Department of Mathematics
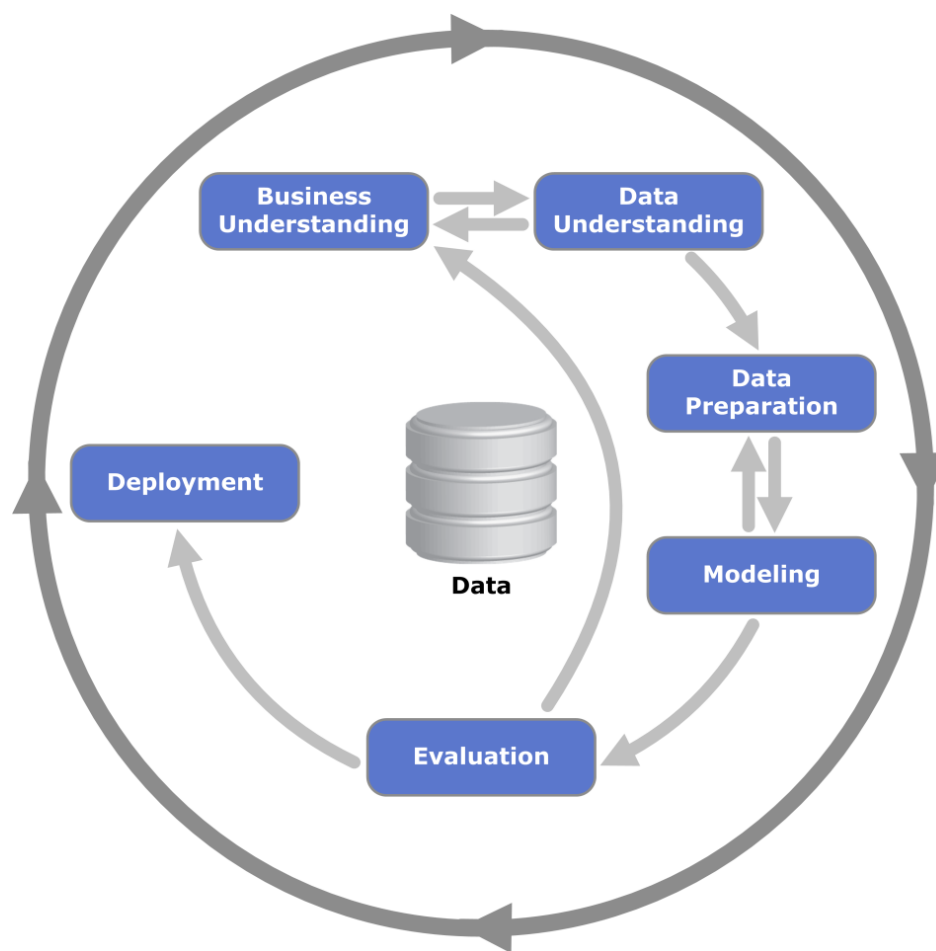University of Calabria

*Data Warehouse and Data Mining*

***Module II – Data Mining***

# Data Preparation

Ph.D. Ettore Ritacco

# The Knowledge Discovery Process (CRISP-DM)

# Business Understanding

- You are a doctor of a medical division, who collected some data from your patients

- All of them contracted the same disease

- The therapy consists in 5 different and exclusive cures, according to the patient condition

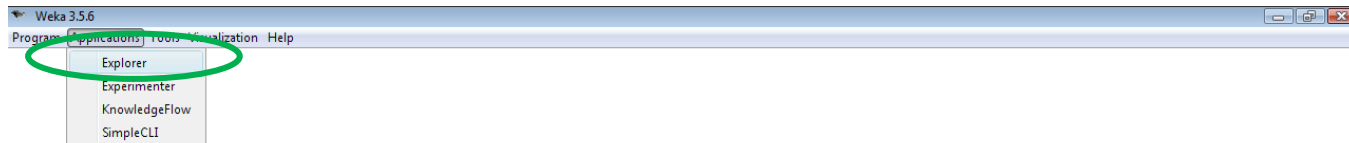- Does an automatic procedure, for the cure assignment, exist?

# Data Understanding

| Attribute | Description |
| --- | --- |
| Instance_number | Incremental tuple ID (INTEGER) |
| ID | Patient's ID (INTEGER) |
| Age | Patient's age (INTEGER) |
| Sex | Patient's gender: F or M |
| BP | Blood Pressure: HIGH, NORMAL or LOW |
| Cholesterol | Concentration of cholesterol in the blood: NORMAL or HIGH |
| Na | Concentration of sodium in the blood (REAL) |
| K | Concentration of potassium in the blood (REAL) |
| Drug | The chosen cure: drugY, drugC, drugX, drugA, drugB |

# Data Understanding

⊙ Data acquisition in Weka:

# Data Understanding

○ Data acquisition in Weka:

# Data Understanding

○ Data acquisition in Weka:

# Data Understanding

- **Target attribute (the class attribute): Drug**

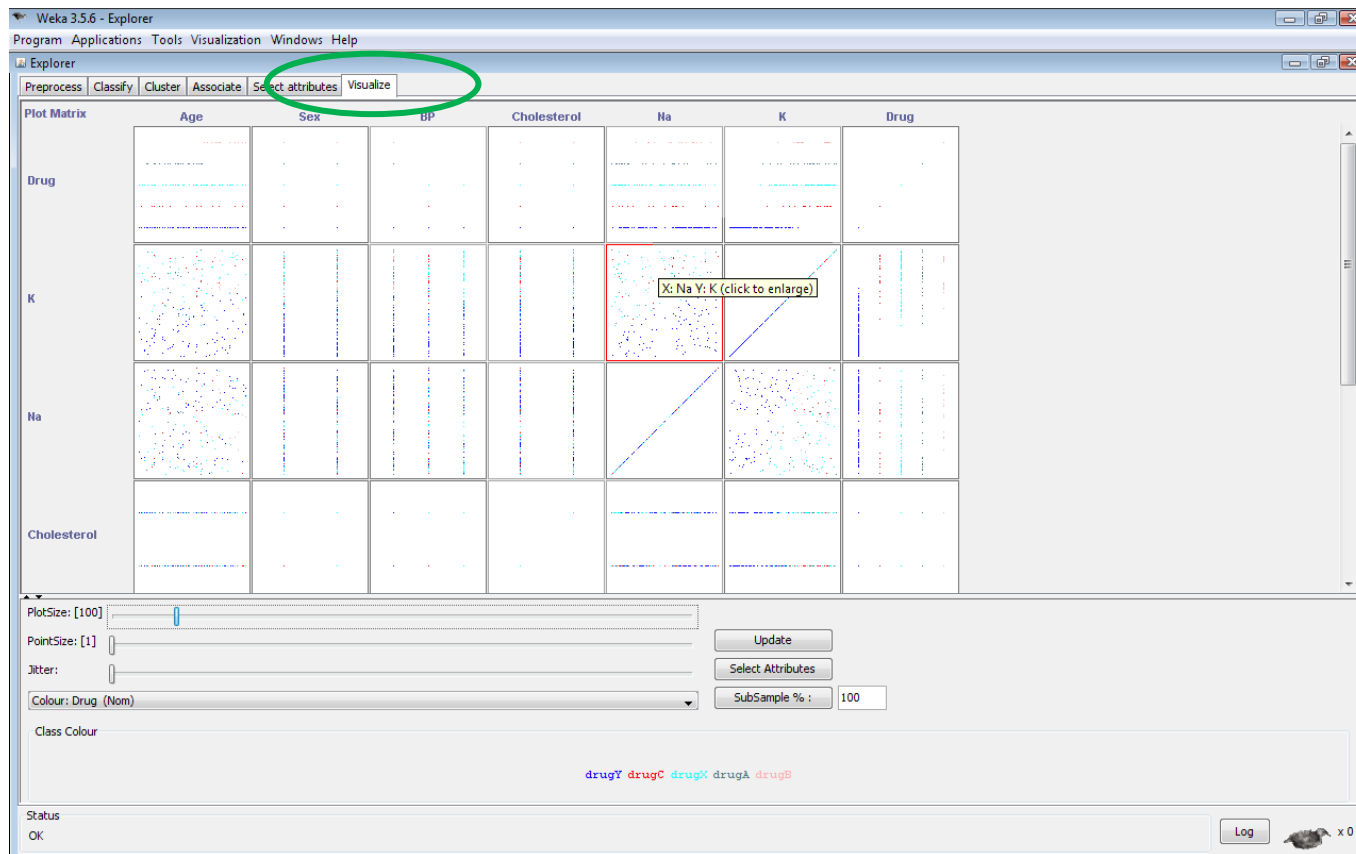# Data Understanding

- Distributions: comments?

# Data Understanding

- When age is greater then (about) 44 years, the class drugB starts to appear

- The sex and the concentration of sodium seem to not affect the class distribution

- Some cures are provided according to the values of blood pressure and cholesterol

- When the concentration of potassium is lower than (about) 0.05 the dominant class is drugY

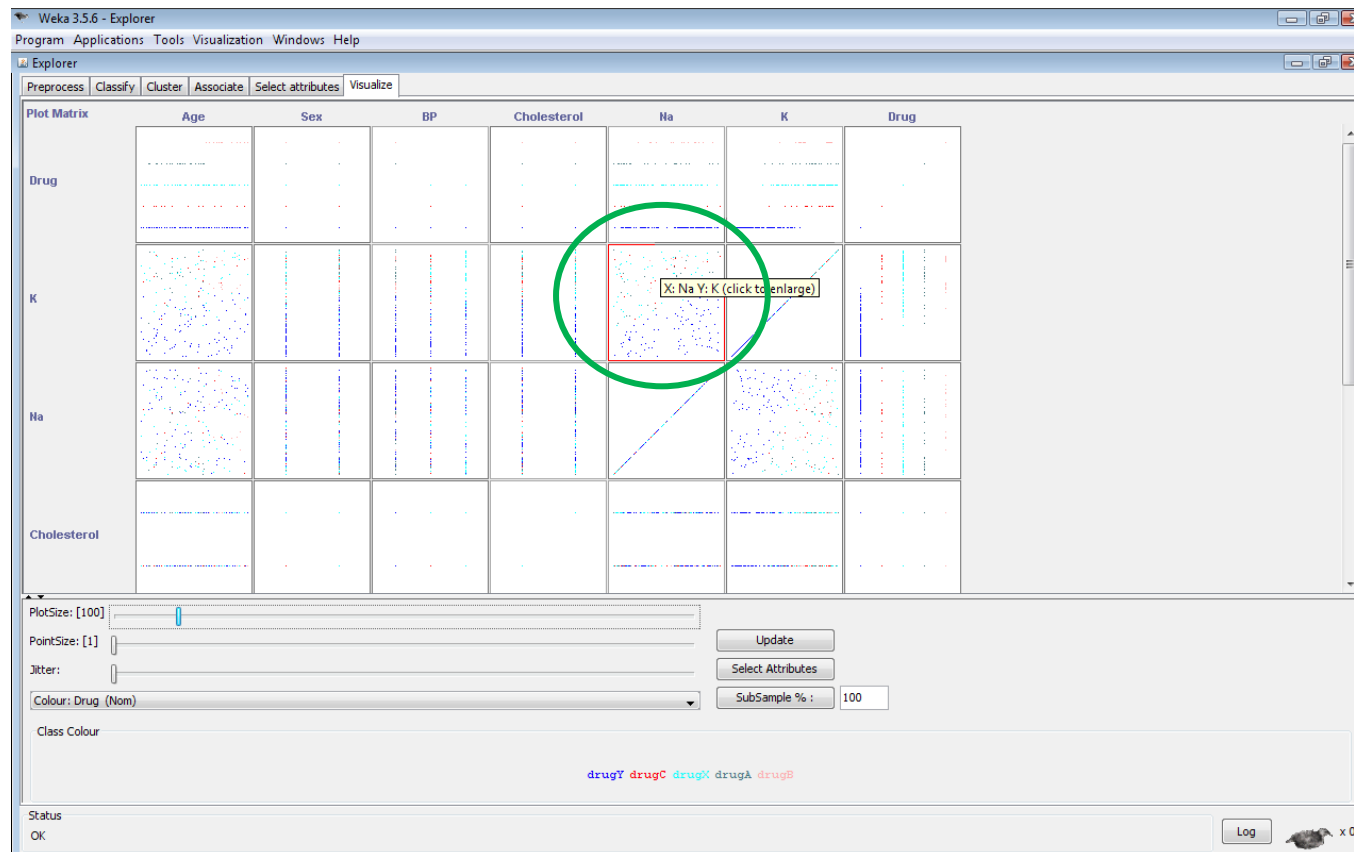- Data are unbalanced:
  - There is one dominant class: drugY (the blue one)
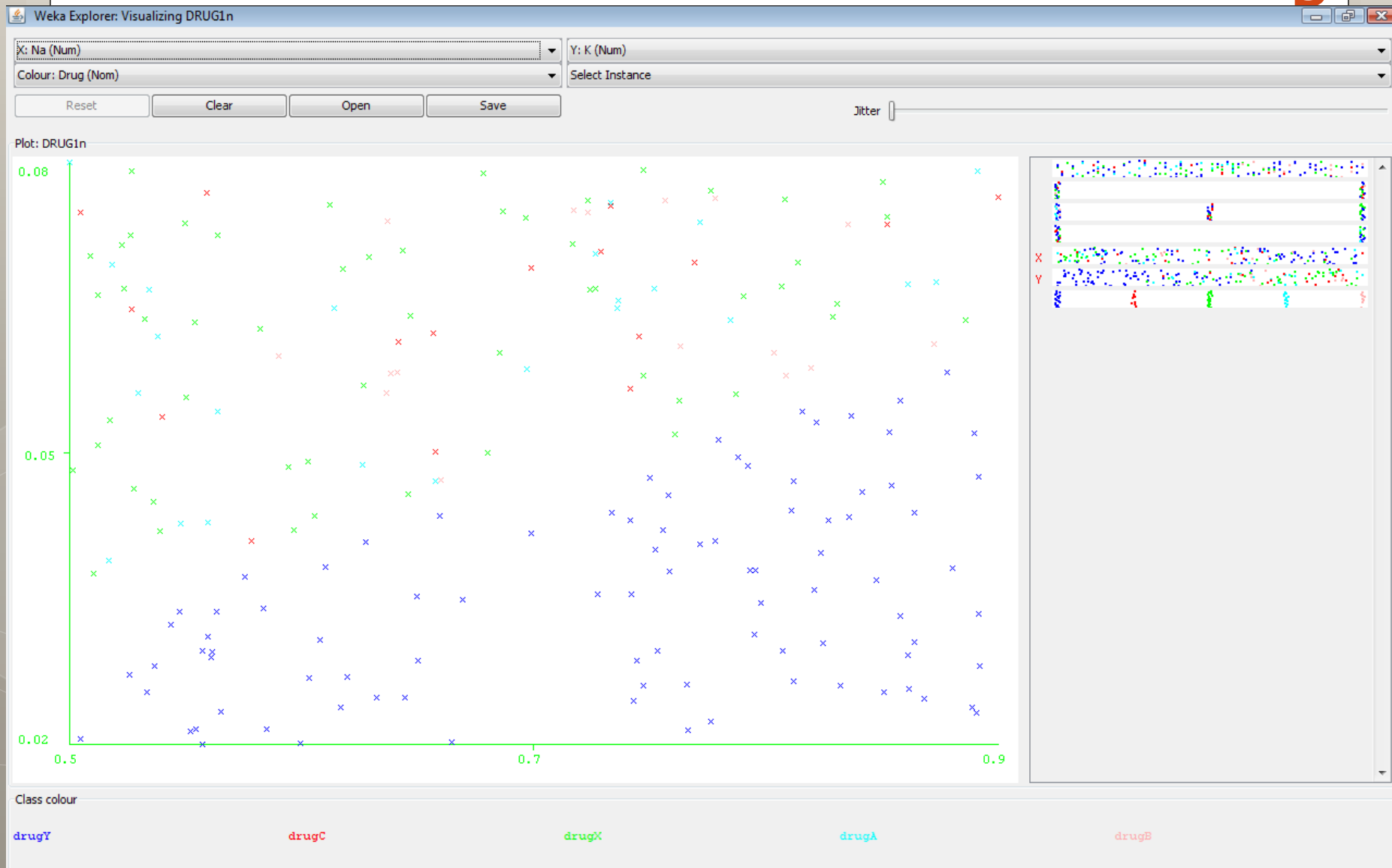
# Data Understanding

- Scatter plot

# Data Understanding

○ Scatter plot

# Data Understanding
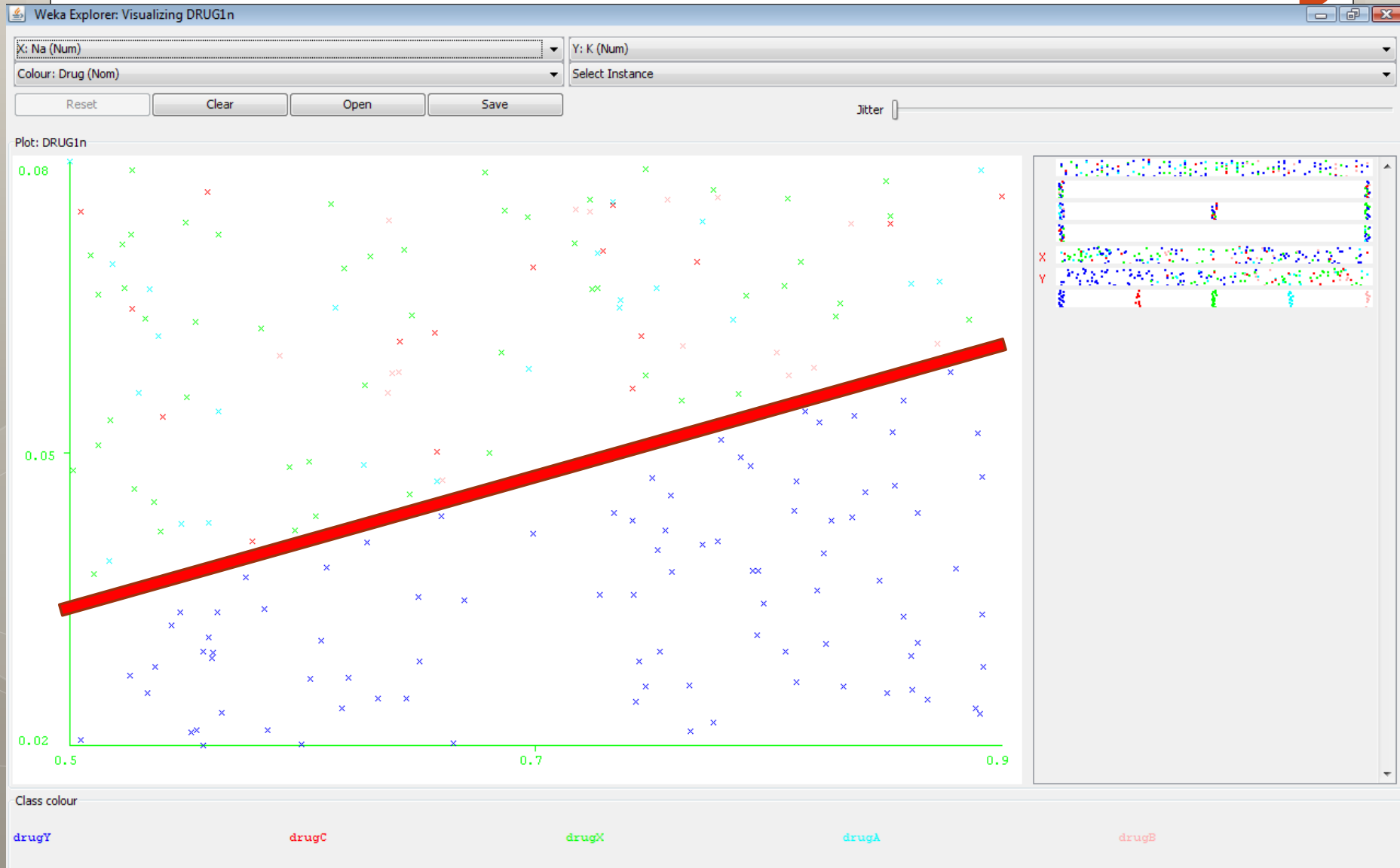
# Data Understanding

- Comments?

# Data Understanding

○ A hyperplane can clearly separate one class (drugY, the blue one) from the others!

# Data Understanding

# Data Preparation

- Hyperplane:
  - K = m * Na + q

- Let's ignore the constant *q*, then
  - m = Na / K

- We can create a new attribute Na_su_K (Na out of K, in *italian*)

- Choose the AddExpession filter

# Data Preparation

- Configure the filter

# Data Preparation

- This filter adds a new attribute to the table, applaying a mathematical expression to the existing attributes (field *expression*). The new attribute will take name from the field *name*.

- Click on *more* and *capabilities* for further information about the filter

Department of Mathematics
University of Calabria

# Data Preparation

# Data Preparation

- Comments?

# Data Preparation

- The attribute Na and K now are redundant. Let's remove them

# Data Preparation

# Modeling

○ Data seems ready to be used in the modeling phase