

# Linear regression with L<sup>A</sup>T<sub>E</sub>X

Battista Benciolini

December 14, 2024

## Contents

|           |   |           |
|-----------|---|-----------|
| <b>1</b>  | <b>Introduction: first description of the problem</b>                       | <b>2</b>  |
| <b>2</b>  | <b>Geometric definition of three optimality criteria</b>                    | <b>2</b>  |
| <b>3</b>  | <b>Information on the realised solution, including limitations</b>          | <b>3</b>  |
| <b>4</b>  | <b>Some comments about the programming of the package</b>                   | <b>4</b>  |
| <b>5</b>  | <b>A ready to use simple user interface</b>                                 | <b>4</b>  |
| <b>6</b>  | <b>A user manual for the package</b>  | <b>4</b>  |
| <b>7</b>  | <b>An example</b>   | <b>6</b>  |
| <b>8</b>  | <b>A package for linear regression and the theory behind it</b>             | <b>6</b>  |
| 8.1       | Math preliminaries and notation . . . . .                                   | 6         |
| 8.2       | Package declaration, required package and definition of variables . . . . . | 8         |
| 8.3       | Preparing data input . . . . .  | 9         |
| 8.4       | Computation . . . . .   | 10        |
| 8.4.1     | Computation of first and second order moments . . . . .                     | 10        |
| 8.4.2     | Classical linear regression . . . . .                                       | 12        |
| 8.4.3     | Symmetric linear regression . . . . .                                       | 13        |
| 8.5       | Print of table of results . . . . .   | 16        |
| 8.6       | Plot of points and lines . . . . .  | 17        |
| 8.7       | Functions for internal use . . . . .  | 18        |
| <b>9</b>  | <b>Versions</b>   | <b>21</b> |
| <b>10</b> | <b>Acknowledgments</b>  | <b>21</b> |
| <b>11</b> | <b>Citations and references</b>   | <b>21</b> |

---

Linear regression with LaTeX - available in CTAN  
Battista Benciolini - contact: benciolinibattista at gmail dot com

# 1 Introduction: first description of the problem

I start with a quote from *ArXiv* (April 2021, number 31, page 73):

The physicist Mario Rossi is investigating a phenomenon, presumably linear, and he performs measurements in his laboratory to verify his hypothesis; he measures the quantity  $x$  which generates the phenomenon and he measures also one of the characteristics  $y$  showed by the phenomenon under the effect of the stimulation  $x$ .

...  
Subsequently Mario graphs the data of the table to judge if the points reasonably follow a linear trend or not; in this regard he computes the parameters of the regression line and he draws this line on the graph in order to judge the quality of the obtained results.

...  
Being a  $\text{\LaTeX}$  user, he thinks to kill two birds with one stone: using  $\text{\LaTeX}$  to draw the graph with the experimental data consisting in the  $x, y$  points and, at the same time, to compute the parameter  $a$  e  $b$  of the regression line  $y = ax + b$ , and finally to draw also this line on the same graph.

A summary description of the problem is therefore the following. A set of data pairs is available and each pair is represented as a point in the plane. A straight line is searched that optimally approximates the points. The first step is therefore the choice of an optimality criterion. This choice is the topic of the next section. From the text we also know that the possible deviation of  $y$  with respect to the model is quite larger than the uncertainty of  $x$ .

After reading the description of the problem of Mario Rossi I tried to produce a solution. In this work I will use  $y_1$  and  $y_2$  instead of  $x$  and  $y$  for the two measured quantities that will become the first and second coordinate, or abscissa and ordinate, in the Cartesian plane.

The problem can be treated as a mere problem of approximation or alternatively as an estimation problem in the frame of a probabilistic description of the uncertainty. The two treatments are conceptually different. The probabilistic treatment produces some more results, but the estimation of the parameters is the same. On the other hand the treatment as an approximation problem is in some sense more immediate and requires a less extended theoretical background. For this reason it will be preferred here. I consider the original problem and also a variation of it based on the assumption that the two variables are known with the same uncertainty. The two considered situations will prove to be quite different.

## 2 Geometric definition of three optimality criteria

For each point given in the plane we can consider the corresponding point with the same abscissa and belonging to the line. Remember that the line is exactly what has to be determined. The distance between the given point and the just defined point on the line is a reasonable measure of the discrepancy between the empirical data and the corresponding theoretical model. The distances we are speaking about are the lengths of the segments shown in the leftmost scheme of figure (1).

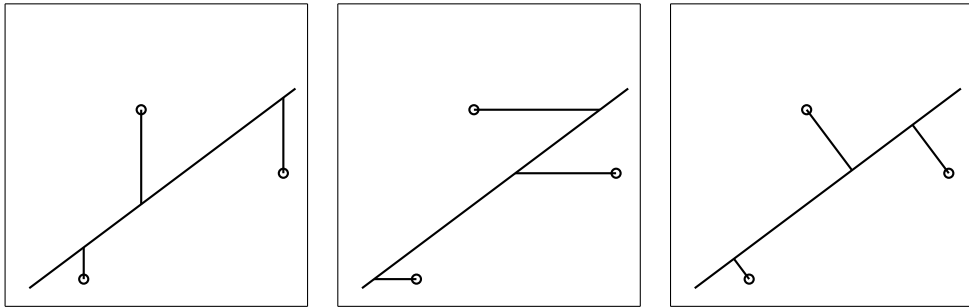


Figure 1: The three kinds of segments used in the definition of the objective function

To obtain a global discrepancy measure that considers all the points at once we perform the sum of the squares of the lengths of the mentioned segments. It is now clear that the two coordinates of the points are treated quite differently and play a different role in the definition of the optimality criterion. This choice is reasonable when the measuring errors only (or mainly) affect the second coordinate. The optimal line is the line that minimize the just defined global discrepancy. The procedure for the determination of the optimal line is named linear regression. In this work it is named *classical linear regression*. We can easily exchange the role of the two quantities, i.e. we can imagine that the first quantity is affected by errors. The problem is not conceptually different. The segments plotted in the central picture of figure (1) represent the discrepancy between the empirical data and the model. This other procedure is named *classical linear regression with inverted role of the coordinates*.

The situation is really different if the two coordinates have to be treated symmetrically. In this case the discrepancy between the empirical data and the model must be defined in a purely geometrical way. Just the line and the points enter in the definition without any special role for any predefined direction. With these requirements it is quite natural to use the distance of each point from the line. Remember that the distance of a point from a line is intended along the shortest path, i.e. measured in the direction orthogonal to the line itself. The rightmost scheme of figure (1) shows the segments that are considered. The global measure of discrepancy is again obtained as the sum of the squares of the lengths of the mentioned orthogonal segments. The procedure that obtain the optimal line that minimize the just defined global discrepancy is named *symmetrical linear regression*.

Some arguments of the present section will be repeated in section 8 from the algebraic and computational point of view.

### 3 Information on the realised solution, including limitations

The code that implements the solution is recorded in two files, that are a package (sty) file and a main interactive document. The file `linearregression.sty` provides several commands that can be used in any document. The file `mainlinearregression.tex` provides a simple interactive user interface. The package described in the sections 6 and 8 (user manual and implementation) provides the functions that execute the various needed operations, i.e. data input,

computations, printing the numerical results and generating a graphic representation of data and results. Some auxiliary functions complete the package. The design of the output (tables and plots) includes some arbitrary choices. The style of the graphic output is quite minimalist (e.g.: no colors, no variations of line styles).

## 4 Some comments about the programming of the package

Large part of the code is written using the `expl3` language. I have tried to be compliant with the various recommendations and prescriptions for a correct use of the language, but I probably only partly succeeded.

Different more elegant and more coherent solutions probably exist both for the general structure of the package and for some specific part of the code, but this is what I have been able to do. Some perhaps problematic aspects are mentioned here after.

Several used variables are global and they are accessed by various functions. This makes the various parts of the package quite connected to each other and creates strong dependencies.

The layered programming style is only partially applied. The partition between document command and lower level functions is present, but part of the low level code is directly in the document commands. Variants are not used.

## 5 A ready to use simple user interface

The main file asks the user for the name of a file containing the data and generates a one page output.

```
1 \documentclass[a4paper]{article}
2 \usepackage{lmodern}
3 \usepackage{linearregression}
4 \begin{document}
5 \pagestyle{empty}
6 \lraskfilename
7 \lrcomputation
8 \lrplotparameters{0.16}{11.0}
9 \lrplot{11.0}{+}{+}{-}{-}
10 \lrprint
11 \end{document}
```

The plot just includes the points and the line obtained with the classical linear regression. Lines 8 and 9 of the code must be modified if a different result is desired. See the next section for details.

## 6 A user manual for the package

The various analysis of a data set and the representation of the data and of the results is obtained with a sequence of several commands. The main operations are: (i) selection of the data file, (ii) data input and computation, (iii) printing

of a table, (iv) printing of a picture (that can be repeated with different parameters). It is generally convenient to put the table and the picture(s) in a proper floating environment. The commands for the four mentioned operations are described here after. The first needed operation is to set the name of the data file.

`\lrfilename` This is done with the command `\lrfilename{file}` that has a mandatory argument. The argument is the name of the data file. As an alternative the command

`\lraskfilename` `\lraskfilename` can be used. It asks the user to type the name of the data file in the terminal.

`\lrcomputation` The macro `\lrcomputation` reads the data and performs all the computations. The results of the computations remain available in internal variables and are then used by the macro that print them or generates a plot.

`\lrprint` The macro `\lrprint` generates a table with all the estimated parameters and some information about the data. The computed numbers are printed with four decimal digit maximum by default, but this number can be set to any desired

value with the command `\lrnumdigit{number}`. The mandatory argument is the desired maximum number of digits.

`\lrplot` The macro `\lrplot{imagewidth}{key1}{key2}{key3}{key4}` really generates the plot. The first argument is the width of the plot in centimeters, while the height is computed according to the distribution of the points or it is recorded with the command `\lrplotparameters` (see below). The other four arguments are referred to the data points, to the lines determined with classical regression, with classical regression with inverted role of the coordinates and with symmetric regression. The four items, i.e. the set of points and the three lines, are drawn or not according to the corresponding string found in `keyi`. Each item is not plotted if the string is the character -, it is plotted in any other case. Furthermore the lines are accompanied by a label made by the corresponding `key`, unless it is just a + or a -.

`\lrplotparameters` The command `\lrplotparameters{diameter}{imageheight}` is used to record some more parameters that are related to the generation of the plot. The first argument `{diameter}` is the diameter in centimeters of the discs that represent the points. It must be a real positive number. The default value is 0.2. The second argument `{imageheight}` is the height of the plot in centimeters. It must be a real number with a special meaning of any non positive value. If the number is positive it is really used as the height and the variables are independently scaled to fill the width and the height of the plot. If the value is not positive or if the command is not used the two variables are scaled with the same factor and the height is computed accordingly. This can generate reasonable results only if the ranges of the two variables are not too much different. All the arguments are mandatory.

The use of `\lrnumdigit` and `\lrplotparameters` is a convenient solution to introduce some flexibility in the output and to completely preserve compatibility with the previous version of the package. Some more elegant solution could be realized with the use of optional arguments and perhaps of key-value pairs, but the implemented solution is preferred because it is extremely simple.

Few words are necessary about the format of the data file. Each record of the file hold the two values related to a point. The two values must be separated by any number (one is needed as a minimum) of space and comma characters. No character different from space can be accepted before the first value and after the second value.

## 7 An example

The data reported here after will be available in `sampledata.txt` and will be used in the examples presented in this section .

|    |        |        |    |        |        |    |        |        |    |        |        |
|----|--------|--------|----|--------|--------|----|--------|--------|----|--------|--------|
| 12 | -0.546 | 0.107  | 18 | -0.203 | -0.292 | 24 | 0.181  | -2.616 | 30 | -0.931 | -1.613 |
| 13 | 1.093  | -0.510 | 19 | 1.517  | 0.779  | 25 | 0.619  | 1.859  | 31 | -1.070 | 0.592  |
| 14 | 1.440  | 1.995  | 20 | 0.559  | -1.341 | 26 | -0.223 | -1.915 | 32 | 2.341  | 0.413  |
| 15 | 1.414  | 0.991  | 21 | -0.462 | -0.437 | 27 | 0.629  | -0.534 | 33 | 1.993  | -0.111 |
| 16 | 0.735  | 1.585  | 22 | -0.785 | -0.661 | 28 | -1.989 | -2.300 | 34 | -2.357 | -0.312 |
| 17 | -1.848 | -0.235 | 23 | -0.558 | 0.397  | 29 | -0.241 | 1.098  | 35 | -1.975 | 0.140  |

The analysis of the sample data and the generation of a numeric table is operated by a code similar to the following (see table 1).

```
\lrfilename{sampledata.txt}
\lrcomputation
\lrnumdigit{5}
\begin{table}
  \lrprint
  \caption{Analysis of ... }
\label{tab:sampledata}\end{table}
```

The generation of some different graphical representation of the data and of the results is operated by a code similar to the following (see figures 2 ).

```
\begin{figure}
\lrplot{10.}{-}{AA}{BB}{S}
\lrplot{10.}{+}{-}{-}{+}
\caption{LEFT The three lines are obtained with the three optimality criteria.
(AA) classical linear regression; (BB) classical linear regression with inverted role
of the coordinates; (S) symmetric linear regression. RIGHT Data points and line
estimated with symmetric linear regression.}
\label{fig:sampledataB} \end{figure}
```

The same data are used to show the effect of the command `\lrplotparameters`. The commands:

```
...
\lrplotparameters{0.12}{5.} \lrplot{4.}{+}{-}{-}{+} \hfill
\lrplotparameters{0.10}{3.} \lrplot{4.}{+}{-}{-}{+} \hfill
\lrplotparameters{0.11}{-1.} \lrplot{4.}{+}{-}{-}{+}
```

...  
generate the figure (3).

## 8 A package for linear regression and the theory behind it

### 8.1 Math preliminaries and notation

The coordinates of a set of  $m$  points on the plane are available. A straight line is searched that optimally approximates the points.

The coordinates of a generic point are  $y_1$  and  $y_2$  and they are collected in the vector  $\underline{y}$ . Any given point is identified with the index  $i$ . Explicit indices  $(\dots)_1$  or  $(\dots)_2$  always refer to the first or second coordinate of a point or to the first or

|  |                     |         |
|--|---------------------|---------|
| Data File  | sampledata.txt      |         |
| Number of points   | 24                  |         |
| Mean values of the coordinates   | -0.02779<br>-0.1217 |         |
| Minimum values of the coordinates  | -2.357<br>-2.616    |         |
| Maximum values of the coordinates  | 2.341<br>1.995      |         |
| Second order moments   | abscissa            | 1.62892 |
|  | mixed               | 0.55492 |
|  | ordinate            | 1.4467  |
| Slope and intercept of optimal line<br>(estimated with errors in ordinate)   | 0.34066<br>-0.11224 |         |
| Slope and intercept of optimal line<br>(estimated with errors in abscissa)   | 2.60702<br>-0.04925 |         |
| Unit vector along the line   | 0.76223<br>0.64729  |         |
| Slope and intercept of optimal line<br>(estimated with symmetric regression) | 0.8492<br>-0.0981   |         |

Table 1: Analysis of the sample data

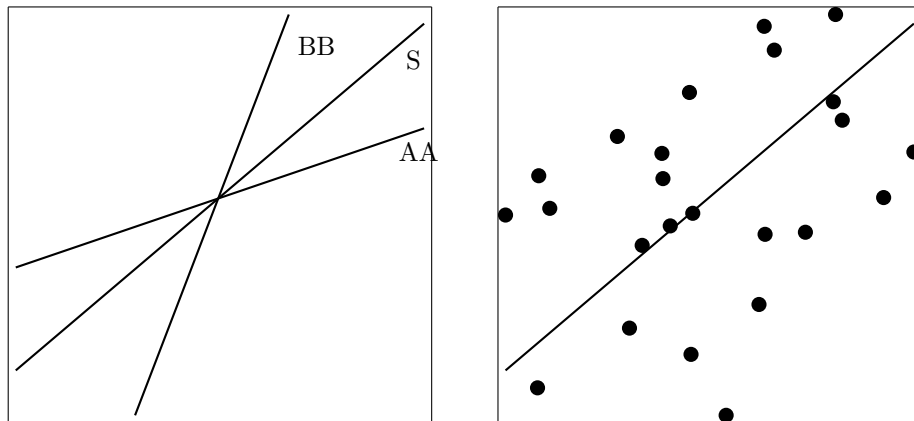


Figure 2: LEFT The three lines are obtained with the three optimality criteria. (AA) classical linear regression; (BB) classical linear regression with inverted role of the coordinates; (S) symmetric linear regression. RIGHT Data points and line estimated with symmetric linear regression.

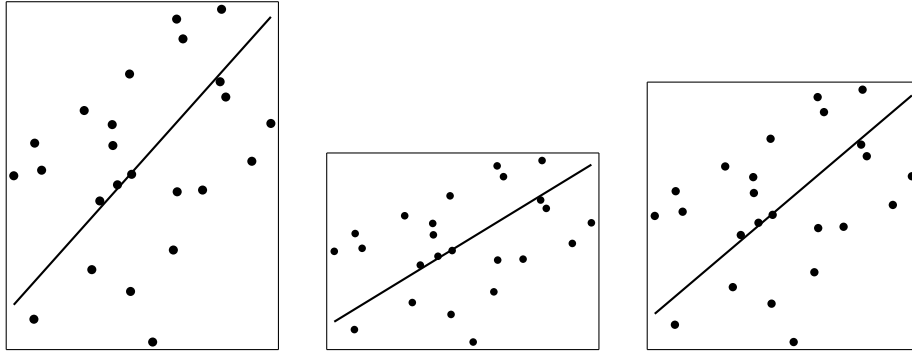


Figure 3: Three different version of a plot to demonstrate the use of the command `\lrplotparameters`

second component of a vector in the plane. Symbolic index  $(\dots)_i$  always refers to the different points. Few formulas require both indices  $(\dots)_{1i}$ ,  $(\dots)_{2i}$ .

With more then two points a criterion of best approximation is needed to select the optimal line that describes the data.

Lower case symbols are used for scalars. Lower case underlined symbols are used for vectors in the plane. Upper case symbols are used for matrices.

It is possible that certain data generate an ambiguity or a singularity in the computation. The following mathematical treatment of the problem do not mention these situations and the code does not deal with them.

## 8.2 Package declaration, required package and definition of variables

The various macro will be provided in a package file that is introduced as usual. Most of the macros require the L<sup>A</sup>T<sub>E</sub>X3 syntax.

```
36 \ProvidesPackage{linearregression}[2024-12-14]
37 \RequirePackage{pict2e}
38 \ExplSyntaxOn
```

The variables used in the package are defined hereafter.

```
39 \ior_new:N \g_BBLR_file_ior
40 \tl_new:N \g_BBLR_file_name_tl
41 \int_new:N \g_BBLR_number_of_points_int
42 \fp_new:N \g_BBLR_abscissa_fp
43 \fp_new:N \g_BBLR_ordinate_fp
44 \fp_new:N \g_BBLR_mean_abscissa_fp
45 \fp_new:N \g_BBLR_mean_ordinate_fp
46 \fp_new:N \g_BBLR_abscissa_SecOrdMoment_fp
47 \fp_new:N \g_BBLR_ordinate_SecOrdMoment_fp
48 \fp_new:N \g_BBLR_mixed_SecOrdMoment_fp
49 \fp_new:N \g_BBLR_slope_A_fp
50 \fp_new:N \g_BBLR_slope_B_fp
51 \fp_new:N \g_BBLR_slope_S_fp
52 \fp_new:N \g_BBLR_intercept_A_fp
53 \fp_new:N \g_BBLR_intercept_B_fp
```



```

54 \fp_new:N \g_BBLR_intercept_S_fp
55 \fp_new:N \g_BBLR_cos_fp
56 \fp_new:N \g_BBLR_sin_fp
57 \fp_new:N \g_BBLR_sig_sin_fp
58 \fp_new:N \g_BBLR_eig_diff_fp
59 \fp_new:N \g_BBLR_diag_diff_fp
60 \tl_new:N \g_BBLR_file_line_tl
61 \fp_new:N \g_BBLR_min_abscissa_fp
62 \fp_new:N \g_BBLR_min_ordinate_fp
63 \fp_new:N \g_BBLR_max_abscissa_fp
64 \fp_new:N \g_BBLR_max_ordinate_fp
65 \fp_new:N \g_BBLR_min_draw_abscissa_fp
66 \fp_new:N \g_BBLR_max_draw_abscissa_fp
67 \bool_new:N \g_BBLR_data_eof_bool
68 \int_new:N \g_BBLR_record_length_int
69 \int_new:N \g_BBLR_rec_count_int
70 \int_new:N \g_BBLR_first_separator_int
71 \int_new:N \g_BBLR_last_separator_int
72 \str_const:Nn \c_BBLR_space_str {-}
73 \str_const:Nn \c_BBLR_comma_str {,}
74 \str_const:Nn \c_BBLR_plus_str {+}
75 \str_const:Nn \c_BBLR_minus_str {-}
76 \bool_new:N \g_BBLR_plot_points_bool
77 \bool_new:N \g_BBLR_plot_lineA_bool
78 \bool_new:N \g_BBLR_plot_lineB_bool
79 \bool_new:N \g_BBLR_plot_lineS_bool
80 \bool_new:N \g_BBLR_two_scale_bool
81 \bool_gset_false:N \g_BBLR_two_scale_bool
82 \fp_new:N \g_BBLR_base_fp
83 \fp_new:N \g_BBLR_height_fp
84 \fp_new:N \g_BBLR_Xbase_fp
85 \fp_new:N \g_BBLR_Xheight_fp
86 \fp_new:N \g_BBLR_XXheight_fp
87 \fp_new:N \g_BBLR_Dabscissa_fp
88 \fp_new:N \g_BBLR_Dordinate_fp
89 \fp_new:N \g_BBLR_diameter_fp
90 \fp_gset:Nn \g_BBLR_diameter_fp{0.2}
91 \fp_new:N \g_BBLR_line_base_length_fp
92 \fp_new:N \g_BBLR_scale_factor_fp
93 \fp_new:N \g_BBLR_scale_factor_H_fp
94 \str_new:N \g_BBLR_point_code_str
95 \str_new:N \g_BBLR_labelA_str
96 \str_new:N \g_BBLR_labelB_str
97 \str_new:N \g_BBLR_labelS_str
98 \int_new:N \g_BBLR_num_dec_dig_int
99 \int_gset:Nn \g_BBLR_num_dec_dig_int{4}

```

### 8.3 Preparing data input

`\lrfilename` The command `\lrfilename` records the file name passed as argument.

```

100 \NewDocumentCommand{\lrfilename}{m}{
101 \tl_gset:Nn \g_BBLR_file_name_tl {#1}
102 }

```

`\lraskfilename` The command `\lraskfilename` asks for the data file name from the terminal.

```

103 \NewDocumentCommand{\lraskfilename}{-}{
104 \ior_get_term:nN {filename ? } \g_BBLR_file_name_tl
105 \tl_trim_spaces:N \g_BBLR_file_name_tl
106 }

```

## 8.4 Computation

`\lrcomputation` The command `\lrcomputation` reads the data file and performs all the relevant computations to solve the proposed problem.

```

107 \NewDocumentCommand{\lrcomputation}{-}{%

```

### 8.4.1 Computation of first and second order moments

In the sequel it will results that the first and second order moments of the data provide everything needed to solve the problem. The barycenter of the data is defined as

$$\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i. \quad (1)$$

It is convenient to scan the data to accumulate the sum that appears in (1). The coordinates of each point are read from the file and they are immediately used. The data are not globally recorded.

```

108 \bool_gset_false:N \g_BBLR_data_eof_bool
109 \int_zero:N \g_BBLR_number_of_points_int
110 \fp_zero:N \g_BBLR_mean_abscissa_fp
111 \fp_zero:N \g_BBLR_mean_ordinate_fp
112 \ior_open:Nn \g_BBLR_file_ior \g_BBLR_file_name_tl
113 \bool_until_do:Nn \g_BBLR_data_eof_bool {
114 \ior_str_get:NN \g_BBLR_file_ior \g_BBLR_file_line_tl
115 \if_eof:w \g_BBLR_file_ior
116 \bool_gset_true:N \g_BBLR_data_eof_bool
117 \else:
118 \int_incr:N \g_BBLR_number_of_points_int
119 \BBLR_decode_data:
120 \fp_gset:Nn \g_BBLR_mean_abscissa_fp
121 {\g_BBLR_mean_abscissa_fp + \g_BBLR_abscissa_fp}
122 \fp_gset:Nn \g_BBLR_mean_ordinate_fp
123 {\g_BBLR_mean_ordinate_fp + \g_BBLR_ordinate_fp}
124 \fi:
125 }

```

Loop ended. Now close the file and divide by the number of points.

```

126 \ior_close:N \g_BBLR_file_ior
127 \fp_gset:Nn \g_BBLR_mean_abscissa_fp
128 {\g_BBLR_mean_abscissa_fp / \g_BBLR_number_of_points_int}
129 \fp_gset:Nn \g_BBLR_mean_ordinate_fp
130 {\g_BBLR_mean_ordinate_fp / \g_BBLR_number_of_points_int}

```

The barycentric coordinates are defined for each point

$$v_i = y_i - \bar{y} \quad (2)$$

and the empirical dispersion matrix is defined as:

$$C = \frac{1}{m} \sum_{i=1}^m v_i v_i^T. \quad (3)$$

Superscript as in  $()^T$  means transpose. The elements of  $C$  are the second order central moments and they are denoted as:

$$C = \begin{pmatrix} k_{11} & k_{12} \\ k_{12} & k_{22} \end{pmatrix}. \quad (4)$$

A second scan of the data is performed to compute the sums that appears in (3) and to determine the the extremal values of the coordinates. Record scan can be regulated by a record counter, because the number of points is now known.

```

131 \fp_zero:N \g_BBLR_abcissa_SecOrdMoment_fp
132 \fp_zero:N \g_BBLR_ordinate_SecOrdMoment_fp
133 \fp_zero:N \g_BBLR_mixed_SecOrdMoment_fp
134 \fp_gset_eq:NN \g_BBLR_min_abcissa_fp \g_BBLR_mean_abcissa_fp
135 \fp_gset_eq:NN \g_BBLR_min_ordinate_fp \g_BBLR_mean_ordinate_fp
136 \fp_gset_eq:NN \g_BBLR_max_abcissa_fp \g_BBLR_mean_abcissa_fp
137 \fp_gset_eq:NN \g_BBLR_max_ordinate_fp \g_BBLR_mean_ordinate_fp
138 \ior_open:Nn \g_BBLR_file_ior \g_BBLR_file_name_tl
139 \int_zero:N \g_BBLR_rec_count_int
140 \int_do_until:nn
141 {\g_BBLR_rec_count_int = \g_BBLR_number_of_points_int}
142 {
143 \ior_str_get:NN \g_BBLR_file_ior \g_BBLR_file_line_tl
144 \int_incr:N \g_BBLR_rec_count_int
145 \BBLR_decode_data:
146 \fp_gset:Nn \g_tmpa_fp
147 {\g_BBLR_abcissa_fp - \g_BBLR_mean_abcissa_fp}
148 \fp_gset:Nn \g_tmpb_fp
149 {\g_BBLR_ordinate_fp - \g_BBLR_mean_ordinate_fp}
150 \fp_gset:Nn \g_BBLR_abcissa_SecOrdMoment_fp
151 {\g_BBLR_abcissa_SecOrdMoment_fp + \g_tmpa_fp * \g_tmpa_fp}
152 \fp_gset:Nn \g_BBLR_mixed_SecOrdMoment_fp
153 {\g_BBLR_mixed_SecOrdMoment_fp + \g_tmpa_fp * \g_tmpb_fp}
154 \fp_gset:Nn \g_BBLR_ordinate_SecOrdMoment_fp
155 {\g_BBLR_ordinate_SecOrdMoment_fp + \g_tmpb_fp * \g_tmpb_fp}
156 \fp_gset:Nn \g_BBLR_min_abcissa_fp
157 {\min(\g_BBLR_min_abcissa_fp, \g_BBLR_abcissa_fp)}
158 \fp_gset:Nn \g_BBLR_min_ordinate_fp
159 {\min(\g_BBLR_min_ordinate_fp, \g_BBLR_ordinate_fp)}
160 \fp_gset:Nn \g_BBLR_max_abcissa_fp
161 {\max(\g_BBLR_max_abcissa_fp, \g_BBLR_abcissa_fp)}
162 \fp_gset:Nn \g_BBLR_max_ordinate_fp
163 {\max(\g_BBLR_max_ordinate_fp, \g_BBLR_ordinate_fp)}
164 }
165 \ior_close:N \g_BBLR_file_ior
166 \fp_gset:Nn \g_BBLR_abcissa_SecOrdMoment_fp
167 {\g_BBLR_abcissa_SecOrdMoment_fp / \g_BBLR_number_of_points_int}
168 \fp_gset:Nn \g_BBLR_mixed_SecOrdMoment_fp
169 {\g_BBLR_mixed_SecOrdMoment_fp / \g_BBLR_number_of_points_int}

```

```

170 \fp_gset:Nn \g_BBLR_ordinate_SecOrdMoment_fp
171 {\g_BBLR_ordinate_SecOrdMoment_fp / \g_BBLR_number_of_points_int}
172 \fp_gset:Nn \g_BBLR_Dabscissa_fp
173 {\g_BBLR_max_abscissa_fp - \g_BBLR_min_abscissa_fp }
174 \fp_gset:Nn \g_BBLR_Dordinate_fp
175 {\g_BBLR_max_ordinate_fp - \g_BBLR_min_ordinate_fp }

```

A single pass algorithm exists, but it is numerically less stable.

### 8.4.2 Classical linear regression

A line in the plane is described by the equation

$$y_2 = ay_1 + b \quad (5)$$

that contains the parameters  $a$  and  $b$ . For each point it is possible to define a distance or a discrepancy of the experimental data with respect to the model. In the given problem the second coordinate is much more affected by errors than the first coordinate. It is therefore reasonable to define the approximation error of each point as

$$e_i = y_{2i} - ay_{1i} - b \quad (6)$$

i.e. the difference between the empirical value  $y_{2i}$  and its model counterpart  $ay_{1i} + b$ . This is the already mentioned classical linear regression. In fact the absolute value of  $e_i$  expressed in (6) is the length of the segments shown in the leftmost scheme of figure (1). The global discrepancy between the data and the model is measured by the least square objective function defined by:

$$\psi = \sum_{i=1}^m e_i^2 \quad (7)$$

and the parameters  $a$  and  $b$  will be determined just by the minimization of the function  $\psi$  defined in (7).

In the present treatment of the regression problem as a pure approximation problem the definition of  $\psi$  in (7) seems quite arbitrary. It is anyway a convenient choice.

Expression (6) can be rewritten in the different form

$$e_i = v_{2i} - av_{1i} + \bar{y}_2 - a\bar{y}_1 - b \quad (8)$$

so that the function to be minimized can be expressed as the sum of two quadratic functions:

$$\psi = \sum_{i=1}^m (v_{2i} - av_{1i})^2 + m(\bar{y}_2 - a\bar{y}_1 - b)^2 \quad (9)$$

and the minimum can be attained considering the two terms one at a time. The second term in the right-hand side of (9) vanishes if the choice of  $b$  is:

$$b = \bar{y}_2 - a\bar{y}_1. \quad (10)$$

The first term in the right-hand side of (9) becomes:

$$\psi_{(a)} = m(k_{22} - 2ak_{12} + a^2k_{11}). \quad (11)$$

Searching the minimum of  $\psi$  w.r.t.  $a$  is therefore the search of the abscissa of the vertex of a parabola with axis parallel to the second coordinated axis. The result is:

$$a = k_{12}/k_{11} \quad (12)$$

Now the slope  $a$  and the intercept  $b$  can be actually computed.

```
176 \fp_gset:Nn \g_BBLR_slope_A_fp
177 {\g_BBLR_mixed_SecOrdMoment_fp / \g_BBLR_abscissa_SecOrdMoment_fp }
178 \fp_gset:Nn \g_BBLR_intercept_A_fp
179 {\g_BBLR_mean_ordinate_fp - \g_BBLR_slope_A_fp * \g_BBLR_mean_abscissa_fp}
```

The empirical data and the estimated values of  $a$  and  $b$  can be used to compute the value actually attained by the residuals  $e_i$  and by the function  $\psi$ . Then the index

$$\hat{\sigma}_0^2 = \psi/(m - 2) \quad (13)$$

can be used to evaluate the general quality of the data and of the model. This claim is clearly quite generic. A complete understanding of this evaluation would require to treat the linear regression problem in the framework of the probabilistic estimation theory. The used notation is derived from that theory.

If the role of the two coordinates is exchanged the result for  $a$  becomes (still with reference to (5))

$$a = k_{22}/k_{12}. \quad (14)$$

The slope and the intercept can be computed accordingly.

```
180 \fp_gset:Nn \g_BBLR_slope_B_fp
181 {\g_BBLR_ordinate_SecOrdMoment_fp / \g_BBLR_mixed_SecOrdMoment_fp}
182 \fp_gset:Nn \g_BBLR_intercept_B_fp
183 {\g_BBLR_mean_ordinate_fp - \g_BBLR_slope_B_fp * \g_BBLR_mean_abscissa_fp}
```

### 8.4.3 Symmetric linear regression

If both the coordinates of the experimental points are affected by the same uncertainty it is advisable to use a more symmetric optimality criterion and it is convenient to use a different model equation.

The same line can be described by a different equation, i.e.

$$x_1y_1 + x_2y_2 = f \quad (15)$$

or in vector form:

$$\underline{x}^T \underline{y} = f. \quad (16)$$

The parameters in (16) are the scalar  $f$  and the elements of the vector  $\underline{x}$ , i.e.  $x_1$  and  $x_2$ . The line described by (16) is obviously invariant when the three parameters are simultaneously scaled by a constant. The normalization condition

$$\underline{x}^T \underline{x} = 1, \quad (17)$$

supplemented by  $f \geq 0$ , is quite convenient because the parameters will assume a significant geometrical meaning:  $\underline{x}$  is a unit vector orthogonal to the line and  $f$  is the distance of the line from the origin. The expression

$$d = f - \underline{x}^T \underline{y} \quad (18)$$

is the distance of the generic point  $\underline{y}$  from the line with a sign that is positive for points on the same side of the origin.

The distance of each given point from the desired optimal line is denoted by  $d_i$ . It has a clear intrinsic geometrical meaning and it does not privileges one coordinate w.r.t. the other. The function to be minimized by the optimal line is

$$\phi = \frac{1}{m} \sum_{i=1}^m d_i^2. \quad (19)$$

The parameters of (16) are determined by the minimization of the function  $\phi$  that can be expressed as:

$$\phi = \frac{1}{m} \sum_{i=1}^m (\underline{x}^T \underline{y}_i - f)^2 \quad (20)$$

and then, after some algebraic manipulations:

$$\phi = \underline{x}^T C \underline{x} + (f - \underline{x}^T \bar{\underline{y}})^2. \quad (21)$$

The function  $\phi$  is composed (as it was the function  $\psi$ ) by the sum of two parts. The second term in the right-hand side of (21) vanishes if the choice of  $f$  is:

$$f = \underline{x}^T \bar{\underline{y}}. \quad (22)$$

Both (10) and (22) means that the optimal line includes the mean point. Then it is necessary to minimize the function

$$\phi_{(\underline{x})} = \underline{x}^T C \underline{x} \quad (23)$$

with the constrain  $\underline{x}^T \underline{x} = 1$ . It can be proved that the function  $\phi_{(\underline{x})}$  is stationary if  $\underline{x}$  is an eigenvector of  $C$ .

The function  $\phi_{(\underline{x})}$  and the constrain must be combined using a Lagrange multiplier:

$$\Phi = \underline{x}^T C \underline{x} + \lambda(1 - \underline{x}^T \underline{x}). \quad (24)$$

Then the stationarity points of  $\Phi$  must be determined. Equating to zero the derivatives of  $\Phi$  gives

$$C \underline{x} = \lambda \underline{x} \quad (25)$$

i.e.  $\underline{x}$  is an eigenvector of  $C$ .

The same result is obtained with the following argument. The function  $\phi_{(\underline{x})}$  is stationary if its first variation is zero. The variation of  $\underline{x}$  is named  $\underline{\delta}$ . It must respect the constrain, that becomes  $\underline{\delta}^T \underline{x} = 0$ . The first variation of  $\phi_{(\underline{x})}$  is  $2\underline{\delta}^T C \underline{x}$ , and it is zero if and only if the following implication is valid:  $\underline{\delta}^T \underline{x} = 0 \implies \underline{\delta}^T C \underline{x} = 0$ , and the implication is valid if and only if the vector  $C \underline{x}$  has the same direction of  $\underline{x}$ , i.e. if  $\underline{x}$  is an eigenvector of  $C$ .

The result on the optimal line can be described geometrically in the following way: (i) the optimal line includes the barycenter of the data; (ii) the optimal line is orthogonal to the eigenvector of  $C$  corresponding to the minimum eigenvalue. The obtained result can be generalized in  $\mathbb{R}^n$ . A set of points in  $\mathbb{R}^n$  must be approximated by an  $(n - 1)$ -dimensional affine subspace. (Other more general situations can be considered.)

The trace of the matrix  $C$ , denoted as  $\text{tr}(C)$ , is a measure of the global dispersion of the set of points. The minimum eigenvalue  $\lambda_{\min}$  of  $C$  is a measure of the dispersion of the set of points with respect to the optimal affine subspace. Therefore the index

$$\frac{n\lambda_{\min}}{\text{tr}(C)} \quad (26)$$

can be used as an indicator of the relative residual dispersion of the data around the optimal affine subspace. The defined index is dimensionless and it is always between 0 and 1.

For the actual computation of  $\underline{x}$  it is convenient to consider the spectral factorization of the matrix  $C$ , i.e.  $C = X\Lambda X^T$  where  $\Lambda$  is a diagonal matrix whose diagonal elements are the eigenvalues of  $C$  and  $X$  is an orthonormal matrix whose columns are the eigenvectors of  $C$ . The spectral factorization exists for any symmetric matrix, but it is specially simple for a  $2 \times 2$  matrix.

$$\begin{pmatrix} k_{11} & k_{12} \\ k_{12} & k_{22} \end{pmatrix} = \begin{pmatrix} c & -s \\ s & c \end{pmatrix} \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \quad (27)$$

The eigenvalues can be easily obtained because their sum is the trace of  $C$

$$\lambda_1 + \lambda_2 = k_{11} + k_{22} \quad (28)$$

and their product is the determinant of the same matrix. Therefore after some manipulations it results:

$$\lambda_1 - \lambda_2 = \sqrt{(k_{11} - k_{22})^2 + 4k_{12}^2} \quad (29)$$

and the two eigenvalues are then immediately obtained.

It is convenient to compute the difference of the two diagonal elements of the dispersion matrix and the difference of its eigenvalues.

```

184 \fp_gset:Nn \g_BBLR_diag_diff_fp
185 {\g_BBLR_abscissa_SecOrdMoment_fp - \g_BBLR_ordinate_SecOrdMoment_fp}
186 \fp_gset:Nn \g_BBLR_eig_diff_fp
187 {\sqrt{\g_BBLR_diag_diff_fp * \g_BBLR_diag_diff_fp +
188 4 * \g_BBLR_mixed_SecOrdMoment_fp * \g_BBLR_mixed_SecOrdMoment_fp}}

```

The computation of  $c$  and  $s$  is obtained from (27) taking into account that  $c^2 + s^2 = 1$ . From (27) it results:

$$k_{11} - k_{22} = (\lambda_1 - \lambda_2)(c^2 - s^2) \quad (30)$$

and also

$$k_{12} = (\lambda_1 - \lambda_2)cs \quad (31)$$

that is only used to determine the sign of  $cs$ . The expression for the parameters  $c$  and  $s$  are:

$$c = \sqrt{\frac{1}{2} + \frac{k_{11} - k_{22}}{2(\lambda_1 - \lambda_2)}} \quad (32)$$

$$s = \text{sgn}(k_{12})\sqrt{\frac{1}{2} - \frac{k_{11} - k_{22}}{2(\lambda_1 - \lambda_2)}} \quad (33)$$

The parameters  $s$  and  $c$  are the sine and cosine of the angle between the axis of  $y_1$  and the eigenvector corresponding to the maximum eigenvalue. They are computed using the already defined elements.

```

189 \fp_gset:Nn \g_BBLR_cos_fp%
190 {sqrt((1 + \g_BBLR_diag_diff_fp / \g_BBLR_eig_diff_fp) / 2)}
191 \fp_gset:Nn \g_BBLR_sig_sin_fp {\fp_sign:n {\g_BBLR_mixed_SecOrdMoment_fp}}
192 \fp_gset:Nn \g_BBLR_sin_fp
193 {\g_BBLR_sig_sin_fp*sqrt((1-\g_BBLR_diag_diff_fp / \g_BBLR_eig_diff_fp) / 2)}

```

The vector  $\underline{x}$  is :

$$\underline{x} = \text{sgn}(-s\bar{y}_1 + c\bar{y}_2) \begin{pmatrix} -s \\ c \end{pmatrix}. \quad (34)$$

The parameter  $a$  of model (5) can be obtained as:

$$a = s/c \quad (35)$$

Now the slope and the intercept of the optimal line corresponding to the symmetric criterion can be computed.

```

194 \fp_gset:Nn \g_BBLR_slope_S_fp
195 {\g_BBLR_sin_fp / \g_BBLR_cos_fp }
196 \fp_gset:Nn \g_BBLR_intercept_S_fp
197 {\g_BBLR_mean_ordinate_fp -
198 \g_BBLR_slope_S_fp * \g_BBLR_mean_abscissa_fp}
199 }

```

The theoretical treatment of the proposed problem and the implementation of its numerical solution end here.

## 8.5 Print of table of results

`\lrprint` The command `\lrprint` prints some info on the data and the results of the computations in tabular form.

```

200 \NewDocumentCommand{\lrprint}{-}{
201 \begin{center}
202 \begin{tabular}{| l | r |} \hline
203 Data-File & \g_BBLR_file_name_tl & \\
204 \hline
205 Number-of-points & \int_use:N\g_BBLR_number_of_points_int & \\
206 \hline
207 Mean-values-of-the-coordinates & & \%
208 $\BBLR_print_fp:N \g_BBLR_mean_abscissa_fp$ \& \\
209 $\BBLR_print_fp:N \g_BBLR_mean_ordinate_fp$ \& \hline
210 Minimum-values-of-the-coordinates & & \%
211 $\BBLR_print_fp:N \g_BBLR_min_abscissa_fp$ \& \\
212 $\BBLR_print_fp:N \g_BBLR_min_ordinate_fp$ \& \hline
213 Maximum-values-of-the-coordinates & & \%
214 $\BBLR_print_fp:N \g_BBLR_max_abscissa_fp$ \& \\
215 $\BBLR_print_fp:N \g_BBLR_max_ordinate_fp$ \& \hline
216 {Second-order-moments}\phantom{xxxxxxxx}{abscissa} & & \%
217 $\BBLR_print_fp:N \g_BBLR_abscissa_SecOrdMoment_fp$ \& \\
218 \multicolumn{1}{|r|}{mixed} & & \%
219 $\BBLR_print_fp:N \g_BBLR_mixed_SecOrdMoment_fp$ ~ \& \\
220 \multicolumn{1}{|r|}{ordinate} & & \%

```



```

221 $\BBLR_print_fp:N \g_BBLR_ordinate_SecOrdMoment_fp$ \\ \hline
222 Slope-and-intercept-of-optimal-line & $\BBLR_print_fp:N
223 \g_BBLR_slope_A_fp$ \\
224 (estimated-with-errors-in-ordinate)&$\BBLR_print_fp:N
225 \g_BBLR_intercept_A_fp$\\ \hline
226 Slope-and-intercept-of-optimal-line & $\BBLR_print_fp:N
227 \g_BBLR_slope_B_fp$ \\
228 (estimated-with-errors-in-abscissa)&$\BBLR_print_fp:N
229 \g_BBLR_intercept_B_fp$\\ \hline
230 Unit-vector-along-the-line & $\BBLR_print_fp:N
231 \g_BBLR_cos_fp$ \\
232 & $\BBLR_print_fp:N \g_BBLR_sin_fp$ \\
233 Slope-and-intercept-of-optimal-line &$\BBLR_print_fp:N
234 \g_BBLR_slope_S_fp$ \\
235 (estimated-with-symmetric-regression) &
236 $\BBLR_print_fp:N \g_BBLR_intercept_S_fp$\\ \hline
237 \end{tabular}
238 \end{center}
239 }

```

`\lrnumdigit` The command `\lrnumdigit` records the maximum number of digits used in the table for computed numbers. The default value is 4.

```

240 \NewDocumentCommand{\lrnumdigit}{m}{
241 \int_gset:Nn \g_BBLR_num_dec_dig_int{#1}
242 }

```

## 8.6 Plot of points and lines

`\lrplot` The command `\lrplot` produce a framed plot of the data and of the regression line(s). The size of the plot and its actual content are determined by the arguments.

```

243 \NewDocumentCommand{\lrplot}{m}{}

```

The plotting area is divided into a main plotting area for the representation of points and line(s) and a small surrounding free space. The height is computed taking into account the distribution of the points, if it is not explicitly given..

```

244 \fp_gset:Nn \g_BBLR_base_fp {#1}
245 \fp_gset:Nn \g_BBLR_Xbase_fp {\g_BBLR_base_fp - 0.6}
246 \fp_gset:Nn \g_BBLR_scale_factor_fp{\g_BBLR_Xbase_fp /
247 \g_BBLR_Dabscissa_fp}
248 \bool_if:NTF \g_BBLR_two_scale_bool
249 {\fp_gset:Nn \g_BBLR_height_fp {\g_BBLR_XXheight_fp}
250 \fp_gset:Nn \g_BBLR_Xheight_fp {\g_BBLR_height_fp - 0.6}
251 \fp_gset:Nn \g_BBLR_scale_factor_H_fp{\g_BBLR_Xheight_fp /
252 \g_BBLR_Dordinate_fp}
253 }
254 {\fp_gset:Nn \g_BBLR_Xheight_fp {\g_BBLR_Dordinate_fp *
255 \g_BBLR_scale_factor_fp}
256 \fp_gset:Nn \g_BBLR_height_fp {\g_BBLR_Xheight_fp + 0.6}
257 \fp_gset:Nn \g_BBLR_scale_factor_H_fp {\g_BBLR_scale_factor_fp }
258 }

```

The information about the items to be plotted is in the remaining arguments.

```

259 \str_gset:Nn \g_BBLR_point_code_str {#2}

```

```

260 \str_gset:Nn \g_BBLR_labelA_str      {#3}
261 \str_gset:Nn \g_BBLR_labelB_str      {#4}
262 \str_gset:Nn \g_BBLR_labelS_str      {#5}
263 \bool_gset:Nn \g_BBLR_plot_points_bool
264 {!(\str_if_eq_p:NN \g_BBLR_point_code_str \c_BBLR_minus_str)}
265 \bool_gset:Nn \g_BBLR_plot_lineA_bool
266 {!(\str_if_eq_p:NN \g_BBLR_labelA_str \c_BBLR_minus_str)}
267 \bool_gset:Nn \g_BBLR_plot_lineB_bool
268 {!(\str_if_eq_p:NN \g_BBLR_labelB_str \c_BBLR_minus_str)}
269 \bool_gset:Nn \g_BBLR_plot_lineS_bool
270 {!(\str_if_eq_p:NN \g_BBLR_labelS_str \c_BBLR_minus_str)}
The unit of length is 1 centimeter. The plotting area is framed.
271 \setlength{\unitlength}{1.0cm}
272 \fp_gset:Nn \g_tmpa_fp {\g_BBLR_Xbase_fp +0.2}
273 \fp_gset:Nn \g_tmpb_fp {\g_BBLR_Xheight_fp +0.1}
274 \begin{picture}(\fp_use:N\g_BBLR_base_fp,\fp_use:N\g_BBLR_height_fp)(-0.3,-0.3)
275 \put(-0.1,-0.1){\line(1,0){\fp_use:N\g_tmpa_fp}}
276 \put(-0.1,\fp_use:N\g_tmpb_fp){\line(1,0){\fp_use:N\g_tmpa_fp}}
277 \fp_gset:Nn \g_tmpa_fp {\g_tmpa_fp -0.1}
278 \fp_gset:Nn \g_tmpb_fp {\g_tmpb_fp +0.1}
279 \put(-0.1,-0.1){\line(0,1){\fp_use:N\g_tmpb_fp}}
280 \put(\fp_use:N\g_tmpa_fp,-0.1){\line(0,1){\fp_use:N\g_tmpb_fp}}
The plot of points and line(s) is obtained using auxiliary functions.
281 \thicklines
282 \bool_if:nT {\g_BBLR_plot_points_bool}{\BBLR_plot_points:}
283 \bool_if:nT {\g_BBLR_plot_lineA_bool}{
284 \BBLR_draw_line:NNN \g_BBLR_slope_A_fp\g_BBLR_intercept_A_fp
285 \g_BBLR_labelA_str}
286 \bool_if:nT {\g_BBLR_plot_lineB_bool}{
287 \BBLR_draw_line:NNN \g_BBLR_slope_B_fp\g_BBLR_intercept_B_fp
288 \g_BBLR_labelB_str}
289 \bool_if:nT {\g_BBLR_plot_lineS_bool}{
290 \BBLR_draw_line:NNN \g_BBLR_slope_S_fp\g_BBLR_intercept_S_fp
291 \g_BBLR_labelS_str}
292 \end{picture}
293 }%

```

`\lrplotparameters` The command `\lrplotparameters` records two parameters related to the generation of the plot and sets a boolean variable.

```

294 \NewDocumentCommand{\lrplotparameters}{mm}{
295 \fp_gset:Nn \g_BBLR_diameter_fp{#1}
296 \fp_gset:Nn \g_BBLR_XXheight_fp{#2}
297 \bool_gset:Nn \g_BBLR_two_scale_bool
298 {\fp_compare_p:nNn {\g_BBLR_XXheight_fp}>{0.}}
299 }

```

## 8.7 Functions for internal use

The functions listed here after are for internal use and are just minimally documented.

`\BBLR_decode_data:` The function `\BBLR_decode_data:` extract two numeric values from the string

read from the file. Some clever actions are necessary because a so called csv file sometime do not contains the separating commas.

```

300 \cs_new_protected:Nn \BBLR_decode_data: {
301 \tl_trim_spaces:N \g_BBLR_file_line_tl
302 \int_gzero:N \g_tmpa_int
303 \int_gzero:N \g_BBLR_first_separator_int
304 \int_gzero:N \g_BBLR_last_separator_int
305 \int_gset:Nn \g_BBLR_record_length_int {
306 \str_count:N \g_BBLR_file_line_tl}
307 \str_map_variable:Nnn \g_BBLR_file_line_tl \g_tmpa_str {
308 \int_gincr:N \g_tmpa_int
309 \bool_lazy_or:nnTF
310 {\str_if_eq_p:NN \g_tmpa_str \c_BBLR_comma_str}
311 {\str_if_eq_p:NN \g_tmpa_str \c_BBLR_space_str}
312 {\int_gset_eq:NN \g_BBLR_last_separator_int \g_tmpa_int
313 \int_if_zero:nTF {\g_BBLR_first_separator_int}
314 {\int_gset_eq:NN \g_BBLR_first_separator_int \g_tmpa_int
315 }{\prg_do_nothing:}
316 }{\prg_do_nothing:}
317 }
318 \int_gincr:N \g_BBLR_last_separator_int
319 \int_gdecr:N \g_BBLR_first_separator_int
320 \fp_gset:Nn \g_BBLR_abscissa_fp{
321 \str_range:Nnn \g_BBLR_file_line_tl{1}{\g_BBLR_first_separator_int}}
322 \fp_gset:Nn \g_BBLR_ordinate_fp{
323 \str_range:Nnn \g_BBLR_file_line_tl
324 {\g_BBLR_last_separator_int}{\g_BBLR_record_length_int}}
325 }

```

`\BBLR_plot_points:` The function `\BBLR_plot_points:` scans the data file to read the coordinates and it draws a circle for each point.

```

326 \cs_new_protected:Nn \BBLR_plot_points: {
327 \ior_open:Nn \g_BBLR_file_ior \g_BBLR_file_name_tl
328 \int_zero:N \g_BBLR_rec_count_int
329 \int_do_until:nn
330 {\g_BBLR_rec_count_int = \g_BBLR_number_of_points_int}
331 {
332 \ior_str_get:NN \g_BBLR_file_ior \g_BBLR_file_line_tl
333 \int_incr:N \g_BBLR_rec_count_int
334 \BBLR_decode_data:
335 \fp_gset:Nn \g_tmpa_fp{(\g_BBLR_abscissa_fp-\g_BBLR_min_abscissa_fp)*
336 \g_BBLR_scale_factor_fp}
337 \fp_gset:Nn \g_tmpb_fp{(\g_BBLR_ordinate_fp-\g_BBLR_min_ordinate_fp)*
338 \g_BBLR_scale_factor_H_fp}
339 \put(\fp_use:N\g_tmpa_fp, \fp_use:N\g_tmpb_fp){
340 {\circle*{\fp_use:N\g_BBLR_diameter_fp}}}
341 }
342 \ior_close:N \g_BBLR_file_ior
343 }

```

`\BBLR_draw_line:NNN` The function `\BBLR_draw_line:NNN` draws the line. The first two parameters given as arguments are the slope and the intercept. The third parameter is a label.

The next code finds the intersection of the line with the plotting area.

```

344 \cs_new_protected:Nn \BBLR_draw_line:NNN {
345 \fp_gset:Nn \fp_tmpa_fp {#1 * \g_BBLR_min_abcissa_fp + #2 }
346 \fp_compare:nTF{\fp_tmpa_fp > \g_BBLR_max_ordinate_fp}{
347 \fp_gset:Nn \g_BBLR_min_draw_abcissa_fp {(\g_BBLR_max_ordinate_fp -#2) / #1}
348 }{
349 \fp_compare:nTF{\fp_tmpa_fp < \g_BBLR_min_ordinate_fp}{
350 \fp_gset:Nn \g_BBLR_min_draw_abcissa_fp {(\g_BBLR_min_ordinate_fp - #2) / #1}
351 }{
352 \fp_gset:Nn \g_BBLR_min_draw_abcissa_fp { \g_BBLR_min_abcissa_fp }
353 }}
354 \fp_gset:Nn \fp_tmpa_fp {#1 * \g_BBLR_max_abcissa_fp + #2 }
355 \fp_compare:nTF{\fp_tmpa_fp > \g_BBLR_max_ordinate_fp}{
356 \fp_gset:Nn \g_BBLR_max_draw_abcissa_fp {(\g_BBLR_max_ordinate_fp -#2) / #1}
357 }{
358 \fp_compare:nTF{\fp_tmpa_fp < \g_BBLR_min_ordinate_fp}{
359 \fp_gset:Nn \g_BBLR_max_draw_abcissa_fp { (\g_BBLR_min_ordinate_fp - #2) / #1}
360 }{
361 \fp_gset:Nn \g_BBLR_max_draw_abcissa_fp { \g_BBLR_max_abcissa_fp }
362 }}

```

Some parameters (coordinates of starting point, base-length, scaled slope) are computed and the line is drawn.

```

363 \fp_gset:Nn \fp_tmpa_fp {(\g_BBLR_min_draw_abcissa_fp -
364 \g_BBLR_min_abcissa_fp)* \g_BBLR_scale_factor_fp}
365 \fp_gset:Nn \fp_tmpb_fp {(#1 * \g_BBLR_min_draw_abcissa_fp + #2 -
366 \g_BBLR_min_ordinate_fp)* \g_BBLR_scale_factor_H_fp}
367 \fp_gset:Nn \fp_BBLR_line_base_length_fp{(\g_BBLR_max_draw_abcissa_fp -
368 \g_BBLR_min_draw_abcissa_fp) * \g_BBLR_scale_factor_fp}
369 \fp_gset:Nn \fp_scaled_slope_fp
370 {(#1 * \g_BBLR_scale_factor_H_fp / \g_BBLR_scale_factor_fp)}
371 \put(\fp_use:N\fp_tmpa_fp, \fp_use:N\fp_tmpb_fp){
372 \line(1., \fp_use:N\fp_scaled_slope_fp){\fp_use:N\fp_BBLR_line_base_length_fp}}

```

The third parameter is used as a label, if it is not a +.

```

373 \bool_if:nF {\str_if_eq_p:NN #3 \c_BBLR_plus_str}{
374 \fp_gset:Nn \fp_tmpa_fp
375 {0.08 * \g_BBLR_min_draw_abcissa_fp + 0.92 * \g_BBLR_max_draw_abcissa_fp}
376 \fp_gset:Nn \fp_tmpb_fp {#1 * \fp_tmpa_fp + #2 }
377 \fp_gset:Nn \fp_tmpa_fp
378 {(\fp_tmpa_fp-\g_BBLR_min_abcissa_fp)*\g_BBLR_scale_factor_fp
379 + 0.3 * #1 /sqrt(1.+#1*#1)}
380 \fp_gset:Nn \fp_tmpb_fp
381 {(\fp_tmpb_fp-\g_BBLR_min_ordinate_fp)* \g_BBLR_scale_factor_fp
382 - 0.3 /sqrt(1.+#1*#1)}
383 \put(\fp_use:N\fp_tmpa_fp, \fp_use:N\fp_tmpb_fp){#3}
384 }
385 }

```

`\BBLR_print_fp:N` The function `\BBLR_print_fp:N` is used to print floating point numbers with a limited number of decimal digit.

```

386 \cs_new_protected:Nn \BBLR_print_fp:N {%
387 {\fp_eval:n{trunc(#1 ,\g_BBLR_num_dec_dig_int)}}
388 }

```

```

389 \ExplSyntaxOff

```

## 9 Versions

2024-06-10: first post on CTAN.

2024-11-23: small changes in the text, different format of numbers in the table.

2024-12-14: more flexibe output, some changes in the text.

## 10 Acknowledgments

The colleagues Paolo Zatelli, Alfonso Vitti and Giulia Graldi read some preliminary version of this text and suggested several improvements. Claudio Beccari usefully commented about the first version when I was preparing the second version.

## 11 Citations and references

### Mathematics

The books by Lang [7] and by Strang [14] give all the background on linear algebra. The texts by Sansò [12,13] (in italian) treat the teory of probability and its application to metrology. The paper by Pearson [11] is the oldest text that I have found on the symmetric regression, or total regression.

### Programming

The two documents [9,10] are the fountamental and official guide for L<sup>A</sup>T<sub>E</sub>X<sub>3</sub> programming. The books by Donald Knuth [6] and Leslie Lamport [8] are still essential references. The papers by Enrico Gregorio [1,2,3,4,5] explain some general and some special aspect of L<sup>A</sup>T<sub>E</sub>X<sub>3</sub> programming.

### References

- [1] Enrico Gregorio, *L<sup>A</sup>T<sub>E</sub>X<sub>3</sub>: un nuovo gioco per i maghi e per diventarlo*, ArsTeXnica **14** (2012), 41–47.
- [2] Enrico Gregorio, *Liste, cicli, L<sup>A</sup>T<sub>E</sub>X<sub>3</sub>*, ArsTeXnica **22** (2016), 69–77.
- [3] Enrico Gregorio, *Condizionali in L<sup>A</sup>T<sub>E</sub>X*, ArsTeXnica **24** (2017), 37–44.
- [4] Enrico Gregorio, *Funzioni e exp13*, ArsTeXnica **30** (2020), 36–45.
- [5] Enrico Gregorio, *Functions and exp13*, TUGboat **41** (2020), no. 3, 299–307.
- [6] Donald Knuth, *The TeXbook*, American Mathematical Society and Addison-Wesley, 1986.
- [7] Serge Lang, *Linear Algebra*, Springer-Verlag, Berlin Heidelberg, 1987.
- [8] Leslie Lamport, *LaTeX - A document preparation system (2nd ed. )*, Addison-Wesley, 1994. something interesting in the fist edition, too.
- [9] The LaTeX project team, *The exp13 package and LaTeX3 programming* (2024). file: `exp13.pdf` available in CTAN in l3kernel.
- [10] The LaTeX project team, *The L<sup>A</sup>T<sub>E</sub>X<sub>3</sub> interface* (2024). file: `interface3.pdf` available in CTAN in l3kernel.

- [11] Karl Pearson, *On lines and planes of closest fit to systems of points in space*, Philosophical Magazine **2** (1901), no. 11, 559–572.
- [12] Fernando Sansò, *Elementi di teoria della probabilità*, Città-Studi, Milano, 1996. see: <http://www.geolab.polimi.it/text-books/>.
- [13] Fernando Sansò, *La teoria della stima*, Città-Studi, Milano, 1996. see: <http://www.geolab.polimi.it/text-books/>.
- [14] Gilbert Strang, *Introduction to linear algebra*, Wellesley-Cambridge press, 2009.

\*\*\*