

Artificial Intelligence II

Some supplementary notes on probability

Sean B. Holden

February 2010

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

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

and 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

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

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

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

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

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

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.