



Trinity College Dublin
Coláiste na Tríonóide, Baile Átha Cliath
[The University of Dublin](#)

School of Computer Science and Statistics

CS7CS4 Machine Learning **Week 2 Assignment**

Boris Flesch
20300025

October 22, 2020

MSc Computer Science, Intelligent Systems

1 Downloaded Dataset

```
# id:1-1-1
```

2 Questions

- (a) (i) Read data from CSV file (code taken from the assignment sheet) and visualise it using Matplotlib (Figure 1).

```

1      df = pd.read_csv("week2.csv", comment="#")
2      X1 = df.iloc[:,0]
3      X2 = df.iloc[:,1]
4      X = np.column_stack((X1,X2))
5      y = df.iloc[:,2]
6
7      X_m1 = X[np.where(y == -1)]
8      X_p1 = X[np.where(y == 1)]
9      plt.scatter(X_m1[:, 0], X_m1[:, 1], c='r', marker='+',
10                  label="y = -1")
11      plt.scatter(X_p1[:, 0], X_p1[:, 1], c='b', marker='+',
12                  label="y = 1")
```

Note: instead of creating `X_m1` and `X_p1` variables, it would have been possible to scatter all points in a single line as following:

```

1      plt.scatter(X1, X2, c=np.where(y == -1, 'r', 'b'),
2                  marker='+')
```

However this method does not allow to plot the label of each different color (i.e. what red crosses and blue crosses refer to).

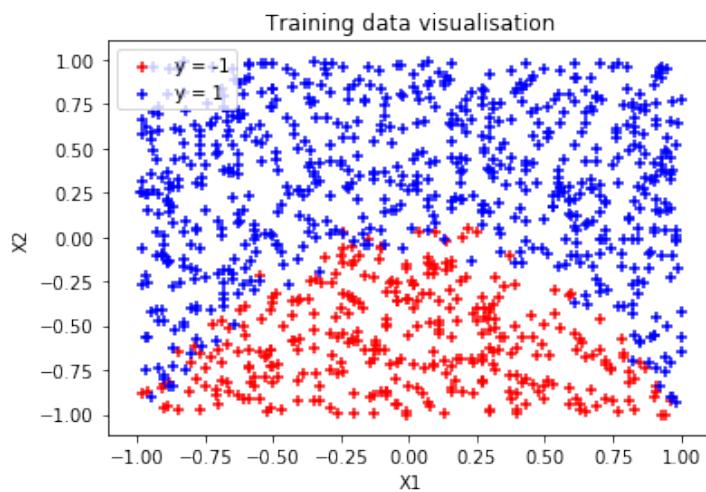


Figure 1

- (ii) Train a logistic regression classifier on the data:

```

1   from sklearn.linear_model import LogisticRegression
2   model = LogisticRegression().fit(X, y)
3   theta = [model.intercept_[0], model.coef_[0][0],
4           ↵ model.coef_[0][1]]

```

$$\theta = \begin{bmatrix} 1.6713466864527504 \\ -0.03177179625806771 \\ 5.045678045815783 \end{bmatrix} \quad (1)$$

An issue with this model is that it is of dimension 2 and will therefore define a line to separate outputs. However, our training data is not linearly separable as seen on Figure 1.

With the parameter values obtained, another issue is the neglect of feature X_1 due to the value of θ_1 which is close to zero. This indicates that the prediction will essentially rely on the value of feature X_2 . Considering $\theta_1 = 0$, we roughly obtain :

$$\text{sign}(\theta^T X) = +1 \quad (2)$$

$$\Leftrightarrow 5.04 \cdot X_2 > -1.67 \quad (3)$$

$$\Leftrightarrow X_2 > -0.33 \quad (4)$$

- (iii) Figure 2 below shows logistic regression classifier predictions and its decision boundary over the training data. Our model is $h_\theta(x) = \text{sign}(\theta^T X)$, the decision boundary can therefore be defined by the equation $\theta^T X = 0$, so that all +1 predictions appear above this line and all -1 predictions appear below it. To plot this decision boundary on the graph with X_1 and X_2 respectively on x- and y-axis, we have to express X_2 as a function of X_1 :

$$\theta^T X = 0 \quad (5)$$

$$\Leftrightarrow \theta_0 + \theta_1 X_1 + \theta_2 X_2 = 0 \quad (6)$$

$$\Leftrightarrow X_2 = \frac{-1}{\theta_2} (\theta_0 + \theta_1 X_1) \quad (7)$$

This can be plotted using the following code:

```

1   # Training data
2   plt.scatter(X_m1[:, 0], X_m1[:, 1], c='r', marker='+', s=500,
3               ↵ label="Training y = -1")
3   plt.scatter(X_p1[:, 0], X_p1[:, 1], c='b', marker='+', s=500,
4               ↵ label="Training y = 1")
4
5   # Predictions
6   X_pred_m1 = X[np.where(model.predict(X) == -1)]
7   X_pred_p1 = X[np.where(model.predict(X) == 1)]
8   plt.scatter(X_pred_m1[:, 0], X_pred_m1[:, 1], marker='o',
9               ↵ facecolor='none', edgecolor='r', s=500, label="Prediction
10              ↵ y = -1")

```

```

9     plt.scatter(X_pred_p1[:, 0], X_pred_p1[:, 1], marker='o',
10    ↳ facecolor='none', edgecolor='b', s=500, label="Prediction
11    ↳ y = 1")
12
13    # Decision boundary
14    plt.plot(X1, -1/theta[2] * (theta[0] + theta[1]*X1),
15    ↳ linewidth=5, c='g', label='Decision boundary')

```

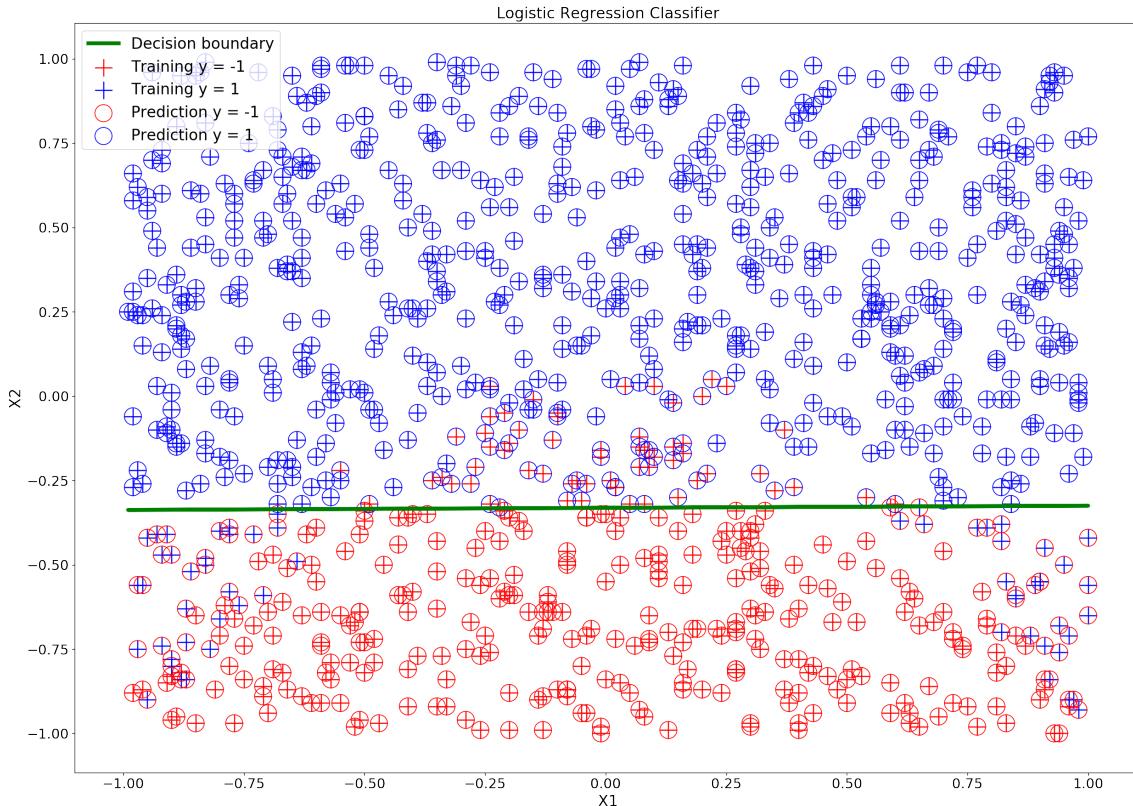


Figure 2

This representation of predictions over training data is quite interesting as it allow us to easily distinguish correct and missed predictions: basically, when a cross is not of the same color as its surrounding circle, the prediction failed. The color then allow us to determine if the missed prediction is a false positive or false negative.

- (iv) The trained logistic regression classifier gives a straight decision boundary line close to -0.33 on y-axis, as expected and explained in (a)(ii). We can therefore see that the predictions obtained are very approximate. Predictions that are very distant from the decision boundary are quite accurate, unlike the ones close to the decision boundary and in the edges. The linear shape of the prediction model does not follow the quadratic behaviour of the data: this is **under-fitting**.
- (b) (i) The LinearSVC function from sklearn allows us to easily train SVM classifiers given a C value as following:

```

1   from sklearn.svm import LinearSVC
2   Cs = [0.001, 1, 1000]

```

```

3   for Ci in Cs:
4       model = LinearSVC(C=Ci).fit(X, y)
5       theta = [model.intercept_[0], model.coef_[0][0],
6       ↪ model.coef_[0][1]]
7       print("C = %.3f" % Ci, "\n theta =", theta)

```

The following parameter values have been obtained (and will be discussed in (b)(iii)):

- For $C = 0.001$ (see Figure 3):

$$\theta = \begin{bmatrix} 0.2189386198207775 \\ 0.007272583598406376 \\ 0.48180974288823664 \end{bmatrix} \quad (8)$$

- For $C = 1$ (see Figure 4):

$$\theta = \begin{bmatrix} 0.5899032183960152 \\ -0.01362192025379305 \\ 1.8196765208475205 \end{bmatrix} \quad (9)$$

- For $C = 1000$ (see Figure 5):

$$\theta = \begin{bmatrix} -0.060980250342034914 \\ -0.6821202341260354 \\ 2.5697912854521707 \end{bmatrix} \quad (10)$$

- (ii) To plot the data, we can use a code similar to the one used in (a)(iii), which gives the following figures:

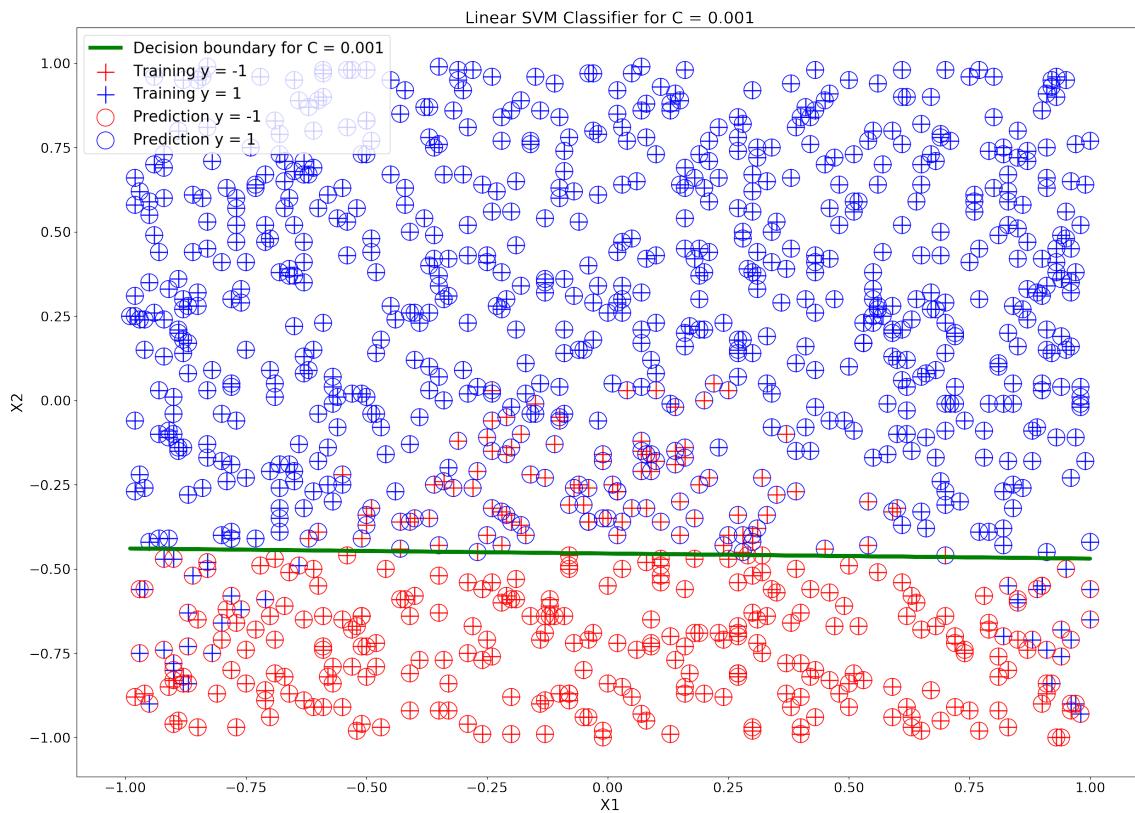


Figure 3

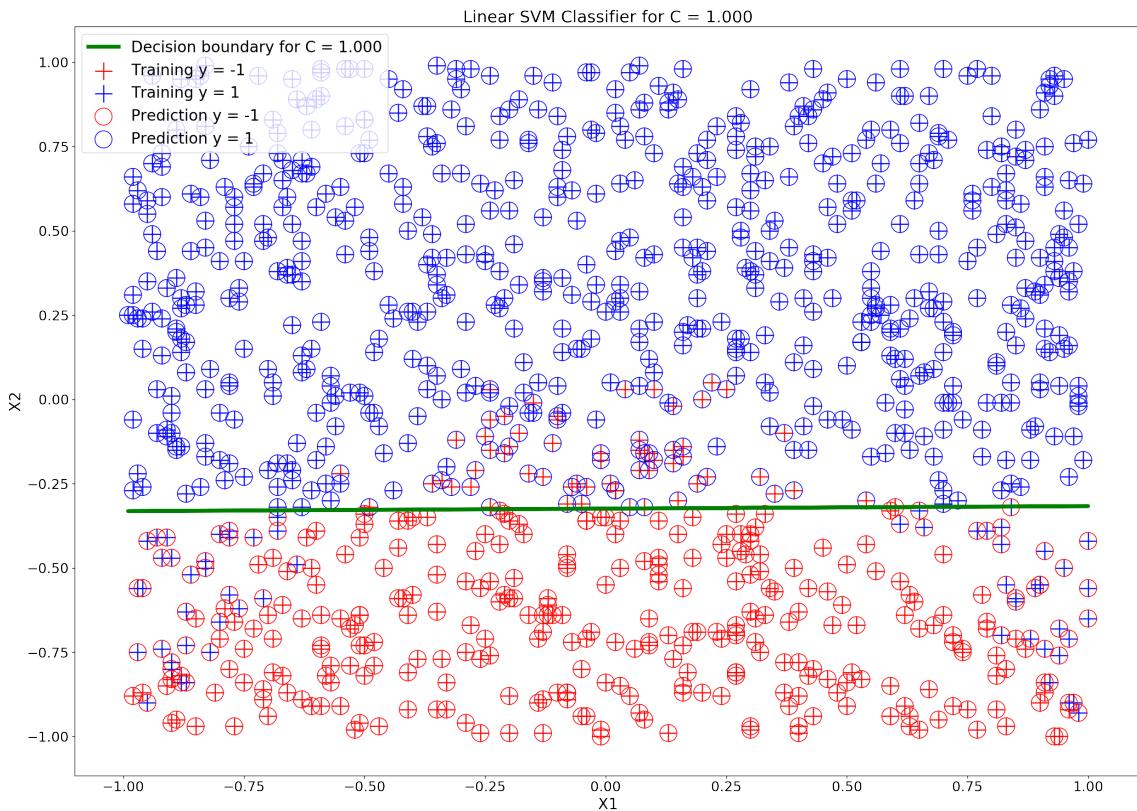


Figure 4

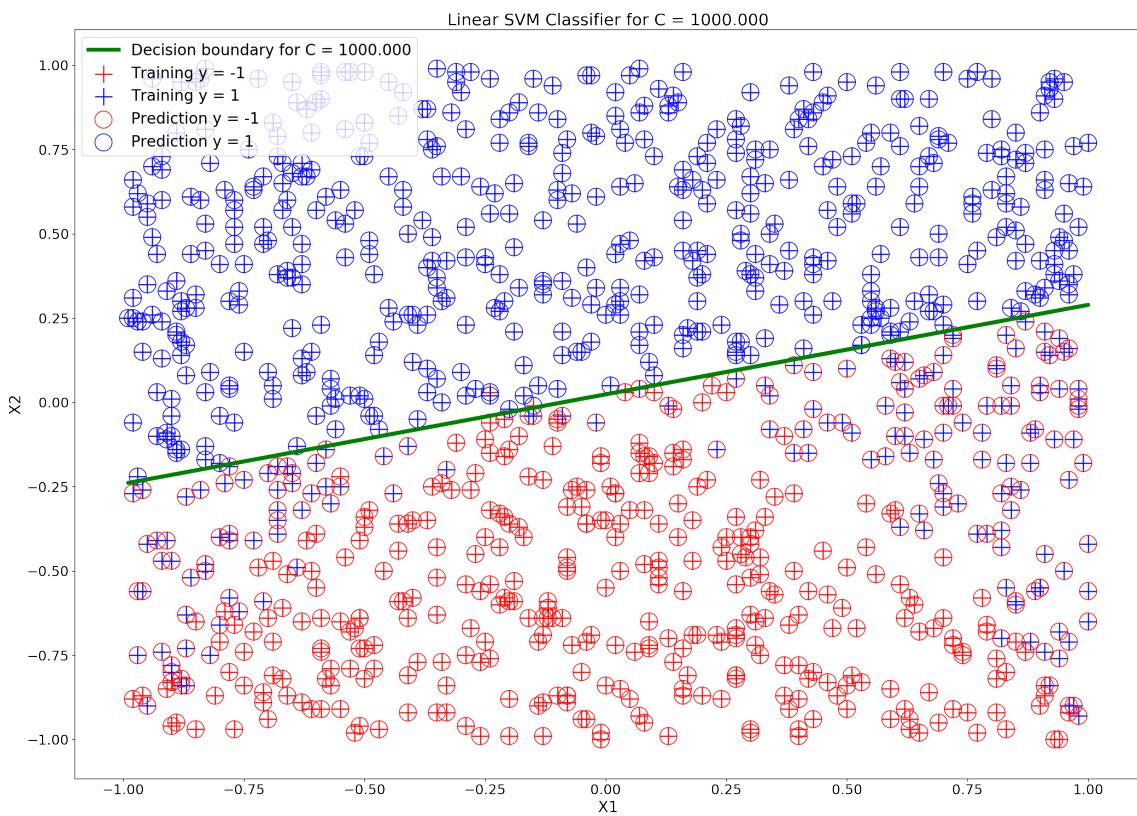


Figure 5

- (iii) Changing the weighting parameter C influences the impact of the penalty on the hinge-loss function of the SVM model.

This penalty should only penalise large values of θ ; however with a too small C value, θ will systematically get penalised, leading to an inaccurate model. That is why we obtain so small parameter values for $C = 0.001$ (especially for θ_3) compared to $C = 1$ and $C = 1000$.

On the contrary, a too big value (i.e. 1000) can lead to an insignificant penalty and therefore some difficulties for the cost function to converge (i.e. *ConvergenceWarning* generated by sklearn).

This explains why the trained model for $C = 1$ (i.e. not an extreme value) is the most accurate; although it remains a linear model and can not fit well the training data especially in the edges.

- (c) (i) Generating additional features and training a logistic regression model can be done with the following Python code:

```

1      X3 = np.power(X1, 2)
2      X4 = np.power(X2, 2)
3      X = np.column_stack((X1,X2,X3,X4))
4
5      model = LogisticRegression().fit(X, y)
6      theta = np.array([model.intercept_[0], model.coef_[0][0],
    ↴ model.coef_[0][1], model.coef_[0][2], model.coef_[0][3]])

```

The trained logistic regression classifier gives the following parameter values:

$$\theta = \begin{bmatrix} 0.29322948 \\ -0.09209744 \\ 6.72896399 \\ 6.48226436 \\ -0.84717209 \end{bmatrix} \quad (11)$$

This 4-dimension model generated thanks to feature engineering allows us to obtain a curved prediction boundary rather than the straight lines obtained until now. In addition to that, the additional feature $X3$ (i.e. $X1^2$) balances the use of feature $X1$ in our prediction as it has been systematically neglected in our previous models. Thus all features are now nearly equally used in this new logistic regression model.

- (ii) The plot showing predicted values and training data can be seen on Figure 6. Adding the square of each initial feature leads to a more reliable model that fits well the training data and therefore produces more accurate predictions.
- (iii) Our training data is composed of 335 outputs -1 and 663 outputs 1. Considering a baseline predictor that always predicts the most common class (i.e. 1 in this case) is therefore a good idea as it gives an accuracy of 66.4% on the training data. Its parameter values are the following:

$$\theta = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad (12)$$

The decision boundary of this baseline model is also shown on Figure 6.

To evaluate and compare these models, we can compare their respective cost function value after calculating them using Python:

```

1     def calcCost(X, y, theta):
2         m = X.shape[0]
3         Xtmp = np.insert(X, 0, 1, axis=1).T
4         cost = 1/m * sum(np.log(1 + np.exp(-y * np.dot(theta,
5             ↵ Xtmp)))/np.log(2))
6         return cost

```

Note: this cost function value calculator automatically adds feature $X_0 = 1$ to the given X data in order to calculate the cost function value properly. That prevents from altering the original value of X in the rest of the program.

The value of the cost function for the logistic regression model and the baseline model are the following:

$$J(\theta) = 0.179947 \quad (13)$$

$$J(\theta_{\text{baseline}}) = 0.936212 \quad (14)$$

Although our baseline model does not seem to be totally inaccurate (as it predicts the most common value of our training data), we can see that $J(\theta_{\text{baseline}}) > J(\theta)$, thus confirming the better accuracy of the logistic regression model.

- (iv) It is possible to use the same principle as in (a)(iii) to obtain the decision boundary equation for this logistic regression model. Starting from $\theta^T X = 0$, we have to get X_2 as a function of the other features. Furthermore features X_3 and X_4 can be substituted by X_1^2 and X_2^2 respectively. This gives us the following equation:

$$\theta^T X = 0 \quad (15)$$

$$\Leftrightarrow \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \theta_3 X_3 + \theta_4 X_4 = 0 \quad (16)$$

$$\Leftrightarrow X_2 = -\frac{1}{\theta_2}(\theta_0 + \theta_1 X_1 + \theta_3 X_3 + \theta_4 X_4) \quad (17)$$

$$\Leftrightarrow X_2 = -\frac{1}{\theta_2}(\theta_0 + \theta_1 X_1 + \theta_3 X_1^2 + \theta_4 X_2^2) \quad (18)$$

$$\Leftrightarrow X_2 = -\frac{\theta_3}{\theta_2} X_1^2 - \frac{\theta_1}{\theta_2} X_1 - \frac{\theta_0}{\theta_2} - \frac{\theta_4}{\theta_2} X_2^2 \quad (19)$$

$$\Leftrightarrow X_2 = ax^2 + bx + c \quad (20)$$

With $a = -\frac{\theta_3}{\theta_2}$, $b = -\frac{\theta_1}{\theta_2}$ and $c = -\frac{\theta_0}{\theta_2} - \frac{\theta_4}{\theta_2} X_2^2$.

The following Python code has been written to plot the decision boundaries of the logistic regression model and the baseline model:

```

1      X1_sorted = np.sort(X1)
2      X2_sorted = np.sort(X2)
3      a = -theta[3]/theta[2]
4      b = -theta[1]/theta[2]
5      c = -theta[0]/theta[2] - theta[4] / theta[2] *
       ↵ np.power(X2_sorted, 2)
6      boundary = a * np.power(X1_sorted, 2) + b * X1_sorted + c
7
8      plt.plot(X1_sorted, boundary, linewidth=5, c='g',
       ↵ label="Decision boundary")
9      plt.axhline(-baselineTheta[0], c='y', linewidth=5,
       ↵ label='Baseline model')

```

Note: sorting X_1 and X_2 features has been required to avoid lines crossing from one side of the curve to the other on the figure.

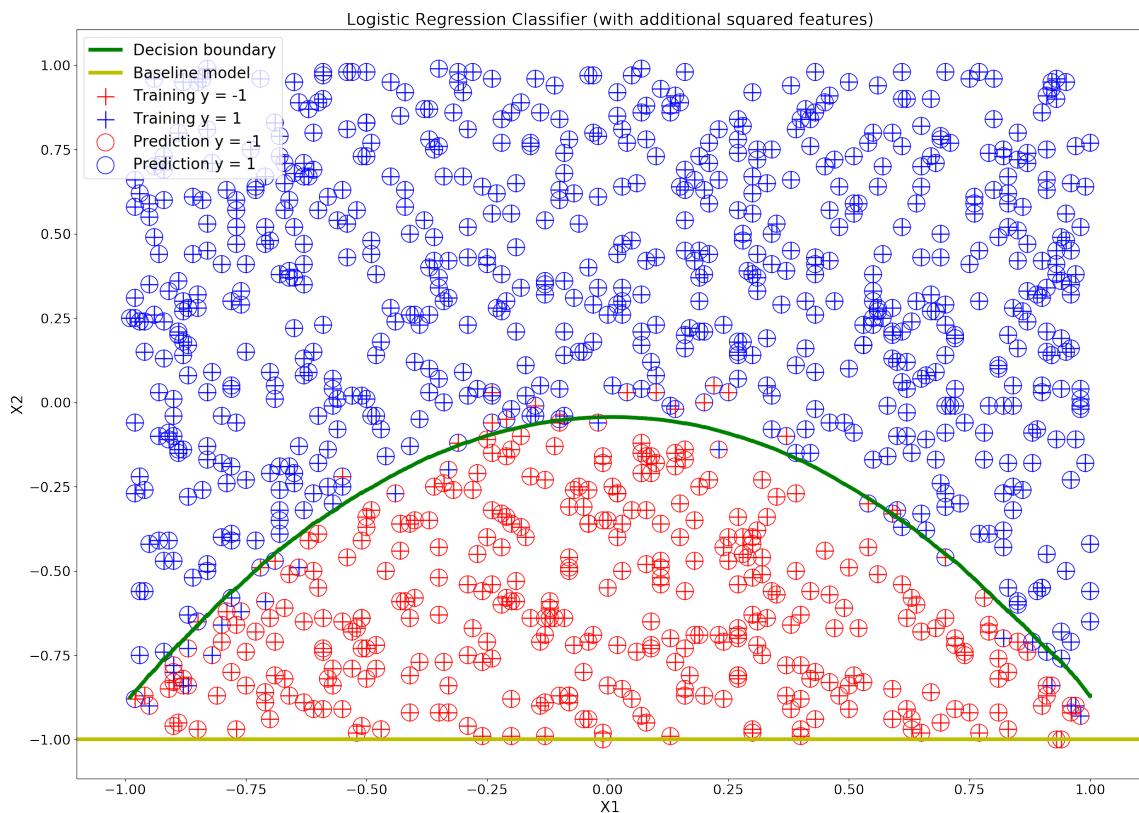


Figure 6

A Appendix

A.1 Python Code

```

1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # CS7CS4/CSU44061 Machine Learning
5  # Week 2 Assignment
6  # Boris Flesch (20300025)
7  #
8  # Downloaded dataset
9  # id:1-1-1
10
11 # Read data
12 import numpy as np
13 import pandas as pd
14 df = pd.read_csv("week2.csv", comment="#")
15 print(df.head())
16 X1 = df.iloc[:,0]
17 X2 = df.iloc[:,1]
18 X = np.column_stack((X1,X2))
19 y = df.iloc[:,2]
20
21 # (a)(i)
22 import matplotlib.pyplot as plt
23
24 X_m1 = X[np.where(y == -1)]
25 X_p1 = X[np.where(y == 1)]
26
27 plt.scatter(X_m1[:, 0], X_m1[:, 1], c='r', marker='+', label="y = -1")
28 plt.scatter(X_p1[:, 0], X_p1[:, 1], c='b', marker='+', label="y = 1")
29 plt.title("Training data visualisation")
30 plt.legend()
31 plt.xlabel("X1")
32 plt.ylabel("X2")
33 plt.show()
34
35 # (a)(ii) Use sklearn to train a logistic regression classifier on the
36 # → data and report the parameter values of the tra
37 from sklearn.linear_model import LogisticRegression
38 model = LogisticRegression().fit(X, y)
39 theta = [model.intercept_[0], model.coef_[0][0], model.coef_[0][1]]
40 print("theta = ", theta)
41
42 # (a)(iii)
43 plt.figure(figsize=(25,18), dpi=120)
44 plt.rcParams.update({'font.size': 18})

```

```

44
45 # Training data
46 plt.scatter(X_m1[:, 0], X_m1[:, 1], c='r', marker='+', s=500,
   ↪ label="Training y = -1")
47 plt.scatter(X_p1[:, 0], X_p1[:, 1], c='b', marker='+', s=500,
   ↪ label="Training y = 1")
48
49 # Predictions
50 X_pred_m1 = X[np.where(model.predict(X) == -1)]
51 X_pred_p1 = X[np.where(model.predict(X) == 1)]
52 plt.scatter(X_pred_m1[:, 0], X_pred_m1[:, 1], marker='o',
   ↪ facecolor='none', edgecolor='r', s=500, label="Prediction y = -1")
53 plt.scatter(X_pred_p1[:, 0], X_pred_p1[:, 1], marker='o',
   ↪ facecolor='none', edgecolor='b', s=500, label="Prediction y = 1")
54
55 # Decision boundary
56 plt.plot(X1, -1/theta[2] * (theta[0] + theta[1]*X1), linewidth=5, c='g',
   ↪ label='Decision boundary')
57
58 plt.legend(scatterpoints=1, fontsize=20)
59 plt.title("Logistic Regression Classifier", fontsize=20)
60 plt.xlabel("X1", fontsize=20)
61 plt.ylabel("X2", fontsize=20)
62 plt.show()
63
64 # (b)(i)
65 from sklearn.svm import LinearSVC
66 Cs = [0.001, 1, 1000]
67
68 for Ci in Cs:
69     model = LinearSVC(C=Ci).fit(X, y)
70     theta = [model.intercept_[0], model.coef_[0][0], model.coef_[0][1]]
71     print("C = %.3f" % Ci, "\n theta =", theta)
72
73     plt.figure(figsize=(25,18), dpi=120)
74     plt.rcParams.update({'font.size': 18})
75
76     # Training data
77     plt.scatter(X_m1[:, 0], X_m1[:, 1], c='r', marker='+', s=500,
   ↪ label="Training y = -1")
78     plt.scatter(X_p1[:, 0], X_p1[:, 1], c='b', marker='+', s=500,
   ↪ label="Training y = 1")
79
80     # Predictions
81     X_pred_m1 = X[np.where(model.predict(X) == -1)]
82     X_pred_p1 = X[np.where(model.predict(X) == 1)]

```

```

83     plt.scatter(X_pred_m1[:, 0], X_pred_m1[:, 1], marker='o',
84                  facecolor='none', edgecolor='r', s=500, label="Prediction y =
85                  -1")
86
87     plt.scatter(X_pred_p1[:, 0], X_pred_p1[:, 1], marker='o',
88                  facecolor='none', edgecolor='b', s=500, label="Prediction y = 1")
89
90     # Decision boundary
91     plt.plot(X1, -1/theta[2] * (theta[0] + theta[1]*X1), linewidth=5,
92               c='g', label="Decision boundary for C = %.3f"%Ci)
93
94     plt.legend(scatterpoints=1, fontsize=20)
95     plt.title("Linear SVM Classifier for C = %.3f"%Ci)
96     plt.xlabel("X1", fontsize=20)
97     plt.ylabel("X2", fontsize=20)
98     plt.show()
99
100
101    #(c)(i)
102    # Additional features
103    X3 = np.power(X1, 2)
104    X4 = np.power(X2, 2)
105    X = np.column_stack((X1,X2,X3,X4))
106
107    model = LogisticRegression().fit(X, y)
108    theta = np.array([model.intercept_[0], model.coef_[0][0],
109                     model.coef_[0][1], model.coef_[0][2], model.coef_[0][3]])
110    print("theta =", theta)
111
112    #(c)(ii)
113    plt.figure(figsize=(25,18), dpi=120)
114    plt.rcParams.update({'font.size': 18})
115
116    # Training data
117    plt.scatter(X_m1[:, 0], X_m1[:, 1], c='r', marker='+', s=500,
118                label="Training y = -1")
119    plt.scatter(X_p1[:, 0], X_p1[:, 1], c='b', marker='+', s=500,
120                label="Training y = 1")
121
122    # Predictions
123    X_pred_m1 = X[np.where(model.predict(X) == -1)]
124    X_pred_p1 = X[np.where(model.predict(X) == 1)]
125
126    plt.scatter(X_pred_m1[:, 0], X_pred_m1[:, 1], marker='o',
127                facecolor='none', edgecolor='r', s=500, label="Prediction y = -1")
128    plt.scatter(X_pred_p1[:, 0], X_pred_p1[:, 1], marker='o',
129                facecolor='none', edgecolor='b', s=500, label="Prediction y = 1")
130
131    # (c)(iii)
132    def calcCost(X, y, theta):
133        """
134
```

```

122     Calculates the cost function value
123
124     :param X: Features X (without  $X_0 = 1$  for each row)
125     :param y: Outputs y
126     :param theta: Theta parameter values
127     :return: Cost function
128     """
129
130     m = X.shape[0]
131     Xtmp = np.insert(X, 0, 1, axis=1).T
132     cost = 1/m * sum(np.log(1 + np.exp(-y * np.dot(theta,
133         ↴ Xtmp))))/np.log(2))
134     return cost
135
136
137 # Cost function for the logistic regression model
138 print("J(theta) = %f\n"%calcCost(X, y, theta))
139
140 # Baseline model
141 (values,counts) = np.unique(y, return_counts=True)
142 mostCommonValue = values[np.argmax(counts)] # prints the most frequent
143     ↴ element
144 baselineTheta = [mostCommonValue, 0, 0, 0, 0]
145 print("theta_baseline =", baselineTheta)
146 print("J(theta_baseline) = %f"%calcCost(X, y, baselineTheta))
147
148 # (c)(iv)
149 # Decision boundary
150 X1_sorted = np.sort(X1)
151 X2_sorted = np.sort(X2)
152 a = -theta[3]/theta[2]
153 b = -theta[1]/theta[2]
154 c = -theta[0]/theta[2] - theta[4] / theta[2] * np.power(X2_sorted, 2)
155 boundary = a * np.power(X1_sorted, 2) + b * X1_sorted + c
156
157 plt.plot(X1_sorted, boundary, linewidth=5, c='g', label="Decision
158     ↴ boundary")
159 plt.axhline(-baselineTheta[0], c='y', linewidth=5, label='Baseline
160     ↴ model')
161 plt.legend(scatterpoints=1, fontsize=20)
162 plt.title("Logistic Regression Classifier (with additional squared
163     ↴ features)")
164 plt.xlabel("X1", fontsize=20)
165 plt.ylabel("X2", fontsize=20)
166 plt.show()

```