

# Machine Learning

Notes taken by Runqiu Ye  
Carnegie Mellon University

Spring 2025

## Contents

<b>1</b>	<b>Probability and statistical inference</b>	<b>3</b>
1.1	Probability . . . . .	3
1.2	Statistical inference . . . . .	4
1.3	PAC learning . . . . .	5
<b>2</b>	<b>Mixture models and EM</b>	<b>9</b>
2.1	Kullback-Leibler (KL) Divergence . . . . .	9
2.2	The EM Algorithm in General . . . . .	11
<b>3</b>	<b>Approximate inference</b>	<b>12</b>

# 1 Probability and statistical inference

## 1.1 Probability

**Definition** (Types of convergence). Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of random variables and  $X$  be another random variable. Let  $F_n$  be the CDF of  $X_n$  for each  $n \in \mathbb{N}$  and  $F$  be the CDF of  $X$ .

1.  $X_n$  converges to  $X$  **in probability** and write  $X_n \xrightarrow{P} X$  if for arbitrary  $\varepsilon > 0$ ,

$$\mathbb{P}[|X_n - X| > \varepsilon] \rightarrow 0$$

as  $n \rightarrow \infty$ .

2.  $X_n$  converges to  $X$  **in distribution** and write  $X_n \rightsquigarrow X$  if

$$\lim_{n \rightarrow \infty} F_n(t) = F(t)$$

for all  $t$  where  $F$  is continuous.

3.  $X_n$  converges to  $X$  in  $L^p$  if

$$\mathbb{E}[|X_n - X|^p] \rightarrow 0$$

as  $n \rightarrow \infty$ . In particular, say  $X_n$  converges to  $X$  in **quadratic mean** and write  $X_n \xrightarrow{qm} X$  if  $X_n$  converges to  $X$  in  $L^2$ .

4.  $X_n$  converges to  $X$  **almost surely** and write  $X_n \xrightarrow{as} X$  if

$$\mathbb{P}\left[\lim_{n \rightarrow \infty} X_n = X\right] = 1.$$

**Theorem.** The following implication holds:

1. If  $X_n$  converges to  $X$  almost surely, then  $X_n$  converges to  $X$  in probability.
2. If  $X_n$  converges to  $X$  in  $L^p$ , then  $X_n$  converges to  $X$  in probability.

*Proof.* 1. If  $X_n$  converges to  $X$  almost surely, the set of points  $O = \{\omega : \lim_{n \rightarrow \infty} X_n(\omega) \neq X(\omega)\}$  has measure zero. Now fix  $\varepsilon > 0$  and consider the sequence of sets

$$A_n = \bigcup_{m=n}^{\infty} \{|X_m - X| > \varepsilon\}.$$

Note that  $A_n \supset A_{n+1}$  for each  $n \in \mathbb{N}$  and let  $A_{\infty} = \bigcap_{n=1}^{\infty} A_n$ . Now show  $\mathbb{P}[A_{\infty}] = 0$ . If  $\omega \notin O$ , then  $\lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)$  and thus  $|X_n(\omega) - X(\omega)| < \varepsilon$  for some  $n \in \mathbb{N}$ . Therefore,  $\omega \notin A_{\infty}$ . It follows that  $A_{\infty} \subset O$  and  $\mathbb{P}[A_{\infty}] = 0$ .

By monotone continuity, we have  $\lim_{n \rightarrow \infty} \mathbb{P}[A_n] = \mathbb{P}[A_{\infty}]$ . It follows that

$$\mathbb{P}[|X_n - X| > \varepsilon] \leq \mathbb{P}[A_n] \rightarrow 0$$

as  $n \rightarrow \infty$ . This completes the proof.

2. From Chebyshev's inequality, we have

$$\mathbb{P}[|X - X_n| > \varepsilon] \leq \frac{1}{\varepsilon^p} \mathbb{E}[|X - X_n|^p].$$

The claim follows directly.

□

**Theorem** (Central Limit Theorem). Let  $X_1, \dots, X_n$  be i.i.d. with mean  $\mu$  and variance  $\sigma^2$ . Let  $S_n = \frac{1}{n} \sum_{i=1}^n X_i$ . Then

$$Z_n = \frac{S_n - \mu}{\sqrt{\text{Var } S_n}} = \frac{\sqrt{n}(S_n - \mu)}{\sigma} \rightsquigarrow Z,$$

where  $Z \sim \mathcal{N}(0, 1)$ . In other words,

$$\lim_{n \rightarrow \infty} \mathbb{P}[Z_n < z] = \Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx.$$

Also write  $Z_n \approx \mathcal{N}(0, 1)$ .

## 1.2 Statistical inference

**Definition.** Let  $X_1, \dots, X_n$  be  $n$  i.i.d. data points observed from some distribution  $F$  with respect to parameter  $\theta$ . A point estimator  $\hat{\theta}_n$  of the parameter  $\theta$  is some function of  $X_1, \dots, X_n$ :

$$\hat{\theta}_n = g(X_1, \dots, X_n).$$

The bias of an estimator is defined as

$$\text{bias}(\hat{\theta}_n) = \mathbb{E}_\theta[\hat{\theta}_n] - \theta.$$

The mean squared error is defined as

$$\text{MSE} = \mathbb{E}_\theta(\hat{\theta}_n - \theta)^2.$$

**Definition** (Consistent point estimator). A point estimator  $\hat{\theta}_n$  of a parameter  $\theta$  is **consistent** if  $\hat{\theta}_n \xrightarrow{\text{P}} \theta$ .

Next we an important relation between bias, variance, and MSE. This is a more rigorous way to express the **bias-variance tradeoff** of point estimators.

**Theorem.** The MSE can be written as

$$\text{MSE} = \text{bias}^2(\hat{\theta}_n) + \text{Var}_\theta(\hat{\theta}_n).$$

*Proof.* Let  $\bar{\theta}_n = \mathbb{E}_\theta(\hat{\theta}_n)$ . Then we have

$$\begin{aligned} \mathbb{E}_\theta(\theta - \hat{\theta}_n)^2 &= \mathbb{E}_\theta(\theta - \bar{\theta}_n + \bar{\theta}_n - \hat{\theta}_n)^2 \\ &= \mathbb{E}_\theta(\theta - \bar{\theta}_n)^2 - 2(\theta - \bar{\theta}_n)\mathbb{E}_\theta(\bar{\theta}_n - \hat{\theta}_n) + \mathbb{E}_\theta(\bar{\theta}_n - \hat{\theta}_n)^2 \\ &= (\theta - \bar{\theta}_n)^2 + \mathbb{E}_\theta(\bar{\theta}_n - \hat{\theta}_n)^2 \\ &= \text{bias}^2(\hat{\theta}_n) + \text{Var}_\theta(\hat{\theta}_n), \end{aligned}$$

where we have used the fact that  $\mathbb{E}_\theta(\bar{\theta}_n - \hat{\theta}_n) = \bar{\theta}_n - \mathbb{E}_\theta(\hat{\theta}_n) = \bar{\theta}_n - \bar{\theta}_n = 0$ . □

Below is the definition of a confidence set/interval.

**Definition.** A  $1 - \alpha$  interval for a parameter  $\theta$  is an interval  $C_n = (a, b)$  where  $a = a(X_1, \dots, X_n)$  and  $b = b(X_1, \dots, X_n)$  are functions of data such that

$$\mathbb{P}_\theta[\theta \in C_n] \geq 1 - \alpha \text{ for all } \theta \in \Theta.$$

In other word,  $(a, b)$  traps  $\theta$  with probability  $1 - \alpha$ .

**Warning!** In the above definition,  $C_n$  is random and  $\theta$  is fixed.

### 1.3 PAC learning

PAC learning is short for Probably Approximate Correct learning, and the setting of PAC learning is as follows:

- We have data  $x \in \mathbb{R}^d$  and label  $y \in \{-1, +1\}$ .
- We collect features  $x_1, \dots, x_n$  iid from some distribution  $D$ . Note that we make no assumption on the distribution  $D$  of the features.
- We collect corresponding labels  $y_1, \dots, y_n$ .
- Assume there exists some true classifier  $h^*$ .
- Let  $\mathcal{H}$  be the set of all hypotheses.

The goal of PAC learning is  $(\varepsilon, \delta)$ -PAC, which is defined as follows.

**Definition.** An  $(\varepsilon, \delta)$ -PAC learning algorithm refers to an algorithm that picks an hypothesis  $h \in \mathcal{H}$  after observing training data  $\{x_i, y_i\}_{i=1}^n$ , where the hypothesis  $h \in \mathcal{H}$  satisfies

$$\text{err}_D(\hat{h}) = \mathbb{P}_{x \sim D} [\hat{h}(x) \neq h^*(x)] < \varepsilon.$$

with probability at least  $1 - \delta$ , where the probability is with respect to the randomness of the training set.

Putting it more concretely, suppose  $g$  is a point estimator (the algorithm), then we want

$$\mathbb{P}_{x_i \sim D} [\text{err}_D(g(x_1, \dots, x_n)) < \varepsilon] \geq 1 - \delta.$$

Through the  $\varepsilon$ - $\delta$  definition of limit, it is not hard to notice the similarity between this definition and **convergence in probability**. As we observe more and more data ( $n \rightarrow \infty$ ), we want the hypothesis produce by the algorithm to converge to the true classifier **in probability**. This is exactly the definition of a **consistent** point estimator in the previous section.

We also say the algorithm is a PAC learner if it uses  $n$  samples and the running time is at most  $\text{poly}(d, \frac{1}{\varepsilon}, \log \frac{1}{\delta}, \text{bits}(h^*))$ , but in this section we do not focus on the runtime of the algorithm.

There are two types of PAC learning – realizable PAC learning, in which case  $h^* \in \mathcal{H}$ , and agnostic PAC learning, in which case  $h^* \notin \mathcal{H}$ . We first discuss realizable PAC learning, and we will find out agnostic PAC learning is a more general setting but than PAC learning but a natural extension.

#### Realizable learning

For realizable PAC learning, we present a simple algorithm:

**Algorithm** (Consistent Learner). Pick any  $\hat{h} \in \mathcal{H}$  such that  $h(x_i) = y_i$  for all  $1 \leq i \leq n$ .

Now we analyze this algorithm in terms of the sample size needed to produce a desired hypothesis.

**Theorem.** Over the dataset of  $n$  iid samples, the consistent learning algorithm produces  $\hat{h}$  such that  $\text{err}_D(\hat{h}) > \varepsilon$  with probability at most  $|\mathcal{H}| e^{-n\varepsilon}$ .

*Proof.* Suppose  $h \in \mathcal{H}$  is such that  $\text{err}_D(h) = \mathbb{P}_{x \sim D} [h(x) \neq h^*(x)] > \varepsilon$ . For such an  $h$  and some data  $x_i$ , we have

$$\mathbb{P}_{x_i \sim D} [h(x_i) = y_i = h^*(x_i)] < 1 - \varepsilon.$$

Since  $\{x_i, y_i\}_{i=1}^n$  are iid, we have

$$\mathbb{P}_{x_{1:n} \sim D} [h(x_i) = h^*(x_i) \text{ for all } 1 \leq i \leq n] < (1 - \varepsilon)^n.$$

Note that our consistent learner do not make any mistake on the training set  $\{x_i, y_i\}_{i=1}^n$

$$\mathbb{P}_{x_{1:n} \sim D} [h = \hat{h}] < (1 - \varepsilon)^n.$$

It follows that

$$\begin{aligned}\mathbb{P}_{x_{1:n} \sim D}[\text{err}_D(\hat{h}) > \varepsilon] &\leq \sum_{h \in \mathcal{H} : \text{err}_D(h) > \varepsilon} \mathbb{P}[\hat{h} = h] \\ &\leq |\mathcal{H}| (1 - \varepsilon)^n \\ &\leq |\mathcal{H}| e^{-n\varepsilon},\end{aligned}$$

where in the last step we used the inequality  $(1 - u)^n \leq e^{-nu}$ . This completes the proof.  $\square$

**Corollary.** If we want  $\mathbb{P}_{x_{1:n} \sim D}[\text{err}_D(\hat{h}) > \varepsilon] \leq \delta$ , we must have

$$n \geq \frac{\ln(|\mathcal{H}|/\delta)}{\varepsilon}.$$

Technically speaking the implication should be the other direction, but this bound for sample size  $n$  **guarantees**  $(\varepsilon, \delta)$ -PAC.

To illustrate how we should think of  $|\mathcal{H}|$ , we present an example.

**Example.** Consider the binary half-spaces hypotheses:

$$\mathcal{H} = \{h_w(x) = \text{sign}(\langle w, x \rangle) : w_i \in \{-1, 1\}\}.$$

In this case  $|\mathcal{H}| = 2^d$ , so  $\ln |\mathcal{H}| = \Theta(d)$ .

We should always think of  $|\mathcal{H}|$  as exponential with respect to the dimension, so  $\ln |\mathcal{H}|$  is linear with respect to dimension.

Now we move on to the setting of agnostic learning.

### Agnostic learning

In agnostic learning, again we collect features  $x_1, \dots, x_n \sim D$  iid, and labels  $y_1, \dots, y_n$ . This forms a data set  $S_n = \{x_i, y_i\}_{i=1}^n$ . The goal is to use this dataset  $S_n$  to produce an hypothesis  $\hat{h}$  such that

$$\text{reg}_{D, \mathcal{H}}(\hat{h}) = \text{err}_D(\hat{h}) - \min_{h \in \mathcal{H}} \text{err}_D(h) < \varepsilon$$

with probability at least  $1 - \delta$ , where the probability is with respect to the randomness of the dataset  $S_n$ .

For this setting, we present an algorithm called the empirical loss minimizer (ERM).

**Algorithm** (empirical loss minimizer). Choose  $\hat{h} \in \mathcal{H}$  by minimizing the  $\{0, 1\}$  loss:

$$\hat{h} = \underset{h \in \mathcal{H}}{\text{argmin}} \sum_{i=1}^n 1\{h(x_i) \neq y_i\}.$$

In a more general setting, suppose  $\ell(\cdot, \cdot)$  be any loss function bounded in  $[0, 1]$ , choose  $\hat{h} \in \mathcal{H}$  by minimizing the loss:

$$\hat{h} = \underset{h \in \mathcal{H}}{\text{argmin}} \sum_{i=1}^n \ell(h(x_i), y_i).$$

We next prove a similar relation between the sample size and the performance of the hypothesis, evaluated in terms of risk. We first present the definition of risk.

**Definition** (Risk). The risk of a hypothesis  $h$  is defined as

$$\text{risk}_D(h) = \mathbb{E}_{x \sim D}[\ell(h(x), y)].$$

**Theorem.** Over the dataset of  $n$  iid samples, the ERM algorithm produces  $\hat{h}$  such that

$$\text{reg}_{D,\mathcal{H}}(\hat{h}) = \text{risk}_D(\hat{h}) - \min_{h \in \mathcal{H}} \text{risk}_D(h) \leq \sqrt{\frac{\ln(|\mathcal{H}|/\delta)}{n}}.$$

with probability at least  $1 - \delta$ .

*Proof.* First we define

$$h^{\text{opt}} = \underset{h \in \mathcal{H}}{\text{argmin}} \text{risk}_D(h).$$

Note that  $\text{reg}_{D,\mathcal{H}}(h^{\text{opt}}) = 0$ . Define also the estimated risk of hypothesis  $h$  for data set  $S$ :

$$\widehat{\text{risk}}_S(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i).$$

Note that  $\hat{h} = \underset{h \in \mathcal{H}}{\text{argmin}} \widehat{\text{risk}}_S(h)$  and  $\mathbb{E}[\widehat{\text{risk}}_S(h)] = \text{risk}_D(h)$ , where the expected value is with respect to the randomness of the dataset. Note that

$$\begin{aligned} \text{reg}_{D,\mathcal{H}}(\hat{h}) &= \text{risk}_D(\hat{h}) - \text{risk}_D(h^{\text{opt}}) \\ &= (\text{risk}_D(\hat{h}) - \widehat{\text{risk}}_S(\hat{h})) + (\widehat{\text{risk}}_S(\hat{h}) - \widehat{\text{risk}}_S(h^{\text{opt}})) - (\text{risk}_D(h^{\text{opt}}) - \widehat{\text{risk}}_S(h^{\text{opt}})) \quad (*) \\ &\leq (\text{risk}_D(\hat{h}) - \widehat{\text{risk}}_S(\hat{h})) - (\text{risk}_D(h^{\text{opt}}) - \widehat{\text{risk}}_S(h^{\text{opt}})), \end{aligned}$$

where the inequality is because our algorithm always gives  $\widehat{\text{risk}}_S(\hat{h}) - \widehat{\text{risk}}_S(h^{\text{opt}}) \leq 0$ .

For the first term, we have for any  $h \in \mathcal{H}$ ,

$$\begin{aligned} \widehat{\text{risk}}_S(h) - \text{risk}_D(h) &= \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i) - \text{risk}_D(h) \\ &= \frac{1}{n} \sum_{i=1}^n (\ell(h(x_i), y_i) - \mathbb{E}_{x,y \sim D}[\ell(h(x), y)]). \end{aligned}$$

Define now random variable  $Z_i := \ell(h(x_i), y_i) - \mathbb{E}_{x,y \sim D}[\ell(h(x), y)]$ . Notice that  $Z_i$  iid,  $\mathbb{E}[Z_i] = 0$ , and  $|Z_i| \leq 1$ . Therefore, we can utilize the Hoeffding bound to get

$$\mathbb{P} \left[ \frac{1}{n} \sum_{i=1}^n Z_i > \varepsilon \right] \leq e^{-2n\varepsilon^2} \quad \text{and} \quad \mathbb{P} \left[ \frac{1}{n} \sum_{i=1}^n Z_i < -\varepsilon \right] \leq e^{-2n\varepsilon^2}$$

This implies that for any fixed  $h \in \mathcal{H}$ ,

$$\mathbb{P} \left[ \left| \widehat{\text{risk}}_S(h) - \text{risk}_D(h) \right| > \varepsilon \right] \leq 2e^{-2n\varepsilon^2}.$$

It follows from the union bound that

$$\mathbb{P} \left[ \exists h \in \mathcal{H} : \left| \widehat{\text{risk}}_S(h) - \text{risk}_D(h) \right| > \varepsilon \right] \leq 2|\mathcal{H}| e^{-2n\varepsilon^2}.$$

Notice that if  $\left| \widehat{\text{risk}}_S(h) - \text{risk}_D(h) \right| \leq \varepsilon$  for all  $h \in \mathcal{H}$ , we must have  $\text{reg}_{D,\mathcal{H}}(\hat{h}) \leq 2\varepsilon$  by inequality (\*). above. Therefore,  $\text{reg}_{D,\mathcal{H}}(\hat{h}) \leq 2\varepsilon$  with probability at least  $1 - 2|\mathcal{H}| e^{-2n\varepsilon^2}$ .  $\square$

**Corollary.** Let  $\delta > 0$  be given. Then, setting  $n = \frac{1}{2\varepsilon^2} \ln(2|\mathcal{H}|/\delta)$ , we obtain that  $\text{reg}_{D,\mathcal{H}}(\hat{h}) \leq 2\varepsilon$  with probability at least  $1 - \delta$ . Therefore, if we want  $\text{reg}_{D,\mathcal{H}}(\hat{h}) \leq \varepsilon$  with probability at least  $1 - \delta$ , we must have

$$n \geq \frac{2 \ln(2|\mathcal{H}|/\delta)}{\varepsilon^2}.$$

### Comparison

As a comparison between realizable learning and agnostic learning, we can see the sample size needed to achieve the same  $(\varepsilon, \delta)$ -PAC bound is

$$n = \Theta\left(\frac{\ln(|\mathcal{H}|/\delta)}{\varepsilon}\right)$$

for realizable learning, and

$$n = \Theta\left(\frac{\ln(|\mathcal{H}|/\delta)}{\varepsilon^2}\right)$$

for agnostic learning. Since  $\varepsilon < 1$ , we can deduce that agnostic learning needs **more sample** for the estimator to achieve the same  $(\varepsilon, \delta)$ -PAC bound.



## 2 Mixture models and EM

### 2.1 Kullback-Leibler (KL) Divergence

In this section, we adopt the convention that  $0 \log 0 = 0$ .

**Definition** (KL divergence). The Kullback-Leibler (KL) divergence of two distributions  $P(X)$  and  $Q(X)$  over the outcome space  $X$  is defined as follows:

$$\text{KL}(P\|Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)}.$$

To understand the significance of KL-divergence better, we first discuss some related concepts in information theory, starting with the definition of **entropy**.

**Definition.** The entropy of a distribution  $P(X)$  is defined as

$$H(P) = - \sum_{x \in X} P(x) \log P(x)$$

Intuitively, entropy measures how dispersed a probability distribution is. For example, a uniform distribution is considered to have very high entropy (i.e. a lot of uncertainty), whereas a distribution that assigns all its mass on a single point is considered to have zero entropy (i.e. no uncertainty). Notably, it can be shown that among continuous distributions over  $\mathbb{R}$ , the Gaussian distribution  $\mathcal{N}(\mu, \sigma^2)$  has the highest entropy (highest uncertainty) among all possible distributions that have the given mean  $\mu$  and variance  $\sigma^2$ .

To further solidify our intuition, we present motivation from communication theory. Suppose we want to communicate from a source to a destination, and our messages are always (a sequence of) discrete symbols over space  $X$  (for example,  $X$  could be letters  $\{a, b, \dots, z\}$ ). We want to construct an encoding scheme for our symbols in the form of sequences of binary bits that are transmitted over the channel. Further, suppose that in the long run the frequency of occurrence of symbols follow a probability distribution  $P(X)$ . This means, in the long run, the fraction of times the symbol  $x$  gets transmitted is  $P(x)$ .

A common desire is to construct an encoding scheme such that the average number of bits per symbol transmitted remains as small as possible. Intuitively, this means we want very frequent symbols to be assigned to a bit pattern having a small number of bits. Likewise, because we are interested in reducing the average number of bits per symbol in the long term, it is tolerable for infrequent words to be assigned to bit patterns having a large number of bits, since their low frequency has little effect on the long term average. The encoding scheme can be as complex as we desire, for example, a single bit could possibly represent a long sequence of multiple symbols (if that specific pattern of symbols is very common). The entropy of a probability distribution  $P(X)$  is its optimal bit rate, i.e., the lowest average bits per message that can possibly be achieved if the symbols  $x \in X$  occur according to  $P(X)$ . It does not specifically tell us how to construct that optimal encoding scheme. It only tells us that no encoding can possibly give us a lower long term bits per message than  $H(P)$ .

To see a concrete example, suppose our messages have a vocabulary of  $K = 32$  symbols, and each symbol has an equal probability of transmission in the long term (i.e, uniform probability distribution). An encoding scheme that would work well for this scenario would be to have  $\log_2 K$  bits per symbol, and assign each symbol some unique combination of the  $\log_2 K$  bits. In fact, it turns out that this is the most efficient encoding one can come up with for the uniform distribution scenario.

It may have occurred to you by now that the long term average number of bits per message depends only on the frequency of occurrence of symbols. The encoding scheme of scenario A can in theory be reused in scenario B with a different set of symbols (assume equal vocabulary size for simplicity), with the same long term efficiency, as long as the symbols of scenario B follow the same probability distribution as the symbols of scenario A. It might also have occurred to you, that reusing the encoding scheme designed to be optimal for scenario A, for messages in scenario B having a different probability of symbols, will always be suboptimal for scenario B. To be clear, we do not need know what the specific optimal schemes are

in either scenarios. As long as we know the distributions of their symbols, we can say that the optimal scheme designed for scenario A will be suboptimal for scenario B if the distributions are different.

Concretely, if we reuse the optimal scheme designed for a scenario having symbol distribution  $Q(X)$ , into a scenario that has symbol distribution  $P(X)$ , the long term average number of bits per symbol achieved is called the cross entropy, denoted by  $H(P, Q)$ :

**Definition.** The cross-entropy of two distributions  $P(X)$  and  $Q(X)$  is defined as

$$H(P, Q) = - \sum_{x \in X} P(x) \log Q(x)$$

To recap, the entropy  $H(P)$  is the best possible long term average bits per message (optimal) that can be achieved under a symbol distribution  $P(X)$  by using an encoding scheme (possibly unknown) specifically designed for  $P(X)$ . The cross entropy  $H(P, Q)$  is the long term average bits per message (suboptimal) that results under a symbol distribution  $P(X)$ , by reusing an encoding scheme (possibly unknown) designed to be optimal for a scenario with symbol distribution  $Q(X)$ .

Now, KL divergence is the penalty we pay, as measured in average number of bits, for using the optimal scheme for  $Q(X)$ , under the scenario where symbols are actually distributed as  $P(X)$ . It is straightforward to see this:

$$\begin{aligned} \text{KL}(P\|Q) &= \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} \\ &= \sum_{x \in X} P(x) \log P(x) - \sum_{x \in X} P(x) \log Q(x) \\ &= H(P, Q) - H(P). \quad (\text{difference in average number of bits}) \end{aligned}$$

If the cross entropy between  $P$  and  $Q$  is zero (and hence  $\text{KL}(P\|Q) = 0$ ) then it necessarily means  $P = Q$ . In Machine Learning, it is a common task to find a distribution  $Q$  that is “close” to another distribution  $P$ . To achieve this, we use  $\text{KL}(P\|Q)$  to be the loss function to be optimized. As we will see in this below, Maximum Likelihood Estimation, which is a commonly used optimization objective, turns out to be equivalent minimizing KL divergence between the training data (i.e. the empirical distribution over the data) and the model.

Now we present some useful properties of KL divergence.

**Theorem** (Non-negativity). For any distribution  $P$  and  $Q$ , we have

$$\text{KL}(P\|Q) \geq 0,$$

and  $\text{KL}(P\|Q) = 0$  if and only if  $P = Q$  (almost everywhere).

*Proof.* By definition,

$$\text{KL}(P\|Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} = - \sum_{x \in X} P(x) \log \frac{Q(x)}{P(x)}.$$

Since  $-\log x$  is strictly convex, by Jensen’s inequality, we have

$$\text{KL}(P\|Q) = - \sum_{x \in X} P(x) \log \frac{Q(x)}{P(x)} \geq - \log \sum_{x \in X} P(x) \frac{Q(x)}{P(x)} = 0.$$

When the equality holds,

$$\log \frac{Q(x)}{P(x)} = 0$$

almost everywhere. That is,  $Q = P$  almost everywhere. This completes the proof.  $\square$

**Definition** (KL divergence of conditional distribution). The KL divergence between 2 conditional distributions  $P(X | Y)$ ,  $Q(X | Y)$  is defined as follows:

$$\text{KL}(P(X | Y) \| Q(X | Y)) = \sum_y P(y) \left( \sum_x P(x | y) \log \frac{P(x | y)}{Q(x | y)} \right).$$

This can be thought of as the expected KL divergence between the corresponding conditional distributions on  $x$ . That is, between  $P(X | Y = y)$  and  $Q(X | Y = y)$ , where the expectation is taken over the random  $y$ .

**Theorem** (Chain rule for KL divergence). The following equality holds:

$$\text{KL}(P(X, Y) \| Q(X, Y)) = \text{KL}(P(X) \| Q(X)) + \text{KL}(P(Y | X) \| Q(Y | X)).$$

*Proof.*

$$\begin{aligned} \text{LHS} &= \sum_x \sum_y P(x, y) \log \frac{P(x, y)}{Q(x, y)} \\ &= \sum_x \sum_y P(y | x) P(x) \left[ \log \frac{P(y | x)}{Q(y | x)} + \log \frac{P(x)}{Q(x)} \right] \\ &= \sum_x \sum_y P(y | x) P(x) \log \frac{P(y | x)}{Q(y | x)} + \sum_x P(x) \log \frac{P(x)}{Q(x)} \sum_y P(y | x) \\ &= \sum_x \sum_y P(y | x) P(x) \log \frac{P(y | x)}{Q(y | x)} + \sum_x P(x) \log \frac{P(x)}{Q(x)} \\ &= \text{KL}(P(X) \| Q(X)) + \text{KL}(P(Y | X) \| Q(Y | X)) \\ &= \text{RHS}. \end{aligned}$$

□

## 2.2 The EM Algorithm in General

### 3 Approximate inference

A central task in the application of probabilistic models is the evaluation of the posterior distribution  $p(\mathbf{Z}|\mathbf{X})$  of the latent variables  $\mathbf{Z}$  given the observed (visible) data variables  $\mathbf{X}$ , and the evaluation of expectations computed with respect to this distribution. For many models of practical interest, it will be infeasible to evaluate the posterior distribution or indeed to compute expectations with respect to this distribution.

In this chapter, we introduce a range of deterministic approximation schemes, some of which scale well to large applications. These are based on analytical approximations to the posterior distribution, for example by assuming that it factorizes in a particular way or that it has a specific parametric form such as a Gaussian. As such, they can never generate exact results, and so their strengths and weaknesses are complementary to those of sampling methods.