

MTAT.03.227 MACHINE LEARNING

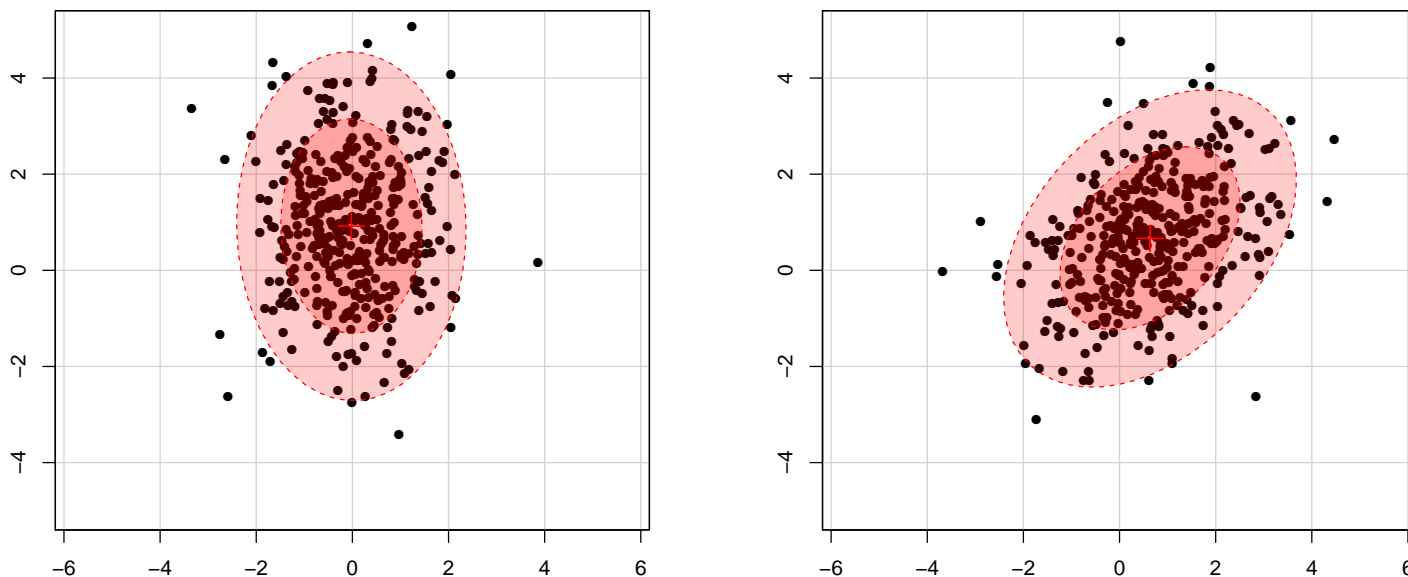
# **Principal Component Analysis**

**Fitting Multivariate Normal Distributions**

Sven Laur  
University of Tartu

## Two-dimensional normal distribution

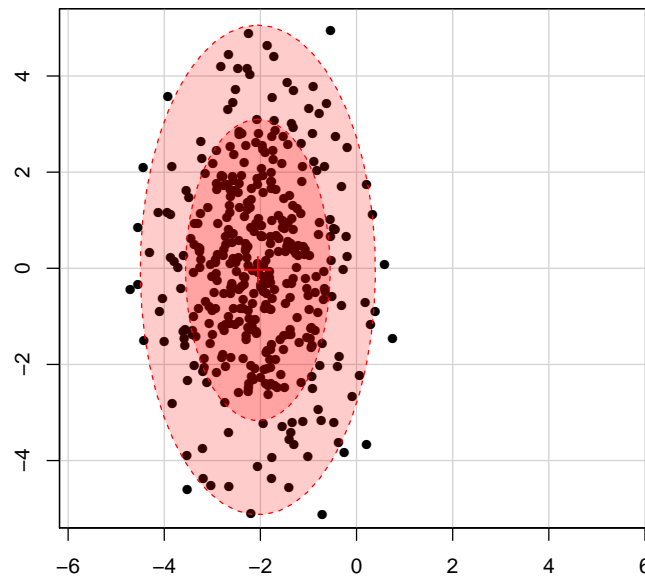
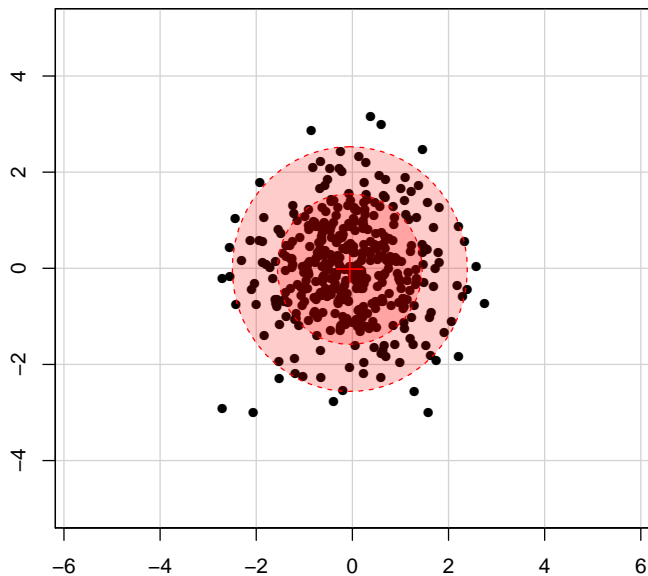
We can form a two-dimensional normal distribution by considering two independent quantities that have a normal distribution.



As the choice of coordinate axis is sometimes arbitrary, e.g., bullet holes in shooting targets, there are also other ways to form a normal distribution.

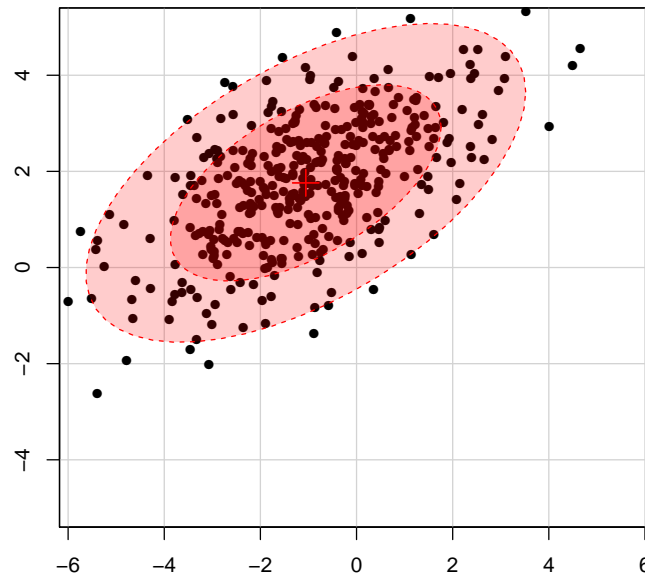
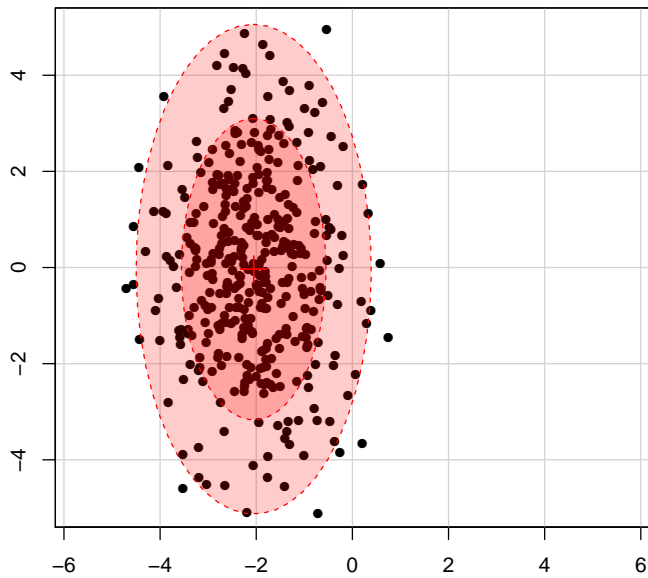
## Scaling and shifting

A *univariate normal distribution*  $\mathcal{N}(\mu, \sigma^2)$  can be obtained from a *standard normal distribution*  $\mathcal{N}(0, 1)$  by shifting and scaling. Hence, the distribution  $y_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $y_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$  can be obtained by scaling and shifting source distribution  $x_1 \sim \mathcal{N}(0, 1)$  and  $x_2 \sim \mathcal{N}(0, 1)$ .



# Rotations

Besides rescaling and shifting we can also rotate the coordinate axis.



## Affine transformations

Let  $x_1 \sim \mathcal{N}(0, 1)$  and  $x_2 \sim \mathcal{N}(0, 1)$  be the quantities generated by independent normal distributions. Let  $y_1$  and  $y_2$  quantities that are observed after the scaling, translation and rotation of the coordinate plane.

Then we can express  $\mathbf{x}$  and  $\mathbf{y}$  in terms of an affine transformation

$$\begin{aligned}\mathbf{y} &= A\mathbf{x} + \boldsymbol{\mu} \text{ ,} \\ \mathbf{x} &= A^{-1}(\mathbf{y} - \boldsymbol{\mu}) \text{ .}\end{aligned}$$

**Observation.** Affine transformations are closed with respect to composition, i.e., applying two affine transformations yields a new affine transformation.

**Remark.** Not all affine transformations are invertible.

## What is density?

Recall that density assigns probability to small enough regions  $\mathcal{R}$ :

$$\Pr \left[ \begin{array}{l} x_1^* \leftarrow \mathcal{N}(0, 1) : x_1 \leq x_1^* \leq x_1 + \Delta x_1 \\ x_2^* \leftarrow \mathcal{N}(0, 1) : x_2 \leq x_2^* \leq x_2 + \Delta x_2 \end{array} \right] = p(x_1, x_2) \cdot \underbrace{\Delta x_1 \Delta x_2}_S + \varepsilon$$

where  $\varepsilon = o(\Delta x_1 \cdot \Delta x_2)$  in the process  $\Delta x_1 \rightarrow 0$  and  $\Delta x_2 \rightarrow 0$ .

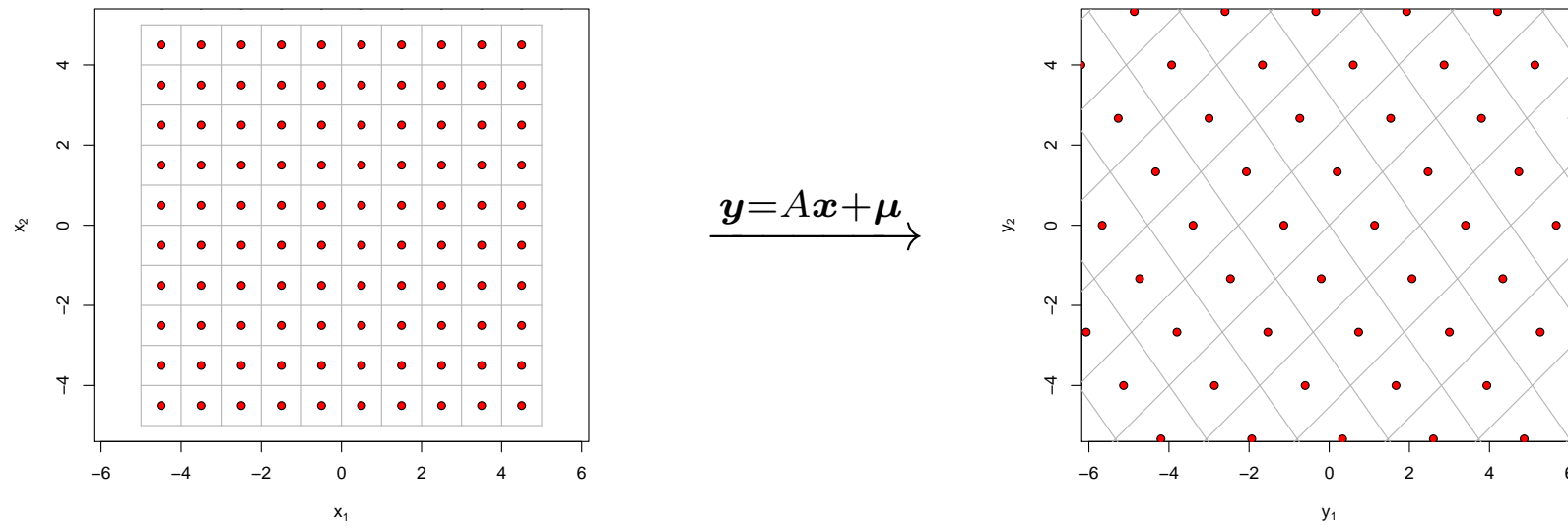
**Remark.** Regions  $\mathcal{R}$  do not have to be rectangular as long as:

- ▷ The area  $S(\mathcal{R})$  of a region can be computed.
- ▷ Probability can be assigned to the region  $\mathcal{R}$  and its scalings.

Then  $\varepsilon = o(S)$  when we rescale the region  $\mathcal{R}$  around the point  $(x_1, x_2)$ .

# Density recalibration

Any affine transformation changes a square grid into parallelograms.



As a result, the area of the regions is different on the left and on the right:

$$p(x_1, x_2) \cdot S_1 \approx q(y_1, y_2) \cdot S_2 \quad \implies \quad q(y_1, y_2) = \frac{S_1}{S_2} \cdot p(x_1, x_2)$$

Fortunately, the ratio between areas are constant over the entire plane!

## Density of two-variate normal distribution

The density of  $(x_1, x_2)$  pairs can be computed based on independence:

$$p(x_1, x_2) = p(x_1) \cdot p(x_2) = \frac{1}{2\pi} \cdot \exp\left(-\frac{x_1^2 + x_2^2}{2}\right) .$$

To estimate density  $q(y_1, y_2)$ , we must find the corresponding  $(x_1, x_2)$ :

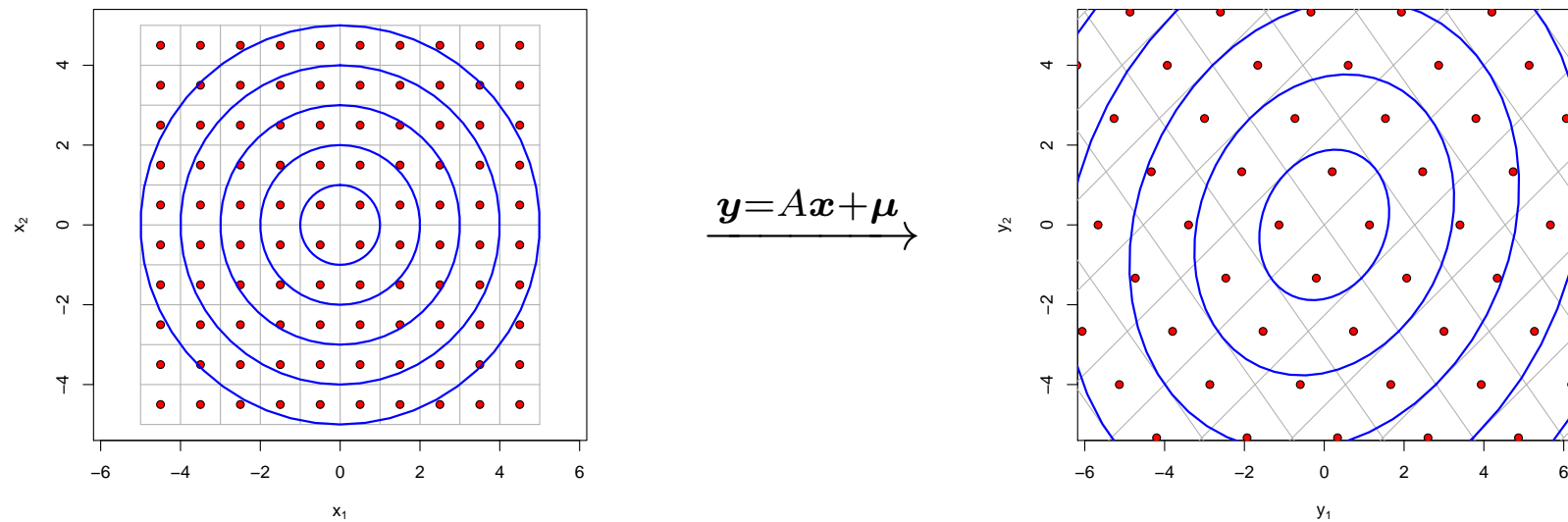
$$\mathbf{y} = A\mathbf{x} + \boldsymbol{\mu} \quad \Leftrightarrow \quad \mathbf{x} = A^{-1}(\mathbf{y} - \boldsymbol{\mu}) .$$

Thus we get

$$\begin{aligned} q(y_1, y_2) &= \frac{S_1}{S_2} \cdot \frac{1}{2\pi} \cdot \exp\left(-\frac{(\mathbf{y} - \boldsymbol{\mu})^T A^{-T} A^{-1}(\mathbf{y} - \boldsymbol{\mu})}{2}\right) \\ &= \frac{1}{\sqrt{\det(\Sigma)}} \cdot \frac{1}{2\pi} \cdot \exp\left(-\frac{(\mathbf{y} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{y} - \boldsymbol{\mu})}{2}\right) . \end{aligned}$$



## Illustrative example



- ▷ Affine transformation changes the square grid into parallelograms.
- ▷ Affine transformation changes circular equiprobability lines into ellipses.
- ▷ The axes of the ellipses may intersect with the sides of parallelograms.

## Generalisation to multivariate case

If observed quantities  $\mathbf{y}$  are generated by applying the affine transformation

$$\mathbf{y} = A\mathbf{x} + \boldsymbol{\mu} \quad \Leftrightarrow \quad \mathbf{x} = A^{-1}(\mathbf{y} - \boldsymbol{\mu})$$

to the *independent source signals*  $x_1, \dots, x_n \sim \mathcal{N}(0, 1)$ , then the resulting distribution is *a multivariate normal distribution* with the density:

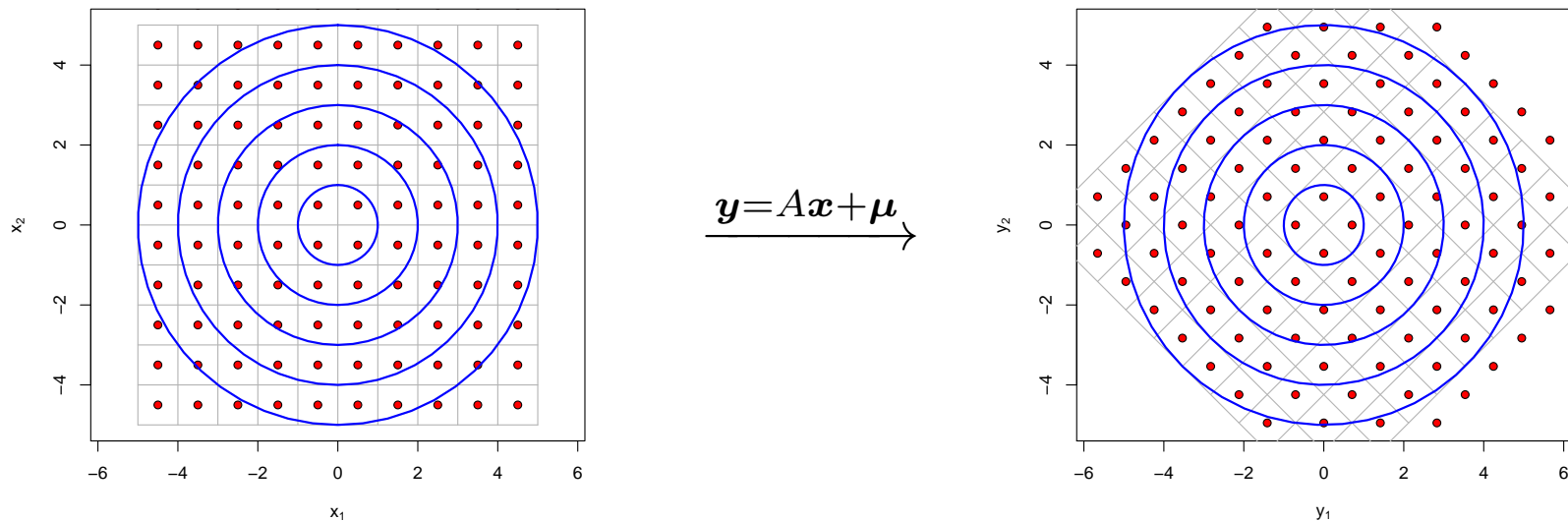
$$q(\mathbf{y}) = \frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{\sqrt{\det(\Sigma)}} \cdot \exp\left(-\frac{(\mathbf{y} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu})}{2}\right)$$

where  $\Sigma^{-1} = A^{-T} A^{-1}$  is *a positively definite symmetric matrix*.

# Distribution reconstruction task

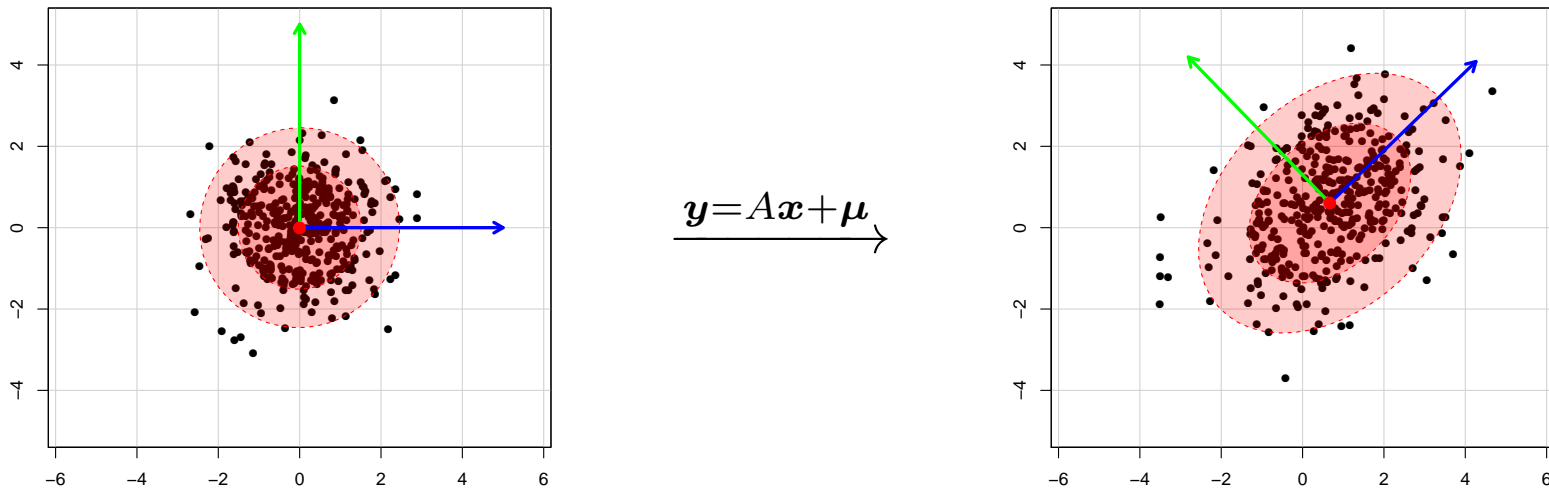
**Original goal.** Given the set of observations  $\mathbf{y}_1, \dots, \mathbf{y}_m$  determine the affine transformation  $\mathbf{y} = A\mathbf{x} + \boldsymbol{\mu}$  and original source signals  $\mathbf{x}_1, \dots, \mathbf{x}_m$ .

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



## Simplified distribution reconstruction task

**Achievable goal.** Given the set of observations  $\mathbf{y}_1, \dots, \mathbf{y}_m$  determine the affine transformation by fixing the centre and axis of the ellipsoid.



- ▷ We need to find the origin and semi-axes  $\mathbf{a}_1, \dots, \mathbf{a}_n$  of the ellipsoid.
- ▷ Unit vectors  $\mathbf{e}_1, \dots, \mathbf{e}_n$  are mapped to semi-axes  $\mathbf{a}_1, \dots, \mathbf{a}_n$  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?

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

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?

$$\mathbf{E}(A\mathbf{x} + \boldsymbol{\mu}) = \mathbf{E}(A\mathbf{x}) + \mathbf{E}(\boldsymbol{\mu}) = \boldsymbol{\mu} .$$

## 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:

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

## Principal components as eigenvectors

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

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

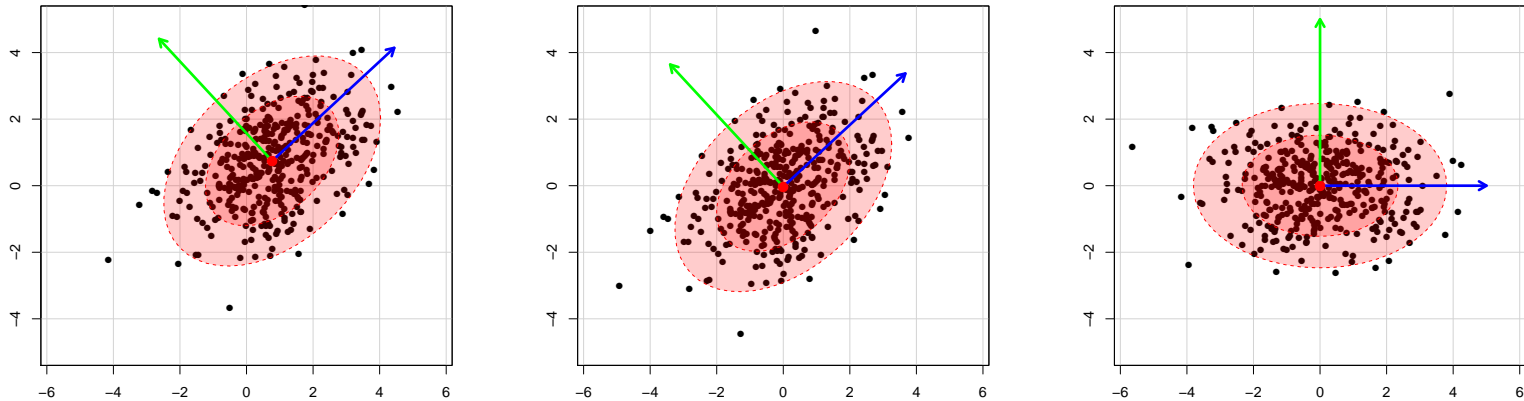
**Fact.** If  $n \times n$  matrix is symmetric and positively definite then there exists  $n$  orthogonal eigenvectors  $\mathbf{w}_1, \dots, \mathbf{w}_n$  with *eigenvalues*  $\lambda_1 \geq \dots \geq \lambda_n > 0$ .

**Corollary.** Principal components corresponding to observations  $\mathbf{y}_1, \dots, \mathbf{y}_m$  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):

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

## 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

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

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  $\boldsymbol{\mu}$ :

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

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

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

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

## Maximum likelihood estimates for parameters

The shift must be the mean of all observations

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

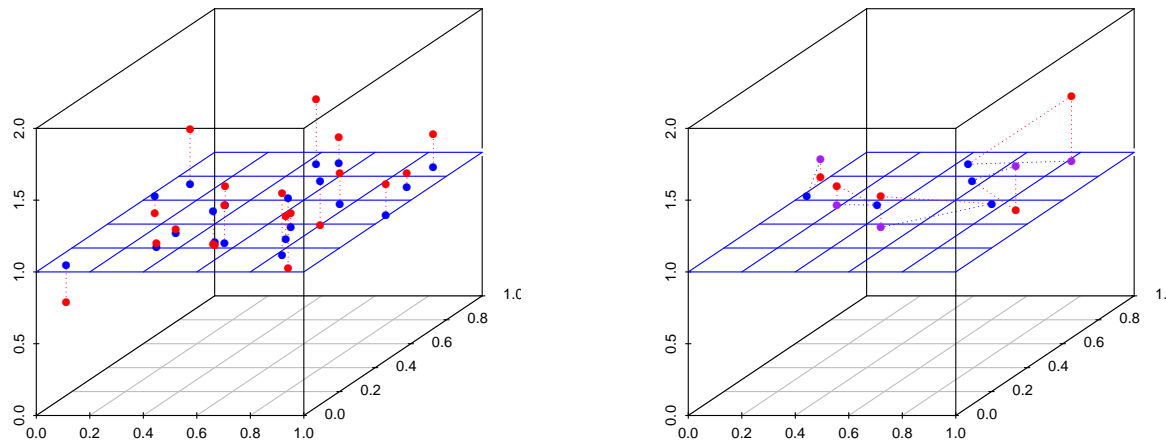
The covariance matrix

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

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

# 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

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

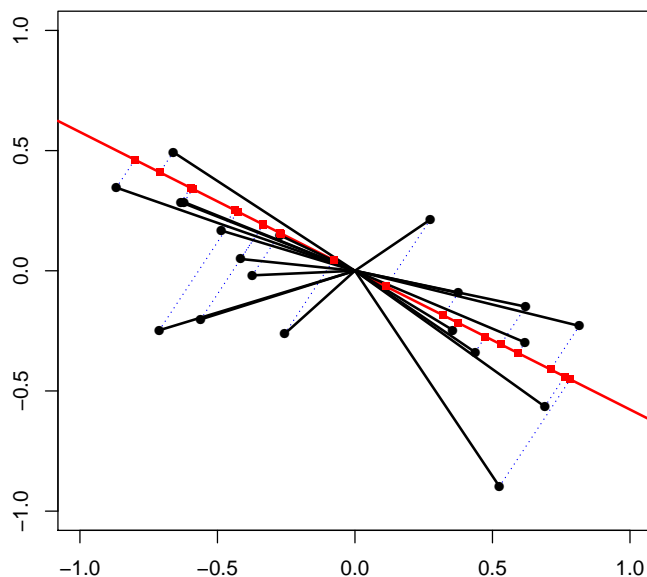
where  $d_i$  is the distance between the plane  $\mathcal{H}$  and the point  $\mathbf{y}_i$ . Thus

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

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.

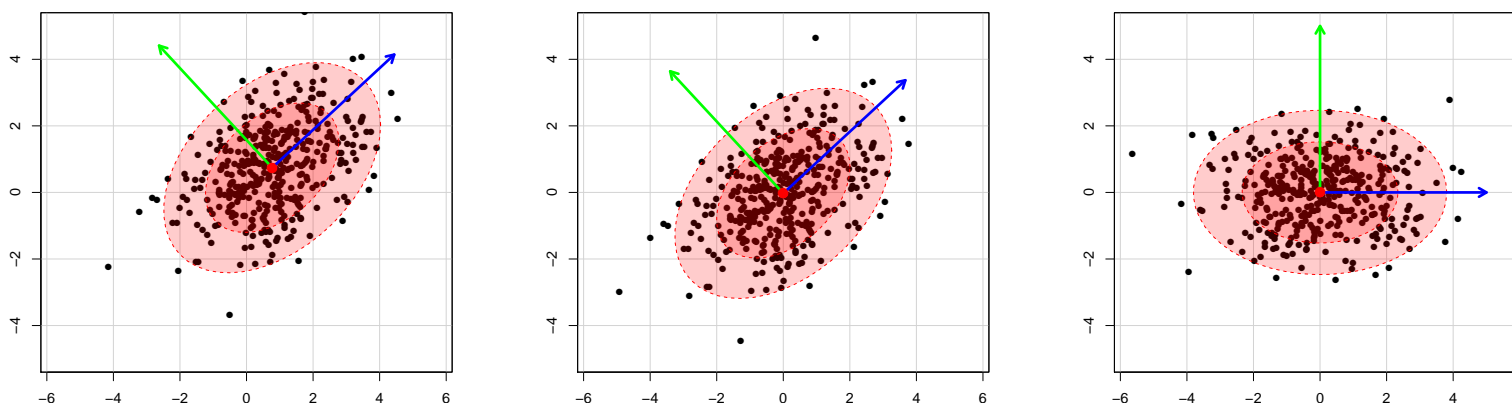


**Corollary.** PCA chooses directions  $w_1, \dots, w_n$  such that the sum of distance squares from the hyperplane formed by  $w_1, \dots, w_k$  is minimal.



## 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.

# Going beyond PCA

Weighted Principal Component Analysis:

- ▷ Sometimes data contains potential outliers.
- ▷ Sometimes we can assign reliability scores to the data points.

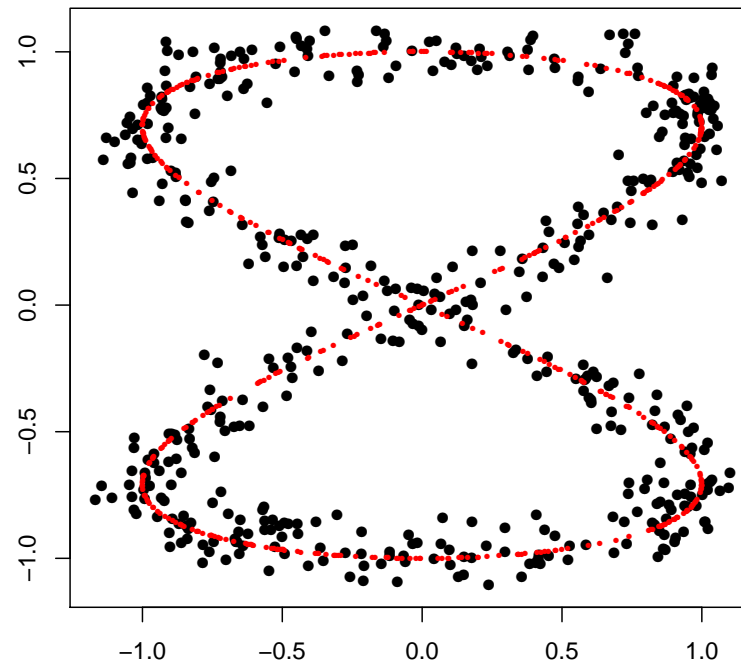
Principal curves and manifolds

- ▷ The original data might be on a low dimensional manifold.
- ▷ The observed data is corrupted by additive white gaussian noise.
- ▷ The task is to reconstruct the manifold and ML estimate for the data.

Independent Component Analysis

- ▷ What if the source components are non-gaussian?
- ▷ Then the reconstruction is possible up to scaling!

# Principal curves and manifolds

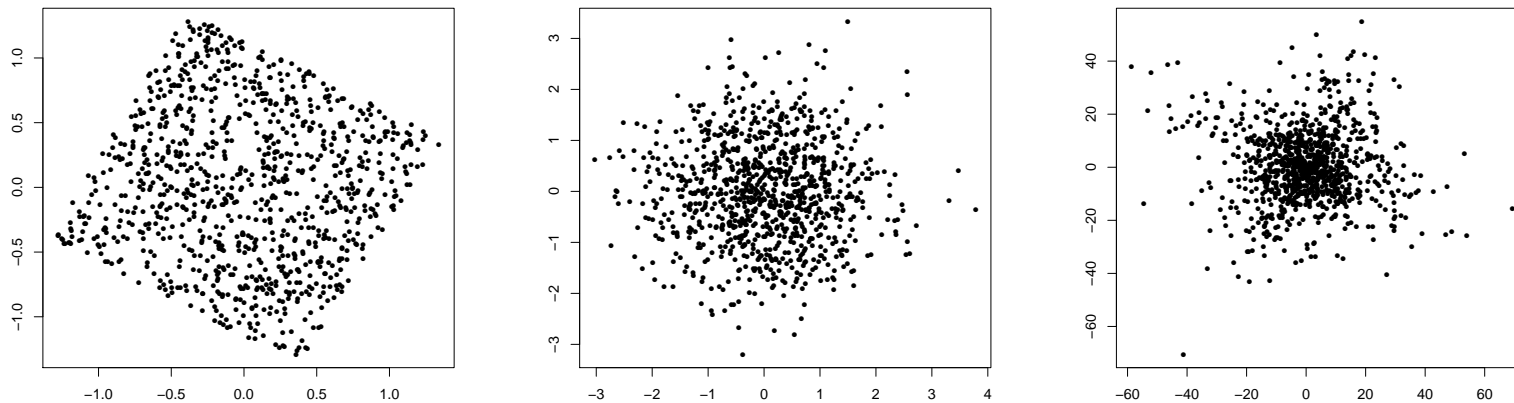


Reconstruction of the underlying curve is much more difficult.

- ▷ We must fix a curve parametrisation
- ▷ The task is different from regression since we have only outputs.

# Independent Component Analysis

Assume that the components of the source data  $\mathbf{x}_1, \dots, \mathbf{x}_m$  are independent but an unknown affine transformation  $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\mu}$  disturbs observations.



It is possible to recover the translation and rotation only if independent components are sufficiently different from the normal distribution.