

# MSA101/MVE187 2020 Lecture 9.1

## From Information to Kullback-Leibler divergence

Petter Mostad

Chalmers University

September 23, 2020

# The information of an event

We assume given a probability mass function  $\pi(x)$  on a finite set.

- ▶ We want to define the “information”  $h(x)$  in an event  $x$ .

Requirements:

- ▶ An event with probability 1 should have zero information.
- ▶ The information should increase with decreasing probability  $\pi(x)$ .
- ▶ The information in two independent events should be the sum of the information in each.
- ▶ We define  $h(x) = -\log(\pi(x))$ .
- ▶ When using the base 2 logarithm  $\log_2$ , information is measured in “bits”. We however use the natural logarithm.

# Expected information: Entropy

- ▶ Define the entropy  $H[X]$  of the random variable  $X$  as the expected information:

$$H[X] = \sum_x h(x)\pi(x) = - \sum_x \pi(x) \log(\pi(x))$$

- ▶ Note:  $H[X]$  is always non-negative.
- ▶ Example: A uniform distribution on  $n$  values has entropy  $\log n$ . This is the largest entropy possible for a distribution on  $n$  values.
- ▶ Shannon's coding theorem: The entropy (using  $\log_2$ ) is a lower bound on the expected number of bits needed to transfer the information from  $X$ .

# (Differential) entropy for continuous distributions

- ▶ For any random variable  $X$ , its (differential) entropy is defined as

$$H[X] = \mathbb{E}[-\log(\pi(x))] = - \int_x \log(\pi(x))\pi(x) dx$$

- ▶  $H[X]$  may now be negative.
- ▶ Example: Assume  $X \sim \text{Normal}(\mu, \sigma^2)$ . Then

$$\begin{aligned} \mathbb{E}[-\log(\pi(x))] &= \mathbb{E}\left[-\log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \frac{1}{2\sigma^2}(x - \mu)^2\right] \\ &= \frac{1}{2}\log(2\pi\sigma^2) + \frac{1}{2\sigma^2}\mathbb{E}[(x - \mu)^2] = \frac{1}{2}\log(2\pi\sigma^2) + \frac{1}{2}. \end{aligned}$$

- ▶ In fact, among all random variables  $X$  with  $\mathbb{E}[X] = \mu$  and  $\text{Var}[X] = \sigma^2$ , the normal has the largest entropy.

# Conditional entropy and mutual information

- ▶ The conditional entropy is defined as

$$H[Y|X] = \int \left[ \int \pi(y | x) (-\log(\pi(y | x))) dy \right] \pi(x) dx$$

- ▶ Show that

$$H[X, Y] = H[Y|X] + H[X].$$

- ▶ The mutual information is defined as

$$I[X, Y] = - \int \int \pi(x, y) \log \left( \frac{\pi(x)\pi(y)}{\pi(x, y)} \right) dx dy$$

- ▶ Show that

$$I[X, Y] = H[X] + H[Y] - H[X, Y]$$

# The Kullback-Leibler divergence (relative entropy)

- ▶ For two densities  $p(x)$  and  $q(x)$  we define the Kullback-Leibler divergence from  $p$  to  $q$  as

$$\text{KL}[p||q] = - \int p(x) \log \left( \frac{q(x)}{p(x)} \right) dx$$

- ▶ Note that  $\text{KL}[p||q]$  is generally different from  $\text{KL}[q||p]$ .
- ▶ However, it has the distance property that  $\text{KL}[p||q] \geq 0$  always, while  $\text{KL}[p||q] = 0$  if and only if  $p = q$ .
- ▶ The standard proof uses *Jensen's inequality*.
- ▶ Jensen's inequality: If a function  $\psi$  is *convex*, then  $\psi(\text{E}[X]) \leq \text{E}[\psi(X)]$ .

# The KL divergence

- ▶ Note that

$$\text{KL}(\pi(x, y) || \pi(x)\pi(y)) = I[X, Y]$$

- ▶ Note that

$$\text{KL}[p || q] = \mathbb{E}_p[-\log(q(x))] - H_p[X]$$

where  $X$  is a random variable with density  $p(x)$ .

- ▶ EXAMPLE: Assume  $X \sim \text{Normal}(\mu_X, \sigma_X^2)$  and  $Y \sim \text{Normal}(\mu_Y, \sigma_Y^2)$ .

Show by direct computation that

$$\text{KL}[\pi_X || \pi_Y] = \frac{1}{2} \log(2\pi\sigma_Y^2) + \frac{\sigma_X^2}{2\sigma_Y^2} + \frac{1}{2\sigma_Y^2}(\mu_X - \mu_Y)^2 - \frac{1}{2} \log(2\pi\sigma_X^2) - \frac{1}{2}.$$

We see how the result is zero when the two distributions are identical.

We see how  $\text{KL}[\pi_X || \pi_Y] \neq \text{KL}[\pi_Y || \pi_X]$  in general.

# MSA101/MVE187 2020 Lecture 9.2

## The EM algorithm

Petter Mostad

Chalmers University

September 24, 2020



# Maximum posterior (MAP)

- ▶ The advantages and disadvantages of only finding the Maximal APosteriori (MAP): The value  $\hat{\theta}$  that maximizes the posterior  $\pi(\theta \mid \text{data})$ .
- ▶ When the prior is flat,  $\pi(\theta) \propto 1$ , this corresponds to finding the maximum likelihood (ML) estimate for  $\theta$ .
- ▶ This should be easy to compute when  $\theta$  consists of all unknown variables: Just differentiate  $\log(\pi(\theta \mid \text{data}))$ , i.e. differentiate  $\log(\pi(\text{data} \mid \theta)\pi(\theta))$ .
- ▶ Much harder if the model also contains other unknown variables  $Z$ : Then  $\pi(\theta \mid \text{data})$  is the marginal of  $\pi(\theta, Z \mid \text{data})$  and much harder to maximize.
- ▶ The Expectation-Maximization (EM) algorithm comes to the rescue...

# The EM algorithm

- ▶ We want to find the  $\theta$  maximizing the posterior  $\pi(\theta | x)$ ; i.e., maximizing

$$\log(\pi(x | \theta)\pi(\theta)) = \log(\pi(x | \theta)) + \log(\pi(\theta))$$

- ▶ Assume we have a joint model  $\pi(x, z | \theta)$  which includes augmented data  $z$ . We may then write, for any density  $q(z)$ ,

$$\log(\pi(x | \theta)) + \log(\pi(\theta)) = \text{KL}(q || \pi_z) + \mathcal{L}(q, \theta) + \log(\pi(\theta)) \quad (1)$$

where

$$\mathcal{L}(q, \theta) = \int q(z) \log \left( \frac{\pi(x, z | \theta)}{q(z)} \right) dz$$

and

$$\text{KL}(q || \pi_z) = - \int q(z) \log \left( \frac{\pi_z(z | x, \theta)}{q(z)} \right) dz$$

# The EM algorithm, cont.

- ▶ Fix  $q(z) = \pi_z(z \mid x, \theta^{old})$  for some value  $\theta^{old}$ .
- ▶ With this  $q(z)$ ,  $KL(q \parallel \pi_z)$  will be zero when  $\theta = \theta^{old}$  and positive for other  $\theta$ 's. THUS: If we find  $\theta^{new}$  maximizing  $\mathcal{L}(q, \theta) + \log(\pi(\theta))$ , so that  $\mathcal{L}(q, \theta^{new}) + \log(\pi(\theta^{new})) > \mathcal{L}(q, \theta^{old}) + \log(\pi(\theta^{old}))$ , replacing  $\theta^{old}$  with  $\theta^{new}$  will increase the right side of Equation 1, and thus also the left side.
- ▶ Set  $\theta^{old}$  to the value  $\theta^{new}$  and start again from the first step above. Continue until convergence.
- ▶ Note that maximizing  $\mathcal{L}(q, \theta) + \log(\pi(\theta))$  is the same as maximizing

$$\int q(z) \log(\pi(x, z \mid \theta)) dz + \log(\pi(\theta))$$

where the left term is the expected full loglikelihood, taking the expectation over the density  $q(z) = \pi_z(z \mid x, \theta^{old})$ .

- ▶ E-step: Computing the expectation above. M-step: Maximizing.

# MSA101/MVE187 2020 Lecture 9.3

## The EM algorithm: A toy example

Petter Mostad

Chalmers University

October 8, 2020

## A toy example

We have data  $x_1, \dots, x_n$ , where we assume the following model, with a single parameter  $\mu$ : With probability 0.5,  $x_i \sim \text{Normal}(0, 1)$  and with probability 0.5,  $x_i \sim \text{Normal}(\mu, 1)$ . We assume a flat prior on  $\mu$ .

- ▶ The likelihood can be written as

$$\pi(x_1, \dots, x_n \mid \mu) = \prod_{i=1}^n (0.5 \cdot \text{Normal}(x_i; 0, 1) + 0.5 \cdot \text{Normal}(x_i; \mu, 1))$$

- ▶ With the loglikelihood programmed numerically, we may for example
  - ▶ Optimize to find the maximum likelihood estimate  $\hat{\mu}$  for  $\mu$ .
  - ▶ Simulate from the posterior, using, e.g., discretization.
- ▶ Instead, we may introduce *augmented* data  $z_1, \dots, z_n$ , where each  $z_i$  has value 0 or 1, so that  $z_i \sim \text{Bernoulli}(0.5)$  and  $x_i \mid z_i \sim \text{Normal}(z_i \cdot \mu, 1)$ . The full joint density may be written as

$$\pi(x_1, \dots, x_n, z_1, \dots, z_n, \mu) \propto \prod_{i=1}^n \pi(x_i \mid z_i, \mu) = \prod_{i=1}^n \text{Normal}(x_i; z_i \cdot \mu, 1)$$

- ▶ The augmented model may be used both for simulation (using Gibbs sampling (inefficient)) and for finding the  $\mu$  maximizing the posterior using the EM-algorithm.

## A toy example: Using the EM algorithm

- First, find the complete data logposterior (which in our case is the same as the loglikelihood). It is (up to a constant)

$$l(\mu) = \sum_{i=1}^n -\frac{1}{2}(x_i - z_i \cdot \mu)^2$$

- Then, for a fixed value  $\mu = \mu^{old}$ , find the distribution  $z_i \mid x_i, \mu^{old}$ :

$$\pi(x_1, \dots, x_n, \dots, z_i = 0, \dots, \mu^{old}) = K \cdot \text{Normal}(x_i; 0, 1)$$

$$\pi(x_1, \dots, x_n, \dots, z_i = 1, \dots, \mu^{old}) = K \cdot \text{Normal}(x_i; \mu^{old}, 1)$$

Normalizing the distribution, we get

$$z_i \mid x_i, \mu^{old} \sim \text{Bernoulli}(p_i), \text{ where}$$

$$p_i = \frac{\text{Normal}(x_i; \mu^{old}, 1)}{\text{Normal}(x_i; 0, 1) + \text{Normal}(x_i; \mu^{old}, 1)}$$

- E step: Compute  $E_Z[l(\mu)]$ . M step: Set  $\mu^{new}$  as the parameter maximizing this function.

## A toy example continued

- ▶ The E step becomes

$$\begin{aligned} \mathbb{E}_Z[l(\mu)] &= \mathbb{E}_Z \left[ \sum_{i=1}^n -\frac{1}{2} (x_i - z_i \mu)^2 \right] \\ &= \mathbb{E}_Z \left[ -\frac{1}{2} \sum_{i=1}^n x_i^2 - 2x_i z_i \mu + z_i^2 \mu^2 \right] \\ &= -\frac{1}{2} \sum_{i=1}^n x_i^2 - 2x_i \mathbb{E}_Z[z_i] \mu + \mathbb{E}_Z[z_i^2] \mu^2 \\ &= -\frac{1}{2} \sum_{i=1}^n x_i^2 - 2x_i p_i \mu + p_i \mu^2 \end{aligned}$$

- ▶ The M step becomes

$$\frac{\partial}{\partial \mu} \mathbb{E}_Z[l(\mu)] = -\frac{1}{2} \sum_{i=1}^n (-2x_i p_i + 2p_i \mu) = \sum_{i=1}^n x_i p_i - \mu \sum_{i=1}^n p_i = 0$$

resulting in  $\mu^{\text{new}} = (\sum_{i=1}^n x_i p_i) / (\sum_{i=1}^n p_i)$ .