

MLPR Assignment 2 2018

s1879797

November 9, 2018

1 Question 1

1.1 Q1a

The standard error verification:

```
1 def compute_standard_error(y_train, y_val):
2     mean_y_train = np.mean(y_train)                # -9.13868774539957e-15
3     print("mean_y_train: ", mean_y_train)
4     if round(mean_y_train, 10) == 0:
5         print("The mean is approximately 0!")
6     else:
7         print("The mean is not approximately 0!")
8
9     mean_y_val = np.mean(y_val)                    # -0.2160085093241599
10    print("mean_y_val: ", mean_y_val)
11    if round(mean_y_val, 10) == 0:
12        print("The mean is approximately 0!")
13    else:
14        print("The mean is not approximately 0!")
15
16    std_y_train = np.std(y_train[:len(y_val)])
17    print("std_y_train: ", std_y_train)
18    print("std_y_train for entire y_train: ", np.std(y_train))
19    std_y_val = np.std(y_val)
20    print("std_y_val: ", std_y_val)
21
22    sem_y_train = std_y_train / np.sqrt(len(y_val))
23    print("sem_y_train: ", sem_y_train)
24    sem_y_val = std_y_val / np.sqrt(len(y_val))
25    print("sem_y_val: ", sem_y_val)
```

It is true that these standard error bars sometimes could not reliably indicate the mean for future data because the actual mean of some sample is within 1 standard error 68% of time, within 2 standard errors 95% of time and within 3 standard errors 99.7% of time (**68–95–99.7 rule**). For a good estimation you have to choose good confidence intervals and I think this is the biggest problem when using standard error bars.

MEAN	STANDARD ERROR	MEAN	STANDARD ERROR
Y_TRAIN	Y_TRAIN	Y_VAL	Y_VAL
$-9.139 * 10^{-15}$	0.012	-0.216	0.013

Table 1: The mean and the standard error for training and validation

The results are presented in table 1. We can clearly see that the mean for y_train can be easily approximate to 0 and the mean for y_val is different than 0.

1.2 Q1b

The method for removing bad columns:

```

1 def remove_unnecessary_features(X_train, X_val, X_test):
2     # Remove constant features
3     X_train_delete_1 = remove_constant_features(X=X_train)
4     X_val_delete_1 = remove_constant_features(X=X_val)
5     X_test_delete_1 = remove_constant_features(X=X_test)
6
7     # Remove duplicate features
8     X_train_delete_2 = remove_duplicates_features(X=X_train)
9     X_val_delete_2 = remove_duplicates_features(X=X_val)
10    X_test_delete_2 = remove_duplicates_features(X=X_test)
11
12    X_train_delete = X_train_delete_1 + X_train_delete_2
13    X_val_delete = X_val_delete_1 + X_val_delete_2
14    X_test_delete = X_test_delete_1 + X_test_delete_2
15
16    # Get only the columns that are to be removed from all three datasets
17    to_delete = reduce(np.intersect1d, (X_train_delete, X_val_delete,
18    X_test_delete))
19
20    # Filter categories
21    X_train_first, X_train_second = filter_remove_categories(to_delete,
22    X_train_delete_1, X_train_delete_2)
23    print("X_train first: ", X_train_first)
24    print("X_train second: ", X_train_second)
25    X_val_first, X_val_second = filter_remove_categories(to_delete,
26    X_val_delete_1, X_val_delete_2)
27    print("X_val first: ", X_val_first)
28    print("X_val second: ", X_val_second)
29    X_test_first, X_test_second = filter_remove_categories(to_delete,
30    X_test_delete_1, X_test_delete_2)
31    print("X_test first: ", X_test_first)
32    print("X_test second: ", X_test_second)
33
34    # Delete bad columns
35    X_train = np.delete(X_train, to_delete, axis=1) # (40754, 373)
36    X_val = np.delete(X_val, to_delete, axis=1) # (5785, 373)
37    X_test = np.delete(X_test, to_delete, axis=1) # (6961, 373)
38    return X_train, X_val, X_test

```

For this question I want to state that every set of data had a different number of columns that had to be removed. In order to keep the number of feature the same for every set of data I removed only the columns that had an equivalent to be removed in the other two sets. For this operation I used the function **reduce** from **functools** in order to apply **numpy intersect1d** over all three list. I had to mention this and I hope it will not be a problem since it does not influence the solution of the assignment and it does not provide a great help.

Small functions used in the above method:

```

1 # Remove constant input features (CATEGORY 1)
2 def remove_constant_features(X):
3     X_delete_1 = []
4     for i in range(X.shape[1]):
5         if (len(set(X[:, i])) == 1):
6             X_delete_1.append(i)
7     return X_delete_1
8
9 # Remove duplicate input features (CATEGORY 2)
10 def remove_duplicates_features(X):
11     unique, train_indices = np.unique(X, return_index=True, axis=1)
12     X_delete_2 = list(set(range(X.shape[1])) - set(train_indices))
13     return X_delete_2
14
15 # Classify removed features
16 def filter_remove_categories(to_delete, X_delete_1, X_delete_2):
17     X_first = []
18     X_second = []
19     for i in to_delete:
20         if i in X_delete_1:
21             X_first.append(i)
22         if i in X_delete_2:
23             X_second.append(i)
24     return X_first, X_second

```

The results are presented in table 2.

DATA SET	CONSTANT FEATURES	DUPLICATE FEATURES
TRAINING SET	[59, 69, 179, 189, 351]	[69, 78, 79, 179, 188, 189, 199, 287, 351, 359]
VALIDATION SET	[59, 69, 179, 189, 351]	[69, 78, 79, 179, 188, 189, 199, 287, 351, 359]
TEST SET	[59, 69, 179, 188, 189, 351]	[59, 69, 78, 79, 179, 188, 189, 199, 287, 351, 359]

Table 2: The columns that had to be remove from the data sets

2 Question 2

The main function for this question:

```
1 def set_linear_regression_baseline(X_train, X_val, y_train, y_val):
2     # For training
3     w_ls, b_ls = least_squares(X_train, y_train)
4     err = compute_err(X=X_train, yy=y_train, ww=w_ls, bb=b_ls)
5     print("ERROR Least Squares:", err) # 0.35524169481074713
6     X_prime, y_prime, w_prime = fit_linreg(X_train, y_train, 10)
7     err = compute_err(X=X_prime, yy=y_prime, ww=w_prime)
8     print("ERROR Regression with Regularization:", err) # 0.36915529643475953
9     w_grad, b_grad = fit_linreg_gradopt(X_train, y_train, 10)
10    err = compute_err(X=X_train, yy=y_train, ww=w_grad, bb=b_grad)
11    print("ERROR Gradient-Based Optimizer:", err) # 0.35575973762757745
12
13    # For validation
14    err = compute_err(X=X_val, yy=y_val, ww=w_ls, bb=b_ls)
15    print("ERROR Least Squares:", err) # 0.4182210942058271
16    X_prime, y_prime, dummy = fit_linreg(X_val, y_val, 10)
17    err = compute_err(X=X_prime, yy=y_prime, ww=w_prime)
18    print("ERROR Regression with Regularization:", err) # 0.47376619192358777
19    err = compute_err(X=X_val, yy=y_val, ww=w_grad, bb=b_grad)
20    print("ERROR Gradient-Based Optimizer:", err) # 0.42060375407081924
```

The `fit_linreg` function:

```
1 def fit_linreg(X, yy, alpha):
2     D = X.shape[1]
3     yy_prime = np.concatenate([yy, np.zeros(D)])
4     alphaI = alpha * np.identity(D)
5     X_prime = np.concatenate([X, alphaI], axis=0)
6     w_prime = np.linalg.lstsq(X_prime, yy_prime, rcond=0)[0]
7
8     return X_prime, yy_prime, w_prime
```

I also wanted to see the results for **least squares** method:

```
1 def least_squares(X, yy):
2     X_bias = np.concatenate([np.ones((X.shape[0], 1)), X], axis=1)
3     w_bias = np.linalg.lstsq(X_bias, yy, rcond=0)[0];
4
5     return w_bias[1:], w_bias[0]
```

The function which compute the predicted output:

```
1 def compute_err(X, yy, ww, bb=0):
2     y_predicted = X.dot(ww) + bb
3     return root_mean_square_error(y-expected=yy, y-predicted=y_predicted)
```

The function which computes the root mean square error:

```
1 def root_mean_square_error(y_expected, y_predicted):
2     sum = 0
3     for i in range(len(y_expected)):
4         sum += ((y_expected[i] - y_predicted[i]) ** 2)
5
6     return np.sqrt(sum/len(y_expected))
```

The result are displayed below in table 3.

DATA SET	LEAST SQUARES	LINEAR REGRESSION WITH REGULARIZATION	GRADIENT-BASED OPTIMIZER
TRAINING SET	0.355	0.369	0.356
VALIDATION SET	0.418	0.474	0.421

Table 3: Least Squares vs. Linear Regression with Regularization vs. Gradient-Based Optimizer

The results obtained from the Gradient-Based Optimizer function for root mean square error were better than the ones obtained using regularization. More than that, the results from the Gradient-Based Optimizer were similar with the ones obtained using least squares without regularization. Using a Gradient-Based Optimizer produce a lower error than Linear Regression using Regularization because for the first one the parameters are update at each iteration which can provide a faster convergence.

3 Question 3

3.1 Q3a

Overall function for this question:

```

1 def decrease_and_increase_input(X_train, X_val, y_train, y_val):
2     # Ex 3.a
3
4     # For training
5     V = pca(X=X_train, yy=y_train, alpha=10, K=10)
6     X_reduced = X_train.dot(V)
7     X_prime, yy_prime, w_prime = fit_linreg(X_reduced, y_train, 10)
8     err = compute_err(X=X_prime, yy=yy_prime, ww = w_prime)
9     print("Training error for K = 10: ", err) # 0.5756376489484005
10    X_reduced = X_val.dot(V)
11    X_prime, yy_prime, dummy = fit_linreg(X_reduced, y_val, 10)
12    err = compute_err(X=X_prime, yy=yy_prime, ww = w_prime)
13    print("Validation error for K = 10: ", err)# 0.5757941612096046
14
15    # For validation
16    V = pca(X=X_train, yy=y_train, alpha=10, K=100)
17    X_reduced = X_train.dot(V)
18    X_prime, yy_prime, w_prime = fit_linreg(X_reduced, y_train, 10)
19    err = compute_err(X=X_prime, yy=yy_prime, ww = w_prime)
20    print("Training error for K = 100: ", err) # 0.4153790239683063
21    X_reduced = X_val.dot(V)
22    X_prime, yy_prime, dummy = fit_linreg(X_reduced, y_val, 10)
23    err = compute_err(X=X_prime, yy=yy_prime, ww = w_prime)
24    print("Validation error for K = 100: ", err)# 0.4605743588041712
25

```

```

26 # Ex 3.b
27 X_prime, yy_prime, w_prime = histogram(X=X_train, yy=y_train, alpha=10)
28 err = compute_err(X=X_prime, yy=yy_prime, ww=w_prime)
29 print("Training error: ", err) # 0.3260341808564488
30 X_prime, yy_prime, dummy = histogram(X=X_val, yy=y_val, alpha=10)
31 err = compute_err(X=X_prime, yy=yy_prime, ww=w_prime)
32 print("Validation error: ", err) # 0.38817888465014794

```

The function which uses the reduced inputs to compute the error:

```

1 def pca(X, yy, alpha, K):
2     X_mu = np.mean(X, 0)
3     X_centred = X - X_mu
4
5     return pca_zm_proj(X=X_centred, K=K)

```

The result are displayed below in table 4.

DATA SET	LINEAR REGRESSION	PCA K=10	PCA K=100
TRAINING SET	0.369	0.576	0.415
VALIDATION SET	0.474	0.576	0.461

Table 4: Linear Regression with a normal matrix and with a reduced matrix

We use PCA in order to reduce the dimension of data so that the computation to be done faster. Normally, PCA should improve speed, accuracy and reduce error. I think that the reason that I got worse results using PCA with $K = 10$ than just doing Linear Regression on the normal data inputs is that PCA has some limitations and changing the number of features from 373 to 10 is too much and the results lose accuracy (for $K = 100$ PCA works really good and the results on validation set are better than using Linear Regression with Regularization).

3.2 Q3b

The histogram plotting function:

```

1 def histogram(X, yy, alpha):
2     print("46th feature")
3     count_0, count_25, count_0_per, count_25_per = get_percentage(X[:,46])
4     print('Count: 0({0}); -0.25({1})'.format(count_0, count_25))
5     print('Percentage: 0({0}); -0.25({1})'.format(count_0_per, count_25_per))
6     plt.clf()
7     plt.hist(X[:,45], bins= 20)
8     plt.show()
9
10    print("all data")
11    count_0, count_25, count_0_per, count_25_per = get_percentage(np.ravel(X))
12    print('Count: 0({0}); -0.25({1})'.format(count_0, count_25))
13    print('Percentage: 0({0}); -0.25({1})'.format(count_0_per, count_25_per))

```

```

14 # plt.clf()
15 # plt.hist(np.ravel(X), bins= 20)
16 # plt.show()
17
18 aug_fn = lambda X: np.concatenate([X, X==0, X<0], axis=1)
19 X_prime = aug_fn(X)
20
21 print("all data after adding extra binary features")
22 count_0, count_25, count_0_per, count_25_per = get_percentage(np.ravel(
X_prime))
23 print('Count: 0({0}); -0.25({1})'.format(count_0, count_25))
24 print('Percentage: 0({0}); -0.25({1})'.format(count_0_per, count_25_per))
25 # plt.clf()
26 # plt.hist(np.ravel(X_prime), bins= 20)
27 # plt.show()
28
29 return fit_linreg(X_prime, yy, alpha)

```

The function used for counting the number of values of 1 and -0.25:

```

1 def get_percentage(X):
2     count_0 = 0
3     count_25 = 0
4     for i in X:
5         if i == 0:
6             count_0 += 1
7         elif i == -0.25:
8             count_25 += 1
9
10    count_0_per = count_0 * 100 / len(X)
11    count_25_per = count_25 * 100 / len(X)
12    return count_0, count_25, count_0_per, count_25_per

```

The result fractions of -0.25 and 0 regarding the training set are presented in table 5, respectively in table 6 for validation set.

MODE	-0.25	0
	PERCENTAGE	PERCENTAGE
46TH FEATURE	34.38%	60.51%
ALL DATA	13.58%	66.78%
ALL DATA WITH EXTRA FEATURES	4.52%	62.14%

Table 5: The fractions of -0.25 and 0 from training data

MODE	-0.25	0
	PERCENTAGE	PERCENTAGE
46TH FEATURE	35.31%	56.72%
ALL DATA	14.25%	65.93%
ALL DATA WITH EXTRA FEATURES	4.75%	61.92%

Table 6: The fractions of -0.25 and 0 from validation data

The new errors for training and validation sets are in table 7.

DATA SET	LINEAR REGRESSION	EXTRA BINARY FEATURES
TRAINING SET	0.369	0.326
VALIDATION SET	0.474	0.388

Table 7: The results after adding extra binary features to input data

Adding more input features normally should decrease the error (more examples, accuracy should increase), but giving that our extra features are binary features I am not really sure of what should happen with the error. Given that we add 1 for all data that are equal to 0 and for all data lower than 0, as we can see in table 5 and table 6, the percentage of zeros from the input data is almost the same, while the percentage of -0.25 (and implicitly of all negative numbers) decreased and of course a huge number of ones has been added to the data set. Taking into account that now most of the values from the data set are 0 and 1 I think that the output has stabilized and the accuracy increased greatly which brought us a lower value for error.

The representation of the percentage for -0.25 and 0 from 46th feature is displayed in figure 2.

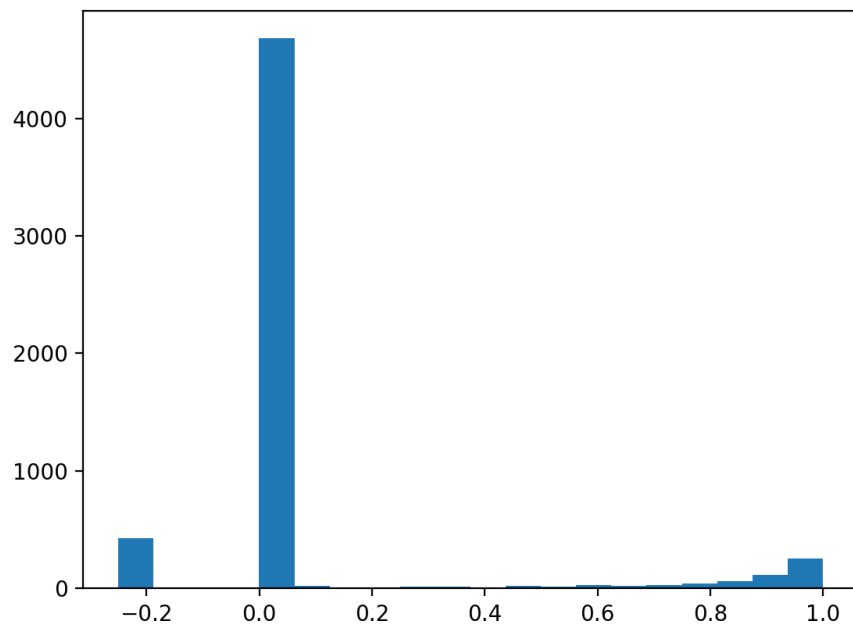
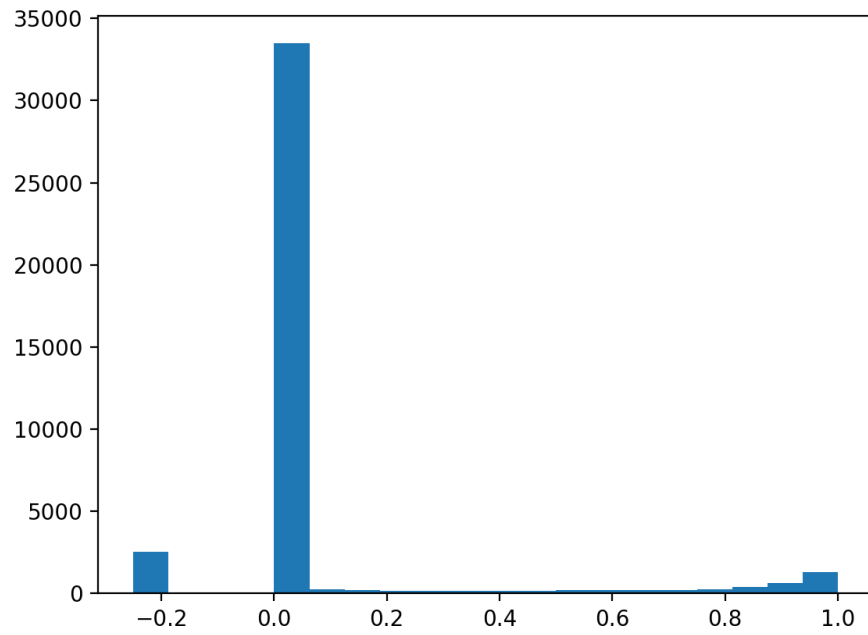


Figure 1: Histogram for training data (up) and for validation data (down)

4 Question 4

Overall function for this question:

```
1 def invent_classification(X_train, X_val, y_train, y_val):
2     K = 10          # number of thresholded classification problems to fit
3     # For training
4     ww_log, bb_log = use_logreg(X=X_train, yy=y_train, alpha=10, K=10)
5     X_reduced = X_train.dot(ww_log) + bb_log
6     X_prime, yy_prime, w_prime = fit_linreg(X=X_reduced, yy=y_train, alpha=10)
7     err = compute_err(X=X_prime, yy=yy_prime, ww = w_prime)
8     print("ERROR Logistic Regression: ", err) # 0.42230977042363277
9
10    # For validation
11    X_reduced = X_val.dot(ww_log) + bb_log
12    X_prime, yy_prime, dummy = fit_linreg(X=X_reduced, yy=y_val, alpha=10)
13    err = compute_err(X=X_prime, yy=yy_prime, ww = w_prime)
14    print("ERROR Logistic Regression: ", err) # 0.4477788975944975
```

Function added in **ct_support_code.py** in order to use the cost function for logistic regression.

```
1 def fit_logreg(X, yy, alpha):
2     D = X.shape[1]
3     args = (X, yy, alpha)
4     init = (np.zeros(D), np.array(0))
5     ww, bb = minimize_list(logreg_cost, init, args)
6
7     return ww, bb
```

Apply logistic regression to reduce the input:

```
1 def use_logreg(X, yy, alpha, K):
2     mx = np.max(yy)
3     mn = np.min(yy)
4     hh = (mx-mn)/(K+1)
5     thresholds = np.linspace(mn+hh, mx-hh, num=K, endpoint=True)
6     ww_array = []
7     bb_array = []
8     for kk in range(K):
9         labels = yy > thresholds[kk]
10        ww, bb = fit_logreg(X, labels, alpha)
11        ww_array.append(ww)
12        bb_array.append(bb)
13
14    ww = np.column_stack(ww_array)
15    bb = np.column_stack(bb_array)
16
17    return ww, bb
```

I choose to compare the results for Linear Regression with Regularization, PCA with $K = 10$ and Logistic Regression (the input will have 10 columns so we can assume $K = 10$ too). As we can see in table 8, the results using Logistic Regression are better than both, PCA and Linear Regression. The PCA method will produce a new set of features based on linear combinations of the previous features (we may lose important information sometimes

and this will produce very bad results) while the Logistic Regression produce a new set of parameters based on multiple probability classes which can action on the input data changing his dimension but keeping the all the important information and providing better results, but taking more time than PCA to compute.

DATA SET	LINEAR REGRESSION	PCA K=10	LOGISTIC REGRESSION
TRAINING SET	0.369	0.576	0.422
VALIDATION SET	0.474	0.576	0.447

Table 8: Linear Regression vs. PCA (K = 10) vs. Logistic Regression (K = 10)

5 Question 5

Overall function for this question:

```

1 def neural_network(X_train, X_val, y_train, y_val):
2     # Random Initialization
3     params = random_initialization(X_train.shape[1], 10)
4     new_params = fit_nn(params=params, X=X_train, yy=y_train, alpha=10)
5     y_predicted = nn_cost(params=new_params, X=X_train, alpha=10)
6     err = root_mean_square_error(y_expected=y_train, y_predicted=y_predicted)
7     print("ERR Training Random:", err) # 0.10109716170454178
8     y_predicted = nn_cost(params=new_params, X=X_val, alpha=10)
9     err = root_mean_square_error(y_expected=y_val, y_predicted=y_predicted)
10    print("ERR Validation Random:", err) # 0.24988936555719002
11
12    # Q4 Initialization
13    params = q4_initialization(X_train, y_train, 10, 10)
14    new_params = fit_nn(params=params, X=X_train, yy=y_train, alpha=10)
15    y_predicted = nn_cost(params=new_params, X=X_train, alpha=10)
16    err = root_mean_square_error(y_expected=y_train, y_predicted=y_predicted)
17    print("ERR Training Ex4:", err) # 0.1034439502578436
18    y_predicted = nn_cost(params=new_params, X=X_val, alpha=10)
19    err = root_mean_square_error(y_expected=y_val, y_predicted=y_predicted)
20    print("ERR Validation Ex4:", err) # 0.2651730161911833

```

Function which randomly initializes the parameters

```

1 def random_initialization(D, K):
2     V = np.random.randn(K, D)
3     bk = np.random.randn(K)
4     ww = np.random.randn(K)
5     bb = np.random.randn(1)[0]
6
7     return (ww, bb, V, bk)

```

Function which initializes the parameters with the values from question 4.

```

1 def q4_initialization(X, yy, alpha, K):
2     V, bk = use_logreg(X=X, yy=yy, alpha=alpha, K=K)
3     X_reduced = X.dot(V) + bk
4     X_prime, yy_prime, w_prime = fit_linreg(X=X_reduced, yy=yy, alpha=alpha)
5     ww = w_prime
6     bb = 0
7
8     return (ww, bb, V.T, bk.T.reshape(K,))

```

Function added in `ct_support_code.py` in order to use the cost function for this neural network.

```

1 def fit_nn(params, X, yy=None, alpha=None):
2     args = (X, yy, alpha)
3     ww_bar, bb_bar, V_bar, bk_bar = minimize_list(nn_cost, params, args)
4
5     return (ww_bar, bb_bar, V_bar, bk_bar)

```

The neural network using Question 4 parameters initialization produced very good results. Regarding using random initialization of parameters (from a normal distribution), sometimes we can get the best results, as we can see in table 9 and sometimes not, but I think it is worth the risk because even when the results are not the best, they are still comparable with using initialization from Question 4.

DATA SET	LINEAR REGRESSION	PCA K=10	LOGISTIC REGRESSION	NN RANDOM	NN Q4 INITIALIZATION
TRAINING SET	0.369	0.576	0.422	0.101	0.103
VALIDATION SET	0.474	0.576	0.447	0.250	0.265

Table 9: Linear Regression Regularization, PCA, Logistic Regression and Neural Network

6 Question 6

For this question I tried to implement a simple Gradient Descent algorithm because is obviously simple to implement, could converge very fast with a proper learning rate (0.05 in this case) and, as we can see in table 10, produces better results than all other methods, except for Neural Network from Question 5. Even the error value for the test set is good: 0.430. I have calculated the gradients with respect to the weights and bias using the following formulas:

$$\frac{\partial}{\partial m} = \frac{2}{N} \sum_{i=1}^N -x_i(y_i - (mx_i + b))$$

$$\frac{\partial}{\partial b} = \frac{2}{N} \sum_{i=1}^N -(y_i - (mx_i + b))$$

Figure 2: Partial derivatives

DATA SET	LINEAR REG	PCA K=10	PCA K=100	LOGISTIC REG	NN RANDOM	NN Q4	GRADIENT DESCENT
TRAINING SET	0.369	0.576	0.415	0.422	0.101	0.103	0.370
VALIDATION SET	0.474	0.576	0.461	0.447	0.250	0.265	0.432

Table 10: All models

The results using Gradient Descent are display in figure 3. I have also plotted the error for each iteration to make sure that the Gradient it is working correctly.

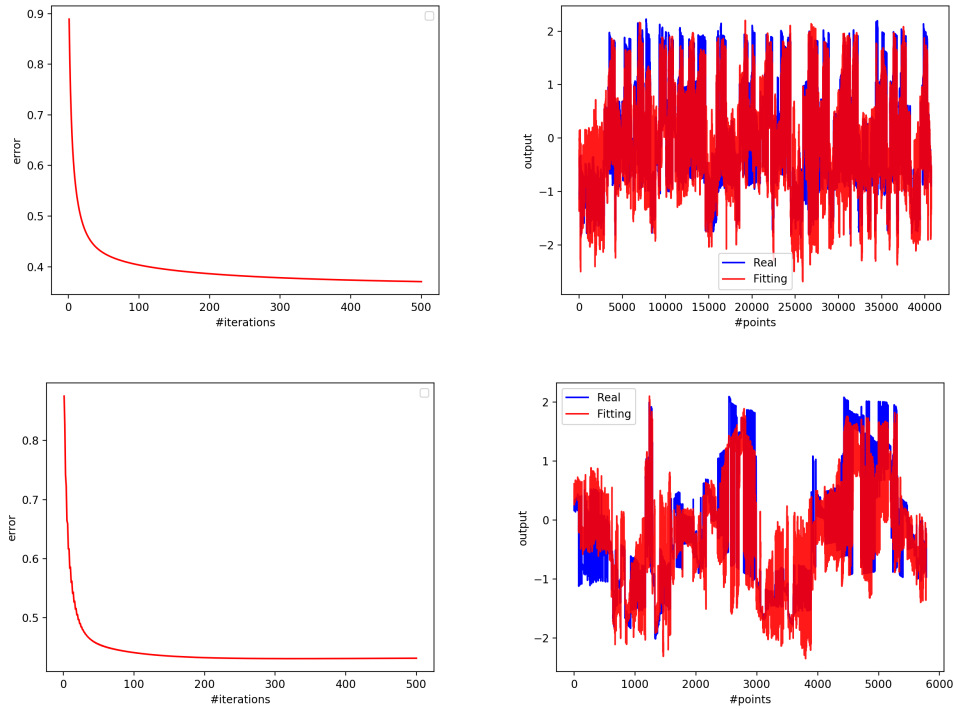


Figure 3: Error Training (top-left), Output Training (top-right), Error Validation (bottom-left), Output Validation (bottom-right)

Overall function for this question:

```
1 def gradient_descent(X_train, X_val, y_train, y_val, X_test, y_test):
2     init = (np.zeros(X_train.shape[1]), np.array(0))
3     ww, bb, params_arr = fit_gd(init, X=X_train, yy=y_train, iter=500,
4     learning_rate=0.05)
5
6     errors = plot_errors(params_arr = params_arr, X=X_train, yy=y_train)
7     plot_output(X=X_train, yy=y_train, ww=ww, bb=bb)
8     print("ERROR GD Training:", errors[-1])           # 0.3704747114536849
9
10    errors = plot_errors(params_arr = params_arr, X=X_val, yy=y_val)
11    plot_output(X=X_val, yy=y_val, ww=ww, bb=bb)
12    print("ERROR GD Validation:", errors[-1])         # 0.43155709194018754
13
14    errors = plot_errors(params_arr = params_arr, X=X_test, yy=y_test)
15    plot_output(X=X_test, yy=y_test, ww=ww, bb=bb)
16    print("ERROR GD Test:", errors[-1])              # 0.4297184947308366
```

The errors plotting function:

```
1 def plot_errors(params_arr, X, yy):
2     errors = []
3     for ww, bb in params_arr:
4         y_predicted = X.dot(ww) + bb
5         err = root_mean_square_error(y_expected=yy, y_predicted=y_predicted)
6         errors.append(err)
7
8     grid = np.linspace(1, len(errors), len(errors))
9
10    plt.clf()
11    plt.plot(grid, errors, 'r-')
12    plt.legend()
13    plt.ylabel("error")
14    plt.xlabel("#iterations")
15    plt.show()
16
17    return errors
```

The output plotting function:

```
1 def plot_output(X, yy, ww, bb=0):
2     y_predicted = X.dot(ww) + bb
3     grid = np.linspace(1, len(yy), len(yy))
4
5     plt.clf()
6     plt.plot(grid, yy, 'b-', label='Real')
7     plt.plot(grid, y_predicted, 'r-', label='Fitting', alpha=0.9)
8     plt.legend()
9     plt.ylabel("output")
10    plt.xlabel("#points")
11    plt.show()
```

These functions were added in `ct_support_code.py`:

The Gradient Descent algorithm:

```
1 def gradient_descent(X, yy, ww, bb, learning_rate):
2     bb_gradient = np.array(0.)
3     ww_gradient = np.zeros(X.shape[1])
4     N = X.shape[0]
5     for i in range(0, N):
6         x = X[i][:]
7         y = yy[i]
8         bb_gradient += -(2/float(N)) * (y - (x.dot(ww) + bb))
9         ww_gradient += -(2/float(N)) * x * (y - (x.dot(ww) + bb))
10    bb_bar = bb - (learning_rate * bb_gradient)
11    ww_bar = ww - (learning_rate * ww_gradient)
12
13    return [ww_bar, bb_bar]
```

The fitting function for Gradient Descent:

```
1 def fit_gd(params, X, yy=None, iter=500, learning_rate=0.01):
2     ww = params[0]
3     bb = params[1]
4     params_arr = []
5     for i in range(iter):
6         print(i)
7         ww, bb = gradient_descent(X, yy, ww, bb, learning_rate)
8         params_arr.append((ww, bb))
9
10    return [ww, bb, params_arr]
```