# Prototype selection for nearest neighbor

**Qimin Chen**
A53284263
qic003@ucsd.edu

## 1    Prototype selection

I firstly applied PCA (Principal Components Analysis) on original MNIST dataset to reduce the dimensionality from $28 \times 28$ to 25, then each class of MNIST (roughly 5000 images per class) was clustered by K-means into $M/10$ clusters and the center of clusters are used to represent the corresponding classes.

## 2    Pseudocode

As described above, PCA and K-means are used to implement prototype selection. The original MNIST training set is read as $60000 \times 784$ matrix so each image is converted to one-hot.

---
**Algorithm 1** Prototype selection

---
**Input:** original MNIST training set $T$ with $60000 \times 784$ dimension and integer $M$.
**Output:** reduced training set $R$ with $M \times 25$ dimension.
  1: apply PCA on training set and testing set to reduce dimensionality;
  2: reduced_training = [ ]
  3: **for** $label = 0 \cdots 9$ **do**
  4:      model = KMeans(T, M/10)
  5:      reduced_training += model.cluster_centers
  6: **end for**

---

## 3    Experimental results

The confidence intervals is computed by equation (1):

$$(\bar{x} - z^* \frac{\sigma}{\sqrt{n}}, \bar{x} + z^* \frac{\sigma}{\sqrt{n}}) \tag{1}$$

where $\bar{x}$ is accuracy mean, $z^*$ is critical value which is 1.96 in 95% confidence intervals, $\sigma$ is the standard deviation of accuracy and $n$ is the number of accuracy.

Table 1 shows the comparison of average accuracy and confidence intervals along with accuracy standard deviations in parentheses over 10 times for each value of $M = 100, 500, 1000, 5000, 10000$ between Prototype selection and Random selection. It is apparent that the performance of Prototype selection is better than Random selection. Also, PCA helps speeding up nearest neighbor classification because it reduces the dimensionality of dataset.

## 4    Critical evaluation

My method is a clear improvement over random selection as seen from Table 1. Since I have already applied PCA on dataset so that it helps improving the accuracy and saving training time as well, I would like to try on other methods like k-NN where k is greater than 1 or neural networks.

1

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

| Experiment | Prototype selection | Random selection |
|------------|---------------------|------------------|
| M = 100 | $92.19 \pm 0.0011$ % (0.0012) | $73.83 \pm 0.0162$ % (0.0184) |
| M = 500 | $94.97 \pm 0.0010$ % (0.0011) | $87.35 \pm 0.0036$ % (0.0041) |
| M = 1000 | $95.69 \pm 0.0007$ % (0.0008) | $90.40 \pm 0.0022$ % (0.0025) |
| M = 5000 | $96.59 \pm 0.0012$ % (0.0014) | $94.37 \pm 0.0012$ % (0.0013) |
| M = 10000 | $96.81 \pm 0.0009$ % (0.0010) | $95.58 \pm 0.0013$ % (0.0015) |

Table 1: Averaged accuracy and confidence intervals along with accuracy standard deviations in parentheses over 10 times experiments per $M$.