

Mathematical Foundations of Machine Learning.
Homework 3

Teresa Morales

October, 28th 2019

Homework 3

Teresa Morales

$$1. a) \begin{cases} \hat{r} = y - X\hat{w} \\ r = y - Xw \end{cases} \quad \begin{cases} y = \hat{r} + X\hat{w} \\ r = (\hat{r} + X\hat{w}) - Xw = \hat{r} + X(\hat{w} - w) \end{cases}$$

$$b) \|r\|^2 = (\hat{r} + X(\hat{w} - w))^T (\hat{r} + X(\hat{w} - w)) \\ = \hat{r}^T \hat{r} + \hat{r}^T X(\hat{w} - w) + (\hat{w} - w)^T X^T \hat{r} + (\hat{w} - w)^T X^T X (\hat{w} - w)$$

$$c) \hat{r} \text{ is orthogonal to the columns of } X \text{ which implies that} \\ X^T \hat{r} = \hat{r}^T X = 0 \rightarrow \|r\|^2 = \hat{r}^T \hat{r} + (\hat{w} - w)^T X^T X (\hat{w} - w)$$

If $X^T X$ is positive definite, then $(\hat{w} - w)^T X^T X (\hat{w} - w) > 0 \forall (\hat{w} - w)$
This would therefore imply that $\|r\|^2 > \|\hat{r}\|^2 \Leftarrow \|r\| > \|\hat{r}\|$
and \hat{w} is the least squares solution.

2. a) We can use Gram-Schmidt orthogonalization:

$$X_1 = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} \quad X_2 = \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix}$$

$$(1^{st}) \text{ we initialize } u_1 = X_1 / \|X_1\|_2 = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} / 3 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

(2nd) We write X_2 as a weighted sum of u_1 plus a residual

$$X_2 = w \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \text{resid}$$

$$w = \argmin_w \|X_2 - u_1 w\|_2^2$$

$$w = (\underbrace{u_1^T u_1}_1)^{-1} u_1^T X_2 = u_1^T X_2$$

$$\text{resid} = \underline{X_2} - \underline{u_1} (u_1^T X_2)$$

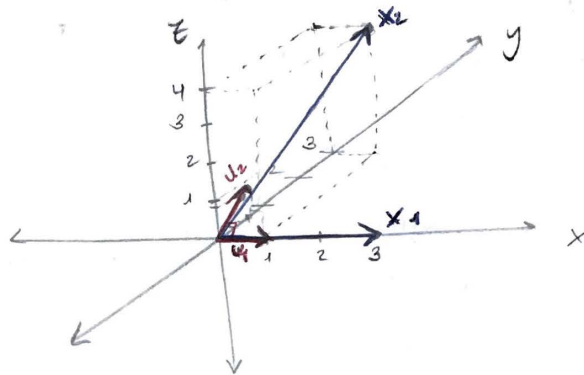
$$= \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} [1 \ 0 \ 0] \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix} - \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 3 \\ 4 \end{bmatrix}$$

$$(8^o) u_2 = \frac{\text{resid}}{\| \text{resid} \|_2} = \begin{bmatrix} 0 \\ 3 \\ 4 \end{bmatrix} / \sqrt{25} = \begin{bmatrix} 0 \\ 3/5 \\ 4/5 \end{bmatrix}$$

$$\Rightarrow \text{the orthonormal vectors are } \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ and } \begin{bmatrix} 0 \\ 3/5 \\ 4/5 \end{bmatrix}$$

b)



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

I can express x as a weighted sum of columns of u , which means that $\hat{y} = Xw = u\tilde{w}$, where

$$\tilde{w} = \underset{w}{\text{argmin}} \| y - u\tilde{w} \|_2^2 = \underbrace{(u^T u)^{-1}}_I u^T y = u^T y$$

$$\begin{aligned} \hat{y} &= u u^T y = \begin{bmatrix} 1 & 0 \\ 0 & 3/5 \\ 0 & 4/5 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 3/5 & 4/5 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 9/25 & 12/25 \\ 0 & 12/25 & 16/25 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1.56 \\ 2.08 \end{bmatrix} \end{aligned}$$

Problem 3

The code to perform Gram-Schmidt orthogonalization is the following:

```
1 import numpy as np
2 import math
3
4 def eliminate_zeros(X):
5     X_nz = X[:, ~np.all(X == 0, axis=0)]
6     return X_nz
7
8
9 def initialize(X, rows): # Calculate Uj for j = 1.
10    U_1 = X[:, 0]/math.sqrt(sum([X[i, 0]**2 for i in range(rows)]))
11    return U_1
12
13
14 def calc_Uj(X, U, rows, j): # Calculate Uj for j = 2, 3...p.
15    X_j = X[:, j-1] - \
16        sum([np.dot(U[:, i], np.dot(np.transpose(U[:, i]), X[:, j
17    -1])) for i in range(j-1)])
18    X_jr = np.array([round(X_j[i], 10) for i in range(rows)])
19    if np.all(X_jr == 0):
20        U_j = X_jr
21    else:
22        U_j = X_j/math.sqrt(sum([X_j[i]**2 for i in range(rows)]))
23    return U_j
24
25 def run_GS(X):
26    X_np = eliminate_zeros(X)
27    n_rows = len(X_np)
28    U_1 = initialize(X=X_np, rows=n_rows)
29    U = np.transpose(np.array([U_1])) # creates structure to
30    append Ujs.
31    n_cols = len(X_np[0])
32    for j in range(2, n_cols+1): # loop to append Ujs for j = 2...
33    p.
34        U_j = calc_Uj(X_np, U, n_rows, j)
35        U_c = np.transpose(U)
36        U = np.transpose(np.vstack([U_c, U_j]))
37    U_final = eliminate_zeros(U)
38    print('U is', U_final)
39    print('The rank of X is', len(U_final[0]))
```

Here are some examples (X1 and X2 correspond to examples given in class):

```
1 X_1 = np.array([[1, 1],
2                 [1, 3],
3                 [0, 0]])
4
5 X_2 = np.array([[1, -1, 2],
6                 [-1, 1, 2]])
7
8 X_3 = np.array([[0, 1, -1, 2, 0, 8],
9                 [1, 3, -1, 1.8, 0, 2],
10                [4, 6, -8, 0, 0, 2],
```

```

11         [2, 1, 0, 3.5, 0, 2],
12         [1, 3, -2.3, 1, 0, 7]])
13
14 run_GS(X_1)
15 U is [[ 0.70710678 -0.70710678]
16        [ 0.70710678  0.70710678]
17        [ 0.          0.          ]]
18 The rank of X is 2
19
20 run_GS(X_2)
21 U is [[ 0.70710678  0.70710678]
22        [-0.70710678  0.70710678]]
23 The rank of X is 2
24
25 run_GS(X_3)
26 U is [[ 0.          0.32522182 -0.06847489  0.92229249 -0.19727777]
27        [ 0.21320072  0.50261554  0.62675823 -0.2381579  -0.50237239]
28        [ 0.85280287  0.05913124 -0.48948081 -0.08874738 -0.14752358]
29        [ 0.42640143 -0.62087802  0.57368199  0.29095535  0.13757275]
30        [ 0.21320072  0.50261554  0.18380103  0.01123671  0.81732123]]
31 The rank of X is 5

```

4. Suppose $A = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$ what is its SVD?

By inspection, we can see that A is a diagonal matrix with positive entries and with $a_{11} \geq a_{22} \geq 0$

Then, the simplest SVD would be the following

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$U = I$$

$$V = I$$

$$\Sigma = A$$

5. $A = \begin{bmatrix} -2 & 0 \\ 0 & -1 \end{bmatrix}$

In this case, A is also a diagonal matrix, but its diagonal entries are negative values. $a_{ii} \leq 0$. However, absolute value of those entries does satisfy $|a_{11}| \geq |a_{22}|$. We could therefore represent A as

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

So the SVD would be $U = I$, $V = -I$ and

$$\Sigma = -A$$

6. $A \in \mathbb{R}^{n \times d}$

a) let $r = \text{rank}(A)$, we can represent A as

$$A = U \Sigma V^T = \begin{bmatrix} u_1 & u_2 & \dots & u_r \\ \vdots & \vdots & & \vdots \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_r \end{bmatrix} \begin{bmatrix} -v_1- \\ -v_2- \\ \vdots \\ -v_r- \end{bmatrix}$$

$(n \times r) \quad (r \times r) \quad (r \times d)$

This would represent the economy singular value decomposition. We can now use the outer-product representation and express this as a sum of rank-1 matrices:

$$A = \begin{bmatrix} u_1 \\ \vdots \\ u_r \end{bmatrix} \begin{bmatrix} \sigma_1 \end{bmatrix} \begin{bmatrix} -v_1- \end{bmatrix} + \begin{bmatrix} u_2 \\ \vdots \\ u_r \end{bmatrix} \begin{bmatrix} \sigma_2 \end{bmatrix} \begin{bmatrix} -v_2- \end{bmatrix} + \dots + \begin{bmatrix} u_r \\ \vdots \\ u_r \end{bmatrix} \begin{bmatrix} \sigma_r \end{bmatrix} \begin{bmatrix} -v_r- \end{bmatrix}$$

$(n \times 1) \quad (1 \times 1) \quad (1 \times d)$

$$A = \begin{bmatrix} u_{11} \sigma_1 v_{11} & u_{11} \sigma_1 v_{12} & \dots & u_{11} \sigma_1 v_{1d} \\ u_{21} \sigma_1 v_{11} & u_{21} \sigma_1 v_{12} & \dots & u_{21} \sigma_1 v_{1d} \\ \vdots & \vdots & & \vdots \\ u_{r1} \sigma_1 v_{11} & u_{r1} \sigma_1 v_{12} & \dots & u_{r1} \sigma_1 v_{1d} \end{bmatrix} + \begin{bmatrix} u_{12} \sigma_2 v_{21} & u_{12} \sigma_2 v_{22} & \dots & u_{12} \sigma_2 v_{2d} \\ u_{22} \sigma_2 v_{21} & u_{22} \sigma_2 v_{22} & \dots & u_{22} \sigma_2 v_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ u_{r2} \sigma_2 v_{21} & u_{r2} \sigma_2 v_{22} & \dots & u_{r2} \sigma_2 v_{2d} \end{bmatrix}$$

$(n \times d) \quad (n \times d)$

$$+ \dots + \begin{bmatrix} u_{1r} \sigma_r v_{r1} & u_{1r} \sigma_r v_{r2} & \dots & u_{1r} \sigma_r v_{rd} \\ u_{2r} \sigma_r v_{r1} & u_{2r} \sigma_r v_{r2} & \dots & u_{2r} \sigma_r v_{rd} \\ \vdots & \vdots & \ddots & \vdots \\ u_{rr} \sigma_r v_{r1} & u_{rr} \sigma_r v_{r2} & \dots & u_{rr} \sigma_r v_{rd} \end{bmatrix}$$

$(n \times d)$

b) Recall rank- k approximations A_k , where $k < r$. Express A_k with the sum of rank one matrices

The subspace approximation theorem tells us that if $A \in \mathbb{R}^{n \times d}$ with rank r , and $k < r$, the best rank k approximation to A is:

$$\min_{Z \in \mathbb{R}^{n \times d}} \|A - Z\|_F^2 = U_k \Sigma_k V_k^T$$

where U_k is formed by the first k columns of U , Σ_k is a k by k diagonal matrix formed by the first k diagonal elements of Σ and V_k^T is formed by the first k rows of V^T .

$$A_k = \begin{bmatrix} u_1 & u_2 & \dots & u_k \end{bmatrix} \begin{bmatrix} \beta_1 & 0 & \dots & 0 \\ 0 & \beta_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \beta_k \end{bmatrix} \begin{bmatrix} -v_1 \\ -v_2 \\ \vdots \\ -v_k \end{bmatrix}$$

$(n \times k) \quad (k \times k) \quad (k \times d)$

We can use the same type of representation and express A as a sum of rank one matrices:

$$A_k = \begin{bmatrix} 1 \\ u_1 \\ 1 \end{bmatrix} \beta_1 \begin{bmatrix} -v_1 \end{bmatrix} + \begin{bmatrix} u_2 \end{bmatrix} \beta_2 \begin{bmatrix} -v_2 \end{bmatrix} + \dots + \begin{bmatrix} u_k \end{bmatrix} \beta_k \begin{bmatrix} -v_k \end{bmatrix}$$

$$A_k = \begin{bmatrix} u_{11} \beta_1 v_{11} & u_{11} \beta_1 v_{12} & \dots & u_{11} \beta_1 v_{1d} \\ u_{21} \beta_1 v_{11} & u_{21} \beta_1 v_{12} & \dots & u_{21} \beta_1 v_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ u_{n1} \beta_1 v_{11} & u_{n1} \beta_1 v_{12} & \dots & u_{n1} \beta_1 v_{1d} \end{bmatrix} + \dots + \begin{bmatrix} u_{1k} \beta_k v_{k1} & u_{1k} \beta_k v_{k2} & \dots & u_{1k} \beta_k v_{kd} \\ u_{2k} \beta_k v_{k1} & u_{2k} \beta_k v_{k2} & \dots & u_{2k} \beta_k v_{kd} \\ \vdots & \vdots & \ddots & \vdots \\ u_{nk} \beta_k v_{k1} & u_{nk} \beta_k v_{k2} & \dots & u_{nk} \beta_k v_{kd} \end{bmatrix}$$

7- a) What we would do is estimate three different models of least squares for each flower species.
for serotia we would have

$$y_s \in \{-1, 1, 4\} \quad \text{and} \quad \hat{y}_s = X w_s \approx y_s$$

$$\hat{w}_s = \underset{w_s}{\text{argmin}} \|y_s - X w_s\|_2^2$$

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

for versicolor we would have

$$y_{vc} \in \{-1, 1, 4\} \quad \text{and} \quad \hat{y}_{vc} = X w_{vc} \approx y_{vc}$$

$$\hat{w}_{vc} = \underset{w_{vc}}{\text{argmin}} \|y_{vc} - X w_{vc}\|_2^2$$

$$\hat{w}_{vc} = (X^T X)^{-1} X^T y_{vc}$$

for virginica we would have

$$y_{vir} \in \{-1, 1, 4\} \quad \text{and} \quad \hat{y}_{vir} = X w_{vir} \approx y_{vir}$$

$$\hat{w}_{vir} = \underset{w_{vir}}{\text{argmin}} \|y_{vir} - X w_{vir}\|_2^2$$

Problem 7

The code for sections b and c is the following:

```
1 import numpy as np
2 import math
3 import random
4 import pandas as pd
5 import scipy.io as sio
6 import matplotlib.pyplot as plt
7
8 file_path = '/Users/teresamorales/Documents/Harris/MFML/Homework 3/
9   fisheriris.mat'
10 fisher_data = sio.loadmat(file_path)
11 Species = fisher_data['species']
12 X = fisher_data['meas']
13 len(X[0]) # 4 columns
14 len(X) # 150 rows
15
16 def create_dummies(name_species, n):
17     y = np.array([1 if Species[i, 0] == name_species else -1 for i
18   in range(n)])
19     return y
20
21 def calc_weights(X, y):
22     Xt = np.transpose(X)
23     XtX = np.matmul(Xt, X)
24     Inverse_XtX = np.linalg.inv(XtX)
25     w = np.matmul(np.matmul(Inverse_XtX, Xt), y)
26     return w
27
28 def divide_randomly(sample_size, training_size, seed):
29     sample_size_index = list(range(sample_size))
30     random.Random(seed).shuffle(sample_size_index)
31     return [sample_size_index[0:training_size], sample_size_index[
32   training_size:]]
33
34 def sign(num):
35     return -1 if num < 0 else 1
36
37 def calc_error_rate(splited_index, features_matrix, labels_vector):
38     X_reserved = features_matrix[splited_index[1]]
39     y_reserved = labels_vector[splited_index[1]]
40     X_training = features_matrix[splited_index[0]]
41     y_training = labels_vector[splited_index[0]]
42     w_training = calc_weights(X_training, y_training)
43     y_hat = np.matmul(X_reserved, w_training)
44     y_tilda = [sign(i) for i in y_hat]
45     return np.sum([y_tilda[i] != y_reserved[i] for i in range(len(
46   splited_index[1]))])/len(splited_index[1])
47
48 #To run the program, experimenting with different training sample
49 sizes, we do the following:
```

```

49
50 n_repetitions = 50
51 training_size = [120, 115, 110, 105, 100, 95, 90, 85, 80, 75,
52                  70, 65, 60, 55, 50, 45, 40, 35, 30, 25, 20, 15,
                    10, 5, 3]
53
54 # (1 ) Setosa:
55 y_setosa = create_dummies('setosa', 150)
56
57 error_rates_setosa = {i: [] for i in training_size}
58
59 for size in training_size:
60     for i in range(n_repetitions):
61         index = divide_randomly(150, size, seed=i)
62         error = calc_error_rate(index, X, y_setosa)
63         error_rates_setosa[size].append(error)
64
65 mean_error_setosa = {i: np.mean(error_rates_setosa[i]) for i in
                    training_size}
66
67 # (2 ) Versicolor:
68
69 y_versicolor = create_dummies('versicolor', 150)
70
71 error_rates_versicolor = {i: [] for i in training_size}
72
73 for size in training_size:
74     for i in range(n_repetitions):
75         index = divide_randomly(150, size, seed=i)
76         error = calc_error_rate(index, X, y_versicolor)
77         error_rates_versicolor[size].append(error)
78
79 mean_error_versicolor = {i: np.mean(error_rates_versicolor[i]) for
                    i in training_size}
80
81 # (3 ) Virginica:
82
83 y_virginica = create_dummies('virginica', 150)
84
85 error_rates_virginica = {i: [] for i in training_size}
86
87 for size in training_size:
88     for i in range(n_repetitions):
89         index = divide_randomly(150, size, seed=i)
90         error = calc_error_rate(index, X, y_virginica)
91         error_rates_virginica[size].append(error)
92
93 mean_error_virginica = {i: np.mean(error_rates_virginica[i]) for i
                    in training_size}
94
95 # Making a graph
96
97 setosa = {'training_size': list(mean_error_setosa.keys()),
98          'mean_error_setosa': list(mean_error_setosa.values())}
99 flours_df = pd.DataFrame(data=setosa)
100
101 flours_df['mean_error_versicolor'] = list(mean_error_versicolor.

```

```

    values())
102
103 flours_df['mean_error_virginica'] = list(mean_error_virginica.
    values())
104
105 fig, ax = plt.subplots(figsize=(12, 6))
106 plt.plot('training_size', 'mean_error_setosa', data=flours_df,
    color='purple')
107 plt.plot('training_size', 'mean_error_versicolor', data=flours_df,
    color='red')
108 plt.plot('training_size', 'mean_error_virginica', data=flours_df,
    color='green')
109 plt.show()

```

b) What is the average test error (number of mistakes divided by 30)? We can print the average test error for our three models, with sample sizes of 120, for our 50 repetitions as follows:

```

1 >>> flours_df.loc[0,:]
2
3 training_size          120.000000
4 mean_error_setosa      0.000000
5 mean_error_versicolor  0.288667
6 mean_error_virginica   0.123333
7 Name: 0, dtype: float64

```

We observe that the error rate is very different for each type of flour with a much smaller error for setosa, a slightly bigger error for virginica and a bigger error rate for versicolor.

c) Experiment with even smaller sized training sets. Clearly we need at least one training example from each type of flower. Make a plot of average test error as a function of training set size.

We observe that when we use smaller training sets the error tends to get bigger for the three species and the order that we just described is maintained so that the error rate for setosa is always the smallest and for versicolor the largest as we can see in the following graph.

Figure 1:



d) Now design a classifier using only the first three measurements (sepal length, sepal width, and petal length). What is the average test error in this case?

We repeat the process for a different selection of features as follows:

```

1 X_sel = X[:, 0:3]
2
3 # (1 ) Setosa:
4
5 error_rates_setosa_sel = {i: [] for i in training_size}
6
7 for size in training_size:
8     for i in range(n_repetitions):
9         index = divide_randomly(150, size, seed=i)
10        error = calc_error_rate(index, X_sel, y_setosa)
11        error_rates_setosa_sel[size].append(error)
12
13 mean_error_setosa_sel = {i: np.mean(error_rates_setosa_sel[i]) for
14 i in training_size}
15
16 # (2 ) Versicolor:
17
18 error_rates_versicolor_sel = {i: [] for i in training_size}
19
20 for size in training_size:
21     for i in range(n_repetitions):
22         index = divide_randomly(150, size, seed=i)
23         error = calc_error_rate(index, X_sel, y_versicolor)
24         error_rates_versicolor_sel[size].append(error)
25
26 mean_error_versicolor_sel = {i: np.mean(error_rates_versicolor_sel[
27 i]) for i in training_size}
28

```

```

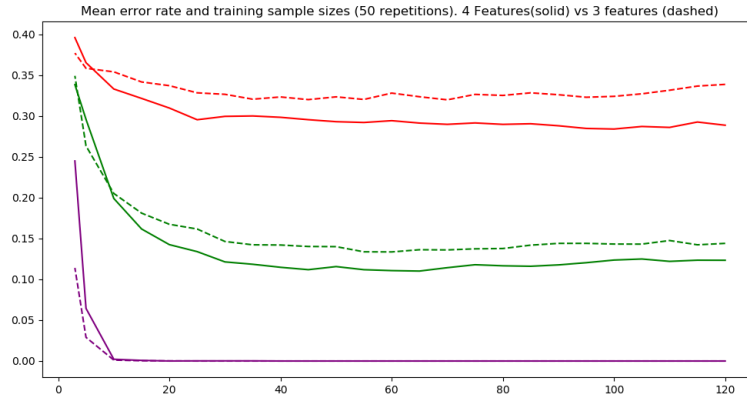
27 # (3 ) Virginica:
28
29 error_rates_virginica_sel = {i: [] for i in training_size}
30
31 for size in training_size:
32     for i in range(n_repetitions):
33         index = divide_randomly(150, size, seed=i)
34         error = calc_error_rate(index, X_sel, y_virginica)
35         error_rates_virginica_sel[size].append(error)
36
37 mean_error_virginica_sel = {i: np.mean(error_rates_virginica_sel[i
38 ]) for i in training_size}
39
40 # Making a graph
41
42 flours_df['mean_error_setosa_sel'] = list(mean_error_setosa_sel.
43 values())
44 flours_df['mean_error_versicolor_sel'] = list(
45     mean_error_versicolor_sel.values())
46 flours_df['mean_error_virginica_sel'] = list(
47     mean_error_virginica_sel.values())
48
49 fig, ax = plt.subplots(figsize=(12, 6))
50 plt.plot('training_size', 'mean_error_setosa', data=flours_df,
51         color='purple')
52 plt.plot('training_size', 'mean_error_versicolor', data=flours_df,
53         color='red')
54 plt.plot('training_size', 'mean_error_virginica', data=flours_df,
55         color='green')
56 plt.plot('training_size', 'mean_error_setosa_sel', data=flours_df,
57         color='purple', linestyle='--')
58 plt.plot('training_size', 'mean_error_versicolor_sel', data=
59     flours_df, color='red', linestyle='--')
60 plt.plot('training_size', 'mean_error_virginica_sel', data=
61     flours_df, color='green', linestyle='--')
62 plt.show()

```

We observe that the error rates increase when we use 3 features instead of three. In particular for the flour setosa, with the same training sample size (120) we obtain a mean error rate of 0 again, for virginica we go from 0.12 with 4 features to 0.14 with 3 features and for versicolor we go from 0.29 to 0.34.

As we can see in the graph, this increase happens for all flours except setosa and for most training sample sizes. However, when the training sample sizes get smaller, the error rates tend to get closer and we even observe points where the average error rate with 3 features is smaller than the average error rate with 4 features.

Figure 2:



The following tables show some selected values of this graph. For each species we show the error rates using 4 features (second column) and 3 features (third column):

```

1 >>> pd.set_option('display.max_columns', None)
2 >>> flours_df.loc[[0,4,12, 18, 24],['training_size',
3     mean_error_setosa', 'mean_error_setosa_sel']]
4     training_size  mean_error_setosa  mean_error_setosa_sel
5 0                120              0.000000              0.000000
6 4                100              0.000000              0.000000
7 12               60              0.000000              0.000000
8 18               30              0.000167              0.000000
9 24                3              0.244898              0.114014
10 >>> flours_df.loc[[0,4,12, 18, 24],['training_size',
11     mean_error_virginica', 'mean_error_virginica_sel']]
12     training_size  mean_error_virginica  mean_error_virginica_sel
13 0                120              0.123333              0.144000
14 4                100              0.123600              0.143200
15 12               60              0.110667              0.133556
16 18               30              0.121333              0.146333
17 24                3              0.338367              0.349116
18 >>> flours_df.loc[[0,4,12, 18, 24],['training_size',
19     mean_error_versicolor', 'mean_error_versicolor_sel']]
20     training_size  mean_error_versicolor  mean_error_versicolor_sel
21 0                120              0.288667              0.338667
22 4                100              0.284000              0.324000
23 12               60              0.294222              0.328000
24 18               30              0.299500              0.326500
25 24                3              0.395918              0.377007
26 >>>

```