# Digit recognition

## Project Proposal

Presented by Daria Morkis, Alex Kohlmann and Karolina Walach

Supervisor: Dr. Leonid Kostrykin, PD Dr. Karl Rohr

Tutor: Hannah Winter

# What is "Digit recognition"?

"Digit recognition is the ability of computers to recognize human handwritten digits"
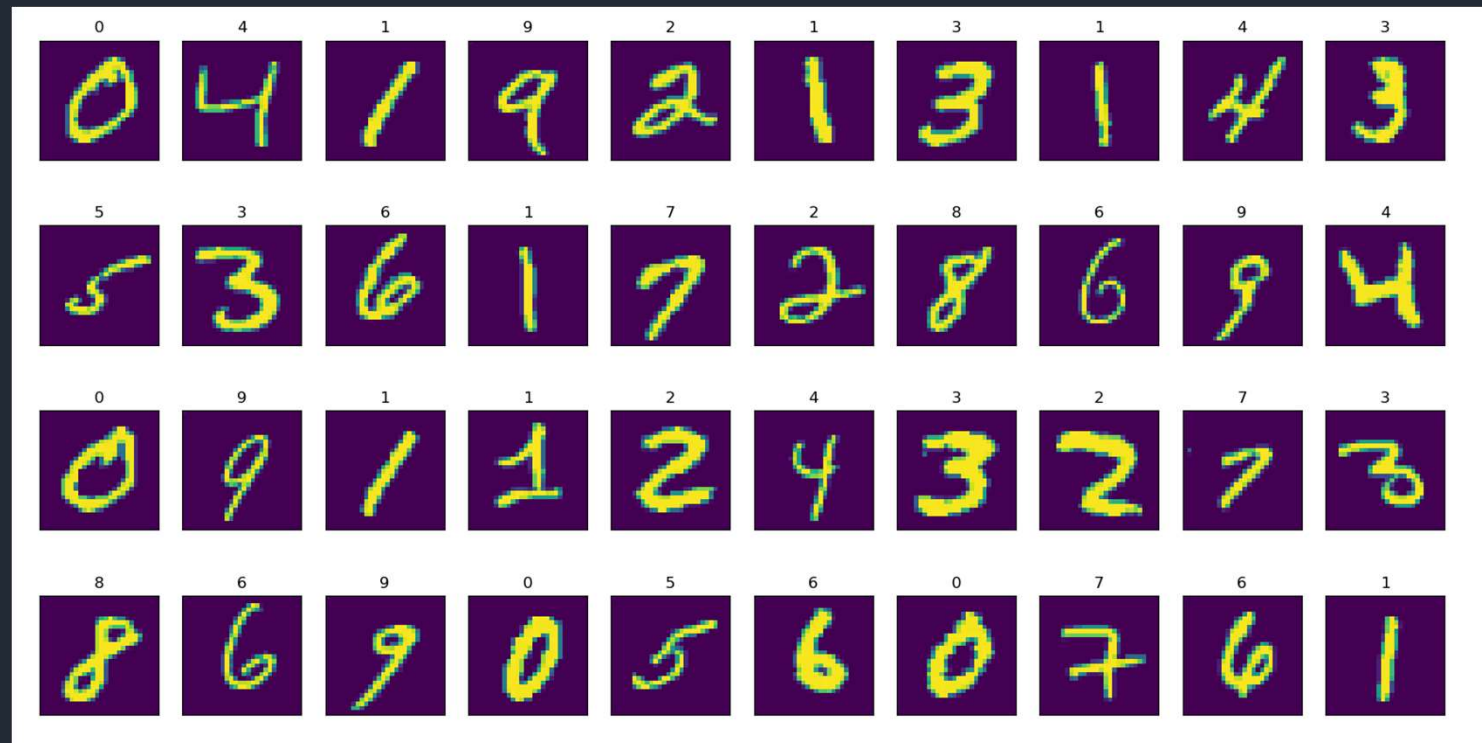
(data-flair.training)

Pixel Values

Label

38,222,225

# Conversion with matplotlib

# List of planned analysis steps

### Principal Component Analysis (PCA)
- Z-transform, covariance, eigenvalues and -vectors

### K-Nearest Neighbour (KNN)
- Euclidian distance, distance sorting, most common neighbours, prediction of class
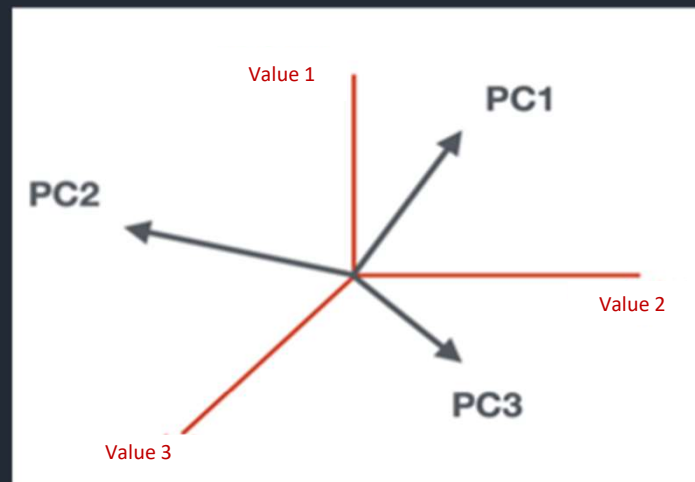
### Evaluation
- Accuracy, optimizing PC/K-value, error analysis

### Additional:
- Implementation of a neural network / comparing outcomes

# Principal Component Analysis (PCA)

60.000 dimensions ⟶ n principal components (reduction of dimensions)

# Principal Component Analysis (PCA)

1. Import libraries and dataset

   

2. z-transform

$$z = \frac{x - \mu}{\sigma}$$

3. Implement PCA
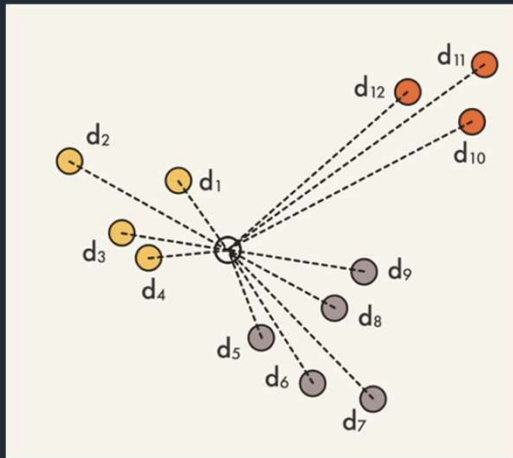
- Covariance

- Eigenvalues and -vectors   $Xv = \lambda v$
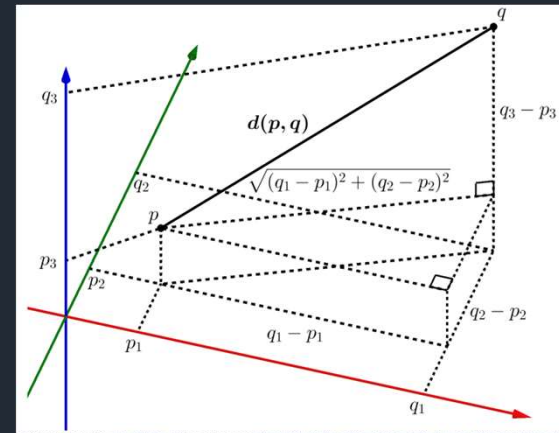
  eigenvector   eigenvalue
  of X          of X

# K-Nearest Neighbour (KNN) algorithm

Euclidean distance

- length of line segment between two points
- $d(p, q) = \sqrt{\sum_{i=1}^{n}(q_i - p_i)^2}$



https://de.wikipedia.org/wiki/Euklidischer_Abstand



https://youtu.be/0p0o5cmgLdE

- calculating distance between test data image and train data images
- Euclidean distance: square of the differences between their coordinates in n-dimensional space (set by PCA)

# K-Nearest Neighbour (KNN) algorithm
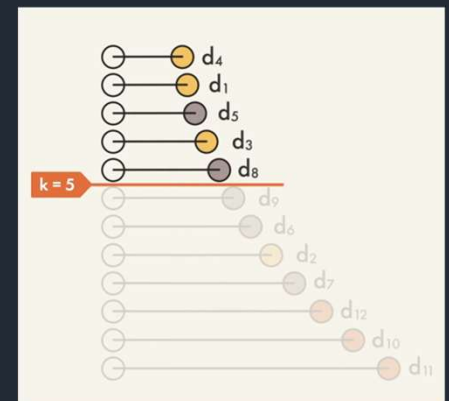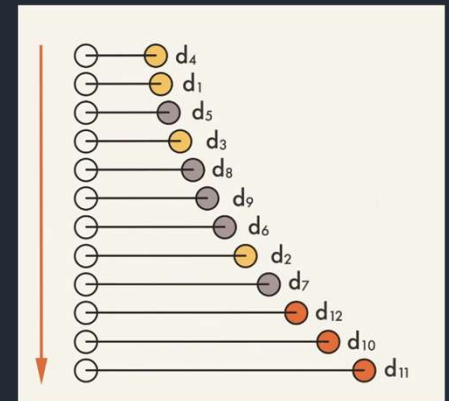
## Sort distances in ascending order

- sorting labeled neighbors (train data points) by ascending distance
- Get the k nearest neighbors by taking top k rows from sorted array

## Most common neighbours

- select the most common labels (digits) for these rows by majority vote
- predicting class of new data point
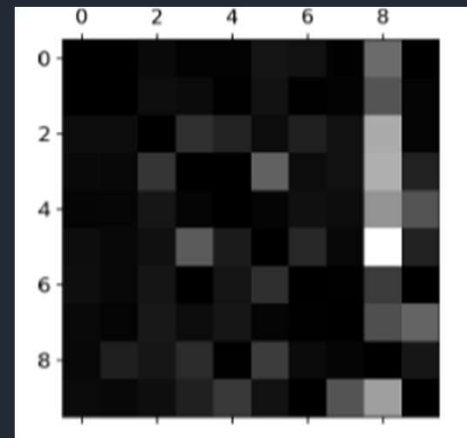
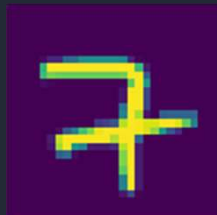## Evaluation

- computing mean accuracy



https://youtu.be/0p0o5cmgLdE

# Problem Nr.1

1. Some numbers are difficult to distinguish

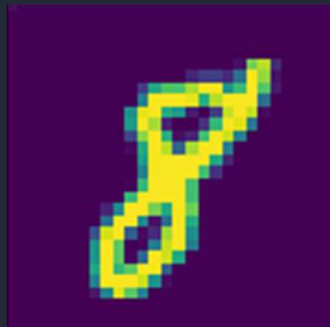 1 or 7?  $\longrightarrow$ Estimating error frequency with e.g. confusion matrix



*Géron, A. (2017). Hands-on machine learning with Scikit-Learn and TensorFlow (Sebastopol, CA: O'Reilly Media).*
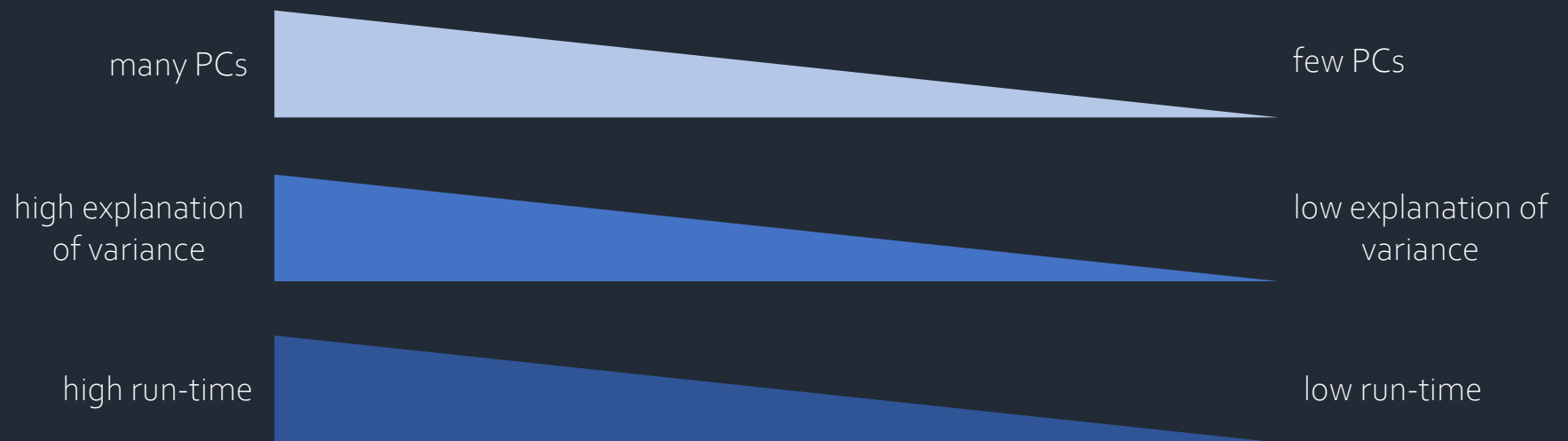
# Solution for Problem Nr.1

1. Some numbers are difficult to distinguish

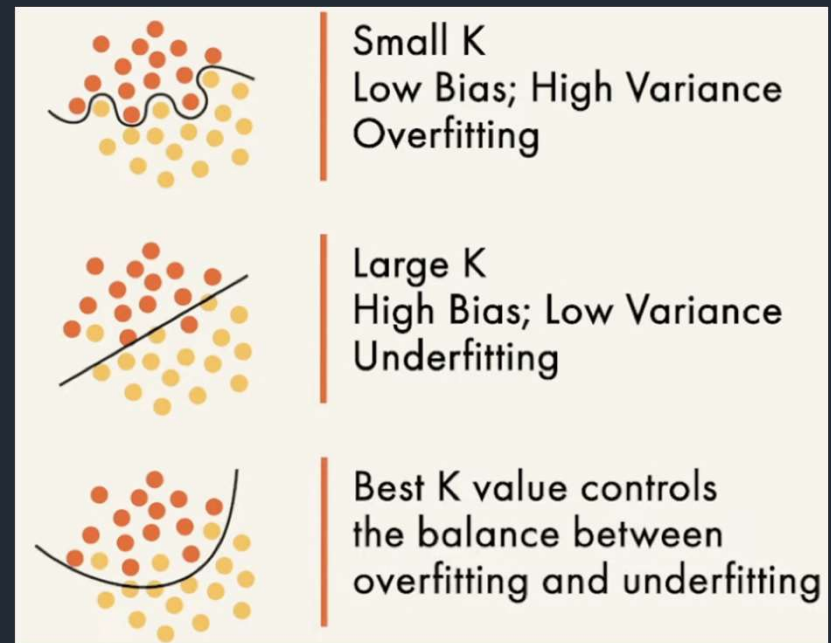Solution: implementation of an algorithm that recognizes closed loops

# Problem Nr. 2

2. Principal components: Tradeoff run-time vs. variance

many PCs                                                                few PCs

high explanation                                                        low explanation of
of variance                                                             variance

high run-time                                                           low run-time

# Problem Nr. 3

3. over- and underfitting of k-value

→ write an algorithm which
   determines best k-value



https://youtu.be/0p0o5cmgLdE

# Thank you for listening

## Questions?