

Statistical Data Analysis

Dr. Jana de Wiljes

18. Januar 2022

Universität Potsdam

Principal Component Analysis

Principal Component Analysis (PCA)

Snapshot information:

1. Dimensionality reduction, i.e., represent it in a more tractable, lower-dimensional form, without losing too much information

- Data Compression (Save computation/memory) ~~Noise Reduction~~ \mathbb{R}^{30} \mathbb{R}^7
- Noise Reduction/ avoid overfitting to noise
- Data Visualization (e.g., in two dimensions) $\mathbb{R}^{10} \rightsquigarrow \mathbb{R}^2, \mathbb{R}^3$

2. unsupervised learning algorithm (Pearson 1901)

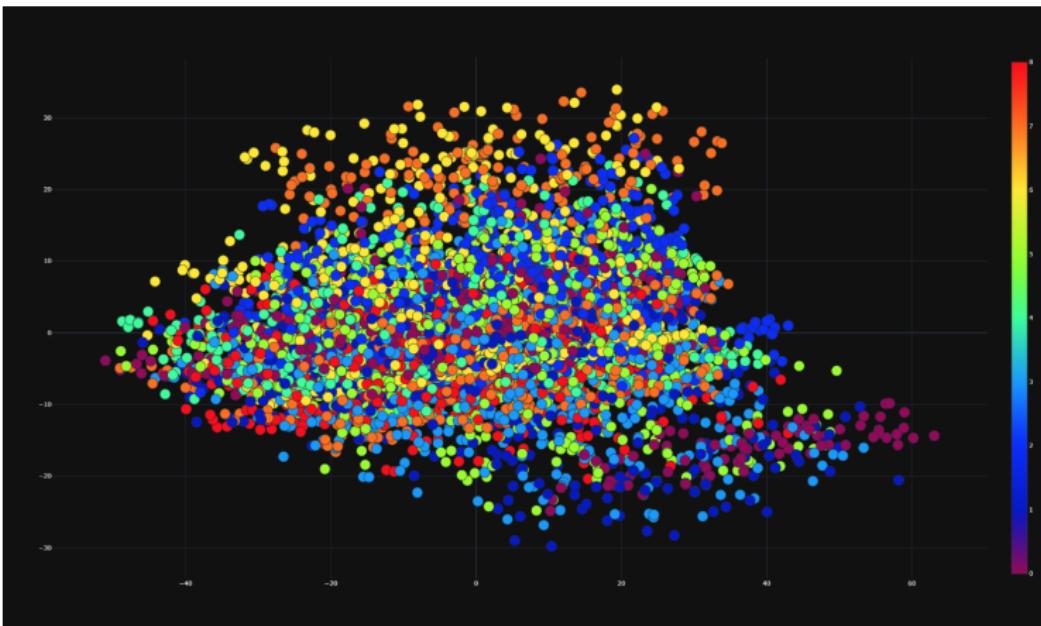
3. Idea: uses an orthogonal transformation to convert a set of observations of correlated variables into a set of linearly uncorrelated variables (called principal components)

4. Known under many different names:

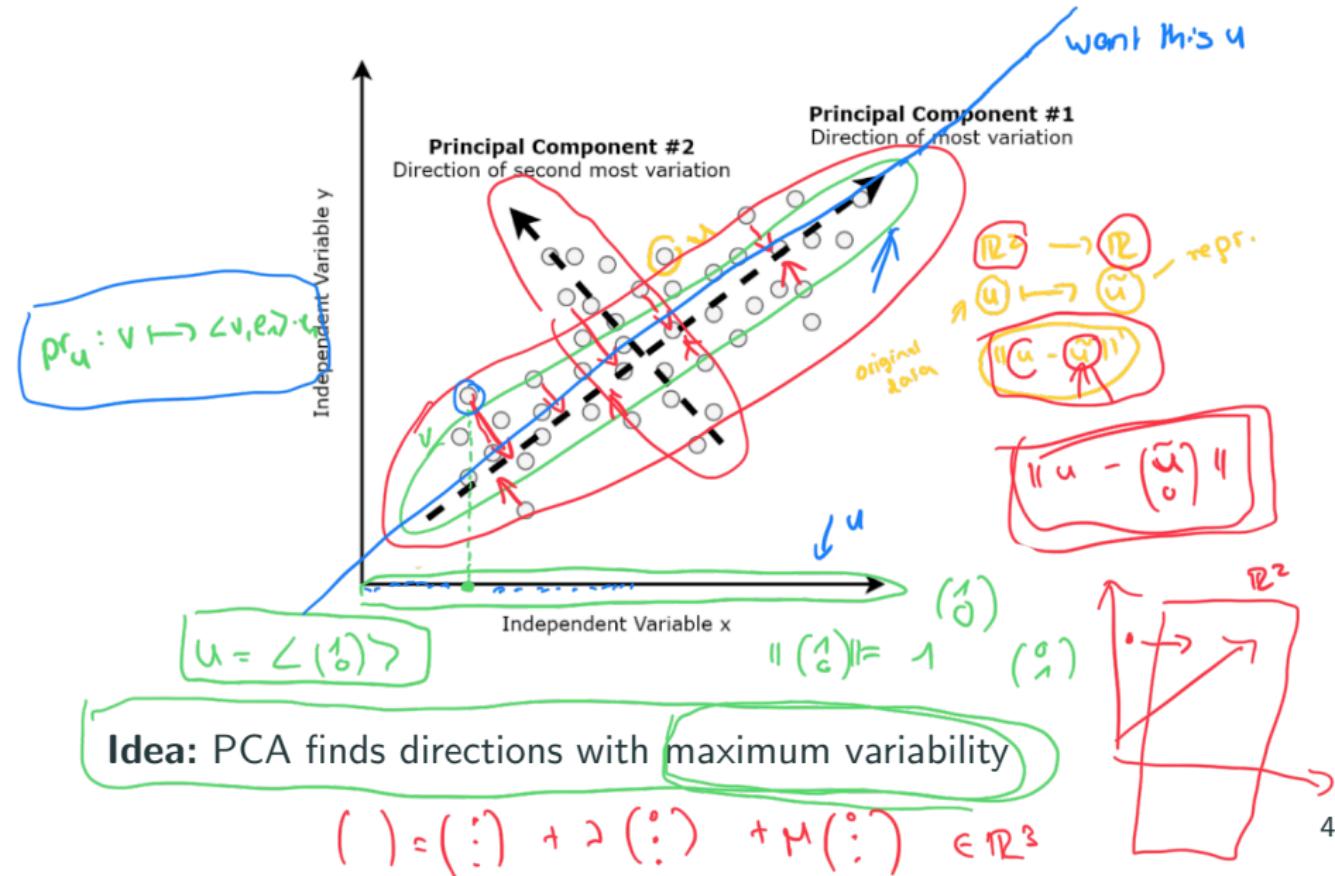
- Karhunen-Loeve transformation
- Hotelling transformation
- empirical orthogonal functions



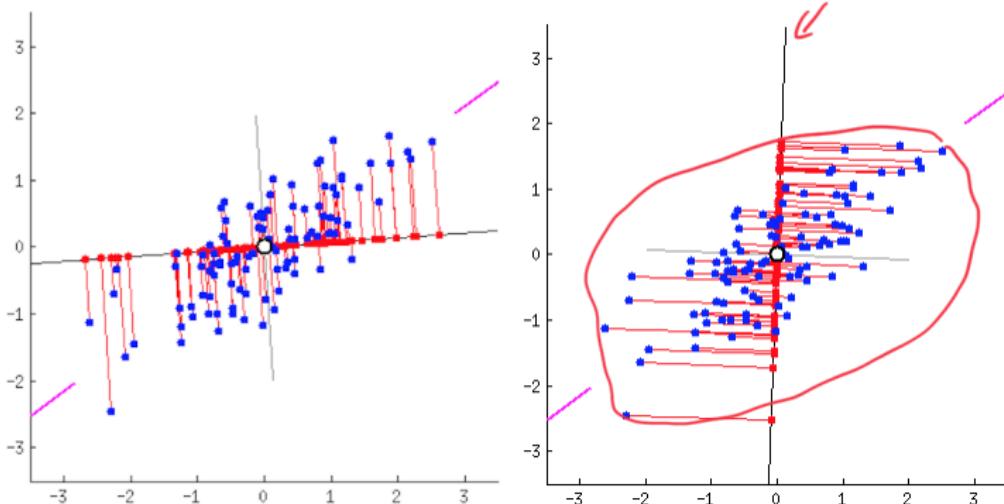
Visualisation



Motivation behind PCA



Motivation behind PCA



PCA:

- Identify a Hyperplane that lies closest to the data
- Project the data onto the hyperplane.

Background information

Orthogonal basis

Def: Let V be a vector space with scalar product $\langle \cdot, \cdot \rangle$ and $\{v_j\}_{j \in J}$ a family of vectors.

$$\begin{aligned} \langle x + y, z \rangle &= \langle x, z \rangle + \langle y, z \rangle \\ \langle x + y, x \rangle &= \langle x, x \rangle + \langle y, x \rangle \\ \text{if } \langle v_j, v_k \rangle = 0 \forall j \neq k \in J \quad \text{then} \quad \langle v_j, v_j \rangle &= 1 \end{aligned}$$

- $\{v_j\}_{j \in J}$ are a **orthogonal system**, if $\langle v_j, v_k \rangle = 0 \forall j \neq k \in J$ and $v_j \neq 0 \forall j \in J$.
- $\{v_j\}_{j \in J}$ is an **orthonormal system**, if additionally:
 $\langle v_j, v_i \rangle = 1 \forall j \quad (\Leftrightarrow \|v_j\| = 1)$, in other words:
 $\langle v_j, v_k \rangle = \delta_{jk} \forall j, k \in J$.
- An Orthogonal- respectively. -normalsystem is called **orthogonal basis** bzw. **orthonormal basis**, if the die vectors of the systems form a basis.

$$\left(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \right)$$

$$\langle x + y, x \rangle = \langle x, x \rangle + \langle y, x \rangle$$

Representation with respect to an orthonormal basis

Theorem

Let $\{v_j\}_{j=1,\dots,n}$ be an orthonormal basis of the vector space V and $w \in V$ a second vector. Then the following holds:

$$w = \langle v_1, w \rangle v_1 + \cdots + \langle v_n, w \rangle v_n.$$
$$w = \langle v_1, w \rangle v_1 + \cdots + \langle v_n, w \rangle v_n.$$

Proof

- Compute the difference $u = w - \sum_{j=1}^n \langle v_j, w \rangle v_j$
- apply the scalar product $\langle v_k, \cdot \rangle$ to it

$$\begin{aligned}\langle v_k, u \rangle &= \left\langle v_k, w - \sum_{j=1}^n \langle v_j, w \rangle v_j \right\rangle \\ &= \langle v_k, w \rangle - \sum_{j=1}^n \underbrace{\langle v_k, \underbrace{\langle v_j, w \rangle v_j} \rangle}_{\text{scalar}}\end{aligned}$$

$$\begin{aligned}&= \langle v_k, w \rangle - \left(\sum_{j=1}^n \langle v_j, w \rangle \underbrace{\langle v_k, v_j \rangle}_{\text{scalar}} \right) \\ &= \langle v_k, w \rangle - \langle v_k, w \rangle \underbrace{\sum_{j=1}^n \langle v_k, v_j \rangle}_{\text{scalar}}\end{aligned}$$

$$= \langle v_k, w \rangle - \langle v_k, w \rangle = 0 \quad \begin{cases} 1 & \text{if } k=j \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow \langle u, v_k \rangle = 0 \quad \forall k \quad \Rightarrow \left\langle u, \sum_{i=1}^n \lambda_i v_i \right\rangle = 0$$

$$u = \left\langle \sum_{i=1}^n \lambda_i v_i \right\rangle$$

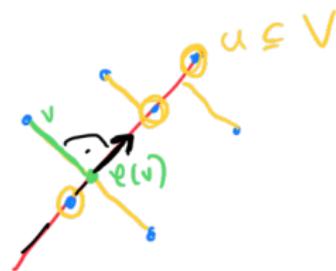
↓ can find Basis

$$\Rightarrow \langle u, u \rangle = 0 \Rightarrow \|u\| = 0 \Rightarrow u = 0 \quad \square$$

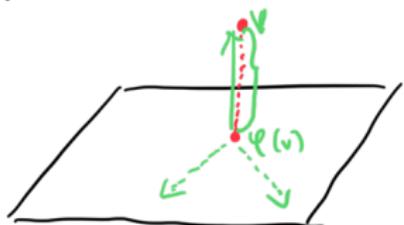
Orthogonal Projection

Def: Let $U \subseteq V$ be a subspace. A map $\varphi: V \rightarrow U$ is called **projection of V onto U** , if for every $u \in U$ gilt: $\varphi(u) = u$. A projection is called **orthogonal projection onto subspace U** , if for every vector $v \in V$ holds:

$$(\varphi(v) - v) \perp U.$$



Example: Hyperspace



Orthonormalbasis

Def: Let V be a K -vector space with scalar product and $U \subseteq V$ a finite-dimensional subspace. Furthermore let $\{u_1, \dots, u_k\}$ be an orthonormal basis of U . Then the map

$$\text{pr}_U: V \rightarrow U, \quad v \mapsto \sum_{j=1}^k \langle u_j, v \rangle u_j$$

$\{u_1, \dots, u_k\}$
linearly independent
then one can always
extend to a
basis of V

is an orthogonal projection.

$$\left(\sum_{j=1}^n \langle u_j, v \rangle u_j \right) \in U?$$

(Scalar) $\in U$

$$u, v \in U \Rightarrow u+v \in U$$
$$u \in U, \lambda \in \mathbb{R} \quad \lambda \cdot u \in U$$

$$\left\langle \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \right\rangle = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = e_3$$

$$= U$$

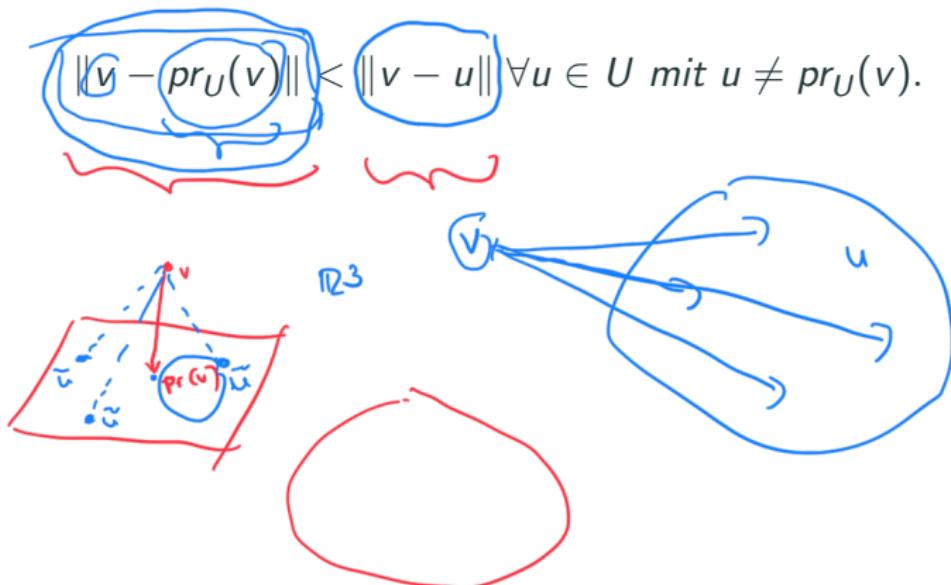
$$v = \sum_{i=1}^3 \langle e_i, v \rangle e_i$$

$$v \approx \sum_{i=1}^3 \langle e_i, v \rangle \cdot e_i$$

Approximation theorem

Theorem

Let V be a \mathbb{R} -vector space, with a scalar product and the corresponding norm $\|\cdot\|$. Let U be a subspace of V . Then for every $v \in V$ $pr_U(v)$ is the best approximation of v in U , i.e.,



PCA

Given: data set of p dimensional vectors

(X) $n \times p$ dimensional

$x_i \in \mathbb{R}^p$ want subspace
 $U \cong \mathbb{R}^q$

Goal: Want to project them to q -dimensional subspace ($q \ll p$)

$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$

Principal components: q uncorrelated, orthogonal directions formed by projecting the original data

$$v_1 = \langle e_1 \rangle$$

$$u = \langle w \rangle$$

$$\|w\| = 1$$



Derivation

Special case : Data $x_i \in \mathbb{R}^P$ and there are n data points

Assume : Data is normalised , i.e. $\frac{1}{n} \sum_{i=1}^n x_i = 0$

Want to project to a line (i.e. $g=1$)

$$\begin{aligned} \|x_i - \underbrace{\langle w, x_i \rangle w}_{\text{Data point } p_w(x_i)}\|^2 &= \underbrace{\langle x_i - \langle w, x_i \rangle w, x_i - \langle w, x_i \rangle w \rangle}_{\text{scalar}} \\ &= \underbrace{\langle x_i, x_i \rangle}_{\text{scalar}} - \underbrace{\langle x_i, \langle w, x_i \rangle w \rangle}_{\text{scalar}} - \underbrace{\langle \langle w, x_i \rangle w, x_i \rangle}_{\text{scalar}} \\ &= \|x_i\|^2 - 2 \underbrace{\langle w, x_i \rangle^2}_{\text{scalar}} + \underbrace{\langle w, x_i \rangle^2}_{\text{scalar}} \cdot \underbrace{\langle w, w \rangle}_{\text{scalar}} \end{aligned}$$

Orthogonal basis vector

$$= \|x_i\|^2 - 2 \underbrace{\langle w, x_i \rangle^2}_{\text{scalar}} + \underbrace{\langle w, x_i \rangle^2}_{\text{scalar}} \cdot \underbrace{\langle w, w \rangle}_{\text{scalar}} \quad (1)$$

$$= \underbrace{\|x_i\|^2}_{\text{MSE } (w)} - \underbrace{\langle w, x_i \rangle^2}_{\text{scalar}}$$

$$\text{MSE } (w) = \frac{1}{n} \sum_{i=1}^n (\|x_i\|^2 - \underbrace{\langle w, x_i \rangle^2}_{\text{scalar}})$$

↑ want a w that maximizes this term here

Derivation

want to maximise $\frac{1}{n} \sum_{i=1}^n \langle w, x_i \rangle^2$

$$\frac{1}{n} \sum_{i=1}^n \langle w, x_i \rangle^2 \Leftrightarrow$$

$$\left(\frac{1}{n} \sum_{i=1}^n \langle x_i, w \rangle \right)^2 + \text{Var}[\langle w, x_i \rangle]$$

$$\text{Var}(x) = E[x^2] - E[x]^2 = 0$$

is zero

$$\frac{1}{n} \sum_{i=1}^n \langle x_i, w \rangle \cdot w$$

$$= \left\langle \frac{1}{n} \sum_{i=1}^n x_i, w \right\rangle \cdot w = 0$$

= 0

Mean of projected data points
still zero

$\therefore x_i = \dots$

$$x_i \sim X_i$$

x_i have the
same distribution

Compute empirical estimate of
 $\text{Var}(x_i)$

Derivation

Maximizes: $\text{Var} [\langle w, x_i \rangle]$ $\forall i$

$\Rightarrow \frac{1}{n} \sum_{i=1}^n \langle w, x_i \rangle^2$

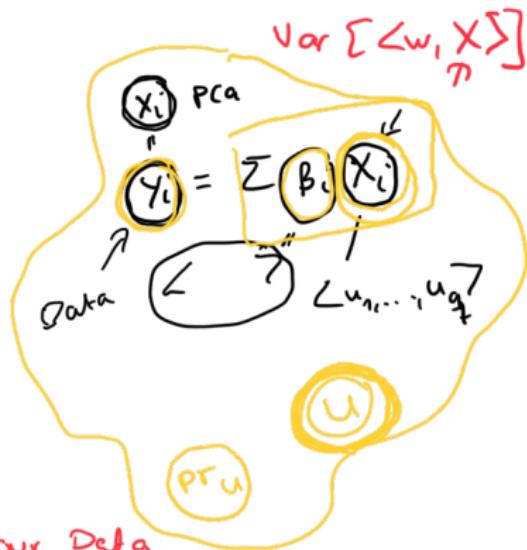
$\times n \times p$

$x_i \sim X_i$ but all x_i have distribution

$$\begin{aligned}\text{Var}(w, X) &= \frac{1}{n} \sum \langle x_i, w \rangle^2 \\ &= \frac{1}{n} (X^T w)^T (X^T w) \\ &= \frac{1}{n} w^T X^T X w\end{aligned}$$

$$= \boxed{\frac{w^T \cancel{X^T X} w}{n}}$$

emp. Covariance of our Data



$\boxed{\|w\|=1}$

$$\hat{\sigma}_{\text{var}}^2 = w^T \cancel{\frac{X^T X}{n}} w - \hat{\sigma}^2 (w^T w - 1)$$

$w^T w = 1$; $\boxed{\frac{X^T}{n} w = \hat{\sigma}^2 w}$ larger value

Derivation

Derivation
