

# Statistics for Machine Learning

## Foundational Statistical Concepts

## Foundational Statistical Concepts

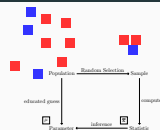


Figure 1: Statistical inference means inferring something about a population using information from samples [Hao22]

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.

Statistics: 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.

Machine learning and statistics are both concerned with the construction of models to explain data. In statistics data is observed and we try to figure out a process that explains the data. In machine learning we look for a best-fit [DFO20].

## Statistics for Machine Learning

## └ Foundational Statistical Concepts

## └ Probability and random variables [DFO20]

**Sample space**  $\Omega$ 

The sample space contains all possible outcomes of an experiment. A coin toss, for example, can have two outcomes: heads (h) or tails (t). Which leads to the set  $\{h, t\}$ . Two successive tosses generate the larger space  $\{hh, tt, ht, th\}$ .

**Event space**  $\mathcal{A}$ 

A set of events, an event is a set of outcomes from the sample space.

**Probability**  $P$ 

With each event  $A$  we associate a number  $P(A)$ . This number measures the probability that the event will occur.

Lets consider a hypothetical bag. The bag contains dollar \$ and euro € coins. Of ten coins in the bag four are euro coins. If we draw twice the probability space will be the set  $\{\text{€€}, \text{\$€}, \text{€\$}, \text{\$\$}\}$ . We can now construct a discrete function for the random variable  $X$ .  $X(\text{€}, \text{€}) = 2$ ,  $X(\text{€}, \$) = 1$ ,  $X(\$, \text{€}) = 1$ ,  $X(\$, \$) = 0$

$$P(X = 2) = P((\text{€}, \text{€})) = P(\text{€}) \cdot P(\text{€}) = 0.4 \cdot 0.4 = 0.16 \quad (1)$$

$$\begin{aligned} P(X = 1) &= P((\text{€}, \$) \cup P((\$, \text{€})) = P((\text{€}, \$) + P((\$, \text{€})) \\ &= 0.4 \cdot (1 - 0.4) + (1 - 0.4) \cdot 0.4 = 0.48 \quad (2) \end{aligned}$$

$$P(X = 0) = P(\$) \cdot P(\$) = 0.6 \cdot 0.6 = 0.36 \quad (3)$$

# Statistics for Machine Learning

## Foundational Statistical Concepts

### Marginal and conditional probability

#### Marginal and conditional probability

We can compute marginal probabilities by summing rows or columns.

$$P(X = x_i) = \frac{c_i}{N} = \frac{\sum_{j=1}^J a_{ji}}{N} \quad (5)$$

$$P(Y = y_j) = \frac{r_j}{N} = \frac{\sum_{i=1}^I a_{ij}}{N} \quad (6)$$

The marginal probabilities allow us to define conditional probability:

$$P(Y = y_j | X = x_i) = \frac{a_{ji}}{c_i} \quad (7)$$

$$P(X = x_i | Y = y_j) = \frac{a_{ji}}{r_j} \quad (8)$$

Made up data:

		$\leq 1 \text{ h}$	$> 1 \text{ h}$
$y_1$	Bonn	0.24	0.08
$y_2$	Cologne	0.04	0.32
$y_3$	Siegburg	0.16	0.16
		$x_1$	$x_2$

The table entries contain normalized frequencies. Row and column frequency sums are the so-called marginal probabilities.

$$\Rightarrow P(X = x_1) = 0.44. \quad (9)$$

$$\Rightarrow P(Y = y_2) = 0.36. \quad (10)$$

Probability for more than one hour from cologne:

$$P(X = x_2 | Y = y_2) = 0.32 / (0.32 + 0.04) = 0.89 \quad (11)$$

## Statistics for Machine Learning

## └ Foundational Statistical Concepts

## └ Bayes Law [DFO20]

Sometimes, we have no direct way of observing a property. We are forced to infer knowledge indirectly. In such cases, Bayes law helps. Bayes states

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} \quad (26)$$

The law is a consequence of our ability to factorize distributions as  $p(x, y) = p(x|y)p(y)$ . If we can't observe  $x$  directly, we may have expectations of its distribution  $p(x)$ , and the likelihood  $p(y|x)$ . Bayes allows us to find a posterior  $p(x|y)$  given evidence  $p(y)$ .

Say 50 in 100k people of a population have a given illness.

We have  $P(S) = 0.0005$ , and  $P(H) = 1 - (50/100000) = 0.9995$ . We have a test that detects the disease with an accuracy of 98%. In other words,  $P(T|S) = 0.98$ . Unfortunately, it also yields a positive result for 1% of healthy people  $P(T|H) = 0.01$ . What happens if we use the test to look for the disease in the general population?

$$P(S|T) = \frac{P(T|S)P(S)}{P(T|S)P(S) + P(T|H)P(H)} \quad (27)$$

$$= \frac{0.98 \cdot 0.0005}{0.98 \cdot 0.0005 + 0.01 \cdot 0.9995} = 0.05 \quad (28)$$

In this case, testing is probability not a great idea. Note, total probability for exclusive events:  $P(T) = P(T|S)P(S) + P(T|H)P(H)$  if  $H$  is not  $S$ .

## Statistics for Machine Learning

## └ Gaussian mixture models

## └ Gaussian mixture models

A Gaussian mixture model has the density [McN16]

$$f(\mathbf{x}|\theta) = \sum_{g=1}^G \pi_g \phi(\mathbf{x}|\mu_g, \Sigma_g). \quad (35)$$

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 Gaussian function. Parameters  $\mu_g$  and  $\Sigma_g$  are mean vector and covariance matrix.

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  $\mathbf{x} \in \mathbb{R}^p$ . Sampling  $\phi(\mathbf{x})$  tells us how likely it was to see the point we have. Big values mean it was likely small mean it was not.

## 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]

$$\mathcal{L}_n(\theta) = \prod_{i=1}^n \prod_{g=1}^G [\rho_g \phi(\mathbf{x}_i | \mu_g, \Sigma_g)]^{z_{ig}} \quad (36)$$

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  $\hat{\mu}_g$  and  $\hat{\Sigma}_g$  [McN16],

$$z_{ig} = \frac{\phi_g(x_i | \hat{\mu}_g, \hat{\Sigma}_g)}{\sum_{k=1}^K \phi_k(x_i | \hat{\mu}_k, \hat{\Sigma}_k)} \quad (38)$$

tells us the probability with which point  $x_i$  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 the maxima in 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.