# ▾ Assignment 3 (From Scratch)

## ▾ Penalized Logistic Ridge Regression CV with Batch Gradient Decent

- **Programmers:**
  - Shaun Pritchard
  - Ismael A Lopez
- **Date:** 11-15-2021
- **Assignment:** 3
- **Prof:** M.DeGiorgio

---

## Overview: Assignment 3

- In this assignment you will still be analyzing human genetic data from $N$ = 183 training observations (individuals) sampled across the world. The goal is to fit a model that can predict (classify) an individual's ancestry from their genetic data that has been projected along $p$ = 10 top principal components (proportion of variance explained is 0.2416) that we use as features rather than the raw genetic data

- Using ridge regression, fit a penalized (regularized) logistic (multinomial) regression with model parameters obtained by batch gradient descent. Based on K = 5 continental ancestries (African, European, East Asian, Oceanian, or Native American), predictions will be made. Ridge regression will permit parameter shrinkage (tuning parameter $\lambda \geq 0$) to mitigate overfitting. In order to infer the bestfit model parameters on the training dataset, the tuning parameter will be selected using five-fold cross validation. After training, the model will be used to predict new test data points.

## ▾ Imports

> Import libaries and data

```
1   #Math libs
2   from math import sqrt
3   from scipy import stats
4   from numpy import median
5   from decimal import *
```

```
 6   import os
 7   # Data Science libs
 8   import numpy as np
 9   import pandas as pd
10   # Graphics libs
11   import seaborn as sns
12   import matplotlib.pyplot as plt
13   %matplotlib inline
14   #Timers
15   # !pip install pytictoc
16   # from pytictoc import TicToc
```

```
1   # Import training and test datasets
2   train_df = pd.read_csv('TrainingData_N183_p10.csv')
3   test_df = pd.read_csv('TestData_N111_p10.csv')
```

```
1   # Validate traning data import correclty
2   train_df.head(2)
```

|   | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | |
|---|---|---|---|---|---|---|---|---|---|
| **0** | -10.901171 | 0.798743 | -1.143301 | -1.070960 | 11.856396 | -2.265965 | 4.536405 | 1.519959 | -2.2 |
| **1** | -9.990054 | 1.416821 | -0.729626 | -0.443621 | 10.418594 | 0.443514 | 2.640659 | -4.637746 | 3.3 |

```
1   # Validate testing data import correclty
2   test_df.head(2)
```

|   | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 2.517683 | 5.464283 | 9.067873 | -4.965928 | -0.741937 | 0.039785 | 0.573279 | -0.216918 | 2.454 |
| **1** | 6.077012 | 1.032867 | -5.795883 | -3.490064 | -0.600204 | -0.120803 | 1.243767 | 1.821390 | -1.173 |

## Data Pre-Proccessing

- Pre-proccess test and training datasets
- Impute categorical variables in features
- Validate correct output of test data

```
1   # recode the categories
2   data = train_df['Ancestry'].unique().tolist()
3   num_features = len(data)
4   train_df['Ancestry2'] = train_df['Ancestry'].apply(lambda x: data.index(x))
```

5

```
1   # Validate training data set
2   train_df.head(2)
```

|   | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 |
|---|-----|-----|-----|-----|-----|-----|-----|-----|
| **0** | -10.901171 | 0.798743 | -1.143301 | -1.070960 | 11.856396 | -2.265965 | 4.536405 | 1.519959 | -2.2 |
| **1** | -9.990054 | 1.416821 | -0.729626 | -0.443621 | 10.418594 | 0.443514 | 2.640659 | -4.637746 | 3.3 |

```
1   # Shape trianing data
2   train_df.shape
```

```
(183, 12)
```

## Seperate X Y Predictors and Responses

- Seperate predictors from responses
- Validate correct output

```
1   # Seperate dependant categorical feature data for training and test data set
2   Y_train_names = train_df['Ancestry'].tolist()
3   Y_test_names = test_df['Ancestry'].tolist()
```

```
1   # Separate training feature predictors from responses
2   X_train = np.float32(train_df.to_numpy()[:, :-2])
3   Y_train = train_df['Ancestry2'].to_numpy()
```

```
1   # Separate test feature predictors from responses
2   X_test = np.float32(test_df.to_numpy()[:, :-1])
```

```
1   X_train.shape
```

```
(183, 10)
```

```
1
```

## Set Global Vairibles

- $\lambda$ = tunning parameters

- α = learning rate
- k = number of folds
- n_itters = nu mber of itterations
- X_p = predictor vlaues from training data
- Y_p = response values from training data

```
1   # Set local variables
2
3   # Tuning Parms
4   λ  =  10 ** np.arange(-4., 4.)
5
6   # learning rate
7   α =  1e-4
8
9   # K-folds
10  k = 5
11
12
13  # Itterations
14  n_iters = 10000
15
16  # Set n x m matrix predictor variable
17  X_p = X_train
18
19  # Set n vector response variable
20  Y_p  = Y_train
```

## Instantiate Data

- Handle logic and set variables needed for calculating ridge logistic regression
- Handle logic and set variables for batch gradient descent
- Handle logic and set variables for cross-validation

```
1   # Encode response variable from CV design matrix for cross vlaidation
2   def imputeResponse(response_vector, num_features):
3       response_vector = np.int64(response_vector)
4       X1 = response_vector.shape[0]
5       response_mat = np.zeros([X1, num_features])
6       response_mat[np.arange(X1), response_vector] = 1
7       return response_mat
```

```
1 # Method to handle randomization of training data predictors and responses
```

```
2 def randomizeData(X_p, Y_p):
3     data = np.concatenate((X_p, Y_p[:, None]), 1)
4     np.random.shuffle(data)
5     return data[:, :-1], data[:, -1]
```

```
1 # Randomize predictors and responses into new vairables
2 x, y = randomizeData(X_p, Y_p)
```

```
1 # Set Global variable for samples and number of X = N X M features
2 X1 = x.shape[0]
3 X2 = x.shape[1]
```

```
1 # Get number of training feature classes = 5
2 num_features = np.unique(y).size
```

```
1  # Call method imputation method on training response variables
2  y = imputeResponse(y, num_features)
```

```
1 # Store 5 K-fold cross validation results in symetric matrices
2 CV = np.zeros([k, len(λ)])
```

```
1 # Number of validation sample index values based on k-folds
2 val_samples = int(np.ceil(X1 / k))
3 test_i = list(range(0, X1, val_samples))
```

```
1 # Create a β zero matrix to store the trained predictors
2 β = np.zeros([k, len(λ), X2 + 1, num_features])
```

## ▾ Implement logic

- Main functions to handle logic within the preceding algorithms

```
1 # Standardize X coefficients
2 def standardize(x, mean_x, std_x):
3     return (x - mean_x) / std_x
```

```
1 # Concatenate ones column matrix with X coefficiants
2 def intercept(x):
3     col = np.ones([x.shape[0], 1])
4     return np.concatenate((col, x), 1)
```

```
1 # Predict standardize expotential X values from intercepts
2 def predict(x):
3   x = standardize(x, mean_x, std_x)
4   x = intercept(x)
5
6   X_p = np.exp(np.matmul(x, βx))
7   return X_p / np.sum(X_p, 1)[:, None]
```

```
1 # Splitting the data into k groups resampling method
2 def cv_folds(i_test):
3     if i_test + val_samples <= X1:
4         i_tests = np.arange(i_test, i_test + val_samples)
5     else:
6         i_tests = np.arange(i_test, X1)
7
8     x_test = x[i_tests]
9     x_train = np.delete(x, i_tests, axis = 0)
10
11
12     y_test = y[i_tests]
13     y_train = np.delete(y, i_tests, axis = 0)
14     return x_train, x_test, y_train, y_test
```

```
1 # Calculate model CV score
2 def score(x, y, βx):
3     # Compute exponent values of X coef and BGD unnormilized probality matrix
4     U = np.exp(np.matmul(x, βx))
5     # Calculate sum unnormilized probality / sum unnormilized matrix by 1
6     P = U / np.sum(U, 1)[:, None]
7     # Calulate to cost error score
8     err = -(1 / x.shape[0]) * np.sum(np.sum(y * np.log10(P), 1))
9     return err
```

## Batch Gradient Descent

> Alorithm 1 used for this computation

$$\mathbf{B} := \mathbf{B} + \alpha[\mathbf{X}^T(\mathbf{Y} - \mathbf{P}) - 2\lambda(\mathbf{B} - \mathbf{Z})]$$

```
1
2 def BGD(x, y, βx, lamb):
```

```
 3      # Unormalized class probability matrix
 4      U = np.exp(np.dot(x, βx))
 5      # Normalized class probability matrix
 6      P = U / np.sum(U, 1)[:, None]
 7      # K intercept matrix
 8      Z = βx.copy()
 9      Z[1:] = 0
10      # Update parameter matrix
11      βx = βx + α * (np.matmul(np.transpose(x), y - P) - 2 * lamb * (βx - Z))
12      return βx
```

```
1
```

## ▾ CV Ridge Penlized Logistic Regression

- Compute ridge-penalized logistic regression with cross vlaidation
- Performing a ridge-penalized logistic regression fit to training data $\{(x1, y1), (x2, y2), \ldots, (xN, yN)\}$ is to minimize the cost function

$$J(\mathbf{B}, \lambda) = -\log \mathcal{L}(\mathbf{B}) + \lambda \sum_{k=1}^{K} \sum_{j=1}^{p} \beta_{jk}^{2}$$

```
 1 # Compute ridge-penalized logistic regression with cross vlaidation
 2 for i_lambda, lamb in enumerate(λ):
 3     for i_fold, i_test in zip(range(k), test_i):
 4
 5         # Validates and trains the CV iteration based on the validation and tr
 6         x_train, x_test, y_train, y_test = cv_folds(i_test)
 7
 8         # Standardize x and center y  5 K-fold trianing and test data
 9         mean_x, std_x = np.mean(x_train, 0), np.std(x_train, 0)
10
11         # implement standardize X training and test sets
12         x_train = standardize(x_train, mean_x, std_x)
13         x_test = standardize(x_test, mean_x, std_x)
14
15         # Add training and test intercept column to the design matrix
16         x_train = intercept(x_train)
```

```
17          x_test = intercept(x_test)
18
19          # initialize Beta coef for lambdas and fold
20          βx =  np.zeros([X2 + 1, num_features])
21
22          # Loop through beta and lambdas with batch gradient decent
23          for iter in range(n_iters):
24              βx = BGD(x_train, y_train, βx, lamb)
25
26          # Score CV cost error tp the model and store the values
27          CV[i_fold, i_lambda] = score(x_test, y_test, βx)
28
29          # Save the updated coefficient vectors βx
30          β[i_fold, i_lambda] = βx
```

## ▾ Deliverable 1

> Illustrate the effect of the tuning parameter on the inferred ridge regression
> coefficients by generating five plots (one for each of the $K$ = 5 ancestry classes) of
> 10 lines (one for each of the $p$ = 10 features)

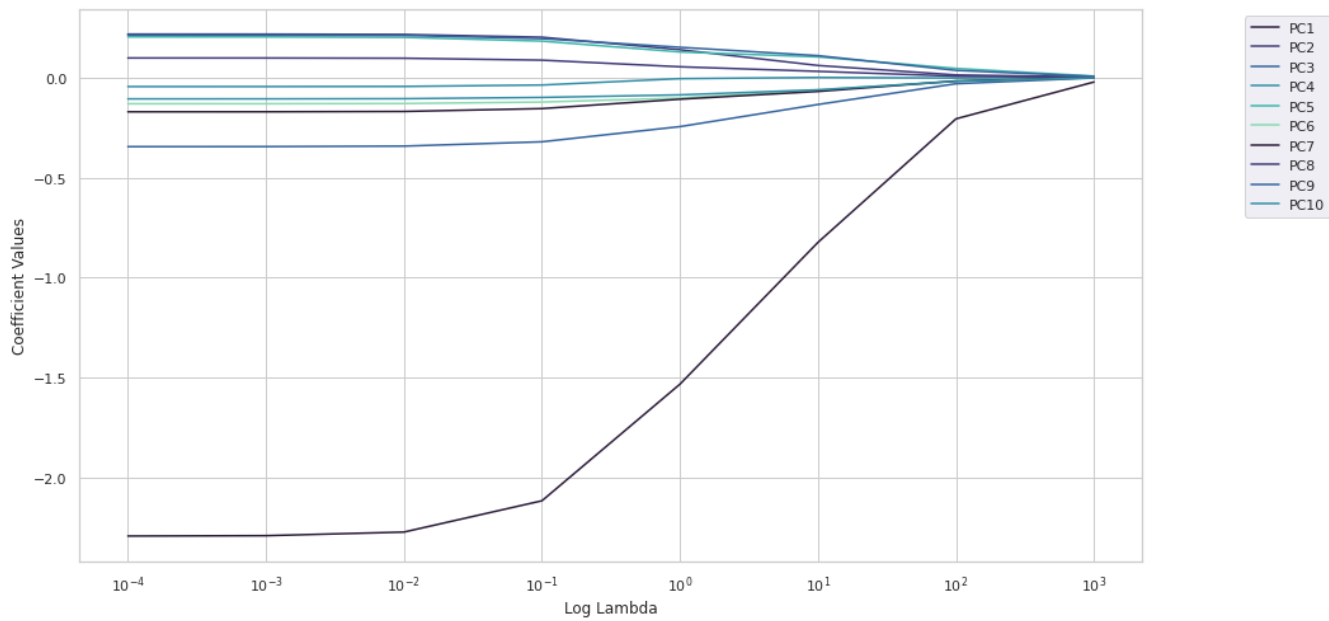$$CV_{(5)} = \frac{1}{5} \sum_{m=1}^{5} CategoricalCrossEntropy_m$$

```
 1 # Plot tuning parameter on the inferred ridge regression coefficients
 2 βμ = np.mean(β, 0)
 3 sns.set(rc = {'figure.figsize':(15,8)})
 4 for i, c in enumerate(data):
 5     βμk = βμ[..., i]
 6     sns.set_theme(style="whitegrid")
 7     sns.set_palette("mako")
 8     for j in range(1, 1 + X2):
 9         sns.lineplot( x=λ, y=βμk[:, j], palette='mako',    label = 'PC{}'.forma
10         sns.set()
11     plt.xscale('log')
12     plt.legend(bbox_to_anchor=(1.09, 1), loc='upper left')
13     plt.xlabel('Log Lambda')
14     plt.ylabel('Coefficient Values')
15     plt.suptitle('Inferred Ridge Regression Coefficient Tuning Parameters of'
```
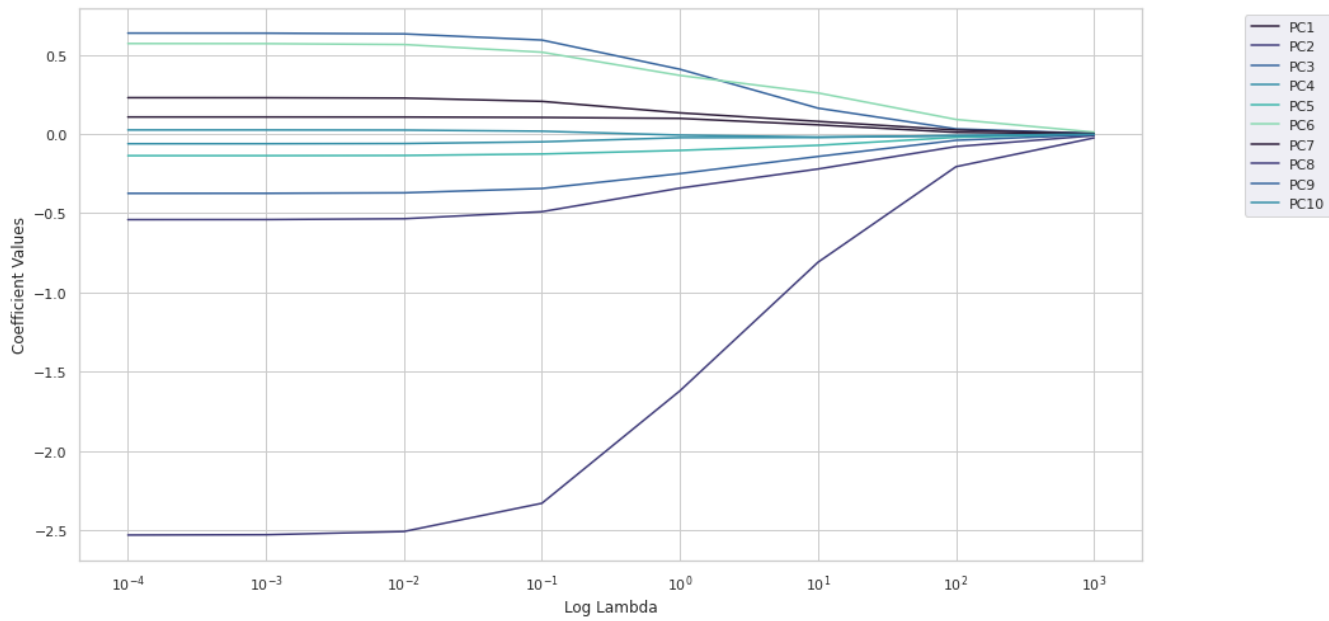
```
16      for l in range(i):
17        # Output Deliverable 1
18        plt.savefig("Assignment3_Deliverable1.{}.png".format(l))
19    plt.show()
20
21
22
23
```
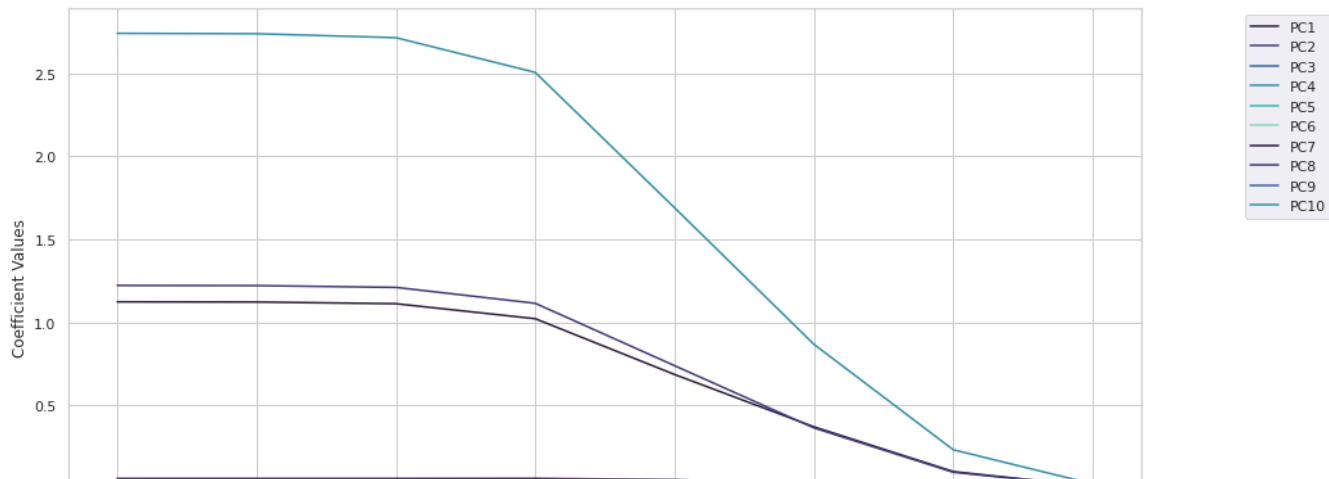
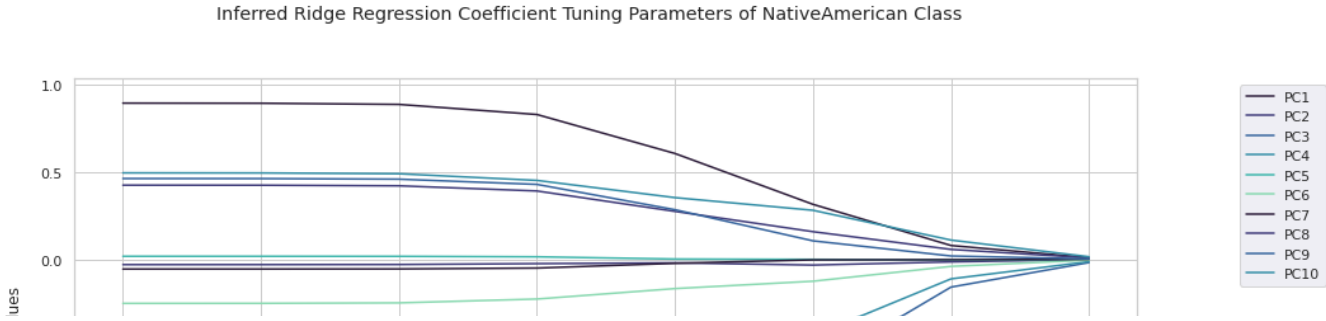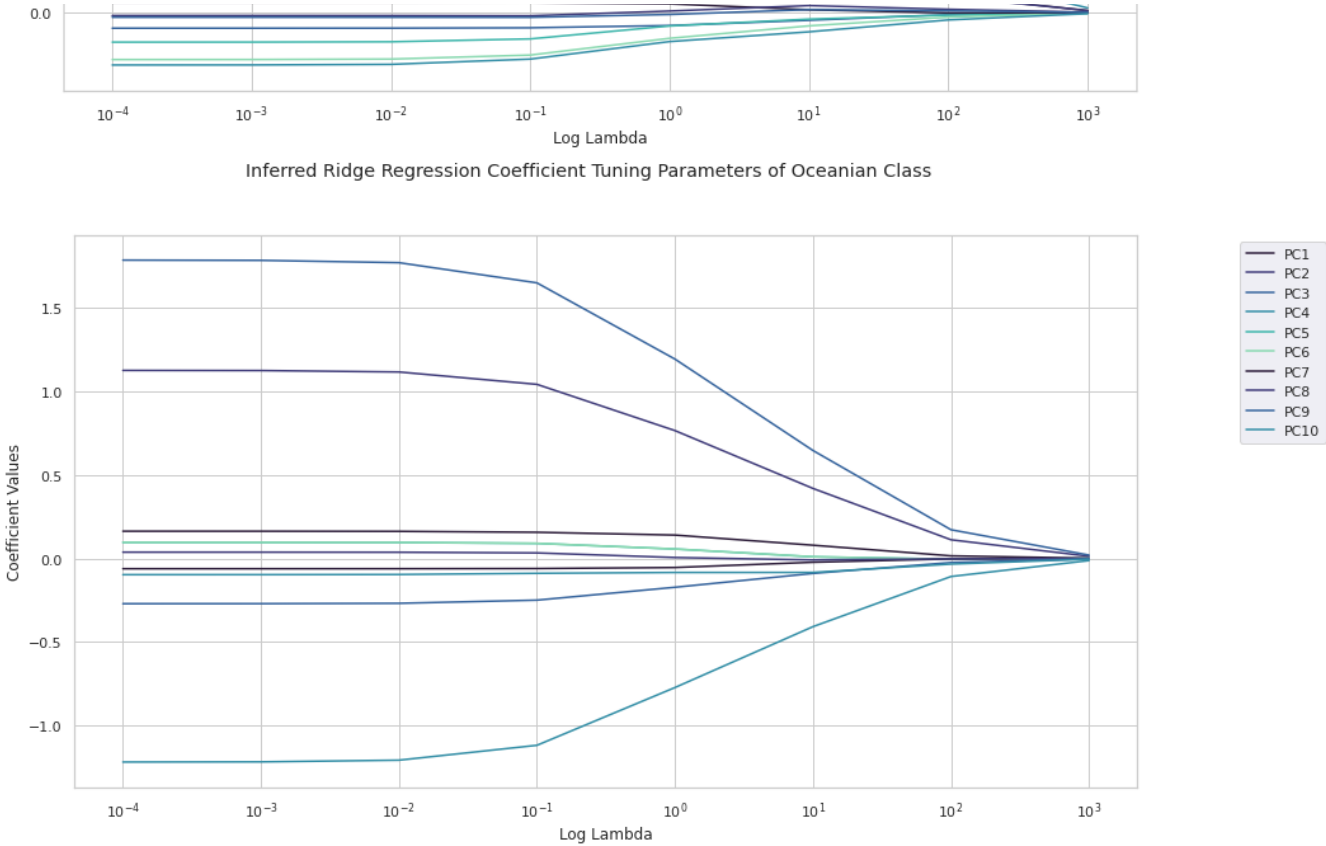Inferred Ridge Regression Coefficient Tuning Parameters of African Class



Inferred Ridge Regression Coefficient Tuning Parameters of European Class



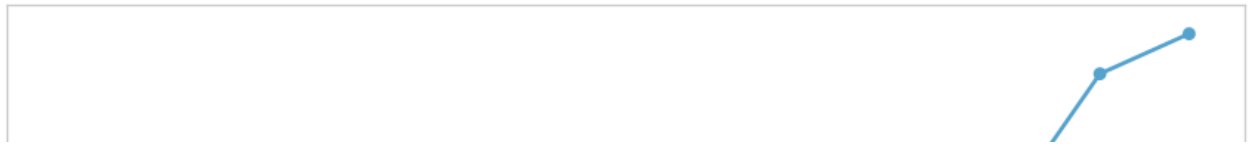Inferred Ridge Regression Coefficient Tuning Parameters of EastAsian Class

Inferred Ridge Regression Coefficient Tuning Parameters of Oceanian Class



Inferred Ridge Regression Coefficient Tuning Parameters of NativeAmerican Class

## ▾ Deliverable 2

> Illustrate the effect of the tuning parameter on the cross validation error by generating a plot with the $y$-axis as CV(5) error, and the $x$-axis the corresponding log-scaled tuning parameter value $\log10(\lambda)$ that generated the particular CV(5) error.

```
1 # Compute tuning parameter on the cross validation error
2 err = np.std(CV, 0) / np.sqrt(CV.shape[0])
3 sns.set(rc = {'figure.figsize':(15,8)})
4 sns.set_theme(style="whitegrid")
5 sns.set_palette("icefire")
6 sns.pointplot(x=λ, y=np.mean(CV, 0),yerr = err)
7 sns.set()
8 plt.xlabel('log10(lambda)')
9 plt.ylabel('CV(5) error')
10 plt.xscale('log')
11 plt.yscale('log')
12 plt.suptitle('Effect of the uning parameter on the cross validation error log1
13 plt.savefig("Assignment3_Deliverable2.png")
14 plt.show()
15
```

Effect of the uning parameter on the cross validation error log10(lambda)

## Retrain model with best lambda

```
1 # Set array of indices into the best lambda
2 best_λ = λ[np.argmin(np.mean(CV, 0))]
```

```
1 # Set standaridzed variables
2 mean_x, std_x = np.mean(x, 0), np.std(x, 0)
```

```
1 # Implement standarization of predictors and copy response variables
2 x = standardize(x, mean_x, std_x)
3 x = intercept(x)
4 y = y.copy()
```

```
1 # Set zeros matrix to coef and retiran model on batch gradient decent
2 βx = np.zeros([X2 + 1, num_features])
3 for iter in range(n_iters):
4     βx = BGD(x, y, βx, best_λ)
```

## Deliverable 3

Indicate the value of $\lambda$ value that generated the smallest CV(5) error

### Optimal lambda

```
1 # Plot lowest optimal lambda
2 palette = sns.color_palette('mako')
3 sns.set(rc = {'figure.figsize':(15,8)})
4 sns.set_theme(style="whitegrid")
5 ak = Decimal(best_λ)
6 plt.plot(λ)
7 sns.set_palette("icefire")
8 sns.countplot(data=λ)
9 plt.xscale('log')
```
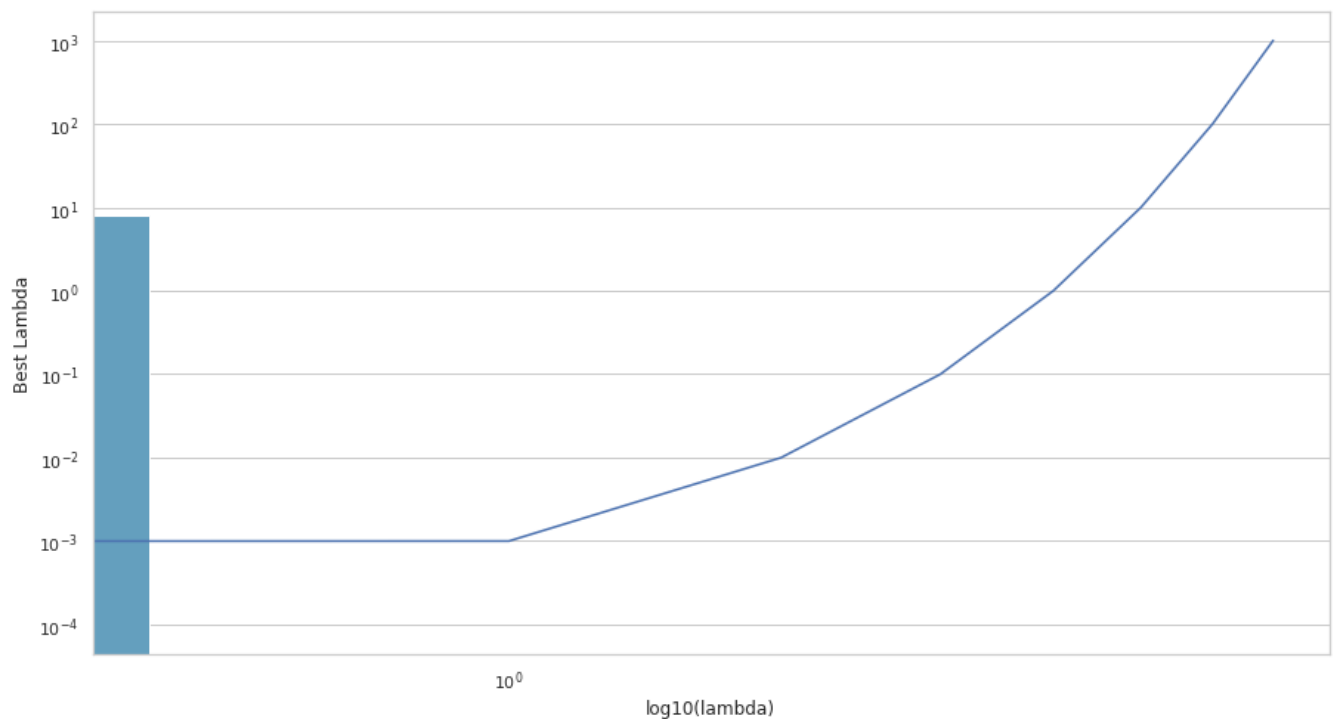
```
10 plt.yscale('log')
11 plt.xlabel('log10(lambda)')
12 plt.ylabel('Best Lambda')
13 plt.suptitle('Lowest optimal Lamda value:= log_1e{:.1f} = {}'.format(ak.log10(
14 print('Optimal lambda value:= {}'.format(best_λ))
15 plt.savefig("Assignment3_Deliverable3-1.png")
```

Optimal lambda value:= 0.0001

Lowest optimal Lamda value:= log_1e-4.0 = 0.0001
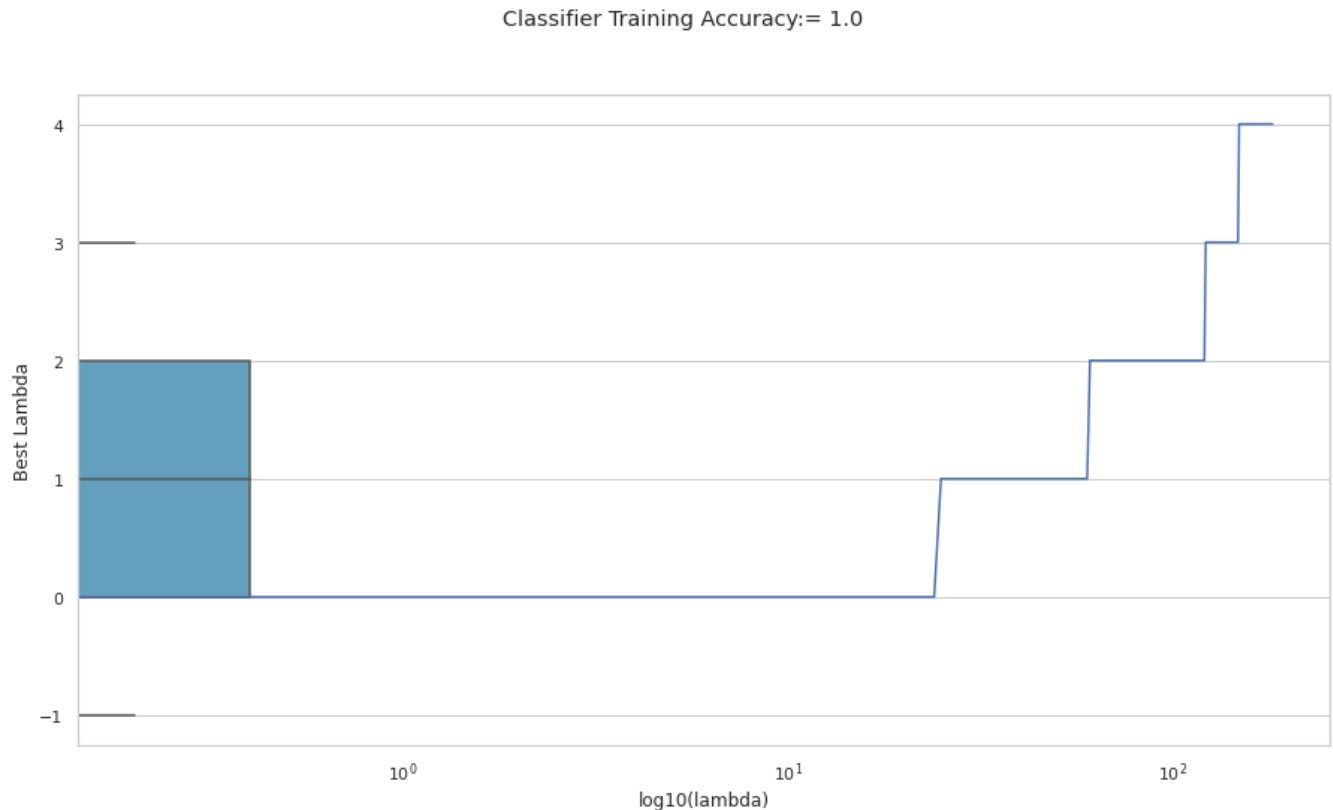


## Accuracy on training classifier

```
1
2  #Implement Prediction function
3  ŷ_p = predict(X_train)
4  # Return the maximum value along a given y axis
5  ŷ0 = np.argmax(ŷ_p, 1)
6  # Return mean traning accuaracy
7  μ = np.mean(ŷ0 == Y_train)
8  sns.set(rc = {'figure.figsize':(15,8)})
9  sns.set_theme(style="whitegrid")
10 plt.plot(ŷ0)
```

```
10    plt.plot(y0)
11    sns.set_palette("icefire")
12    sns.boxplot(data=ŷ0 -μ**2)
13    plt.xscale('log')
14    plt.xlabel('log10(lambda)')
15    plt.ylabel('Best Lambda')
16    plt.suptitle('Classifier Training Accuracy:= {}'.format(μ))
17    plt.savefig("Assignment3_Deliverable3-2.png")
18    # print('Classifier Training Accuracy: {}'.format(μ))
```



Classifier Training Accuracy:= 1.0

## ▾ Retrain model on the entire dataset for optimal $\lambda$

- Given the optimal $\lambda$, retrain your model on the entire dataset of $N$ = 183 observations to obtain an estimate of the $(p + 1) \times K$ model parameter matrix as $\mathbf{B}$ and make predictions of the probability for each of the $K$ = 5 classes for the 111 test individuals located in TestData_N111_p10.csv.
- Add probability predictions to the test dataframe

```
1 # Create new test predicotr and response variables ŷ
2 ŷ_test = predict(X_test)
3 Y_class = np.argmax(ŷ_test, 1)
```

```
1 # Re-lable feature headers and add new class prediction index column
2 new_colNames = ['{}_Probability'.format(c_name) for c_name in data] + ['ClassP
```

```
1 # Implemnt index array of probabilities
2 i_prob = np.concatenate((ŷ_test, Y_class[:, None]), 1)
```

```
1 # Create New dataframe for probality indeces
2 df2 = pd.DataFrame(i_prob, columns = new_colNames)
```

```
1 # Concat dependant Ancestory features to dataframe
2 dep_preds = pd.concat([test_df['Ancestry'], df2], axis = 1)
```

```
1 # Add new
2 dep_preds['ClassPredName'] = dep_preds['ClassPredInd'].apply(lambda x: data[in
```

```
1 # Validate Probability predictions dataframe
2 dep_preds.head()
```

| | Ancestry | African_Probability | European_Probability | EastAsian_Probability | Oceania |
|---|---|---|---|---|---|
| 0 | Unknown | 0.000348 | 0.000265 | 0.000142 | |
| 1 | Unknown | 0.000057 | 0.000600 | 0.000100 | |
| 2 | Unknown | 0.000884 | 0.992759 | 0.005403 | |
| 3 | Unknown | 0.992049 | 0.000587 | 0.000214 | |
| 4 | Unknown | 0.000066 | 0.000024 | 0.999733 | |

```
1 # Slice prediction and set new feature vector column variable
2 prob_1 = dep_preds.loc[:, 'Ancestry':'NativeAmerican_Probability']
```

```
1 # Unpivot convert dataFrame to long format
2 prob_2 = pd.melt(prob_1, id_vars = ['Ancestry'], var_name = 'Ancestry_Predicti
```

```
1 # Test for true probability
2 prob_2['Ancestry_Predictions'] = prob_2['Ancestry_Predictions'].apply(lambda x
```

```
1 # Validate dataframe
2 prob_2.head(5)
```

|   | Ancestry | Ancestry_Predictions | Probability |
|---|----------|---------------------|-------------|
| 0 | Unknown  | African_            | 0.000348    |
| 1 | Unknown  | African_            | 0.000057    |
| 2 | Unknown  | African_            | 0.000884    |
| 3 | Unknown  | African_            | 0.992049    |
| 4 | Unknown  | African_            | 0.000066    |

```
1 # Validate dataframe features
2 print('Describe Columns:=', prob_2.columns, '\n')
3 print('Data Index values:=', prob_2.index, '\n')
4 print('Describe data:=', prob_2.describe(), '\n')
```

```
Describe Columns:= Index(['Ancestry', 'Ancestry_Predictions', 'Probability'], dtype='ob

Data Index values:= RangeIndex(start=0, stop=555, step=1)

Describe data:=        Probability
count  555.000000
mean     0.200000
std      0.340312
min      0.000020
25%      0.002991
50%      0.012157
75%      0.183692
max      0.999733
```
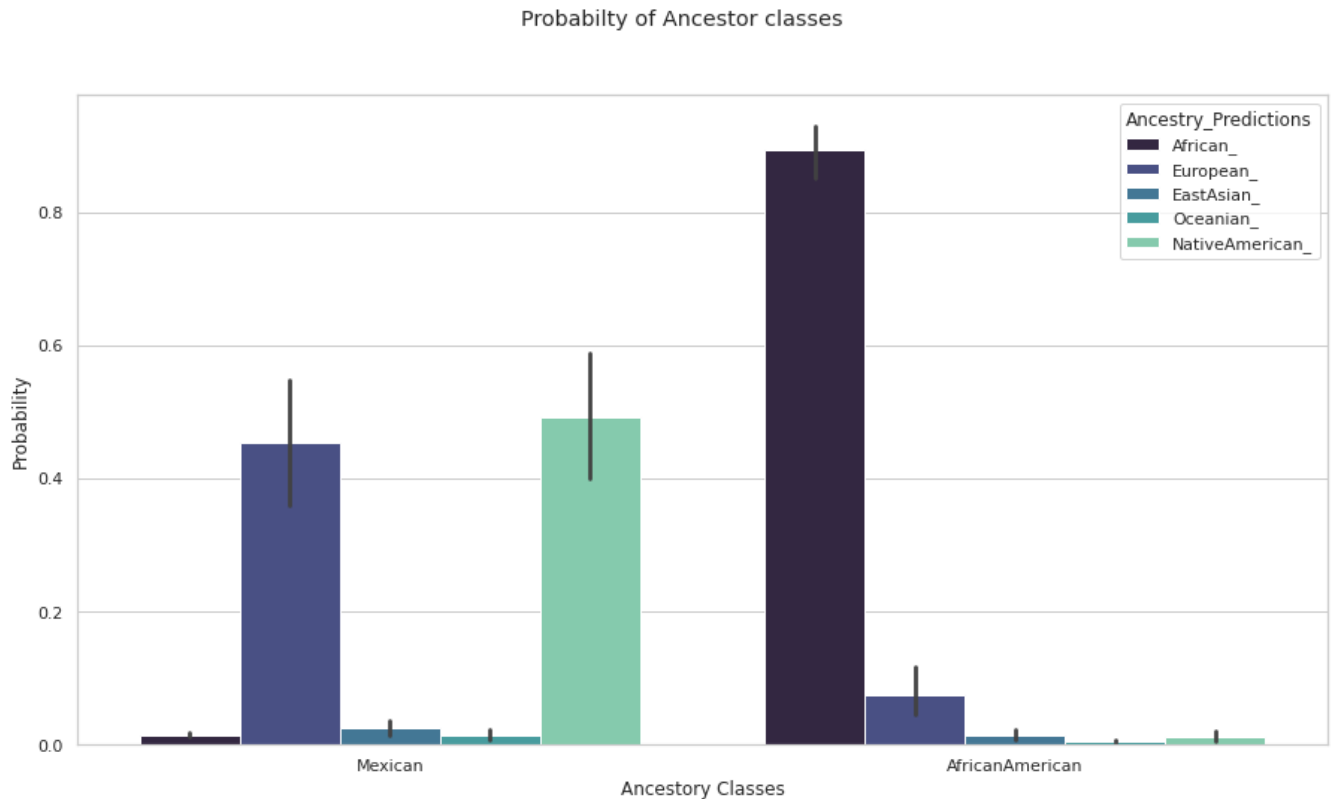
## Deliverable 4

Given the optimal $\lambda$, retrain your model on the entire dataset of $N$ = 183 observations to obtain an estimate of the $(p + 1) \times K$ model parameter matrix as $\mathbf{B}$ and make predictions of the probability for each of the $K = 5$ classes for the 111 test individuals located in TestData_N111_p10.csv.

$$\hat{Y}(X) = \underset{k \in \{1,2,\ldots,K\}}{\arg\max} \; p_k(X; \widehat{\mathbf{B}})$$

```
 1 # Plot Probality prediction matrix
 2 sns.set(rc = {'figure.figsize':(15,8)})
 3 sns.set_theme(style="whitegrid")
 4 fig, ax = plt.subplots()
 5 sns.barplot(data = prob_2[prob_2['Ancestry'] != 'Unknown'],color = 'r', x = 'A
 6 plt.xlabel('Ancestory Classes')
 7 plt.ylabel('Probability')
 8 plt.suptitle('Probabilty of Ancestor classes')
 9 plt.savefig("Assignment3_Deliverable4.png")
10 plt.show()
```



Probabilty of Ancestor classes

# Deliverable 5

**How do the class label probabilities differ for the Mexican and African American samples when compared to the class label probabilities for the unknown samples?**

> In comparison to the class label probabilities for the unknown samples, those with unknown ancestry show a probability close to or equal to one while the other classes