

WEEK 1

## PREDICTION PROBLEMS

### Nearest neighbor classification

#### Topics we'll cover

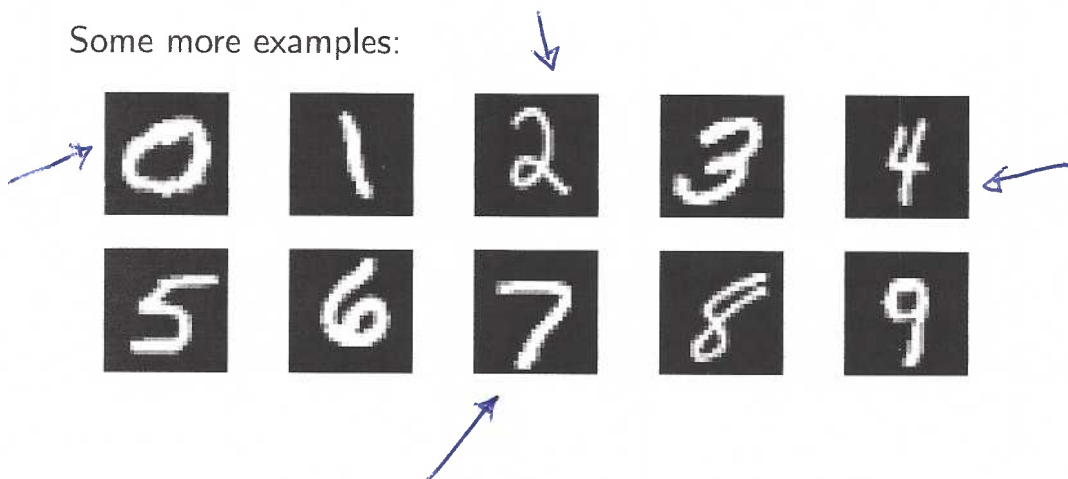
- ① What is a classification problem?
- ② The training set and test set
- ③ Representing data as vectors
- ④ Distance in Euclidean space
- ⑤ The 1-NN classifier
- ⑥ Training error versus test error
- ⑦ The error of a random classifier

## 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:

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

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

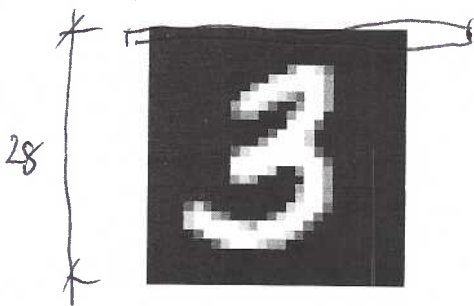


How to **classify** a new image  $x$ ?

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

## The data space

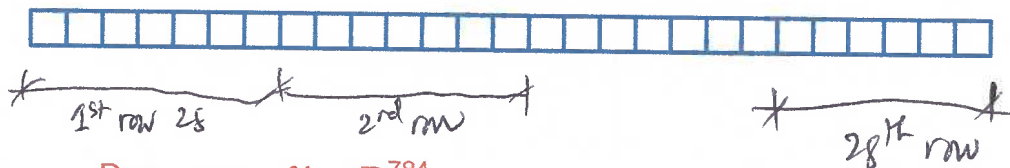
How to measure the distance between images?



MNIST images:

- Size  $28 \times 28$  (total: 784 pixels)
- Each pixel is grayscale: 0-255

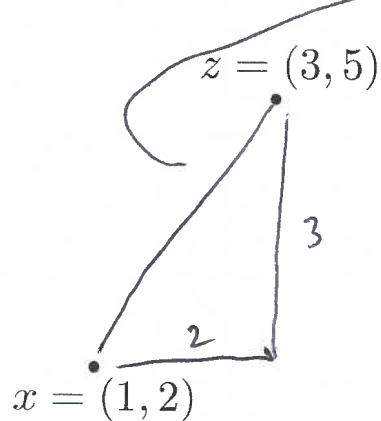
Stretch each image into a vector with 784 coordinates:



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

## The distance function

Remember Euclidean distance in two dimensions?



$$\sqrt{2^2 + 3^2} = \sqrt{13}$$

## Euclidean distance in higher dimension

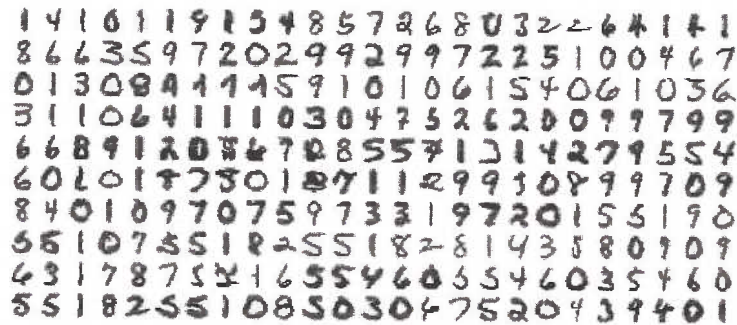
Euclidean distance between 784-dimensional vectors  $x, z$  is

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

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?

## 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:

