

# A Primer on Probability, Expectation, and Variance

Phil Schniter

January 26, 2022

## 1 Random variables

A *random variable* (RV) is a quantity that can take on many possible values. We never know what value the RV will take on, but we do know which values are possible and how often they occur.

A good analogy is a giant bag of marbles, where each marble has a number on it. We never know the particular number will appear when we draw a marble from the bag. These drawn values are called *realizations*, and the RV is the entire bag. The information contained by the RV includes the complete set of values written on the marbles, as well as the fraction of the bag that is taken by each of those values.

In this document, we will use capital letters for RVs (e.g., “ $A$ ”) and lowercase letters for deterministic variables (e.g., “ $a$ ”) in an attempt to be clear. But we can’t do this throughout all of our lecture derivations, because that would conflict with many other notational conventions that we (and the machine learning community) have adopted.

When we write things out mathematically, we have to make a distinction between discrete RVs (i.e., those that take on values from a countable set) and continuous RVs (i.e., those that take on values from an uncountably infinite set).

### 1.1 Discrete random variables

A discrete RV  $A$  is often described using the *probability mass function* (pmf)  $\{p_A[k]\}_{k=1}^K$ , where integer  $k$  indexes over the  $K$  possible values in the set. For example, suppose that  $A \in \{a^{(1)}, a^{(2)}, \dots, a^{(K)}\}$ . Then  $p_A[k] = \Pr\{A = a^{(k)}\}$ , i.e., the probability that  $A$  takes on the value  $a^{(k)}$ . The pmf must be non-negative and sum to one, i.e.,

$$p_A[k] \geq 0 \text{ for all } k, \text{ and } \sum_{k=1}^K p_A[k] = 1. \quad (1)$$

### 1.2 Continuous random variables

A continuous RV  $A$  is often described using the *probability density function* (pdf)  $p_A(\cdot)$ , which is defined for any real-valued input. The pdf must be non-negative and integrate to one, i.e.,

$$p_A(a) \geq 0 \text{ for all } a \in \mathbb{R}, \text{ and } 1 = \int_{-\infty}^{\infty} p_A(a) da. \quad (2)$$

Since it's a density, the pdf can be used to determine the probability of the event that  $A$  takes on values in a range, say  $[a_1, a_2]$ :

$$\Pr\{A \in [a_1, a_2]\} = \int_{a_1}^{a_2} p_A(a) da. \quad (3)$$

But notice that, when  $p_A(a)$  is finite-valued for all  $a$ , there is zero probability that  $A$  takes on any particular value, like  $a_1$ . This can be seen by taking the limit  $a_2 \rightarrow a_1$  in (3): the probability  $\Pr\{A = a_1\} = \lim_{a_2 \rightarrow a_1} \Pr\{A \in [a_1, a_2]\}$  vanishes.

Note: it is possible to describe discrete a RV using a pdf as follows:

$$p_A(a) = \sum_{k=1}^K p_A[k] \delta(a - a^{(k)}), \quad (4)$$

where  $\delta(\cdot)$  is the “point mass” or Dirac delta function. In words,  $\delta(\cdot)$  is an infinitely tall, infinitely narrow, unit-area spike centered at the origin (a strange beast!). Mathematically,

$$\delta(0) = \infty, \quad \delta(a) = 0 \text{ for any } a \neq 0, \text{ and } \int_{-\infty}^{\infty} \delta(a) da = 1. \quad (5)$$

### 1.3 Joint random variables

Often we encounter expressions with multiple RVs that may be statistically related. We can describe this collection of RVs using a *joint pmf* (if they are all discrete) or a *joint pdf* more generally. In the case of two RVs  $A$  and  $B$ , we would have the joint pdf  $p_{A,B}(\cdot, \cdot)$ . Not surprisingly,  $p_{A,B}$  must be non-negative and integrate to one, i.e.,

$$p_{A,B}(a, b) \geq 0 \text{ for all } a, b \in \mathbb{R}, \text{ and } 1 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p_{A,B}(a, b) da db. \quad (6)$$

Also,  $p_{A,B}(a, b) = p_{B,A}(b, a)$  for all  $a, b \in \mathbb{R}$ .

We say that  $A$  and  $B$  are *statistically independent* when

$$p_{A,B}(a, b) = p_A(a)p_B(b) \text{ for all } a, b \in \mathbb{R}. \quad (7)$$

As we will see below, that means you can't predict anything about  $A$  from  $B$  (or vice versa).

### 1.4 Conditional probability

One way to quantify the dependence among RVs is through *conditional probability*. Given two RVs  $A$  and  $B$ , the *conditional pdf of  $A$  given  $B$*  is written as  $p_{A|B}(\cdot|\cdot)$ . The quantity  $p_{A|B}(a|b)$  describes the pdf of  $A$  evaluated at  $a$  given the event that  $B = b$ .

*Example:* Suppose that you model temperature as a continuous random variable  $A$  and whether or not it is snowing as a binary random variable  $B$  (i.e.,  $B = 1$  means it is snowing). Then  $p_{A|B}(a|1)$  would be the pdf of temperatures given that it is snowing. And  $p_{B|A}(1|a)$  would be probability that it is snowing given that the temperature equals  $a$ .

The conditional and joint distributions can be connected as follows:

$$p_{A|B}(a|b) p_B(b) = p_{A,B}(a, b). \quad (8)$$

Also, the order of conditioning can be “swapped” using *Bayes rule*:

$$p_{A|B}(a|b) = \frac{p_{B|A}(b|a)p_A(a)}{p_B(b)}, \quad (9)$$

which follows from (8) and the equality  $p_{A,B}(a, b) = p_{B,A}(b, a)$ .

In the special case that  $A$  and  $B$  are statistically independent, (7) and (8) imply

$$p_{A|B}(a|b)p_B(b) = p_A(a)p_B(b) \quad (10)$$

$$\Leftrightarrow p_{A|B}(a|b) = p_A(a) \quad (11)$$

which says that the density of  $A$  is invariant to the value of  $B$ . Thus, you can’t predict anything about  $A$  from  $B$ . By exchanging  $A$  and  $B$  in (8), you can similarly prove that  $p_{B|A}(b|a) = p_B(b)$ , i.e., you can’t predict anything about  $B$  from  $A$ .

## 2 Expectation

The *expected value*  $E\{A\}$  of a RV  $A$  is the average or “mean” value of  $A$ . Mathematically, it has the following definition:

$$\text{discrete : } E\{A\} = \sum_{k=1}^K a^{(k)} p_A[k] \quad (12)$$

$$\text{continuous : } E\{A\} = \int_{-\infty}^{\infty} a p_A(a) da \quad (13)$$

It’s important to understand the distinction between this statistical average of  $A$  and the *sample average* of a set of realizations  $\{a_1, \dots, a_n\}$  of  $A$ , which we write as

$$\bar{a} \triangleq \frac{1}{n} \sum_{i=1}^n a_i. \quad (14)$$

For example, when the statistical average is zero (i.e.,  $E\{A\} = 0$ ), the sample average may still be non-zero. Concretely, if  $A$  is a zero-mean Gaussian random variable, then  $E\{A\} = 0$  but  $\bar{a} \neq 0$  for any finite value of  $n$ . (You can easily check this in Python or Matlab.)

There are two key properties of expectation that we will use: For any deterministic constants  $c$  and  $d$ , deterministic functions  $f(\cdot)$  and  $g(\cdot)$ , and independent RVs  $A$  and  $B$ , we have

$$1. \text{ Linearity: } E\{c + d f(A)\} = c + d E\{f(A)\}$$

$$2. \text{ Independence: } E\{f(A)g(B)\} = E\{f(A)\} E\{g(B)\}$$

## 3 Variance

The *variance*  $\text{var}\{A\}$  of a RV  $AA$  is the mean squared deviation from the mean of  $A$ . Mathematically, it has the following definition:

$$\text{discrete : } \text{var}\{A\} = E\{(A - E\{A\})^2\} = \sum_{k=1}^K (a^{(k)} - E\{A\})^2 p_A[k] \quad (15)$$

$$\text{continuous : } \text{var}\{A\} = E\{(A - E\{A\})^2\} = \int_{-\infty}^{\infty} (a - E\{A\})^2 p_A(a) da \quad (16)$$

Note that the statistical variance is not equal to the *sample variance* of set of realizations  $\{a_1, \dots, a_n\}$  of  $A$ , which we write (using the sample mean  $\bar{a}$ ) as

$$s_a^2 \triangleq \frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})^2. \quad (17)$$

One very useful identity is

$$\mathbb{E}\{A^2\} = \text{var}\{A\} + (\mathbb{E}\{A\})^2. \quad (18)$$

A short proof is now provided using the shorthand notation  $\mu \triangleq \mathbb{E}\{A\}$ :

$$\mathbb{E}\{A^2\} = \mathbb{E}\{(A - \mu) + \mu\}^2 \quad (19)$$

$$= \mathbb{E}\{(A - \mu)^2 + 2(A - \mu)\mu + \mu^2\} \quad (20)$$

$$= \underbrace{\mathbb{E}\{(A - \mu)^2\}}_{\text{var}\{A\}} + 2\underbrace{(\mathbb{E}\{A\} - \mu)\mu}_{=0} + \underbrace{\mu^2}_{(\mathbb{E}\{A\})^2}. \quad (21)$$

## 4 Conditional expectation and variance

Given several RVs, we sometimes want to describe the statistical average of one variable given fixed values of the rest. We can do this using *conditional expectation*. For example, given two RVs  $A$  and  $B$ , we define the *conditional mean of  $A$  given the event  $B = b$*  as follows:

$$\mathbb{E}\{A | B = b\} = \int_{-\infty}^{\infty} a p_{A|B}(a|b) da. \quad (22)$$

Note that this quantity is a deterministic function of the deterministic constant  $b$ , i.e.,

$$\mathbb{E}\{A | B = b\} = f(b) \quad (23)$$

for some function  $f(\cdot)$ . Sometimes we encounter the closely related quantity  $\mathbb{E}\{A | B\}$ . This is nothing more than the same function  $f(\cdot)$  from (23) applied to the random variable  $B$ , i.e.,

$$\mathbb{E}\{A | B\} = f(B). \quad (24)$$

Thus  $\mathbb{E}\{A | B\}$  is itself random, i.e., a random variable.

Similarly, we may want to know how far  $A$  deviates from its mean, on average, when  $B = b$ . We thus define the *conditional variance of  $A$  given the event  $B = b$*  as follows:

$$\text{var}\{A | B = b\} = \mathbb{E}\{(A - \mathbb{E}\{A | B\})^2 | B = b\} \quad (25)$$

$$= \int_{-\infty}^{\infty} (a - \mathbb{E}\{A | B = b\})^2 p_{A|B}(a|b) da. \quad (26)$$

Note that  $\text{var}\{A | B = b\}$  is a deterministic function of  $b$ . We could also consider  $\text{var}\{A | B\}$ , which is a function of the random variable  $B$ .

Next we consider some useful applications of conditional mean and variance.

## 4.1 Law of total expectation

Consider evaluating the quantity  $E\{f(A, B, C)\}$ , which is the average value of a function  $f(\cdot)$  involving the three RVs  $A$ ,  $B$ , and  $C$ . In computing this quantity, we are “averaging out” the random contributions from  $A$ ,  $B$ , and  $C$ . In some cases, when computing  $E\{f(A, B, C)\}$ , it is more convenient to “average out” *one random variable at a time*. We can do this by *nesting the expectations*. For example, we could first remove the randomness due to  $C$  alone by holding  $A$  and  $B$  fixed via

$$g(A, B) \triangleq E\{f(A, B, C) \mid A, B\}, \quad (27)$$

which gives a function  $g$  of RVs  $A$  and  $B$ . Next we could remove the randomness due to  $B$  by holding  $A$  fixed:

$$h(A) \triangleq E\{g(A, B) \mid A\}. \quad (28)$$

And finally we could remove the randomness due to  $A$ :

$$E\{h(A)\}. \quad (29)$$

The key point is that this latter quantity is the identical to what we would get from averaging  $A, B, C$  all at once, i.e.,

$$E\{f(A, B, C)\} = E\{h(A)\} \quad (30)$$

$$= E\{ E\{g(A, B) \mid A\} \} \quad (31)$$

$$= E\{ E\{ E\{f(A, B, C) \mid A, B\} \mid A\} \}. \quad (32)$$

This is known as the *law of total expectation*. It’s typically stated in the 2-variable case, and without the function  $f(\cdot)$ , as follows:

$$E\{A\} = E\{ E\{A \mid B\} \}. \quad (33)$$

A short proof is now provided:

$$E\{ E\{A \mid B\} \} = \int \left[ \int a p_{A|B}(a|b) da \right] p_B(b) db \quad (34)$$

$$= \int a \int \underbrace{p_{A|B}(a|b)p_B(b)}_{p_{A,B}(a,b)} db da \quad (35)$$

$$= \int a p_A(a) da \quad (36)$$

$$= E\{A\}. \quad (37)$$

## 4.2 MMSE estimation

Another application of condition expectation is *minimum mean-squared error* (MMSE) estimation. Say that our goal is to predict a target variable  $y$  from vector of features  $\mathbf{x} = [x_1, \dots, x_d]^\top$  using a function  $f(\mathbf{x})$ . We’d like that, on average, the squared error  $(y - f(\mathbf{x}))^2$  is minimized. We can formalize our notion of “average” using statistical expectation if we model the the target and features as random variables  $Y$  and  $\mathbf{X} = [X_1, \dots, X_d]^\top$ .

In this statistical formulation, we want that  $f(\cdot)$  minimizes the *mean-squared error*

$$\text{MSE} \triangleq \text{E} \{ (Y - f(\mathbf{X}))^2 \}. \quad (38)$$

We can derive the MMSE  $f(\cdot)$  by first using the law of total expectation to write

$$\text{MSE} = \text{E} \{ \underbrace{\text{E} \{ (Y - f(\mathbf{X}))^2 \mid \mathbf{X} \}}_{\text{MSE}(\mathbf{X})} \}. \quad (39)$$

If we can find the  $f(\mathbf{X})$  that minimizes  $\text{MSE}(\mathbf{X})$  for *any*  $\mathbf{X}$ , then this same function will minimize MSE when we average over  $\mathbf{X}$  in (39). Notice that  $\text{MSE}(\mathbf{X})$  can be written as

$$\text{MSE}(\mathbf{X}) = \text{E} \{ (Y - f(\mathbf{X}))^2 \mid \mathbf{X} \} \quad (40)$$

$$= \text{E} \{ Y^2 - 2Yf(\mathbf{X}) + f(\mathbf{X})^2 \mid \mathbf{X} \} \quad (41)$$

$$= \text{E} \{ Y^2 \mid \mathbf{X} \} - 2\text{E}\{Y \mid \mathbf{X}\}f(\mathbf{X}) + f(\mathbf{X})^2. \quad (42)$$

The optimal  $f(\mathbf{X})$ , which we'll call  $\hat{f}(\mathbf{X})$ , is simply find the value of  $f(\mathbf{X})$  that sets the derivative of  $\text{MSE}(\mathbf{X})$  to zero. In other words,

$$0 = \left. \frac{\partial \text{MSE}(\mathbf{X})}{\partial f(\mathbf{X})} \right|_{f(\mathbf{X})=\hat{f}(\mathbf{X})} = -2\text{E}\{Y \mid \mathbf{X}\} + 2\hat{f}(\mathbf{X}) \quad (43)$$

which implies

$$\hat{f}(\mathbf{X}) = \text{E}\{Y \mid \mathbf{X}\}. \quad (44)$$

So, the MMSE estimator of  $Y$  from  $\mathbf{X}$  is the conditional mean,  $\hat{f}(\mathbf{X}) = \text{E}\{Y \mid \mathbf{X}\}$ .

Suppose that we would like to know the smallest possible value of  $\text{MSE}(\mathbf{X})$ , i.e., the value attained by the MMSE estimator  $\hat{f}(\mathbf{X})$ . We can simply plug  $\hat{f}(\mathbf{X})$  from (44) into (40) and obtain

$$\text{MSE}(\mathbf{X})|_{\min} = \text{E} \{ (Y - \hat{f}(\mathbf{X}))^2 \mid \mathbf{X} \} \quad (45)$$

$$= \text{E} \{ (Y - \text{E}\{Y \mid \mathbf{X}\})^2 \mid \mathbf{X} \}, \quad (46)$$

which we recognize as the *conditional variance* of  $Y$ . So, the minimum value of  $\text{MSE}(\mathbf{X})$  is the conditional variance,  $\text{MSE}(\mathbf{X})|_{\min} = \text{var}\{Y \mid \mathbf{X}\}$ .