

LECTURE NOTES

---

# PROBABILISTIC GENERATIVE MODELS

---

This version: 16th January 2026

Latest version: [github.com/felipe-tobar/Probabilistic-Generative-Models](https://github.com/felipe-tobar/Probabilistic-Generative-Models)

Felipe Tobar  
Department of Mathematics  
Imperial College London

[f.tobar@imperial.ac.uk](mailto:f.tobar@imperial.ac.uk)  
<https://www.ma.ic.ac.uk/~ft410/>

# Preface

This notes are under development for 2026.

Felipe Tobar,  
London,  
January 2026.

# Contents

<b>1. Foundations</b>	<b>5</b>
1.1. Introduction . . . . .	5
1.2. Discriminative versus generative . . . . .	6
1.3. The pushforward measure . . . . .	7
1.4. Likelihood-based training . . . . .	9
1.5. A brief intro to information theory . . . . .	12
1.6. KL divergence as a way to compare distributions . . . . .	15
1.7. Concluding remarks . . . . .	19
1.7.1. Suggested exercises . . . . .	19
<b>2. Latent variable models and Expectation Maximisation</b>	<b>21</b>
2.1. Gaussian mixtures . . . . .	21
2.2. The Gaussian mixture model . . . . .	22
2.3. Expectation Maximisation for GMMs . . . . .	23
2.4. An interpretation of EM . . . . .	24
2.5. EM in its general form . . . . .	26
2.6. Concluding remarks . . . . .	27
2.6.1. Suggested exercises . . . . .	27
<b>References</b>	<b>29</b>



# Week 1

## Foundations

### 1.1 Introduction

A Probabilistic generative model (PGM), or simply, a GM, is a methodology for generating data. In general, the PGM is constructed and adjusted using observations with the aim to synthesise samples with the same statistical properties of the available observations. The *probabilistic* nature of the PGMs studied in this course follows from the fact that the available data will be considered to be realisations of an underlying random variable (RV), e.g.,  $X$ .

In this sense, the probability distribution of  $X$ , denoted  $P_X(x)$ , as well as its probability density function (pdf)  $p_X(x)$  will be central to the study of PGMs. In particular, targeting the pdf is one way of constructing PGMs, in which case the whole PGM paradigm becomes equivalent to the classical statistical modelling approach. However, as we will see in the course, enforcing the sought-after PGM to have an explicit parametric pdf can be rather restrictive.

Throughout the course, we will consider a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , with 3 RVs given by the following measurable maps:

$$\begin{array}{lll} X : \Omega \rightarrow \mathcal{X} & Y : \Omega \rightarrow \mathcal{Y} & Z : \Omega \rightarrow \mathcal{Z} \\ \text{(observed input)} & \text{(observed output)} & \text{(latent variable)} \end{array}$$

**Remark 1.1.**

Not all three RVs will be present in all our settings. For instance, in classification there is no justification for the latent variable  $Z$  (in general), while in clustering, there is no need for  $X$ . However, we build the general setup here for formality.

We equip  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$  with their Borel  $\sigma$ -algebras  $\mathcal{B}(\mathcal{X}), \mathcal{B}(\mathcal{Y}), \mathcal{B}(\mathcal{Z})$ , and consider the product measurable space

$$(\mathcal{X} \times \mathcal{Y} \times \mathcal{Z}, \mathcal{B}(\mathcal{X}) \otimes \mathcal{B}(\mathcal{Y}) \otimes \mathcal{B}(\mathcal{Z})).$$

Then the joint random variable  $(X, Y, Z) : \Omega \rightarrow \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$  is measurable with respect to  $\mathcal{F}$  and the above product  $\sigma$ -algebra.

Furthermore, we will assume that the law of  $(X, Y, Z)$  is absolutely continuous with respect to the product base measure on  $\mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$  (e.g. Lebesgue measure when  $\mathcal{X}, \mathcal{Y}, \mathcal{Z} \subseteq \mathbb{R}^d$ ), and hence admits a joint density  $p(x, y, z)$ . That is, for all  $A \in \mathcal{B}(\mathcal{X}), B \in \mathcal{B}(\mathcal{Y}), C \in$

$\mathcal{B}(\mathcal{Z})$ ,

$$\mathbb{P}(X \in A, Y \in B, Z \in C) = \int_{A \times B \times C} p(x, y, z) dx dy dz. \quad (1.1)$$

We will also assume that all relevant marginals and conditionals admit densities (with respect to the corresponding base measures), e.g.  $p(x, y)$ ,  $p(y|x)$ ,  $p(z|x, y)$ , etc.

## 1.2 Discriminative versus generative

The generative approach aims to characterise the complete generative distribution  $p(x, y, z)$ , whereas, in some application-specific cases, only the discriminative model, e.g.,  $p(y|x)$ , is needed. Let us examine the following example.

**Example 1.1** (Generative and discriminative views of binary classification).

Consider the binary classification problem, where, given an observation  $X = x$ , one needs to estimate its label  $Y$ . A discriminative model would directly parametrise  $\mathbb{P}(Y|X = x)$ . Since this is a binary classification case, without loss of generality, we can assume  $Y \in \{0, 1\}$ , and model  $\mathbb{P}(Y = 1|X = x)$ , since  $\mathbb{P}(Y = 0|X = x) = 1 - \mathbb{P}(Y = 1|X = x)$ . A model for this probability only needs to map  $x \in \mathbb{R}^d \rightarrow \mathbb{P}(Y = 1|X = x) \in [0, 1]$ . For instance, a reasonable candidate for this is

$$\mathbb{P}(Y = 1|X = x) = \frac{1}{1 + e^{-\theta^\top x}} \quad (1.2)$$

which is known as the logistic regression.

Conversely, in a generative approach, we aim to model the joint probability  $p(Y = y, X = x)$ . Modelling this distribution is not easy, however, observe that we can factorise it as

$$p(Y = y, X = x) = p(X = x|Y = y)p(Y = y), \quad (1.3)$$

which yields a pair of much more intuitive distributions to model:

- the class probability  $p(Y = y) = (\pi, 1 - \pi)$ ,  $\pi \in [0, 1]$ , and
- the class-conditional probability  $p(X = x|Y = y)$ , given by a two distributions over  $\mathcal{X}$ , denoted  $f_{\theta_0}$  and  $f_{\theta_1}$ .

Therefore, the classifier is

$$\begin{aligned} p(Y = 1|X = x) &= \frac{p(X = x|Y = 1)p(Y = 1)}{p(X = x)} \\ &= \frac{1}{1 + e^{-\log\left(\frac{\pi}{1-\pi} \frac{f_{\theta_1}(x)}{f_{\theta_0}(x)}\right)}}. \end{aligned} \quad (1.4)$$

**Exercise 1.1.**

Evaluate eq. (1.4) for  $f_{\theta_0} = \mathcal{N}(\mu_0, \Sigma_0)$  and  $f_{\theta_1} = \mathcal{N}(\mu_1, \Sigma_1)$ . What happens when  $\Sigma_0 = \Sigma_1$ ?

## 1.3 The pushforward measure

Despite the abundant collection of well-studied statistical models, in some scenarios we can construct a more ad hoc model by applying an appropriate transformation.

**Definition 1.1.**

Consider a RV  $X \in \mathcal{X}$  with measure  $P_X$ , and a nonlinear map  $T : \mathcal{X} \rightarrow \mathcal{X}$ . The measure of the transformed RV  $T(X)$  is known as the *push forward measure* of  $P_X$  through  $T$ , and it is denoted by  $T_{\#}P_X$

**Remark 1.2.**

The transformations considered in the course will be such that the pushforward measure has a density. With a slight abuse of notation, we will denote this density as  $T_{\#}p_X$ .

**Example 1.2** (Discrete pushforward).

Let  $X$  be a discrete random variable taking values in  $\{1, 2, 3\}$  with

$$\mathbb{P}(X = 1) = 0.2, \quad \mathbb{P}(X = 2) = 0.5, \quad \mathbb{P}(X = 3) = 0.3.$$

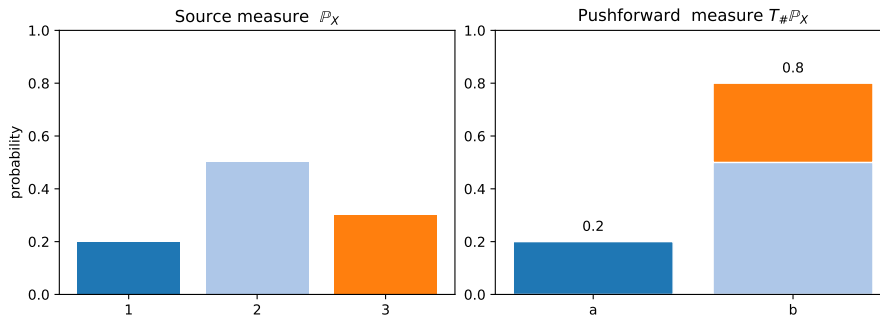
Define the map  $T : \{1, 2, 3\} \rightarrow \{a, b\}$  by

$$T(1) = a, \quad T(2) = b, \quad T(3) = b.$$

Then the pushforward  $T_{\#}\mathbb{P}$  satisfies

$$(T_{\#}\mathbb{P})(\{a\}) = 0.2, \quad (T_{\#}\mathbb{P})(\{b\}) = 0.8.$$

For an illustration see Fig. 1.1.



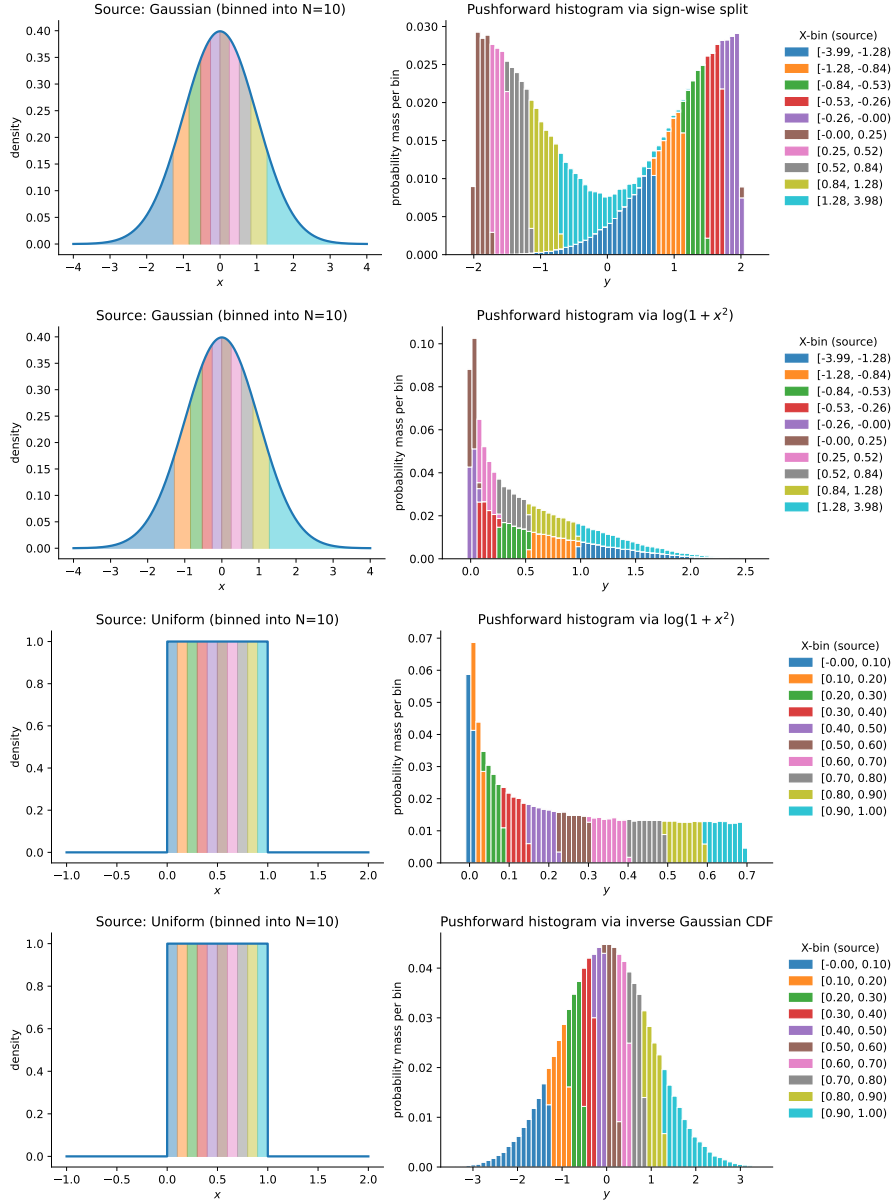
**Figure 1.1:** Source and pushforward distributions: discrete example.

**Example 1.3** (Continuous pushforward).

Let  $X \sim \mathcal{N}(0, 1)$  on  $\mathbb{R}$  and define  $T(x) = x^2$ . The pushforward  $T_{\#}\mathbb{P}$  is the law of  $Y = T(X)$ , supported on  $\mathbb{R}_+$ . Its density is given by

$$p_Y(y) = \frac{1}{\sqrt{2\pi y}} \exp\left(-\frac{y}{2}\right), \quad y > 0.$$

For an illustration of this case and other related examples, see Fig. 1.2.



**Figure 1.2:** Source and target distributions: continuous examples

In general, for arbitrary source distributions and maps it is difficult to compute the target density in closed form, at least in the continuous case. For the specific case of differentiable and invertible maps  $T$ , the following theorem gives a recipe to compute  $p_{T(X)}$ .

**Theorem 1.1** (Change of variable).

Consider two RVs  $X, Y \in \mathbb{R}^d$ , such that  $Y = T(X)$ , where  $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a  $C^1$  diffeomorphism. If  $X$  and  $Y$  have densities  $p_X$  and  $p_Y$  respectively, then

$$p_Y(y) = p_X(T^{-1}(y)) \left| \det \nabla_y T^{-1}(y) \right|, \quad (1.5)$$

where  $\nabla_y T^{-1}(y)$  is the Jacobian of the inverse map.

**Remark 1.3.**



Though the above result provides a closed-form expression for the pushforward measure only when  $T$  is a  $C^1$  diffeomorphism (continuously differentiable with an the inverse having the same property), we can transform a source RV  $X$  into a target RV  $T$  with any measurable map. This is because

$$(T_{\#}P_X)(A) = P_X(T^{-1}(A)), \quad \forall A \in \mathcal{B}(\mathcal{X}). \quad (1.6)$$

Though in general the pdf of  $T$  will not be available in closed form.

## 1.4 Likelihood-based training

Maximum likelihood (ML) is going to be the canonical methodology for training our PGMs, and, as we will see next, it will recover other forms of training criteria in particular cases.

Consider a PGM for the RV  $Y$ , with density  $p_{\theta}(y)$ , where  $\theta \in \Theta$  denotes the model parameter. Also, consider the realisations of  $Y$  given by  $y_1, y_2, \dots, y_N$ .

**Definition 1.2** (Likelihood function).

The likelihood of the parameter  $\theta$  is the function  $L : \Theta \rightarrow \mathbb{R}_+$  given by the probability density function of  $Y$  evaluated on the observations. That is,

$$L(\theta) = p_{\theta}(y_1, y_2, \dots, y_N). \quad (1.7)$$

**NB:** We abused notation above stating the joint pdf for the observations.

**Remark 1.4.**

Very important: the likelihood function is not a probability/density function, as it is a function of the parameter. In particular, it is not true that  $\int_{\Theta} L(\theta) d\theta$  is one.

**Definition 1.3** (Maximum likelihood estimator).

The ML estimator is given by

$$\theta_{ML} = \arg \max_{\theta} L(\theta). \quad (1.8)$$

**Remark 1.5.**

In general (but, importantly, not always) we will consider i.i.d observations, in which case the likelihood factorises as  $L(\theta) = \prod_{n=1}^N p_{\theta}(y_n)$ . Furthermore, when optimising the likelihood we will consider the log-likelihood instead; in the i.i.d. case, this is

$$l(\theta) = \log L(\theta) = \sum_{n=1}^N \log p_{\theta}(y_n). \quad (1.9)$$

**Example 1.4** (Gaussian linear regression).

Let us consider the PGM given by

$$Y|x \sim \mathcal{N}(ax, \sigma^2), a, x \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+. \quad (1.10)$$

This is equivalent to  $Y = ax + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$ . The parameters in this setting are  $\theta = (a, \sigma^2)$ . Now consider the observations  $\{(x_n, y_n)\}_{n=1}^N$ .

Since  $p(y_n|x_n) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-1}{2\sigma^2}(y_n - ax_n)^2\right)$ , we can write the log-likelihood as

$$l(\theta) = \sum_{n=1}^N \frac{-1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2}(y_n - ax_n)^2. \quad (1.11)$$

The optimal  $(a, \sigma^2)$  can be found in closed form using the first order optimality conditions.

**Remark 1.6.**

Observe that optimising eq. (1.11) recovers the least squares solution.

**Example 1.5** (Binary classification).

Consider observations  $\{(x_n, y_n)\}_{n=1}^N \subset \mathbb{R}^d \times \{0, 1\}$  from a binary classification setting. Model the classifier as

$$p_\theta(y = 1|x) = \sigma(s(x)), \quad (1.12)$$

where  $\sigma(s(x)) = \frac{1}{1+e^{-s(x)}}$ , and  $s : \mathbb{R}^d \rightarrow \mathbb{R}$  is a feature extractor (e.g.,  $s(x) = a^\top x + b$ ). Assuming that the observations are i.i.d., we have

$$L(\theta) = \prod_{n=1}^N p(y_n|x_n) = \prod_{n=1}^N \sigma(s(x_n))^{y_n} (1 - \sigma(s(x_n)))^{1-y_n}, \quad (1.13)$$

and equivalently

$$l(\theta) = \sum_{n=1}^N y_n \log \sigma(s(x_n)) + (1 - y_n) \log(1 - \sigma(s(x_n))). \quad (1.14)$$

Does this expression seem familiar? If not, we will find out soon what this is.

**Example 1.6** (Clustering).

Consider a set of observations  $\{x_n\}_{n=1}^N \in \mathbb{R}^d$  and implement a clustering algorithm. We will assume that there are  $K \in \mathbb{N}$  clusters, each specified by a density  $p_k, k = 1, \dots, K$ ; this means that the probability of the RV  $X$  coming from the  $k$ -th cluster is  $\mathbb{P}(X \in C_k) = \pi_k$ , where  $\forall k, 0 \leq \pi_k \leq 1$  and  $\sum_{k=1}^K \pi_k = 1$ .

This is a mixture model, with density  $p(x) = \sum_{k=1}^K \pi_k p_k(x)$ , and parameters given by the cluster probabilities  $\pi_k$  and the parameters of the densities  $p_k = p_{\theta_k}$ . The log-likelihood is

$$l(\theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k p_k(x_n) \quad (1.15)$$

Note that there are two issues associated to optimising eq. (1.15).

- We do not recover the cluster assignments.
- The problem is ill-posed. E.g., if  $p_k = \mathcal{N}(\mu_k, \Sigma_k)$ , which is the usual choice, we can set  $\mu_k = x_n, \Sigma_k = 0$  which gives  $l = \infty$ .

We can overcome this drawback by introducing a collection of latent random variables  $Z_{nk} \in \{0, 1\}$ , that represents the cluster assignments. That is,

$$Z_{nk} = 1 \iff x_n \in C_k. \quad (1.16)$$

This allows us to write the conditional densities  $p(x_n|z_{nk}) = \prod_{k=1}^K p_k^{z_{nk}}$ , and thus to express the **complete-data likelihood** given by

$$l(\theta) = \log \prod_{n=1}^N \prod_{k=1}^K p_k^{z_{nk}}(x_n) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \log p_k(x_n). \quad (1.17)$$

Good and bad news: this objective is now theoretically feasible to optimise but impractical since we do not have access to the latent cluster assignments  $\{z_{nk}\}_{nk}$ .

A workaround to this is to estimate the cluster assignments, via its conditional expectation wrt the observations. That is,

$$\mathbb{E}(z_{nk}|x_{1:N}) = 1 * \mathbb{P}(z_{nk} = 1|x_n) + 0 * \mathbb{P}(z_{nk} = 0|x_n) = \mathbb{P}(z_{nk} = 1|x_n), \quad (1.18)$$

which can be computed explicitly using Bayes theorem in terms of the model parameters. Then, we can perform an iterative procedure by: i) optimising  $l(\theta)$  using  $\mathbb{E}(z_{nk}|x_{1:N})$ , and ii) computing  $\mathbb{E}(z_{nk}|x_{1:N})$  using  $\theta_{ML} = \arg \max l(\theta)$ .

This means that exact ML cannot be performed in this case. Also, does this procedure seem familiar?

The maximum likelihood estimator (MLE) satisfies several important theoretical properties:

- **Consistency:** Under the assumption that the statistical model is *identifiable*—i.e., different parameter values correspond to different probability distributions—the MLE converges to the true parameter as the number of observations grows. Intuitively, maximising the likelihood asymptotically minimises the Kullback–Leibler divergence between the true distribution and the distribution induced by a candidate parameter.
- **Equivariance:** If  $\hat{\theta}_{MLE}$  is the MLE of  $\theta$ , then for any transformation  $g$ , the MLE of  $g(\theta)$  is  $g(\hat{\theta}_{MLE})$ . This property allows us to compute MLEs under reparametrisations directly.
- **Asymptotic normality:** For large sample sizes, the MLE is approximately normally distributed around the true parameter with covariance matrix given by the inverse Fisher information. Formally,

$$\sqrt{n}(\hat{\theta}_{MLE} - \theta) \xrightarrow{d} \mathcal{N}(0, I(\theta)^{-1}),$$

where  $I(\theta)$  is the Fisher information matrix.

- **Asymptotic efficiency:** As a consequence of asymptotic normality, the MLE achieves the Cramér–Rao lower bound for the variance in the limit of large  $n$ , making it asymptotically optimal among unbiased estimators.

In practice, these properties justify the widespread use of the MLE: it not only converges to the true parameter under mild assumptions, but also allows for straightforward reparametrisations and provides an estimator with minimal asymptotic variance.

## 1.5 A brief intro to information theory

**NB:** This section is based on Chapter 6 of (Murphy, 2022)

**Motivation.** Let us consider a discrete RV  $X \in \{1, 2, \dots, K\}$  with pmf  $p_X$ . Observe that  $-\log p_X(x)$  represents a measure of *information* gained from obtaining the value  $a$  as a sample of  $X$ . Now consider a communication channel  $A \rightarrow B$ , where  $A$  is transmitting samples of  $X$  to  $B$ . When  $B$  received the samples, its *average information* can be expressed as

$$H(X) = - \sum_{x=1}^K p_X(x) \log p_X(x). \quad (1.19)$$

This quantity is known as *entropy* and—in connection with thermodynamics—it represent a measure of disorder or un-predictability of  $X$ .

**NB:** We will use  $H(X)$  and  $H(p_X)$  interchangeably.

**NB:** We will usually denote  $H(X) = -\mathbb{E}(\log p_X) = \mathbb{E}\left(\log \frac{1}{p_X}\right)$ .

Clearly,  $H(X) \geq 0$  with equality achieved for an RV that has always the same outcome with  $p_X(x) = 1$ . This is the deterministic, predictable, case. To revise further properties, let us recall the following result.

**Jensen's Inequality** Let  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  be a *convex* function and let  $X$  be a random variable such that  $\mathbb{E}[|X|] < \infty$ . Then

$$\phi(\mathbb{E}[X]) \leq \mathbb{E}[\phi(X)]. \quad (1.20)$$

since  $\log(\dots)$  is *concave*, the inequality is reversed and we have:

$$\mathbb{E}[\log X] \leq \log \mathbb{E}[X]. \quad (1.21)$$

**Remark 1.7.**

Equality in eq. (1.20) is only achieved when either  $\phi$  is affine, or  $X$  is constant almost surely, that is,  $\mathbb{P}(X = c) = 1$ . As a consequence, equality in eq. (1.21) is only achieved when  $X$  is constant almost surely (when the argument of the logarithm does not depend on  $x$ )

Keep this result in mind, as it will be used throughout the module.

Using Jensen on the definition of the entropy, we have

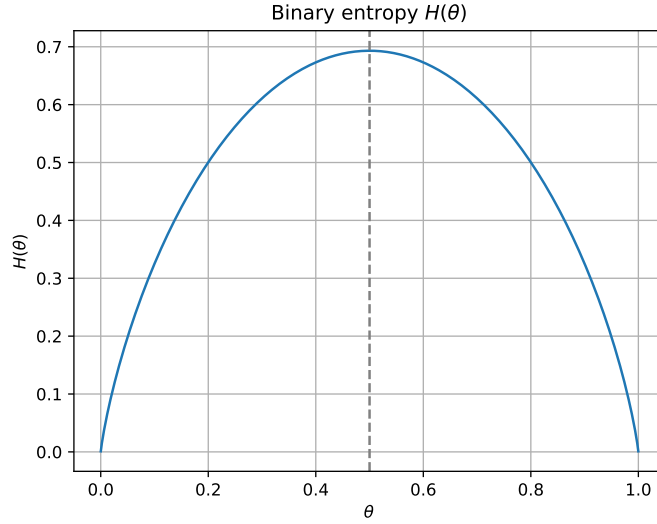
$$H(X) = \mathbb{E}\left(\log \frac{1}{p}\right) \leq \log \mathbb{E}\left(\frac{1}{p}\right) = \log \sum \frac{1}{p} = \log K. \quad (1.22)$$

This directly implies that the uniform distribution over  $\{1, 2, \dots, K\}$  has the largest entropy, since

$$H(U_{1:K}) = \mathbb{E} \left( \log \frac{1}{1/K} \right) = \log K. \quad (1.23)$$

**Example 1.7** (Bernoulli distribution).

Consider  $X \sim p_X(x) = \theta^x(1-\theta)^{1-x}$ ,  $\theta \in [0, 1]$ . The entropy is given by  $H(X) = -\theta \log \theta - (1-\theta) \log(1-\theta)$ . Figure 1.3 shows this function.



**Figure 1.3:** Entropy for Bernoulli.

The entropy, in addition to being a measure of disorder, can be understood as the cost of the optimal for of compression. Here, think of a compression strategy using symbols  $s_1, s_2, \dots$  with increasing size (or storage cost). For instance, think of storage using logical gates, meaning that these symbols are (equivalent to)

$$s_1 = 0, s_2 = 1, s_3 = 10, s_4 = 11, \dots \quad (1.24)$$

where storing more and more symbols becomes increasingly expensive. The compression strategy is then to assign each outcome of  $X$  with a symbol  $s_i$ . Intuitively, the optimal compression would assign  $s_i$  to the  $i$ -th most frequent value in  $\{1, 2, \dots, K\}$ , meaning that the *cost of storage* precisely grows precisely with  $\log \frac{1}{p_X}$ .

As a consequence, the average message size is the entropy  $H(X)$ , and thus can be understood as the cost of this compression strategy.

Now let us go back to our communication channel  $A \rightarrow B$ , where the receiver in  $B$  now mistakenly believes that  $X \sim q_X$ .  $B$ 's estimated entropy, or averaged information, would be

$$H_{CE}(p_X, q_X) = - \sum_{x=1}^K p_X(x) \log q_X(x). \quad (1.25)$$

Following the same rationale as above, this can be interpreted as the cost of compressing the sequence  $x_1, x_2, x_3, \dots$  using  $q_x$ .

This quantity is known as the cross-entropy between  $p_X$  and  $q_X$ . Note that this quantity is not symmetric.

It is relevant to study how  $H_{CE}(p_X, q_X)$  and  $H(P) = H_{CE}(p_X, p_X)$  relate to one another, in particular if one of them is (always) larger than the other.

Let us see:

$$H(p) - H_{CE}(p_X, q_X) = \sum_x p(x) \log \frac{q(x)}{p(x)} \quad (1.26)$$

$$\stackrel{\text{Jensen's}}{\leq} \log \sum_x p(x) \frac{q(x)}{p(x)} \quad (1.27)$$

$$= \log 1 = 0. \quad (1.28)$$

Therefore,  $H(p) \leq H(p, q)$ , with equality only achieved when  $p = q$ , as per Remark 1.7.

**Remark 1.8.**

Minimising the cross-entropy wrt to one of its arguments is precisely an attempt to match  $p = q$ .

The use of the entropy/cross-entropy that is going to be more relevant in our case is via its application to continuous RVs. This extension is

$$H(p) = - \int_{\mathcal{X}} p(x) \log p(x) dx \quad (1.29)$$

$$H(p, q) = - \int_{\mathcal{X}} p(x) \log q(x) dx \quad (1.30)$$

$$(1.31)$$

**Remark 1.9.**

Unlike its discrete formulation,  $H(p)$  can be positive, negative or zero for continuous RVs.

**Example 1.8** (Continuous uniform distribution).

Consider a RV  $X \sim U_{[a,b]}$ , its entropy is

$$H(U_{[a,b]}) = - \int_a^b \frac{1}{b-a} \log \frac{1}{b-a} dx = \log(b-a).$$

and can be zero (resp. negative) if  $b-a = 1$  (resp.  $b-a \leq 1$ ).

The difference between the entropy and crossentropy is of critical relevance in this module (and life). We will recall a relevant definition first.

**Definition 1.4** (Absolute Continuity).

Let  $(\Omega, \mathcal{F})$  be a measurable space and let  $P$  and  $Q$  be probability measures on it. We say that  $P$  is *absolutely continuous* with respect to  $Q$ , denoted  $P \ll Q$ , if for every  $A \in \mathcal{F}$ ,

$$Q(A) = 0 \implies P(A) = 0.$$

**Remark 1.10.**

If  $P \ll Q$ , and  $Q$  admits a density  $q$ ,  $P$  admits a density  $p$  satisfying  $p(x) = 0$  whenever  $q(x) = 0$ .

**Definition 1.5** (Kullback–Leibler Divergence).

Let  $P$  and  $Q$  be probability measures on a measurable space  $(\Omega, \mathcal{F})$  such that  $P \ll Q$ . If  $P$  and  $Q$  admit densities  $p$  and  $q$  with respect to a common base measure (e.g. Lebesgue measure), then

$$\text{KL}(p \parallel q) = \mathbb{E}_p \left[ \log \frac{p(X)}{q(X)} \right] = \int p(x) \log \frac{p(x)}{q(x)} dx.$$

**Definition 1.6** (Discrete KL Divergence).

For discrete distributions,

$$\text{KL}(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)}.$$

**Remark 1.11.**

Notice that the KL divergence is always positive:

$$\text{KL}(p \parallel q) = H(p, q) - H(p) \geq 0. \quad (1.32)$$

The KL is a *divergence*, i.e., a function that quantifies how far  $p$  is from  $q$  that is i) always positive, and ii)  $\text{KL}(p \parallel q) = 0 \iff p = q$  (identify of the indiscernible). However, note that the KL is not a distance, since

- is not symmetric
- does not have triangle inequality.

Critically, the  $\text{KL}(p \parallel q)$  is only defined when  $P \ll Q$ .

## 1.6 KL divergence as a metric to compare $p$ and $q$

In the continuous case, we are interested in understanding what type of convergence KL gives. Let us consider, other two divergences:

- $L_1(p \parallel q) = \int_{\mathcal{X}} |p(x) - q(x)| dx$
- $\chi^2(p \parallel q) = \int_{\mathcal{X}} \frac{|p(x) - q(x)|^2}{q(x)} dx.$

**Example 1.9** (KL versus  $L_1$ ).

Consider  $p(x) = \text{Uniform}(0, 1)$  and  $q_n(x) = \mathbb{1}_{x \in [0, 1/n]} e^{-n} + \mathbb{1}_{x \in [1/n, 1]} c_n$ , with  $c_n \geq 0$  so that  $q$  integrates 1 ( $n > 1$ ).

Let us first compute  $c_n$  explicitly. Since  $q_n$  must integrate to one, we require

$\int_0^{1/n} e^{-n} dx + \int_{1/n}^1 c_n dx = 1$ , which yields

$$\frac{1}{n} e^{-n} + \left(1 - \frac{1}{n}\right) c_n = 1.$$

Solving for  $c_n$ , we obtain

$$c_n = \frac{1 - \frac{1}{n} e^{-n}}{1 - \frac{1}{n}} = \frac{n - e^{-n}}{n - 1}. \quad (1.33)$$

Note that here, we have

$$L_1(p||q_n) = \int_0^{1/n} |e^{-n} - 1| dx + \int_{1/n}^1 |c_n - 1| dx = \frac{|e^{-n} - 1|}{n} + \frac{n|c_n - 1|}{n - 1} \quad (1.34)$$

$$\text{KL}(p||q_n) = \int_0^{1/n} -\log e^{-n} dx + \int_{1/n}^1 -\log c_n dx = 1 + \frac{-n}{n - 1} \log c_n. \quad (1.35)$$

Now take  $n \rightarrow \infty$ . From eq. (1.33) we can see that  $c_n$  converges to 1. Therefore,  $L_1(p||q) \rightarrow 0$ . However, note that  $\text{KL}(p||q) \rightarrow 1$ .

**Example 1.10** (KL versus  $\chi^2$ ).

Consider now  $p_\epsilon = (1 - \epsilon, \epsilon)$  and  $q_\epsilon = (1 - \epsilon^2, \epsilon^2)$  two Bernoulli distribution with different parameters. Again, we have:

$$\chi^2(p_\epsilon||q_\epsilon) = \frac{\|1 - \epsilon - 1 + \epsilon^2\|^2}{1 - \epsilon^2} + \frac{\|\epsilon - \epsilon^2\|^2}{\epsilon^2} = \frac{\|\epsilon^2 - \epsilon\|^2}{1 - \epsilon^2} + \|1 - \epsilon\|^2 \quad (1.36)$$

$$\text{KL}(p_\epsilon||q_\epsilon) = (1 - \epsilon) \log \frac{1 - \epsilon}{1 - \epsilon^2} + \epsilon \log \frac{\epsilon}{\epsilon^2} = (1 - \epsilon) \log \frac{1}{1 + \epsilon} + \epsilon \log \frac{1}{\epsilon}. \quad (1.37)$$

This time, taking  $\epsilon \rightarrow 0$ , we have  $\text{KL}(p_\epsilon||q_\epsilon) \rightarrow 0$  (l'Hôpital's rule), but  $\chi^2(p_\epsilon||q_\epsilon) \rightarrow 1$

**Remark 1.12.**

The objective of these examples is to show that under different divergences, one can have different criteria of convergence. In the first case,  $q_n$  converges to  $p$  under  $L_1$ , but not under KL. In the second case,  $p_\epsilon$  converges to  $q_\epsilon$  under KL but not under  $\chi^2$ . This give a sense of *hierarchy* across divergences, where some are said to induce stronger topologies than others. The stronger the topology, the more demanding the conditions for convergence (or fewer sequences are admitted to converge). In general, we consider KL as one of the stronger divergences (but there are some that are even stronger as we just saw).

**Direct versus reverse KL.** Since  $\text{KL}(p||q)$  is not symmetric, we are interested in studying the *reverse* divergence  $\text{KL}(q||p)$  and understanding how it relates its *direct* counterpart.

Since  $P \ll Q$  is needed for  $\text{KL}(p||q)$ , it is required that  $P \gg Q$  for  $\text{KL}(q||p)$ . This gives intuition of  $\text{KL}(p||q)$  as a metric assessing how well  $q$  approximates  $p$  (and not viceversa); this is because if there is a set  $A \subset \mathcal{X}$  such that  $Q(A) = 0$  and  $P(A) > 0$  is strongly penalised, unlike the opposite case.



Let us see a numerical example.

**Example 1.11** (Assymetry of the KL between two Gaussians).

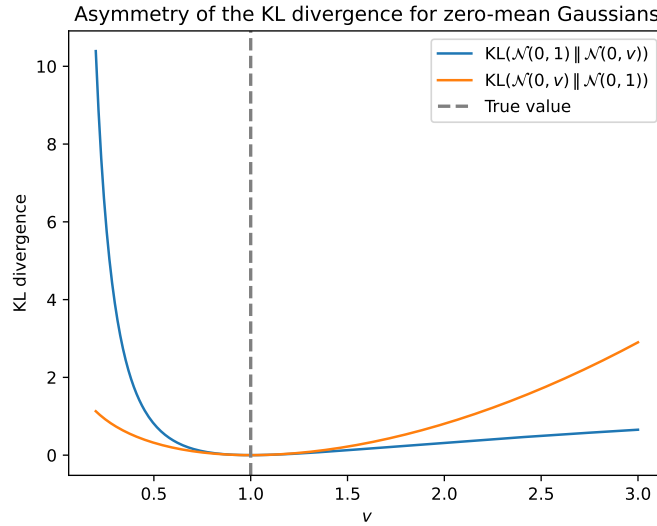
The KL divergence between two Gaussians is

$$\text{KL}(\mathcal{N}(\mu_0, \sigma_0) \parallel \mathcal{N}(\mu_1, \sigma_1)) = \log \frac{\sigma_1}{\sigma_0} + \frac{\sigma_0^2 + (\mu_0 - \mu_1)^2}{2\sigma_1^2} - \frac{1}{2}. \quad (1.38)$$

Let us consider  $p = \mathcal{N}(0, 1)$  and  $q = \mathcal{N}(0, v^2)$ , and evaluate

- $\text{KL}(p \parallel q) = \frac{1}{2}(\log v^2 + v^{-2} - 1)$
- $\text{KL}(q \parallel p) = \frac{1}{2}(-\log v^2 + v^2 - 1)$ .

Fig. 1.4 shows these functions, note how the penalisation strength depends on the direction.



**Figure 1.4:** Direct and reverse KL for zero mean Gaussians as a function of the variance.

**Example 1.12** (KL gradient flow).

Let us now find the approximating  $q$  via optimisation for the above example. We can do this via optimisation. Differentiating eq. (1.38) wrt to  $\mu_1$  and  $\sigma_1$ , we have

$$\nabla_{\mu_1} \text{KL}(\mathcal{N}(\mu_0, \sigma_0) \parallel \mathcal{N}(\mu_1, \sigma_1)) = \frac{(\mu_1 - \mu_0)}{\sigma_1^2} \quad (1.39)$$

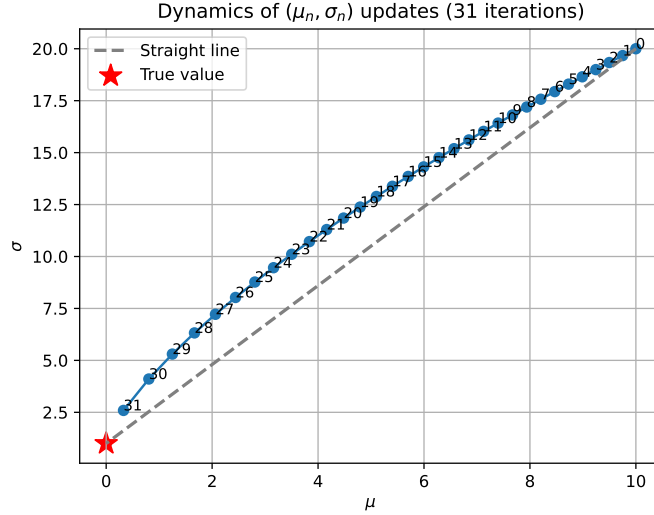
$$\nabla_{\sigma_1} \text{KL}(\mathcal{N}(\mu_0, \sigma_0) \parallel \mathcal{N}(\mu_1, \sigma_1)) = \frac{1}{\sigma_1} - \frac{\sigma_0^2}{\sigma_1^3} = \frac{\sigma_1^2 - \sigma_0^2}{\sigma_1^3} \quad (1.40)$$

Where it is clear the this is minimised for  $q = p$ . Additionally, we can build the gradient descent rule:

$$\mu_n \rightarrow \mu_n - \eta_\mu \frac{(\mu_n - \mu_0)}{\sigma_n^2} \quad (1.41)$$

$$\sigma_n \rightarrow \sigma_n - \eta_\sigma \frac{\sigma_n^2 - \sigma_0^2}{\sigma_n^3} \quad (1.42)$$

Figure 1.5 implements these recursions.



**Figure 1.5:** KL gradient flow between two Gaussians.

**KL and maximum likelihood.** Let us now return to the setting of the learning problem. Consider a true model given by  $p$  and iid observations  $x_1, x_2, \dots, x_N \sim p(x)$ . We could use the KL to look for the best approximator of  $p$  with a given family of candidate models  $\{q_\theta, \theta \in \Theta\}$ . That is,

$$\theta^* = \arg \min \text{KL}(p \| q_\theta). \quad (1.43)$$

Though it sounds good, this is unfeasible in practice since  $p$  is unknown. However, note that there is a workaround to that. We can write the above expression as

$$\theta^* = \arg \min \int_{\mathcal{X}} p(x) \log p(x) dx - \int_{\mathcal{X}} p(x) \log q_\theta(x) dx \quad (1.44)$$

$$= \arg \max \mathbb{E}_p[\log q_\theta(x)] \quad (1.45)$$

$$\approx \arg \max \sum_{x_i} \log q_\theta(x_i), \quad (1.46)$$

where the last approximation is due to Monte Carlo. This reveals that maximum likelihood is (asymptotically) equivalent to minimising the KL divergence between the model candidate and the (true) empirical distribution.

## 1.7 Concluding remarks

As we will throughout the module, when designing/choosing a PGM, we will be faced with the following scenarios.

- **Explicit-likelihood models:** These include classical statistical models such as Gaussians, exponentials, Bernoulli,  $X^2$ , log-normal, but also combinations or transformations that we construct as long as they have an explicit likelihood. These are Gaussian mixtures, piece-wise defined distributions, and any pushforward model constructed through an invertible transformation so that its density can be calculated via the change of variable theorem. These last models are referred to as normalising flows, as they assume a Gaussian source measure.
- **Implicit models:** As the name suggests, these models are only implicitly defined via a data-generating mechanism, usually involving sampling. For instance, take a RV  $Z \in \mathbb{R}^d \sim \mathcal{N}(0, I_d)$  and construct  $X = T_\theta(Z)$ , where  $T_\theta$  is collection of neural networks which are sequentially applied to  $Z$  with the aim to replicate learnt dynamics that make  $X$  flow towards the desired distribution. Depending on the parametrisation, these models are known as score-based models, diffusion models, or flow matching.

### 1.7.1 Suggested exercises

1. **Generative vs discriminative modelling (theory).** Let  $(X, Y)$  be random variables with joint distribution  $p(x, y)$ .
  - (a) Define what is meant by a *generative model* and a *discriminative model*.
  - (b) Show how a generative model can be used to construct a classifier.
  - (c) Discuss one advantage and one limitation of generative modelling relative to discriminative modelling.
2. **Information-theoretic objectives (theory).** Let  $p$  be a data-generating distribution and  $q_\theta$  a parametric model.
  - (a) Define the entropy  $H(p)$ , cross-entropy  $H(p, q_\theta)$ , and KL divergence  $\text{KL}(p||q_\theta)$ .
  - (b) Show that maximising the log-likelihood of data sampled from  $p$  is equivalent to minimising  $\text{KL}(p||q_\theta)$ .
  - (c) Explain why minimising  $\text{KL}(p||q_\theta)$  and  $\text{KL}(q_\theta||p)$  lead to qualitatively different approximations.
3. **Maximum likelihood density estimation (coursework).** You are given samples from a one-dimensional distribution.
  - (a) Fit a Gaussian model by maximum likelihood.
  - (b) Fit a mixture of Gaussians by maximum likelihood.
  - (c) Empirically compare the learned models using log-likelihood and visual inspection.

4. **Forward and reverse KL divergence (coursework).** Consider approximating a multimodal target distribution using a unimodal Gaussian.
- (a) Numerically minimise  $\text{KL}(p\|q)$  and  $\text{KL}(q\|p)$ .
  - (b) Visualise the resulting solutions.
  - (c) Explain the observed behaviour using the geometry of the KL divergence.

## Week 2

# Latent variable models and Expectation Maximisation

**NB:** This is based on Chapter 9 of (Bishop, 2006).

## 2.1 Gaussian mixtures

Consider a dataset  $\{x_1, x_2, \dots, x : N\} \subset \mathbb{R}^d$ . Our task is to partition this set into  $K \in \mathbb{N}$  subsets; we will consider  $K$  known for now. Intuitively, each subset—referred to as *cluster* from now on—of points should share some common or similar patterns; a formal definition of similarity in this case will be ignored for now.

A natural solution for this is to define  $K$  prototypes  $\{\mu_1, \mu_2, \dots, \mu_K\} \subset \mathbb{R}^d$  and determine the assignment of each datapoint  $x_n$  to each prototype  $\mu_k$ , according to a given criterion.

To solve this optimisation problem, we can define a set of binary variables  $\{r_{nk}\}_{nk} \subset \{0, 1\}$ , where

$$r_{nk} = 1 \iff x_n \text{ is assigned to } \mu_k. \quad (2.1)$$

Then, using the Euclidean distance as similarity criterion, the objective can be

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2. \quad (2.2)$$

The solution to the clustering problem obtained from the optimisation of eq. (2.2) is

$$r_{nk} = \begin{cases} 1, & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2, \\ 0, & \text{if not.} \end{cases} \quad (2.3)$$

$$\mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}. \quad (2.4)$$

The solution can be found via iteratively implementing these equations, which is known as the  $k$ -means algorithm.

**Remark 2.1.**

Observe that the  $k$ -means recursion ensure convergence in a finite number of steps:

this is because eq. (2.3) define a discrete number of solutions, and (2.4) is the global optima for a given  $\{r_{nk}\}_{nk}$ .

There are some known drawbacks of  $k$ -means, for instance

- Speed: computing the assignment variables has a cost  $\mathcal{O}(NK)$ .
- It depends on the Euclidean distance that might not be robust to outliers
- It only provide hard assignments, not a degree of *responsibility*.

## 2.2 The Gaussian mixture model

Let us consider the following PGM:

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k), \quad (2.5)$$

where  $0 \ll \pi_k \ll 1$ ,  $\sum_{k=1}^K \pi_k = 1$ ,  $\mu_k \in \mathbb{R}^d$  and  $\Sigma_k \in \mathbb{R}^{d \times d}$ .

This model suggests to be an improved clustering model over  $K$ -means, since it—at least—allows for learning the shape (variance) of each cluster, admits the definition of a soft assignment variable.

However, note that the likelihood of this models is ill posed. Denoting the parameters by  $\theta = \{\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K}\}$  and the i.i.d. data  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$  the log-likelihood is given by

$$l(\theta) = \log p(\mathbf{x}|\theta) = \log \prod_{n=1}^N p(x_n|\theta) = \sum_{n=1}^N \log p(x_n|\theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k). \quad (2.6)$$

This objective can reach an infinite value if a component is assigned to a single data-point with a vanishing variance. Additionally, for each possible assignment, there are  $K!$  different solution that provide such assignment.

We will derive an equivalent formulation to the PGM above that admits an easier training procedure. To this end, let us introduce a set of  $K$  latent variables  $\{z_k\} \subset \{0, 1\}$ ,  $\sum_{k=1}^K z_k = 1$ . We can write

$$p(x, z) = p(x|z)p(z). \quad (2.7)$$

Also, defining  $p(z_k = 1) = \pi_k$ , we can express the pmf/pdf:

$$p(z) = \prod_{k=1}^K \pi_k^{z_k} \quad (2.8)$$

$$p(x|z) = \prod_{k=1}^K \mathcal{N}(\mu_k, \Sigma_k)^{z_k}, \quad (2.9)$$

with the marginal pdf over  $x$  as

$$p(x) = \sum_{k=1}^K p(z_k) p(x|z_k) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k), \quad (2.10)$$

Thus, showing that the formulations are equivalent.

In this formulation, let us define the *responsibilities* of the  $k$ -th component to explain the observation  $x$  given by

$$\gamma(z_k) \stackrel{\text{def}}{=} p(z_k = 1|x) = \frac{p(x|z_k = 1)p(z_k = 1)}{\sum_{j=1}^K p(x|z_j = 1)p(z_j = 1)} \quad (2.11)$$

$$= \frac{\pi_k \mathcal{N}(x; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x; \mu_j, \Sigma_j)}. \quad (2.12)$$

**Remark 2.2.**

The latent-variable formulation of the GMM allows for direct sampling from that PGM: first sample  $z \sim p(z) = \prod_{k=1}^K \pi_k$ , and then sample  $x \sim p(x|z) = \prod_{k=1}^K \mathcal{N}(x; \mu_k, \Sigma_k)^{z_k}$ . This is known as *ancestral sampling*.

## 2.3 Expectation Maximisation for GMMs

We will introduce a learning approach for PGMs that feature a latent variable called Expectation Maximisation (EM). First we will present in in the particular case of the GMM model, and then in the general case.

The first order optimality conditions for the log-likelihood in eq. (2.6) give

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n \quad (2.13)$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^\top \quad (2.14)$$

$$\pi_k = \frac{N_k}{N}, \quad (2.15)$$

where we have defined the effective number of samples per component as  $N_k = \sum_{n=1}^N \gamma(z_{nk})$ .

**Exercise 2.1.**

Derive eqs. (2.13)-(2.15)

**Remark 2.3.**

Observe how the optimal mean and variance of each component is a weighted average of all the data points, where the weights are proportional to the responsibility (contribution) of that component to generation of the sample. Also, note that eqs. (2.13)-(2.15) can be considered as the soft-assignment version of the  $K$ -means solutions, with the additional flexibility of having an learnable expression for the shape of the clusters.

Eqs. (2.13)-(2.15) do not provide a direct closed-form solution, since they depend on the responsibilities  $\gamma(z_{nk})$  which are functions of the latent variable. However, they can still be implemented in following the steps:

E Compute  $\gamma(z_{nk}) = p(z_{nk} = 1|\mathbf{x})$

M use  $\gamma(z_{nk})$  to compute eqs. (2.13)-(2.15).

## 2.4 An interpretation of EM

Let us now leave the GMM aside. In more general, perhaps abstract, terms, the goal of the EM algorithm is to find maximum likelihood solutions for latent variable models (LVMs) by breaking down the the optimisation problem into a functional approximation of the likelihood, and the the (simpler) optimisation of the approximation.

Recall our notation involving a an observed variable  $x$  and a latent variable  $z$ . In general LVMs, the log-likelihood is given by

$$\log p(x|\theta) = \int_{\mathcal{X}} p(x|\theta, z)p(z|\theta)d\theta. \quad (2.16)$$

**NB:** We treat both the discrete and continuous equivalently.

In general, direct optimisation of eq.(2.16) is difficult. In fact, even calculating the above expression is only possible in limited cases, since mixtures do not mix well with the logarithm. In fact, even for likelihood in the exponential family, the mixture is not longer exponential and thus the application of the logarithm does not remove the exponential as in the single-Gaussian case.

Let us then consider the hypothetical scenario where we have access to the values of  $z$  alongside the observed  $x$ .

### Definition 2.1.

We will refer to the  $\{\mathbf{x}, \mathbf{z}\}$  as the complete dataset, while  $\mathbf{x}$  will be the *observed* or *incomplete* dataset.

The related likelihood to the complete dataset would be

$$\log p(\mathbf{x}, \mathbf{z}|\theta) = \log \prod_{n=1}^N p(x_n, z_n|\theta) \quad (2.17)$$

$$= \log \prod_{n=1}^N p(x_n|z_n, \theta)p(z_n|\theta) \quad (2.18)$$

$$= \sum_{n=1}^N \log p(x_n|z_n, \theta) + \log p(z_n|\theta) \quad (2.19)$$

which we will assume are simpler to evaluate and optimise.

This, however, is impractical since  $\mathbf{z}$  is unknown. An interesting interpretation of this optimisation problem is presented next.



**Remark 2.4.**

Since  $z_n$  is not observed, the complete-data log-likelihood in eq. (2.19) can be interpreted as a random function. Therefore, an alternative optimisation strategy is estimate its expectation (E step) and optimise this deterministic expression (M step). We will formally justify why taking the expectation is more than an intuition.

In more detail, the 2-step optimisation procedure can be implemented to move from a candidate solution  $\theta^{\text{old}}$  by first computing  $p(\mathbf{z}|\mathbf{x}, \theta^{\text{old}})$  and then the expectation of the complete-data log-likelihood in eq. (2.19) given by

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{z}} \log p(\mathbf{x}, \mathbf{z}|\theta) p(\mathbf{z}|\mathbf{x}, \theta^{\text{old}}), \quad (2.20)$$

to finally implement

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}). \quad (2.21)$$

**Remark 2.5.**

This procedure can also be used for maximum a posteriori estimation, in which case  $Q(\theta, \theta^{\text{old}}) \rightarrow Q(\theta, \theta^{\text{old}}) + \log p(\theta)$  incorporates the prior over the parameter.

Now let us return to the GMM case and feed back these observations. For the GMM, the complete-data log-likelihood is

$$p(\mathbf{x}, \mathbf{z}|\theta) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(x_n|\mu_k, \Sigma_k)^{z_{nk}}, \quad (2.22)$$

and thus the log-likelihood issues

$$l(\theta) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} (\log \pi_k + \log \mathcal{N}(x_n|\mu_k, \Sigma_k)), \quad (2.23)$$

which is tractable.

**Remark 2.6.**

The objective in eq. (2.23) is straightforward to optimise: since only one term in the  $k$ -sum is non-zero, which means the that mean and covariances are simple estimates as in the case of the single Gaussian. Furthermore, imposing the first order optimality condition and enforcing the  $\sum_{k=1}^K \pi_k = 1$ , gives  $\pi_k = \sum_{n=1}^N x_n / N$  directly.

Recall that this optima is a function of  $\mathbf{z}$ , so we need to compute the expectation wrt to it. We have

$$p(\mathbf{z}|\mathbf{x}, \theta) \propto p(\mathbf{z}, \mathbf{x}|\theta) = \prod_{n=1}^N \prod_{k=1}^K \underbrace{\pi_k^{z_{nk}} \mathcal{N}(x_n|\mu_k, \Sigma_k)^{z_{nk}}}_{\propto p(z_{nk}|\mathbf{x}_n, \theta)}, \quad (2.24)$$

which means that  $p(\mathbf{z}|\mathbf{x}, \theta)$  factorises wrt to  $n$  and  $k$ , and thus all the  $z_n$  are independent. This is reasonable, since the cluster responsibilities over one samples should affect the rest (due to the i.i.d. assumption).

The expectation of  $\mathbf{z}$  can be computed as follows,

$$\mathbb{E}(z_{nk}) = 1 \cdot p(z_{nk} = 1 | \mathbf{x}, \theta) + 0 \cdot p(z_{nk} = 0 | \mathbf{x}, \theta) \quad (2.25)$$

$$= p(z_{nk} = 1 | \mathbf{x}, \theta) \quad (2.26)$$

$$= \gamma(z_{nk}) \quad (2.27)$$

$$\stackrel{\text{def}}{=} \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)} \quad (2.28)$$

Note fm eq.(2.23) that the objective is linear in  $z_{nk}$ , which makes the computation of its expectation direct:

$$\mathbf{E}_z \log p(\mathbf{x}, \mathbf{z} | \theta) = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) (\log \pi_k + \log \mathcal{N}(x_n | \mu_k, \Sigma_k)). \quad (2.29)$$

### Exercise 2.2.

Show that GMM recovers  $K$ -means. For that, choose  $p(x | \mu_k, \Sigma_k) = \mathcal{N}(x | \mu_k, \epsilon I)$  with  $\epsilon > 0$  fixed to show that you recover a *soft-assignment* version of  $K$ -means. Then, take  $\epsilon \rightarrow 0$  to recover vanilla  $K$ -means.

### Exercise 2.3.

See the applications of EM to mixtures of Bernoulli and Bayesian linear regression in [Bishop].

## 2.5 EM in its general form

Recall:

$$\underbrace{p(\mathbf{x} | \theta)}_{\text{difficult}} = \sum_{\mathbf{z}} \underbrace{p(\mathbf{x}, \mathbf{z} | \theta)}_{\text{easier}}. \quad (2.30)$$

We are interested in deriving EM as a model-approximation approach. To that end, let us consider a distribution over the latent variable  $q(z)$ ; intuitively, this distribution will approximate  $p(z | x, \theta)$ . For any choice of  $q$ , the following holds:

$$\begin{aligned} \log p(x | \theta) &= \sum_z q(z) \log p(x | \theta) && \leftarrow \log p(x | \theta) \text{ is constant wrt } z \\ &= \sum_z q(z) \log \left( p(x | \theta) \frac{p(z | x, \theta) q(z)}{p(z | x, \theta) q(z)} \right) && \leftarrow \text{multiply by 1} \\ &= \sum_z q(z) \log \left( \frac{p(x, z | \theta)}{q(z)} \cdot \frac{q(z)}{p(z | x, \theta)} \right) && \leftarrow \text{arrange} \\ &= \underbrace{\sum_z q(z) \log \left( \frac{p(x, z | \theta)}{q(z)} \right)}_{\mathcal{L}(q, \theta)} + \underbrace{\sum_z q(z) \log \left( \frac{q(z)}{p(z | x, \theta)} \right)}_{\text{KL}(q(z) \| p(z | x, \theta))} && \leftarrow \text{split} \end{aligned}$$

**Remark 2.7.**

Recall that the KL is always non-negative, meaning that  $\text{KL}(q(z)||p(z|x, \theta))$  is a lower bound for  $\log p(x|\theta)$ .

EM breaks the ML problem into two simpler problems related to the computation of a lower bound and its optimisation. The objective above can be minimised in two stages:

- First, the distribution  $q$  is chosen in order to minimise the term  $\text{KL}(q(z)||p(z|x, \theta))$ , where, the optimal solution is  $q(z) = p(z|x, \theta)$ .
- Then, using that choice for  $q$ , the term  $\mathcal{L}(q, \theta)$  can be optimised. Notice that this term is an expectation

**Remark 2.8.**

This view formalises the intuition that we presented for the GMM case. The choice of taking the expectation of the complete log-likelihood is not arbitrary, but follows from finding the optimal approximating distribution in the KL sense.

**To do:** A few more things to include: Monotonicity of EM its ability to fix the ill-posedness of the objective

## 2.6 Concluding remarks

**To do:** to be completed

### 2.6.1 Suggested exercises

- 1. Latent variable models and incomplete data (theory).**
  - (a) Define a latent variable model and distinguish between complete and incomplete data.
  - (b) Explain why direct maximisation of the marginal likelihood is often intractable.
  - (c) Describe how the introduction of latent variables simplifies modelling but complicates inference.
- 2. Expectation–Maximisation algorithm (theory).** Consider a latent variable model with parameters  $\theta$ .
  - (a) Derive the EM algorithm starting from the marginal log-likelihood.
  - (b) Define the  $Q$ -function and explain the role of the E-step and the M-step.
  - (c) Prove that each EM iteration does not decrease the log-likelihood.
- 3. Gaussian mixture models (coursework).**
  - (a) Implement the EM algorithm for Gaussian mixture models.
  - (b) Investigate the effect of initialisation on convergence.

- (c) Illustrate the relationship between  $k$ -means and GMMs by varying the covariance structure.

**4. Likelihood degeneracy and regularisation (coursework).**

- (a) Demonstrate empirically the likelihood degeneracy of Gaussian mixture models.
- (b) Propose and implement at least one regularisation strategy.
- (c) Analyse how regularisation affects the learned parameters and likelihood.

## References

- Bishop, C. M. (2006). *Pattern recognition and machine learning*. New York: Springer.
- Murphy, K. P. (2022). *Probabilistic machine learning: An introduction* (1st ed.). Cambridge, MA, USA: The MIT Press.