

# Statistics for Machine Learning

## Foundational Statistical Concepts

## Foundational Statistical Concepts

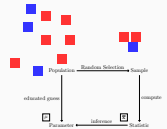


Figure 1: Statistical inference means inferring something about a population using information from samples [Hao22]

Population includes all of the elements from a set of data. Sample consists of one or more observations from the population.

Parameter Characteristic of a distribution describing a population, such as the mean or standard deviation of a normal distribution. Often notated using Greek letters.

Statistic A numerical value that represents a property of a random sample.

Examples of statistics are

- the mean value of the sample data.
- the range of the sample data.
- deviation of the data from the sample mean.

The covariance matrix of a multidimensional variables is filled with individual variables, consider the two-dimensional case:

$$\Sigma = \begin{pmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{yx} & \sigma_{yy} \end{pmatrix} \quad (12)$$

$$\Sigma = (\mathbf{x}_i - \hat{\mu}_g)(\mathbf{x}_i - \hat{\mu}_g)^T \quad (13)$$

$$\Sigma = (\mathbf{x}_i - \hat{\mu}_g)(\mathbf{x}_i - \hat{\mu}_g)^T \quad (14)$$

$$= \begin{pmatrix} (\mathbf{x}_1 - \mu_1)^2 & (\mathbf{x}_1 - \mu_1)(\mathbf{x}_2 - \mu_2) & \dots & (\mathbf{x}_1 - \mu_1)(\mathbf{x}_n - \mu_n) \\ (\mathbf{x}_2 - \mu_2)(\mathbf{x}_1 - \mu_1) & (\mathbf{x}_2 - \mu_2)^2 & \dots & (\mathbf{x}_2 - \mu_2)(\mathbf{x}_n - \mu_n) \\ \vdots & \vdots & \ddots & \vdots \\ (\mathbf{x}_n - \mu_n)(\mathbf{x}_1 - \mu_1) & (\mathbf{x}_n - \mu_n)^2 & \dots & (\mathbf{x}_n - \mu_n)(\mathbf{x}_n - \mu_n) \end{pmatrix} \quad (15)$$

$$= \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{12} & \sigma_2^2 & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix} \quad (16)$$

## Statistics for Machine Learning

## └ Gaussian mixture models

## └ Gaussian mixture models

A Gaussian mixture model has the density [McN16]

$$f(\mathbf{x}|\theta) = \sum_{g=1}^G \pi_g \phi(\mathbf{x}|\mu_g, \Sigma_g). \quad (20)$$

With the normal distribution  $\phi$  defined as before,  $\pi_g$  denotes the global probability with which a data value could originate from gaussian  $g$ . The  $g$ 's number the gaussians, and  $G$  is the total number of Gaussians in the mix. We will use two.  $\phi$  denotes the parameters  $\mu_g$  and  $\Sigma_g$ .

Typically we want as many  $g$  as we have classes in the data. I.e. one for healthy and one for diabetic. The data vectors are  $p$  dimensional  $\mathbf{x} \in \mathbb{R}^p$

## Statistics for Machine Learning

## └ Gaussian mixture models

## └ Likelihood

Likelihood models the probability of data originating from a distribution as a function of the parameters. The gaussian case is modelled by [McN16]

$$\mathcal{L}_n(\theta) = \prod_{i=1}^n \prod_{g=1}^G [\rho_g \phi(\mathbf{x}_i | \mu_g, \Sigma_g)]^{z_{ig}} \quad (21)$$

We want to maximize the likelihood.

In other words, we want to transform the balls in such a way, that they explain the points as plausible as possible.

To maximize  $\phi$  it needs to sit on top of the points it labels. When a gaussian sits on top of many points it's  $\rho_g$  should be large. Finally, when this works well we want a big weight from  $z_{ig}$ .

## Statistics for Machine Learning

## └ Gaussian mixture models

## └ Clustering using a GMM

After guessing an initial choice for all  $\hat{\mu}_k$  and  $\hat{\Sigma}_k$  [McN16],

$$z_{ig} = \frac{\phi_g(\mathbf{x}_i | \hat{\mu}_g, \hat{\Sigma}_g)}{\sum_{k=1}^K \phi_k(\mathbf{x}_i | \hat{\mu}_k, \hat{\Sigma}_k)} \quad (23)$$

tells us the probability with which point  $\mathbf{x}_i$  came from gaussian  $g$ . It creates an association between the data points and the Gaussians. Numerically evaluation results in a matrix  $\mathbf{Z} \in \mathbb{R}^{G \times n}$ . Use it's output to select the points which belong to each class.

The  $z_{ig}$  are the true labels,  $\hat{z}_{ig}$  is our estimation. The  $\hat{z}_{ig}$  are the expected value of the complete data log-likelihood. Why?  $\phi$  is a pdf. A pdf can be interpreted as providing a relative likelihood that the value of the random variable would be equal to that sample <sup>1</sup>. We ask for all gaussians and every point and normalize.