

# Note Set 1: Review of Basic Concepts in Probability

Padhraic Smyth,  
Department of Computer Science  
University of California, Irvine  
January 2019

This set of notes is intended as a brief refresher on probability. As a student reading these notes you will likely have seen (in other classes) most or all of the ideas discussed below. Nonetheless, even though these ideas are relatively straightforward, they are the key building blocks for more complex ideas in probabilistic learning. Thus, it is important to be familiar with these ideas in order to understand the material we will be discussing later in the course.

## 1 Discrete Random Variables and Distributions

Consider a variable  $A$  that can take a finite set of values  $\{a_1, \dots, a_m\}$ . We can define a **probability distribution** for  $A$  by specifying a set of numbers  $\{P(A = a_1), \dots, P(A = a_m)\}$ , where  $0 \leq P(A = a_i) \leq 1$  and where  $\sum_{i=1}^m P(A = a_i) = 1$ . We can think of  $A = a_i$  as an **event** that either is or is not true, and  $P(A = a_i)$  is the probability of this particular event being true. We can also think of  $A = a_i$  as a **proposition** or **logical statement** about the world that is either true or false. Probability expresses our uncertainty about whether the proposition  $A = a_i$  is true or is false.

Notational comment: for convenience, we follow a typical convention in probability notation and will often shorten  $P(A = a_i)$  to just  $P(a_i)$ , and will also often use  $a$  (and  $P(a)$ ) to indicate a generic value for  $A$ , i.e., one of the possible  $a_i$  values.

We call  $A$  a *random variable* as long as the set of values  $\{a_1, \dots, a_m\}$  are **mutually exclusive and exhaustive**. *Mutually exclusive* means that the random variable  $A$  can only take one value at a time—another way to state this is that the  $P(A = a_i \text{ AND } A = a_j) = 0, \forall i, j, i \neq j$ , i.e., the probability of the event occurring where  $A$  takes two different values is 0 (i.e., it is deemed impossible)<sup>1</sup>. *Exhaustive* means

---

<sup>1</sup>the symbol  $\forall$  means “for all”

that the variable  $A$  always takes one of the values in the set  $\{a_1, \dots, a_m\}$ , i.e., there are no other values it can take.

EXAMPLE 1: In medical diagnosis we often want to be able to predict whether a patient has a particular disease or not, given other measurements—this is a particular type of prediction problem known as **classification** that we will discuss later in more detail. Let  $C$  be a variable representing what disease a particular patient has. Consider for example  $C$  taking values in  $\{c_1, c_2\}$ , where  $c_1 = \textit{has the flu}$  and  $c_2 = \textit{does not have the flu}$ . In this case  $C$  is a random variable since  $c_1$  and  $c_2$  are clearly mutually exclusive and exhaustive, and we will be uncertain in general about the actual true value of  $C$  for any particular patient.

Now consider another set of events:  $d_1 = \textit{has the flu}$  and  $d_2 = \textit{has malaria}$  and  $d_3 = \textit{healthy}$ . Could  $\{d_1, d_2, d_3\}$  be used as the set of events to define a random variable? The answer is no since the set of events are not mutually exclusive: a person could have both the flu and malaria (this have a very small probability of happening, but this probability is not zero, i.e., it is possible). Nor are the events exhaustive: a person could be in none of these states since they could in general have some condition other than flu or malaria.

EXAMPLE 2: Consider modeling the probability distribution of English words in a particular set of text documents. Let  $W$  be the word variable and let  $w_i$  be a particular word. For example, we might want to know the probability that a particular word  $w_i$  will be typed by a user in a text editor or will be spoken by a user. Such probabilities are used in practice in a variety of applications such as speech recognition, automated machine translation, and so forth. In order to treat  $W$  as a random variable we need to ensure that the words being spoken or written by a user are *mutually exclusive* and *exhaustive* by defining  $W$  appropriately. For example, for mutual exclusion, we could define  $W$  to be the next word in a sequence of words, which by definition means there will always be a single unique “next word”—and we would need to include in our vocabulary a special symbol for the end of a sequence. The requirement for exhaustivity is more difficult to satisfy. We could define the set of possible words  $\{w_1, \dots, w_m\}$  as the set of all words in an English dictionary—but which dictionary should we use? and what about words not in the dictionary such as regional variations of words, and proper nouns such as “California”? This is a problem commonly faced with text modeling, since it is impossible to know all future words that might occur in practice. One way to get around this is to limit the set of words being modeled to (say) the  $m = 20,000$  most frequent words in a particular set of text and then an additional symbol  $w_{20,001}$  is defined to represent “all other words,” ensuring that  $\sum_i P(w_i) = 1$ . A more interesting problem, that we will not discuss at this point, is how much probability mass we should assign to the event “all other words,” or equivalently, how likely are we to encounter previously unseen words in the future?

Before we finish with discrete random variables we note that not all discrete random variables necessarily take values from a finite set, but instead could take values from a countably infinite set. For example, we can define a random variable  $A$  taking values from the set of positive integers  $\{1, 2, 3, \dots\}$ . The sum  $\sum_{i=1}^{\infty} P(A = i)$  must converge to 1 of course.

For random variables taking values in a countably infinite set it is impossible to explicitly represent a probability distribution directly by a table of numbers—and even in the finite case it may be often inconvenient or inefficient. In such cases we use a *parametric model* for the distribution, i.e., a function that describes how  $P(i)$  varies as a function of  $i$  and as a function of one or more **parameters** of the model.

EXAMPLE 3: An example of a probability distribution defined on the positive integers is the geometric distribution,

$$P(A = i) = (1 - \alpha)^{i-1}\alpha, \quad i = 1, 2, \dots$$

where  $\alpha$  is a parameter of the model and  $0 < \alpha < 1$ . This can be used for example to describe the distribution of the number of consecutive “tails” that we will see before we see the first “heads event” in a coin-tossing experiment, where  $\alpha$  is the probability of a head occurring on each coin toss.

## 2 Continuous Random Variables and Density Functions

The variables we discussed above such as  $A$  and  $C$  take values that can be put in one-to-one correspondence with sets of integers, and are often referred to as **discrete random variables**. It is also useful to be able to build and use probability models for **real-valued or continuous random variables**, e.g., a random variable  $X$  that can take values  $x$  anywhere on the real line.

Generally speaking, where we have sums for discrete random variables we usually will have integrals for continuous random variables. So, instead of the requirement that  $\sum_i P(x_i) = 1$  we have  $\int p(x)dx = 1$ . Here  $p(x)$  is the **probability density function** for the variable  $X$ , where  $p(x) \geq 0, \forall x$ . We can calculate the probability that the variable  $X$  lies between 2 values  $a$  and  $b$  by integrating  $p(x)$  between  $a$  and  $b$ , i.e.,  $P(a \leq X \leq b) = \int_a^b p(x)dx$ .

EXAMPLE 4: The **uniform density** is defined as  $p(x) = \frac{1}{b-a}$  for  $a \leq x \leq b$  and 0 otherwise. We sometimes use the notation  $X \sim U(a, b)$  to denote that the random variable  $X$  has a particular type of distribution (in this case the uniform distribution  $U$ ). Here  $a$  and  $b$  are the **parameters** of the uniform density function.

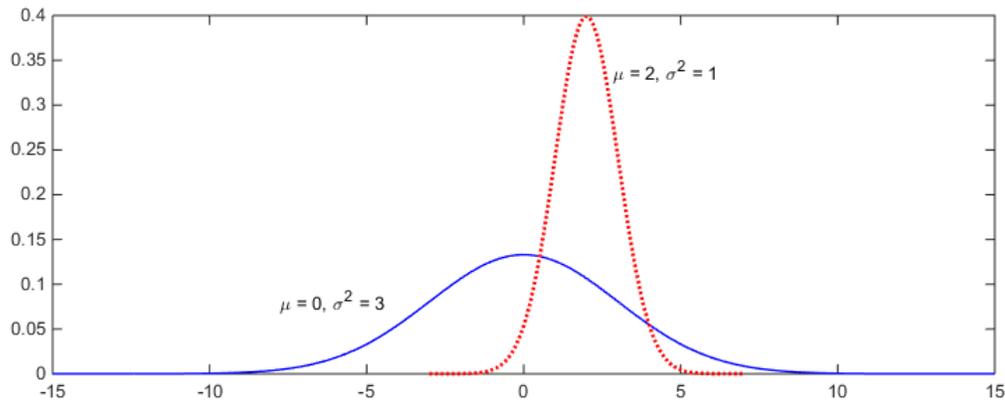


Figure 1: An example of two Gaussian density functions with different parameter values.

Note that while we had the option of describing a discrete probability distribution (for a finite set) with a list or table of probabilities, we can't define a continuous density function this way—instead we must parametrize the function and describe it via its functional form and its parameters.

Note also that although  $\int p(x)dx = 1$  this does not imply that  $p(x)$  needs to be less than 1! For example, for the uniform density if  $a = 0$  and  $b = 0.1$ , then  $p(x) = \frac{1}{0.1-0.0} = 10$  for  $a \leq x \leq b$ . The key point is that the **area** under the density function,  $p(x)$ , is constrained to be 1, but the height of the function can be any non-negative quantity.

EXAMPLE 5: The Gaussian or Normal density function is defined as:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-[\frac{1}{2\sigma^2}(x-\mu)^2]}.$$

Here the parameters are the mean  $\mu$  and the variance  $\sigma^2$ . We can also say  $X \sim N(\mu, \sigma^2)$  for shorthand. Figure 1 shows an example of two Gaussians, each with different parameters—the dotted one is much narrower and concentrated than the solid one. We will examine the Gaussian density function in more detail later in this class.

### 3 Multiple Random Variables

#### 3.1 Conditional Probabilities

Modeling single random variables is a useful starting point but the real power in probabilistic modeling comes from modeling sets of random variables. Assume for the moment that we have two variables  $A$  and

$B$ . We define the **conditional probability**  $P(A = a|B = b) = P(a|b)$  as the probability that  $A$  takes value  $a$  given that  $B$  takes value  $b$ . The vertical “bar”  $|$  means that the probability of the proposition to the left of the bar (in this case  $A = a$ ) is conditioned on knowing or assuming that the proposition to the right of the bar (in this case  $B = b$ ) is true.

We can define a *conditional probability distribution* over all values for  $A$ , i.e., the list of probabilities  $P(a_1|b), P(a_2|b), \dots, P(a_m|b)$ , for some value  $b$  of the variable  $B$ . A key point is that a conditional probability distribution is a standard probability distribution for the variable that is the argument of the distribution function (here  $A$ ), but now conditioned on some event (here  $B = b$ ). In particular,  $\sum_{i=1}^m P(a_i|b) = 1$ .

There are two obvious interpretations of what a conditional probability means. In the first case,  $P(a|b)$  could be interpreted as “having measured  $B = b$  this is the probability of  $A = a$  now that we know  $B = b$ ”. The second case allows for hypothetical reasoning, i.e.,  $P(a|b)$  can be interpreted as the probability of  $A = a$  **if** hypothetically  $B = b$ . The value of the conditional probability  $P(a|b)$  and the associated mathematics will be the same for both cases, but the interpretation is a little different.

In many respects conditional probabilities are the key central concept in probabilistic information processing. A conditional probability distribution  $P(A|B)$  gives us a quantitative way to represent how  $B$  provides information about  $A$ . We will see later for example, that when a variable  $B$  provides no information about another variable  $A$ , i.e.,  $p(A|B) = p(A)$ , that this statement is equivalent to saying that the two variables  $A$  and  $B$  are independent. Information and dependence are closely intertwined concepts in probability modeling.

### 3.2 Joint Probability Models

We can define the joint probability  $P(a, b)$ , which is short for  $P(A = a \text{ AND } B = b)$ , to represent the probability that variable  $A$  takes value  $a$  and variable  $B$  takes value  $b$ . Again note that the terms inside the parentheses, i.e.  $a, b$  in  $P(a, b)$ , represent a logical statement about the world, namely that  $A$  takes value  $a$  and variable  $B$  takes value  $b$ . We can think of the “joint variable”  $AB$  as taking values from the cross-product of the sets values of  $A$  and  $B$ —it is sometimes conceptually useful to think of  $AB$  as a “super-variable” defined as the cross-product of individual variables  $A$  and  $B$ . By definition, since  $A$  and  $B$  are random variables, then one can easily show that the values of the joint variable  $AB$  must also be mutually exclusive and exhaustive, and we have that

$$\sum_{i=1}^m \sum_{j=1}^n P(a_i, b_j) = 1.$$

Note that if  $A$  takes  $m$  values and  $B$  takes  $n$  values, then there are  $m \times n$  different possible combinations of  $a$  and  $b$  values, and the joint probability table will contain  $mn$  numbers, which could be very large for large values of  $m$  and  $n$ . This hints at a combinatorial issue that arises when we model multiple variables, namely that the number of entries in the joint distribution will grow exponentially fast with the number of variables in the model.

### 3.3 Relating Conditional and Joint Probabilities

From the basic axioms of probability one can show straightforwardly that the joint probability  $P(a, b)$  is related to the conditional probability  $P(a|b)$  and the probability  $P(b)$  (often referred to as the **marginal probability** of  $b$ ) in the following manner:

$$P(a, b) = P(a|b)P(b)$$

This is always true and is one of the basic rules of probability. There is a simple intuitive interpretation in words, namely, the probability that  $A$  takes value  $a$  and  $B$  takes value  $b$  can be decomposed into (1) the probability that  $A$  takes value  $a$  given that  $B$  takes value  $b$ , (2) times the probability that  $B = b$  is true. This argument makes intuitive sense (although its not a formal proof).

### 3.4 More than 2 Variables

We can directly extend our definitions above beyond just 2 variables, to 3, 4, and indeed thousands of variables. As an example, say we have 4 random variables  $A, B, C, D$ .

- We can define the conditional probability  $P(a|b, c, d)$ , which is interpreted as the conditional probability that  $A$  takes value  $a$ , given that we know (or assume) that  $B = b$  and  $C = c$  and  $D = d$ . (Note again that the comma “,” notation is used as shorthand for the conjunction AND).
- We can define a joint distribution on all 4 variables,  $P(a, b, c, d)$  where this “joint variable” takes values from the cross-product of the individual value sets of each variable.
- Similarly we could define  $P(a, b|c, d)$ . This is the conditional distribution of  $A$  and  $B$  given that  $C = c$  and  $D = d$ . If  $A$  has  $m$  possible values and  $B$  has  $n$  possible values, then  $P(a, b|c, d)$  is a table  $m \times n$  numbers for fixed values of  $c$  and  $d$ , and

$$\sum_{a,b} P(a, b|c, d) = \sum_i \sum_j P(A = a_i, B = b_j|c, d) = 1.$$

Probability models with multiple variables can be referred to as **joint probability models** or **multivariate probability models** (as opposed to **univariate probability models** for a single variable). As mentioned earlier, one issue with multivariate models is that the size of the table required to specify the joint distribution (for discrete random variables) grows exponentially. For example, if we have  $K$  random variables each taking  $m$  values, then the  $K$ -dimensional table for the joint distribution will contain  $m^K$  entries (each of which is a joint probability for some particular instantiation of the values of  $K$  variables).

## 4 Computing with Probabilities

We are often interested in computing the conditional probability of some proposition of interest (e.g.,  $A = a$ ), given a joint probability table (e.g.,  $P(a, b, c, d)$ ) and given some observed evidence (e.g.,  $B = b$ ).

### 4.1 The Law of Total Probability: Computing Marginals from Joint Probabilities

We can relate the probability distribution of any random variable  $A$  to that of any other random variable  $B$  as follows:

$$P(a) = \sum_b P(a, b) = \sum_b P(a|b)P(b)$$

This is known as the **law of total probability** and it tells us how to compute a marginal probability  $P(a)$  from a joint distribution or from a conditional and another marginal<sup>2</sup>. The problem here is that we would like to know the distribution of  $A$ , but don't have it directly: instead we have the joint distribution  $P(a, b)$ . So, in order to be able to make use of this joint distribution, we introduce the variable  $B$  into the problem and express  $P(a)$  in terms of what we know about  $A$  and  $B$  together. This “summing out of the other variable,” also known as **marginalization**, is a very useful manipulation.

We can of course generalize this idea. For example, to get  $P(a, b)$  we have

$$P(a, b) = \sum_{c,d} P(a, b, c, d) = \sum_{c,d} P(a, b|c, d)P(c, d).$$

(Make sure you can convince yourself that the equation above makes sense: it may help to think of  $AB$  as a “super variable” that is conditioned on some value of  $CD$ .)

---

<sup>2</sup>Note below that the term “marginal” refers to a univariate distribution, e.g.,  $P(a)$  or  $P(b)$ .

## 4.2 Bayes' Rule

We have seen already that we can express  $P(a, b)$  as  $P(a|b)P(b)$ . Since there is nothing special about which argument  $a$  or  $b$  comes first in  $P(a, b)$  then it is clear that we can also write  $P(a, b) = P(b, a) = P(b|a)P(a)$ . In fact if we equate these two different expressions for  $P(a, b)$  we get

$$P(a|b)P(b) = P(b|a)P(a)$$

and if we divide each side by  $P(b)$  we can derive **Bayes' rule**, i.e.,

$$P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

Furthermore, from the law of total probability we know that we can re-express the denominator on the right hand side to give

$$P(a|b) = \frac{P(b|a)P(a)}{\sum_j P(b|a_j)P(a_j)}.$$

Bayes' rule expresses how we can “reason in reverse”, i.e., given a forward model connecting  $b$  to  $a$ , namely  $P(b|a)$ , and given a marginal distribution on  $a$  (i.e.,  $P(a)$ ), we can make inferences about  $a$  given  $b$ .

**EXAMPLE 6:** Consider a medical diagnosis problem where a random variable  $A$  can take two values, 0 meaning a patient *does not have the disease* and 1 meaning a patient *has the disease*. The random variable  $T$  represents the outcome of a test for the disease, where  $T$  can take values  $t = 0$  (negative) and  $t = 1$  (positive). Assume we know (based on past medical data or prior experience) the following:

$$P(T = 1|A = 0) = 0.01, \quad P(T = 1|A = 1) = 0.9,$$

from which we can deduce that  $P(T = 0|A = 0) = 0.99$  and  $P(T = 0|A = 1) = 0.1$ . In words this tells us that for healthy people ( $A = 0$ ) the test is negative 99% of the time—or equivalently there is a 1% chance of a false positive. And for people with the disease ( $A = 1$ ), the test is negative only 10% of the time.

To use Bayes rule to calculate  $P(A = a|T = t)$  we will also need to know what  $P(A)$  is. Suppose we know that  $P(A = 1) = 0.001$  (i.e., only 1 in a thousand people on average have the disease). Given this information we can now use Bayes rule to compute the conditional probability that a person has the disease given (a) that the test outcome is negative (has value 0), and (b) given that the test outcome is positive (has value 1). (Calculation of the actual conditional probabilities for this problem is left as a homework exercise).

### 4.3 Factorization and the Chain Rule

We have seen earlier that we can write  $P(a, b)$  as  $P(a|b)P(b)$ . In fact one can do a similar decomposition of a joint distribution defined on  $K$  variables into a product of conditional probabilities and a marginal, e.g.,

$$\begin{aligned}
 P(a, b, c, d) &= P(a|b, c, d)P(b, c, d) \\
 &\quad \text{treating } b, c, d \text{ as a conjunctive value from } B \times C \times D \\
 &= P(a|b, c, d)P(b|c, d)P(c, d) \\
 &\quad \text{repeating the same trick with } P(b, c, d) \\
 &= P(a|b, c, d)P(b|c, d)P(c|d)P(d). \\
 &\quad \text{now factorizing } P(c, d)
 \end{aligned}$$

Note that this works for any ordering, i.e., there is nothing special about the ordering  $a, b, c, d$ —we could have just as easily have decomposed as  $P(a, b, c, d) = P(d|a, b, c)P(b|a, c)P(c|a)P(a)$  or using any other ordering. Also note that this works for any number of variables. There is no assumption being made here, this factoring always holds (as long as we are working with random variables). These types of factored representations can often be useful when we are working with a joint distribution and we wish to break it down into simpler factors that may be easier to manipulate and interpret.

## 5 Real-valued Variables

All of the properties and equations above extend naturally to real-valued variables and density functions, i.e., we just replace our distributions  $P(\cdot)$  with density functions  $p(\cdot)$ . For example,

- **Conditional density functions:** Given two real-valued random variables  $X$  and  $Y$ , we can define the conditional density  $p(x|y)$ , interpreted as the density of  $X$  conditioned on  $Y = y$  being true. As with a conditional probability distribution, a conditional density function is itself a density function and obeys the laws of density functions, i.e.,  $p(x|y) \geq 0$  and  $\int p(x|y)dx = 1$ . We can also define conditional density functions of the form  $p(x|a)$  where the conditioning is now on a discrete random variable taking a particular value,  $A = a$ .
- **Joint density functions:** Naturally, we can define joint probability density functions over multiple real-valued variables. For example  $p(x, y)$  is a joint density function for two real-valued random

variables  $X$  and  $Y$ . Again, we can think of this as a “super-variable” taking values (a two-component vector) in the cross-product of the values of  $X$  and the values of  $Y$ . If  $X$  and  $Y$  can each take values anywhere on the real line, then the cross-product of the two defines a two-dimensional plane. The density function  $p(x, y)$  can then be thought of as a scalar-valued function (or surface, pointing into the 3rd dimension) that is defined over the 2d-plane, where  $p(x, y) \geq 0$  and  $\int_x \int_y p(x, y) dx dy = 1$ . We could plot  $p(x, y)$  as a contour plot or heatmap in two dimensions. We can extend this to any number of variables: with  $K$  real-valued variables we have a joint density defined as a scalar function of  $K$  real-valued variables (and, of course, although we can define such density functions for higher-dimensional sets of variables, we won’t be able to plot or visualize them).

- Other properties such as the law of total probability, Bayes rule, and factorization, are extended to real-valued variables in the manner one might expect, e.g., extending the law of total probability to probability densities we have  $p(x) = \int p(x, y) dy$  with integration replacing summation.

## 6 Mixing Discrete and Continuous Random Variables

Up to this point we have been discussing sets of variables that either all discrete-valued or real-valued. In practice of course we will frequently encounter mixed sets of random variables where some are discrete-valued and some are real-valued. The general theory and principles above can be extended in a natural way, as long we are careful to interpret which parts of our expressions are distributions and which parts are densities.

To illustrate the ideas consider an example where  $A$  is a discrete random variable and  $X$  is a real-valued random variable. For example,  $A$  might represent the binary events  $a_1 = \textit{has the flu}$  and  $a_2 = \textit{does not have the flu}$ , and  $X$  could be a patient’s temperature. We can define various conditional and joint distributions and densities:

- $P(a_1|x)$  is a conditional probability between 0 and 1 that is a function of  $x$ , i.e., it is a function defined over the real-line (if  $x$  is defined on the real-line) where the function takes values between 0 and 1. And  $P(a_2|x) = 1 - P(a_1|x)$  by definition. For example, in a very simple model, for high-values of  $x$  (high temperature) we might have that  $P(a_1|x)$  (the probability of flu) is quite high (close to 1), and the probability of flu could decrease monotonically as the temperature  $x$  decreases.
- $p(x|a_1)$  and  $p(x|a_2)$  are two **conditional density functions**. For example,  $p(x|a_1)$  might be a Gaussian model with a mean temperature of 102 degrees, conditioned on the patient having the flu ( $A =$

$a_1$ ), and  $p(x|a_2)$  could be another Gaussian model with a mean temperature of 98 degrees, conditioned on the patient not having the flu ( $A = a_2$ ).

- We can define the marginal distribution of  $p(x)$  as  $p(x) = p(x|a_1)P(a_1) + p(x|a_2)P(a_2)$ : this is known as a **finite mixture model** since the marginal distribution on temperature  $x$  can be expressed as a weighted combination of two conditional distributions (corresponding to the two subpopulations, people with the flu and people without the flu). We can also view this definition of  $p(x)$  as just another application of the law of total probability, but where one variable is discrete and the other is real-valued.
- Bayes rule also works in this context, e.g.,

$$\begin{aligned} P(a_1|x) &= \frac{p(x|a_1)P(a_1)}{p(x|a_1)P(a_1) + p(x|a_2)P(a_2)} \\ &= \frac{p(x|a_1)P(a_1)}{p(x)} \end{aligned}$$

EXAMPLE 7: In the example above with  $A$  and  $x$ , say we assume that  $p(x|a_1) = N(\mu_1, \sigma^2)$  and  $p(x|a_2) = N(\mu_2, \sigma^2)$ , i.e., that the conditional densities on patients' temperatures are both Gaussian, with different means and a common variance. Assume that the value of  $p(a_1)$  is also known. With a little algebraic manipulation and the use of Bayes rule, one can show (homework problem) that

$$P(a_1|x) = \frac{1}{1 + e^{-(\alpha_0 + \alpha x)}}$$

where  $\alpha_0$  and  $\alpha$  are scalars that are functions of  $p(a_1)$  and of the parameters of the two Gaussian models. This functional form for  $P(a_1|x)$  is known as the **logistic model** and we will encounter it again in the future when we discuss classification and regression.

## 7 Expectation

The **expected value** of a discrete-valued random variable is defined as

$$E[a] = \sum_i P(a_i)a_i.$$

This only makes sense of course if the values  $a_i$  are numerical, e.g., the integers 1, 2, 3, . . . For example, if  $A$  represents the variable *job type* with possible values such as *engineer*, *teacher*, *cook*, etc., then taking

the expected value of these categories is not meaningful. These types of discrete random variables, with values that cannot meaningfully be ordered or mapped to numerical values, are referred to as **categorical variables**.

We can also define the expected value of a conditional random variable, e.g.,

$$E_{P(a|b)}[a] = \sum_i P(a_i|b)a_i.$$

Note that the averaging here occurs with respect to  $P(a_i|b)$  rather than  $P(a_i)$ . Strictly speaking, wherever the term  $E[a]$  is used it should always have a subscript indicating what distribution the expectation is being taken with respect to, e.g.,  $E_{P(a)}[a]$  or  $E_{P(a|b)}[a]$ . In practice when this subscript is not indicated explicitly one should be careful to make sure that one knows what distribution the expectation is being taken with respect to (it is common in many papers and textbooks to not explicitly state what the expectation is with respect to and to assume the reader knows from the context what distribution is being used).

We can similarly define the expected value for real-valued random variables, e.g.,

$$E[x] = \int p(x)x dx \quad \text{and} \quad E_{p(x|y)}[x] = \int p(x|y)x dx.$$

The expected value in an intuitive sense represents the “center of mass” of the density function, i.e., the point on the real line where it would be balanced if it consisted of actual mass.

A useful property to know about expectations is that they have the property of linearity, i.e.,

$$E[ax + b] = aE[x] + b$$

where  $a$  and  $b$  are arbitrary constants and where the random variable  $X$  can be either real-valued or discrete. (Proof left as an exercise).

Finally, we can also take **expectations of functions of random variables**,  $g(x)$ , i.e.,  $E[g(x)] = \int_x g(x)p(x)dx$ . Note that in general  $E[g(x)] \neq g(E[X])$ , i.e., the expectation of a function is not the same as the function of the expectation. A well known example of  $g(x) = (x - E[X])^2$  is the squared difference between  $X$  and the expected value of  $X$ , i.e.,

$$E[g(x)] = E[(x - E[X])^2] = \int_x p(x)(x - E[x])^2 dx$$

which is the well-known expression for the variance of a density  $p(x)$ , measuring how “spread out” a density function is (and is also known as a “scale parameter” in statistics).

## 8 Summary of Key Points

- The ideas above can be summarized in terms of a few relatively general principles: definition and properties of a random variable, distributions and densities, conditional and joint probabilities, Bayes rule, the law of total probability, factorization, expectation. These basic ideas are extremely powerful and well worth knowing: we will use them extensively in probabilistic learning.
- Note that all of the above principles hold *in general* for all distributions and density functions, irrespective of the functional form or numerical specification of the distributions and densities. Thus, we can reason with and manipulate distributions and densities at an abstract level, independently of the detailed specifications. Note also that we have made no assumptions above in stating the general principles and properties (apart from some specific assumptions in our illustrative examples)—so the properties stated above are very general.
- Two open problems remain that we have not discussed:
  1. Where do the numbers come from? the probabilities in the tables for discrete distributions and the parameters in the density functions? This is where probabilistic learning comes in—we can frequently learn these numbers (or parameters) from observed data.
  2. As we include more variables in our multivariate models, how will we avoid the problems of computation in high-dimensional spaces? e.g., for marginalization, how can we avoid computationally expensive high-dimensional integration for real-valued densities and exponential increase in the number of terms in the sums for discrete distributions? We will see that there ways to impose *independence structure* on our probability models that allow us to simplify models in a systematic way.

## 9 The Semantic Interpretation of Probability

We conclude this section with a brief discussion on the semantic interpretation of probability. Consider an event  $a$  and a probability  $P(a)$ . As we discussed earlier, we can think of  $a$  as a logical statement about the world that is either true or false, and  $P(a)$  expresses our uncertainty about  $a$  being true. But what exactly do we mean when we say for example that  $P(a) = 0.3$ ? Of interest here is our interpretation of the number 0.3.

By probability we mean that the proposition  $a$  is assigned a number between 0 and 1 representing our uncertainty. We can agree that it makes sense to interpret  $P(a) = 0$  as meaning that  $a$  is logically false or

impossible, and  $P(a) = 1$  as stating that  $a$  is logically true.

Numbers between 0 and 1 are a little more open to interpretation. The classical and traditional viewpoint, that most of us learned as undergraduates, is that  $P(a)$  represents the relative frequency with which  $a$  is true, in the infinite limit over a series of experiments or trials. This interpretation works well for problems like tossing a coin or throwing a pair of dice—we can imagine repeating such experiments until the observed frequency can be trusted as being a good estimate of the true probability of occurrence of a particular event (e.g., “heads” on the coin). This is known as the **frequentist interpretation of probability**.

However, there are other propositions  $a$  for which it does not make sense conceptually to imagine an infinite series of experiments. For example, let  $a = \textit{the US soccer team will win the World Cup within the next 20 years}$ . This is a proposition for which we can't easily imagine conducting repeated trials. Similarly imagine propositions such as  $a = \textit{life exists on other planets}$  or  $a = \textit{Alexander the Great played the harp}$ . We could come up with subjective estimates (our best guesses) of probabilities for any of these propositions, even though there is clearly no notion of repeated trials. This leads to the **subjective or Bayesian interpretation of probability**, which can perhaps best be summarized as thinking of  $P(a)$  as the **degree of belief** that an agent (a human or a computer program) attaches to the likelihood that  $a$  is true. Note that there is no notion of repeated trials being required: this is just a number that reflects some agent's degree of belief. More formally, degree of belief can be stated as a conditional probability  $P(a|I)$  where  $I$  is the background information available to the agent (although we usually don't bother including  $I$  in practice). So, for example,  $I$  could be a person's model or data on repeated coin tossing: so in a sense our “definition” includes the frequentist approach as a special case, in the sense that  $I$  represents whatever assumptions or background knowledge we use to construct  $P$ .

While the Bayesian view of probability as a degree of belief might seem somewhat informal and non-precise, one can in fact make the concept of “degree of belief” quite precise. For example, one way to do this is to consider the degree of belief for a proposition as the probability value for which an individual will change their minds about making a bet about whether the proposition is true, by relating degrees of belief to what kind of odds a rational individual would be willing to accept.

We will in general use the Bayesian “degree of belief” interpretation of probability since it covers a broader range of situations than the frequentist approach, and can be thought of as a generalization. The Bayesian interpretation of probability also maps well to artificial intelligence and machine learning: the probabilities computed by our agents (computer programs) correspond to “subjective estimates” of how likely it is that particular events are true, conditioned on the information available to the agent.

One very important point, however, is that whichever interpretation we use, the rules and properties of

probability are the same, i.e., we use the same equations. Only the semantic interpretation of the probabilities changes, not the mathematics of how we combine them and work with them.