

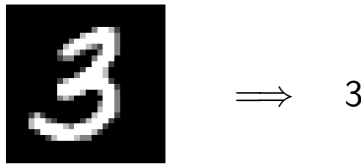
Nearest neighbor classification

Topics we'll cover

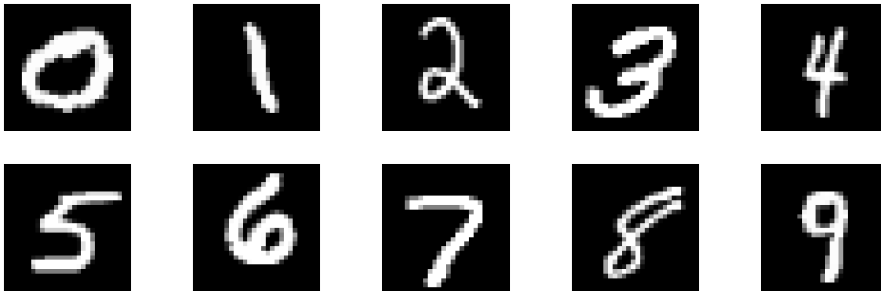
- ① What is a classification problem?
- ② The training set and test set
- ③ Representing data as vectors
- ④ Distance in Euclidean space
- ⑤ The 1-NN classifier
- ⑥ Training error versus test error
- ⑦ The error of a random classifier

The problem we'll solve today

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



Some 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

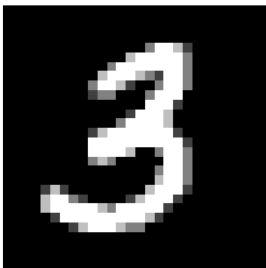


How to **classify** a new image x ?

- 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

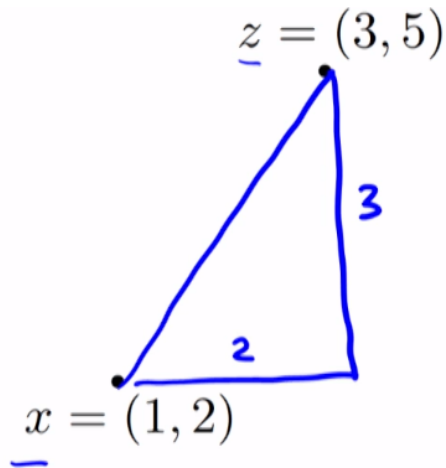
Stretch each image into a vector with 784 coordinates:



- Data space $\mathcal{X} = \mathbb{R}^{784}$ a 784 dimensional vector space then, which we're gonna denote by script X
- Label space $\mathcal{Y} = \{0, 1, \dots, 9\}$ is 784 dimensional Euclidean space R to the 784th.

The distance function

Remember Euclidean distance in two dimensions?



$$\|x - z\| = \sqrt{2^2 + 3^2}$$

common, or default distance

function is perhaps just Euclidean distance.

When you have two points, the Euclidean distance

between them is just the length of the line connecting them.

So it's the length of this line.

And what is that length?

Well, if you look at these two points, X and Z ,

along the first coordinate, they differ by two

and along the second coordinate, they differ by three.

So the length of the line, the distance from X to Z

is simply the square root of two squared plus three squared

which is the square root of 13.

That's the Euclidean distance between X and Z

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 $y^{(1)}, \dots, y^{(60000)}$



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?

So we have these 60,000 training images.
For any training point, its nearest neighbor in the training set is itself.
So it'll definitely get the right label.
So the error rate on the training set is zero.
What that means is that training error
is not a good predictor of future performance.

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					