

An investigation into use of machine learning techniques
as a tool in the study of the factors that affect life
expectancy on a global scale

Adriaan Louw (53031377)

December 27, 2019

Abstract

The pursuit of increasing ones life expectancy is a very human endeavour. In order to study human life expectancy, demographers and other scientists, study the factors that affect life expectancy. Those include socio-economic indicators like income, educational attainment and per capita spending on healthcare. Understanding these relationships can lead a country to better a understand of which socio-economic areas to focus on, in order to increase their life expectancy. A popular techinique when studying the relationship between life expectancy and these factors is linear regression. This study will use datasets, from the World Bank, that descibe these indicators on a global scale and compare how well linear regression fits the data in comparrison with various machine learning techniques. These techniques are k-nearest neighbour, Support Vector Machines and clustering. All of these techniques will then be compared using Mean Sqared Error.

Contents

1	Introduction	5
2	Literature Review	5
2.1	What do we mean by life expectancy?	5
2.2	Multi indicator studies	6
2.3	Determinants of life expectancy	6
2.3.1	Income	6
2.3.2	Education attainment	6
2.3.3	Spending on health	8
2.3.4	Unemployment	8
2.4	Modelling Techniques	9
2.4.1	Linear Regression	9
2.4.2	Logistic regression	9
2.4.3	k-Nearest Neighbour	9
2.4.4	Support Vector Machines	10
2.4.5	Clustering	10
2.4.6	Neural Networks	11
2.4.7	Decision Trees	11
2.5	Datasets	11
2.5.1	Life Expectancy	11
2.5.2	Income	12
2.5.3	Educational Attainment	12
2.5.4	Per Capita spending on health	12
2.5.5	Unemployment	12
3	Methodology/Procedure	13
3.1	Choice of algorithms	13
3.2	Cross-Validation	13
4	Results	14
4.1	Linear Regression Results	15
4.2	k-Nearest Neighbour Results	15
4.3	Clustering Results	18
5	Discussion	20
6	Conclusions and limitations	20
7	Recommendations	21
	References	21

List of Figures

1	The original Preston curve from Preston (1975)	7
2	Matrix scatter plot of data	14
3	Mean Error in Predicted Life Expectancy per Fold for the kNN Algorithm . . .	16
4	kNN Error Data Plot with standard deviation over all folds	18

List of Tables

1	OLS Regression Results	15
2	Unweighted kNN Results for best value of k	17
3	Weighted kNN Results for best value of k	17
4	Clustering algorithm results per randomly selected number of centroids	19
5	Clustering algorithm results with 1296 evenly spaced centroids	19

1 Introduction

Human beings have always had a fascination with longevity. Myths like the fountain of youth or the Holy Grail are a testament to this fact. Today, longevity and causes of mortality are studied by professionals like Demographers and Actuaries. Trying to determine why some people or group live long lives.

This study investigates the use of Machine Learning techniques in studying determinants of life expectancy for countries. Indicators that have shown to have some form of correlation with life expectancy will be selected. Their relationship with life expectancy will be investigated using various techniques from Machine Learning and contrasted to various forms of regression whose use is ubiquitous in the literature. This analysis will seek out to prove the appropriateness of using these machine learning algorithms for use in research to find the exact correlation between these indicators and life expectancy. It is the hypothesis of this study that machine learning techniques like k-Nearest Neighbour and Support Vector Machines will model the indicator/life expectancy relationship better than regression techniques can. Also that these techniques can create more accurate models. In the hope that the causes of long life expectancies in certain countries can be better understood (Chen & Asch 2017). This study does not aim to prove causation between the indicators chosen and life expectancy, but rather the usefulness of machine learning algorithms as tools.

Statistical regression techniques are predominantly used algorithm for data analysis. Linear regression (Section 2.4.1) for instance assumes a linear relationship between the independent variables and the dependent variable. Which might not be the case. This we will discuss in Literature Review (Section 2).

2 Literature Review

Life expectancy and mortality are 2 related terms. Mortality is generally expressed as a mortality rate. It describes the rate at which people die under certain circumstances. Life expectancy (Section 2.1) is the amount of years an individual or group of people are expected to live. If the amount of people who are dying increases, the mortality rate increases and the life expectancy for people in that group decreases and visa versa.

2.1 What do we mean by life expectancy?

A life table is a table given for a specific year that contains the probability that a person of a certain age will die in that specific year. Actuaries and Demographers use life tables in the insurance industry and the study of demographics respectively. Life expectancy is one element of a life table. Both countries and the United Nations create life tables for use in policy creation.

There are 2 types of life tables, namely period and cohort life tables. A cohort is a group of people who were born in the same year. A cohort life table will follow a cohort over its lifetime until every member of the cohort has died. A cohort life table requires the mortality information of the cohort over many years. This information is often unavailable. While for a period life table, a hypothetical cohort is created and subjected to current mortality rates. This gives the user of the period timetable a window to see mortality rates at that point in time. This makes period life tables the most common type (Arias et al. 2017).

A life expectancy figure from a period life table is called a period life expectancy.

2.2 Multi indicator studies

In this section, studies that used various indicators will be discussed.

Kabir (2008) investigated how well the life expectancy of 93 developing countries were predicted by indicators like income, education and fertility (among others). It classed a countries life expectancy into 3 categories. Then used a probit model where the input variables have a linear relationship. Multiple Ordinary Least Squares Regression was then applied to study indicators' influences.

The study Hu et al. (2015), also used a linear regression model of GDP per capita, Gini indeces, ect with respect to life expectancy. The intention of the study was to link income inequality to mortality rates and life expectancy.

While Shaw et al. (2005) investigated factors like smoking, pharماسutical spending, amount of butter consumed and amount of fruit and vegetables consumed. Then putting them in a linear model and applying regression.

In an examination of 108 countries, Hassan et al. (2016) investigated indicators like education, GDP and health spending as it relates to life expectancy. They used Grossmans model (Grossman 2000) to model their data and Vector Error Correction Model to anlyse the data.

By using the componets of the Human Development index and Pearsons r with multivariate regression, Bulled & Sosis (2010) investigated the link between life expectancy, education and reproduction.

This study aims to add machiine learning algorithms as a tool to life expectancy studies, like these, that look at multiple factors that affect life expectancy.

2.3 Determinants of life expectancy

This section summarises the relationship some indicators have with life expectancy.

2.3.1 Income

The relationship between life expectancy and income has been given a lot of attention in academic circles (Preston 1975, Hu et al. 2015, Chetty et al. 2016, Oeppen 2019).

Preston (1975) was the first to show the correlation between life expectancy and per capita income. His original curve can be seen in Figure 1. As we can see from Figure 1, for low income countries, life expectancy increases rapidly with per capita income. Whereas in high income countries a small increase in per capita income does not have a large effect on life expectancy.

This relationship has also been shown in more recent studies (Chetty et al. 2016, Oeppen 2019). Even though Shkolnikov et al. (2019) found that in Russia, the Preston curve is not an accurate predictor of life expectancy, they found that the actual life expectancy should be "substantially higher" when comparing to the Preston curve predicted value.

Studies in first world countries involving mortality rather that life expectancy have also found a relationship with income level (Blakely et al. 2004, Kalwij et al. 2013, von Gaudecker & Scholz 2007).

Just 16% of the improvement in life expectancy between the 1930s and 1960s could be explained by rising income levels Preston (1975). Which seems to indicate that a countrie's life expectancy is dependant on more than income levels.

2.3.2 Education attainment

Kaplan et al. (2015) investigated the relationship between educational attainment and life expectancy in eight states in the United States. They found that even when controlling for

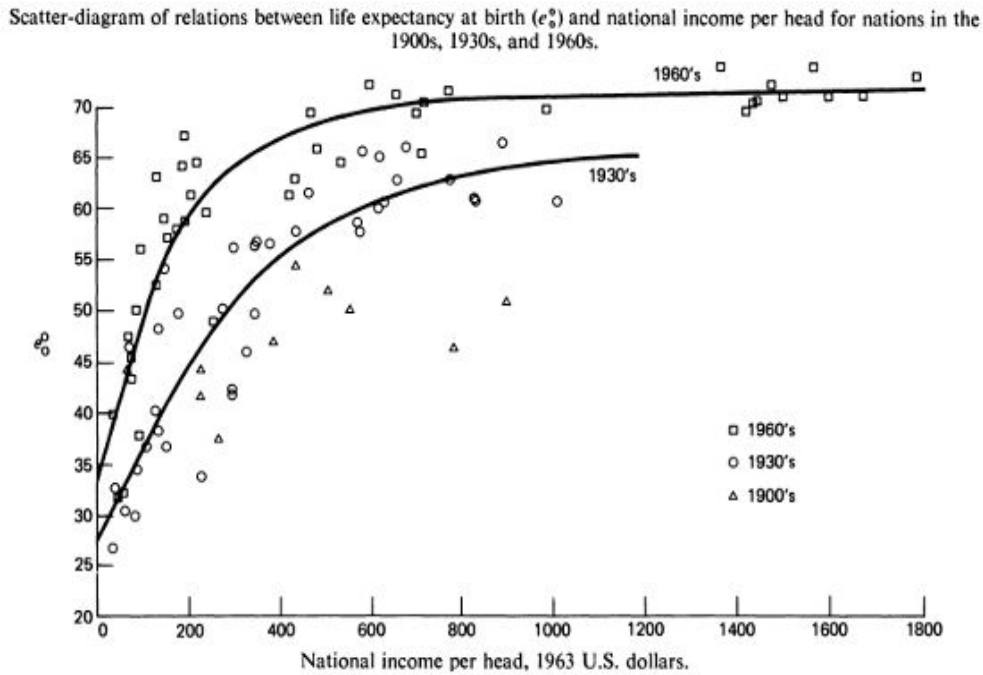


Figure 1: The original Preston curve from Preston (1975)

variables like income, race, sex, and common medical issues like cardiovascular disease, the relationship between educational attainment and life expectancy remains statistically significant.

Luy et al. (2019) studied 3 developed nations, namely the United States, Italy and Denmark. They have also found a strong correlation between education levels and longevity.

But what is the nature of this correlation? According to Deary & Gottfredson (2004), Intelligence Quotient or IQ could explain the association. While Hayward et al. (2015) does not believe in a “causal relationship”, but rather that it depends on factors like “time, place, and the social environment”.

In an attempt to find a causal relationship between education and life expectancy, van Kippersluis et al. (2011) investigated the result of the Netherlands increasing the mandatory number of years a child had to attend school for to 7 years. It was 6 years previously. van Kippersluis et al. (2011) found a decrease in mortality of 3% for 81-year old males who had the additional year of schooling.

This relationship appears strongest in more developed countries where the life expectancy is already above 60 years (Bulled & Sosis 2010). In these countries, any educational investment leads to greater compensation for the learner than they would get in a less developed country Bulled & Sosis (2010), Handwerker (1986). In addition, Kabir (2008) also studied this relationship, among others, with regards to developing countries and did not find a correlation.

The literature appears not to be in agreement.

The question remains, which educational indicators should be used when investigating the relationship between education and life expectancy?

Various educational indicators have been used in the literature for comparing to life expectancy. One approach is to use the International Classification of Education (ISCED) system (UNESCO Institute for Statistics 2012). The ISCED 2011 standard consists of 9 levels ranging from ISCED level 0 (Early childhood education) to ISCED level 8 (Doctoral or equivalent level).

Luy et al. (2019) used the United Nations ISCED-97 (consisting of 7 levels) scale to break education attainment down into 3 levels namely Low (None to Lower Secondary), Medium (Upper secondary) and High (Tertiary education). In van Kippersluis et al. (2011) the Dutch

SOI system (Standaard Onderwijs Indeling). Which, according to van Kippersluis et al. (2011), is similar to the ISCED system. While in Deboosere et al. (2009), educational attainment was broken into 5 levels, also ranging from no education to Tertiary education.

Kaplan et al. (2015) broke educational attainment into 4 levels ranging from less than high school to college graduate.

In the study Bulled & Sosis (2010), the relationship between educational investments and fertility against life expectancy, over 193 countries, was investigated. They used adult literacy and the enrolment ratios for primary, secondary and tertiary schooling.

2.3.3 Spending on health

Healthcare spending and life expectancy in the United States, between 1960 and 2000, was compared in Cutler et al. (2006). They found that increased spending on health per capita, controlling for inflation, is positively correlated to US life expectancy for the time period in question.

Most Eastern European countries, who have joined the European Union, have seen an increase in healthcare spending. This has generally been accompanied by an improvement in life expectancy (Jakovljevic et al. 2016). This has to be seen in the context of the so called “Russian Mortality Crisis” where former Soviet Union countries faced a sudden drop in life expectancy after the fall of the Berlin wall (Brainerd & Cutler 2005). Jakovljevic et al. (2016) found that the best metric to use when comparing health spending of countries, is to use their total per capita health spending (in US dollars).

The same relationship was found in Canada. When spending on healthcare is decreased, life expectancy follows (Crémieux et al. 1999).

It is well known that life expectancy in Sub-Saharan Africa is low. Here spending on health care can also be correlated to increases in life expectancy. Even though poor governance can undo some of the effects of increased spending (Makuta & O’Hare 2015).

A country’s per capita healthcare is not necessarily in proportion to its per capita income. In 2005, the United States spent 50% more on healthcare per capita than its income per capita would suggest (Anderson & Frogner 2008).

2.3.4 Unemployment

According to Bonamore et al. (2015), the literature has 2 main views on the relationship between unemployment and life expectancy. The first view states that during an economic downturn, people suffer from more stress and depression. This leads to more unhealthy lifestyle choices like smoking and alcohol. Which in turn lowers life expectancy. Bonamore et al. (2015) cites the works of Lundin et al. (2014), Montgomery et al. (2013), Garcy & Vågerö (2012), Browning & Heinesen (2012), Dávalos et al. (2012), Backhans & Hemmingsson (2011), Deb et al. (2011) and Strully (2009), who take this view. The second view focusses on times when there is economic growth, i.e. less unemployment. This period of economic growth also can lead to stress eg. burning out and having less time for activities that benefit one’s health. Like going to the gym. This view is held by Tapia Granados & Ionides (2011), Ruhm (2005), Tapia Granados (2005), Neumayer (2004) and Ruhm (2000). Then there are also studies express the view that no connection can be established (Bonamore et al. 2015). The view of Bonamore et al. (2015) is that this relationship is non-linear.

2.4 Modelling Techniques

2.4.1 Linear Regression

Linear Regression is a popular technique, used to find relationships in data. As the name suggests Linear Regression assumes a linear relationship between the input variables and the result (Murphy 2012). This might not be the case for the target function. The target function could be any potential function. In the case of life expectancy modelling, we know that according to the Preston curve (Section 2.3.1), the relationship between income and life expectancy is not linear. Thus using Linear Regression should return a sub-optimal result unless the data is transformed using some non-linear function into a linear space. Linear regression was used in the following studies: Chetty et al. (2016), Jakovljevic et al. (2016), Hu et al. (2015), Mackenbach & Looman (2013), Bulled & Sosis (2010).

The Ordinary Least Square (OLS) method of estimation will be used. This algorithm will be implemented in Python and used to analyse how well this technique models the problem.

2.4.2 Logistic regression

This regression is appropriate when the dependant variable is discrete, e.g. a yes/no answer. It also assumes a linear relationship between inputs. In contrast to linear regression, this linear sum is passed through the sigmoid function (Murphy 2012). All of this makes it inappropriate for non-linear target functions and this study.

2.4.3 k-Nearest Neighbour

The k-Nearest Neighbour algorithm (kNN) is an instance based form of machine learning. It uses the classification of those datapoints closest to the data point to be classified to determine its classification. The kNN-algorithm allows for non-linear problem spaces to be classified, because it does not make an assumption on the nature of the problem space. It just sees a data point as a function of its closes neighbours. This is usefull for indicators like income that are highly non-linear as described in Section 2.3.1. Additionally, how the algorithm determined its output value is transparent and can be used to study how various components affect the end result. This is important , because if the algorithm was not transparent it can not be seen as a replacement for other algorithms that are transparent like regression alorithms. In this study, the standard kNN-algorithm will be altered to accomodate a real valued output and not just a class classification. This will be accomplished by taking the mean life expectancy for all the data points determined to be closest to the target point. Care will have to be taken to reduce the number of features of the data, because this algorithm is sensitive to the so-called “curse of dimentionality” (Mitchell 1997). As discussed in Section 2.5.3, educational attainment will comprise of many indicators. This could lead to educational attainment being weighted more than other indicators. Thus each educational indicator will be given a fractional weight. If there are 4 educational attainment then each indicator will have a 4th of the weight of other indicators.

The value of k will have to be found experimentally and will be affected by the density of the data points. One potential issue could be if the data tends to cluster around certain points. This could happen if developing countries cluster together and developed countries cluster together. That could make developing countries, that are on the higher end of the development scale, more difficult to classify, since not a lot of data points will be very close to them.

2.4.4 Support Vector Machines

The classic Support Vector Machine (SVM) is used in classification tasks. It involves determining the decision surface with regards to the data points closest to the surface. This study will use a modified SVM algorithm that makes it useful for regression tasks. This is called Support Vector Regression (SVR). The SVR algorithm is described in Smola & Schölkopf (2004). The final model can be described as:

$$f(x) = \sum_{i=1}^l (\alpha_i - \alpha_i^*) k(x_i, x) + b \quad (1)$$

where α_i and α_i^* are Lagrange multipliers, b is the bias and the term $k(x_i, x)$ is the kernel function. The kernel is a function to determine the similarity between 2 input vectors. Kernel functions are used to map a non-linear problem space into a smaller subspace where the features are linearly separable. A kernel has to map from the dimension of the problem space down to 1 dimension. Instead of computing the dot product of 2 datapoints, we can replace it with a kernel that requires less computation. Various kernels will be investigated, including the Radial-basis function

$$k(x_i, x) = \exp\left(-\frac{\|x - x_i\|^2}{2\sigma^2}\right) \quad (2)$$

and the polynomial kernel

$$k(x_i, x) = (x_i \cdot x + 1)^p \quad (3)$$

Unfortunately, SVR models are difficult to analyse. Making it hard to determine how the input features of the model affect the output. In order for SVR to be seen as a viable alternative to multivariate regression, an appropriate technique has to be found that can assist in this dilemma. For purpose, this study will use the techniques described in Üstün et al. (2007) to understand which input feature the SVR model deems most important for the modelling of our life expectancy problem space.

In the algorithm described in Üstün et al. (2007), a kernel matrix or Gram matrix has to be computed. This matrix K will be of size $N \times N$ and contains in each cell the kernel function corresponding to the input vectors of the same index such that

$$K_{ij} = k(x_i, x_j) \quad (4)$$

The vector of input data(I) is of size $M \times N$ where M is the number of input features and N is the size of the dataset. The Correlation matrix(R) is calculated by correlating each row of I with each column in K. This gives a resultant vector R of size $N \times M$. Now we have a measure of the importance of each input to the kernel matrix. This resultant matrix is then converted to an image so that the relationships can be seen visually.

2.4.5 Clustering

Clustering techniques are a form of unsupervised learning where the algorithm tries to group the data in a number of groups. In essence, classifying the data into categories. This cannot be easily compared to regression techniques. This study will rather follow the approach of Clustered linear regression(CLR) (Ari & Güvenir 2002). In CLR, clustering techniques are used to break the dataset into clusters. Then apply linear regression to the datapoints in each cluster. In this way one linear regression model does not have to account for all data the data in the data set. Each regression model can be tailored to data similar to it, thereby increasing

the chance of getting a good linear fit on the data. This technique also allows the researcher to gain information about the dataset. If the data is very clustered, it could mean that the problem space is not very linear and that data tends to “jump” between clusters. Comparisons will be made to how the number of clusters affect the results from the regression analyses.

Various clustering algorithms exist. They can be classified into 2 broad groups. The first group is flat clustering. In this group, the algorithms break the data into various groups that have no hierarchy. In other words, there are no group of groups. The second type is hierarchical clustering (Murphy 2012).

Agglomerative clustering and divisive clustering are the 2 types of hierarchical clustering. Agglomerative clustering starts with each data point being a cluster and then combining clusters until some condition is reached. Divisive clustering adds all the datapoints to one cluster and then breaks them into smaller sub-clusters until some condition is met (Murphy 2012).

Each cluster has an initial start point, i.e. the initial position of its centroid. The centroid will be randomly placed in the hypothesis space. Each data point is then assigned to the nearest centroid, determined by Euclidean distance. The centroid's position will then be updated to be the mean of all its newly assigned datapoints. The algorithms then repeatedly update the position of each centroid based on the closest data points, until the centroid stops moving. It is worth noting that the algorithm might not converge to the optimum solution. Therefore multiple runs of this algorithm will be necessary (Murphy 2012).

2.4.6 Neural Networks

Even though Neural networks are capable of representing non-linear hypothesis spaces (Mitchell 1997), they are not appropriate for this study for a couple of reasons. Firstly, the amount of processing power and processing time required, will not be available to this study. Secondly, the results of neural networks are hard to interpret. How the Neural Network came to its conclusion is not clear to the researcher. Which makes it unsuitable as a tool to study the relationship between life expectancy and its various indicators.

2.4.7 Decision Trees

Traceability and understandability are some of the hallmarks of Decision Trees. These algorithms are suited to problem spaces where the target function and the input attributes are discrete values. It is possible to approximate continuous input attributes by making a branch in the tree when a value is smaller or greater than some value, or is between some value. For functions where input attributes span over large ranges, this leads to very large and sub-optimum trees (Mitchell 1997). Many decision tree algorithms exist, like ID3 and C4.5. The problem of determining life expectancy from socio-economic indices has a continuous target function output and continuous input attributes. Therefore, Decision Trees will be excluded from this study.

2.5 Datasets

2.5.1 Life Expectancy

For life expectancy data, this study will use the indicator “Life expectancy at birth, total(years) {SP.DYN.LE00.IN}” from the World Bank’s Development Indicators Database (World Bank Group 2019h). This is a weighted average combining both male and female life expectancy and is calculated in a period life table (see Section 2.1). Only the data that is available for the last 30 years (1981-2010) will be used. This is in keeping with other studies where the amount of

years that their studies look back on is limited to relatively recently (Luy et al. 2019, Hu et al. 2015, Tarkiainen et al. 2012, Kabir 2008, Shaw et al. 2005).

2.5.2 Income

This study will use GDP per capita as was used in Oeppen (2019), Shkolnikov et al. (2019), Mackenbach & Looman (2013) and De Vogli et al. (2005). The source will be GDP per capita, PPP (constant 2011 international \$){NY.GDP.PCAP.PP.KD} from the World Bank (World Bank Group 2019b).

2.5.3 Educational Attainment

This study will use the same indicators and was used in Bulled & Sosis (2010). For an adult literacy indicator, SE.ADT.LITR.ZS will be used from the world bank website (World Bank Group 2019c). This indicator describes the percentage of adults, from the age of 15, who can read and write to a certain level of proficiency. Determining the literacy level of a country is a difficult endeavour. The definition of literacy might be exactly the same between countries, because this indicator uses a lot of data that each country determines on its own. This also might include some measure of competency with numbers.

The indicator SE.PRM.ENRR will serve as the enrollment ratio for primary education (World Bank Group 2019d). As the name suggests, it is the percentage of students in primary school. It is calculated as the fraction of students in primary school over the amount of students that should be in primary school based on population figures. This can lead to a value greater than 100%. This is due to older students or adults that should have completed primary school that are now undergoing primary school education. This can be a sign of a poorly performing school system. This indicator also does not take issues like truancy into account.

For secondary school enrolment, SE.SEC.ENRR (World Bank Group 2019e) will be used. This indicator can be higher than expected or above 100% for the same reasons as the primary school enrollment indicator.

SE.TER.ENRR (World Bank Group 2019f) will be used for tertiary enrollment. A high number here tends to indicate that it is a number from a developed country, because tertiary education is much more of a luxury in developing countries.

It is expected that the enrollment at tertiary and secondary will be highly correlated to industrialised countries and thus with life expectancy. The world bank credits UNESCO Institute of Statistics (2019) for all of these educational indicator data.

2.5.4 Per Capita spending on health

To represent Per capita spending on health in this study's model, SH.XPD.CHEX.PP.CD (World Bank Group 2019a) which describes Current health expenditure per capita, will be used.

2.5.5 Unemployment

Unemployment will be ignored for this study because sufficient data on unemployment per country is not available. This can be seen by looking at indicators like SL.UEM.TOTL.NE.ZS (World Bank Group 2019g) which describe percentage unemployment per country for both sexes.

3 Methodology/Procedure

This study will attempt to create a model that can predict life expectancy for a country based on various socio-economic conditions in the country over a 30-year period. This will be done by implementing the various algorithms discussed in Section 2.4 in Python. Then these programs will be applied to the dataset. Python is a common tool used in machine learning applications, with a wide variety of libraries that can be used. The algorithms will be newly implemented for this project but libraries like Numpy will be used to help with reading data from files and writing data to files.

Then by using these programs we will evaluate how well these machine learning techniques model the relationship between Life expectancy and the chosen indicators. Unlike Shaw et al. (2005), this study will not take into account the age distribution of each country.

The philosophical standpoint of this study is Positivism. By using the scientific method, this study will comprise of an experiment to inductively determine whether machine learning techniques can provide more accurate life expectancy models than those created using regression. This cross-sectional study will use life expectancy indicators shown, from the literature, to have some correlation to life expectancy.

Firstly appropriate data sets will be chosen. Taking into account that some indicators might have sparse information. This will impact the size of the dataset. Looking at available indicators like World Bank Group (2019c), we can see there are approximately 200 countries. If all the indicators have data for each year in the study there will be around 6000 data points. Looking at indicators like World Bank Group (2019c) and World Bank Group (2019a), it is clear that the data is quite sparse, especially for developing countries. The intention is to have the study encompass developed and developing countries, but this will depend on data availability.

Once the dataset is finalised, the algorithms and their permutations described in Section 3.1 will be applied. Three models will be created for each algorithm and permutation thereof. The first model will be using all the data. The second will only use data from developed countries and the third will only use data from developing countries, data permitting. Cross-validation will be applied to the created models as described in Section 3.2. The k data subsets will be determined before starting the validation. Each algorithm will then be run on the same data sets to facilitate comparison.

The mean squared error will be calculated for each of the k Cross-Validation datasets for each of the algorithms. The average of each of these mean squared errors will be calculated for each algorithm. This number will be used to determine the fitness of each algorithm and enables comparison. The mean squared error was chosen as the method of comparison due to its ease of use and that it can be applied to all the algorithms.

3.1 Choice of algorithms

As discussed in Section 2.4 of the Literature Review, this study will compare k -Nearest Neighbour, Support Vector Machines and Clustered linear regression with Linear regression. Linear Regression will be used as a baseline for comparison on the dataset.

3.2 Cross-Validation

By using stratified k -fold cross validation, this study will aim to reduce the impact of the relative small dataset that will be analysed. This form of cross validation will ensure that when the validation set is chosen, no important data points are ignored for training. The data will be broken down randomly into k subsets of equal size. Each data subset will also

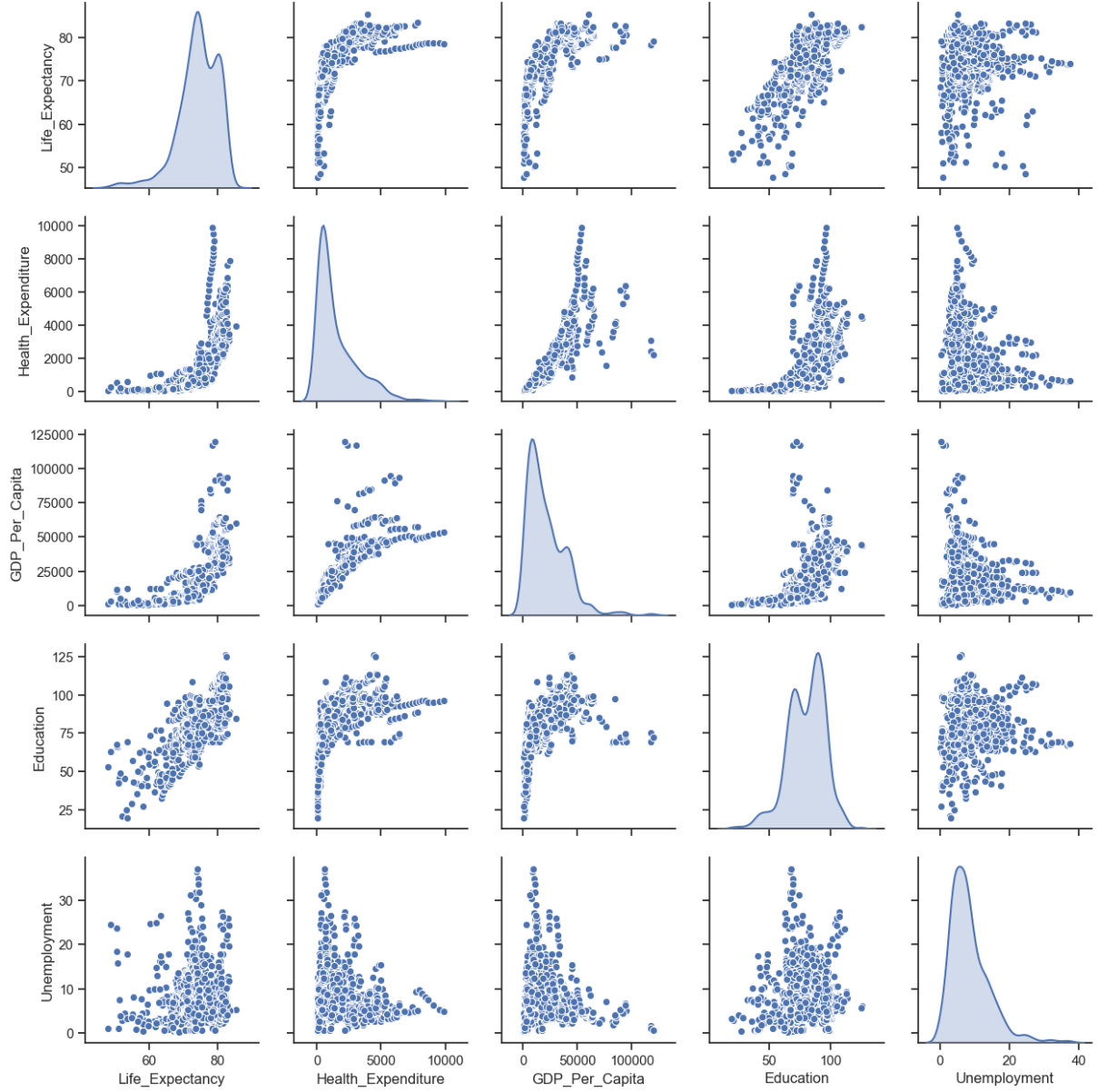


Figure 2: Matrix scatter plot of data

contain equal amounts of data points with low and high life expectancies, so that no dataset is completely towards one end of the data range. One data subset is chosen to be the validation subset and the remaining $k - 1$ subsets are combined into the training set. The model is then trained on the training dataset and its performance is measured against the validation subset. This is done k times in order for each subset to be the validation subset. For each of the training runs the mean of the error will be calculated (Mitchell 1997, Murphy 2012). The value of k will be dependant on the final dataset.

4 Results

The data set for each indicator was downloaded separately. The data was then combined based on country and year. Sparse data became a problem. If at least one indicator did not have data for a specific year and country, then the entry for that that year and country was removed from the dataset. The indicator for Adult Literacy was also removed from all entries due to that indicator having very sparse data. The total number of accepted data points were 1062.

Run	Average Mean Error	Average Standard Deviation of Error
1	9.598848355	7.594930051
2	9.626031548	7.633113769
3	9.609315926	7.621318362
4	9.605118586	7.635851148
5	9.606997134	7.65353075
6	9.6332431	7.622023874
7	9.590347552	7.621950019
8	9.607112393	7.611794799
9	9.598032704	7.608716795
10	9.619144548	7.590407725
Average:	9.609419185	7.619363729

Table 1: OLS Regression Results

For each of those data points, every indicator (except adult literacy) was defined.

The data for each indicatr was normalised between 0 and 1. For each indicator, the largest value was assigned 1 and the smallest 0. The rest was then scaled proportionally between 0 and 1. This was done so that algorithms like kNN do not favour one indicator/dimention over another due to its relative size.

The values for Primary-, Secondary-, and Tertiary Enrollment was added together and divided by 3 to create a new Indicator namely “Education”. This was done in order to reduce the dimentionality of this 7-dimentional data set to a more manageable 5-dimentional data set. The dataset can be visually seen in a matrix scatter plot in Figure 2. It is worth noting that the Preston curve can be seen in this catter plot where Life Expectancy and GDP Per Capita intersect.

The following procedure was followed in order to get 10 stratified folds. The dataset was then sorted according to decending life expectancy. The dataset was then split in half such that the life expectancy of each data point in the “top” data set is higher than the life expectancy of each data point in the “bottom” data set. Then 10 folds of equal size was created from the “top” and “bottom” data sets, by taking 53 data points from each set. Thus each fold has 106 data points, with 53 high life expectancy points and 53 low life expectancy points.

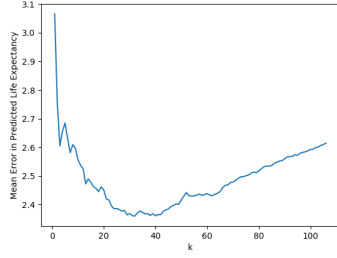
This process was done 10 times to create 10 runs of 10 folds each. Each dataset contains the same data but the data is spread around the different folds randomly.

4.1 Linear Regression Results

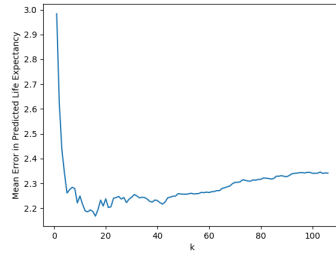
As described in Section 2.4.1 the OLS algorithim was implemented in Python and run for each fold each dataset. The results can be see in Table 1. The table shows the Average Mean Error per dataset. This is done by calculating the mean error of each fold in a dataset and then computing the average over all the folds in the dataset. The same process was followed for the Standard Deviation. The average Error over all the runs and thus all the runs of the algorithm was 9.609 years. While the Standard Deviation was 7.619.

4.2 k-Nearest Neighbour Results

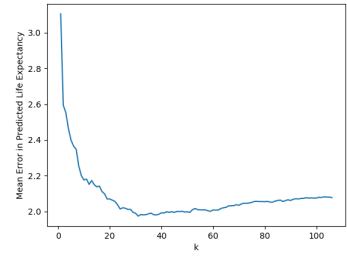
There were 2 versions of the kNN-algorithm. The first was an unweighted implementation where the predicted Life Expectancy was calculated as an average of the life expectancy of the k nearest neighbours. The second version was a weighted nearest neighbour algorithm. The predicted life expectancy(LE_p) was calculated with the following formula:



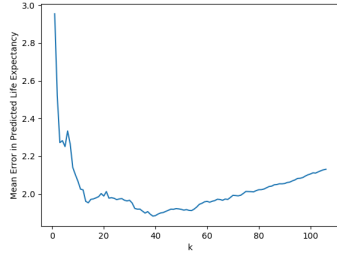
(a) Fold 1



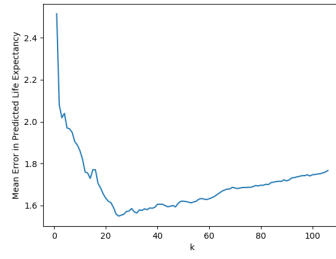
(b) Fold 2



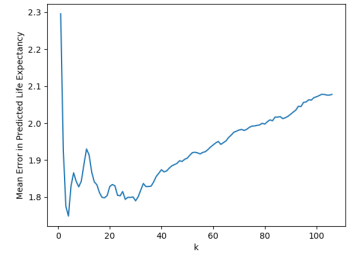
(c) Fold 3



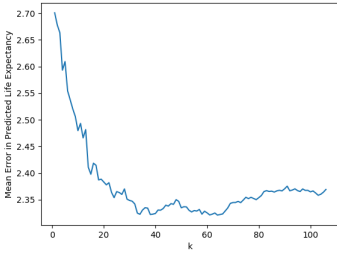
(d) Fold 4



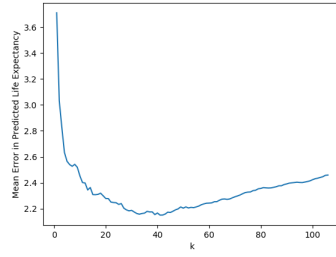
(e) Fold 5



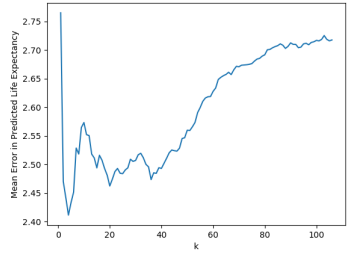
(f) Fold 6



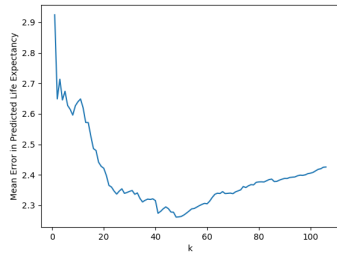
(g) Fold 7



(h) Fold 8



(i) Fold 9



(j) Fold 10

Figure 3: Mean Error in Predicted Life Expectancy per Fold for the kNN Algorithm

Run	k	Average Mean Error	Average Standard Deviation of Error
1	41	2.120386929	2.389743605
2	35	2.118507644	2.330946813
3	39	2.120214137	2.343546292
4	39	2.122089962	2.345210974
5	39	2.122707998	2.339231525
6	39	2.122234215	2.336469553
7	39	2.120867119	2.339882294
8	39	2.11824741	2.340979544
9	39	2.119204169	2.339061524
10	39	2.118872887	2.338973832
Average:	38.8	2.120333247	2.344404596

Table 2: Unweighted kNN Results for best value of k

Run	k	Average Mean Error	Average Standard Deviation of Error
1	41	2.121160816	2.381492154
2	28	2.110539349	2.262304954
3	39	2.11851791	2.314154022
4	41	2.127764586	2.358604755
5	40	2.124401498	2.315950553
6	41	2.119369173	2.327937177
7	38	2.112441118	2.344628653
8	39	2.101194006	2.343523783
9	41	2.125380836	2.322854298
10	37	2.091687809	2.275117546
Average:	38.5	2.11524571	2.324656789

Table 3: Weighted kNN Results for best value of k

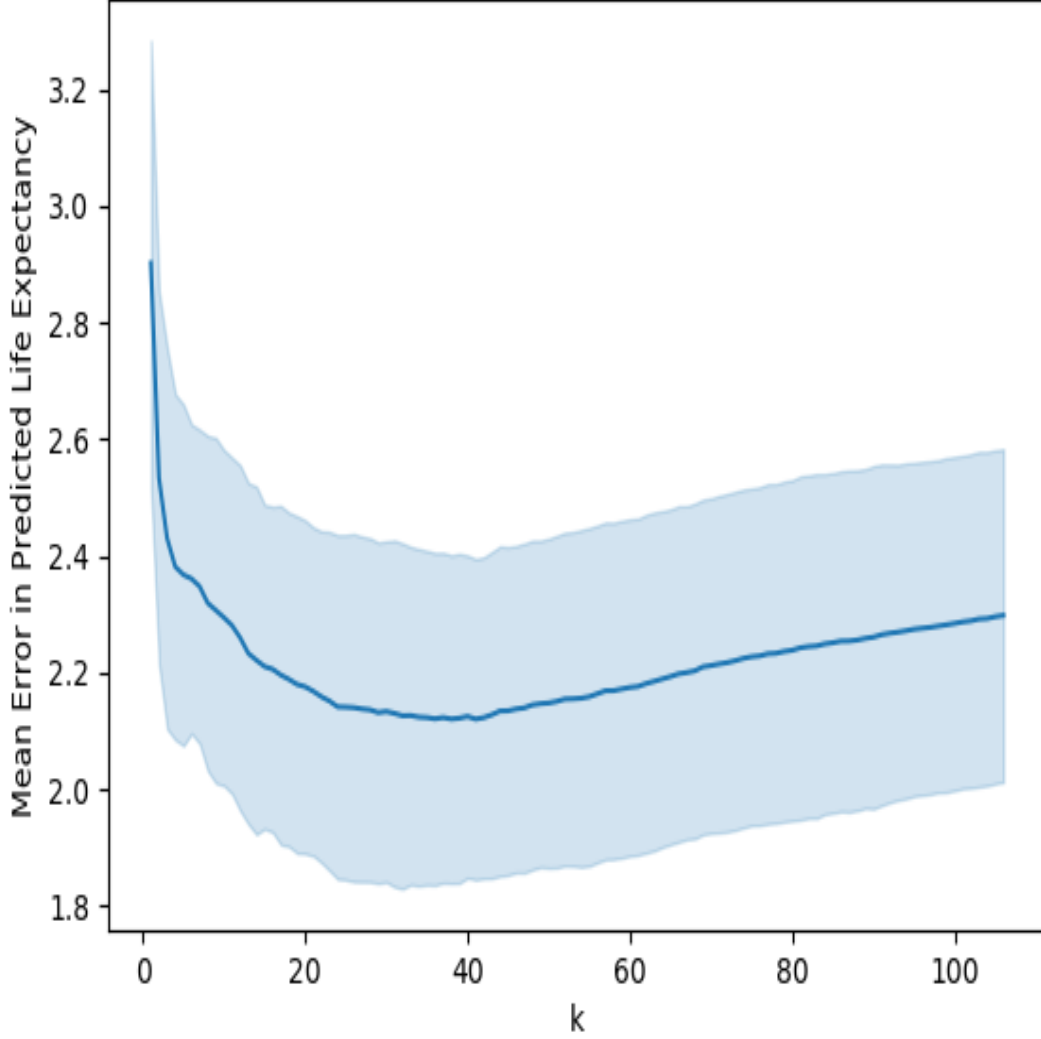


Figure 4: kNN Error Data Plot with standard deviation over all folds

$$LE_p = \frac{\sum_i^k LE_i e^{d_i}}{\sum_i^k e^{d_i}} \quad (5)$$

where LE_i is the life expectancy of the i th data point and d_i is the distance from the i th data point to the point in question. The results of the unweighted kNN-algorithm can be seen in Table 2 and the results for the weighted kNN-algorithm can be found in Table 3.

4.3 Clustering Results

At first, random centroids were selected for each run of the clustering algorithm on each fold. This was also done for various numbers of centroids. The runs ran for up to 15 centroids per fold. The amount of centroids that could be used was limited by the amount of time the algorithm took to run. It took approximately 12 hours to run through 2 to 15 centroids for each fold. 4 of these total runs were done. The algorithm terminated when all the datapoints assigned to each centroid was within a euclidean distance of 0.001. This number was determined

# of centroids	Average Mean Error	Average Standard Deviation of Error
2	7.762052171	5.96294992
3	7.629710319	5.844772366
4	7.005638769	5.599632559
5	6.631308225	5.422260091
6	6.325310938	5.132051508
7	6.167973105	5.048466276
8	6.12395204	4.873890668
9	6.925555773	4.823754093
10	6.022997788	4.71522236
11	6.495814931	4.626752313
12	5.857179844	4.626444289
13	7.591622412	4.492052798
14	5.712254947	4.549610122
15	7.823888766	4.400055159

Table 4: Clustering algorithm results per randomly selected number of centroids

Run	Average Mean Error	Average Standard Deviation of Error
1	9.740346466	3.921416011
2	31.32104966	2.332399296
3	5.883284899	2.711148601
4	8.483479693	3.029508192
5	3.02759275	1.949908716
6	6.221257613	2.338482549
7	6.909572466	2.351402757
8	9.928680799	2.519936579
9	9.152888827	2.537911275
10	8.425976639	2.496948897
Average:	9.909412981	2.618906287

Table 5: Clustering algorithm results with 1296 evenly spaced centroids

experimentally. That number was the smallest distance the datapoints could be from the centroids while the algorithm would still terminate within a reasonable amount of time. The number is dimensionless, since it is calculated using the normalised values of each indicator or dimension. Which in turn is also dimensionless. The averaged results of these random runs can be seen in Table 4.

5 Discussion

Since Linear regression is used often in the literature (Section 2.4.1), the result for linear regression will be regarded as the baseline. From Table 1 we can see that the Error value averaged 9.609 years over all the data sets. The large error is to be expected since the data is not linear, as can be seen in Figure 2. More specifically, the Preston Curve discussed in Section 2.3.1 can be seen in Figure 2 where Life Expectancy and GDP per Capita intersect. Thus Linear Regression when done over the whole dataset, gives a result with such high error.

The kNN-algorithm fared much better. The kNN-algorithm was first implemented without weighting the nearest neighbours. Table 2 shows that result. This table only takes the Average Mean Error into account for the value of k for which the Average Mean Error is the smallest. The value of k does change per run but averages around 39. Over the 10 the best value of the Average Mean Error averaged 2.12 years.

The weighted version of the kNN algorithm was then implemented in order to improve on the numbers found in Table 2. The results for the 10 runs of the Weighted kNN-algorithm is shown in Table 3.

The weighted version of the algorithm did improve on the results of the unweighted algorithm. The weighted algorithm had an average Mean Error of 2.115 years. Which is an improvement of 0.0051 years. The standard deviation of both kNN-Algorithms also improved more than 3 fold over the standard deviation of the errors for the OLS-Regression algorithm.

Unlike OLS Linear regression, the kNN-algorithm only takes datapoints into account that are local to the point in question. Therefore, not caring about the general structure of the whole dataset. One disadvantage of this algorithm was the need to run the data for so many different values of k . The max value of k was 106. This meant that for each fold the kNN-algorithm was run for 1 Nearest Neighbour up to all the Neighbours in the fold. Given that a fold only contained 106 data points.

The first results from the clustering algorithm was where the optimum number of clusters was being determined using random positions of the centroids. As discussed in Section 4.3, the number of centroids were varied between 2 and 15.

6 Conclusions and limitations

Multi-dimensional data is hard to visualise and thus hard to understand and study. The situation gets worse with each additional dimension. Combining Primary-, Secondary- and Tertiary School Enrollment, meant that the amount of dimensions studied reduced from 7 to 5. Which is unfortunately still difficult to analyse, mostly because it is impossible to visualise mentally. Techniques like a matrix scatter plot from Figure 2 helped but still require interpretation.

The source code and raw results can be found on GitHub (Louw 2019).

The greatest limitation for the analysis of the centroid algorithm was the computation time. Running through all the folds of a dataset could take more than a week, not including external factors like load shedding or technical issues on the chosen computer. Python is not an ideal language to scientific computation when all the computation, especially long loops, are not being done in a more performant language like C or C++ (Cai et al. 2005). It would have been

more efficient to implement the clustering algorithm in a more performant language and then calling those libraries from python. If that was the case then the algorithm could have been tested against more clusters. The threshold distance that the datapoints had to be from the centroids, for the algorithm to terminate, could also have been reduced.

7 Recommendations

References

- Anderson, G. F. & Frogner, B. K. (2008), ‘Health spending in OECD countries: Obtaining value per dollar’, *Health Affairs* **27**(6), 1718–1727.
- Ari, B. & Güvenir, H. A. (2002), ‘Clustered linear regression’, *Knowledge-Based Systems* **15**(3), 169–175.
URL: www.elsevier.com/locate/knosys
- Arias, E., Heron, M. & Xu, J. (2017), ‘United States Life Tables, 2013.’, *National vital statistics reports : from the Centers for Disease Control and Prevention, National Center for Health Statistics, National Vital Statistics System* **66**(3), 1–64.
- Backhans, M. C. & Hemmingsson, T. (2011), ‘Unemployment and mental healthwho is (not) affected?’, *The European Journal of Public Health* **22**(3), 429–433.
- Blakely, T., Kawachi, I., Atkinson, J. & Fawcett, J. (2004), ‘Income and mortality: The shape of the association and confounding New Zealand Census - Mortality study, 1981-1999’, *International Journal of Epidemiology* **33**(4), 874–883.
- Bonamore, G., Carmignani, F. & Colombo, E. (2015), ‘Addressing the unemployment-mortality conundrum: Non-linearity is the answer’, *Social Science and Medicine* **126**, 67–72.
- Brainerd, E. & Cutler, D. M. (2005), Autopsy on an Empire: Understanding Mortality in Russia and the Former Soviet Union, Technical Report 1.
- Browning, M. & Heinesen, E. (2012), ‘Effect of job loss due to plant closure on mortality and hospitalization’, *Journal of Health Economics* **31**(4), 599–616.
- Bulled, N. L. & Sosis, R. (2010), ‘Examining the Relationship between Life Expectancy, Reproduction, and Educational Attainment’, *Human Nature* **21**(3), 269–289.
- Cai, X., Langtangen, H. P. & Moe, H. (2005), ‘On the performance of the Python programming language for serial and parallel scientific computations’, *Scientific Programming* **13**(1), 31–56.
- Chen, J. H. & Asch, S. M. (2017), ‘Machine Learning and Prediction in Medicine Beyond the Peak of Inflated Expectations’, *New England Journal of Medicine* **376**(26), 2507–2509.
- Chetty, R., Stepner, M., Abraham, S., Lin, S., Scuderi, B., Turner, N., Bergeron, A. & Cutler, D. (2016), ‘The association between income and life expectancy in the United States, 2001-2014’, *JAMA - Journal of the American Medical Association* **315**(16), 1750–1766.
- Crémieux, P.-Y., Ouellette, P. & Pilon, C. (1999), ‘Health care spending as determinants of health outcomes’, *Health Economics* **8**(7), 627–639.

- Cutler, D. M., Rosen, A. B. & Vijan, S. (2006), 'The Value of Medical Spending in the United States, 1960-2000'.
- Dávalos, M. E., Fang, H. & French, M. T. (2012), 'EASING THE PAIN OF AN ECONOMIC DOWNTURN: MACROECONOMIC CONDITIONS AND EXCESSIVE ALCOHOL CONSUMPTION', *Health Economics* **21**(11), 1318–1335.
URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/hec.1788>
- De Vogli, R., Mistry, R., Gnesolto, R. & Cornia, G. A. (2005), 'Has the relation between income inequality and life expectancy disappeared? Evidence from Italy and top industrialised countries', *Journal of Epidemiology and Community Health* **59**(2), 158–162.
URL: <http://jech.bmj.com/>
- Deary, I. J. & Gottfredson, L. S. (2004), 'Intelligence Predicts Health and Longevity, but Why?', *Current Directions in Psychological Science* **13**(1), 1–4.
- Deb, P., Gallo, W. T., Ayyagari, P., Fletcher, J. M. & Sindelar, J. L. (2011), 'The effect of job loss on overweight and drinking', *Journal of Health Economics* **30**(2), 317–327.
- Deboosere, P., Gadeyne, S. & Van Oyen, H. (2009), 'The 1991-2004 evolution in life expectancy by educational level in Belgium based on linked census and Population register data', *European Journal of Population* **25**(2), 175–196.
- Garcy, A. M. & Vågerö, D. (2012), 'The length of unemployment predicts mortality, differently in men and women, and by cause of death: A six year mortality follow-up of the Swedish 1992-1996 recession', *Social Science & Medicine* **74**(12), 1911–1920.
- Grossman, M. (2000), THE HUMAN CAPITAL MODEL, in 'Handbook of Health Economics, Volume 1', Vol. 1, pp. 348–407.
- Handwerker, W. P. (1986), 'The Modern Demographic Transition: An Analysis of Subsistence Choices and Reproductive Consequences', *American Anthropologist* **88**(2), 400–417.
- Hassan, F. A., Minato, N., Ishida, S. & Mohamed Nor, N. (2016), 'Social Environment Determinants of Life Expectancy in Developing Countries: A Panel Data Analysis', *Global Journal of Health Science* **9**(5), 105–117.
- Hayward, M. D., Hummer, R. A. & Sasson, I. (2015), 'Trends and Group Differences in the Association between Educational Attainment and U.S. Adult Mortality: Implications for Understanding Education's Causal Influence *', *Soc Sci Med* **127**, 8–18.
- Hu, Y., van Lenthe, F. J. & Mackenbach, J. P. (2015), 'Income inequality, life expectancy and cause-specific mortality in 43 European countries, 1987-2008: a fixed effects study', *European Journal of Epidemiology* **30**(8), 615–625.
- Jakovljevic, M. B., Vukovic, M. & Fontanesi, J. (2016), 'Life expectancy and health expenditure evolution in Eastern Europe: DiD and DEA analysis', *Expert Review of Pharmacoeconomics and Outcomes Research* **16**(4), 537–546.
- Kabir, M. (2008), 'Determinants of Life Expectancy in Developing Countries', *The Journal of Developing Areas* **41**(2), 185–204.
- Kalwij, A. S., Alessie, R. J. M., Knoef, M. G., Kalwij, A. S., Alessie, R. J. M. & Knoef, M. G. (2013), 'The Association Between Individual Income and Remaining Life Expectancy at the Age of 65 in the Netherlands', *Demography* **50**(1), 181–206.

- Kaplan, R. M., Howard, V. J., Safford, M. M. & Howard, G. (2015), ‘Educational attainment and longevity: Results from the REGARDS U.S. national cohort study of blacks and whites’, *Annals of Epidemiology* **25**(5), 323–328.
- Louw, A. (2019), ‘Source Code and Raw Data’.
URL: <https://github.com/ajlouw7/Unisa/tree/master/Hons%20Project>
- Lundin, A., Falkstedt, D., Lundberg, I. & Hemmingsson, T. (2014), ‘Unemployment and coronary heart disease among middle-aged men in Sweden: 39 243 men followed for 8years’, *Occupational and Environmental Medicine* **71**(3), 183 LP – 188.
- Luy, M., Zannella, M., Wegner-Siegmundt, C., Minagawa, Y., Lutz, W. & Caselli, G. (2019), ‘The impact of increasing education levels on rising life expectancy: a decomposition analysis for Italy, Denmark, and the USA’, *Genus* **75**(1), 11.
- Mackenbach, J. P. & Looman, C. W. (2013), ‘Life expectancy and national income in Europe, 1900-2008: An update of Preston’s analysis’, *International Journal of Epidemiology* **42**(4), 1100–1110.
- Makuta, I. & O’Hare, B. (2015), ‘Quality of governance, public spending on health and health status in Sub Saharan Africa: a panel data regression analysis’, *BMC Public Health* **15**(1), 932.
- Mitchell, T. M. (1997), *Machine Learning*, McGraw-Hill international editions - computer science series, McGraw-Hill.
- Montgomery, S., Udumyan, R., Magnuson, A., Osika, W., Sundin, P.-O. & Blane, D. (2013), ‘Mortality following unemployment during an economic downturn: Swedish register-based cohort study’, *BMJ Open* **3**(7), e003031.
- Murphy, K. P. (2012), *Machine Learning: A Probabilistic Perspective*.
- Neumayer, E. (2004), ‘Recessions lower (some) mortality rates:: evidence from Germany’, *Social Science & Medicine* **58**(6), 1037–1047.
- Oeppen, J. (2019), Life Expectancy Convergence Among Nations Since 1820: Separating the Effects of Technology and Income, in T. Bengtsson & N. Keilman, eds, ‘Old and New Perspectives on Mortality Forecasting’, Springer International Publishing, Cham, pp. 197–219.
- Preston, S. H. (1975), ‘The Changing Relation between Mortality and level of Economic Development’, *Population Studies* **29**(2), 231–248.
- Ruhm, C. (2000), ‘Are Recessions Good for Your Health?’, *The Quarterly Journal of Economics* **115**(2), 617–650.
- Ruhm, C. J. (2005), ‘Healthy living in hard times’, *Journal of Health Economics* **24**(2), 341–363.
- Shaw, J. W., Horrace, W. C. & Vogel, R. J. (2005), The Determinants of Life Expectancy: An Analysis of the OECD, Technical Report 4.
- Shkolnikov, V. M., Andreev, E. M., Tursun-zade, R. & Leon, D. A. (2019), ‘Patterns in the relationship between life expectancy and gross domestic product in Russia in 200515: a cross-sectional analysis’, *The Lancet Public Health* **4**(4), e181–e188.

- Smola, A. J. & Schölkopf, B. (2004), ‘A tutorial on support vector regression’, *Statistics and Computing* **14**(3), 199–222.
- Strully, K. W. (2009), ‘Job loss and health in the U.S. labor market.’, *Demography* **46**(2), 221–246.
- Tapia Granados, J. A. (2005), ‘Increasing mortality during the expansions of the US economy, 1900-1996.’, *International journal of epidemiology* **34**(6), 1194–1202.
- Tapia Granados, J. A. & Ionides, E. L. (2011), ‘Mortality and Macroeconomic Fluctuations in Contemporary Sweden’, *European Journal of Population / Revue européenne de Démographie* **27**(2), 157–184.
- Tarkiainen, L., Martikainen, P., Laaksonen, M. & Valkonen, T. (2012), ‘Trends in life expectancy by income from 1988 to 2007: decomposition by age and cause of death’, *Journal of Epidemiology and Community Health* **66**(7), 573 LP – 578.
- UNESCO Institute for Statistics (2012), *International Standard Classification of Education ISCED 2011*.
- UNESCO Institute of Statistics (2019), ‘UNESCO Institute of Statistics’.
URL: <http://uis.unesco.org/>
- Üstün, B., Melssen, W. J. & Buydens, L. M. C. (2007), ‘Visualisation and interpretation of Support Vector Regression models’, *Analytica Chimica Acta* **595**(1), 299–309.
- van Kippersluis, H., O’Donnell, O. & van Doorslaer, E. (2011), ‘Long Run Returns to Education: Does Schooling Lead to an Extended Old Age?’, *Journal of Human Resources* **46**(4), 695–721.
- von Gaudecker, H. M. & Scholz, R. D. (2007), ‘Differential mortality by lifetime earnings in Germany’, *Demographic Research* **17**, 83–108.
- World Bank Group (2019a), ‘Current health expenditure per capita, PPP (current international \$)’.
URL: <https://data.worldbank.org/indicator/SH.XPD.CHEX.PP.CD>
- World Bank Group (2019b), ‘GDP per capita, PPP (constant 2011 international \$)’.
URL: <https://data.worldbank.org/indicator/NY.GDP.PCAP.PP.KD>
- World Bank Group (2019c), ‘Literacy rate, adult total (% of people ages 15 and above)’.
URL: <https://data.worldbank.org/indicator/SE.ADT.LITR.ZS>
- World Bank Group (2019d), ‘School enrollment, primary (% gross)’.
URL: <https://data.worldbank.org/indicator/SE.PRM.ENRR>
- World Bank Group (2019e), ‘School enrollment, secondary (% gross)’.
URL: <https://data.worldbank.org/indicator/SE.SEC.ENRR>
- World Bank Group (2019f), ‘School enrollment, tertiary (% gross)’.
URL: <https://data.worldbank.org/indicator/SE.TER.ENRR>
- World Bank Group (2019g), ‘Unemployment, total (% of total labor force) (national estimate)’.
URL: <https://data.worldbank.org/indicator/SL.UEM.TOTL.NE.ZS>
- World Bank Group (2019h), ‘World Development Indicators’.
URL: <http://datatopics.worldbank.org/world-development-indicators/>