# CIS 520 - Project Report

## Section 1: Data Processing

### 1.1 Skew Correction

The distribution of the labels is skewed. We know that regression works better on normal distribution as regression models the mean. In a skewed distribution, the mean is not a good indicator or measure of central tendency, so in order to perform regression, we skew corrected it.

The table below shows the original skewness of each label and the skewness after calling the function transformLabels().

| Outcome Feature | Original Skew | Skew of transformed labels |
|---|---|---|
| health_aamort | 0.4557 | 0 |
| health_fairpoor | 0.6989 | 0.0625 |
| health_mentunh | 0.0336 | 0.0336 |
| Health_pcdiab | 0.4615 | 0.0513 |
| Health_pcexcdrin | 0.0761 | 0.0761 |
| Health_pcinact | -0.0886 | -0.0886 |
| health_pcsmoker | 0.3072 | -0.1104 |
| health_physunh | 0.3191 | 0.0348 |
| Heath_pcobese | -0.4445 | 0.0563 |

After we predict the labels for the test set, we reversed the skew correction.

## 1.2 Dimensionality Reduction:

We performed Principal Component Analysis on our 2000 Twitter features as we wanted to capture the variation in our data whilst removing any possible correlation. The table below shows the percentage of the total variance explained by selecting certain numbers of PCA components.

| Number of PCA Components | Percentage of Variance Explained |
|---|---|
| 30 | 89.2919 |
| 40 | 90.7310 |
| 80 | 93.7676 |
| 90 | 94.2255 |

As we can see, out of 2000 Twitter features, 90 principal components are enough to explain 94.2255% of the variance in the data.

# Section 2: Leaderboard Method

We used the same method described in our generative model for our leaderboard submission.

# Section 3: Other ML Techniques

## 3.1 Generative Model: Gaussian Process Regression

This is the model that we submitted to the leaderboard.

We used the inbuilt MATLAB function fitrgp() for implementing Gaussian Process Regression.

```
fitrgp(X,Y,'KernelFunction','rationalquadratic', 'Standardize',1, 'Sigma',
        obj.sigmaParams{index}, 'BasisFunction', 'none', 'FitMethod', 'sr')
```

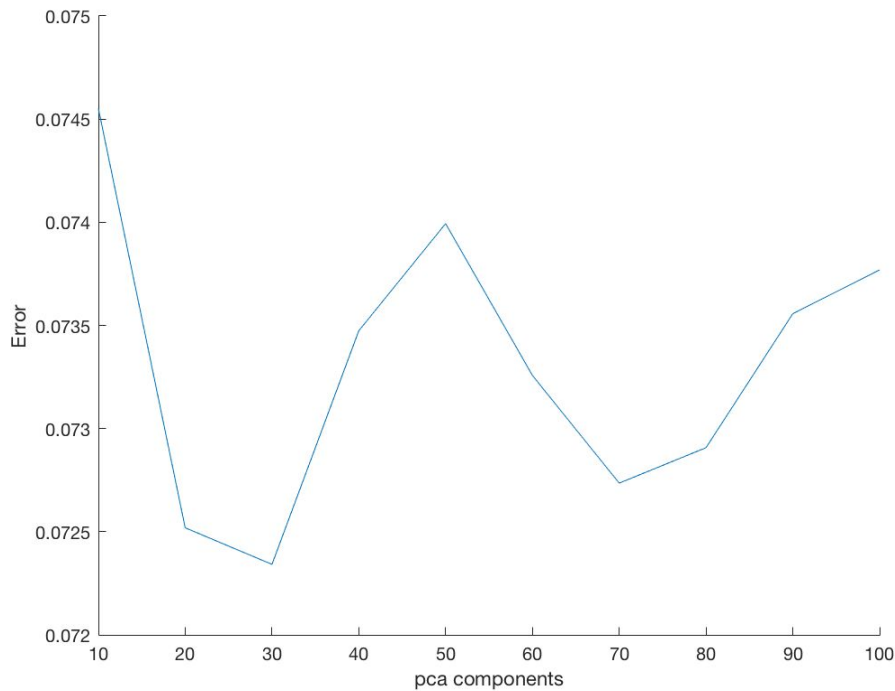We tried 'exponential', 'squaredexponential', 'matern32', 'matern52', and 'rationalquadratic' kernel functions.

| Kernels | Training Error | Test Error |
|---|---|---|
| exponential | 0.0025 | 0.0733 |
| squaredexponential | 0.0430 | 0.0739 |
| matern32 | 0.0327 | 0.0731 |
| matern52 | 0.0385 | 0.0731 |
| rationalquadratic | 0.0382 | 0.729 |

As we can observe from the above table, "rationalquadratic" gives a good tradeoff between the training and test error, so we chose that as our kernel function.

For regularization, we passed the FitMethod as "sr" (takes the standard deviation of the labels) and we tried using exact but that was overfitting the data.

To improve the accuracy of our model, we analyzed the PCA component vs Error graph for every outcome label and chose the number of principal components that gave us the least error. Below are the plots for the same.
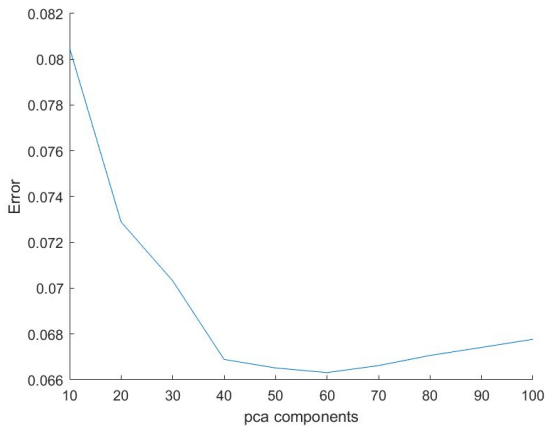
LABEL 1 :

As we can see from the above plot, we get the minimum error for 30 PCA components.
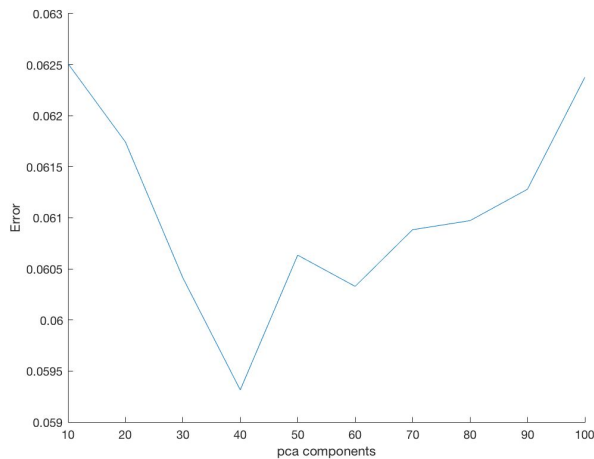Hence, number of PCA components chosen for Label 1 = 30.

LABEL 2 :



As we can see from the above plot, we get the minimum error for 70 PCA components.
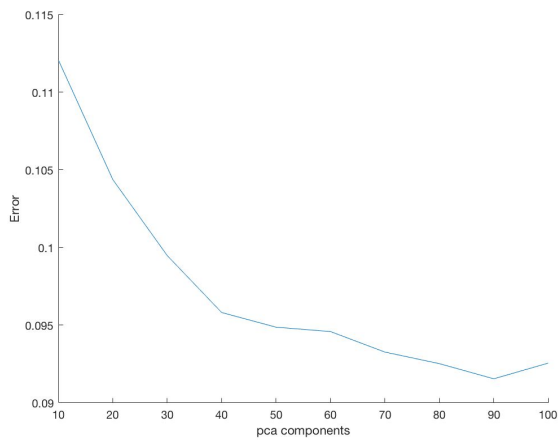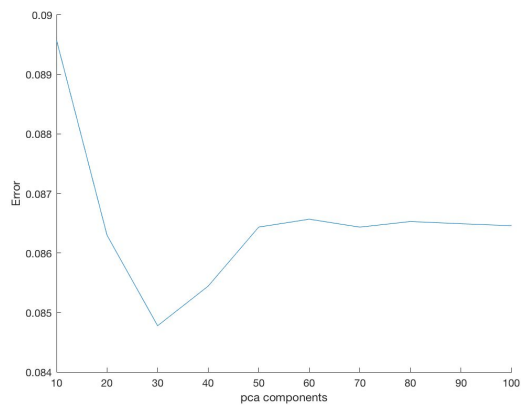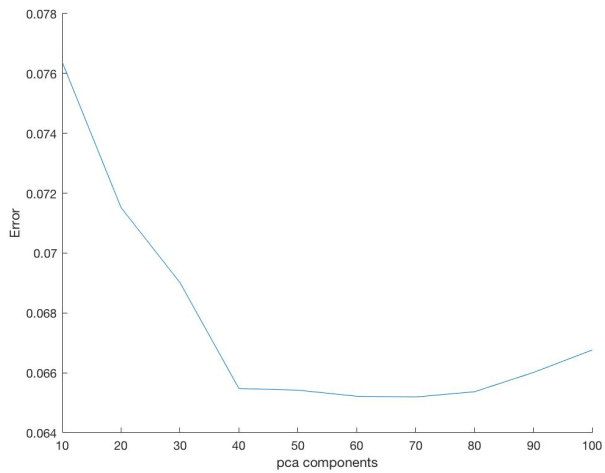Hence, number of PCA components chosen for label 2 = 70.
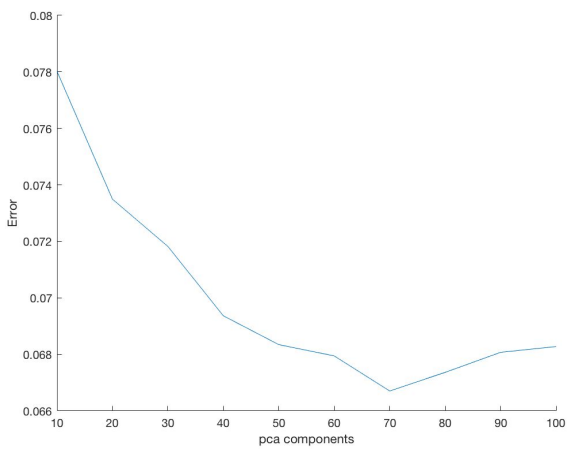
LABEL 3 :

As we can see from the above plot, we get the minimum error for 60 PCA components. Hence, number of PCA components chosen for label 3 = 60.

LABEL 4 :



As we can see from the above plot, we get the minimum error for 40 PCA components. Hence, number of PCA components chosen for label 4 = 40.

LABEL 5 :

As we can see from the above plot, we get the minimum error for 90 PCA components. Hence, number of PCA components chosen for label 5= 90.

LABEL 6 :



As we can see from the above plot, we get the minimum error for 30 PCA components. Hence, number of PCA components chosen for label 6 = 30.
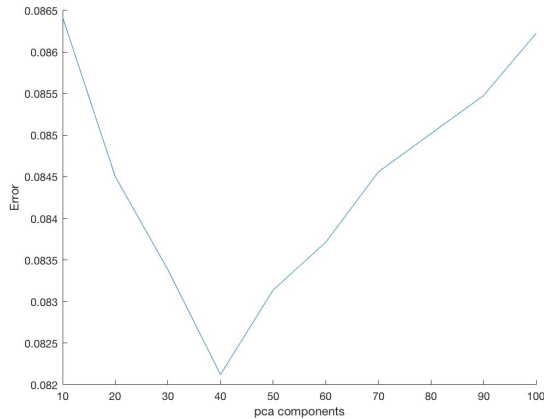
LABEL 7 :

As we can see from the above plot, we get the minimum error for 40 PCA components. Hence, number of PCA components chosen for label 7 = 40.

LABEL 8 :



As we can see from the above plot, we get the minimum error for 70 PCA components. Hence, number of PCA components chosen for label 8 = 70.

LABEL 9 :



As we can see from the above plot, we get the minimum error for 40 PCA components. Hence, number of PCA components chosen for label 9 = 40.

# 3.2 Discriminative Model: ElasticNet Regression

## 3.2.1 Model

We created a class implementing the ElasticNet model with the following parameter values:
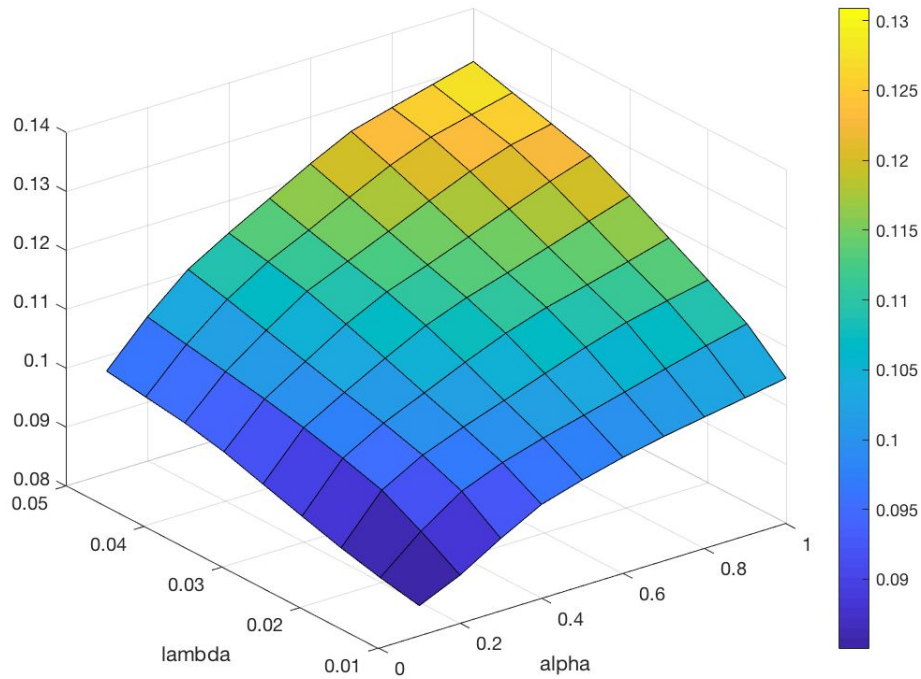- $\alpha = 1.0$
- $\lambda$ for each outcome:
    - health_aamort = 0.0085
    - health_fairpoor = 0.0028
    - health_mentunh = 6.8172e-05
    - Health_pcdiab = 9.4965e-4
    - Health_pcexcdrin = 3.053e-04
    - Health_pcinact = 0.0012
    - health_pcsmoker = 0.0026
    - health_physunh = 0.0014
    - Heath_pcobese = 5.4241e-04

where $\alpha$ is the ratio of weight of lasso regularization vs weight of ridge regularization.

If we try to optimize same alpha/lambda for all the labels,
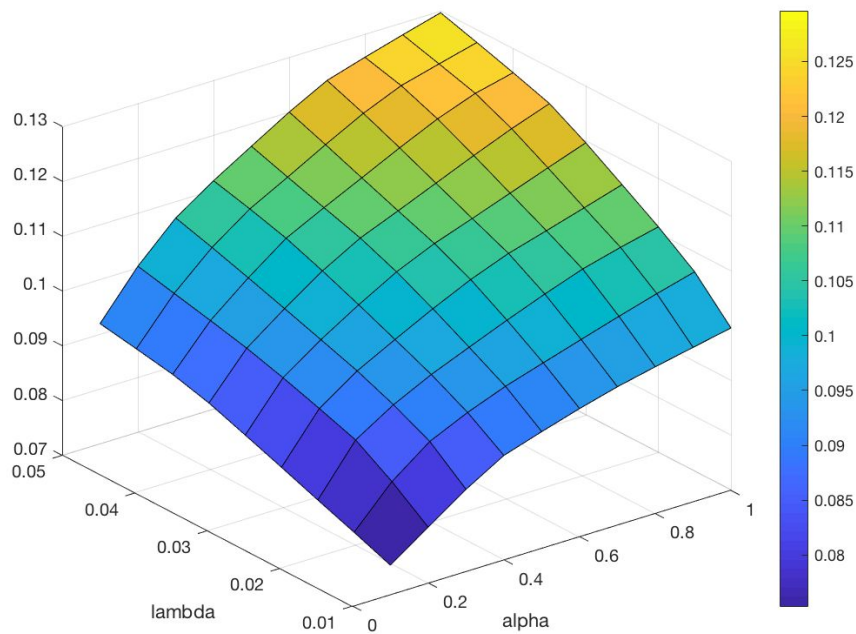
Surface Plot

for Cross Validation Error:



Cross Validation Errors :

| Alpha\Lambda | 0.01 | 0.015 | 0.02 | 0.025 | 0.03 | 0.035 | 0.04 | 0.045 | 0.05 |
|---|---|---|---|---|---|---|---|---|---|
| 0.1 | 0.085081 79614 | 0.086404 05856 | 0.088006 9137 | 0.089952 89962 | 0.091955 39853 | 0.093999 22298 | 0.095466 90855 | 0.096429 12038 | 0.097380 80436 |
| 0.2 | 0.088150 3512 | 0.092156 57685 | 0.095700 34999 | 0.097673 55518 | 0.099414 09353 | 0.100845 4686 | 0.102045 8972 | 0.103240 3966 | 0.104465 7499 |
| 0.3 | 0.092262 23227 | 0.096849 86441 | 0.099519 16615 | 0.101524 3262 | 0.103290 1482 | 0.105103 1835 | 0.106943 3602 | 0.108760 4017 | 0.110419 1769 |
| 0.4 | 0.095922 64261 | 0.099591 9611 | 0.102175 9103 | 0.104537 8866 | 0.106976 0507 | 0.109333 0851 | 0.111470 5623 | 0.113044 6431 | 0.114472 1466 |
| 0.5 | 0.097952 97238 | 0.101638 7204 | 0.104573 4667 | 0.107615 202 | 0.110421 5179 | 0.112695 6378 | 0.114442 2434 | 0.116274 902 | 0.118233 7574 |
| 0.6 | 0.099683 | 0.103411 | 0.107028 | 0.110436 | 0.113033 | 0.115151 | 0.117413 | 0.119749 | 0.122001 |

| | 68887 | 6628 | 0139 | 5513 | 3548 | 4452 | 7528 | 0571 | 6477 |
|---|---|---|---|---|---|---|---|---|---|
| 0.7 | 0.101118 3559 | 0.105224 6354 | 0.109386 1624 | 0.112707 9006 | 0.115149 238 | 0.117787 4925 | 0.120475 555 | 0.123042 1194 | 0.125384 5647 |
| 0.8 | 0.102308 4786 | 0.107063 7357 | 0.111503 3638 | 0.114440 0264 | 0.117393 3585 | 0.120447 4314 | 0.123363 4129 | 0.125752 7089 | 0.127303 5479 |
| 0.9 | 0.103478 3141 | 0.108860 3349 | 0.113057 1346 | 0.116256 8076 | 0.119673 9133 | 0.122979 0703 | 0.125708 0726 | 0.127430 708 | 0.129130 739 |
| 1.0 | 0.104680 1742 | 0.110480 7134 | 0.114456 4816 | 0.118157 4123 | 0.121882 5604 | 0.125253 8003 | 0.127200 9037 | 0.129072 8832 | 0.130909 6888 |

Surface Plot for Training Error:



| Alpha\Lambda | 0.01 | 0.015 | 0.02 | 0.025 | 0.03 | 0.035 | 0.04 | 0.045 | 0.05 |
|---|---|---|---|---|---|---|---|---|---|
| 0.1 | 0.075342 65156 | 0.077907 07241 | 0.080354 4146 | 0.082885 70178 | 0.085335 68463 | 0.087744 25991 | 0.089537 62724 | 0.090758 29905 | 0.091913 42568 |
| 0.2 | 0.080160 06482 | 0.085051 07536 | 0.089198 31957 | 0.091507 28504 | 0.093679 18961 | 0.095632 96823 | 0.097244 06553 | 0.098757 70409 | 0.100272 3768 |
| 0.3 | 0.084958 81994 | 0.090259 30432 | 0.093520 87925 | 0.096259 19532 | 0.098518 72303 | 0.100745 9644 | 0.102924 0595 | 0.105032 3405 | 0.106973 2313 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0.4 | 0.089031 76374 | 0.093441 60444 | 0.096928 55964 | 0.099873 10235 | 0.102768 7638 | 0.105516 4562 | 0.108038 0755 | 0.109954 005 | 0.111655 6275 |
| 0.5 | 0.091265 03126 | 0.096099 74951 | 0.099793 78189 | 0.103378 7129 | 0.106665 2607 | 0.109405 5268 | 0.111505 3476 | 0.113638 5311 | 0.115830 6141 |
| 0.6 | 0.093362 83332 | 0.098277 06148 | 0.102612 2586 | 0.106588 8937 | 0.109726 8743 | 0.112255 538 | 0.114833 8407 | 0.117451 3038 | 0.119957 7087 |
| 0.7 | 0.095235 47761 | 0.100425 1732 | 0.105281 0371 | 0.109253 8222 | 0.112182 1732 | 0.115189 8095 | 0.118204 5651 | 0.121056 5972 | 0.123617 9955 |
| 0.8 | 0.096767 22908 | 0.102533 9482 | 0.107734 0053 | 0.111282 076 | 0.114691 5935 | 0.118125 4263 | 0.121368 1472 | 0.124000 7869 | 0.125738 3199 |
| 0.9 | 0.098195 11951 | 0.104560 5036 | 0.109576 3964 | 0.113353 4455 | 0.117211 7895 | 0.120909 0645 | 0.123920 4545 | 0.125854 4597 | 0.127725 824 |
| 1.0 | 0.099630 17321 | 0.106436 8628 | 0.111208 0125 | 0.115462 2467 | 0.119658 5242 | 0.123378 6853 | 0.125574 6742 | 0.127645 5678 | 0.129631 9316 |

**Note:** However, for each alpha, optimization leads to unique values of lambda for each outcome label. In such cases, working with to each feature, and optimizing lambda for each alpha, resulted same cross validation error for all alpha's. We pick alpha as 1.0 as we prefer sparse vector

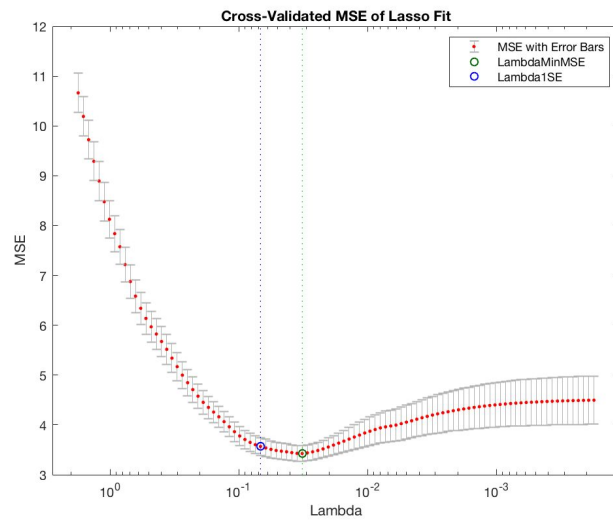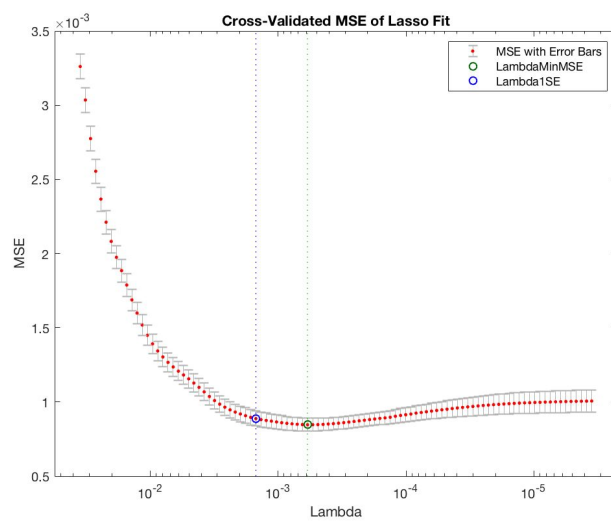Optimizing specific to each labels,

LABEL 1 :



LABEL 2 :



LABEL 3 :

Cross-Validated MSE of Lasso Fit

LABEL 4 :



Cross-Validated MSE of Lasso Fit
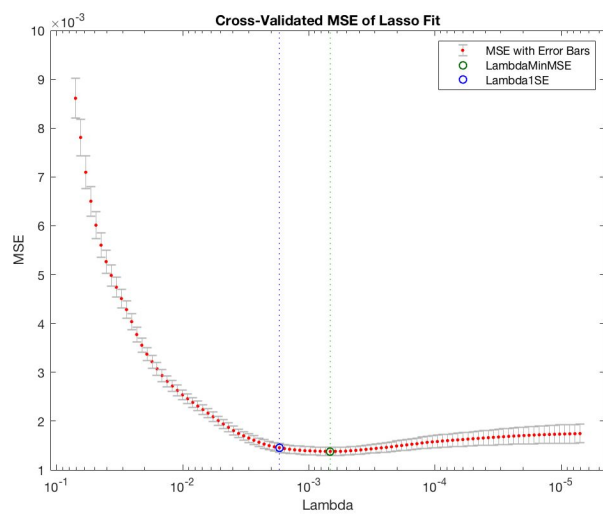
LABEL 5 :

LABEL 6 :



LABEL 7 :

LABEL 8 :



LABEL 9 :

Cross-Validated MSE of Lasso Fit
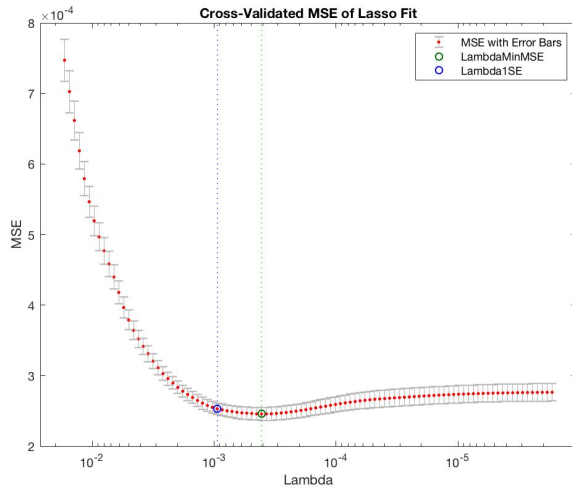
In order to actually create the ElasticNet Regression model with the parameters specified, we called the inbuilt MATLAB function lasso() and passed our data and parameters for each outcome feature.

```
lasso(X, Y, 'Alpha', obj.alpha,'Lambda', lambda)
```

### 3.2.2 Development

To develop a discriminative model for our data, we first tried using ridge regression, as we had a large number of features and wanted to prevent overfitting. However, the least value of cross validation error after performing 10-fold cross validation using Ridge regression was 0.0808. Clearly, simply adding a penalty was not good enough for the model.

However, as we have already performed PCA, the features that we are using will not be correlated, so instead of performing Lasso regression, we decided to try the ElasticNet regression model which led to a mean cross validation error of 0.0799 after performing 10-fold cross validation.

### 3.2.3 Results

Mean Cross Validation Error = 0.0811 (Optimizing for each label)
Train Error = 0.0702

# 3.3 Instance Based Model: K-Nearest Neighbors

## 3.3.1 Model

We have created two versions of KNN based on how we developed and tuned the models.

### 3.3.1.1 Model 1:

| Outcome Feature | K | Distance Metric |
|---|---|---|
| health_aamort | 187 | Correlation |
| health_fairpoor | 55 | Euclidean |
| health_mentunh | 58 | Seuclidean |
| Health_pcdiab | 6 | Seuclidean |
| Health_pcexcdrin | 6 | Correlation |
| Health_pcinact | 40 | Euclidean |
| health_pcsmoker | 335 | Euclidean |
| health_physunh | 1 | Seuclidean |
| Heath_pcobese | 95 | Euclidean |

### 3.3.1.2 Model 2:

K-Nearest Neighbors has hyperparameters K = 1, and Distance Metric = Euclidean.

In order to create both these models, we used the inbuilt MATLAB function fitcknn(), where we passed these hyperparameters.

## 3.3.2 Development

To develop Model 1, we used the "OptimizeHyperParameters" option with the value "auto" which means that we only optimize K and the distance metric. This performed 5-fold cross validation on the training data and found the hyperparameters that minimized loss. As we used

more than 10 feature columns, the "exhaustive" search method was used. This means that for every point, we find the distance from every other point to identify the K-Nearest Neighbors. The other option is to use the "kdtree" option in which the algorithm creates a tree like structure for storing and finding the k-nearest neighbors of any new point.

To develop Model 2 on the other hand, we manually found the lowest cross validation and train set errors for values from K=1 to K=15. We performed 10-Fold Cross validation in each case. The lowest cross validation error was for K=1, so our model performs 1-Nearest Neighbor classification.



### 3.3.3 Results:

Model 1:

Cross Validation Error = 0.1441
Train Error = 0.1226

Model 2:
Cross Validation Error = 0.1390

Train Error = 0

## 3.4 Novel Method: Ensemble - Stacking



Stacking requires weak learners so \lambda is 0.0002. Will not work with strong classifiers

This was our 3rd submission to the leaderboard., Ensemble Error on Leader Board - 0.0775

As we can observe from the above table, for Lasso, though the error is decreasing as we increase the number of classifiers, but we chose 50 classifiers as increasing the number of classifiers increases the time complexity.

For Lasso, we selected 200 features out of 2000 features because the probability of missing out a feature comes out to be 0.00051 which is very less and affordable.

For Ridge, we selected 200 features out of 2000 features because the probability of missing a feature comes out to be $7*(10^{-10})$.

The reason why we chose 200 features out of 2000 features, 50 and 200 classifiers in case of lasso and ridge respectively because, increasing any of the values would make it computationally time consuming and with the chosen values, the probability of missing a feature is very less.

We did some experiments on choosing the number of classifiers for lasso and Ridge. Keeping the feature count limit to 200.

Cross Validation Errors:
Elastic Net Config(0.001, 0.002):



Ridge Config(0.001, 0.002):

**Results:**

Cross Validation Error: 0.0799
Training Error: 0.0684

# Section 5: Interpretation

## 5.1 Correlation Between The Labels

The heatmap showing the correlation between the different outcome variables is given below.

In this heatmap, the more yellow the block, the more correlation between the two labels.

The most surprising fact is that outcome label 5, 'health_pcexcdrin', which is the percent of adults that report excessive drinking seems to have no correlation to any other outcome label. This label was also the hardest to classify as it gave us the maximum error. According to the heatmap, label 5 had the most correlation to label 9, which indicates the percent of adults who reported being obese. This makes sense, as there is medical evidence that alcohol causes weight gain.

There also appears to be a correlation between label 8, average number of reported physically unhealthy days, and labels 2 and 3, which are the percent of adults reporting fair or poor health and the percent of adults reporting mentally unhealthy days in a month. This mean that more physical activity results in better health and better mental health too, or at least, there is some correlation between these variables.

## 5.2 Comparison between Feature and Labels

**health_aamort**

**Elastic Net Weights**



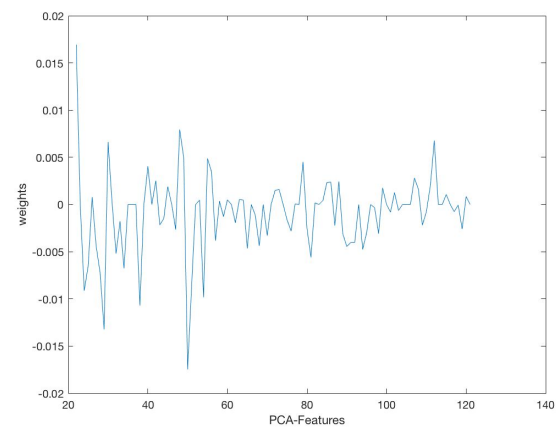**Regression Tree Weights (fraction representing the importance of feature.)**



We see that for **health_aamort,** which predicts the years of potential life lost, Demographic Feature, ses_log_hhinc, which is the median household income, has more importance. This could

be because a lower median income results in lack of facilities and perhaps inability to afford healthcare. Principal Component 1 has more weight. Followed by 9, 17.

| health_fairpoor |
| --- |

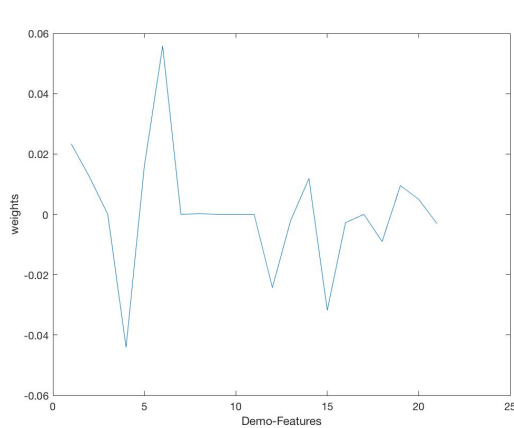| Elastic Net Weights |
| --- |



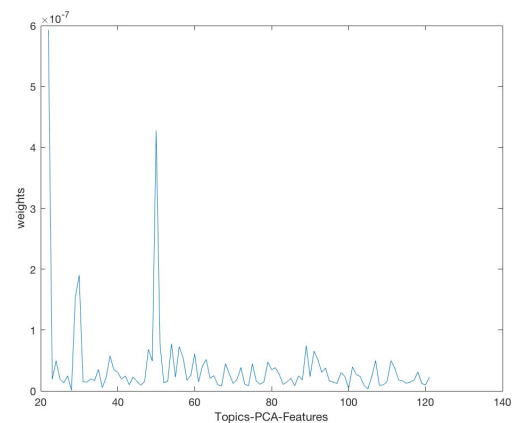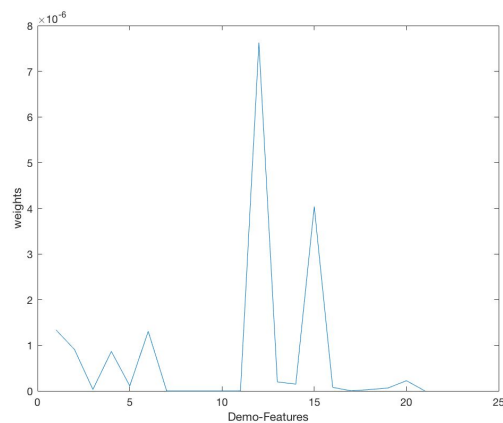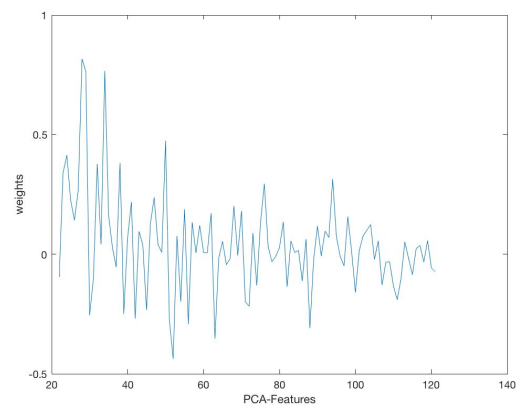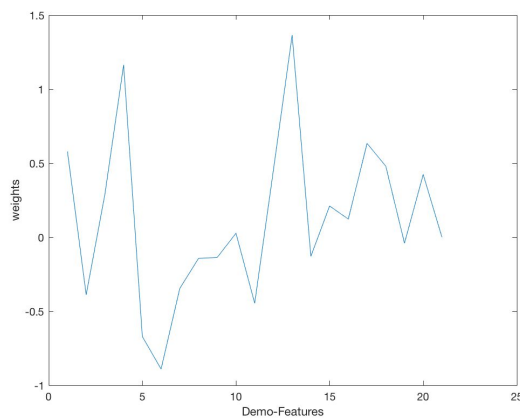| Regression Tree Weights (fraction representing the importance of feature.) |
| --- |



We see that for **health_fairpoor,** Demographic Features, household income, percentage of population with a college degree, and the percentage of the population that is white have more

importance. Income has been given the highest weight. Clearly, income plays a large role in determining the health statistics of a county. And the Principal Component 1 has more weight. Followed by 9, 7.

**health_mentunh**

**Elastic Net Weights**



**Regression Tree Weights (fraction representing the importance of feature.)**

We see that for **health_mentunh,** Demographic Features, household income, the percentage of population with a college degree and the percentage of unemployed people has more importance. And the Principal Component (7,9, 17) has more weight.

| Health_pcdiab |
| --- |

| Elastic Net Weights |
| --- |



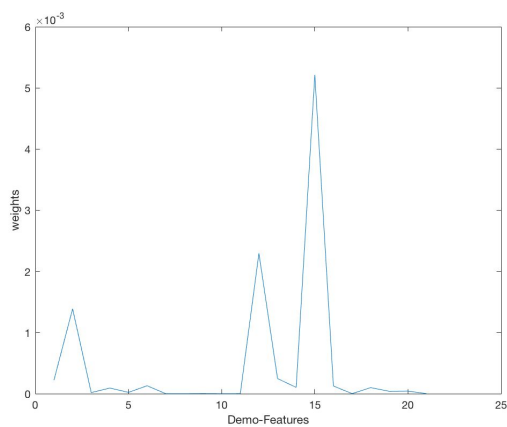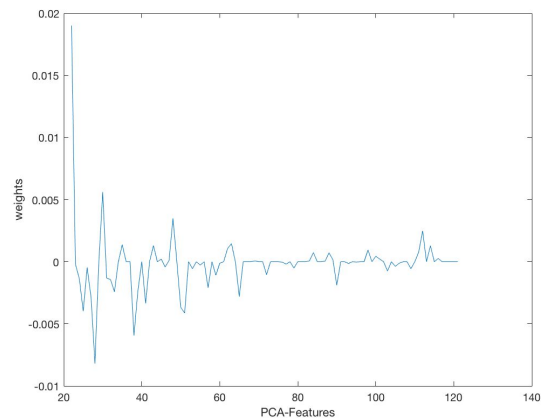| Regression Tree Weights (fraction representing the importance of feature.) |
| --- |

We see that for **Health_pcdiab,** Demographic Features, the percentage of white population, percentage of people over 65, income and percent of college degree holders have more importance. And the Principal Component (1, 29, 9) has more weight.
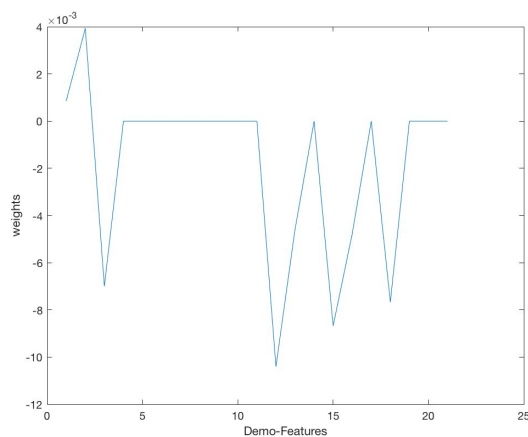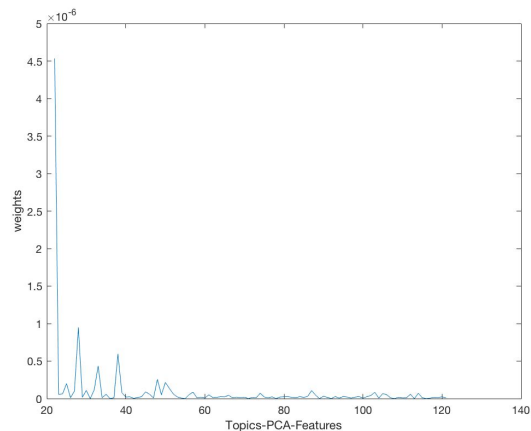
**Health_pcexcdrin**

**Elastic Net Weights**



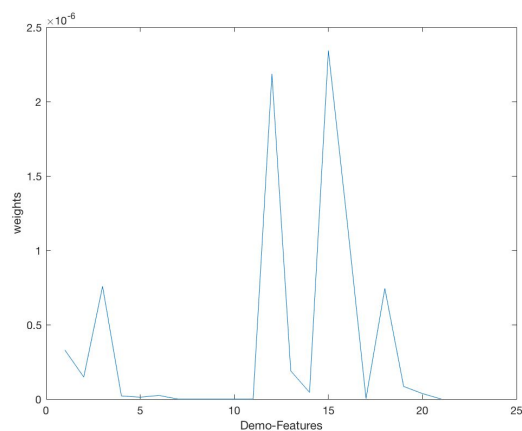**Regression Tree Weights (fraction representing the importance of feature.)**

We see that for **Health_pcexcdrin,** Demographic Features, the percentage of white population, percentage of people over 65, income and percent of college degree holders have more importance. And the Principal Component (7, 9, 13) has more weight.

## Health_pcinact
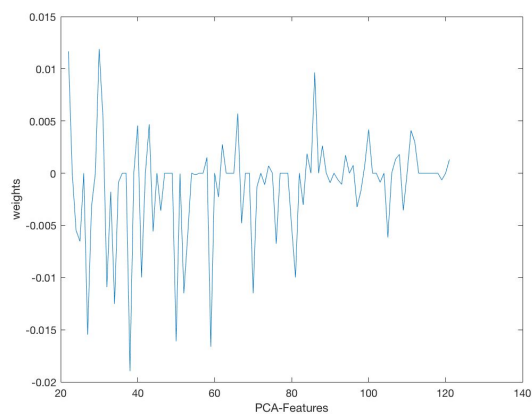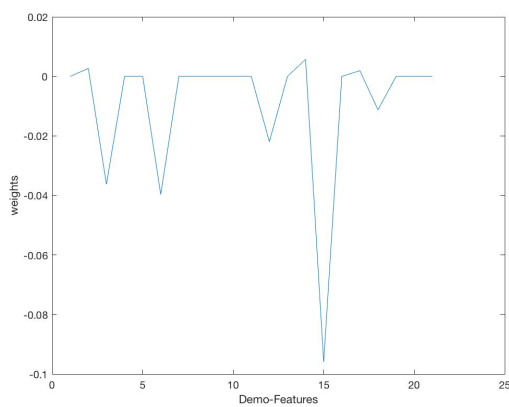
## Elastic Net Weights



## Regression Tree Weights (fraction representing the importance of feature.)
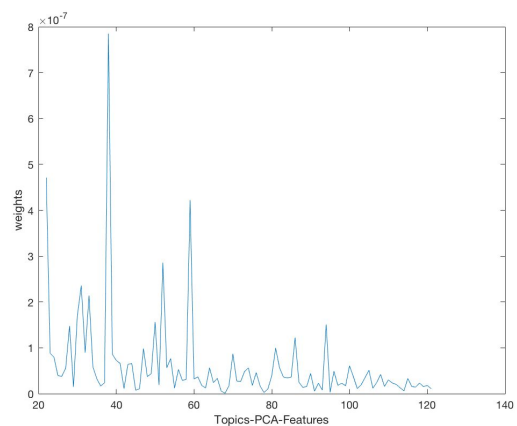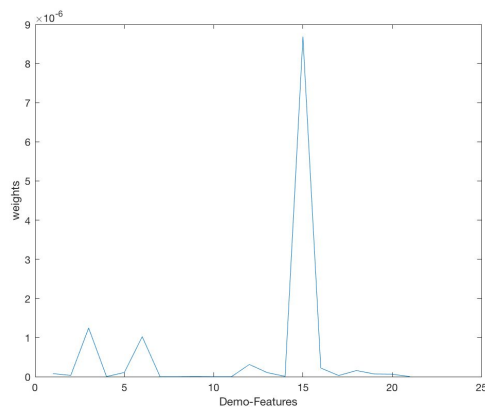
We see that for **Health_pcinact,** Demographic Features, household income and percentage of population with a college degree has more importance. And the Principal Component (1,7) has more weight.

**health_pcsmoker**

**Elastic Net Weights**



**Regression Tree Weights (fraction representing the importance of feature.)**

We see that for **health_pcsmoker,** Demographic Feature, household income, percentage of hispanic population, and percentage of people under the age of 18 have more importance. And the Principal Component (17, 38) has more weight.
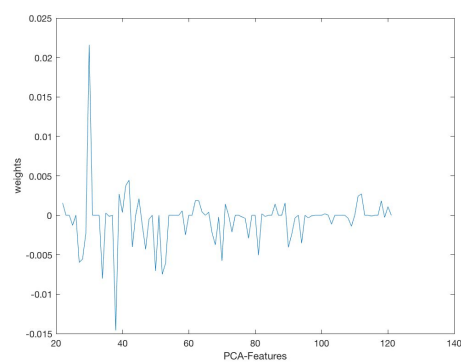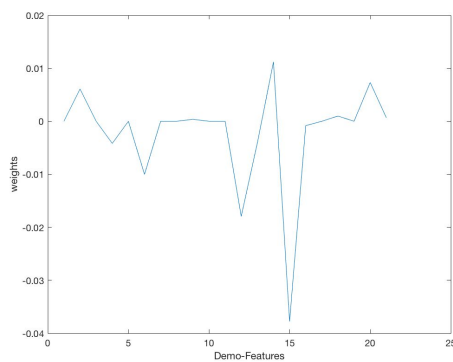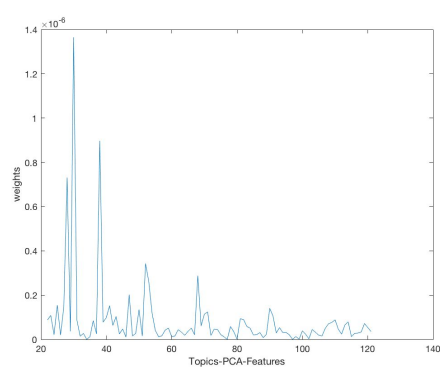
**health_physunh**

**Elastic Net Weights**



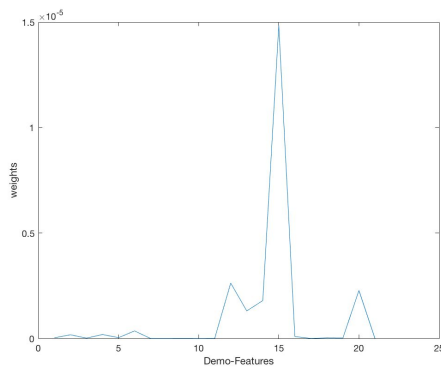**Regression Tree Weights (fraction representing the importance of feature.)**



We see that for **health_physunh,** Demographic Features household income and percentage of college degree holders have more importance. And the Principal Component (9, 7, 17) has more weight.

| Heath_pcobese |
| --- |



We see that for **Heath_pcobese,** Demographic Features, household income, percentage of college educated people and percentage of black population have more importance. And the Principal Component (1, 2) has more weight.

# Conclusion

## Model Comparison

| Model | Cross-Validation Error |
| --- | --- |

| Gaussian Process Regression, Fit-Grp with non-linear kernel | 0.725 |
| --- | --- |
| Stacking | 0.0799 |
| ElasticNet | 0.0811 |
| KNN | 0.1441 |

Clearly, non-linear models work better on the data. Stacking and ElasticNet are linear models (as we have used Lasso and Ridge Regression models in our stacking algorithm). While KNN can model non-linear relationships, in this case, it is hard to say what metric is suitable as a distance metric and how we should calculate the K-nearest neighbors.

Therefore, the Gaussian Process Regression Model gave us the lowest error.

## Improving Accuracy:

To improve the accuracy of our model, we analyzed the PCA component vs Error graph for every outcome label and chose different number of principal components for every label that gave us the least error.

Skew correction of labels has helped to increase the accuracy. Used Log, square root and cubic-root for skew correction of labels. Linear regression gives better accuracies when the data distribution is similar to normal distribution. We also used PCA to reduce the dimensionality of the data and

DIFFERENT MODELS TRIED AND WHY IT DIDN'T WORK :

K nearest neighbours uses distance metric for learning, but this is a regression problem, so standardization is not a good thing as we do not know the metrics of the features. That is why KNN performed bad.

Neural Networks didn't work because the data that was given to us is very less ( around 1000 samples), and the neural net had to learn a lot of parameters(weights), but since the data was not enough it couldn't learn the optimal parameters.

For Boosting Trees the error was 0.08. It wasn't working bad but gaussian regression process was working better. Since data can be modeled properly with a kernel, fitgrp finds kernel parameters such that it gives an optimal fit. Whereas for boosting trees, there is a trade-off between over-fitting and under-fitting and it was difficult to generalize as we approximate too

much it was leading to over-fitting. Even after finding a trade-off, gaussian regression process was giving better results as compared to boosting trees.

The data given to us is non-linear and stacking and elastic net are linear regression methods, whereas gaussian regression process is non linear. So gaussian regression process worked the best among the models that we tried.

DATA PROCESSING THAT DIDN'T WORK:
Skew correction on input data didn't work because we are doing skew correction independent of every column feature . Even though the overall skew was less, it could have been that it could have missed some correlation and so we were getting better training error without performing skew correction on the input data.

We also tried CCA for dimensionality reduction as performing only PCA on the input set may lead to loss of information that is correlated with our labels. However, that did not work because the rank of our labels is 9 and CCA reduces our Twitter data to a 9-dimensional subspace from 2000 dimensions. Therefore, it is not able to capture the correlation between the features and the labels very well.

# Future Work

In the future, finding an approximate non-linear transformation could help model the data better.