

# Computational Intelligence Lat Zusammenfassung

## - Linear Autoencoder

- $C \hat{=} \text{Encoder}$

- $D \hat{=} \text{Decoder}$

- Loss:  $L(x, \theta) = \frac{1}{2} \|x - \hat{x}(\theta)\|^2$

- $\hat{x}(\theta) = D \cdot C \cdot x$

- Parameters:  $D$  and  $C$

- $z$  is bottleneck layer / latent representation

- $z$  is bottleneck  $\Rightarrow$  reduces rank to  $k \Rightarrow$  loss of information

- $\text{rank}(C \cdot D) \leq \min\{\text{rank}(C), \text{rank}(D)\} \leq k$

## - Eckart-Young Theorem

- $\arg \min \|x - \hat{x}\|_F^2 = U \cdot \Sigma_k \cdot V^T$

- You get the best rank  $k$  approximation of a matrix when you take the SVD and the first  $k$  dominant singular values

## - Singular Value Decomposition (Recommender System)

- $U/V$ : - Left/right singular vectors

- Orthogonal

- $U \cdot U^T = I$

- $U^T = U^{-1}$

- $\Sigma$ : - Singular values

- $\text{diag}(\sigma_1, \dots, \sigma_s) \quad (s \leq \min(m, n))$

$$\begin{matrix} A \\ m \times n \end{matrix} = \begin{matrix} U \\ m \times m \end{matrix} \cdot \begin{matrix} \Sigma \\ m \times n \end{matrix} \cdot \begin{matrix} V^T \\ n \times n \end{matrix}$$

## - Optimal Autoencoder

- $C^* = U_K^T$
- $D^* = U_K$

$$\begin{aligned} \hat{x} &= D \cdot C \cdot x = D \cdot C \cdot (U \cdot \Sigma \cdot V^T) \\ &= U_K \cdot U_K^T \cdot U \cdot \Sigma \cdot V^T \\ &\approx U \cdot I_K \cdot \Sigma \cdot V^T \\ &= U \cdot \Sigma_K \cdot V^T \end{aligned}$$

## - Principle Eigenvector Decomposition (image compression)

- Principle eigenvector has the highest trend direction

- $A$  must be squared matrix

- $u$  is eigenvector of  $A$  iff  $A \cdot u = \lambda \cdot u$

$$\hookrightarrow A \cdot u = \lambda \cdot u$$

- $\lambda$  is called eigenvalue, stored in  $\Lambda$

- We preserve largest data variance

- $\Sigma = U \cdot \Lambda \cdot U^T$

- $\Sigma$  is symmetric:  $\Sigma = \Sigma^T$

$\hookrightarrow$  Variance covariance matrix:  $\Sigma = \frac{1}{n} \cdot \bar{x} \cdot \bar{x}^T = \frac{1}{n} \sum_{i=1}^n x_i \cdot x_i^T$

- diag: variance between same feature

- off-diag: variance between two features

- $\Sigma$  is p.s.d  $\Rightarrow$  all eigenvalues are positive

1. Remove mean  $\tilde{x} = x - \bar{x}$  ( $m \times n$ )

2. Get  $\Sigma = \frac{1}{n} \cdot \tilde{x} \cdot \tilde{x}^T$  ( $m \times m$ )

3. Pick  $d$  eigenvectors  $u_d$  from  $\Sigma$ :  $\Sigma = U \cdot \Lambda \cdot U^T$

4. Reduced matrix  $z = \underbrace{U_d^T}_{d \times m} \cdot \underbrace{\tilde{x}}_{m \times n}$

4.  $\tilde{x} = U_d \cdot z$

## - Power Method

$$\cdot v_{t+1} = \frac{A \cdot v_t}{\|A \cdot v_t\|}$$

$$\cdot \lim_{t \rightarrow \infty} v_t = u_1$$

## - Iterative View

- Find eigenvector of

$$(\Sigma - \lambda \cdot U U^T) \quad (= 2nd EV of \Sigma)$$

- Iterate to find lower EV

- Detect knee in eigenspectrum to choose  $d$

- PCA vs. Autoencoder: - PCA is unique and interpretable

- From SVD to PCA:  $A \cdot A^T = (U \cdot \Sigma \cdot V^T) \cdot (U \cdot \Sigma \cdot V^T)^T$

$$\Rightarrow U \cdot \Sigma \cdot V^T \cdot V \cdot \Sigma^T \cdot U^T$$

$$= U \cdot \underbrace{\Sigma \cdot \Sigma^T}_{\Lambda} \cdot U^T = U \cdot \Sigma^2 \cdot U^T$$

$$\hookrightarrow \Sigma^2 = \Lambda$$

## - Lagrange Multiplier

- can be used for PCA

- $L(u_1, \lambda_1) = u_1^T \Sigma \cdot u_1 - \lambda_1 (u_1^T u_1 - 1)$   
 $\underbrace{u_1 \text{ needs to be unit length}}$

$$\nabla_{u_1} L(u_1, \lambda_1) = 2(\Sigma \cdot u_1 - \lambda_1 \cdot u_1) = 0$$

$$\Sigma \cdot u_1 - \lambda_1 \cdot u_1 = 0$$

$$\Sigma \cdot u_1 = \lambda_1 \cdot u_1 \quad (= \text{Definition of Eigenvector})$$

- If we want to get second eigenvector we need to add that  $u_1$  and  $u_2$  must be orthogonal:  $u_2^T \cdot u_1 = 0$

## - Recommender System

- Matrix completion
- Find best low-rank approximation
- Frobenius norm:  $\|A\|_F = \sqrt{\sum \sum a_{ij}^2}$
- Sum of all entries
- $k$ -dimensional: number of latent factors
- SVD of rating matrix:
- $U$ : User-to-factor associations
- $\Sigma$ : Dominance of each factor
- $V$ : items-to-factor associations

18,000 movies						
480,000 users	x	1	1	x	...	x
	x	x	x	5	...	x
	x	x	3	x	...	x
	x	4	3	x	...	2
	...	x	x	x	...	x
	x	5	x	1	...	x
	x	x	3	3	...	x
	x	1	x	x	...	2

Approach over: Only Eigenvalues that are not included

$$\begin{array}{c}
 \text{Cremators} \quad \text{Evil spawn} \quad \text{Fatal justice} \quad \text{Clerks} \quad \text{American pie} \\
 \downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow \\
 \left( \begin{array}{ccccc} 5 & 5 & 5 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 4 & 4 \end{array} \right) = \left( \begin{array}{cc} 0.57 & 0 \\ 0.46 & 0 \\ 0.57 & 0 \\ 0.34 & 0 \\ 0 & 0.52 \\ 0 & 0.66 \\ 0 & 0.52 \end{array} \right) \times \left( \begin{array}{cc} 15 & 0 \\ 0 & 10.67 \end{array} \right) \times \left( \begin{array}{ccccc} 0.57 & 0.57 & 0.57 & 0 & 0 \\ 0 & 0 & 0 & 0.70 & 0.70 \end{array} \right)
 \end{array}$$

Horror      Comedy  
 ↓            ↓  
 Strength of Horror      Strength of Comedy

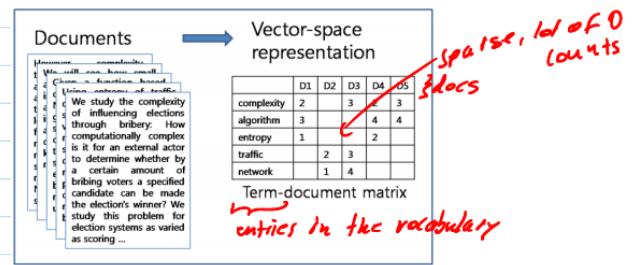
Movie Cremators

U            Σ      V

- For 0 elements we take the mean
- SVD not convex optimization problem
- $\hookrightarrow f(u, v) = 1/I \sum (a_{i,j} - \langle u_i, v_j \rangle)^2 \Rightarrow$  Not convex
- $U \leftarrow \arg \min_u f(u, v) \Rightarrow$  convex
- $V \leftarrow \arg \min_v f(u, v) \Rightarrow$  convex      } fix other variable and alternate
- $A \cdot A^T$ : Similarity between users      } Based on ratings
- $A^T \cdot A$ : Similarity between movies      }
- Frobenius norm only depends on singular values  $\Sigma$
- $\hookrightarrow \|F\|_F^2 = \sum \sigma_i^2$
- Euclidean 2-norm only depends on largest sing. value
- $\hookrightarrow \|A\|_2 = \sigma_1$

## - Document Representation

- Vocabulary: All meaningful words in a language
- Term filtering: Remove misspellings
  - Exclude stop words (the, is, at, ...)
- Term normalization: Reduce word to stem
  - argue (argued, arguing, ...)
- Bag of words: Ignore structure of words
  - ↳ Only co-occurrence counts



## - Probabilistic LSA (= pLSA)

- Get main topics of text
- Document has distribution of topics
- Topic: distribution over terms
  - Soccer topic: high "teams" word probability

1. For each document get the topic

2. For each topic get the words

$$p(w|d) = \sum_z p(w,z|d) = \sum_z p(w|d,z) \cdot p(z|d) \stackrel{z=1}{=} \underbrace{\sum_{\text{topics}}}_{\# \text{topics}} \underbrace{\underbrace{p(w|z)}_{\text{predefined}} \cdot \underbrace{p(z|d)}_{\text{conditional independence assump.}}}_{\text{①}}$$

in doc d & conditional independence assumption.

$X = X_{i,j}$  (→ # of occurrences of  $w_j$  in document  $d_i$ )

$$\ell(U, V) = \sum_{i,j} X_{i,j} \cdot \log p(w_j|d_i) = \sum_{i,j} \log \sum_{z=1}^K p(w_j|z) \cdot p(z|d_i)$$

$Q_{z,i,j} \in [0,1]$ , probability that word  $w_j$  in doc  $d_i$  is generated

$v_{z,j}$        $u_{z,i}$   
 $\underbrace{v_{z,j} \quad u_{z,i}}_{\text{probabilities, - therefore non-negative matrices}}$

from topic  $z$

- Expectation step (solve for optimal  $q$ )

$$\hookrightarrow q_{z,i,j} = \frac{u_{zi} \cdot v_{zj}}{\sum_{k=1}^K u_{ki} \cdot v_{kj}}$$

$\begin{cases} \text{word comes from topic } z \\ \text{word comes from other topics} \end{cases}$

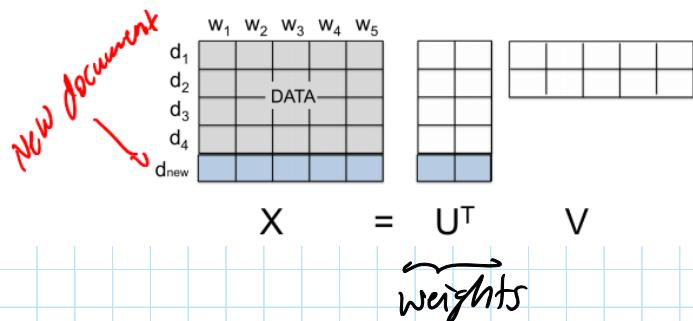
- Maximization step (solve for optimal parameters)

$\hookrightarrow U$ : reestimate probability of word  $w$  for topic  $z$

$U$ : reestimate probability of doc  $i$  covering topic  $z$

### - Latent Dirichlet Allocation

- Generate a new document



- Need probability of  $v$ : doc  $d_{\text{new}}$  covering topic  $z$
- Use dirichlet distribution for generating topic weights  $v$
- $U_i$  contains probability of each topic being in document  $i$
- LDA tries to use an appropriate prior for the language rather than just the one fits all simple 2nd order correlation that pLSA loses.
- Dirichlet prior acts as regularizer
- Need to have good hyperparameter  $\alpha$

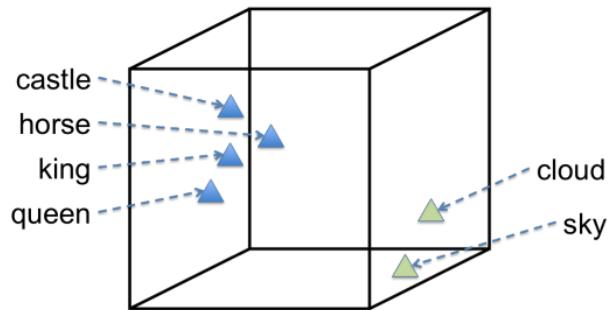
## Word Embeddings

- Predict context word given recent word (=skip-gram)
- $p_\theta(w|w') = \text{probability that } w \text{ occurs in context window } w'$
- Objective:  $\mathcal{L}(\theta; w) = \sum_{t=1}^T \sum_{\Delta \in \Sigma} \log p_\theta(w^{+\Delta} | w^{(t)})$
- Each word is represented in vector space (=Embedding)

$$\cdot p_\theta(w|w') = \langle x_w, x_{w'} \rangle + b_w \quad w \mapsto (x_w, b_w) \in \mathbb{R}^{d+1}, \text{ (vector + bias)}$$

### Context vector:

- Distinguish output vocabulary  $V$  and input vocabulary  $C$
- Output embedding  $x_w$
- Input embedding  $y_w$
- $p_\theta(w|w') = \langle x_w, y_{w'} \rangle + b_w$



### Negative Sampling

- Reduce to binary classification
- Negative samples don't belong together
- $\mathcal{L}(\theta) = \sum_{(i,j) \in \Delta^+} \log \sigma(\langle x_i, y_j \rangle) + \sum_{(i,j) \in \Delta^-} \log \sigma(-\langle x_i, y_j \rangle)$

### GloVe (=try to match co-occurrence count of two words)

- Summarize data in co-occurrence matrix
- $n_{ij} = \# \text{occurrences of } \underbrace{w_i}_{\text{Output}} \text{ in context of } \underbrace{w_j}_{\text{Input}}$
- $n_{\text{castle,king}} = \# \text{occurrences of castle in context of word king}$

- $N$  can be computed in one pass
- Sparse matrix

- Weighted least squares

$$\cdot H(\theta, N) = \sum f(n_{ij}) \cdot \left( \underbrace{\log n_{ij}}_{\text{Target}} - \underbrace{\log p_\theta(w_i) w_j}_{\text{Model to train}} \right)^2$$

- $f(n)$ :
  - downweights word-pairs with small counts
  - cut off at certain threshold  $\Rightarrow$  limit influence

## - Data Clustering

- Assign each data point a cluster
  - (clusters represented by centroids  $u$ )
  - Mapping via nearest centroid rule

$\hookrightarrow \pi(x) = \arg \min_{j=1 \dots k} \|u_j - x\|$

## • Optimization problem

- Find clusters  $U$  and assignments  $\pi$
  - Indicator variable  $z_{i,j} = \begin{cases} 1 & \text{if } \pi(x_i) = j \\ 0 & \text{otherwise} \end{cases}$
  - $\sum_{j=1}^K z_{i,j} = 1$
  - $J(U, Z) = \sum_{i=1}^N \sum_{j=1}^K z_{i,j} \cdot \|x_i - u_j\|^2$

↳ Alternating minimization

• Alternating optimization (=K-means)

- Compute  $Z$ , given  $U$  (=cluster assignment / E-step)
 
$$\hookrightarrow Z_{ij}^* = \begin{cases} 1 & \text{if } j = \arg \min_k \|x_i - U_k\|^2 \\ 0 & \text{otherwise} \end{cases}$$
  - Compute  $U$ , given  $Z$  (=centroid update / M-step)
 
$$\hookrightarrow U_k^+ = \frac{\sum_{n=1}^N z_{kn}^+ \cdot x_n}{\sum_{n=1}^N z_{kn}^+} \quad \begin{cases} \text{Mean of all } x \text{ values} \\ \text{Normalization e.g. } 1/n \end{cases}$$

initialize  $\mathbf{U}$  on  $K$  distinct random data points  
initialize  $\mathbf{Z} \leftarrow \mathbf{Z}^*(\mathbf{U})$

**repeat**

## repeat

$\mathbf{U} \leftarrow \mathbf{U}^*(\mathbf{Z})$  (see above)

$$\mathbf{Z}^{\text{new}} \leftarrow \mathbf{Z}^*(\mathbf{U}) \text{ (se)}$$

$$\text{same} = (\mathbf{Z}^{\text{new}} == \mathbf{Z})$$

$\mathbf{z} \leftarrow \mathbf{z}^{\text{new}}$

**until** (same)

**until** (same)

- because 1 +

= No random initialization

= pay more attention to data

far away

far away

Converges to local minimum, but global minimum is NP-hard

## • Probabilistic clustering

- Instead of hard assignments, use probabilities
- $p(x, \theta) = \sum_{j=1}^K \pi_j \cdot p(x, \mu_j, \Sigma_j)$ 
  - $\hookrightarrow \sum_{j=1}^K \pi_j = 1$  (Mixing coefficient)  
( $\rightarrow$  like prior)
- $p(x, \theta_j) = \text{location \& shape of cluster}$
- $\hat{\theta} = \arg \max_{\theta} \sum_{i=1}^N \log \left[ \sum_{j=1}^K \pi_j \cdot p(x_i, \theta_j) \right]$ 
  - $\curvearrowleft$  Data points  $\curvearrowright$  Clusters
  - $\curvearrowleft$  Expectation maximization
    - $\curvearrowleft$  Points to clusters
    - $\curvearrowright$  Cluster location and shape

## • Generate data point

1. Sample cluster index
2. Sample data point from sampled cluster

## Comparison with K-means

- ▶ Assignments
  - ▶ K-means algorithm: hard assignment points to clusters
  - ▶ EM algorithm: soft assignment based on posteriors
- ▶ Shapes
  - ▶ K-means: spherical cluster shapes, uniform spread
  - ▶ EM algorithm: can learn covariance matrix
- ▶ K-means as a special case
  - ▶ Gaussian mixture model with (fixed) covariances  $\Sigma_j = \sigma^2 \mathbf{I}$
  - ▶ in the limit of  $\sigma \rightarrow 0$ , recover K-means (hard assignments)
  - ▶ can be more formally derived (EM objective  $\rightarrow$  K-means objective)

## - Variational Autoencoder

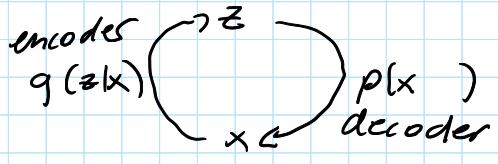
- Generative model: Get  $x$  from  $z$

- Inference model: Get  $z$  from  $x$

$$1. \log p(x)$$

$$2. p(x) = \int p(x|z) \cdot p(z) \cdot dz$$

$$3. p(x) = \int p(x|z) \cdot \frac{p(z)}{q(z|x)} \cdot q(z|x)$$



we get this  
from encoder

$$4. \log p(x) \geq \text{ELBO}(\phi, \theta) = \mathbb{E}_{q\phi} [\underbrace{\log p_{\theta}(x|z)}_{\text{Reconstruction quality}} - \text{KL}(q_{\phi}(z|x) || p(z))]$$

- Maximize  $\theta$ : generative model

- Maximize  $\phi$ : inference model

- Inference model suggests  $z$  for generative model

## - Generative Adversarial Network

$$\cdot \theta^* := \underset{\phi}{\operatorname{argmin}} \left\{ \sup_{\mathcal{G}} \mathcal{L}(\theta, \phi) \right\}$$

$\mathcal{G}$  Generator  $\mathcal{C}$  Classifier

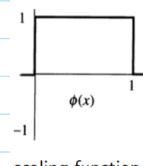
## - Autoregressive models

- Justified by chain rule:  $p(x) = \prod_{t=1}^m p(x_t | x_{1:t-1})$

## - Haar Wavelets

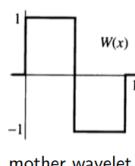
- Original signal  $x$  and orthogonal transform matrix  $U$
- $z = U^T \cdot x = \langle U, x \rangle / \|x\|$
- Energy & distance basis is preserved:  $\|U^T x\|^2 = \|x\|^2$
- Signal compression: Truncate small values of  $z$
- No information is lost
- Can denoise signals

- Haar wavelets good for localized signals



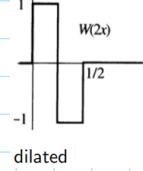
scaling function

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$



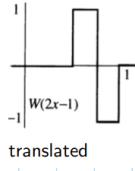
mother wavelet

$$\begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}$$



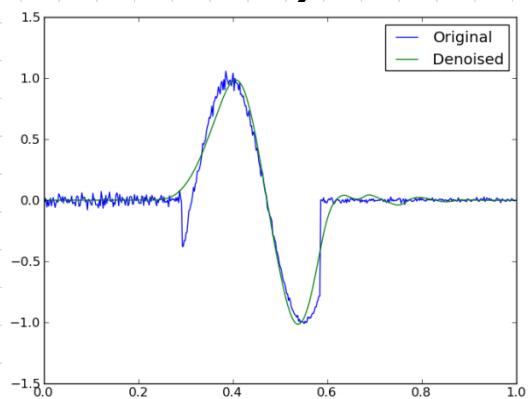
dilated

$$\begin{bmatrix} 1 \\ -1 \\ 0 \\ 0 \end{bmatrix}$$



translated

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ -1 \end{bmatrix}$$



- Fourier basis:
  - Good for sine like signals
  - Global support
  - Not good for localized signals

Fixed transform basis  $U$

- sender/receiver agrees on basis beforehand
- Transmit only non-zero elements of  $\hat{z}$

PCA transform basis  $U$

- Called Hoteling transform
- Use  $K$  largest eigenvectors
- Optimal for given  $\epsilon$

- Basis:
  - Linearly independent vectors
  - Spans the space

## - Overcomplete Dictionaries

- More general dictionaries
- $U \in \mathbb{R}^{D \times L}$ , where  $L > D$  (More atoms than dim.)
- More atoms (dictionary elements) than dimensions
- $U$  is overcomplete, more unknowns than equations
  - ↪ find sparsest  $z \in \mathbb{R}^L$  such that  $x = Uz$
  - ↪  $z^* \in \operatorname{argmin}_z \|z\|_0$
- s.t.  $x = Uz$
- Column  $\equiv$  Atom
- All atoms  $\leq$  dictionary

## - Matching Pursuit

- chooses one atom at a time to greedily reduce approximation error
- Want to minimize residual  $r_{n+1} = r_n - \hat{x}_n$
- Initialize  $r_0 = x$  and  $\hat{x}_0 = 0$
- 1. Find  $j^* = \operatorname{argmax}_j |\langle r_i, u_j \rangle|$
- 2. Compute better approximation  $\hat{x}_{i+1} \leftarrow \hat{x}_i + \langle r_i, u_j \rangle \cdot u_j$
- 3. Update residual  $r_{i+1} \leftarrow r_i - \underbrace{\langle r_i, u_j \rangle}_{\text{Save these in } z} \cdot u_j$

## - Compressive Sensing

- compressing while collecting the data
- $y = W \cdot x = W \cdot U \cdot z$ 
  - ↑ Gaussian random matrix
- $z^* \in \operatorname{argmin}_z \|z\|_0$  such that  $y = \Theta z$   
where  $\Theta = W \cdot U$

## - Dictionary learning ( $\Rightarrow$ Matrix Factorization)

- Mixture of different characteristics in  $U$
- We adapt dictionary to signal characteristics
- $(U, Z) = \underset{U, Z}{\operatorname{arg\,min}} \|X - U \cdot Z\|_F^2$ 
  - ↪ Not jointly convex
- Shared dictionary among many samples
- Coding step:  $z_n \in \operatorname{arg\,min} \|z\|_0$  for  $n=1, \dots, N$ 
  - s.t.  $\|x_n - U^T z\|_2 \leq \sigma \cdot \|x_n\|_2$ 
    - ↑ approx error
- Dictionary step: 1. Fix all atoms except  $u_i$ 
  - 2. zero out atom  $u_i$
  - 3. Look at all points that used  $u_i$
  - 4. Update atom as first singular vector from residual

