

---

# Kullback Matching and Maximum Likelihood Estimation

---

SEUNGJIN CHOI  
 DEPARTMENT OF COMPUTER SCIENCE  
 POSTECH, KOREA

## Outline

This note shows that maximum likelihood estimation is identical to the minimization of Kullback-Leibler divergence between the empirical distribution and model distribution.

## Details

The Kullback-Leibler divergence (KL-divergence) is a popular measure for a similarity between two probability distributions, defined by

$$KL[p||q] = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}, \quad (1)$$

where  $p(\mathbf{x})$  and  $q(\mathbf{x})$  are probability distributions.

Let denote by  $\tilde{p}(\mathbf{x})$  and  $p(\mathbf{x}|\boldsymbol{\theta})$ , the empirical distribution and model distribution, respectively. Two fundamental properties of KL-divergence are:

- $KL[p||q] \geq 0$  (Gibb's inequality) with equality holding if and only if  $p(\mathbf{x}) = q(\mathbf{x})$ .
- The KL-divergence is asymmetric, i.e.,  $KL[p||q] \neq KL[q||p]$ .

Given a set of data points,  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  drawn from the underlying distribution  $p(\mathbf{x})$ , let  $\tilde{p}(\mathbf{x})$  be the empirical distribution which puts probability  $\frac{1}{N}$  on each data point, leading to

$$\tilde{p}(\mathbf{x}) = \frac{1}{N} \sum_{t=1}^N \delta(\mathbf{x} - \mathbf{x}_t). \quad (2)$$

We consider the KL-divergence from the empirical distribution  $\tilde{p}(\mathbf{x})$  to the model distribution  $p(\mathbf{x}|\boldsymbol{\theta})$

$$\begin{aligned} KL[\tilde{p}(\mathbf{x})||p(\mathbf{x}|\boldsymbol{\theta})] &= \int \tilde{p}(\mathbf{x}) \log \frac{\tilde{p}(\mathbf{x})}{p(\mathbf{x}|\boldsymbol{\theta})} d\mathbf{x} \\ &= -H(\tilde{p}) - \int \tilde{p}(\mathbf{x}) \log p(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x}, \end{aligned} \quad (3)$$

where  $H(\tilde{p}) = -\int \tilde{p}(\mathbf{x}) \log \tilde{p}(\mathbf{x}) d\mathbf{x}$  is the entropy of  $\tilde{p}$ .

It follows from (3) that

$$\arg \min_{\boldsymbol{\theta}} KL[\tilde{p}(\mathbf{x})||p(\mathbf{x}|\boldsymbol{\theta})] \equiv \arg \max_{\boldsymbol{\theta}} \langle \log p(\mathbf{x}|\boldsymbol{\theta}) \rangle_{\tilde{p}}, \quad (4)$$

where  $\langle \cdot \rangle_{\tilde{p}}$  represents the expectation with respect to the distribution  $\tilde{p}$ . Plugging (2) into the righthand side of (4), leads to

$$\begin{aligned} \langle \log p(\mathbf{x}|\boldsymbol{\theta}) \rangle_{\tilde{p}} &= \frac{1}{N} \int \sum_{t=1}^N N \delta(\mathbf{x} - \mathbf{x}_t) \log p(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} \\ &= \frac{1}{N} \sum_{t=1}^N \log p(\mathbf{x}_t|\boldsymbol{\theta}). \end{aligned} \quad (5)$$

Apart from the scaling factor  $\frac{1}{N}$ , this is just the log-likelihood function. In other words, maximum likelihood estimation is obtained from the minimization of (3) (*Kullback matching*).