

Probabilistic modeling

Motivation: many machine learning problems can be conveniently framed into a probabilistic framework

Quick refresher on probabilities

- A **random variable** X represents the uncertainty of a certain event
- X can be either **discrete** (finite number of possible values) or **continuous** (real values)
- X is described by:
 - (discrete) its probability mass function $P(X = x)$
 - (continuous) its probability density function p such that
$$P(a \leq X \leq b) = \int_a^b p(x) dx$$
- A **multivariate** random variable \mathbf{X} in dimension D is a vector of random variables, taking values $\mathbf{x} = [x_1, \dots, x_D]^T$

Sum rule, product rule, Bayes theorem

- **Joint** probability¹ of two random variables: $p(x, y)$
- **Sum rule:** the marginal probability is given by

$$p(x) = \begin{cases} \sum_y p(x, y) & \text{if } y \text{ is discrete} \\ \int p(x, y) dy & \text{if } y \text{ is continuous} \end{cases}$$

- **Product rule:** relates joint and conditional probabilities

$$\begin{aligned} p(x, y) &= p(y | x) p(x) \\ &= p(x | y) p(y) \end{aligned}$$

- From the product rule, we can derive **Bayes theorem**:

$$\underbrace{p(y | x)}_{\text{posterior}} = \frac{\overbrace{p(x | y)}^{\text{likelihood}} \overbrace{p(y)}^{\text{prior}}}{\underbrace{p(x)}_{\text{evidence}}}$$

¹for continuous variables, we actually refer to the probability density function

Expected value and variance

For a random variable X ,

- The **expected value** of a function f over x is defined as:

$$\mathbb{E}_x[f] = \sum_x p(x)f(x) \quad (\text{in the discrete case})$$

$$\mathbb{E}_x[f] = \int p(x)f(x) \, dx \quad (\text{in the continuous case})$$

- The **variance** of f is then defined as:

$$\text{Var}_x[f] = \mathbb{E}_x[(f(x) - \mathbb{E}_x[f])^2]$$

For a multivariate random variable \mathbf{X} in dimension D ,

- The expected value is a vector $\mathbb{E}_{\mathbf{X}}[\mathbf{x}] = [\mathbb{E}[x_1], \dots, \mathbb{E}[x_D]]^T$
- The covariance matrix $\text{Cov}_{\mathbf{X}}[\mathbf{x}]$ is the symmetric matrix of co-variances of its components

Bernoulli distribution

A discrete random variable X with states $x \in \{0, 1\}$ such that:

$$p(x) = \mu^x (1 - \mu)^{1-x}$$

where $\mu \in [0, 1]$ is the parameter of the distribution.

We have:

$$\mathbb{E}[x] = \mu$$

$$\text{Var}[x] = \mu(1 - \mu)$$

Example: flip of a coin.

Multivariate case: **categorical** or **multinoulli** distribution over k states, parameterized by vector $\mathbf{p} \in [0, 1]^k$ where $\sum_i p_i = 1$.

Normal distribution

Normal distribution (gaussian)

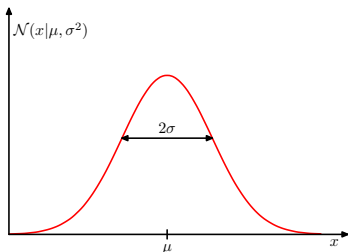
A continuous random variable x following a normal distribution has the following probability density function:

$$p(x) = \mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Parameterized by its mean μ and variance σ^2 :

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x \, dx = \mu$$

$$\text{Var}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) (x - \mu)^2 \, dx = \sigma^2$$



Multivariate normal distribution

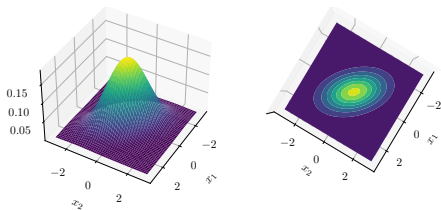
Multivariate normal distribution

For continuous random vector \mathbf{x} in dimension d , the probability density function is given by:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

Parameterized by its mean vector $\boldsymbol{\mu}$
and its covariance matrix $\boldsymbol{\Sigma}$:

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} \quad \text{Var}[\mathbf{x}] = \boldsymbol{\Sigma}$$

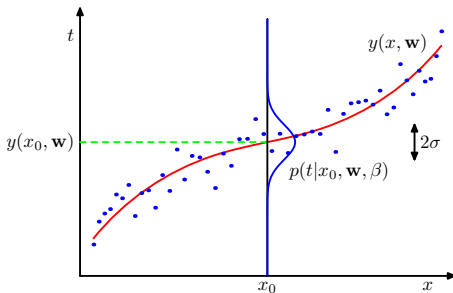


Linear regression as a probabilistic model

Recall the linear regression problem: $y = f(x; \mathbf{w})$

Hypothesis: assume target y is generated from a normal distribution with a mean equal to $f(x; \mathbf{w})$

$$p(y | x, \mathbf{w}, \sigma^2) = \mathcal{N}(y | f(x; \mathbf{w}), \sigma^2)$$



Also assume y_i are drawn independently from each other (conditional independence).

Maximum likelihood estimation

With target vector \mathbf{y} and data matrix \mathbf{X} , we can write the conditional likelihood of the training data:

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{i=1}^N \mathcal{N}(y_i | (f(\mathbf{x}_i, \mathbf{w}), \sigma^2))$$

Maximum likelihood (ML) estimation: find \mathbf{w} by maximizing the conditional (log-)likelihood of the training data

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2)$$

with:

$$\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) = -\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

→ **Equivalent to minimizing the mean squared error**

Maximum a posteriori estimation

Hypothesis 2: parameters \mathbf{w} follow a prior normal distribution

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, b^2 \mathbf{I})$$

Maximum a posteriori (MAP) estimation:

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} \log p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) = \arg \max_{\mathbf{w}} \log \frac{p(\mathbf{y} \mid \mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}, \mathbf{X})}$$

We have:

$$\begin{aligned} \log p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) &= \log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) + \log p(\mathbf{w}) + \text{const} \\ &= -\frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 - \frac{1}{2b^2} \|\mathbf{w}\|^2 + \text{const} \end{aligned}$$

→ **Equivalent to minimizing the regularized mean squared error**