

i Note: This assignment makes use of the **CodeGrade** tool that allows you to receive instantaneous feedback on your code while you work on this assignment. You will also use this tool to hand in your final solution. In the lecture on *Tuesday, 21st November*, a short demonstration will be provided on how you can use CodeGrade to receive feedback and submit your solution.

Instructions

- Download the file `period2_assignment.zip` from Canvas and extract it. In this zip-file, you can find three files: `assignment.py`, `Q2_data.csv` and `Q3_data.csv`. While you work on the assignment, keep these files in the same folder and don't rename any of them. When you submit on CodeGrade, you should only submit `assignment.py`, don't upload the datasets.
- The questions in this assignment will ask you to complete the code of different functions provided in the file `assignment.py`. **You must work in this file. Do not rename the file.** You can edit this file using an editor of your choice (e.g. Spyder, VS Code etc.).
- To answer the questions in this assignment, add your code in between the answer blocks in each function (see below). **Do not modify any of the lines outside these blocks**, otherwise you risk your code not running and you might not receive points for the questions for which your code crashes.

```
## BEGIN ANSWER
# ...
# TODO: Add your code here
# ...
## END ANSWER
```

- To answer the questions in this assignment, you are only allowed to use basic Python functions, as well as the libraries imported in `assignment.py`. As a hint, some useful functions from these libraries have already been imported (see the first lines of `assignment.py`). In addition, the exercises below might specify some restrictions on which functions/libraries you are allowed to use for each question. If you use libraries or functions not permitted for the exercise, you will not receive any points for the questions for which you used them.
- **CodeGrade instantaneous feedback:** You can receive feedback on whether your code is working correctly by copy and pasting your code into the editor in CodeGrade and clicking the **Run** button. This will run your code through a series of functionality checks which check whether your code computes the right answers for each question. You will see for which questions your code passes the checks and which questions you still need to work on. *You can hand-in your code as many times as you want*, so make sure that your code passes as many of the checks as possible in the final version that you submit.
- **Questions/doubts:** If you have any doubts about the questions or you believe that your code is working correctly, but it is not passing the functionality checks, please send a Canvas message to the TA (*Boaz van Bokhorst*), attaching your current version of `assignment.py`. Make sure you clearly state your group name and *specific (!) question*. General questions on whether your code is correct, how to solve the questions etc. will be ignored.

1 Random Number Generation

❗ Note: The questions below will ask you to generate pseudo-random numbers from different distributions. In your code, you should use the `np.random.rand()` or `np.random.uniform()` functions to generate random numbers from the uniform distribution (i.e. don't code a uniform RNG yourself). However, you may not use any other random number generators (e.g. for different distributions, from other libraries etc.).

The *Cauchy distribution* has support $x \in \mathbb{R}$ and has parameters location $\mu \in \mathbb{R}$ and scale $\sigma > 0$. It has the probability density function (pdf)

$$f_X(x) = \frac{1}{\pi\sigma \left(1 + \left(\frac{x-\mu}{\sigma}\right)^2\right)} \quad (1)$$

The *Laplace distribution* has support $x \in \mathbb{R}$ and has parameters location $\mu \in \mathbb{R}$ and scale $b > 0$. It has the probability density function (pdf)

$$f_X(x) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right) \quad (2)$$

- (a) Complete the code for the function `pdf_cauchy(x, mu, sigma)` which accepts a scalar input x and returns the scalar density value $f_X(x)$.
- (b) Complete the code for the function `pdf_laplace(x, mu, b)` which accepts a scalar input x and returns the scalar density value $f_X(x)$.
- (c) Complete the code for the function `rng_cauchy(n, mu, sigma, seed)`. Your code should return a vector of n pseudo-random numbers from the $\text{Cauchy}(\mu, \sigma)$ distribution. To generate these numbers, you should use the *Inverse CDF method*.
- (d) Complete the code for the function `rng_std_laplace(n, seed)`. Your code should return a vector of n pseudo-random numbers from the $\text{Laplace}(0, 1)$ distribution. To generate these numbers, you should use the *Rejection-Sampling method* using a $\text{Cauchy}(0, 1)$ proposal density. To draw Cauchy random numbers, you should use your function from (c).

Note: It should be clear from your code what value of M you used, i.e. define `M=...`

- (e) Complete the code for the function `hist_std_laplace()` such that it draws a histogram of 50,000 draws from the $\text{Laplace}(0, 1)$ distribution using your function from (d). To generate the plot, use `plt.hist()`. The histogram should show the density of the Laplace random draws using 200 bins. The x limits of the plot should be fixed to $[-7, 7]$. The histogram should be overlaid with a plot of the theoretical Laplace density in the appropriate range using your pdf function from (b).

2 Numerical Optimization

❗ Note: Some of the questions below will ask you to optimize objective functions using numerical optimization methods. In question (e), you should code up the optimization procedure yourself following the steps laid out in the exercise. Only in question (g) are you allowed to use the `minimize()` function from the `scipy` library.

This question concerns the simple linear regression model you discussed in Econometrics I.

$$y_i = \beta_0 + \beta_1 x_i + u_i, \quad i = 1, 2, \dots, n \quad (3)$$

Let \mathbf{y} be the n -dimensional vector with entries y_i and \mathbf{x} be the n -dimensional vector with entries x_i . In addition, let \mathbf{X} be the $(n \times 2)$ matrix whose first column contains 1 in every row and whose second column is \mathbf{x} .

The **Ordinary Least Squares (OLS)** estimator in this model $\hat{\beta}_{OLS}$, which minimizes the sum of squared residuals (SSR) of the model, is the 2-dimensional vector given by

$$\hat{\beta}_{OLS} := \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \left(\mathbf{X}^\top \mathbf{y} \right) \quad (4)$$

In contrast, the **Least Absolute Deviation (LAD)** estimator in this model $\hat{\beta}_{LAD}$ is defined as the solution to the following optimization problem

$$\hat{\beta}_{LAD} := \arg \min_{\mathbf{b} \in \mathbb{R}^2} SAR(\mathbf{b}, \mathbf{y}, \mathbf{x})/n \quad (5)$$

$$SAR(\mathbf{b}, \mathbf{y}, \mathbf{x}) = \sum_{i=1}^n |y_i - b_0 - b_1 x_i| \quad (6)$$

For the questions below, you will work with the dataset `Q2_data.csv`. Make sure this dataset is placed in the same folder as the file `assignment.py` that you are editing. For testing your functions while you develop them, you can load the data using the function supplied in the template by running `y, x = load_data_q2()`.

- (a) Using the formula in eq. (4), complete the code for the function `ols_estimator(y, x)` which should compute and return $\hat{\beta}_{OLS}$.
- (b) Complete the function `ols_scatterplot()` such that it plots a scatter plot of the data and overlays it with the fitted regression line given by $\hat{\beta}_{OLS}$ (black solid line).
- (c) Complete the function `sar(b,y,x)` such that it computes the sum of absolute residuals given in equation (6).
- (d) Complete the function `sar_grad(b,y,x)` such that it computes the gradient of the sum of absolute residuals with respect to \mathbf{b} i.e. $\nabla_{\mathbf{b}} SAR(\mathbf{b}, \mathbf{y}, \mathbf{x})$.

- (e) Consider the implementation of a gradient-based optimization algorithm called **Gradient-descent** that is described below. This algorithm is similar to the Newton-Raphson procedure with line search we discussed in the lecture, but does not make use of Hessian information. It also uses a different backtracking procedure than the one we discussed in the lecture.

GRADIENT-DESCENT

- Start with an initial guess $\mathbf{b}^{(0)}$.
- For $k = 0, 1, \dots$, loop until the convergence criterion $|f(\mathbf{b}^{(k-1)}) - f(\mathbf{b}^{(k)})| < f_{tol}$ holds or the maximum number of iterations is exceeded.
 - Compute the gradient $G(\mathbf{b}^{(k)})$ and updating direction $\mathbf{h}^{(k)} = -G(\mathbf{b}^{(k)})$.
 - Update $\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} + \alpha_k \cdot \mathbf{h}^{(k)}$ where α_k is determined using the following *backtracking procedure* for $c_1 = 1e - 4$ and $\rho = 0.95$.
 - (a) Start with $\alpha_k = 1.0$
 - (b) Check if the Armijo rule $f(\mathbf{b}^{(k)} + \alpha_k \mathbf{h}^{(k)}) \leq f(\mathbf{b}^{(k)}) + c_1 \alpha_k \cdot \mathbf{h}^{(k)T} G(\mathbf{b}^{(k)})$ holds.
 - (c) If yes, use this value of α_k . If not, update $\alpha_k = \rho \cdot \alpha_k$ and go to (b).
- Return $\mathbf{b}^{(k+1)}$.

Complete the function `gradient_descent(f, grad, b0, y, x, max_iter, f_tol)` so that it implements this algorithm.

- (f) Complete the function `lad_scatterplot()` such that it plots a scatter plot of the data and overlays it with the fitted regression line given by $\hat{\beta}_{OLS}$ (black solid line) and $\hat{\beta}_{LAD}$ (red dashed line). To compute $\hat{\beta}_{LAD}$, you should use the `gradient_descent()` function from (e) with a starting value of $\mathbf{b}_0 = (0.0, 0.0)^T$.
- (g) Compute $\hat{\beta}_{LAD}$ again for the dataset by completing the code for the function `lad_nelder_mead()`, but this time use the `minimize()` function from `scipy.optimize` to minimize the sum of absolute residuals using the **Nelder-Mead algorithm** with a starting value of $\mathbf{b}_0 = (0, 0)^T$.

3 Solving Linear Equations & Eigenvalues

❗ Note: Some of the questions below will ask you to solve linear equation systems or to compute eigenvalues and eigenvectors using numerical methods. For these exercises, you should use the functions `np.linalg.solve()` and `np.linalg.eig()`, instead of coding up the algorithms that were discussed in the lecture yourself.

This exercise concerns the **Ridge estimator** in the multiple linear regression model (without constant) given by

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + u_i, \quad i = 1, 2, \dots, n \quad (7)$$

Let \mathbf{y} be the n -dimensional vector with entries y_1, y_2, \dots, y_n and \mathbf{x}_j be the n -dimensional vector with entries $x_{1j}, x_{2j}, \dots, x_{nj}$. Further, let \mathbf{X} be the $(n \times k)$ matrix with columns corresponding to \mathbf{x}_j for $j = 1, \dots, k$. As you discussed in Econometrics I, the Ridge estimator is the solution to the following optimization problem for some choice of $\alpha \geq 0$.

$$\hat{\beta}_{ridge} := \arg \min_{\mathbf{b} \in \mathbb{R}^k} Q(\mathbf{b}, \mathbf{y}, \mathbf{X}) \quad (8)$$

$$\begin{aligned} Q(\mathbf{b}, \mathbf{y}, \mathbf{X}) &:= (\mathbf{y} - \mathbf{X}\mathbf{b})^\top (\mathbf{y} - \mathbf{X}\mathbf{b}) + \alpha \cdot (\mathbf{b}^\top \mathbf{b}) \\ &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{b}^\top \mathbf{X}^\top \mathbf{y} + \alpha \cdot \mathbf{b}^\top \mathbf{b} \end{aligned} \quad (9)$$

- Analytically derive an expression for $\nabla_{\mathbf{b}} Q(\mathbf{b}, \mathbf{y}, \mathbf{X}) = 0$. Then, bring this equation into the form $\mathbf{M}\mathbf{b} = \mathbf{z}$. Complete the code for the functions `ridge_M(X,alpha)` and `ridge_z(y,X)` so that they implement the expressions for \mathbf{M} and \mathbf{z} .
- Complete the code for the function `ridge_estimator(y,X,alpha)` so that it computes $\hat{\beta}_{ridge}$ by solving the equation from (a) using `np.linalg.solve()`.
- Let \mathbf{w} be the k -dimensional vector that collects the eigenvalues of $\mathbf{X}^\top \mathbf{X}$. Further, let \mathbf{V} be the $(k \times k)$ matrix whose columns are the eigenvectors associated with eigenvalues w_1, w_2, \dots, w_k . Complete the code for the function `ridge_ev_decomp_XtX(X)` so that it returns \mathbf{w} and \mathbf{V} .
- Note that using the equation in (a), we can also compute the ridge estimator as $\mathbf{b} = \mathbf{M}^{-1}\mathbf{z}$. In this question, we will directly compute \mathbf{M}^{-1} via an eigenvalue decomposition of the matrix $\mathbf{X}^\top \mathbf{X}$. Complete the code for the function `ridge_Minv(X,alpha)` so that it computes \mathbf{M}^{-1} via the following steps.
 - Use your function from (c) to compute the vector of eigenvalues \mathbf{w} and the matrix of eigenvectors \mathbf{V} of the matrix $\mathbf{X}^\top \mathbf{X}$.
 - Compute the modified vector of eigenvalues $\tilde{\mathbf{w}}$ where $\tilde{w}_j = 1/(w_j + \alpha)$ for $j = 1, 2, \dots, k$.
 - Compute $\mathbf{M}^{-1} = \mathbf{V}\tilde{\mathbf{W}}\mathbf{V}^\top$ where $\tilde{\mathbf{W}}$ is a diagonal matrix with $\tilde{\mathbf{w}}$ on the diagonal.
- Complete the code for the function `ridge_estimator_via_inv(y,X,alpha)` so that it computes the ridge estimator as $\mathbf{M}^{-1}\mathbf{z}$. You may not compute \mathbf{M}^{-1} using numpy, but you should use your function from (d). For computing \mathbf{z} , use your function from (a).