

## Topic 2: Review of Probability Theory

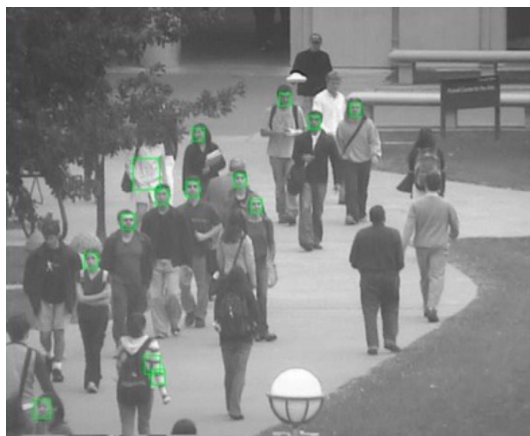
INSTRUCTOR: DANIEL L. PIMENTEL-ALARCÓN

© COPYRIGHT 2017

## 2.1 Why Probability?

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

**Example 2.1.** 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:

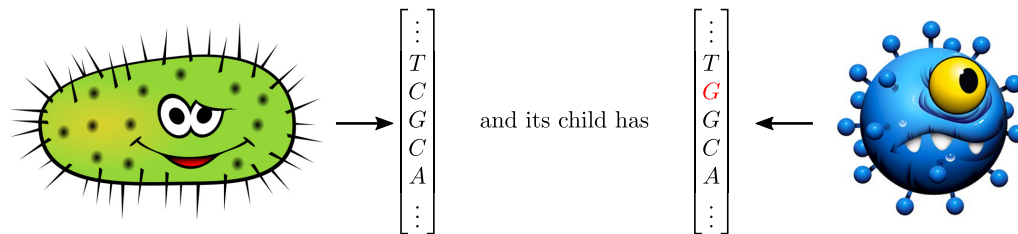


You can predict where some of these people will be in a few seconds (based on their current locations and directions). However, you cannot be certain. What if someone trips? What if the wind moves the leaves of a tree and produces occlusions?) All of this can be modeled probabilistically in a precise 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 2.2.** All organisms have a DNA sequence (genome), i.e., a very long vector of nucleotides ( $A, C, G, T$ ). For reasons that nobody knows *yet*, sometimes organisms *mutate*. That is, when organisms

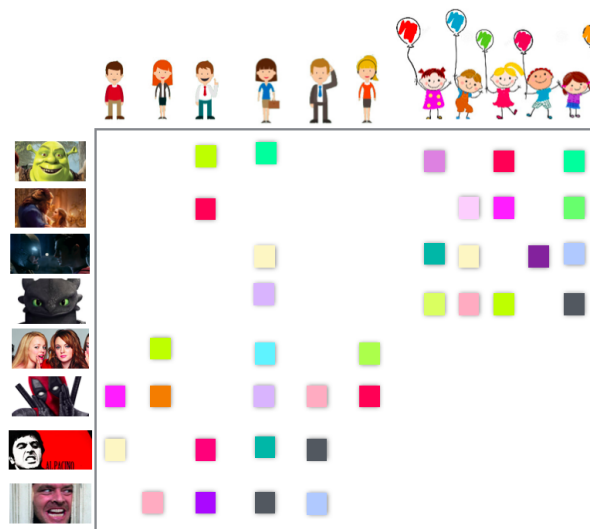
reproduce, some of their nucleotides get changed. For example, it is possible that a parent has the sequence



The location of these mutations are random, but more common in certain regions of the genome (called genes). This can be modeled probabilistically. For instance, we can say that the mutation may appear uniformly at random in each gene, but some genes are more likely to present a mutation 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 problem. The reasons are not as apocalyptic as a zombi virus, but close. 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.

**Example 2.3.** 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.

## 2.2 The Basics

**Definition 2.1.** There are three elemental concepts in basic probability theory:

$\Omega$  := Sample space = set of all possible outcomes.

$\mathcal{A}$  := Set of all possible events.

$P$  := Probability measure.

**Example 2.4.** Consider a fair die. Then

$$\Omega = \{1, 2, 3, 4, 5, 6\}.$$

$$\mathcal{A} = \left\{ \{1\}, \dots, \{6\}, \{1, 2\}, \dots, \{5, 6\}, \dots, \{1, 2, 3, 4, 5, 6\} \right\}$$

$$P(k) = 1/6 \text{ for every } k = 1, \dots, 6.$$

**Definition 2.2** (Probability measure). A mapping  $P : \mathcal{A} \rightarrow [0, 1]$  is a *probability measure* if it satisfies the next properties:

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

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

## 2.3 Conditional Probability

**Definition 2.3** (Conditional probability). Let  $A, B \in \mathcal{A}$ . The *conditional probability* that  $A$  occurred given  $B$  occurred is

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

**Example 2.5.** Continuing with Example 2.4, let  $A = \{1, 2\}$ ,  $B = \{2, 3\}$ . The probability that  $A$  occurs is  $P(A) = 1/3$ . However, if you already know that  $B$  occurred, then the probability that  $A$  also occurred increases to

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

## 2.4 Independence

**Definition 2.4** (Independent events). Let  $A, B \in \mathcal{A}$ . We say  $A$  and  $B$  are *independent* if  $P(A|B) = P(A)$ .

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

**Example 2.6.** 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

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

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

## 2.5 Bayes Rule

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

**Definition 2.5** (Bayes rule). Let  $A, B \in \mathcal{A}$ . Then

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

Bayes rule plays a crucial role in modern applications.

**Example 2.7.** Geneticists have determined that 90% of the people with disease  $B$  have gene  $A$  active, i.e.,  $P(A|B) = 0.9$ . 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  $P(B|A)$ ? At first glance you might think it is very likely that you will develop disease  $B$ . However, to determine this you need to know  $P(A)$  and  $P(B)$ . Of the whole population, if only 5% have disease  $B$ , while 45% have gene  $A$  active, what is  $P(B|A)$ ? This is a simple application of Bayes rule:

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

## 2.6 Random Variables

Sometimes the sample space  $\Omega$  contains elements that are cumbersome to handle. For example, imagine making a list of all animals,  $\Omega = \{elephant, giraffe, \dots\}$ . Other times,  $\Omega$  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 2.6** (Random variable). A random variable is a mapping  $x : \Omega \rightarrow \mathbb{R}$ .

For example, we could define a mapping  $x$  that maps  $elephant \mapsto 1$ ,  $giraffe \mapsto 2$ , etc. Since  $P$  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  $P(x \leq 0) = P(\{\omega \in \Omega : x(\omega) \leq 0\})$ . Continuing with our example,  $P(x \leq 2) = P(\{elephant, giraffe\})$ .

## 2.7 Densities

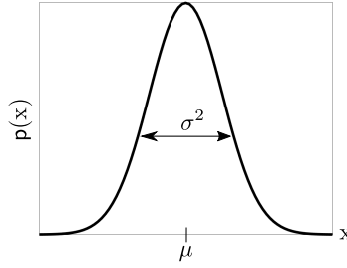
Intuitively, a probability measure  $P$  is the rule that assigns probability to the events in  $\mathcal{A}$ . For instance, in Example 2.4,  $P$  assigns an equal probability of  $1/6$  to each of the possible outcomes  $k = 1, \dots, 6$ . To calculate the probability of an event, all we would have to do is compute

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

This was easily done because  $x$  was discrete. In this case,  $P$  is called a *mass function*.

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*  $p$  that indicates the probability that  $x$  falls within certain intervals. Then

$$P(x \in A) = \int_A p(x) dx.$$

Figure 2.1: Gaussian density function  $p(x)$  with mean  $\mu$  and variance  $\sigma^2$ .

**Example 2.8.** The height of a person can be modeled as a gaussian random variable with mean  $\mu = 5'5''$  for females and  $\mu = 5'10''$  for males (see Section 2.9 and Figure 2.1). 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  $(5'8'', 5'12'')$ .

## 2.8 Expectation

**Definition 2.7** (Expectation). For a continuous random variable  $x$  and an arbitrary function  $f(x)$ ,

$$E[f(x)] := \int f(x)p(x)dx.$$

If  $x$  is a discrete random variable,

$$E[f(x)] := \sum_k f(x_k)P(x = x_k).$$

**Example 2.9.** Special cases of expectations:

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

**Definition 2.8** (Conditional expectation). The *conditional expectation* of a continuous random variable  $f(x)$  given a random variable  $y$  takes the value  $y$  is defined as

$$E[f(x)|y = y] := \int f(x)p(x|y)dx,$$

and if  $x$  is a discrete random variable,

$$\mathbb{E}[f(x)|y = y] := \sum_k f(x_k)P(x = x_k|y),$$

## 2.9 Common Probability Measures

Tables 2.1 and 2.2 give examples of common probability measures, also known as *distributions*.

## 2.10 Multivariate distributions

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

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

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,  $\mathbf{p}(\mathbf{x}) = \mathbf{p}(x_1, x_2, x_3)$  would simply factor into  $\mathbf{p}(x_1)\mathbf{p}(x_2)\mathbf{p}(x_3)$ . However, if the random variables in the vector are dependent (as is actually the case with height, weight and cholesterol level), then  $\mathbf{p}(\mathbf{x})$  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 Gaussian  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  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:

$$\mathbf{p}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})}.$$

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

Table 2.1: Examples of common probability mass functions.

Continuous	$p(x)$	$E[x]$	$\text{var}(x)$
Uniform[a, b]	$p(x) = \frac{1}{b-a}, a \leq x \leq b$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
Exponential( $\lambda$ )	$p(x) = \lambda e^{-\lambda x}, x \geq 0$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
Laplace( $\lambda$ )	$p(x) = \frac{\lambda}{2} e^{-\lambda x }$	0	$\frac{2}{\lambda^2}$
Gaussian or Normal( $\mu, \sigma^2$ )	$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$	$\mu$	$\sigma^2$
Gamma( $\alpha, \beta$ )	$p(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$ $\Gamma(\alpha) := \int_0^\infty x^{\alpha-1} e^{-x} dx$	$\frac{\alpha}{\beta}$	$\frac{\alpha}{\beta^2}$
Beta( $\alpha, \beta$ )	$p(x) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$	$\frac{\alpha}{\alpha+\beta}$	$\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$
$\chi^2(k)$	Gamma( $k/2, 1/2$ )		
$F(\alpha, \beta)$	$p(x) = \frac{\sqrt{\frac{(\alpha x)^\alpha \beta^\beta}{(\alpha x + \beta)^{\alpha+\beta}}}}{x \mathcal{B}(\alpha/2, \beta/2)}$ $\mathcal{B}(\alpha, \beta) := \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$	$\frac{\beta}{\beta-2}, \beta > 2$	$\frac{2\beta^2(\alpha+\beta-2)}{\alpha(\beta-2)^2(\beta-4)}, \beta > 4$

Table 2.2: Examples of common probability density functions.



## 2.11 Sums of Independent Random Variables

In many applications we want to know the distribution of the sum of independent random variables. Table 2.3 gives a few examples.

**Example 2.10.** Suppose there is an epidemic in a city with  $N$  habitants. The  $i^{\text{th}}$  person will independently contract the disease with probability  $p$ . We can model this as  $x_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$ ,  $i = 1, \dots, N$ . Let  $m = \sum_{i=1}^N x_i$  be the number of people that get infected. The Center for Disease Control wants to determine  $P(m > m)$ . 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  $\{0, \dots, N\}$ , 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?

## 2.12 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  $\mathbf{x} \in \mathbb{R}^D$ ,  $\text{cov}(\mathbf{Ax}) = \mathbf{Acov}(\mathbf{x})\mathbf{A}^\top$ .
- **Squared gaussians:** If  $x \sim \mathcal{N}(0, 1)$ , then  $x^2 \sim \chi^2$ .
- **Ratio of  $\chi^2$ 's:** If  $x \sim \chi^2(k)$  is independent of  $y \sim \chi^2(\ell)$ , then  $\frac{\ell x}{ky} \sim F(k, \ell)$ .

$x_i$	$\sum_{i=1}^N x_i$
Bernoulli(p)	<b>Exercise</b>
Binomial( $n_i, p$ )	Binomial $\left(\sum_{i=1}^N n_i, p\right)$
Poisson( $\lambda_i$ )	Poisson $\left(\sum_{i=1}^N \lambda_i\right)$
$\exp(\lambda)$	Gamma( $N, \lambda$ )
Gamma( $\alpha_i, \beta$ )	Gamma $\left(\sum_{i=1}^N \alpha_i, \beta\right)$
$\mathcal{N}(\mu_i, \sigma_i^2)$	$\mathcal{N}\left(\sum_{i=1}^N \mu_i, \sum_{i=1}^N \sigma_i^2\right)$
$\chi^2(k_i)$	$\chi^2\left(\sum_{i=1}^N k_i\right)$

Table 2.3: Examples of distributions of sums of **independent** random variables.