96, 418.95

Class: dog with probability: 0.823  
Coord: 83.82, 206.64, 179.14, 476.79

# Groovy Community Information



- [groovy.apache.org](http://groovy.apache.org)
- [groovy-lang.org](http://groovy-lang.org)
- [github.com/apache/groovy](https://github.com/apache/groovy)
- [groovycommunity.com](http://groovycommunity.com)
- [dev@groovy.apache.org](mailto:dev@groovy.apache.org)
- [users@groovy.apache.org](mailto:users@groovy.apache.org)
- [@ApacheGroovy](https://twitter.com/ApacheGroovy)
- [objectcomputing.com/training/catalog/groovy-programming/](http://objectcomputing.com/training/catalog/groovy-programming/)

# THANK YOU



Find me on twitter  
[@paulk\\_asert](https://twitter.com/paulk_asert)



OBJECT  
COMPUTING

## CONNECT WITH US



1+ (314) 579-0066



[@objectcomputing](https://twitter.com/objectcomputing)



[objectcomputing.com](https://objectcomputing.com)

# References

- <https://www.analyticsvidhya.com/blog/2017/02/introductory-guide-on-linear-programming-explained-in-simple-english/>
- Computational Mixed-Integer Programming  
<https://www.birs.ca/cmo-workshops/2018/18w5208/files/GleixnerAmbros.pdf>
- Integer Linear Programming: Graphical Introduction  
<https://www.youtube.com/watch?v=RhHhy-8sz-4>
- Integer Linear Programming: Binary Constraints  
<https://www.youtube.com/watch?v=-3my1TkyFiM>  
<https://www.youtube.com/watch?v=B3biWsBLeCw>  
<https://www.youtube.com/watch?v=MO8uQnIch6I>
- Linear Programming: Alternate solutions, Infeasibility, Unboundedness, & Redundancy  
<https://www.youtube.com/watch?v=eMA0LWsRQQ>