



Cairo University
Faculty of Computers and Artificial Intelligence



Supervised Learning

Prepared By:

Christina Montasser Saad

20190382

Content

1. About data:

1.1. Description

1.2. Sub-Conclusion

2. Dealing with Features

2.1. Binarization

2.2. Feature reduction

2.2.1. Chosen Size

2.2.2. Moment algorithm

3. KNN Classifier

3.1. Algorithm

4. Different tests

5. Conclusion

1. About data:

Data with the title of “**Minst Dataset**”.

1.1 Description

Minst is a handwriting image dataset with 60,000 training points and 10,000 test points. Each image has a size of 28x28. They're RGB images which means their values lay between 0 -255.

Each image has a label number which is an integer from 0 -9.

I downloaded it from Keras library.

1.2 Sub-Conclusion

- It's a supervised classification problem.
- Must use a feature reduction technique as it's too large to have a feature vector of 784x1.
- It's recommended to invert it from RGB to gray-scale images.

2. Dealing with Features

2.1 Binarization

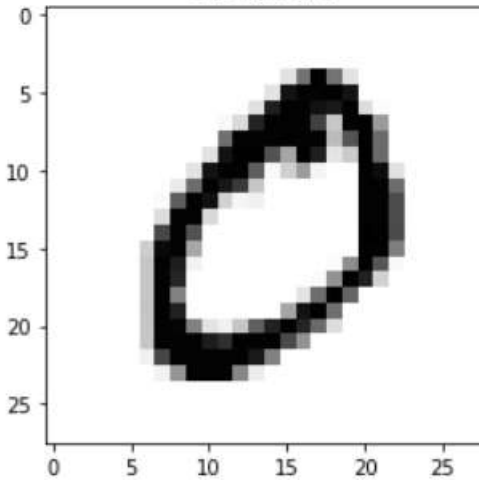
Binarization¹ seems to be an important step especially in cases of document images due to the problems of shadows, non-uniform illumination and blur. Binarization helps to decrease computational load which leads to increase efficiency. Besides, binarization helps to update Optimizely KNN classifier by calculating dis-similarities instead of Euclidean distance.

I faced some of the binarization challenges, like noise and varying in gray levels of the written number. **Global binarization method**² is used with Minst data set. Dependently on how the numbers are written, I choose different threshold 200 and 150. Here are plotting samples after and before binarization, with the two chosen different thresholds.

¹ Binarization computes the threshold value that differentiate object and background pixels.

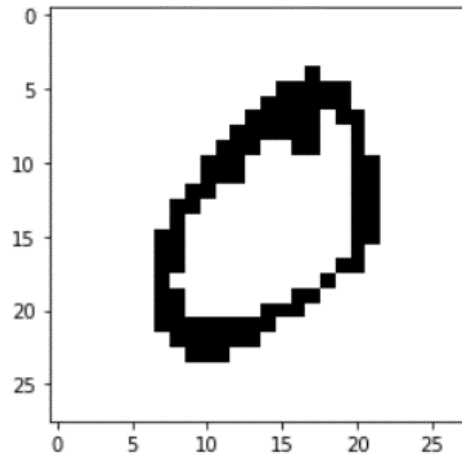
² Choosing one threshold value over the whole image.

true label: 0



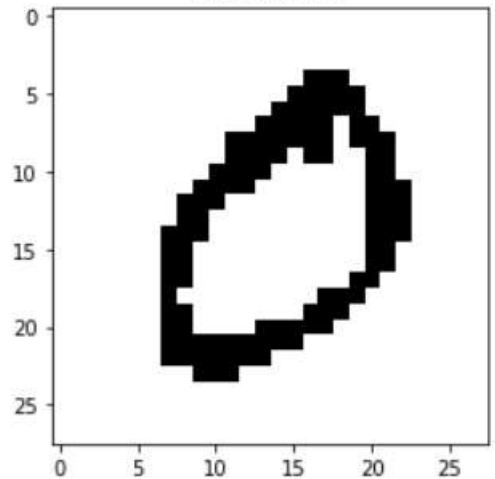
Without binarization

true label: 0



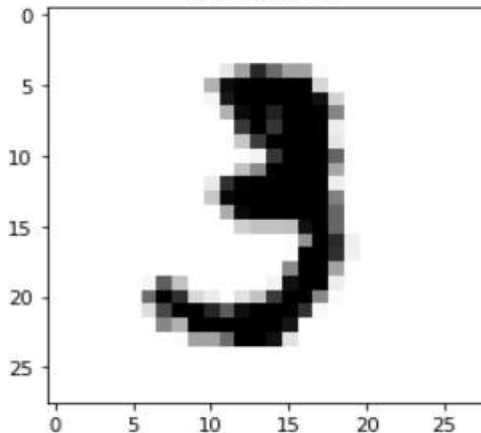
With binarization threshold 200

true label: 0



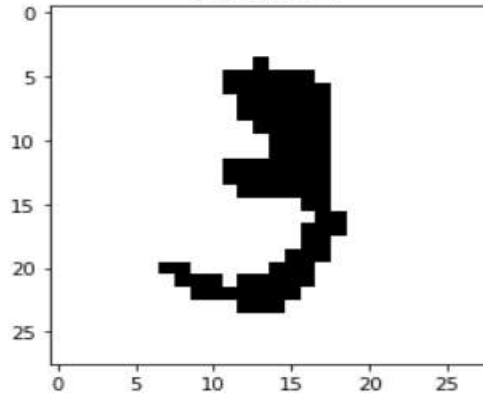
With binarization threshold 150

true label: 3



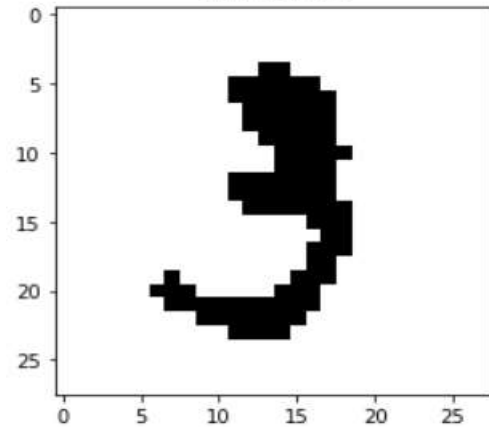
Without binarization

true label: 3



With binarization threshold 200

true label: 3



With binarization threshold 150

Both works fine, but somehow, I see “200” as a threshold fits more to not let noise and blur being considered as important bits. Because it may lead to matched a number can be written in small on-bits with another written in large on-bit, just because the first is noisier and blurrier. Especially I used here the “Centroid”³ technique.

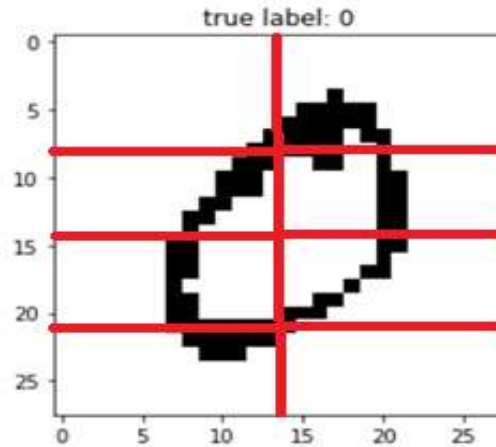
Let this discussion open till the classifier says what it sees...

³ the weighted average of all the pixels of the number.

2.2 Feature reduction

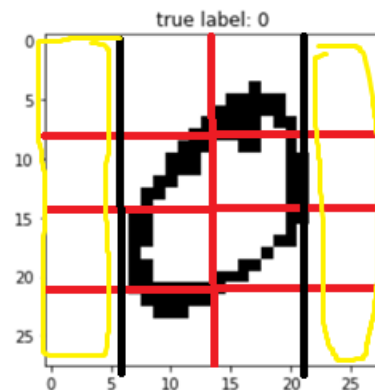
2.2.1 Chosen Size

I split all images to 8 parts, 2x4.



This chosen size depends on where most of the objects lay. Most numbers in Minst dataset lay in the middle, so it would be redundant to split a part of the most-left or right x axis of any image because those are totally white background parts so their x, y moments will be zeros.

Those yellow colored spaces are meaningless. Even if they contained pixels of the written numbers, ignoring them doesn't affect anything. Ignoring them can be described as different binarization threshold and those bits were noise.



2.2.2 Moment algorithm

Centroid is the weighted average of all pixels. “Weighted” means there are pixels contribute more than others so it’s multiplied by a factor of its importance.

$$C_x = \frac{\sum_{x=1}^m x f(x,y)}{f(x,y)} \quad C_y = \frac{\sum_{y=1}^m y f(x,y)}{f(x,y)}$$

$$C_x = \frac{M10}{M00} \quad C_y = \frac{M01}{M00}$$

$$M_{ij} = \sum_{x=1}^m \sum_{y=1}^n x^i y^j I(x, y)$$

Note: M_{00} which refers to the image function.

I applied this algorithm on each block of pixels of the split parts. It’s 2x4, the returning was an ordered pair (x, y) -size 2-. It ended with each image has *16 features*.

Progress: from 728 features to only 16.

Now, our data with a size of ***60,000x16***

3. KNN Classifier

3.1 Algorithm

3.1.1 Dis-similarities (Euclidean distance)

$$d(x_1, x_2) = \sum_{m=1}^D I(x_1 \neq x_2)$$

Because features’ value is binary. Calculate dis-similarities instead of Euclidean distance is better.

3.1.2 Neighbors

Ordered them ascendingly and choose the first k neighbors.

3.1.3 Voting

Find the dominant class over the k returned labels.

4. Different tests

Test#	Binarization	Feature size	Accuracy	KNN distance
Test1	Without	2x4	92.74	Euclidean distance high computational function
Test2	With Threshold = 150	2x4	92.53	Dis-similarity
Test3	With Threshold = 200	2x4	91.74	Dis-similarity

Conclusion: It seems with threshold 150 is the best overall.