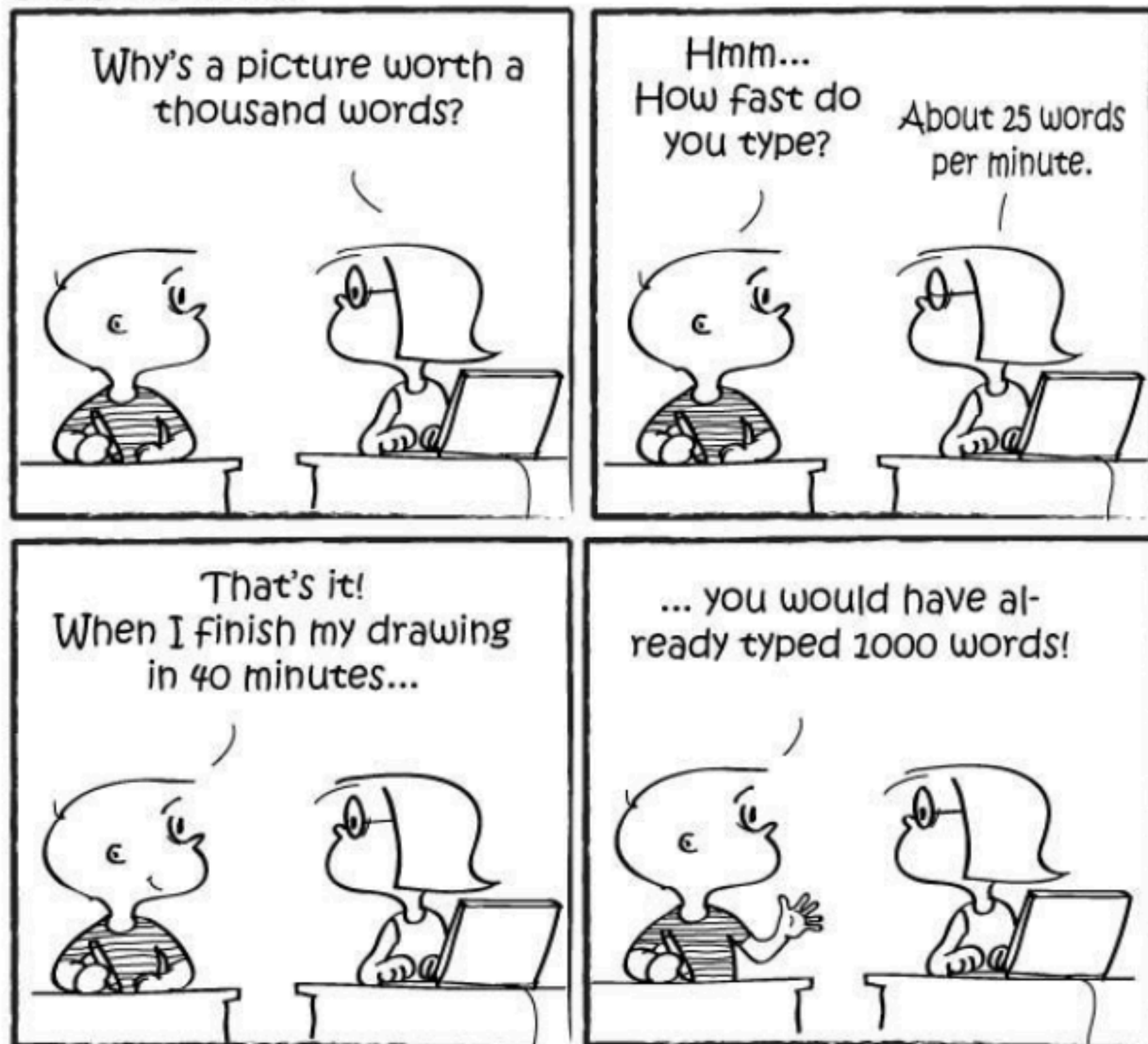# Data Mining
## (Mining Knowledge from Data)

## Data Visualization

Magda Friedjungová

# Exercises organization

- Attending the practice is compulsory.

- 4 blocks of exercises, dates on EDUX.

- 6 homework assignments, 2 of them assigned at every practice (excluding the last one), for 8 points each.

- 2 points for compulsory activity at last practice.

- The required minimum is 25 points, the maximum is 50.

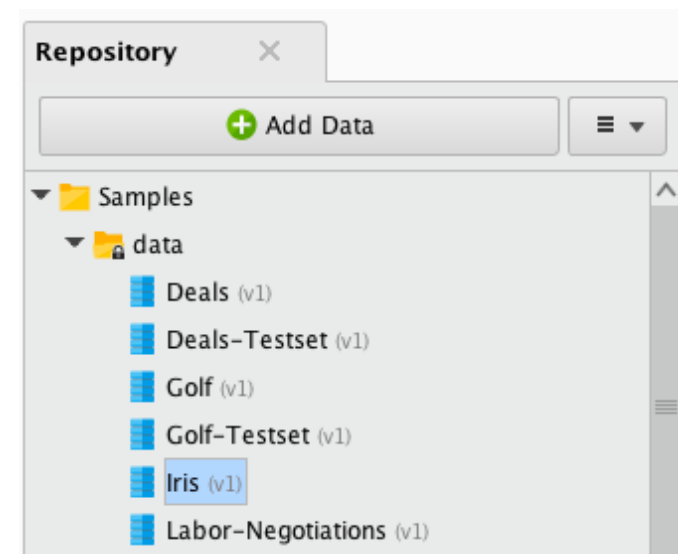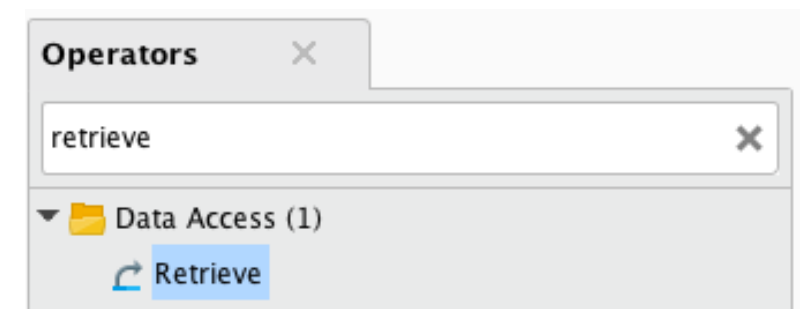# Importance of Visualization

- Helps understand data

- Extremely powerful tool

- Bypass language centers, go directly to the visual cortex

- Ability to recognize patterns

- Animations, live data processing

"As the Chinese say, 1001 words is worth more than a picture."
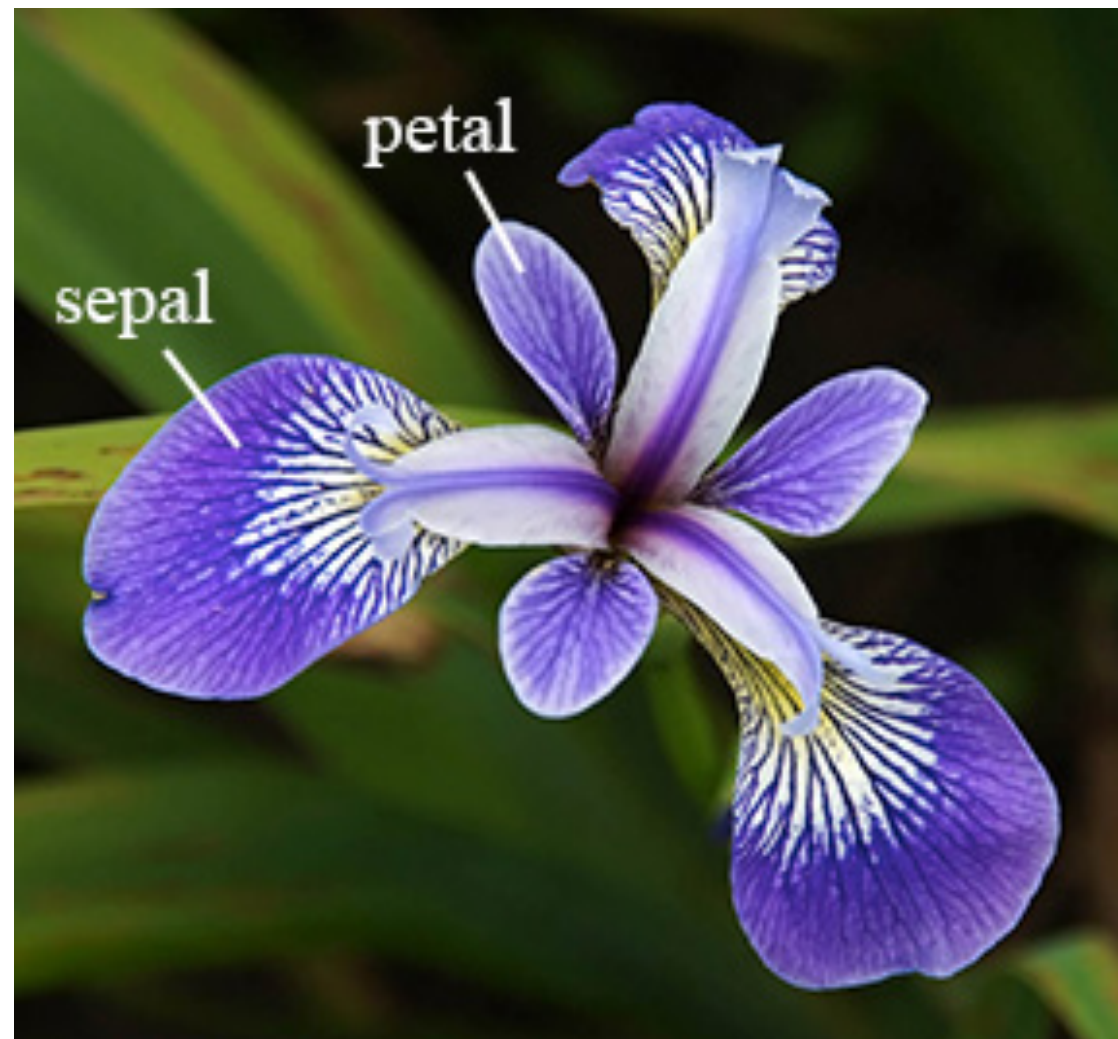
# RapidMiner Studio

- Start RapidMiner Studio.

- Create new blank process.

- Find the block "Retrieve" and drag to your process (at MacOS you can skip this step).

- Find the Iris dataset in Samples and set it to the Retrieve block.

- Run the model.

# Iris Dataset

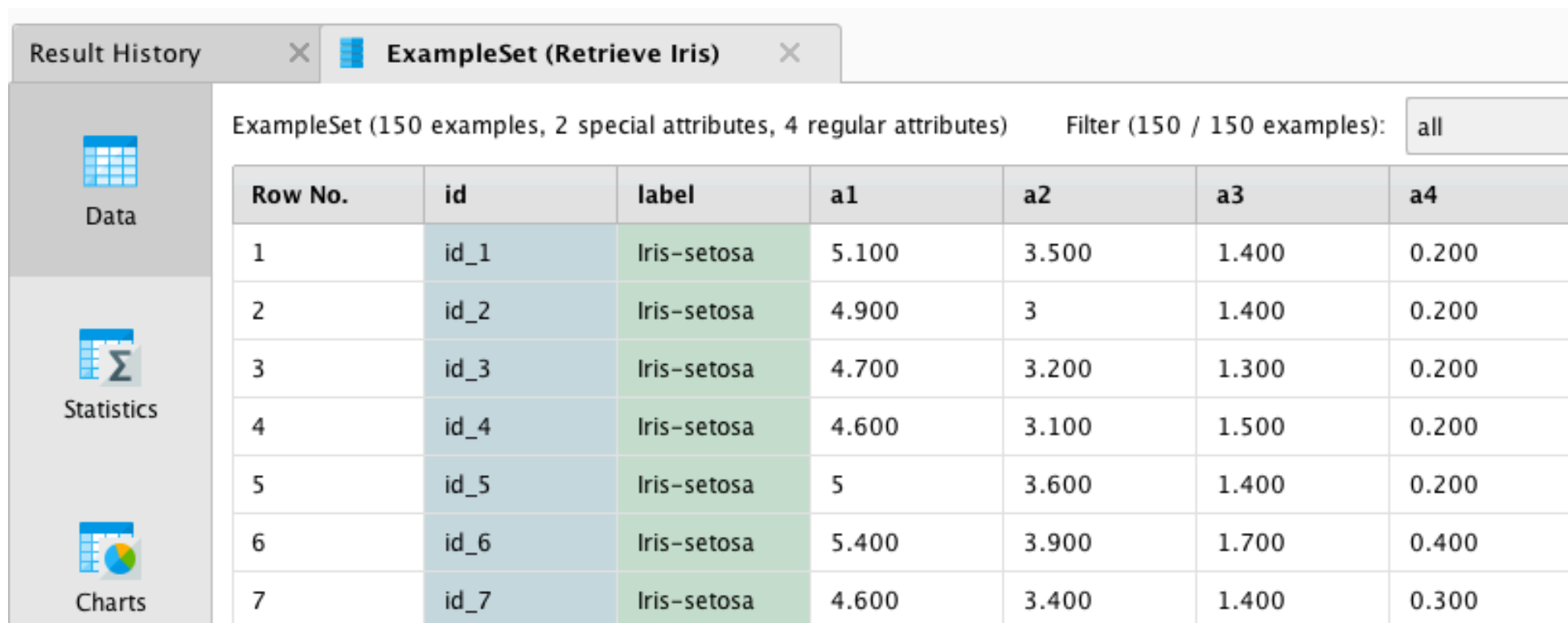| Data Set Characteristics: | Multivariate | Number of Instances: | 150 | Area: | Life |
|---|---|---|---|---|---|
| Attribute Characteristics: | Real | Number of Attributes: | 4 | Date Donated | 1988-07-01 |
| Associated Tasks: | Classification | Missing Values? | No | Number of Web Hits: | 1104221 |

# Questions 1/2

- Switch to Data.

  - How many attributes are in the dataset?

  - How many samples are in the dataset?

  - How many classes are in the dataset?

  - How many samples are in each class?

# Answers 1/2

- There are 4 attributes in the dataset.

- The dataset has 150 examples.

- There are 3 classes in the dataset.

- Each class has 50 samples.



| Row No. | id | label | a1 | a2 | a3 | a4 |
|---------|------|-------------|-------|-------|-------|-------|
| 1 | id_1 | Iris–setosa | 5.100 | 3.500 | 1.400 | 0.200 |
| 2 | id_2 | Iris–setosa | 4.900 | 3 | 1.400 | 0.200 |
| 3 | id_3 | Iris–setosa | 4.700 | 3.200 | 1.300 | 0.200 |
| 4 | id_4 | Iris–setosa | 4.600 | 3.100 | 1.500 | 0.200 |
| 5 | id_5 | Iris–setosa | 5 | 3.600 | 1.400 | 0.200 |
| 6 | id_6 | Iris–setosa | 5.400 | 3.900 | 1.700 | 0.400 |
| 7 | id_7 | Iris–setosa | 4.600 | 3.400 | 1.400 | 0.300 |

ExampleSet (150 examples, 2 special attributes, 4 regular attributes)    Filter (150 / 150 examples): all

# Questions 2/2

- Switch to Charts -> Plot View.

  - Is there any attribute which differentiates among the Irises?

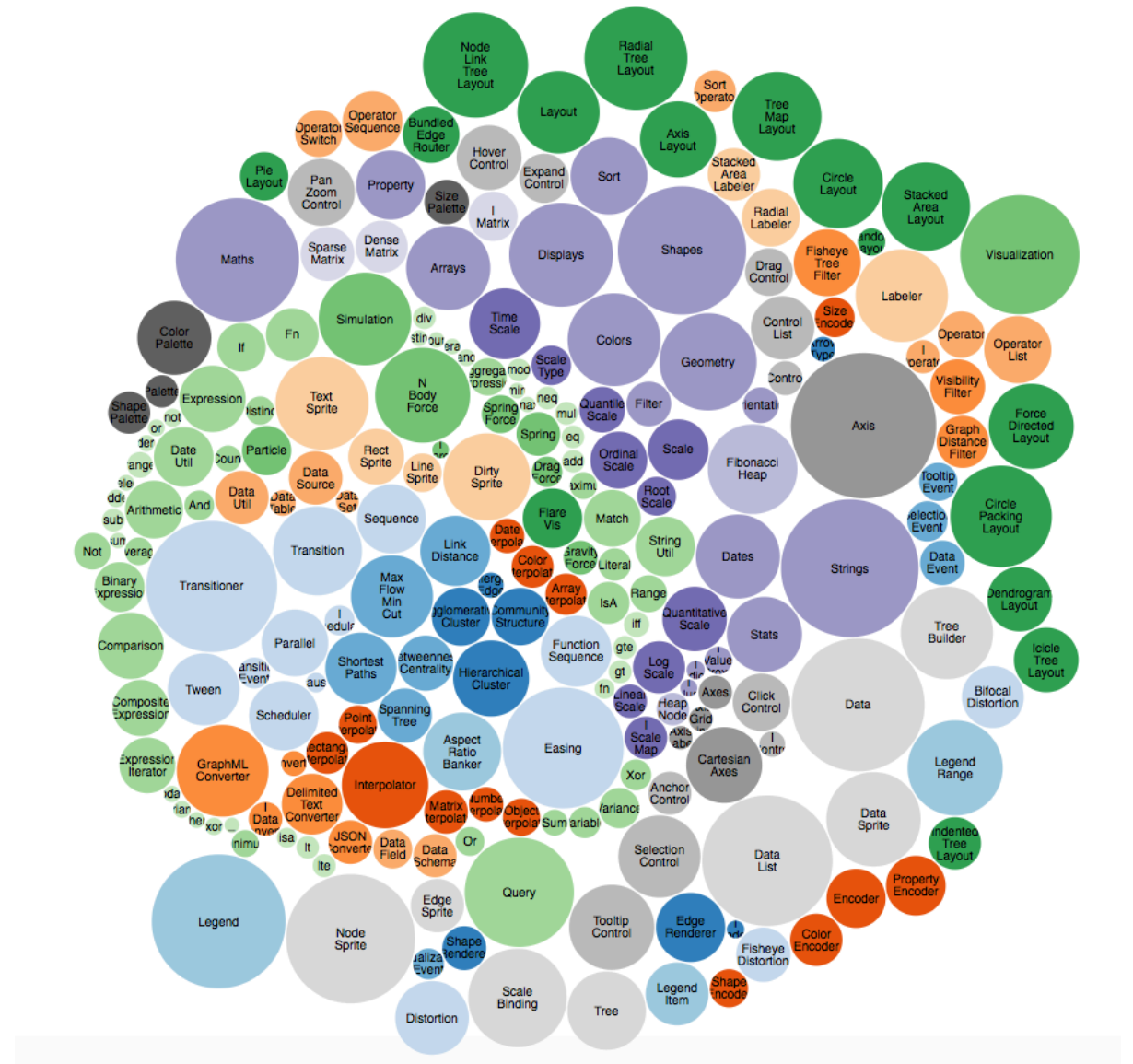  - Which attribute is the best one for differentiating among the classes?

# Answers 2/2

- The best attribute which differentiates among classes does not exist (versicolor mixes with virginica).

- The best attribute for separation is a3 or a4.

# Some Popular Tools

- Tableau Software (student license for free)
  http://www.tableau.com

- Google Data Studio (beta version)
  https://www.google.com/analytics/data-studio/

- Chart.js
  http://www.chartjs.org

- D3.js
  https://d3js.org

- Visualize Free
  http://visualizefree.com

- Gephi
  https://gephi.org

- Microsoft Power BI (Cloud version for free, register with another domain than google, yahoo etc.)
  https://powerbi.microsoft.com/en-us/

- Tag Cloud

  - Create Tag Cloud for an arbitrary document.

  - http://www.wordle.net

  - Meaningless word (a, the, for, etc.) can be removed by right-clicking.
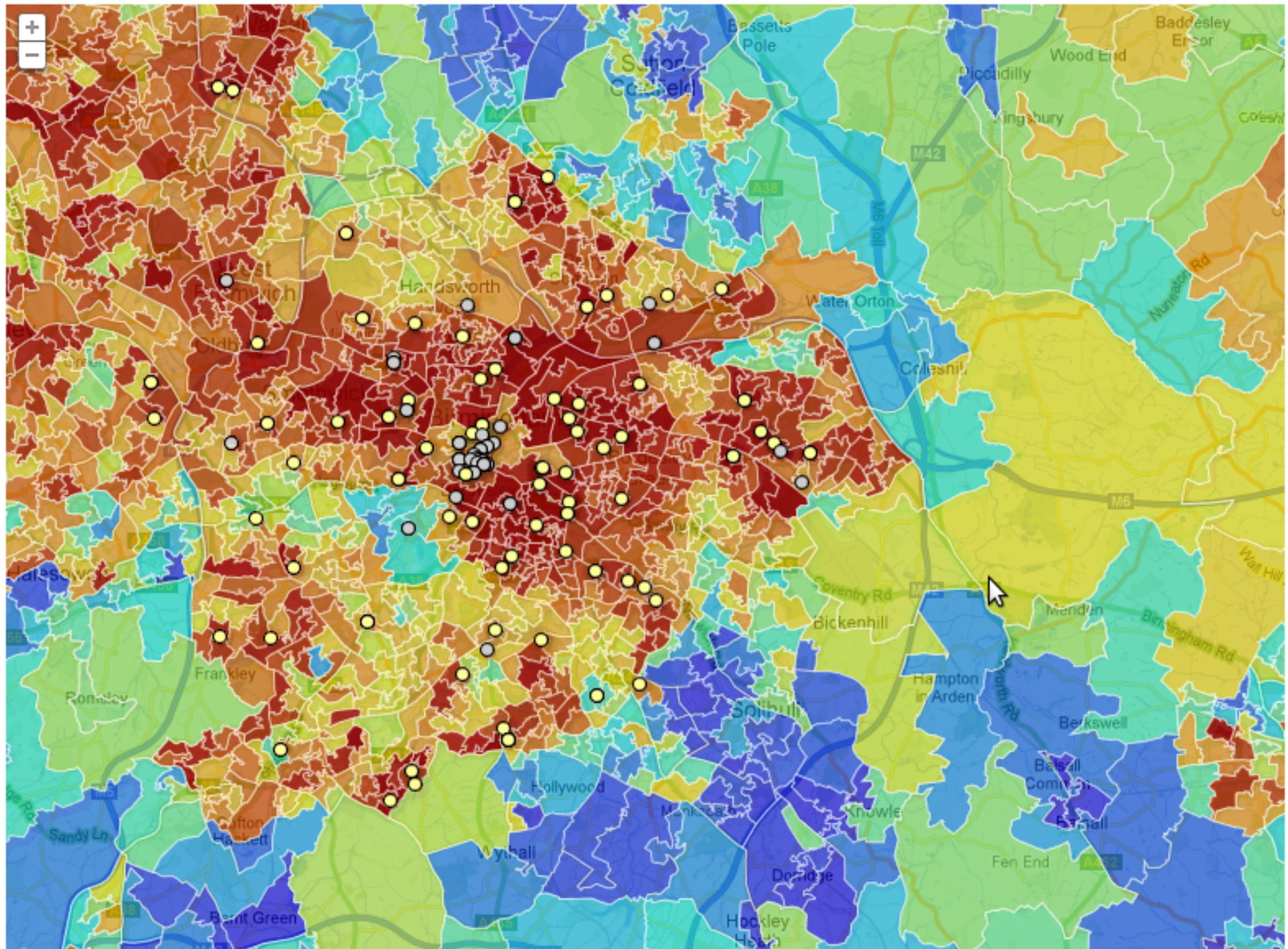
- Google Fusion Tables
  http://www.google.com/fusiontables

  - Search for interesting data.

  - Visualize an example, using the "Map".

Key ⬭ Offence ⬭ Accused address **Deprivation** *Poorer* ■ ■ ■ ■ ■ ■ *Richer*
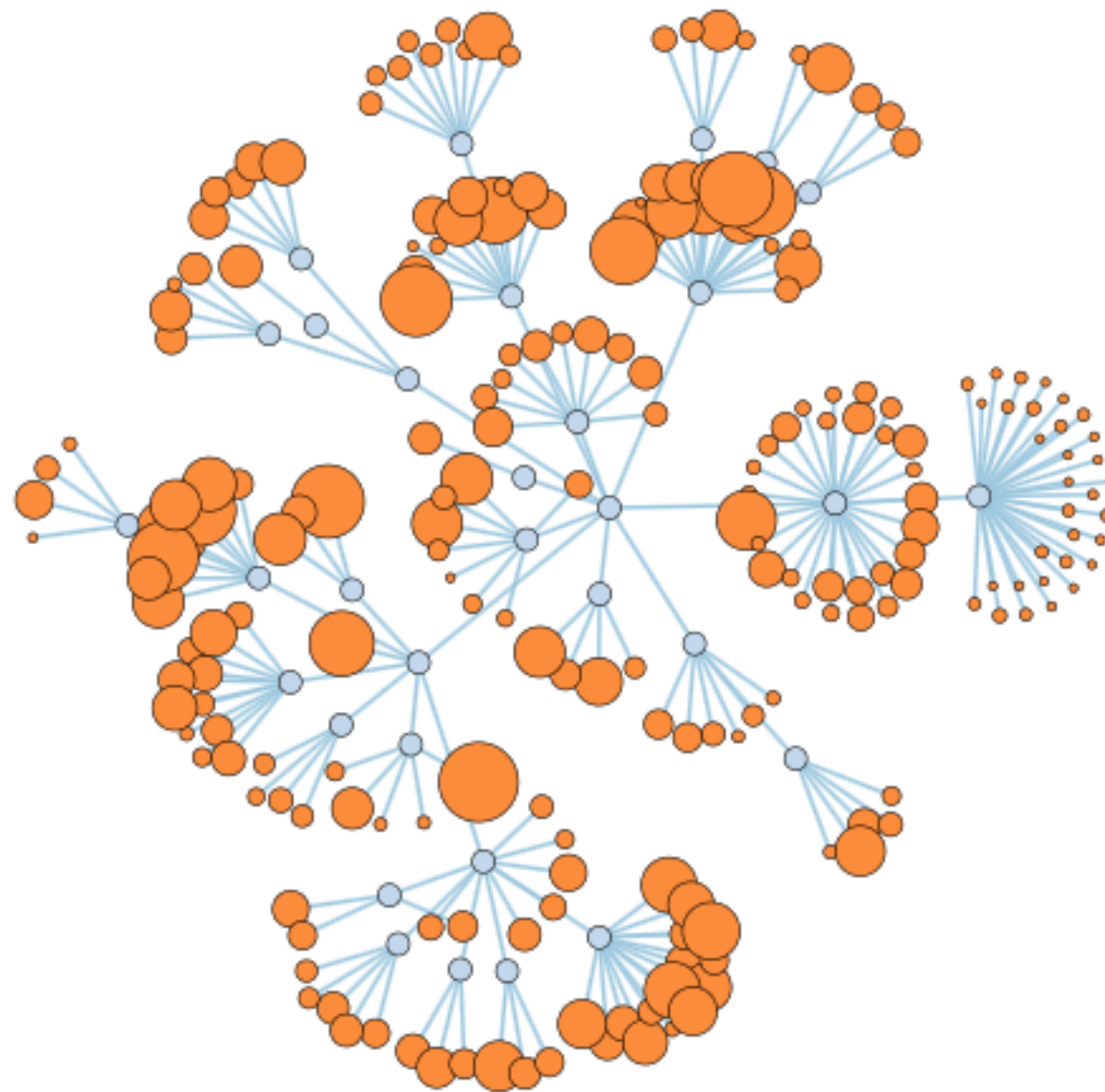
- mapsdata
  http://www.mapsdata.co.uk/

  - The Busiest Stations on the London Underground

- D3.js http://mbostock.github.io/d3/talk/20111116/force-collapsible.html

- Linux Kernel Development Visualization
  https://www.youtube.com/watch?v=P_02QGsHzEQ