

Nearest neighbor classification

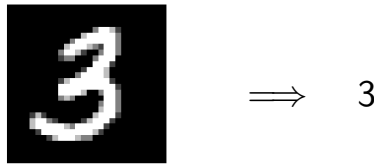
DSE 220

Outline

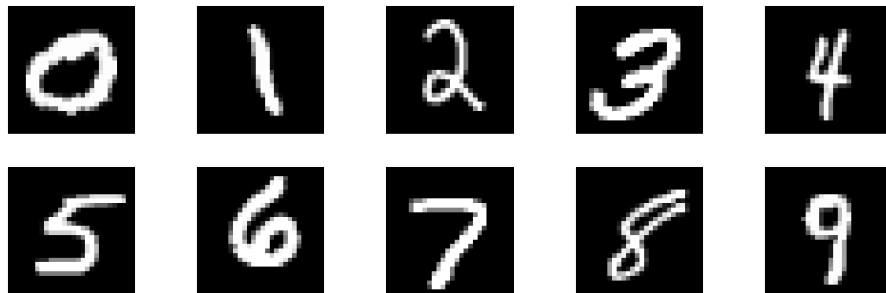
- ① Nearest neighbor classification
- ② k -nearest neighbor
- ③ Choosing the features and distance function

The problem we'll solve today

Given an image of a handwritten digit, say which digit it is.



More examples:



The machine learning approach

Assemble a data set:



The MNIST data set of handwritten digits:

- **Training set** of 60,000 images and their labels.
- **Test set** of 10,000 images and their labels.

And let the machine figure out the underlying patterns.

Nearest neighbor classification

Training images $x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(60000)}$

Labels $y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(60000)}$ are numbers in the range 0 – 9

1 4 1 6 1 1 9 1 3 4 8 5 7 2 6 8 0 3 2 2 6 4 1 4 1
8 6 6 3 5 9 7 2 0 2 9 9 2 9 9 7 2 2 5 1 0 0 4 6 7
0 1 3 0 8 4 1 1 1 5 9 1 0 1 0 6 1 5 4 0 6 1 0 3 6
3 1 1 0 6 4 1 1 1 0 3 0 4 7 5 2 6 2 0 0 9 9 7 9 9
6 6 8 9 1 2 0 8 6 7 0 8 5 5 7 1 3 1 4 2 7 9 5 5 4
6 0 1 0 1 8 7 3 0 1 8 7 1 1 2 9 9 1 0 8 9 9 7 0 9
8 4 0 1 0 9 7 0 7 5 9 7 3 3 1 9 7 2 0 1 5 5 1 9 0
6 6 1 0 7 5 5 1 8 2 5 5 1 8 2 8 1 4 3 5 8 0 9 0 9
6 3 1 7 8 7 5 2 1 6 5 5 4 6 0 5 5 4 6 0 3 5 4 6 0
5 5 1 8 2 5 5 1 0 8 5 0 3 0 4 7 5 2 0 4 3 9 4 0 1

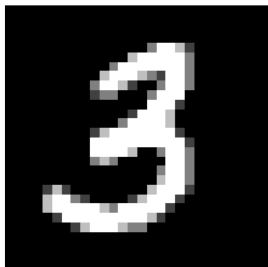


How to **classify** a new image x ?

- Find its nearest neighbor amongst the $x^{(i)}$
- Return $y^{(i)}$

The data space

How to measure the distance between images?



MNIST images:

- Size 28×28 (total: 784 pixels)
- Each pixel is grayscale: 0-255

Stretch each image into a vector with 784 coordinates:



- Data space $\mathcal{X} = \mathbb{R}^{784}$
- Label space $\mathcal{Y} = \{0, 1, \dots, 9\}$

The distance function

Remember Euclidean distance in two dimensions?

$$z = \underset{\bullet}{(3, 5)}$$

$$x = \overset{\bullet}{(1, 2)}$$

Euclidean distance in higher dimension

Euclidean distance between 784-dimensional vectors x, z is

$$\|x - z\| = \sqrt{\sum_{i=1}^{784} (x_i - z_i)^2}$$

Here x_i is the i th coordinate of x .

Nearest neighbor classification

Training images $x^{(1)}, \dots, x^{(60000)}$, labels $y^{(1)}, \dots, y^{(60000)}$

1 4 1 0 1 1 9 1 5 4 8 5 7 2 6 8 0 3 2 2 6 4 1 4 1
8 6 6 3 5 9 7 2 0 2 9 9 2 9 9 7 2 2 5 1 0 0 4 6 7
0 1 3 0 8 4 1 1 1 5 9 1 0 1 0 6 1 5 4 0 6 1 0 3 6
3 1 1 0 6 4 1 1 1 0 3 0 4 7 5 2 6 2 0 0 9 9 7 9 9
6 6 8 9 1 2 0 8 6 7 8 8 5 5 7 1 3 1 4 2 7 9 5 5 4
6 0 6 0 1 8 7 5 0 1 8 7 1 1 2 9 9 1 0 8 9 9 7 0 9
8 4 0 1 0 9 7 0 7 5 9 7 3 3 1 9 7 2 0 1 5 5 1 9 0
6 5 1 0 7 5 5 1 8 2 5 5 1 8 2 8 1 4 3 5 8 0 9 0 9
4 3 1 7 8 7 5 2 1 6 5 5 4 6 0 5 5 4 6 0 3 5 4 6 0
5 5 1 8 2 5 5 1 0 8 5 0 3 0 4 7 5 2 0 4 3 9 4 0 1



To classify a new image x :

- Find its nearest neighbor amongst the $x^{(i)}$ using **Euclidean distance in \mathbb{R}^{784}**
- Return $y^{(i)}$

How accurate is this classifier?

Accuracy of nearest neighbor on MNIST

Training set of 60,000 points.

- What is the error rate on training points? **Zero.**
In general, **training error** is an overly optimistic predictor of future performance.
- A better gauge: separate test set of 10,000 points.
Test error = fraction of test points incorrectly classified.
- What test error would we expect for a *random classifier*?
(One that picks a label 0 – 9 at random?) **90%.**
- Test error of nearest neighbor: **3.09%.**

Examples of errors

Test set of 10,000 points:

- 309 are misclassified
- Error rate 3.09%

Examples of errors:

Query					
NN					

Ideas for improvement: (1) k -NN (2) better distance function.

Outline

- ① Nearest neighbor classification
- ② k -nearest neighbor
- ③ Choosing the features and distance function

K-nearest neighbor classification

Classify a point using the labels of its k nearest neighbors among the training points.

MNIST:

k	1	3	5	7	9	11
Test error (%)	3.09	2.94	3.13	3.10	3.43	3.34

In real life, there's no test set. How to decide which k is best?

① Hold-out set.

- Let S be the training set.
- Choose a subset $V \subset S$ as a *validation set*.
- What fraction of V is misclassified by the k -nearest neighbors in $S \setminus V$?

② Leave-one-out cross-validation.

- For each point $x \in S$, find the k -nearest neighbors in $S \setminus \{x\}$.
- What fraction are misclassified?

Cross-validation

How to estimate the error of k -NN for a particular k ?

10-fold cross-validation

- Divide the training set into 10 equal pieces.
Training set (call it S): 60,000 points
Call the pieces S_1, S_2, \dots, S_{10} : 6,000 points each.
- For each piece S_i :
 - Classify each point in S_i using k -NN with training set $S - S_i$
 - Let ϵ_i = fraction of S_i that is incorrectly classified
- Take the average of these 10 numbers:

$$\text{estimated error with } k\text{-NN} = \frac{\epsilon_1 + \dots + \epsilon_{10}}{10}$$

Outline

- ① Nearest neighbor classification
- ② k -nearest neighbor
- ③ Choosing the features and distance function

Another improvement: better distance functions

The Euclidean (ℓ_2) distance between these two images is very high!



Much better idea: distance measures that are invariant under:

- Small translations and rotations. e.g. *tangent distance*.
- A broader family of natural deformations. e.g. *shape context*.

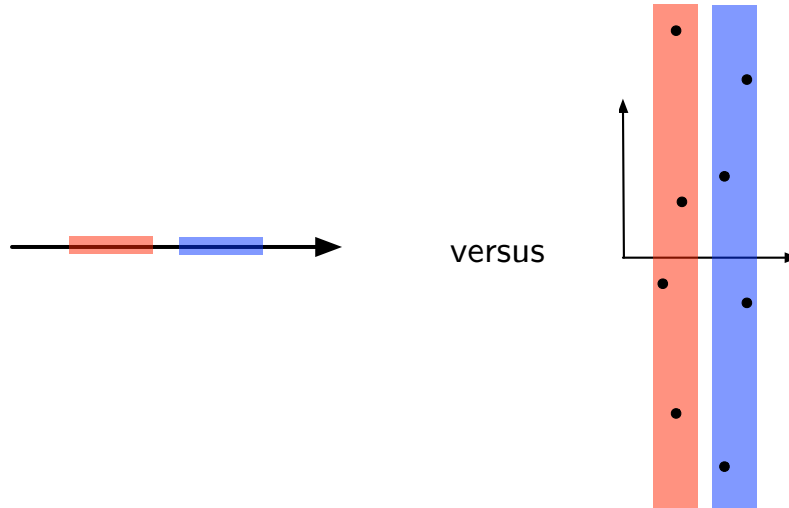
Test error rates:

ℓ_2	tangent distance	shape context
3.09	1.10	0.63

Related problem: feature selection

Feature selection/reweighting is part of picking a distance function.

And, one noisy feature can wreak havoc with nearest neighbor!



Algorithmic issue: speeding up NN search

Naive search takes time $O(n)$ for training set of size n : slow!

There are data structures for speeding up nearest neighbor search, like:

- ① Locality sensitive hashing
- ② Ball trees
- ③ K -d trees

These are part of standard Python libraries for NN, and help a lot.

Example: k -d trees for NN search

A hierarchical, rectilinear spatial partition.

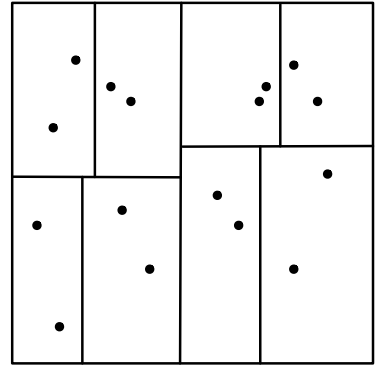
For data set $S \subset \mathbb{R}^d$:

- Pick a coordinate $1 \leq i \leq d$.
- Compute $v = \text{median}(\{x_i : x \in S\})$.
- Split S into two halves:

$$S_L = \{x \in S : x_i < v\}$$

$$S_R = \{x \in S : x_i \geq v\}$$

- Recurse on S_L, S_R



Two types of search, given a query $q \in \mathbb{R}^d$:

- *Defeatist search*: Route q to a leaf cell and return the NN in that cell. This might not be the true NN.
- *Comprehensive search*: Grow the search region to other cells that cannot be ruled out using the triangle inequality.