# Applied Statistics (ECS764P) - Lab 4

## Fredrik Dahlqvist

## 6 Dec 2023

# 1 Theory

1. Consider the measure $\mathbb{P}$ on $\{1, 2, 3\} \times \{1, 2, 3\}$ defined by

$$
\begin{array}{ccc}
\mathbb{P}(1,1) = {}^1\!/_{10} & \mathbb{P}(1,2) = {}^2\!/_{10} & \mathbb{P}(1,3) = {}^1\!/_{10} \\
\mathbb{P}(2,1) = {}^1\!/_{10} & \mathbb{P}(2,2) = {}^1\!/_{10} & \mathbb{P}(2,3) = {}^2\!/_{10} \\
\mathbb{P}(3,1) = {}^1\!/_{10} & \mathbb{P}(3,2) = 0 & \mathbb{P}(3,3) = {}^1\!/_{10}
\end{array}
$$

(where I've written $\mathbb{P}(x, y)$ for $\mathbb{P}(\{(x, y)\})$ in order to keep things readable.)

(a) Is $\mathbb{P}$ a probability measure?

(b) Prove that $\mathbb{P}$ cannot be written as a product measure. *Hint: prove it by contradiction.*

(c) Compute the two marginals of $\mathbb{P}$.

(d) Compute the covariance and correlation of $\mathbb{P}$.

# 2 Practice

*You can assume that* `numpy`, `matplotlib`, `scipy` *and* `pandas-datareader` *are installed on the machine of the person who will run and mark your notebook. There is no need to force an install with the* `!` *command. For textual answers please use a markdown cell.*

1. **(5 marks)** You will first download the world GDP data from the World Bank using `pandas_datareader`. The following code will download and plot the entire world GDP time series. Do NOT make any local copies of your data!

```
1    from pandas_datareader import wb
2    import matplotlib.pyplot as plt
3    import numpy as np
4
5    gdp_data = wb.downloa\P(indicator='NY.GDP.MKTP.CD', country='WLD',
     start='1960', end='2021')
6    time = np.arange(1960,2022)
7    gdp = gdp_data.iloc[:,0].astype(float).to_numpy()
8    # Data is returned in inverse chronological order, so reverse order
9    gdp = np.flip(gdp)
10   # Plot world GDP data against time
11   plt.plot(time,gdp,label='US GDP')
12   plt.legen\P()
13   plt.show()
14
```

(you can ignore the warning about the code 'WLD'). You will try to estimate the long-term annual growth rate of the world using a regression.

(a) If the growth rate was a constant $r$, then the world's GDP would grow as

$$
GDP_k = GDP_0(1 + r)^k
$$

where $k$ is the number of years since 1960 and $GDP_0$ is the world's GDP in 1960. This is clearly not a linear relationship between time ($k$, in years) and $GDP$. However, we can get a linear relationship by applying a simple transformation $f(-)$ on both side of the equation. What is this transformation? (*Hint: we used this transformation in the context of MLE, it turns products into sums.*)

(b) Apply this transformation $f(-)$ to the GDP data, and perform a regression against the time variable. On the same plot, display your regression line, a scatter-plot of the (transformed) data points, and your $R^2$ value.

(c) Compute the residuals of your regression (i.e. the difference between the model and the observations), and print their mean and their standard deviation $\hat{\sigma}$. Perform a KS-test to determine whether we can reject the null hypothesis that the residuals are sampled from a normal distribution with mean 0 and standard deviation $\hat{\sigma}$. Take $\alpha = 99\%$.

(d) You will now apply the inverse of the transformation $f(-)$ to your linear model in order to get a non-linear model for the GDP. On the same plot, display your (non-linear) model and a scatter-plot of the (original) data points.

(e) What is the relationship between the slope of the regression and the long-term growth rate of the world GDP? Compute the long-term growth rate of the world GDP.

(f) What do you observe since approximately 2015?

2. **(5 marks)** In this question you will study the distribution of the slope and intercept parameters of a linear model. Consider the following model

$$y_i = ax_i + b + \varepsilon_i \qquad \text{where} \qquad a = \frac{1}{2}, b = 2, \varepsilon_i \sim \text{Normal}\left(0, \frac{1}{5}\right), 1 \le i \le N \qquad (1)$$

For the purpose of this exercise you will take $N = 200$ and generate the $x_i$s by

$$x = \text{np.linspace}(-5, 5, 200)$$

(a) Generate 10000 sets of error vectors $\varepsilon_i$ and use them to perform 10000 linear regression of the $N$-dimensional vectors $(y_i)$ against $(x_i)$, where $y_i$ is given by (1).

(b) Collect the slopes and the intercepts of these 10000 linear regressions and plot their histograms against their respective theoretical densities given in the lecture. What do you observe?

(c) For each of the 10000 regression, compute the test statistic for the slope and for the intercept (given in the lecture) and plot their histograms against their theoretical density (also given in the lecture). What do you observe?

(d) Take the last of your regressions and perform the following two tests with $\alpha = 99\%$ (you may use either $p$-values or critical regions but make sure you think about whether this is a one-sided or two-sided test).

$$H_0 : a = \frac{1}{2} \qquad \text{(assuming } b = 2\text{)}$$

$$H_1 : b = 2 \qquad \text{(assuming } a = \frac{1}{2}\text{)}$$

(e) Change the model to

$$y_i = ax_i + b + \varepsilon_i \qquad \text{where} \qquad a = \frac{1}{2}, b = 2, \varepsilon_i \sim \text{Cauchy}\left(0, \frac{1}{5}\right), 1 \le i \le N \qquad (2)$$

Perform another 10000 regressions based on this model. Collect the slopes and intercepts of these regressions as well as the associated statistics. Plot their histograms. What do you observe?