# NYCU Pattern Recognition, Homework 2

## Part. 1, Coding (70%):

In this coding assignment, you are requested to implement 1) logistic Regression and 2) Fisher's Linear Discriminant by using only Numpy, then train your model on the provided dataset and finally evaluate the performance on testing data. Please train your logistic regression model using Gradient Descent, not the closed-form solution.

### (20%) Logistic Regression Model

Requirements:
- Use Gradient Descent
- Use CE (Cross-Entropy) as your loss function.
- Use Softmax for this multiclass classification task.

Criteria:

1. (0%) Show the learning rate, epoch, and batch size that you used.

```python
# For Q1
lr = 0.1
batch_size = 16
epoch = 1000

logistic_reg = MultiClassLogisticRegression()
logistic_reg.fit(X_train, y_train, lr=lr, batch_size=batch_size, epochs=epoch)
✓  4.7s
```
```
[[ 0.60340734 -0.01977367  2.41636633]
 [ 1.2290092  -0.97770969  2.74870049]
 [ 3.01373641  0.9910362  -1.00477261]]
```

2. (5%) What's your training accuracy?

```python
# For Q2
print('Training acc: ', logistic_reg.evaluate(X_train, y_train))
✓  0.0s
```
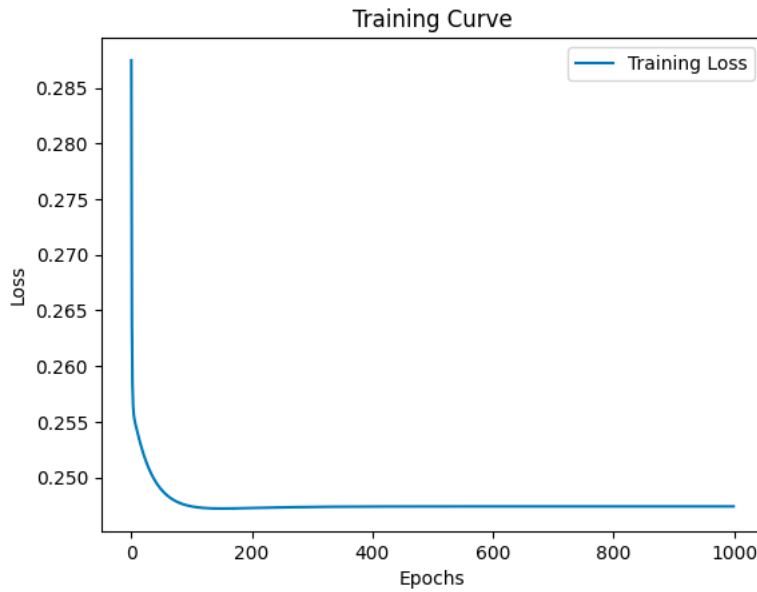```
ans : 0.896
Training acc:  0.896
```

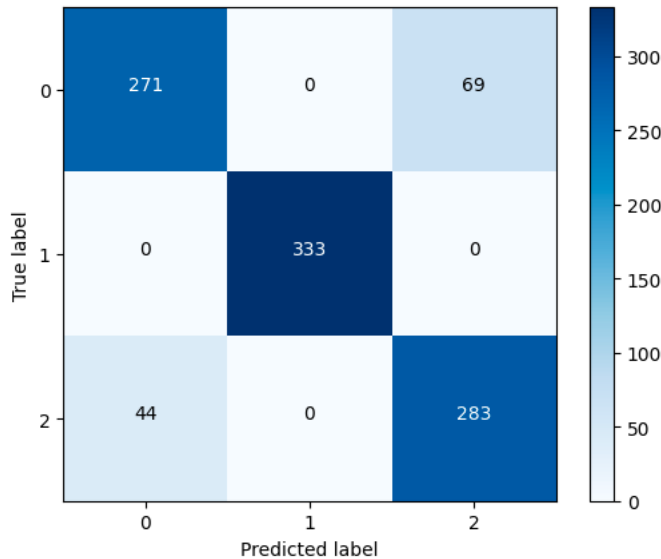3. (5%) What's your testing accuracy?

```python
# For Q3
print('Testing acc: ', logistic_reg.evaluate(X_test, y_test))
✓  0.0s
```
```
ans : 0.887
Testing acc:  0.887
```

4. (5%) Plot the learning curve of the training. (x-axis=epoch, y-axis=loss)

Training Curve

5. (5%) Show the confusion matrix on testing data.



# (30%) Fisher's Linear Discriminant (FLD) Model

Requirements:

- Use FLD to reduce the dimension of the data from 2 to 1.

Criteria:

6. (2%) Compute the mean vectors $m_i$ (i=1, 2, 3) of each class on training data.

```
# For Q6
print("Class mean vector: ", fld.mean_vectors)
✓ 0.0s
Class mean vector:  [array([-4.17505764,  6.35526804]), array([-9.43385176, -4.87830741]), array([-2.54454008,  7.53144179])]
```

7. (2%) Compute the within-class scatter matrix $S_W$ on training data.

```
# For Q7
print("Within-class scatter matrix SW: ", fld.sw)
✓ 0.0s

Within-class scatter matrix SW:  [[1052.70745046  -12.5828441 ]
 [ -12.5828441    971.29686189]]
```

8. (2%) Compute the between-class scatter matrix $S_B$ on <u>training data.</u>

```
# For Q8
print("Between-class scatter matrix SB: ", fld.sb)
✓ 0.0s

Between-class scatter matrix SB:  [[ 8689.12907035 16344.86572983]
 [16344.86572983 31372.93949414]]
```
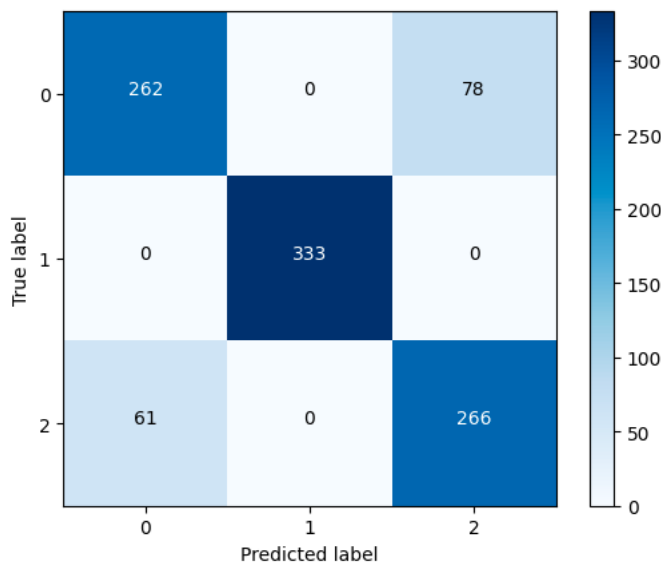
9. (4%) Compute the Fisher's linear discriminant $w$ on <u>training data.</u>

```
# For Q9
print("W: ", fld.w)

✓ 0.0s

W:  [[-0.44115384]
 [-0.8974315 ]]
```

10. (8%) Project the <u>testing data</u> to get the prediction using the shortest distance to the class m ean. Report the accuracy score and draw the confusion matrix on <u>testing data.</u>

FLD using class mean, accuracy:  0.861



11. (8%) Project the <u>testing data</u> to get the prediction using <u>K-Nearest-Neighbor</u>. Compare the accuracy score on the <u>testing data</u> with K values from 1 to 5.

```
# For Q11
y_pred_k1 = fld.predict_using_knn(X_train, y_train, X_test, k=1)
print("FLD using knn (k=1), accuracy: ", fld.accuracy_score(y_test, y_pred_k1))

y_pred_k2 = fld.predict_using_knn(X_train, y_train, X_test, k=2)
print("FLD using knn (k=2), accuracy: ", fld.accuracy_score(y_test, y_pred_k2))

y_pred_k3 = fld.predict_using_knn(X_train, y_train, X_test, k=3)
print("FLD using knn (k=3), accuracy: ", fld.accuracy_score(y_test, y_pred_k3))

y_pred_k4 = fld.predict_using_knn(X_train, y_train, X_test, k=4)
print("FLD using knn (k=4), accuracy: ", fld.accuracy_score(y_test, y_pred_k4))

y_pred_k5 = fld.predict_using_knn(X_train, y_train, X_test, k=5)
print("FLD using knn (k=5), accuracy: ", fld.accuracy_score(y_test, y_pred_k5))
✓ 7.2s

FLD using knn (k=1), accuracy:  0.822
FLD using knn (k=2), accuracy:  0.819
FLD using knn (k=3), accuracy:  0.843
FLD using knn (k=4), accuracy:  0.84
FLD using knn (k=5), accuracy:  0.862
```
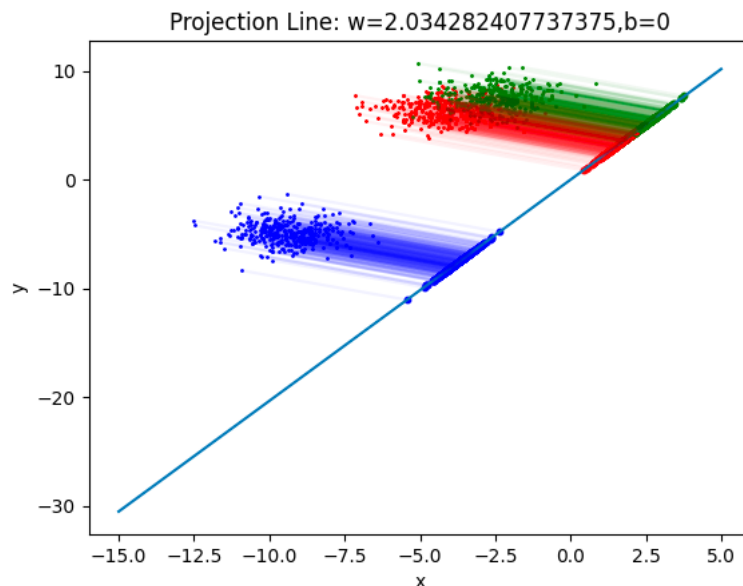
12. (                                    4                        %                              )

   1) Plot the best projection line on the <u>training data</u> and <u>show the slope and intercept on the title</u> *(you can choose any value of intercept for better visualization)*
   2) c o l o r i z e   t h e   t r a i n i n g   d a t a   w i t h   e a c h   c l a s s
   3) project all training data points on your projection line. Your result should look like the below image (This image is for reference, not the answer)



Projection Line: w=2.034282407737375,b=0

## (20%) Train your own model

Requirements:

- Using another dataset that we provided (a real-world dataset).
- Train your model (FLD or Logistics Regression model that you implemented above).

- Try different parameters and feature engineering to beat the baseline.
- Save your testing predictions in the CSV file.

**Criteria:**

13. Explain how you chose your model and what feature processing you have done in detail. Otherwise, no points will be given.

| Point | Accuracy |
|-------|----------|
| 20 | testing acc > 0.921 |
| 15 | 0.91 < testing acc <= 0.921 |
| 8 | 0.9 < testing acc <= 0.91 |
| 0 | testing acc <= 0.9 |

batch_size: 8

learning rate : 0.9

epoch : 10000000

Use feature :

```python
def varible_extend( df, expend = ['Feature1','Feature2','Feature3','Feature4'] ) :
    X = df.drop('Target', axis=1)

    for i in range(len(expend)) :
        for j in range(i,len(expend)) :
            name = f'{expend[i]}*{expend[j]}'
            X[name] = X[expend[i]]*X[expend[j]]

    return X
```
Python

I used logistic regression as my classifier and did feature expansion because I did not know the effect of feature on this Target, so I left all the expansion features, which can be learned from the last homework that more features are more helpful to the regression performance.

# Part. 2, Questions (30%):

(6%) 1. Discuss and analyze the performance

a) between Q10 and Q11, which approach is more suitable for this dataset. Why?

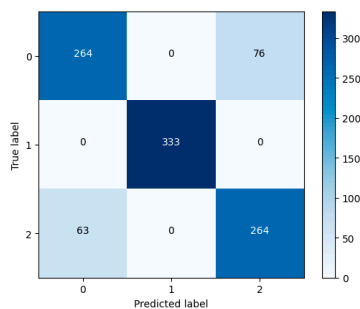b) between different values of k in <u>Q11</u>. (Which is better, a larger or smaller k? Does this always hold?)

(a) According to Q12, it can be found that two classes are very close to each other, so when using knn, k is set a little larger, it is easy to be biased towards a certain class, but it is not necessarily correct, we can learn from Q11, when k=4 the accuracy rate will be lower than k=3, in this data set should take the class mean to predict better
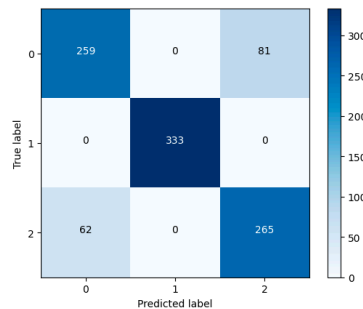
(b)

| k | 4 | 5 | 10 | 20 | 100 | 327 | all |
|---|---|---|----|----|-----|-----|-----|
| accuracy | 0.84 | 0.862 | 0.854 | 0.861 | 0.854 | 0.857 | 0.327 |

The number of this dataset class is 339 334 337, so I tried the numbers in the table above to experiment and found that the larger the number, the better the result.



k=20 :                                    k=327

Because this time the dataset has two classes overlapping each other, so when the k is set a little too large, but will be estimated wrong
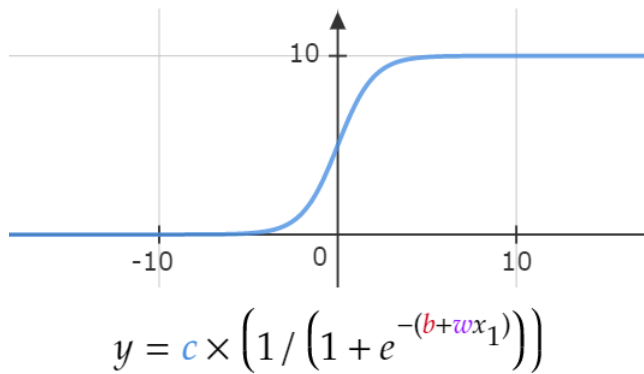
(6%) 2. Compare the sigmoid function and softmax function.
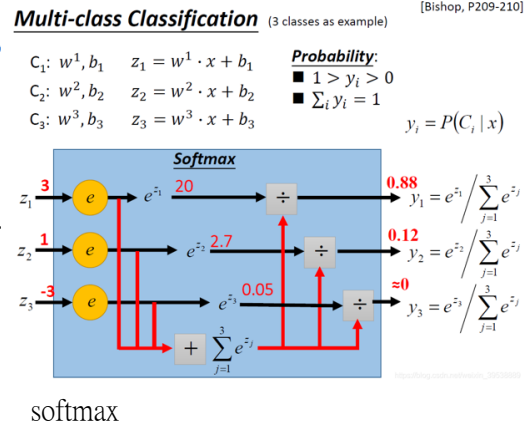
(1). Range of values:

The sigmoid function maps any real-valued number to a value between 0 and 1. This makes it suitable for binary classification problems where the output is either 0 or 1. In contrast, the softmax function maps a vector of real numbers to a probability distribution, meaning that the output values are between 0 and 1 and sum up to 1. This makes it suitable for multiclass classification problems.

(2). Shape of the curve:

The sigmoid function has an S-shaped curve and is monotonically increasing, which means that the output increases as the input increases. In contrast, the softmax function has a curved shape, and the output is not necessarily monotonically increasing.

$$y = c \times \left(1 / \left(1 + e^{-(b + wx_1)}\right)\right)$$

sigmoid                                               softmax

(6%) 3. Why do we use cross entropy for classification tasks and mean square error for regression tasks?

Cross entropy for classification tasks:

Cross entropy is a loss function that measures the difference between the predicted probability distribution and the actual probability distribution for a classification problem. In other words, it measures the distance between the predicted class probabilities and the true class probabilities. Cross entropy is used in classification tasks because it is more suited to handling probability distributions and penalizes more strongly for larger prediction errors.

Mean square error for regression tasks:

Mean square error is a loss function that measures the average squared difference between the predicted values and the true values. It is used in regression tasks because it is well-suited for continuous and numeric data. In regression tasks, the goal is to predict a continuous value, such as a real number, and mean square error provides a measure of the average squared difference between the predicted values and the true values.

In summary, cross entropy is more suited to classification tasks because it can handle probability distributions and penalizes more strongly for larger prediction errors.

(6%) 4. In Q13, we provide an imbalanced dataset. Are there any methods to improve Fisher Linear Discriminant's performance in handling such datasets?
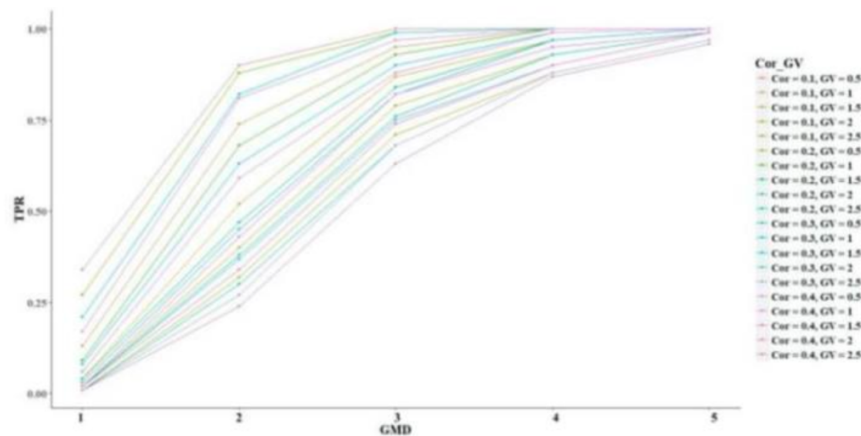
Yes! Here are a few ways to solve the problem
Sampling Techniques :

Sampling techniques are a group of methods used to address imbalanced datasets in machine learning. These techniques aim to balance the distribution of the classes in the dataset by either undersampling the majority class or oversampling the minority class.

According to "provide an imbalanced dataset. are there any methods to improve Fisher Linear Discriminant's performance in handling such datasets? "The authors of this paper propose three methods to solve this problem

LDA :



(c) LDA

Random Undersampling: (RUS)

This technique involves randomly selecting a subset of examples from the majority class to balance the class distribution. While it is simple to implement, it can lead to loss of important information as some examples from the majority class may be discarded.
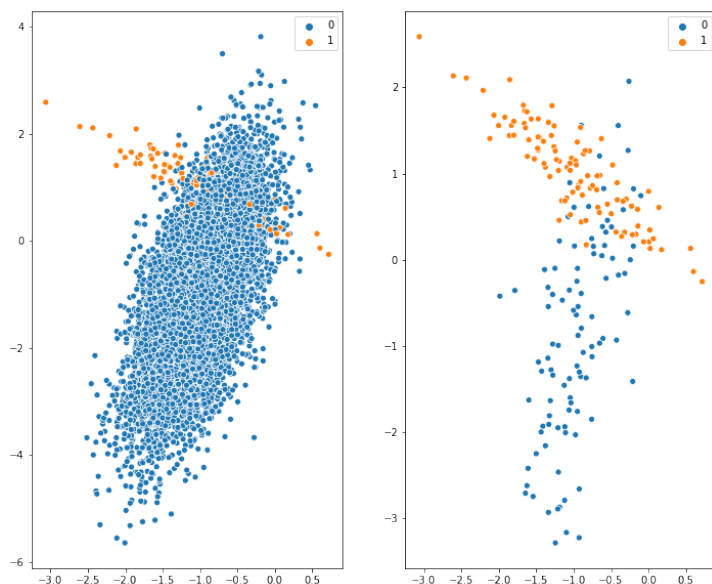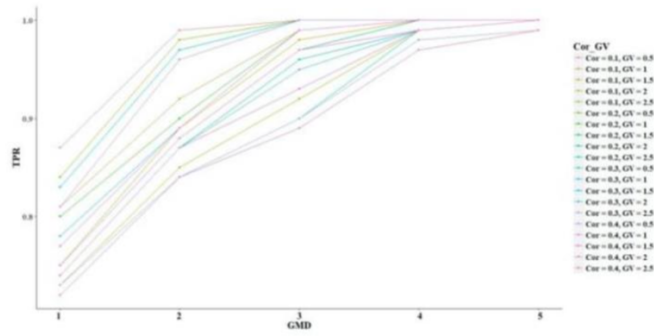


Figure : Training data (left), Undersampling (right)

(b) RUS-LDA

## Random Oversampling: (ROS)

This technique involves randomly replicating examples from the minority class until the class distribution is balanced. While this method is easy to implement and can prevent loss of information, it can also lead to overfitting and increased computational costs.
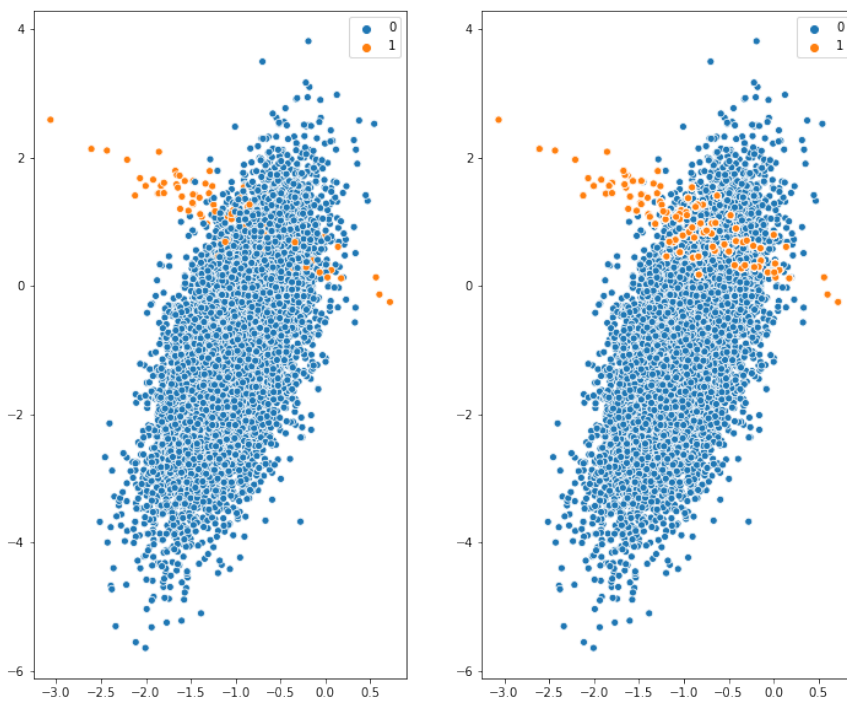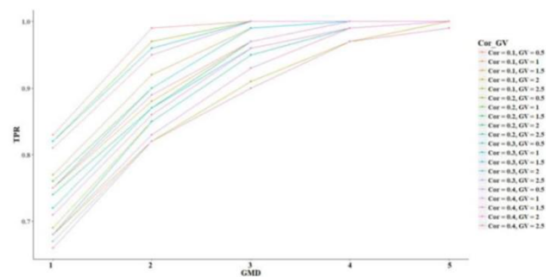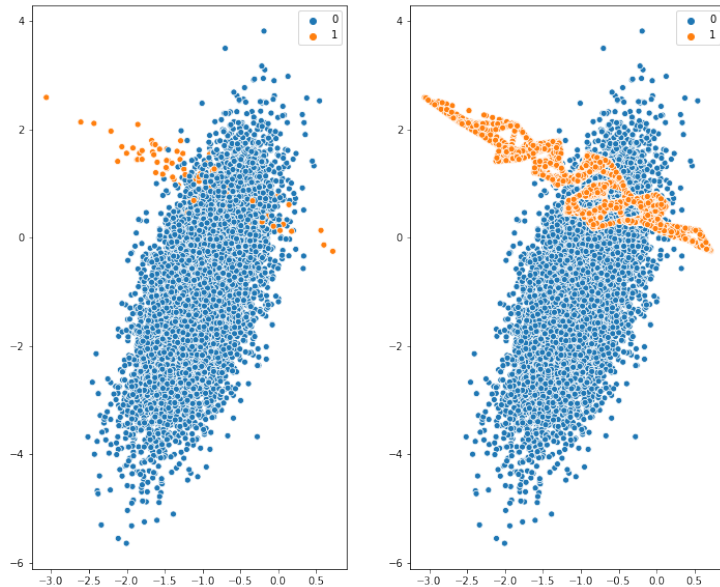


Figure : Training data (left), Oversampling (right)



(a) ROS-LDA

Synthetic Minority Over-sampling Technique (SMOTE):

SMOTE is an oversampling technique that generates synthetic examples for the minority class by interpolating between existing examples. It works by selecting a minority class example and finding its k-nearest neighbors. It then creates synthetic examples by randomly interpolating between the selected example and its neighbors. SMOTE can help to prevent overfitting and can lead to improved generalization.



(6%) 5. Calculate the results of the partial derivatives for the following equations. (The first one is binary cross-entropy loss, and the second one is mean square error loss followed by a sigmoid function.)

$$\frac{\partial}{\partial x}\left(y * \ln(\sigma(x)) + (1 - y) * \ln(1 - \sigma(x))\right)$$

$$\frac{\partial}{\partial x}\left((y - \sigma(x))^2\right)$$

sigmoid derivatives



(1)

$$\frac{\partial}{\partial x}\left(y * \ln(G(x)) + (1-y)\ln(1-G(x))\right)$$

$$\Rightarrow y * \frac{1}{G(x)} G'(x) + (1-y) * \frac{1}{1-G(x)} \cdot (1-G(x))'$$

$$\Rightarrow y * \frac{1}{G(x)} * (1-G(x)) G(x) + (1-y) * \frac{1}{1-G(x)} \cdot G(x) \cdot (1-G(x))$$

$$= y(1-G(x)) + (1-y) G(x) = G(x) - 2y G(x) + y.$$

(2)

$$\frac{\partial}{\partial x}\left((y-G(x))^2\right) \Rightarrow \frac{\partial}{\partial x}\left(y^2 - 2y G(x) + [G(x)]^2\right).$$

$$\Rightarrow -2y G(x)' + 2 G(x)' \Rightarrow (-2y+2) G(x)' \Rightarrow G(x)(1-G(x)) [-2y+2]$$