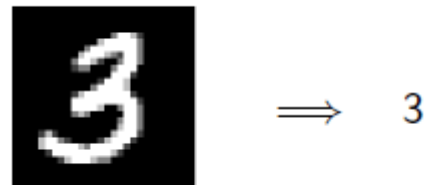


NEAREST NEIGHBOR CLASSIFICATION (najstasza i najprostsza)

Zadanie: **Dany jest obraz ręcznie napisanej cyfry, określ jaka to jest cyfra**



Próbowano rozwiązać problem zapisując zestaw zasad (np. liczba pętli, odcinków itp.), ale to nie zdawało egzaminu (bo: pismo ręczne jest bardzo nieregularne i zmienne).

ML – (sposób podejścia za pomocą ML)

Utworzyć zbiór danych:

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

Zlecić maszynie znalezienie „wzoru” rozpoznawania.

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

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

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

Nearest neighbor = najbardziej podobny = najmniej się różni = najmniejszy dystans

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 $\mathcal{Y} = \{0, 1, \dots, 9\}$

EUCLIDIAN DISTANCE

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 .

$\| \quad \|$ - oznacza Euclidian Distance



To classify a new image x :

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

Ocena dokładności:

- 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%**.

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

K – nearest neighbours

To classify a new point:

- Find the k nearest neighbors in the training set.
- Return the most common label amongst them.

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

W rzeczywistości często nie dysponujemy zbiorem testowym – operujemy wyłącznie zbiorem treningowym. Stosuje się metodę **Cross-validacji** do wyznaczania k w metodzie **k-NN**, jak i wiele innych parametrów w innych metodach.

Dzielimy zbiór treningowy na n – kawałków (folds)