



Computer Vision Project – Summer Term 2021

Face Recognition

Dr.-Ing. Thomas Köhler

Pattern Recognition Lab, Friedrich-Alexander-Universität (FAU) Erlangen-Nürnberg

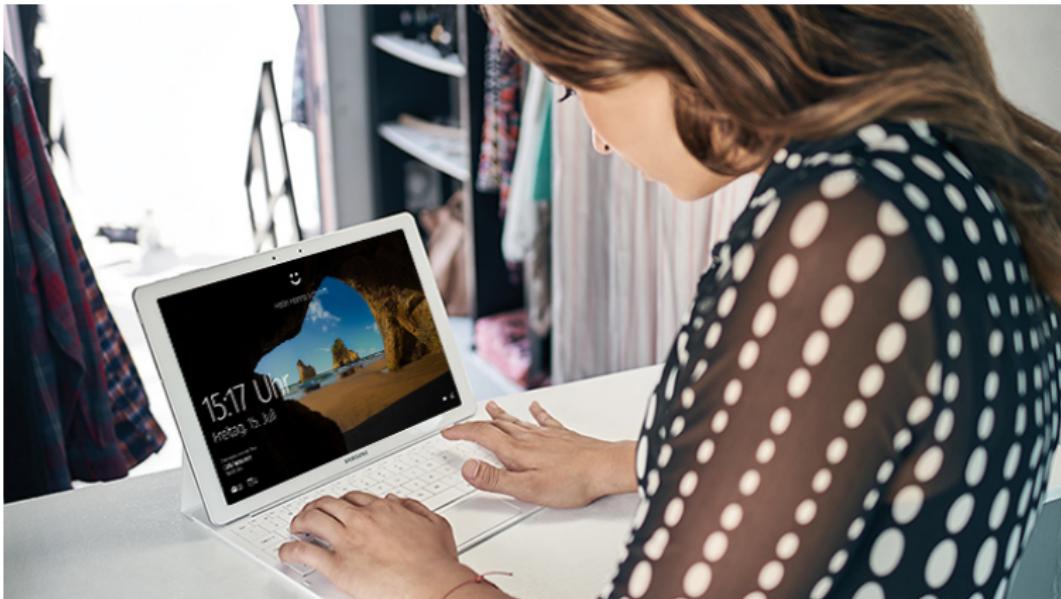
Augmented Reality & Cognitive Services, e.solutions GmbH, Erlangen

May 31, 2021



Face Recognition Use Cases

Microsoft Windows Hello: Windows login via face recognition¹



¹ <https://www.microsoft.com/de-de/windows/windows-hello>

Face Recognition Use Cases

Apple Face ID: unlock smart phone via face recognition²



²<https://www.apple.com/iphone-x/#face-id>

Face Recognition Use Cases

Identification of Vehicle Occupants and Personalization of Vehicle Settings



Scope of this Project

In the lecture:

- Learn how to design and evaluate the core of current facial recognition systems from a technical point of view
- Overview of modern machine learning methods in this field
- Discussion of strengths and limitations of such methods

In the exercise:

- Implementation of a simple system comprising basic functionality for face verification, identification, and clustering from webcam videos
- Evaluation of face recognition algorithms

Outline

Introduction

Face Representation

Eigenfaces

Fisherfaces

Deep Features

Selected Topics in Face Recognition

Distance Measures and Face Verification

Face Identification

Face Clustering

Evaluating Face Recognition Systems

Summary



FAU

FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

FACULTY OF ENGINEERING

Introduction



The Face as a Biometric Marker

Biometric markers for identification³:

- Fingerprint
- Iris
- Speech
- Face

Face vs. other markers:

- Face the only marker that can be captured at large distances
- Only low-cost hardware required

³Jain, Anil K., and Stan Z. Li. Handbook of face recognition. New York, Springer, 2011.

Challenges

Intrinsic conditions:

- High intra-subject variation: facial expressions, eye glasses, changes in facial hair, changes in pose, aging
- Low inter-subject variation: similar skin or hair color, similar eye glasses

Extrinsic conditions:

- Varying illumination conditions: images captured at day and at night
- Image quality: images/videos processed with different codecs

Some Difficult Examples

Frontal and profile views⁴:



⁴ Example from CMU Multi-PIE database: <http://www.cs.cmu.edu/afs/cs/project/PIE/MultiPie/Multi-Pie/Home.html>

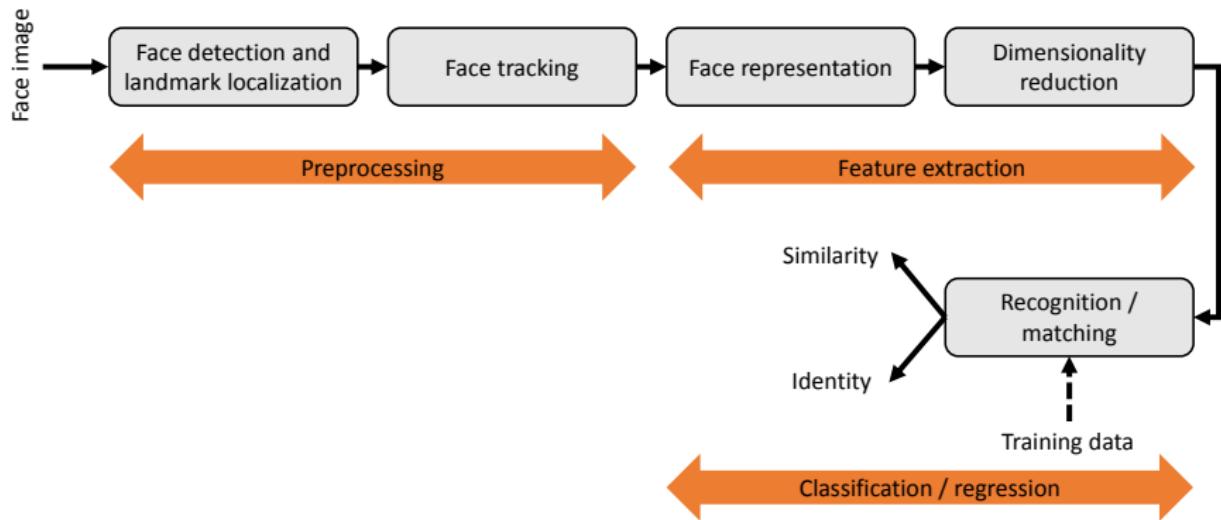
Some Difficult Examples

Close family relationships (e. g. father and son, twins):



Face Recognition Pipeline

Approaching face recognition via machine learning techniques:

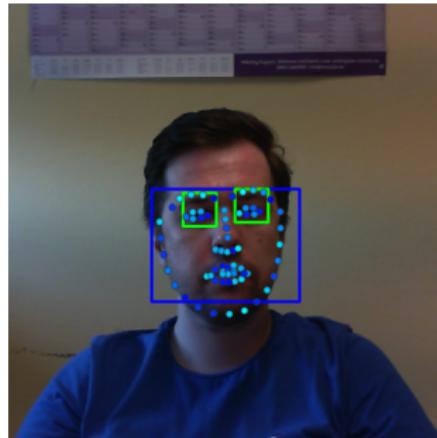


Adopting common pattern recognition pipeline: preprocessing, feature extraction, classification/regression

Face Detection and Facial Landmark Localization

Initial steps of the face recognition pipeline:

- Detection of face region
- Detection of eye regions
- Extracting landmarks to model pose and facial expression
- Alignment of face image using detected landmarks for data normalization
 - Compensate for different head poses
 - Eyes, nose, and mouth at predefined positions



Face Tracking

Continuous monitoring of a human face over time:



- Track face bounding box and landmark positions over time in video data
- 2-D approach: template matching and related motion estimation techniques
- 3-D approach: constrained local model (cf. active shape models (AAM))
 - Describes shape variations of the face
 - Learned from annotated face exemplars
 - Enables the estimation of the head pose

Methodologies of Face Recognition Systems

Geometric face representations:

- Distances, areas, or angles between salient points (eyes, nose, mouth)
- Obtained by feature detection algorithms
- Requires engineering of meaningful features (hand-crafted features)



Data-driven face representations:

- Image intensities form raw features
- Learn suitable representation from exemplars (face manifold)
- Requires training from large datasets

Problem Statements in Face Recognition

Verification (one-to-one matching):

- Given one probe image and one gallery image
- Check if both images show the same identity

Identification (one-to-many matching):

- Given one probe image and a set of gallery images with identity labels
- Retrieve identity of probe image

Clustering (many-to-many matching):

- Given a set of unlabeled face images
- Cluster all images according to the identities captured in the data



FAU

FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

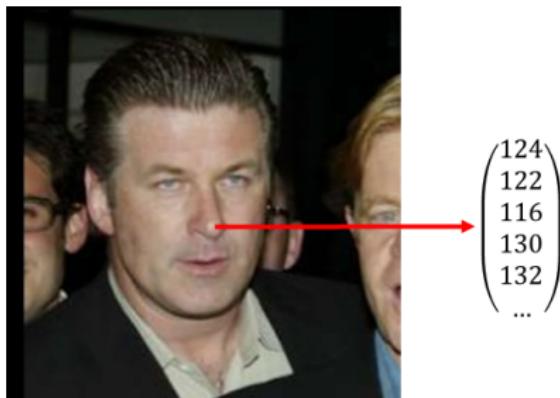
FACULTY OF ENGINEERING

Face Representation



Data-Driven Face Representations

How to represent faces in digital images?

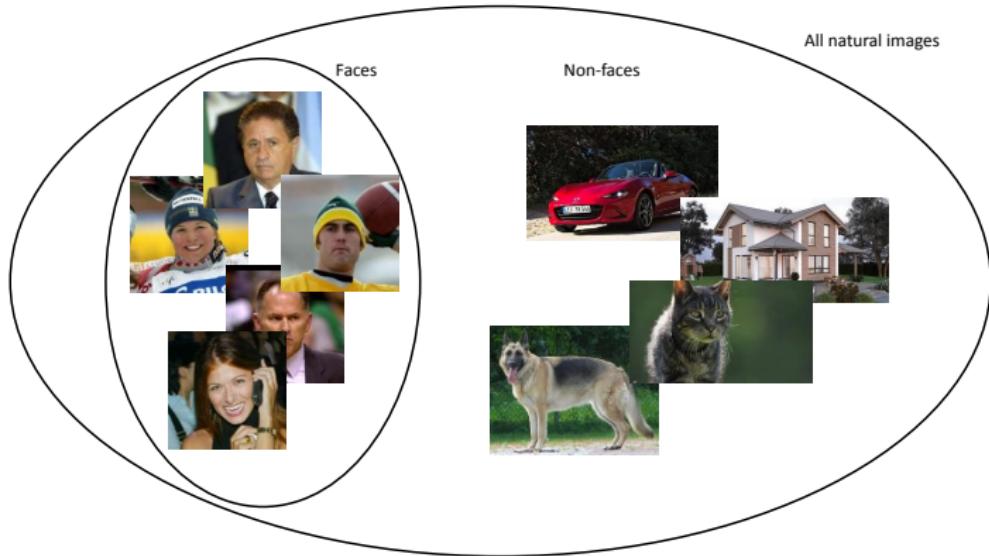


- Grayscale image with $M \times N$ pixels and b bits per pixel (e.g. 64×64 , 8 bit)
 $\rightarrow (2^b)^{M \cdot N}$ different images (e.g. $256^{4096} \gg 10^{1000}$)
- Only a small fraction of images corresponds to valid faces

Manifold of Face Images

Faces live on a low-dimensional manifold

- Learn the manifold form exemplar data
- Face detection: distinguish faces and non-faces
- Face recognition: distinguish faces of different subjects



Eigenfaces

- Let $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_i \in \mathbb{R}^D$ be a set of n face images with mean:

$$\mu_x = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (1)$$

and covariance:

$$\Sigma_x = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mu_x)(\mathbf{x}_i - \mu_x)^\top \quad (2)$$

- Project each \mathbf{x}_i to a M -dimensional space ($M \ll D$):

$$\mathbf{y}_i = \mathbf{W}\mathbf{x}_i \in \mathbb{R}^M \quad (3)$$

- Seek linear transform $\mathbf{W} \in \mathbb{R}^{M \times D}$ that maximizes the variance of $\mathbf{y}_1, \dots, \mathbf{y}_n$:

$$\begin{aligned} \Sigma_y &= \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \mu_y)(\mathbf{y}_i - \mu_y)^\top \\ &= \mathbf{W}^\top \Sigma_x \mathbf{W} \end{aligned} \quad (4)$$

→ Principal Component Analysis (Principal components ≡ Eigenfaces)

Fisherfaces

- Eigenfaces ignore class labels to construct the feature transform
- Fisherfaces exploit two constraints (one unique face \equiv one class):
 1. The between-class scatter is maximized
 2. The within-class scatter is minimized
- Seek the transform $\mathbf{W} \in \mathbb{R}^{M \times D}$ that maximizes Fisher's linear discriminant:

$$J(\mathbf{W}) = \frac{\mathbf{W}^\top \Sigma_{\text{inter}} \mathbf{W}}{\mathbf{W}^\top \Sigma_{\text{intra}} \mathbf{W}} \quad (5)$$

Intra-class and inter-class scatter for c classes:

$$\Sigma_{\text{intra}} = \sum_{i=1}^c \frac{1}{|\mathcal{X}_i|} \sum_{\mathbf{x}_k \in \mathcal{X}_i} (\mathbf{x}_k - \boldsymbol{\mu}_i)(\mathbf{x}_k - \boldsymbol{\mu}_i)^\top \quad (6)$$

$$\Sigma_{\text{inter}} = \sum_{i=1}^c (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top \quad (7)$$

$\boldsymbol{\mu}_i$ is the mean of all faces \mathcal{X}_i in the i -th class and $\boldsymbol{\mu}$ is the overall mean face

Deep Features

Limitations of the methods discussed so far:

- Determine face representations under a linear model
- Linear models feature limited robustness against pose or illumination variations

Extension:

- Non-linear face representation $f_\theta(\mathbf{x})$:

$$\mathbf{y} = f_\theta(\mathbf{x}) \tag{8}$$

- The transform $f_\theta(\mathbf{x})$ is implemented by a deep neural network
- Parameters θ are learned from exemplars

Convolutional Neural Networks (CNNs)

Neural network:

- Computation graphs comprising neurons
- Propagation of input signals through the network to obtain outputs

Network design with input, output, and hidden layers:

- Convolutional layer: convolution of input neuron activations with filter kernel
- Pooling layer: fusion of clusters of input neuron activations
- Locally/fully connected layer: weighted sum of all input neuron activations

θ : Parameters of all layers in the network

Face Representation via Classification

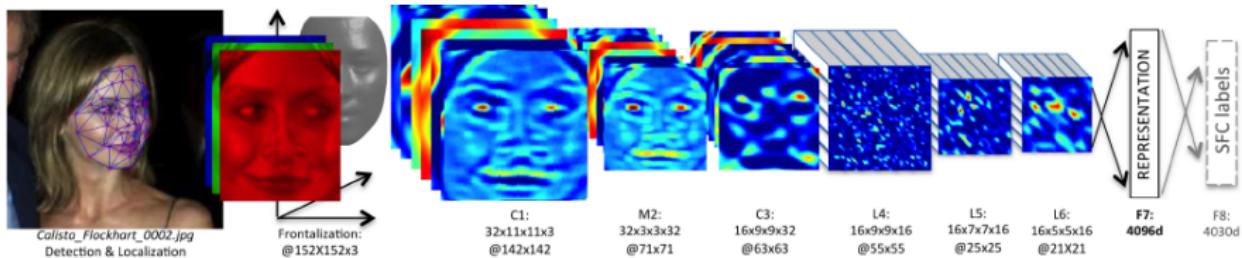
- Design neural network to classify face images according to their identity
- Fully connected output layer (\mathbf{W}, \mathbf{b}) to model probability distribution over c classes:

$$p_i = \frac{\exp(\mathbf{w}_i^\top \mathbf{x} + b_i)}{\sum_{j=1}^c \exp(\mathbf{w}_j^\top \mathbf{x} + b_j)}, \quad i = 1, \dots, c \quad (9)$$

→ Softmax activation

- Train network parameters on face exemplars by minimizing misclassification error (e.g. cross entropy loss)
- Activation of fully connected layer \equiv face representation (aka. embedding)

Example: Facebooks DeepFace (trained on 4M faces)⁵



- Face frontalization for preprocessing
- Eight layer network: convolutional (C1, C3), max-pooling (M2), locally connected (L4, L4, L6), and fully connected layer (F7, F8)
- Feature maps model face from high-level to low-level
- Activation of F7 fully connected layer used as face representation

⁵Taigman, Yaniv, et al. "Deepface: Closing the gap to human-level performance in face verification." Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. 2014.

Face Representation via Regression

Discussion of the classification-based approach:

- Classification requires class labels (identities)
- Models discriminative features for face recognition only implicitly
- Learned features are not necessarily optimal

Extension to regression-based face representation:

- Training of deep neural networks for classification
- Fine-tuning by regression-based loss

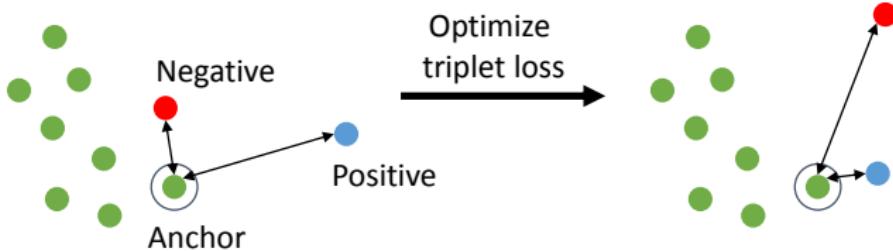
Triplet Loss

- Let \mathbf{x}_i^a be an anchor face, \mathbf{x}_i^p a face of the same subject (positive) and \mathbf{x}_i^n be a face of different subject (negative)
- Ensure for all triplets $(\mathbf{x}_i^a, \mathbf{x}_i^p, \mathbf{x}_i^n)$ with margin α :

$$\underbrace{\|\mathbf{f}_\theta(\mathbf{x}_i^a) - \mathbf{f}_\theta(\mathbf{x}_i^p)\|_2^2}_{\text{anchor-to-positive distance}} + \alpha < \underbrace{\|\mathbf{f}_\theta(\mathbf{x}_i^a) - \mathbf{f}_\theta(\mathbf{x}_i^n)\|_2^2}_{\text{anchor-to-negative distance}} \quad (10)$$

- Penalize distances of positive and negative samples w.r.t. the anchor

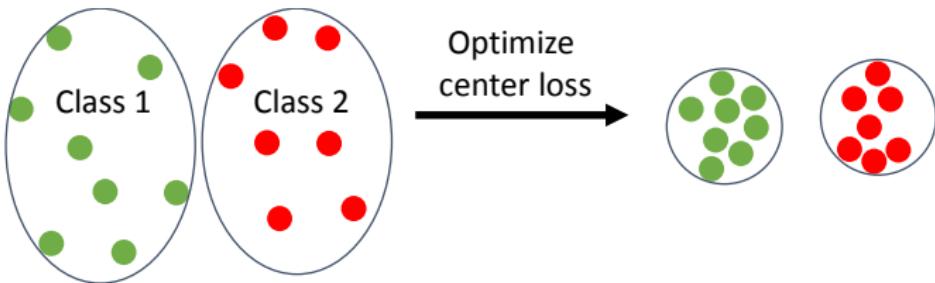
$$\mathcal{L}_{\text{triplet}}(\theta) = \sum_{i=1}^n \|\mathbf{f}_\theta(\mathbf{x}_i^a) - \mathbf{f}_\theta(\mathbf{x}_i^p)\|_2^2 - \|\mathbf{f}_\theta(\mathbf{x}_i^a) - \mathbf{f}_\theta(\mathbf{x}_i^n)\|_2^2 + \alpha \quad (11)$$



Center Loss

- Penalize variations of deep features \mathbf{x}_i around their class center (mean) μ_i

$$\mathcal{L}_{\text{center}}(\theta) = \sum_{i=1}^n \|\mathbf{f}_\theta(\mathbf{x}_i) - \mu_i\|_2^2 \quad (12)$$



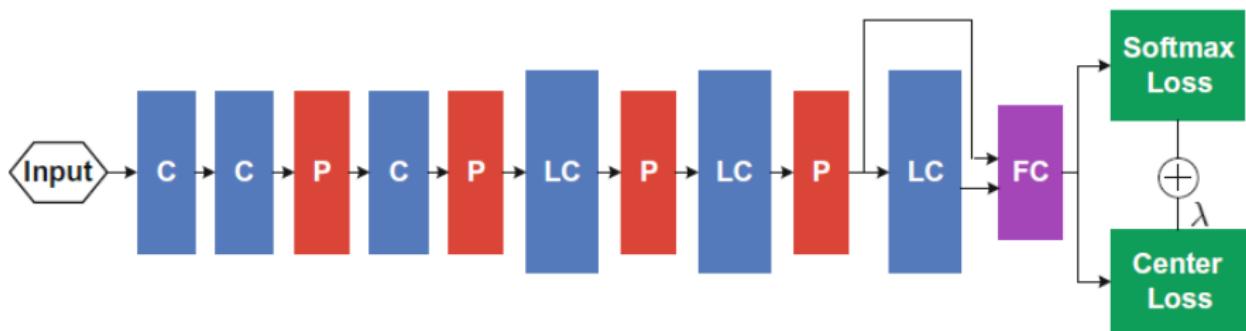
- Similar to Fisherfaces but used with a non-linear model
- Does not require the handling of triplets (combinatorial explosion)

Combining Classification and Regression Approaches

Joint supervision with cross entropy and center loss:

$$\mathcal{L}(\theta) = \mathcal{L}_{\text{ce}}(\theta) + \lambda \mathcal{L}_{\text{center}}(\theta) \quad (13)$$

Example architecture of Wen et al.⁶:

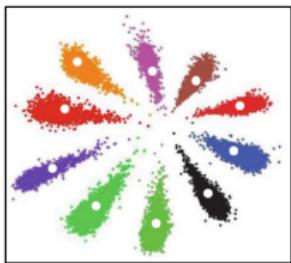


Network with convolutional (C), max-pooling (P), local connected (LC), and fully connected layer (FC)

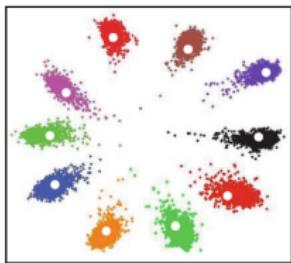
⁶Wen, Yandong, et al. "A discriminative feature learning approach for deep face recognition." European Conference on Computer Vision. Springer, 2016.

Combining Classification and Regression Approaches

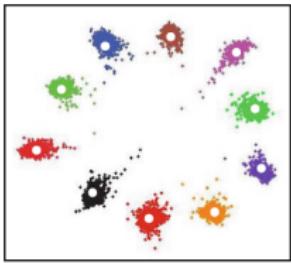
Controlling discriminative power of deep features in the method of Wen et al.:



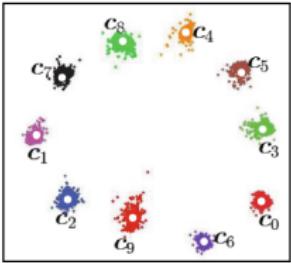
(a) $\lambda = 0.001$



(b) $\lambda = 0.01$



(c) $\lambda = 0.1$



(d) $\lambda = 1$

Larger center loss weight (large λ) \Rightarrow higher discriminative power of deep features

Practical Considerations on Deep Features

Learning:

- Learn $f_\theta(\mathbf{x})$ on large datasets (Facebook: 4M faces, Google: 200M faces)
- Computationally very demanding
- Implemented on graphics processing units (GPU)

Inference:

- Inference of face representation $f_\theta(\mathbf{x})$ by forward pass
- Efficient to compute using additional optimizations (weight quantization)
- Use face representation within lightweight classification or clustering models
- Can be implemented on embedded devices



FAU

FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

FACULTY OF ENGINEERING

Selected Topics in Face Recognition



Face Verification

- Given two face images \mathbf{x}_1 and \mathbf{x}_2 , is the image pair of the same person?
- Define a suitable distance measure $d(\mathbf{x}_1, \mathbf{x}_2)$
 - Use face representation instead of the image space
 - Example: cosine distance

$$d(\mathbf{x}_1, \mathbf{x}_2) = \frac{||f_{\theta}(\mathbf{x}_1) - f_{\theta}(\mathbf{x}_2)||_2^2}{||f_{\theta}(\mathbf{x}_1)||_2^2 + ||f_{\theta}(\mathbf{x}_2)||_2^2} \quad (14)$$

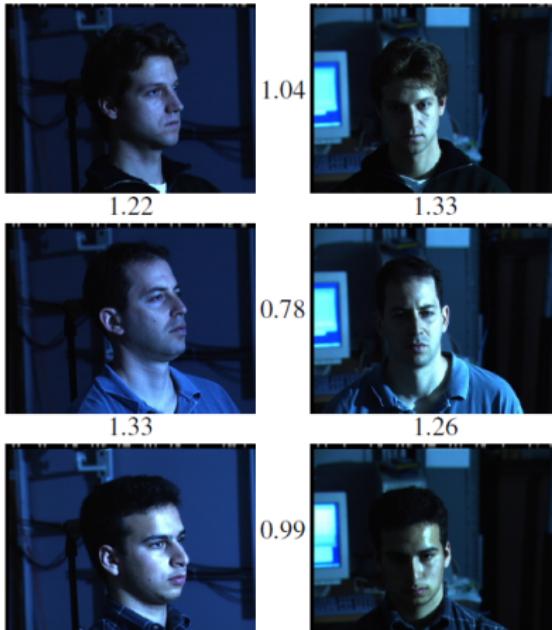
- Verification using the distance measure and threshold $\tau > 0$:

$$v(\mathbf{x}_1, \mathbf{x}_2) = \begin{cases} \text{same} & \text{if } d(\mathbf{x}_1, \mathbf{x}_2) \leq \tau \\ \text{not same} & \text{if } d(\mathbf{x}_1, \mathbf{x}_2) > \tau \end{cases} \quad (15)$$

- Alternatively, use similarity measure, e. g. $s(\mathbf{x}_1, \mathbf{x}_2) = -d(\mathbf{x}_1, \mathbf{x}_2)$

Face Verification

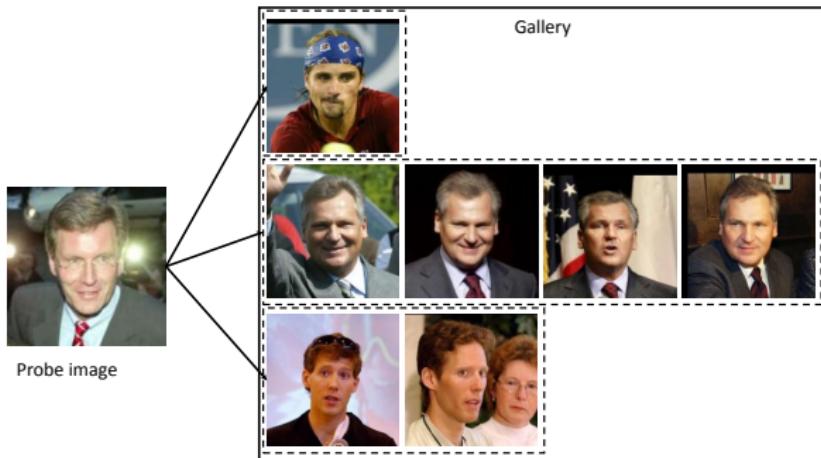
Example: Face identification using Google FaceNet embeddings⁷



⁷ Schroff, Florian, Dmitry Kalenichenko, and James Philbin. "Facenet: A unified embedding for face recognition and clustering." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.

Face Identification

- Given a probe image x_{probe} of unknown identity
- Determine identity $I(x_{\text{probe}})$ from labeled images in a gallery



Two different protocols:

- Closed-set: all possible identities of probes are contained in the gallery
- Open-set: identities of some probes are missing in the gallery

Face Identification with Closed-Set Protocol

Repeated face verification using k-nearest neighbors (k-NN) classifier:

- For the probe image $\mathbf{x}_{\text{probe}}$, find the k closest gallery images $\mathbf{x}_1, \dots, \mathbf{x}_k$ according to a distance $d(\mathbf{x}_{\text{probe}}, \cdot)$
- Identity of $\mathbf{x}_{\text{probe}}$ is the identity of the majority (mode) in $\mathbf{x}_1, \dots, \mathbf{x}_k$:

$$I(\mathbf{x}_{\text{probe}}) = \text{mode}(I(\mathbf{x}_1), \dots, I(\mathbf{x}_k)) \quad (16)$$

- Alternatively, we can consider similarities $s(\mathbf{x}_{\text{probe}}, \cdot) = -d(\mathbf{x}_{\text{probe}}, \cdot)$

Other discriminative classification models:

- Support vector machine (SVM)
- Random forests
- Boosting methods

Face Identification with Open-Set Protocol

Handle open-set space with distance-based approach (thresholded NN):

- In addition, extract the best matching gallery image $\mathbf{x}_{\text{match}}$ with minimum distance to $\mathbf{x}_{\text{probe}}$
- Possibly assign "unknown" identity to $\mathbf{x}_{\text{probe}}$:

$$I(\mathbf{x}_{\text{probe}}) = \begin{cases} \text{mode}(I(\mathbf{x}_1), \dots, I(\mathbf{x}_k)) & \text{if } d(\mathbf{x}_{\text{probe}}, \mathbf{x}_{\text{match}}) \leq \tau \\ \text{unknwon} & \text{if } d(\mathbf{x}_{\text{probe}}, \mathbf{x}_{\text{match}}) > \tau \end{cases} \quad (17)$$

- τ is the face verification threshold

Handle open-set space via extreme value theory⁸:

- Statistical model for inclusion probability of probe images in classes contained in the gallery
- Thresholding of inclusion probabilities instead of using distances

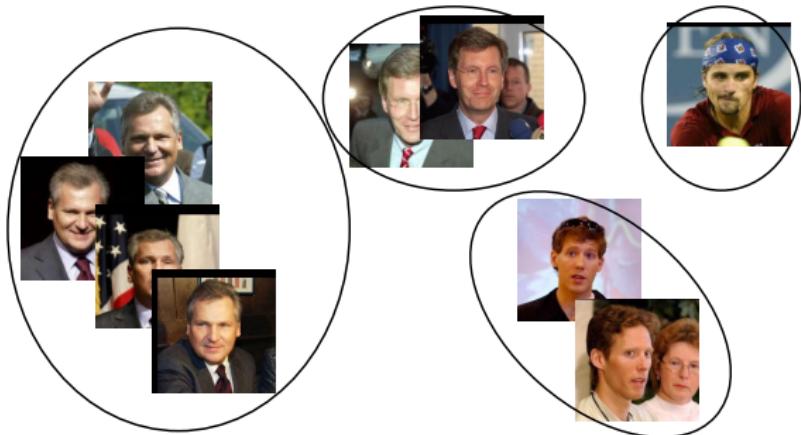
⁸Günther, Manuel, et al. "Toward open-set face recognition." Conference on Computer Vision and Pattern Recognition (CVPR) Workshops, 2017.

Face Clustering

- Given n face images $\mathbf{x}_1, \dots, \mathbf{x}_n$, cluster them into $k \leq n$ clusters $\mathcal{C}_1, \dots, \mathcal{C}_k$
- Minimize cluster scatter with cluster centers μ_i and distance measure $d(\cdot, \cdot)$:

$$(\mathcal{C}_1, \dots, \mathcal{C}_k) = \operatorname{argmin}_{\mathcal{C}_1, \dots, \mathcal{C}_k} \sum_{i=1}^k \sum_{\mathbf{x} \in \mathcal{C}_i} d(\mathbf{x}, \mu_i) \quad (18)$$

- k -means clustering: $d(\mathbf{x}, \mu_i) = \|\mathbf{x} - \mu_i\|_2^2$ and $\mu_i \equiv$ cluster mean



Clustering using k -Means Algorithm

Iterative algorithm:

1. Assignment: Assign each face to the cluster with closest center

$$\mathcal{C}_i^t = \{\mathbf{x} : \|\mathbf{x} - \mu_i\|_2^2 \leq \|\mathbf{x} - \mu_j\|_2^2 \text{ for all } i \neq j\} \quad (19)$$

2. Update: Re-calculate the cluster centers from current clustering $\mathcal{C}_1^t, \dots, \mathcal{C}_k^t$

$$\mu_i^{t+1} = \frac{1}{|\mathcal{C}_i^t|} \sum_{\mathbf{x} \in \mathcal{C}_i^t} \mathbf{x} \quad (20)$$

Other clustering approaches⁹:

- Soft clustering: membership degrees instead of assignment to single clusters
- Agglomerative clustering: adaptive selection of the number of clusters

⁹Otto, C., Wang, D., and Jain, A. K. (2018). Clustering millions of faces by identity. IEEE Transactions on Pattern Analysis and Machine Intelligence, 40(2), 289-303.

Evaluating Face Recognition Systems

Face identification with closed-set protocol:

- Rank of a probe image $\mathbf{x}_{\text{probe}}$ for gallery \mathcal{G} with true matching image $\mathbf{x}_{\text{match}}$:

$$\text{Rank}(\mathbf{x}_{\text{probe}}) = \left| \left\{ \mathbf{x}' \in \mathcal{G} : s(\mathbf{x}', \mathbf{x}_{\text{probe}}) \geq s(\mathbf{x}_{\text{match}}, \mathbf{x}_{\text{probe}}) \right\} \right| \quad (21)$$

$\text{Rank}(\mathbf{x}_{\text{probe}}) = 1$ if $\mathbf{x}_{\text{probe}}$ is correctly associated with $\mathbf{x}_{\text{match}}$

- Rank- k Identification rate on a test set \mathcal{T} :

$$\text{IR}(r) = \frac{\left| \left\{ \mathbf{x}' \in \mathcal{T} : \text{Rank}(\mathbf{x}') \geq r \right\} \right|}{|\mathcal{T}|} \quad (22)$$

For $r = 1$ it is equivalent to the accuracy

- Depict the identification rate $\text{IR}(r)$ versus the rank r

Evaluating Face Recognition Systems

Face identification with open-set protocol:

- Trade-off between true identifications and unknowns that are incorrectly detected as knowns (false alarms) depending on similarity threshold θ
- Detection and identification rate for test set of knowns \mathcal{K} at rank r :

$$\text{DIR}(\theta) = \frac{|\{\mathbf{x}' \in \mathcal{K} : s(\mathbf{x}', \mathbf{x}_{\text{gallery}}) \geq \theta \text{ and } \text{Rank}(\mathbf{x}') \geq r\}|}{|\mathcal{K}|} \quad (23)$$

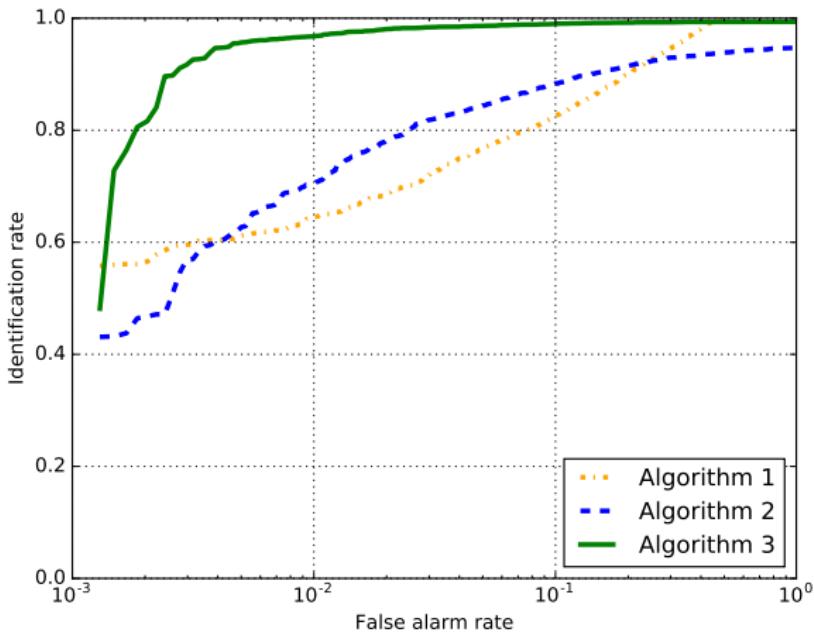
- False alarm rate for complementary test set of unknowns \mathcal{U} :

$$\text{FAR}(\theta) = \frac{|\{\mathbf{x}' \in \mathcal{U} : s(\mathbf{x}', \mathbf{x}_{\text{gallery}}) \geq \theta \text{ for any } \mathbf{x}_{\text{gallery}} \in \mathcal{G}\}|}{|\mathcal{U}|} \quad (24)$$

- Depict $\text{DIR}(\theta)$ at different $\text{FAR}(\theta)$ for a given rank r (typically $r = 1$)

Example: Comparing Different Face Recognition Algorithms

DIR curve with semi-logarithmic axes for three algorithms (shown for rank $r = 1$):





FAU

FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

FACULTY OF ENGINEERING

Summary



Take Home Messages

- Widespread application domains of face recognition
- Still a hard problem under uncontrolled conditions (e. g. difficult poses)
- Face representation is a key component for modern face recognition systems
 - Learned on large training datasets
 - Different methodologies: Eigenfaces, Fisherfaces, deep features
- Recognition tasks (verification, identification, clustering) solved by common machine learning algorithms
 - Based on suitable face representation
 - Today also applicable on embedded devices

Further Readings

Overview on face image analysis and recognition techniques:

Anil K. Jain and Stan Z. Li. "Handbook of face recognition". Springer, 2011

Eigenfaces, Fisherfaces, and other classical methods:

Peter N. Belhumeur, João P. Hespanha, and David J. Kriegman. "Eigenfaces vs. Fisherfaces: Recognition using class specific linear projection." IEEE Transactions on Pattern Analysis and Machine Intelligence 19(7), 1997, 711-720.

Deep learning based methods:

- Yaniv Taigman *et al.* "Deepface: Closing the gap to human-level performance in face verification." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2014
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. "Facenet: A unified embedding for face recognition and clustering." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015
- Yandong Wen *et al.* "A discriminative feature learning approach for deep face recognition." European Conference on Computer Vision, 2016.

Open Project Topics (Collaboration with e.solutions Erlangen)

Improving face representation learning using adversarial samples (5/10 ECTS)

- Adversarial attack: fool face recognition system by manipulating input data
- Project goal: investigate benefit of adversarial samples for data augmentation
- Begin: February 2021

Learning image formation models for super-resolution (5/10 ECTS)

- Learning-based super-resolution aims at inferring mappings from low-resolution to high-resolution images
- Classical reconstruction-based methods use reversed mappings for inverse problem formulations
- Project goal: learn reversed mapping for reconstruction algorithms
- Begin: February 2021

All topics are extendable to a master's thesis

Thanks for listening.
Any questions?