

**Project: Population Regression**

**Instructions:** Please work in your preassigned groups to complete and submit your work to the appropriate folder in LumiNUS. Project submissions are due on

Please submit all the following documents as a single zip file named Group-X-Project.zip:

- (i) Powerpoint slides named as Group-X-Project.pptx (20 slides max)
- (ii) Completed Word file named as Group-X-Project.docx (with all results)
- (iii) Print preview of ipynb file named as Group-X-Project.pdf (with all results)
- (iv) Your working ipynb file named as Group-X-Project.ipynb
- (v) Your data files (either csv or excel).

**1. Introduction and Reading Assignment**

In this project, we will look at the human population statistics collected by the various national governments and build a machine learning model to make population predictions.

Please read the following article from Nature Education.

An Introduction to Population Growth

By: Sunny B. Snider (College of Agriculture, California State University, Chico) & Jacob N. Brimlow (College of Agriculture, California State University, Chico)

<https://www.nature.com/scitable/knowledge/library/an-introduction-to-population-growth-84225544/>

**2. Data Source: Singapore Department of Statistics (SingStat)**

Let's start with looking at the population statistics of Singapore. Download Singapore population data from 1950 to 2019 from: <https://www.singstat.gov.sg/>

- a. Graph the total population vs year.
- b. Use linear regression to build an estimator of the total population of Singapore in the future. Use the data for years 2013 and earlier as training data.
- c. Performance metrics:
  - i. What are the slope and y-intercept of the best fit line? Plot the best fit line over the empirical data.
  - ii. What is the  $R^2$  coefficient for the best fit line? See Appendix for definition of the  $R^2$  coefficient.
  - iii. What is the mean squared error (MSE) of the estimator on the training data?
  - iv. Use years greater than 2013 as test data and predict the population for those years.
  - v. What is the MSE of the estimator on the test data? Hint: you may want to normalize the mean squared error for it to be meaningful.
- d. What is your estimate of Singapore's population in 2030 and 2050? Do you think these estimates are reasonable? Explain your answer.
- e. What pattern do you expect for human population growth in Singapore?
- f. How could you improve your estimates of the future population?

**3. Data Source: The World Bank (<https://www.worldbank.org/>)**

Download population data from: <https://www.worldbank.org/>. Note that you will have to search for the data from the World Bank website (sp.pop.totl).

- You should be able to get an excel file with the population of every country from 1960 to 2019. First, verify that the data from the World Bank matches the Singapore population data you previously downloaded from SingStat.gov.sg.
- Use the total population data for China. Graph the total population vs year.
- Use linear regression to build an estimator of the total population of China in the future. Use the data for years 2013 and earlier as training data.
- Performance metrics:
  - What are the slope and y-intercept of the best fit line? Plot the best fit line over the empirical data.
  - What is the  $R^2$  coefficient for the best fit line?
  - What is the mean squared error (MSE) of the estimator on the training data?
  - Use years greater than 2013 as test data and predict the population for those years.
  - What is the MSE of the estimator on the test data? Hint: you may want to normalize the mean squared error for it to be meaningful.
- What is your estimate of China's population in 2030 and 2050? Do you think these estimates are reasonable? Explain your answer.
- What pattern do you expect for human population growth in China?
- How could you improve your estimates of the future population?

**4. Appendix: On the  $R^2$  Coefficient**

The coefficient of determination, or  $R^2$ , is a measure that provides some information about the goodness of fit of a model. In the context of regression, it is a statistical measure of how well the regression line approximates the actual data. It is therefore important when a statistical model is used either to predict future outcomes or in the testing of hypotheses. The most widely used expression for  $R^2$  is shown below.

$$SS_{\text{res}} = \sum_i (y_i - f_i)^2$$

$$SS_{\text{tot}} = \sum_i (y_i - \bar{y})^2$$

$$R^2 \equiv 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

The better the linear regression fits the data in comparison to the simple average, the closer the value of  $R^2$  is to 1.

See the Wikipedia entry on  $R^2$ : [https://en.wikipedia.org/wiki/Coefficient\\_of\\_determination](https://en.wikipedia.org/wiki/Coefficient_of_determination)

See video on  $R^2$  from Khan Academy: <https://youtu.be/Ing4ZgConCM>