

# Principal components analysis

Victor Kitov

[v.v.kitov@yandex.ru](mailto:v.v.kitov@yandex.ru)

# Table of Contents

- 1 Linear algebra reminder
- 2 Dimensionality reduction intro
- 3 Principal component analysis

## Scalar product reminder

- Here we will assume  $\langle a, b \rangle = a^T b$
- $\|a\| = \sqrt{\langle a, a \rangle}$
- Signed projection of  $x$  on  $a$  is equal to  $\langle x, a \rangle / \|a\|$
- Unsigned projection (length) of  $x$  onto  $a$  is equal to  $|\langle x, a \rangle| / \|a\|$

# Eigenvectors, eigenvalues

- If for some  $A \in \mathbb{R}^{D \times D}$  there exist scalar  $\lambda$  and  $D$ -dimensional vector  $v$  such that  $Av = \lambda v$  then
  - $v$  is called eigenvector of  $A$
  - $\lambda$  is called eigenvalue of  $A$ , corresponding to eigenvector  $v$ .
- $\exists v \neq 0 : Av = \lambda v \Leftrightarrow (A - \lambda I)v = 0 \Leftrightarrow \det(A - \lambda I) = 0$ . So all eigenvalues satisfy  $\det(A - \lambda I) = 0$  which
  - is a polynomial equation of order  $D$
  - so has  $D$  solutions<sup>1</sup> (accounting for their multiplicity, possibly complex)

---

<sup>1</sup>According to Fundamental theorem of algebra.

# Symmetric matrices

- Matrix  $A \in \mathbb{R}^{D \times D}$  is called *symmetric* if  $A^T = A$ .
- Properties:
  - All eigenvalues of symmetric matrix are real.
  - Eigenvectors, corresponding to different eigenvalues of symmetric matrix  $B$  are orthogonal to each other.
  - If  $\tilde{\lambda}$  is a repeated root of  $\det(A - \lambda I) = 0$  for some symmetric  $A \in \mathbb{R}^{D \times D}$  with multiplicity  $m$  then there exist  $m$  orthonormal eigenvectors of  $A$ , corresponding to  $\tilde{\lambda}$ .
  - For any symmetric matrix  $A \in \mathbb{R}^{D \times D}$  there exists orthonormal basis of eigenvectors of this matrix.

# Spectral decomposition

## Theorem 1 (Spectral decomposition.)

*Every symmetric  $A \in \mathbb{R}^{D \times D}$  can be factorized*

$$A = P \Lambda P^T$$

*where  $P \in \mathbb{R}^{D \times D}$  is orthogonal matrix whose columns  $p_1, \dots, p_D$  are eigenvectors of  $A$  and  $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_D\}$  is diagonal matrix with corresponding eigenvalues on the diagonal.*

**Intuition:** transformation  $Ax$  by symmetric matrix is equivalent to

- 1 rotation of  $x$  to ortonormal basis formed by eigenvectors of  $A$
- 2 scaling coordinates in this basis by eigenvalues  $\lambda_1, \dots, \lambda_D$ .
- 3 reverse rotation to initial basis.

# Positivity of matrices

## Definition

Symmetric matrix  $A \in \mathbb{R}^{D \times D}$  is called *positive semi-definite* when

$$\forall x \in \mathbb{R}^D : \langle x, Ax \rangle = x^T Ax \geq 0$$

- Positive semi-definiteness of  $A$  is denoted as  $A \succcurlyeq 0$ .
- Are the following matrix positive

semi-definite:  $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ ?<sup>2</sup>

## Theorem

*Symmetric matrix  $A$  is positive semi-definite  $\iff$  all its eigenvalues are non-negative.*

<sup>2</sup>Are these matrices  $\succcurlyeq 0$  for  $D > 2$  dimensional case?

## Distribution properties

Let  $x_1, \dots, x_N \in \mathbb{R}^D$  be observations of some vector random variable  $x \sim F$ . Group these observations into matrix

$$X = [x_1^T, \dots, x_N^T]^T \in \mathbb{R}^{N \times D}$$

- Expectation:  $\mu = \mathbb{E}x$
- Covariance matrix  $\Sigma = \mathbb{E}(x - \mu)(x - \mu)^T$ .
- Sample mean  $\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n$
- Sample covariance matrix  
$$\hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu})(x_n - \hat{\mu})^T = \frac{1}{N} X^T X$$

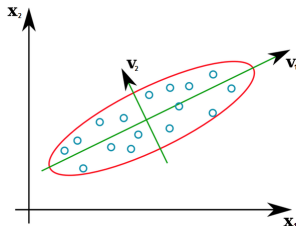


## Properties of covariance matrix

- If random vector  $x \in \mathbb{R}^D$  has covariance  $\Sigma$ , then random variable  $\alpha^T x$  for any  $\alpha \in \mathbb{R}^D$  has variance  $\alpha^T \Sigma \alpha$ .
- Covariance matrix is symmetric and positive semi-definite.
- For any matrix  $X \in \mathbb{R}^{N \times D}$   $X^T X \in \mathbb{R}^{D \times D}$  is symmetric and positive semi-definite.
  - So all eigenvalues of  $X^T X$  are non-negative
- Sample covariance matrix is symmetric and positive semi-definite.

## Estimating scatter by covariance matrix

- For different  $\alpha \in \mathbb{R}^D$ ,  $\alpha^T \alpha = 1$  estimate  $\text{var}(\alpha^T x) = \alpha^T \Sigma \alpha = \alpha^T P \Lambda P^T \alpha = (\Lambda^{1/2} P^T \alpha)^T (\Lambda^{1/2} P^T \alpha)$ .
- $\alpha$  gets rotated to new orthonormal basis  $P$ , then stretched along axes with factors  $\sqrt{\lambda_1}, \dots, \sqrt{\lambda_D}$ .



- We can evaluate scatter by looking at trace  $\Sigma = \lambda_1 + \dots + \lambda_D$  or  $\det \Sigma = \lambda_1 \cdot \dots \cdot \lambda_D$ .
  - This is similar to arithmetic and geometric averaging.

## Vector derivatives

- Suppose  $x = [x^1, \dots, x^D]$  and  $f(x) = f(x^1, \dots, x^D)$ . Vector derivative

$$\frac{\partial f(x)}{\partial x} := \begin{pmatrix} \frac{\partial f(x)}{\partial x^1} \\ \frac{\partial f(x)}{\partial x^2} \\ \dots \\ \frac{\partial f(x)}{\partial x^D} \end{pmatrix}$$

- For any  $x, b \in \mathbb{R}^D$  it holds that<sup>3</sup>:

$$\frac{\partial [b^T x]}{\partial x} = b$$

- For any  $x \in \mathbb{R}^D$  and symmetric  $B \in \mathbb{R}^{D \times D}$  it holds that<sup>4</sup>:

$$\frac{\partial [x^T B x]}{\partial x} = 2Bx$$

---

<sup>3</sup>Prove it.

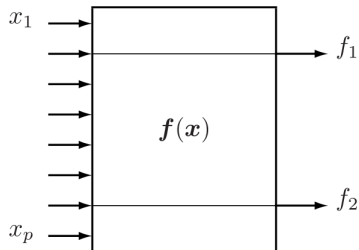
<sup>4</sup>Prove it. How will the formula change for non-symmetric  $B$ ?

# Table of Contents

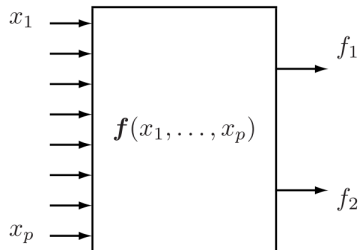
- 1 Linear algebra reminder
- 2 Dimensionality reduction intro
- 3 Principal component analysis

# Dimensionality reduction

## Feature selection / Feature extraction



(a) feature selector



(b) feature extractor

**Feature extraction:** find transformation of original data which extracts most relevant information for machine learning task.

# Applications of dimensionality reduction

## Applications:

- visualization in 2D or 3D
- reduce operational costs on data storage, transfer and processing
  - memory
  - disk
  - CPU usage
- remove multi-collinearity to improve performance of some machine-learning models

# Categorization of dimensionality reduction methods

Supervision:

- supervised
- unsupervised

Mapping to reduced space:

- linear
- non-linear

Principal components analysis - linear unsupervised method of dimensionality reduction.

# Table of Contents

- 1 Linear algebra reminder
- 2 Dimensionality reduction intro
- 3 Principal component analysis
  - Definition
  - Applications of PCA
  - Application details
  - Construction of principal components
  - Proof of optimality of principal components



### 3 Principal component analysis

- Definition
- Applications of PCA
- Application details
- Construction of principal components
- Proof of optimality of principal components

# Projections, orthogonal complements

- For point  $x$  and subspace  $L$  denote:
  - $p$ : the projection of  $x$  on  $L$
  - $h$ : orthogonal complement
  - $x = p + h$ ,  $\langle p, h \rangle = 0$ .
- For training set  $x_1, x_2, \dots, x_N$  and subspace  $L$  find:
  - projections:  $p_1, p_2, \dots, p_N$
  - orthogonal complements:  $h_1, h_2, \dots, h_N$ .

# Best subspace fit<sup>5</sup>

## Definition 1

Best-fit  $k$ -dimensional subspace for a set of points  $x_1, x_2, \dots, x_N$  is a subspace, spanned by  $k$  vectors  $v_1, v_2, \dots, v_k$ , solving

$$\sum_{n=1}^N \|h_n\|^2 \rightarrow \min_{v_1, v_2, \dots, v_k}$$

## Proposition 1

Vectors  $v_1, v_2, \dots, v_k$ , solving

$$\sum_{n=1}^N \|p_n\|^2 \rightarrow \max_{v_1, v_2, \dots, v_k}$$

also define best-fit  $k$ -dimensional subspace.

<sup>5</sup>Prove 1 using that  $\|x\|^2 = \|p\|^2 + \|h\|^2$  for  $x = p + h$  and  $\langle p, h \rangle = 0$ .

# Definition of PCA

## Definition 2

Principal components  $a_1, a_2, \dots, a_k$  are vectors, forming orthonormal basis in the  $k$ -dimensional subspace of best fit.

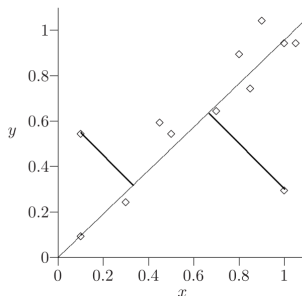
- Properties:
  - Not invariant to translation:
    - center data before PCA:

$$x \leftarrow x - \mu \text{ where } \mu = \frac{1}{N} \sum_{n=1}^N x_n$$

- Not invariant to scaling:
  - scale features to have unit variance before PCA

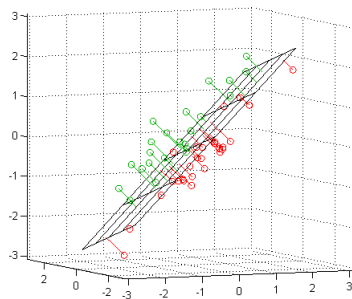
## Example: line of best fit

- In PCA the sum of squared perpendicular distances to line is minimized:



- *What is the difference with least squares minimization in regression?*

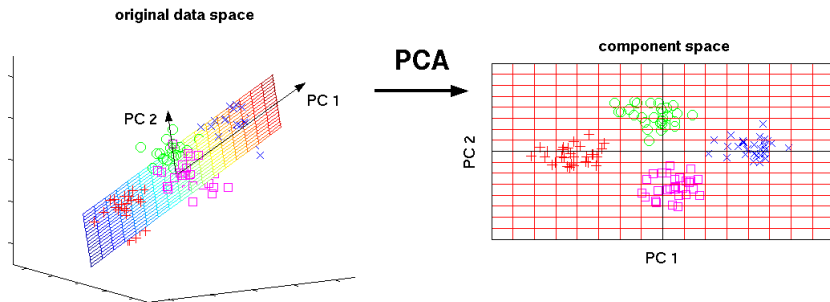
## Example: plane of best fit



### 3 Principal component analysis

- Definition
- Applications of PCA
- Application details
- Construction of principal components
- Proof of optimality of principal components

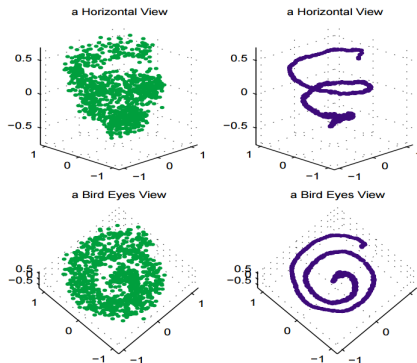
# Visualization





# Data filtering

Remove noise to get a cleaner picture of data distribution:



X. Huo and Jihong Chen (2002). Local linear projection (LLP). First IEEE Workshop on Genomic Signal Processing and Statistics (GENSIPS), Raleigh, NC, October.  
<http://www.gensips.gatech.edu/proceedings/>.

## Economic description of data

Faces database:



## Eigenvectors (called eigenfaces)

Projections on first several eigenvectors describe most of face variability.



# Text analysis

- Objects=text files
- Binary, TF, TF-IDF representations have huge  $D$ .
  - math operations with  $X$  - inefficient
    - ML methods work longer
- Sparsity induces complications with query matching
  - consider query “automobile”
  - simple cosine-metric matching won’t match documents with “car”, “bus”, etc.

# Latent semantic analysis (LSA)

## Latent semantic analysis (LSA)

Get economical document representations with coordinates of most important PCA components found without centering.

Comments:

- usually 200-300 components are sufficient.
- Do *not* center  $X$  before computing PCA
  - otherwise will lose sparsity of  $X$
  - $\mu \approx 0$  anyway, because most features are 0.
- Technically done with truncated SVD of  $X$ .

### 3 Principal component analysis

- Definition
- Applications of PCA
- **Application details**
- Construction of principal components
- Proof of optimality of principal components

## Quality of approximation

Consider vector  $x$ . Since all  $D$  principal components form a full orthonormal basis,  $x$  can be written as

$$x = \langle x, a_1 \rangle a_1 + \langle x, a_2 \rangle a_2 + \dots + \langle x, a_D \rangle a_D$$

Let  $p^K$  be the projection of  $x$  onto subspace spanned by first  $K$  principal components:

$$p^K = \langle x, a_1 \rangle a_1 + \langle x, a_2 \rangle a_2 + \dots + \langle x, a_K \rangle a_K$$

Error of this approximation is

$$h^K = x - p^K = \langle x, a_{K+1} \rangle a_{K+1} + \dots + \langle x, a_D \rangle a_D$$

## Quality of approximation

Using that  $a_1, \dots, a_D$  is an orthonormal set of vectors, we get

$$\begin{aligned}\|x\|^2 &= \langle x, x \rangle = \langle x, a_1 \rangle^2 + \dots + \langle x, a_D \rangle^2 \\ \|p^K\|^2 &= \langle p^K, p^K \rangle = \langle x, a_1 \rangle^2 + \dots + \langle x, a_K \rangle^2 \\ \|h^K\|^2 &= \langle h^K, h^K \rangle = \langle x, a_{K+1} \rangle^2 + \dots + \langle x, a_D \rangle^2\end{aligned}$$

We can measure how well first  $K$  components describe our dataset  $x_1, x_2, \dots, x_N$  using relative loss

$$L(K) = \frac{\sum_{n=1}^N \|h_n^K\|^2}{\sum_{n=1}^N \|x_n\|^2} \quad (1)$$

or relative score

$$S(K) = \frac{\sum_{n=1}^N \|p_n^K\|^2}{\sum_{n=1}^N \|x_n\|^2} \quad (2)$$

Evidently  $L(K) + S(K) = 1$ .



## Contribution of individual component

Contribution of  $a_k$  for explaining  $x$  is  $\langle x, a_k \rangle^2$ .

Contribution of  $a_k$  for explaining  $x_1, x_2, \dots, x_N$  is:

$$\sum_{n=1}^N \langle x_n, a_k \rangle^2$$

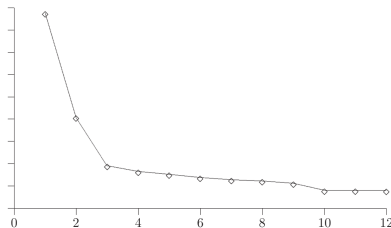
Explained variance ratio:

$$E(a_k) = \frac{\sum_{n=1}^N \langle x_n, a_k \rangle^2}{\sum_{d=1}^D \sum_{n=1}^N \langle x_n, a_d \rangle^2} = \frac{\sum_{n=1}^N \langle x_n, a_k \rangle^2}{\sum_{n=1}^N \|x_n\|^2}$$

- Explained variance ratio measures relative contribution of component  $a_k$  to explaining our dataset  $x_1, \dots, x_N$ .
- Note that  $\sum_{k=1}^K E(a_k) = S(K)$ .

## How many principal components to select?

- Data visualization: 2 or 3 components.
- Take most significant components until their explained variance ratio falls sharply down:



- Or take minimum  $K$  such that  $L(K) \leq t$  or  $S(K) \geq 1 - t$ , where typically  $t = 0.95$ .

# Transformation $\xi \rightleftharpoons x$

Dependence between original and transformed features:

$$\xi = A^T(x - \mu), \quad x = A\xi + \mu,$$

where  $\mu = \frac{1}{N} \sum_{n=1}^N x_n$ .

Taking first  $r$  components -  $A_r = [a_1|a_2|\dots|a_r]$ , we get the image of the reduced transformation:

$$\xi_r = A_r^T(x - \mu)$$

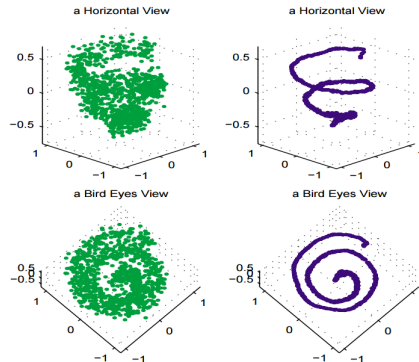
$\xi_r$  will correspond to

$$x_r = A \begin{pmatrix} \xi_r \\ 0 \end{pmatrix} + \mu = A_r \xi_r + \mu$$

$$x_r = A_r A_r^T(x - \mu) + \mu$$

$A_r A_r^T$  is projection matrix with rank  $r$   
(follows from the property  $\text{rank}[A_r A_r^T] = \text{rank}[A_r^T A_r]$  for arbitrary  $A_r$ ).

# Local linear projection



X. Huo and Jihong Chen (2002). Local linear projection (LLP). First IEEE Workshop on Genomic Signal Processing and Statistics (GENSIPS), Raleigh, NC, October.  
<http://www.gensips.gatech.edu/proceedings/>.

## Local linear projection

Local linear projection method makes denoised version of original data by locally projecting it onto hyperplane of small rank.

**INPUT:**

p-local dimensionality of data  
K-number of nearest neighbours

for each  $x_i$  in  $X$ :

- 1) find K nearest neighbours of  $x_i$ :  $x_{j(i,1)}, \dots, x_{j(i,K)}$
- 2) find linear hyperplane  $L_p$  of dimensionality  $p$ ,  
describing  $x_{j(i,1)}, \dots, x_{j(i,K)}$  # hyperplane-subspace with offset
- 3) let  $\hat{x}_i$  be the projection of  $x_i$  onto this hyperplane

**OUTPUT:**

denoised version of objects  $\hat{x}_1, \hat{x}_2, \dots, \hat{x}_K$ .

- Projection is made on hyperplane, not subspace!

## PCA solution

- Center  $x_1, \dots, x_N$  to have zero mean.
- Scale  $x_1, \dots, x_N$  to have equal variance.
- Form  $X = [x_1^T; \dots, x_N^T]^T \in \mathbb{R}^{N \times D}$
- Estimate sample covariance matrix of  $x$ :  $\hat{\Sigma} = \frac{1}{N} X^T X$
- Find eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D \geq 0$  and corresponding eigenvectors  $a_1, a_2, \dots, a_D$ .
- $a_1, a_2, \dots, a_k$  are first  $k$  principal components,  $k = 1, 2, \dots, D$ .
- Sum of squared projections onto  $a_i$  is  $\|X a_i\|^2 = \lambda_i$ .
- *Explained variance ratio* by component  $a_i$  is equal to

$$\frac{\lambda_i}{\sum_{d=1}^D \lambda_d}$$

### 3 Principal component analysis

- Definition
- Applications of PCA
- Application details
- Construction of principal components
- Proof of optimality of principal components

# Constructive definition of PCA

- Principal components  $a_1, a_2, \dots, a_D \in \mathbb{R}^D$  are found such that
$$\langle a_i, a_j \rangle = \begin{cases} 1, & i = j \\ 0 & i \neq j \end{cases}$$
- $Xa_i$  is a vector of projections of all objects onto the  $i$ -th principal component.
- For any object  $x$  its projections onto principal components are equal to:

$$p = A^T x = [\langle a_1, x \rangle, \dots, \langle a_D, x \rangle]^T$$

where  $A = [a_1; a_2; \dots, a_D] \in \mathbb{R}^{D \times D}$ .



# Constructive definition of PCA

- ①  $a_1$  is selected to maximize  $\|Xa_1\|$  subject to  $\langle a_1, a_1 \rangle = 1$
- ②  $a_2$  is selected to maximize  $\|Xa_2\|$  subject to  $\langle a_2, a_2 \rangle = 1$ ,  
 $\langle a_2, a_1 \rangle = 0$
- ③  $a_3$  is selected to maximize  $\|Xa_3\|$  subject to  $\langle a_3, a_3 \rangle = 1$ ,  
 $\langle a_3, a_1 \rangle = \langle a_3, a_2 \rangle = 0$   
etc.
- It turns out that:
  - $a_1, \dots, a_k$  form  $k$ -dimensional subspace of best fit.
  - $a_1, a_2, \dots$  are first, second, ... eigenvectors of  $X^T X$  (ordered by decreasing eigenvalue).

### 3 Principal component analysis

- Definition
- Applications of PCA
- Application details
- Construction of principal components
- Proof of optimality of principal components

# Componentwise optimization leads to best fit subspace

## Theorem 2

*Let  $L_k$  be the subspace spanned by  $a_1, a_2, \dots, a_k$ . Then for each  $k$   $L_k$  is the best-fit  $k$ -dimensional subspace for  $X$ .*

Proof: use induction. For  $k = 1$  the statement is true by definition since projection maximization is equivalent to distance minimization.

Suppose theorem holds for  $k - 1$ . Let  $L_k$  be the plane of best-fit of dimension with  $\dim L = k$ . We can always choose an orthonormal basis of  $L_k$   $b_1, b_2, \dots, b_k$  so that

$$\begin{cases} \|b_k\| = 1 \\ b_k \perp a_1, b_k \perp a_2, \dots, b_k \perp a_{k-1} \end{cases} \quad (3)$$

by setting  $b_k$  perpendicular to projections of  $a_1, a_2, \dots, a_{k-1}$  on  $L_k$ .

## Componentwise optimization leads to best fit subspace

Consider the sum of squared projections:

$$\|Xb_1\|^2 + \|Xb_2\|^2 + \dots + \|Xb_{k-1}\|^2 + \|Xb_k\|^2$$

By induction proposition  $L[a_1, a_2, \dots, a_{k-1}]$  is space of best fit of rank  $k-1$  and  $L[b_1, \dots, b_{k-1}]$  is some space of same rank, so sum of squared projections on it is smaller:

$$\|Xb_1\|^2 + \|Xb_2\|^2 + \dots + \|Xb_{k-1}\|^2 \leq \|Xa_1\|^2 + \|Xa_2\|^2 + \dots + \|Xa_{k-1}\|^2$$

and

$$\|Xb_k\|^2 \leq \|Xa_k\|^2$$

since  $b_k$  by (3) satisfies constraints of optimization problem (??) and  $a_k$  is its optimal solution.

# Summary

- Every symmetric matrix  $A$  can be decomposed into rotation, scaling and backward rotation:

$$A = P\Lambda P^T$$

- Sample covariance matrix  $\frac{1}{N}X^T X$  is symmetric and  $\succcurlyeq 0$ .
  - so it has non-negative eigenvalues  $\lambda_1 \geq \dots \geq \lambda_D \geq 0$  with corresponding eigenvectors  $a_1, \dots, a_D$ .
  - spread of distribution is characterized by eigenvalues.

## Summary

- Dimensionality reduction - common preprocessing step for efficiency and numerical stability.
- Subspace of best fit of rank  $k$  for training set  $x_1, \dots, x_N$  is  $k$ -dimensional subspace  $\mathcal{L}(b_1, \dots, b_k)$ , minimizing:

$$\|h_1\|^2 + \dots + \|h_N\|^2 \rightarrow \min_{b_1, \dots, b_k}$$

- Solution vectors are called top  $k$  principal components.
- Principal component analysis - expression of  $x$  in terms of first  $k$  principal components.
- It is unsupervised linear dimensionality reduction.
- Solution: principal components  $a_1, \dots, a_k$  are top  $k$  eigenvectors of  $X^T X$ .