# LIFE EXPECTANCY PREDICTION

By Trent Casillas

---

## OUTLINE

**Generalized Capstone Timeline:**

- Motivation
- Dataset Introduction
- Exploratory Data Analysis
- Feature Correlation
  - Top Positive and Negative correlated features to Life Expectancy and correlated pairs overall
- Regression Methods and Evaluation
  - Models include: Ridge Regression, Gradient Boosting , Robust (Thiel-sen) Regression
- Clustering and Evaluation
  - K-means, Mean-Shift
- Mixed Effect Model Evaluation
  - Check ICC
  - Run Random Intercepts, Random Slopes, Random Intercepts and Slopes Models
  - Compare best model to Regression Methods
- Conclusions
  - Takeaways and further studies

---

## MOTIVATION AND PURPOSE

- Predict life expectancy by looking at the positive and negatively correlated factors to improve life quality

- Try to discern any differences between countries groupings

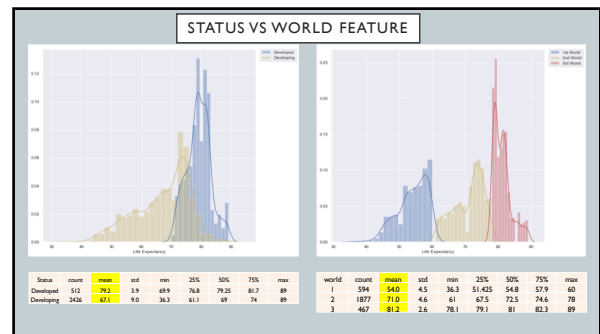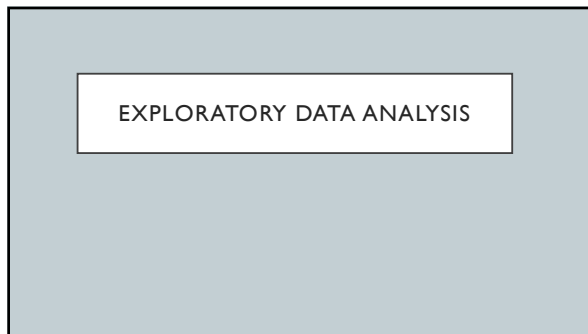- Serves as an example for countries to assess to improve life expectancy for their citizens.

---

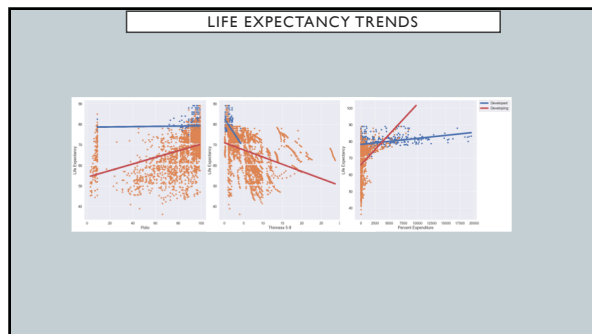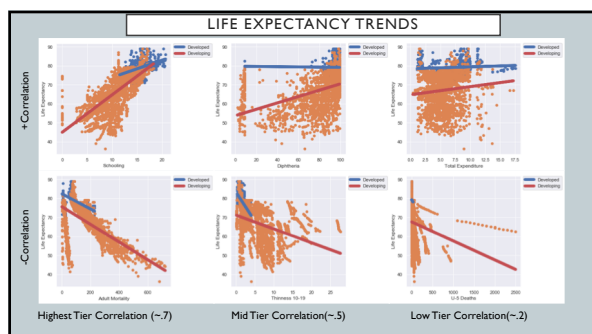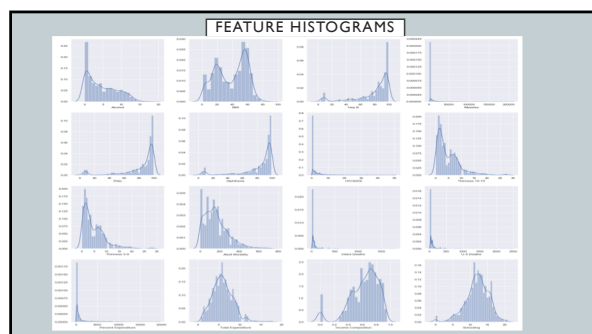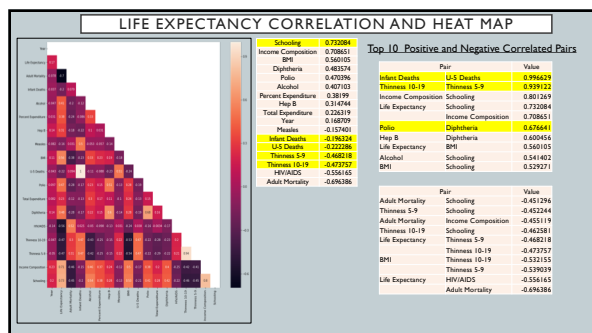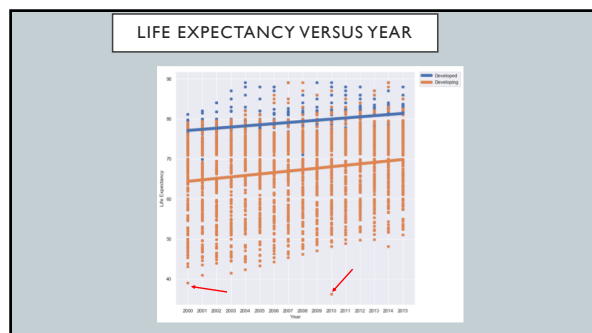## DATA SET

---

## DATASET INFORMATION AND CLEANING

- My data is from The World Health Organization and United Nations website, Life Expectancy expresses the results of 193 countries for life expectancy with 2938 rows for 20 different features columns spanning from 2000-2015 with 2563 missing values.

- The features vary from whether country demographics such as Population, GDP, Total Expenditure spend on health to population statistics such immunizations and mortality rates along with BMI, Alcohol Consumption, and years of schooling.

- The data cleaning included removed 1100 removed between GDP(448) and Population,(652).

- 818 values replaced with missing respective country means.

- 645 replaced by the mean related to the status of the country.

| | Country | Year | Status | Life Expectancy | Adult Mortality | Infant Deaths | Alcohol | Percent Expenditure | Hep B | Measles | BMI | U-5 Deaths | Polio | Total Expenditure | Diphtheria | HIV/AIDS | GDP | Population | Thinness 10-19 | Thinness 5-9 | Income Composition | Schooling |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 2015 | Developing | 65 | 263 | 62 | 0.01 | 71.279624 | 65 | 1154 | 19.1 | 83 | 6 | 8.16 | 65 | 0.1 | 584.2591 | 33736494 | 17.2 | 17.3 | 0.479 | 10.1 |
| 1 | Afghanistan | 2014 | Developing | 59.9 | 271 | 64 | 0.01 | 73.523582 | 62 | 492 | 18.6 | 86 | 58 | 8.18 | 62 | 0.1 | 612.696514 | 327582 | 17.5 | 17.5 | 0.476 | 10 |
| 2 | Afghanistan | 2013 | Developing | 59.9 | 268 | 66 | 0.01 | 73.219243 | 64 | 430 | 18.1 | 89 | 62 | 8.13 | 64 | 0.1 | 631.744976 | 31731688 | 17.7 | 17.7 | 0.47 | 9.9 |
| 3 | Afghanistan | 2012 | Developing | 59.5 | 272 | 69 | 0.01 | 78.184215 | 67 | 2787 | 17.6 | 93 | 67 | 8.52 | 67 | 0.1 | 669.959 | 3696958 | 17.9 | 18 | 0.463 | 9.8 |
| 4 | Afghanistan | 2011 | Developing | 59.2 | 275 | 71 | 0.01 | 7.097109 | 68 | 3013 | 17.2 | 97 | 68 | 7.87 | 68 | 0.1 | 63.537231 | 2978599 | 18.2 | 18.2 | 0.454 | 9.5 |
| 5 | Afghanistan | 2010 | Developing | 58.8 | 279 | 74 | 0.01 | 79.679367 | 66 | 1989 | 16.7 | 102 | 66 | 9.2 | 66 | 0.1 | 553.32894 | 2883167 | 18.4 | 18.4 | 0.448 | 9.2 |
| 6 | Afghanistan | 2009 | Developing | 58.6 | 281 | 77 | 0.01 | 56.762217 | 63 | 2861 | 16.2 | 106 | 63 | 9.42 | 63 | 0.1 | 445.893298 | 284331 | 18.6 | 18.7 | 0.434 | 8.9 |
| 7 | Afghanistan | 2008 | Developing | 58.1 | 287 | 80 | 0.03 | 25.873925 | 64 | 1599 | 15.7 | 110 | 64 | 8.33 | 64 | 0.1 | 373.361116 | 2729941 | 18.8 | 18.9 | 0.433 | 8.7 |
| 8 | Afghanistan | 2007 | Developing | 57.5 | 295 | 82 | 0.02 | 10.910156 | 63 | 1141 | 15.2 | 113 | 63 | 6.73 | 63 | 0.1 | 369.835796 | 26616792 | 19 | 19.1 | 0.415 | 8.4 |
| 9 | Afghanistan | 2006 | Developing | 57.3 | 295 | 84 | 0.03 | 17.171518 | 64 | 1990 | 14.7 | 116 | 58 | 7.43 | 58 | 0.1 | 272.56377 | 298946 | 19.2 | 19.3 | 0.405 | 8.1 |

---

## FEATURE LIST

**Country- Country**
**Year- Year**
**Status- Developed or Developing status**
**Life Expectancy- Age(years)**
**Adult Mortality- Adult Mortality Rates of both sexes(probability of dying between 15&60 years per 1000 population)**
**Infant Deaths- Number of Infant Deaths per 1000 population**
**Alcohol- Alcohol, recorded per capita (15+) consumption (in litres of pure alcohol)**
**Percent Expenditure- Expenditure on health as a percentage of Gross Domestic Product per capita(%)**
**Hep B- Hepatitis B (HepB) immunization coverage among 1-year-olds(%)**
**Measles- number of reported measles cases per 1000 population**
**BMI- Average Body Mass Index of entire population**
**U-5 Deaths- Number of under-five deaths per 1000 population**
**Polio- Polio(Pol3) immunization coverage among 1-year-olds(%)**
**Total Expenditure- General government expenditure on health as a percentage of total government expenditure(%)**
**Diphtheria- Diphtheria tetanus toxoid and pertussis (DTP3) immunization coverage among 1-year-olds(%)**
**HIV/AIDS- Deaths per 1000 live births HIV/AIDS(0-4 years)**
**GDP- Gross Domestic Product per capita(in USD)**
**Population- Population**
**Thinness 10-19- Prevalence of thinness among children and adolescents for Age 10 to 19**
**Thinness 5-9(%)- Prevalence of thinness among children for Age 5 to 9(%)**
**Income Composition-Human Development Index in terms of income composition of resources(0-1)**
**Schooling- Number of years of Schooling**

## DATASET



- Population and GDP will be dropped due to discrepancies with the data.
- Fill in remaining missing with country mean where applicable
- In other cases, fill in data with associated country status mean.

## 3 SIGMA OUTLIERS

- Adult Mortality 99 Percentile 675.2 Max 723.0
- Infant Deaths 99th Percentile 1307.0 Max 1800
- Alcohol 99th Percentile 15.07 Max 17.87
- Percent Expenditure 99th Percentile 15357.4 Max 19479.9
- Percent Expenditure Measles 99th Percentile 111472.7 Max 212183
- BMI 99th Percentile 76.7 Max 87.3
- U-5 Deaths 99th Percentile 1807.0 Max 2500
- Total Expenditure 99th Percentile 16.2 Max 17.6
- HIV/AIDS 99th Percentile 42.1 Max 50.6
- Thinness 10-19 99th Percentile 27.0 Max 27.7
- Thinness 5-9 99th Percentile 27.9 Max 28.6
- Schooling 99th Percentile 20.3 Max 20.7

## EXPLORATORY DATA ANALYSIS

## LIFE EXPECTANCY



- Maximum value 89.0 years
- Mean 69.2 years
- Minimum 36.3 years
- Standard Deviation 9.5 years

- Normal Fit : 0.9566084742546082,
- P-value: 9.622531605232346e-29

## COUNTRY CATEGORIES



Country Status:
Developing: 82.6% Developed: 17.4%

World Feature:
3rd World:16.0% 2nd World:63.8% 1st World: 20.3%

3rd World is below 60 and 1st world is above 78

## STATUS VS WORLD FEATURE



| Status | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| Developed | 512 | 79.2 | 3.9 | 69.9 | 76.8 | 79.25 | 81.7 | 89 |
| Developing | 2426 | 67.1 | 9.0 | 36.3 | 61.1 | 69 | 74 | 89 |

| world | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| 1 | 594 | 54.0 | 4.5 | 36.3 | 51.425 | 54.8 | 57.9 | 60 |
| 2 | 1877 | 71.0 | 4.6 | 61 | 67.5 | 72.5 | 74.6 | 78 |
| 3 | 467 | 81.2 | 2.6 | 78.1 | 79.1 | 81 | 82.3 | 89 |

## LIFE EXPECTANCY VERSUS YEAR



## LIFE EXPECTANCY CORRELATION AND HEAT MAP



## FEATURE HISTOGRAMS



## LIFE EXPECTANCY TRENDS



Highest Tier Correlation (~.7)   Mid Tier Correlation(~.5)   Low Tier Correlation(~.2)

## LIFE EXPECTANCY TRENDS
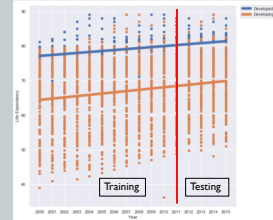


## LIFE EXPECTANCY TRENDS

# REGRESSION

## PREPARATION AND REGRESSION METHODS

- 781 Outliers were removed from data and the data was scaled due to varying feature values.
- Thiel-sen (Robust Regression), Ridge , and Gradient Boosting Regression methods.
- The data was split into 5 different statuses: Developed, Developing, 1st World, 2nd World, 3rd World, a training and test set with all statuses included, and a full data set.
- Training data was from 2000- 2011, Testing after 2011 to 2015.

| count | 2157 |
|---|---|
| mean | 70.663418 |
| sd | 8.447265 |
| min | 41 |
| 25% | 65.9 |
| 50% | 72.9 |
| 75% | 76 |
| max | 89 |



## COMPILED MEAN RESULTS BY METHOD AND STATUS



## THIEL AND RIDGE RESULTS



## GRADIENT BOOSTING RESULTS



## GRADIENT BOOSTING FEATURE IMPORTANCE

| Developed | Developing | 1st World | 2nd World | 3rd world | Xtrain | X | AVG |
|---|---|---|---|---|---|---|---|
| Adult Mortality | HIV/AIDS | Adult Mortality | AdultMortality | AdultMortality | Income Composition | Income Composition | Income Composition |
| Income Composition | Income Composition | Income Composition | Income Composition | HIV/AIDS | HIV/AIDS | HIV/AIDS | Adult Mortality |
| Thinness 5-9 | Adult Mortality | Thinness 10-19 | Schooling | country_code | Adult Mortality | AdultMortality | HIV/AIDS |
| TotalExpenditure | Schooling | Year | HIV/AIDS | Income Composition | Schooling | Schooling | Schooling |
| Alcohol | Polio | Schooling | Thinness 5-9 | Schooling | Thinness 5-9 | U5 Deaths | Thinness 5-9 |

Income Composition ,Adult Mortality, HIV/AIDS, Schooling and Thinness 5-9 most important in LE gradient boosting across all models.

## CLUSTERING

---

### CLUSTERING METHODS AND EVALUATION

- Mean-Shift and K-Means Clustering with scaling of the full data set.
- K-Means cluster chosen by comparing scores calculated below.
- Silhouette and Calinski-Harabaz scores were used as non-ground truth scores.
- Completeness, Homogeneity, and ARI scores were used as ground truth scores related to the world feature.
- Percentages of each cluster were calculated to choose the appropriate clustering methods.

---

### CLUSTERING EVALUATION SCORES

**Calinski-Harabaz Index-**
- Additionally known as the Variance Ratio Criterion where a higher score means a better defined cluster.
- It compares the ratio of the between-clusters dispersion mean and the within-cluster dispersion.
- Scores are higher when dense and separated from other clusters.
- The score is normalized with respect to the others scores for comparison in one chart.

**Silhouette Score-**
- the ratio of difference between the mean nearest-cluster distance and mean intra-cluster distance over the maximum between both scores.
- +1 indicates a highly dense cluster while scores close to 0 indicates overlapping clusters and -1 indicates incorrect clustering.
  s=(b-a)/max(b,a)
  mean intra-cluster distance (a)
  mean nearest-cluster distance (b)

Source: Caliński, T., & Harabasz, J. (1974). "A dendrite method for cluster analysis". Communications in Statistics-theory and Methods 3: 1-27. doi:10.1080/03610926.2011.560741.
Source: Peter J. Rousseeuw (1987). "Silhouettes: a Graphical Aid to the Interpretation and Validation of Cluster Analysis". Computational and Applied Mathematics 20: 53–65. doi:10.1016/0377-0427(87)90125-7.

---

### CLUSTERING EVALUATION SCORES

- Homogeneity Score- Each cluster contains only members of a single class.
- Completeness Score-All members of a given class are assigned to the same cluster.

Adjusted Rand Score-The Rand Index compares how pairs of datapoints relate in the ground truth and in the post-clustering assignment. There are four possible types of pair relationships:
a=Members of the same cluster in the ground truth match same cluster in the new solution.
b=Members of the same cluster in the ground truth match different clusters in the new solution.
c=Members of different clusters in the ground truth match the same cluster in the new solution.
d=Members of different clusters in the ground truth match different clusters in the new solution.
E(RI)- expected RI
RI=(a+c)/sum(a,b,c,d)
ARI=(RI-E(RI)/(max(RI)-E(RI))

---

### COMPARISON OF K-MEANS AND MEAN SHIFT

| Mean-Shift | |
|---|---|
| Number of Estimated Clusters: | 3 |
| Calinski Harabaz Score | 0.09 |
| Silhouette Score | 0.35 |
| Homogenity Score | 0.15 |
| Completeness Score | 0.02 |
| Mean Shift Cluster Percentages | |
| 0 | 97.5 |
| 1 | 1.9 |
| 2 | 0.6 |

| K-Means | |
|---|---|
| Number of Clusters | 3 |
| Calinski Harabaz Score | 0.640 |
| Silhouette Score | 0.215 |
| Homogenity Score | 0.362 |
| Completeness Score | 0.460 |
| K-Means Cluster Percentages | |
| 0 | 47.5 |
| 1 | 29.3 |
| 2 | 23.3 |

---

### K-CLUSTER DETERMINATION

### APPLIED K-MEANS CLUSTERING



### APPLIED K-MEANS CLUSTERING



### APPLIED K-MEANS CLUSTERING



### MIXED EFFECT MODELS

### MIXED EFFECT MODEL ASSUMPTIONS

- Mixed Effect Models have the following assumptions:
- The features are related linearly to the outcome.
- The errors have constant variance: In other words, the model fits equally well for all values of the outcome and features within a level.
- The errors are independent: The fit of the model within a group(level 1) is uncorrelated with the fit of the model at between a group (level 2).
- The errors are normally distributed.
- Observations within a person/group are correlated with one another(high ICC).

### ICC AND GROUP FEATURE SELECTION

```
Model Country
         Mixed Linear Model Regression Results
===================================================
Model:            MixedLM  Dependent Variable: Life_Expectancy
No. Observations: 2157     Method:             REML
No. Groups:       180      Scale:              5.2217
Min. group size:  1        Likelihood:         -5290.2892
Max. group size:  16       Converged:          Yes
Mean group size:  12.0
---------------------------------------------------
              Coef.  Std.Err.    z     P>|z|  [0.025  0.975]
---------------------------------------------------
Intercept    69.776   0.634   110.102  0.000  68.534  71.018
Group Var    71.585   3.488
===================================================

The Intraclass Correlation is: 0.932
```

| Column | ICC | Group Var |
|---|---|---|
| Year | 0.004 | 0.308 |
| Country | 0.932 | 71.58 |
| Status | 0.469 | 49.42 |
| Adult_Mortality | 0.905 | 94.02 |
| Infant_Deaths | 0.408 | 27.72 |
| Alcohol | 0.24 | 17.12 |
| Percent_Expenditure | 0.008 | 0.58 |
| Hep_B | 0.582 | 58.37 |
| Measles | 0.46 | 42.51 |
| BMI | 0.611 | 40.48 |
| U_5_Deaths | 0.431 | 28.92 |
| Polio | 0.503 | 43.37 |
| Total_Expenditure | 0.162 | 11.6 |
| Diphtheria | 0.528 | 47.27 |
| HIV_AIDS | 0.43 | 20.75 |
| Thinness_10_19 | 0.504 | 37.14 |
| Thinness_5_9 | 0.504 | 36.81 |
| Income_Composition | 0.82 | 78 |
| Schooling | 0.78 | 90.07 |
| world | 0.908 | 172.93 |

Sample Calculation
ICC=71.585/71.585+5.22=.932

group var: between-group variance
resid: within-group variance
ICC= group var/(group var +resid)

## MIXED EFFECT MODELS

Random Intercepts
RE=Null

Random Slopes
RE=~0+Feature Name

Random Slopes and Intercepts
RE=~Feature Name

## RI MODEL RESULTS

Running Random Intercepts Model
Mixed Linear Model Regression Results

- 2nd best performing model type!
- RMSE Range from 1.72 to 2.18.

## RS MODEL RESULTS

Running Random Slopes Model
Mixed Linear Model Regression Results

- Worst performing model type!
- RMSE Range from 1.77 to 7.28

## RIS MODEL RESULTS

- Best performing model type!
- RMSE Range from 1.32 to 1.99.

## MODEL RESULTS

Residual histogram

Shapiro Wilk: (0.7576553225517273, 0.0)

The incremental difference between developing and developed is 11.3 years with an average of 71.5 years. Year doesn't vary all that much.

## MODEL RESULTS

Between Variance

Residual

## FINDINGS



| Method | RMSE | R²ICC | LE Min | LE Max | Mean LE | LE Std | LE Var |
|---|---|---|---|---|---|---|---|
| Gradient Boosting | 1.688 | 0.9601 | 42.0 | 85.0 | 70.7 | 8.147 | 66.366 |
| RIS | 1.370 | 0.983 | 41.9 | 85.9 | 70.7 | 8.295 | 68.809 |
| Data | - | 1 | 41.0 | 89.0 | 70.7 | 8.447 | 71.356 |

## TAKEAWAYS

| Developed | Developing | 1st World | 2nd World | 3rd world | Xtrain | X | AVG |
|---|---|---|---|---|---|---|---|
| Adult Mortality | HIV/AIDS | Adult Mortality | Adult Mortality | Adult Mortality | Income Composition | Income Composition | Income Composition |
| Income Composition | Income Composition | Income Composition | Income Composition | HIV/AIDS | HIV/AIDS | HIV/AIDS | Adult Mortality |
| Thinness 5-9 | Adult Mortality | Thinness 10-19 | Schooling | country_code | Adult Mortality | Adult/Mortality | HIV/AIDS |
| Total Expenditure | Schooling | Year | HIV/AIDS | Income Composition | Schooling | Schooling | Schooling |
| Alcohol | Polio | Schooling | Thinness 5-9 | U-5 Deaths | Thinness 5-9 | Thinness 5-9 | Thinness 5-9 |
| | | | | | | U-5 Deaths | |

- Life Expectancy Ranges:
  - Developed :70 to 89
  - Developing: 41 to 86 years
  - 1st World: 78 to 88 years
  - 2nd World: 65 to 78
  - 3rd World: 41 to 65

- The developed and developing status don't fully cover the different categories of countries.

- Testing Data averaged much higher across all groupings than the training groups. Most generally, life span is increasing every year about .3 years of added life.

- Disease and hunger relief are an universal key to improving life expectancy globally.

- Mixed Models are effective , but still need to keep improving.

## FUTURE WORK

- Look at class within a particular country and see if these same factors are same in determining life expectancy for an individual.
- Use the Twitter API to incorporate NLP analysis for a country to see how it relates to Life Expectancy.
- Increase the dataset size with continuing UN and Global Data to incorporate new added features like population, GDP, environmental, and etc in order to test and clarify country groupings.
- Mental Health versus Life Expectancy