CS 760: Machine Learning

Spring 2024

Topic 3: Review of Probability Theory

INSTRUCTOR: DANIEL L. PIMENTEL-ALARCÓN

© Copyright 2024

DO NOT POLLUTE! AVOID PRINTING, OR PRINT 2-SIDED MULTIPAGE.

3.1 Why Probability?

Many (if not all) applications of machine learning involve randomness.

Example 3.1. In computer vision we often want to track objects or people. The location of these objects over time is random. For example, look at this frame of a video:



You can predict where some of these people will be in a few seconds (based on their current locations and directions). However, you cannot be certain. What if someone trips? What if the wind moves the leaves of a tree and produces occlusions? All of this can be modeled probabilistically in a precise mathematical way. For instance, we can say that your predictions will be accurate with a probability p close to 1 (if nothing unusual happens), and inaccurate with probability 1 - p (if something unusual happens). Furthermore, it is harder to predict where these people will be later in time. So we can refine our probability model and say your predictions will be accurate with probability p/t (where t is the amount of time), and inaccurate with probability 1 - p/t.

You can see that *probability theory is nothing but common sense reduced to calculation* — Pierre Laplace, 1812.

Example 3.2. All organisms have a DNA sequence (genome), i.e., a very long vector of nucleotides (A, C, G, T). Every now and then, organisms *mutate*. That is, when organisms reproduce, some of their nucleotides get changed. For example, it is possible that a parent has the sequence



The location of these mutations can be modeled probabilistically. For instance, it is known that certain mutations are approximately distributed uniformly at random within each gene (delimited region of the genome), and that some well-identified genes are more likely to present mutations than others.

Mutations are more common in smaller organisms, like bacteria and viruses, that reproduce very rapidly. That is why pharmaceuticals, insurance companies, the National Health Service (NHS), the National Security Agency (NSA), and even the Department of Defense (DoD), among many others, are very interested in this phenomenon. The reasons are not as apocalyptic as a zombi virus, but close (does this sound familiar/relevant right now in 2020?). For example, some bacteria have become immune to most antibiotics. Similarly, it is harder to produce vaccines for viruses that mutate quickly. In fact these mutations are the reason why we don't have a definitive flu vaccine yet. What about Corona virus?

Example 3.3. Let's now consider Netflix. Some users have rated some movies, some users have rated others. However, nobody has rated all of them. This produces an *incomplete* data matrix like this (colors indicate how much each person liked a movie)



The goal is to predict which users will like which items, in order to make good recommendations. Again, this can be modeled probabilistically. The movies that each user has rated (and hence the samples in this matrix) are somewhat random. For example, adults are more likely to watch (and enjoy) adult movies, while kids and parents are more likely to watch (and enjoy) children movies. This can be modeled probabilistically.

We could construct similar models for songs, shoes, clothes, restaurants, groceries, etc. So in fact this also applies to Amazon, Pandora, Spotify, Pinterest, Yelp, Apple, etc. If these companies recommend you an item you will like, you are more likely to buy it. You can see why all these companies have a great interest in this problem, and they are paying *a lot* of money to people who work on this.

I hope these few examples help convincing you of the ubiquitousness of randomness in machine learning. Probability theory allows us to model randomness in a precise mathematical way. Furthermore, despite the uncertainty produced by randomness, probability theory allows us to draw *likely* conclusions in a sensible manner, and quantify how certain we are about these conclusions.

3.2 The Basics

Definition 3.1. There are three elemental concepts in basic probability theory:

 Ω := Sample space = set of all possible outcomes.

- $\mathcal{A} := \sigma$ -Algebra = Set of all possible events.
- \mathbb{P} := Probability measure.

Example 3.4. Consider a fair die. Then

$$\begin{split} \Omega &= \{1, 2, 3, 4, 5, 6\}. \\ \mathcal{A} &= \left\{\{1\}, \dots, \{6\}, \{1, 2\}, \dots, \{5, 6\}, \dots, \{1, 2, 3, 4, 5, 6\}\right\} \\ \mathbb{P} &(\mathbf{x}) &= \frac{1}{6} \text{ for every } \mathbf{x} = 1, \dots, 6. \end{split}$$

Definition 3.2 (Probability measure). A mapping $\mathbb{P} : \mathcal{A} \to [0,1]$ is a *probability measure* if it satisfies the next properties:

- (i) $\mathbb{P}(A) \ge 0$ for every $A \in \mathcal{A}$.
- (ii) $\mathbb{P}(\Omega) = 1$.
- (iii) $A \cap B = \emptyset$ for some $A, B \in \mathcal{A}$, then $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$.

Condition (iii) implies that $\mathbb{P}(A \cup B) \leq \mathbb{P}(A) + \mathbb{P}(B)$, which is often known as the *union bound*.

3.3 Conditional Probability

Definition 3.3 (Conditional probability). Let $A, B \in \mathcal{A}$. The *conditional probability* that A occurs given B occurred is

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

Example 3.5. Continuing with Example 3.4, let $A = \{1, 2\}$, $B = \{2, 3\}$. The probability that A occurs is $\mathbb{P}(A) = \frac{1}{3}$. However, if you already know that B occurred, then the probability that A also occurs increases to

$$\mathbb{P}(A|B) = \frac{1/6}{1/3} = \frac{1}{2}$$

3.4 Independence

Definition 3.4 (Independent events). Let $A, B \in \mathcal{A}$. We say A and B are *independent* if $\mathbb{P}(A|B) = \mathbb{P}(A)$.

In words, two events are independent if they provide no information of one an other.

Example 3.6. Consider two fair dice. Let A be the event that the first die is 1; let B be the event that the second die is 1. Then

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A=1 \cap B=1)}{\mathbb{P}(B=1)} = \frac{1/36}{1/6} = \frac{1}{6} = \mathbb{P}(A).$$

Hence the events A and B are independent. This matches our intuition that one die has no influence on the outcome of the other.

3.5 Bayes Rule

Given the conditional probability $\mathbb{P}(A|B)$, Bayes rule gives us a formula for the *inverse* probability, $\mathbb{P}(B|A)$.

Proposition 3.1 (Bayes rule). Let $A, B \in \mathcal{A}$. Then

$$\mathbb{P}(B|A) = \frac{\mathbb{P}(A|B)\mathbb{P}(B)}{\mathbb{P}(A)}$$
(3.1)

Proof. On one hand:

 $\mathbb{P}(A \cap B) \ = \ \mathbb{P}(A|B)\mathbb{P}(B).$

By symmetry,

$$\mathbb{P}(A \cap B) = \mathbb{P}(B|A)\mathbb{P}(A)$$

It follows that

$$\mathbb{P}(A|B)\mathbb{P}(B) = \mathbb{P}(B|A)\mathbb{P}(A),$$

and solving for $\mathbb{P}(B|A)$ we obtain (3.1), as desired.

3-4

Bayes rule plays a crucial role in modern applications.

Example 3.7. Geneticists have determined that 90% of the people with disease *B* have gene *A* active, i.e., $\mathbb{P}(A|B) = 0.9$. If you sequence your genome and find out that your gene *A* is active, what is the probability that you develop disease *B*? In other words, what is $\mathbb{P}(B|A)$? At first glance you might think it is very likely that you will develop disease *B*. However, to determine this you need to know $\mathbb{P}(A)$ and $\mathbb{P}(B)$. Of the whole population, if only 5% have disease *B*, while 45% have gene *A* active, what is $\mathbb{P}(B|A)$? This is a simple application of Bayes rule:

$$\mathbb{P}(B|A) \ = \ \frac{\mathbb{P}(A|B)\mathbb{P}(B)}{\mathbb{P}(A)} \ = \ \frac{(0.9)(0.05)}{0.45} \ = \ 0.1$$

3.6 Random Variables

Sometimes the sample space Ω contains elements that are cumbersome to handle. For example, imagine making a list of all animals, $\Omega = \{elephant, giraffe, ...\}$. Other times, Ω is not explicitly identified. Hence we want to translate from the world of outcomes to a more familiar and measurable space, like \mathbb{R} . That is essentially why we use random variables.

Definition 3.5 (Random variable). A random variable is a mapping $x : \Omega \to \mathbb{R}$.

For example, we could define a mapping x that maps $elephant \mapsto 1$, $giraffe \mapsto 2$, etc. Since \mathbb{P} specifies a probability for every $A \in \mathcal{A}$, it also induces a probability in terms of x. For instance, the event $\{x \leq 0\}$ is equivalent to the event $\{\omega \in \Omega : x(\omega) \leq 0\}$, and $\mathbb{P}(x \leq 0) = \mathbb{P}(\{\omega \in \Omega : x(\omega) \leq 0\})$. Continuing with our example, $P(x \leq 2) = \mathbb{P}(\{elephant, giraffe\})$.

3.7 Densities

Intuitively, a probability measure \mathbb{P} is the rule that assigns probability to the events in \mathcal{A} . For instance, in Example 3.4, \mathbb{P} assigns an equal probability of 1/6 to each of the possible outcomes $x = 1, \ldots, 6$. To calculate the probability of an event, all we would have to do is compute

$$\mathbb{P}(x \in A) = \sum_{\mathbf{x} \in A} \mathbb{P}(\mathbf{x}).$$

This was easily done because x was discrete. In this case, \mathbb{P} is called a *mass function*. Notice that x is used for a random variable, whereas x is used for a specific value that x may take. For discrete random variables, $\mathbb{P}(x)$ is essentially shorthand for $\mathbb{P}(x = x)$. If x were continuous, the probability that it takes a particular value is zero. So instead of a mass function, we use a *density function* \mathbb{P} that indicates the probability that x falls within certain intervals. Then

$$\mathbb{P}(x \in A) = \int_A \mathbb{P}(\mathbf{x}) d\mathbf{x}.$$

In general, probability mass functions and densities depend on certain parameters $\boldsymbol{\theta}$. Whenever this want to be made noticed explicitly, we use the notation $\mathbb{P}(\cdot|\boldsymbol{\theta})$.



Figure 3.1: Normal density function $\mathbb{P}(\mathbf{x}|\boldsymbol{\theta}) = \mathcal{N}(\mu, \sigma^2)$. In this case, $\boldsymbol{\theta}$ denotes the parameters corresponding to the mean μ and the variance σ^2 .

Example 3.8. The height of a person can be modeled as a Normal random variable with mean $\mu = 5'5"$ for females and $\mu = 5'10"$ for males (see Figure 3.1). However, the probability that your height is *exactly* 5'5" or 5'10" is zero. You are more likely to be somewhere in between (5'3", 5'7") or (5'8", 5'12").

3.8 Expectation

Definition 3.6 (Expectation). For a continuous random variable x and an arbitrary function f(x),

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

If x is a discrete random variable,

$$\mathbb{E}[f(x)] := \sum_{\mathbf{x}} f(\mathbf{x}) \mathbb{P}(\mathbf{x}).$$

Example 3.9. Special cases of expectations:

- **Probability:** $\mathbb{P}(x \in A) = \mathbb{E}[\mathbb{1}_{x \in A}].$
- Mean: $\mu := \mathbb{E}[x]$.
- Variance: $\sigma^2 := \mathbb{E}[(x \mu)^2].$

Definition 3.7 (Conditional expectation). The *conditional expectation* of a continuous random variable f(x) given a random variable y is defined as

$$\mathbb{E}[f(x)|y] := \int f(\mathbf{x})\mathbb{P}(\mathbf{x}|y)d\mathbf{x},$$

and if x is a discrete random variable,

$$\mathbb{E}[f(x)|y] := \sum_{\mathbf{x}} f(\mathbf{x})\mathbb{P}(\mathbf{x}|y).$$

Notice that $\mathbb{E}[f(x)|y]$ is a function of the random variable y, and hence it is also a random variable. Notice the difference between $\mathbb{E}[f(x)|y]$ and $\mathbb{E}[f(x)|y]$, which is no longer a random variable, because the random variable y is already known to have taken the value y.

3.9 Common Probability Measures

Tables 3.1 and 3.2 give examples of common probability measures, also known as distributions.

Discrete	Parameters $(\boldsymbol{\theta})$	$\mathbb{P}(x = \mathbf{x})$	$\mathbb{E}[x]$	var(x)
Bernoulli	р	$\mathbb{P}(x=1) = \mathbf{p}, \mathbb{P}(x=0) = 1 - \mathbf{p}$	р	p(1 - p)
Binomial	\mathbf{n},\mathbf{p}	$\mathbb{P}(x = \mathbf{x}) = \binom{\mathbf{n}}{\mathbf{x}} \mathbf{p}^{\mathbf{x}} (1 - \mathbf{p})^{\mathbf{n} - \mathbf{x}}, \mathbf{x} = 0, \dots, \mathbf{n}$	np	np(1-p)
Poisson	λ	$\mathbb{P}(x = \mathbf{x}) = \frac{\lambda^{\mathbf{x}} e^{-\lambda}}{\mathbf{x}!}, \mathbf{x} = 0, 1, \dots$	λ	λ

Table 3.1: Examples of common probability mass functions.

3.10 Multivariate distributions

In many modern applications it is convenient to arrange random variables in vectors. For example, we might consider a *random vector*

$$oldsymbol{x} = egin{bmatrix} x_1 \ x_2 \ x_3 \end{bmatrix}$$

containing the information of a person's height, weight and cholesterol level. If the random variables in the vector are independently distributed, then its *joint* distribution is just the product of the univariate distributions of each component. In our example, $\mathbb{P}(\mathbf{x}) = \mathbb{P}(x_1, x_2, x_3)$ would simply factor into $\mathbb{P}(x_1)\mathbb{P}(x_2)\mathbb{P}(x_3)$. However, if the random variables in the vector are dependent (as is actually the case with height, weight and cholesterol level), then $\mathbb{P}(\mathbf{x})$ does not factor in this simple way.

Multivariate densities model dependent random variables. The one we will use most in this course is the multivariate Normal $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with mean vector $\boldsymbol{\mu} \in \mathbb{R}^{D}$ and covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$, which has the following form:

$$\mathbb{P}(\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi)^{\mathrm{D}}|\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})}.$$

Continuous	Parameters $(\boldsymbol{\theta})$	$\mathbb{P}(\mathrm{x})$	$\mathbb{E}[x]$	var(x)	
Uniform	$\mathrm{a,b}$	$\mathbb{P}(x) = \frac{1}{b-a}, a \leq x \leq b$	$\frac{a+b}{2}$	$\frac{(\mathrm{b}-\mathrm{a})^2}{12}$	
Exponential	λ	$\mathbb{P}(\mathbf{x}) = \lambda e^{-\lambda \mathbf{x}}, \mathbf{x} \ge 0$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	
Laplace	λ	$\mathbb{P}(\mathrm{x}) = rac{\lambda}{2} e^{-\lambda \mathrm{x} }$	0	$rac{2}{\lambda^2}$	
Normal	μ,σ^2	$\mathbb{P}(\mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2}\left(\frac{\mathbf{x}-\mu}{\sigma}\right)^2}$	μ	σ^2	
Gamma	lpha,eta	$\mathbb{P}(\mathbf{x}) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} \mathbf{x}^{\alpha-1} e^{-\beta \mathbf{x}}$ $\Gamma(\alpha) := \int_0^\infty \mathbf{x}^{\alpha-1} e^{-\mathbf{x}} d\mathbf{x}$	$\frac{lpha}{eta}$	$\frac{lpha}{eta^2}$	
Beta	lpha,eta	$\mathbb{P}(\mathbf{x}) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \mathbf{x}^{\alpha - 1} (1 - \mathbf{x})^{\beta - 1}$	$\frac{\alpha}{\alpha+\beta}$	$\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$	
χ^2	k	Gamma(^k /2, ¹ /2)			
F	lpha,eta	$\mathbb{P}(\mathbf{x}) = \frac{\sqrt{\frac{(\alpha \mathbf{x})^{\alpha} \beta^{\beta}}{(\alpha \mathbf{x} + \beta)^{\alpha + \beta}}}}{\mathbf{x} \mathcal{B}(\alpha/2, \beta/2)}$ $\mathcal{B}(\alpha, \beta) := \int_{0}^{1} \mathbf{x}^{\alpha - 1} (1 - \mathbf{x})^{\beta - 1} d\mathbf{x}$	$\frac{\beta}{\beta-2},\beta>2$	$\frac{2\beta^2(\alpha+\beta-2)}{\alpha(\beta-2)^2(\beta-4)}, \ \beta > 4$	

Table 3.2: Examples of common probability density functions.

3.11 Probability and Likelihood

A probability (mass or distribution) $\mathbb{P}(x|\boldsymbol{\theta})$ determines the frequency with which that a random variable x takes each value, given some parameter $\boldsymbol{\theta}$. For example, if $x \sim \text{Bernoulli}(\boldsymbol{\theta})$, with $\boldsymbol{\theta} = 1/2$, then the probability that x takes the value 1 is $\mathbb{P}(x = 1|\boldsymbol{\theta}) = \boldsymbol{\theta} = 1/2$.

Conversely, the *likelihood* $\mathbb{P}(\mathbf{x}|\boldsymbol{\theta})$ determines the probability that a parameter $\boldsymbol{\theta}$ was the one that generated a sample x. We emphasize this distinction using x instead of x, to indicate that x is already known, i.e., observed data that has already taken a specific value. Under the same Bernoulli example, if we observe $\mathbf{x} = 1$, then the likelihood of the parameter $\boldsymbol{\theta}$ is $\mathbb{P}(\mathbf{x} = 1|\boldsymbol{\theta}) = \boldsymbol{\theta}$.

The probability and the likelihood may *look* a lot alike. The difference is very subtle, and mainly conceptually: the probability $\mathbb{P}(x|\theta)$ is a function where x is the variable, and θ is fixed. In contrast, the likelihood $\mathbb{P}(x|\theta)$

is a function where θ is the variable, and x is fixed. We use $\mathbb{P}(x|\theta)$ when we know θ and want to guess x; we use $\mathbb{P}(x|\theta)$ when we have already observed data with the specific value x, and we want to guess the parameter θ that generated it.

Example 3.10. Suppose x_1, \ldots, x_6 are *independently and identically distributed* (i.i.d.) according to a Bernoulli(¹/₄) distribution. Then the probability that $x_1 = x_2 = x_3 = 1$, and $x_4 = x_5 = x_6 = 0$ is:

$$\mathbb{P}(x_1 = x_2 = x_3 = 1, x_4 = x_5 = x_6 = 0|\theta) = \prod_{i=1}^3 \mathbb{P}(x_i = 1|\theta) \cdot \prod_{i=4}^6 \mathbb{P}(x_i = 0|\theta)$$
$$= \theta^3 (1-\theta)^3 = (1/4)^3 (3/4)^3.$$

Instead, suppose that we observe $x_1 = x_2 = x_3 = 1$, and $x_4 = x_5 = x_6 = 0$. Then the likelihood of θ under this sample is:

$$\begin{split} \mathbb{P}(x_1 = x_2 = x_3 = 1, x_4 = x_5 = x_6 = 0 | \theta) &= \prod_{i=1}^{3} \mathbb{P}(x_i = 1 | \theta) \cdot \prod_{i=4}^{6} \mathbb{P}(x_i = 0 | \theta) \\ &= \theta^3 (1 - \theta)^3. \end{split}$$

Based on this sample, which would be your intuitive best guess at the value of θ ? Is this the same value that maximizes the likelihood $\mathbb{P}(\mathbf{x}_1, \ldots, \mathbf{x}_6 | \theta)$?

3.12 Sums of Independent Random Variables

In many applications we want to know the distribution of the sum of independent random variables. Table 3.3 gives a few examples.

Example 3.11. Suppose there is an epidemic in a city with N habitants. The ith person will independently contract the disease with probability p. We can model this as $x_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$, i = 1, ..., N. Let $m = \sum_{i=1}^{N} x_i$ be the number of people that get infected. The Center for Disease Control wants to determine $\mathbb{P}(m > m)$. This raises the question: what is the distribution of m? Notice that m is the sum of i.i.d. Bernoulli random variables. However, m is clearly not Bernoulli. To begin with, m can take values in $\{0, \ldots, N\}$, while a Bernoulli random variable can only take values in $\{0, 1\}$. So the question is: what is the distribution of a sum of N i.i.d. Bernoulli(p) random variables?

3.13 Other Common Functions of Random Variables

In addition to sums, other common functions of random variables include

- Linear multiplication: For a constant matrix $\mathbf{A} \in \mathbb{R}^{M \times D}$ and a random vector $\boldsymbol{x} \in \mathbb{R}^{D}$, $\mathsf{cov}(\mathbf{A}\boldsymbol{x}) = \mathbf{A}\mathsf{cov}(\boldsymbol{x})\mathbf{A}^{\mathsf{T}}$.
- Squared Normal: If $x \sim \mathcal{N}(0, 1)$, then $x^2 \sim \chi^2$.
- Ratio of χ^2 's: If $x \sim \chi^2(\mathbf{k})$ is independent of $y \sim \chi^2(\ell)$, then $\frac{\ell x}{\mathbf{k}y} \sim F(\mathbf{k}, \ell)$.

x_{i}	$\sum_{\mathrm{i}=1}^{\mathrm{N}} x_{\mathrm{i}}$
Bernoulli(p)	Exercise
$Binomial(n_i,p) \\$	$Binomial \left(\sum_{i=1}^N n_i, p \right)$
$Poisson(\lambda_i)$	$\operatorname{Poisson}\left(\sum_{i=1}^N \lambda_i\right)$
$\exp(\lambda)$	$\operatorname{Gamma}(\mathrm{N},\lambda)$
$\operatorname{Gamma}(\alpha_{\mathrm{i}},\beta)$	$Gamma\left(\sum_{i=1}^{N} \alpha_{i}, \beta\right)$
$\mathcal{N}(\mu_i,\sigma_i^2)$	$\mathcal{N}\left(\sum_{i=1}^{N} \mu_{i}, \sum_{i=1}^{N} \sigma_{i}^{2} ight)$
$\chi^2({ m k_i})$	$\chi^2\left(\sum_{i=1}^N k_i\right)$

Table 3.3: Examples of distributions of sums of **independent** random variables.

3.14 Simulations

Often it is useful to generate random numbers to simulate data. For example, if I want to simulate a dataset containing the weight of a population of 1000 individuals, I can generate normal numbers centered at 178 pounds, with a standard deviation of 25 pounds. I can do this very easily in Matlab using:

```
1 data = 178+5*randn(1000,1);
2 histogram(data);
3 set(gca,'FontSize',18);
4 ylabel('Frequency');
5 xlabel('Weight');
6 title('Normal(178,25)');
```

which should produce the following histogram:



Here the function randn(m,n) will generate an $m \times n$ matrix with normal numbers (he *n* in *randn* stands for *normal*). You can similarly generate random numbers of other distributions, like uniform or exponential:

```
subplot(1,2,1);
1
2
  data = rand(1000, 1);
3 histogram(data);
  set(gca, 'FontSize', 18);
4
   ylabel('Frequency');
\mathbf{5}
   xlabel('Value');
6
   title('Uniform(0,1)');
7
8
9
   subplot(1,2,2);
10
   data = exprnd(1,1000,1);
11 histogram(data);
  set(gca, 'FontSize', 18);
12
  xlabel('Value');
13
14
   title('Exponential(1)');
```

which should produce the following histogram:

