

Machine Learning

Notes taken by Runqiu Ye
Carnegie Mellon University

Spring 2025

Contents

1	Probability and statistical inference	3
1.1	Probability	3
1.2	Statistical inference	4
1.3	PAC learning	5
1.3.1	Realizable learning	5
1.3.2	Agnostic learning	6
1.4	Infinite hypotheses space and VC dimension	8
2	Mixture models and EM	9
2.1	Kullback-Leibler (KL) Divergence	9
2.2	The EM Algorithm in General	11
3	Approximate inference	13

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 μ and variance σ^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 θ . A point estimator $\hat{\theta}_n$ of the parameter θ 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 θ 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 θ 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

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

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

Warning! In the above definition, C_n is random and θ 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 i.i.d. 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 (ε, δ) -PAC, which is defined as follows.

Definition. An (ε, δ) -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 ε - δ 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.

1.3.1 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 i.i.d. 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 i.i.d., 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

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

Technically speaking the implication should be the other direction, but this bound for sample size n **guarantees** (ε, δ) -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 $\log |\mathcal{H}| = \Theta(d)$.

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

Now we move on to the setting of agnostic learning.

1.3.2 Agnostic learning

In agnostic learning, again we collect features $x_1, \dots, x_n \sim D$ i.i.d., 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 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:

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

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 i.i.d. samples, the ERM algorithm produces \hat{h} such that

$$\text{reg}_{D, \mathcal{H}}(\hat{h}) = \text{risk}_D(\hat{h}) - \min_{h \in \mathcal{H}} \text{risk}_D(h) \leq 2\varepsilon.$$

with probability at least $1 - 2|\mathcal{H}|e^{-2n\varepsilon^2}$.

Proof. First we define

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

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

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

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

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

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

For any $h \in \mathcal{H}$, we have

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

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 i.i.d., $\mathbb{E}[Z_i] = 0$, and $|Z_i| \leq 1$. Therefore, we can utilize the **Hoeffding bound** to get

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

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

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

It follows from the union bound that

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

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

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

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

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

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

for realizable learning, and

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

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

1.4 Infinite hypotheses space and VC dimension

In last section, we discussed PAC learning with finite hypotheses space \mathcal{H} . However, we often need to deal with situations where the hypotheses space is infinite. Consider the following example.

Example. Consider the binary half space hypotheses:

$$\mathcal{H} = \{h_w(x) = \text{sign}(\langle w, x \rangle) : w \in \mathbb{R}^d\}.$$

Then, $|\mathcal{H}| = \infty$ since each $w \in \mathbb{R}^d$ gives a different decision boundary.

However, on some fixed dataset S , not all of the hypotheses in the hypotheses space \mathcal{H} gives a different prediction. This leads to the following definition and theorem, which characterize the “effective size” of the hypotheses space.

Definition. Let S be a dataset with n points $\{x_1, \dots, x_n\}$. Define $|\mathcal{H}(S)|$ as the number of distinct labeling $\{(x_i, y_i)\}_{i=1}^n$, where each $y_i = h(x_i)$ for a fixed $h \in \mathcal{H}$.

With this definition, we can derive a similar confidence bound for our learned hypotheses.

Theorem. Define the effective size

$$\mathcal{H}[n] = \sup_{x_1, \dots, x_n} |\mathcal{H}(\{x_1, \dots, x_n\})|,$$

and let \hat{h}_n denote the consistency learner from n examples. Then we have

$$\mathbb{P}_{x \sim D}[\text{err}_D(\hat{h}) \geq \varepsilon] \leq \delta,$$

where

$$\varepsilon = O\left(\frac{\log \mathcal{H}[n] + \log(1/\delta)}{n}\right).$$

Example. Consider the threshold hypotheses space:

$$\mathcal{H} = \{h_w(x) = 2 \cdot \mathbf{1}\{x > w\} - 1 : w \in \mathbb{R}\}.$$

Then it is easy to see that $\mathcal{H}[n] = n$. It follows that

$$\varepsilon = O\left(\frac{\log \mathcal{H}[n] + \log(1/\delta)}{n}\right) = O\left(\frac{\log n + \log(1/\delta)}{n}\right),$$

which approaches to 0 as $n \rightarrow \infty$.

Another way to characterize the “effective size” of the hypotheses space is **VC dimension**.

Definition (VC dimension). For a hypotheses space \mathcal{H} , we say \mathcal{H} **shatters** a dataset $\{x_1, \dots, x_n\}$ if for any choice of labels $\{-1, 1\}^n$, there exists $h \in \mathcal{H}$ such that $h(x_i) = y_i$ for all $1 \leq i \leq n$.

The **VC dimension** of a hypotheses space \mathcal{H} is defined as the largest n such that there exists a dataset S of size n such that \mathcal{H} shatters S . Denote this as $\text{VC}(\mathcal{H}) = n$.

The VC dimension of a hypotheses space \mathcal{H} is related to $\mathcal{H}[n]$ by the following lemma by Sauer:

Theorem (Sauer’s Lemma). For any hypotheses class \mathcal{H} , we have

$$\log \mathcal{H}[n] \leq O(\text{VC}(\mathcal{H}) \cdot \log n).$$

Recalling our bound for realizable learning, we get (ε, δ) -PAC with

$$\varepsilon = O\left(\frac{\log \mathcal{H}[n] + \log(1/\delta)}{n}\right) = O\left(\frac{\text{VC}(\mathcal{H}) \cdot \log n + \log(1/\delta)}{n}\right).$$

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 μ and variance σ^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

Consider a probabilistic model in which we collectively denote all of the observed variables by X and all of the hidden variables by Z . The joint distribution $p(X, Z | \theta)$ is governed by a set of parameters denoted θ . Our goal is to maximize the likelihood function that is given by

$$p(X | \theta) = \sum_Z p(X, Z | \theta).$$

Assume direct optimization of $p(X | \theta)$ is difficult but optimization of $p(X, Z | \theta)$ is easy. Introduce distribution $q(Z)$ for the latent variable Z , we then have

$$\log p(X | \theta) = \mathcal{L}(q, \theta) + \text{KL}(q \| p),$$

where

$$\begin{aligned} \mathcal{L}(q, \theta) &= \sum_Z q(Z) \log \frac{p(X, Z | \theta)}{q(Z)}, \\ \text{KL}(q \| p) &= - \sum_Z q(Z) \log \frac{p(Z | X, \theta)}{q(Z)}. \end{aligned}$$

Note that $\text{KL}(q \| p)$ is the KL divergence of $q(Z)$ and $p(Z | X, \theta)$ so $\text{KL}(q \| p) \geq 0$. Therefore, $\mathcal{L}(q, \theta) \leq \log p(X | \theta)$. In other words, $\mathcal{L}(q, \theta)$ is a lower bound for $\log p(X | \theta)$.

Therefore, for the EM algorithm, suppose currently the model parameter is θ_t . In the E step, the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to $q(Z)$ while holding θ fixed. Note that $\log p(X | \theta_t)$ is not

dependent on $q(Z)$, so the lower bound will be maximized when $\text{KL}(q\|p) = 0$. This is equivalent to when $q(Z) = p(Z | X, \theta)$. Therefore, in the E step we simply set

$$q(Z) = p(Z | X, \theta_t).$$

In the subsequent M step, the distribution $q(Z)$ is fixed and the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to θ to give some new model parameter θ_{t+1} . Because the distribution q is determined using the old parameter values rather than the new values and is held fixed during the M step, it will not equal the new posterior distribution $p(Z | X, \theta_{t+1})$, and hence there will be a nonzero KL divergence. Substituting the result from E step, we have the following expression for the lower bound of M step:

$$\begin{aligned} \mathcal{L}(q, \theta) &= \sum_Z p(Z | X, \theta_t) \log p(X, Z | \theta) - \sum_Z p(Z | X, \theta_t) \log p(Z | X, \theta_t) \\ &= Q(\theta_t, \theta) + \text{const}, \end{aligned}$$

where

$$Q(\theta_t, \theta) = \sum_Z p(Z | X, \theta_t) \log p(X, Z | \theta)$$

and the constant is simply the negative entropy of the q distribution. Note that the variable θ over which we are optimizing appears only inside the logarithm.

3 Approximate inference

A central task in the application of probabilistic models is the evaluation of the posterior distribution $p(\mathbf{Z} \mid \mathbf{X})$ 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.