# PYU44C01 Linear Algebra Assignment 2 2021

You are provided with four data sets for a set of 29 EU and EFTA countries from
https://ec.europa.eu/eurostat which are: hours worked per week (hours), employment rates (%),
Labour mobility (000's of workers working outside their country of citizenship by citizenship) and
working population (000's of workers).

(1) Prove that $A_1 = \sigma_1 u_1 v_1^T$ is a rank 1 approximation to $A$ where $u_1$ and $v_1$ are columns of the $Q_1$
and $Q_2$ matrices in the SVD of $A$.

(2) Write a Python script which reads the four data sets provided into a 4 x 29 array $A$ and plots them
*in the order given in the hours worked data set*. Note that order varies by data set. You will find the
arrays in *script.py* useful for this.

(3) Modify the data so that each row of $A$ has zero mean.

(4) Find a diagonal matrix $D$ which scales the rows of $A$ so that each row has unit magnitude.

(5) Form the correlation matrix $C = (D \cdot A) \cdot (D \cdot A)^T$

(6) Perform a principal component analysis of the data sets using the correlation matrix, i.e. project
the data points in $D \cdot A$ onto the two most important principal axes of the correlation matrix.

(7) Put on your economist hat and discuss the meaning of the principal component analysis. You may
find the *Nature* article on Blackboard on genetic variation throughout Europe as well as material
covered in the lecture on principal component analysis useful for this.