# Java as a fundamental working tool of the Data Scientist

Speaker : Alexey Zinoviev

# About



- I am a <graph theory, machine learning, traffic jams prediction, BigData algorythms> scientist
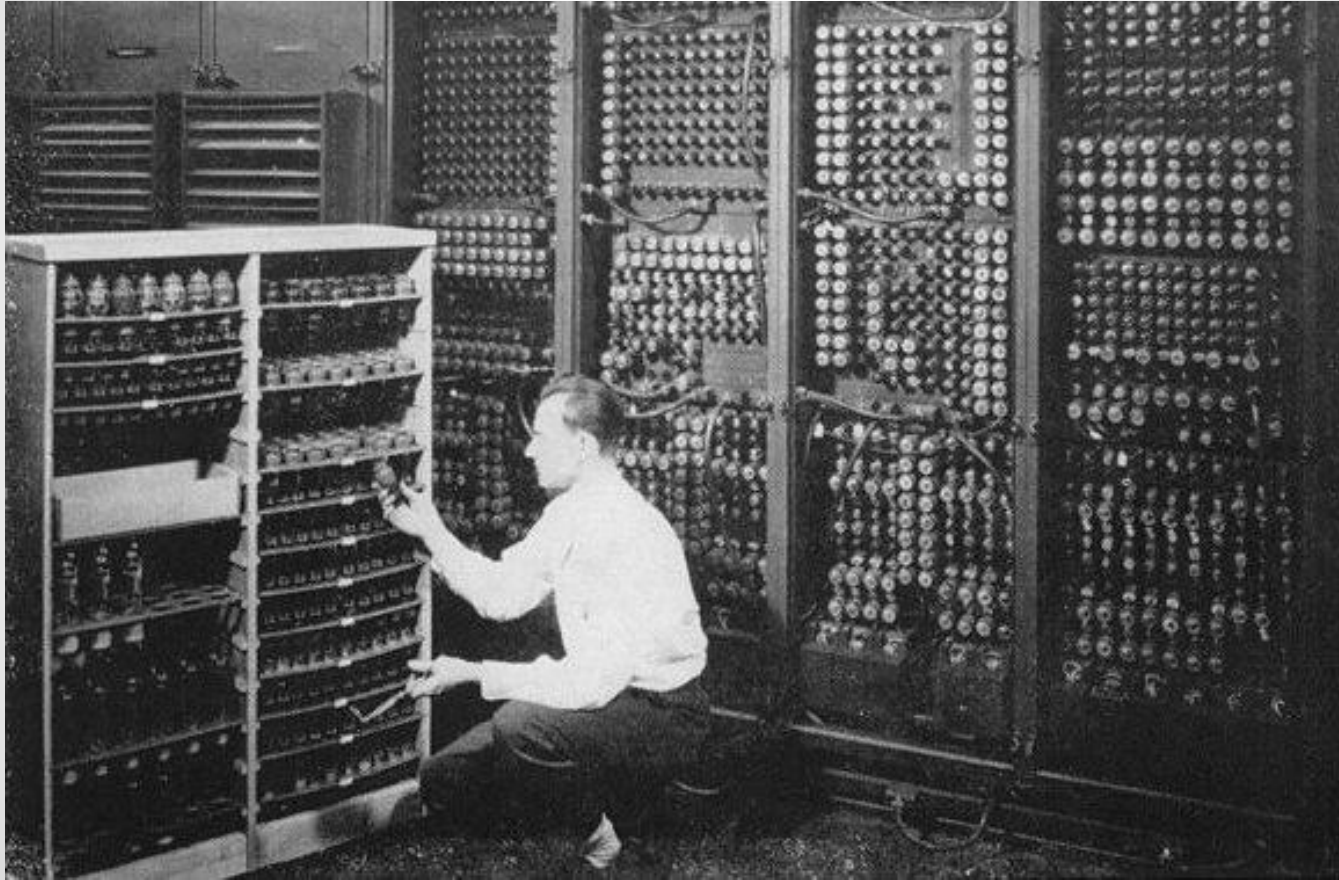- But I'm a <Java, JavaScript, Android, NoSQL, Hadoop, Spark> programmer

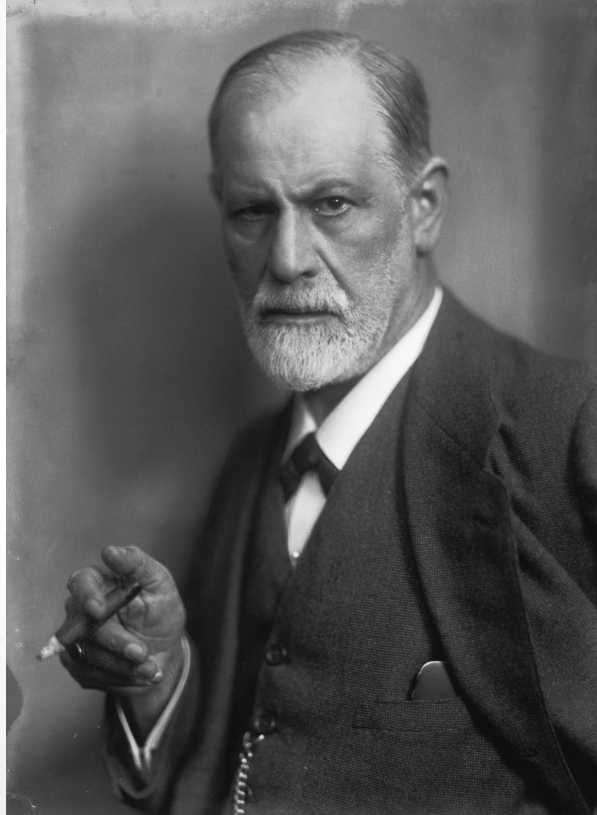# One of these fine days...

# We need in Python dev 'cause Data Mining

# You're a programmer, not an analyst

# Write your backends!

# Let's talk about it, Java-boy…

# Data mining



Mining coal
in your data

# Hey, man, predict me something!

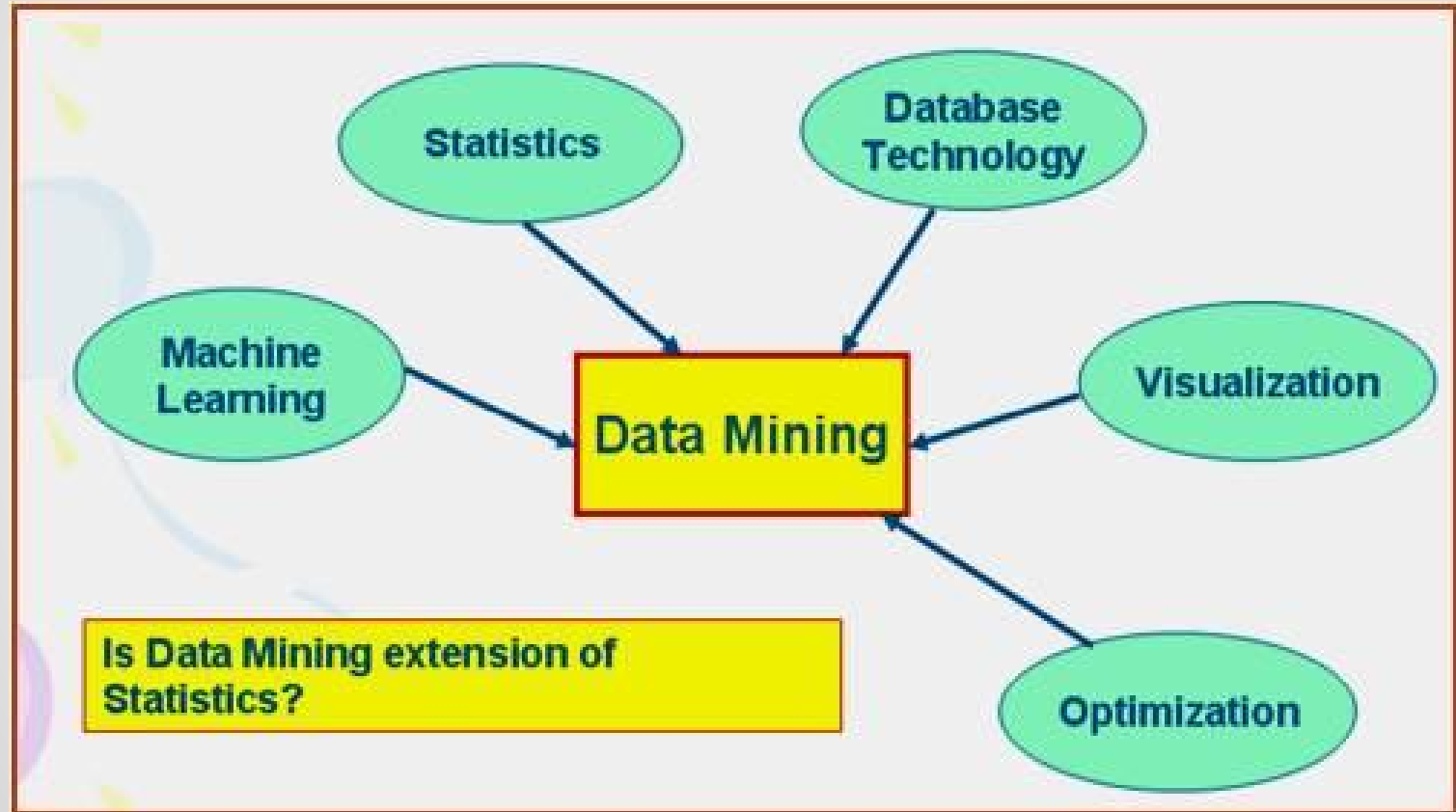# Man or sofa?

# Typical questions for DM

- Which loan applicants are high-risk?

- How do we detect phone card fraud?

- Which customers do prefer product A over product B?

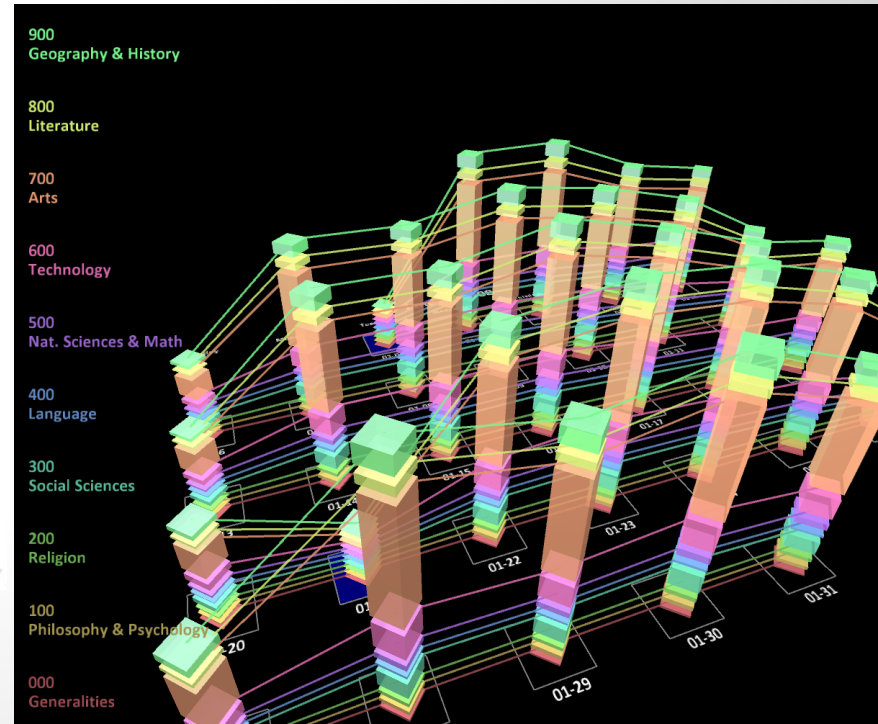- What is the revenue prediction for next year?
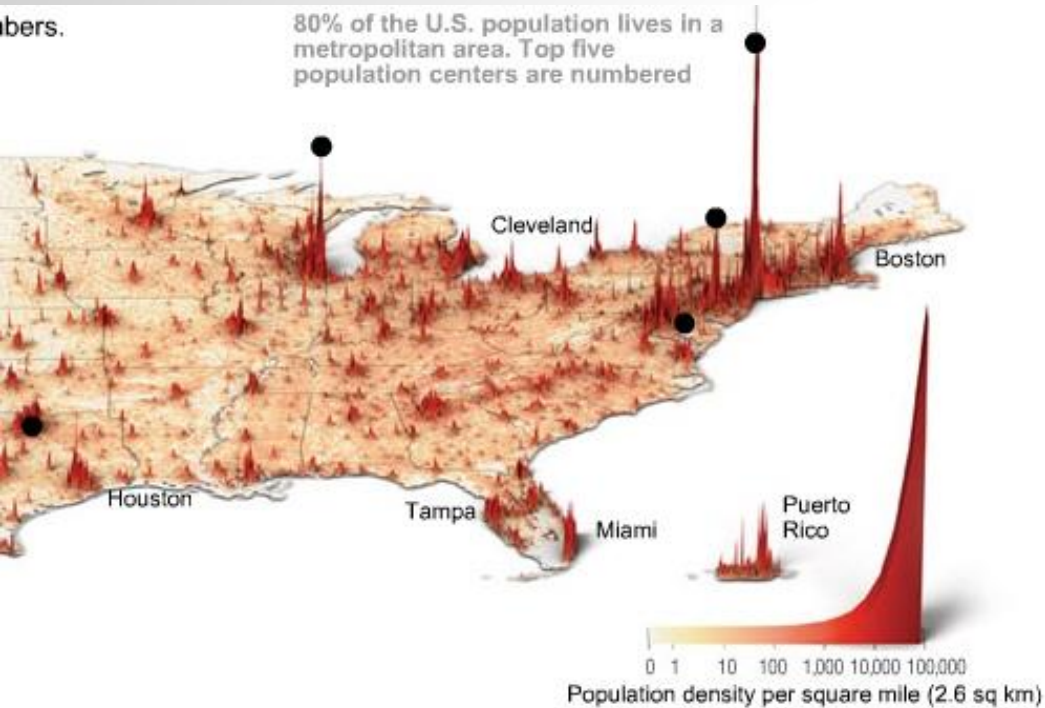
# What is Data Mining?

# Statistics?

# Tag cloud?

# Data visualization?

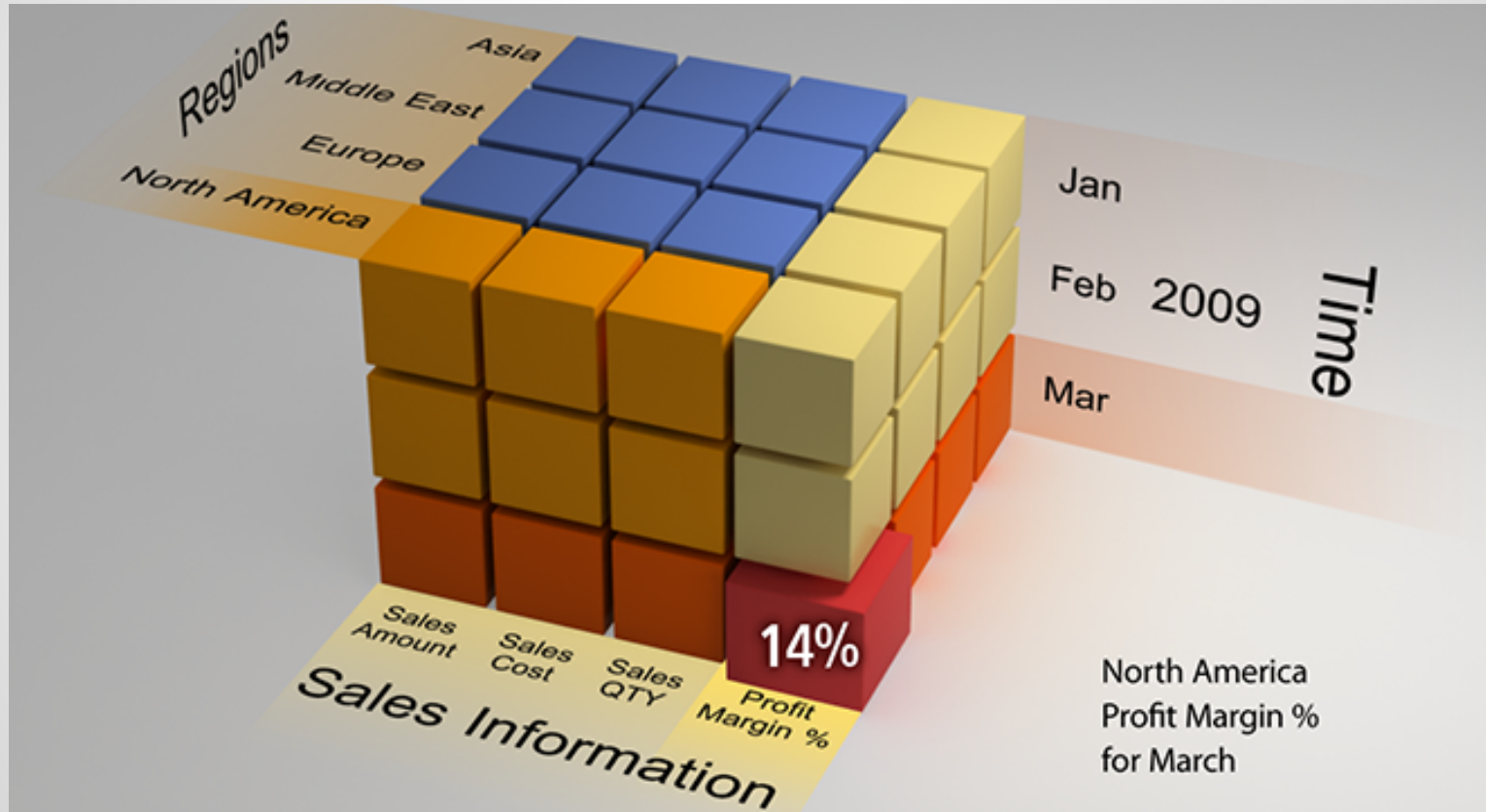# Not OLAP, 100%

# Magic part of KDD (Knowledge Discovery in Databases)

**Figure 1.** Overview of the steps constituting the KDD process

Selection → Target Data → Pre-processing → Preprocessed Data → Trans-formation → Transformed Data → Data Mining → Patterns → Interpretation/ Evaluation → Knowledge

Data

1. Selection
2. Pre-processing
3. Transformation
4. *Data Mining*
5. Interpretation/Evaluation

# How it really works



1. Share your date with us

2. Our magic manipulations

3. Building an answering machine

4. PROFIT!!!

# Data

20/72

# Data examples



- Facebook users, tweets

- Weather

- Sea routes

- Trade transactions

- Goverment

- Medicine (genomic data)

- Telecommuncations (phone call records)

# Data sources



- Relational Databases

- Data warehouses (Historical data)

- Files in CSV or in binary format

- Internet or electronic mails

- Scientific, research (R, Octave, Matlab)

# Target Data & Personal Data

## 23/72

# Pay with your personal data



- All your personal data (PD) are being deeply mined
- The industry of collecting, aggregating, and brokering PD is "database marketing."
- 1.1 billion browser cookies, 200 million mobile profiles, and an average of 1,500 pieces of data per consumer in Acxiom

# Preprocessing

26/72

Data Preparation

- Select small pieces
- Define default values for missed data
- Remove strange signals from data
- Merge some tables in one if required

# Pattern mining

28/72

# Association rule learning

# What is Cluster Analysis?



It is the process of finding model of function that describes and distinguishes data class to predict the class of objects whose class label is unknown.

# Regression



Simple Linear Regression
(with a *continuous* dependent [Y] variable)

- Statistical process for estimating the relationships among variables
- The estimation target is function (it can be probability distribution)
- Can be linear, polynomial, nonlinear and etc.

# Classification





- Training set of classified examples (supervised learning)
- Test set of non-classified items
- Main goal: find a function (classifier) that maps input data to a category
- Computer vision, drug discovery, speech recognition, biometric indentification, credit scoring

# Decision trees

# Decision trees

# kNN (k-nearest neighbor)

- There are two classes of objects A & B

- Define the class of new object, based on information about its neighbors

- Changing the boundaries of an new object area, we form a set of neighbors.

- New object is B becuase majority of the neighbors is a B.

15-Nearest Neighbor Classifier

# Skills & Tools

## 36/72

# Big Data Landscape



**Vertical Apps**
PREDICTIVE POLICING
bloomreach GET FOUND.
MYRRIX

**Log Data Apps**
splunk> loggly sumologic

**Ad/Media Apps**
rocketfuel
bluefin
Media Science
TURN
collective[i]
Recorded Future
LuckySort
DataXu
Data. Insight. Action.

**Data As A Service**
factual.
GNIP DATASIFT Windows Azure Marketplace INRIX LexisNexis SPACE CURVE
kaggle
knoema beta
LOQATE
Everything Location

**Business Intelligence**
ORACLE | Hyperion
SAP Business Objects RJMetrics
Microsoft | Business Intelligence
IBM COGNOS birst
MicroStrategy
Autonomy
QlikView bime DOMO
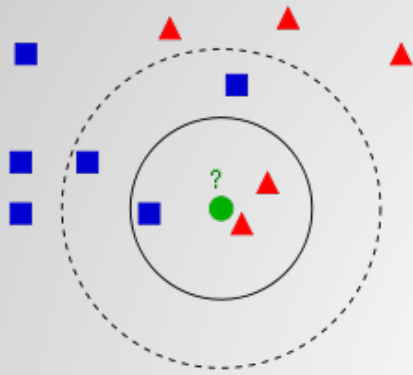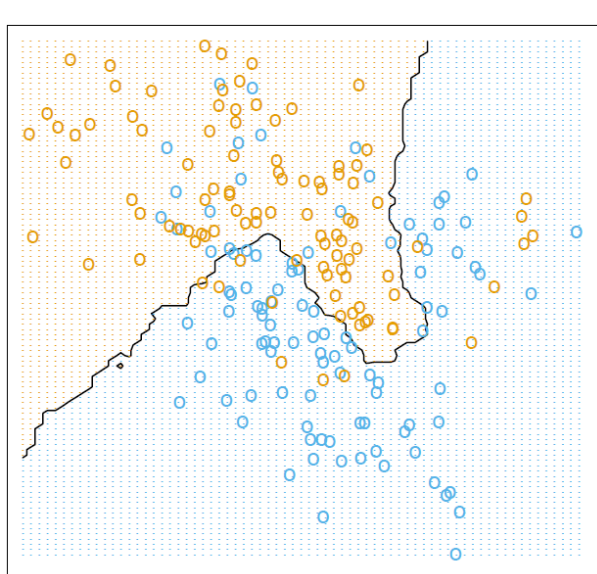Chart.io GoodData

**Analytics and Visualization**
tableau Palantir
OPERA SOLUTIONS metaLayer
METAMARKETS dataspora centrifuge
TERADATA ASTER
SAS TIBCO KARMASPHERE
panopticon Real-Time Visual Data Analysis
Datameer pentaho
platfora ClearStory CIRRO
alteryx visual.ly AYATA

**Analytics Infrastructure**
Hortonworks VERTICA An HP Company MAPR TECHNOLOGIES
Cloudera INFOBRIGHT PARACCEL
EMC² GREENPLUM
NETEZZA kognitio
DATASTAX EXASOL calpont

**Operational Infrastructure**
COUCHBASE 10gen The MongoDB company
TERADATA HADAPT
TERRACOTTA VoltDB
MarkLogic INFORMATICA

**Infrastructure As A Service**
amazon web services
Windows Azure
infochimps
Google BigQuery

**Structured Databases**
ORACLE MySQL
Microsoft SQL Server PostgreSQL
IBM DB2. SYBASE
memsql

**Technologies**
hadoop
hadoop MapReduce
mahout
APACHE HBASE
Cassandra

dave@vcdave.com

Bar chart of data mining technique usage:

- Decision Trees: 69%
- Regression: 68%
- Cluster Analysis: 60%
- Time Series: 32%
- Neural Nets: 31%
- Factor Analysis: 27%
- Text Mining: 26%
- Association Rules: 25%
- Ensemble Models: 22%
- Support Vector: 21%
- Bayesian: 21%
- Anomoly Detection: 16%
- Survival Analysis: 14%
- Rule Induction: 13%
- Social Network Analysis: 12%
- Genetic Algorithms: 11%
- Link Analysis: 9%
- Uplift Modeling: 9%
- MARS: 8%

**Consultants are more likely to use Ensemble Models**

| Corporate | Consultants | Academic | NGO / Gov't |
|---|---|---|---|
| 21% | 27% | 20% | 18% |

**Consultants and corporate data miners are more likely to use Uplift Modeling**

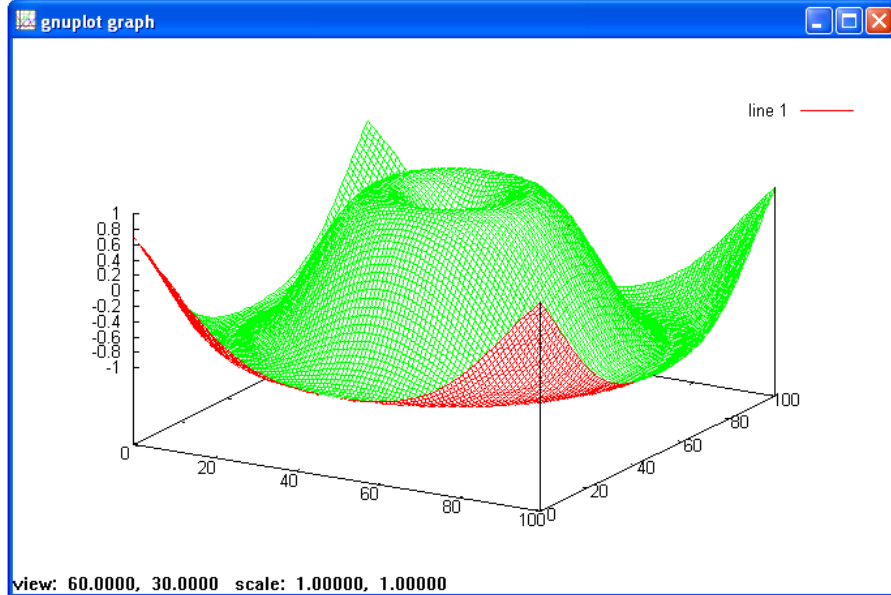| Corporate | Consultants | Academic | NGO / Gov't |
|---|---|---|---|
| 10% | 12% | 4% | 5% |

# Fashion Languages

File    Help

./
../
MinGW/
build/
gnuplot/
msys/

```
> [x,y]=meshgrid(-5:0.1:5);

> z=sin((x.^2+y.^2).^(1/2));

> mesh(z)

> A = [1 2 3; -2 1 5; 4 -1 1]
A =

    1    2    3
   -2    1    5
    4   -1    1


> b = 1:3
b =

    1    2    3


> x = A/b
x =

   1.00000
   1.07143
   0.35714

>
```
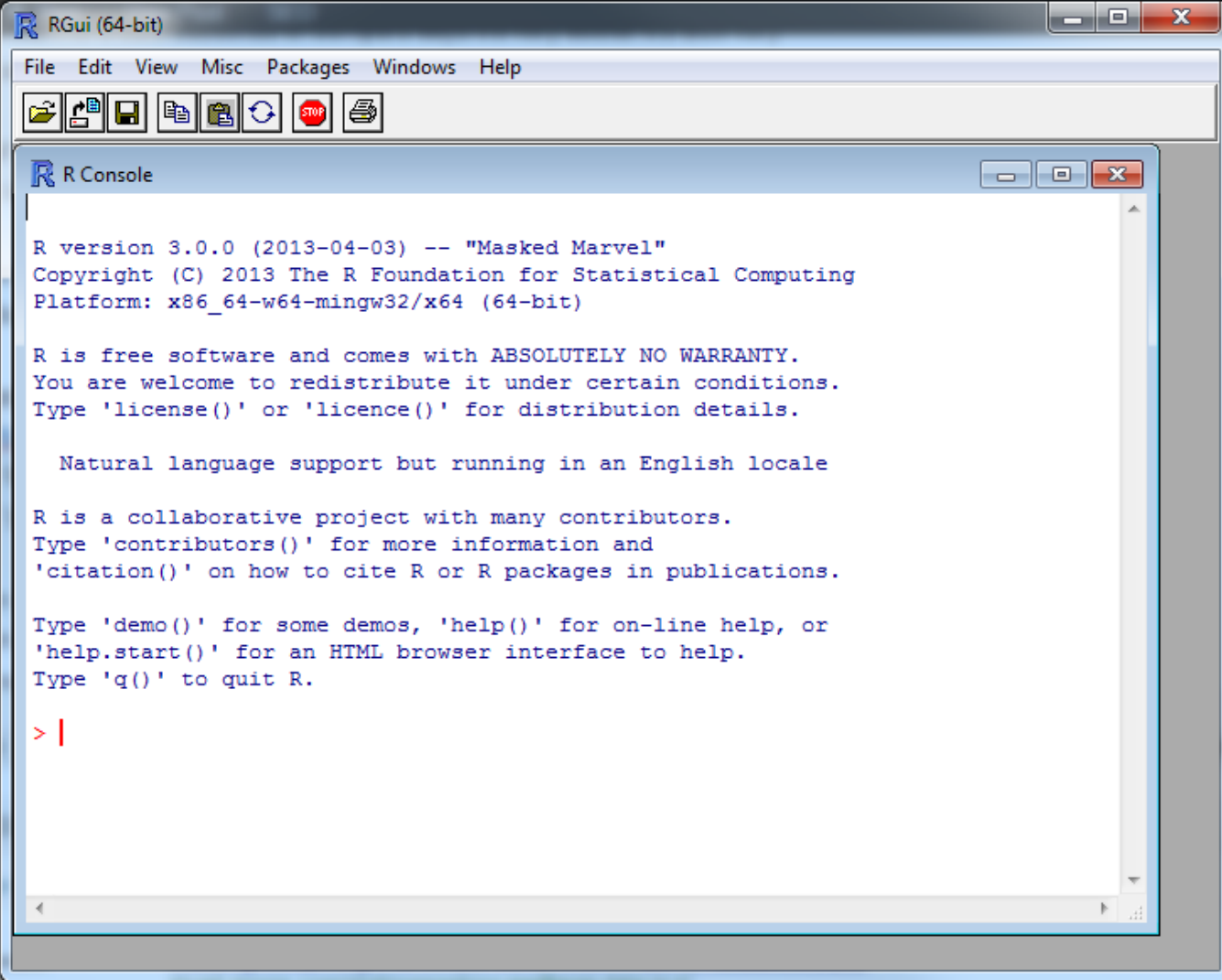
| Name | Dim | Size |
|------|-----|------|
| __nargin__ | 1x1 | 8 |
| __octfiledir__ | 1x77 | 77 |
| A | 3x3 | 72 |
| b | 1x3 | 24 |
| x | 3x1 | 24 |
| y | 101x101 | 81608 |
| z | 101x101 | 81608 |



gnuplot graph

line 1

view: 60.0000, 30.0000   scale: 1.00000, 1.00000

# Why not Octave?

- It's free
- Not full implemented stack of ML algorythms
- All your matrix are belong to us!
- Single thread model
- Java support

File   Edit   View   Misc   Packages   Windows   Help

R Console

```
R version 3.0.0 (2013-04-03) -- "Masked Marvel"
Copyright (C) 2013 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

  Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

>
```

# Why not R?

- 25% of R packages are written in Java

- Syntax is too sweet

- You should read 1000 lines in docs to write 1 line of code

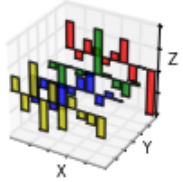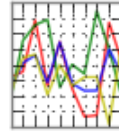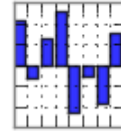- Single thread model for 95% algorythms

# Why not Python?

- Now Python is an idol for young scientists due to the low barrier to entry
- We are not Python developers
- High-level language
- Have you ever heard about a Jython?
- Long long way to real Highload production

# DM libraries in Python

# Java ecosystem

46/72

# JDM

- Java API for Data mining, JSR 73 and JSR 247

- **javax.datamining.supervised** defines the supervised function-related interfaces

- **javax.datamining.algorithm** contains all mining algorithm subclass packages

- JDM 2.0 adds Text Mining, Time series and so on..

# **Weka**



- Connectors to R, Octave, Matlab, Hadoop, NoSQL/SQL databases

- Source code of all algorythms in Java

- Preprocessing tools: discretization, normalization, resampling, attribute selection, transforming and combining

**Weka Explorer**

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

Classifier

Choose | J48 -C 0.25 -M 2

Test options
- Use training set
- Supplied test set | Set...
- Cross-validation | Folds | 10
- Percentage split

More op

(Nom) class

Start

Result list (right-click
13:54:50 - trees.J48

Classifier output

```
=== Stratified cross-validation ===
=== Summary ===
```

96    %
4    %

**Weka Classifier Tree Visualizer: 13:54:50 - trees.J48 (iris)**

Tree View

petalwidth
<= 0.6        > 0.6
Iris-setosa (50.0)        petalwidth
<= 1.7        > 1.7
petallength        Iris-virginica (46.0/1.0)

Iris-ver

**Weka Classifier Visualize: 13:54:50 - trees.J48 (iris)**

X: Instance_number (Num)        Y: sepallength (Num)
Colour: class (Nom)        Select Instance

Reset | Clear | Open | Save        Jitter

Plot: iris_predicted

7.9

6.1

4.3

0        74.5        149

| ROC Area | Class |
|---|---|
| 0.99 | Iris-setosa |
| 0.952 | Iris-versicolor |
| 0.961 | Iris-virginica |
| 0.968 | |

Class colour

Iris-setosa Iris-versicolor Iris-virginica

Status
OK

Log        x 0

# Weka + Hadoop

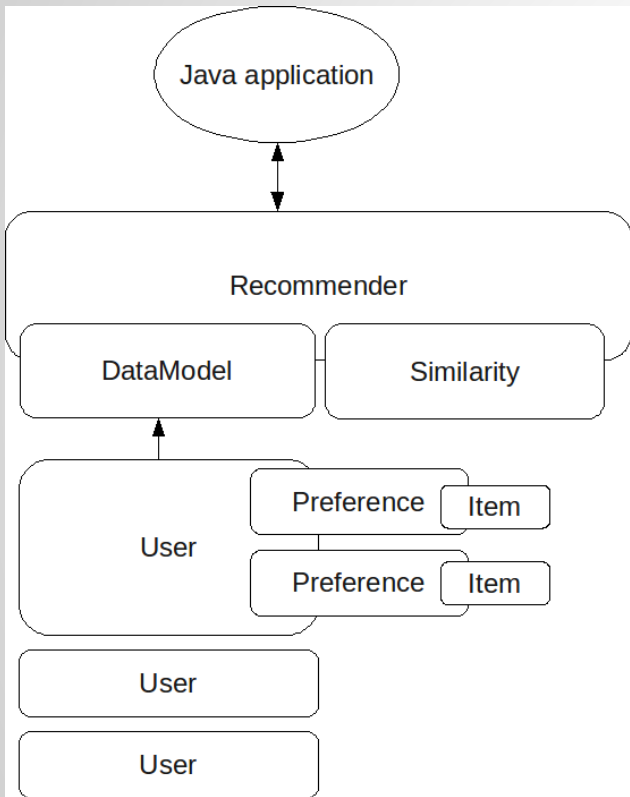# SPMF



Visitor count on SPMF website

- It's codebase of algorythms in pattern mining field
- It has [cool examples](#) and implementation of [78 algorythms](#)
- Cool performance [results](#) in specific area
- Codebase grows very fast

# **Mahout**

- Driven by Ng et al.'s [paper](#) "MapReduce for Machine Learning on Multicore"

- Next algorythms were adopted: Locally Weighted Linear Regression(LWLR), Naive Bayes (NB), k-means, Logistic Regression, Neural Network (NN), Principal Components Analysis (PCA), Support Vector Machine (SVM) and so on..

- The complexity was reduced in $n$ times for $n$ processors.
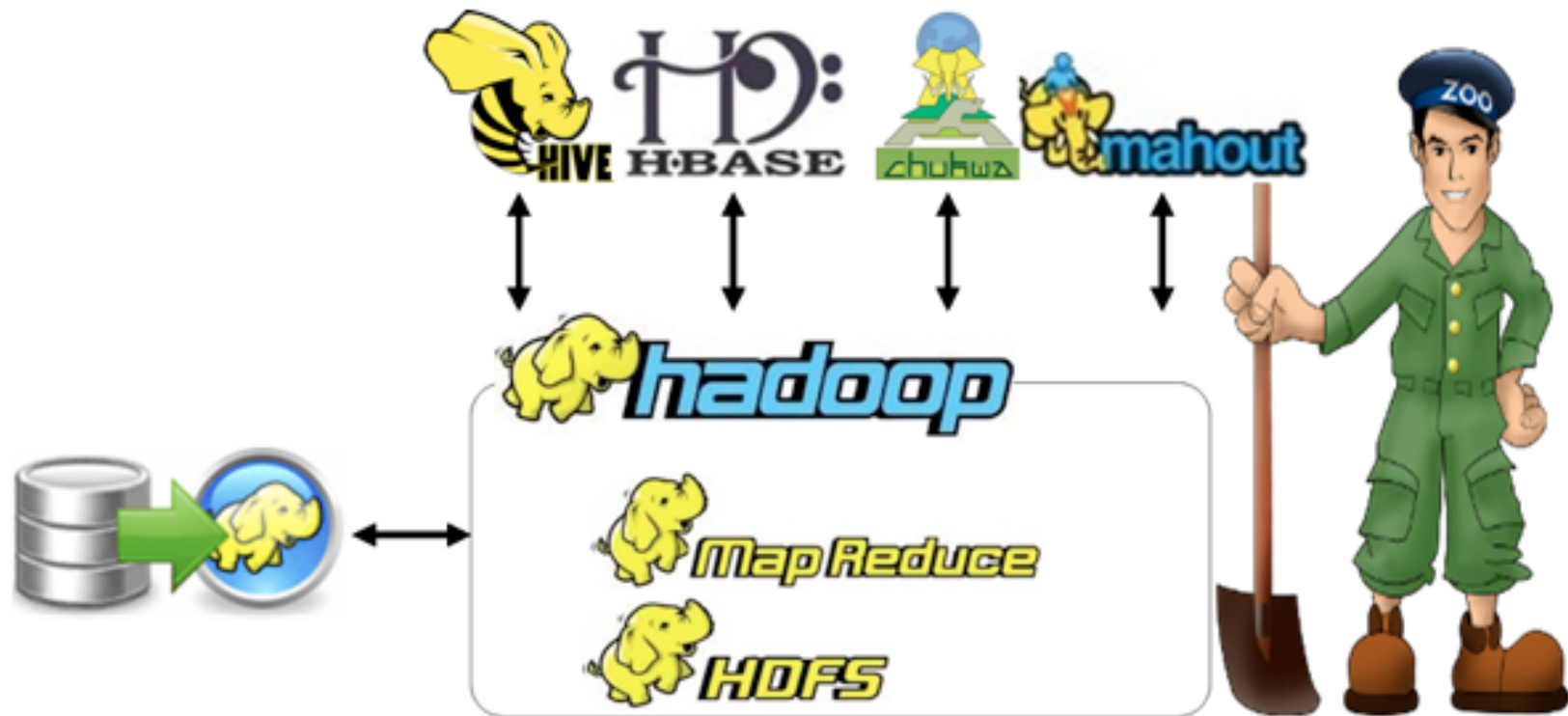
# Mahout



- DataModel (File, MySQL, PostgreSQL, Mongo, Cassandra)

- UserSimilarity
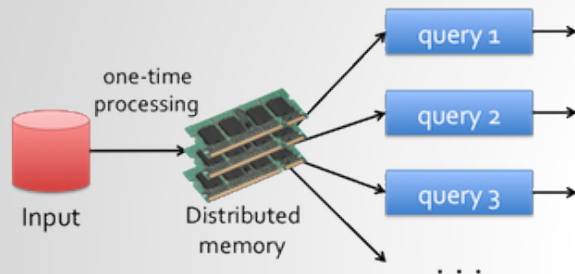
- ItemSimilarity
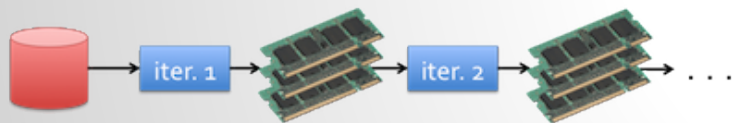
- UserNeighborhood

- Recommender

# **Mahout**

- Advanced Implementations of Java's Collections Framework for better Performance.

- Very close to Apache Giraph

- New algorythms will build on Spark platform

- [Spark shell](#)

- [Spring + Mahout demo](#)

- Collaborative Filtering, Classification, Clustering, Dimensionality Reduction, Miscellaneous are supported

# Spark



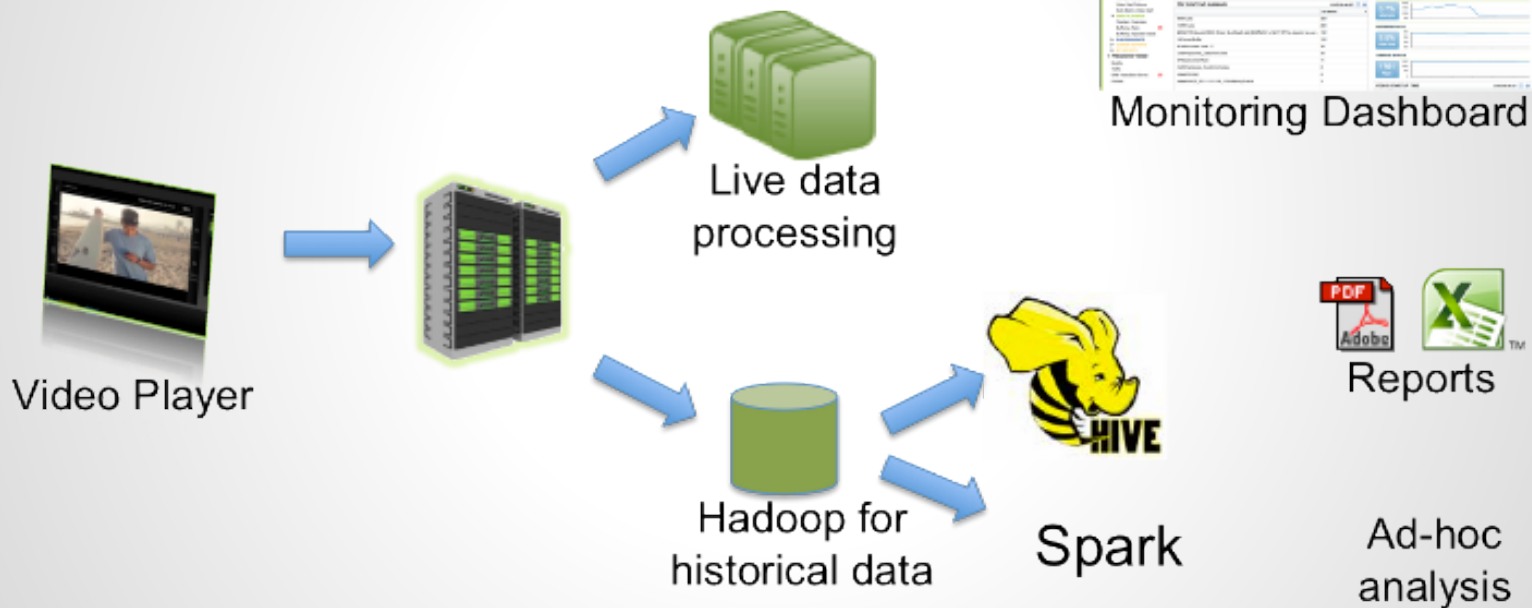(a) Low-latency computations (queries)

(b) Iterative computations

- MapReduce in memory

- Up to 50x faster than Hadoop

- Support for Shark (like Hive), MLlib (Machine learning), GraphX (graph processing)

- RDD is a basic building block (immutable distributed collections of objects)

# Spark



Video Player

Live data processing

Hadoop for historical data

Spark

Monitoring Dashboard

Reports

Ad-hoc analysis

# Spark + Java 8

**Java 7 search example:**

```
JavaRDD<String> lines = sc.textFile("hdfs://log.txt").filter(
  new Function<String, Boolean>() {
    public Boolean call(String s) {
      return s.contains("Tomcat");
    }
});
long numErrors = lines.count();
```

**Java 8 search example:**

```
JavaRDD<String> lines = sc.textFile("hdfs://log.txt")
                .filter(s -> s.contains("Tomcat"));
long numErrors = lines.count();
```
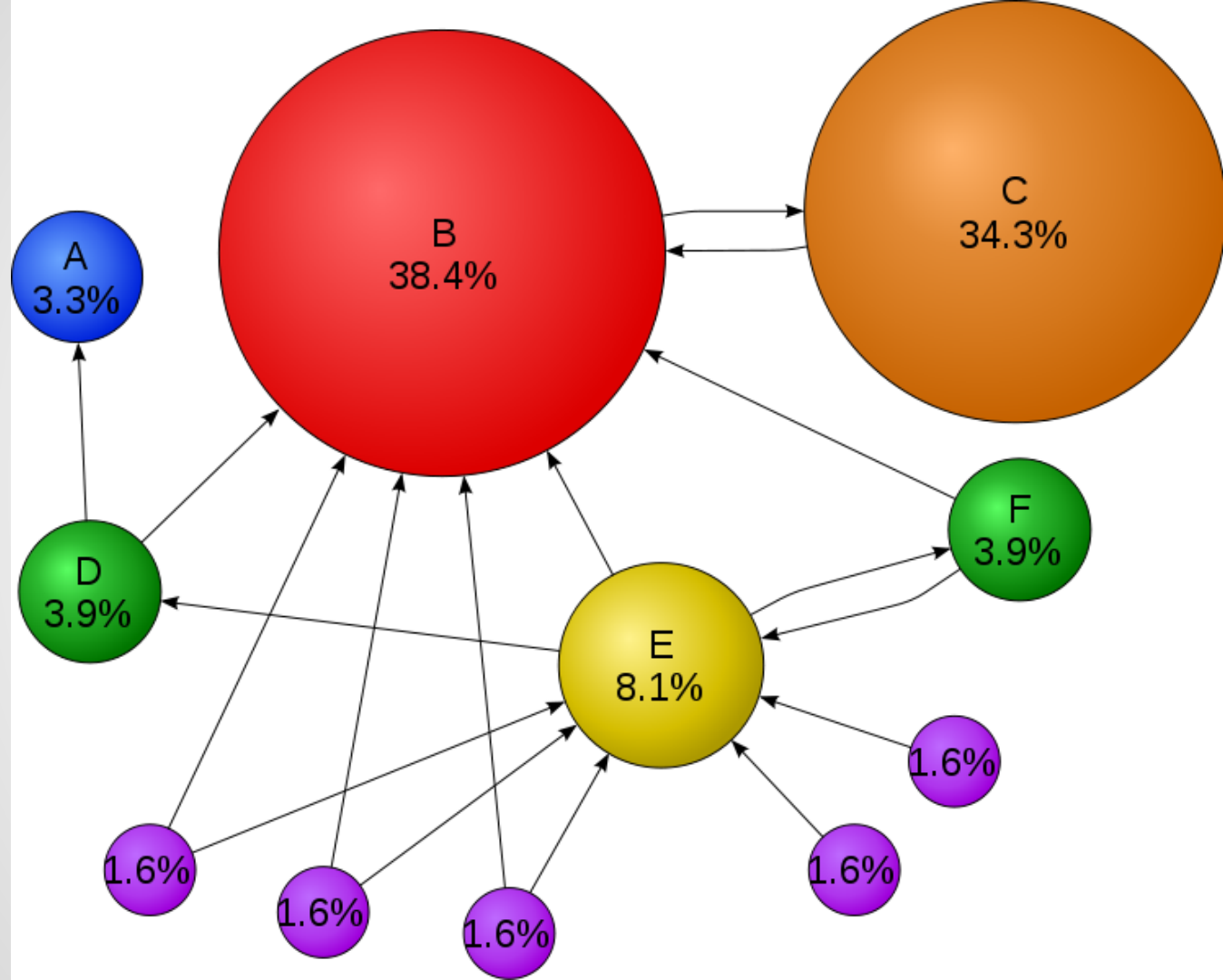
# Mahout's killer

# MLlib

- Classification and regression. collaborative filtering and clustering, Dimensionality reduction and Optimization are supported

- It extends scikit-learn (Python lib) and Mahout and run on Spark

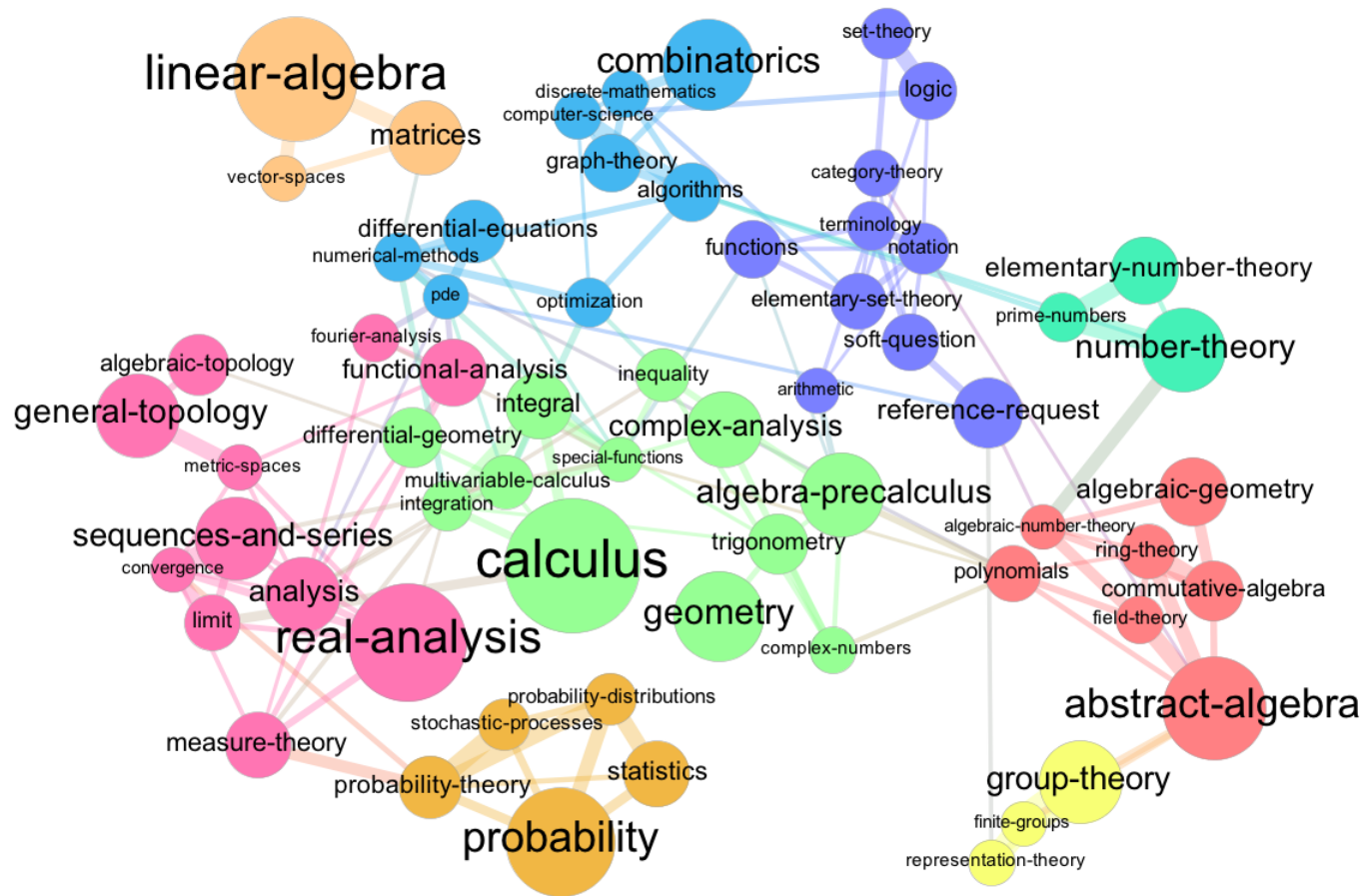- Well documented and integrated with many Java solutions

| Size | Classification | Tools |
|------|----------------|-------|
| **Lines**<br>Sample Data | Analysis and Visualization | Whiteboard, bash |
| **KBs - low MBs**<br>Prototype Data | Analysis and Visualization | Matlab, Octave, R |
| **MBs - low GBs**<br>Online Data | Storage | MySQL (DBs) |
| **MBs - low GBs**<br>Online Data | Analysis | NumPy, SciPy, Weka, BLAS/LAPACK |
| **GBs - TBs - PBs**<br>BigData | Storage | HDFS, HBase, Cassandra |
| **GBs - TBs - PBs**<br>Big Data | Analysis | Hive, Mahout, Hama, Giraph,MLlib |

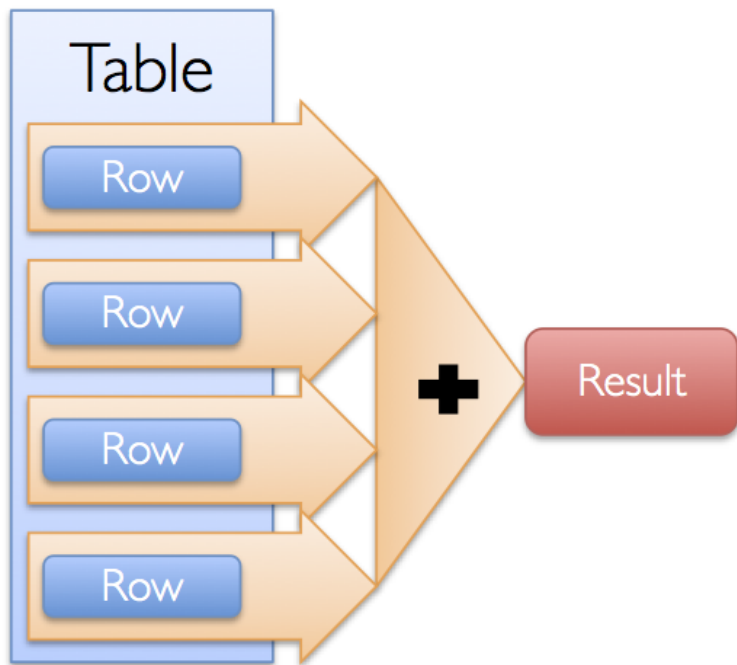# Large graph processing tools

63/72

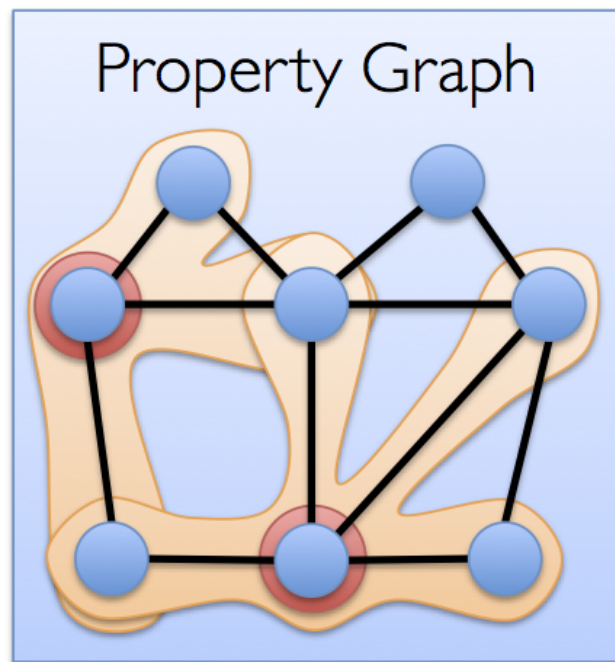| Graph | Number of vertexes | Number of edges | Volume | Data/per day |
|---|---|---|---|---|
| Web-graph | 1,5 * 10^12 | 1,2 * 10^13 | 100 PB | 300 TB |
| Facebook (friends graph) | 1,1 * 10^9 | 160 * 10^9 | 1 PB | 15 TB |
| Road graph of EU | 18 * 10^6 | 42 * 10^6 | 20 GB | 50 MB |
| Road graph of this city | 250 000 | 460 000 | 500 MB | 100 KB |

# MapReduce for iterative calculations

- High complexity of graph problem reduction to key-value model
- Iteration algorythms, but multiple chained jobs in M/R with full saving and reading of each state

*Think like a vertex...*

# C++ API

```cpp
template <typename VertexValue,
          typename EdgeValue,
          typename MessageValue>
class Vertex {
 public:
  virtual void Compute(MessageIterator* msgs) = 0;

  const string& vertex_id() const;
  int64 superstep() const;

  const VertexValue& GetValue();
  VertexValue* MutableValue();
  OutEdgeIterator GetOutEdgeIterator();

  void SendMessageTo(const string& dest_vertex,
                     const MessageValue& message);
  void VoteToHalt();
};
```
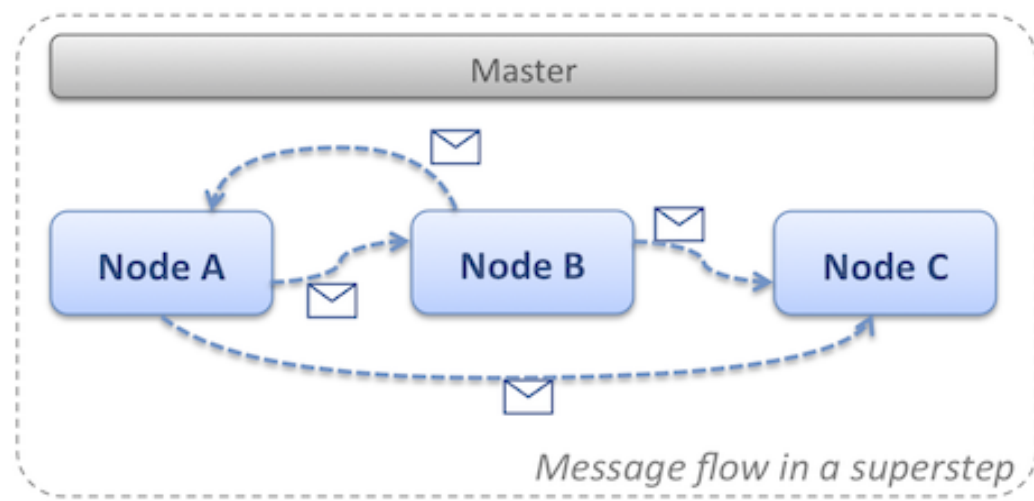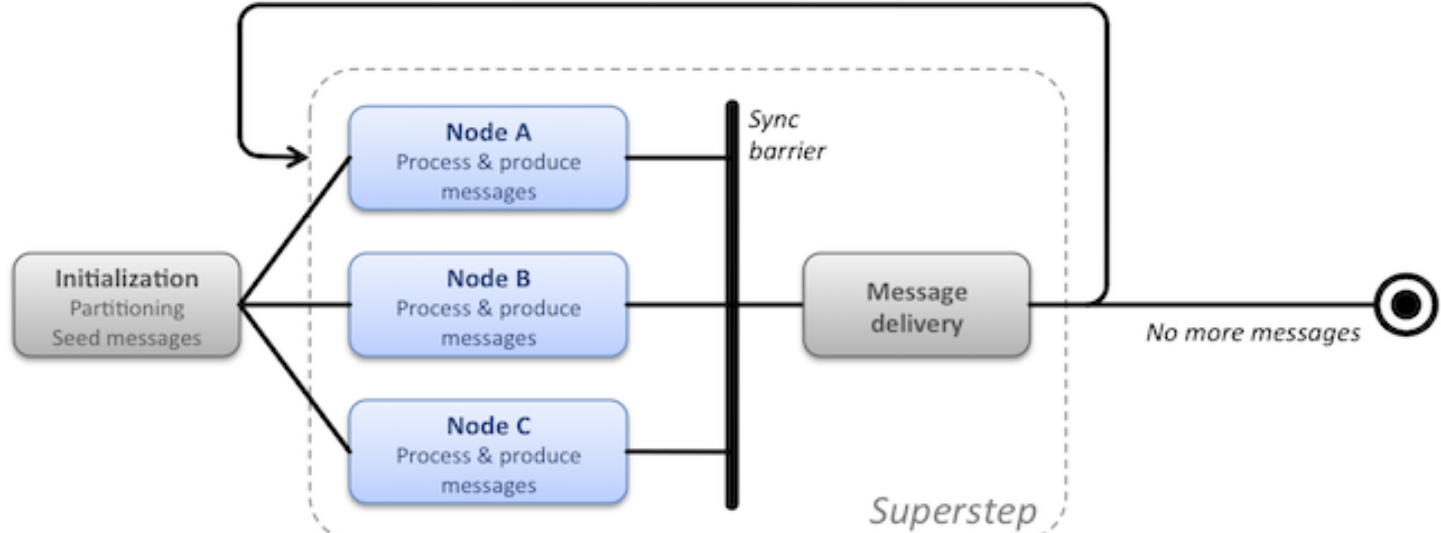
Superstep

Message flow in a superstep

# Books and papers

- "Mahout in Action", Owen et. al., Manning Pub.

- "Pattern Recognition and Machine Learning", Christopher Bishop, Springer Pub.

- "Elements of Statistical Learning: Data Mining, Inference, and Prediction", Hastie et. al., Springer Pub.

- "Collective Intelligence in Action" Satnam Alag et. al., Manning Pub.

# Your questions?