



Gatsby Bridging Program

Probability: Discrete Distributions

Cameron Stewart

Gatsby Computational Neuroscience Unit

Table of contents

1. Random Variables, Probability Mass Functions, and Cumulative Distribution Functions
2. Expectation and Variance
3. Common Discrete Distributions
4. Introduction to Stochastic Processes

Random Variables, Probability Mass Functions, and Cumulative Distribution Functions

What is a Random Variable?

Random variables are functions which map the possible outcomes of an experiment to numerical values. How we define these functions is up to us. In general, for random variable X and sample space Ω , we have that $X : \Omega \rightarrow \mathbb{R}$.

What is a Random Variable?

Random variables are functions which map the possible outcomes of an experiment to numerical values. How we define these functions is up to us. In general, for random variable X and sample space Ω , we have that $X : \Omega \rightarrow \mathbb{R}$.

Example 1

Consider a bag containing 3 red balls (R) and 5 green balls (G). In our experiment, we are going to draw 2 balls from the bag, with replacement. The sample space is $\Omega = \{RR, RG, GR, GG\}$.

We are interested in determining the probabilities of drawing various numbers of red balls. To do this, we could start by defining a random variable $X : \Omega \rightarrow \{0, 1, 2\}$ such that

$$X(\omega) = \begin{cases} 0 & \text{if } \omega = GG \\ 1 & \text{if } \omega \in \{RG, GR\} \\ 2 & \text{if } \omega = RR \end{cases} .$$

What is a Random Variable?

Random variables are often utilised without any explicit reference to the sample space. Instead of writing $X(\omega)$, we will typically just write X . Instead of writing $\mathbb{P}(E)$ for some event E , we typically write $\mathbb{P}(X \in A)$ for some set of numerical values A .

What is a Random Variable?

Random variables are often utilised without any explicit reference to the sample space. Instead of writing $X(\omega)$, we will typically just write X . Instead of writing $\mathbb{P}(E)$ for some event E , we typically write $\mathbb{P}(X \in A)$ for some set of numerical values A .

Example 1 Cont.

What is the probability of drawing one red ball? We could express this as $\mathbb{P}(\{RG, GR\})$, but it is common and accepted notation to instead write $\mathbb{P}(X \in \{1\})$, or $\mathbb{P}(0 < X < 2)$, or most preferably in this case $\mathbb{P}(X = 1)$.

$$\begin{aligned}\mathbb{P}(X = 1) &= \frac{3}{8} \frac{5}{8} + \frac{5}{8} \frac{3}{8} \\ &= \frac{15}{32}.\end{aligned}$$

The Chain Rule and Independence

In a previous lecture, we covered the following relationship for events E and F :

$$\begin{aligned}\mathbb{P}(E \cap F) &= \mathbb{P}(E \mid F) \mathbb{P}(F) \\ &= \mathbb{P}(F \mid E) \mathbb{P}(E) .\end{aligned}$$

The Chain Rule and Independence

In a previous lecture, we covered the following relationship for events E and F :

$$\begin{aligned}\mathbb{P}(E \cap F) &= \mathbb{P}(E \mid F) \mathbb{P}(F) \\ &= \mathbb{P}(F \mid E) \mathbb{P}(E) .\end{aligned}$$

This also applies to random variables. Here we are looking at the probability of random variables X and Y taking values in sets A and B respectively:

Chain Rule for Two Random Variables

$$\begin{aligned}\overbrace{\mathbb{P}(X \in A, Y \in B)}^{\text{Joint Prob.}} &= \overbrace{\mathbb{P}(X \in A \mid Y \in B)}^{\text{Conditional Prob.}} \overbrace{\mathbb{P}(Y \in B)}^{\text{Marginal Prob.}} \\ &= \mathbb{P}(Y \in B \mid X \in A) \mathbb{P}(X \in A)\end{aligned}$$

The Chain Rule and Independence

In a previous lecture, we covered the following relationship for events E and F :

$$\begin{aligned}\mathbb{P}(E \cap F) &= \mathbb{P}(E \mid F) \mathbb{P}(F) \\ &= \mathbb{P}(F \mid E) \mathbb{P}(E) .\end{aligned}$$

This also applies to random variables. Here we are looking at the probability of random variables X and Y taking values in sets A and B respectively:

Chain Rule for Two Random Variables

$$\begin{aligned}\overbrace{\mathbb{P}(X \in A, Y \in B)}^{\text{Joint Prob.}} &= \overbrace{\mathbb{P}(X \in A \mid Y \in B)}^{\text{Conditional Prob.}} \overbrace{\mathbb{P}(Y \in B)}^{\text{Marginal Prob.}} \\ &= \mathbb{P}(Y \in B \mid X \in A) \mathbb{P}(X \in A)\end{aligned}$$

A future lecture will cover joint, conditional, and marginal distributions and the chain rule, independence, and marginalisation in more detail.

The Chain Rule and Independence

In this lecture, we will only consider independent random variables. For two random variables to be independent, the realisation of one must have no effect on the distribution of the other. E.g. a coin flip is independent of the outcome of the previous coin flip.

The Chain Rule and Independence

In this lecture, we will only consider independent random variables. For two random variables to be independent, the realisation of one must have no effect on the distribution of the other. E.g. a coin flip is independent of the outcome of the previous coin flip. Mathematically, we write this as

$$\mathbb{P}(X \in A, Y \in B) = \mathbb{P}(X \in A) \mathbb{P}(Y \in B) ,$$

taking note that independence implies

$$\mathbb{P}(X \in A \mid Y \in B) = \mathbb{P}(X \in A) .$$

The Chain Rule and Independence

In this lecture, we will only consider independent random variables. For two random variables to be independent, the realisation of one must have no effect on the distribution of the other. E.g. a coin flip is independent of the outcome of the previous coin flip. Mathematically, we write this as

$$\mathbb{P}(X \in A, Y \in B) = \mathbb{P}(X \in A) \mathbb{P}(Y \in B) ,$$

taking note that independence implies

$$\mathbb{P}(X \in A \mid Y \in B) = \mathbb{P}(X \in A) .$$

For N independent random variables, this generalises to:

Independence of N Random Variables

$$\mathbb{P}(X_1 \in A_1, \dots, X_N \in A_N) = \prod_{n=1}^N \mathbb{P}(X_n \in A_n)$$

The Chain Rule and Independence

Example 1 Cont.

In our red ball example, we could also use separate random variables for the outcomes of each draw from the bag. Let $X_1 \in \{0, 1\}$ be the random variable representing the first draw and $X_2 \in \{0, 1\}$ be the random variable representing the second draw (0 for a green ball and 1 for red ball). Then it is straightforward to see that $X = X_1 + X_2$.

The Chain Rule and Independence

Example 1 Cont.

In our red ball example, we could also use separate random variables for the outcomes of each draw from the bag. Let $X_1 \in \{0, 1\}$ be the random variable representing the first draw and $X_2 \in \{0, 1\}$ be the random variable representing the second draw (0 for a green ball and 1 for red ball). Then it is straightforward to see that $X = X_1 + X_2$.

Are X_1 and X_2 independent?

The Chain Rule and Independence

Example 1 Cont.

In our red ball example, we could also use separate random variables for the outcomes of each draw from the bag. Let $X_1 \in \{0, 1\}$ be the random variable representing the first draw and $X_2 \in \{0, 1\}$ be the random variable representing the second draw (0 for a green ball and 1 for red ball). Then it is straightforward to see that $X = X_1 + X_2$.

Are X_1 and X_2 independent? Yes, they are independent random variables, as the first draw has no effect on the second draw.

The Chain Rule and Independence

Example 1 Cont.

In our red ball example, we could also use separate random variables for the outcomes of each draw from the bag. Let $X_1 \in \{0, 1\}$ be the random variable representing the first draw and $X_2 \in \{0, 1\}$ be the random variable representing the second draw (0 for a green ball and 1 for red ball). Then it is straightforward to see that $X = X_1 + X_2$.

Are X_1 and X_2 independent? Yes, they are independent random variables, as the first draw has no effect on the second draw.

What if we had the same experimental setup, but without replacing the ball after the first draw?

The Chain Rule and Independence

Example 1 Cont.

In our red ball example, we could also use separate random variables for the outcomes of each draw from the bag. Let $X_1 \in \{0, 1\}$ be the random variable representing the first draw and $X_2 \in \{0, 1\}$ be the random variable representing the second draw (0 for a green ball and 1 for red ball). Then it is straightforward to see that $X = X_1 + X_2$.

Are X_1 and X_2 independent? Yes, they are independent random variables, as the first draw has no effect on the second draw.

What if we had the same experimental setup, but without replacing the ball after the first draw? In this case they are not independent, as the colour of the first drawn ball will dictate the probabilities of the second drawn ball.

Discrete and Continuous Random Variables

Discrete random variables can take a countable number of values. For example:

- $X \in \{0, 1\}$
- $X \in \{0, 1, 2, \dots\}$
- X representing the number of buses arriving within an hour.

Discrete and Continuous Random Variables

Discrete random variables can take a countable number of values. For example:

- $X \in \{0, 1\}$
- $X \in \{0, 1, 2, \dots\}$
- X representing the number of buses arriving within an hour.

Continuous random variables can take values in continuous ranges. For example:

- $X \in [0, 1]$
- $X \in \mathbb{R}$
- X representing the waiting time until the next bus.

Probability Mass Functions

Discrete distributions can be defined by their probability mass function (PMF). The PMF of random variable X is often denoted by f_X or p_X , and is defined as:

Probability Mass Functions

$$f_X(x) = \mathbb{P}(X = x)$$

Probability Mass Functions

Discrete distributions can be defined by their probability mass function (PMF). The PMF of random variable X is often denoted by f_X or p_X , and is defined as:

Probability Mass Functions

$$f_X(x) = \mathbb{P}(X = x)$$

f_X is simply a function name. It is fine to use a different name, as long as it is clear how the function is defined. Occasionally, the same name is used for PMFs if it is clear from context how these are defined, but I'd advise against this practice for the sake of clarity. E.g. it is clearer to write $p_X(x)$ and $p_Y(y)$ than $p(x)$ and $p(y)$ if these correspond to 2 different PMFs.

Probability Mass Functions

Probabilities can't be negative and must sum to 1 over the set of all possible values \mathcal{X} , so we have the following constraints:

$$f_X(x) \geq 0 \text{ for all } x$$

and

$$\sum_{x \in \mathcal{X}} f_X(x) = 1.$$

\mathcal{X} is referred to as the support of X . $f_X(x) = 0$ for $x \notin \mathcal{X}$.

Example 1 Cont.

Let's continue with our red ball example. What is the PMF of X ? First, we recognise that $\mathcal{X} = \{0, 1, 2\}$. You can verify for yourself that

$$f_X(x) = \begin{cases} \frac{25}{64} & \text{if } x = 0 \\ \frac{15}{32} & \text{if } x = 1 \\ \frac{9}{64} & \text{if } x = 2 \\ 0 & \text{otherwise} \end{cases},$$

and that this function satisfies the conditions for a PMF on the previous slide.

Often, we use a \sim to mean "is distributed as". It is very common to see the notation

$$X \sim f_x ,$$

which, in the discrete case, means that X has the PMF f_x . Variations of this notation exist, but there should never be any ambiguity about the distribution of a random variable when using a \sim .

Probability Mass Functions

Often, we use a \sim to mean "is distributed as". It is very common to see the notation

$$X \sim f_x,$$

which, in the discrete case, means that X has the PMF f_x . Variations of this notation exist, but there should never be any ambiguity about the distribution of a random variable when using a \sim .

Sometimes you will see

$$X_1, \dots, X_N \stackrel{\text{i.i.d.}}{\sim} f_x,$$

which implies that all N random variables are independent and identically distributed (i.i.d.) with PMF f_x .

Cumulative Distribution Functions

A probability distribution can also be defined in terms of its cumulative distribution function (CDF). The CDF of a random variable X is a monotonically increasing function defined as:

Cumulative Distribution Functions

$$F_X(x) = \mathbb{P}(X \leq x)$$

Cumulative Distribution Functions

A probability distribution can also be defined in terms of its cumulative distribution function (CDF). The CDF of a random variable X is a monotonically increasing function defined as:

Cumulative Distribution Functions

$$F_X(x) = \mathbb{P}(X \leq x)$$

For discrete random variables, we can relate this definition to the PMF as follows:

$$F_X(x) = \sum_{y \leq x} f_X(y) .$$

Cumulative Distribution Functions

As probabilities must sum to 1 over the support, we have the following emergent properties for CDFs:

$$\lim_{x \rightarrow -\infty} F_X(x) = 0$$

$$\lim_{x \rightarrow \infty} F_X(x) = 1$$

Cumulative Distribution Functions

As probabilities must sum to 1 over the support, we have the following emergent properties for CDFs:

$$\lim_{x \rightarrow -\infty} F_X(x) = 0$$

$$\lim_{x \rightarrow \infty} F_X(x) = 1$$

We can also see that the following is true:

$$\mathbb{P}(a < X \leq b) = F_X(b) - F_X(a) .$$

Cumulative Distribution Functions

Example 1 Cont.

Back to the red ball example. What is the CDF of X ? Previously, we found that

$$f_X(x) = \begin{cases} \frac{25}{64} & \text{if } x = 0 \\ \frac{15}{32} & \text{if } x = 1 \\ \frac{9}{64} & \text{if } x = 2 \\ 0 & \text{otherwise} \end{cases}.$$

Hence, the CDF is

$$F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ \frac{25}{64} & \text{if } 0 \leq x < 1 \\ \frac{55}{64} & \text{if } 1 \leq x < 2 \\ 1 & \text{if } x \geq 2 \end{cases}.$$

Example 1 Cont.

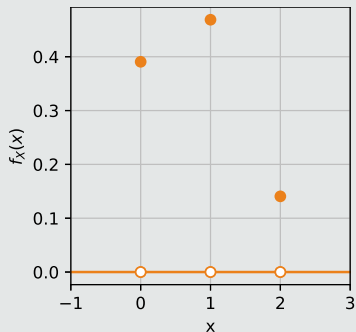


Figure 1: Probability Mass Function

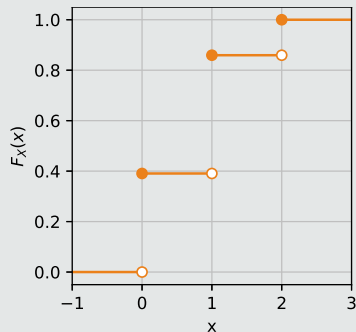


Figure 2: Cumulative Distribution Function

Sampling

Similarly to how we refer to realisations of random variables, we can also talk about samples of distributions. Sampling from a discrete distribution with PMF f_x simply implies that we randomly generate a number x with probability $f_x(x)$.

Sampling

Similarly to how we refer to realisations of random variables, we can also talk about samples of distributions. Sampling from a discrete distribution with PMF f_x simply implies that we randomly generate a number x with probability $f_x(x)$.

To be entirely correct, drawing N samples from f_x involves running N random experiments with N associated random variables $X_1, \dots, X_N \stackrel{\text{i.i.d.}}{\sim} f_x$ and generating N realisations x_1, \dots, x_N .

Sampling

Similarly to how we refer to realisations of random variables, we can also talk about samples of distributions. Sampling from a discrete distribution with PMF f_x simply implies that we randomly generate a number x with probability $f_x(x)$.

To be entirely correct, drawing N samples from f_x involves running N random experiments with N associated random variables $X_1, \dots, X_N \stackrel{\text{i.i.d.}}{\sim} f_x$ and generating N realisations x_1, \dots, x_N .

As $N \rightarrow \infty$, we will observe that

$$\frac{\sum_{n=1}^N [x_n = x]}{N} \rightarrow f_x(x)$$

Intermission

Expectation and Variance

Expectations

Suppose we want to know the average value of a function $g(x)$ evaluated at samples from a distribution. More precisely, we wish to draw realisations of $X_1, \dots, X_N \stackrel{\text{i.i.d.}}{\sim} f_x$ and then compute the mean of $\{g(x_1), \dots, g(x_N)\}$. As N increases, this will converge to a value which we call the expectation (or expected value). This theorem is referred to as the law of large numbers.

We can estimate the expectation for finite N as follows:

Empirical Estimates of Expectations

$$\mathbb{E}_{X \sim f_X} [g(X)] \approx \frac{1}{N} \sum_{n=1}^N g(x_n) \quad \text{for large } N$$

If X is discrete with support \mathcal{X} and PMF f_X , we can also define this expectation as a weighted average:

Expectations on Discrete Distributions

$$\mathbb{E}_{X \sim f_X} [g(X)] = \sum_{x \in \mathcal{X}} f_X(x) g(x)$$

If the distribution on which the expectation is taken is clear from context, the text under the \mathbb{E} is often omitted.

Rules for Expectations

- Linearity:

$$\mathbb{E} [ag(X) + bh(X)] = a\mathbb{E} [g(X)] + b\mathbb{E} [h(X)] ,$$

for constants a, b and functions g, h .

- Non-multiplicativity: In general

$$\mathbb{E} [g(X) h(X)] \neq \mathbb{E} [g(X)] \mathbb{E} [h(X)] .$$

- Multiplicativity under independence: If X_1, \dots, X_N are independent, then

$$\mathbb{E} \left[\prod_{n=1}^N g_n(X_n) \right] = \prod_{n=1}^N \mathbb{E} [g_n(X_n)] ,$$

for functions g_n . The definition of a multivariate expectation will make more sense in later lectures when multivariate distributions are covered.

Example 2

Suppose we roll 2 dice. Let X_1 and X_2 represent the first and second roll respectively.

Intermission

Common Discrete Distributions

Intermission

Introduction to Stochastic Processes
