# Machine Learning, Fall 2019: Project 1

Xin Zhang (G32507616)

Header: List the major resources you used to complete this project and the programming language you used.

Language: Python

(30 points) Dataset details:

Data set 1 - MNIST:

The data size in train.csv: 42000

1000 from 95% (39900) are used for training

5% are used for validation: 2099

Test data are read from test.csv: 28000
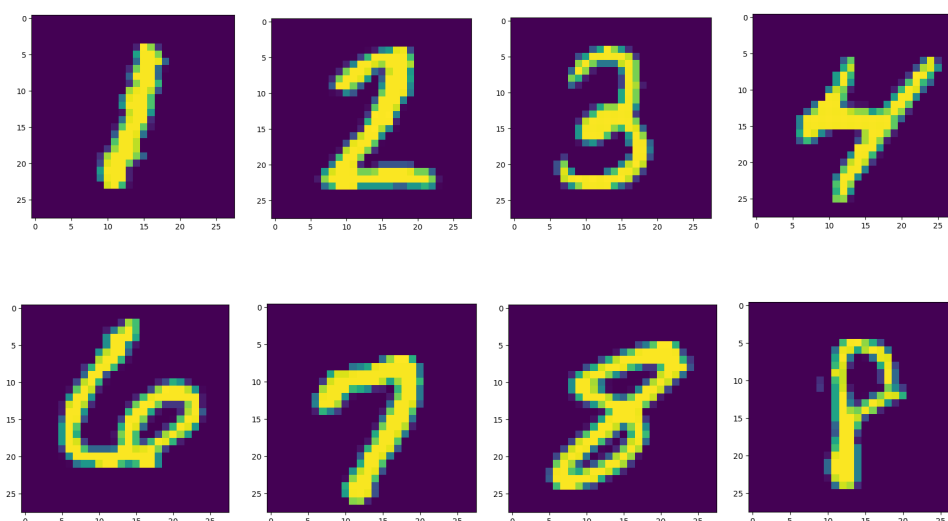
Figure 1-1 shows some simple visualizations:



Figure 1-1 simple visualizations from data set 1

(15 points) Algorithm Description: K-NN is a very clear algorithm, so here describe any data pre-processing, feature scaling, distance metrics, or otherwise that you did.

MNIST:

1. shuffle the data and label with same random seed.
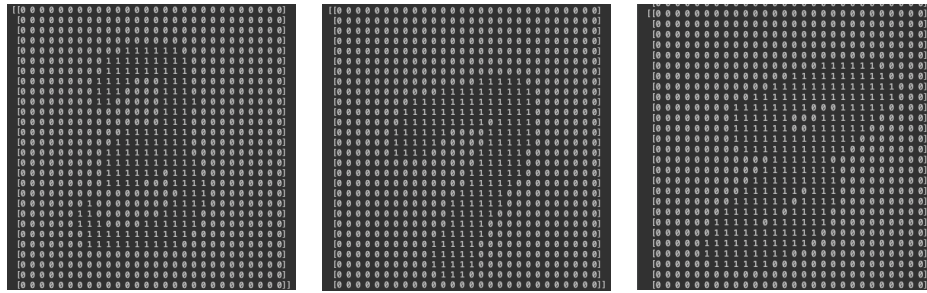2. Transfer the data matrix to 0,1 matrix, as shown in figure 1-2:

Fig 1-2: 0,1 matrix of the handwritten digits

3. Set K = 7. Then step into adaptive training phase for data model:

4. Select 1000 random images as to build the KNN training model. Where include 100 images for each digit. That is to say, 100 '0', 100 '1',..,100 '9'.

5. Using 1000 random images from validation data set to calculate the accuracy of the model.

6. Repeat step 3 and step 4 epoch = 100 times, and choose the one with highest accuracy as final data set for KNN then save the data set and label set to csv files.

7. Adaptive tuning K from 1 to 10 to choose best K for this problem.

KNN for hand written digits recognition:

7 nearest recognition is used to classify the digits.

Euclidian distance is adapted to calculate the distance between two pictures.

KNN Algorithm pseudocode:

For data in test_data:

 Calculate the distance of data to all imgs in train_data

 Get the 7 items with lowest distance

 If only one img has the largest number:

  Choose the label of img with most amount as prediction

 Else multiple imgs have the same largest number:

  Choose the one with smallest distance as prediction

Data set 2 - Pima:

Data size is 768,

Train data size is 518 (70%*95%)

Validation data size is 26 (70%*5%)

Test data size is 230 (30%)

In this data set first do the normalization of data, next Euclidian distance is adapted to calculate the distance between two data.
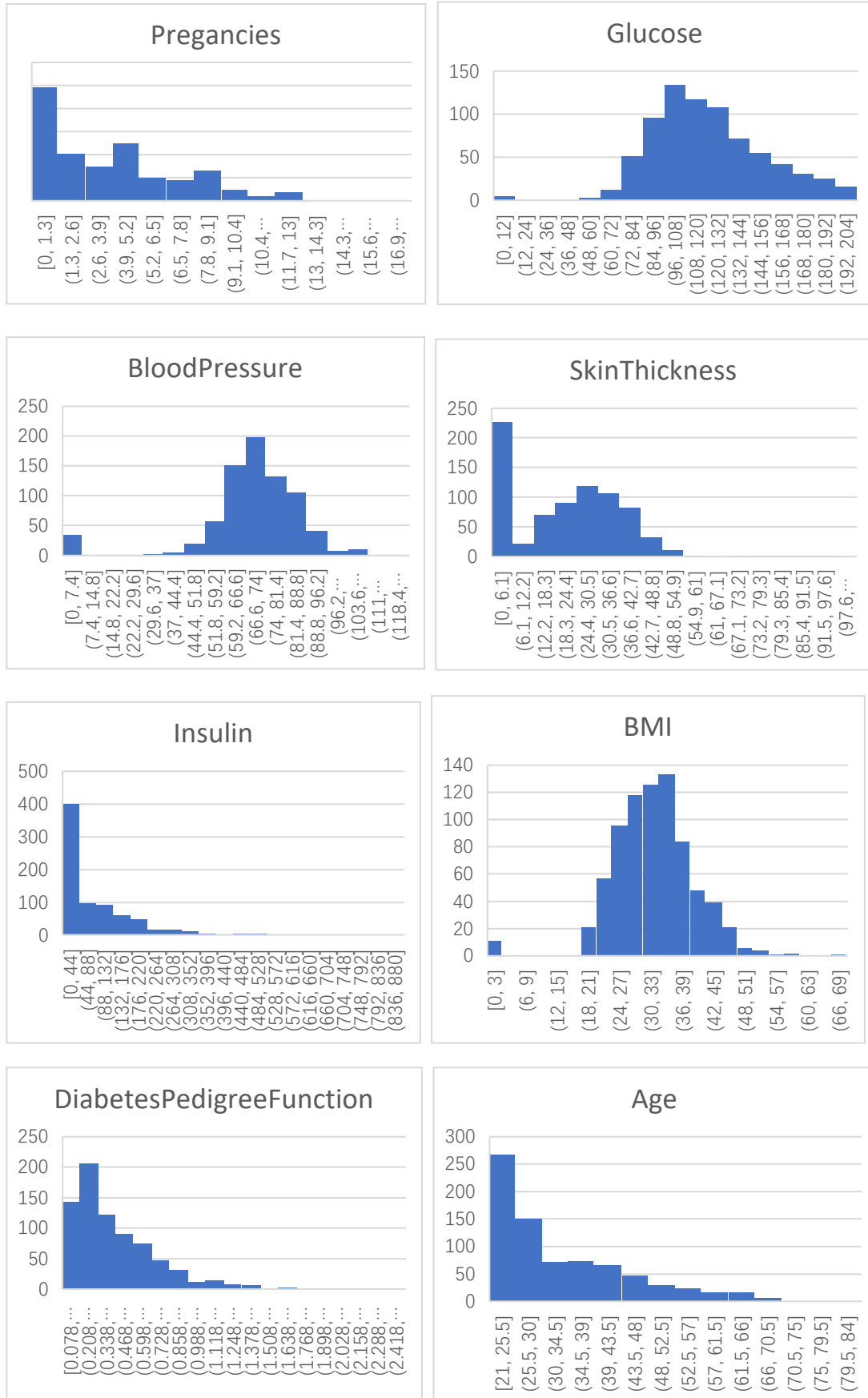
Distribution of the data is shown in Fig 1-3:

Fig 1-3: Distribution of the data set Pima

(45 points) Algorithm Results: Show the accuracy of your algorithm — in the case of the Pima Dataset, show accuracy with tables showing false positive, false negative, true positive and true negatives. For the Pima Dataset, use three different distance metrics and compare the results.

In the case of the MNIST digits show the complete confusion matrix. Choose a single digit to measure accuracy and show how that number varies as a function of K.

The accuracy of validation data 2099 is 89.3 %.
Take digit 1 as an example:
Confusion matrixes with distance are show as below:

| | | Actual Class | | | | | | | | | | Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | |
| | 0 | **209** | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 1 | |
| | 1 | 0 | **231** | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | |
| | 2 | 0 | 15 | **190** | 1 | 1 | 1 | 2 | 10 | 7 | 0 | |
| | 3 | 1 | 5 | 0 | **182** | 0 | 7 | 2 | 2 | 10 | 6 | |
| Predicted | 4 | 0 | 8 | 1 | 0 | **134** | 0 | 6 | 0 | 0 | 42 | 89.3% |
| Class | 5 | 3 | 3 | 0 | 11 | 1 | **144** | 7 | 0 | 4 | 6 | |
| | 6 | 6 | 6 | 0 | 0 | 0 | 0 | **203** | 0 | 1 | 0 | |
| | 7 | 0 | 14 | 0 | 0 | 1 | 0 | 0 | **200** | 0 | 21 | |
| | 8 | 5 | 9 | 3 | 10 | 1 | 4 | 0 | 2 | **170** | 5 | |
| | 9 | 1 | 2 | 0 | 4 | 5 | 2 | 0 | 4 | 1 | **162** | |

From the figure we could see:
Choose digit 5 as an example and let k varies from 1 to 100:
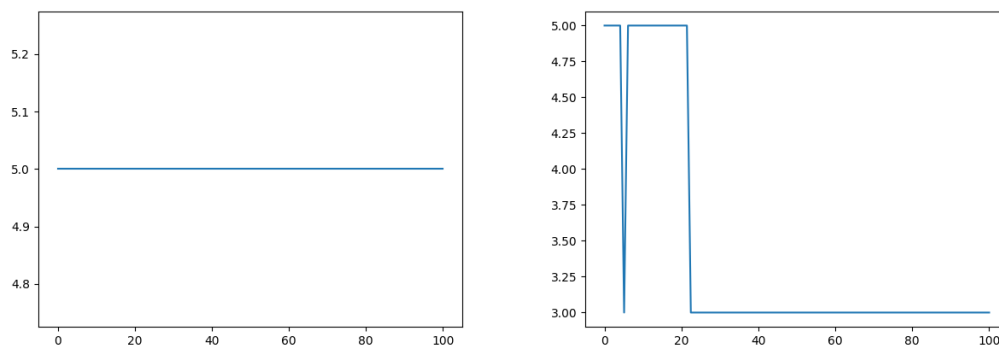The predicted result is shown in figure 2-1:



Fig 2-1 predicted result with different k values

From Fig 2-1, for some figures like figure in the right, prediction result varies with k.

In the case of Pima:
1) Euclidean distance:
For validation data:

| | Actual Class | | | Accuracy |
|---|---|---|---|---|
| Predicted Class | | 0 | 1 | 100% |
| | 0 | **12** | 0 | |
| | 1 | 0 | **13** | |

For test data:

| | Actual Class | | | Accuracy |
|---|---|---|---|---|
| Predicted Class | | 0 | 1 | 100% |
| | 0 | **147** | 0 | |
| | 1 | 0 | **82** | |

Run time: 1.725 s with 228 test data.

```
228
Run time: 1.725228342

Process finished with exit code 0
```

2) Manhattan distance:

| | Actual Class | | | Accuracy |
|---|---|---|---|---|
| Predicted Class | | 0 | 1 | 100% |
| | 0 | **147** | 0 | |
| | 1 | 0 | **82** | |

Run time: 1.592 s with 228 test data.


(10 points) Runtime: Describe the run-time of your algorithm and also share the actual "wall-clock" time that it took to compute your results.

MNIST:

Run time is number of images in model * number of images to predict = 1000*28000 =2.8 * 10^7 cycles

```
28000
Run time: 10690.181816266

Process finished with exit code 0
```

Run time: 10690 s

Pima:

Run time is number of images in model * number of images to predict = 518*228 =1.2 * 10^5 cycles

Run time: 1.592 s


To view the code and final output go:

Github: https://github.com/zhangxin0/KNN