

# LTAT.02.004 MACHINE LEARNING II

## **Affine data projections**

**based on normal distribution**

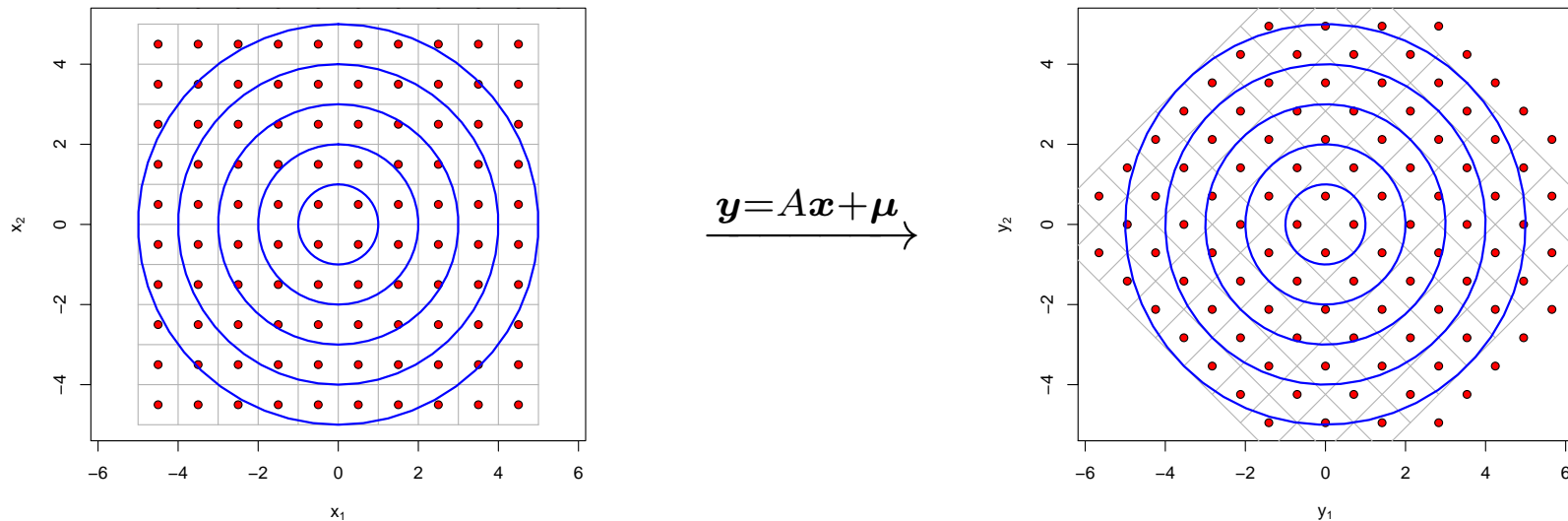
Sven Laur  
University of Tartu

# Principal component analysis

# Distribution reconstruction task

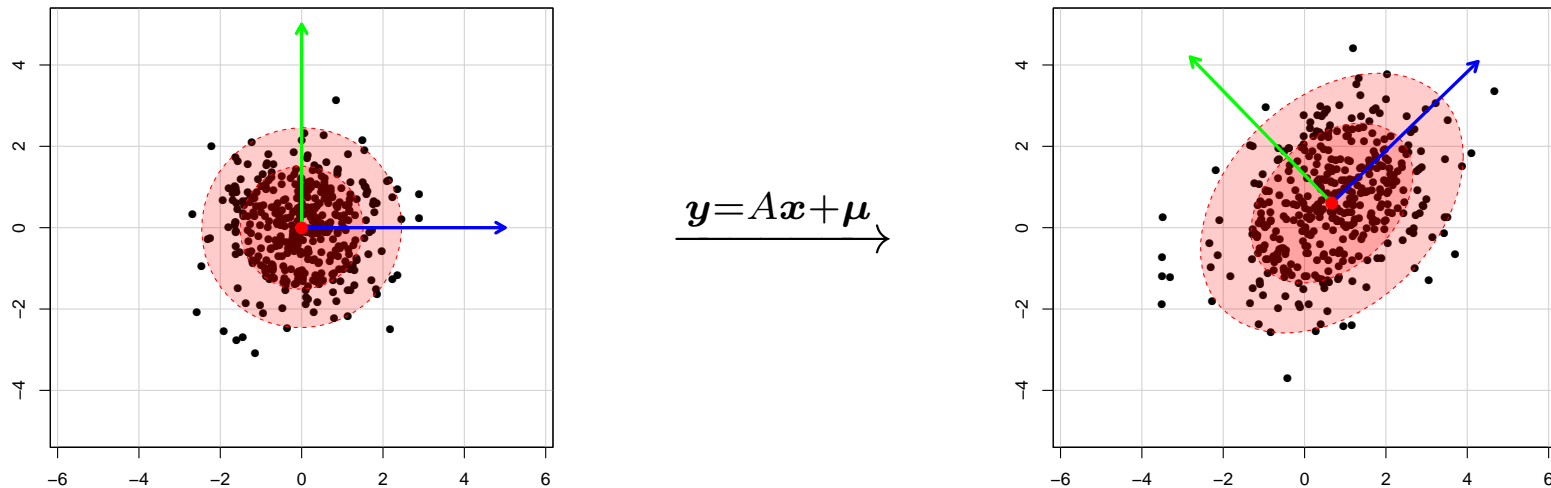
**Original goal.** Given the set of observations  $y_1, \dots, y_m$  determine the affine transformation  $y = Ax + \mu$  and original source signals  $x_1, \dots, x_m$ .

**Impossibility result.** The matrix  $A$  can be recovered *only* up to rotations.



# Simplified distribution reconstruction task

**Achievable goal.** Given the set of observations  $y_1, \dots, y_m$  determine the affine transformation by fixing the centre and axis of the ellipsoid.



- ▷ We need to find the origin and semi-axes  $a_1, \dots, a_n$  of the ellipsoid.
- ▷ Unit vectors  $e_1, \dots, e_n$  are mapped to semi-axes  $a_1, \dots, a_n$  of ellipsoid.

## Variance for a fixed direction

**Fact.** Orthogonal projection onto a unit vector  $w$  is given by scalar product.

**Question.** What is the direction  $w$  that maximises the variance for ellipsoid?

$$\mathbf{Var}(w^T \text{diag}(a)x) = \mathbf{Var}\left(\sum_{i=1}^n w_i a_i x_i\right) = \sum_{i=1}^n w_i^2 a_i^2 .$$

The variance is maximised in the direction of the longest ellipse axis  $a_1$ .

**Question.** How is the center of the ellipsoid and mean values connected?

$$\mathbf{E}(Ax + \mu) = \mathbf{E}(Ax) + \mathbf{E}(\mu) = \mu .$$

## Principal component analysis

- ▷ Compute the average value of the observations  $\mathbf{y}_1, \dots, \mathbf{y}_m$ :

$$\hat{\boldsymbol{\mu}} \leftarrow \frac{\mathbf{y}_1 + \dots + \mathbf{y}_m}{m} .$$

- ▷ Centre the data by substituting  $\hat{\boldsymbol{\mu}}$ :

$$\mathbf{y}_i \leftarrow \mathbf{y}_i - \hat{\boldsymbol{\mu}}, \quad i \in \{1, \dots, m\} .$$

- ▷ Find the unit direction  $\mathbf{w}_1$  that has *a maximal empirical* variance:

$$F(\mathbf{w}) = \text{Var}(\mathbf{w}^T \mathbf{y}_1, \dots, \mathbf{w}^T \mathbf{y}_n) = \frac{(\mathbf{w}^T \mathbf{y}_1)^2 + \dots + (\mathbf{w}^T \mathbf{y}_m)^2}{m} .$$

- ▷ Find unit directions  $\mathbf{w}_i$  orthogonal to previous directions that maximise the empirical variance of the corresponding the projection onto  $\mathbf{w}_i$ .

## Covariance matrix and optimisation goal

We can use matrix algebra to simplify the variance estimate

$$\begin{aligned} F(\mathbf{w}) &= \frac{1}{m} \cdot \left( \mathbf{w}^T \mathbf{y}_1 \mathbf{y}_1^T \mathbf{w} + \cdots + \mathbf{w}^T \mathbf{y}_m \mathbf{y}_m^T \mathbf{w} \right) \\ &= \mathbf{w}^T \left( \frac{\mathbf{y}_1 \mathbf{y}_1^T + \cdots + \mathbf{y}_m \mathbf{y}_m^T}{m} \right) \mathbf{w} \end{aligned}$$

The  $n \times n$  matrix in the middle is known as a *covariance matrix*  $\Sigma$ .

Due to the restriction  $\|\mathbf{w}\|_2^2 = \mathbf{w}^T \mathbf{w} = 1$ , we have to use Lagrange' trick:

$$F_*(\mathbf{w}) = \mathbf{w}^T \Sigma \mathbf{w} - 2\lambda \mathbf{w}^T \mathbf{w} \quad \Rightarrow \quad \frac{\partial F_*(\mathbf{w})}{\partial \mathbf{w}} = 2\Sigma \mathbf{w} - 2\lambda \mathbf{w} = \mathbf{0}.$$

## Principal components as eigenvectors

The  $F_*(\boldsymbol{w})$  is maximised only if the direction  $\boldsymbol{w}$  is an *eigenvector* of  $\Sigma$ :

$$\Sigma \boldsymbol{w} = \lambda \boldsymbol{w} \quad \Rightarrow \quad \boldsymbol{w}^T \Sigma \boldsymbol{w} = \boldsymbol{w}^T \lambda \boldsymbol{w} = \lambda .$$

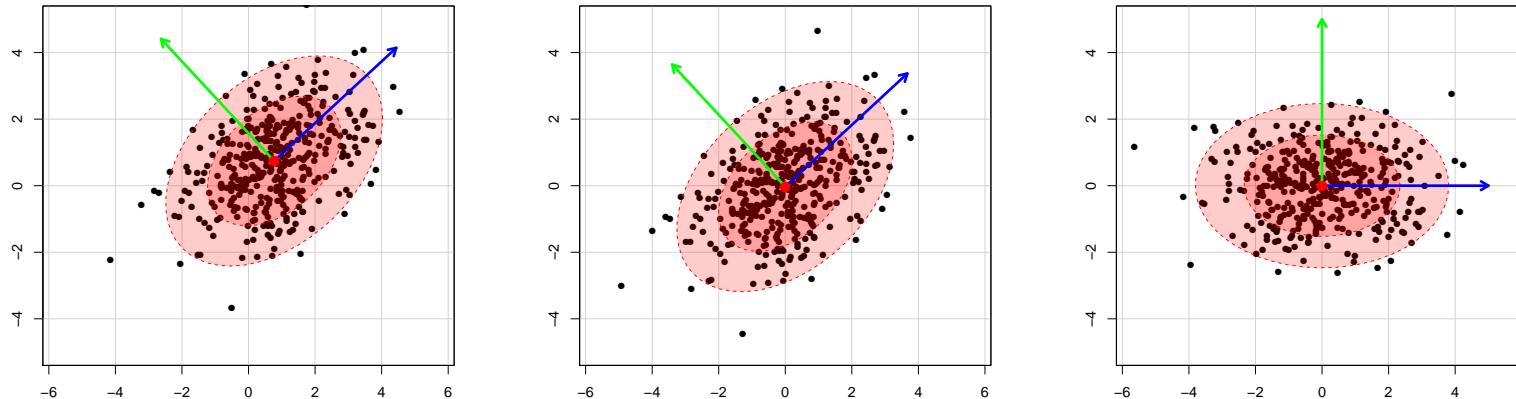
**Fact.** If  $n \times n$  matrix is symmetric and positively definite then there exists  $n$  orthogonal eigenvectors  $\boldsymbol{w}_1, \dots, \boldsymbol{w}_n$  with *eigenvalues*  $\lambda_1 \geq \dots \geq \lambda_n > 0$ .

**Corollary.** Principal components corresponding to observations  $\boldsymbol{y}_1, \dots, \boldsymbol{y}_m$  are the eigenvectors of the covariance matrix  $\Sigma$ .



# Principal component analysis as a rotation

Reconstruction of the source signal can be viewed as a *translation* followed by a *rotation* to orientate the ellipsoid wrt coordinate axis.



As vectors  $w_1, \dots, w_n$  are orthogonal, the rotation can be done through computing projections (read scalar products):

$$\hat{x}_i = (w_1 || \dots || w_n)^T (y_i - \hat{\mu}_0) = W(y_i - \hat{\mu}) \quad .$$

## Maximum likelihood estimate

The algorithm formulated above was based on *ad hoc* reasoning:

- ▷ Empirical estimates for the mean and variance are not precise!

Theoretically correct way to handle the problem is

- ▷ obtain the maximum likelihood estimate on the model parameters,
- ▷ determine the translation and rotation based on the model parameters.

What are the model parameters?

- ▷ Parameters of the density formula  $\Sigma$  and  $\mu$ .
- ▷ Parameters of the affine transformation  $A$  and  $\mu$ .

## Likelihood function under iid assumption

If all observations  $\mathbf{y}_1, \dots, \mathbf{y}_m$  are independent then

$$p[\mathbf{y}_1, \dots, \mathbf{y}_m | \Sigma, \boldsymbol{\mu}] = \prod_{i=1}^m p[\mathbf{y}_i | \Sigma, \boldsymbol{\mu}]$$

where

$$p[\mathbf{y}_i | \Sigma, \boldsymbol{\mu}] = \frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{\sqrt{\det(\Sigma)}} \cdot \exp \left( -\frac{(\mathbf{y}_i - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{y}_i - \boldsymbol{\mu})}{2} \right)$$

The *log-likelihood* of the data  $\ln p[\mathbf{y}_1, \dots, \mathbf{y}_m | \Sigma, \boldsymbol{\mu}]$  can be expressed

$$\mathcal{L}(\Sigma, \boldsymbol{\mu}) = \text{const} + \frac{m}{2} \cdot \ln \det(\Sigma^{-1}) - \sum_{i=1}^m \frac{(\mathbf{y}_i - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{y}_i - \boldsymbol{\mu})}{2}$$

Now we have to find the arrangement  $(\Sigma, \boldsymbol{\mu})$  that maximises  $\mathcal{L}(\Sigma, \boldsymbol{\mu})$ .

## Gradients of the log-likelihood function

Gradient with respect to the shift  $\mu$ :

$$\frac{\partial \mathcal{L}}{\partial \mu} = - \sum_{i=1}^m \frac{\partial}{\partial \mu} \frac{(\mathbf{y}_i - \mu)^T \Sigma^{-1} (\mathbf{y}_i - \mu)}{2} = - \sum_{i=1}^m \frac{\Sigma^{-1} (\mathbf{y}_i - \mu)}{2} \cdot (-1)$$

Gradient with respect to the inverse matrix  $\Sigma^{-1}$ :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial (\Sigma^{-1})} &= \frac{m}{2} \cdot \frac{\partial}{\partial (\Sigma^{-1})} \ln \det(\Sigma^{-1}) - \sum_{i=1}^m \frac{\partial}{\partial (\Sigma^{-1})} \frac{(\mathbf{y}_i - \mu)^T \Sigma^{-1} (\mathbf{y}_i - \mu)}{2} \\ &= \frac{m}{2} \cdot \Sigma^T - \sum_{i=1}^m \frac{(\mathbf{y}_i - \mu)^T (\mathbf{y}_i - \mu)}{2} \end{aligned}$$

As  $\Sigma$  is symmetric and  $\Sigma^{-1}$  exists we can derive closed form solutions.

## Maximum likelihood estimates for parameters

The shift must be the mean of all observations

$$\boldsymbol{\mu} = \frac{1}{m} \cdot \sum_{i=1}^m \mathbf{y}_i \ .$$

The covariance matrix

$$\boldsymbol{\Sigma} = \frac{1}{m} \cdot \sum_{i=1}^m (\mathbf{y}_i - \boldsymbol{\mu})^T (\mathbf{y}_i - \boldsymbol{\mu})$$

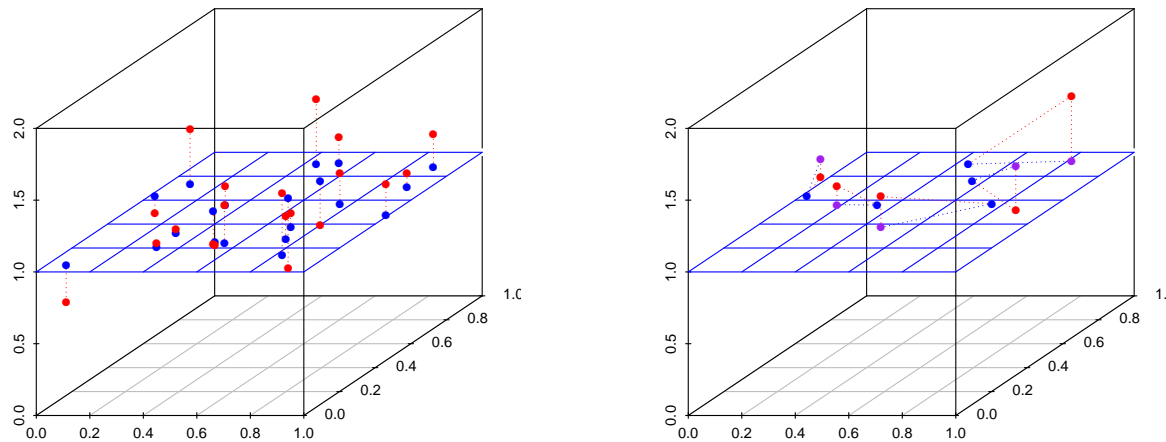
**Correctness of PCA.** As ML estimates are exactly the same we used in principal component analysis, the method is theoretically justified!

# Principal component analysis

## Alternative formalisations

# Dimensionality reduction

What if the actual data  $\mathbf{x}_1, \dots, \mathbf{x}_m$  lies in a lower-dimensional plane and the observation  $\mathbf{y}_1, \dots, \mathbf{y}_m$  are obtained by random shifts?



The shifts can be either orthogonal to the plane or just random. The first model is easier to analyse while the second is more plausible.

## Maximum likelihood estimate

Let  $\mathcal{H}$  be the plane. Assume that the random shifts  $\varepsilon_i$  are orthogonal to the plane and have a normal distribution  $\mathcal{N}(0, \sigma I)$ . Then

$$p[\mathbf{y}_i | \mathcal{H}, \sigma] = \text{const} \cdot \exp \left( -\frac{d_i^2}{2\sigma^2} \right)$$

where  $d_i$  is the distance between the plane  $\mathcal{H}$  and the point  $\mathbf{y}_i$ . Thus

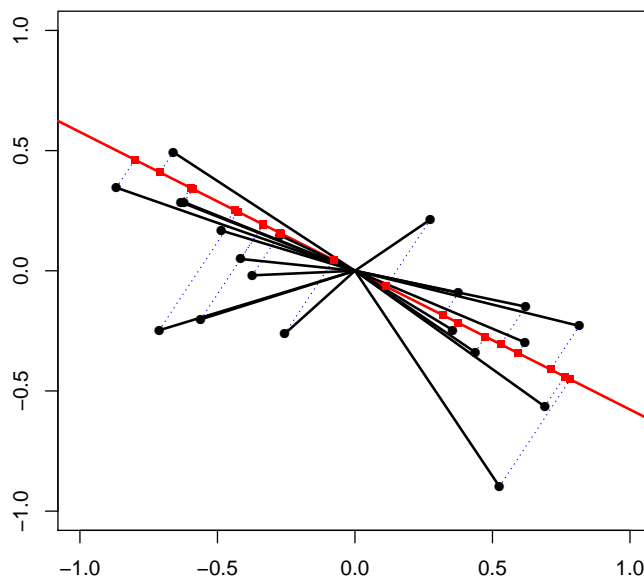
$$p[\mathbf{y}_1, \dots, \mathbf{y}_m | \mathcal{H}, \sigma] = \text{const} \cdot \exp \left( -\sum_{i=1}^m \frac{d_i^2}{2\sigma^2} \right)$$

and the maximum likelihood estimate of the plane minimises sum of the distance squares. Corresponding estimates of  $\mathbf{x}_1, \dots, \mathbf{x}_m$  are projections of  $\mathbf{y}_1, \dots, \mathbf{y}_m$  to the plane  $\mathcal{H}$ .



## Another characterisation of PCA

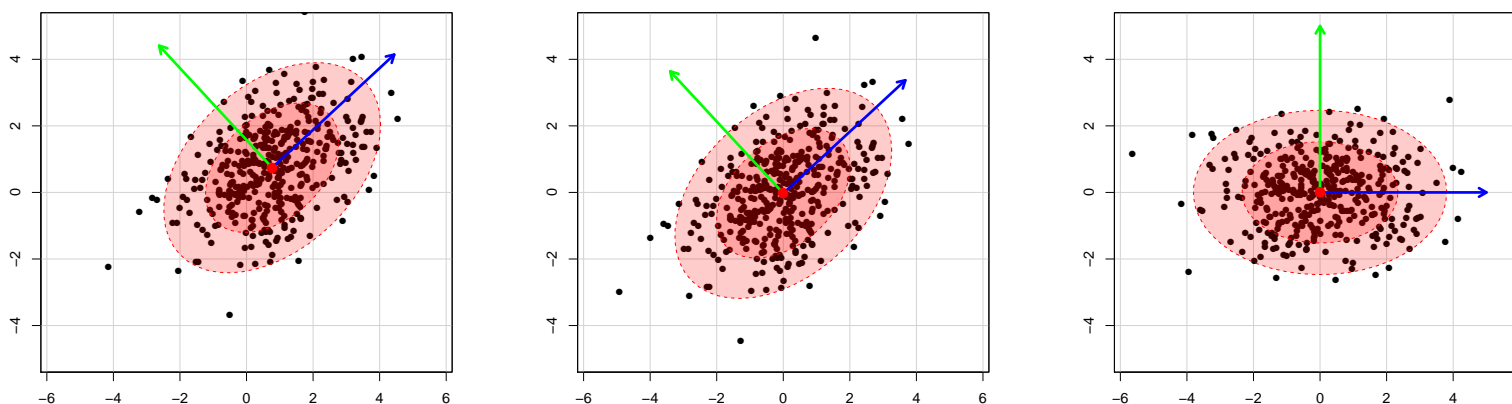
**Fact.** If the data is centred then PCA chooses the direction  $w_1$  such that the sum of squares of the projections  $w_1^T y_i$  is maximal.



**Corollary.** PCA chooses directions  $w_1, \dots, w_n$  such that the sum of distance squares from the hyperplane formed by  $w_1, \dots, w_k$  is minimal.

# PCA as a dimensionality reduction tool

**Corollary.** PCA rotates the data such way that first  $k$  coordinates of the rotated data correspond to maximum likelihood reconstructions of original vectors corrupted with white Gaussian noise  $\mathcal{N}(0, \sigma I)$ .



Alternatively, we can view the last components of the source signal  $x$  as the uninformative noise. The overall noise component should be small.