

Mathematical Foundations of Machine Learning.

Homework 2

Teresa Morales

October 2019

1. Let

$$X = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}.$$

- (a) What is the rank of X ? **The rank of this matrix is 3.** rank of a matrix is the number of linearly independent columns, which has to be equal to the number of linearly independent rows. A collection of vectors $v_1, v_2 \dots v_p \in R^n$ are linearly independent if $\sum_{i=1}^p a_i v_i = 0$ if and only if $a_i = 0$ for all i .

In this case, we have that

$$a_1 + a_2 + a_3 = 0$$

$$a_1 + a_2 = 0$$

$$a_1 = 0$$

This system of equations has only one solution in which $a_1 = 0$, $a_2 = 0$ and $a_3 = 0$, which means that the three vectors are linearly independent. Thus, the rank of X is 3.

- (b) Suppose that $y = Xw$. Derive an explicit formula for w in terms of y .

$$w = X^{-1}y$$

We have to calculate the inverse of matrix X , using the following steps:

First, we estimate the matrix of minors:

$$MM = \begin{bmatrix} 0 & 0 & -1 \\ 0 & -1 & -1 \\ -1 & -1 & 0 \end{bmatrix}.$$

Second, we estimate the matrix of co-factors which coincides with the adjugate because it is symmetrical (equal to its transpose):

$$Adjugate = \begin{bmatrix} 0 & 0 & -1 \\ 0 & -1 & 1 \\ -1 & 1 & 0 \end{bmatrix}.$$

Finally, we divide by the determinant of X which is equal to -1, so we have that:

$$w = X^{-1}y = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & -1 \\ 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} y_3 \\ y_1 - y_3 \\ y_1 - y_2 \end{bmatrix}$$

2. Consider trying to estimate how many calories there are per gram of carbohydrates, fat, and protein in breakfast cereals. The data matrix for this problem is X. Each row contains the grams/serving of carbohydrates, fat, and protein, and each row corresponds to a different cereal (Frosted Flakes, Grape-Nuts, Teenage Mutant Ninja Turtles). The total calories for each cereal are y.

$$X = \begin{bmatrix} 25 & 0 & 1 \\ 20 & 1 & 2 \\ 40 & 1 & 6 \end{bmatrix} \quad y = \begin{bmatrix} 110 \\ 110 \\ 210 \end{bmatrix}$$

- (a) Write a small program (e.g., in Matlab or Python) that solves the system of equations $Xw = y$. Recall the solution w gives the calories/gram of carbohydrate, fat, or protein. What is the solution?

The formula to calculate the estimated weights for these features would be the following:

$$\hat{w} = (X^T X)^{-1} X^T y$$

The code is the following:

```

1 import numpy as np
2 def calc_weights(X, y):
3     Xt = np.transpose(X)
4     XtX = np.matmul(Xt, X)
5     Inverse_XtX = np.linalg.inv(XtX)
6     w = np.matmul(np.matmul(Inverse_XtX, Xt), y)
7     print('Weights are', w)
8     return w
9
10 X0 = [[25, 0, 1],
11        [20, 1, 2],
12        [40, 1, 6]]
13
14 y0 = [110, 110, 210]
15
16 calc_weights(X0, y0)

```

The estimated weights are 4.25 calories per gram of carbohydrates, 17.5 calories per gram of fat and 3.75 calories per gram of protein, as shown in our results:

```
1 Weights are [ 4.25 17.5 3.75]
```

- (b) The solution may not agree with the known calories/gram, which are 4 for carbs, 9 for fat and 4 for protein. We suspect this may be due to rounding the grams to integers, especially the fat grams. Assuming the true value for calories/gram is:

$$w = \begin{bmatrix} 4 \\ 9 \\ 4 \end{bmatrix}$$

and that the total calories, grams of carbs, and grams of protein are correctly reported above, determine the “correct” grams of fat in each cereal.

We know that,

$$y = Xw$$

$$\begin{bmatrix} 110 \\ 110 \\ 210 \end{bmatrix} = \begin{bmatrix} 25 & x & 1 \\ 20 & y & 2 \\ 40 & z & 6 \end{bmatrix} \cdot \begin{bmatrix} 4 \\ 9 \\ 4 \end{bmatrix}$$

Solving this system of equations we have that $x = 2/3 = 0.\overline{66}$, $y = 22/9 = 2.\overline{44}$ **and** $z = 26/9 = 2.\overline{88}$.

- (c) Now suppose that we predict total calories using a more refined breakdown of carbohydrates, into total carbohydrates, complex carbohydrates and sugars (simple carbs). So now we will have 5 features to predict calories (the three carb features + fat and protein). So let's suppose we measure the grams of these features in 5 different cereals to obtain this data matrix:

$$X = \begin{bmatrix} 25 & 15 & 10 \\ 0 & 1 & & & \\ 20 & 12 & 8 & & \\ 1 & 2 & & & \\ 40 & 30 & 10 & & \\ 1 & 6 & & & \\ 30 & 15 & 15 & & \\ 0 & 3 & & & \\ 35 & 20 & 15 & & \\ 2 & 4 & & & \end{bmatrix}$$

$$y = \begin{bmatrix} 104 \\ 97 \\ 193 \\ 132 \\ 174 \end{bmatrix}$$

Can you solve $Xw = y$ Carefully examine the situation in this case. Is there a solution that agrees with the true calories/gram?

If we observe matrix X we can see that we cannot solve this problem since there is some information that is redundant. That is to say, **the columns of X are not linearly independent**. In particular we observe that the first column is the sum of the second and third columns. In order to solve this problem we need to eliminate one of these columns. We eliminate the first one since it would make sense to estimate the weights for each type of carbohydrates (complex and sugars).

Using the same program as above, the additional code to estimate this new problem would be:

```

1
2 X1 = [[15, 10, 0, 1],
3       [12, 8, 1, 2],
4       [30, 10, 1, 6],
5       [15, 15, 0, 3],
6       [20, 15, 2, 4]]
7
8 Y1 = [104, 97, 193, 132, 174]
9
10 calc_weights(X1, Y1)
```

As we can see, in this case, **the solution agrees with the true weights and we have 4 calories per gram for each type of carbohydrate, 9 calories per gram for fats and 4 calories per gram for proteins.**

```

1 Weights are [4. 4. 9. 4.]
```

3. Gradients Calculate the gradients of the following functions with respect to w:

(a) $f(w) = w^T(2x)$

$$\nabla_w f = 2x$$

(b) $f(w) = 3w^T x - 0.5x^T w$

$$\nabla_w f = 3x - 0.5x = 2.5x$$

$$(c) f(w) = w^T \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} w$$

$$\begin{aligned} \nabla_w f &= Qw + Q^T w \\ &= \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \\ &= \begin{bmatrix} 2w_1 + 5w_2 \\ 5w_1 + 8w_2 \end{bmatrix} \end{aligned}$$

$$(d) f(w) = w^T \begin{bmatrix} 1 & 2.5 \\ 2.5 & 4 \end{bmatrix} w$$

$$\nabla_w f = Qw + Q^T w = 2Qw = \begin{bmatrix} 2w_1 + 5w_2 \\ 5w_1 + 8w_2 \end{bmatrix}$$

4. Design classifier to detect if a face image is happy.

Consider the two faces below. It is easy for a human, like yourself, to decide which is happy and which is not. Can we get a machine to do it?

The key to this classification task is to find good features that may help to discriminate between happy and mad faces. What features do we pay attention to? The eyes, the mouth, maybe the brow?

The image below depicts a set of points or “landmarks” that can be automatically detected in a face image (notice there are points corresponding to the eyes, the brows, the nose, and the mouth). The distances between pairs of these points can indicate the facial expression, such as a smile or a frown. We chose $p = 9$ of these distances as features for a classification algorithm. The features extracted from $n = 128$ face images (like the two shown above) are stored in the $n \times p$ matrix \mathbf{X} in the Matlab file `face_emotion_data.mat`. This file also includes an $n \times 1$ binary vector \mathbf{y} ; happy faces are labeled $+1$ and mad faces are labeled -1 . The goal is to find a set of weights for the features in order to predict whether the emotion of a face image is happy or mad.

- (a) Use the training data \mathbf{X} and \mathbf{y} to find an good set of weights.

The python code used to find a set of weights is the following:

```
1 import scipy.io as sio
2 import numpy as np
3
4 file_path = '/Users/teresamorales/Documents/Harris/MFML/
    Homework 2/face_emotion_data.mat'
5 face_data = sio.loadmat(file_path)
6 X = face_data['X']
7 y = face_data['y']
8
9
10 def calc_weights(X, y):
```

```

11 Xt = np.transpose(X)
12 XtX = np.matmul(Xt, X)
13 Inverse_XtX = np.linalg.inv(XtX)
14 w = np.matmul(np.matmul(Inverse_XtX, Xt), y)
15 return w
16
17 w = calc_weights(X, y)

```

The weights obtained are the following:

$$\hat{w} = \begin{bmatrix} 0.944 \\ 0.214 \\ 0.266 \\ -0.392 \\ -0.005 \\ -0.017 \\ -0.166 \\ -0.082 \\ -0.166 \end{bmatrix}$$

- (b) How would you use these weights to classify a new face image as happy or mad?

Let

$$X_{new} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \end{bmatrix}$$

be a vector containing the 9 specific distances of a given new face that was not part of the training sample. Then, **our best prediction of whether the face is happy or mad would be determined by the following equations, using the estimated weights:**

$$\tilde{y} = \text{sign}(\hat{y})$$

$$\hat{y} = \hat{w}X_{new}$$

The following code shows how we would build our classifier and an example that results in a face classified as being mad:

```

1 def sign(num):
2     return -1 if num < 0 else 1
3
4 def classify_face(training_features, labels, new_features)
5     :

```

```

5     w_training= calc_weights(training_features, labels)
6     y_hat = np.matmul(new_features, w_training)
7     y_tilda = [sign(i) for i in y_hat]
8     if y_tilda == 1:
9         print('The face is happy')
10    else:
11        print('The face is mad')
12
13    # Example:
14    x_new = [0.2, 0.8, 0.3,0.2, 0.8, 0.3,0.2, 0.8, 0.3]
15    classify_face(X, y, x_new)
16

```

- (c) Which features seem to be most important? Justify your answer. Given the fact that all features are measuring distances in a face we could assume that these features have the same unit of measure. Similarly, we could assume that they have a similar variance.

Given these assumptions we could pick features by just selecting the ones with the highest absolute value, since they seem to be explaining a bigger share of the outcome (being happy or mad). **The features that have the highest absolute value are 1, 4, 3 and 2 (in that order).**

- (d) Can you design a classifier based on just 3 of the 9 features? Which 3 would you choose? How would you build a classifier?

As explained in the previous section we could pick features with the highest weights (absolute values). We would therefore choose features 1, 4 and 3. Given this selection, **we would need to rebuild our classifier to estimate weights when we are only considering those three features.**

Using the same program as in section a, we would now apply our function to our subset of features.

```

1 X_sel = X[:, [0, 2, 3]]
2 w_sel = calc_weights(X_sel, y)

```

Our new weights are $w_1 = 0.705$, $w_4 = -0.788$, $w_3 = 0.874$. Note that while the sign of the weights has been kept, the relative importance of each of these weights has changed when we have excluded the other features. This is probably due to the fact that our selected features are correlated to other features that we were taking into account so, when we exclude the other distances, part of that information is captured by our new weights.

The new classifier gives the same result for the example given above (the face is classified as being mad).

```

1 X_sel = X[:, [0, 2, 3]]
2 x_new_sel= [0.2, 0.3,0.2]
3 classify_face(X_sel, y, x_new_sel)

```

- (e) A common method for estimating the performance of a classifier is cross-validation (CV). CV works like this. Divide the data set into 8 equal sized subsets (e.g., examples 1 – 16, 17 – 32, etc). Use 7 sets of the data to chose your weights, then use the weights to predict the labels of the remaining “hold-out” set. Compute the number of mistakes made on this hold-out set and divide that number by 16 (the size of the set) to estimate the error rate. Repeat this process 8 times (for the 8 different choices of the hold-out set) and average the error rates to obtain a final estimate.

The python code for cross validation is the following:

```

1 import random
2 import scipy.io as sio
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 def calc_weights(X, y):
7     Xt = np.transpose(X)
8     XtX = np.matmul(Xt, X)
9     Inverse_XtX = np.linalg.inv(XtX)
10    w = np.matmul(np.matmul(Inverse_XtX, Xt), y)
11    return w
12
13 def random_partition(sample_size, n_chunks, chunk_size):
14    sample_size_index = list(range(sample_size))
15    random.Random(222).shuffle(sample_size_index) #
16    randomizes selection setting seed 222
17    return [sample_size_index[round(chunk_size * i):round(
18    chunk_size * (i + 1))]] for i in range(n_chunks)]
19
20 def sign(num):
21    return -1 if num < 0 else 1
22
23 def calc_error_rate(reserved_index, n_chunks, chunk_size,
24    features_matrix, labels_vector):
25    X_reserved = features_matrix[sliced_index[
26    reserved_index]]
27    y_reserved = labels_vector[sliced_index[reserved_index
28    ]]
29    X_training = features_matrix[sum(sliced_index[0:
30    reserved_index] +
31    sliced_index[
32    reserved_index+1:n_chunks], [])]
33    y_training = labels_vector[sum(sliced_index[0:
34    reserved_index] +
35    sliced_index[
36    reserved_index+1:n_chunks], [])]
37    w_training = calc_weights(X_training, y_training)
38    y_hat = np.matmul(X_reserved, w_training)
39    y_tilda = [sign(i) for i in y_hat]
40    return np.sum([y_tilda[i] != y_reserved[i] for i in
41    range(int(chunk_size))])/chunk_size
42
43 # Cross validation with 9 features:
44

```



```

35 sample_size = 128
36 number_of_chunks = 8
37 chunk_size = sample_size/number_of_chunks
38
39 sliced_index = random_partition(sample_size,
40                                 number_of_chunks, chunk_size)
41
42 error_rates = [calc_error_rate(i, number_of_chunks,
43                                chunk_size, X, y)
44                 for i in range(number_of_chunks)] #loop for
45                 8 alternative training samples
46
47 np.mean(error_rates)
48
49 # Cross validation with 3 features (same randomization as
50 # with 9 features):
51
52 X_sel = X[:, [0, 2, 3]]
53
54 error_rates_3features = [calc_error_rate(
55     i, number_of_chunks, chunk_size, X_sel, y) for i in
56     range(number_of_chunks)]
57
58 np.mean(error_rates_3features)

```

- (f) What is the estimated error rate using all 9 features? What is it using the 3 features you chose in (d) above?

The estimated error rate using all 9 features is 0.03125.

The estimated error rate using only 3 features is 0.0547.

The error rate is therefore smaller when we use the whole set of features to predict whether a face is happy or mad.

5. **Polynomial fitting.** Suppose we observe pairs of scalar points (z_i, y_i) , $i = 1, \dots, n$. Imagine these points are measurements from a scientific experiment. The variables z_i are the experimental conditions and the y_i correspond to the measured response in each condition. Suppose we wish to fit a degree $d < n$ polynomial to these data. In other words, we want to find the coefficients of a degree d polynomial p so that $p(z_i) \approx y_i$ for $i = 1, 2, \dots, n$. We will set this up as a least-squares problem.

- (a) Suppose p is a degree d polynomial. Write the general expression for $p(z) = y$.

The general expression would be the following: $y_i = w_0 z_i^0 + w_1 z_i^1 + w_2 z_i^2 + \dots + w_d z_i^d$ for $i = 1, \dots, n$.

- (b) Express the $i = 1, \dots, n$ equations as a system in matrix form $Xw = y$. Specifically, what is the form/structure of X in terms of the given x_i .

X would have the form of the Vandermonde matrix:

$$\begin{bmatrix} z_1^0 & z_1^1 & z_1^2 \\ \dots & z_2^d & \dots \\ z_2^0 & z_2^1 & z_2^2 \\ \dots & z_2^d & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \\ z_n^0 & z_n^1 & z_n^2 \\ \dots & z_n^d & \dots \end{bmatrix} \cdot \begin{bmatrix} w_0 \\ w_1 \\ \dots \\ w_d \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix}$$

$$(n * (d + 1)) * ((d + 1) * 1) = n * 1$$

- (c) Write a Matlab or Python script to find the least-squares fit to the $n = 30$ data points in `polydata.mat`. Plot the points and the polynomial fits for $d = 1, 2, 3$.

The script is the following:

```

1 import numpy as np
2 import scipy.io as sio
3 import matplotlib.pyplot as plt
4
5 file_path = '/Users/teresamorales/Documents/Harris/MFML/
    Homework 2/polydata.mat'
6 poly_data = sio.loadmat(file_path)
7 x = poly_data['x']
8 y = poly_data['y']
9
10
11 def calc_weights(X, y):
12     Xt = np.transpose(X)
13     XtX = np.matmul(Xt, X)
14     Inverse_XtX = np.linalg.inv(XtX)
15     w = np.matmul(np.matmul(Inverse_XtX, Xt), y)
16     return w
17
18
19 V_Matrix_b = np.fromfunction(lambda i, j: x[i] ** j, (30,
    30), dtype=int)
20
21 Weights = {d: calc_weights([V_Matrix_b[i][0][0:d+1] for i
    in range(30)], y)
    for d in (1, 2, 3)}
22
23
24 for d in (1, 2, 3):
25     print('Weights for a polynomial of degree', d, 'are',
    Weights[d])
26
27 #If estimated y's are needed we would calculate them and
    store them in our data set as follows:
28
29 y_est = {d: np.matmul([V_Matrix_b[i][0][0:d+1] for i in
    range(30)], Weights[d])

```

```

30         for d in (1, 2, 3)}
31 poly_data['y_est_1'], poly_data['y_est_2'], poly_data['
    y_est_3'] = y_est[1], y_est[2], y_est[3]
32
33 # Graph data and fitted functions:
34
35 x_graph = np.linspace(0, 1, num=100)
36
37 fx_2 = []
38
39 for i in range(len(x_graph)):
40     fx_2.append(Weights[2][0]+Weights[2][1]*x_graph[i]+
        Weights[2][2]*x_graph[i]**2)
41
42 fx_3 = []
43
44 for i in range(len(x_graph)):
45     fx_3.append(Weights[3][0]+Weights[3][1]*x_graph[i]+
        Weights[3][2]*x_graph[i]**2 + Weights
        [3][3]*x_graph[i]**3)
47
48 fig, ax = plt.subplots(figsize=(12, 6))
49 plt.scatter(x, y)
50 plt.plot([min(x), max(x)], [min(y_est[1]), max(y_est[1])],
    color='red')
51 plt.plot(x_graph, fx_2, color='green')
52 plt.plot(x_graph, fx_3, color='pink')
53 plt.show()

```

Our results are:

```

1 Weights for a polynomial of degree 1 are [[-0.11914368]
2   [ 0.21420653]]
3 Weights for a polynomial of degree 2 are [[ 0.00876394]
4   [-0.53816355]
5   [ 0.70458083]]
6 Weights for a polynomial of degree 3 are [[-0.04418309]
7   [ 0.03303464]
8   [-0.58563589]
9   [ 0.796239   ]]

```

Figure 1 shows graph of our sample (dots) and the three fitted lines. The red line represents a polynomial of degree $d=1$, the green line represents a polynomial of degree $d = 2$ and the pink line a polynomial of degree $d = 3$.

Figure 1:

