# Simple missing functionalities

Anton Antonov
MathematicaForPrediction at GitHub
MathematicaVsR project at GitHub
November, 2016

## Introduction

This notebook belongs to the Mathematica-part of MathematicaVsR at GitHub project DataWrangling.

This notebook illustrates commands that are (in my opinion) are missing from Mathematica but are present and used often in R. See these corresponding R-part HTML file or RMarkdown file.

In this notebook functionalities of those missing commands are obtained by the functions `RecordsSummary`, `VariableDependenceGrid`, `CrossTabulate`, and `MosaicPlot`.

### Load packages

The following commands load the packages used in this notebook.

```
Import[
 "https://raw.githubusercontent.com/antononcube/MathematicaForPrediction/master/
   MathematicaForPredictionUtilities.m"]
```

```
Import[
 "https://raw.githubusercontent.com/antononcube/MathematicaForPrediction/master/
   MosaicPlot.m"]
```

## Data load and rudimentary analysis

### Titanic data

Here is the summary of the Titanic data used below:

```
titanicData =
   (Flatten@*List) @@@ ExampleData[{"MachineLearning", "Titanic"}, "Data"];
columnNames = (Flatten@*List) @@
    ExampleData[{"MachineLearning", "Titanic"}, "VariableDescriptions"];
titanicData = DeleteCases[titanicData, {___, _Missing, ___}];
RecordsSummary[titanicData, columnNames]
```

|  | 2 passenger age | | | |
|---|---|---|---|---|
| 1 passenger class | Min 0.1667 | | | |
| | 1st Qu 21. | 3 passenger sex | 4 passenger survival | |
| 3rd 501 | Median 28. | male 658 | died 619 | |
| 1st 284 | Mean 29.8811 | female 388 | survived 427 | |
| 2nd 261 | 3rd Qu 39. | | | |
| | Max 80. | | | |

This variable dependence grid shows the relationships between the variables.

```
Magnify[#, 0.7] &@VariableDependenceGrid[titanicData, columnNames]
```

## Employee attitude dataset

Here is the summary of the Titanic data used below:

```
eaData = ExampleData[{"Statistics", "EmployeeAttitude"}];
eaColumnNames = (Flatten@*List) @@
    ExampleData[{"Statistics", "EmployeeAttitude"}, "ColumnHeadings"];
Multicolumn[RecordsSummary[N@eaData, eaColumnNames], 3, Dividers → All]
```
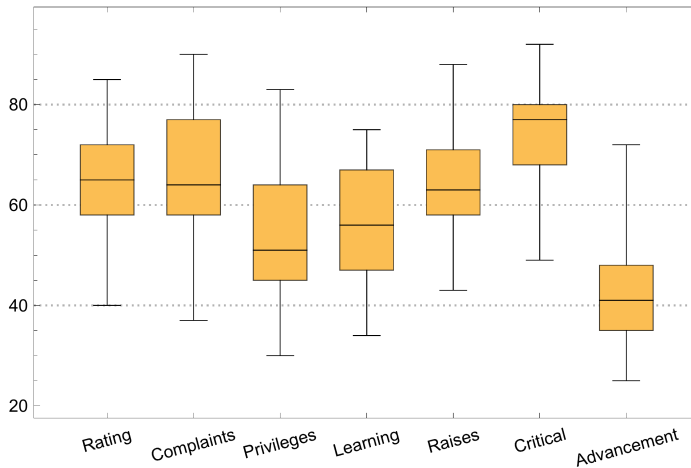
| 1 Rating | | 4 Learning | | 7 Advancement | |
|---|---|---|---|---|---|
| Min | 40. | Min | 34. | Min | 25. |
| 1st Qu | 58. | 1st Qu | 47. | 1st Qu | 35. |
| Mean | 64.6333 | Mean | 56.3667 | Median | 41. |
| Median | 65.5 | Median | 56.5 | Mean | 42.9333 |
| 3rd Qu | 72. | 3rd Qu | 67. | 3rd Qu | 48. |
| Max | 85. | Max | 75. | Max | 72. |
| 2 Complaints | | 5 Raises | | | |
| Min | 37. | Min | 43. | | |
| 1st Qu | 58. | 1st Qu | 58. | | |
| Median | 65. | Median | 63.5 | | |
| Mean | 66.6 | Mean | 64.6333 | | |
| 3rd Qu | 77. | 3rd Qu | 71. | | |
| Max | 90. | Max | 88. | | |
| 3 Privileges | | 6 Critical | | | |
| Min | 30. | Min | 49. | | |
| 1st Qu | 45. | 1st Qu | 68. | | |
| Median | 51.5 | Mean | 74.7667 | | |
| Mean | 53.1333 | Median | 77.5 | | |
| 3rd Qu | 64. | 3rd Qu | 80. | | |
| Max | 83. | Max | 92. | | |

It is a good idea to get an impression of the numerical variables distributions in a given dataset.

There are several approaches for doing this (in Mathematica and in general.)
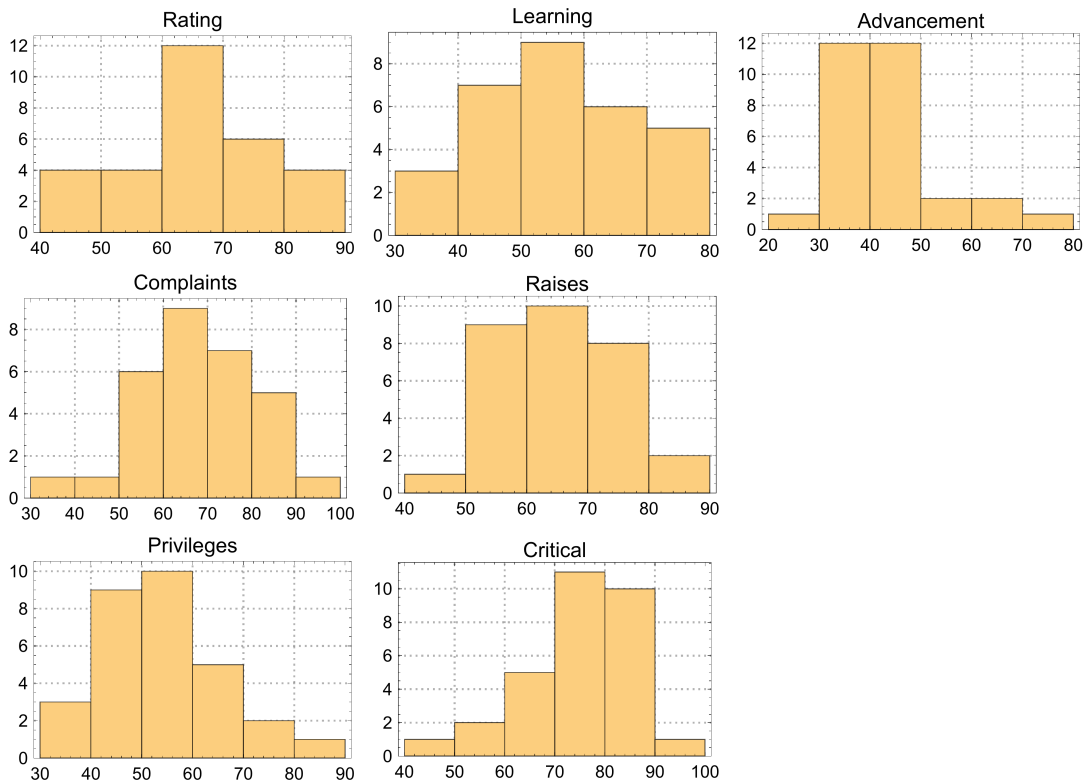
## Box-and-whisker diagrams

```
DistributionChart[Transpose[eaData], ChartElementFunction → "BoxWhisker",
 ChartLabels → Map[Rotate[#, π / 12] &, eaColumnNames], PlotTheme -> "Detailed"]
```



## Panel of histograms

```
Multicolumn[
 MapThread[Histogram[#1, PlotLabel → #2, PlotRange → All, PlotTheme -> "Detailed"] &,
  {Transpose[eaData], eaColumnNames}], 3]
```

# Cross tabulation and mosaic plots

## Cross tabulation

In statistics contingency tables are matrices used to show the co-occurrence of variable values of multi-dimensional data. They are fundamental in many types of research. Below are some examples of cross-tabulation. For a detailed discussion see the Markdown file "Contingency-tables-creation-examples.md" of this project or the corresponding PDF file.

```
CrossTabulate[titanicData[All, {1, 3}]] // MatrixForm
```

$$
\begin{pmatrix}
 & \text{female} & \text{male} \\
\hline
\text{1st} & 133 & 151 \\
\text{2nd} & 103 & 158 \\
\text{3rd} & 152 & 349
\end{pmatrix}
$$

A generalization of `CrossTabulate` is the function `CrossTensorate` implemented in Mathematica-ForPredictionUtilities.m that takes a "formula" argument similar to R's xtabs.

```
CrossTensorate[Count == "passenger class" + "passenger sex" + "passenger survival",
   titanicData, columnNames] // MatrixForm
```

$$
\begin{pmatrix}
 & \text{female} & \text{male} \\
\hline
\text{1st} & \begin{pmatrix} \text{died} & 5 \\ \text{survived} & 128 \end{pmatrix} & \begin{pmatrix} \text{died} & 98 \\ \text{survived} & 53 \end{pmatrix} \\
\text{2nd} & \begin{pmatrix} \text{died} & 11 \\ \text{survived} & 92 \end{pmatrix} & \begin{pmatrix} \text{died} & 135 \\ \text{survived} & 23 \end{pmatrix} \\
\text{3rd} & \begin{pmatrix} \text{died} & 80 \\ \text{survived} & 72 \end{pmatrix} & \begin{pmatrix} \text{died} & 290 \\ \text{survived} & 59 \end{pmatrix}
\end{pmatrix}
$$
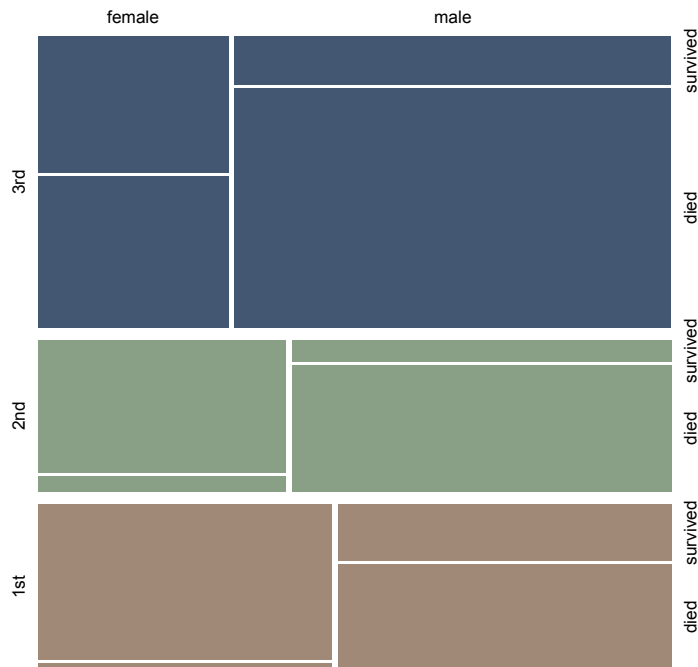
```
CrossTensorate["passenger age" == "passenger class" + "passenger sex",
   titanicData, columnNames] // MatrixForm
```

$$
\begin{pmatrix}
 & \text{female} & \text{male} \\
\hline
\text{1st} & 4926. & 6195.42 \\
\text{2nd} & 2832.42 & 4868.83 \\
\text{3rd} & 3372.17 & 9060.83
\end{pmatrix}
$$

## Mosaic plots

Mosaic plots can illustrate fairly well the (conditional) dependencies between the values of the categorical variables in a dataset.

```
MosaicPlot[titanicData[All, {1, 3, 4}]]
```



In contrast with R (and RStudio) in Mathematica's FrontEnd we can have tooltips showing the exact conditional values. (Hover with the mouse pointer over the rectangles in the plot above.)



| condition | event | probability |
|---|---|---|
| | 3rd | 0.478967 |
| | 3rd ∩ male | 0.333652 |
| 3rd | male | 0.696607 |
| | 3rd ∩ male ∩ died | 0.277247 |
| 3rd | male ∩ died | 0.578842 |
| 3rd ∩ male | died | 0.830946 |