

# Estimation of country's population using Lasso Regression

Michael Bass , Samarth Mistry and Samyuktha Sankaran

**Abstract**—This challenge is to estimate the population of the given 259 countries between the years 2000 and 2016 based on the data given for the years 1960 to 1999. The lasso regression was used to carry out the estimate.

## I. INTRODUCTION

LASSO (least absolute shrinkage and selection operator) is a regression analysis method that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretability of the statistical model it produces. Lasso was originally formulated for least squares models and this simple case reveals a substantial amount about the behavior of the estimator, including its relationship to ridge regression and best subset selection and the connections between lasso coefficient estimates and so-called soft thresholding. It also reveals that (like standard linear regression) the coefficient estimates need not be unique if covariates are collinear.

Though originally defined for least squares, lasso regularization is easily extended to a wide variety of statistical models including generalized linear models, generalized estimating equations, proportional hazards models, and M-estimators, in a straightforward fashion. Lasso's ability to perform subset selection relies on the form of the constraint and has a variety of interpretations including in terms of geometry, Bayesian statistics, and convex analysis.

## II. DATA SET

The data set used is from the world bank data, and we are using the population of about 260 countries. The data from 1960 to 1999 is used to train the model, and the model is tested on its accuracy to predict the population from 2000 to 2016. One of the requirements is to predict a test country's population as a linear combination of other countries. However, the linear combination must have a sparsity of five.

## III. REGULARIZATION STRATEGY

Lasso regularization achieves its goal by adding to the least squares minimization an  $\alpha \cdot \text{norm}_1$  of Beta. When  $\alpha$  approaches 0, the solution is equal to that of least squares. As  $\alpha$  increases, the solution becomes sparser. However, for equations, different  $\alpha$  values will yield varying degrees of sparsity. This lead us to realize that we must search the  $\alpha$  space for the  $\alpha$  value that meets our constraints and yields the most optimal solution.

Our approach was to perform an initial fast coarse search followed by an iterative fine-grained search. For each country we created a line space of alphas from 0 to 5000, with 25 points. We performed lasso regression with this set of alphas and evaluated the sparsity of the results. Given our

constraint of using only 5 countries we must find an  $\alpha$  that correlates with 5 non-zero coefficients. If there were no acceptable solutions from the alphas between 0 to 5000, we then exponentially increase our search space as follows:

- 0 to 5000
- 5000 to 10000
- 10000 to 20000
- 20000 to 40,000
- etc.

Once we find sparsity values of 5, we evaluate each  $\alpha$  that results in 5 coefficients, and determine which yields than minimum mean squared error (MSE) and call this our current  $\alpha$ . Next, we proceed to the fine grain search.

For the fine-grained search, we refer to our line space, and select the  $\alpha$  to the left of the current  $\alpha$  as the leftmost bound of the next line space, and the  $\alpha$  to the right of the current  $\alpha$  as the rightmost bound of the next line space. Then, we create a line space between these points with 100 points. Again, we apply the lasso regression, and select the optimal  $\alpha$ . We perform this fine-grained search iteratively and can do this as many times as needed. For our results, we applied the fine-grained search twice.

We had to modify the coarse search slightly as sometimes no results had a sparsity of 5. As the search space increases, and since the coarse search uses only 25 points per line space, the distance between points increases. As a result, it may be possible for the result to include sparsity greater than 5 and sparsity less than 5, but no 5. If this happens, we narrow in on the region that surpassed 5 until we find an  $\alpha$  with sparsity of 5.

### A. Alternate Approaches

We had to modify the coarse search slightly as sometimes no results had a sparsity of 5. As the search space increases, and since the coarse search uses only 25 points per line space, the distance between points increases. As a result, it may be possible for the result to include sparsity greater than 5 and sparsity less than 5, but no 5. If this happens, we narrow in on the region that surpassed 5 until we find an  $\alpha$  with sparsity of 5.

A tuning parameter ( $\lambda$ ) controls the strength of the penalty term. When  $\lambda = 0$ , ridge regression equals least squares regression. If  $\lambda = \infty$ , all coefficients are shrunk to zero. The ideal penalty is therefore somewhere in between 0 and  $\infty$ .

## IV. RESULTS

### A. Co-efficients

For each country, the five best countries were selected, i.e., the other coefficients, except these five were reduced to zero.

A sample of the first five countries are presented in Table I.

### B. Prediction

The test data set is fed to the trained model, to predict the population of the given countries between years 2000 and 2016. A sample of the first five countries between 2000 and 2004 are shown in Table II.

## V. VISUALIZATION

Data visualization is viewed by many disciplines as a modern equivalent of visual communication. It involves the creation and study of the visual representation of data.

To communicate information clearly and efficiently, data visualization uses statistical graphics, plots, information graphics and other tools. Numerical data may be encoded using dots, lines, or bars, to visually communicate a quantitative message. Effective visualization helps users analyze and reason about data and evidence. It makes complex data more accessible, understandable and usable. Users may have particular analytical tasks, such as making comparisons or understanding causality, and the design principle of the graphic (i.e., showing comparisons or showing causality) follows the task. Tables are generally used where users will look up a specific measurement, while charts of various types are used to show patterns or relationships in the data for one or more variables.

Here, each node represents a country and each country is connected to its five best features (i.e., countries). Figure 2 and 3 shows the distribution for countries in the African and European continents respectively. Figure 1 shows the distribution for the entire data set distributed geographically.

We added the latitudes and longitudes for every country given in the data set to see if we could get any inference. We plotted the graph using the GeoLayout in Gephi but could not say anything definitively as there are connections all over the place. To see how it follows for all countries of a particular continent, we added an attribute continent to every node to filter the nodes in Gephi. Where,

- Avg. Degree is the average no. of arcs incident to a node in the graph.
- Network Diameter is the average graph distance between all pairs of nodes. Connected nodes have a graph distance of 1. The diameter is the longest graph distance between any two nodes in the network. (i.e. How far apart are the two most distant nodes).
- Network Density measures how close the network is to 'complete'. A complete network has all possible connections and its density is 1.
- Avg. Clustering Coefficient gives an overall indication of the clustering in the network. It indicates how nodes are embedded in their neighborhood.

## VI. CONCLUSION

The lasso has a major advantage over ridge regression, in that it produces simpler and more easily interpreted models that involve only a subset of the predictors. Lasso leads to qualitatively similar behavior to ridge regression, in that as  $\lambda$

Country	Coefficients	Respective Countries
Aruba	0.0026 -0.0033 -0.0007 0.5825 0.0196	Afghanistan Bulgaria Bosnia and Herzegovina Grenada Malta
Afghanistan	215.44 80.584 1.5730 -148.589 9.8225	Aruba Antigua and Barbuda Austria Grenada Luxembourg
Angola	0.2622 0.5245 0.3133 0.9837 1.8412	United Arab Emirates Benin Burkina Faso Central African Republic Comoros
Albania	1.5508 0.00016 0.6787 0.0056 0.0121	Bhutan Europe & Central Asia Lithuania Poland Russian Federation
Andorra	0.00042 0.11576 0.1050 0.0610 0.00216	Algeria Micronesia, Fed. Sts. Guam Marshall Islands Tajikistan

TABLE I  
COEFFICIENTS EXAMPLES

Country	Estimate	Year
Aruba	90853 92898 94992 97017 98737	2000 2001 2002 2003 2004
Afghanistan	20093756 20966463 21979923 23064851 24118979	2000 2001 2002 2003 2004
Angola	16440924 16983266 17572649 18203369 18865716	2000 2001 2002 2003 2004
Albania	3089027 3060173 3051010 3039616 3026939	2000 2001 2002 2003 2004
Andorra	65390 67341 70049 73182 76244	2000 2001 2002 2003 2004

TABLE II  
PREDICTION EXAMPLES

X	Avg. degree	Network diameter	Graph Density	Avg. clustering coefficient
Entire graph	5	11	0.019	0.206
Africa	2.321	10	0.045	0.183
Europe	2.060	15	0.042	0.093
Asia	1.022	6	0.023	0.115
North America	0.939	4	0.019	1.797
Oceania	0.368	2	0.020	0.000
South America	1.250	4	0.114	0.235
Asia + South America	3.147	15	0.025	0.205
Africa + Oceania				
Asia + South America + Africa + Oceania with Groupings	3.897	18	0.022	0.244
North America + Europe	2.554	11	0.031	0.084
North America + Europe with Groupings	3.047	23	0.024	0.193
Groupings only	1.891	18	0.042	0.330

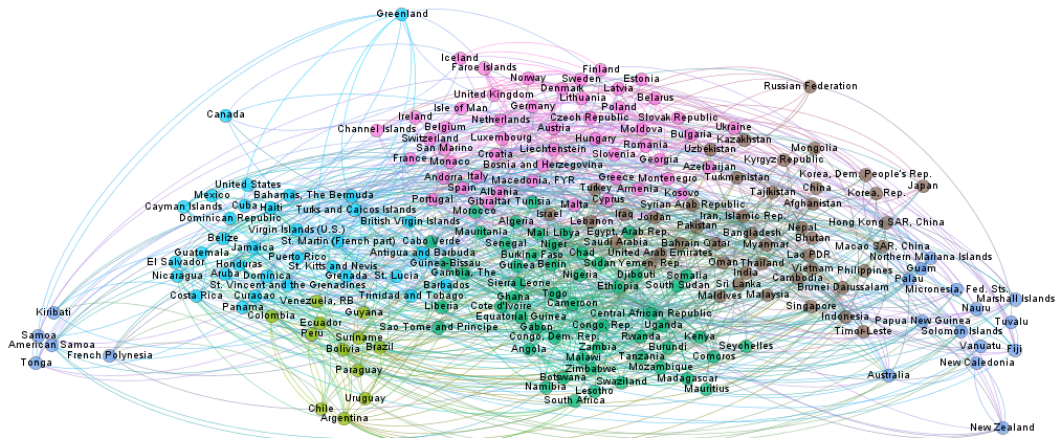
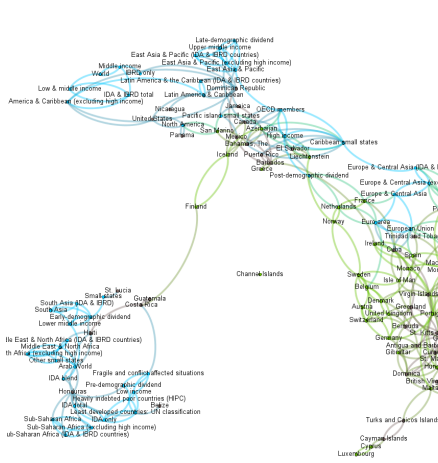
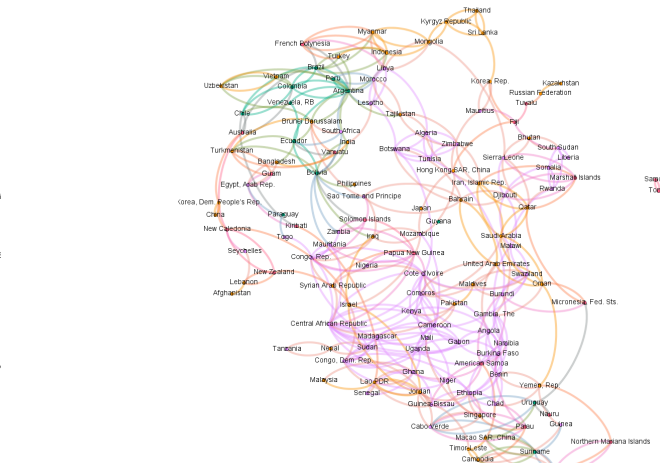


Fig. 1. World



America Europe with groupings.png



South America, Africa and Oceania.png

Fig. 2. North America Europe with Groupings

Fig. 4. Asia, South America, Africa and Oceania

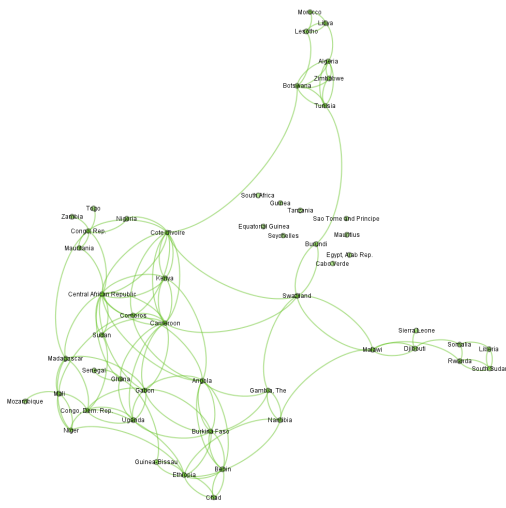


Fig. 3. Africa

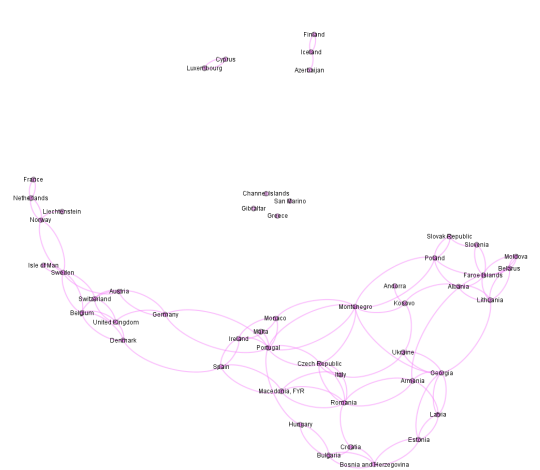


Fig. 5. Europe

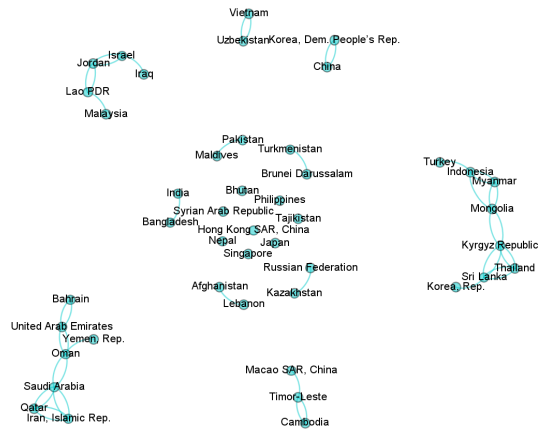


Fig. 6. Asia

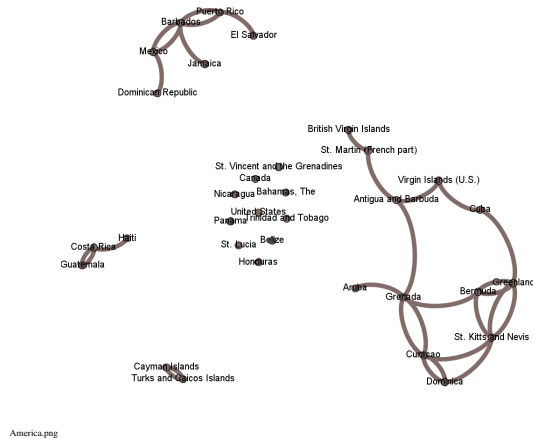


Fig. 7. North America

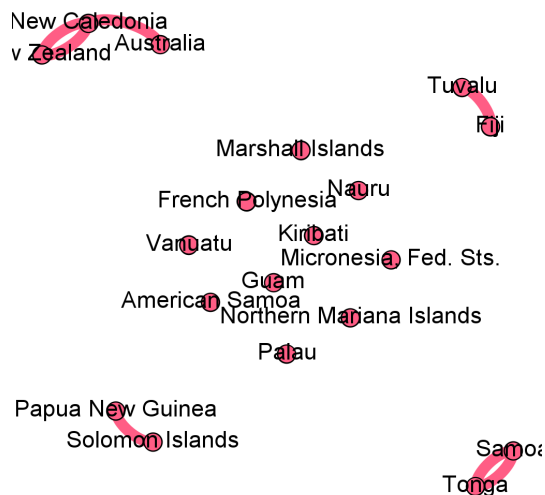


Fig. 8. Oceania

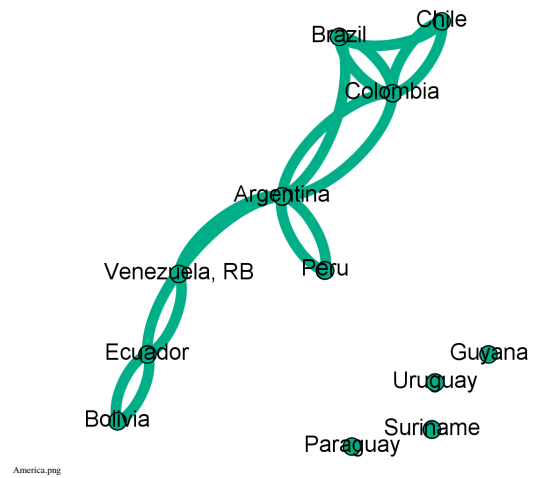


Fig. 9. South America

increases, the variance decreases and the bias increases. The lasso implicitly assumes that a number of the coefficients truly equal zero. Consequently, it is not surprising that ridge regression outperforms the lasso in terms of prediction error in this setting. In general, one might expect the lasso to perform better in a setting where a relatively small number of predictors have substantial coefficients, and the remaining predictors have coefficients that are very small or that equal zero. Ridge regression will perform better when the response is a function of many predictors, all with coefficients of roughly equal size. However, the number of predictors that is related to the response is never known a priori for real data sets. As with ridge regression, when the least squares estimates have excessively high variance, the lasso solution can yield a reduction in variance at the expense of a small increase in bias, and consequently can generate more accurate predictions. Unlike ridge regression, the lasso performs variable selection, and hence results in models that are easier to interpret. There are very efficient algorithms for fitting both ridge and lasso models; in both cases the entire coefficient paths can be computed with about the same amount of work as a single least squares fit.

## REFERENCES

- [1] An Introduction to Statistical Learning with Applications in R by Gareth James and Daniela Witten.
- [2] Understanding Machine Learning: From Theory to Algorithms by Shai Ben-David and Shai Shalev-Shwartz.