

Advanced Probabilistic Machine Learning and Applications

Caterina De Bacco and Isabel Valera

1 Notations

We try to be consistent with notations throughout the classes.

Here we give a summary of the main used. These are the same as in the [Bishop \(2006\)](#).

- $\mathbf{x} = (x_1, \dots, x_D)$: a column D -dim vector, i.e. lower case bold Roman letter
- \mathbf{x}^T : the row vector transposed of \mathbf{x} .
- $\mathbf{x}_1, \dots, \mathbf{x}_N$: N samples of a D -dim vector
- \mathbf{X} : data matrix with n -th row the row vector \mathbf{x}_n^T ; i.e. the n, i entry of \mathbf{X} is the i -th entry of the n -th observation \mathbf{x}_n .
- \mathbf{x} : data matrix for one-dimensional variables, i.e. this is a column vector whose n -th element is x_n .
- $\mathbf{x} \neq \mathbf{x}$: we use two different typefaces. \mathbf{x} denotes a N -dim vector (i.e. N samples of one-dim variables); \mathbf{x} denotes a D -dim vector (i.e. 1 sample of a D -dim variable).
- $\mathbb{E}_x[f(x, y)]$: expected value with respect to the random variable x of the function $f(x, y)$.
- $\mathbb{E}[x]$: expected value of x when there is no ambiguity as to which variable is being averaged over.

2 One variable

Suppose we have one random variable X , one-dimensional.

Assume that this is distributed according to some probability distribution, for instance Gaussian:

$$X \sim \mathcal{N}(\mu, \sigma^2) \quad , \quad (1)$$

where μ and σ^2 are the parameters.

Now, we observe N samples of this variable x_1, \dots, x_N and assume they are independent and identically distributed as a Gaussian.

Question: given this sample, how do we estimate μ and σ ?

2.1 Maximum Likelihood Estimation

One possibility is to find the pair (μ, σ) that maximizes the likelihood of observing the data. In general, the likelihood of the data \mathbf{x} for a probability distribution of parameters θ is defined as $p(\mathbf{x}|\theta)$.

Often it is mathematically more convenient to consider the logarithm of this quantity, also called log-likelihood. This is also useful numerically, as there is less risk of underflow when taking a sum

instead of a product of a very small numbers. We denote this as $\mathcal{L}(\mu, \sigma^2)$ and we omit from the notation the explicit dependence on \mathbf{x} .

Given the logarithm is monotonically increasing, the maximum of these two coincide. Then we have:

$$\mu_{MLE}, \sigma_{MLE}^2 = \arg \max_{\mu, \sigma^2} \{\mathcal{L}(\mu, \sigma^2)\} \quad (2)$$

We have N independent observations, denoted as a vector \mathbf{x} , so we can sum the log-likelihood of each of them:

$$\mathcal{L}(\mu, \sigma^2) = p(\mathbf{x}|\mu, \sigma^2) = \sum_i^N \log(\mathcal{N}(x_i|\mu, \sigma^2)) \quad (3)$$

$$= \sum_i^N \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}\right)\right) \quad (4)$$

$$= -\frac{1}{2} \sum_i^N \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_i^N (x_i - \mu)^2 + \text{const} \quad (5)$$

where the term *const* does not depend on the parameters. Now we want to extract the maximum of that expressions:

$$\begin{cases} \frac{\partial \mathcal{L}(\mu, \sigma^2)}{\partial \mu} \equiv 0 \\ \frac{\partial \mathcal{L}(\mu, \sigma^2)}{\partial \sigma^2} \equiv 0 \end{cases} \quad (6)$$

After some algebra, we can find the analytical expressions of the Maximum Likelihood estimation:

$$\mu_{MLE} = \frac{1}{N} \sum_{n=1}^N x_n \quad (7)$$

$$\sigma_{MLE}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2 \quad (8)$$

These correspond to the sample *empirical* mean and variance respectively.

Obs1: this was easy, because we could split the terms inside the sum as the samples were *independent* (and identically) distributed.

Obs2: from this estimation, we do not have idea of any measure of error or uncertainty in estimating the parameters, i.e. μ_{MLE} and σ_{MLE}^2 are point-estimates but no extra information was obtained from the calculations.

2.2 Posterior and Maximum a Posteriori estimates (MAP)

To address Obs2 we consider also μ, σ^2 as random variables, each distributed according to some distribution.

For instance, assume that:

$$p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2) \quad , \quad (9)$$

and keep σ^2 fixed. $p(\mu)$ is called the *Prior* of the parameter μ . This distribution should be chosen in order to use some information (if any) that we know about the parameters. For instance, based on domain knowledge or by looking at the observed data; example: if the samples are all positive and discrete, perhaps it is more appropriate to set the prior as Poisson instead of Gaussian.

With this notion in mind and using Bayes' rule, we can now define the *Posterior* of μ as:

$$p(\mu|\mathbf{x}) = \frac{p(\mathbf{x}|\mu)p(\mu)}{p(\mathbf{x})} \propto p(\mathbf{x}|\mu)p(\mu) \quad , \quad (10)$$

where we made explicit the proportionality of the left-most term because in general the marginal likelihood $p(\mathbf{x})$ is not accessible (it involves complex calculations of integrals). Hence, in practice, one estimates $p(\mu|\mathbf{x})$ by writing $p(\mathbf{x}|\mu)p(\mu)$ and attempting to then normalize the resulting expression. This is not generally possible analytically, but in this Gaussian example everything is doable. The result, after some algebra, is that the posterior is also Gaussian distributed, $p(\mu|\mathbf{x}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$, with parameters:

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{MLE} \quad (11)$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \quad (12)$$

Obs1: μ_N is a combination of the prior's mean μ_0 and the MLE's estimate μ_{MLE} . In particular, when $N = 0$ we recover the prior, i.e. $\mu_N = \mu_0$; on the opposite scenario $N \rightarrow \infty$, we get $\mu_N \rightarrow \mu_{MLE}$, i.e. the MLE estimate and the mean of the Posterior coincide.

Obs2: $\frac{1}{\sigma_N^2}$, also called *precision*, is additive; which means that for $N \rightarrow \infty$ the σ_N^2 decreases, i.e. the more data we observe, the more the posterior variance decreases and the posterior distribution becomes peaked around μ_{MLE} .

Obs3: for *finite* N , if we have $\sigma_0^2 \rightarrow \infty$ (prior with infinite variance, also called *non informative prior*), then also in this case $\mu_N \rightarrow \mu_{MLE}$, but this time with $\sigma_N^2 = \sigma^2/N$.

These analysis and results can be observed by playing around with the [jupyter notebook](#)¹ included in L1's folder.

Being able to characterize the Posterior distribution allows not only to measure the variance of the estimated parameters (i.e. σ_N^2), but one can perform other tasks such as sampling from that distribution.

However, if one is interested to get a point-estimate from the Posterior, one common choice is taking the value at the peak of this distribution (also called *mode*), the Maximum a posteriori estimate (MAP):

$$\mu_{MAP} = \arg \max_{\mu} p(\mu|\mathbf{x}) \quad . \quad (13)$$

For the example considered above of a Gaussian Posterior, we have $\mu_{MAP} \equiv \mu_N$. But this is not the case in general.

Obs1: if the Prior is uniform, then $\mu_{MAP} \equiv \mu_N$.

Obs2: given that the (usually hard-to-compute) term $p(\mathbf{x})$ in the denominator of the exact posterior does not depend on μ , to calculate μ_{MAP} one can simply consider the maximization over $p(\mathbf{x}|\mu)p(\mu)$. This quantity is easier to access because does not require the hard task of normalizing as for the Posterior.

All the results of this section can be generalized for D -dimensional variables.

3 More than one variable

Suppose now that we have more than one random variable. For simplicity, we consider the case of having two of them X_1 and X_2 and both of them one-dimensional.

¹https://github.com/APMLA-2021/APMLA-2021_material/blob/master/L1/L1.ipynb

Also in this case, assume that these are distributed according to some probability distribution, for instance Gaussian:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{bmatrix}\right), \quad (14)$$

this is also called *Multivariate Gaussian* with parameters $\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ and $\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{bmatrix}$.

Σ is called *covariance matrix* and to be valid it has to be *positive semi-definite* (all its eigenvalues should be nonnegative) and *symmetric*.

In general, X_1 and X_2 are not *independent*, the way they depend to each other is regulated by the cross terms of the covariance matrix σ_{12}^2 .

3.1 Conditional and marginal distributions

When two variables are correlated with each other, a relevant question is how does one vary when the other takes a given fixed value. This answer is provided by the *conditional* distribution of x_1 given x_2 :

$$p(x_1|x_2) = \frac{p(x_1, x_2)}{p(x_2)} \quad (15)$$

The term in the denominator is the marginal probability of x_2 and is defined as:

$$p(x_2) = \int p(x_1, x_2) dx_1 \quad (16)$$

Given that this integral is not always easy to calculate, one can instead derive an expression for the conditional distribution rewriting as:

$$p(x_1|x_2) = \frac{p(x_1, x_2)}{Z_1} \quad (17)$$

where Z_1 is a normalization constant that one should find to make $p(x_1|x_2)$ a valid probability distribution. Usually finding Z_1 is easier than calculating the integral in (16). This can be done rewriting the Multivariate distribution isolating the terms depending on x_1 :

$$(\mathbf{x} - \boldsymbol{\mu})^T \Sigma (\mathbf{x} - \boldsymbol{\mu}) = \frac{1}{\sigma_1^2} \left[x_1 - \left(\mu_1 + (x_2 - \mu_2) \frac{\sigma_{12}^2}{\sigma_2^2} \right) \right]^2 \quad (18)$$

This means that $p(x_1|x_2)$ is also a Gaussian, i.e. $\mathcal{N}(\mu_{x_1|x_2}, \sigma_{x_1|x_2}^2)$, with parameters:

$$\mu_{x_1|x_2} = \mu_1 + (x_2 - \mu_2) \frac{\sigma_{12}^2}{\sigma_2^2} \quad (19)$$

$$\sigma_{x_1|x_2}^2 = \sigma_1^2 \quad (20)$$

These calculations can be done explicitly, this is a good practice exercise for handling Gaussian distribution (e.g. completing the square). Similar thinking can be applied to derive the expression of the marginal $p(x_1)$, but in this case calculations are much more tedious!

The result is that also $p(x_1)$ is a Gaussian distribution, namely $\mathcal{N}(x_1|\mu_1, \sigma_1^2)$.

You can find interactive plots visualizing these distributions inside the [jupyter notebook](https://github.com/APMLA-2021/APMLA-2021_material/blob/master/L1/L1.ipynb)² included in L1's folder.

²https://github.com/APMLA-2021/APMLA-2021_material/blob/master/L1/L1.ipynb

Obs1: the number of parameters for the 2-variable example is 2 for the mean plus 4 for the covariance matrix. More generally, it grows as $N + N(N + 1)/2$, i.e. quadratically in the number of variables (and we are only considering one-dimensional ones!).

Obs2: the Multivariate Gaussian has a single maximum (it is *unimodal*), therefore it may be too limited to capture distributions with many maxima.

Obs3: if x_1 and x_2 come from different types of distributions, e.g. one is Gaussian and the other is Gamma, and they are correlated, there is no such thing as a Multivariate distribution with mixed types of variables. In these cases, how do we capture correlations? As an example, think of the variables (*salary, balance*). These are two correlated variables, if you increase your salary, then it is more likely that also your bank account's balance increases. However, one is always positive (salary), the other can go negative (balance).

These three observations show that a Multivariate Gaussian suffers to major problems: it might not capture complex situations like non-unimodal or mixed-type distributions. And it has too many parameters, growing quadratic in the number of variables.

4 Latent variables and Generative models

Question: how do we capture correlations in these complex scenarios?

The answer is given by *latent* variables. The purpose of these is to capture correlations of a complicated distributions via simpler *conditional* distributions. This is formalized by De Finetti's theorem as follows. Consider two D -dimensional variables $\mathbf{x}_1, \mathbf{x}_2$ and a K -dimensional variable \mathbf{z} (our *latent* variable). The we have:

$$p(\mathbf{x}_1, \mathbf{x}_2) = \int p(\mathbf{x}_1|\mathbf{z}) p(\mathbf{x}_2|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \quad (21)$$

The quantity in the right-hand-side is a joint distribution, in general this can be complicated as generally $p(\mathbf{x}_1, \mathbf{x}_2) \neq p(\mathbf{x}_1)p(\mathbf{x}_2)$. The terms inside the integral are instead *factorized*, thanks to the introduction of the auxiliary latent variable \mathbf{z} . This is the key point of the theorem.

We can interpret $p(\mathbf{x}_1, \mathbf{x}_2, \mathbf{z})$ as a model expressing the process that generated the observed data, i.e. we call this a *generative model*.

Obs1: even though we can factorize $p(\mathbf{x}_1, \mathbf{x}_2, \mathbf{z}) = p(\mathbf{x}_1|\mathbf{z}) p(\mathbf{x}_2|\mathbf{z}) p(\mathbf{z})$, we can still capture complex correlations between $\mathbf{x}_1, \mathbf{x}_2$ because when we marginalize out over \mathbf{z} as in (21), we get a non-factorized joint $p(\mathbf{x}_1, \mathbf{x}_2)$.

4.1 Example: Gaussian Factor Analysis

An example of this is found in the following model (see Chapter 12 of (Murphy, 2012) for further details), where for simplicity we consider one-dimensional variables:

$$p(x_1|z) = \mathcal{N}(x_1|w_1 z, \sigma_1^2) \quad (22)$$

$$p(x_2|z) = \mathcal{N}(x_2|w_2 z, \sigma_2^2) \quad (23)$$

$$p(z) = \mathcal{N}(z|\mu_0, \sigma_0^2) \quad (24)$$

Let's calculate the joint $p(x_1, x_2) = \int \mathcal{N}(x_1|w_1 z, \sigma_1^2) \mathcal{N}(x_2|w_2 z, \sigma_2^2) \mathcal{N}(z|\mu_0, \sigma_0^2) dz$. Without loss of generality, we can set $\sigma_0^2 = 1$. Rewriting everything explicitly, we obtain that the marginal distribution of $(x_1, x_2)^T$ is a Multivariate Gaussian distribution with mean and covariance:

$$\boldsymbol{\mu} = (w_1\mu_0, w_2\mu_0)^T \quad (25)$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} w_1^2 & w_1w_2 \\ w_1w_2 & w_2^2 \end{bmatrix} + \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} . \quad (26)$$

The expression for $\boldsymbol{\Sigma}$ highlights that we now have correlations between the two variables, thanks to the terms w_1w_2 in the non-diagonal entries.

Obs1: even though x_1 and x_2 were *conditionally* independent given \mathbf{z} , their joint distribution is not factorized anymore, they have indeed a non-diagonal covariance matrix $\boldsymbol{\Sigma}$.

Obs2: In general models, the integral in $p(x_1, \dots, x_D) = \int \prod_{d=1}^D p(x_d|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ does not have closed-form solution.

References

C. M. Bishop, *Pattern recognition and machine learning* (Springer, 2006).

K. P. Murphy, *Machine learning: a probabilistic perspective* (MIT press, 2012).