

Artificial Intelligence II

Some supplementary notes on probability

Sean B. Holden, 2010-12

1 Introduction

These notes provide a reminder of some simple manipulations that turn up a great deal when dealing with probabilities. The material in this handout—assuming you know it well—should suffice for getting you through most of the AI material on uncertain reasoning. In particular, the boxed results are the really important ones.

Random variables (RVs) are by convention given capital letters. Say we have the RVs X_1, \dots, X_n . Their values are given using lower case. So for example X_1 might be a binary RV taking values `true` and `false`, and X_2 might be the outcome of rolling a die and therefore taking values `one`, `two`, `...`, `six`.

The use of probability in AI essentially reduces to representing in some usable way the joint distribution $P(X_1, \dots, X_n)$ of all the RVs our agent is interested in, because if we can do that then in principle we can compute *any* probability that might be of interest. (This is explained in full below.)

To be clear, the joint distribution is talking about the *conjunction* of the RVs. We'll stick to the convention that a comma-separated list of RVs (or a set of RVs) represents a conjunction. Also, the notation

$$\sum_{x_i \in X_i} (\dots x_i \dots)$$

denotes the sum over all *values* of a random variable. So for example if X_1 is binary then

$$\sum_{x_1 \in X_1} P(x_1, X_2) = P(\text{true}, X_2) + P(\text{false}, X_2). \quad (1)$$

This extends to summing over *sets* of RVs. Let's define

$$\mathbf{X} = \{X_1, \dots, X_n\}$$

and

$$\mathbf{X}' = \{X'_1, \dots, X'_m\}.$$

Then for any sets \mathbf{X} and $\mathbf{X}' \subseteq \mathbf{X}$ of RVs define $\mathbf{X} \setminus \mathbf{X}'$ to be the set \mathbf{X} with the elements of \mathbf{X}' removed

$$\mathbf{X} \setminus \mathbf{X}' = \{X \in \mathbf{X} \mid X \notin \mathbf{X}'\}.$$

We'll always be assuming that $\mathbf{X}' \subseteq \mathbf{X}$. Finally

$$\sum_{x' \in \mathbf{X}'} (\dots, x'_1, \dots, x'_m, \dots)$$

means

$$\sum_{x'_1 \in X'_1} \sum_{x'_2 \in X'_2} \dots \sum_{x'_m \in X'_m} (\dots, x'_1, \dots, x'_m, \dots).$$

2 Standard trick number 1: marginalising

Marginalising is the process of getting rid of RVs that we don't want to have to think about—although in some cases it's used the other way around to introduce variables. In general, say we want to ignore X_i . Then

$$P(\mathbf{X} \setminus \{X_i\}) = \sum_{x_i \in X_i} P(\mathbf{X}).$$

So for example, equation 1 is actually telling us that with $\mathbf{X} = \{X_1, X_2\}$

$$\begin{aligned} P(X_2) &= P(\mathbf{X} \setminus \{X_1\}) \\ &= \sum_{x_1 \in X_1} P(x_1, X_2) \\ &= P(\text{true}, X_2) + P(\text{false}, X_2). \end{aligned}$$

This can obviously be iterated for as many RVs as we like, so if \mathbf{X}' is the set of random variables we're not interested in then

$$\boxed{P(\mathbf{X} \setminus \mathbf{X}') = \sum_{x' \in \mathbf{X}'} P(\mathbf{X})}.$$

These notes assume for the most part that RVs are discrete. Everything still applies when continuous RVs are involved, but sums are then replaced by integrals. For example, we can marginalise the two-dimensional Gaussian density

$$p(x_1, x_2) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x_1^2 + x_2^2)\right)$$

as follows

$$p(x_1) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2}(x_1^2 + x_2^2)\right) dx_2.$$

3 Standard trick number 2: you can treat a conjunction of RVs as an RV

When we consider events such as $X_1 = \text{true}$ and $X_2 = \text{four}$, the *conjunction* of the events is also an event. This goes for any number of events, and any number of RVs as well. Why is that interesting? Well, Bayes' theorem usually looks like this

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}.$$

However as a conjunction of RVs can be treated as a RV we can also write things like

$$P(X_1, X_5 | X_2, X_3, X_{10}) = \frac{P(X_2, X_3, X_{10} | X_1, X_5) P(X_1, X_5)}{P(X_2, X_3, X_{10})}$$

and Bayes' theorem still works.

4 Standard trick number 3: conditional distributions are still distributions

This is perhaps the point I want to make that's most often missed: *a conditional probability distribution is still a probability distribution*. Consequently the first two tricks extend to them without any extra work—you simply apply them while leaving the conditioning RVs (the ones on the right hand side of the $|$ in $P(\dots | \dots)$) alone. So, for instance, we can write

$$P(X_1|X_3) = \sum_{x_2 \in X_2} P(X_1, X_2|X_3)$$

or in general for sets of RVs

$$P(\mathbf{X}|\mathbf{Z}) = \sum_{\mathbf{y} \in \mathbf{Y}} P(\mathbf{X}, \mathbf{Y}|\mathbf{Z}).$$

Quite often this trick is used to *introduce* extra RVs in \mathbf{Y} rather than eliminate them. The reason for this is that you can then try to re-arrange the contents of the sum to get something useful. In particular you can often use the following further tricks.

Just as marginalisation still works for conditional distributions, so do Bayes' theorem and related ideas. For example, the definition of a conditional distribution looks like this

$$P(X|Y) = \frac{P(X, Y)}{P(Y)} \tag{2}$$

so

$$P(X, Y) = P(X|Y)P(Y).$$

As the left hand side of this equation is a joint probability distribution, and conjunctions of RVs act like RVs, we can extend this to arbitrary numbers of RVs to get, for example

$$\begin{aligned} P(X_1, X_2, X_3) &= P(X_1|X_2, X_3)P(X_2, X_3) \\ &= P(X_1|X_2, X_3)P(X_2|X_3)P(X_3). \end{aligned}$$

What's more useful however is to note that Bayes' theorem is obtained from equation 2 and its twin

$$P(Y|X) = \frac{P(X, Y)}{P(X)}$$

by a simple re-arrangement. How might this work if we have conjunctions of random variables? Consider

$$P(X|Y, Z) = \frac{P(X, Y, Z)}{P(Y, Z)}$$

and its twin

$$P(Y|X, Z) = \frac{P(X, Y, Z)}{P(X, Z)}$$

both of which follow from the definition of conditional probability. Re-arranging to eliminate the $P(X, Y, Z)$ gives

$$P(X|Y, Z) = \frac{P(Y|X, Z)P(X, Z)}{P(Y, Z)}.$$

We now have two smaller joint distributions $P(Y, Z)$ and $P(X, Z)$ which we can split to give

$$\begin{aligned} P(X|Y, Z) &= \frac{P(Y|X, Z)P(X|Z)P(Z)}{P(Y|Z)P(Z)} \\ &= \frac{P(Y|X, Z)P(X|Z)}{P(Y|Z)} \end{aligned}$$

or in general, with sets of RVs

$$\boxed{P(\mathbf{X}|\mathbf{Y}, \mathbf{Z}) = \frac{P(\mathbf{Y}|\mathbf{X}, \mathbf{Z})P(\mathbf{X}|\mathbf{Z})}{P(\mathbf{Y}|\mathbf{Z})}}. \quad (3)$$

5 How to (in principle) compute absolutely anything

Say you want to compute a conditional probability $P(\mathbf{X}|\mathbf{Z})$. By definition

$$P(\mathbf{X}|\mathbf{Z}) = \frac{P(\mathbf{X}, \mathbf{Z})}{P(\mathbf{Z})}$$

and if the complete collection of all the RVs our agent is interested in is $\{\mathbf{X}, \mathbf{Y}, \mathbf{Z}\}$ then both the numerator and the denominator can be computed by marginalising the joint distribution $P(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$. In fact as the denominator serves essentially just to make the left hand side sum to 1 (when we sum over \mathbf{X}) so that it's a proper probability distribution, we often treat it just as a constant and write

$$\boxed{P(\mathbf{X}|\mathbf{Z}) = \frac{1}{Z} \sum_{\mathbf{y} \in \mathbf{Y}} P(\mathbf{X}, \mathbf{Y}, \mathbf{Z})}.$$

The quantity Z is called the *partition function* if you're a physicist or *evidence* if you're a computer scientist, for reasons that will become clear during the lectures.

6 Why multiplication of factors works

This section is really about an algorithm rather than probabilities. We provide a simple explanation of why the multiplication of *factors* in the manner suggested in the lecture notes works, and why it avoids duplicating the computation of sub-expressions.

6.1 Why it works

Let's drop any reference to probabilities for the moment and just look at general summations involving functions. Say you have three finite sets $X = \{x_1, \dots, x_p\}$, $Y = \{y_1, \dots, y_q\}$ and $Z = \{z_1, \dots, z_r\}$, and you want to compute a summation like

$$f(x, y) = \sum_{z \in Z} g(x, y, z)h(y, z) \quad (4)$$

for values $x \in X$ and $y \in Y$. The sum will look like

$$f(x, y) = g(x, y, z_1)h(y, z_1) + \dots + g(x, y, z_r)h(y, z_r).$$

In other words, the products you need to compute are the ones *for which values of z coincide*. This, in a nutshell, is what the process of combining factors achieves. In this example, we would write the factors in the sum as

x	y	z	$F_x(x, y, z)$
x_1	y_1	z_1	$g(x_1, y_1, z_1)$
x_1	y_1	z_2	$g(x_1, y_1, z_2)$
\vdots	\vdots	\vdots	\vdots
x_p	y_q	z_r	$g(x_p, y_q, z_r)$

and

y	z	$F_y(y, z)$
y_1	z_1	$h(y_1, z_1)$
y_1	z_2	$h(y_1, z_2)$
\vdots	\vdots	\vdots
y_q	z_r	$h(y_q, z_r)$

When the factors are multiplied we match up and multiply the table entries for which variables common to both have matching values. So for example

x	y	z	$F_{x,y}(x, y, z)$
x_1	y_1	z_1	$g(x_1, y_1, z_1)h(y_1, z_1)$
x_1	y_1	z_2	$g(x_1, y_1, z_2)h(y_1, z_2)$
\vdots	\vdots	\vdots	\vdots
x_p	y_q	z_r	$g(x_p, y_q, z_r)h(y_q, z_r)$

In order to deal with the summation to form the factor $F_{x,y}(x, y)$ we now form sums of the entries in $F_{x,y}(x, y, z)$ over all values for z , so

x	y	$F_{x,y}(x, y)$
x_1	y_1	$\sum_{z \in Z} F_{x,y}(x_1, y_1, z)$
x_1	y_2	$\sum_{z \in Z} F_{x,y}(x_1, y_2, z)$
\vdots	\vdots	\vdots
x_p	y_q	$\sum_{z \in Z} F_{x,y}(x_p, y_q, z)$

Expanding out one of these summations results in something like

$$F_{x,y}(x_1, y_2) = \sum_{z \in Z} g(x_1, y_2, z)h(y_2, z). \quad (5)$$

Comparing equations 4 and 5 we see that the entries in the factor $F_{x,y}(x, y)$ are just the values of the summation for each possible pair of values x and y .

6.2 The relationship to probabilities

So: the tabular process using factors is just a way of keeping track of the values needed to compute the sum. In the probabilistic inference algorithm the functions f , g , h and so on are all just (conditional) probability distributions, and because we're dealing with the decomposition

$$\Pr(X_1, \dots, X_n) = \prod_{i=1}^n \Pr(X_i | \text{parents}(X_i))$$

on a directed, acyclic graph we start off with a factor for each RV, and each time we get to a summation we sum out the corresponding variable. (The above example does not have this structure, which is why the summing out notation $F_{x,y,\bar{z}}$ does not appear, but the process is identical.)

6.3 Why it avoids duplication

Reverting now to probabilities, say we have a Bayes network that represents the decomposition

$$\Pr(X, Y_1, Y_2, E_1, E_2) = \Pr(E_1|Y_1, Y_2) \Pr(E_2|Y_2) \Pr(Y_2|X) \Pr(X) \Pr(Y_1).$$

(Exercise: draw it.) Now we attempt to compute the inference

$$\Pr(X|e_1, e_2) = \frac{1}{Z} \Pr(X) \sum_{y_1 \in Y_1} \Pr(Y_1) \sum_{y_2 \in Y_2} \Pr(Y_2|X) \Pr(e_1|Y_1, Y_2) \Pr(e_2|Y_2).$$

The repetition of computations arises in summations like this because—handled in the naive way using recursive depth-first evaluation—the summation

$$\sum_{y_2 \in Y_2} \Pr(Y_2|X) \Pr(e_1|Y_1, Y_2) \Pr(e_2|Y_2)$$

will involve computing the product $\Pr(Y_2|X) \Pr(e_2|Y_2)$ for each value of Y_1 . This problem will repeat itself for each value of X . By computing and storing each of the products needed only once the method based on factors avoids this.

7 Further tricks

We now look at some further simple manipulations that are needed to understand the application of Bayes' theorem to supervised learning. Once again, random variables are assumed to be discrete, but all the following results still hold for continuous random variables, with sums replaced by integrals where necessary.

7.1 Some (slightly) unconventional notation

In the machine learning literature there is a common notation intended to make it easy to keep track of which random variables and which distributions are relevant in an expression. While this notation is common within the field, it's rarely if ever seen elsewhere; it is however very useful.

A statistician would define the *expected value* of the random variable X as

$$\mathbb{E}[X] = \sum_{x \in X} xP(x)$$

or when we're interested in the expected value of a function of a random variable

$$\mathbb{E}[f(X)] = \sum_{x \in X} f(x)P(x)$$

where f is some function defined on X . Here, it is implicit that the probability distribution for X is P . With complex expressions involving combinations of functions defined on random variables with multiple underlying distributions it can be more tricky to keep track of which distributions are relevant. Thus the notation

$$\mathbb{E}_{x \sim P(X)} [f(X)]$$

is intended to indicate explicitly that the distribution of X is P , in situations where we don't write out the full definition

$$\mathbb{E}_{x \sim P(X)} [f(X)] = \sum_{x \in X} f(x)P(x).$$

to make it clear. The same notation is also often applied to statements about probabilities rather than expected values.

7.2 Expected value and conditional expected value

The standard definition of the expected value of a function f of a random variable X is

$$\mathbb{E}_{x \sim P(X)} [f(X)] = \sum_{x \in X} f(x)P(x)$$

as already noted. We can also define the *conditional expected value* of $f(X)$ given Y as

$$\mathbb{E}_{x \sim P(X|Y)} [f(X)|Y] = \sum_{x \in X} f(x)P(x|Y).$$

Now here's an important point: *the value of this expression depends on the value of Y* . Thus, the conditional expected value is itself a function of the random variable Y . What is its expected value? Well

$$\begin{aligned} \mathbb{E}_{y \sim P(Y)} [\mathbb{E}_{x \sim P(X|Y)} [f(X)|Y]] &= \sum_{y \in Y} \mathbb{E}_{x \sim P(X|Y)} [f(X)|Y] P(y) \\ &= \sum_{y \in Y} \sum_{x \in X} f(x)P(x|y)P(y) \\ &= \sum_{y \in Y} \sum_{x \in X} f(x)P(x, y) \\ &= \sum_{x \in X} f(x) \sum_{y \in Y} P(x, y) \\ &= \sum_{x \in X} f(x)P(x) \\ &= \mathbb{E}_{x \sim P(X)} [f(X)] \end{aligned}$$

or in the more usual notation

$$\mathbb{E} [\mathbb{E} [f(X)|Y]] = \mathbb{E} [f(X)].$$

7.3 Expected value of the indicator function

For any $b \in \{\text{true}, \text{false}\}$ the *indicator function* \mathbb{I} is defined as

$$\mathbb{I}(b) = \begin{cases} 1 & \text{if } b = \text{true} \\ 0 & \text{if } b = \text{false} \end{cases}.$$

Let f be a Boolean-valued function on a random variable X . Then

$$\begin{aligned}\mathbb{E}_{x \sim P(X)} [\mathbb{I}(f(x))] &= \sum_{x \in X} \mathbb{I}(f(x))P(x) \\ &= \sum_{x \in X, f(x) \text{ is true}} \mathbb{I}(f(x))P(x) + \sum_{x \in X, f(x) \text{ is false}} \mathbb{I}(f(x))P(x) \\ &= \sum_{x \in X, f(x) \text{ is true}} P(x) \\ &= P_{x \sim P(x)} [f(x) = \text{true}]\end{aligned}$$

In other words, *the probability of an event is equal to the expected value of its indicator function*. This provides a standard method for calculating probabilities by evaluating expected values. So for example if we roll a fair die and consider $f(X)$ to be true if and only if the outcome is even then

$$P(\text{outcome is even}) = \mathbb{E} [\mathbb{I}(f(X))] = 1/6 + 1/6 + 1/6 = 1/2$$

as expected.