Syllabus



Fri. 21.10.	(1)	0. Introduction
		A. Supervised Learning: Linear Models & Fundamentals
Fri. 27.10.	(2)	A.1 Linear Regression
Fri. 3.11.	(3)	A.2 Linear Classification
Fri. 10.11.	(4)	A.3 Regularization
Fri. 17.11.	(5)	A.4 High-dimensional Data
		B. Supervised Learning: Nonlinear Models
Fri. 24.11.	(6)	B.1 Nearest-Neighbor Models
Fri. 1.12.	(7)	B.2 Neural Networks
Fri. 8.12.	(8)	B.3 Decision Trees
Fri. 15.12.	(9)	B.4 Support Vector Machines
Fri. 12.1.	(10)	B.5 A First Look at Bayesian and Markov Networks
		C. Unsupervised Learning
Fri. 19.1.	(11)	C.1 Clustering
Fri. 26.1.	(12)	C.2 Dimensionality Reduction
Fri. 2.2.	(13)	C.3 Frequent Pattern Mining

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Outline

- 1. Principal Components Analysis
- 2. Probabilistic PCA & Factor Analysis
- 3. Non-linear Dimensionality Reduction
- 4. Supervised Dimensionality Reduction



Outline



1. Principal Components Analysis

- 2. Probabilistic PCA & Factor Analysis
- 3. Non-linear Dimensionality Reduction
- 4. Supervised Dimensionality Reduction

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 1. Principal Components Analysis

The Dimensionality Reduction Problem

Given

- ▶ a set \mathcal{X} called **data space**, e.g., $\mathcal{X} := \mathbb{R}^m$,
- a set $X \subseteq \mathcal{X}$ called data,
- ► a function

$$D: \bigcup_{X\subseteq \mathcal{X}, K\in \mathbb{N}} (\mathbb{R}^K)^X \to \mathbb{R}_0^+$$

called **distortion** where D(P) measures how bad a low dimensional representation $P: X \to \mathbb{R}^K$ for a data set $X \subseteq \mathcal{X}$ is, and

▶ a number $K \in \mathbb{N}$ of latent dimensions,

find a low dimensional representation $P : X \to \mathbb{R}^K$ with K dimensions with minimal distortion D(P).



Distortions for Dimensionality Reduction (1/2)

Let $d_{\mathcal{X}}$ be a distance on \mathcal{X} and d_{Z} be a distance on the latent space \mathbb{R}^{K} usually just the Euclidean distance

$$d_Z(v,w) := ||v-w||_2 = (\sum_{i=1}^K (v_i - w_i)^2)^{\frac{1}{2}}$$

Multidimensional scaling aims to find latent representations *P* that **reproduce the distance measure** $d_{\mathcal{X}}$ as good as possible:

$$egin{aligned} D(P) &:= rac{2}{|X|(|X|-1)} \sum_{x,x' \in X top x
eq x
eq x'} (d_{\mathcal{X}}(x,x') - d_{Z}(P(x),P(x')))^2 \ &= rac{2}{n(n-1)} \sum_{i=1}^n \sum_{j=1}^{i-1} (d_{\mathcal{X}}(x_i,x_j) - ||z_i - z_j||)^2, \quad z_i := P(x_i) \end{aligned}$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 1. Principal Components Analysis

Distortions for Dimensionality Reduction (2/2)

Feature reconstruction methods aim to find latent representations P and reconstruction maps $r : \mathbb{R}^K \to \mathcal{X}$ from a given class of maps that **reconstruct features** as good as possible:

$$egin{aligned} D(P,r) &:= rac{1}{|X|} \sum_{x \in X} d_{\mathcal{X}}(x,r(P(x))) \ &= rac{1}{n} \sum_{i=1}^n d_{\mathcal{X}}(x_i,r(z_i)), \quad z_i := P(x_i). \end{aligned}$$





Singular Value Decomposition (SVD)

Theorem (Existence of SVD) For every $A \in \mathbb{R}^{n \times m}$ there exist matrices

$$U \in \mathbb{R}^{n \times k}, V \in \mathbb{R}^{m \times k}, \Sigma := diag(\sigma_1, \dots, \sigma_k) \in \mathbb{R}^{k \times k}, \qquad k := \min\{n, m\}$$

$$\sigma_1 \ge \sigma_2 \ge \dots \ge \sigma_r > \sigma_{r+1} = \dots = \sigma_k = 0, \qquad r := \operatorname{rank}(A)$$

$$U, V \text{ orthonormal, i.e., } U^T U = I, V^T V = I$$

with

$$A = U \Sigma V^T$$

 σ_i are called singular values of A.

Note: $I := \text{diag}(1, \ldots, 1) \in \mathbb{R}^{k \times k}$ denotes the unit matrix.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 1. Principal Components Analysis

Singular Value Decomposition (SVD; 2/2)

It holds:

a)
$$\sigma_i^2$$
 are eigenvalues and V_i eigenvectors of $A^T A$:
 $(A^T A)V_i = \sigma_i^2 V_i, \quad i = 1, ..., k, V = (V_1, ..., V_k)$
b) σ_i^2 are eigenvalues and U_i eigenvectors of AA^T :
 $(AA^T)U_i = \sigma_i^2 U_i, \quad i = 1, ..., k, U = (U_1, ..., U_k)$
proof:

a)
$$(A^T A)V_i = V\Sigma^T U^T U\Sigma V^T V_i = V\Sigma^2 e_i = \sigma_i^2 V_i$$

b) $(AA^T)U_i = U\Sigma^T V^T V\Sigma^T U^T U_i = U\Sigma^2 e_i = \sigma_i^2 U_i$





Truncated SVD



Let $A \in \mathbb{R}^{n \times m}$ and $U \Sigma V^T = A$ its SVD. Then for $k' \leq \min\{n, m\}$ the decomposition

$$A \approx U' \Sigma' V'^T$$

with

$$U' := (U_{,1}, \ldots, U_{,k'}), V' := (V_{,1}, \ldots, V_{,k'}), \Sigma' := \operatorname{diag}(\sigma_1, \ldots, \sigma_{k'})$$

is called **truncated SVD** with rank k'.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 1. Principal Components Analysis

Low Rank Approximation

Let $A \in \mathbb{R}^{n \times m}$. For $k \leq \min\{n, m\}$, any pair of matrices

 $U \in \mathbb{R}^{n \times k}, V \in \mathbb{R}^{m \times k}$

is called a **low rank approximation of** A with rank k. The matrix

 UV^T

is called the **reconstruction of** A by U, V and the quantity

$$||A - UV^{T}||_{F} = \sum_{i=1}^{n} \sum_{j=1}^{m} (A_{i,j} - U_{i}^{T}V_{j})^{2}$$

the L2 reconstruction error.

Note: $||A||_F$ is called Frobenius norm. (Do not confuse this with the L2 norm $|| \cdot ||_2$ for matrices.) Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany



Optimal Low Rank Approximation is Truncated SVD

· ·

Theorem (Low Rank Approximation; Eckart-Young theorem) Let $A \in \mathbb{R}^{n \times m}$. For $k' \leq \min\{n, m\}$, the optimal low rank approximation of rank k' (i.e., with smallest reconstruction error)

 $(U^*, V^*) := \operatorname*{arg\,min}_{U \in \mathbb{R}^{n \times k'}, V \in \mathbb{R}^{m \times k'}} ||A - UV^T||_F^2$

is the truncated SVD.

Note: As U, V do not have to be orthonormal, one can take $U := U'\Sigma', V := V'$ for the SVD $A = U'\Sigma'V'^T$. Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 1. Principal Components Analysis

Principal Components Analysis (PCA)

Let $X := \{x_1, \ldots, x_n\} \subseteq \mathbb{R}^m$ be a data set and $K \in \mathbb{N}$ the number of latent dimensions $(K \leq m)$.

PCA finds

- *K* principal components $v_1, \ldots, v_K \in \mathbb{R}^m$ and
- ▶ latent weights $z_i \in \mathbb{R}^K$ for each data point $i \in \{1, ..., n\}$,
- such that the linear combination of the principal components

$$x_i \approx \sum_{k=1}^K z_{i,k} v_k$$

reconstructs the original features x_i as good as possible:

$$\arg\min_{\substack{v_1,...,v_K\\z_1,...,z_n}} \sum_{i=1}^n ||x_i - \sum_{k=1}^K z_{i,k} v_k||^2$$

= $\sum_{i=1}^n ||x_i - Vz_i||^2$, $V := (v_1, ..., v_K)^T$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany



PCA Algorithm



10 /

- 1: procedure DIMRED-PCA($\mathcal{D} := \{x_1, \ldots, x_N\} \subseteq \mathbb{R}^M, K \in \mathbb{N}$) 2: $X := (x_1, x_2, \ldots, x_N)^T$ 3: $(U, \Sigma, V) := svd(X)$

- 4:
- $Z := U_{.,1:K} \cdot \Sigma_{1:K,1:K}$ return $\mathcal{D}^{\text{dimred}} := \{Z_{1,.}, \dots, Z_{N,.}\}$ 5:

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 1. Principal Components Analysis

Principal Components Analysis (Example 1)



Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Principal Components Analysis (Example 2)



12 /



[HTFF05, p. 538] Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 2. Probabilistic PCA & Factor Analysis

Outline

1. Principal Components Analysis

2. Probabilistic PCA & Factor Analysis

- 3. Non-linear Dimensionality Reduction
- 4. Supervised Dimensionality Reduction

Probabilistic Model



Probabilistic PCA provides a probabilistic interpretation of PCA.

It models for each data point

- ► a multivariate normal distributed latent factor z,
- that influences the observed variables linearly:

$$p(z) := \mathcal{N}(z; 0, I)$$

 $p(x \mid z; \mu, \sigma^2, W) := \mathcal{N}(x; \mu + Wz, \sigma^2 I)$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 2. Probabilistic PCA & Factor Analysis

Probabilistic PCA Loglikelihood

$$\ell(X, Z; \mu, \sigma^{2}, W) = \sum_{i=1}^{n} \ln p(x_{i} \mid z_{i}; \mu, \sigma^{2}, W) + \ln p(z_{i}) \\ = \sum_{i} \ln \mathcal{N}(x_{i}; \mu + Wz_{i}, \sigma^{2}I) + \ln \mathcal{N}(z_{i}; 0, I) \\ \propto \sum_{i} -\frac{1}{2} \log \sigma^{2} - \frac{1}{2\sigma^{2}} (x_{i} - \mu - Wz_{i})^{T} (x_{i} - \mu - Wz_{i}) - \frac{1}{2} z_{i}^{T} z_{i} \\ \propto -\sum_{i} \log \sigma^{2} + \frac{1}{\sigma^{2}} (\mu^{T} \mu + z_{i}^{T} W^{T} Wz_{i} - 2x_{i}^{T} \mu - 2x_{i}^{T} Wz_{i} + 2\mu^{T} Wz_{i}) \\ + z_{i}^{T} z_{i}$$

remember: $\mathcal{N}(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^m} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)\Sigma^{-1}(x-\mu)}$. Lars Schmidt-Thieme, Information $\sqrt[4]{(2\pi)^m} |\Sigma|^{\frac{1}{2}}$ Machine Learning Lab (ISMLL), University of Hildesheim, Germany



PCA vs Probabilistic PCA

$$\ell(X, Z; \mu, \sigma^2, W)$$

$$\propto \sum_i -\frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (x_i - \mu - Wz_i)^T (x_i - \mu - Wz_i) - \frac{1}{2} z_i^T z_i$$

as PCA: Decompose with minimal L2 loss

$$egin{aligned} x_i &pprox \mu + \sum_{k=1}^{\mathcal{K}} z_{i,k} oldsymbol{v}_k \ & ext{with } oldsymbol{v}_k &:= W_{\cdot,k} \end{aligned}$$

- ▶ different from PCA: L2 regularized row features z. ► cannot be solved by SVD. Use EM as learning algorithm!
- ► additionally also regularization of column features W possible (through a prior on W).

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 2. Probabilistic PCA & Factor Analysis

EM / Block Coordinate Descent: Outline

$$\ell(X, Z; \mu, \sigma^2, W)$$

$$\propto -\sum_{i} \log \sigma^2 + \frac{1}{\sigma^2} (\mu^T \mu + z_i^T W^T W z_i - 2x_i^T \mu - 2x_i^T W z_i + 2\mu^T W z_i) + z_i^T z_i$$

1. expectation step: $\forall i$

$\frac{\partial \ell}{\partial z_i} \stackrel{!}{=} 0$	$\rightsquigarrow z_i = \dots$	(0)
$0 \mathbf{z}_{\mathbf{i}}$		

2. minimization step:

$rac{\partial \ell}{\partial \mu} \stackrel{!}{=} 0$	$\rightsquigarrow \mu = \dots$	(1)
$rac{\partial \ell}{\partial \sigma^2} \stackrel{!}{=} 0$	$\rightsquigarrow \sigma^2 = \dots$	(2)
$\frac{\partial \ell}{\partial W} \stackrel{!}{=} 0$	$\rightsquigarrow {\it W}=\ldots$	(3)

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany





EM / Block Coordinate Descent



17

$$\ell(X, Z; \mu, \sigma^2, W)$$

$$\propto -\sum_{i} \log \sigma^2 + \frac{1}{\sigma^2} (\mu^T \mu + z_i^T W^T W z_i - 2x_i^T \mu - 2x_i^T W z_i + 2\mu^T W z_i) + z_i^T z_i$$

$$\frac{\partial \ell}{\partial z_i} = -\frac{1}{\sigma^2} (2z_i^T W^T W - 2x_i^T W + 2\mu^T W) - 2z_i^T \stackrel{!}{=} 0$$
$$(W^T W + \sigma^2 I) z_i = W^T (x_i - \mu)$$
$$z_i = (W^T W + \sigma^2 I)^{-1} W^T (x_i - \mu)$$
(0)

$$\frac{\partial \ell}{\partial \mu} = -\frac{1}{\sigma^2} \sum_{i} 2\mu^T - 2x_i^T + 2z_i^T W^T \stackrel{!}{=} 0$$

Note: As $\mathbb{E}(z_i) = \underbrace{\mu}_{-} \underbrace{\mu}_{-$

Machine Learning 2. Probabilistic PCA & Factor Analysis

EM / Block Coordinate Descent: Summary

alternate until convergence:

- 1. expectation step: $\forall i$ $z_i = (W^T W + \sigma^2 I)^{-1} W^T (x_i - \mu) \qquad (0)$
- 2. minimization step:

$$\mu = \frac{1}{n} \sum_{i} x_i - W z_i \tag{1}$$

$$\sigma^{2} = \frac{1}{n} \sum_{i} (x_{i} - \mu - Wz_{i})^{T} (x_{i} - \mu - Wz_{i})$$
(2)

$$W = \sum_{i} (x_{i} - \mu) z_{i}^{T} (\sum_{i} z_{i} z_{i}^{T})^{-1}$$
(3)

Probabilistic PCA Algorithm (EM)



1: procedure DIMRED-PPCA($\mathcal{D} := \{x_1, \dots, x_N\} \subseteq \mathbb{R}^M, K \in \mathbb{N}, \epsilon \in \mathbb{R}^+$) 2: allocate $z_1, \dots, z_N := 0 \in \mathbb{R}^K, \mu := 0 \in \mathbb{R}^M, W := 0 \in \mathbb{R}^{N \times K}, \sigma^2 := 1 \in \mathbb{R}$ 3: repeat 4: $\sigma_{old}^2 := \sigma^2, z^{old} := z$ 5: for $n := 1, \dots, N$ do 6: $z_n := (W^T W + \sigma^2 I)^{-1} W^T (x_n - \mu)$ 7: $\mu_{old} := \mu$ 8: $\mu := \frac{1}{n} \sum_i x_i - W z_i$ 9: $\sigma^2 := \frac{1}{n} \sum_i (x_i - \mu_{old} - W z_i)^T (x_i - \mu_{old} - W z_i)$ 10: $W := \sum_i (x_i - \mu_{old}) z_i^T (\sum_i z_i z_i^T)^{-1}$ 11: until $\frac{1}{N} \sum_{n=1}^N ||z_n - z_n^{old}|| < \epsilon$ 12: return $\mathcal{D}^{dimred} := \{z_1, \dots, z_N\}$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 2. Probabilistic PCA & Factor Analysis

EM / Block Coordinate Descent: Example





Regularization of Column Features W



Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 2. Probabilistic PCA & Factor Analysis

Bayesian PCA: Example







Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany



Factor Analysis



$$p(z) := \mathcal{N}(z; 0, I)$$

 $p(x \mid z; \mu, \Sigma, W) := \mathcal{N}(x; \mu + Wz, \Sigma), \quad \Sigma ext{ diagonal}$

$$\ell(X, Z; \mu, \Sigma, W)$$

 $\propto \sum_{i} -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu - Wz_i)^T \Sigma^{-1} (x_i - \mu - Wz_i) - \frac{1}{2} z_i^T z_i$

EM:

$$z_i = (W^T \Sigma^{-1} W + I)^{-1} W^T \Sigma^{-1} (x_i - \mu)$$
 (0')

$$\mu = \frac{1}{n} \sum_{i} x_i - W z_i \tag{1}$$

Note: See appendix for derivation of EM formulas:

$$\sum_{i} j = \frac{1}{n} \sum_{j} ((x_i - \mu_j - W_{Z_j})_j)^2 \qquad (2')$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 3. Non-linear Dimensionality Reduction

Outline

1. Principal Components Analysis

2. Probabilistic PCA & Factor Analysis

3. Non-linear Dimensionality Reduction

4. Supervised Dimensionality Reduction

23 /

Linear Dimensionality Reduction



Dimensionality reduction accomplishes two tasks:

compute lower dimensional representations for given data points x_i
 ▶ for PCA:

$$u_i = \Sigma^{-1} V^T x_i, \quad U := (u_1, \ldots, u_n)^T$$

 compute lower dimensional representations for new data points x (often called "fold in")

► for PCA:

$$u := \arg\min_{u} ||x - V\Sigma u||^2 = \Sigma^{-1} V^T x$$

PCA is called a **linear dimensionality reduction technique** because the latent representations *u* depend linearly on the observed representations *x*.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 3. Non-linear Dimensionality Reduction

Kernel Trick

Represent (conceptionally) non-linearity by linearity in a higher dimensional embedding

$$\phi: \mathbb{R}^m \to \mathbb{R}^{\tilde{m}}$$

but compute in lower dimensionality for methods that depend on x only through a scalar product

$$ilde{x}^T ilde{ heta} = \phi(x)^T \phi(heta) = k(x, heta), \quad x, heta \in \mathbb{R}^m$$

if k can be computed without explicitly computing ϕ .



Kernel Trick / Example Example:



$$\phi : \mathbb{R} \to \mathbb{R}^{1001}, \\ x \mapsto \left(\left(\begin{array}{c} 1000 \\ i \end{array} \right)^{\frac{1}{2}} x^{i} \right)_{i=0,\dots,1000} = \begin{pmatrix} 1 \\ 31.62 x \\ 706.75 x^{2} \\ \vdots \\ 31.62 x^{999} \\ x^{1000} \end{pmatrix}$$

$$\tilde{x}^T \tilde{\theta} = \phi(x)^T \phi(\theta) = \sum_{i=0}^{1000} \left(\begin{array}{c} 1000 \\ i \end{array} \right) x^i \theta^i = (1+x\theta)^{1000} =: k(x,\theta)$$

Naive computation:

► 2002 binomial coefficients, 3003 multiplications, 1000 additions.

Kernel computation:

▶ 1 multiplication, 1 addition, 1 exponentiation.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 3. Non-linear Dimensionality Reduction

Kernel PCA

$$\phi : \mathbb{R}^{m} \to \mathbb{R}^{\tilde{m}}, \quad \tilde{m} \gg m$$
$$\tilde{X} := \begin{pmatrix} \phi(x_{1}) \\ \phi(x_{2}) \\ \vdots \\ \phi(x_{n}) \end{pmatrix}$$
$$\tilde{X} \approx U \Sigma \tilde{V}^{T}$$

We can compute the columns of U as eigenvectors of $\tilde{X}\tilde{X}^T \in \mathbb{R}^{n \times n}$ without having to compute $\tilde{V} \in \mathbb{R}^{\tilde{m} \times k}$ (which is large!):

$$\tilde{X}\tilde{X}^{\mathsf{T}}U_i=\sigma_i^2U_i$$



Kernel PCA / Removing the Mean

Juiversit

Issue 1: The $\tilde{x}_i := \phi(x_i)$ may not have zero mean and thus distort PCA.

$$\begin{split} \tilde{x}'_{i} &:= \tilde{x}_{i} - \frac{1}{n} \sum_{i=1}^{n} \tilde{x}_{i} \\ &= \tilde{X}^{T} \left(I - \frac{1}{n} \mathbb{1} \right) \\ \tilde{X}' &:= \left(\tilde{x}'_{1}, \dots, \tilde{x}'_{n} \right)^{T} = \left(I - \frac{1}{n} \mathbb{1} \right) \tilde{X}^{T} \\ \mathcal{K}' &:= \tilde{X}' \tilde{X}'^{T} = \left(I - \frac{1}{n} \mathbb{1} \right) \tilde{X}^{T} \tilde{X} \left(I - \frac{1}{n} \mathbb{1} \right) \\ &= H \mathcal{K} \mathcal{H}, \quad \mathcal{H} := \left(I - \frac{1}{n} \mathbb{1} \right) \text{ centering matrix} \end{split}$$

Thus, the kernel matrix K' with means removed can be computed from the kernel matrix K without having to access coordinates. Note: $1 := (1)_{i=1,...,n,j=1,...,n}$ vector of ones, $I := (\delta(i = j))_{i=1},...,n, j=1,...,n}$ unity matrix. Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 3. Non-linear Dimensionality Reduction

Kernel PCA / Fold In

Issue 2: How to compute projections u of new points x (as \tilde{V} is not computed)?

$$u := \arg\min_{u} ||x - \tilde{V}\Sigma u||^2 = \Sigma^{-1}\tilde{V}^T x$$

With

$$\begin{split} \tilde{V} &= \tilde{X}^{T} U \Sigma^{-1} \\ u &= \Sigma^{-1} \tilde{V}^{T} x = \Sigma^{-1} \Sigma^{-1} U^{T} \tilde{X} x = \Sigma^{-2} U^{T} (k(x_{i}, x))_{i=1,...,n} \end{split}$$

u can be computed with access to kernel values only (and to U, Σ).



Kernel PCA / Summary

Given:

- data set $X := \{x_1, \ldots, x_n\} \subseteq \mathbb{R}^m$,
- kernel function $k : \mathbb{R}^m \times \mathbb{R}^m \to R$.

task 1: Learn latent representations U of data set X:

$$K := (k(x_i, x_j))_{i=1,...,n, j=1,...,n}$$
(0)

$$K' := HKH, \quad H := \left(I - \frac{1}{n}\mathbb{1}\right) \tag{1}$$

$$(U, \Sigma)$$
 :=eigen decomposition $K'U = U\Sigma$ (2)

task 2: Learn latent representation u of new point x:

$$u := \Sigma^{-2} U^{T}(k(x_{i}, x))_{i=1,...,n}$$
(3)

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Eigenvalue=20.936

Machine Learning 3. Non-linear Dimensionality Reduction

Eigenvalue=22.558

Kernel PCA: Example 1



Eigenvalue=4.648



Eigenvalue=3.988

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany





Kernel PCA: Example 2





Machine Learning 4. Supervised Dimensionality Reduction

Outline

- 1. Principal Components Analysis
- 2. Probabilistic PCA & Factor Analysis
- 3. Non-linear Dimensionality Reduction

4. Supervised Dimensionality Reduction

38

Dimensionality Reduction as Pre-Processing

Given a prediction task and

a data set $\mathcal{D}^{\mathsf{train}} := \{(x_1, y_1), \dots, (x_n, y_n)\} \subseteq \mathbb{R}^m \times \mathcal{Y}.$

- 1. compute latent features $z_i \in \mathbb{R}^K$ for the objects of a data set by means of dimensionality reduction of the predictors x_i .
 - e.g., using PCA on $\{x_1, \ldots, x_n\} \subseteq \mathbb{R}^m$
- 2. learn a prediction model

$$\hat{y}: \mathbb{R}^{K} \to \mathcal{Y}$$

on the latent features based on

$$\mathcal{D}'^{\mathsf{train}} := \{(z_1, y_1), \dots, (z_n, y_n)\}$$

3. treat the number K of latent dimensions as hyperparameter.
▶ e.g., find using grid search.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 4. Supervised Dimensionality Reduction

Dimensionality Reduction as Pre-Processing

Advantages:

- ► simple procedure
- ► generic procedure
 - works with any dimensionality reduction method and any prediction method as component methods.
- usually fast

Disadvantages:

- dimensionality reduction is unsupervised, i.e., not informed about the target that should be predicted later on.
 - leads to the very same latent features regardless of the prediction task.
 - likely not the best task-specific features are extracted.





Supervised PCA



$$p(z) := \mathcal{N}(z; 0, 1)$$
 $p(x \mid z; \mu_x, \sigma_x^2, W_x) := \mathcal{N}(x; \mu_x + W_x z, \sigma_x^2 I)$
 $p(y \mid z; \mu_y, \sigma_y^2, W_y) := \mathcal{N}(y; \mu_y + W_y z, \sigma_y^2 I)$

- ► like two PCAs, coupled by shared latent features *z*:
 - one for the predictors *x*.
 - one for the targets *y*.
- Iatent features act as information bottleneck.
- also known as Latent Factor Regression or Bayesian Factor Regression.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 4. Supervised Dimensionality Reduction

Supervised PCA: Discriminative Likelihood

A simple likelihood would put the same weight on

- reconstructing the predictors and
- reconstructing the targets.

A weight $\alpha \in \mathbb{R}_0^+$ for the reconstruction error of the predictors should be introduced (**discriminative likelihood**):

$$L_{\alpha}(\Theta; x, y, z) := \prod_{i=1}^{n} p(y_i \mid z_i; \Theta) p(x_i \mid z_i; \Theta)^{\alpha} p(z_i; \Theta)$$

 α can be treated as hyperparameter and found by grid search.



Supervised PCA: EM



- The M-steps for μ_x, σ_x^2, W_x and μ_y, σ_y^2, W_y are exactly as before.
- ► the coupled E-step is:

$$z_i = \left(\frac{1}{\sigma_y^2} W_y^T W_y + \alpha \frac{1}{\sigma_x^2} W_x^T W_x\right)^{-1} \left(\frac{1}{\sigma_y^2} W_y^T (y_i - \mu_y) + \alpha \frac{1}{\sigma_x^2} W_x^T (x_i - \mu_x)\right)$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning 4. Supervised Dimensionality Reduction

Conclusion (1234/4)

- Dimensionality reduction aims to find a lower dimensional representation of data that preserves the information as much as possible. — "Preserving information" means
 - to preserve pairwise distances between objects (multidimensional scaling).
 - to be able to reconstruct the original object features (feature reconstruction).
- The truncated Singular Value Decomposition (SVD) provides the best low rank factorization of a matrix in two factor matrices.
 - SVD is usually computed by an algebraic factorization method (such as QR decomposition).
- Principal components analysis (PCA) finds latent object features and latent variable features that provide the best linear reconstruction (in L2 error).
 - PCA is a truncated SVD of the data matrix.



Readings



- Principal Components Analysis (PCA)
 - ▶ [HTFF05], ch. 14.5.1, [Bis06], ch. 12.1, [Mur12], ch. 12.2.
- Probabilistic PCA
 - ▶ [Bis06], ch. 12.2, [Mur12], ch. 12.2.4.
- ► Factor Analysis
 - ▶ [HTFF05], ch. 14.7.1, [Bis06], ch. 12.2.4.
- ► Kernel PCA
 - ▶ [HTFF05], ch. 14.5.4, [Bis06], ch. 12.3, [Mur12], ch. 14.4.4.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Further Readings

- (Non-negative) Matrix Factorization
 - ► [HTFF05], ch. 14.6
- Independent Component Analysis, Exploratory Projection Pursuit
 - ▶ [HTFF05], ch. 14.7 [Bis06], ch. 12.4 [Mur12], ch. 12.6.
- Nonlinear Dimensionality Reduction
 - ▶ [HTFF05], ch. 14.9, [Bis06], ch. 12.4



Factor Analysis: Loglikelihood



$$\ell(X, Z; \mu, \Sigma, W) = \sum_{i=1}^{n} \ln p(x \mid z; \mu, \Sigma, W) + \ln p(z) = \sum_{i} \ln \mathcal{N}(x; \mu + Wz, \Sigma) + \ln \mathcal{N}(z; 0, I) \propto \sum_{i} -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_{i} - \mu - Wz_{i})^{T} \Sigma^{-1} (x_{i} - \mu - Wz_{i}) - \frac{1}{2} z_{i}^{T} z_{i} \propto -\sum_{i} \log |\Sigma| + (x_{i}^{T} \Sigma^{-1} x_{i} + \mu^{T} \Sigma^{-1} \mu + z_{i}^{T} W^{T} \Sigma^{-1} Wz_{i} - 2x_{i}^{T} \Sigma^{-1} \mu - 2x_{i}^{T} \Sigma^{-1} Wz_{i} + 2\mu^{T} \Sigma^{-1} Wz_{i}) + z_{i}^{T} z_{i}$$

remember: $\mathcal{N}(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^m} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)\Sigma^{-1}(x-\mu)}$. Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Factor Analysis: EM / Block Coordinate Descent



$$\ell(X, Z; \mu, \Sigma, W)$$

$$\propto -\sum_{i} \log |\Sigma| + (x_i^T \Sigma^{-1} x_i + \mu^T \Sigma^{-1} \mu + z_i^T W^T \Sigma^{-1} W z_i - 2x_i^T \Sigma^{-1} \mu$$

$$-2x_i^T \Sigma^{-1} W z_i + 2\mu^T \Sigma^{-1} W z_i) + z_i^T z_i$$

$$\begin{aligned} \frac{\partial \ell}{\partial z_i} &= -(2z_i^T W^T \Sigma^{-1} W - 2x_i^T W \Sigma^{-1} + 2\mu^T \Sigma^{-1} W) - 2z_i^T = \\ (W^T \Sigma^{-1} W + I) z_i &= W^T \Sigma^{-1} (x_i - \mu) \\ z_i &= (W^T \Sigma^{-1} W + I)^{-1} W^T \Sigma^{-1} (x_i - \mu) \end{aligned}$$

$$\frac{\partial \ell}{\partial \mu} = -\sum_{i} 2\mu^{T} \Sigma^{-1} - 2x_{i}^{T} \Sigma^{-1} + 2z_{i}^{T} W^{T} \Sigma^{-1} \stackrel{!}{=} 0$$

Note: As $\mathbb{E}(z_i) = \underline{\mathbb{D}}$, whethere is fixed to $\mu := \frac{1}{n} \sum_i x_i$. Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany (1')

References





Christopher M. Bishop. *Pattern recognition and machine learning*, volume 1. springer New York, 2006.

Trevor Hastie, Robert Tibshirani, Jerome Friedman, and James Franklin. *The elements of statistical learning: data mining, inference and prediction*, volume 27. Springer, 2005.



Kevin P. Murphy.

Machine learning: a probabilistic perspective. The MIT Press, 2012.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Matrix Trace

The function

$$\mathsf{tr}:\bigcup_{n\in\mathbb{N}}\mathbb{R}^{n\times n}\to\mathbb{R}$$

$$A\mapsto \operatorname{tr}(A):=\sum_{i=1}^n a_{i,i}$$

is called matrix trace. It holds:

a) invariance under permutations of factors:

$$\mathsf{tr}(AB) = \mathsf{tr}(BA)$$

b) invariance under basis change:

$$\operatorname{tr}(B^{-1}AB) = \operatorname{tr}(A)$$

proof:

a)
$$tr(AB) = \sum_{i} \sum_{j} A_{i,j} B_{j,i} = \sum_{i} \sum_{j} B_{i,j} A_{j,i} = tr(BA)$$

b) $tr(B^{-1}AB) = tr(BB^{-1}A) = tr(A)$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

43 / 38

Machine Learning

Frobenius Norm



The function
$$||\cdot||_F : \bigcup_{n,m\in\mathbb{N}} \mathbb{R}^{n\times m} \to \mathbb{R}_0^+$$

 $A \mapsto ||A||_F := (\sum_{i=1}^n \sum_{j=1}^m a_{i,j}^2)^{\frac{1}{2}}$

is called Frobenius norm. It holds:

a) trace representation:

$$||A||_F = (\operatorname{tr}(A^T A))^{\frac{1}{2}}$$

b) invariance under orthonormal transformations:

$$tr(UAV^T) = tr(A), \quad U, V \text{ orthonormal}$$

proof:

a)
$$tr(A^{T}A) = \sum_{i} \sum_{j} A_{j,i}A_{j,i} = ||A||_{2}^{2}$$

b) $||UAV||_{F}^{2} = tr(VA^{T}U^{T}UAV^{T}) = tr(VA^{T}AV^{T})$
 $= tr(A^{T}AV^{T}V) = tr(A^{T}A) = ||A||_{F}^{2}$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Machine Learning

Frobenius Norm (2/2)

c) representation as sum of squared singular values:

$$||A||_F = \sum_{i=1}^{\min\{m,n\}} \sigma_i^2$$

proof:

c) let
$$A = U\Sigma V^T$$
 the SVD of A

$$||A||_F = ||U\Sigma V^T||_F = ||\Sigma||_F = \operatorname{tr}(\Sigma^T \Sigma) = \sum_{i=1}^{\min\{m,n\}} \sigma_i^2$$

