

CS-E4891 Deep Generative Models

Lecture 1: Statistical methods

Harri Lähdesmäki

Department of Computer Science
Aalto University

March 17, 2025

Outline

- Divergence measures (Sections 2.7, 5.1 from (Murphy, 2023))
- Monte Carlo (Sec. 11.2)
- Importance sampling (Sec. 11.5)
- Neural network likelihood models (Sec. 16.3)

- Reading: parts of Sections 2.7, 5.1, 11.2, 11.5, and 16.3 from (Murphy, 2023)

Commonly used training objectives in probabilistic modeling

- Likelihood, maximum likelihood: e.g.
 - Bernoulli distribution for binary classification (cross-entropy loss)
 - Categorical distribution for multi-class classification (cross-entropy loss)
 - Gaussian distribution for regression (a special case: MSE loss)
 - Generalized models with link function, typically for exponential family
- Bayesian inference:
 - Likelihood, prior, Bayes theorem, posterior distribution
- Divergence measures between probability distributions, e.g.:
 - f -divergences: e.g. Kullback-Leibler
 - Integral probability metrics: e.g. maximum mean discrepancy (MMD)

Divergence measures between probability distributions

- Goal: compare two probability distributions, P and Q , defined on the same space \mathcal{X}
- Two settings: compute a divergence measure $D(P, Q)$ assuming distributions are defined by
 - ➊ Probability densities $p(\mathbf{x})$ and $q(\mathbf{x})$, $\mathbf{x} \in \mathcal{X}$
 - ➋ Samples from the distributions

$$X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, \text{ where } \mathbf{x}_i \sim P$$
$$X' = \{\mathbf{x}'_1, \dots, \mathbf{x}'_M\}, \text{ where } \mathbf{x}'_j \sim Q$$

Divergence measures between probability distributions

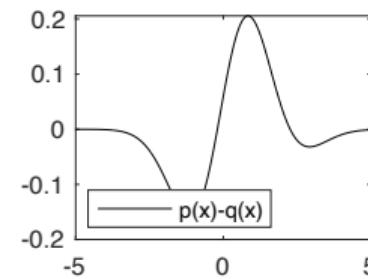
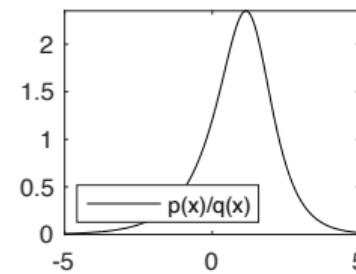
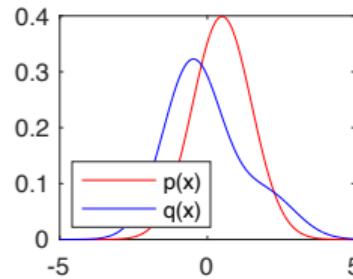
- Goal: compare two probability distributions, P and Q , defined on the same space \mathcal{X}
- Two settings: compute a divergence measure $D(P, Q)$ assuming distributions are defined by

- ➊ Probability densities $p(\mathbf{x})$ and $q(\mathbf{x})$, $\mathbf{x} \in \mathcal{X}$
- ➋ Samples from the distributions

$$X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, \text{ where } \mathbf{x}_i \sim P$$

$$X' = \{\mathbf{x}'_1, \dots, \mathbf{x}'_M\}, \text{ where } \mathbf{x}'_j \sim Q$$

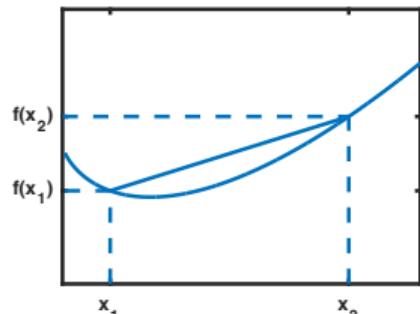
- Two main ways to compute the divergence between a pair of distributions:
 - ➊ In terms of their ratio: P/Q
 - ➋ In terms of their difference: $P - Q$



Convex and concave functions

A real-valued function $f : \mathcal{X} \rightarrow \mathbb{R}$ is^a **convex** if for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ and for all $0 \leq t \leq 1$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



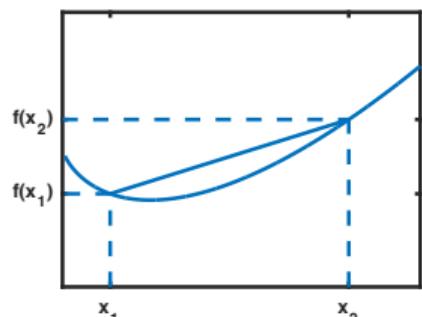
$$f(x) = x \log(x)$$

^aAssume \mathcal{X} is a convex (sub)set, e.g. $\mathcal{X} = \mathbb{R}$, \mathcal{X} is an interval in \mathbb{R} , $\mathcal{X} = \mathbb{R}^2$, etc.

Convex and concave functions

A real-valued function $f : \mathcal{X} \rightarrow \mathbb{R}$ is^a **convex** if for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ and for all $0 \leq t \leq 1$:

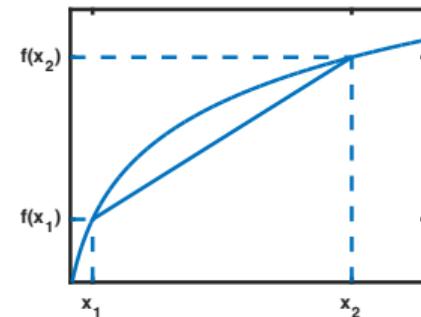
$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



$$f(x) = x \log(x)$$

A real-valued function $f : \mathcal{X} \rightarrow \mathbb{R}$ is **concave** if for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ and for all $0 \leq t \leq 1$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \geq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



$$f(x) = \log(x)$$

If f is convex, then $-f$ is concave

^aAssume \mathcal{X} is a convex (sub)set, e.g. $\mathcal{X} = \mathbb{R}$, \mathcal{X} is an interval in \mathbb{R} , $\mathcal{X} = \mathbb{R}^2$, etc.

f -divergence

f -divergence

f -divergence between two probability distributions P and Q is defined via the density ratio $r(\mathbf{x}) = \frac{p(\mathbf{x})}{q(\mathbf{x})}$

$$D_f(p \parallel q) = \int q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x},$$

where $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ is a convex function, and $f(1) = 0$

Special attention needed such that the ratio $\frac{p(\mathbf{x})}{q(\mathbf{x})}$ and $f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right)$ are defined

f-divergence

f-divergence

f-divergence between two probability distributions P and Q is defined via the density ratio $r(\mathbf{x}) = \frac{p(\mathbf{x})}{q(\mathbf{x})}$

$$D_f(p \parallel q) = \int q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x},$$

where $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ is a convex function, and $f(1) = 0$

Special attention needed such that the ratio $\frac{p(\mathbf{x})}{q(\mathbf{x})}$ and $f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right)$ are defined

Important special cases of the *f*-divergence

- Kullback-Leibler divergence: $f(r) = r \log(r)$
- Alpha divergence: $f(r) = \frac{4}{1-\alpha^2}(1 - r^{(1+\alpha)/2})$, $\alpha \neq \pm 1$
- Hellinger distance: $f(r) = (\sqrt{r} - 1)^2$
- Chi-squared distance: $f(r) = (r - 1)^2$

Jensen's inequalities

Jensen's inequality for concave functions

If x is a random variable with probability distribution P and f is a concave function, then

$$f(\mathbb{E}_{p(x)}[x]) \geq \mathbb{E}_{p(x)}[f(x)],$$

where $\mathbb{E}_{p(x)}[\cdot]$ is the expectation w.r.t. probability density $p(x)$

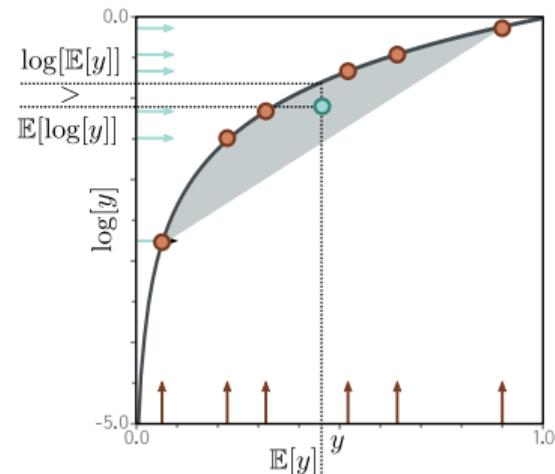


Figure 17.4 from (Prince, 2023)
Note: y has the role of x

Jensen's inequalities

Jensen's inequality for concave functions

If \mathbf{x} is a random variable with probability distribution P and f is a concave function, then

$$f(\mathbb{E}_{p(\mathbf{x})}[\mathbf{x}]) \geq \mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})],$$

where $\mathbb{E}_{p(\mathbf{x})}[\cdot]$ is the expectation w.r.t. probability density $p(\mathbf{x})$

Jensen's inequality for convex functions

If \mathbf{x} is a random variable with probability distribution P and f is a convex function, then

$$f(\mathbb{E}_{p(\mathbf{x})}[\mathbf{x}]) \leq \mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})]$$

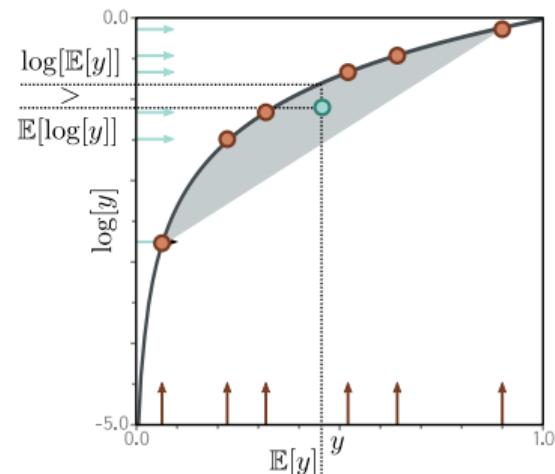


Figure 17.4 from (Prince, 2023)
Note: y has the role of \mathbf{x}

f -divergence: properties

- f -divergence can be written as

$$D_f(p \parallel q) = \int q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x} = \mathbb{E}_{q(\mathbf{x})} \left[f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) \right]$$

f -divergence: properties

- f -divergence can be written as

$$D_f(p \parallel q) = \int q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x} = \mathbb{E}_{q(\mathbf{x})} \left[f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) \right]$$

- Jensen's inequality for convex functions gives

$$\mathbb{E}_{q(\mathbf{x})} \left[f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) \right] \geq f\left(\mathbb{E}_{q(\mathbf{x})} \left[\frac{p(\mathbf{x})}{q(\mathbf{x})} \right]\right) = f\left(\int q(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}\right) = f\left(\int p(\mathbf{x}) d\mathbf{x}\right) = f(1) = 0$$

f-divergence: properties

- *f*-divergence can be written as

$$D_f(p \parallel q) = \int q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x} = \mathbb{E}_{q(\mathbf{x})} \left[f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) \right]$$

- Jensen's inequality for convex functions gives

$$\mathbb{E}_{q(\mathbf{x})} \left[f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) \right] \geq f\left(\mathbb{E}_{q(\mathbf{x})} \left[\frac{p(\mathbf{x})}{q(\mathbf{x})} \right]\right) = f\left(\int q(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}\right) = f\left(\int p(\mathbf{x}) d\mathbf{x}\right) = f(1) = 0$$

- Thus, *f*-divergence satisfies

- $D_f(p \parallel q) \geq 0$
- $D_f(p \parallel p) = 0$

- *f*-divergence is not a valid distance measure because

- it is not generally symmetric
- it does not satisfy triangle inequality

Kullback-Leibler (KL) divergence

Kullback-Leibler divergence

The Kullback-Leibler divergence between two probability distributions P and Q (divergence of P from Q) is defined as

$$\begin{aligned} D_{\text{KL}}(p \parallel q) &= \int p(\mathbf{x}) \log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) d\mathbf{x} \\ &= \mathbb{E}_{p(\mathbf{x})} \left[\log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right]. \end{aligned}$$

In other words, the KL divergence is the f -divergence with $f(r) = r \log(r)$.

- KL divergence satisfies
 - ① $D_{\text{KL}}(p \parallel q) \geq 0$ (since KL is an f -divergence)
 - ② $D_{\text{KL}}(p \parallel q) = 0$ iff $p = q$ for all \mathbf{x}
- KL divergence is not symmetric, i.e., generally $D_{\text{KL}}(p \parallel q) \neq D_{\text{KL}}(q \parallel p)$

KL divergence for independent random variables

For probability distributions of D -dimensional random variables where dimensions are independent,

$$q(\mathbf{x}) = \prod_{i=1}^D q_i(x_i) \quad \text{and} \quad p(\mathbf{x}) = \prod_{i=1}^D p_i(x_i),$$

the KL factorizes across dimensions

$$\begin{aligned} D_{\text{KL}}(p \parallel q) &= \mathbb{E}_{p(\mathbf{x})} \left[\log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right] = \mathbb{E}_{p(\mathbf{x})} \left[\log \left(\frac{\prod_{i=1}^D p(x_i)}{\prod_{i=1}^D q(x_i)} \right) \right] = \mathbb{E}_{p(\mathbf{x})} \left[\log \left(\prod_{i=1}^D \frac{p(x_i)}{q(x_i)} \right) \right] \\ &= \mathbb{E}_{p(\mathbf{x})} \left[\sum_{i=1}^D \log \left(\frac{p(x_i)}{q(x_i)} \right) \right] = \sum_{i=1}^D \mathbb{E}_{p(\mathbf{x})} \left[\log \left(\frac{p(x_i)}{q(x_i)} \right) \right] \\ &= \sum_{i=1}^D \mathbb{E}_{p(x_i)} \left[\log \left(\frac{p(x_i)}{q(x_i)} \right) \right] \\ &= \sum_{i=1}^D D_{\text{KL}}(p_i \parallel q_i), \end{aligned}$$

KL divergence: chain rule

Chain rule for probability distributions: $p(\mathbf{x}) = p(x_1)p(x_2 \mid x_1) \dots p(x_D \mid x_1, \dots, x_{D-1})$

KL divergence satisfies a natural chain rule

$$\begin{aligned} D_{\text{KL}}(p(x_1, x_2) \parallel q(x_1, x_2)) &= \iint p(x_1, x_2) \log \left(\frac{p(x_1, x_2)}{q(x_1, x_2)} \right) dx_1 dx_2 \\ &= \iint p(x_1) p(x_2 \mid x_1) \left[\log \left(\frac{p(x_1)}{q(x_1)} \right) + \log \left(\frac{p(x_2 \mid x_1)}{q(x_2 \mid x_1)} \right) \right] dx_1 dx_2 \\ &= D_{\text{KL}}(p(x_1) \parallel q(x_1)) + \mathbb{E}_{p(x_1)} [D_{\text{KL}}(p(x_2 \mid x_1) \parallel q(x_2 \mid x_1))] \end{aligned}$$

KL divergence: two Gaussian distributions

The KL divergence between two univariate Gaussian distributions

$$D_{\text{KL}}(\mathcal{N}(x | \mu_1, \sigma_1^2) || \mathcal{N}(x | \mu_2, \sigma_2^2)) = \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

KL divergence: two Gaussian distributions

The KL divergence between two univariate Gaussian distributions

$$D_{\text{KL}}(\mathcal{N}(x | \mu_1, \sigma_1^2) || \mathcal{N}(x | \mu_2, \sigma_2^2)) = \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

The KL divergence between two multivariate Gaussian distributions

$$\begin{aligned} D_{\text{KL}}(\mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) || \mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) &= \frac{1}{2} \left[\text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right. \\ &\quad \left. - D + \log\left(\frac{\det(\boldsymbol{\Sigma}_2)}{\det(\boldsymbol{\Sigma}_1)}\right) \right] \end{aligned}$$

KL divergence: two Gaussian distributions

The KL divergence between two univariate Gaussian distributions

$$D_{\text{KL}}(\mathcal{N}(x | \mu_1, \sigma_1^2) || \mathcal{N}(x | \mu_2, \sigma_2^2)) = \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

The KL divergence between two multivariate Gaussian distributions

$$\begin{aligned} D_{\text{KL}}(\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) || \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) &= \frac{1}{2} \left[\text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right. \\ &\quad \left. - D + \log\left(\frac{\det(\boldsymbol{\Sigma}_2)}{\det(\boldsymbol{\Sigma}_1)}\right) \right] \end{aligned}$$

If $\boldsymbol{\Sigma}_1 = \text{diag}(\sigma_{11}^2, \dots, \sigma_{1D}^2)$ and $\boldsymbol{\Sigma}_2 = \text{diag}(\sigma_{21}^2, \dots, \sigma_{2D}^2)$, then due to the independence property

$$D_{\text{KL}}(\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) || \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) = \sum_{i=1}^D D_{\text{KL}}(\mathcal{N}(x_i | \mu_{1i}, \sigma_{1i}^2) || \mathcal{N}(x_i | \mu_{2i}, \sigma_{2i}^2))$$

KL divergence: two Gaussian distributions

The KL divergence between two univariate Gaussian distributions

$$D_{\text{KL}}(\mathcal{N}(x | \mu_1, \sigma_1^2) || \mathcal{N}(x | \mu_2, \sigma_2^2)) = \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

The KL divergence between two multivariate Gaussian distributions

$$\begin{aligned} D_{\text{KL}}(\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) || \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) &= \frac{1}{2} \left[\text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right. \\ &\quad \left. - D + \log\left(\frac{\det(\boldsymbol{\Sigma}_2)}{\det(\boldsymbol{\Sigma}_1)}\right) \right] \end{aligned}$$

If $\boldsymbol{\Sigma}_1 = \text{diag}(\sigma_{11}^2, \dots, \sigma_{1D}^2)$ and $\boldsymbol{\Sigma}_2 = \text{diag}(\sigma_{21}^2, \dots, \sigma_{2D}^2)$, then due to the independence property

$$D_{\text{KL}}(\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) || \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) = \sum_{i=1}^D D_{\text{KL}}(\mathcal{N}(x_i | \mu_{1i}, \sigma_{1i}^2) || \mathcal{N}(x_i | \mu_{2i}, \sigma_{2i}^2))$$

These are useful results for developing deep generative models because

- Distributions are often assumed to be Gaussians
- Often $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_D^2)$

KL divergence has closed-form solution for any two exponential family distributions from the same family

KL divergence and maximum likelihood

Suppose $p_{\mathcal{D}}$ is the empirical data distribution of data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n),$$

where $\delta(\mathbf{x})$ is the delta distribution

Suppose we want to find distribution q that is as close as possible to $p_{\mathcal{D}}$

$$q^* = \arg \min_q D_{\text{KL}}(p_{\mathcal{D}} \parallel q)$$

KL divergence and maximum likelihood

Suppose $p_{\mathcal{D}}$ is the empirical data distribution of The objective function can be written as data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n),$$

where $\delta(\mathbf{x})$ is the delta distribution

Suppose we want to find distribution q that is as close as possible to $p_{\mathcal{D}}$

$$q^* = \arg \min_q D_{\text{KL}}(p_{\mathcal{D}} \parallel q)$$

$$\begin{aligned} D_{\text{KL}}(p_{\mathcal{D}} \parallel q) &= \underbrace{\int p_{\mathcal{D}}(\mathbf{x}) \log p_{\mathcal{D}}(\mathbf{x}) d\mathbf{x} - \int p_{\mathcal{D}}(\mathbf{x}) \log q(\mathbf{x}) d\mathbf{x}}_{\text{constant } C} \\ &= C - \int \left[\frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n) \right] \log q(\mathbf{x}) d\mathbf{x} \\ &= C - \frac{1}{N} \sum_{n=1}^N \int \delta(\mathbf{x} - \mathbf{x}_n) \log q(\mathbf{x}) d\mathbf{x} \\ &= C - \frac{1}{N} \sum_{n=1}^N \log q(\mathbf{x}_n) \\ &= C - \frac{1}{N} \log \prod_{n=1}^N q(\mathbf{x}_n) \end{aligned}$$

Minimizing $D_{\text{KL}}(p_{\mathcal{D}} \parallel q)$ maximizes the likelihood

KL divergence: minimization

Given a true distribution P , minimize the KL divergence w.r.t. an approximate distribution Q

Two optimization approaches due to the asymmetry of the KL divergence

- ① Forwards KL: $q = \arg \min_q D_{\text{KL}}(p \parallel q)$
 - Results in mode covering, i.e., a wide (under-confident) distribution q
- ② Reverse KL: $q = \arg \min_q D_{\text{KL}}(q \parallel p)$
 - Results in mode seeking, i.e., a narrow (over-confident) distribution q

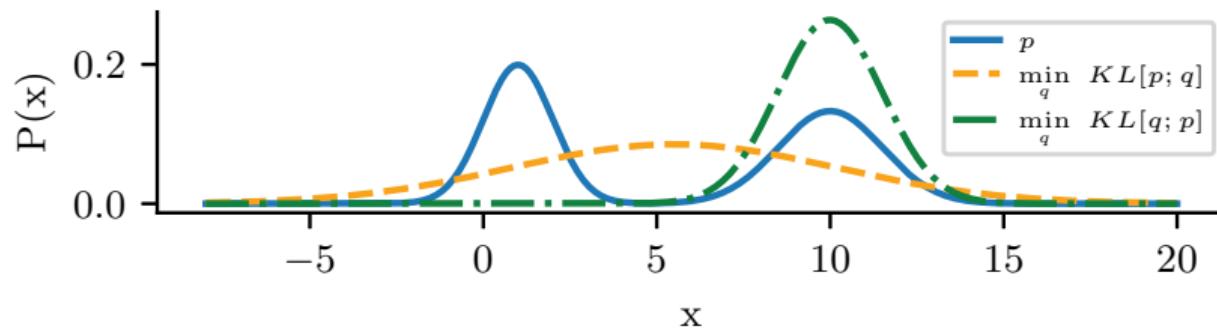


Figure 5.1 from (Murphy, 2023)

Forwards KL minimization for two Gaussians

Assume the true distribution is a 2-D correlated Gaussian $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$

Assume that $q(\mathbf{x})$ is a mean-field Gaussian $q(\mathbf{x} \mid \mathbf{m}, \mathbf{V}) = \mathcal{N}(x_1 \mid m_1, v_1) \mathcal{N}(x_2 \mid m_2, v_2)$

Forwards KL minimization for two Gaussians

Assume the true distribution is a 2-D correlated Gaussian $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$

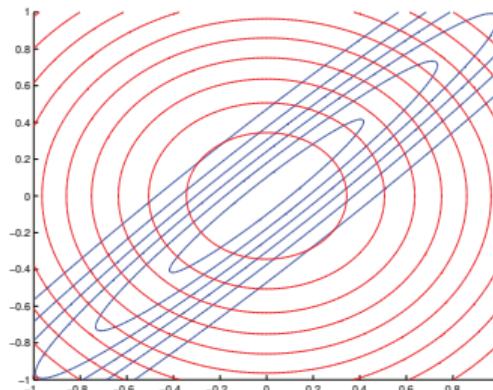
Assume that $q(\mathbf{x})$ is a mean-field Gaussian $q(\mathbf{x} | \mathbf{m}, \mathbf{V}) = \mathcal{N}(x_1 | m_1, v_1) \mathcal{N}(x_2 | m_2, v_2)$

It can be shown that

- Optimizing the forwards KL $D_{\text{KL}}(p || q)$ w.r.t. q results in moment matching, where

$$q(\mathbf{x}) = N(x_1 | \mu_1, \Sigma_{11})N(x_2 | \mu_2, \Sigma_{22})$$

- This holds more generally for exponential family
- Forwards KL is called moment projection



(a)

Figure 5.2 a) from (Murphy, 2023)

Reverse KL minimization for two Gaussians

Assume the true distribution is a 2-D correlated Gaussian $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$

Assume that $q(\mathbf{x})$ is a mean-field Gaussian $q(\mathbf{x} \mid \mathbf{m}, \mathbf{V}) = \mathcal{N}(x_1 \mid m_1, v_1) \mathcal{N}(x_2 \mid m_2, v_2)$

Reverse KL minimization for two Gaussians

Assume the true distribution is a 2-D correlated Gaussian $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$

Assume that $q(\mathbf{x})$ is a mean-field Gaussian $q(\mathbf{x} | \mathbf{m}, \mathbf{V}) = \mathcal{N}(x_1 | m_1, v_1) \mathcal{N}(x_2 | m_2, v_2)$

It can be shown that

- Optimizing the reverse KL $D_{\text{KL}}(q || p)$ w.r.t. q results in
$$q(\mathbf{x}) = N(x_1 | \mu_1, \Lambda_{11})N(x_2 | \mu_2, \Lambda_{22})$$
- Reverse KL is called information projections
- Reverse KL is typically easier to optimize because expectations are w.r.t. q which we can choose

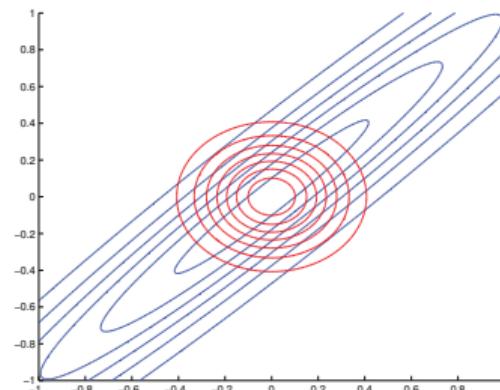


Figure 5.2 b) from (Murphy, 2023)

Jensen-Shannon divergence

- Jensen-Shannon divergence is a symmetric variant of the KL

$$\text{JSD}(p, q) = \frac{1}{2}D_{\text{KL}}(p \parallel m) + \frac{1}{2}D_{\text{KL}}(q \parallel m) = \text{JSD}(q, p),$$

where $m(\mathbf{x}) = \frac{1}{2}p(\mathbf{x}) + \frac{1}{2}q(\mathbf{x})$

- This divergence is important when training implicit generative models, such as generative adversarial networks (GANs)

Integral probability metric (IPM)

- f -divergences are defined in terms of density ratio P/Q
- Many methods use the KL divergence as their objective function
- Divergence measures can also be defined in terms of $P - Q$
- Integral probability metric is defined as

$$D_{\mathcal{F}}(P, Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})] - \mathbb{E}_{q(\mathbf{x}')}[f(\mathbf{x}')]|,$$

where \mathcal{F} is a class of smooth functions

- Depending on the choice of \mathcal{F} we get different divergence measures
- This class of divergences is important when training implicit generative models, such as generative adversarial networks (GANs)

Total variation distance

- The total variation (TV) distance between two probability distributions P and Q for a continuous random variable x

$$D_{\text{TV}} = \frac{1}{2} \int |p(x) - q(x)| dx$$

- TV satisfies the properties of a distance: symmetric, positive, definite, triangle inequality
- TV is the (only) divergence that is both f -divergence (with $f(r) = \frac{1}{2}|r - 1|$) and IPM

Monte Carlo integration

Monte Carlo methods are a stochastic approach to solve numerical integration problems

- A random variable $\boldsymbol{x} \in \mathcal{X} = \mathbb{R}^n$ with probability density $p(\boldsymbol{x})$, and a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
- Interested in computing the expected value of function $f(\boldsymbol{x})$

$$\mathbb{E}[f(\boldsymbol{x})] = \int f(\boldsymbol{x})p(\boldsymbol{x})d\boldsymbol{x}$$

Monte Carlo integration

Monte Carlo methods are a stochastic approach to solve numerical integration problems

- A random variable $\boldsymbol{x} \in \mathcal{X} = \mathbb{R}^n$ with probability density $p(\boldsymbol{x})$, and a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
- Interested in computing the expected value of function $f(\boldsymbol{x})$

$$\mathbb{E}[f(\boldsymbol{x})] = \int f(\boldsymbol{x})p(\boldsymbol{x})d\boldsymbol{x}$$

Monte Carlo integration

Monte Carlo integration estimate of $\mathbb{E}[f(\boldsymbol{x})]$ is

$$\mathbb{E}[f(\boldsymbol{x})] \approx \bar{f}_N = \frac{1}{N} \sum_{s=1}^N f(\boldsymbol{x}_s), \quad \text{where } \boldsymbol{x}_s \stackrel{\text{i.i.d.}}{\sim} p(\boldsymbol{x})$$

and N is the number of samples

- By design, function $f(\boldsymbol{x})$ is only evaluated at points \boldsymbol{x} that have non-negligible probability

Monte Carlo integration: example

- The area of a circle with radius r is $A = \pi r^2$
- Alternatively, A equals

$$A = \int_{-r}^r \int_{-r}^r \mathbb{I}(x^2 + y^2 \leq r^2) dx dy$$

where $\mathbb{I}(\cdot)$ is the indicator function (value is 1 inside the circle, 0 outside)

Monte Carlo integration: example

- The area of a circle with radius r is $A = \pi r^2$
- Alternatively, A equals

$$A = \int_{-r}^r \int_{-r}^r \mathbb{I}(x^2 + y^2 \leq r^2) dx dy$$

where $\mathbb{I}(\cdot)$ is the indicator function (value is 1 inside the circle, 0 outside)

- Define e.g. $p(x, y) = p(x)p(y)$, where $p(x)$ and $p(y)$ are uniform over $[-r, r]$
- Monte Carlo estimate of A

$$\begin{aligned} A &= (2r)(2r) \int \int \mathbb{I}(x^2 + y^2 \leq r^2) p(x)p(y) dx dy \\ &\approx 4r^2 \frac{1}{N} \sum_{s=1}^N \mathbb{I}(x_s^2 + y_s^2 \leq r^2), \quad x_s, y_s \sim p(x, y) \end{aligned}$$

Monte Carlo integration: example

- The area of a circle with radius r is $A = \pi r^2$
- Alternatively, A equals

$$A = \int_{-r}^r \int_{-r}^r \mathbb{I}(x^2 + y^2 \leq r^2) dx dy$$

where $\mathbb{I}(\cdot)$ is the indicator function (value is 1 inside the circle, 0 outside)

- Define e.g. $p(x, y) = p(x)p(y)$, where $p(x)$ and $p(y)$ are uniform over $[-r, r]$
- Monte Carlo estimate of A

$$\begin{aligned} A &= (2r)(2r) \int \int \mathbb{I}(x^2 + y^2 \leq r^2) p(x)p(y) dx dy \\ &\approx 4r^2 \frac{1}{N} \sum_{s=1}^N \mathbb{I}(x_s^2 + y_s^2 \leq r^2), \quad x_s, y_s \sim p(x, y) \end{aligned}$$

- An estimate of π is then $\hat{\pi} = \hat{A}/(r^2)$

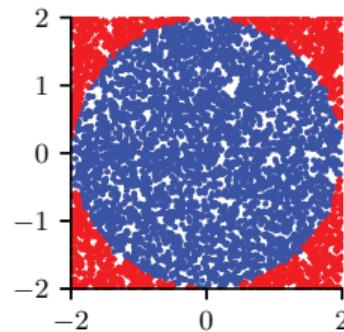


Figure 11.1 from (Murphy, 2023)

Monte Carlo integration: accuracy

Monte Carlo estimate \bar{f}_N of $\mathbb{E}[f(\mathbf{x})]$

$$\bar{f}_N = \frac{1}{N} \sum_{s=1}^N f(\mathbf{x}_s), \quad \text{where } \mathbf{x}_s \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

Assuming the true mean $\mu = \mathbb{E}[f(\mathbf{x})]$ exists, by the law of large numbers

$$\bar{f}_N \rightarrow \mu \quad \text{when} \quad N \rightarrow \infty$$

Monte Carlo integration: accuracy

Monte Carlo estimate \bar{f}_N of $\mathbb{E}[f(\mathbf{x})]$

$$\bar{f}_N = \frac{1}{N} \sum_{s=1}^N f(\mathbf{x}_s), \quad \text{where } \mathbf{x}_s \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

Assuming the true mean $\mu = \mathbb{E}[f(\mathbf{x})]$ exists, by the law of large numbers

$$\bar{f}_N \rightarrow \mu \quad \text{when } N \rightarrow \infty$$

Assuming a scalar-valued function f and finite variance $\sigma^2 = \mathbb{V}[f(\mathbf{x})] = \mathbb{E}[(f(\mathbf{x}) - \mu)^2] < \infty$, by the central limit theorem

$$(\bar{f}_N - \mu) \rightarrow \mathcal{N}\left(0, \frac{\sigma^2}{N}\right)$$

Monte Carlo integration: accuracy

Monte Carlo estimate \bar{f}_N of $\mathbb{E}[f(\mathbf{x})]$

$$\bar{f}_N = \frac{1}{N} \sum_{s=1}^N f(\mathbf{x}_s), \quad \text{where } \mathbf{x}_s \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

Assuming the true mean $\mu = \mathbb{E}[f(\mathbf{x})]$ exists, by the law of large numbers

$$\bar{f}_N \rightarrow \mu \quad \text{when} \quad N \rightarrow \infty$$

Assuming a scalar-valued function f and finite variance $\sigma^2 = \mathbb{V}[f(\mathbf{x})] = \mathbb{E}[(f(\mathbf{x}) - \mu)^2] < \infty$, by the central limit theorem

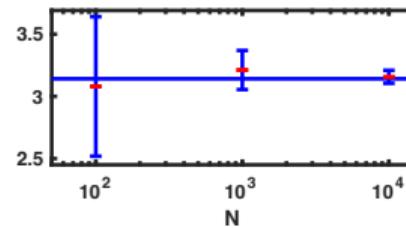
$$(\bar{f}_N - \mu) \xrightarrow{\text{d}} \mathcal{N}\left(0, \frac{\sigma^2}{N}\right)$$

Variance σ^2 can be estimated

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{s=1}^N (f(\mathbf{x}_s) - \bar{f}_N)^2$$

Thus, for large N

$$P\left(\bar{f}_N - 1.96 \frac{\hat{\sigma}}{\sqrt{N}} \leq \mu \leq \bar{f}_N + 1.96 \frac{\hat{\sigma}}{\sqrt{N}}\right) \approx 0.95$$



An example of π estimates and empirical standard errors for varying N

Monte Carlo integration: accuracy

To summarize, accuracy of Monte Carlo integration

- is independent of dimensionality of x
- increases with increasing sample size N

Importance sampling

- Consider again using Monte Carlo method to approximate integrals of the form

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x},$$

where $f(\mathbf{x})$ is the target function and $\pi(\mathbf{x})$ is now the target distribution

- If $\pi(\mathbf{x})$ is difficult to sample, we can instead sample from a proposal distribution $q(\mathbf{x})$
- We need to adjust the Monte Carlo method for this inaccuracy by associating weights with each sample

$$\mathbb{E}[f(\mathbf{x})] \approx \frac{1}{N} \sum_{s=1}^N W_n f(\mathbf{x}_n)$$

- Two different importance sampling variants:
 - Direct importance sampling: assume that, albeit difficult to sample, we can still evaluate the normalized target distribution $\pi(\mathbf{x})$
 - Self-normalized importance sampling: we can only evaluate unnormalized target distribution $\tilde{\pi}(\mathbf{x})$

Direct importance sampling

- Write the target integral as

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} = \int f(\mathbf{x})\frac{\pi(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x}$$

assuming $q(\mathbf{x}) > 0$ whenever $\pi(\mathbf{x}) > 0$, i.e., $\text{sup } \pi \subseteq \text{sup } q$

Direct importance sampling

- Write the target integral as

$$\mathbb{E}[f(\mathbf{x})] = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} = \int f(\mathbf{x})\frac{\pi(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x}$$

assuming $q(\mathbf{x}) > 0$ whenever $\pi(\mathbf{x}) > 0$, i.e., $\text{sup } \pi \subseteq \text{sup } q$

Direct importance sampling

Direct importance sampling estimate of $\mathbb{E}[f(\mathbf{x})]$ is

$$\mathbb{E}[f(\mathbf{x})] \approx \frac{1}{N} \sum_{s=1}^N \frac{\pi(\mathbf{x}_s)}{q(\mathbf{x}_s)} f(\mathbf{x}_s) = \frac{1}{N} \sum_{s=1}^N \tilde{w}_s f(\mathbf{x}_s), \quad \text{where } \mathbf{x}_s \stackrel{\text{i.i.d.}}{\sim} q(\mathbf{x})$$

and $\tilde{w}_s = \frac{\pi(\mathbf{x}_s)}{q(\mathbf{x}_s)}$ are the importance weights

Direct importance sampling: proposal distribution

Performance of the importance sampling depends crucially on the quality of $q(\mathbf{x})$

Ideally, the proposal distribution should

- cover the target distribution,
 $\pi(\mathbf{x}) > 0 \Rightarrow q(\mathbf{x}) > 0$
- be similar with $\pi(\mathbf{x})$
- take into account properties of f

One good strategy is to learn a proposal

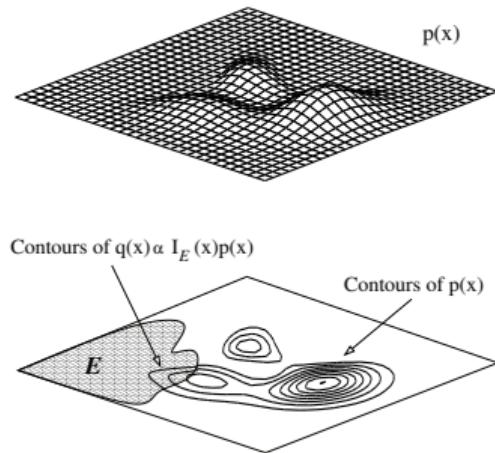


Figure 11.6 from Murphy (2023) illustrating the proposal $q(\mathbf{x})$

Self-normalized importance sampling

- It is often much easier to evaluate the unnormalized distribution $\tilde{\pi}(\mathbf{x}) = Z\pi(\mathbf{x})$, where $Z = \int \tilde{\pi}(\mathbf{x})d\mathbf{x}$

Self-normalized importance sampling

- It is often much easier to evaluate the unnormalized distribution $\tilde{\pi}(\mathbf{x}) = Z\pi(\mathbf{x})$, where $Z = \int \tilde{\pi}(\mathbf{x})d\mathbf{x}$

Self-normalized importance sampling

Self-normalized importance sampling estimate of $\mathbb{E}[f(\mathbf{x})]$ is

$$\begin{aligned}\mathbb{E}[f(\mathbf{x})] &= \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x} = \int f(\mathbf{x})\frac{1}{Z}\tilde{\pi}(\mathbf{x})d\mathbf{x} = \frac{\int f(\mathbf{x})\tilde{\pi}(\mathbf{x})d\mathbf{x}}{\int \tilde{\pi}(\mathbf{x})d\mathbf{x}} = \frac{\int \left[\frac{\tilde{\pi}(\mathbf{x})}{q(\mathbf{x})}f(\mathbf{x})\right]q(\mathbf{x})d\mathbf{x}}{\int \left[\frac{\tilde{\pi}(\mathbf{x})}{q(\mathbf{x})}\right]q(\mathbf{x})d\mathbf{x}} \\ &\approx \frac{\frac{1}{N} \sum_{s=1}^N \tilde{w}_s f(\mathbf{x}_s)}{\frac{1}{N} \sum_{s=1}^N \tilde{w}_s}, \quad \text{where } \mathbf{x}_s \stackrel{\text{i.i.d.}}{\sim} q(\mathbf{x})\end{aligned}$$

$\tilde{w}_s = \frac{\tilde{\pi}(\mathbf{x}_s)}{q(\mathbf{x}_s)}$ are the unnormalized importance weights, and assuming $\sup \tilde{\pi} \subseteq \sup q$

Multilayer perceptron (MLP)

- Multilayer perceptron (MLP) $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^C$ consists of L layers of linear computations and element-wise non-linearities φ (here $L = 2$)

$$\begin{aligned}\mathbf{y} &= \mathbf{f}(\mathbf{x}; \boldsymbol{\theta}) \\ &= \varphi(\mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2)\end{aligned}$$

where $\boldsymbol{\theta} = (\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L)$ and

$$\mathbf{W}_1 \in \mathbb{R}^{D \times K_1}$$

$$\mathbf{W}_l \in \mathbb{R}^{K_{l-1} \times K_l}$$

$$\mathbf{W}_L \in \mathbb{R}^{K_{L-1} \times C}$$

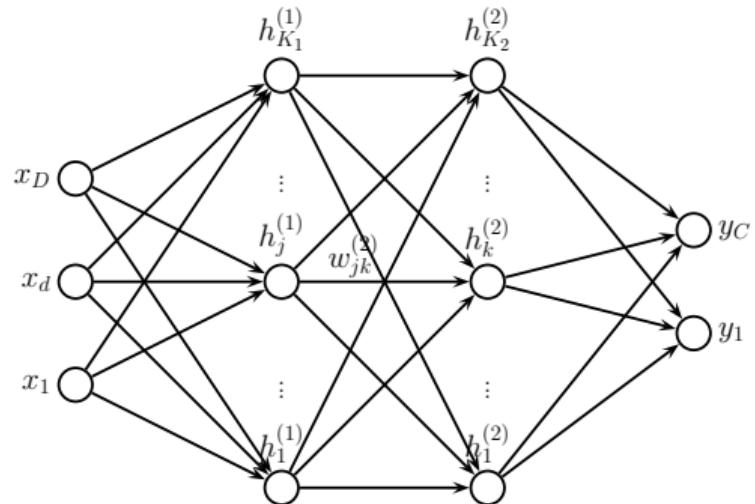


Figure 16.10 from (Murphy, 2023)

MLPs as likelihood models

Throughout this course we will use neural nets to parameterize probabilistic models

Example: regression with homoscedastic Gaussian noise

- Univariate response variable $y \in \mathbb{R}$ with MLP

$$f : \mathbb{R}^D \rightarrow \mathbb{R}$$

$$p_{\theta}(y | \mathbf{x}) = \mathcal{N}(y | f(\mathbf{x}; \theta), \sigma_{\text{noise}}^2),$$

where $f(\mathbf{x}; \theta) = \mathbb{E}(y | \mathbf{x})$

MLPs as likelihood models

Throughout this course we will use neural nets to parameterize probabilistic models

Example: regression with homoscedastic Gaussian noise

- Univariate response variable $y \in \mathbb{R}$ with MLP

$$f : \mathbb{R}^D \rightarrow \mathbb{R}$$

$$p_{\theta}(y | \mathbf{x}) = \mathcal{N}(y | f(\mathbf{x}; \theta), \sigma_{\text{noise}}^2),$$

where $f(\mathbf{x}; \theta) = \mathbb{E}(y | \mathbf{x})$

- Multivariate response variable $\mathbf{y} \in \mathbb{R}^D$ with MLP

$$f : \mathbb{R}^D \rightarrow \mathbb{R}^C$$

$$p_{\theta}(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{f}(\mathbf{x}; \theta), \Sigma),$$

where $\mathbf{f}(\mathbf{x}; \theta) = \mathbb{E}(\mathbf{y} | \mathbf{x})$

MLPs as likelihood models

Throughout this course we will use neural nets to parameterize probabilistic models

Example: regression with homoscedastic Gaussian noise

- Univariate response variable $y \in \mathbb{R}$ with MLP
 $f : \mathbb{R}^D \rightarrow \mathbb{R}$

$$p_{\theta}(y | \mathbf{x}) = \mathcal{N}(y | f(\mathbf{x}; \theta), \sigma_{\text{noise}}^2),$$

where $f(\mathbf{x}; \theta) = \mathbb{E}(y | \mathbf{x})$

- Multivariate response variable $\mathbf{y} \in \mathbb{R}^D$ with MLP
 $f : \mathbb{R}^D \rightarrow \mathbb{R}^C$

$$p_{\theta}(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{f}(\mathbf{x}; \theta), \Sigma),$$

where $\mathbf{f}(\mathbf{x}; \theta) = \mathbb{E}(\mathbf{y} | \mathbf{x})$

Example: multiclass classification $y \in \{1, \dots, C\}$

$$p_{\theta}(y | \mathbf{x}) = \text{Cat}(y | \text{softmax}(\mathbf{f}(\mathbf{x}; \theta))),$$

where

$$\text{softmax}(\mathbf{a}) = \left(\frac{e^{a_1}}{\sum_{c=1}^C e^{a_c}}, \dots, \frac{e^{a_C}}{\sum_{c=1}^C e^{a_c}} \right)$$

MLPs as likelihood models

Throughout this course we will use neural nets to parameterize probabilistic models

Example: regression with homoscedastic Gaussian noise

- Univariate response variable $y \in \mathbb{R}$ with MLP
 $f : \mathbb{R}^D \rightarrow \mathbb{R}$

$$p_{\theta}(y | \mathbf{x}) = \mathcal{N}(y | f(\mathbf{x}; \theta), \sigma_{\text{noise}}^2),$$

where $f(\mathbf{x}; \theta) = \mathbb{E}(y | \mathbf{x})$

- Multivariate response variable $\mathbf{y} \in \mathbb{R}^D$ with MLP
 $f : \mathbb{R}^D \rightarrow \mathbb{R}^C$

$$p_{\theta}(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{f}(\mathbf{x}; \theta), \Sigma),$$

where $\mathbf{f}(\mathbf{x}; \theta) = \mathbb{E}(\mathbf{y} | \mathbf{x})$

Example: multiclass classification $y \in \{1, \dots, C\}$

$$p_{\theta}(y | \mathbf{x}) = \text{Cat}(y | \text{softmax}(\mathbf{f}(\mathbf{x}; \theta))),$$

where

$$\text{softmax}(\mathbf{a}) = \left(\frac{e^{a_1}}{\sum_{c=1}^C e^{a_c}}, \dots, \frac{e^{a_C}}{\sum_{c=1}^C e^{a_c}} \right)$$

We use interchangeably $p_{\theta}(y | \mathbf{x}) = p(y | \mathbf{x}, \theta)$

All neural network architectures can be converted to likelihood models with appropriate link functions

If neural network parameters θ are random variables, then these are called Bayesian neural networks

MLP for heteroscedastic regression

- Heteroscedastic nonlinear regression = nonlinear regression with input-dependent noise variance
- Assume univariate response variable $y \in \mathbb{R}$ with Gaussian noise
- MLP has two outputs ($C = 2$) which approximate $f_\mu(\mathbf{x}) \approx \mathbb{E}[y|\mathbf{x}, \theta]$ and $f_\sigma(\mathbf{x}) \approx \mathbb{V}[y|\mathbf{x}, \theta]$ such that

$$p_{\theta}(y | \mathbf{x}) = \mathcal{N}(y | \mathbf{w}_\mu^\top f(\mathbf{x}; \mathbf{w}_{\text{shared}}), \sigma_+(\mathbf{w}_\sigma^\top f(\mathbf{x}; \mathbf{w}_{\text{shared}}))),$$

where e.g. $\sigma_+(a) = \log(1 + e^a)$

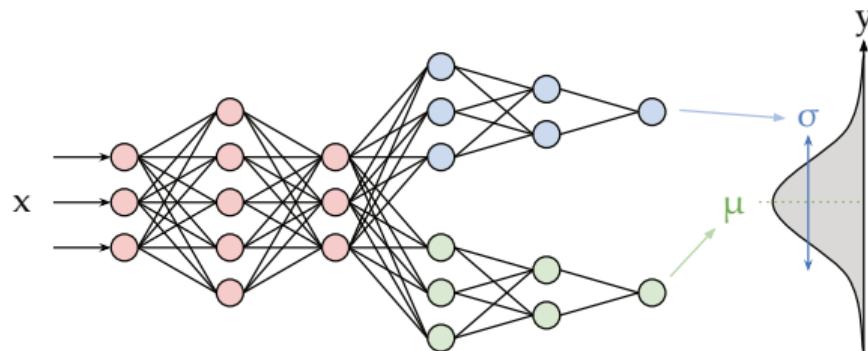


Figure 16.11 from (Murphy, 2023)

References

- Murphy K, Probabilistic Machine Learning: Advanced Topics, 2023.
- Prince SJD, Understanding Deep Learning, The MIT Press, 2023.