# Binary Classification Problem of Breast Cancer by Logistic Regression and KNN (midterm-project)

Name: WangYexiang
SID: 12012529

*Abstract*—In this project, the Logistic Regression and KNN are used to solve a binary classfication problem of breast cancer. By the researching, the correct method of pre-process of data is proved important and influence the performance of the direction a lot. How to use less data to train a better model with effective algorithm is the direction of future.

*Keywords*—Dataset, binary classification, Logistic Regression, KNN, sklearn, pre-process, metrics

## I. INTRODUCTION

To deal with binary classification problem of breast cancer by logistic regression and KNN, there are two datasets for training models with different data, "origin_breast_cancer_data.csv" containing 357 benign and 212 malignant samples and "breast_cancer_data_357B_100M.csv" containing 357 benign and 100 malignant samples, which is unbalanced for positive and negative samples.

For the original dataset "origin_breast_cancer_data.csv", use Logistic Regression method and KNN method for binary classification problem (to predict whether it is malignant or benign) with the two datasets mentioned above. Use pre-process the data and tune possible hyper-parameters, to get the best binary classification results for each method. The metrics (recall, precision, and F1 score) will be calculated for both training set and validation set. Compare these two methods and discuss for the conclusions.

For the unbalanced dataset "breast_cancer_data_357B-_100M.csv", if the performance of the previous ways is degraded, then try some modification for these two methods to see if there is improvement.

In this project, the sklearn is used to help training the prediction models and calculate relative parameters.

## II. PROBLEM FORMUALTION

### A. Logistic Regression with original dataset

Firstly, use all data(with 30 dimensions in this project) to train the Logistic Regression model and observe the weights. If there are some weights very small, which means the relative data will not influence the results too much, and even worsen the prediction. What's more, the standardization and normalization of data will be considered. After testing, adjust the model and find whether the optimization works. For the hyper-parameters, it will be shown by drawing the graph of learning rate.

### B. KNN with original dataset

To begin with, use all data(with 30 dimensions in this project) to train the KNN model and use the standardization and normalization to pre-process the data. Then, traverse many values of k to find the best k according to different dataset used to train. By drawing the graph of different k(hyper-parameters here) to indicate the process. Finally, observe the results about different ways to pre-process.

### C. Try to deal with unbalanced dataset

After two cases above, the unbalanced dataset will cause that too few malignant samples, whose results is high error rate about the malignant cases. Therefore, when selecting the samples to train the model randomly, try to keep as much as possible malignant cases to promise the necessary training samples. The pre-process methods of standardization and normalization will also be used and observe the results.

## III. METHOD AND ALGORITHMS

In this project, the Logistic Regression method and KNN method are used to figure out the binary classification problem. Relative knowledge and algorithms will be introduced briefly below.

### A. Logistic Regression

In simple terms, logistic regression is a machine learning method used to solve binary classification (0 or 1) problems to estimate the likelihood of something. For example, the possibility of a user purchasing a certain product, the possibility that a patient has a certain disease, and the possibility of an advertisement being clicked by the user. Note that "probability" is used here, not mathematical "probability", and the result of logisitic regression is not a probability value in the mathematical definition and cannot be used directly as a probability value. This result is often used for weighted summation with other eigenvalues rather than direct multiplication.

- Sigmoid function (or Logistic function)

$$g(z) = \frac{1}{1 + e^{-z}} \tag{1}$$

- Cost Function
  In logistic regression, the most commonly used cost function is Cross Entropy:

$$J(\theta) = -\frac{1}{m} \left( \sum_{i=1}^{m} \left( y^{(i)} \log h_\theta \left( x^{(i)} \right) + \right. \right.$$
$$\left. \left. \left( 1 - y^{(i)} \right) \log \left( 1 - h_\theta \left( x^{(i)} \right) \right) \right) \right) \tag{2}$$

- Optimize

$$\min_{\theta} J(\theta) \qquad (3)$$

## B. KNN (K Nearest Neighbors)

The full name of KNN is K Nearest Neighbors, which means K nearest neighbors. From this name, we can see some clues of the KNN algorithm. K nearest neighbors, there is no doubt that the value of K must be crucial. In fact, the principle of KNN is that when predicting a new value x, it determines which category x belongs to based on what class it is the nearest K points.

- Distance calculation

$$d(x, y) := \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}. \qquad (4)$$

- Selection of k

Start by picking a smaller K value, increase the value of K, then calculate the variance of the validation set, and finally find a suitable K value.

## C. Pre-process of data

- Normalization

In this project, the Min-Max Normalization was used.

$$x^* = \frac{x - x_{min}}{x_{max} - x_{min}} \qquad (5)$$

- Standardization

In this project, Z-score normalization (standard deviation normalization / zero mean normalization) was used.

$$x^* = \frac{x - \bar{x}}{\sigma} \qquad (6)$$

## D. sklearn

SciKit-Learn, also known as sklearn, is an open source machine learning toolkit based on the Python language. It enables efficient algorithm applications through python numerical computing libraries such as NumPy, SciPy and Matplotlib, and covers almost all mainstream machine learning algorithms. In this project, the sklearn.linear_model, sklearn.neighbors, sklearn.metrics, sklearn.preprocessing and so on will be used.

## IV. EXPERIMENT RESULTS AND ANALYSIS

### A. Logistic Regression model trained by original dataset

First of all, all data are used to train the model, and the parameters of Logistic Regression model is used with l1 regulation. It is obvious that many parameters are very small, which means they even don't influence the prediction of the model at all. the result is shown below as Figure 1:

In other side, the l2 regulation, Normalization and Standardization are uesed beacuse of the difference of the 30 sets of data. The effects are better, results shown as Figure 2, Figure 3 and Figure 4:

By selecting the effective data to train the model, the result is getting better according to the case using all data, shown as Figure 5: It is obvious that after selecting the proper data, the effect of model is better according to the higher f1_score.



Fig. 1. Figure1:all data with regulation l1 of LR



Fig. 2. Figure2:all data with regulation l2 of LR

### B. KNN model trained by original dataset

To begin with, the all data without pre-process was used to train the KNN model, whose score graph of different k was shown as Figure 6. After finding the best k, the result of this best k was shown as Figure 7.

Then the Normalization and Standardization are uesed. For Standardization, the score graph of different k was shown as Figure 8. After finding the best k, the result of this best k was shown as Figure 9.

For Normalization, the score graph of different k was shown as Figure 10. After finding the best k, the result of this best k was shown as Figure 11.

### C. Logistic Regression model trained by unbalanced dataset

First of all, all data are used to train the model, and the parameters of Logistic Regression model is used with l1 regulation. It is obvious that many parameters are very small, which means they even don't influence the prediction of the model at all. the result is shown below as Figure 12: In other side, the l2 regulation, Normalization and Standardization are



Fig. 3. Figure3:all data with standardization of LR



Fig. 4. Figure4:all data with normalization scale of LR

Fig. 5. Figure5:after selecting data of LR



Fig. 6. Figure6:find best k without pre-process of KNN



Fig. 7. Figure7:result without pre-process of KNN



Fig. 8. Figure8:find best k with Standardization of KNN



Fig. 9. Figure9:result with Standardization of KNN



Fig. 10. Figure10:find best k with Normalization of KNN



Fig. 11. Figure11:result with Normalization of KNN

uesed beacuse of the difference of the 30 sets of data. However, in this time, the effects are uncertain. Some of them even become worse. The reason inferred is unbalanced data and improper pre-process operation of dataset. The results shown as Figure 13, Figure 14 and Figure 15:

By selecting the effective data to train the model, the result is also uncertain according to the case using all data. Some even get wore a lot, shown as Figure 16:

### D. KNN model trained by unbalanced dataset

To begin with, the all data without pre-process was used to train the KNN model, whose score graph of different k was shown as Figure 17. After finding the best k, the result of this best k was shown as Figure 18.

Then the Normalization and Standardization are uesed. For Standardization, the score graph of different k was shown as



Fig. 12. Figure12:all data with regulation l1 of LR



Fig. 13. Figure13:all data with regulation l2 of LR

Fig. 14. Figure14:all data with standardization of LR



Fig. 15. Figure15:all data with normalization scale of LR

Figure 19. After finding the best k, the result of this best k was shown as Figure 20.

For Normalization, the score graph of different k was shown as Figure 21. After finding the best k, the result of this best k was shown as Figure 22.

It is clear that the pre-process of data which used Normalization and Standardization is effective and helpful for KNN model in the condition unbalanced data, which means reduce the influence caused by one-result samples

## V. CONCLUSION AND FUTURE PROBLEMS

For the machine learning, the pre-process of data is very important. How to select proper data to train the model decide the result of the model.

In this project, the pre-process operation of data almost implement the f1_score of the model, which means the rate of correct prediction is lifted. For instance, as KNN model, the Normalization and Standardization are very notable. However,



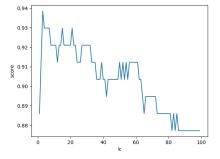Fig. 16. Figure16:after selecting data of LR



Fig. 17. Figure17:find best k with Normalization of KNN



Fig. 18. Figure18:result with Normalization of KNN



Fig. 19. Figure19:find best k with Normalization of KNN



Fig. 20. Figure20:result with Normalization of KNN



Fig. 21. Figure10:find best k with Normalization of KNN



Fig. 22. Figure11:result with Normalization of KNN

some times it don't perform very well, like the Logistic Regression in this project.

About the future problems, when the number of data is small, just like this project, how to use the more effective algorithms to reduce the dependence of dataset, which means the model can find effective features and laws is a direction deserving consideration.

## REFERENCES

[1] Python—KNN https://zhuanlan.zhihu.com/p/143092725
[2] http://t.csdn.cn/WuV4W
[3] http://t.csdn.cn/dTS0C
[4] https://www.knowledgedict.com/tutorial/sklearn-lr.html
[5] https://zhuanlan.zhihu.com/p/462192060
[6] https://developer.aliyun.com/article/692322: :text=Python
[7] http://t.csdn.cn/ANXzE

.