

Computational Neuroscience

EEE 482/582

Can Kocagil

21602218

Homework-3



Department of Electric & Electronics Engineering

Bilkent University

Ankara, Turkey

1.04.2021

TABLE OF CONTENTS

List of Figures	ii
1. Question 1	1
1.1. Part A	1
1.2. Part B	8
1.3. Part C	10
2. Question 2	12
2.1. Part A	12
2.2. Part B	15
2.3. Part C	18
2.4. Part D	19
2.5. Part E	21
3. Source Code	23
References	37

LIST OF FIGURES

1	Ridge regression model fits for different tuning parameters λ	4
2	Ridge regression model fits for different tuning parameters λ in the logarithmic sense ..	7
3	Ridge regression model weights for bootstrapped iterations and its 95% confidence interval	10
4	Ridge regression model weights for bootstrapped iterations and its 95% confidence interval	11
5	Difference of Means of Bootstrapped Populations 1 and 2.....	14
6	Difference of Means of Bootstrapped Populations 1 and 2.....	15
7	Sampling distribution of correlation between vox1 and vox2	16
8	Sampling distribution of correlation between vox1 and vox2	18
9	Sampling distribution of difference of means of building and face	20
10	Sampling distribution of difference of means of building and face	21

1. QUESTION 1

In this question, Blood-oxygen level dependent (BOLD) responses of a neural population in human visual cortex are provided in the file hw3.data2.mat that consist of a variable Y_n that represents 1000 response samples and variable X_n that represent 100 regressors that may explain the responses.

1.1. **Part A.** In part a, we are asked to use the ridge regression method to fit regularized linear models to predict noisy BOLD responses as a weighted sum of given regressors. Then, to tune the ridge parameter $\lambda \in [0, 10^{12}]$, we will perform 10-fold cross-validation. For each cross-validation fold, I will do a three-way split of the data, i.e., select a validation set of 100 contiguous samples, a testing set of 100 samples and a training set of length 800 samples. Ultimately, we will fit the each model separately to estimate model performance based proportion of explained variance R^2 .

Ridge regression is a extension of linear regression to reduce model complexity and prevent over-fitting which may result from simple linear regression so it shrinks the coefficients and it helps to reduce the model complexity and multi-collinearity in the context of statistical estimation. In ridge regression, the aim is to minimize sum of squared error plus additional penalty term/weight decay to regularize the network so that the model hopefully is not over-fitted. This constraint on the coefficient of ridge (W), penalize large values in the gradient sense and norm equation so that it is effective algorithm to prevent overfitting. This model solves a regression model where the loss function is the linear least squares function and regularization is given by L_{2-norm} . Hence, our main objective is to minimize the linear least squares with additional L_{2-norm} penalty term as follows

$$(1) \quad W^* = \min_W \frac{1}{2} \|W^T X - y\|_2^2 + \frac{\lambda}{2} \|W\|^2$$

where W^* is estimated value W such that the equation in the above is minimized, X is the n x k multivariate input matrix and y is the dependant variable to be estimated. Hence, the loss function is

$$(2) \quad L(W, y) = \frac{1}{2} \|W^T X - y\|_2^2 + \frac{\lambda}{2} \|W\|^2$$

We can get norm equation by setting the derivative of the loss function $L(W, y)$ equal to 0 w.r.t. model parameter W so that we can find the W^* such that loss function is minimized.

$$(3) \quad L(W, y) = \frac{1}{2} \|W^T X - y\|_2^2 + \frac{\lambda}{2} \|W\|^2$$

$$(4) \quad = \frac{1}{2} (XW - y)^T \cdot (XW - y) + \frac{\lambda}{2} \|W\|^2$$

$$(5) \quad = \frac{1}{2} (W^T X^T X w - 2y^T X w + y^T y) + \frac{\lambda}{2} \|W\|^2$$

To find the norm equation in the context of optimization of $L(W, y)$, we simply set $\frac{\partial L(W, y)}{\partial W} = 0$ so that resulting W will be equal to W^* .

$$(6) \quad \frac{\partial L(W, y)}{\partial W} = X^T \cdot (XW - y) + \lambda W = 0$$

After the manipulation of the W , we can find W^* as follows

$$(7) \quad W^* = (X^T X + \lambda I_k)^{-1} X^T \cdot y$$

where I_k is $k \times k$ identity matrix. We successfully derive the norm equation of the Ridge regression. One can observe that this is the extension of Ordinary Least Squares equation with extra λI_k term. This term also helps us to invert the matrix $X^T X$ since the inverse of it may not exist all the time with ease. Then, we can move on the construction of Ridge regression in the Python. Before, note that the following libraries are imported to compute necessary actions in the following parts of the questions.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import h5py
4 import scipy
5 import pandas as pd
6 import scipy.special as special
7 import random
```

Now, we can move on the construction of Ridge regression in the Python.

```
1 class RidgeRegression(object):
2     """
3         Ridge regression is a method of estimating the coefficients of
4         ↪ multiple-regression models in scenarios where independent variables are
5         ↪ highly correlated.
6         """
7     def __init__(self, Lambda:float=1):
8         """
9             Constructor method for initialization of ridge regression model.
10            Arguments:
11            - Lambda (float): is the parameter which balances the amount
12            ↪ of emphasis given to minimizing RSS vs minimizing sum of square of
13            ↪ coefficients
14            """
15     self.Lambda = Lambda
```

After the initialization of the Ridge instances, we can fit our data via calling fit function as follows

```
1 def fit(self, X:np.ndarray, y:np.ndarray) -> None:
2     """
3         Given the pair of X,y, fit the data, i.e., find parameter W such
4         ↪ that sum of square error is minimized.
5     """
6     Arguments:
7     - X (np.ndarray) : Regressor data
8     - y (np.ndarray) : Ground truths for regressors
9     Returns:
10    - None
11    """
12
```

```

13         I = np.eye(X.shape[1])
14
15         self.W = np.linalg.inv(
16             X.T.dot(X) + self.Lambda * I
17             ).dot(X.T).dot(y)
18
19         return self

```

At this point, we construct Ridge regression model for initializing & training part. In the inference part, i.e., when we predict the dependant variable, we need to compute $X \cdot W$ as a estimation of its target variable. Hence, we can compute that estimation as follows

```

1  def predict(self,X:np.ndarray) -> np.ndarray :
2      """
3          Given the test data X, we predict the target variable.
4
5          Arguments:
6              - X (np.ndarray) : The independant variable (regressor)
7
8          Returns:
9              - Y_hat (np.ndarray) : Estimated value of y
10         """
11
12         return X.dot(self.W)

```

Here, the fundamental context of Ridge regression is completed, we can move on the following parts. In this question, for all parts, the proportion of explained variance (R^2) should be calculated as the square of Pearson's correlation coefficient between measured and predicted responses. Explained variance (also called explained variation) is used to measure the discrepancy between a model and actual data. In other words, it's the part of the model's total variance that is explained by factors that are actually present and isn't due to error variance. Higher percentages of explained variance indicates a stronger strength of association. Also, we can calculate the R^2 term as squared of Pearson correlation. But the main idea is same and can be explained as follows

$$(8) \quad R^2 = 100 * \left(1 - \frac{\text{unexplained variance}}{\text{total variance}} \right)$$

Let's compute R^2 in the Python as follows

```

1  def eval_r2(self,y_true:np.ndarray, y_pred:np.ndarray) -> np.float:
2      """
3          Given the true dependant variable and estimated variable, computes
4          ↪ proportion of explained variance R^2 by square the Pearson correlation
5          ↪ between true dependant variable and estimated variabl
6
7          Arguments:
8              - y_true (np.ndarray) : true dependant variable
9              - y_pred (np.ndarray) : estimated variable
10
11          Returns:
12              - r_squared (np.float) : Proportion of explained variance
13         """
14
15         _pearson = np.corrcoef(y_true,y_pred)

```

```

14     pearson = _pearson[1][0]
15     r_squared = np.square(pearson)
16     return r_squared

```

Then, we perform 10-fold cross-validation to tune the ridge parameter ($\lambda \in [0, 10^{12}]$) based on model performance. From now on, I fit a separate model for each λ using the training then I find R^2 of each model on the testing set. After that, I separately estimate R^2 of each model on the validation set with the plot of the average R^2 across cross-validation folds as a function of λ . After, I find the optimal ridge parameter λ_{opt} that maximizes average R^2 based on the validation performance. As a final step, I will find the model performance by calculating the average R^2 across cross-validation folds, measured on the validation set for λ_{opt} with corresponding plots. So before the 10-fold cross validation, let's start with the analysis of Ridge regression with varying λ parameter as a first intuition of the problem.

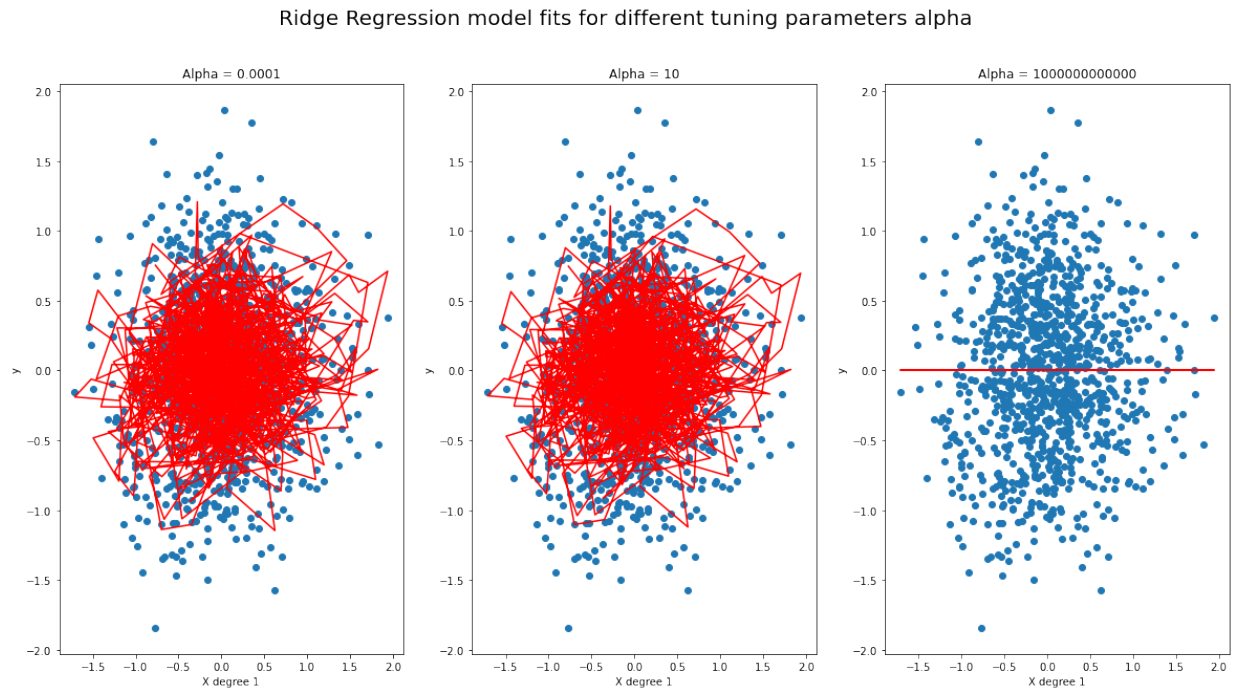


FIGURE 1. Ridge regression model fits for different tuning parameters λ

Then, we can move on the actual parameter and performance estimation part via cross-validation technique. Cross-validation is a technique to evaluate predictive models by partitioning the original sample into a training set to train the model, and a test/val set to evaluate it.

In k-fold cross-validation, the original sample is randomly partitioned into k equal size subsamples. Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k-1 subsamples are used as training data. The cross-validation process is then repeated k times (the folds), with each of the k subsamples used exactly once as the validation data. The k results from the folds can then be averaged (or otherwise combined) to produce a single estimation. The advantage of this method is that all observations are used for both training and validation, and each observation is used for validation exactly once. But note that in our case, we are asked to do a three-way split of the data, i.e., select a validation set of 100 contiguous samples, a testing set of 100 samples and a training set of length 800 samples.

To accomplish that, the following code snippet performs K-fold cross validation as a three-way split.

```

1 class K_fold(object):
2     """
3     Cross-validation, sometimes called rotation estimation or out-of-sample
    ↪ testing, is any of various similar model validation techniques for assessing
    ↪ how the results of a statistical analysis will generalize to an independent
    ↪ data set
4     """
5     def __init__(self, sample_size: int = y.shape[0], folds: int = 10):
6         """
7         Constructor method for initializing the sample size and the number
    ↪ of folds
8
9         Arguments:
10            - sample_size (int) : How many samples are in the dataset
11            - folds (int) : the number of folds
12        """
13
14        self.sample_size = sample_size
15        self.folds = folds
16        self.fold_size = int(sample_size / folds)
17
18    def split(self) -> tuple:
19        """
20        Generator function for splitting data as validation (10%), testing
    ↪ (10%) and training (80%) as K-fold cross validation based resampling
21        """
22
23        for idx in range(self.folds):
24            _val_idx = idx * self.fold_size
25            _test_idx = (idx + 1) * self.fold_size
26            _train_idx = (idx + 2) * self.fold_size
27
28            val_idx = np.arange(_val_idx, _test_idx) % self.sample_size
29            test_idx = np.arange(_test_idx, _train_idx) % self.sample_size
30            train_idx = np.arange(_train_idx, self.sample_size + _val_idx) %
    ↪ self.sample_size
31
32            yield val_idx, test_idx, train_idx

```

Then, let's instantiate the K-fold class and perform 10-fold cross validation as follows.

```

1 dict_inference = {
2     'test' : dict(),
3     'val'  : dict()
4 }
5
6
7 phases = [
8     'train',
9     'val',
10    'test'

```



```

11
12 ]
13
14 log_lambda_arr = np.logspace(
15     start = 0,
16     stop  = 12,
17     num   = 500,
18     base  = 10
19 )
20
21 cv = K_fold(folds = 10)
22
23 for val_idx, test_idx, train_idx in cv.split():
24
25     X_list = [
26         X[train_idx],
27         X[val_idx],
28         X[test_idx]
29     ]
30
31     y_list = [
32         y[train_idx],
33         y[val_idx],
34         y[test_idx]
35     ]
36
37
38     for _lambda in log_lambda_arr:
39
40         for phase, X_phase, y_phase in zip(phases, X_list, y_list):
41             if phase == 'train':
42                 model = RidgeRegression(_lambda)
43                 model.fit(X_phase, y_phase)
44
45             else:
46                 preds = model.predict(X_phase)
47                 r2_score = model.eval_r2(y_phase, preds)
48                 dict_inference[phase].setdefault(
49                     _lambda, list()).append(r2_score)
50
51 inference_r2 = {
52     phase : {
53         _lambda : np.mean(r2_score) for _lambda, r2_score in
54         ↪ dict_inference[phase].items()
55     }
56     for phase in ['val', 'test']
57 }

```

At this point, I successfully estimate the R^2 on both validation and test. Additionally, I found the λ_{opt} as follows.

```

1 best_r2 = 0
2 for _lambda, r_2 in inference_r2['val'].items():
3     if r_2 > best_r2:
4         best_r2 = r_2
5         best_lambda = _lambda
6
7
8 print(f'Best lambda parameter that maximizes the R^2 is : {best_lambda}')
9 print('Best R^2 along the testing :', inference_r2['test'][best_lambda])
10 print('Best R^2 along the validation :', inference_r2['val'][best_lambda])

```

Best lambda parameter that maximizes the R^2 is : 395.5436244734702

Best R^2 along the testing : 0.16042061044928463

Best R^2 along the validation : 0.15259887784859996

Hence, we found $\lambda_{opt} = 395.543$. Let's visualize the inference progress through the cross-validation as a better intuition of how our ridge model estimate target variable w.r.t. varying parameter λ .

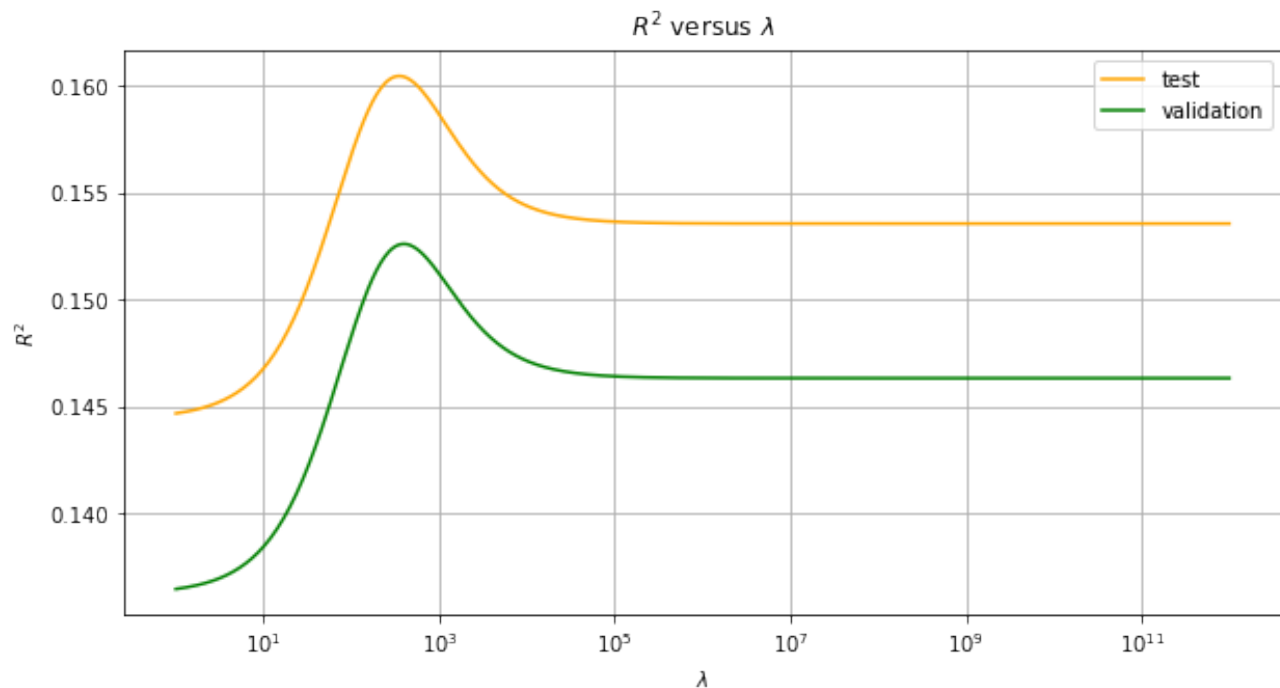


FIGURE 2. Ridge regression model fits for different tuning parameters λ in the logarithmic sense

We can see from the figure that, there is a optimal λ , λ_{opt} that maximizes the explained variance in the validation as well as testing phase. But, since validation data is used for estimating the model parameters, we choose the optimal parameter of the ridge model λ_{opt} according to validation samples.

1.2. **Part B.** In this part of the question, we will determine confidence intervals for parameters of the OLS model from part a. OLS model corresponds to $\lambda = 0$ case in ridge regression so OLS model minimize the sum of squared error in the sense of statistical estimation. To do that, I generate bootstrap samples from the 1000 samples in the original data and perform 500 bootstrap iterations, and refit a separate model at each iteration.

Bootstrapping is any test or metric that uses random sampling with replacement (e.g. mimicking the sampling process), and falls under the broader class of resampling methods [1]. Bootstrapping assigns measures of accuracy (bias, variance, confidence intervals, prediction error, etc.) to sample estimates [1]. This technique allows estimation of the sampling distribution of almost any statistic using random sampling methods [1]

In our case, we generate bootstrap samples to estimate the 95% confidence interval and to find statistically significant weights. So, let's see the Python code for bootstrapping. But before that, for reproducibility purposes, I create random seed operation as follows:

```

1 def random_seed(seed:int = 42) -> None :
2     """ Random seeding for reproducibility
3
4         Arguments:
5             - seed (int) : random state
6
7         Returns:
8             - None
9     """
10    np.random.seed(seed)
11    random.seed(seed)

```

Then, here is code for bootstrapping.

```

1 random_seed(10)
2
3 bootstrap_iters = range(500)
4 sample_idx = np.arange(X.shape[0])
5 parameters = list()
6
7 for idx in bootstrap_iters:
8
9     bootstrap_idx = np.random.choice(sample_idx, size = 1000, replace = True)
10    y_bootstrap = y[bootstrap_idx]
11    X_bootstrap = X[bootstrap_idx]
12    ridge = RidgeRegression(Lambda = 0)
13    ridge.fit(X_bootstrap,y_bootstrap)
14    parameters.append(ridge.parameters())
15
16 w_bootstrap = np.array(parameters)
17 w_mean = np.mean(w_bootstrap, axis=0)
18 w_std = np.std(w_bootstrap, axis=0)

```

So, what we done here is that we generate random indexes with replacement, then sample from our actual data, fit OLS model. After the completion of bootstrapping, we estimate model reliability by averaging the means of each bootstrap iteration. Then, we also calculate the standard deviation of sampling distribution of model parameter to compute test statistics in the following parts.

Then, let's move on the statistical test, confidence intervals and hypothesis testing part. Since, this assignment focused on the concept of hypothesis testing, I will briefly describe the concept for future use and easy understanding by focusing on the concept not equations since every statistical test requires similar but different derivations to calculate confidence interval, etc.

A confidence interval is a range of values that is likely to contain an unknown population parameter. If you draw a random sample many times, a certain percentage of the confidence intervals will contain the population mean. This percentage is the confidence level. In our case, we will use confidence interval for OLS regression coefficients.

The relationship between p-values and confidence intervals can be interpreted as follows. In the context of hypothesis testing, one can use either p-values or confidence intervals to determine whether your results are statistically significant. If a hypothesis test produces both, these results will agree. As summary of relationship between p-value and confidence interval. We set type-I probability error term α and confidence level is equivalent to $1 - \alpha$ level. Hence, if we predetermine significance level, for example 0.05, the corresponding confidence level yields 95%. Then, we can end up with

- $p\text{-value} < \alpha$ then the hypothesis test is statistically significant,
- If the confidence interval does not contain the null hypothesis value, the results are statistically significant
- $p\text{-value} < \alpha$ the confidence interval will not contain the null hypothesis value

Hence, we can interpret p-value as a significance level so that calculated value of a test statistic is borderline between rejection and acceptance region. Smaller p-values imply the rejection of null assumptions so they indicates research hypothesis is statistically significant. In other words, smaller p-values indicates strong evidence against the null hypothesis and also we can say that we get strong evidence for behoof of alternative hypothesis. Even the significance level changes study to study, researcher generally set 0.05 as the significance level so $p\text{-value} < 0.05$ is small enough to say that alternative hypothesis is statistically significant.

Finally, we can move on our case. Then, here is code for plotting its mean and 95% confidence interval in the sense of error plot.

```
1 plt.figure(figsize = (10,5))
2 plt.errorbar(np.arange(1, 101),
3             w_mean,
4             yerr= w_std,
5             ecolor='red',
6             elinewidth=1,
7             capsize=1)
8 plt.title('Ridge Model OLS Weights')
9 plt.xlabel('i')
10 plt.ylabel('$W_i$')
11 plt.show()
```

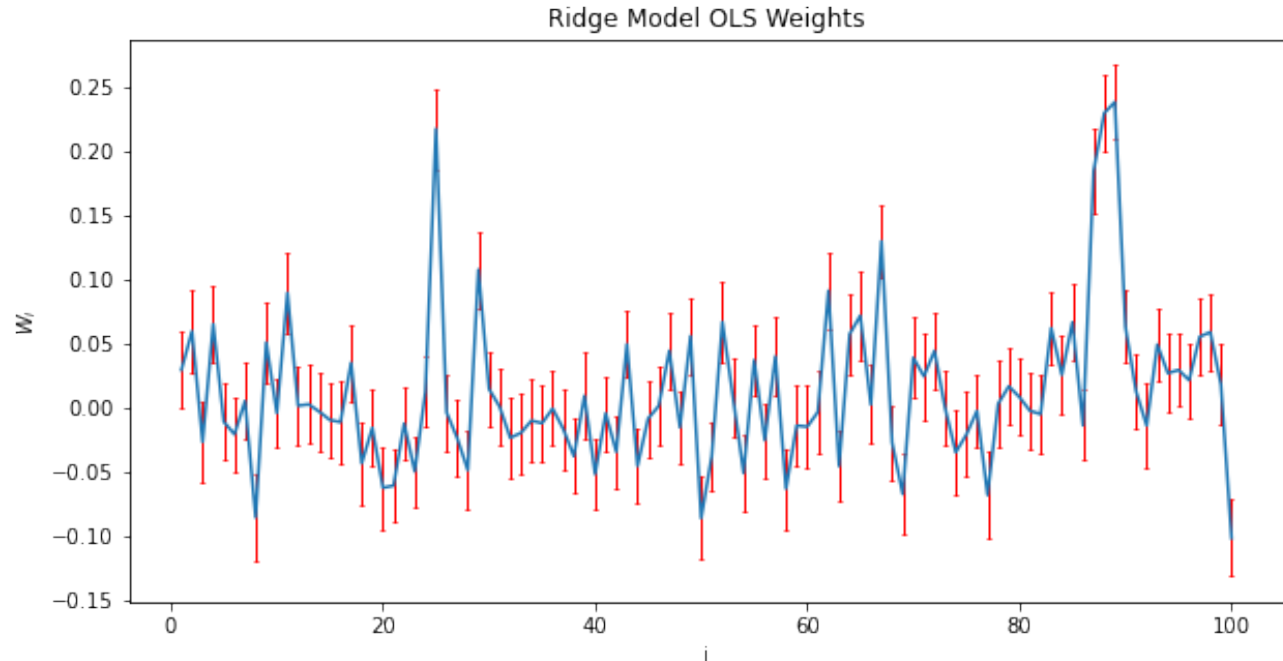


FIGURE 3. Ridge regression model weights for bootstrapped iterations and its 95% confidence interval

Hence, to calculate p-value, we need to standardize our sampling distribution. So, here is the standardization and calculation of p-value with the indexes of statistically significant weights different than 0.

```
1 two_sided = 2
2 p_values = special.ndtr(- w_mean / w_std) * two_sided
3 alpha_level = 0.05
4 significant = np.argwhere(p_values < alpha_level).flatten()
5 print(f' Index of the parameters that are significantly different than 0: \n
   ↪ {significant}')
```

Index of the parameters that are significantly different than 0:

3 10 24 28 51 61 64 66 82 84 86 87 88 89 97

1.3. Part C. In this part of the question, we will determine confidence intervals for parameters of the regularized linear model from part a, i.e., the model obtained for λ_{opt} . To do that, as we done in the previous part, I generate bootstrap samples from the 1000 samples in the original data and perform 500 bootstrap iterations, and refit a separate model at each iteration using λ_{opt} . Then, I will plot the mean and 95% confidence intervals of the parameters in the same graph with providing the model regressors which have weights that are significantly, different than 0.

Here is the Python code for generating bootstrapped samples.

```
1 random_seed(10)
2
3 bootstrap_iters = range(500)
4 sample_idx = np.arange(X.shape[0])
5 parameters = list()
6
```

```

7 for idx in bootstrap_iters:
8
9     bootstrap_idx = np.random.choice(sample_idx, size = 1000, replace = True)
10    y_bootstrap = y[bootstrap_idx]
11    X_bootstrap = X[bootstrap_idx]
12    ridge = RidgeRegression(Lambda = best_lambda)
13    ridge.fit(X_bootstrap,y_bootstrap)
14    parameters.append(ridge.parameters())
15
16 w_bootstrap = np.array(parameters)
17 w_mean = np.mean(w_bootstrap, axis=0)
18 w_std = np.std(w_bootstrap, axis=0)

```

Then, here is code for plotting its mean and 95% confidence interval in the sense that error plot.

```

1 plt.figure(figsize = (10,5))
2 plt.errorbar(np.arange(1, 101),
3             w_mean,
4             yerr= w_std,
5             ecolor='red',
6             elinewidth=1,
7             capsize=1)
8 plt.title('Ridge Model  $\lambda_{\text{optimal}}$  Weights')
9 plt.xlabel('i')
10 plt.ylabel('$W_i$')
11 plt.show()

```

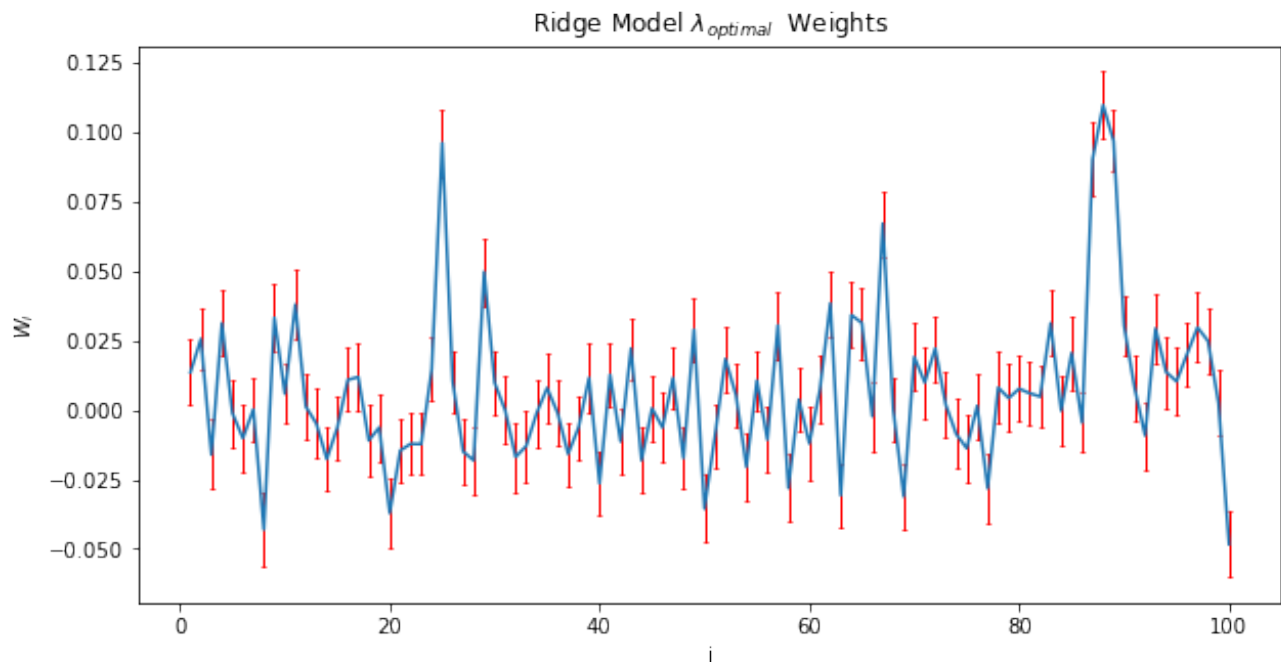


FIGURE 4. Ridge regression model weights for bootstrapped iterations and its 95% confidence interval

Hence, to calculate p-value, we need to standardize our sampling distribution. So, here is the standardization and calculation of p-value with the indexes of statistically significant weights different than 0.

```
1 p_values = scipy.special.ndtr(- w_mean / w_std) * two_sided
2 significants = np.argwhere(p_values < alpha_level).flatten()
3 print(f' Index of the parameters that are significantly different than 0: \n
   ↪ {significants}')
```

Index of the parameters that are significantly different than 0:

1 3 8 10 24 28 42 48 56 61 63 64 66 82 86 87 88 89 92 96 97

Hence, we can see that with the optimal Ridge model the number of indexes of the model regressors which have weights that are significantly different than 0 (at a significance level of $p < 0.05$) is increased when compared to OLS. The reason behind this, we added penalty L_{2-norm} term to the OLS estimator which stricts the parameters within smaller numbers w.r.t. OLS case with the increasing number of bootstrapped samples. So, as we truly estimate the coefficient of the model, the number of parameters that are significantly different than 0 (with alpha level = 0.05) is increased.

2. QUESTION 2

In this question, a series of neural response measurements are provided in the file hw3_data3.mat.

2.1. Part A. In this part, responses from two separate populations of neurons are stored in the variables pop1 and pop2. We would like to examine whether the mean responses of the two populations are significantly different. In other words, we examine if they are drawn from the same distribution or not with alpha level = 0.05. As a prior information to the researcher, the first population contains 7 neurons, whereas the second population contains 5 neurons.

To test our null hypothesis that states that two datasets follow the same distribution versus the alternative that indicates the just the opposite of null hypothesis (two-tailed) with alpha level = 0.05. Let's construct the hypothesis testing formulation as follows. Let μ_1 and μ_2 represents the means of pop1 and pop2 respectively such that

$$\text{Null hypothesis } H_o : \mu_1 - \mu_2 = 0 \text{ and researcher hypothesis } H_a : \mu_1 - \mu_2 \neq 0$$

So, we priory believe that our responses are coming from the same distribution so that there is no such underlying mean difference in the sense of statistical estimation. Now, we will test our hypothesis versus alternative one as follows.

In the question, there is a hint that says if the two datasets come from a common distribution, is there any need to separate them. So, we priory assume that there is no separation between the the pop1 and pop2 so we can combine in a single population as first step of hypothesis testing. Then, we need to bootstrap the data with 10 000 iteration to estimate its sampling distribution and characterization of its test-statistics. Here is the Python code for bootstrapping.

```

1 def bootstrap(sample:np.ndarray, bootstrap_iters:iter = range(10000),
  ↪ random_state:int = 11) -> np.ndarray:
2     """
3
4     Generate bootstrap samples using random sampling with replacement.
5
6     Arguments:
7         - sample (np.ndarray) : Sample to be bootstrapped
8         - bootstrap_iters (iterator object) : Specification of bootstrap
  ↪ iterations
9         - random_state (int) : Random seed for reproducibility
10
11     Returns:
12         - bootstrap_samples (np.ndarray) : Bootstrapped array
13
14     """
15     random_seed(random_state)
16     size = sample.shape[0]
17     bootstrap_samples = list()
18
19     for idx in bootstrap_iters:
20         bootstrap_idx = np.random.choice(np.arange(sample.shape[0]), size =
  ↪ size, replace = True)
21         bootstrap_samples.append(sample[bootstrap_idx])
22
23     return np.array(bootstrap_sample

```

Then, I as mentioned, I aggregated the pop1 and pop2 in column wise. Then, I bootstrapped from that aggregated samples. Next step is select 7 values from sample and represent the values as first populations responses. In the same manner, we interpret the remaining 5 values of each sample as second populations responses. Since, as a priory, we believe that they come from the same underlying distribution so there is no discrepancy to claim in this manner. Hence, we compute difference of means by subtracting the representation of first population with bootstrapped samples from second one. Then, here is the code for computation of difference of means and its visualization.

```

1 pop = np.vstack([pop1,pop2])
2 pop_bootstrap = bootstrap(pop)
3 sample_1 = pop_bootstrap[:,len(pop1)].squeeze(2)
4 sample_2 = pop_bootstrap[:,len(pop1):].squeeze(2)
5 sample_1_bootstrap_mean = sample_1.mean(axis = 1)
6 sample_2_bootstrap_mean = sample_2.mean(axis = 1)
7 sample_diff_means = sample_1_bootstrap_mean - sample_2_bootstrap_mean
8 sample_mean_dist = pd.DataFrame()
9 sample_mean_dist['Mean Difference'] = sample_diff_means.flatten()
10 fig, ax = plt.subplots(figsize = (10,5))
11 sample_mean_dist.plot.kde(ax=ax, title='Difference of Means of Bootstrapped
  ↪ Populations 1 and 2')
12 sample_mean_dist.plot.hist(density=True, ax = ax, bins = 15)
13 ax.set_ylabel('Probability $P_X(x)$')
14 ax.set_xlabel('Difference in means (x)')
15 ax.grid(axis='y')
16 ax.set_yticks([])

```

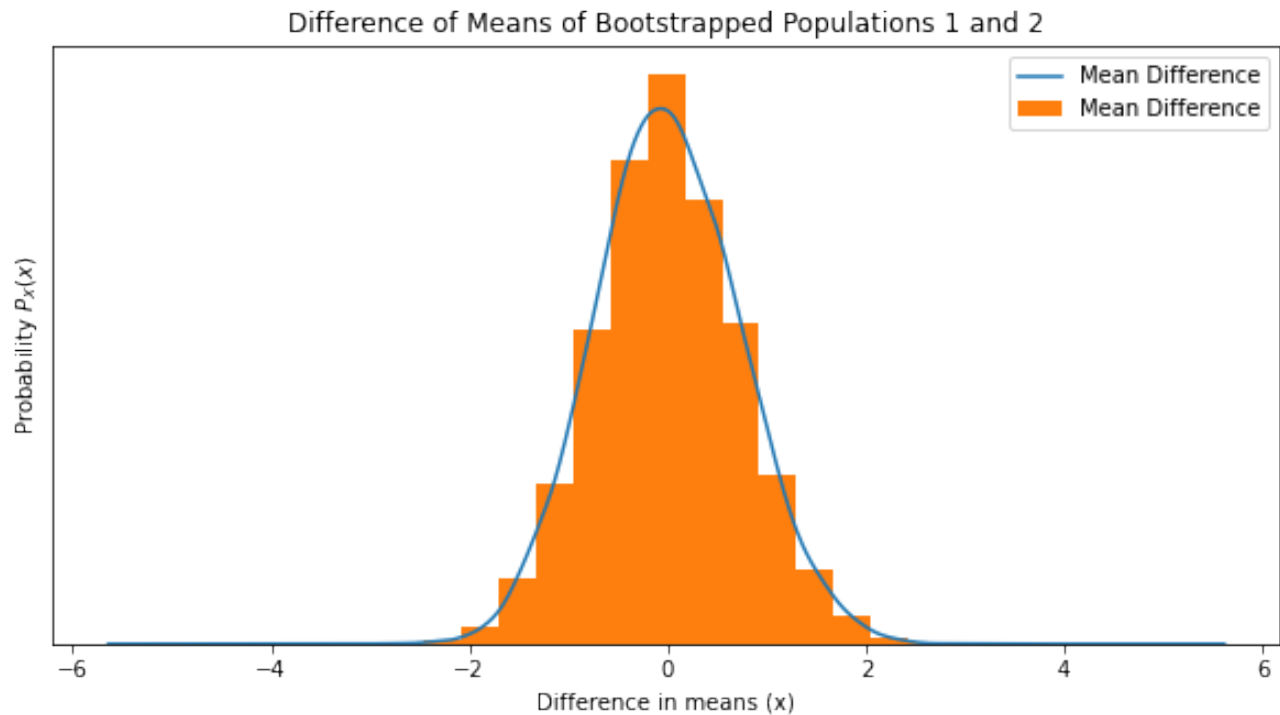



FIGURE 5. Difference of Means of Bootstrapped Populations 1 and 2

Here, the figure represents our sampling distribution of difference of means. As we can see, it is almost zero centered Gaussian. Let's analyze the sampling distribution of the pop1 and pop2 separately to compare above figure as it is very intuitive way of understanding of the problem. Let's see the corresponding code and the figures.

```

1 pop1_bootstrap = bootstrap(pop1)
2 pop2_bootstrap = bootstrap(pop2)
3 pop1_bootstrap_mean = np.mean(pop1_bootstrap, axis = 1)
4 pop2_bootstrap_mean = np.mean(pop2_bootstrap, axis = 1)
5
6 mean_dist = pd.DataFrame()
7 mean_dist['pop1 Mean'] = pop1_bootstrap_mean.flatten()
8 mean_dist['pop2 Mean'] = pop2_bootstrap_mean.flatten()
9 mean_dist['Mean Difference'] = pop1_bootstrap_mean - pop2_bootstrap_mean
10
11 fig, ax = plt.subplots(figsize = (10,5))
12 mean_dist.plot.kde(ax=ax, title='Difference of Means of Bootstrapped Populations
   ↪ 1 and 2')
13 mean_dist.plot.hist(density=True, ax = ax, bins = 15)
14 ax.set_ylabel('Probability $P_X(x)$')
15 ax.set_xlabel('Difference in means (x)')
16 ax.grid(axis='y')
17 ax.set_yticks([])

```

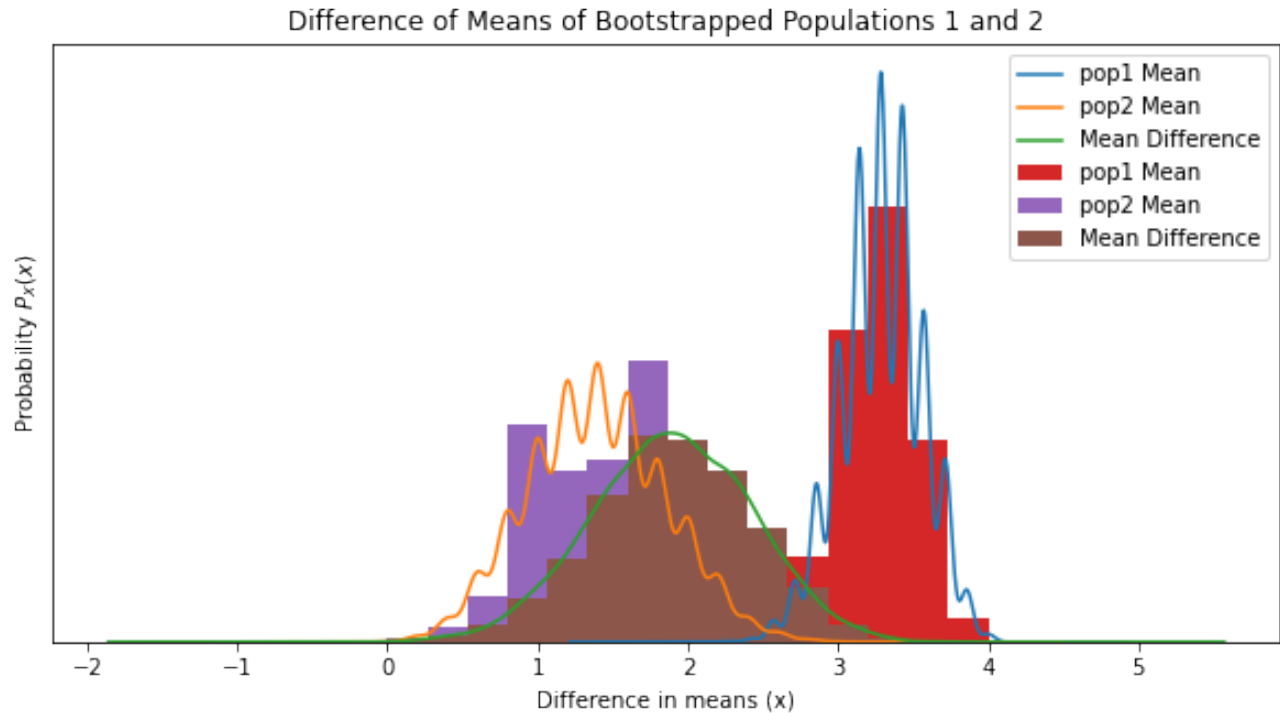


FIGURE 6. Difference of Means of Bootstrapped Populations 1 and 2

Hence, before the calculation of the p-value and test statistic, we get the intuition that the actual bootstrapped difference of means centered around 2 whereas the null hypothesis statements centered around 0.

Then, let's calculate test statistic and p-value to test our hypothesis as follows.

```

1 actual_diff_means = pop1.mean() - pop2.mean()
2 std_test = sample_mean_dist['Mean Difference'].std()
3 mean_test = sample_mean_dist['Mean Difference'].mean()
4
5 z_cal = (mean_test - actual_diff_means) / std_test
6 p_values = scipy.special.ndtr(z_cal) * two_sided
7
8 print('The two sided p-value is:', p_values)

```

The two sided p-value is: 0.01003287406404004

As we can see, p-value < alpha level that indicates the rejection of our prior believes about the population, i.e., null hypothesis. Hence, researcher find a strong evidence to indicate the rejection of prior believe that pop1 and pop2 are drawn from the same distribution. Hence, researcher hypothesis is statistically significant.

2.2. Part B. In this part, BOLD responses recorded in two voxels in the human brain are stored in the variables vox1 and vox2 and we would like to examine whether the voxel responses are similar to each other, by calculating their correlation. As done in the previous parts, I will find the mean and 95% confidence interval by bootstrapping 10 000 iterations. Then, I will find the percentile of the bootstrap distribution, corresponding to a correlation value of 0.

Here, we bootstrapped the vox1 and vox2 separately as our prior believe they are not correlated in the sense of Pearson. Hence, our null hypothesis implies that the distributions vox1 and vox2 are not correlated.

In the context of hypothesis testing, researcher hypothesis implies the correlation between vox1 and vox2. Let's look at the sampling distribution of the correlation between vox1 and vox2, then calculate test statistic to test our prior believe under the null assumption as follows. Note that since the null hypothesis is not given directly, either ways of computing the end result should agree.

```

1 def corr(X:np.ndarray,Y:np.ndarray) -> list:
2     """
3         Given the X,Y distributions, computes the correlation element wise.
4
5         Arguments:
6             - X (list or np.ndarray) : First distribution
7             - Y (list or np.ndarray) : Second distribution
8
9         Returns:
10            - pearson_corrs (list[float]) : Correlations element wise
11    """
12    return [scipy.stats.pearsonr(X[i], Y[i])[0] for i in range(X.shape[0])]
13
14 vox1_bootstrap = bootstrap(vox1)
15 vox2_bootstrap = bootstrap(vox2)
16 corr_bootstrap = corr(vox1_bootstrap,vox2_bootstrap)
17 fig, ax = plt.subplots(figsize = (10,5))
18 pd.Series(corr_bootstrap).plot.kde(ax=ax, legend = False, title='Sampling
19 ↪ Distribution of Correlation between vox1 and vox2')
20 pd.Series(corr_bootstrap).plot.hist(density=True, ax = ax, bins = 20, alpha =
21 ↪ 0.8,color = 'red')
22 ax.set_ylabel('Probability $P_Y(y)$')
23 ax.set_xlabel('Pearson Correlation y')
24 ax.grid(axis='y')
25 ax.set_yticks([])

```

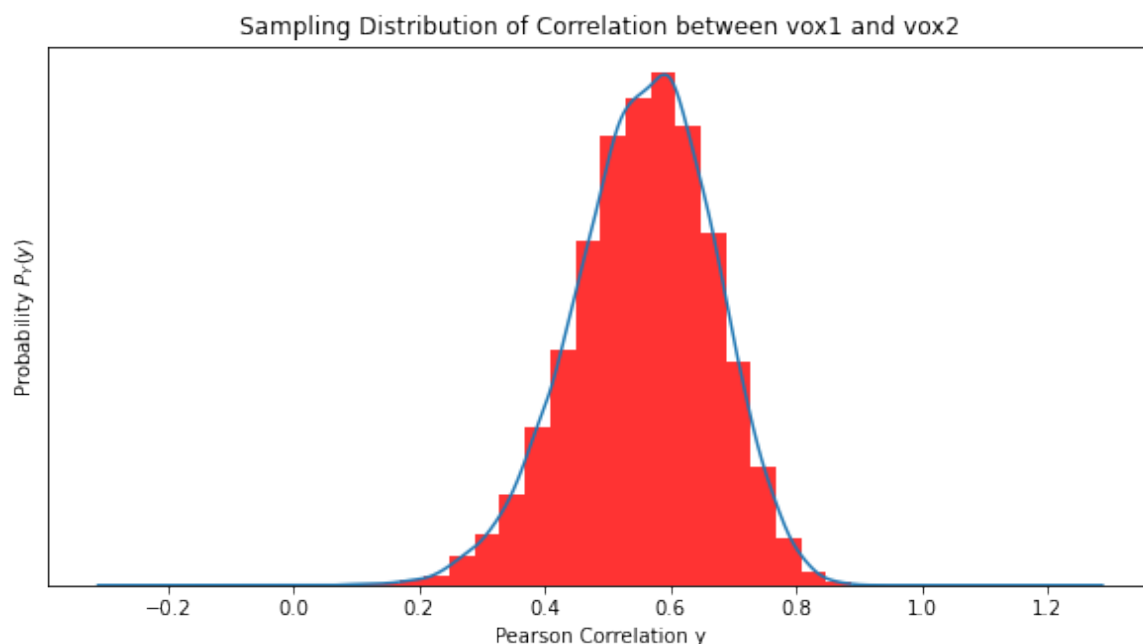


FIGURE 7. Sampling distribution of correlation between vox1 and vox2

As we can see from the figure, sampling distribution of correlation between vox1 and vox2 have correlation centered around 0.55 as a first intuition of our test progress. It seems we will find enough evidence against our null assumptions such that rejection of null hypothesis is likely. Then, let's calculate mean and 95% confidence interval as follows.

```

1  # Thanks to
2  # https://stackoverflow.com/questions/15033511/
3  # compute-a-confidence-interval-from-sample-data
4  def confidence_interval(data:np.ndarray, confidence:float=0.95) -> tuple:
5      """
6          Given the distribution and confidence level, computes the confidence
        ↪ interval.
7
8          Arguments:
9              - data (list or np.ndarray) : Input distribution
10             - confidence (float) : confidence level in the range [0,1]
11
12          Returns:
13              - confidence_level (tuple[np.ndarray]) : lower, upper limits
        ↪ respectively
14          """
15      a = 1.0 * np.array(data)
16      n = len(a)
17      m, se = np.mean(a), scipy.stats.sem(a)
18      h = se * scipy.stats.t.ppf((1 + confidence) / 2., n-1)
19      return m-h, m+h
20
21 corr_mean = np.mean(corr_bootstrap)
22 lower, upper = confidence_interval(corr_bootstrap,confidence=0.95)
23 print('Mean correlation value:', corr_mean)
24 print(f'95% confidence interval of the correlation values: {lower, upper}')

```

Mean correlation value: 0.5571443738567518

95% confidence interval of the correlation values: (0.5549693272197517, 0.5593194204937518)

As we previously analyze the distribution of sampling correlation and conclude that expected correlation is around 0.55 and our actual mean correlation of sampling distribution value is found as 0.557.

After that, we are asked to find the percentile of the bootstrap distribution, corresponding to a correlation value of 0. But, since we already got enough evidence against our null assumptions and analyzed the bootstrapped and mean correlation and conclude that vox1 and vox2 is correlated and this correlation is centered around 0.55, we expect little or no percentile of the bootstrap distribution, corresponding to a correlation value of 0. Let's find direct percentage through the Python as follows.

```

1  is_corr_zero = np.argwhere(corr_bootstrap == 0)
2  corr_zero_percentage = 100 * is_corr_zero.shape[0] / 10000
3  print('Percentage of zero correlation values:', corr_zero_percentage)

```

Percentage of zero correlation values: 0.0

As our intuition, and direct calculation of the percentile of the bootstrap distribution, corresponding to a correlation value of 0 indicates that its value is zero.

2.3. **Part C.** As the question implies estimation of confidence intervals and hypothesis testing are dual problems, we already have good enough intuition about this question and the evidence against the null hypothesis. But, for the sake of the question indicates the computation of break apart correlation of bootstrapped samples, let's visualize the sampling distribution and compute p-value for two voxels having zero or negative correlation. Let's see the code and visualization for part C in one shot.

```

1 vox1_ind = bootstrap(vox1, range(10000), random_state=42)
2 vox2_ind = bootstrap(vox2, range(10000), random_state=21)
3
4 _corr_ind = corr(vox1_ind,vox2_ind)
5 corr_ind = pd.Series(_corr_ind)
6
7 fig, ax = plt.subplots(figsize = (10,5))
8 corr_ind.plot.kde(ax=ax, legend = False, title='Sampling Distribution of
  ↳ Correlation between vox1 and vox2')
9 corr_ind.plot.hist(density=True, ax = ax, bins = 20, alpha = 0.8,color = 'red')
10 ax.set_ylabel('Probability $P_Y(y)$')
11 ax.set_xlabel('Pearson Correlation y')
12 ax.grid(axis='y')
13 ax.set_yticks([])
14
15 actual_corr, _ = scipy.stats.pearsonr(vox1,vox2)
16 mean_corr = corr_ind.mean()
17 std_corr = corr_ind.std()
18 z_score = mean_corr - actual_corr
19 z_score /= std_corr
20 p_value = scipy.special.ndtr(z_score)
21 print('The one sided p-value is:', p_value)

```

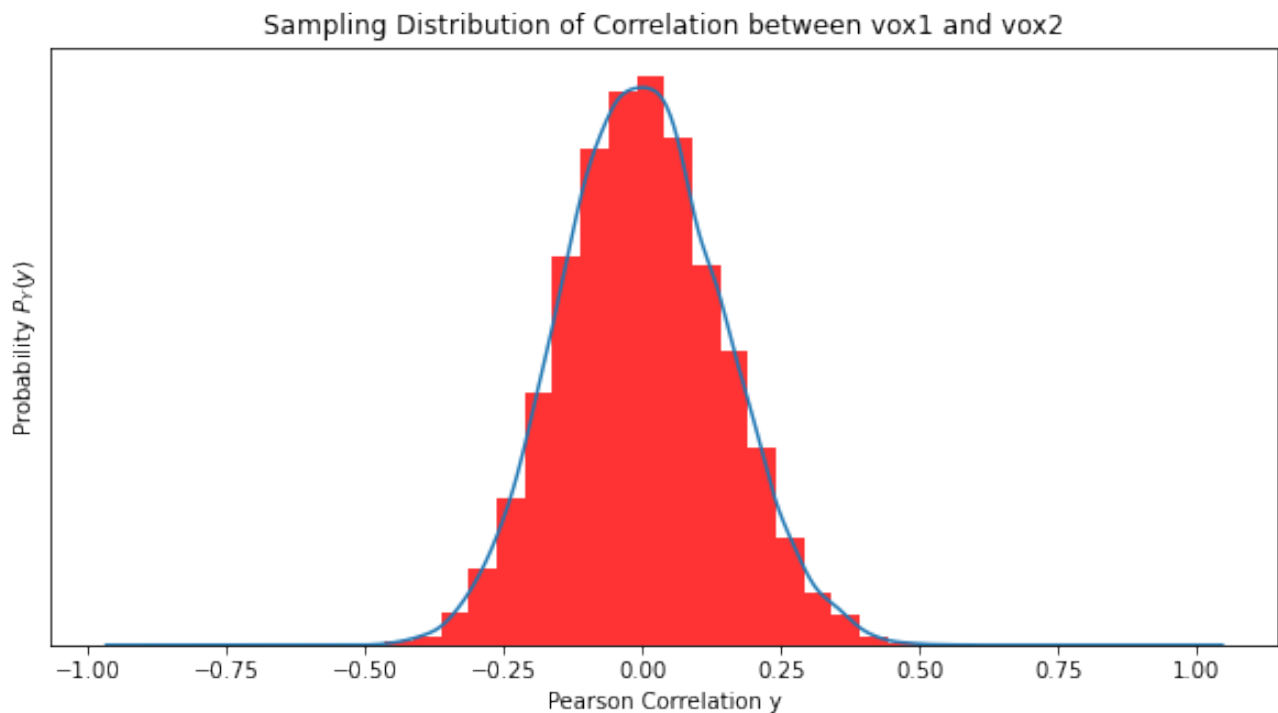


FIGURE 8. Sampling distribution of correlation between vox1 and vox2

The one sided p-value is: 3.812065700707869e-05

As our intuition and calculated p-value indicates that ($p\text{-value} < \alpha\text{ level}$) we have a strong evidence against the null assumption that states that two voxel responses have zero or negative correlation so that we reject our null hypothesis.

2.4. Part D. In this part of the question, the average BOLD responses in a face-selective region of the human brain have been recorded in two separate experiments. The responses of this region to building images (1st experiment) and face images (2nd experiment) are stored in the variables `building` and `face` for 20 subjects. Additionally, there is prior assumption that the same subject population was recruited in both experiments. We need to use bootstrapping (10000 iterations) to calculate the two-tailed p-value for the null hypothesis that there is no difference between the building and face responses.

In the underlying conditions, our prior believe is that there is no difference between the building and face as it is our null hypothesis. Also, researcher/alternative hypothesis is just the opposite so that our hypothesis procedure requires two-tailed p-value calculation. However, in these case, the situation is bit different than the previous ones as we have have an assumption that the subjects are the same for both experiments. We need to generate single difference of means. Our subject can be faced with $2^2 = 4$ cases which are

- Face, Face
- Face, Building
- Building, Face
- Building, Building

Hence, to generate unbiased experiment in statistical sense we need to generate samples by choosing 4 cases with the number of bootstrap iteration. Things will be more clear with the code part as it is placed below.

```
1 random_seed(31)
2
3 assert building.shape[0] == face.shape[0], 'Dimensionality Mismatch!'
4
5 sample_size = np.arange(building.shape[0])
6
7 _mean_diff = list()
8 bootstrap_iters = np.arange(10000)
9
10 for ii in bootstrap_iters:
11     resample = []
12
13     for jj in sample_size:
14
15         bootstrap_idx = np.random.choice(np.arange(building.shape[0]), replace =
16             ↪ True)
17         options = [0] * 2
18         _option = building[jj] - face[jj]
19         options.append(_option)
20         _option = face[jj] - building[jj]
21         options.append(_option)
22         resample.append(np.random.choice(options))
23
24     _mean_diff.append(np.mean(resample))
```

Here, I bootstrapped samples according to our underlying situation and calculate the difference of means in. Let's look at the distribution of difference of means of follows.

```

1 mean_diff = pd.Series(_mean_diff)
2 fig, ax = plt.subplots(figsize = (10,5))
3 mean_diff.plot.kde(ax=ax, legend = False, title='Difference in means of building
  ↳ and face')
4 mean_diff.plot.hist(density=True, ax = ax, bins = 40, alpha = 0.8, color =
  ↳ 'red')
5 ax.set_ylabel('Probability $P_X(x)$')
6 ax.set_xlabel('Difference in means (x)')
7 ax.grid(axis='y')
8 ax.set_yticks([])

```

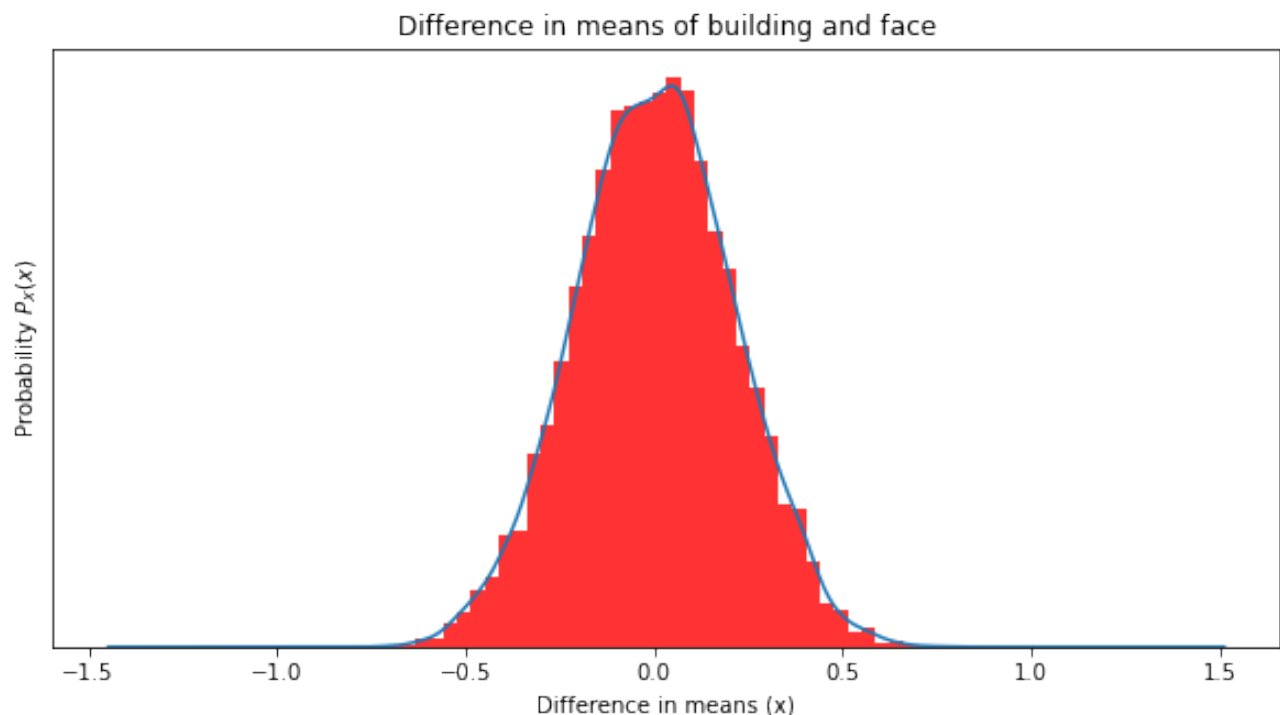


FIGURE 9. Sampling distribution of difference of means of building and face

The shape of the distribution implies zero centered Gaussian. Let's test our hypothesis by calculating test statistic and p-value as follows.

```

1 x_actual = np.mean(building) - np.mean(face)
2 mean = mean_diff.mean()
3 std = mean_diff.std()
4 z_score = mean - x_actual
5 z_score /= std
6 p_value = scipy.special.ndtr(- z_score) * two_sided
7 print('The two sided p-value is:', p_value)

```

The two sided p-value is: 0.000356527118064005

As p-value indicates rejection of our prior believe that there is no difference between the building and face since it is much smaller than our alpha level (0.05). Hence, we kindly reject our null hypothesis as we found strong evidence against it and conclude that the hypothesis that states there is difference between the building and face responses are statistically significant.

2.5. **Part E.** In this part, we are asked to repeat the exercise in part d, but this time assuming that the subject populations recruited for the two experiments are distinct. To do that, we will bootstrap (10000 iterations) to calculate the two-tailed p-value for the null hypothesis that there is no difference between the building and face responses.

Here is Python code for computation of its sampling distribution, visualization of its discretized density and p-value.

```
1 arr_stack = np.hstack((building, face))
2 arr_bootstrap = bootstrap(arr_stack)
3 samples1 = arr_bootstrap[:, :len(building)]
4 samples2 = arr_bootstrap[:, len(building):]
5 means1 = np.mean(samples1, axis=1)
6 means2 = np.mean(samples2, axis=1)
7 sample_diff_means = means1 - means2
```

At this point, difference of means are calculated, let's visualize the distribution as follows.

```
1 sample_mean_dist = pd.DataFrame()
2 sample_mean_dist['Mean Difference'] = sample_diff_means.flatten()
3 fig, ax = plt.subplots(figsize = (10,5))
4 sample_mean_dist.plot.kde(ax=ax, title='Difference of Means of Bootstrapped
  ↳ Populations building and face')
5 sample_mean_dist.plot.hist(density=True, ax = ax, bins = 50)
6 ax.set_ylabel('Probability $P_X(x)$')
7 ax.set_xlabel('Difference in means (x)')
8 ax.grid(axis='y')
9 ax.set_yticks([])
```

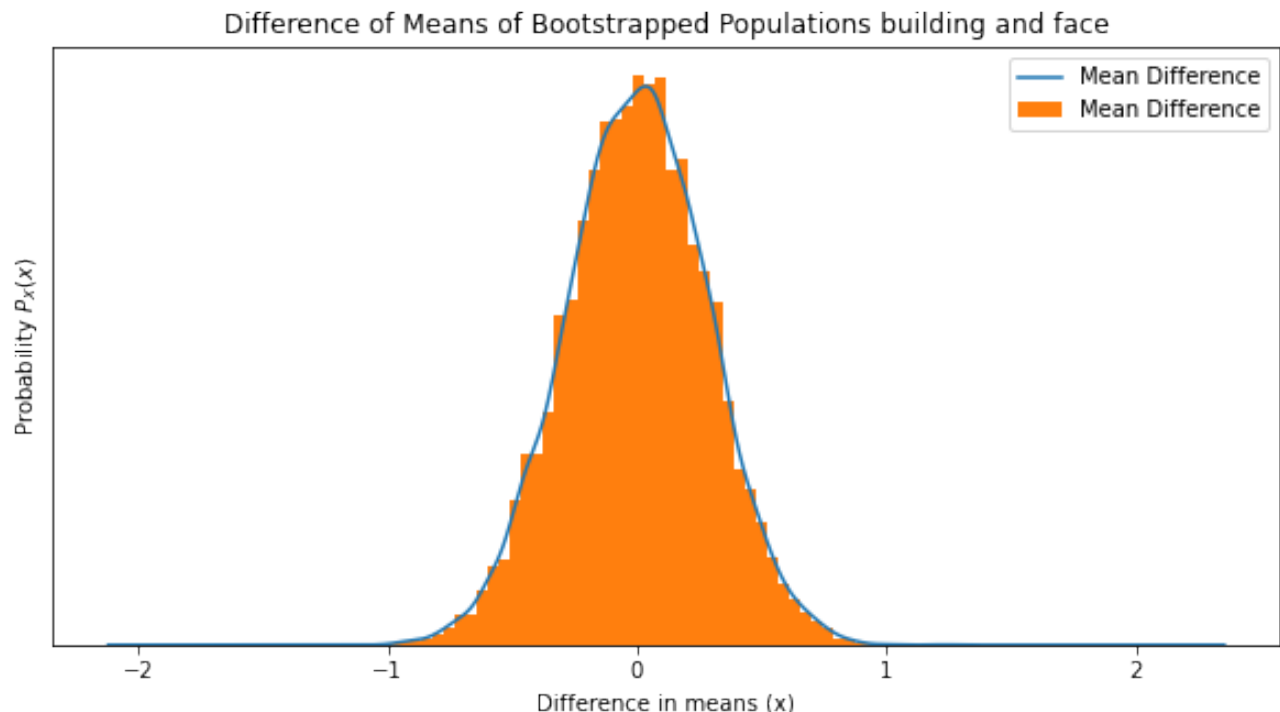


FIGURE 10. Sampling distribution of difference of means of building and face


```
1 x_actual = np.mean(building) - np.mean(face)
2 mean = sample_mean_dist.mean()
3 std = sample_mean_dist.std()
4 z_score = mean - x_actual
5 z_score /= std
6 p_value = scipy.special.ndtr(- z_score) * two_sided
7 print('The two sided p-value is:', p_value)
```

The one sided p-value is: Mean Difference 0.00844

Hence, corresponding p-value is calculated as 0.00844 that implies strong evidence against the null hypothesis since it is smaller than alpha level so that we reject our prior belief that there is no difference between the building and face responses.

3. SOURCE CODE

```
1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # In[19]:
5
6
7  import numpy as np
8  import matplotlib.pyplot as plt
9  import h5py
10 import scipy
11 import pandas as pd
12 import scipy.special as special
13 import random
14
15
16 # In[20]:
17
18
19 #cd D:\ThisSemester\CompNeuro\Homeworks\Hw3\HW3_Can_Kocagil\Assignment
20
21
22 # ### Question 1
23
24 # In[21]:
25
26
27 f = h5py.File('hw3_data2.mat', 'r')
28
29 X = np.array(f.get('Xn')).T
30
31 y = np.array(f.get('Yn')).flatten()
32
33 print(X.shape, y.shape)
34
35
36 # In[22]:
37
38
39 def random_seed(seed:int = 42) -> None :
40     """ Random seeding for reproducibility
41
42         Arguments:
43             - seed (int) : random state
44
45         Returns:
46             - None
47
48     """
49     np.random.seed(seed)
50     random.seed(seed)
51
```

```

52
53 # In[23]:
54
55
56 class RidgeRegression(object):
57
58     """
59     Ridge regression is a method of estimating the coefficients of
    ↪ multiple-regression models in
60     scenarios where independent variables are highly correlated.
61
62     """
63     def __init__(self, Lambda:float=1):
64         """
65         Constructor method for initialization of ridge regression model.
66
67         Arguments:
68         - Lambda (float): is the parameter which balances the amount
69           of emphasis given to minimizing RSS vs minimizing sum of
70     ↪ square of coefficients
71
72         """
73
74         self.Lambda = Lambda
75
76     def fit(self, X:np.ndarray, y:np.ndarray) -> None:
77         """
78
79         Given the pair of X,y, fit the data, i.e., find parameter W such
80     ↪ that sum of square error
81         is minimized.
82
83         Arguments:
84         - X (np.ndarray) : Regressor data
85         - X (np.ndarray) : Ground truths for regressors
86
87         Returns:
88         - None
89
90         """
91
92         I = np.eye(X.shape[1])
93
94         self.W = np.linalg.inv(
95             X.T.dot(X) + self.Lambda * I
96             ).dot(X.T).dot(y)
97
98         return self
99
100     def predict(self,X:np.ndarray) -> np.ndarray :
101         """
102

```

```

103         Given the test data X, we predict the target variable.
104
105         Arguments:
106         - X (np.ndarray) : The independant variable (regressor)
107
108         Returns:
109         - Y_hat (np.ndarray) : Estimated value of y
110
111         """
112
113         return X.dot(self.W)
114
115
116     def parameters(self) -> None:
117         """
118         Returns the estimated parameter W of the Ridge Regression
119
120         """
121         return self.W
122
123     def eval_r2(self, y_true: np.ndarray, y_pred: np.ndarray) -> np.float:
124         """
125         Given the true dependant variable and estimated variable, computes
126         ↪ proportion of
127         explained variance R^2 by square the Pearson correlation between
128         ↪ true dependant
129         variable and estimated variabl
130
131         Arguments:
132         - y_true (np.ndarray) : true dependant variable
133         - y_pred (np.ndarray) : estimated variable
134
135         Returns:
136         - r_squared (np.float) : Proportion of explained variance
137
138         """
139         _pearson = np.corrcoef(y_true, y_pred)
140         pearson = _pearson[1][0]
141         r_squared = np.square(pearson)
142         return r_squared
143
144     @staticmethod
145     def R2(y_true: np.ndarray, y_pred: np.ndarray) -> np.float:
146         r_squared = (1 - (sum((y_true - (y_pred))**2) / ((len(y_true) - 1) *
147         ↪ np.var(y_true.T, ddof=1)))) * 100
148         return r_squared
149
150     def __str__(self):
151         model = RidgeRegression().__class__.__name__
152         model += f" with parameter \n"
153         model += f"{self.Lambda}"
154         return model

```

```

154
155
156     def __repr__(self):
157         model = RidgeRegression().__class__.__name__
158         model += f" with parameter \n"
159         model += f"{self.Lambda}"
160         return model
161
162
163 # In[24]:
164
165
166 class K_fold(object):
167     """
168     Cross-validation, sometimes called rotation estimation or out-of-sample
↪ testing,
169     is any of various similar model validation techniques for assessing how the
↪ results
170     of a statistical analysis will generalize to an independent data set
171
172
173     """
174     def __init__(self, sample_size: int = y.shape[0], folds: int = 10):
175         """
176         Constructor method for initializing the sample size and the number
↪ of folds
177
178         Arguments:
179             - sample_size (int) : How many samples are in the dataset
180             - folds (int) : the number of folds
181
182         """
183
184         self.sample_size = sample_size
185         self.folds = folds
186         self.fold_size = int(sample_size / folds)
187
188     def split(self):
189         """
190
191         Generator function for splitting data as validation (10%), testing
↪ (10%) and
192         training (80%) as K-fold cross validation based resampling
193
194         """
195
196         for idx in range(self.folds):
197             _val_idx = idx * self.fold_size
198             _test_idx = (idx + 1) * self.fold_size
199             _train_idx = (idx + 2) * self.fold_size
200
201             val_idx = np.arange(_val_idx, _test_idx) % self.sample_size
202             test_idx = np.arange(_test_idx, _train_idx) % self.sample_size

```

```

203         train_idx = np.arange(_train_idx, self.sample_size + _val_idx) %
           ↪ self.sample_size
204
205         yield val_idx, test_idx, train_idx
206
207
208     # In[26]:
209
210
211     dict_inference = {
212         'test' : dict(),
213         'val'  : dict()
214     }
215
216
217     phases = [
218         'train',
219         'val',
220         'test'
221     ]
222 ]
223
224     log_lambda_arr = np.logspace(
225         start = 0,
226         stop  = 12,
227         num   = 500,
228         base  = 10
229     )
230
231     cv = K_fold(folds = 10)
232
233     for val_idx, test_idx, train_idx in cv.split():
234
235         X_list = [
236             X[train_idx],
237             X[val_idx],
238             X[test_idx]
239         ]
240
241         y_list = [
242             y[train_idx],
243             y[val_idx],
244             y[test_idx]
245         ]
246
247
248         for _lambda in log_lambda_arr:
249
250             for phase, X_phase, y_phase in zip(phases, X_list, y_list):
251                 if phase == 'train':
252                     model = RidgeRegression(_lambda)
253                     model.fit(X_phase, y_phase)
254
255                 else:
256                     preds = model.predict(X_phase)

```

```

257         r2_score = model.eval_r2(y_phase, preds)
258         dict_inference[phase].setdefault(
259             _lambda, list()).append(r2_score)
260
261 inference_r2 = {
262     phase : {
263         _lambda : np.mean(r2_score) for _lambda, r2_score in
264             ↪ dict_inference[phase].items()
265     }
266     for phase in ['val', 'test']
267 }
268
269 # In[27]:
270
271
272 best_r2 = 0
273 for _lambda, r_2 in inference_r2['val'].items():
274     if r_2 > best_r2:
275         best_r2 = r_2
276         best_lambda = _lambda
277
278
279 print(f'Best lambda parameter that maximizes the R^2 is : {best_lambda}')
280 print('Best R^2 along the testing :', inference_r2['test'][best_lambda])
281 print('Best R^2 along the validation :', inference_r2['val'][best_lambda])
282
283
284 # In[28]:
285
286
287 lists1 = sorted(inference_r2['val'].items())
288 x1, y1 = zip(*lists1)
289 lists2 = sorted(inference_r2['test'].items())
290 x2, y2 = zip(*lists2)
291 plt.figure(figsize = (10,5))
292 plt.plot(x2, y2, color='orange')
293 plt.plot(x1, y1, color='g')
294 plt.legend(['test', 'validation'])
295 plt.ylabel('$R^2$')
296 plt.xlabel('$\lambda$')
297 plt.title('$R^2$ versus $\lambda$')
298 plt.xscale('log')
299 plt.grid()
300 plt.show()
301
302
303 # In[29]:
304
305
306 random_seed(10)
307
308 bootstrap_iters = range(500)
309 sample_idx = np.arange(X.shape[0])
310 parameters = list()

```

```

311
312 for idx in bootstrap_iters:
313
314     bootstrap_idx = np.random.choice(sample_idx, size = 1000, replace = True)
315     y_bootstrap = y[bootstrap_idx]
316     X_bootstrap = X[bootstrap_idx]
317     ridge = RidgeRegression(Lambda = 0)
318     ridge.fit(X_bootstrap,y_bootstrap)
319     parameters.append(ridge.parameters())
320
321 w_bootstrap = np.array(parameters)
322 w_mean = np.mean(w_bootstrap, axis=0)
323 w_std = np.std(w_bootstrap, axis=0)
324
325
326 # In[30]:
327
328
329 plt.figure(figsize = (10,5))
330 plt.errorbar(np.arange(1, 101),
331             w_mean,
332             yerr= w_std,
333             ecolor='red',
334             elinewidth=1,
335             capsize=1)
336 plt.title('Ridge Model OLS Weights')
337 plt.xlabel('i')
338 plt.ylabel('$W_i$')
339 plt.show()
340
341
342 # In[31]:
343
344
345 two_sided = 2
346 p_values = special.ndtr(- w_mean / w_std) * two_sided
347 alpha_level = 0.05
348 significants = np.argwhere(p_values < alpha_level).flatten()
349 print(f' Index of the parameters that are significantly different than 0: \n
350 ↪ {significants}')
351
352
353 # In[32]:
354
355
356 random_seed(10)
357
358 bootstrap_iters = range(500)
359 sample_idx = np.arange(X.shape[0])
360 parameters = list()
361
362 for idx in bootstrap_iters:
363
364     bootstrap_idx = np.random.choice(sample_idx, size = 1000, replace = True)

```



```

365     y_bootstrap = y[bootstrap_idx]
366     X_bootstrap = X[bootstrap_idx]
367     ridge = RidgeRegression(Lambda = best_lambda)
368     ridge.fit(X_bootstrap, y_bootstrap)
369     parameters.append(ridge.parameters())
370
371 w_bootstrap = np.array(parameters)
372 w_mean = np.mean(w_bootstrap, axis=0)
373 w_std = np.std(w_bootstrap, axis=0)
374
375
376 # In[33]:
377
378
379 plt.figure(figsize = (10,5))
380 plt.errorbar(np.arange(1, 101),
381             w_mean,
382             yerr= w_std,
383             ecolor='red',
384             elinewidth=1,
385             capsize=1)
386 plt.title('Ridge Model  $\lambda_{\text{optimal}}$  Weights')
387 plt.xlabel('i')
388 plt.ylabel('$W_i$')
389 plt.show()
390
391
392 # In[34]:
393
394
395 p_values = scipy.special.ndtr(- w_mean / w_std) * two_sided
396 significant = np.where(p_values < alpha_level).flatten()
397 print(f' Index of the parameters that are significantly different than 0: \n
398     ↪ {significant}')
399
400
401 # ### Question 2
402
403 # ## Part A
404
405 # In[44]:
406
407
408 f = h5py.File('hw3_data3.mat', 'r')
409
410 pop1 = np.array(
411     f.get('pop1')
412 )
413
414 pop2 = np.array(
415     f.get('pop2')
416 )
417
418

```

```

419 # In[45]:
420
421
422 def bootstrap(sample:np.ndarray, bootstrap_iters:iter = range(10000),
↳ random_state:int = 11) -> np.ndarray:
423     """
424
425     Generate bootstrap samples using random sampling with replacement.
426
427     Arguments:
428         - sample (np.ndarray) : Sample to be bootstrapped
429         - bootstrap_iters (iterator object) : Specification of bootstrap
↳ iterations
430         - random_state (int) : Random seed for reproducibility
431
432     Returns:
433         - bootstrap_samples (np.ndarray) : Bootstrapped array
434
435     """
436     random_seed(random_state)
437     size = sample.shape[0]
438     bootstrap_samples = list()
439
440     for idx in bootstrap_iters:
441         bootstrap_idx = np.random.choice(np.arange(sample.shape[0]), size =
↳ size, replace = True)
442         bootstrap_samples.append(sample[bootstrap_idx])
443
444     return np.array(bootstrap_samples)
445
446
447 # In[46]:
448
449
450 pop = np.vstack([pop1,pop2])
451 pop_bootstrap = bootstrap(pop)
452 sample_1 = pop_bootstrap[:,len(pop1)].squeeze(2)
453 sample_2 = pop_bootstrap[:,len(pop1):].squeeze(2)
454 sample_1_bootstrap_mean = sample_1.mean(axis = 1)
455 sample_2_bootstrap_mean = sample_2.mean(axis = 1)
456 sample_diff_means = sample_1_bootstrap_mean - sample_2_bootstrap_mean
457 sample_mean_dist = pd.DataFrame()
458 sample_mean_dist['Mean Difference'] = sample_diff_means.flatten()
459 fig, ax = plt.subplots(figsize = (10,5))
460 sample_mean_dist.plot.kde(ax=ax, title='Difference of Means of Bootstrapped
↳ Populations 1 and 2')
461 sample_mean_dist.plot.hist(density=True, ax = ax, bins = 15)
462 ax.set_ylabel('Probability $P_X(x)$')
463 ax.set_xlabel('Difference in means (x)')
464 ax.grid(axis='y')
465 ax.set_yticks([])
466
467
468 # In[47]:

```

```

469
470
471 pop1_bootstrap = bootstrap(pop1)
472 pop2_bootstrap = bootstrap(pop2)
473 pop1_bootstrap_mean = np.mean(pop1_bootstrap, axis = 1)
474 pop2_bootstrap_mean = np.mean(pop2_bootstrap, axis = 1)
475
476 mean_dist = pd.DataFrame()
477 mean_dist['pop1 Mean'] = pop1_bootstrap_mean.flatten()
478 mean_dist['pop2 Mean'] = pop2_bootstrap_mean.flatten()
479 mean_dist['Mean Difference'] = pop1_bootstrap_mean - pop2_bootstrap_mean
480
481 fig, ax = plt.subplots(figsize = (10,5))
482 mean_dist.plot.kde(ax=ax, title='Difference of Means of Bootstrapped Populations
↳ 1 and 2')
483 mean_dist.plot.hist(density=True, ax = ax, bins = 15)
484 ax.set_ylabel('Probability $P_X(x)$')
485 ax.set_xlabel('Difference in means (x)')
486 ax.grid(axis='y')
487 ax.set_yticks([])
488
489
490 fig, ax = plt.subplots(figsize = (10,5))
491 mean_dist['Mean Difference'].plot.kde(ax=ax, legend = True, title='Difference of
↳ Means of Bootstrapped Populations 1 and 2')
492 mean_dist['Mean Difference'].plot.hist(density=True, ax = ax, bins = 15)
493 ax.set_ylabel('Probability $P_X(x)$')
494 ax.set_xlabel('Difference in means (x)')
495 ax.grid(axis='y')
496 ax.set_yticks([])
497
498
499 # In[48]:
500
501
502 actual_diff_means = pop1.mean() - pop2.mean()
503 std_test = sample_mean_dist['Mean Difference'].std()
504 mean_test = sample_mean_dist['Mean Difference'].mean()
505
506 z_cal = (mean_test - actual_diff_means) / std_test
507 p_values = scipy.special.ndtr(z_cal) * two_sided
508
509 print('The two sided p-value is:', p_values)
510
511
512 # ## Part B
513
514 # In[49]:
515
516
517 vox1 = np.array(
518     f.get('vox1')
519     ).flatten()
520

```

```

521 vox2 = np.array(
522     f.get('vox2')
523     ).flatten()
524
525
526 print(
527     vox1.shape,
528     vox2.shape
529 )
530
531 vox1_bootstrap = bootstrap(vox1)
532 vox2_bootstrap = bootstrap(vox2)
533
534 def corr(X: list or np.ndarray, Y: list or np.ndarray) -> list:
535     """
536
537     Given the X,Y distributions, computes the Pearson Correlation element
538     ↪ wise.
539
540     Arguments:
541     - X (list or np.ndarray) : First distribution
542     - Y (list or np.ndarray) : Second distribution
543
544     Returns:
545     - pearson_corrs (list[float]) : Computed correlations element
546     ↪ wise
547
548     """
549     assert X.shape == Y.shape, 'Dimension Mismatch!'
550     return [scipy.stats.pearsonr(X[i], Y[i])[0] for i in range(X.shape[0])]
551
552 corr_bootstrap = corr(vox1_bootstrap, vox2_bootstrap)
553
554 fig, ax = plt.subplots(figsize = (10,5))
555 pd.Series(corr_bootstrap).plot.kde(ax=ax, legend = False, title='Sampling
556     ↪ Distribution of Correlation between vox1 and vox2')
557 pd.Series(corr_bootstrap).plot.hist(density=True, ax = ax, bins = 20, alpha =
558     ↪ 0.8,color = 'red')
559 ax.set_ylabel('Probability $P_Y(y)$')
560 ax.set_xlabel('Pearson Correlation y')
561 ax.grid(axis='y')
562 ax.set_yticks([])
563 # Thanks to
564 ↪ https://stackoverflow.com/questions/15033511/compute-a-confidence-interval-from-sample-
565 def confidence_interval(data: list or np.ndarray, confidence:float=0.95) ->
566     ↪ tuple:
567     """
568
569     Given the distribution and confidence level, computes the confidence
570     ↪ interval.
571
572     Arguments:
573     - data (list or np.ndarray) : Input distribution

```

```

568         - confidence (float) : confidence level in the range [0,1]
569
570
571     Returns:
572         - confidence_level (tuple[np.ndarray]) : lower, upper limits
    ↪     respectively
573
574
575
576     """
577     a = 1.0 * np.array(data)
578     n = len(a)
579     m, se = np.mean(a), scipy.stats.sem(a)
580     h = se * scipy.stats.t.ppf((1 + confidence) / 2., n-1)
581     return m-h, m+h
582
583 def _confidence_interval(data, confidence=0.95):
584     return scipy.stats.t.interval(confidence, len(data)-1, loc=np.mean(data),
    ↪     scale=st.sem(data))
585
586 corr_mean = np.mean(corr_bootstrap)
587 lower, upper = confidence_interval(corr_bootstrap, confidence=0.95)
588 print('Mean correlation value:', corr_mean)
589 print(f'95% confidence interval of the correlation values: {lower, upper}')
590
591 is_corr_zero = np.argwhere(corr_bootstrap == 0)
592 corr_zero_percentage = 100 * is_corr_zero.shape[0] / 10000
593 print('Percentage of zero correlation values:', corr_zero_percentage)
594
595
596 # ## Part C
597
598 # In[50]:
599
600
601 vox1_ind = bootstrap(vox1, range(10000), random_state=42)
602 vox2_ind = bootstrap(vox2, range(10000), random_state=21)
603
604 _corr_ind = corr(vox1_ind, vox2_ind)
605
606 corr_ind = pd.Series(_corr_ind)
607
608
609 fig, ax = plt.subplots(figsize = (10,5))
610 corr_ind.plot.kde(ax=ax, legend = False, title='Sampling Distribution of
    ↪     Correlation between vox1 and vox2')
611 corr_ind.plot.hist(density=True, ax = ax, bins = 20, alpha = 0.8, color = 'red')
612 ax.set_ylabel('Probability $P_Y(y)$')
613 ax.set_xlabel('Pearson Correlation y')
614 ax.grid(axis='y')
615 ax.set_yticks([])
616
617
618 actual_corr, _ = scipy.stats.pearsonr(vox1, vox2)

```

```

619 mean_corr = corr_ind.mean()
620 std_corr = corr_ind.std()
621 z_score = mean_corr - actual_corr
622 z_score /= std_corr
623 p_value = scipy.special.ndtr(z_score)
624 print('The one sided p-value is:', p_value)
625
626
627 # ## Part D
628
629 # In[52]:
630
631
632 building = np.array(f.get('building')).flatten()
633
634 face = np.array(f.get('face')).flatten()
635
636 print(
637     building.shape,
638     face.shape
639 )
640
641 random_seed(31)
642
643 assert building.shape[0] == face.shape[0], 'Dimensionality Mismatch!'
644
645 sample_size = np.arange(building.shape[0])
646
647 _mean_diff = list()
648 bootstrap_iters = np.arange(10000)
649
650 for ii in bootstrap_iters:
651
652     resample = []
653
654     for jj in sample_size:
655
656         bootstrap_idx = np.random.choice(np.arange(building.shape[0]), replace =
        ↪ True)
657         options = [0] * 2
658         _option = building[jj] - face[jj]
659         options.append(_option)
660         _option = face[jj] - building[jj]
661         options.append(_option)
662         resample.append(np.random.choice(options))
663
664     _mean_diff.append(np.mean(resample))
665
666
667 mean_diff = pd.Series(_mean_diff)
668 fig, ax = plt.subplots(figsize = (10,5))
669 mean_diff.plot.kde(ax=ax, legend = False, title='Difference in means of building
        ↪ and face')

```

```

670 mean_diff.plot.hist(density=True, ax = ax, bins = 40, alpha = 0.8, color =
    ↪ 'red')
671 ax.set_ylabel('Probability $P_X(x)$')
672 ax.set_xlabel('Difference in means (x)')
673 ax.grid(axis='y')
674 ax.set_yticks([])
675
676
677 x_actual = np.mean(building) - np.mean(face)
678 mean = mean_diff.mean()
679 std = mean_diff.std()
680 z_score = mean - x_actual
681 z_score /= std
682 p_value = scipy.special.ndtr(- z_score) * two_sided
683 print('The two sided p-value is:', p_value)
684
685
686 # ## Part E
687
688 # In[53]:
689
690
691 arr_stack = np.hstack((building, face))
692 arr_bootstrap = bootstrap(arr_stack)
693 samples1 = arr_bootstrap[:, :len(building)]
694 samples2 = arr_bootstrap[:, len(building):]
695 means1 = np.mean(samples1, axis=1)
696 means2 = np.mean(samples2, axis=1)
697 sample_diff_means = means1 - means2
698
699
700 sample_mean_dist = pd.DataFrame()
701 sample_mean_dist['Mean Difference'] = sample_diff_means.flatten()
702 fig, ax = plt.subplots(figsize = (10,5))
703 sample_mean_dist.plot.kde(ax=ax, title='Difference of Means of Bootstrapped
    ↪ Populations building and face')
704 sample_mean_dist.plot.hist(density=True, ax = ax, bins = 50)
705 ax.set_ylabel('Probability $P_X(x)$')
706 ax.set_xlabel('Difference in means (x)')
707 ax.grid(axis='y')
708 ax.set_yticks([])
709
710
711
712 x_actual = np.mean(building) - np.mean(face)
713 mean = sample_mean_dist.mean()
714 std = sample_mean_dist.std()
715 z_score = mean - x_actual
716 z_score /= std
717 p_value = scipy.special.ndtr(- z_score) * two_sided
718 print('The two sided p-value is:', p_value)
719
720
721 # In[ ]:

```

REFERENCES

- [1] Wikipedia contributors. *Bootstrapping (statistics)* — *Wikipedia, The Free Encyclopedia*. [https://en.wikipedia.org/w/index.php?title=Bootstrapping_\(statistics\)&oldid=1014701581](https://en.wikipedia.org/w/index.php?title=Bootstrapping_(statistics)&oldid=1014701581). [Online; accessed 2-April-2021]. 2021.