



Bookmarks



Bookmark

► Introduction

▼ 1. Probability and Inference

Introduction to Probability (Week 1)

Exercises due Sep 21, 2016 at 21:00 UTC

Probability Spaces and Events (Week 1)

Exercises due Sep 21, 2016 at 21:00 UTC

Random Variables (Week 1)

Exercises due Sep 21, 2016 at 21:00 UTC

Jointly Distributed Random Variables (Week 2)

Exercises due Sep 28, 2016 at 21:00 UTC

Conditioning on Events (Week 2)

Exercises due Sep 28, 2016 at 21:00 UTC

Homework 1 (Week 2)

Homework due Sep 28, 2016 at 21:00 UTC

Inference with Bayes' Theorem for Random Variables (Week 3)

Exercises due Oct 05, 2016 at 21:00 UTC

Independence Structure (Week 3)

Exercises due Oct 05, 2016 at 21:00 UTC

Homework 2 (Week 3)

Homework due Oct 05, 2016 at 21:00 UTC

1. Probability and Inference > Inference with Bayes' Theorem for Random Variables (Week 3) > The Product Rule for Random Variables

The Product Rule for Random Variables (Also Called the Chain Rule)

6.008.1x - The Product Rule for Random Variables (Also Called the Ch



▶ 0:00 / 0:00

▶ 1.25x



Video

[Download video file](#)

Transcripts

[Download SubRip \(.srt\) file](#)[Download Text \(.txt\) file](#)

These notes cover roughly the same content as the video:

THE PRODUCT RULE FOR RANDOM VARIABLES (COURSE NOTES)

In many real world problems, we aren't given what the joint distribution of two random variables is although we might be given other information from which we can compute the joint distribution. Often times, we can compute

Notation Summary (Up Through Week 3)

Mini-project 1:

Movie

Recommendations (Week 3)

Mini-projects due Oct 12, 2016 at 21:00 UTC

out the joint distribution using what's called the *product rule* (often also called the chain rule). This is precisely the random variable version of the product rule for events.

As we saw from before, we were able to derive Bayes' theorem for events using the product rule for events: $\mathbb{P}(\mathcal{A} \cap \mathcal{B}) = \mathbb{P}(\mathcal{A})\mathbb{P}(\mathcal{B} \mid \mathcal{A})$. The random variable version of the product rule is derived just like the event version of the product rule, by rearranging the equation for the definition of conditional probability. For two random variables \mathbf{X} and \mathbf{Y} (that take on values in sets \mathcal{X} and \mathcal{Y} respectively), the *product rule* for random variables says that

$$p_{\mathbf{X},\mathbf{Y}}(x, y) = p_{\mathbf{Y}}(y)p_{\mathbf{X}|\mathbf{Y}}(x \mid y) \quad \text{for all } x \in \mathcal{X}, y \in \mathcal{Y} \text{ such that } p_{\mathbf{Y}}(y) > 0.$$

Interpretation: If we have the probability table for \mathbf{Y} , and separately the probability table for \mathbf{X} conditioned on \mathbf{Y} , then we can come up with the joint probability table (i.e., the joint distribution) of \mathbf{X} and \mathbf{Y} .

What happens when $p_{\mathbf{Y}}(y) = 0$? Even though $p_{\mathbf{X}|\mathbf{Y}}(x \mid y)$ isn't defined in this case, one can readily show that $p_{\mathbf{X},\mathbf{Y}}(x, y) = 0$ when $p_{\mathbf{Y}}(y) = 0$.

To see this, think about what is happening computationally: Remember how $p_{\mathbf{Y}}(y)$ is computed from joint probability table $p_{\mathbf{X},\mathbf{Y}}$? In particular, we have $p_{\mathbf{Y}}(y) = \sum_x p_{\mathbf{X},\mathbf{Y}}(x, y)$, so $p_{\mathbf{Y}}(y)$ is the sum of either a row or a column in the joint probability table (whether it's a row or column just depends on how you write out the table and which random variable is along which axis—along rows or columns). So if $p_{\mathbf{Y}}(y) = 0$, it must mean that the individual elements being summed are 0 (since the numbers we're summing up are nonnegative).

We can formalize this intuition with a proof:

Claim: Suppose that random variables \mathbf{X} and \mathbf{Y} have joint probability table $p_{\mathbf{X},\mathbf{Y}}$ and take on values in sets \mathcal{X} and \mathcal{Y} respectively. Suppose that for a specific choice of $y \in \mathcal{Y}$, we have $p_{\mathbf{Y}}(y) = 0$. Then

$$p_{\mathbf{X},\mathbf{Y}}(x, y) = 0 \quad \text{for all } x \in \mathcal{X}.$$

Proof: Let $y \in \mathcal{Y}$ satisfy $p_{\mathbf{Y}}(y) = 0$. Recall that we relate marginal distribution $p_{\mathbf{Y}}$ to joint distribution $p_{\mathbf{X},\mathbf{Y}}$ via marginalization:

$$0 = p_Y(y) = \sum_{x \in \mathcal{X}} p_{X,Y}(x, y).$$

Next, we use a crucial mathematical observation: If a sum of nonnegative numbers (such as probabilities) equals 0, then each of the numbers being summed up must also be 0 (otherwise, the sum would be positive!). Hence, it must be that each number being added up in the right-hand side sum is 0, i.e.,

$$p_{X,Y}(x, y) = 0 \quad \text{for all } x \in \mathcal{X}.$$

This completes the proof. \square

Thus, in general:

$$p_{X,Y}(x, y) = \begin{cases} p_Y(y)p_{X|Y}(x | y) & \text{if } p_Y(y) > 0, \\ 0 & \text{if } p_Y(y) = 0. \end{cases}$$

Important convention for this course: For notational convenience, throughout this course, we will often just write

$p_{X,Y}(x, y) = p_Y(y)p_{X|Y}(x | y)$ with the understanding that if $p_Y(y) = 0$, even though $p_{X|Y}(x | y)$ is not actually defined, $p_{X,Y}(x, y)$ just evaluates to 0 anyways.

The product rule is symmetric: We can use the definition of conditional probability with X and Y swapped, and rearranging factors, we get:

$$p_{X,Y}(x, y) = p_X(x)p_{Y|X}(y | x) \quad \text{for all } x \in \mathcal{X}, y \in \mathcal{Y} \text{ such that } p_X(x) > 0,$$

and so similarly we could show that

$$p_{X,Y}(x, y) = \begin{cases} p_X(x)p_{Y|X}(y | x) & \text{if } p_X(x) > 0, \\ 0 & \text{if } p_X(x) = 0. \end{cases}$$

Again for notational convenience, we'll typically just write

$p_{X,Y}(x, y) = p_X(x)p_{Y|X}(y | x)$ with the understanding that the expression is 0 when $p_X(x) = 0$.

Interpretation: If we're given the probability table for X and, separately, the probability table for Y conditioned on X , then we can come up with the joint probability table for X and Y .

Importantly, for any two jointly distributed random variables \mathbf{X} and \mathbf{Y} , the product rule is always true, without making any further assumptions! Also, as a recurring theme that we'll see later on as well, we are decomposing the joint distribution into the product of factors (in this case, the product of two factors).

Many random variables: If we have many random variables, say, $\mathbf{X}_1, \mathbf{X}_2$, up to \mathbf{X}_N where N is not a random variable but is a fixed constant, then we have

$$\begin{aligned} p_{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \\ = p_{\mathbf{X}_1}(\mathbf{x}_1) p_{\mathbf{X}_2 | \mathbf{X}_1}(\mathbf{x}_2 | \mathbf{x}_1) p_{\mathbf{X}_3 | \mathbf{X}_1, \mathbf{X}_2}(\mathbf{x}_3 | \mathbf{x}_1, \mathbf{x}_2) \\ \cdots p_{\mathbf{X}_N | \mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_{N-1}}(\mathbf{x}_N | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N-1}). \end{aligned}$$

Again, we write this to mean that this holds for every possible choice of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ for which we never condition on a zero probability event. Note that the above factorization always holds without additional assumptions on the distribution of $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$.

Note that the product rule could be applied in arbitrary orderings. In the above factorization, you could think of it as introducing random variable \mathbf{X}_1 first, and then \mathbf{X}_2 , and then \mathbf{X}_3 , etc. Each time we introduce another random variable, we have to condition on all the random variables that have already been introduced.

Since there are N random variables, there are $N!$ different orderings in which we can write out the product rule. For example, we can think of introducing the last random variable \mathbf{X}_N first and then going backwards until we introduce \mathbf{X}_1 at the end. This yields the, also correct, factorization

$$\begin{aligned} p_{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) \\ = p_{\mathbf{X}_N}(\mathbf{x}_N) p_{\mathbf{X}_{N-1} | \mathbf{X}_N}(\mathbf{x}_{N-1} | \mathbf{x}_N) p_{\mathbf{X}_{N-2} | \mathbf{X}_{N-1}, \mathbf{X}_N}(\mathbf{x}_{N-2} | \mathbf{x}_{N-1}, \mathbf{x}_N) \\ \cdots p_{\mathbf{X}_1 | \mathbf{X}_2, \mathbf{X}_3, \dots, \mathbf{X}_N}(\mathbf{x}_1 | \mathbf{x}_2, \dots, \mathbf{x}_N). \end{aligned}$$

© All Rights Reserved



© 2016 edX Inc. All rights reserved except where noted. EdX, Open edX and the edX and Open EdX logos are registered trademarks or trademarks of edX Inc.

POWERED BY
OPENedX®

