

Probability Review for Machine Learning

Chris J. Maddison¹

University of Toronto

¹Slides adapted from CSC 311.

Uncertainty comes from:

- Noisy measurements
- Variability between samples
- Finite size of data sets

Probability theory provides a consistent formalism for the quantification and manipulation of uncertainty.

Sample Space

The data comes from a measurement of the real world, which we can think of as an experiment:

- **Sample space** Ω is the set of all possible outcomes of the experiment.
- **Observations** $\omega \in \Omega$ are points in the space also called sample outcomes, realizations, or elements.
- **Events** $E \subset \Omega$ are subsets of the sample space.

For example, if we flip a coin twice:

Sample space All outcomes $\Omega = \{HH, HT, TH, TT\}$

Observation $\omega = HT$ valid sample since $\omega \in \Omega$

Event Both flips same $E = \{HH, TT\}$ valid event since $E \subset \Omega$

The probability of an event E , $P(E)$, satisfies three axioms:

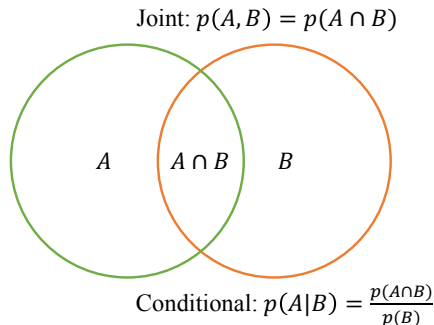
- 1: $P(E) \geq 0$ for every E
- 2: $P(\Omega) = 1$
- 3: If E_1, E_2, \dots are disjoint then

$$P\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_{i=1}^{\infty} P(E_i)$$

Joint and Conditional Probabilities

Joint Probability of A and B is denoted $P(A, B)$.

Conditional Probability of A given B is denoted $P(A|B)$.



$$p(A, B) = p(A|B)p(B) = p(B|A)p(A)$$

Conditional Example

Probability of passing the midterm is 60% and probability of passing both the final and the midterm is 45%.

What is the probability of passing the final given the student passed the midterm?

$$\begin{aligned}P(F|M) &= P(M, F)/P(M) \\ &= 0.45/0.60 \\ &= 0.75\end{aligned}$$

Independence

Events A and B are **independent** if $P(A, B) = P(A)P(B)$.

Suppose $P(A) = P(B) = 0.5$.

- Independent: A : first toss is HEAD; B : second toss is HEAD;

$$P(A, B) = 0.5 * 0.5 = P(A)P(B)$$

- Not Independent: A : first toss is HEAD; B : first toss is HEAD;

$$P(A, B) = 0.5 \neq P(A)P(B)$$

Independence

Events A and B are **conditionally independent** given C if

$$P(A, B|C) = P(B|C)P(A|C) \quad (1)$$

Consider two coins ²: A regular coin and a coin which always outputs HEAD or always outputs TAIL.

Now consider the following events.

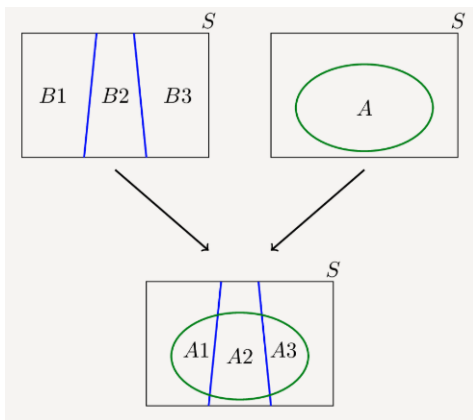
- A =The first toss is HEAD.
- B =The second toss is HEAD.
- C =The regular coin is used.
- D =The other coin is used.

Then A and B are conditionally independent given C , but A and B are NOT conditionally independent given D .

²www.probabilitycourse.com/chapter1/1_4_4_conditional_independence.php

Marginalization and Law of Total Probability³

$$P(X) = \sum_Y P(X, Y) = \sum_Y P(X|Y)P(Y)$$



³www.probabilitycourse.com/chapter1/1_4_2_total_probability.php

Bayes' Rule:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(\theta|x) = \frac{P(x|\theta)P(\theta)}{P(x)}$$

$$\text{Posterior} = \frac{\text{Likelihood} * \text{Prior}}{\text{Evidence}}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

Bayes' Example

Suppose you have tested positive for a disease. What is the probability you actually have the disease?

This depends on the prior probability of the disease:

- $P(T = 1|D = 1) = 0.95$ (likelihood)
- $P(T = 1|D = 0) = 0.10$ (likelihood)
- $P(D = 1) = 0.1$ (prior)

So $P(D = 1|T = 1) = ?$

Bayes' Example

Suppose you have tested positive for a disease. What is the probability you actually have the disease?

$$P(T = 1|D = 1) = 0.95 \text{ (true positive)}$$

$$P(T = 1|D = 0) = 0.10 \text{ (false positive)}$$

$$P(D = 1) = 0.1 \text{ (prior)}$$

So $P(D = 1|T = 1) = ?$

Use Bayes' Rule:

$$P(D = 1|T = 1) = \frac{P(T = 1|D = 1)P(D = 1)}{P(T = 1)} = \frac{0.95 * 0.1}{P(T = 1)} = 0.51$$

$$\begin{aligned} P(T = 1) &= P(T = 1|D = 1)P(D = 1) + P(T = 1|D = 0)P(D = 0) \\ &= 0.95 * 0.1 + 0.1 * 0.90 = 0.185 \end{aligned}$$

How do we connect sample spaces and events to data?

A **random variable** is a mapping which assigns a real number $X(\omega)$ to each observed outcome $\omega \in \Omega$

For example, let's flip a coin 10 times. $X(\omega)$ counts the number of Heads we observe in our sequence. If $\omega = HHTHTHHTHT$ then $X(\omega) = 6$.

Discrete and Continuous Random Variables

Discrete Random Variables

- Takes countably many values, e.g., number of heads
- Distribution defined by probability mass function (PMF) p_X .
- The probability that $X(\omega) \in B$ is given by:

$$P(X(\omega) \in B) = \sum_{x \in B} p_X(x). \quad (2)$$

Continuous Random Variables

- Takes uncountably many values, e.g., time to complete task
- Distribution defined by probability density function (PDF) f_X .
- The probability that $X(\omega) \in B$ is given by:

$$P(X(\omega) \in B) = \int_B f_X(x) dx \quad (3)$$

Joint Distributions

Two random variables $X(\omega)$, $Y(\omega)$ have a **joint distribution**.

- If both are discrete, then they have a joint PMF $p_{X,Y}(x,y)$.
- If both are continuous, then they have a joint PDF $f_{X,Y}(x,y)$.

We can marginalize the joint to get the **marginal distributions** of X or Y :

- $\sum_y p_{X,Y}(x,y) = p_X(x)$.
- $\int f_{X,Y}(x,y) dx = f_Y(y)$.

p_X and f_Y are often just called **marginals**.

Conditional Distributions

The **conditional distribution** of X given Y can be described using conditional PMF or PDF:

- If both are discrete, then the conditional PMF of X given Y is

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$$

- If both are continuous, then the conditional PMF of X given Y is

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$$

X is **independent** of Y if $p_{X|Y}(x|y) = p_X(x)$ or $f_{X|Y}(x|y) = f_X(x)$.

Random variables are said to be **independent and identically distributed** (i.i.d.) if they are sampled from the same probability distribution and are mutually independent.

This is a common assumption for observations. For example, coin flips are assumed to be iid.

A Note on Notation

The machine learning discipline has a variety of different notation norms for probability, which range from very formal to very informal. We will lean informal, because it's efficient and important to learn the norms.

We use lower case for random variables, use p for all PMFs / PDFs, and sometimes even omit the ranges of integration.

$$P(X(\omega) \in B) = \int_B f_X(x) dx$$

$$f_{X,Y}(x,y) = f_X(x)f_Y(y)$$

$$f_{X|Y}(x|y)$$

becomes \rightarrow

$$P(x \in B) = \int_B p(x) dx$$

$$p(x,y) = p(x)p(y)$$

$$p(x|y)$$

This seems like a bit of a mess, but with practice it's always clear what is intended from context.

Mean: First Moment, μ

$$\mathbb{E}[X] = \sum_{i=1}^{\infty} x_i p(x_i) \quad (\text{univariate discrete r.v.})$$

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} xp(x)dx \quad (\text{univariate continuous r.v.})$$

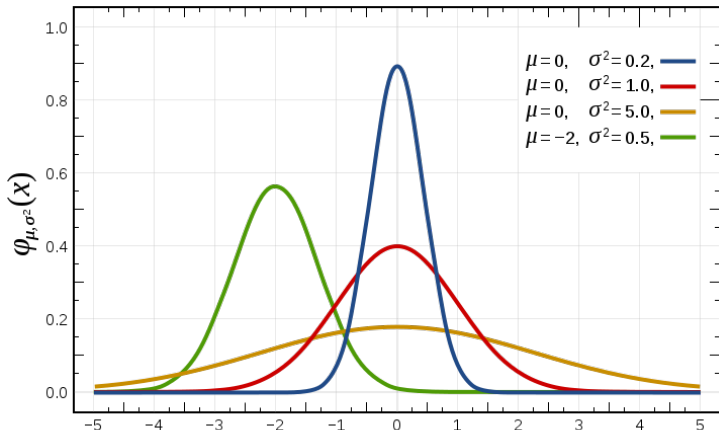
Variance: Second (central) Moment, σ^2

$$\begin{aligned} \text{Var}[X] &= \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx \\ &= \mathbb{E}[(X - \mu)^2] \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \end{aligned}$$

Univariate Gaussian Distribution

Also known as the **Normal Distribution**, $\mathcal{N}(\mu, \sigma^2)$

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



Mixed Discrete and Continuous

It is very common in machine learning to have a joint distribution over a mixture of discrete and continuous random variables. In this case, we can define the joint distribution in terms of marginals and conditionals.

Suppose x is discrete and y is continuous. Then we can define the joint distribution by

$$p(x, y) = p(x)p(y|x)$$

More specifically, suppose $x \in \{1, \dots, n\}$ with probabilities given by p_x . Suppose also we have n real numbers μ_i , then we can define a joint distribution between x and a continuous y via

$$p(x, y) = p_x \mathcal{N}(y|\mu_x, 1)$$

Multivariate Gaussian Distribution

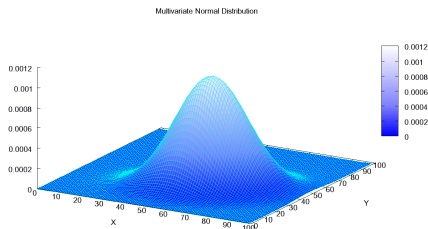
\mathbf{x} is a D -dimensional vector

$\boldsymbol{\mu}$ is a D -dimensional mean vector

Σ is a $D \times D$ covariance matrix with determinant $|\Sigma|$.

When Σ is invertible, the Gaussian has a density:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$



Covariance Matrix

Mean $\boldsymbol{\mu}$ is the D-dimensional vector whose i entry is the expected value

$$\begin{aligned}\boldsymbol{\mu}_i &= \mathbb{E}[\mathbf{x}_i] \\ &= \int \mathbf{x}_i \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}\end{aligned}$$

Covariance matrix $\boldsymbol{\Sigma}$ is a matrix whose (i, j) entry is the covariance

$$\begin{aligned}\Sigma_{ij} &= \text{Cov}(\mathbf{x}_i, \mathbf{x}_j) \\ &= \mathbb{E}[(\mathbf{x}_i - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_j)] \\ &= \mathbb{E}[(\mathbf{x}_i \mathbf{x}_j)] - \boldsymbol{\mu}_i \boldsymbol{\mu}_j\end{aligned}$$

so notice that the diagonal entries are the variance of each elements.

The covariance matrix has the property that it is symmetric and positive-semidefinite.

Inferring Parameters

We have data \mathbf{x} and we assume it is sampled from some distribution $p(\mathbf{x}|\theta)$ with parameters θ , which are unknown to us.

How do we figure out the θ that 'best' fit that distribution?

One classical approach is **maximum likelihood estimation** (MLE),

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \log p(\mathbf{x}|\theta)$$

We will discuss this in more detail in the later part of the course. For now, just assume that this makes sense to do.

MLE for Univariate Gaussian Distribution

For example, we are trying to infer the parameters for a Univariate Gaussian Distribution, mean (μ) and variance (σ^2).

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

The **likelihood** that our observations $\mathbf{x} = (x_1, \dots, x_N)$ were generated by a univariate Gaussian with parameters μ and σ^2 is

$$\text{Likelihood} = p(x_1 \dots x_N | \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x_i - \mu)^2\right\}$$

In this case $\theta = (\mu, \sigma^2)$.

MLE for Univariate Gaussian Distribution

For MLE we want to maximize this likelihood, which is difficult because it is represented by a product of terms

$$\text{Likelihood} = p(x_1 \dots x_N | \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x_i - \mu)^2\right\}$$

So we take the log of the likelihood so the product becomes a sum

$$\begin{aligned} \text{Log Likelihood} &= \log p(x_1 \dots x_N | \mu, \sigma^2) \\ &= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x_i - \mu)^2\right\} \end{aligned}$$

Since log is monotonically increasing $\max L(\theta) = \max \log L(\theta)$

MLE for Univariate Gaussian Distribution

The log Likelihood simplifies to

$$\begin{aligned}\mathcal{L}(\mu, \sigma) &= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x_i - \mu)^2\right\} \\ &= -\frac{1}{2}N \log(2\pi\sigma^2) - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2}\end{aligned}$$

Which we want to maximize. How?

MLE for Univariate Gaussian Distribution

To maximize we take the derivatives, set equal to 0, and solve:

$$\mathcal{L}(\mu, \sigma) = -\frac{1}{2}N \log(2\pi\sigma^2) - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2}$$

Derivative w.r.t. μ , set equal to 0, and solve for $\hat{\mu}$

$$\frac{\partial \mathcal{L}(\mu, \sigma)}{\partial \mu} = 0 \implies \hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

Therefore the $\hat{\mu}$ that maximizes the likelihood is the average of the data points.

Derivative w.r.t. σ^2 , set equal to 0, and solve for $\hat{\sigma}^2$

$$\frac{\partial \mathcal{L}(\mu, \sigma)}{\partial \sigma^2} = 0 \implies \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$