

CS-233 Theoretical Exercise

1. Logistic regression can be used to

- a) predict tomorrow's temperature from previous daily temperature observation.
- b) predict whether the object is glass or stone given its reflectance.
- c) predict someone's weight based on their height.
- d) predict the probability of raining given the humidity, wind force and temperature.
- e) tell apart a dog from a cat given height and weight measurements.

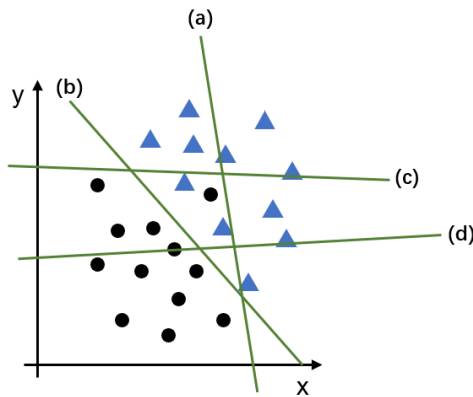
Solution: b), d) and e).

2. Which of the following statements are true:

- a) Linear regression cannot be used for classification problems.
- b) Logistic regression is a linear classifier and can only separate classes using linear decision boundaries.
- c) The output of logistic regression is the estimated probability of the sample belonging to a specific class.

Solution: b) and c).

3. In the following figure, the dots and triangles are samples from two different classes. Which line is the most likely decision boundary obtained by logistic regression?



Solution: (b)

4. Consider a dataset with the four data points shown in Table 1. Assume $x^{(1)}$ and $x^{(2)}$ are two measured biochemical indicators of patients, and $y = 0$ and $y = 1$ indicate the patients without and with a specific symptom, respectively. We want to build a logistic regression model to predict whether a patient has the symptom based on the input features $x^{(1)}$ and $x^{(2)}$. The prediction model is expressed as

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{x}), \quad (1)$$

| Data Point | $x^{(1)}$ | $x^{(2)}$ | y |
|------------|-----------|-----------|-----|
| 1 | -10 | -20 | 0 |
| 2 | 0 | -10 | 0 |
| 3 | 10 | 0 | 1 |
| 4 | 20 | 10 | 1 |

Table 1: Data points.

Algorithm 1: Gradient Descent

Given the step size η ;
Initialize \mathbf{w}_0 ;
while *not converged* **do**
 $\Delta \mathbf{w} = \sum_{i=1}^4 (\hat{y}_i - y_i) \mathbf{x}_i$;
 $\mathbf{w}_k = \mathbf{w}_{k-1} - \eta \Delta \mathbf{w}$;
end

where σ is the sigmoid function, $\mathbf{w} = [w^{(0)}, w^{(1)}, w^{(2)}]^T$, $\mathbf{x} = [1, x^{(1)}, x^{(2)}]^T$.

(1) Write down the algorithm (pseudo-code) that uses gradient descent to compute the optimal \mathbf{w} .

Solution:

(2) Perform one iteration of the previous algorithm with an initialization of $\mathbf{w} = [1, 1, 1]^T$ and a step size of 0.1.

Solution:

(2)

$$\mathbf{w}_0 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad (3)$$

$$\hat{y}_1 = \sigma(\mathbf{w}_0^T \mathbf{x}_1) = \sigma(-29) \approx 0, \quad (4)$$

$$\hat{y}_2 = \sigma(\mathbf{w}_0^T \mathbf{x}_2) = \sigma(-9) \approx 0.00012, \quad (5)$$

$$\hat{y}_3 = \sigma(\mathbf{w}_0^T \mathbf{x}_3) = \sigma(11) \approx 0.99998, \quad (6)$$

$$\hat{y}_4 = \sigma(\mathbf{w}_0^T \mathbf{x}_4) = \sigma(31) \approx 1, \quad (7)$$

$$\Delta \mathbf{w} = 0.00012 \begin{bmatrix} 1 \\ 0 \\ -10 \end{bmatrix} - 0.00002 \begin{bmatrix} 1 \\ 10 \\ 0 \end{bmatrix} \quad (8)$$

$$= \begin{bmatrix} 0.0001 \\ -0.0002 \\ -0.0012 \end{bmatrix} \quad (9)$$

$$\mathbf{w}_1 = \mathbf{w}_0 - 0.1 \Delta \mathbf{w} \quad (10)$$

$$= \begin{bmatrix} 0.99999 \\ 1.00002 \\ 1.00012 \end{bmatrix} \quad (11)$$

(3) Assume \mathbf{w}^* is the optimal solution obtained after the full gradient descent algorithm (not just a single step) in Question (2). Now, we divide $x^{(1)}$ and $x^{(2)}$ by 100 and perform the gradient descent algorithm again with the scaled data. This results in a model with different parameters \mathbf{w}' . Given the test data of new patients $\{\mathbf{x}_5, \dots, \mathbf{x}_N\}$, will the two classifiers defined by \mathbf{w}^* and \mathbf{w}' produce different results for them? Justify your answer mathematically. Note that you need to also scale the test data when using the classifier defined by \mathbf{w}' .

Solution: No. Let $\mathbf{w}^* = [w^{(0)*}, w^{(1)*}, w^{(2)*}]^T$ be the optimal solution for the original data from Question (2), i.e., $\mathbf{x} = [1, x^{(1)}, x^{(2)}]^T$ without scaling, obtained after the full gradient descent algorithm. That is, we have

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} R(\mathbf{w}, \mathbf{X}) , \quad (12)$$

$$= \arg \min_{\mathbf{w}} - \sum_{i=1}^4 y_i \ln(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1 - y_i) \ln(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) . \quad (13)$$

Let us then define the minimum loss $R^* = R(\mathbf{w}^*, \mathbf{X})$, obtained by evaluating the training objective with the optimal parameters.

Furthermore, let $\mathbf{x}' = [1, x^{(1)}/100, x^{(2)}/100]^T$ be the data scaled by a factor 100. Then, we can define a parameter vector $\mathbf{w}' = [w^{(0)*}, 100w^{(1)*}, 100w^{(2)*}]^T$ such that $R(\mathbf{w}', \mathbf{X}') = R(\mathbf{w}^*, \mathbf{X}) = R^*$. This means that \mathbf{w}' is the optimal solution for the scaled data.

For $\mathbf{x}_j \in \{\mathbf{x}_5, \dots, \mathbf{x}_N\}$, we will have

$$\hat{y} = \sigma(\mathbf{w}^{*T} \mathbf{x}_j) = \sigma(\mathbf{w}'^T \mathbf{x}'_j), \quad (14)$$

which indicates that the classifiers \mathbf{w}^* and \mathbf{w}' will produce the same result for \mathbf{x}_j .

(4) If we switch the meaning of the y value, i.e., 0 and 1 now indicating with and without the symptom, respectively, and train a logistic regression model on the resulting data, will the model produce different results for $\{\mathbf{x}_5, \dots, \mathbf{x}_N\}$ than before switching?

Solution: Although switching the value of y will lead to a model with different parameters, the prediction of whether the patient has the symptom or not will not change (i.e., correct predictions will remain correct, and incorrect predictions will remain incorrect).