**Data Mining**

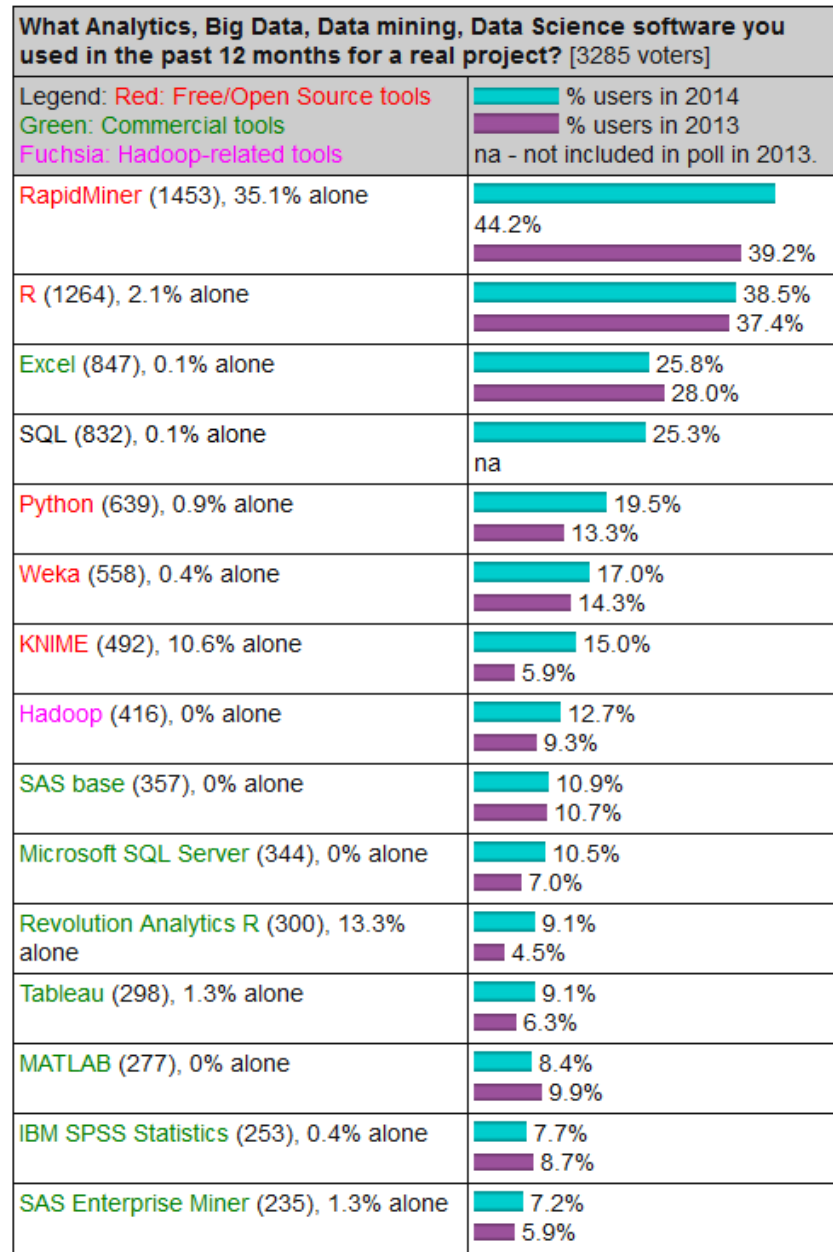# Introduction to RapidMiner

# Outline

1. What is RapidMiner?

2. RapidMiner Terminology

3. RapidMiner User Interface

    1. The Design Perspective

    2. The Results Perspective

4. The RapidMiner Repository

5. Preprocessing Operators

6. Data Visualization

7. RapidMiner Resources

# RapidMiner

- A very comprehensive open-source data mining tool

  - The data mining process is visually modeled as an operator chain

  - RapidMiner has over 400 build in data mining operators

  - RapidMiner provides broad collection of charts for visualizing data

- Project started in 2001 by Ralf Klinkenberg, Ingo Mierswa, and Simon Fischer at University of Dortmund, Germany

- Today: Maintained by commercial company plus open-source developers

- RapidMiner Editions

  - Community Edition: Free
    (= Second Last Edition)

  - Enterprise Edition: Commercial
    (= Last Edition plus professional support)

# KDnuggets Poll 2014 (and 2013)

| What Analytics, Big Data, Data mining, Data Science software you used in the past 12 months for a real project? [3285 voters] | | |
|---|---|---|
| Legend: Red: Free/Open Source tools<br>Green: Commercial tools<br>Fuchsia: Hadoop-related tools | % users in 2014<br>% users in 2013<br>na - not included in poll in 2013. | |
| RapidMiner (1453), 35.1% alone | 44.2% | 39.2% |
| R (1264), 2.1% alone | 38.5% | 37.4% |
| Excel (847), 0.1% alone | 25.8% | 28.0% |
| SQL (832), 0.1% alone | 25.3% | na |
| Python (639), 0.9% alone | 19.5% | 13.3% |
| Weka (558), 0.4% alone | 17.0% | 14.3% |
| KNIME (492), 10.6% alone | 15.0% | 5.9% |
| Hadoop (416), 0% alone | 12.7% | 9.3% |
| SAS base (357), 0% alone | 10.9% | 10.7% |
| Microsoft SQL Server (344), 0% alone | 10.5% | 7.0% |
| Revolution Analytics R (300), 13.3% alone | 9.1% | 4.5% |
| Tableau (298), 1.3% alone | 9.1% | 6.3% |
| MATLAB (277), 0% alone | 8.4% | 9.9% |
| IBM SPSS Statistics (253), 0.4% alone | 7.7% | 8.7% |
| SAS Enterprise Miner (235), 1.3% alone | 7.2% | 5.9% |

**Source**:
http://www.kdnuggets.com/2014/06/kdnuggets-annual-software-poll-rapidminer-continues-lead.html

# RapidMiner Terminology

– Concept: Class of examples, e.g. *person, flower*

– Example: Instance of a class, e.g. *Franz Müller*

– Example Set: The collection of all examples, that are mined.

– Attributes: Properties, e.g. *name, address, age, size*

   • Attributes have Value Types (see next slides)

   • Attributes have Roles (see next slides)

– Process: The complete data transformation and mining process.

– Model: The result of applying a process, e.g. *a classification model or set of rules*.

# Value Types in RapidMiner

| Value Type | RapidMiner Name | Description |
| --- | --- | --- |
| Nominal | nominal | Categorical non-numerical values, usually used for finite quantities of different characteristics |
| 2-value nominal | binominal | Special case of nominal, where only two different values are permitted |
| Multi-value nominal | polynominal | Special case of nominal, where more than two different values are permitted |
| Numerical values | numeric | For numerical values in general |
| Integers | integer | Whole numbers, positive and negative |
| Real numbers | real | Real numbers, positive and negative |
| Date time | date_time | Date as well as time |
| Date | date | Only date |
| Time | time | Only time |
| Text | text | Random free text without structure |

# Attribute Roles in RapidMiner

- Special attributes

  - Id: A unique identifier of examples, e.g. *matriculation number*

  - Cluster: The assigned cluster of an example, e.g. *after clustering*

  - Label: The class attribute, e.g. *for classification tasks*

  - Prediction: The predicted label, e.g. *after applying a classification model*

- Regular Attributes

  - Regular attributes are used to describe characteristics of an example

# The Design Perspective

# Specifying a Process by Chaining Operators

# The Results Perspective

– Perspective for displaying and visualizing Example Sets and Models.



Data Statistics View

Raw Data View

Charts

# The Repository

- Repository types

  - Local (files on your machine)

  - Remote (in the Cloud or own RapidMiner server)

- Location to store:

  - Datasets and its meta-data

  - Processes

- Data Import/Export Formats

  - CSV and XML Files

  - Excel Spreadsheets

  - Microsoft Access Tables

  - ARFF Files (Weka file format)

  - ...

# (Some common) Preprocessing Operators

- Type and Role Conversions

    - Type: Provides for converting between different attribute types

    - Role: Provides for setting attribute roles (e.g. label, id)

- Attribute Set Reduction and Transformation

    - Select Attributes: Selects a sub set of the available attributes

    - Generate Attribute: Generates new attribute from existing attributes

- Value Transformations

    - Normalize: Normalizes attribute to certain range

- Filtering

    - Filter Examples: Filters the examples for a give criterion

- Aggregation

    - Performs a SQL-like aggregation (count, sum, ...)
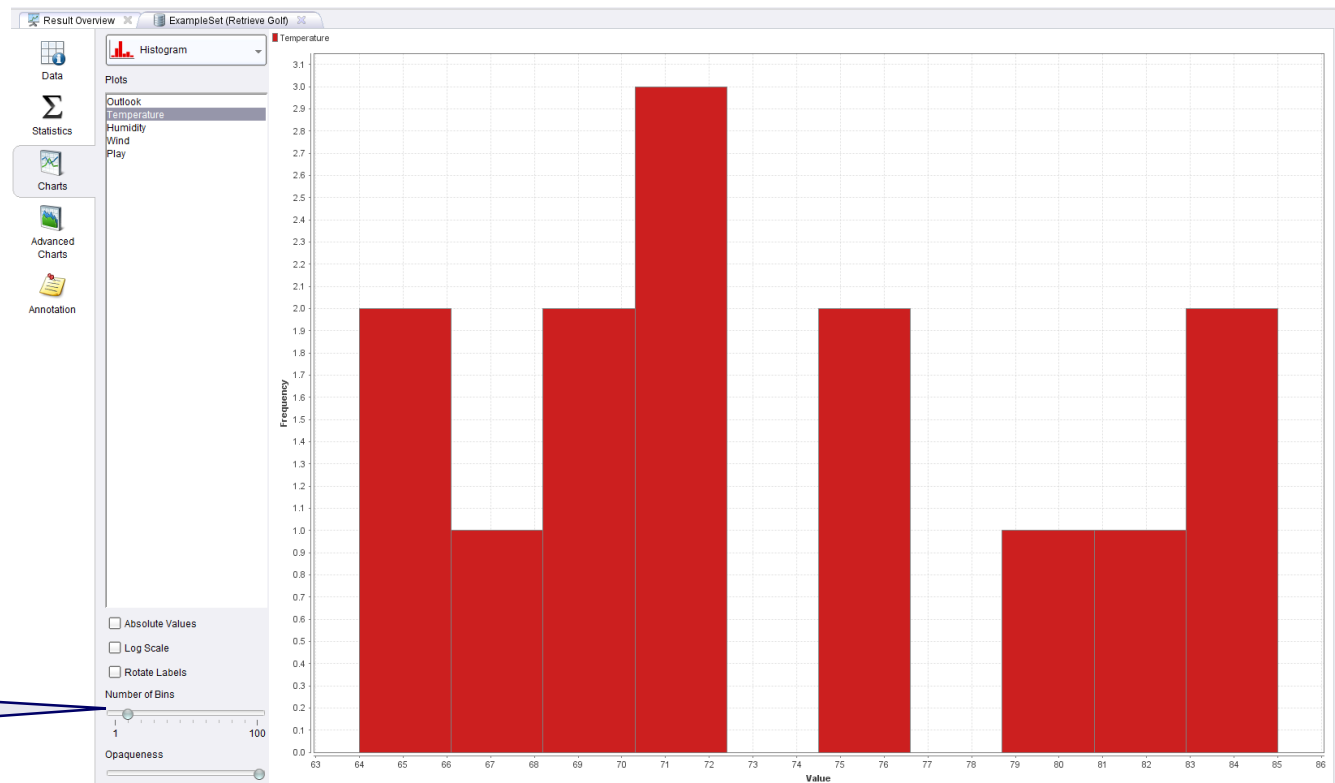
# Data Visualization

- Visualization of data is one of the most powerful and appealing techniques for data exploration

  - Humans have a well developed ability to analyze large amounts of information that is presented visually

  - Can detect general patterns and trends

  - Can detect outliers and unusual patterns

**Visualization is the conversion of data into a visual format so that the characteristics of the data and the relationships among data items or attributes can be analyzed.**

# Visualization Techniques: Histogram

- Usually used to display the distribution of values of a single attribute

  - Divide the values into bins and show a bar plot of the number of objects in each bin

  - The height of each bar indicates the number of objects per bin

  - Shape of histogram depends on the number of bins
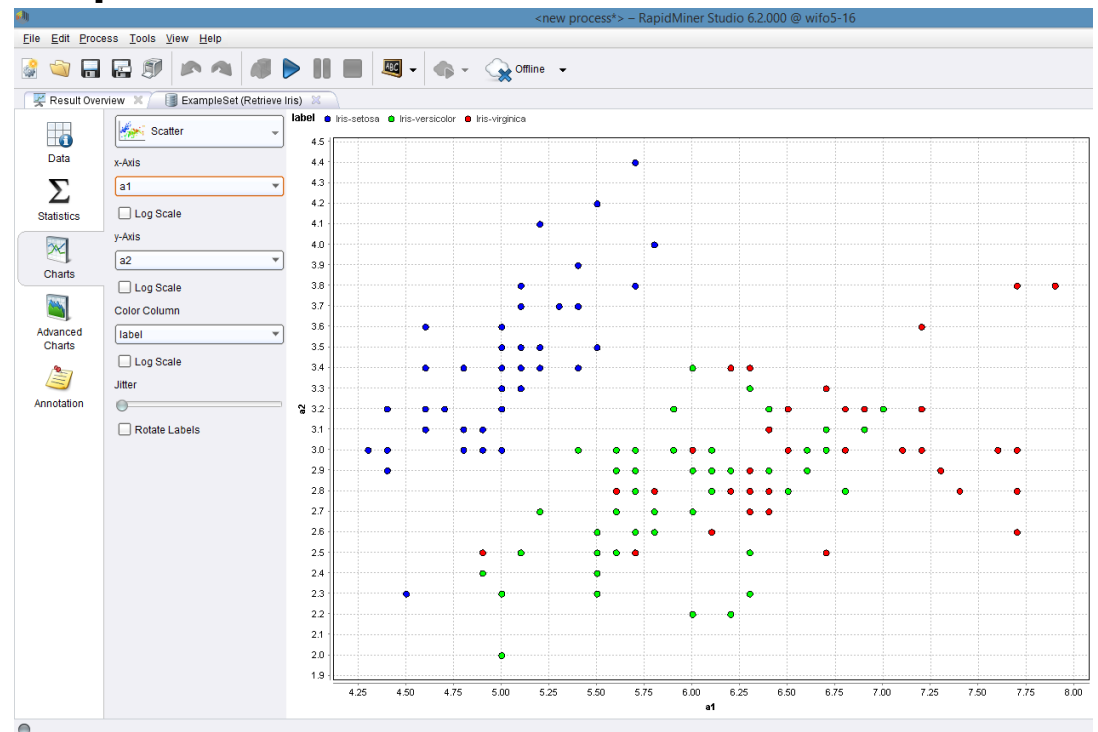
**RapidMiner Chart: Histogram**
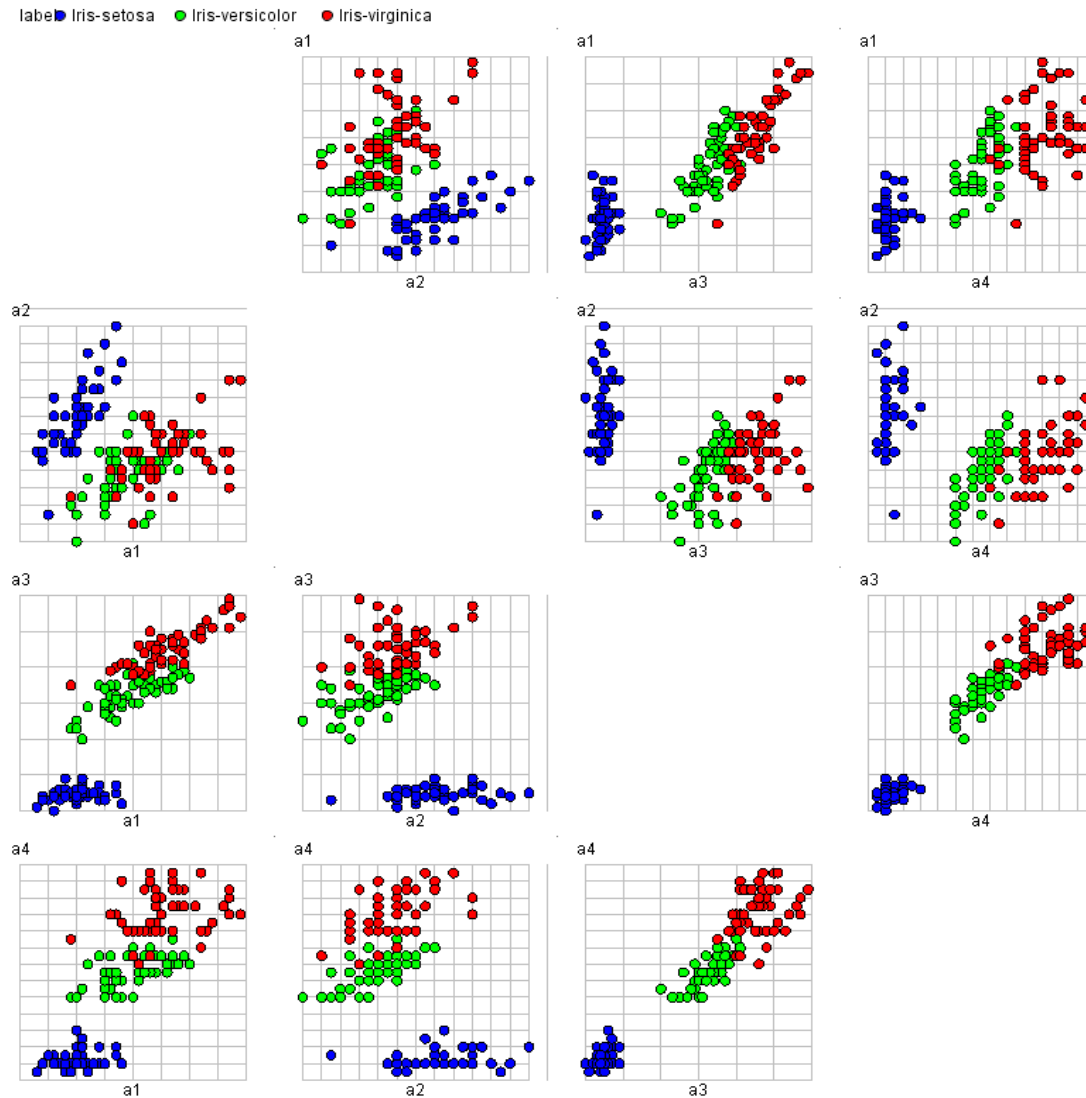


**Number of Bins**

# Visualization Techniques: Scatter Charts

– Two-dimensional scatter charts are most commonly used

– Often additional attributes/dimensions are displayed by using the size, shape, and color of the markers that represent the objects

– It is useful to have arrays of scatter charts that can compactly summarize the relationships of several pairs of attributes

**RapidMiner Plotter: Scatter**



– RapidMiner Scatter Charts

- Scatter (single chart)

- Scatter Multiple

- Scatter Matrix

- Scatter 3D

# RapidMiner Chart: Scatter Matrix

# RapidMiner Resources

- RapidMiner 7:
  - https://my.rapidminer.com/nexus/account/index.html#downloads
  - **License Key** can be found in ILIAS

- Rapidminer User Manuals: http://rapidminer.com/documentation/

- Open Access Book covering RapidMiner

  - Matthew North: Data Mining For The Masses: http://dl.dropbox.com/u/31779972/DataMiningForTheMasses.pdf

- Operator Documentation: http://rapid-i.com/wiki/

- RapidMiner Forum and Discussion Groups: http://forum.rapid-i.com/

- Video Tutorials

  - by Rapid-I: https://www.youtube.com/user/RapidIVideos

  - by NDLR: https://dspace.ndlr.ie/jspui/handle/10633/2353

  - by Neutral Market Trends: http://www.neuralmarkettrends.com/tutorials/

- MyExperiment: process repository: http://www.myexperiment.org/