

## FINAL project part B: linear regression

### Play by the rules:

1. The use of any form of cycles or symbolic toolbox functions in your implementation is prohibited. Using the vector features of MATLAB is mandatory.
2. It is prohibited to use the function `polyfit()` or the `\` operator when implementing the `linreg()` function. It is mandatory to implement it using the equations derived in the lecture, potentially using the built-in functions `mean()`, `var()`, `cov()`.

### Exercise 1: Preparation

Read the documentation for the following MATLAB functions: `rand`, `randn`, `randi`, `mean`, `var`, `cov`, `polyfit`. Read the short help provided by the `>>help function` command, and also the more detailed help provided by the `>>doc function` command. Except where explicitly stated otherwise, implement all code into a single .m file using cell mode (use the cell separator: `%` ). For further examples look at the MATLAB linear regression documentation online (in english): [https://www.mathworks.com/help/matlab/data\\_analysis/linear-regression.html](https://www.mathworks.com/help/matlab/data_analysis/linear-regression.html)

Further documentation in slovak available here: <https://kurzy.kpi.fei.tuke.sk/nm/student/13.html>

Implement exercises 5 (copy from part A), 7, 8, 10 directly into `lin_regression.m`.

### Exercise 6: Implement linear regression coefficient calculation algorithm

Open the provided .m file `linreg.m` with a function prototype given as :

`function [alpha, beta] = linreg( x, y )`. Implement the body of this function to perform the calculation of simple linear regression coefficients  $\alpha$ ,  $\beta$ . The necessary formulas are provided below in the “Equations” section. Use, however, the MATLAB built-in functions `mean()`, `var()`, `cov()` in your implementation.

Note that when using linear interpolation, we denoted the coefficients of the linear function:  $f(x) = ax + b$ . Then, when describing regression, we denoted the coefficients:  $f(x) = \alpha + \beta x$ . Furthermore, various MATLAB functions may use another notation, such as:  $f(x) = p_1 x + p_2$ . It is our duty to not be confused by this fact of life :)

### Exercise 7: Perform linear regression – find the actual coefficients

For the  $x$  coordinates and both noisy datasets  $Y_1$ ,  $Y_2$  from Exercise 5 find the coefficients  $\alpha$ ,  $\beta$  of the simple linear regression:

- a) using your own function `linreg()`
- b) using the MATLAB function `polyfit()`
- c) using the `\` operator

Observe the results.

Note: The result for each of the points a) b) c) will be the tuple  $(\alpha, \beta)$  of the linear regression coefficients. Together 6 values. The results obtained in a) may be very slightly different from results of b) and c). This is a known property of floating point calculations and number representation. To compare the equality of these results, the `==` operator is not an appropriate tool.

### Exercise 8: Perform linear regression - plot fitted data

Use the `polyval()` function to find the data  $Y_{f1}$  and  $Y_{f2}$  for both noisy datasets  $Y_1$  and  $Y_2$  using the previously found linear regression coefficients  $\alpha$ ,  $\beta$  for the  $x$  vector defined in Exercise 5.

Use the `plot()`, `subplot()` and `scatter()` functions to plot three plots side by side in one figure window:

1. The noisy data  $Y_1$  along with the fitted linear function values  $Y_{f1}$ .
2. The noisy data  $Y_2$  along with the fitted linear function values  $Y_{f2}$ .
3. The linear fitted values (vectors)  $Y_{f1}$  and  $Y_{f2}$ .

### Exercise 9: Calculate coefficient of determination

Implement the body of the function `r_squared()`, provided in a separate .m-file `r_squared.m`, so that it first calculates the residuals as defined by eq. (7) and uses these to perform the calculation of the coefficient of determination  $R^2$  as defined by eq. (6).

Note: The symbol  $Y_f$  is used here in the text to denote the fitted data while in MATLAB source code the symbol `yf` is used (potentially further indexed as: `yf1`, `yf2`). Equation (7) uses the symbol  $\hat{y}$  (y hat).

### Exercise 10: Plot the residuals

For the fits  $Y_{f1}$  and  $Y_{f2}$ , and know values  $Y_1$  and  $Y_2$ , calculate residuals for  $r_1$  and  $r_2$  and plot them using the `stem()` (for  $r_1$ ) and `scatter()` (for  $r_2$ ) functions. Which of these functions appears to be more suitable for this data visualization ?

Also use the `hist()` function to calculate and plot the histogram of one the residuals vectors of your choice.

Note: The vectors of residuals  $r_1$  and  $r_2$  are vectors of the same length as vectors  $Y_1$  and  $Y_2$  (which is the same length as  $Y_{f1}$  and  $Y_{f2}$ ).

### Exercise 11: Change the resolution

Change the value of the  $x$  step in exercise 4 from 0.5 to 0.1 and run the exercises 4 to 9 again. Also try step 0.001. Observe results.

## FINAL projekt časť B: lineárna regresia

### Pravidlá:

1. Pri implementácii úloh je zakázané používať akékoľvek cykly, alebo funkcie symbolického toolboxu. Využitie vektorových operácií systému MATLAB je povinné.
2. Pri implementácii funkcie `linreg()` je zakázané používať funkciu `polyfit()` a operátor `\`, implementovať ju nutné pomocou vzorcov odvodených na prednáške s možným využitím funkcií `mean()`, `var()`, `cov()`.
3. Cvičenia 5 – 10 musia používať krok nastavený ako premennú (nie literál).

### Príprava

Prečítajte si dokumentáciu knasledujúcim funkciám MATLABu: `rand`, `randn`, `randi`, `mean`, `var`, `cov`, `polyfit`. Prečítajte si nie len krátku verziu dokumentácie poskytovanú pomocou príkazu: `>>help function`, ale aj detailnejšiu dokumentáciu, ktorú zobrazuje príkaz: `>>doc function`. Okrem úloh, kde je explicitne uvedený iný postup, implementujte všetky úlohy do jedného spoločného .m súboru, pričom použite bunkový režim (oddelenie pomocou: `%%`) pre oddelenie implementácie jednotlivých úloh. Prečítajte si podrobnú dokumentáciu k implementácii lineárnej regresie v angličtine: [https://www.mathworks.com/help/matlab/data\\_analysis/linear-regression.html](https://www.mathworks.com/help/matlab/data_analysis/linear-regression.html).

Dokumentácia v slovenčine je dostupná tu: <https://kurzy.kpi.fei.tuke.sk/nm/student/13.html>

Implementujte cvičenia 5 (skopírujte z časti A), 7, 8, 10 priamo do súboru `lin_regression.m`.

### Cvičenie 6: Implementujte algoritmus výpočtu koeficientov lineárnej regresie

Otvorte si priložený .m-súbor `linreg.m` obsahujúci funkčný prototyp daný ako:

`function [alpha, beta] = linreg( x, y )`. Implementujte telo tejto funkcie tak aby táto vypočítavala koeficienty jednoduchej lineárnej regresie  $f(x) = \alpha + \beta x$ . Potrebné vzorce sú uvedené dole v sekcii "Rovnice". Tieto vzorce ale nemusíte sami implementovať, vo Vašej implementácii použite hotové funkcie MATLABu: `mean()`, `var()`, `cov()`. **Vo Vašej implementácii funkcie `linreg()` je zakázané používať funkciu `polyfit()` alebo operátor `\`.**

Všimnite si kontrast pri označení: pri lineárnej interpolácii sme si označili koeficienty lineárnej funkcie  $f(x) = ax + b$ , zatiaľ čo pri regresii používame označenie  $f(x) = \alpha + \beta x$ , pričom rôzne funkcie MATLABu používajú ešte aj iné označenia (napr.  $f(x) = p_1x + p_2$ ). Je našou úlohou si to nepoplietť :)

### Cvičenie 7: Vykonajte lineárnu regresiu – nájdite koeficienty

Pre hodnoty  $x$  a obidve zašumené pozorovania  $Y_1$ ,  $Y_2$  z Cvičenia 5 nájdite koeficienty jednoduchej lineárnej regresie  $\alpha$ ,  $\beta$ :

- a) pomocou vlastnej funkcie `linreg()`
- b) pomocou funkcie MATLABu `polyfit()`
- c) pomocou operátora `\`

Pozorujte výsledky.

Poznámka: Pre každý z horeuvedených bodov a) b) c) bude výsledkom dvojica koeficientov  $(\alpha, \beta)$  lineárnej funkcie. Spolu teda 6 hodnôt. Výsledky dosiahnuté v bode a) sa môžu nepatrne líšiť od výsledkov z bodu b) a c). Toto je známa vlastnosť výpočtov a reprezentácie čísel s pohyblivou desatinnou čiarkou. Pre porovnanie týchto výsledkov nie je operátor `==` vhodným nástrojom.

### Cvičenie 8: Vykonajte lineárnu regresiu - vykreslite lineárne priebehy

Použite funkciu `polyval()` pre nájdanie hodnôt  $Y_{f1}$  a  $Y_{f2}$ , získaných dosadením do lineárnych funkcií s koeficientami  $\alpha$ ,  $\beta$  nájdenej pomocou jednoduchej lineárnej regresie pre všetky hodnoty vektora  $x$ , ako je definovaný v Cvičení č. 5.

Použite funkcie `plot()`, `subplot()` a `scatter()` na vykreslenie troch grafov v jednom okne vedľa seba:

1. Zašumené hodnoty  $Y_1$  zároveň s preloženými hodnotami regresiou získanej lineárnej funkcie  $Y_{f1}$ .
2. Zašumené hodnoty  $Y_2$  zároveň s preloženými hodnotami regresiou získanej lineárnej funkcie  $Y_{f2}$ .
3. Hodnoty regresiou získaných lineárnych funkcií (vektorov)  $Y_{f1}$  a  $Y_{f2}$ .

### Cvičenie 9: Vypočítajte koeficient determinácie

Implementujte telo funkcie `r_squared()`, dodanej v samostatnom m.-súbore `r_squared.m`, tak aby táto vypočítala najprv reziduá a tieto potom použila na výpočet koeficientu determinácie  $R^2$  podľa vzorca (6).

Poznámka: Symbol  $Y_f$  je v texte používaný pre linearizované dáta, zatiaľ čo v zdrojovom MATLAB súbore je použitý symbol `yf` (prípadne ďalej indexovaný ako: `yf1`, `yf2`). V rovnici (7) je pre túto istú veličinu použitý symbol  $\hat{y}$  (y so strieškou).

### Cvičenie 10: Vykreslite reziduá

Pre regresiou získané hodnoty  $Y_{f1}$  a  $Y_{f2}$  a známe hodnoty  $Y_1$  a  $Y_2$  vypočítajte reziduá  $r_1$  a  $r_2$  a tieto vykreslite pomocou funkcií `stem()` (pre  $r_1$ ) a `scatter()` (pre  $r_2$ ). Ktorá z týchto funkcií sa Vám javí ako vhodnejšia pre vizualizáciu týchto dát ?

Použite funkciu `hist()` pre výpočet a vykreslenie histogramu Vami zvoleného vektora rezidiu.

Poznámka: Vektory rezidiu  $r_1$  a  $r_2$  sú vektory rovnakej dĺžky ako vektory  $Y_1$  a  $Y_2$  (ktoré sú rovnakej dĺžky ako  $Y_{f1}$  a  $Y_{f2}$ ).

### Cvičenie 11: Zmeňte rozlíšenie

Zmeňte hodnotu kroku pri vektore  $x$  v cvičení č. 4 z 0.5 na 0.1 a spustite kód cvičení 4 až 9 znova. Skúste aj krok 0.001. Pozorujte výsledky.

## Equations / Rovnice

### Basic statistics / Základné štatistiky

Sample mean, výberový priemer

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (1)$$

Sample variance, výberová variancia

$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (2)$$

Sample covariance, výberová kovariancia

$$s_{x,y} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \quad (3)$$

### Calculating simple linear regression coefficients

#### Výpočet koeficientov jednoduchéj lineárnej regresie

$$\beta = \frac{s_{x,y}}{s_x^2} \quad (4)$$

$$\alpha = \bar{y} - \beta \bar{x} \quad (5)$$

### Calculating residuals and coefficient of determination

#### Výpočet reziduí a koeficientu determinácie

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (6)$$

$$r_i = y_i - \hat{y}_i \quad (7)$$