

CS-233 Theoretical Exercise 4

1 When will CFF trains break down?

You are given the task to predict whether a CFF train will break down or not under certain weather conditions. The dataset, represented by $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$, has $N = 1000$ data entries, and is explained as follows: For each data entry $(\mathbf{x}^{(i)}, y^{(i)})$,

- $\mathbf{x}^{(i)}$ has 5 features: Train line, time of the day, temperature of the day, precipitation of the day, and maximum wind speed of the day.
- $y^{(i)}$ is either 0, which means that the train reaches its final destination with a delay of less than 3 minutes (negative - not broken down), or 1 if the train is delayed by more than 3 minutes or is cancelled (positive - broken down).

The dataset has 900 cases where $y^{(i)} = 0$ and 100 cases where $y^{(i)} = 1$. On this dataset, your model has the following performance:

- Among the 900 cases where the train doesn't break down, i.e., $y^{(i)} = 0$, you successfully predicted 800 cases;
- Among the 100 cases where the train breaks down, i.e., $y^{(i)} = 1$, you successfully predicted 10 cases.

Now let us evaluate your prediction model!

Question 1: Draw the confusion matrix.

Question 2: What is the accuracy of the model?

Question 3: What is the precision of the model?

Question 4: What is the recall of the model?

Question 5: What is the F1 score of the model?

Question 6: Do you think this model is a good model or a bad model ? Explain your reasoning.

Solution: We have the following confusion matrix:

| | Pos | Neg |
|-----|-------------|-------------|
| Yes | TP=10 | FP=100 |
| No | FN=90 | TN=800 |
| | P=TP+FN=100 | N=FP+TN=900 |

As a result, accuracy $ACC = (TP + TN) / (P + N) = 810 / 1000 = 0.810$, $Precision = TP / (TP + FP) = 0.091$, $Recall = TP / P = 0.100$, $FP_rate = FP / N = 0.110$, $F1_score = 2(Precision * Recall) / (Precision + Recall) = 0.095$.

The model is not a good model because it fails to recognize most positive cases. Furthermore the algorithm can be outperformed by simply random guessing.

2 Multiclass classification

Your boss wants to differentiate a canceled train from a delayed train. Therefore, the dataset is now labeled with three classes:

- 0 - a train on time;
- 1 - a train delayed;
- 2 - a train canceled.

For this task, you decide to use a least-square classifier with a one-hot encoding of the label. The prediction for the entire dataset (same dataset as in the previous question) can be formalized in matrix form as

$$\hat{\mathbf{Y}} = \mathbf{X} \cdot \mathbf{W}.$$

Question 1. Write down the one-hot encoding of each class.

Question 2. What are the shapes of $\hat{\mathbf{Y}}$, \mathbf{X} and \mathbf{W} ? Note that we append a 1 to each input when performing classification to account for the bias.

Question 3. Write down the loss function for the MSE loss. Given a learning rate η , what is the update of a single gradient descent step? When the gradient descent algorithm converges (i.e., the gradient goes to 0), what is the final \mathbf{W} ? (Hint: Solve for each column vector of \mathbf{W} .)

Solution: The representations for {train on time, train delayed, train canceled} are, respectively, $[1, 0, 0]^\top, [0, 1, 0]^\top, [0, 0, 1]^\top$.

The shapes of $\hat{\mathbf{Y}}$, \mathbf{X} and \mathbf{W} are, respectively, 1000×3 , 1000×6 and 6×3 .

The loss function is ($N = 1000$)

$$L = \frac{1}{N} \sum_{i=1}^N \|\mathbf{W}^\top \mathbf{x}_i - \mathbf{y}_i\|_2^2.$$

For each column vector $j = 1, 2, 3$, the gradient is

$$\nabla_{\mathbf{w}_{(j)}} L(\mathbf{w}_{(j)}) = \frac{2}{N} \sum_{i=1}^N (\mathbf{w}_{(j)}^\top \mathbf{x}_i - y_{i(j)}) \cdot \mathbf{x}_i = \frac{2}{N} \mathbf{X}^\top \mathbf{X} \mathbf{w}_{(j)} - \frac{2}{N} \mathbf{X}^\top \mathbf{Y}_{(j)}.$$

As a result, the update step for $j = 1, 2, 3$ is

$$\mathbf{w}_{(j)} \leftarrow \mathbf{w}_{(j)} - \alpha \nabla_{\mathbf{w}_{(j)}} L(\mathbf{w}_{(j)}) = \mathbf{w}_{(j)} - \alpha \left(\frac{2}{N} \mathbf{X}^\top \mathbf{X} \mathbf{w}_{(j)} - \frac{2}{N} \mathbf{X}^\top \mathbf{Y}_{(j)} \right).$$

This can also be written in matrix form as

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \left(\frac{2}{N} \mathbf{X}^\top \mathbf{X} \mathbf{W} - \frac{2}{N} \mathbf{X}^\top \mathbf{Y} \right).$$

When the algorithm converges, the gradient is zero. Solving

$$\frac{2}{N} \mathbf{X}^\top \mathbf{X} \mathbf{W} - \frac{2}{N} \mathbf{X}^\top \mathbf{Y} = 0$$

gives us the optimal weight matrix

$$\mathbf{W} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}.$$

Note that this is identical to the solution on Slide 35 in the lecture. The converged weight matrix is the optimal solution, because the MSE loss is a convex function.

3 Connection between binary and multi-class logistic regression

Consider a binary logistic regression task. Your model (model 1) using a sigmoid function has a weight vector \mathbf{w} . If you now use another model (model 2) with a one-hot representation and a soft-max function, is there a weight matrix \mathbf{W} such that model 2 has the same decision boundary as model 1? If not, explain your reasoning; if yes, compute \mathbf{W} .

Solution: Yes. $\mathbf{W} = [\mathbf{0}, \mathbf{w}]$. This can be verified by comparing the softmax function and the sigmoid function: For model 2 with a one-hot encoding, we have

$$\hat{y}^{(0)}(x) = \frac{e^{\mathbf{W}_0^\top x}}{e^{\mathbf{W}_0^\top x} + e^{\mathbf{W}_1^\top x}} = \frac{e^{0^\top x}}{e^{0^\top x} + e^{\mathbf{w}^\top x}} = 1 - \sigma(\mathbf{w}^\top x)$$

and

$$\hat{y}^{(1)}(x) = \frac{e^{\mathbf{W}_1^\top x}}{e^{\mathbf{W}_0^\top x} + e^{\mathbf{W}_1^\top x}} = \frac{e^{\mathbf{w}^\top x}}{e^{0^\top x} + e^{\mathbf{w}^\top x}} = \sigma(\mathbf{w}^\top x).$$

This is exactly the same output as for model 1.