

Stochastic Processes 1

STAT 333

Fall 2021 (1219)¹

Cameron Roopnarine²

Steve Drekić³

Mirabelle Huynh⁴

13th October 2021

¹Online Course

²TeXer

³Main Instructor

⁴Assisting Instructor

Contents

1	Review of Elementary Probability	2
2	Conditional Distributions and Conditional Expectation	16
2.1	Definitions and Construction	16
2.2	Computing Expectation by Conditioning	27
2.3	Computing Probabilities by Conditioning	32
2.4	Some Further Extensions	36

Chapter 1

Review of Elementary Probability

WEEK 1
8th to 15th September

Fundamental Definition of a Probability Function

Probability Model: A probability model consists of 3 essential components: a *sample space*, a collection of *events*, and a *probability function (measure)*.

- **Sample Space:** For a random experiment in which all possible outcomes are known, the set of all possible outcomes is called the sample space (denoted by Ω).
- **Event:** Every subset A of a sample space Ω is an event.
- **Probability Function:** For each event A of Ω , $\mathbb{P}(A)$ is defined as the *probability of an event* A , satisfying 3 conditions:
 - (i) $0 \leq \mathbb{P}(A) \leq 1$,
 - (ii) $\mathbb{P}(\Omega) = 1$, or equivalently, $\mathbb{P}(\emptyset) = 0$, where \emptyset is the *null event*,
 - (iii) For $n \in \mathbb{Z}^+$ (in fact, $n = \infty$ as well), $\mathbb{P}(\bigcup_{i=1}^n A_i) = \sum_{i=1}^n \mathbb{P}(A_i)$ if the sequence of events $\{A_i\}_{i=1}^n$ is *mutually exclusive* (i.e., $A_i \cap A_j = \emptyset \forall i \neq j$).

As a result of conditions (ii) and (iii), and noting that A^c is the complement of A , it follows that

$$1 = \mathbb{P}(\Omega) = \mathbb{P}(A \cup A^c) = \mathbb{P}(A) + \mathbb{P}(A^c) \implies \mathbb{P}(A^c) = 1 - \mathbb{P}(A).$$

Conditional Probability

Conditional Probability: The *conditional probability of event* A *given event* B *occurs* is defined as

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)},$$

provided that $\mathbb{P}(B) > 0$.

Remarks:

- (1) When $B = \Omega$, $\mathbb{P}(A | \Omega) = \mathbb{P}(A \cap \Omega) / \mathbb{P}(\Omega) = \mathbb{P}(A) / 1 = \mathbb{P}(A)$, as one would expect.

- (2) Rewriting the above formula, $\mathbb{P}(A \cap B) = \mathbb{P}(A | B) \mathbb{P}(B)$, which is often referred to as the basic “multiplication rule.” For a sequence of events $\{A_i\}_{i=1}^n$, the generalized multiplication rule is given by

$$\mathbb{P}(A_1 \cap A_2 \cap \cdots \cap A_n) = \mathbb{P}(A_1) \mathbb{P}(A_2 | A_1) \cdots \mathbb{P}(A_n | A_1 \cap A_2 \cap \cdots \cap A_{n-1}).$$

Example 1.1. Suppose that we roll a fair six-sided die once (i.e., $\Omega = \{1, 2, 3, 4, 5, 6\}$). Let A denote the event of rolling a number less than 4 (i.e., $A = \{1, 2, 3\}$), and let B denote the event of rolling an odd number (i.e., $B = \{1, 3, 5\}$). Given that the roll is odd, what is the probability that number rolled is less than 4?

Solution: Since the die is fair, it immediately follows that $\mathbb{P}(A) = 3/6 = 1/2$ and $\mathbb{P}(B) = 3/6 = 1/2$. Moreover,

$$\mathbb{P}(A \cap B) = \mathbb{P}(\{1, 2, 3\} \cap \{1, 3, 5\}) \quad (1.1)$$

$$= \mathbb{P}(\{1, 3\}) \quad (1.2)$$

$$= \frac{2}{6} \quad (1.3)$$

$$= \frac{1}{3}. \quad (1.4)$$

Therefore,

$$\mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{1/3}{1/2} = \frac{2}{3}.$$

Independence of Events

Independence of Events: Two events A and B are *independent* if and only if (iff)

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$$

In general, if an experiment consists of a sequence of independent trials, and A_1, A_2, \dots, A_n are events such that A_i depends only on the i^{th} trial, then A_1, A_2, \dots, A_n are independent events and

$$\mathbb{P}(\cap_{i=1}^n A_i) = \prod_{i=1}^n \mathbb{P}(A_i).$$

Law of Total Probability

Law of Total Probability: For $n \in \mathbb{Z}^+$ (and even $n = \infty$), suppose that $\Omega = \cup_{i=1}^n B_i$, where the sequence

of events $\{B_i\}_{i=1}^n$ is mutually exclusive. Then,

$$\begin{aligned}
 \mathbb{P}(A) &= \mathbb{P}(A \cap \Omega) \\
 &= \mathbb{P}(A \cap \{\cup_{i=1}^n B_i\}) \\
 &= \mathbb{P}(\cup_{i=1}^n \{A \cap B_i\}) \\
 &= \sum_{i=1}^n \mathbb{P}(A \cap B_i) \\
 &= \sum_{i=1}^n \mathbb{P}(A | B_i) \mathbb{P}(B_i),
 \end{aligned}$$

where the second last equality follows from the fact that the sequence of events $\{A \cap B_i\}_{i=1}^n$ is also mutually exclusive.

Bayes' Formula

Bayes' Formula: Under the same assumptions as in the previous slide,

$$\mathbb{P}(B_j | A) = \frac{\mathbb{P}(A \cap B_j)}{\mathbb{P}(A)} = \frac{\mathbb{P}(A | B_j) \mathbb{P}(B_j)}{\sum_{i=1}^n \mathbb{P}(A | B_i) \mathbb{P}(B_i)}.$$

Definition of a Random Variable

Definition: A *random variable* (rv) X is a real-valued function which maps a sample space Ω onto a state space $\mathcal{S} \subseteq \mathbb{R}$ (i.e., $X: \Omega \rightarrow \mathcal{S}$).

Discrete type: \mathcal{S} consists of a finite or countable number of possible values. Important functions include:

$$\begin{aligned}
 p(a) &= \mathbb{P}(X = a) && \text{(pmf),} \\
 F(a) &= \mathbb{P}(X \leq a) = \sum_{x \leq a} p(x) && \text{(cdf),} \\
 \bar{F}(a) &= \mathbb{P}(X > a) = 1 - F(a) && \text{(tpf),}
 \end{aligned}$$

where pmf stands for *probability mass function*, cdf stands for *cumulative distribution function*, and tpf stands for *tail probability function*.

Remark: If X takes on values in the set $\mathcal{S} = \{a_1, a_2, a_3, \dots\}$ where $a_1 < a_2 < a_3 < \dots$ such that $p(a_i) > 0 \forall i$, then we can recover the pmf from knowledge of the cdf via

$$\begin{aligned}
 p(a_1) &= F(a_1), \\
 p(a_i) &= F(a_i) - F(a_{i-1}), \quad i = 2, 3, 4, \dots
 \end{aligned}$$

Discrete Distributions

Special Discrete Distributions:

1. **Bernoulli:** If we consider a *Bernoulli trial*, which is a random trial with probability p of being a “success” (denoted by 1) and a probability $1 - p$ of being a “failure” (denoted by 0), then X is *Bernoulli* (i.e., $X \sim \text{BERN}(p)$) with pmf

$$p(x) = p^x(1 - p)^{1-x}, \quad x = 0, 1.$$

2. **Binomial:** If X denotes the number of successes in $n \in \mathbb{Z}^+$ independent Bernoulli trials, each with probability p of being a success, then X is *Binomial* (i.e., $X \sim \text{BIN}(n, p)$) with pmf

$$p(x) = \binom{n}{x} p^x (1 - p)^{n-x}, \quad x = 0, 1, \dots, n,$$

where

$$\binom{n}{x} = \frac{n!}{(n-x)!x!} = \frac{(n)_x}{x!} = \frac{n(n-1) \cdots (n-x+1)}{x!}$$

is the number of distinct groups of x objects chosen from a set of n objects.

Remarks:

- (1) A $\text{BIN}(1, p)$ distribution simplifies to become the $\text{BERN}(p)$ distribution.
- (2) The binomial pmf is even defined for $n = 0$, in which case $p(x) = 1$ for $x = 0$. Such a distribution is said to be degenerate at 0.
- (3) Note that $\binom{n}{x} = 0$ if $n, x \in \mathbb{N}$ with $n < x$.

3. **Negative Binomial:** If X denotes the number of Bernoulli trials (each with success probability p) required to observe $k \in \mathbb{Z}^+$ successes, then X is *Negative Binomial* (i.e., $X \sim \text{NB}_t(k, p)$) with pmf

$$p(x) = \binom{x-1}{k-1} p^k (1-p)^{x-k}, \quad x = k, k+1, k+2, \dots$$

Remarks:

- (1) In the above pmf, $\binom{x-1}{k-1}$ appears rather than $\binom{x}{k}$ since the final trial must always be a success.
- (2) Sometimes, a negative binomial distribution is alternatively defined as the number of failures observed to achieve k successes. If Y denotes such a rv and $X \sim \text{NB}_t(k, p)$, then we clearly have the relationship $X = Y + k$, which immediately leads to the following pmf for Y :

$$p_Y(y) = \mathbb{P}(Y = y) = \mathbb{P}(X = y + k) = \binom{y+k-1}{k-1} p^k (1-p)^y, \quad y = 0, 1, 2, \dots$$

To refer to this negative binomial distribution, we will write $Y \sim \text{NB}_f(k, p)$.

4. **Geometric:** If $X \sim \text{NB}_t(1, p)$, then X is *Geometric* (i.e., $X \sim \text{GEO}_t(p)$) with pmf

$$p(x) = p(1-p)^{x-1}, \quad x = 1, 2, 3, \dots$$

In other words, the geometric distribution models the number of Bernoulli trials required to observe the first success.

Remark: Similarly, if $X \sim \text{NB}_f(1, p)$ then we obtain an alternative geometric distribution (denoted by $X \sim \text{GEO}_f(p)$) which models the number of failures observed prior to the first success.

5. **Discrete Uniform:** If X is equally likely to take on values in the (finite) set $\{a, a+1, \dots, b\}$ where $a, b \in \mathbb{Z}$ with $a \leq b$, then X is *Discrete Uniform* (i.e., $X \sim \text{DU}(a, b)$) with pmf

$$p(x) = \frac{1}{b-a+1}, \quad x = a, a+1, \dots, b.$$

6. **Hypergeometric:** If X denotes the number of success objects in n draws without replacement from a finite population of size N containing exactly r success objects, then X is *Hypergeometric* (i.e., $X \sim \text{HG}(N, r, n)$) with pmf

$$p(x) = \frac{\binom{r}{x} \binom{N-r}{n-x}}{\binom{N}{n}}, \quad x = \max\{0, n-N+r\}, \dots, \min\{n, r\}.$$

7. **Poisson:** A rv X is *Poisson* (i.e., $X \sim \text{POI}(\lambda)$) with parameter $\lambda > 0$ if its pmf is one of the form

$$p(x) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

Remark: The pmf is even defined for $\lambda = 0$ (if we use the standard convention that $0^0 = 1$), in which case $p(x) = 1$ for $x = 0$ (i.e., X is degenerate at 0).

Example 1.2. Show that when n is large and p is small, the $\text{BIN}(n, p)$ distribution may be approximated by a $\text{POI}(\lambda)$ distribution where $\lambda = np$.

Solution: Recall $e^z = \lim_{n \rightarrow \infty} (1 + z/n)^n$, $z \in \mathbb{R}$. Letting $X \sim \text{BIN}(n, p)$, we have

$$\begin{aligned} \mathbb{P}(X = x) &= \binom{n}{x} p^x (1-p)^{n-x} \\ &= \frac{n(n-1) \cdots (n-x+1)}{x!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \\ &= \frac{n}{n} \frac{n-1}{n} \cdots \frac{n-x+1}{n} \frac{\lambda^x}{x!} \frac{(1-\lambda/n)^n}{(1-\lambda/n)^x} \\ &\simeq (1)(1) \cdots (1) \frac{\lambda^x e^{-\lambda}}{x!} \quad \text{when } n \text{ is large} \\ &= \frac{e^{-\lambda} \lambda^x}{x!} \end{aligned}$$

Continuous Random Variables

Continuous type: A rv X takes on a continuum of possible values (which is uncountable) with cdf

$$F(x) = \mathbb{P}(X \leq x) = \int_{-\infty}^x f(y) dy,$$

where $f(x)$ denotes the *probability density function* (pdf) of X , which is a non-negative real-valued function that satisfies

$$\mathbb{P}(X \in B) = \int_{x \in B} f(x) dx,$$

where B is the set of real numbers (e.g., an interval).

Remarks:

- (1) If $F(x)$ (or the tpf $\bar{F}(x) = 1 - F(x)$) is known, we can recover the pdf using the relation

$$f(x) = \frac{d}{dx} F(x) = F'(x) = -\bar{F}'(x),$$

which holds by the *Fundamental Theorem of Calculus*.

- (2) When working with pdfs in general, it is usually not necessary to be precise about specifying whether a range of numbers includes the endpoints. This is quite different from the situation we encounter with discrete rvs. Throughout this course, however, we will adopt the convention of **not including** the endpoints when specifying the range of values for pdfs.

Continuous Distributions

Special Continuous Distributions:

1. **Uniform:** A rv X is *Uniform* on the real interval (a, b) (i.e., $X \sim U(a, b)$) if it has pdf

$$f(x) = \frac{1}{b-a}, \quad a < x < b,$$

where $a, b \in \mathbb{R}$ with $a < b$.

Remark: The choice of name is because X takes on values in (a, b) with all subintervals of a fixed length being equally likely.

2. **Beta:** A rv X is *Beta* with parameters $m \in \mathbb{Z}^+$ and $n \in \mathbb{Z}^+$ (i.e., $X \sim \text{Beta}(m, n)$) if it has pdf

$$f(x) = \frac{(m+n-1)!}{(m-1)!(n-1)!} x^{m-1} (1-x)^{n-1}, \quad 0 < x < 1.$$

Remark: A $\text{Beta}(1, 1)$ distribution simplifies to become the $U(0, 1)$ distribution.

3. **Erlang:** A rv X is *Erlang* with parameters $n \in \mathbb{Z}^+$ and $\lambda > 0$ (i.e., $X \sim \text{Erlang}(n, \lambda)$) if it has pdf

$$f(x) = \frac{\lambda^n x^{n-1} e^{-\lambda x}}{(n-1)!}, \quad x > 0.$$

Remark: The $\text{Erlang}(n, \lambda)$ distribution is actually a special case of the more general Gamma distribution in which n is extended to be any positive real number.

4. **Exponential:** A rv X is *Exponential* with parameter $\lambda > 0$ (i.e., $X \sim \text{EXP}(\lambda)$) if it has pdf

$$f(x) = \lambda e^{-\lambda x}, \quad x > 0.$$

Remark: An $\text{Erlang}(1, \lambda)$ distribution actually simplifies to become the $\text{EXP}(\lambda)$ distribution.

Expectation

Expectation: If $g(\cdot)$ is an arbitrary real-valued function, then

$$\mathbb{E}[g(X)] = \begin{cases} \sum_x g(x)p(x) & , \text{ if } X \text{ is a discrete rv,} \\ \int_{-\infty}^{\infty} g(x)f(x) dx & , \text{ if } X \text{ is a continuous rv.} \end{cases}$$

Special choices of $g(\cdot)$:

1. $g(X) = X^n, n \in \mathbb{N} \implies \mathbb{E}[g(X)] = \mathbb{E}[X^n]$ is the n^{th} moment of X . In general, moments serve to describe the shape of a distribution. If $n = 0$, then $\mathbb{E}[X^0] = 1$. If $n = 1$, then $\mathbb{E}[X] = \mu_X$ is the mean of X .

2. $g(X) = (X - \mathbb{E}[X])^2 \implies \mathbb{E}[g(X)] = \mathbb{E}[(X - \mathbb{E}[X])^2]$ is the variance of X . Note that

$$\text{Var}(X) = \sigma_X^2 = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2,$$

or equivalently

$$\sigma_X^2 = \mathbb{E}[X(X-1)] + \mathbb{E}[X] - \mathbb{E}[X]^2.$$

Related to this quantity, the *standard deviation* of X is $\sqrt{\text{Var}(X)} = \sigma_X$.

3. $g(X) = aX + b, a, b \in \mathbb{R}$ (i.e., $g(X)$ is a linear function of X). Note that

$$\mu_{aX+b} = \mathbb{E}[aX + b] = a\mu_X + b,$$

$$\sigma_{aX+b}^2 = \text{Var}(aX + b) = a^2 \sigma_X^2,$$

$$\sigma_{aX+b} = \sqrt{\text{Var}(aX + b)} = |a| \sigma_X.$$

Moment Generating Function

4. $g(X) = e^{tX}$, $t \in \mathbb{R} \implies E[g(X)] = E[e^{tX}]$ is the *moment generating function* (mgf) of X . This quantity is a function of t and is denoted by

$$\phi_X(t) = E[e^{tX}].$$

First, $\phi_X(0) = E[e^{0X}] = E[1] = 1$. Moreover, making use of the linearity property of the expected value operator, note that

$$\begin{aligned} \phi_X(t) &= E[e^{tX}] \\ &= E\left[\sum_{n=0}^{\infty} \frac{(tX)^n}{n!}\right] \\ &= E\left[\frac{t^0 X^0}{0!} + \frac{t^1 X^1}{1!} + \frac{t^2 X^2}{2!} + \cdots + \frac{t^n X^n}{n!} + \cdots\right] \\ &= E[X^0] \frac{t^0}{0!} + E[X] \frac{t^1}{1!} + E[X^2] \frac{t^2}{2!} + \cdots + E[X^n] \frac{t^n}{n!} + \cdots, \end{aligned}$$

implying that the n^{th} moment of X is simply the coefficient of $t^n/n!$ in the above series expansion.

We have: $\phi_X(t) = E[t^X] = E[X^0] \frac{t^0}{0!} + E[X] \frac{t^1}{1!} + E[X^2] \frac{t^2}{2!} + \cdots + E[X^n] \frac{t^n}{n!} + \cdots$.

Remarks:

- (1) Given the mgf of X , we can extract its n^{th} moment via

$$E[X^n] = \phi_X^{(n)}(0) = \left. \frac{d^n}{dt^n} \phi_X(t) \right|_{t=0} = \lim_{t \rightarrow 0} \frac{d^n}{dt^n} \phi_X(t), \quad n \in \mathbb{N}.$$

Note that the 0^{th} derivative of a function is simply the function itself.

- (2) A mgf uniquely characterizes the probability distribution of a rv (i.e., there exists a one-to-one correspondence between the mgf and the pmf/pdf of a rv). In other words, if two rvs X and Y have the same mgf, then they must have the same probability distribution (which we denote by $X \sim Y$). Thus, by finding the mgf of a rv, one has indeed determined its probability distribution.

Example 1.3. Suppose that $X \sim \text{BIN}(n, p)$. Find the mgf of X and use it to find $E[X]$ and $\text{Var}(X)$.

Solution: Recall the binomial series formula

$$(a + b)^m = \sum_{x=0}^m \binom{m}{x} a^x b^{m-x}, \quad a, b \in \mathbb{R}, \quad m \in \mathbb{N}.$$

Using this formula, we obtain

$$\begin{aligned}\phi_X(t) &= \mathbb{E}[e^{tX}] \\ &= \sum_{x=0}^n e^{tx} \binom{n}{x} p^x (1-p)^{n-x} \\ &= \sum_{x=0}^n \binom{n}{x} (pe^t)^x (1-p)^{n-x} \\ &= (pe^t + 1 - p)^n, \quad t \in \mathbb{R}.\end{aligned}$$

Then,

$$\phi'_X(t) = n(pe^t + 1 - p)^{n-1} pe^t \quad \text{and} \quad \phi''_X(t) = n(pe^t + 1 - p)^{n-1} pe^t + npe^t(n-1)(pe^t + 1 - p)^{n-2} pe^t.$$

Thus,

$$\begin{aligned}\mathbb{E}[X] &= \phi'_X(0) = n(pe^0 + 1 - p)^{n-1} pe^0 = np, \\ \text{Var}(X) &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \phi''_X(0) - n^2 p^2 = np + np(n-1)p - n^2 p^2 = np.\end{aligned}$$

Joint Distributions

Joint Distributions: The following results are presented for the bivariate case mostly, but these ideas extend naturally to an arbitrary number of rvs.

Definition: The *joint cdf* of X and Y is

$$\begin{aligned}F(a, b) &= \mathbb{P}(X \leq a, Y \leq b) \\ &= \mathbb{P}(\{X \leq a\} \cap \{Y \leq b\}), \quad a, b \in \mathbb{R}.\end{aligned}$$

Remark: If the joint cdf is known, then we can recover their marginal counterparts as follows:

$$\begin{aligned}F_X(a) &= \mathbb{P}(X \leq a) = F(a, \infty) = \lim_{b \rightarrow \infty} F(a, b), \\ F_Y(a) &= \mathbb{P}(Y \leq b) = F(\infty, b) = \lim_{a \rightarrow \infty} F(a, b).\end{aligned}$$

Jointly Discrete Case:

Joint pmf:

$$p(x, y) = \mathbb{P}(X = x, Y = y)$$

Marginals:

$$\begin{aligned}p_X(x) &= \mathbb{P}(X = x) = \sum_y p(x, y) \\ p_Y(y) &= \mathbb{P}(Y = y) = \sum_x p(x, y)\end{aligned}$$

Multinomial Distribution: Consider an experiment which is repeated $n \in \mathbb{Z}^+$ times, with one of $k \geq 2$ distinct outcomes possible each time. Let p_1, p_2, \dots, p_k denote the probabilities of the k types of outcomes (with $\sum_{i=1}^k p_i = 1$). If $X_i, i = 1, 2, \dots, k$, counts the number of type- i outcomes to occur,

then (X_1, X_2, \dots, X_k) is *Multinomial* (i.e., $(X_1, X_2, \dots, X_k) \sim \text{MN}(n, p_1, p_2, \dots, p_k)$) with joint pmf

$$p(x_1, x_2, \dots, x_k) = \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}, \quad x_i = 0, 1, \dots, n \quad \forall i \text{ and } \sum_{i=1}^k x_i = n$$

Remark: A $\text{MN}(n, p_1, 1 - p_1)$ distribution simplifies to become the $\text{BIN}(n, p_1)$ distribution.

Jointly Continuous Case:

Joint pdf: The joint pdf $f(x, y)$ is a non-negative real-valued function which enables one to calculate probabilities of the form

$$\mathbb{P}(X \in A, Y \in B) = \int_B \int_A f(x, y) \, dx \, dy = \int_A \int_B f(x, y) \, dx \, dy$$

where A and B are sets of real numbers (e.g., intervals). As a result,

$$F(a, b) = \int_{-\infty}^b \int_{-\infty}^a f(x, y) \, dx \, dy = \int_{-\infty}^a \int_{-\infty}^b f(x, y) \, dy \, dx$$

Marginals:

$$\begin{aligned} f_X(x) &= \int_{-\infty}^{\infty} f(x, y) \, dy \\ f_Y(y) &= \int_{-\infty}^{\infty} f(x, y) \, dx \end{aligned}$$

Jointly Continuous Case:

Important Relationship:

$$f(x, y) = \frac{\partial^2}{\partial x \partial y} F(x, y)$$

Transformations: Let (X, Y) be jointly continuous with joint pdf $f(x, y)$ and region of support $\mathcal{S}(X, Y)$. Suppose that the rvs V and W are given by $V = b_1(X, Y)$ and $W = b_2(X, Y)$, where the functions $v = b_1(x, y)$ and $w = b_2(x, y)$ defined a one-to-one transformation that maps the set $\mathcal{S}(X, Y)$ onto the set $\mathcal{S}(V, W)$. If x and y are expressed in terms of v and w (i.e., $x = h_1(v, w)$ and $y = h_2(v, w)$), then the joint pdf of V and W is given by

$$g(v, w) = \begin{cases} f(h_1(v, w), h_2(v, w)) |J| & , \text{ if } (v, w) \in \mathcal{S}(V, W), \\ 0 & , \text{ elsewhere,} \end{cases}$$

where J is the *Jacobian* of the transformation given by

$$J = \frac{\partial x}{\partial v} \frac{\partial y}{\partial w} - \frac{\partial x}{\partial w} \frac{\partial y}{\partial v}.$$

Expectation

Expectation: If $g(\cdot, \cdot)$ denotes an arbitrary real-valued function, then

$$\mathbb{E}[g(X, Y)] = \begin{cases} \sum_x \sum_y g(x, y) p(x, y) & , \text{ if } X \text{ and } Y \text{ are jointly discrete,} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f(x, y) dy dx & , \text{ if } X \text{ and } Y \text{ are jointly continuous.} \end{cases}$$

Remark: The order of summation/integration is irrelevant and can be interchanged.

Special choices of $g(\cdot, \cdot)$:

1. $g(X, Y) = (X - \mathbb{E}[X])(Y - \mathbb{E}[Y]) \implies \mathbb{E}[g(X, Y)] = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$ is the *covariance* of X and Y . Note that

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

and $\text{Cov}(X, X) = \text{Var}(X)$.

2. $g(X, Y) = aX + bY$, $a, b \in \mathbb{R}$ (i.e., $g(X, Y)$ is a linear combination of X and Y). Note that:

$$\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y],$$

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y).$$

3. $g(X, Y) = e^{sX+tY}$, $s, t \in \mathbb{R} \implies \mathbb{E}[g(X, Y)] = \mathbb{E}[e^{sX+tY}]$ is the joint mgf of X and Y . A joint mgf (denoted by $\phi(s, t)$) also uniquely characterizes a joint probability distribution and can be used to calculate joint moments of X and Y via the formula

$$\mathbb{E}[X^m Y^n] = \phi^{(m,n)}(0, 0) = \left(\frac{\partial^{m+n}}{\partial s^m \partial t^n} \phi(s, t) \right)_{s=0, t=0} = \lim_{s \rightarrow 0, t \rightarrow 0} \frac{\partial^{m+n}}{\partial s^m \partial t^n} \phi(s, t), \quad m, n \in \mathbb{N}$$

Independence of Random Variables

Formal Definition: If X and Y are *independent* rvs if

$$\begin{aligned} F(a, b) &= \mathbb{P}(X \leq a, Y \leq b) \\ &= \mathbb{P}(X \leq a) \mathbb{P}(Y \leq b) \\ &= F_X(a) F_Y(b) \quad \forall a, b \in \mathbb{R}. \end{aligned}$$

Equivalently, independence exists iff $p(x, y) = p_X(x)p_Y(y)$ (in the jointly discrete case) or $f(x, y) = f_X(x)f_Y(y)$ (in the jointly continuous case) $\forall x, y \in \mathbb{R}$.

Important Property: For arbitrary real-valued functions $g(\cdot)$ and $h(\cdot)$, if X and Y are independent, then

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

Remark: As a consequence of this property, $\text{Cov}(X, Y) = 0$ if X and Y are independent, implying that $\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y)$. However, if $\text{Cov}(X, Y) = 0$, we cannot conclude that X and Y are independent (we can only say that X and Y are *uncorrelated*).

Example 1.4. Suppose that X and Y have joint pmf (and corresponding marginals) of the form

		y		$p_X(x)$
		0	1	
x	$p(x, y)$	0.2	0	0.2
	0	0	0.6	0.6
	1	0.2	0	0.2
$p_Y(y)$		0.4	0.6	1

Show that $\text{Cov}(X, Y) = 0$ holds, but X and Y are not independent.

Solution: Recall that $\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$. Note that

$$\begin{aligned} E[XY] &= \sum_x \sum_y xyp(x, y) \\ &= (0)(0)(0.2) + (0)(1)(0) + (1)(0)(0.2) + (1)(1)(0.6) + (2)(0)(0.2) + (2)(1)(0) \\ &= 0.6, \end{aligned}$$

$$E[X] = \sum_x xp_X(x) = (0)(0.2) + (1)(0.6) + (2)(0.2) = 1,$$

$$E[Y] = \sum_y yp_Y(y) = (0)(0.4) + (1)(0.6) = 0.6.$$

Thus,

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y] = 0.6 - (1)(0.6) = 0.$$

However, from the given table, it is clear that $p(2, 0) = 0.2 \neq 0.08 = (0.2)(0.4) = p_X(2)p_Y(0)$. Thus, we conclude that while $\text{Cov}(X, Y) = 0$, X and Y are not independent.

Theorem 1.1. If X_1, X_2, \dots, X_n are independent rvs where $\phi_{X_i}(t)$ is the mgf of X_i , $i = 1, 2, \dots, n$, then $T = \sum_{i=1}^n X_i$ has mgf $\phi_T(t) = \prod_{i=1}^n \phi_{X_i}(t)$.

Proof: Note that the mgf of T given by

$$\begin{aligned} \phi_T(t) &= E[e^{tT}] \\ &= E[e^{t(X_1 + X_2 + \dots + X_n)}] \\ &= E[e^{tX_1} e^{tX_2} \dots e^{tX_n}] \\ &= E[e^{tX_1}] E[e^{tX_2}] \dots E[e^{tX_n}] && \text{by independence of } \{X_i\}_{i=1}^n \\ &= \phi_{X_1}(t) \phi_{X_2}(t) \dots \phi_{X_n}(t) \\ &= \prod_{i=1}^n \phi_{X_i}(t). \end{aligned}$$

Remarks:

- (1) Simply put, Theorem 1.1 states that the mgf of a sum of independent rvs is just the product of their individual mgfs.
- (2) As a special case of the above result, note that $\phi_T(t) = \phi_{X_1}(t)^n$ if X_1, X_2, \dots, X_n is an independent and identically distributed (iid) sequence of rvs.

Example 1.5. Let X_1, X_2, \dots, X_m be an independent sequence of rvs where $X_i \sim \text{BIN}(n_i, p)$, $i = 1, 2, \dots, m$. Find the distribution of $T = \sum_{i=1}^m X_i$.

Solution: Looking at the mgf of T , note that

$$\begin{aligned} \phi_T(t) &= \prod_{i=1}^m \phi_{X_i}(t) && \text{by Theorem 1.1} \\ &= \prod_{i=1}^m (pe^t + 1 - p)^{n_i} && \text{using the result of Example of 1.3} \\ &= (pe^t + 1 - p)^{\sum_{i=1}^m n_i}, \quad t \in \mathbb{R}. \end{aligned}$$

By the mgf uniqueness property we recognize that $T = \sum_{i=1}^m X_i \sim \text{BIN}(\sum_{i=1}^m n_i, p)$.

Remark: As a special case of the above example, if X_1, X_2, \dots, X_m are iid $\text{BERN}(p)$ rvs, then $T = \sum_{i=1}^m X_i \sim \text{BIN}(m, p)$.

Convergence of Random Variables

Modes of Convergence: If $X_n, n \in \mathbb{Z}^+$, and X are rvs, then

1. $X_n \rightarrow X$ in distribution iff

$$\lim_{n \rightarrow \infty} \mathbb{P}(X_n \leq x) = \mathbb{P}(X \leq x), \quad \forall x \in \mathbb{R} \text{ at which } \mathbb{P}(X \leq x) \text{ is continuous,}$$

2. $X_n \rightarrow X$ in probability, iff $\forall \varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0,$$

3. $X_n \rightarrow X$ almost surely (a.s.) iff

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} X_n = X\right) = 1.$$

Remarks:

- (1) In probability theory, an event is said to happen a.s. if it happens with probability 1.
- (2) The following implications hold true in general:

$$X_n \rightarrow X \text{ a.s.} \implies X_n \rightarrow X \text{ in probability} \implies X_n \rightarrow X \text{ in distribution.}$$

Strong Law of Large Numbers

Strong Law of Large Numbers (SLLN): If X_1, X_2, \dots, X_n is an iid sequence of rvs with common mean μ and $E[|X_1|] < \infty$, then

$$\bar{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n} \rightarrow \mu \text{ a.s. as } n \rightarrow \infty.$$

Remark: The SLLN is one of the most important results in probability and statistics, indicating that the sample mean will, with probability 1, converge to the true mean of the underlying distribution as the sample size approaches infinity. In other words, if the same experiment or study is repeated independently many times, the average of the results of the trials must be close to the mean. The result gets closer to the mean as the number of trials is increased.

Chapter 2

Conditional Distributions and Conditional Expectation

WEEK 2
15th to 22nd September

2.1 Definitions and Construction

Jointly Discrete Case

Formulation: If X_1 and X_2 are both discrete rvs with joint pmf $p(x_1, x_2)$ and marginal pmfs $p_1(x_1)$ and $p_2(x_2)$, respectively, then the conditional distribution of X_1 given $X_2 = x_2$, denoted by $X_1 | (X_2 = x_2)$, is defined via its *conditional pmf*

$$p_{1|2}(x_1 | x_2) = \mathbb{P}(X_1 = x_1 | X_2 = x_2) = \frac{\mathbb{P}(X_1 = x_1, X_2 = x_2)}{\mathbb{P}(X_2 = x_2)} = \frac{p(x_1, x_2)}{p_2(x_2)},$$

provided that $p_2(x_2) > 0$. Similarly, the conditional distribution of $X_2 | (X_1 = x_1)$ is defined via its conditional pmf

$$p_{2|1}(x_2 | x_1) = \mathbb{P}(X_2 = x_2 | X_1 = x_1) = \frac{p(x_1, x_2)}{p_1(x_1)}, \text{ provided that } p_1(x_1) > 0.$$

Remarks:

- (1) If X_1 and X_2 are independent, then $p(x_1, x_2) = p_1(x_1)p_2(x_2) \forall x_1, x_2 \in \mathbb{R}$, and so $p_{1|2}(x_1 | x_2) = p_1(x_1)$ and $p_{2|1}(x_2 | x_1) = p_2(x_2)$.
- (2) These ideas extend beyond the simple bivariate case naturally. For example, suppose that X_1, X_2 , and X_3 are discrete rvs. We can define the conditional distribution of (X_1, X_2) given $X_3 = x_3$ via its conditional pmf as follows:

$$p_{12|3}(x_1, x_2 | x_3) = \frac{\mathbb{P}(X_1 = x_1, X_2 = x_2, X_3 = x_3)}{\mathbb{P}(X_3 = x_3)} = \frac{p(x_1, x_2, x_3)}{p_3(x_3)},$$

provided that $p_3(x_3) > 0$. Alternatively, we can define the conditional distribution of X_2 given $(X_1 = x_1, X_3 = x_3)$ via its conditional pmf given by

$$p_{2|13}(x_2 | x_1, x_3) = \frac{p(x_1, x_2, x_3)}{p_{13}(x_1, x_3)}, \text{ provided that } p_{13}(x_1, x_3) > 0,$$

where $p_{13}(x_1, x_3)$ is the joint pmf of X_1 and X_3 .

Conditional Expectation: The *conditional mean* of $X_1 | (X_2 = x_2)$ is

$$\mathbb{E}[X_1 | X_2 = x_2] = \sum_{x_1} x_1 p_{1|2}(x_1 | x_2).$$

More generally, if $w(\cdot, \cdot)$, $h(\cdot)$, and $g(\cdot)$ are arbitrary real-valued functions, then

$$\mathbb{E}[w(X_1, X_2) | X_2 = x_2] = \mathbb{E}[w(X_1, x_2) | X_2 = x_2] = \sum_{x_1} w(x_1, x_2) p_{1|2}(x_1 | x_2)$$

and

$$\mathbb{E}[g(X_1)h(X_2) | X_2 = x_2] = \mathbb{E}[g(X_1)h(x_2) | X_2 = x_2] = h(x_2) \mathbb{E}[g(X_1) | X_2 = x_2].$$

As an immediate consequence, if $a, b \in \mathbb{R}$, then we obtain

$$\mathbb{E}[ag(X_1) + bh(X_1) | X_2 = x_2] = a \mathbb{E}[g(X_1) | X_2 = x_2] + b \mathbb{E}[h(X_1) | X_2 = x_2].$$

Furthermore, if we recall that $\mathbb{E}[X_1 + X_2] = \sum_{x_1} \sum_{x_2} (x_1 + x_2) p(x_1, x_2)$, then it correspondingly follows

that

$$\begin{aligned}
 E[X_1 + X_2 \mid X_3 = x_3] &= \sum_{x_1} \sum_{x_2} (x_1 + x_2) p_{12|3}(x_1, x_2 \mid x_3) \\
 &= \sum_{x_1} \sum_{x_2} (x_1 + x_2) \frac{p(x_1, x_2, x_3)}{p_3(x_3)} \\
 &= \sum_{x_1} \sum_{x_2} x_1 \cdot \frac{p(x_1, x_2, x_3)}{p_3(x_3)} + \sum_{x_1} \sum_{x_2} x_2 \cdot \frac{p(x_1, x_2, x_3)}{p_3(x_3)} \\
 &= \sum_{x_1} \frac{x_1}{p_3(x_3)} \sum_{x_2} p(x_1, x_2, x_3) + \sum_{x_2} \frac{x_2}{p_3(x_3)} \sum_{x_1} p(x_1, x_2, x_3) \\
 &= \sum_{x_1} \frac{x_1}{p_3(x_3)} p_{13}(x_1, x_3) + \sum_{x_2} \frac{x_2}{p_3(x_3)} p_{23}(x_2, x_3) \\
 &= \sum_{x_1} x_1 p_{1|3}(x_1 \mid x_3) + \sum_{x_2} x_2 p_{2|3}(x_2 \mid x_3) \\
 &= E[X_1 \mid X_3 = x_3] + E[X_2 \mid X_3 = x_3].
 \end{aligned}$$

We have: $E[X_1 + X_2 \mid X_3 = x_3] = E[X_1 \mid X_3 = x_3] + E[X_2 \mid X_3 = x_3]$. In other words, the conditional expected value is also a **linear** operator. In fact, more generally, if $a_i \in \mathbb{R}$, $i = 1, 2, \dots, n$, then the same essential approach can be used to show that

$$E\left[\sum_{i=1}^n a_i X_i \mid Y = y\right] = \sum_{i=1}^n a_i E[X_i \mid Y = y].$$

Conditional Variance: If we take $g(X_1) = (X_1 - E[X_1 \mid X_2 = x_2])^2$, then

$$E[g(X_1) \mid X_2 = x_2] = E\left[(X_1 - E[X_1 \mid X_2 = x_2])^2 \mid X_2 = x_2\right] = \text{Var}(X_1 \mid X_2 = x_2)$$

is the *conditional variance* of $X_1 \mid (X_2 = x_2)$.

As with the calculation of variance, the following result provides an alternative (and often times preferred) way to calculate $\text{Var}(X_1 \mid X_2 = x_2)$.

Theorem 2.1. $\text{Var}(X_1 \mid X_2 = x_2) = E[X_1^2 \mid X_2 = x_2] - E[X_1 \mid X_2 = x_2]^2$.

Proof:

$$\begin{aligned}
 \text{Var}(X_1 \mid X_2 = x_2) &= E\left[(X_1 - E[X_1 \mid X_2 = x_2])^2 \mid X_2 = x_2\right] \\
 &= E[X_1^2 - 2X_1 E[X_1 \mid X_2 = x_2] + E[X_1 \mid X_2 = x_2]^2 \mid X_2 = x_2] \\
 &= E[X_1^2] - 2E[X_1 \mid X_2 = x_2]^2 + E[X_1 \mid X_2 = x_2]^2 \\
 &= E[X_1^2 \mid X_2 = x_2] - E[X_1 \mid X_2 = x_2]^2
 \end{aligned}$$

Example 2.1. Suppose that X_1 and X_2 are discrete rvs having joint pmf of the form

$$p(x_1, x_2) = \begin{cases} 1/5 & , \text{ if } x_1 = 1 \text{ and } x_2 = 0, \\ 2/15 & , \text{ if } x_1 = 0 \text{ and } x_2 = 1, \\ 1/15 & , \text{ if } x_1 = 1 \text{ and } x_2 = 2, \\ 1/5 & , \text{ if } x_1 = 2 \text{ and } x_2 = 0, \\ 2/5 & , \text{ if } x_1 = 1 \text{ and } x_2 = 1, \\ 0 & , \text{ otherwise.} \end{cases}$$

Find the conditional distribution of $X_1 \mid (X_2 = 1)$. Also, calculate $E[X_1 \mid X_2 = 1]$ and $\text{Var}(X_1 \mid X_2 = 1)$.

Solution: Note that for problems of this nature, it often helps to create a table summarizing the information:

$p(x_1, x_2)$		x_2			$p_1(x_1)$
		0	1	2	
x_1	0	0	2/15	0	2/15
	1	1/5	2/5	1/15	2/3
	2	1/5	0	0	1/5
$p_2(x_2)$		2/5	8/15	1/15	1

Then,

- $p_{1|2}(0 \mid 1) = \mathbb{P}(X_1 = 0 \mid X_2 = 1) = (2/15)/(8/15) = 1/4$, and
- $p_{1|2}(1 \mid 1) = \mathbb{P}(X_1 = 1 \mid X_2 = 1) = (2/5)/(8/15) = 3/4$.

Thus, the conditional pmf of $X_1 \mid (X_2 = 1)$ can be represented as follows:

x_1	0	1
$p_{1 2}(x_1 \mid 1)$	1/4	3/4

Note that $X_1 \mid (X_2 = 1) \sim \text{BERN}(3/4)$. Thus, $E[X_1 \mid X_2 = 1] = 3/4$ and $\text{Var}(X_1 \mid X_2 = 1) = 3/4(1 - 3/4) = 3/16$.

Example 2.2. For $i = 1, 2$, suppose that $X_i \sim \text{BIN}(n_i, p)$ where X_1 and X_2 are independent. Find the conditional distribution of X_1 given $X_1 + X_2 = m$.

Solution: We want to find the conditional pmf of $X_1 \mid (Y = m)$, where $Y = X_1 + X_2$. Let this conditional pmf be denoted by $p_{X_1|Y}(x_1 \mid m) = \mathbb{P}(X_1 = x_1 \mid Y = m)$. Recall from Example 1.5 that

$$X_1 + X_2 \sim \text{BIN}(n_1 + n_2, p).$$

$$\begin{aligned} p_{X_1|Y}(x_1 | m) &= \frac{\mathbb{P}(X_1 = x_1, Y = m)}{\mathbb{P}(Y = m)} \\ &= \frac{\mathbb{P}(X_1 = x_1, X_1 + X_2 = m)}{\mathbb{P}(X_1 + X_2 = m)} \\ &= \frac{\mathbb{P}(X_1 = x_1, X_2 = m - x_1)}{\binom{n_1+n_2}{m} p^m (1-p)^{n_1+n_2-m}} \\ &= \frac{p_1(x_1) p_2(m-x_1)}{\binom{n_1+n_2}{m} p^m (1-p)^{n_1+n_2-m}} \\ &= \frac{\binom{n_1}{x_1} p^{x_1} (1-p)^{n_1-x_1} \binom{n_2}{m-x_1} p^{m-x_1} (1-p)^{n_2-(m-x_1)}}{\binom{n_1+n_2}{m} p^m (1-p)^{n_1+n_2-m}} \end{aligned}$$

provided that $0 \leq x_1 \leq n_1$, and $0 \leq m - x_1 \leq n_2$ (i.e., $m - n_2 \leq x_1 \leq m$). Simplifying,

$$p_{X_1|Y}(x_1 | m) = \frac{\binom{n_1}{x_1} \binom{n_2}{m-x_1}}{\binom{n_1+n_2}{m}},$$

for $x_1 = \max\{0, m - n_2\}, \dots, \min\{n_1, m\}$.

Remark: Looking at the conditional pmf we just obtained, we recognize that $X_1 | (X_1 + X_2 = m) \sim \text{HG}(n_1 + n_2, n_1, m)$. The result that $X_1 | (X_1 + X_2 = m)$ has a hypergeometric distribution should not be all that surprising. Consider the sequence of $n_1 + n_2$ Bernoulli trials represented visually as follows:

TODO figure

Of these $n_1 + n_2$ trials in which m of them were known to be successes, we want x_1 successes to have occurred among the first n_1 trials (thereby implying that $m - x_1$ successes are obtained during the final n_2 trials). Since any of these trials were equally likely to be a success (i.e., the same success probability p is assumed), the desired result ends up being the obtained hypergeometric probability.

Example 2.3. Let X_1, X_2, \dots, X_m be independent rvs where $X_i \sim \text{POI}(\lambda_i)$, $i = 1, 2, \dots, m$. Define $Y = \sum_{i=1}^m X_i$. Find the conditional distribution of $X_j | (Y = n)$.

Solution: We are interested in the conditional pmf of $X_j | (Y = n)$, to be denoted by

$$\begin{aligned} p_{X_j|Y}(x_j | n) &= \mathbb{P}(X_j = x_j | Y = n) \\ &= \frac{\mathbb{P}(X_j = x_j, Y = n)}{\mathbb{P}(Y = n)} \\ &= \frac{\mathbb{P}\left(X_j = x_j, \sum_{i=1}^m X_i = n\right)}{\mathbb{P}(Y = n)} \end{aligned}$$

First, we investigate the numerator:

$$\begin{aligned} \mathbb{P}\left(X_j = x_j, \sum_{i=1}^m X_i = n\right) &= \mathbb{P}\left(X_j = x_j, X_j + \sum_{i=1, i \neq j}^m X_i = n\right) \\ &= \mathbb{P}\left(X_j = x_j, \sum_{i=1, i \neq j}^m X_i = n - x_j\right) \\ &= \mathbb{P}(X_j = x_j) \mathbb{P}\left(\sum_{i=1, i \neq j}^m X_i = n - x_j\right) \end{aligned}$$

where the last equality follows due to the independence of $\{X_i\}_{i=1}^m$. We are given that $X_j \sim \text{POI}(\lambda_j)$. Due to the result of Exercise 1.1, it follows that

$$\sum_{i=1, i \neq j}^m X_i \sim \text{POI}\left(\sum_{i=1, i \neq j}^m \lambda_i\right).$$

By the same result, we also have that

$$Y = \sum_{i=1}^m X_i \sim \text{POI}\left(\sum_{i=1}^m \lambda_i\right).$$

Therefore,

$$p_{X_j|Y}(x_j | n) = \frac{\frac{e^{-\lambda_j} \lambda_j^{x_j}}{x_j!} \frac{e^{-\sum_{i=1, i \neq j}^m \lambda_i} (\sum_{i=1, i \neq j}^m \lambda_i)^{n-x_j}}{(n-x_j)!}}{\frac{e^{-\sum_{i=1}^m \lambda_i} (\sum_{i=1}^m \lambda_i)^n}{n!}}$$

provided that $x_j \geq 0$ and $n - x_j \geq 0$ which implies $0 \leq x_j \leq n$. Thus,

$$\begin{aligned} p_{X_j|Y}(x_j | n) &= \binom{n}{x_j} \frac{\lambda_j^{x_j} (\lambda_Y - \lambda_j)^{n-x_j}}{\lambda_Y^n} \\ &= \binom{n}{x_j} \left(\frac{\lambda_j}{\lambda_Y}\right)^{x_j} \left(1 - \frac{\lambda_j}{\lambda_Y}\right)^{n-x_j}, \quad x_j = 0, 1, \dots, n \end{aligned}$$

where $\lambda_Y = \sum_{i=1}^m \lambda_i$ and note that $\lambda_Y^{x_j} \lambda_Y^{n-x_j} = \lambda_Y^n$. We see that

$$X_j | (Y = n) \sim \text{BIN}\left(n, \frac{\lambda_j}{\sum_{i=1}^m \lambda_i}\right).$$

Example 2.4. Suppose that $X \sim \text{POI}(\lambda)$ and $Y | (X = x) \sim \text{BIN}(x, p)$. Find the conditional distribution of $X | (Y = y)$.

Solution: We want to calculate the conditional pmf of $X | (Y = y)$, to be denoted by

$$p_{X|Y}(x | y) = \mathbb{P}(X = x | Y = y) = \frac{\mathbb{P}(X = x, Y = y)}{\mathbb{P}(Y = y)}.$$

First, note that

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

which implies that

$$\mathbb{P}(X = x, Y = y) = \mathbb{P}(Y = y | X = x) \mathbb{P}(X = x) = \frac{e^{-\lambda} \lambda^x}{x!} \binom{x}{y} p^y (1-p)^{x-y},$$

for $x = 0, 1, 2, \dots$ and $y = 0, 1, \dots, x$. Note that the range of y depends on the values of x . A graphical display of the region is given below: TODO

We may rewrite this region with the range of x depending on the values of y . Specifically, note that $x = 0, 1, 2, \dots$ and $y = 0, 1, \dots, x$ is equivalent to $y = 0, 1, 2, \dots$ and $x = y, y+1, y+2, \dots$. We use this

alternative region to find the marginal pmf of Y .

$$\begin{aligned}
 \mathbb{P}(Y = y) &= \sum_x \mathbb{P}(X = x, Y = y) \\
 &= \sum_{x=y}^{\infty} e^{-\lambda} \frac{\lambda^x}{x!} \binom{x}{y} p^y (1-p)^{x-y} \\
 &= \sum_{x=y}^{\infty} e^{-\lambda} \frac{\lambda^x}{x!} \frac{x!}{y!(x-y)!} p^y (1-p)^{x-y} \\
 &= \frac{e^{-\lambda}}{y!} p^y \sum_{x=y}^{\infty} \frac{\lambda^x (1-p)^{x-y}}{(x-y)!} \lambda^{-y} \lambda^y \\
 &= \frac{e^{-\lambda} (\lambda p)^y}{y!} \sum_{x=y}^{\infty} \frac{(\lambda(1-p))^{x-y}}{(x-y)!} && \text{let } z = x - y \\
 &= \frac{e^{-\lambda} (\lambda p)^y}{y!} e^{\lambda(1-p)} \\
 &= \frac{e^{-\lambda p} (\lambda p)^y}{y!} && y = 0, 1, 2, \dots
 \end{aligned}$$

In fact, $Y \sim \text{POI}(\lambda p)$. Therefore,

$$\begin{aligned}
 p_{X|Y}(x | y) &= \frac{\frac{e^{-\lambda} \lambda^x}{x!} \frac{x!}{y!(x-y)!} p^y (1-p)^{x-y}}{\frac{e^{-\lambda p} (\lambda p)^y}{y!}} \\
 &= \frac{e^{-\lambda(1-p)} (\lambda(1-p))^{x-y}}{(x-y)!},
 \end{aligned}$$

for $x = y, y + 1, \dots$

Remark: The above conditional pmf is recognized as that of a **shifted** Poisson distribution (y units to the right). Specifically, we have that

$$X | (Y = y) \sim W + y$$

where $W \sim \text{POI}(\lambda(1-p))$.

Formulation: In the jointly discrete case, it was natural to define:

$$p_{X|Y}(x | y) = \mathbb{P}(X = x | Y = y) = \mathbb{P}(X = x, Y = y) / \mathbb{P}(Y = y).$$

Strictly speaking, this no longer makes sense in a continuous context since $f(x, y) \neq \mathbb{P}(X = x, Y = y)$ and $f_Y(y) \neq \mathbb{P}(Y = y)$. However, for small positive values of dy (as the figure below shows), $\mathbb{P}(y \leq Y \leq y + dy) \approx f_Y(y) dy$.

Formally,

$$f_Y(y) = \lim_{dy \rightarrow 0} \frac{\mathbb{P}(y \leq Y \leq y + dy)}{dy}.$$

Similarly,

$$f(x, y) = \lim_{dx \rightarrow 0, dy \rightarrow 0} \frac{\mathbb{P}(x \leq X \leq x + dx, y \leq Y \leq y + dy)}{dx dy},$$

which implies that $\mathbb{P}(x \leq X \leq x + dx, y \leq Y \leq y + dy) \approx f(x, y) dx dy$. For small positive values of dx and

dy, consider now

$$\begin{aligned}\mathbb{P}(x \leq X \leq x + dx \mid y \leq Y \leq y + dy) &= \frac{\mathbb{P}(x \leq X \leq x + dx \mid y \leq Y \leq y + dy)}{\mathbb{P}(y \leq Y \leq y + dy)} \\ &\approx \frac{f(x, y) dx dy}{f_Y(y) dy} \\ &= \frac{f(x, y)}{f_Y(y)} dx.\end{aligned}$$

As a result, we formally define the *conditional pdf* of X given $Y = y$ (again to be denoted by $X \mid (Y = y)$) as

$$f_{X|Y}(x \mid y) = \frac{f(x, y)}{f_Y(y)} = \lim_{dx \rightarrow 0, dy \rightarrow 0} \frac{\mathbb{P}(x \leq X \leq x + dx, y \leq Y \leq y + dy)}{dx}.$$

Remark: In the jointly continuous case, the conditional probability of an event of the form $\{a \leq X \leq b\}$ given $Y = y$ would be calculated as

$$\mathbb{P}(a \leq X \leq b \mid Y = y) = \int_a^b f_{X|Y}(x \mid y) dx = \frac{\int_a^b f(x, y) dx}{f_Y(y)},$$

which we can also express as

$$\mathbb{P}(a \leq X \leq b \mid Y = y) = \frac{\int_a^b f(x, y) dx}{\int_{-\infty}^{\infty} f(x, y) dx}.$$

In other words, we could view this as a way of assigning probability to an event $\{a \leq X \leq b\}$ over a “slice,” $Y = y$, of the (joint) region of support for the pair of rvs X and Y .

Example 2.5. Suppose that the joint pdf of X and Y is given by

$$f(x, y) = \begin{cases} 5e^{-3x-y} & , \text{ if } 0 \leq 2x \leq y < \infty, \\ 0 & , \text{ elsewhere.} \end{cases}$$

Determine the conditional distribution of $Y \mid (X = x)$ where $0 \leq x < \infty$.

Solution: We wish to find the conditional pdf of $Y \mid (X = x)$ given by

$$f_{Y|X}(y \mid x) = \frac{f(x, y)}{f_X(x)}$$

The region of support for this joint distribution looks like: TODO figure
For $0 < x < \infty$:

$$\begin{aligned}f_X(x) &= \int_{-\infty}^{\infty} f(x, y) dy \\ &= \int_{2x}^{\infty} 5e^{-3x-y} dy \\ &= \left[5e^{-3x}(-e^{-y}) \right]_{y=2x}^{y=\infty} \\ &= 5e^{-3x}e^{-2x} \\ &= 5e^{-5x}\end{aligned}$$

Note that $X \sim \text{EXP}(5)$. Finally, we get:

$$f_{Y|X}(y \mid x) = \frac{5e^{-3x-y}}{5e^{-5x}} = e^{-y+2x}, \quad y > 2x.$$

Remark: The conditional pdf of $Y \mid (X = x)$ is recognized as that of a *shifted exponential distribution* ($2x$ units to the right). Specifically, we have that $Y \mid (X = x) \sim W + 2x$, where $W \sim \text{EXP}(1)$.

Conditional Expectation: If X and Y are jointly continuous rvs and $g(\cdot)$ is an arbitrary real-valued function, then the *conditional expectation* of $g(X)$ given $Y = y$ is

$$E[g(X) \mid Y = y] = \int_{-\infty}^{\infty} g(x) f_{X|Y}(x \mid y) dx,$$

and so the conditional mean of $X \mid (Y = y)$ is given by

$$E[X \mid Y = y] = \int_{-\infty}^{\infty} x f_{X|Y}(x \mid y) dx.$$

Example 2.6. Suppose that the joint pdf of X and Y is given by

$$f(x, y) = \begin{cases} \frac{12}{5}x(2 - x - y) & , \text{ if } 0 < x < 1 \text{ and } 0 < y < 1, \\ 0 & , \text{ elsewhere.} \end{cases}$$

Find the conditional distribution of X given $Y = y$ where $0 < y < 1$, and use it to calculate its conditional mean.

Solution: Using our earlier theory, we wish to find the conditional pdf of $X \mid (Y = y)$ given by

$$f_{X|Y}(x \mid y) = \frac{f(x, y)}{f_Y(y)}.$$

The region of support for this joint distribution of X and Y look like: TODO figure.
For $0 < y < 1$,

$$\begin{aligned} f_Y(y) &= \int_{-\infty}^{\infty} f(x, y) dx \\ &= \int_0^1 \frac{12}{5}x(2 - x - y) dx \\ &= \frac{12}{5} \int_0^1 (2x - x^2 - xy) dx \\ &= \frac{12}{5} \left[x^2 - \frac{x^3}{3} - \frac{x^2 y}{2} \right]_{x=0}^{x=1} \\ &= \frac{12}{5} \left(1 - \frac{1}{3} - \frac{y}{2} \right) \\ &= \frac{2(4 - 3y)}{5} \end{aligned}$$

You can verify this by integrating $f_Y(y)$ over the support of Y (to get 1). Thus,

$$f_{X|Y}(x \mid y) = \frac{12/5x(2 - x - y)}{2/5(4 - 3y)} = \frac{6x(2 - x - y)}{4 - 3y}, \quad 0 < x < 1$$

The conditional mean of X given $Y = y$ is:

$$\begin{aligned}
 E[X | Y = y] &= \int_0^1 x \frac{6x(2-x-y)}{4-3y} dx \\
 &= \frac{6}{4-3y} \int_0^1 (2x^2 - x^3 - x^2y) dx \\
 &= \frac{6}{4-3y} \left[\frac{2x^3}{3} - \frac{x^4}{4} - \frac{x^3y}{3} \right]_{x=0}^{x=1} \\
 &= \frac{6}{4-3y} \left(\frac{2}{3} - \frac{1}{4} - \frac{y}{3} \right) \\
 &= \frac{5-4y}{2(4-3y)}
 \end{aligned}$$

Conditional Variance: Likewise, as in the jointly discrete case, we can also consider the notion of conditional variance, which retains the same definition as before:

$$\text{Var}(X | Y = y) = E[(X - E[X | Y = y])^2 | Y = y] = E[X^2 | Y = y] - E[X | Y = y]^2.$$

A fact that is becoming more and more evident is that conditional expectation inherits many of the properties from regular expectation. Moreover, the same properties concerning conditional expectation that held in the jointly discrete case continue to hold true in the jointly continuous case (as we are effectively replacing summation with integration).

Example 2.6. (continued) Calculate $\text{Var}(X | Y = y)$ where $0 < y < 1$ and the joint pdf of X and Y is given by

$$f(x, y) = \begin{cases} \frac{12}{5}x(2-x-5y) & , \text{ if } 0 < x < 1 \text{ and } 0 < y < 1, \\ 0 & , \text{ elsewhere.} \end{cases}$$

Solution: Our earlier results tell us that

$$\begin{aligned}
 E[X^2 | Y = y] &= \int_0^1 x^2 \frac{6x(2-x-y)}{4-3y} dx \\
 &= \frac{6}{4-3y} \int_0^1 (2x^3 - x^4 - x^3y) dx \\
 &= \frac{6}{4-3y} \left[\frac{x^4}{2} - \frac{x^5}{5} - \frac{x^4y}{4} \right]_{x=0}^{x=1} \\
 &= \frac{6}{4-3y} \left(\frac{1}{2} - \frac{1}{5} - \frac{y}{4} \right) \\
 &= \frac{3(6-5y)}{10(4-3y)}.
 \end{aligned}$$

Therefore, this leads to

$$\begin{aligned}
 \text{Var}(X | Y = y) &= E[X^2 | Y = y] - E[X | Y = y]^2 \\
 &= \frac{3(6-5y)}{10(4-3y)} - \frac{(5-4y)^2}{4(4-3y)^2} \\
 &= \frac{19 + 2y(5y-14)}{20(4-3y)^2}.
 \end{aligned}$$

Mixed Case

We can also consider conditional distributions where the rvs are neither jointly continuous nor jointly discrete. To consider such a situation, suppose X is a continuous rv having pdf $f_X(x)$ and Y is a discrete rv having pmf $p_Y(y)$.

If we focus on the conditional distribution of X given $Y = y$, then let us look at the following quantity:

$$\begin{aligned} \frac{\mathbb{P}(x \leq X \leq x + dx \mid Y = y)}{dx} &= \frac{\mathbb{P}(x \leq X \leq x + dx, Y = y)}{dx \mathbb{P}(Y = y)} \\ &= \frac{\mathbb{P}(x \leq X \leq x + dx) \mathbb{P}(Y = y \mid x \leq X \leq x + dx)}{dx \mathbb{P}(Y = y)} \\ &= \frac{\mathbb{P}(Y = y \mid x \leq X \leq x + dx) \mathbb{P}(x \leq X \leq x + dx)}{\mathbb{P}(Y = y) dx}, \end{aligned}$$

where dx is again, a small positive value.

By letting $dx \rightarrow 0$, we can formally define the conditional pdf of $X \mid (Y = y)$ as follows:

$$\begin{aligned} f(x \mid y) &= \lim_{dx \rightarrow 0} \frac{\mathbb{P}(x \leq X \leq x + dx \mid Y = y)}{dx} \\ &= \lim_{dx \rightarrow 0} \frac{\mathbb{P}(Y = y \mid x \leq X \leq x + dx) \mathbb{P}(x \leq X \leq x + dx)}{\mathbb{P}(Y = y) dx} \\ &= \frac{\mathbb{P}(Y = y \mid X = x)}{\mathbb{P}(Y = y)} f_X(x) \\ &= \frac{p(y \mid x) f_X(x)}{p_Y(y)}, \end{aligned}$$

where $p(y \mid x) = \mathbb{P}(Y = y \mid X = x)$ is defined as the conditional pmf of $Y \mid (X = x)$. Note that since $f(x \mid y)$ is a pdf, it follows that

$$\int_{-\infty}^{\infty} f(x \mid y) dx = 1 \implies p_Y(y) = \int_{-\infty}^{\infty} p(y \mid x) f_X(x) dx.$$

Similarly, we can also write

$$p(y \mid x) = \frac{f(x \mid y) p_Y(y)}{f_X(x)}.$$

Since $p(y \mid x)$ is a pmf, we have that

$$\sum_y p(y \mid x) = 1 \implies f_X(x) = \sum_y f(x \mid y) p_Y(y).$$

Example 2.7. Suppose that $X \sim U(0, 1)$ and $Y \mid (X = x) \sim \text{BERN}(x)$. Find the conditional distribution of $X \mid (Y = y)$.

Solution: We wish to find the conditional pdf of $X \mid (Y = y)$ given by

$$f(x \mid y) = \frac{p(y \mid x) f_X(x)}{p_Y(y)}$$

Based on the given information, we have

$$\begin{aligned} f_X(x) &= 1, \quad 0 < x < 1, \\ p(y \mid x) &= x^y (1 - x)^{1-y}, \quad y = 0, 1. \end{aligned}$$

For $y = 0, 1$, note that

$$\begin{aligned} p_Y(y) &= \int_{-\infty}^{\infty} p(y | x) f_X(x) dx \\ &= \int_0^1 x^y (1-x)^{1-y} dx \end{aligned}$$

- For $y = 0 \implies p_Y(0) = \int_0^1 (1-x) dx = [x - x^2/2]_{x=0}^{x=1} = 1/2$.
- For $y = 1 \implies p_Y(1) = \int_0^1 x dx = [x^2/2]_{x=0}^{x=1} = 1/2$.

In other words, we have that

$$p_Y(y) = \frac{1}{2}, \quad y = 0, 1 \implies Y \sim \text{BERN}\left(\frac{1}{2}\right)$$

Thus, for $y = 0, 1$, we ultimately obtain

$$f(x | y) = \frac{x^y (1-x)^{1-y}}{1/2} = 2x^y (1-x)^{1-y}, \quad 0 < x < 1.$$

2.2 Computing Expectation by Conditioning

An Important Observation

As before, let $g(\cdot)$ be an arbitrary real-valued function. In general, we recognize that $E[g(X) | Y = y] = v(y)$, where $v(y)$ is some function of y . With this in mind, let us make the following definition:

$$E[g(X) | Y] = E[g(X) | Y = y]_{y=Y} = v(Y).$$

Functions of rvs are, once again, rvs themselves. Therefore, it makes sense to consider the expected value of $v(Y)$. In this regard, we would obtain:

$$\begin{aligned} E[E[g(X) | Y]] &= E[v(Y)] \\ &= \begin{cases} \sum_y v(y) p_Y(y) & , \text{ if } Y \text{ is discrete,} \\ \int_{-\infty}^{\infty} v(y) f_Y(y) dy & , \text{ if } Y \text{ is continuous,} \end{cases} \\ &= \begin{cases} \sum_y E[g(X) | Y = y] p_Y(y) & , \text{ if } Y \text{ is discrete,} \\ \int_{-\infty}^{\infty} E[g(X) | Y = y] f_Y(y) dy & , \text{ if } Y \text{ is continuous.} \end{cases} \end{aligned}$$

Law of Total Expectation

The following important result is regarded as the *law of total expectation*.

Theorem 2.2. For rvs X and Y , $E[g(X)] = E[E[g(X) | Y]]$.

Proof: Without loss of generality, assume that X and Y are jointly continuous rvs. From above, we have

$$\begin{aligned}
 E[E[g(X) | Y]] &= \int_{-\infty}^{\infty} E[g(X) | Y = y] f_Y(y) dy \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x) f_{X|Y}(x | y) dx f_Y(y) dy \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x) \frac{f(x, y)}{f_Y(y)} f_Y(y) dx dy \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x) f(x, y) dy dx \\
 &= \int_{-\infty}^{\infty} g(x) \int_{-\infty}^{\infty} f(x, y) dy dx \\
 &= \int_{-\infty}^{\infty} g(x) f_X(x) dx \\
 &= E[g(X)]
 \end{aligned}$$

Remark: Using a similar method of proof, the result of Theorem 2.2 can naturally be extended as follows:

$$E[g(X, Y)] = E[E[g(X, Y) | Y]].$$

The usefulness of the law of total expectation is well-demonstrated in the following example.

Example 2.8. Suppose that $X \sim \text{GEO}_t(p)$ with pmf $p_X(x) = (1 - p)^{x-1}p$, $x = 1, 2, 3, \dots$. Calculate $E[X]$ and $\text{Var}(X)$ using the law of total expectation.

Solution: With $X \sim \text{GEO}_t(p)$, recall that X actually models the number of (independent) trials necessary to obtain the first success. Define:

$$Y = \begin{cases} 0 & \text{, if the 1st trial is a failure,} \\ 1 & \text{, if the 1st trial is a success.} \end{cases}$$

We observe that $Y \sim \text{BERN}(p)$, so that $p_Y(0) = 1 - p$ and $p_Y(1) = p$.

Note:

- $X | (Y = 1)$ is degenerate at 1 (i.e., X given $Y = 1$ is equal to 1 with probability 1).
- $X | (Y = 0)$ is equivalent in distribution $1 + X$ (i.e., $X | (Y = 0) \sim 1 + X$).

By the law of total expectation, we obtain:

$$\begin{aligned}
 E[X] &= E[E[X | Y]] \\
 &= \sum_{y=0}^1 E[X | Y = y] p_Y(y) \\
 &= (1 - p) E[X | Y = 0] + p E[X | Y = 1] \\
 &= (1 - p) E[1 + X] + p \\
 &= (1 - p) + (1 - p) E[X] + p \\
 &= 1 + (1 - p) E[X],
 \end{aligned}$$

which implies that $(1 - (1 - p)) E[X] = 1$, or simply $E[X] = 1/p$. Similarly, we use the law of total

expectation to get

$$\begin{aligned}
 E[X^2] &= E[E[X^2 | Y]] \\
 &= \sum_{y=0}^1 E[X^2 | Y = y] p_Y(y) \\
 &= (1-p) E[X^2 | Y = 0] + p E[X^2 | Y = 1] \\
 &= (1-p) E[(1+X)^2] + p \\
 &= (1-p) (E[X^2] + 2E[X] + 1) + p \\
 &= 1 + (1-p) E[X^2] + \frac{2(1-p)}{p},
 \end{aligned}$$

which implies that

$$(1 - (1-p)) E[X] = \frac{p + 2(1-p)}{p}$$

or simply

$$E[X^2] = \frac{p + 2 - 2p}{p^2} = \frac{2-p}{p^2}$$

Finally,

$$\text{Var}(X) = \frac{2-p}{p^2} - \left(\frac{1}{p}\right)^2 = \frac{1-p}{p^2}$$

Remarks:

- (1) Note that the obtained mean and variance agree with known results. Moreover, the above procedure relied only on basic manipulations and did not involve any complicated sums or the differentiation of an mgf.
- (2) As part of the above solution, we claimed that $X | (Y = 0) \sim Z$ where $Z = 1 + X$, and this implied that $E[X^2 | Y = 0] = E[(1 + X)^2]$. To see why this holds true formally, consider first

$$p_{X|Y}(x | 0) = \mathbb{P}(X = x | Y = 0) = \frac{\mathbb{P}(X = x, Y = 0)}{\mathbb{P}(Y = 0)} = \frac{\mathbb{P}(X = x, Y = 0)}{1-p}.$$

Note that

$$\begin{aligned}
 \mathbb{P}(X = x, Y = 0) &= \mathbb{P}(\text{1st trial is a failure and } x \text{ total trials needed to get 1st success}) \\
 &= \mathbb{P}(\text{1st trial is a failure, next } x-2 \text{ trials are failures, and } x^{\text{th}} \text{ trial is a success}) \\
 &= (1-p)(1-p)^{x-2}p \text{ due to independence of trials.}
 \end{aligned}$$

Thus,

$$p_{X|Y}(x | 0) = \frac{(1-p)(1-p)^{x-2}p}{1-p} = (1-p)^{x-2}p, \quad x = 2, 3, 4, \dots$$

On the other hand, note that

$$\begin{aligned}
 p_Z(z) &= \mathbb{P}(Z = z) \\
 &= \mathbb{P}(1 + X = z) \\
 &= \mathbb{P}(X = z - 1) \\
 &= (1-p)^{(z-1)-1}p \\
 &= (1-p)^{z-2}p, \quad z = 2, 3, 4, \dots
 \end{aligned}$$

Since these two pmfs are identical, it follows that $X \mid (Y = 0) \sim Z$. As a further consequence, for an arbitrary real-valued function $g(\cdot)$, we must have that

$$\mathbb{E}[g(X) \mid Y = 0] = \mathbb{E}[g(Z)] = \mathbb{E}[g(1 + X)].$$

Computing Variances by Conditioning

In recognizing that $\mathbb{E}[g(X) \mid Y = y]$ is a function of y , it similarly follows that $\text{Var}(X \mid Y = y)$ is also a function of y . Therefore, we can make the following definition:

$$\text{Var}(X \mid Y) = \text{Var}(X \mid Y = y) \Big|_{y=Y}.$$

Since $\text{Var}(X \mid Y)$ is a function of Y , it is a rv as well, meaning that we could take its expected value. The following result, usually referred to as the *conditional variance formula*, provides a convenient way to calculate variance through the use of conditioning.

Theorem 2.3. For rvs X and Y , $\text{Var}(X) = \mathbb{E}[\text{Var}(X \mid Y)] + \text{Var}(\mathbb{E}[X \mid Y])$.

Proof: First, consider the term $\mathbb{E}[\text{Var}(X \mid Y)]$. Since

$$\text{Var}(X \mid Y = y) = \mathbb{E}[X^2 \mid Y = y] - \mathbb{E}[X \mid Y = y]^2,$$

it follows that

$$\text{Var}(X \mid Y) = \mathbb{E}[X^2 \mid Y] - \mathbb{E}[X \mid Y]^2,$$

which yields (by Theorem 2.2)

$$\begin{aligned} \mathbb{E}[\text{Var}(X \mid Y)] &= \mathbb{E}[\mathbb{E}[X^2 \mid Y] - \mathbb{E}[X \mid Y]^2] \\ &= \mathbb{E}[\mathbb{E}[X^2 \mid Y]] - \mathbb{E}[\mathbb{E}[X \mid Y]^2] \\ &= \mathbb{E}[X^2] - \mathbb{E}[\mathbb{E}[X \mid Y]^2]. \end{aligned}$$

Next, recall

$$\text{Var}(v(Y)) = \mathbb{E}[v(Y)^2] - \mathbb{E}[v(Y)]^2.$$

Applying Theorem 2.2 once more,

$$\begin{aligned} \text{Var}(\mathbb{E}[X \mid Y]) &= \mathbb{E}[\mathbb{E}[X \mid Y]^2] - \mathbb{E}[\mathbb{E}[X \mid Y]]^2 \\ &= \mathbb{E}[\mathbb{E}[X \mid Y]^2] - \mathbb{E}[X]^2. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{E}[\text{Var}(X \mid Y)] + \text{Var}(\mathbb{E}[X \mid Y]) &= \mathbb{E}[X^2] - \mathbb{E}[\mathbb{E}[X \mid Y]^2] + \mathbb{E}[\mathbb{E}[X \mid Y]^2] - \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\ &= \text{Var}(X). \end{aligned}$$

Example 2.9. Suppose that $\{X_i\}_{i=1}^{\infty}$ is an iid sequence of rvs with common mean μ and common variance σ^2 . Let N be a discrete, non-negative integer-valued rv that is independent of each X_i . Find the mean and variance of $T = \sum_{i=1}^N X_i$ (referred to as a *random sum*).

Solution: By the law of total expectation,

$$E[T] = E[E[T | N]].$$

Note that

$$\begin{aligned} E[T | N = n] &= E\left[\sum_{i=1}^N X_i \mid N = n\right] \\ &= E\left[\sum_{i=1}^n X_i \mid N = n\right] \\ &= \sum_{i=1}^n E[X_i | N = n] \\ &= \sum_{i=1}^n E[X_i] && \text{since } N \text{ is independent of } \{X_i\}_{i=1}^{\infty} \\ &= n\mu. \end{aligned}$$

Thus,

$$E[T | N] = E[T | N = n]_{n=N} = N\mu,$$

and so $E[T] = E[N\mu] = \mu E[N]$. To calculate $\text{Var}(T)$, we employ Theorem 2.3 to obtain

$$\begin{aligned} \text{Var}(T) &= E[\text{Var}(T | N)] + \text{Var}(E[T | N]) \\ &= E[\text{Var}(T | N)] + \text{Var}(N\mu) \\ &= E[\text{Var}(T | N)] + \mu^2 \text{Var}(N). \end{aligned}$$

Now,

$$\begin{aligned} \text{Var}(T | N = n) &= \text{Var}\left(\sum_{i=1}^N X_i \mid N = n\right) \\ &= \text{Var}\left(\sum_{i=1}^n X_i \mid N = n\right) \\ &= \text{Var}\left(\sum_{i=1}^n X_i\right) && \text{since } N \text{ is independent of } \{X_i\}_{i=1}^{\infty} \\ &= \sum_{i=1}^n \text{Var}(X_i) \\ &= n\sigma^2. \end{aligned}$$

Thus, $\text{Var}(T | N) = \text{Var}(T | N = n)_{N=n} = N\sigma^2$. Finally,

$$\begin{aligned} \text{Var}(T) &= E[N\sigma^2] + \mu^2 \text{Var}(N) \\ &= \sigma^2 E[N] + \mu^2 \text{Var}(N). \end{aligned}$$

2.3 Computing Probabilities by Conditioning

For any two rvs, recall that

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X | Y]] = \begin{cases} \sum_y \mathbb{E}[X | Y = y] p_Y(y) & , \text{ if } Y \text{ is discrete,} \\ \int_{-\infty}^{\infty} \mathbb{E}[X | Y = y] f_Y(y) dy & , \text{ if } Y \text{ is continuous.} \end{cases} \quad (2.1)$$

Now suppose that A represents some event of interest, and we wish to determine $\mathbb{P}(A)$. Define an indicator rv X such that

$$X = \begin{cases} 0 & , \text{ if event } A^c \text{ occurs,} \\ 1 & , \text{ if event } A \text{ occurs.} \end{cases}$$

Clearly, $\mathbb{P}(X = 1) = \mathbb{P}(A)$ and $\mathbb{P}(X = 0) = 1 - \mathbb{P}(A)$, so that $X \sim \text{BERN}(\mathbb{P}(A))$. Thus,

$$\begin{aligned} \mathbb{E}[X | Y = y] &= \sum_x x \mathbb{P}(X = x | Y = y) \\ &= 0 \mathbb{P}(X = 0 | Y = y) + 1 \mathbb{P}(X = 1 | Y = y) \\ &= \mathbb{P}(X = 1 | Y = y) \\ &= \mathbb{P}(A | Y = y). \end{aligned}$$

Therefore, (2.1) becomes

$$\mathbb{P}(A) = \begin{cases} \sum_y \mathbb{P}(A | Y = y) p_Y(y) & , \text{ if } Y \text{ is discrete,} \\ \int_{-\infty}^{\infty} \mathbb{P}(A | Y = y) f_Y(y) dy & , \text{ if } Y \text{ is continuous,} \end{cases} \quad (2.2)$$

which are analogues of the law of total probability. In other words, the expectation formula (2.1) can also be used to calculate probabilities of interest as indicated by (2.2).

Example 2.10. Suppose that X and Y are independent continuous rvs. Find an expression for $\mathbb{P}(X < Y)$.

Solution: With the event defined as $A = \{X < Y\}$, we apply (2.2) to get

$$\begin{aligned} \mathbb{P}(X < Y) &= \mathbb{P}(A) \\ &= \int_{-\infty}^{\infty} \mathbb{P}(A | Y = y) f_Y(y) dy \\ &= \int_{-\infty}^{\infty} \mathbb{P}(X < Y | Y = y) f_Y(y) dy \\ &= \int_{-\infty}^{\infty} \mathbb{P}(X < y | Y = y) f_Y(y) dy \\ &= \int_{-\infty}^{\infty} \mathbb{P}(X < y) f_Y(y) dy && \text{since } X \text{ and } Y \text{ are independent rvs} \\ &= \int_{-\infty}^{\infty} \mathbb{P}(X \leq y) f_Y(y) dy && \text{since } X \text{ is a continuous rv} \\ &= \int_{-\infty}^{\infty} F_X(y) f_Y(y) dy \end{aligned} \quad (2.3)$$

Remark: If, in addition, X and Y are identically distributed, then the pdf $f_Y(y)$ is equal to $f_X(y)$ and the

result of Example 2.10 simplifies to become

$$\begin{aligned}
 \mathbb{P}(X < Y) &= \int_{-\infty}^{\infty} F_X(y) f_X(y) dy \\
 &= \int_0^1 u du && \text{where } u = F_X(y) \implies \frac{du}{dy} = f_X(y) \implies du = f_X(y) dy \\
 &= \left[\frac{u^2}{2} \right]_{u=0}^{u=1} \\
 &= \frac{1}{2},
 \end{aligned}$$

as one would expect.

Example 2.11. Suppose that $X \sim \text{EXP}(\lambda_1)$ and $Y \sim \text{EXP}(\lambda_2)$ are independent exponential rvs. Show that

$$\mathbb{P}(X < Y) = \frac{\lambda_1}{\lambda_1 + \lambda_2}.$$

Solution: Since X and Y are both exponential rvs, it immediately follows that

$$\begin{aligned}
 f_Y(y) &= \lambda_2 e^{-\lambda_2 y}, \quad y > 0, \\
 F_X(y) &= \int_0^y \lambda_1 e^{-\lambda_1 x} dx \\
 &= \lambda_1 \left[-\frac{1}{\lambda_1} e^{-\lambda_1 x} \right]_{x=0}^{x=y} \\
 &= 1 - e^{-\lambda_1 y}, \quad y \geq 0.
 \end{aligned}$$

Therefore, (2.3) becomes

$$\begin{aligned}
 \mathbb{P}(X < Y) &= \int_0^{\infty} (1 - e^{-\lambda_1 y}) \lambda_2 e^{-\lambda_2 y} dy \\
 &= \int_0^{\infty} \lambda_2 e^{-\lambda_2 y} dy - \lambda_2 \int_0^{\infty} e^{-(\lambda_1 + \lambda_2)y} dy \\
 &= 1 - \frac{\lambda_2}{(\lambda_1 + \lambda_2)} \int_0^{\infty} (\lambda_1 + \lambda_2) e^{-(\lambda_1 + \lambda_2)y} dy \\
 &= 1 - \frac{\lambda_2}{\lambda_1 + \lambda_2} \\
 &= \frac{\lambda_1}{\lambda_1 + \lambda_2}.
 \end{aligned}$$

Remark: As a matter of interest, this particular result will be featured quite prominently in Chapter 4.

Example 2.12. Suppose W , X , and Y are independent continuous rvs on $(0, \infty)$. If $Z = X \mid (X < Y)$, then show that $(W, X) \mid (W < X < Y)$ and $(W, Z) \mid (W < Z)$ are identically distributed.

Solution: Let us first consider the joint conditional cdf of $(W, X) \mid (W < X < Y)$:

$$\begin{aligned} G(w, x) &= \mathbb{P}(W \leq w, X \leq x \mid W < X < Y) \\ &= \frac{\mathbb{P}(W \leq w, X \leq x, W < X < Y)}{\mathbb{P}(W < X < Y)} \\ &= \frac{\mathbb{P}(W \leq w, X \leq x, W < X, X < Y)}{\mathbb{P}(W < X, X < Y)}, \quad w, x \geq 0. \end{aligned}$$

Conditioning on the rv X and noting that W , X , and Y are independent rvs, it follows that

$$\begin{aligned} \mathbb{P}(W < X, X < Y) &= \int_0^\infty \mathbb{P}(W < X, X < Y \mid X = s) f_X(s) ds \\ &= \int_0^\infty \mathbb{P}(W < s, Y > s \mid X = s) f_X(s) ds \\ &= \int_0^\infty \mathbb{P}(W < s, Y > s) f_X(s) ds \\ &= \int_0^\infty \mathbb{P}(W < s) \mathbb{P}(Y > s) f_X(s) ds \end{aligned} \tag{2.4}$$

and

$$\begin{aligned} \mathbb{P}(W \leq w, X \leq x, W < X, X < Y) &= \int_0^\infty \mathbb{P}(W \leq w, X \leq x, W < X, X < Y \mid X = s) f_X(s) ds \\ &= \int_0^\infty \mathbb{P}(W \leq w, s \leq x, W < s, Y > s \mid X = s) f_X(s) ds \\ &= \int_0^\infty \mathbb{P}(W \leq w, s \leq x, W < s, Y > s) f_X(s) ds \\ &= \int_0^x \mathbb{P}(W \leq w, W < s, Y > s) f_X(s) ds \\ &= \int_0^x \mathbb{P}(W \leq \min\{w, s\}, Y > s) f_X(s) ds \\ &= \int_0^x \mathbb{P}(W \leq \min\{w, s\}) \mathbb{P}(Y > s) f_X(s) ds \end{aligned} \tag{2.5}$$

Next, consider the conditional rv $Z = X \mid (X < Y)$.

$$\begin{aligned} \mathbb{P}(Z \leq z) &= \mathbb{P}(X \leq z \mid X < Y) \\ &= \frac{\mathbb{P}(X \leq z, X < Y)}{\mathbb{P}(X < Y)} \\ &= \frac{\int_0^\infty \mathbb{P}(X \leq z, X < Y \mid X = s) f_X(s) ds}{\mathbb{P}(X < Y)} \\ &= \frac{\int_0^\infty \mathbb{P}(s \leq z, s < Y \mid X = s) f_X(s) ds}{\mathbb{P}(X < Y)} \\ &= \frac{\int_0^\infty \mathbb{P}(s \leq z, s < Y) f_X(s) ds}{\mathbb{P}(X < Y)} \\ &= \frac{\int_0^z \mathbb{P}(Y > s) f_X(s) ds}{\mathbb{P}(X < Y)} \end{aligned}$$

and so the pdf of Z is given by

$$\begin{aligned} h_Z(z) &= \frac{d}{dz} \mathbb{P}(Z \leq z) \\ &= \frac{\frac{d}{dz} \int_0^z \mathbb{P}(Y > s) f_X(s) ds}{\mathbb{P}(X < Y)} \\ &= \frac{\mathbb{P}(Y > z) f_X(z)}{\mathbb{P}(X < Y)}, \quad z > 0. \end{aligned}$$

Now, the joint conditional cdf of $(W, Z) \mid (W < Z)$ is given by

$$\mathbb{P}(W \leq w, Z \leq z \mid W < Z) = \frac{\mathbb{P}(W \leq w, Z \leq z, W < Z)}{\mathbb{P}(W < Z)}, \quad w, z \geq 0$$

Due to the independence of W with X and Y ,

$$\begin{aligned} \mathbb{P}(W < Z) &= \int_0^\infty \mathbb{P}(W < Z \mid Z = s) h_Z(s) ds \\ &= \int_0^\infty \mathbb{P}(W < s \mid Z = s) h_Z(s) ds \\ &= \int_0^\infty \mathbb{P}(W < s) h_Z(s) ds \\ &= \int_0^\infty \mathbb{P}(W < s) \frac{\mathbb{P}(Y > s) f_X(s)}{\mathbb{P}(X < Y)} ds \\ &= \frac{\mathbb{P}(W < X, X < Y)}{\mathbb{P}(X < Y)} \quad \text{from (2.4)} \end{aligned}$$

Next,

$$\begin{aligned} \mathbb{P}(W \leq w, Z \leq z, W < Z) &= \int_0^\infty \mathbb{P}(W \leq w, Z \leq z, W < Z \mid Z = s) h_Z(s) ds \\ &= \int_0^\infty \mathbb{P}(W \leq w, s \leq z, W < s) h_Z(s) ds \\ &= \int_0^z \mathbb{P}(W \leq w, W < s) h_Z(s) ds \\ &= \int_0^z \mathbb{P}(W \leq \min\{w, s\}) \frac{\mathbb{P}(Y > s) f_X(s)}{\mathbb{P}(X < Y)} ds \\ &= \mathbb{P}(W \leq w, X \leq z, W < X, X < Y) \quad \text{from (2.5)} \end{aligned}$$

Therefore, we ultimately obtain:

$$\mathbb{P}(W \leq w, Z \leq z, W < Z) = \frac{\mathbb{P}(W \leq w, X \leq z, W < X, X < Y)}{\mathbb{P}(W < X, X < Y)} = G(w, z), \quad w, z \geq 0.$$

This implies that

$$(W, X) \mid (W < X < Y) \sim (W, Z) \mid (W < Z).$$

Remark: It can likewise be shown that if $V = X \mid (W < X)$, then $(X, Y) \mid (W < X < Y)$ and $(V, Y) \mid (V < Y)$ are identically distributed (left as an upcoming exercise).

2.4 Some Further Extensions

If you consider our treatment of the conditional expectation $E[X | Y = y]$, then one detail you should notice is that this kind of expectation behaves *exactly* the same as the regular (i.e., unconditional) expectation *except* that all pmfs/pdfs used now are conditional on the event $Y = y$. In this sense, conditional expectations essentially satisfy all the properties of regular expectation. Thus, for an arbitrary real-valued function $g(\cdot)$, a corresponding analogue of

$$E[g(X)] = \begin{cases} \sum_w E[g(X) | W = w] p_W(w) & , \text{ if } W \text{ is discrete,} \\ \int_{-\infty}^{\infty} E[g(X) | W = w] f_W(w) dw & , \text{ if } W \text{ is continuous,} \end{cases}$$

would be

$$E[g(X) | Y = y] = \begin{cases} \sum_w E[g(X) | W = w, Y = y] p_{W|Y}(w | y) & , \text{ if } W \text{ is discrete,} \\ \int_{-\infty}^{\infty} E[g(X) | W = w, Y = y] f_{W|Y}(w | y) dw & , \text{ if } W \text{ is continuous.} \end{cases}$$

We remark that the above relation makes sense, since if we assume (without loss of generality) that X and Y are discrete rvs, then we obtain (in the case when W is discrete too):

$$\begin{aligned} \sum_w E[g(X) | W = w, Y = y] p_{W|Y}(w | y) &= \sum_w \sum_x g(x) p_{X|WY}(x | w, y) p_{W|Y}(w | y) \\ &= \sum_w \sum_x g(x) \frac{p_{XWY}(x, w, y)}{p_{WY}(w, y)} \frac{p_{WY}(w, y)}{p_Y(y)} \\ &= \sum_x \frac{g(x)}{p_Y(y)} \sum_w p_{XWY}(x, w, y) \\ &= \sum_x g(x) \frac{p_{XY}(x, y)}{p_Y(y)} \\ &= \sum_x g(x) p_{X|Y}(x, y) \\ &= E[g(X) | Y = y]. \end{aligned}$$

Similarly, if one introduces an event of interest A and defines

$$g(X) = \begin{cases} 0 & , \text{ if event } A^c \text{ occurs,} \\ 1 & , \text{ if event } A \text{ occurs,} \end{cases}$$

then we obtain

$$E[A | Y = y] = \begin{cases} \sum_w E[A | W = w, Y = y] p_{W|Y}(w | y) & , \text{ if } W \text{ is discrete,} \\ \int_{-\infty}^{\infty} E[A | W = w, Y = y] f_{W|Y}(w | y) dw & , \text{ if } W \text{ is continuous.} \end{cases}$$

Furthermore, if we now define

$$E[g(X) | W, Y] = E[g(X) | W = w, Y = y] \Big|_{w=W, y=Y},$$

then the law of total expectation extends to become

$$E[g(X)] = E[E[g(X) | Y]] = E[E[E[g(X) | W, Y] | Y]].$$

Example 2.13. Consider an experiment in which independent trials, each having success probability $p \in (0, 1)$, are performed until k consecutive successes are achieved where $k \in \mathbb{Z}^+$. Determine the

expected number of trials needed to achieve k consecutive successes.

Solution: Let N_k represent the number of trials needed to get k consecutive successes. We wish to determine $E[N_k]$. For $k = 1$, note that $N_1 \sim \text{GEO}_t(p)$, therefore $E[N_k] = \frac{1}{p}$. For arbitrary $k \geq 2$, let us consider conditioning on the outcome of the first trial, represented by W , such that

$$W = \begin{cases} 0 & , \text{ if first trial is a failure,} \\ 1 & , \text{ if first trial is a success.} \end{cases}$$

Thus,

$$\begin{aligned} E[N_k] &= E[E[N_k | W]] \\ &= \mathbb{P}(W = 0) E[N_k | W = 0] + \mathbb{P}(W = 1) E[N_k | W = 1] \\ &= (1 - p) E[N_k | W = 0] + p E[N_k | W = 1] \end{aligned}$$

Now, it is clear $N_k | (W = 0) \sim 1 + N_k$, but unfortunately we do not have a nice corresponding result for $N_k | (W = 1)$. It does not hold true that $N_k | (W = 0) \sim 1 + N_{k-1}$. What else can we try?

Idea: Let's try $E[N_k] = E[E[N_k | N_{k-1}]]$, i.e., to get k in a row, we must first get $k - 1$ in a row. Define

$$Y | (N_{k-1} = n) = \begin{cases} 0 & , \text{ if } (n + 1)^{\text{th}} \text{ trial is a failure,} \\ 1 & , \text{ if } (n + 1)^{\text{th}} \text{ trial is a success.} \end{cases}$$

By independence of the trials,

$$\begin{aligned} \mathbb{P}(Y = 0 | N_{k-1} = n) &= 1 - p, \\ \mathbb{P}(Y = 1 | N_{k-1} = n) &= p. \end{aligned}$$

As a result, we get:

$$\begin{aligned} E[N_k | N_{k-1} = n] &= \sum_{y=0}^1 E[N_k | N_{k-1} = n, Y = y] \mathbb{P}(Y = y | N_{k-1} = n) \\ &= (1 - p) E[N_k | N_{k-1} = n, Y = 0] + p E[N_k | N_{k-1} = n, Y = 1]. \end{aligned}$$

Note that $N_k | (N_{k-1} = n, Y = 0) \sim n + 1 + N_k$ (i.e., given that we know it took n trials to get $k - 1$ consecutive successes, and then on the next trial we got a failure, what happens?). Also, $N_k | (N_{k-1} = n, Y = 1)$ is equal to $n + 1$ with probability 1. Therefore,

$$\begin{aligned} E[N_k | N_{k-1} = n] &= (1 - p)(n + 1 + E[N_k]) + p(n + 1) \\ &= n + 1 + (1 - p) E[N_k]. \end{aligned}$$

Therefore,

$$E[N_k | N_{k-1}] = E[N_k | N_{k-1} = n] \Big|_{n=N_{k-1}} = N_{k-1} + 1 + (1 - p) E[N_k].$$

Now, our whole idea was to apply $E[N_k] = E[E[N_k | N_{k-1}]]$, and now we have the inner piece, so

$$\begin{aligned} E[N_k] &= E[N_{k-1} + 1 + (1 - p) E[N_k]] \\ &= E[N_{k-1}] + 1 + (1 - p) E[N_k] \end{aligned}$$

Therefore,

$$(1 - (1 - p)) E[N_k] = 1 + E[N_{k-1}] \implies E[N_k] = \frac{1}{p} + \frac{E[N_{k-1}]}{p}, \quad k \geq 2,$$

which is a recursive equation for $E[N_k]$. Take $k = 2$:

$$E[N_2] = \frac{1}{p} + \frac{E[N_1]}{p} = \frac{1}{p} + \frac{(1/p)}{p} = \frac{1}{p} + \frac{1}{p^2}.$$

Take $k = 3$:

$$E[N_3] = \frac{1}{p} + \frac{E[N_2]}{p} = \frac{1}{p} + \frac{1}{p^2} + \frac{1}{p^3}.$$

Take $k = 4$:

$$E[N_4] = \frac{1}{p} + \frac{E[N_3]}{p} = \frac{1}{p} + \frac{1}{p^2} + \frac{1}{p^3} + \frac{1}{p^4}.$$

Continuing inductively, we actually have

$$E[N_k] = \frac{1}{p} + \frac{1}{p^2} + \cdots + \frac{1}{p^k},$$

which is a finite geometric series, therefore,

$$E[N_k] = \frac{(1/p) - (1/p^{k+1})}{1 - (1/p)} = \frac{p^{-k} - 1}{1 - p}, \quad k \geq 2.$$

Actually, this holds true for $k \in \mathbb{Z}^+$ (try it).