

Aug 24st 2021 ~~Chapter~~ Chapter 1.

**Definition 1.1.1** The set,  $S$ , of all possible outcomes of a particular experiment is called the *sample space* for the experiment.

Ex. flip a coin :  $S = \{H, T\}$

**Definition 1.1.2** An event is any collection of possible outcomes of an experiment, that is, any subset of  $S$  (including  $S$  itself).

**Theorem 1.1.4** For any three events,  $A$ ,  $B$ , and  $C$ , defined on a sample space  $S$ ,

- a. Commutativity  $A \cup B = B \cup A$ ,  
 $A \cap B = B \cap A$ ;
- b. Associativity  $A \cup (B \cup C) = (A \cup B) \cup C$ ,  
 $A \cap (B \cap C) = (A \cap B) \cap C$ ;
- c. Distributive Laws  $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$ ,  
 $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ ;
- d. DeMorgan's Laws  $(A \cup B)^c = A^c \cap B^c$ ,  
 $(A \cap B)^c = A^c \cup B^c$ .

**Definition 1.1.5** Two events  $A$  and  $B$  are disjoint (or *mutually exclusive*) if  $A \cap B = \emptyset$ . The events  $A_1, A_2, \dots$  are pairwise disjoint (or *mutually exclusive*) if  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ .

**Definition 1.1.6** If  $A_1, A_2, \dots$  are pairwise disjoint and  $\cup_{i=1}^{\infty} A_i = S$ , then the collection  $A_1, A_2, \dots$  forms a partition of  $S$ .

**Definition 1.2.1** A collection of subsets of  $S$  is called a sigma algebra (or *Borel field*), denoted by  $\mathcal{B}$ , if it satisfies the following three properties:

- a.  $\emptyset \in \mathcal{B}$  (the empty set is an element of  $\mathcal{B}$ ).
- b. If  $A \in \mathcal{B}$ , then  $A^c \in \mathcal{B}$  ( $\mathcal{B}$  is closed under complementation).
- c. If  $A_1, A_2, \dots \in \mathcal{B}$ , then  $\cup_{i=1}^{\infty} A_i \in \mathcal{B}$  ( $\mathcal{B}$  is closed under countable unions).

Ex.  $\mathcal{B} = \{\emptyset, S\}$  trivial algebra

$\mathcal{B} = \{\text{all subsets of } S, \text{ including } S \text{ itself}\}$ .  $|\mathcal{B}| = 2^n$  (if  $|S|=n$ )

**Example 1.2.3 (Sigma algebra-II)** Let  $S = (-\infty, \infty)$ , the real line. Then  $\mathcal{B}$  is chosen to contain all sets of the form

$$[a, b], \quad (a, b], \quad (a, b), \quad \text{and} \quad [a, b)$$

for all real numbers  $a$  and  $b$ . Also, from the properties of  $\mathcal{B}$ , it follows that  $\mathcal{B}$  contains all sets that can be formed by taking (possibly countably infinite) unions and intersections of sets of the above varieties. ||

The empty set  $\emptyset$  is a subset of any set. Thus,  $\emptyset \subset S$ . Property (a) states that this subset is always in a sigma algebra. Since  $S = \emptyset^c$ , properties (a) and (b) imply that  $S$  is always in  $\mathcal{B}$  also. In addition, from DeMorgan's Laws it follows that  $\mathcal{B}$  is closed under countable intersections. If  $A_1, A_2, \dots \in \mathcal{B}$ , then  $A_1^c, A_2^c, \dots \in \mathcal{B}$  by property (b), and therefore  $\bigcap_{i=1}^{\infty} A_i^c \in \mathcal{B}$ . However, using DeMorgan's Law (as in Exercise 1.9), we have

$$(1.2.1) \quad \left( \bigcup_{i=1}^{\infty} A_i^c \right)^c = \bigcap_{i=1}^{\infty} A_i.$$

Thus, again by property (b),  $\bigcap_{i=1}^{\infty} A_i \in \mathcal{B}$ . → closed under countable intersections

**Definition 1.2.4** Given a sample space  $S$  and an associated sigma algebra  $\mathcal{B}$ , a *probability function* is a function  $P$  with domain  $\mathcal{B}$  that satisfies

1.  $P(A) \geq 0$  for all  $A \in \mathcal{B}$ .
2.  $P(S) = 1$ .
3. If  $A_1, A_2, \dots \in \mathcal{B}$  are pairwise disjoint, then  $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$ .

$$P: \mathcal{B} \rightarrow [0, 1]$$

Prop 1:  $P(\emptyset) = 0$

Pf: Let  $A_i = \emptyset$ , then  $\bigcup_{i=1}^{\infty} A_i = \emptyset$

$$P(\emptyset) = P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i) = \sum_{i=1}^{\infty} P(\emptyset) \Rightarrow P(\emptyset) = 0.$$

2. Axiom of Finite Additivity: If  $A \in \mathcal{B}$  and  $B \in \mathcal{B}$  are disjoint, then

$$P(A \cup B) = P(A) + P(B).$$

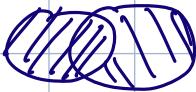
Pf: Let  $A_1 = A$ ,  $A_2 = B$ ,  $A_i = \emptyset$ ,  $\forall i \geq 3$

$$P(A \cup B) = P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i) = P(A) + P(B) + \underbrace{\sum_{i=3}^{\infty} P(\emptyset)}_{=0} = P(A) + P(B)$$

3.  $P(A^c) = 1 - P(A)$ .

Pf:  $1 = P(S) = P(A \cup A^c) = P(A) + P(A^c)$

4.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ ;



Pf:  $P(A \cup B) = P(A) + P(B \cap A^c) = P(A) + P(B) - P(A \cap B)$

$$P(B) = P(\{B \cap A\} \cup \{B \cap A^c\}) = P(B \cap A) + P(B \cap A^c),$$

5. If  $A \subset B$ , then  $P(A) \leq P(B)$

Pf:  $0 \leq P(B \cap A^c) = P(B) - P(A)$

### 6/7. Theorem 1.2.11 If $P$ is a probability function, then

- a.  $P(A) = \sum_{i=1}^{\infty} P(A \cap C_i)$  for any partition  $C_1, C_2, \dots$ ;
- b.  $P(\bigcup_{i=1}^{\infty} A_i) \leq \sum_{i=1}^{\infty} P(A_i)$  for any sets  $A_1, A_2, \dots$ . (Boole's Inequality)

(a) Pf.  $A = A \cap S = A \cap \left( \bigcup_{i=1}^{\infty} C_i \right) = \bigcup_{i=1}^{\infty} (A \cap C_i) = P \left( \bigcup_{i=1}^{\infty} (A \cap C_i) \right) = \sum_{i=1}^{\infty} P(A \cap C_i)$

(b) Pf. To establish (b) we first construct a disjoint collection  $A_1^*, A_2^*, \dots$ , with the property that  $\bigcup_{i=1}^{\infty} A_i^* = \bigcup_{i=1}^{\infty} A_i$ . We define  $A_i^*$  by

Need to check:

$$\textcircled{1} A_i^* \subseteq A_i$$

$$\textcircled{2} \bigcup_{i=1}^{\infty} A_i^* = \bigcup_{i=1}^{\infty} A_i$$

$\textcircled{3} A_i^*$  pairwise disjoint

$$A_1^* = A_1, \quad A_i^* = A_i \setminus \left( \bigcup_{j=1}^{i-1} A_j \right), \quad i = 2, 3, \dots,$$

$$P \left( \bigcup_{i=1}^{\infty} A_i \right) = P \left( \bigcup_{i=1}^{\infty} A_i^* \right) = \sum_{i=1}^{\infty} P(A_i^*) \leq \sum_{i=1}^{\infty} P(A_i)$$

**Theorem 1.2.6** Let  $S = \{s_1, \dots, s_n\}$  be a finite set. Let  $\mathcal{B}$  be any sigma algebra of subsets of  $S$ . Let  $p_1, \dots, p_n$  be nonnegative numbers that sum to 1. For any  $A \in \mathcal{B}$ , define  $P(A)$  by

$$P(A) = \sum_{\{i: s_i \in A\}} p_i.$$

(The sum over an empty set is defined to be 0.) Then  $P$  is a probability function on  $\mathcal{B}$ . This remains true if  $S = \{s_1, s_2, \dots\}$  is a countable set.

**Proof:** We will give the proof for finite  $S$ . For any  $A \in \mathcal{B}$ ,  $P(A) = \sum_{\{i: s_i \in A\}} p_i \geq 0$ , because every  $p_i \geq 0$ . Thus, Axiom 1 is true. Now,

$$P(S) = \sum_{\{i: s_i \in S\}} p_i = \sum_{i=1}^n p_i = 1.$$

Thus, Axiom 2 is true. Let  $A_1, \dots, A_k$  denote pairwise disjoint events. ( $\mathcal{B}$  contains only a finite number of sets, so we need consider only finite disjoint unions.) Then,

$$P \left( \bigcup_{i=1}^k A_i \right) = \sum_{\{j: s_j \in \bigcup_{i=1}^k A_i\}} p_j \stackrel{\text{Ai disjoint}}{=} \sum_{i=1}^k \sum_{\{j: s_j \in A_i\}} p_j = \sum_{i=1}^k P(A_i).$$

The first and third equalities are true by the definition of  $P(A)$ . The disjointedness of the  $A_i$ s ensures that the second equality is true, because the same  $p_j$ s appear exactly once on each side of the equality. Thus, Axiom 3 is true and Kolmogorov's Axioms are satisfied.  $\square$

**Theorem 1.2.14** If a job consists of  $k$  separate tasks, the  $i$ th of which can be done in  $n_i$  ways,  $i = 1, \dots, k$ , then the entire job can be done in  $n_1 \times n_2 \times \dots \times n_k$  ways.

Table 1.2.1. Number of possible arrangements of size  $r$  from  $n$  objects

	Without replacement	With replacement
Ordered	$\frac{n!}{(n-r)!}$	$n^r$
Unordered	$\binom{n}{r}$	$\binom{n+r-1}{r}$

△ Unordered + W.R.

$X_i = \# \text{times object } i \text{ was selected}$

$$\{(X_1, X_2, \dots, X_n) : X_i \in \mathbb{Z}_+^*, \sum_{i=1}^n X_i = r\}$$

$\rightarrow n$  non-negative integers sum to  $r$ , so  $\binom{n+r-1}{r}$

Aug 26st 2021

Very easy. No need to review this example.

**Example 1.2.18 (Poker)** Consider choosing a five-card poker hand from a standard deck of 52 playing cards. Obviously, we are sampling without replacement from the deck. But to specify the possible outcomes (possible hands), we must decide whether we think of the hand as being dealt sequentially (ordered) or all at once (unordered). If we wish to calculate probabilities for events that depend on the order, such as the probability of an ace in the first two cards, then we must use the ordered outcomes. But if our events do not depend on the order, we can use the unordered outcomes. For this example we use the unordered outcomes, so the sample space consists of all the five-card hands that can be chosen from the 52-card deck. There are  $\binom{52}{5} = 2,598,960$  possible hands. If the deck is well shuffled and the cards are randomly dealt, it is reasonable to assign probability  $1/2,598,960$  to each possible hand.

We now calculate some probabilities by counting outcomes in events. What is the probability of having four aces? How many different hands are there with four aces? If we specify that four of the cards are aces, then there are 48 different ways of specifying the fifth card. Thus,

$$P(\text{four aces}) = \frac{48}{2,598,960},$$

less than 1 chance in 50,000. Only slightly more complicated counting, using Theorem 1.2.14, allows us to calculate the probability of having four of a kind. There are 13

ways to specify which denomination there will be four of. After we specify these four cards, there are 48 ways of specifying the fifth. Thus, the total number of hands with four of a kind is  $(13)(48)$  and

$$P(\text{four of a kind}) = \frac{(13)(48)}{2,598,960} = \frac{624}{2,598,960}.$$

To calculate the probability of exactly one pair (not two pairs, no three of a kind, etc.) we combine some of the counting techniques. The number of hands with exactly one pair is

$$(1.2.11) \quad 13 \binom{4}{2} \binom{12}{3} 4^3 = 1,098,240.$$

Expression (1.2.11) comes from Theorem 1.2.14 because

$13 = \# \text{ of ways to specify the denomination for the pair,}$

$\binom{4}{2} = \# \text{ of ways to specify the two cards from that denomination,}$

$\binom{12}{3} = \# \text{ of ways of specifying the other three denominations,}$

$4^3 = \# \text{ of ways of specifying the other three cards from those denominations.}$

Thus,

$$P(\text{exactly one pair}) = \frac{1,098,240}{2,598,960}.$$

Ex.  $\binom{n}{r} = \binom{n-1}{r} + \binom{n-1}{r-1}$

$\boxed{\bullet \circ \dots \circ}$	$\boxed{\bullet \circ \dots \circ}$
object 1 isn't picked	object 1 is picked
$= \binom{n-1}{r}$	$= \binom{n-1}{r-1}$

the *inclusion-exclusion identity* says that

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = P_1 - P_2 + P_3 - P_4 + \dots \pm P_n.$$

where

$$P_1 = \sum_{i=1}^n P(A_i)$$

$$P_2 = \sum_{1 \leq i < j \leq n} P(A_i \cap A_j)$$

$$P_3 = \sum_{1 \leq i < j < k \leq n} P(A_i \cap A_j \cap A_k)$$

:

$$P_n = P(A_1 \cap A_2 \cap \dots \cap A_n).$$

$(-)^{n+1}$

$$\begin{aligned}
\text{Pf. } P(A_1 \cup \dots \cup A_{n+1}) &= P(A_1 \cup \dots \cup A_n) + P(A_{n+1}) - P(A_1 \cup \dots \cup A_n) \cap A_{n+1} \\
&= P(A_1 \cup \dots \cup A_n) + P(A_{n+1}) - P((A_1 \cap A_{n+1}) \cup \dots \cup (A_n \cap A_{n+1})) \\
&= \sum_{i=1}^n P(A_i) - \sum_{1 \leq i_1 < i_2 \leq n} P(A_{i_1} \cap A_{i_2}) + \dots + (-1)^{n+1} P(A_1 \cap \dots \cap A_n) \\
&\quad + P(A_{n+1}) - \sum_{i=1}^n P(A_i \cap A_{n+1}) + \sum_{1 \leq i_1 < i_2 \leq n} P(A_{i_1} \cap A_{i_2} \cap A_{n+1}) \\
&\quad - \dots - (-1)^n \sum_{1 \leq i_1 < \dots < i_{n+1} \leq n} P(A_{i_1} \cap \dots \cap A_{i_{n+1}} \cap A_{n+1}) \\
&\quad - (-1)^{n+1} P(A_1 \cap A_2 \cap \dots \cap A_{n+1}) \\
&= P_1 - P_2 + P_3 - \dots + (-1)^{n+2} P_{n+1}
\end{aligned}$$

Ex.  $N$  men repick their hats randomly  $\rightarrow A$ .

- Calculate the prob. of at least one man selects his own hat

$A_i = i$ th man selects his own hat,  $1 \leq i \leq N$

$$A = \bigcup_{i=1}^N A_i$$

$$P(A) = P\left(\bigcup_{i=1}^N A_i\right) = \sum_{i=1}^N P(A_i) - \sum_{1 \leq i < j \leq N} P(A_i \cap A_j) + \dots$$

$$|S| = N! \quad P(A_i) = \frac{(N-1)!}{N!} = \frac{1}{N}$$

$$i \neq j: \quad P(A_i \cap A_j) = \frac{(N-2)!}{N!} = \frac{1}{N(N-1)}$$

:

$$P(A_{i_1} \cap \dots \cap A_{i_k}) = \frac{(N-k)!}{N!}$$

$$\begin{aligned}
\therefore P\left(\bigcup_{i=1}^N A_i\right) &= N \cdot \frac{1}{N} - \binom{N}{2} \frac{(N-2)!}{N!} + \binom{N}{3} \frac{(N-3)!}{N!} - \dots + (-1)^{N+1} \binom{N}{N} \frac{(N-N)!}{N!} \\
&= 1 - \frac{1}{2!} + \frac{1}{3!} - \dots + (-1)^{N+1} \frac{1}{N!}
\end{aligned}$$

$$(e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}) \rightarrow \text{large } N \quad 1 - e^{-1} \approx 63\%$$

Aug 31st, 2021

**Definition 1.3.2** If  $A$  and  $B$  are events in  $S$ , and  $P(B) > 0$ , then the conditional probability of  $A$  given  $B$ , written  $P(A|B)$ , is

$$(1.3.1) \quad P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

**Theorem 1.3.5 (Bayes' Rule)** Let  $A_1, A_2, \dots$  be a partition of the sample space, and let  $B$  be any set. Then, for each  $i = 1, 2, \dots$ ,

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^{\infty} P(B|A_j)P(A_j)}.$$

Multiplicative Rule:

$$P\left(\bigcap_{i=1}^n A_i\right) = P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2) \cdots P(A_n|A_1 \cap \cdots \cap A_{n-1})$$

**Example 1.3.6 (Coding)** When coded messages are sent, there are sometimes errors in transmission. In particular, Morse code uses “dots” and “dashes,” which are known to occur in the proportion of 3:4. This means that for any given symbol,

$$P(\text{dot sent}) = \frac{3}{7} \quad \text{and} \quad P(\text{dash sent}) = \frac{4}{7}.$$

Suppose there is interference on the transmission line, and with probability  $\frac{1}{8}$  a dot is mistakenly received as a dash, and vice versa. If we receive a dot, can we be sure that a dot was sent? Using Bayes’ Rule, we can write

$$P(\text{dot sent} | \text{dot received}) = P(\text{dot received} | \text{dot sent}) \frac{P(\text{dot sent})}{P(\text{dot received})}.$$

Now, from the information given, we know that  $P(\text{dot sent}) = \frac{3}{7}$  and  $P(\text{dot received} | \text{dot sent}) = \frac{7}{8}$ . Furthermore, we can also write

$$\begin{aligned} P(\text{dot received}) &= P(\text{dot received} \cap \text{dot sent}) + P(\text{dot received} \cap \text{dash sent}) \\ &= P(\text{dot received} | \text{dot sent})P(\text{dot sent}) \\ &\quad + P(\text{dot received} | \text{dash sent})P(\text{dash sent}) \\ &= \frac{7}{8} \times \frac{3}{7} + \frac{1}{8} \times \frac{4}{7} = \frac{25}{56}. \end{aligned}$$

Combining these results, we have that the probability of correctly receiving a dot is

$$P(\text{dot sent} | \text{dot received}) = \frac{(7/8) \times (3/7)}{25/56} = \frac{21}{25} = 84\%$$

**Definition 1.3.7** Two events,  $A$  and  $B$ , are statistically independent if

$$(1.3.8) \quad P(A \cap B) = P(A)P(B). \quad [P(A|B)=P(A)]$$

**Definition 1.3.12** A collection of events  $A_1, \dots, A_n$  are mutually independent if for any subcollection  $A_{i_1}, \dots, A_{i_k}$ , we have

$$P\left(\bigcap_{j=1}^k A_{i_j}\right) = \prod_{j=1}^k P(A_{i_j}).$$

**Theorem 1.3.9** If  $A$  and  $B$  are independent events, then the following pairs are also independent:

- a.  $A$  and  $B^c$ ,
- b.  $A^c$  and  $B$ ,
- c.  $A^c$  and  $B^c$ .

**Proof:** We will prove only (a), leaving the rest as Exercise 1.40. To prove (a) we must show that  $P(A \cap B^c) = P(A)P(B^c)$ . From Theorem 1.2.9a we have

$$\begin{aligned} P(A \cap B^c) &= P(A) - P(A \cap B) \\ &= P(A) - P(A)P(B) \quad (A \text{ and } B \text{ are independent}) \\ &= P(A)(1 - P(B)) \\ &= P(A)P(B^c). \end{aligned}$$

□

**Definition 1.4.1** A random variable is a function from a sample space  $S$  into the real numbers.

$$X: S \rightarrow \mathbb{R}$$

Ex: Toss a fair coin 3 times.  $X := \# \text{heads}$

$$S = \{\text{HHH}, \text{HHT}, \text{HTH}, \text{THH}, \text{HTT}, \text{THT}, \text{TTH}, \text{TTT}\}$$

$$X: \begin{matrix} \downarrow & \downarrow \\ 3 & 2 & 2 & 2 & 1 & 1 & 1 & 0 \end{matrix}$$

$$P(X=2) = P(\{\text{HHT}, \text{HTH}, \text{THH}\}) = 3/8$$

For any set  $A \in \mathcal{X}$ ,

$$P_X(X \in A) = P(\{s \in S : X(s) \in A\}) = P(X^{-1}(A))$$

$$\begin{matrix} \downarrow & \downarrow & \downarrow \\ \mathbb{R} & S & S \end{matrix}$$

events in  $S$

$P_X$  is an induced prob. func on  $\mathbb{R}$  (satisfies axioms of prob.)

**Definition 1.5.1** The cumulative distribution function or cdf of a random variable  $X$ , denoted by  $F_X(x)$ , is defined by

$$F_X(x) = P_X(X \leq x), \quad \text{for all } x.$$

**Example 1.5.2 (Tossing three coins)** Consider the experiment of tossing three fair coins, and let  $X = \text{number of heads observed}$ . The cdf of  $X$  is

$$(1.5.1) \quad F_X(x) = \begin{cases} 0 & \text{if } -\infty < x < 0 \\ \frac{1}{8} & \text{if } 0 \leq x < 1 \\ \frac{1}{2} & \text{if } 1 \leq x < 2 \\ \frac{7}{8} & \text{if } 2 \leq x < 3 \\ 1 & \text{if } 3 \leq x < \infty. \end{cases}$$

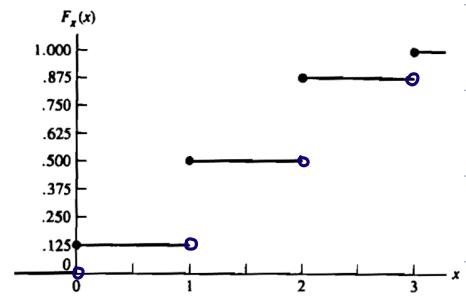


Figure 1.5.1. Cdf of Example 1.5.2

**Theorem 1.5.3** The function  $F(x)$  is a cdf if and only if the following three conditions hold:

- a.  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow \infty} F(x) = 1$ .
- b.  $F(x)$  is a nondecreasing function of  $x$ .
- c.  $F(x)$  is right-continuous; that is, for every number  $x_0$ ,  $\lim_{x \downarrow x_0} F(x) = F(x_0)$ .

Axiom of continuity:

①  $A_1, A_2, \dots$  s.t.  $A_i \subseteq A_{i+1}$ ;  $i=1, 2, \dots$ , then:

$$P(A_n) \uparrow P(A); \quad A = \bigcup_{i=1}^{+\infty} A_i$$

$$\lim_{n \rightarrow +\infty} P(A_n) = \lim_{n \rightarrow +\infty} P(\bigcup_{i=1}^n A_i) \stackrel{?}{=} P\left(\lim_{n \rightarrow +\infty} \bigcup_{i=1}^n A_i\right) = P\left(\bigcup_{i=1}^{+\infty} A_i\right)$$

②  $A_1, A_2, \dots$  s.t.  $A_{i+1} \subseteq A_i$ ;  $i=1, 2, \dots$ , then:

$$P(A_n) \downarrow P(A); \quad A = \bigcap_{i=1}^{+\infty} A_i$$

Proof of 1.5.3.

(Necessity)

$$\text{a. } \lim_{x \rightarrow +\infty} F(x) = 1 : A_n = \{X \leq x_n\} \Rightarrow A_n \uparrow \bigcup_{i=1}^{+\infty} A_i = \{X < +\infty\} = S$$

By the axiom of continuity:  $F(x_n) = P(A_n) \uparrow P(S) = 1$

$$-\lim_{x \rightarrow -\infty} F(x) = 0 : \forall \text{seq } X_n \rightarrow -\infty \text{ as } n \rightarrow +\infty$$

$$A_n = \{X \leq x_n\} \Rightarrow A_n \downarrow \bigcap_{i=1}^{+\infty} A_i = \emptyset$$

By the axiom of continuity:  $\bar{F}_X(x_n) = P(A_n) \downarrow P(\emptyset) = 0$

b.  $\bar{F}_X$  is non-decreasing:

$$\begin{aligned} x \leq y &\Rightarrow \{s : X(s) \leq x\} \subseteq \{s : X(s) \leq y\} \\ &\Rightarrow P(\underset{\substack{\downarrow \\ \bar{F}_X(x)}}{\underset{\parallel}{\cup}}) \leq P(\underset{\substack{\downarrow \\ \bar{F}_X(y)}}{\underset{\parallel}{\cup}}) \end{aligned}$$

c. Right continuous; NTS:  $\lim_{n \rightarrow +\infty} \bar{F}_X(y_n) = \bar{F}_X(x); \forall y_n \downarrow x$

$$A_n = \{X \leq y_n\} \quad A_n \downarrow \bigcap_{i=1}^n A_i = \{X \leq x\}$$

By the axiom of continuity,  $\bar{F}_X(y_n) = P(A_n) \downarrow P(\{X \leq x\}) = \bar{F}_X(x)$

-  $\bar{F}_X(\cdot)$  has left limit:

$$y_n \uparrow x \text{ as } n \rightarrow +\infty; A_n = \{X \leq y_n\} \quad A_n \uparrow \bigcup_{i=1}^{+\infty} A_i = \{X < x\}$$

By the axiom of continuity:  $\bar{F}_X(y_n) = P(A_n) \uparrow P(\{X < x\})$

**Example 1.5.4 (Tossing for a head)** Suppose we do an experiment that consists of tossing a coin until a head appears. Let  $p$  = probability of a head on any given toss, and define a random variable  $X$  = number of tosses required to get a head. Then, for any  $x = 1, 2, \dots$ ,

$$(1.5.2) \quad P(X = x) = (1 - p)^{x-1} p,$$

since we must get  $x - 1$  tails followed by a head for the event to occur and all trials are independent. From (1.5.2) we calculate, for any positive integer  $x$ ,

$$(1.5.3) \quad P(X \leq x) = \sum_{i=1}^x P(X = i) = \sum_{i=1}^x (1 - p)^{i-1} p.$$

The partial sum of the geometric series is

$$(1.5.4) \quad \sum_{k=1}^n t^{k-1} = \frac{1 - t^n}{1 - t}, \quad t \neq 1,$$

a fact that can be established by induction (see Exercise 1.50). Applying (1.5.4) to our probability, we find that the cdf of the random variable  $X$  is

$$\begin{aligned} F_X(x) &= P(X \leq x) \\ &= \frac{1 - (1 - p)^x}{1 - (1 - p)} p \\ &= 1 - (1 - p)^x, \quad x = 1, 2, \dots \end{aligned}$$

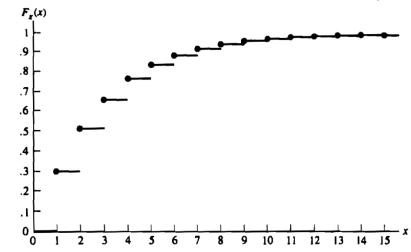


Figure 1.5.2. Geometric cdf,  $p = .3$

**Definition 1.5.7** A random variable  $X$  is *continuous* if  $F_X(x)$  is a continuous function of  $x$ . A random variable  $X$  is *discrete* if  $F_X(x)$  is a step function of  $x$ .

**Definition 1.5.8** The random variables  $X$  and  $Y$  are *identically distributed* if, for every set  $A \in \mathcal{B}^1$ ,  $P(X \in A) = P(Y \in A)$ .

**Example 1.5.9 (Identically distributed random variables)** Consider the experiment of tossing a fair coin three times as in Example 1.4.3. Define the random variables  $X$  and  $Y$  by

$$X = \text{number of heads observed} \quad \text{and} \quad Y = \text{number of tails observed.}$$

The distribution of  $X$  is given in Example 1.4.3, and it is easily verified that the distribution of  $Y$  is exactly the same. That is, for each  $k = 0, 1, 2, 3$ , we have  $P(X = k) = P(Y = k)$ . So  $X$  and  $Y$  are identically distributed. However, for no sample points do we have  $X(s) = Y(s)$ . ||

**Theorem 1.5.10** The following two statements are equivalent:

- a. The random variables  $X$  and  $Y$  are identically distributed.
- b.  $F_X(x) = F_Y(x)$  for every  $x$ .

**Definition 1.6.1** The *probability mass function (pmf)* of a discrete random variable  $X$  is given by

$$f_X(x) = P(X = x) \quad \text{for all } x.$$

**Example 1.6.2 (Geometric probabilities)** For the geometric distribution of Example 1.5.4, we have the pmf

$$f_X(x) = P(X = x) = \begin{cases} (1 - p)^{x-1} p & \text{for } x = 1, 2, \dots \\ 0 & \text{otherwise.} \end{cases}$$

Sep 2nd 2021

**Definition 1.6.3** The *probability density function or pdf*,  $f_X(x)$ , of a continuous random variable  $X$  is the function that satisfies

NOT Unique

$$(1.6.3) \quad F_X(x) = \int_{-\infty}^x f_X(t) dt \quad \text{for all } x.$$

Using the Fundamental Theorem of Calculus, if  $f_X(x)$  is continuous, we have the further relationship

$$(1.6.2) \quad \frac{d}{dx} F_X(x) = f_X(x).$$

**Theorem 1.6.5** A function  $f_X(x)$  is a pdf (or pmf) of a random variable  $X$  if and only if

- a.  $f_X(x) \geq 0$  for ~~all~~  $x$ . A.S.
- b.  $\sum_x f_X(x) = 1$  (pmf) or  $\int_{-\infty}^{\infty} f_X(x) dx = 1$  (pdf).

## Chapter 2: Transformations and Expectations

$$Y = g(X)$$

Def:  $X = \{x : f_X(x) > 0\}$  support of r.v.  $X$ .

(2.1.7)  $\mathcal{Y} = \{y : y = g(x) \text{ for some } x \in X\}$  support of r.v.  $Y$ .

- If  $g$  is monotone:

$g$  increasing:  $\{x \in X : g(x) \leq y\} = \{x \in X : x \leq g^{-1}(y)\}$

$$\Rightarrow F_Y(y) = \int_{\{x \in X | x \leq g^{-1}(y)\}} f_X(x) dx = \int_{-\infty}^{g^{-1}(y)} f_X(x) dx = F_X(g^{-1}(y))$$

$g$  decreasing:  $\{x \in X : g(x) \leq y\} = \{x \in X : x \geq g^{-1}(y)\}$

$$\Rightarrow F_Y(y) = \int_{g^{-1}(y)}^{+\infty} f_X(x) dx = 1 - F_X(g^{-1}(y))$$

**Theorem 2.1.5** Let  $X$  have pdf  $f_X(x)$  and let  $Y = g(X)$ , where  $g$  is a monotone function. Let  $\mathcal{X}$  and  $\mathcal{Y}$  be defined by (2.1.7). Suppose that  $f_X(x)$  is continuous on  $\mathcal{X}$  and that  $g^{-1}(y)$  has a continuous derivative on  $\mathcal{Y}$ . Then the pdf of  $Y$  is given by

$$(2.1.10) \quad f_Y(y) = \begin{cases} f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right| & y \in \mathcal{Y} \\ 0 & \text{otherwise.} \end{cases}$$

**Example 2.1.4 (Uniform-exponential relationship-I)** Suppose  $X \sim f_X(x) = 1$  if  $0 < x < 1$  and 0 otherwise, the uniform(0, 1) distribution. It is straightforward to check that  $F_X(x) = x$ ,  $0 < x < 1$ . We now make the transformation  $Y = g(X) = -\log X$ . Since

$$\frac{d}{dx} g(x) = \frac{d}{dx} (-\log x) = \frac{-1}{x} < 0, \quad \text{for } 0 < x < 1,$$

$g(x)$  is a decreasing function. As  $X$  ranges between 0 and 1,  $-\log x$  ranges between 0 and  $\infty$ , that is,  $\mathcal{Y} = (0, \infty)$ . For  $y > 0$ ,  $y = -\log x$  implies  $x = e^{-y}$ , so  $g^{-1}(y) = e^{-y}$ . Therefore, for  $y > 0$ ,

$$F_Y(y) = 1 - F_X(g^{-1}(y)) = 1 - F_X(e^{-y}) = 1 - e^{-y}. \quad (F_X(x) = x)$$

Of course,  $F_Y(y) = 0$  for  $y \leq 0$ . Note that it was necessary only to verify that  $g(x) = -\log x$  is monotone on  $(0, 1)$ , the support of  $X$ . ||

**Example 2.1.1 (Binomial transformation)** A discrete random variable  $X$  has a *binomial distribution* if its pmf is of the form

$$(2.1.3) \quad f_X(x) = P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n,$$

where  $n$  is a positive integer and  $0 \leq p \leq 1$ . Values such as  $n$  and  $p$  that can be set to different values, producing different probability distributions, are called *parameters*. Consider the random variable  $Y = g(X)$ , where  $g(x) = n - x$ ; that is,  $Y = n - X$ . Here  $\mathcal{X} = \{0, 1, \dots, n\}$  and  $\mathcal{Y} = \{y: y = g(x), x \in \mathcal{X}\} = \{0, 1, \dots, n\}$ . For any  $y \in \mathcal{Y}$ ,  $n - x = g(x) = y$  if and only if  $x = n - y$ . Thus,  $g^{-1}(y)$  is the single point  $x = n - y$ , and

$$\begin{aligned} f_Y(y) &= \sum_{x \in g^{-1}(y)} f_X(x) \\ &= f_X(n - y) \\ &= \binom{n}{n-y} p^{n-y} (1-p)^{n-(n-y)} \\ &= \binom{n}{y} (1-p)^y p^{n-y}. \quad \left( \begin{array}{l} \text{Definition 1.2.17} \\ \text{implies } \binom{n}{y} = \binom{n}{n-y} \end{array} \right) \end{aligned}$$

Thus, we see that  $Y$  also has a binomial distribution, but with parameters  $n$  and  $1 - p$ . ||

**Example 2.1.7 (Square transformation)** Suppose  $X$  is a continuous random variable. For  $y > 0$ , the cdf of  $Y = X^2$  is

$$F_Y(y) = P(Y \leq y) = P(X^2 \leq y) = P(-\sqrt{y} \leq X \leq \sqrt{y}).$$

Because  $x$  is continuous, we can drop the equality from the left endpoint and obtain

$$\begin{aligned} F_Y(y) &= P(-\sqrt{y} < X \leq \sqrt{y}) \\ &= P(X \leq \sqrt{y}) - P(X \leq -\sqrt{y}) = F_X(\sqrt{y}) - F_X(-\sqrt{y}). \end{aligned}$$

The pdf of  $Y$  can now be obtained from the cdf by differentiation:

$$\begin{aligned} f_Y(y) &= \frac{d}{dy} F_Y(y) \\ &= \frac{d}{dy} [F_X(\sqrt{y}) - F_X(-\sqrt{y})] \\ &= \frac{1}{2\sqrt{y}} f_X(\sqrt{y}) + \frac{1}{2\sqrt{y}} f_X(-\sqrt{y}), \end{aligned}$$

where we use the chain rule to differentiate  $F_X(\sqrt{y})$  and  $F_X(-\sqrt{y})$ . Therefore, the pdf is

$$(2.1.11) \quad f_Y(y) = \frac{1}{2\sqrt{y}} (f_X(\sqrt{y}) + f_X(-\sqrt{y})). \quad (\text{ } y > 0)$$

Notice that the pdf of  $Y$  in (2.1.11) is expressed as the sum of two pieces, pieces that represent the intervals where  $g(x) = x^2$  is monotone. In general, this will be the case. ||

**Theorem 2.1.8** Let  $X$  have pdf  $f_X(x)$ , let  $Y = g(X)$ , and define the sample space  $\mathcal{X}$  as in (2.1.7). Suppose there exists a partition,  $A_0, A_1, \dots, A_k$ , of  $\mathcal{X}$  such that  $P(X \in A_0) = 0$  and  $f_X(x)$  is continuous on each  $A_i$ . Further, suppose there exist functions  $g_1(x), \dots, g_k(x)$ , defined on  $A_1, \dots, A_k$ , respectively, satisfying

- i.  $g(x) = g_i(x)$ , for  $x \in A_i$ ,
- ii.  $g_i(x)$  is monotone on  $A_i$ ,
- iii. the set  $\mathcal{Y} = \{y: y = g_i(x) \text{ for some } x \in A_i\}$  is the same for each  $i = 1, \dots, k$ , and
- iv.  $g_i^{-1}(y)$  has a continuous derivative on  $\mathcal{Y}$ , for each  $i = 1, \dots, k$ .

Then

$$f_Y(y) = \begin{cases} \sum_{i=1}^k f_X(g_i^{-1}(y)) \left| \frac{d}{dy} g_i^{-1}(y) \right| & y \in \mathcal{Y} \\ 0 & \text{otherwise.} \end{cases}$$

Sep 7th 2021

**Theorem 2.1.10 (Probability integral transformation)** Let  $X$  have continuous cdf  $F_X(x)$  and define the random variable  $Y$  as  $Y = F_X(X)$ . Then  $Y$  is uniformly distributed on  $(0, 1)$ , that is,  $P(Y \leq y) = y$ ,  $0 < y < 1$ .

Before we prove this theorem, we will digress for a moment and look at  $F_X^{-1}$ , the inverse of the cdf  $F_X$ , in some detail. If  $F_X$  is strictly increasing, then  $F_X^{-1}$  is well defined by

$$(2.1.12) \quad F_X^{-1}(y) = x \Leftrightarrow F_X(x) = y.$$

However, if  $F_X$  is constant on some interval, then  $F_X^{-1}$  is not well defined by (2.1.12), as Figure 2.1.2 illustrates. Any  $x$  satisfying  $x_1 \leq x \leq x_2$  satisfies  $F_X(x) = y$ .

This problem is avoided by defining  $F_X^{-1}(y)$  for  $0 < y < 1$  by

$$(2.1.13) \quad F_X^{-1}(y) = \inf \{x: F_X(x) \geq y\},$$

a definition that agrees with (2.1.12) when  $F_X$  is nonconstant and provides an  $F_X^{-1}$  that is single-valued even when  $F_X$  is not strictly increasing. Using this definition, in Figure 2.1.2b, we have  $F_X^{-1}(y) = x_1$ . At the endpoints of the range of  $y$ ,  $F_X^{-1}(y)$  can also be defined.  $F_X^{-1}(1) = \infty$  if  $F_X(x) < 1$  for all  $x$  and, for any  $F_X$ ,  $F_X^{-1}(0) = -\infty$ .

**Proof of Theorem 2.1.10:** For  $Y = F_X(X)$  we have, for  $0 < y < 1$ ,

$$\begin{aligned} P(Y \leq y) &= P(F_X(X) \leq y) \\ &= P(F_X^{-1}[F_X(X)] \leq F_X^{-1}(y)) \quad (F_X^{-1} \text{ is increasing}) \\ &= P(X \leq F_X^{-1}(y)) \quad (\text{see paragraph below}) \\ &= F_X(F_X^{-1}(y)) \quad (\text{definition of } F_X) \\ &= y. \quad (\text{continuity of } F_X) \end{aligned}$$

At the endpoints we have  $P(Y \leq y) = 1$  for  $y \geq 1$  and  $P(Y \leq y) = 0$  for  $y \leq 0$ , showing that  $Y$  has a uniform distribution.

The reasoning behind the equality

$$P(F_X^{-1}(F_X(X)) \leq F_X^{-1}(y)) = P(X \leq F_X^{-1}(y))$$

is somewhat subtle and deserves additional attention. If  $F_X$  is strictly increasing, then it is true that  $F_X^{-1}(F_X(x)) = x$ . (Refer to Figure 2.1.2a.) However, if  $F_X$  is flat, it may be that  $F_X^{-1}(F_X(x)) \neq x$ . Suppose  $F_X$  is as in Figure 2.1.2b and let  $x \in [x_1, x_2]$ . Then  $F_X^{-1}(F_X(x)) = x_1$  for any  $x$  in this interval. Even in this case, though, the probability equality holds, since  $P(X \leq x) = P(X \leq x_1)$  for any  $x \in [x_1, x_2]$ . The flat cdf denotes a region of 0 probability ( $P(x_1 < X \leq x) = F_X(x) - F_X(x_1) = 0$ ). □

**Definition 2.2.1** The expected value or mean of a random variable  $g(X)$ , denoted by  $E g(X)$ , is

$$E g(X) = \begin{cases} \int_{-\infty}^{\infty} g(x) f_X(x) dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} g(x) f_X(x) = \sum_{x \in \mathcal{X}} g(x) P(X = x) & \text{if } X \text{ is discrete,} \end{cases}$$

provided that the integral or sum exists. If  $E|g(X)| = \infty$ , we say that  $E g(X)$  does not exist. (Ross 1988 refers to this as the “law of the unconscious statistician.”) We

**Example 2.2.3 (Binomial mean)** If  $X$  has a *binomial distribution*, its pmf is given by

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

where  $n$  is a positive integer,  $0 \leq p \leq 1$ , and for every fixed pair  $n$  and  $p$  the pmf sums to 1. The expected value of a binomial random variable is given by

$$E X = \sum_{x=0}^n x \binom{n}{x} p^x (1-p)^{n-x} = \sum_{x=1}^n x \binom{n}{x} p^x (1-p)^{n-x}$$

(the  $x = 0$  term is 0). Using the identity  $x \binom{n}{x} = n \binom{n-1}{x-1}$ , we have

$$\begin{aligned} E X &= \sum_{x=1}^n n \binom{n-1}{x-1} p^x (1-p)^{n-x} \\ &= \sum_{y=0}^{n-1} n \binom{n-1}{y} p^{y+1} (1-p)^{n-(y+1)} \quad (\text{substitute } y = x-1) \\ &= np \sum_{y=0}^{n-1} \binom{n-1}{y} p^y (1-p)^{n-1-y} \\ &= np, \end{aligned}$$

since the last summation must be 1, being the sum over all possible values of a  $\text{binomial}(n-1, p)$  pmf. ||

**Example 2.2.4 (Cauchy mean)** A classic example of a random variable whose expected value does not exist is a *Cauchy random variable*, that is, one with pdf

$$f_X(x) = \frac{1}{\pi} \frac{1}{1+x^2}, \quad -\infty < x < \infty.$$

It is straightforward to check that  $\int_{-\infty}^{\infty} f_X(x) dx = 1$ , but  $E|X| = \infty$ . Write

$$E|X| = \int_{-\infty}^{\infty} \frac{|x|}{\pi} \frac{1}{1+x^2} dx = \frac{2}{\pi} \int_0^{\infty} \frac{x}{1+x^2} dx.$$

For any positive number  $M$ ,

$$\int_0^M \frac{x}{1+x^2} dx = \left. \frac{\log(1+x^2)}{2} \right|_0^M = \frac{\log(1+M^2)}{2}.$$

Thus,

$$E|X| = \lim_{M \rightarrow \infty} \frac{2}{\pi} \int_0^M \frac{x}{1+x^2} dx = \frac{1}{\pi} \lim_{M \rightarrow \infty} \log(1+M^2) = \infty$$

and  $E X$  does not exist. ||

**Theorem 2.2.5** Let  $X$  be a random variable and let  $a, b$ , and  $c$  be constants. Then for any functions  $g_1(x)$  and  $g_2(x)$  whose expectations exist,

- a.  $E(ag_1(X) + bg_2(X) + c) = aE g_1(X) + bE g_2(X) + c$ .
- b. If  $g_1(x) \geq 0$  for all  $x$ , then  $E g_1(X) \geq 0$ .
- c. If  $g_1(x) \geq g_2(x)$  for all  $x$ , then  $E g_1(X) \geq E g_2(X)$ .
- d. If  $a \leq g_1(x) \leq b$  for all  $x$ , then  $a \leq E g_1(X) \leq b$ .

**Proof:** We will give details for only the continuous case, the discrete case being similar. By definition,

$$\begin{aligned} E(ag_1(X) + bg_2(X) + c) \\ &= \int_{-\infty}^{\infty} (ag_1(x) + bg_2(x) + c)f_X(x) dx \\ &= \int_{-\infty}^{\infty} ag_1(x)f_X(x) dx + \int_{-\infty}^{\infty} bg_2(x)f_X(x) dx + \int_{-\infty}^{\infty} cf_X(x) dx \end{aligned}$$

by the additivity of the integral. Since  $a$ ,  $b$ , and  $c$  are constants, they factor out of their respective integrals and we have

$$\begin{aligned} E(ag_1(X) + bg_2(X) + c) \\ &= a \int_{-\infty}^{\infty} g_1(x)f_X(x) dx + b \int_{-\infty}^{\infty} g_2(x)f_X(x) dx + c \int_{-\infty}^{\infty} f_X(x) dx \\ &= aE g_1(X) + bE g_2(X) + c, \end{aligned}$$

establishing (a). The other three properties are proved in a similar manner. □

**Definition 2.3.1** For each integer  $n$ , the  $n$ th moment of  $X$  (or  $F_X(x)$ ),  $\mu'_n$ , is

$$\mu'_n = E X^n.$$

The  $n$ th central moment of  $X$ ,  $\mu_n$ , is

$$\mu_n = E(X - \mu)^n,$$

where  $\mu = \mu'_1 = E X$ .

**Definition 2.3.2** The variance of a random variable  $X$  is its second central moment,  $\text{Var } X = E(X - E X)^2$ . The positive square root of  $\text{Var } X$  is the standard deviation of  $X$ .

$$SD(X) = \sqrt{\text{Var}(X)}$$

**Theorem 2.3.4** If  $X$  is a random variable with finite variance, then for any constants  $a$  and  $b$ ,

$$\text{Var}(aX + b) = a^2 \text{Var } X.$$

$$SD(aX+b)=|a|SD(X)$$

**Proof:** From the definition, we have

$$\begin{aligned} \text{Var}(aX + b) &= E((aX + b) - E(aX + b))^2 \\ &= E(aX - aE X)^2 \quad (E(aX + b) = aE X + b) \\ &= a^2 E(X - E X)^2 \\ &= a^2 \text{Var } X. \end{aligned}$$

□

Another formula for  $\text{Var } X$ :

$$\begin{aligned} \text{Var } X &= E(X - E X)^2 = E[X^2 - 2XE X + (E X)^2] \\ &= E X^2 - 2(E X)^2 + (E X)^2 \\ &= E X^2 - (E X)^2, \Rightarrow (E X)^2 \leq E X^2 \end{aligned}$$

Sep. 9th 2021      (Special case of Jensen's Ineqa)

**Example 2.3.5 (Binomial variance)** Let  $X \sim \text{binomial}(n, p)$ , that is,

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

We have previously seen that  $E X = np$ . To calculate  $\text{Var } X$  we first calculate  $E X^2$ . We have

$$(2.3.2) \quad E X^2 = \sum_{x=0}^n x^2 \binom{n}{x} p^x (1-p)^{n-x}.$$

In order to sum this series, we must first manipulate the binomial coefficient in a manner similar to that used for  $E X$  (Example 2.2.3). We write

$$(2.3.3) \quad x^2 \binom{n}{x} = x \frac{n!}{(x-1)!(n-x)!} = xn \binom{n-1}{x-1}.$$

The summand in (2.3.2) corresponding to  $x = 0$  is 0, and using (2.3.3), we have

$$\begin{aligned}
E X^2 &= n \sum_{x=1}^n x \binom{n-1}{x-1} p^x (1-p)^{n-x} \\
&= n \sum_{y=0}^{n-1} (y+1) \binom{n-1}{y} p^{y+1} (1-p)^{n-1-y} \quad (\text{setting } y = x-1) \\
&= np \sum_{y=0}^{n-1} y \binom{n-1}{y} p^y (1-p)^{n-1-y} + np \sum_{y=0}^{n-1} \binom{n-1}{y} p^y (1-p)^{n-1-y}.
\end{aligned}$$

Now it is easy to see that the first sum is equal to  $(n-1)p$  (since it is the mean of a binomial( $n-1, p$ )), while the second sum is equal to 1. Hence,

$$(2.3.4) \quad E X^2 = n(n-1)p^2 + np.$$

Using (2.3.1), we have

$$\text{Var } X = n(n-1)p^2 + np - (np)^2 = -np^2 + np = np(1-p). \quad \|$$

**Definition 2.3.6** Let  $X$  be a random variable with cdf  $F_X$ . The moment generating function (mgf) of  $X$  (or  $F_X$ ), denoted by  $M_X(t)$ , is

$$M_X(t) = E e^{tX},$$

provided that the expectation exists for  $t$  in some neighborhood of 0. That is, there is an  $h > 0$  such that, for all  $t$  in  $-h < t < h$ ,  $E e^{tX}$  exists. If the expectation does not exist in a neighborhood of 0, we say that the moment generating function does not exist.

More explicitly, we can write the mgf of  $X$  as

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx} f_X(x) dx \quad \text{if } X \text{ is continuous,}$$

or

$$M_X(t) = \sum_x e^{tx} P(X = x) \quad \text{if } X \text{ is discrete.}$$

Theorem. Suppose there exists a  $t^* > 0$  such that  $M_X(t^*) < \infty$ ,

$M_X(t^*) < \infty$ . Then,

$$\forall t_0 \in (-\frac{t^*}{2}, \frac{t^*}{2}): \left. \frac{d^n}{dt^n} M_X(t) \right|_{t=t_0} = E(X^n \cdot e^{t_0 X}), n=1, 2, 3, \dots$$

In particular, when  $t_0 = 0$ :

$$\left. \frac{d^n}{dt^n} M_X(t) \right|_{t=0} = E(X^n)$$

Lemma(Dominated Convergence Theorem)(DCT):

Suppose  $\frac{\partial}{\partial t} f(t, x)$  exists  $\forall t \in (t_0 - \delta, t_0 + \delta)$  &  $\forall x \in \mathbb{R}$

for some  $\delta > 0$ . If :

$$\textcircled{1} \quad \left| \frac{d}{dt} f(t, x) \right| \leq g(x) \quad (\text{uniformly over } t \in (t_0 - \delta, t_0 + \delta))$$

$$\textcircled{2} \quad \int_{\mathbb{R}} g(x) dx < +\infty$$

then  $\left( \frac{d}{dt} \int_{-\infty}^{+\infty} f(t, x) dx \right) \Big|_{t=t_0} = \int_{-\infty}^{+\infty} \left( \frac{d}{dt} f(t, x) \Big|_{t=t_0} \right) dx$

Pf: Start with  $n=1$

$$\begin{aligned} \frac{d}{dt} M_x(t) \Big|_{t=t_0} &= \frac{d}{dt} \left( \int_{\mathbb{R}} e^{tx} f_x(x) dx \right) \Big|_{t=t_0} \\ &\stackrel{\text{DCT}}{=} \int_{-\infty}^{+\infty} \frac{d}{dt} e^{tx} f_x(x) \Big|_{t=t_0} dx \\ &= \int_{-\infty}^{+\infty} x e^{tx} f_x(x) dx \\ &= E(X \cdot e^{tx}) \end{aligned}$$

Verify DCT conditions:  $f(t, x) = e^{tx} f_x(x)$

$$\begin{aligned} \left| \frac{d}{dt} f(t, x) \right| &= \left| x e^{tx} f_x(x) \right| = |x| \cdot e^{tx} f_x(x) \\ &\leq |x| \cdot e^{\frac{|t^*|}{2}|x|} f_x(x) \quad ; \forall t \in (-\frac{t^*}{2}, \frac{t^*}{2}) \end{aligned}$$

Now, since  $t_0 \in (-\frac{t^*}{2}, \frac{t^*}{2})$ , we can choose  $\delta$  small enough such that  $(t_0 - \delta, t_0 + \delta) \subseteq (-\frac{t^*}{2}, \frac{t^*}{2})$

Since  $|y| \leq e^{|y|}$ ,  $\forall y \in \mathbb{R}$ , we have:  $(|\frac{t^*}{2} \cdot x| \leq e^{\frac{|t^*|}{2}|x|})$

$$\begin{aligned} \left| \frac{d}{dt} f(t, x) \right| &\leq \frac{2}{t^*} e^{\frac{|t^*|}{2}|x|} f_x(x) \\ &\leq \frac{2}{t^*} (e^{t^*x} + e^{-t^*x}) f_x(x) =: g(x) \end{aligned}$$

$$\begin{aligned} \int_{\mathbb{R}} g(x) dx &= \frac{2}{t^*} \int_{\mathbb{R}} (e^{t^*x} + e^{-t^*x}) f_x(x) dx \\ &= \frac{2}{t^*} (M_x(t^*) + M_x(-t^*)) < +\infty \end{aligned}$$

Proceed the proof by induction :

Suppose:  $\frac{d^n}{dt^n} M_X(t) \Big|_{t=t_0} = E(X^n e^{t_0 X})$ ;  $\forall t_0 \in (-\frac{t^*}{2}, \frac{t^*}{2})$

Then: 
$$\begin{aligned} \frac{d^{n+1}}{dt^{n+1}} M_X(t) \Big|_{t=t_0} &= \frac{d}{dt} E(X^n e^{t_0 X}) \Big|_{t=t_0} \\ &= \left( \frac{d}{dt} \int_{\mathbb{R}} x^n e^{t_0 x} f_X(x) dx \right) \Big|_{t=t_0} \\ &\quad "f_{t,x} \end{aligned}$$

Let  $\delta$  be chosen such that  $(t_0 - \delta, t_0 + \delta) \subseteq (-\frac{t^*}{2}, \frac{t^*}{2})$ , then:

$$\begin{aligned} \forall t \in (t_0 - \delta, t_0 + \delta), \left| \frac{d}{dt} f_{t,x} \right| &= \left| \frac{d}{dt} x^n e^{t x} f_X(x) \right| \\ &= |x|^{n+1} e^{t x} f_X(x) \leq |x|^{n+1} \cdot e^{\frac{t^*}{2}|x|} f_X(x) \\ (|y| \leq (n+1) \cdot e^{|y|}) &\leq \frac{2^{n+1}}{(t^*)^{n+1}} e^{t^*|x|} \cdot f_X(x) \\ &\leq \frac{2^{n+1}}{(t^*)^{n+1}} (e^{t^*|x|} + e^{-t^*|x|}) f_X(x) =: g(x) \\ \int_{\mathbb{R}} g(x) dx &= \frac{2^{n+1}}{(t^*)^{n+1}} (M_X(t^*) + M_X(-t^*)) < +\infty \end{aligned}$$

Sep 14th. 2021

**Theorem 2.3.11** Let  $F_X(x)$  and  $F_Y(y)$  be two cdfs all of whose moments exist.

- a. If  $X$  and  $Y$  have bounded support, then  $F_X(u) = F_Y(u)$  for all  $u$  if and only if  $E X^r = E Y^r$  for all integers  $r = 0, 1, 2, \dots$ .
- b. If the moment generating functions exist and  $M_X(t) = M_Y(t)$  for all  $t$  in some neighborhood of 0, then  $F_X(u) = F_Y(u)$  for all  $u$ .

**Example 2.3.10 (Nonunique moments)** Consider the two pdfs given by

$$f_1(x) = \frac{1}{\sqrt{2\pi x}} e^{-(\log x)^2/2}, \quad 0 \leq x < \infty,$$

$$f_2(x) = f_1(x)[1 + \sin(2\pi \log x)], \quad 0 \leq x < \infty.$$

(The pdf  $f_1$  is a special case of a *lognormal pdf*.)

It can be shown that if  $X_1 \sim f_1(x)$ , then

$$E X_1^r = e^{r^2/2}, \quad r = 0, 1, \dots,$$

so  $X_1$  has all of its moments. Now suppose that  $X_2 \sim f_2(x)$ . We have

$$\begin{aligned} E X_2^r &= \int_0^\infty x^r f_1(x) [1 + \sin(2\pi \log x)] dx \\ &= E X_1^r + \int_0^\infty x^r f_1(x) \sin(2\pi \log x) dx. = \text{Ex}_1^r \\ &\quad \text{= 0 (see below)} \end{aligned}$$

However, the transformation  $y = \log x - r$  shows that this last integral is that of an odd function over  $(-\infty, \infty)$  and hence is equal to 0 for  $r = 0, 1, \dots$ . Thus, even though  $X_1$  and  $X_2$  have distinct pdfs, they have the same moments for all  $r$ . The two pdfs are pictured in Figure 2.3.2.

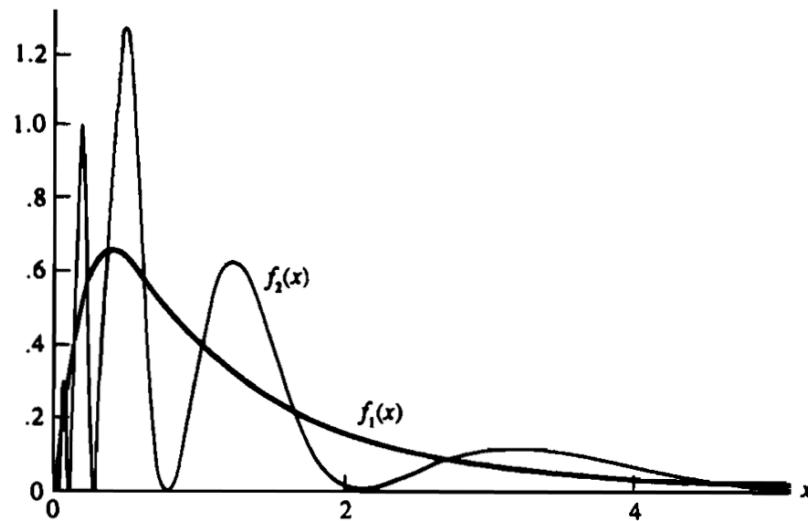


Figure 2.3.2. Two pdfs with the same moments:  $f_1(x) = \frac{1}{\sqrt{2\pi x}} e^{-(\log x)^2/2}$  and  $f_2(x) = f_1(x)[1 + \sin(2\pi \log x)]$

$$\begin{aligned} \int_0^\infty x^r f_1(x) \sin(2\pi \log x) dx &\stackrel{y = \log x - r}{=} \int_{-\infty}^{+\infty} (e^{y+r})^r \frac{1}{\sqrt{2\pi e^{y+r}}} e^{-\frac{(y+r)^2}{2}} \cdot \underbrace{\sin(2\pi(y+r))}_{\text{odd function}} e^{y+r} dy \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} \sin(2\pi y) \cdot e^{y+r^2 - \frac{1}{2}y^2 - py - \frac{1}{2}r^2} dy = \sin(2\pi y) \\ &= \frac{e^{r^2}}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{y^2}{2}} \sin(2\pi y) dy = 0 \end{aligned}$$

odd function & absolute integrable

**Theorem 2.3.12 (Convergence of mgfs)** Suppose  $\{X_i, i = 1, 2, \dots\}$  is a sequence of random variables, each with mgf  $M_{X_i}(t)$ . Furthermore, suppose that

$$\lim_{i \rightarrow \infty} M_{X_i}(t) = M_X(t), \quad \text{for all } t \text{ in a neighborhood of 0,}$$

and  $M_X(t)$  is an mgf. Then there is a unique cdf  $F_X$  whose moments are determined by  $M_X(t)$  and, for all  $x$  where  $F_X(x)$  is continuous, we have

$$\lim_{i \rightarrow \infty} F_{X_i}(x) = F_X(x).$$

That is, convergence, for  $|t| < h$ , of mgfs to an mgf implies convergence of cdfs.

**Theorem 2.3.15** For any constants  $a$  and  $b$ , the mgf of the random variable  $aX+b$  is given by

$$M_{aX+b}(t) = e^{bt} M_X(at).$$

**Proof:** By definition,

$$\begin{aligned} M_{aX+b}(t) &= E(e^{(aX+b)t}) \\ &= E(e^{(aX)t} e^{bt}) \quad (\text{properties of exponentials}) \\ &= e^{bt} E(e^{(at)X}) \quad (e^{bt} \text{ is constant}) \\ &= e^{bt} M_X(at), \quad (\text{definition of mgf}) \end{aligned}$$

proving the theorem.  $\square$

## Chapter 3: Common Families of Distributions

### Discrete Distributions

#### ① Discrete Uniform Distribution

A random variable  $X$  has a discrete uniform  $(1, N)$  distribution if

$$(3.2.1) \quad P(X = x|N) = \frac{1}{N}, \quad x = 1, 2, \dots, N,$$

We then have

$$EX = \sum_{x=1}^N x P(X = x|N) = \sum_{x=1}^N x \frac{1}{N} = \frac{N+1}{2}$$

and

$$EX^2 = \sum_{x=1}^N x^2 \frac{1}{N} = \frac{(N+1)(2N+1)}{6},$$

and so

$$\begin{aligned} \text{Var } X &= EX^2 - (EX)^2 \\ &= \frac{(N+1)(2N+1)}{6} - \left(\frac{N+1}{2}\right)^2 \\ &= \frac{(N+1)(N-1)}{12}. \end{aligned}$$

## (2) Hypergeometric Distribution

Suppose we have a large urn filled with  $N$  balls that are identical in every way except that  $M$  are red and  $N - M$  are green. We reach in, blindfolded, and select  $K$  balls at random (the  $K$  balls are taken all at once, a case of sampling without replacement). What is the probability that exactly  $x$  of the balls are red?

The total number of samples of size  $K$  that can be drawn from the  $N$  balls is  $\binom{N}{K}$ , as was discussed in Section 1.2.3. It is required that  $x$  of the balls be red, and this can be accomplished in  $\binom{M}{x}$  ways, leaving  $\binom{N-M}{K-x}$  ways of filling out the sample with  $K - x$  green balls. Thus, if we let  $X$  denote the number of red balls in a sample of size  $K$ , then  $X$  has a hypergeometric distribution given by

$$(3.2.2) \quad P(X = x|N, M, K) = \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}}, \quad x = 0, 1, \dots, K.$$

Note that there is, implicit in (3.2.2), an additional assumption on the range of  $X$ . Binomial coefficients of the form  $\binom{n}{r}$  have been defined only if  $n \geq r$ , and so the range of  $X$  is additionally restricted by the pair of inequalities

$$M \geq x \quad \text{and} \quad N - M \geq K - x,$$

which can be combined as

$$M - (N - K) \leq x \leq M.$$

In many cases  $K$  is small compared to  $M$  and  $N$ , so the range  $0 \leq x \leq K$  will be contained in the above range and, hence, will be appropriate. The formula for the hypergeometric probability function is usually quite difficult to deal with. In fact, it is not even trivial to verify that

$$\sum_{x=0}^K P(X = x) = \sum_{x=0}^K \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}} = 1.$$

The hypergeometric distribution illustrates the fact that, statistically, dealing with finite populations (finite  $N$ ) is a difficult task.

The mean of the hypergeometric distribution is given by

$$EX = \sum_{x=0}^K x \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}} = \sum_{x=1}^K x \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}}. \quad (\text{summand is 0 at } x=0)$$

To evaluate this expression, we use the identities (already encountered in Section 2.3)

$$x \binom{M}{x} = M \binom{M-1}{x-1},$$

$$\binom{N}{K} = \frac{N}{K} \binom{N-1}{K-1},$$

and obtain

$$EX = \sum_{x=1}^K \frac{M \binom{M-1}{x-1} \binom{N-M}{K-x}}{\frac{N}{K} \binom{N-1}{K-1}} = \frac{KM}{N} \sum_{x=1}^K \frac{\binom{M-1}{x-1} \binom{N-M}{K-x}}{\binom{N-1}{K-1}}.$$

We now can recognize the second sum above as the sum of the probabilities for another hypergeometric distribution based on parameter values  $N - 1$ ,  $M - 1$ , and  $K - 1$ . This can be seen clearly by defining  $y = x - 1$  and writing

$$\begin{aligned} \sum_{x=1}^K \frac{\binom{M-1}{x-1} \binom{N-M}{K-x}}{\binom{N-1}{K-1}} &= \sum_{y=0}^{K-1} \frac{\binom{M-1}{y} \binom{(N-1)-(M-1)}{K-1-y}}{\binom{N-1}{K-1}} \\ &= \sum_{y=0}^{K-1} P(Y = y | N - 1, M - 1, K - 1) = 1, \end{aligned}$$

where  $Y$  is a hypergeometric random variable with parameters  $N - 1$ ,  $M - 1$ , and  $K - 1$ . Therefore, for the hypergeometric distribution,

$$EX = \frac{KM}{N}.$$

A similar, but more lengthy, calculation will establish that

$$\text{Var } X = \frac{KM}{N} \left( \frac{(N-M)(N-K)}{N(N-1)} \right).$$

$$\begin{aligned} E(X(X-1)) &= \sum_{x=2}^K x(x-1) \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}} \\ &= \sum_{x=2}^K x(x-1) \frac{M(M-1)}{x(x-1)} \frac{\binom{M-2}{x-2} \binom{N-M}{K-x}}{\binom{N}{K}} \\ &= \frac{M(M-1)}{\binom{N}{K}} \sum_{y=0}^{K-2} \binom{M-2}{y} \binom{(N-2)-(M-2)}{K-2-y} \\ &= \frac{M(M-1)}{\binom{N}{K}} \binom{N-2}{K-2} = \frac{M(M-1)K(K-1)}{N(N-1)} \\ \text{Var}(X) &= E(X(X-1)) + EX - (EX)^2 = \dots = \frac{KM}{N} \left( 1 - \frac{M}{N} \right) \frac{N-K}{N-1} \end{aligned}$$

### ③ Binomial Distribution

A random variable  $X$  has a *Bernoulli(p)* distribution if

$$(3.2.3) \quad X = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p, \end{cases} \quad 0 \leq p \leq 1.$$

$$EX = 1p + 0(1-p) = p,$$

$$\text{Var } X = (1-p)^2 p + (0-p)^2 (1-p) = p(1-p).$$

In a sequence of  $n$  identical, independent Bernoulli trials, each with success probability  $p$ , define the random variables  $X_1, \dots, X_n$  by

$$X_i = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p. \end{cases}$$

The random variable

$$Y = \sum_{i=1}^n X_i$$

has the binomial( $n, p$ ) distribution.

$$P(Y = y | n, p) = \binom{n}{y} p^y (1-p)^{n-y}, \quad y = 0, 1, 2, \dots, n,$$

$$EX = np, \quad \text{Var } X = np(1-p). \quad M_X(t) = [pe^t + (1-p)]^n.$$

$$\begin{aligned} M_X(t) &= E(e^{tX}) = \sum_{x=0}^n e^{tx} \cdot P(X=x) \\ &= \sum_{x=0}^n \binom{n}{x} e^{tx} \cdot p^x (1-p)^{n-x} \\ &= \sum_{x=0}^n \binom{n}{x} (pe^t)^x (1-p)^{n-x} \\ &= (pe^t + 1-p)^n \end{aligned}$$

$$\frac{dM_X(t)}{dt} \Big|_{t=0} = n \cdot pe^t (pe^t + 1-p)^{n-1} \Big|_{t=0} = np$$

Theorem: Suppose  $X \sim B(n, p)$ ,  $Y \sim B(m, p)$ . Also,  $X$  and  $Y$  are independent. Then,  $Z = X + Y \sim B(n+m, p)$

$$\begin{aligned} \text{Pf. } P(Z=z) &= \sum_{k=0}^n P(X=k, Y=z-k) \\ &= \sum_{k=0}^n P(X=k) P(Y=z-k) \\ &= \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} \binom{m}{z-k} p^{z-k} (1-p)^{m-(z-k)} \\ &= p^z (1-p)^{n+m-z} \sum_{k=0}^z \binom{n}{k} \binom{m}{z-k} \\ &= \binom{n+m}{z} p^z (1-p)^{n+m-z} \end{aligned}$$

④

## Poisson Distribution

A random variable  $X$ , taking values in the nonnegative integers, has a Poisson( $\lambda$ ) distribution if

$$(3.2.5) \quad P(X = x|\lambda) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, \dots$$

$$\begin{aligned} EX &= \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} \\ &= \sum_{x=1}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} \\ &= \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!} \\ &= \lambda e^{-\lambda} \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} \\ &= \lambda \end{aligned}$$

$$\begin{aligned} M_X(t) &= e^{\lambda(e^t - 1)}. \\ M_X(t) &= \sum_{x=0}^{\infty} e^{tx} \frac{e^{-\lambda} \lambda^x}{x!} \\ &= e^{-\lambda} \sum_{x=0}^{\infty} \frac{(xe^t)^x}{x!} \\ &= e^{-\lambda} e^{\lambda e^t} = e^{\lambda(e^t - 1)} \end{aligned}$$

$$\text{Var}(X) = \lambda.$$

$X \sim B(n, p)$ . Suppose  $n \rightarrow +\infty, p \rightarrow 0, n \cdot p \rightarrow \lambda > 0$

$$\begin{aligned} \text{Now: } P(X=x) &= \binom{n}{x} p^x (1-p)^{n-x}; \quad x=0, 1, 2, \dots \\ &= \frac{1}{x!} \cdot \frac{n!}{(n-x)!} \cdot \frac{(np)^x}{n^x} \cdot \left(1 - \frac{np}{n}\right)^{n-x} \\ &= \frac{(np)^x}{x!} \cdot \frac{n(n-1)(n-2)\dots(n-x+1)}{n^x} \cdot \left(1 - \frac{np}{n}\right)^n (1-p)^{-x} \\ &\quad \downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow \\ &\quad \frac{x^x}{x!} \qquad \qquad 1 \qquad \qquad e^{-\lambda} \qquad 1 \\ &\rightarrow e^{-\lambda} \frac{\lambda^x}{x!} \end{aligned}$$

**Example 3.2.4 (Waiting time)** As an example of a waiting-for-occurrence application, consider a telephone operator who, on the average, handles five calls every 3 minutes. What is the probability that there will be no calls in the next minute? At least two calls?

If we let  $X$  = number of calls in a minute, then  $X$  has a Poisson distribution with  $EX = \lambda = \frac{5}{3}$ . So

$$P(\text{no calls in the next minute}) = P(X = 0)$$

$$\begin{aligned} &= \frac{e^{-5/3} \left(\frac{5}{3}\right)^0}{0!} \\ &= e^{-5/3} = .189; \end{aligned}$$

$$P(\text{at least two calls in the next minute}) = P(X \geq 2)$$

$$\begin{aligned} &= 1 - P(X = 0) - P(X = 1) \\ &= 1 - .189 - \frac{e^{-5/3} \left(\frac{5}{3}\right)^1}{1!} \\ &= .496. \end{aligned}$$

\*Sep. 16th 2021

Calculation of Poisson probabilities can be done rapidly by noting the following recursion relation:

$$(3.2.6) \quad P(X = x) = \frac{\lambda}{x} P(X = x - 1), \quad x = 1, 2, \dots$$

This relation is easily proved by writing out the pmf of the Poisson. Similar relations hold for other discrete distributions. For example, if  $Y \sim \text{binomial}(n, p)$ , then

$$(3.2.7) \quad P(Y = y) = \frac{(n - y + 1)}{y} \frac{p}{1 - p} P(Y = y - 1).$$

The recursion relations (3.2.6) and (3.2.7) can be used to establish the Poisson approximation to the binomial, which we have already seen in Section 2.3, where the approximation was justified using mgfs. Set  $\lambda = np$  and, if  $p$  is small, we can write

$$\frac{n - y + 1}{y} \frac{p}{1 - p} = \frac{np - p(y - 1)}{y - py} \approx \frac{\lambda}{y}$$

since, for small  $p$ , the terms  $p(y - 1)$  and  $py$  can be ignored. Therefore, to this level of approximation, (3.2.7) becomes

$$(3.2.8) \quad P(Y = y) = \frac{\lambda}{y} P(Y = y - 1),$$

which is the Poisson recursion relation. To complete the approximation, we need only establish that  $P(X = 0) \approx P(Y = 0)$ , since all other probabilities will follow from (3.2.8). Now

$$P(Y = 0) = (1 - p)^n = \left(1 - \frac{np}{n}\right)^n = \left(1 - \frac{\lambda}{n}\right)^n$$

upon setting  $np = \lambda$ . Recall from Section 2.3 that for fixed  $\lambda$ ,  $\lim_{n \rightarrow \infty} (1 - (\lambda/n))^n = e^{-\lambda}$ , so for large  $n$  we have the approximation

$$P(Y = 0) = \left(1 - \frac{\lambda}{n}\right)^n \approx e^{-\lambda} = P(X = 0),$$

completing the Poisson approximation to the binomial.

The approximation is valid when  $n$  is large and  $p$  is small, which is exactly when it is most useful, freeing us from calculation of binomial coefficients and powers for large  $n$ .

$$X \sim B(n, p) \Rightarrow M_X(t) = (pe^t + 1 - p)^n = (1 + \frac{1}{n}(e^t - 1)np)^n$$

as  $n \rightarrow \infty$   $\xrightarrow{} e^{\lambda(e^t - 1)}$   
 [  $p \rightarrow 0$  ]  $[ np \rightarrow \lambda ]$

**Example 3.2.5 (Poisson approximation)** A typesetter, on the average, makes one error in every 500 words typeset. A typical page contains 300 words. What is the probability that there will be no more than two errors in five pages?

If we assume that setting a word is a Bernoulli trial with success probability  $p = \frac{1}{500}$  (notice that we are labeling an error as a "success") and that the trials are independent, then  $X = \text{number of errors in five pages (1500 words)}$  is binomial( $1500, \frac{1}{500}$ ). Thus

$$P(\text{no more than two errors}) = P(X \leq 2)$$

$$\begin{aligned} &= \sum_{x=0}^2 \binom{1500}{x} \left(\frac{1}{500}\right)^x \left(\frac{499}{500}\right)^{1500-x} \\ &= .4230, \end{aligned}$$

which is a fairly cumbersome calculation. If we use the Poisson approximation with  $\lambda = 1500 \left(\frac{1}{500}\right) = 3$ , we have

$$P(X \leq 2) \approx e^{-3} \left(1 + 3 + \frac{3^2}{2}\right) = .4232.$$

### ⑤ Geometric Distribution

$$P(X = x|p) = p(1-p)^{x-1}, \quad x = 1, 2, \dots, \quad X \sim \text{Geometric}(p)$$

which defines the pmf of a *geometric random variable*  $X$  with success probability  $p$ .  $X$  can be interpreted as the trial at which the first success occurs

The instructor's def of  $G(p)$ :  $P(X=x) = p(1-p)^x, x=0, 1, 2, \dots$

$$\begin{aligned} E[X] &= \sum_{x=1}^{+\infty} x \cdot p \cdot (1-p)^x = p \sum_{x=1}^{+\infty} (\sum_{y=1}^x 1) \cdot (1-p)^x = p \sum_{x=1}^{+\infty} \sum_{y=1}^x (1-p)^x \\ &= p \sum_{y=1}^{+\infty} \sum_{x=y}^{+\infty} (1-p)^x = p \cdot \sum_{y=1}^{+\infty} \frac{(1-p)^y}{1-(1-p)} \\ &= \sum_{y=1}^{+\infty} (1-p)^y = \frac{1-p}{p} \end{aligned}$$

$$\begin{aligned} M_X(t) &= \sum_{x=0}^{+\infty} e^{tx} p(1-p)^x = p \sum_{x=0}^{+\infty} ((1-p)e^t)^x \\ &= \begin{cases} \frac{p}{1-(1-p)e^t}, & \text{if } (1-p)e^t < 1 \text{ or } t < -\log(1-p) \\ +\infty & \text{otherwise} \end{cases} \end{aligned}$$

$$E[X] = \frac{dM_X(t)}{dt} \Big|_{t=0} = \frac{p(1-p)e^t}{(1-(1-p)e^t)^2} \Big|_{t=0} = \frac{p(1-p)}{p^2} = \frac{1-p}{p}$$

$$EX^2 = \frac{d^2/M_x(t)}{dt^2} \Big|_{t=0} = \frac{P(1-P)e^t}{(1-(1-P)e^t)^2} + \frac{2P(1-P)^2e^{2t}}{(1-(1-P)e^t)^3} \Big|_{t=0}$$

$$= \frac{P(1-P)}{P^2} + \frac{2P(1-P)^2}{P^3} = \frac{1-P}{P} + \frac{2(1-P)^2}{P^2}$$

$$\text{Var}(X) = EX^2 - (EX)^2 = \frac{1-P}{P^2}$$

The geometric distribution has an interesting property, known as the "memoryless" property. For integers  $s > t$ , it is the case that

$$(3.2.11) \quad P(X > s | X > t) = P(X > s - t);$$

∴ Sep 21st. 2021

$$P(X \geq t+s | X \geq t) = P(X \geq s), \forall s, t \in \mathbb{N}_0$$

$$\text{Pf. } P(X \geq X) = \sum_{t=\infty}^{+\infty} P(1-P)^t = P \frac{(1-P)^{\infty}}{1-(1-P)} = (1-P)^{\infty}$$

$$P(X \geq t+s | X \geq t) = \frac{P(X \geq t+s, X \geq t)}{P(X \geq t)}$$

$$= \frac{P(X \geq t+s)}{P(X \geq t)}$$

$$= (1-P)^{t+s} / (1-P)^t$$

$$= (1-P)^s = P(X \geq s)$$

Thm. Suppose  $X$  is a r.v. whose range is inside integer numbers, i.e.,

$\text{Range}(X) \subseteq \mathbb{Z}$ , and  $P(X \geq 0) > 0$ . Now, if

$$P(X \geq t+s | X \geq t) = P(X \geq s); \forall s, t \geq 0,$$

then  $X \sim \text{Geometric}(P(X=0))$ .

Pf. Use MP with  $t=s=0$ :  $P(X \geq 0+0 | X \geq 0) = P(X \geq 0)$

$$\frac{P(X \geq 0)}{P(X \geq 0)} \stackrel{||}{=} 1 \quad \therefore P(X \geq 0) = 1$$

Use MP with  $t=s=1$ :

$$\frac{P(X \geq 2)}{P(X \geq 1)} = P(X \geq 2 | X \geq 1) = P(X \geq 1) \Rightarrow P(X \geq 2) = (P(X \geq 1))^2$$

By induction:  $P(X \geq n) = (P(X \geq 1))^n ; \forall n \geq 1$

$$\left( \frac{P(X \geq n+1)}{P(X \geq 1)} \right) = P(X \geq n+1 | X \geq 1) = P(X \geq n) = (P(X \geq 1))^n \Rightarrow P(X \geq n+1) = (P(X \geq 1))^{n+1}$$

$$\begin{aligned} P(X=n) &= P(X \geq n) - P(X \geq n+1) \\ &= (P(X \geq 1))^n - (P(X \geq 1))^{n+1} \\ &= (P(X \geq 1))^n (1 - P(X \geq 1)) \\ &= (1 - P(X=0))^n P(X=0) \end{aligned}$$

## ⑥ Negative Binomial Distribution

Repeat independently a Bernoulli( $p$ )

Define  $X = \# \text{failure before the } r\text{th success}$

$$X \sim \text{Neg Binom}(r, p)$$

$$\text{Range}(X) = \{0, 1, 2, \dots\}$$

$$P(X=x) = p^r (1-p)^{x-r} \binom{x+r-1}{x}$$

$$M_X(t) = \sum_{x=0}^{+\infty} e^{tx} \binom{x+r-1}{x} p^r (1-p)^{x-r}$$

negative binomial formula

$|a| < 1 \text{ & } n \in \mathbb{N}$

$$(1-a)^{-r} = \sum_{x=0}^{+\infty} \binom{x+r-1}{x} a^x \quad = p^r (1 - (1-p)e^t)^{-r}, \text{ if } (1-p)e^t < 1 \text{ or } t < -\log(1-p)$$

$$\begin{aligned} \frac{dM_X(t)}{dt} \Big|_{t=0} &= p^r (-r)(-(1-p)e^t)(1 - (1-p)e^t)^{-r-1} \Big|_{t=0} \\ &= r p^r (1-p)e^t (1 - (1-p)e^t)^{-r-1} \Big|_{t=0} \\ &= r \frac{1-p}{p} = EX \end{aligned}$$

$$\begin{aligned} \frac{d^2 M_X(t)}{dt^2} \Big|_{t=0} &= r p^r (1-p) \left[ e^t (1 - (1-p)e^t)^{-r-1} + e^t (-r-1)(-(1-p)e^t)(1 - (1-p)e^t)^{-r-2} \right] \Big|_{t=0} \\ &= r p^r (1-p) \left( \frac{1}{p^{r+1}} + \frac{(r+1)(1-p)}{p^{r+2}} \right) \\ &= \frac{r^2 (1-p)^2}{p^2} + r \frac{1-p}{p^2} \end{aligned}$$

$$\Rightarrow \text{Var}(X) = EX^2 - (EX)^2 = r \frac{1-p}{p^2}$$

Proof of Negative Binomial Formula:

$$f_m(z) = \sum_{k=0}^{+\infty} \binom{k+m}{m} z^k$$

$$\text{Note that } \binom{k+m}{m} = \binom{k+m-1}{m-1} + \binom{k+m-1}{m}$$

$$\Rightarrow f_m(z) = f_{m-1}(z) + z \sum_{k=1}^{+\infty} \binom{k-1+m}{m} z^{k-1}$$

$$= f_{m-1}(z) + z f_m(z)$$

$$\Rightarrow f_m(z) = \frac{1}{1-z} f_{m-1}(z)$$

$$f_0(z) = \sum_{k=0}^{+\infty} \binom{k+0}{0} z^k = \sum_{k=0}^{+\infty} z^k = \frac{1}{1-z}$$

$$\therefore f_m(z) = (1-z)^{-m+1} \quad \text{Negative Binomial Formula}$$

## Continuous Distributions

### ① Uniform Distribution

The continuous *uniform distribution* is defined by spreading mass uniformly over an interval  $[a, b]$ . Its pdf is given by

$$f(x|a, b) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise.} \end{cases}$$

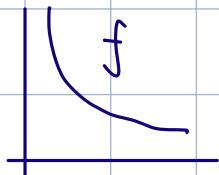
$$Mx(t) = E(e^{tx}) = \int_a^b \frac{e^{tx}}{b-a} dx = \begin{cases} \frac{e^{bt} - e^{at}}{t(b-a)} & ; t \neq 0 \\ 1 & ; t = 0 \end{cases}$$

$$EX = \int_a^b \frac{x}{b-a} dx = \frac{b+a}{2};$$

$$\text{Var } X = \int_a^b \frac{(x - \frac{b+a}{2})^2}{b-a} dx = \frac{(b-a)^2}{12}.$$

### ② Exponential Distribution

$$X \sim \text{Exp}(\theta), \theta > 0, \text{ if } f_X(x) = \begin{cases} \frac{1}{\theta} e^{-\frac{x}{\theta}} & , x > 0 \\ 0 & ; \text{o.w.} \end{cases}$$



$$P(X \in [a-h, a+h]) = e^{-\frac{a}{\theta}} (e^{\frac{h}{\theta}} - e^{-\frac{h}{\theta}})$$

$$EX = \int_0^{+\infty} \frac{x}{\theta} e^{-\frac{x}{\theta}} dx = \theta \Gamma(2) = \theta$$

$$EX^2 = \int_0^{+\infty} \frac{x^2}{\theta} e^{-\frac{x}{\theta}} dx = \theta^2 \Gamma(3) = 2\theta^2$$

$$\text{Var}(X) = EX^2 - (EX)^2 = \theta^2$$

$$\begin{aligned}
 M_X(t) &= \int_0^{+\infty} e^{tx} \frac{1}{\theta} e^{-\frac{x}{\theta}} dx \\
 &= \frac{1}{\theta} \int_0^{+\infty} e^{(t-\frac{1}{\theta})x} dx \\
 &= \begin{cases} +\infty & ; t \geq \frac{1}{\theta} \\ \frac{1/\theta}{1/\theta - t} = \frac{1}{1-\theta t} & ; t < \frac{1}{\theta} \end{cases}
 \end{aligned}$$

Thm:  $X \sim \text{Exp}(\theta) \Rightarrow P(X > t+s | X > t) = P(X > s)$ ;  $\forall t, s \in [0, +\infty)$

$$\text{Pf: } P(X > u) = \int_u^{+\infty} \frac{1}{\theta} e^{-\frac{x}{\theta}} dx = (-e^{-\frac{x}{\theta}}) \Big|_u^{+\infty} = e^{-\frac{u}{\theta}} \quad (u > 0)$$

$$\begin{aligned}
 \text{Now, } P(X > t+s | X > t) &= P(X > t+s) / P(X > t) \\
 &= e^{-\frac{(t+s)}{\theta}} / e^{-\frac{t}{\theta}} \\
 &= e^{-\frac{s}{\theta}} = P(X > s)
 \end{aligned}$$

Thm: Suppose  $X$  is a continuous r.v. with  $P(X > t) > 0$ ;  $\forall t > 0$ , and  
 $P(X > s+t | X > t) = P(X > s)$ ;  $\forall t, s \in [0, +\infty)$ .

Then  $X \sim \text{Exp}(\theta)$  for some  $\theta > 0$ .

Pf:  $t=s=0$ :  $P(X > 0 | X > 0) = P(X > 0) \Rightarrow P(X > 0) = 1 \Rightarrow \text{Range}(X) \in [0, +\infty)$

Define  $G(x) = P(X > x)$ ;  $\forall x \geq 0$  [ $G(\cdot)$  is a continuous function]

$$t=s=1: \quad G(2) = (G(1))^2$$

By induction,  $G(n) = (G(1))^n$ ;  $\forall n \in \mathbb{N}$

Let  $x = \frac{p}{q}$ ;  $p, q \in \mathbb{N}$  using MP.

$$\begin{aligned}
 G(p) &= G\left(\frac{p}{q} \cdot q\right) = \left(G\left(\frac{p}{q}\right)\right)^q \\
 \Rightarrow G\left(\frac{p}{q}\right) &= \left(G(p)\right)^{1/q} = \left((G(1))^p\right)^{1/q} = (G(1))^{p/q}
 \end{aligned}$$

Let  $x \in [0, +\infty)$ , then  $\exists$  sequence of  $\{x_n\}_{n \geq 1}$  of positive rational numbers such that  $x_n \rightarrow x$  as  $n \rightarrow \infty$ . Now,

$$\begin{aligned}
 G(x) &= G\left(\lim_{n \rightarrow \infty} x_n\right) \xrightarrow[\text{continuity of } G]{\text{continuity}} \lim_{n \rightarrow \infty} G(x_n) \\
 &= \lim_{n \rightarrow \infty} (G(1))^{x_n} = (G(1))^x = e^{x \log G(1)} \\
 &= e^{-(-\log G(1))x} \quad (\theta = \frac{-1}{\log G(1)})
 \end{aligned}$$

$F_X(x) = 1 - G(x) = 1 - e^{-\frac{x}{\theta}}$ ;  $\forall x > 0$  (matches the CDF of  $\text{Exp}(\theta)$ )

✓ Sep 23rd. 2021

### ③ Gamma Distribution

$$[\Gamma(1, \beta) = \text{Exp}(\beta)]$$

$$f(x|\alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta}, \quad 0 < x < \infty, \quad \alpha > 0, \quad \beta > 0.$$

The parameter  $\alpha$  is known as the shape parameter, since it most influences the peakedness of the distribution, while the parameter  $\beta$  is called the scale parameter, since most of its influence is on the spread of the distribution.

$$\begin{aligned} EX &= \frac{1}{\Gamma(\alpha)\beta^\alpha} \int_0^\infty x^\alpha e^{-x/\beta} dx \\ &= \frac{1}{\Gamma(\alpha)\beta^\alpha} \Gamma(\alpha+1)\beta^{\alpha+1} \quad \leftarrow \int_0^\infty x^{\alpha-1} e^{-x/\beta} dx = \Gamma(\alpha)\beta^\alpha, \\ &= \frac{\alpha\Gamma(\alpha)\beta}{\Gamma(\alpha)} \\ &= \alpha\beta. \end{aligned}$$

$$\text{Var } X = \alpha\beta^2.$$

$$\begin{aligned} M_X(t) &= E(e^{tx}) = \frac{1}{\Gamma(\alpha)\beta^\alpha} \int_0^{+\infty} e^{tx} x^{\alpha-1} e^{-x/\beta} dx \\ &= \frac{1}{\Gamma(\alpha)\beta^\alpha} \int_0^{+\infty} x^{\alpha-1} e^{-x/(1-\beta t)} dx \\ &= \Gamma(\alpha) \left( \frac{\beta}{1-\beta t} \right)^\alpha / \Gamma(\alpha)\beta^\alpha = \begin{cases} \left( \frac{1}{1-\beta t} \right)^\alpha & ; \text{ if } t < \frac{1}{\beta} \\ +\infty & ; \text{ o.w.} \end{cases} \end{aligned}$$

There are a number of important special cases of the gamma distribution. If we set  $\alpha = p/2$ , where  $p$  is an integer, and  $\beta = 2$ , then the gamma pdf becomes

$$(3.3.10) \quad f(x|p) = \frac{1}{\Gamma(p/2)2^{p/2}} x^{(p/2)-1} e^{-x/2}, \quad 0 < x < \infty, \quad \Gamma(\frac{p}{2}, 2) = \chi_p^2$$

which is the chi squared pdf with  $p$  degrees of freedom. The mean, variance, and mgf of the chi squared distribution can all be calculated by using the previously derived gamma formulas.

④

## Normal Distribution

The normal distribution has two parameters, usually denoted by  $\mu$  and  $\sigma^2$ , which are its mean and variance. The pdf of the normal distribution with mean  $\mu$  and variance  $\sigma^2$  (usually denoted by  $n(\mu, \sigma^2)$ ) is given by

$$(3.3.13) \quad f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/(2\sigma^2)}, \quad -\infty < x < \infty.$$

If  $X \sim n(\mu, \sigma^2)$ , then the random variable  $Z = (X - \mu)/\sigma$  has a  $n(0, 1)$  distribution, also known as the standard normal. This is easily established by writing

$$\begin{aligned} P(Z \leq z) &= P\left(\frac{X - \mu}{\sigma} \leq z\right) \\ &= P(X \leq z\sigma + \mu) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{z\sigma+\mu} e^{-(x-\mu)^2/(2\sigma^2)} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-t^2/2} dt, \quad \left(\text{substitute } t = \frac{x - \mu}{\sigma}\right) \end{aligned}$$

showing that  $P(Z \leq z)$  is the standard normal cdf.

We have not yet established that (3.3.13) integrates to 1 over the whole real line.

By applying the standardizing transformation, we need only to show that

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-z^2/2} dz = 1.$$

Notice that the integrand above is symmetric around 0, implying that the integral over  $(-\infty, 0)$  is equal to the integral over  $(0, \infty)$ . Thus, we reduce the problem to showing

$$(3.3.14) \quad \int_0^{\infty} e^{-z^2/2} dz = \frac{\sqrt{2\pi}}{2} = \sqrt{\frac{\pi}{2}}.$$

The function  $e^{-z^2/2}$  does not have an antiderivative that can be written explicitly in terms of elementary functions (that is, in closed form), so we cannot perform the integration directly. In fact, this is an example of an integration that either you know how to do or else you can spend a very long time going nowhere. Since both sides of (3.3.14) are positive, the equality will hold if we establish that the squares are equal. Square the integral in (3.3.14) to obtain

$$\begin{aligned} \left(\int_0^{\infty} e^{-z^2/2} dz\right)^2 &= \left(\int_0^{\infty} e^{-t^2/2} dt\right) \left(\int_0^{\infty} e^{-u^2/2} du\right) \\ &= \int_0^{\infty} \int_0^{\infty} e^{-(t^2+u^2)/2} dt du. \end{aligned}$$

The integration variables are just dummy variables, so changing their names is allowed. Now, we convert to polar coordinates. Define

$$t = r \cos \theta \quad \text{and} \quad u = r \sin \theta.$$

Then  $t^2 + u^2 = r^2$  and  $dt du = r d\theta dr$  and the limits of integration become  $0 < r < \infty$ ,  $0 < \theta < \pi/2$  (the upper limit on  $\theta$  is  $\pi/2$  because  $t$  and  $u$  are restricted to be positive). We now have

$$\begin{aligned}
\int_0^\infty \int_0^\infty e^{-(t^2+u^2)/2} dt du &= \int_0^\infty \int_0^{\pi/2} r e^{-r^2/2} d\theta dr \\
&= \frac{\pi}{2} \int_0^\infty r e^{-r^2/2} dr \\
&= \frac{\pi}{2} \left[ -e^{-r^2/2} \Big|_0^\infty \right] \\
&= \frac{\pi}{2},
\end{aligned}$$

which establishes (3.3.14).

This integral is closely related to the gamma function; in fact, by making the substitution  $w = \frac{1}{2}z^2$  in (3.3.14), we see that this integral is essentially  $\Gamma(\frac{1}{2})$ . If we are careful to get the constants correct, we will see that (3.3.14) implies

$$(3.3.15) \quad \Gamma\left(\frac{1}{2}\right) = \int_0^\infty w^{-1/2} e^{-w} dw = \sqrt{\pi}.$$

$$\sqrt{\pi} = \int_0^{+\infty} e^{-\frac{z^2}{2}} dz \stackrel{z=\sqrt{y}}{=} \int_0^{+\infty} e^{-y} \frac{1}{\sqrt{2y}} dy = \frac{1}{\sqrt{2}} \Gamma\left(\frac{1}{2}\right) \Rightarrow \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}.$$

$$\begin{aligned}
M_Z(t) &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{tz} e^{-z^2/2} dz \\
&= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z^2-2tz+t^2)} dz \\
&= e^{t^2/2} \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z-t)^2} dz = e^{\frac{t^2}{2}}
\end{aligned}$$

$$\begin{aligned}
M_X(t) &= E(e^{tx}) = E(e^{t(\sigma z + \mu)}) \\
&= e^{t\mu} E(e^{t\sigma z}) = e^{t\mu} M_Z(t\sigma) \\
&= e^{t\mu} e^{\frac{(t\sigma)^2}{2}} = e^{t\mu + t^2\sigma^2/2} \quad ; \forall t \in \mathbb{R}
\end{aligned}$$

✓ Sep 28th 2021

**Example 3.3.2 (Normal approximation)** Let  $X \sim \text{binomial}(25, .6)$ . We can approximate  $X$  with a normal random variable,  $Y$ , with mean  $\mu = 25(.6) = 15$  and standard deviation  $\sigma = ((25)(.6)(.4))^{1/2} = 2.45$ . Thus

$$P(X \leq 13) \approx P(Y \leq 13) = P\left(Z \leq \frac{13 - 15}{2.45}\right) = P(Z \leq -.82) = .206,$$

while the exact binomial calculation gives

$$P(X \leq 13) = \sum_{x=0}^{13} \binom{25}{x} (.6)^x (.4)^{25-x} = .267,$$

showing that the normal approximation is good, but not terrific. The approximation can be greatly improved, however, by a “continuity correction.” To see how this works, look at Figure 3.3.2, which shows the binomial(25, .6) pmf and the  $n(15, (2.45)^2)$  pdf. We have drawn the binomial pmf using bars of width 1, with height equal to the probability. Thus, the areas of the bars give the binomial probabilities. In the approximation, notice how the area of the approximating normal is smaller than the binomial area (the normal area is everything to the left of the line at 13, whereas the binomial area includes the entire bar at 13 up to 13.5). The continuity correction adds this area back by adding  $\frac{1}{2}$  to the cutoff point. So instead of approximating  $P(X \leq 13)$ , we approximate the equivalent expression (because of the discreteness),  $P(X \leq 13.5)$  and obtain

$$P(X \leq 13) = P(X \leq 13.5) \approx P(Y \leq 13.5) = P(Z \leq -.61) = .271,$$

a much better approximation. In general, the normal approximation with the continuity correction is far superior to the approximation without the continuity correction.

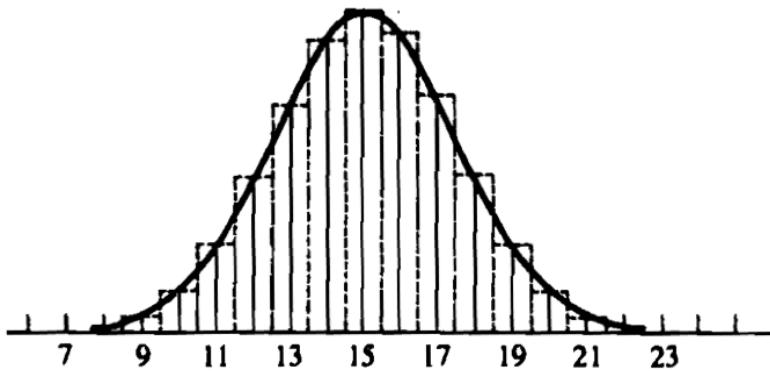


Figure 3.3.2.  $Normal(15, (2.45)^2)$  approximation to the binomial(25, .6)

We also make the correction on the lower end. If  $X \sim \text{binomial}(n, p)$  and  $Y \sim n(np, np(1 - p))$ , then we approximate

$$P(X \leq x) \approx P(Y \leq x + 1/2),$$

$$P(X \geq x) \approx P(Y \geq x - 1/2).$$

||

### Beta Distribution

The beta family of distributions is a continuous family on  $(0, 1)$  indexed by two parameters. The  $\text{beta}(\alpha, \beta)$  pdf is

$$(3.3.16) \quad f(x|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad 0 < x < 1, \quad \alpha > 0, \quad \beta > 0,$$

where  $B(\alpha, \beta)$  denotes the beta function,

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx.$$

The beta function is related to the gamma function through the following identity:

$$(3.3.17) \quad B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}.$$

$\alpha=1, \beta=1$  : Beta(1, 1) = Uniform(0, 1)

$\alpha>1, \beta=1$  : strictly increasing  $f_X(x)$

$\alpha=1, \beta>1$  : strictly decreasing  $f_X(x)$

$\alpha<1, \beta<1$  : V-shaped

$\alpha>1, \beta>1$  : Unimodal

Calculation of moments of the beta distribution is quite easy, due to the particular form of the pdf. For  $n > -\alpha$  we have

$$\begin{aligned} EX^n &= \frac{1}{B(\alpha, \beta)} \int_0^1 x^n x^{\alpha-1} (1-x)^{\beta-1} dx \\ &= \frac{1}{B(\alpha, \beta)} \int_0^1 x^{(\alpha+n)-1} (1-x)^{\beta-1} dx. \\ &= \frac{B(\alpha+n, \beta)}{B(\alpha, \beta)} = \frac{\Gamma(\alpha+n)\Gamma(\alpha+\beta)}{\Gamma(\alpha+\beta+n)\Gamma(\alpha)}. \end{aligned}$$

Using (3.3.3) and (3.3.18) with  $n = 1$  and  $n = 2$ , we calculate the mean and variance of the beta( $\alpha, \beta$ ) distribution as

$$EX = \frac{\alpha}{\alpha + \beta} \quad \text{and} \quad \text{Var } X = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}.$$

## ⑥ Cauchy Distribution

The *Cauchy distribution* is a symmetric, bell-shaped distribution on  $(-\infty, \infty)$  with pdf

$$(3.3.19) \quad f(x|\theta) = \frac{1}{\pi} \frac{1}{1+(x-\theta)^2}, \quad -\infty < x < \infty, \quad -\infty < \theta < \infty.$$

# Exponential Families

A family of pdfs or pmfs is called an exponential family if it can be expressed as

$$(3.4.1) \quad f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta}) \exp\left(\sum_{i=1}^k w_i(\boldsymbol{\theta})t_i(x)\right).$$

Here  $h(x) \geq 0$  and  $t_1(x), \dots, t_k(x)$  are real-valued functions of the observation  $x$  (they cannot depend on  $\boldsymbol{\theta}$ ), and  $c(\boldsymbol{\theta}) \geq 0$  and  $w_1(\boldsymbol{\theta}), \dots, w_k(\boldsymbol{\theta})$  are real-valued functions of the possibly vector-valued parameter  $\boldsymbol{\theta}$  (they cannot depend on  $x$ ). Many common families introduced in the previous section are exponential families. These include the continuous families—normal, gamma, and beta, and the discrete families—binomial, Poisson, and negative binomial.

**Example 3.4.1 (Binomial exponential family)** Let  $n$  be a positive integer and consider the binomial( $n, p$ ) family with  $0 < p < 1$ . Then the pmf for this family, for  $x = 0, \dots, n$  and  $0 < p < 1$ , is

$$\begin{aligned} f(x|p) &= \binom{n}{x} p^x (1-p)^{n-x} \\ &= \binom{n}{x} (1-p)^n \left(\frac{p}{1-p}\right)^x \\ &= \binom{n}{x} (1-p)^n \exp\left(\log\left(\frac{p}{1-p}\right)x\right). \end{aligned}$$

Define

$$h(x) = \begin{cases} \binom{n}{x} & x = 0, \dots, n \\ 0 & \text{otherwise,} \end{cases} \quad c(p) = (1-p)^n, \quad 0 < p < 1,$$

$$w_1(p) = \log\left(\frac{p}{1-p}\right), \quad 0 < p < 1, \quad \text{and} \quad t_1(x) = x.$$

Then we have

$$f(x|p) = h(x)c(p) \exp[w_1(p)t_1(x)],$$

**Theorem 3.4.2** If  $X$  is a random variable with pdf or pmf of the form (3.4.1), then

$$(3.4.4) \quad E\left(\sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X)\right) = -\frac{\partial}{\partial \theta_j} \log c(\boldsymbol{\theta});$$

$$(3.4.5) \quad \text{Var}\left(\sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X)\right) = -\frac{\partial^2}{\partial \theta_j^2} \log c(\boldsymbol{\theta}) - E\left(\sum_{i=1}^k \frac{\partial^2 w_i(\boldsymbol{\theta})}{\partial \theta_j^2} t_i(X)\right).$$

**Example 3.4.3 (Binomial mean and variance)** From Example 3.4.1 we have

$$\begin{aligned}\frac{d}{dp}w_1(p) &= \frac{d}{dp}\log\frac{p}{1-p} = \frac{1}{p(1-p)} \\ \frac{d}{dp}\log c(p) &= \frac{d}{dp}n\log(1-p) = \frac{-n}{1-p}\end{aligned}$$

and thus from Theorem 3.4.2 we have

$$E\left(\frac{1}{p(1-p)}X\right) = \frac{n}{1-p}$$

and a bit of rearrangement yields  $E(X) = np$ . The variance identity works in a similar manner.  $\parallel$

**Example 3.4.4 (Normal exponential family)** Let  $f(x|\mu, \sigma^2)$  be the  $n(\mu, \sigma^2)$  family of pdfs, where  $\theta = (\mu, \sigma)$ ,  $-\infty < \mu < \infty$ ,  $\sigma > 0$ . Then

$$\begin{aligned}(3.4.6) \quad f(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \exp\left(-\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2}\right).\end{aligned}$$

Define

$$h(x) = 1 \text{ for all } x;$$

$$c(\theta) = c(\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-\mu^2}{2\sigma^2}\right), \quad -\infty < \mu < \infty, \sigma > 0;$$

$$w_1(\mu, \sigma) = \frac{1}{\sigma^2}, \quad \sigma > 0; \quad w_2(\mu, \sigma) = \frac{\mu}{\sigma^2}, \quad \sigma > 0;$$

$$t_1(x) = -x^2/2; \quad \text{and} \quad t_2(x) = x.$$

Then

$$f(x|\mu, \sigma^2) = h(x)c(\mu, \sigma) \exp[w_1(\mu, \sigma)t_1(x) + w_2(\mu, \sigma)t_2(x)],$$

From (3.4.1), since the factor  $\exp(\cdot)$  is always positive, it can be seen that for any  $\theta \in \Theta$ , that is, for any  $\theta$  for which  $c(\theta) > 0$ ,  $\{x : f(x|\theta) > 0\} = \{x : h(x) > 0\}$  and this set does not depend on  $\theta$ . So, for example, the set of pdfs given by  $f(x|\theta) = \theta^{-1} \exp(1 - (x/\theta))$ ,  $0 < \theta < x < \infty$ , is not an exponential family even though we can write  $\theta^{-1} \exp(1 - (x/\theta)) = h(x)c(\theta) \exp(w(\theta)t(x))$ , where  $h(x) = e^1$ ,  $c(\theta) = \theta^{-1}$ ,  $w(\theta) = \theta^{-1}$ , and  $t(x) = -x$ . Writing the pdf with indicator functions makes this very clear. We have

$$f(x|\theta) = \theta^{-1} \exp\left(1 - \left(\frac{x}{\theta}\right)\right) I_{[\theta, \infty)}(x).$$

The indicator function cannot be incorporated into any of the functions of (3.4.1) since it is not a function of  $x$  alone, not a function of  $\theta$  alone, and cannot be expressed as an exponential. Thus, this is not an exponential family.

**Definition 3.4.7** A curved exponential family is a family of densities of the form (3.4.1) for which the dimension of the vector  $\theta$  is equal to  $d < k$ . If  $d = k$ , the family is a full exponential family. (See also *Miscellanea 3.8.3*.)

**Example 3.4.8 (A curved exponential family)** The normal family of Example 3.4.4 is a full exponential family. However, if we assume that  $\sigma^2 = \mu^2$ , the family becomes curved. (Such a model might be used in the analysis of variance; see Exercises 11.1 and 11.2.) We then have

$$\begin{aligned} f(x|\mu) &= \frac{1}{\sqrt{2\pi\mu^2}} \exp\left(-\frac{(x-\mu)^2}{2\mu^2}\right) I_{(-\infty, +\infty)}(x) \\ (3.4.9) \quad &= \frac{1}{\sqrt{2\pi\mu^2}} \exp\left(-\frac{1}{2}\right) \exp\left(-\frac{x^2}{2\mu^2} + \frac{x}{\mu}\right) I_{(-\infty, +\infty)}(x). \end{aligned}$$

For the normal family the full exponential family would have parameter space  $(\mu, \sigma^2) = \mathbb{R} \times (0, \infty)$ , while the parameter space of the curved family  $(\mu, \sigma^2) = (\mu, \mu^2)$  is a parabola. ||

**Example 3.4.9 (Normal approximations)** In Chapter 5 we will see that if  $X_1, \dots, X_n$  is a sample from a  $\text{Poisson}(\lambda)$  population, then the distribution of  $\bar{X} = \sum_i X_i/n$  is approximately

$$\bar{X} \xrightarrow{\text{Normal}} n(\lambda, \lambda/n),$$

a curved exponential family.

The  $n(\lambda, \lambda/n)$  approximation is justified by the Central Limit Theorem (Theorem 5.5.14). In fact, we might realize that most such CLT approximations will result in a curved normal family. We have seen the normal binomial approximation (Example 3.3.2): If  $X_1, \dots, X_n$  are iid  $\text{Bernoulli}(p)$ , then

$$\bar{X} \xrightarrow{\text{Normal}} n(p, p(1-p)/n),$$

approximately. For another illustration, see Example 5.5.16. ||

## Location & Scale Families

**Theorem 3.5.1** Let  $f(x)$  be any pdf and let  $\mu$  and  $\sigma > 0$  be any given constants. Then the function

$$g(x|\mu, \sigma) = \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right)$$

is a pdf.

Pf:  $g \geq 0$  a.s. and

$$\int_{-\infty}^{\infty} \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right) dx = \int_{-\infty}^{\infty} f(y) dy \quad \left(\text{substitute } y = \frac{x-\mu}{\sigma}\right)$$

$$= 1, \quad (\text{since } f(y) \text{ is a pdf})$$

**Definition 3.5.2** Let  $f(x)$  be any pdf. Then the family of pdfs  $f(x - \mu)$ , indexed by the parameter  $\mu$ ,  $-\infty < \mu < \infty$ , is called the location family with standard pdf  $f(x)$  and  $\mu$  is called the location parameter for the family.

**Example 3.5.3 (Exponential location family)** Let  $f(x) = e^{-x}$ ,  $x \geq 0$ , and  $f(x) = 0$ ,  $x < 0$ . To form a location family we replace  $x$  with  $x - \mu$  to obtain

$$f(x|\mu) = \begin{cases} e^{-(x-\mu)} & x - \mu \geq 0 \\ 0 & x - \mu < 0 \end{cases}$$

$$= \begin{cases} e^{-(x-\mu)} & x \geq \mu \\ 0 & x < \mu. \end{cases}$$

Graphs of  $f(x|\mu)$  for various values of  $\mu$  are shown in Figure 3.5.2. As in Figure 3.5.1, the graph has been shifted. Now the positive part of the graph starts at  $\mu$  rather than at 0. If  $X$  measures time, then  $\mu$  might be restricted to be nonnegative so that  $X$  will be positive with probability 1 for every value of  $\mu$ . In this type of model, where  $\mu$  denotes a bound on the range of  $X$ ,  $\mu$  is sometimes called a threshold parameter. ||

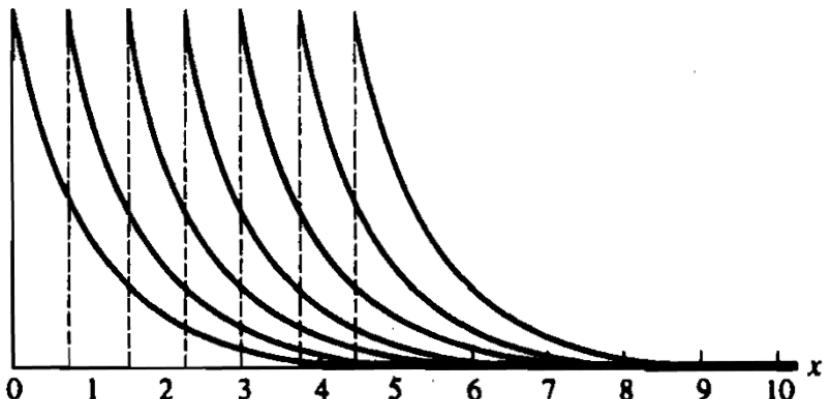


Figure 3.5.2. Exponential location densities

**Definition 3.5.4** Let  $f(x)$  be any pdf. Then for any  $\sigma > 0$ , the family of pdfs  $(1/\sigma)f(x/\sigma)$ , indexed by the parameter  $\sigma$ , is called the scale family with standard pdf  $f(x)$  and  $\sigma$  is called the scale parameter of the family.

The effect of introducing the scale parameter  $\sigma$  is either to stretch ( $\sigma > 1$ ) or to contract ( $\sigma < 1$ ) the graph of  $f(x)$  while still maintaining the same basic shape of the graph.

**Definition 3.5.5** Let  $f(x)$  be any pdf. Then for any  $\mu$ ,  $-\infty < \mu < \infty$ , and any  $\sigma > 0$ , the family of pdfs  $(1/\sigma)f((x - \mu)/\sigma)$ , indexed by the parameter  $(\mu, \sigma)$ , is called the location-scale family with standard pdf  $f(x)$ ;  $\mu$  is called the location parameter and  $\sigma$  is called the scale parameter.

**Theorem 3.5.6** Let  $f(\cdot)$  be any pdf. Let  $\mu$  be any real number, and let  $\sigma$  be any positive real number. Then  $X$  is a random variable with pdf  $(1/\sigma)f((x - \mu)/\sigma)$  if and only if there exists a random variable  $Z$  with pdf  $f(z)$  and  $X = \sigma Z + \mu$ .

**Proof:** To prove the “if” part, define  $g(z) = \sigma z + \mu$ . Then  $X = g(Z)$ ,  $g$  is a monotone function,  $g^{-1}(x) = (x - \mu)/\sigma$ , and  $|(d/dx)g^{-1}(x)| = 1/\sigma$ . Thus by Theorem 2.1.5, the pdf of  $X$  is

$$f_X(x) = f_Z(g^{-1}(x)) \left| \frac{d}{dx} g^{-1}(x) \right| = f\left(\frac{x - \mu}{\sigma}\right) \frac{1}{\sigma}.$$

To prove the “only if” part, define  $g(x) = (x - \mu)/\sigma$  and let  $Z = g(X)$ . Theorem 2.1.5 again applies:  $g^{-1}(z) = \sigma z + \mu$ ,  $|(d/dz)g^{-1}(z)| = \sigma$ , and the pdf of  $Z$  is

$$f_Z(z) = f_X(g^{-1}(z)) \left| \frac{d}{dz} g^{-1}(z) \right| = \frac{1}{\sigma} f\left(\frac{(\sigma z + \mu) - \mu}{\sigma}\right) \sigma = f(z).$$

Also,

$$\sigma Z + \mu = \sigma g(X) + \mu = \sigma \left( \frac{X - \mu}{\sigma} \right) + \mu = X. \quad \square$$

# Probability Inequalities

**Theorem 3.6.1 (Chebychev's Inequality)** Let  $X$  be a random variable and let  $g(x)$  be a nonnegative function. Then, for any  $r > 0$ ,

$$P(g(X) \geq r) \leq \frac{Eg(X)}{r}.$$

**Proof:**

$$\begin{aligned} Eg(X) &= \int_{-\infty}^{\infty} g(x)f_X(x)dx \\ &\geq \int_{\{x:g(x)\geq r\}} g(x)f_X(x)dx \quad (g \text{ is nonnegative}) \\ &\geq r \int_{\{x:g(x)\geq r\}} f_X(x)dx \\ &= rP(g(X) \geq r). \quad (\text{definition}) \end{aligned}$$

Rearranging now produces the desired inequality.  $\square$

**Example 3.6.2 (Illustrating Chebychev)** The most widespread use of Chebychev's Inequality involves means and variances. Let  $g(x) = (x - \mu)^2/\sigma^2$ , where  $\mu = EX$  and  $\sigma^2 = \text{Var } X$ . For convenience write  $r = t^2$ . Then

$$P\left(\frac{(X - \mu)^2}{\sigma^2} \geq t^2\right) \leq \frac{1}{t^2} E\frac{(X - \mu)^2}{\sigma^2} = \frac{1}{t^2}.$$

Doing some obvious algebra, we get the inequality

$$P(|X - \mu| \geq t\sigma) \leq \frac{1}{t^2} \quad P\left(\left|\frac{X - \mu}{\sigma}\right| \geq t\right) \leq \frac{1}{t^2}$$

and its companion

$$P(|X - \mu| < t\sigma) \geq 1 - \frac{1}{t^2},$$

which gives a universal bound on the deviation  $|X - \mu|$  in terms of  $\sigma$ . For example, taking  $t = 2$ , we get

$$P(|X - \mu| \geq 2\sigma) \leq \frac{1}{2^2} = .25,$$

so there is at least a 75% chance that a random variable will be within  $2\sigma$  of its mean (no matter what the distribution of  $X$ ).  $\parallel$

**Example 3.6.3 (A normal probability inequality)** If  $Z$  is standard normal, then

$$(3.6.1) \quad P(|Z| \geq t) \leq \sqrt{\frac{2}{\pi}} \frac{e^{-t^2/2}}{t}, \quad \text{for all } t > 0.$$

Compare this with Chebychev's Inequality. For  $t = 2$ , Chebychev gives  $P(|Z| \geq t) \leq .25$  but  $\sqrt{(2/\pi)}e^{-2}/2 = .054$ , a vast improvement.

To prove (3.6.1), write

$$\begin{aligned} P(Z \geq t) &= \frac{1}{\sqrt{2\pi}} \int_t^{\infty} e^{-x^2/2} dx \\ &\leq \frac{1}{\sqrt{2\pi}} \int_t^{\infty} \frac{x}{t} e^{-x^2/2} dx && \left( \begin{array}{l} \text{since } x/t > 1 \\ \text{for } x > t \end{array} \right) \\ &= \frac{1}{\sqrt{2\pi}} \frac{e^{-t^2/2}}{t} \end{aligned}$$

and use the fact that  $P(|Z| \geq t) = 2P(Z \geq t)$ . A lower bound on  $P(|Z| \geq t)$  can be established in a similar way (see Exercise 3.47). ||

If  $M_X(t)$  exists, then

$$P(X \geq a) \leq e^{-at} M_X(t),$$

~~Oct 7th 2021~~

## ~~Chapter 4: Multiple Random Variables~~

**Definition 4.1.1** An  $n$ -dimensional random vector is a function from a sample space  $S$  into  $\mathbb{R}^n$ ,  $n$ -dimensional Euclidean space.

**Definition 4.1.3** Let  $(X, Y)$  be a discrete bivariate random vector. Then the function  $f(x, y)$  from  $\mathbb{R}^2$  into  $\mathbb{R}$  defined by  $f(x, y) = P(X = x, Y = y)$  is called the joint probability mass function or joint pmf of  $(X, Y)$ . If it is necessary to stress the fact that  $f$  is the joint pmf of the vector  $(X, Y)$  rather than some other vector, the notation  $f_{X,Y}(x, y)$  will be used.

Expectations of functions of random vectors are computed just as with univariate random variables. Let  $g(x, y)$  be a real-valued function defined for all possible values  $(x, y)$  of the discrete random vector  $(X, Y)$ . Then  $g(X, Y)$  is itself a random variable and its expected value  $Eg(X, Y)$  is given by

$$Eg(X, Y) = \sum_{(x,y) \in \mathbb{R}^2} g(x, y)f(x, y).$$

**Theorem 4.1.6** Let  $(X, Y)$  be a discrete bivariate random vector with joint pmf  $f_{X,Y}(x, y)$ . Then the marginal pmfs of  $X$  and  $Y$ ,  $f_X(x) = P(X = x)$  and  $f_Y(y) = P(Y = y)$ , are given by

$$f_X(x) = \sum_{y \in \mathbb{R}} f_{X,Y}(x, y) \quad \text{and} \quad f_Y(y) = \sum_{x \in \mathbb{R}} f_{X,Y}(x, y).$$

**Proof:** We will prove the result for  $f_X(x)$ . The proof for  $f_Y(y)$  is similar. For any  $x \in \mathbb{R}$ , let  $A_x = \{(x, y) : -\infty < y < \infty\}$ . That is,  $A_x$  is the line in the plane with first coordinate equal to  $x$ . Then, for any  $x \in \mathbb{R}$ ,

$$\begin{aligned} f_X(x) &= P(X = x) \\ &= P(X = x, -\infty < Y < \infty) \quad (P(-\infty < Y < \infty) = 1) \\ &= P((X, Y) \in A_x) \quad (\text{definition of } A_x) \\ &= \sum_{(x,y) \in A_x} f_{X,Y}(x, y) \\ &= \sum_{y \in \mathbb{R}} f_{X,Y}(x, y). \end{aligned}$$

□

~~Oct 12th 2021~~

**Definition 4.1.10** A function  $f(x, y)$  from  $\mathbb{R}^2$  into  $\mathbb{R}$  is called a joint probability density function or joint pdf of the continuous bivariate random vector  $(X, Y)$  if, for every  $A \subset \mathbb{R}^2$ ,

$$P((X, Y) \in A) = \int_A \int f(x, y) dx dy.$$

the *expected value of  $g(X, Y)$*  is defined to be

$$\text{E}g(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f(x, y) dx dy.$$

The marginal probability density functions of  $X$  and  $Y$  are also defined as in the discrete case with integrals replacing sums. The marginal pdfs may be used to compute probabilities or expectations that involve only  $X$  or  $Y$ . Specifically, the marginal pdfs of  $X$  and  $Y$  are given by

$$(4.1.3) \quad \begin{aligned} f_X(x) &= \int_{-\infty}^{\infty} f(x, y) dy, \quad -\infty < x < \infty, \\ f_Y(y) &= \int_{-\infty}^{\infty} f(x, y) dx, \quad -\infty < y < \infty. \end{aligned}$$

The joint probability distribution of  $(X, Y)$  can be completely described with the *joint cdf* (cumulative distribution function) rather than with the joint pmf or joint pdf. The joint cdf is the function  $F(x, y)$  defined by

$$F(x, y) = P(X \leq x, Y \leq y)$$

for all  $(x, y) \in \mathbb{R}^2$ . The joint cdf is usually not very handy to use for a discrete random vector. But for a continuous bivariate random vector we have the important relationship, as in the univariate case,

$$F(x, y) = \int_{-\infty}^x \int_{-\infty}^y f(s, t) dt ds.$$

From the bivariate Fundamental Theorem of Calculus, this implies that

$$(4.1.4) \quad \frac{\partial^2 F(x, y)}{\partial x \partial y} = f(x, y)$$

at continuity points of  $f(x, y)$ . This relationship is useful in situations where an expression for  $F(x, y)$  can be found. The mixed partial derivative can be computed to find the joint pdf.

**Definition 4.2.1** Let  $(X, Y)$  be a discrete bivariate random vector with joint pmf  $f(x, y)$  and marginal pmfs  $f_X(x)$  and  $f_Y(y)$ . For any  $x$  such that  $P(X = x) = f_X(x) > 0$ , the conditional pmf of  $Y$  given that  $X = x$  is the function of  $y$  denoted by  $f(y|x)$  and defined by

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

For any  $y$  such that  $P(Y = y) = f_Y(y) > 0$ , the conditional pmf of  $X$  given that  $Y = y$  is the function of  $x$  denoted by  $f(x|y)$  and defined by

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

**Definition 4.2.3** Let  $(X, Y)$  be a continuous bivariate random vector with joint pdf  $f(x, y)$  and marginal pdfs  $f_X(x)$  and  $f_Y(y)$ . For any  $x$  such that  $f_X(x) > 0$ , the conditional pdf of  $Y$  given that  $X = x$  is the function of  $y$  denoted by  $f(y|x)$  and defined by

$$f(y|x) = \frac{f(x, y)}{f_X(x)}.$$

For any  $y$  such that  $f_Y(y) > 0$ , the conditional pdf of  $X$  given that  $Y = y$  is the function of  $x$  denoted by  $f(x|y)$  and defined by

$$f(x|y) = \frac{f(x, y)}{f_Y(y)}.$$

In addition to their usefulness for calculating probabilities, the conditional pdfs or pmfs can also be used to calculate expected values. Just remember that  $f(y|x)$  as a function of  $y$  is a pdf or pmf and use it in the same way that we have previously used unconditional pdfs or pmfs. If  $g(Y)$  is a function of  $Y$ , then the conditional expected value of  $g(Y)$  given that  $X = x$  is denoted by  $E(g(Y)|x)$  and is given by

$$E(g(Y)|x) = \sum_y g(y)f(y|x) \quad \text{and} \quad E(g(Y)|x) = \int_{-\infty}^{\infty} g(y)f(y|x) dy$$

The variance of the probability distribution described by  $f(y|x)$  is called the conditional variance of  $Y$  given  $X = x$ . Using the notation  $\text{Var}(Y|x)$  for this, we have, using the ordinary definition of variance,

$$\text{Var}(Y|x) = E(Y^2|x) - (E(Y|x))^2.$$

**Definition 4.2.5** Let  $(X, Y)$  be a bivariate random vector with joint pdf or pmf  $f(x, y)$  and marginal pdfs or pmfs  $f_X(x)$  and  $f_Y(y)$ . Then  $X$  and  $Y$  are called independent random variables if, for every  $x \in \mathfrak{R}$  and  $y \in \mathfrak{R}$ ,

$$(4.2.1) \quad f(x, y) = f_X(x)f_Y(y).$$

$X \perp\!\!\!\perp Y$

If  $X$  and  $Y$  are independent, the conditional pdf of  $Y$  given  $X = x$  is

$$\begin{aligned} f(y|x) &= \frac{f(x, y)}{f_X(x)} && \text{(definition)} \\ &= \frac{f_X(x)f_Y(y)}{f_X(x)} && \text{(from (4.2.1))} \\ &= f_Y(y), \end{aligned}$$

regardless of the value of  $x$ . Thus, for any  $A \subset \mathfrak{R}$  and  $x \in \mathfrak{R}$ ,  $P(Y \in A|x) = \int_A f(y|x) dy = \int_A f_Y(y) dy = P(Y \in A)$ . The knowledge that  $X = x$  gives us no additional information about  $Y$ .

~~Oct 14th. 2021~~

**Lemma 4.2.7** Let  $(X, Y)$  be a bivariate random vector with joint pdf or pmf  $f(x, y)$ . Then  $X$  and  $Y$  are independent random variables if and only if there exist functions  $g(x)$  and  $h(y)$  such that, for every  $x \in \mathbb{R}$  and  $y \in \mathbb{R}$ ,

$$f(x, y) = g(x)h(y).$$

**Proof:** The “only if” part is proved by defining  $g(x) = f_X(x)$  and  $h(y) = f_Y(y)$  and using (4.2.1). To prove the “if” part for continuous random variables, suppose that  $f(x, y) = g(x)h(y)$ . Define

$$\int_{-\infty}^{\infty} g(x) dx = c \quad \text{and} \quad \int_{-\infty}^{\infty} h(y) dy = d,$$

where the constants  $c$  and  $d$  satisfy

$$\begin{aligned} cd &= \left( \int_{-\infty}^{\infty} g(x) dx \right) \left( \int_{-\infty}^{\infty} h(y) dy \right) \\ (4.2.2) \quad &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x)h(y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy \\ &= 1. \end{aligned} \quad (f(x, y) \text{ is a joint pdf})$$

Furthermore, the marginal pdfs are given by

(4.2.3)

$$f_X(x) = \int_{-\infty}^{\infty} g(x)h(y) dy = g(x)d \quad \text{and} \quad f_Y(y) = \int_{-\infty}^{\infty} g(x)h(y) dx = h(y)c.$$

Thus, using (4.2.2) and (4.2.3), we have

$$f(x, y) = g(x)h(y) = g(x)h(y)cd = f_X(x)f_Y(y),$$

showing that  $X$  and  $Y$  are independent. Replacing integrals with sums proves the lemma for discrete random vectors.  $\square$

**Theorem 4.2.10** Let  $X$  and  $Y$  be independent random variables.

- For any  $A \subset \mathbb{R}$  and  $B \subset \mathbb{R}$ ,  $P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$ ; that is, the events  $\{X \in A\}$  and  $\{Y \in B\}$  are independent events.
- Let  $g(x)$  be a function only of  $x$  and  $h(y)$  be a function only of  $y$ . Then

$$\mathbf{E}(g(X)h(Y)) = (\mathbf{E}g(X))(\mathbf{E}h(Y)).$$

**Proof:** For continuous random variables, part (b) is proved by noting that

$$\mathbf{E}(g(X)h(Y)) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x)h(y)f(x, y) dx dy$$

$$\begin{aligned}
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x)h(y)f_X(x)f_Y(y) dx dy \quad (\text{by (4.2.1)}) \\
&= \int_{-\infty}^{\infty} h(y)f_Y(y) \int_{-\infty}^{\infty} g(x)f_X(x) dx dy \\
&= \left( \int_{-\infty}^{\infty} g(x)f_X(x) dx \right) \left( \int_{-\infty}^{\infty} h(y)f_Y(y) dy \right) \\
&= (\mathbb{E}g(X))(\mathbb{E}h(Y)).
\end{aligned}$$

The result for discrete random variables is proved by replacing integrals by sums. Part (a) can be proved by a series of steps similar to those above or by the following argument. Let  $g(x)$  be the indicator function of the set  $A$ . Let  $h(y)$  be the indicator function of the set  $B$ . Note that  $g(x)h(y)$  is the indicator function of the set  $C \subset \mathbb{R}^2$  defined by  $C = \{(x, y) : x \in A, y \in B\}$ . Also note that for an indicator function such as  $g(x)$ ,  $\mathbb{E}g(X) = P(X \in A)$ . Thus using the expectation equality just proved, we have

$$\mathbb{E}(g(x)) = \mathbb{E}(I_A(x)) = P(X \in A)$$

$$P(X \in A, Y \in B) = P((X, Y) \in C) = \mathbb{E}(g(X)h(Y))$$

$$= (\mathbb{E}g(X))(\mathbb{E}h(Y)) = P(X \in A)P(Y \in B). \quad \square$$

**Example 4.2.11 (Expectations of independent variables)** Let  $X$  and  $Y$  be independent exponential(1) random variables. From Theorem 4.2.10 we have

$$P(X \geq 4, Y < 3) = P(X \geq 4)P(Y < 3) = e^{-4}(1 - e^{-3}).$$

Letting  $g(x) = x^2$  and  $h(y) = y$ , we see that

$$\mathbb{E}(X^2Y) = (\mathbb{E}X^2)(\mathbb{E}Y) = (\text{Var } X + (\mathbb{E}X)^2)\mathbb{E}Y = (1 + 1^2)1 = 2. \quad \parallel$$

**Theorem 4.2.12** Let  $X$  and  $Y$  be independent random variables with moment generating functions  $M_X(t)$  and  $M_Y(t)$ . Then the moment generating function of the random variable  $Z = X + Y$  is given by

$$M_Z(t) = M_X(t)M_Y(t).$$

**Proof:** Using the definition of the mgf and Theorem 4.2.10, we have

$$M_Z(t) = \mathbb{E}e^{tZ} = \mathbb{E}e^{t(X+Y)} = \mathbb{E}(e^{tX}e^{tY}) = (\mathbb{E}e^{tX})(\mathbb{E}e^{tY}) = M_X(t)M_Y(t). \quad \square$$

**Example 4.2.13 (Mgf of a sum of normal variables)** Sometimes Theorem 4.2.12 can be used to easily derive the distribution of  $Z$  from knowledge of the distribution of  $X$  and  $Y$ . For example, let  $X \sim n(\mu, \sigma^2)$  and  $Y \sim n(\gamma, \tau^2)$  be independent normal random variables. From Exercise 2.33, the mgfs of  $X$  and  $Y$  are

$$M_X(t) = \exp(\mu t + \sigma^2 t^2/2) \quad \text{and} \quad M_Y(t) = \exp(\gamma t + \tau^2 t^2/2).$$

Thus, from Theorem 4.2.12, the mgf of  $Z = X + Y$  is

$$M_Z(t) = M_X(t)M_Y(t) = \exp((\mu + \gamma)t + (\sigma^2 + \tau^2)t^2/2).$$

This is the mgf of a normal random variable with mean  $\mu + \gamma$  and variance  $\sigma^2 + \tau^2$ . This result is important enough to be stated as a theorem. ||

**Theorem 4.2.14** Let  $X \sim n(\mu, \sigma^2)$  and  $Y \sim n(\gamma, \tau^2)$  be independent normal random variables. Then the random variable  $Z = X + Y$  has a  $n(\mu + \gamma, \sigma^2 + \tau^2)$  distribution.

~~Oct 19th 2021~~

**Example 4.3.1 (Distribution of the sum of Poisson variables)** Let  $X$  and  $Y$  be independent Poisson random variables with parameters  $\theta$  and  $\lambda$ , respectively. Thus the joint pmf of  $(X, Y)$  is

$$f_{X,Y}(x, y) = \frac{\theta^x e^{-\theta}}{x!} \frac{\lambda^y e^{-\lambda}}{y!}, \quad x = 0, 1, 2, \dots, y = 0, 1, 2, \dots$$

The set  $\mathcal{A}$  is  $\{(x, y) : x = 0, 1, 2, \dots \text{ and } y = 0, 1, 2, \dots\}$ . Now define  $U = X + Y$  and  $V = Y$ . That is,  $g_1(x, y) = x + y$  and  $g_2(x, y) = y$ . We will describe the set  $\mathcal{B}$ , the set of possible  $(u, v)$  values. The possible values for  $v$  are the nonnegative integers. The variable  $v = y$  and thus has the same set of possible values. For a given value of  $v$ ,  $u = x + y = x + v$  must be an integer greater than or equal to  $v$  since  $x$  is a nonnegative integer. The set of all possible  $(u, v)$  values is thus given by  $\mathcal{B} = \{(u, v) : v = 0, 1, 2, \dots \text{ and } u = v, v+1, v+2, \dots\}$ . For any  $(u, v) \in \mathcal{B}$ , the only  $(x, y)$  value satisfying  $x + y = u$  and  $y = v$  is  $x = u - v$  and  $y = v$ . Thus, in this example,  $A_{uv}$  always consists of only the single point  $(u - v, v)$ . From (4.3.1) we thus obtain the joint pmf of  $(U, V)$  as

$$f_{U,V}(u, v) = f_{X,Y}(u - v, v) = \frac{\theta^{u-v} e^{-\theta}}{(u-v)!} \frac{\lambda^v e^{-\lambda}}{v!}, \quad v = 0, 1, 2, \dots, \quad u = v, v+1, v+2, \dots$$

In this example it is interesting to compute the marginal pmf of  $U$ . For any fixed nonnegative integer  $u$ ,  $f_{U,V}(u, v) > 0$  only for  $v = 0, 1, \dots, u$ . This gives the set of  $v$  values to sum over to obtain the marginal pmf of  $U$ . It is

$$f_U(u) = \sum_{v=0}^u \frac{\theta^{u-v} e^{-\theta}}{(u-v)!} \frac{\lambda^v e^{-\lambda}}{v!} = e^{-(\theta+\lambda)} \sum_{v=0}^u \frac{\theta^{u-v}}{(u-v)!} \frac{\lambda^v}{v!}, \quad u = 0, 1, 2, \dots$$

This can be simplified by noting that, if we multiply and divide each term by  $u!$ , we can use the Binomial Theorem to obtain

$$f_U(u) = \frac{e^{-(\theta+\lambda)}}{u!} \sum_{v=0}^u \binom{u}{v} \lambda^v \theta^{u-v} = \frac{e^{-(\theta+\lambda)}}{u!} (\theta + \lambda)^u, \quad u = 0, 1, 2, \dots$$

This is the pmf of a Poisson random variable with parameter  $\theta + \lambda$ . This result is significant enough to be stated as a theorem. ||

**Theorem 4.3.2** If  $X \sim \text{Poisson}(\theta)$  and  $Y \sim \text{Poisson}(\lambda)$  and  $X$  and  $Y$  are independent, then  $X + Y \sim \text{Poisson}(\theta + \lambda)$ .

$$M_X(t) = e^{\theta(e^t - 1)}$$

If  $(X, Y)$  is a continuous random vector with joint pdf  $f_{X,Y}(x, y)$ , then the joint pdf of  $(U, V)$  can be expressed in terms of  $f_{X,Y}(x, y)$  in a manner analogous to (2.1.8). As before,  $\mathcal{A} = \{(x, y) : f_{X,Y}(x, y) > 0\}$  and  $\mathcal{B} = \{(u, v) : u = g_1(x, y) \text{ and } v = g_2(x, y) \text{ for some } (x, y) \in \mathcal{A}\}$ . The joint pdf  $f_{U,V}(u, v)$  will be positive on the set  $\mathcal{B}$ . For the simplest version of this result we assume that the transformation  $u = g_1(x, y)$  and  $v = g_2(x, y)$  defines a one-to-one transformation of  $\mathcal{A}$  onto  $\mathcal{B}$ . The transformation is onto because of the definition of  $\mathcal{B}$ . We are assuming that for each  $(u, v) \in \mathcal{B}$  there is only one  $(x, y) \in \mathcal{A}$  such that  $(u, v) = (g_1(x, y), g_2(x, y))$ . For such a one-to-one, onto transformation, we can solve the equations  $u = g_1(x, y)$  and  $v = g_2(x, y)$  for  $x$  and  $y$  in terms of  $u$  and  $v$ . We will denote this inverse transformation by  $x = h_1(u, v)$  and  $y = h_2(u, v)$ . The role played by a derivative in the univariate case is now played by a quantity called the *Jacobian of the transformation*. This function of  $(u, v)$ , denoted by  $J$ , is the *determinant of a matrix* of partial derivatives. It is defined by

$$J = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} = \frac{\partial x}{\partial u} \frac{\partial y}{\partial v} - \frac{\partial y}{\partial u} \frac{\partial x}{\partial v},$$

where

$$\frac{\partial x}{\partial u} = \frac{\partial h_1(u, v)}{\partial u}, \quad \frac{\partial x}{\partial v} = \frac{\partial h_1(u, v)}{\partial v}, \quad \frac{\partial y}{\partial u} = \frac{\partial h_2(u, v)}{\partial u}, \quad \text{and} \quad \frac{\partial y}{\partial v} = \frac{\partial h_2(u, v)}{\partial v}.$$

We assume that  $J$  is not identically 0 on  $\mathcal{B}$ . Then the joint pdf of  $(U, V)$  is 0 outside the set  $\mathcal{B}$  and on the set  $\mathcal{B}$  is given by

$$(4.3.2) \quad f_{U,V}(u, v) = f_{X,Y}(h_1(u, v), h_2(u, v))|J|,$$

**Example 4.3.3 (Distribution of the product of beta variables)** Let  $X \sim \text{beta}(\alpha, \beta)$  and  $Y \sim \text{beta}(\alpha + \beta, \gamma)$  be independent random variables. The joint pdf of  $(X, Y)$  is

$$f_{X,Y}(x, y) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} \frac{\Gamma(\alpha + \beta + \gamma)}{\Gamma(\alpha + \beta)\Gamma(\gamma)} y^{\alpha+\beta-1} (1-y)^{\gamma-1},$$

$$0 < x < 1, \quad 0 < y < 1.$$

Consider the transformation  $U = XY$  and  $V = X$ . The set of possible values for  $V$  is  $0 < v < 1$  since  $V = X$ . For a fixed value of  $V = v$ ,  $U$  must be between 0 and  $v$  since  $X = V = v$  and  $Y$  is between 0 and 1. Thus, this transformation maps the set  $\mathcal{A}$  onto the set  $\mathcal{B} = \{(u, v) : 0 < u < v < 1\}$ .

For any  $(u, v) \in \mathcal{B}$ , the equations  $u = xy$  and  $v = x$  can be uniquely solved for  $x = h_1(u, v) = v$  and  $y = h_2(u, v) = u/v$ . Note that if considered as a transformation defined on all of  $\mathbb{R}^2$ , this transformation is not one-to-one. Any point  $(0, y)$  is mapped into the point  $(0, 0)$ . But as a function defined only on  $\mathcal{A}$ , it is a one-to-one transformation onto  $\mathcal{B}$ . The Jacobian is given by

$$J = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} = \begin{vmatrix} 0 & 1 \\ \frac{1}{v} & -\frac{u}{v^2} \end{vmatrix} = -\frac{1}{v}.$$

Thus, from (4.3.2) we obtain the joint pdf as

$$(4.3.3) \quad f_{U,V}(u, v) = \frac{\Gamma(\alpha + \beta + \gamma)}{\Gamma(\alpha)\Gamma(\beta)\Gamma(\gamma)} v^{\alpha-1} (1-v)^{\beta-1} \left(\frac{u}{v}\right)^{\alpha+\beta-1} \left(1 - \frac{u}{v}\right)^{\gamma-1} \frac{1}{v},$$

$$0 < u < v < 1.$$

The marginal distribution of  $V = X$  is, of course, a beta( $\alpha, \beta$ ) distribution. But the distribution of  $U$  is also a beta distribution:

$$\begin{aligned} f_U(u) &= \int_u^1 f_{U,V}(u, v) dv \\ &= \frac{\Gamma(\alpha + \beta + \gamma)}{\Gamma(\alpha)\Gamma(\beta)\Gamma(\gamma)} u^{\alpha-1} \int_u^1 \left(\frac{u}{v} - u\right)^{\beta-1} \left(1 - \frac{u}{v}\right)^{\gamma-1} \left(\frac{u}{v^2}\right) dv. \end{aligned}$$

The expression (4.3.3) was used but some terms have been rearranged. Now make the univariate change of variable  $y = (u/v - u)/(1-u)$  so that  $dy = -u/[v^2(1-u)]dv$  to obtain

$$\begin{aligned} f_U(u) &= \frac{\Gamma(\alpha + \beta + \gamma)}{\Gamma(\alpha)\Gamma(\beta)\Gamma(\gamma)} u^{\alpha-1} (1-u)^{\beta+\gamma-1} \int_0^1 y^{\beta-1} (1-y)^{\gamma-1} dy \\ &= \frac{\Gamma(\alpha + \beta + \gamma)}{\Gamma(\alpha)\Gamma(\beta)\Gamma(\gamma)} u^{\alpha-1} (1-u)^{\beta+\gamma-1} \frac{\Gamma(\beta)\Gamma(\gamma)}{\Gamma(\beta + \gamma)} \\ &= \frac{\Gamma(\alpha + \beta + \gamma)}{\Gamma(\alpha)\Gamma(\beta + \gamma)} u^{\alpha-1} (1-u)^{\beta+\gamma-1}, \quad 0 < u < 1. \end{aligned}$$

To obtain the second identity we recognized the integrand as the kernel of a beta pdf and used (3.3.17). Thus we see that the marginal distribution of  $U$  is beta( $\alpha, \beta + \gamma$ ). ||

$$\text{Ex. } f_{X_1, X_2}(x_1, x_2) = \lambda e^{-\lambda x_1} \cdot \lambda e^{-\lambda x_2} \cdot I_{(0, +\infty)}(x_1) \cdot I_{(0, +\infty)}(y_2)$$

$$\text{Define } (Y_1, Y_2) = (g_1(X_1, X_2), g_2(X_1, X_2)) = (X_1 + X_2, \frac{X_1}{X_1 + X_2})$$

$$X_1 = h_1(Y_1, Y_2) = Y_1 \cdot Y_2, \quad X_2 = h_2(Y_1, Y_2) = Y_1 - Y_1 Y_2$$

$$J = \begin{vmatrix} \frac{\partial X_1}{\partial Y_1} & \frac{\partial X_1}{\partial Y_2} \\ \frac{\partial X_2}{\partial Y_1} & \frac{\partial X_2}{\partial Y_2} \end{vmatrix} = \begin{vmatrix} Y_2 & Y_1 \\ 1 - Y_2 & -Y_1 \end{vmatrix} = -Y_1 Y_2 - Y_1 (1 - Y_2) = -Y_1$$

$$\Rightarrow f_{Y_1, Y_2}(y_1, y_2) = \lambda^2 e^{-\lambda y_1} \cdot 1 - y_1 = \lambda^2 y_1 e^{-\lambda y_1}, \quad y_1 > 0, \quad 0 < y_2 < 1$$

$$\Rightarrow \textcircled{1} Y_2 \sim \text{unif}(0, 1) \quad \textcircled{2} Y_1 \sim \text{Gamma}(2, \frac{1}{\lambda}) \quad \textcircled{3} Y_1 \perp\!\!\!\perp Y_2$$

**Theorem 4.3.5** Let  $X$  and  $Y$  be independent random variables. Let  $g(x)$  be a function only of  $x$  and  $h(y)$  be a function only of  $y$ . Then the random variables  $U = g(X)$  and  $V = h(Y)$  are independent.

**Proof:** We will prove the theorem assuming  $U$  and  $V$  are continuous random variables. For any  $u \in \mathbb{R}$  and  $v \in \mathbb{R}$ , define

$$A_u = \{x : g(x) \leq u\} \quad \text{and} \quad B_v = \{y : h(y) \leq v\}.$$

Then the joint cdf of  $(U, V)$  is

$$\begin{aligned} F_{U,V}(u, v) &= P(U \leq u, V \leq v) && \text{(definition of cdf)} \\ &= P(X \in A_u, Y \in B_v) && \text{(definition of } U \text{ and } V\text{)} \\ &= P(X \in A_u)P(Y \in B_v). && \text{(Theorem 4.2.10)} \end{aligned}$$

The joint pdf of  $(U, V)$  is

$$\begin{aligned} f_{U,V}(u, v) &= \frac{\partial^2}{\partial u \partial v} F_{U,V}(u, v) && \text{(by (4.1.4))} \\ &= \left( \frac{d}{du} P(X \in A_u) \right) \left( \frac{d}{dv} P(Y \in B_v) \right), \end{aligned}$$

where, as the notation indicates, the first factor is a function only of  $u$  and the second factor is a function only of  $v$ . Hence, by Lemma 4.2.7,  $U$  and  $V$  are independent.  $\square$

Of course, in many situations, the transformation of interest is not one-to-one. Just as Theorem 2.1.8 generalized the univariate method to many-to-one functions, the same can be done here. As before,  $\mathcal{A} = \{(x, y) : f_{X,Y}(x, y) > 0\}$ . Suppose  $A_0, A_1, \dots, A_k$  form a partition of  $\mathcal{A}$  with these properties. The set  $A_0$ , which may be empty, satisfies  $P((X, Y) \in A_0) = 0$ . The transformation  $U = g_1(X, Y)$  and  $V = g_2(X, Y)$  is a one-to-one transformation from  $A_i$  onto  $\mathcal{B}$  for each  $i = 1, 2, \dots, k$ . Then for each  $i$ , the inverse functions from  $\mathcal{B}$  to  $A_i$  can be found. Denote the  $i$ th inverse by  $x = h_{1i}(u, v)$  and  $y = h_{2i}(u, v)$ . This  $i$ th inverse gives, for  $(u, v) \in \mathcal{B}$ , the unique  $(x, y) \in A_i$  such that  $(u, v) = (g_1(x, y), g_2(x, y))$ . Let  $J_i$  denote the Jacobian computed from the  $i$ th inverse. Then assuming that these Jacobians do not vanish identically on  $\mathcal{B}$ , we have the following representation of the joint pdf,  $f_{U,V}(u, v)$ :

$$(4.3.6) \quad f_{U,V}(u, v) = \sum_{i=1}^k f_{X,Y}(h_{1i}(u, v), h_{2i}(u, v)) |J_i|.$$

**Example 4.3.6 (Distribution of the ratio of normal variables)** Let  $X$  and  $Y$  be independent  $n(0, 1)$  random variables. Consider the transformation  $U = X/Y$  and  $V = |Y|$ . ( $U$  and  $V$  can be defined to be any value, say  $(1, 1)$ , if  $Y = 0$  since  $P(Y = 0) = 0$ .) This transformation is not one-to-one since the points  $(x, y)$  and  $(-x, -y)$  are both mapped into the same  $(u, v)$  point. But if we restrict consideration to either positive or negative values of  $y$ , then the transformation is one-to-one. In the above notation, let

$$A_1 = \{(x, y) : y > 0\}, \quad A_2 = \{(x, y) : y < 0\}, \quad \text{and} \quad A_0 = \{(x, y) : y = 0\}.$$

$A_0, A_1$ , and  $A_2$  form a partition of  $\mathcal{A} = \mathbb{R}^2$  and  $P((X, Y) \in A_0) = P(Y = 0) = 0$ . For either  $A_1$  or  $A_2$ , if  $(x, y) \in A_i$ ,  $v = |y| > 0$ , and for a fixed value of  $v = |y|$ ,  $u = x/y$  can be any real number since  $x$  can be any real number. Thus,  $\mathcal{B} = \{(u, v) : v > 0\}$  is the image of both  $A_1$  and  $A_2$  under the transformation. Furthermore, the inverse transformations from  $\mathcal{B}$  to  $A_1$  and  $\mathcal{B}$  to  $A_2$  are given by  $x = h_{11}(u, v) = uv$ ,  $y = h_{21}(u, v) = v$ , and  $x = h_{12}(u, v) = -uv$ ,  $y = h_{22}(u, v) = -v$ . Note that the first inverse gives positive values of  $y$  and the second gives negative values of  $y$ . The Jacobians from the two inverses are  $J_1 = J_2 = v$ . Using

$$f_{X,Y}(x, y) = \frac{1}{2\pi} e^{-x^2/2} e^{-y^2/2},$$

from (4.3.6) we obtain

$$\begin{aligned} f_{U,V}(u, v) &= \frac{1}{2\pi} e^{-(uv)^2/2} e^{-v^2/2} |v| + \frac{1}{2\pi} e^{-(-uv)^2/2} e^{-(v)^2/2} |v| \\ &= \frac{v}{\pi} e^{-(u^2+1)v^2/2}, \quad -\infty < u < \infty, \quad 0 < v < \infty. \end{aligned}$$

From this the marginal pdf of  $U$  can be computed to be

$$\begin{aligned} f_U(u) &= \int_0^\infty \frac{v}{\pi} e^{-(u^2+1)v^2/2} dv \\ &= \frac{1}{2\pi} \int_0^\infty e^{-(u^2+1)z/2} dz \quad z = v^2 \quad (\text{change of variable}) \\ &= \frac{1}{2\pi} \frac{2}{(u^2+1)} \quad \left( \begin{array}{l} \text{integrand is kernel of} \\ \text{exponential } (\beta = 2/(u^2+1)) \text{ pdf} \end{array} \right) \\ &= \frac{1}{\pi(u^2+1)}, \quad -\infty < u < \infty. \end{aligned}$$

So we see that the ratio of two independent standard normal random variables is a Cauchy random variable. (See Exercise 4.28 for more relationships between normal and Cauchy random variables.) ||

~~Oct 21st 2021~~

#### 4.4 Hierarchical Models and Mixture Distributions

**Example 4.4.1 (Binomial-Poisson hierarchy)** Perhaps the most classic hierarchical model is the following. An insect lays a large number of eggs, each surviving with probability  $p$ . On the average, how many eggs will survive?

The “large number” of eggs laid is a random variable, often taken to be  $\text{Poisson}(\lambda)$ . Furthermore, if we assume that each egg’s survival is independent, then we have Bernoulli trials. Therefore, if we let  $X = \text{number of survivors}$  and  $Y = \text{number of eggs laid}$ , we have

$$X|Y \sim \text{binomial}(Y, p),$$

$$Y \sim \text{Poisson}(\lambda),$$

a hierarchical model. (Recall that we use notation such as  $X|Y \sim \text{binomial}(Y, p)$  to mean that the conditional distribution of  $X$  given  $Y = y$  is  $\text{binomial}(y, p)$ .) ||

The advantage of the hierarchy is that complicated processes may be modeled by a sequence of relatively simple models placed in a hierarchy. Also, dealing with the hierarchy is no more difficult than dealing with conditional and marginal distributions.

**Example 4.4.2 (Continuation of Example 4.4.1)** The random variable of interest,  $X = \text{number of survivors}$ , has the distribution given by

$$\begin{aligned} P(X = x) &= \sum_{y=0}^{\infty} P(X = x, Y = y) \\ &= \sum_{y=0}^{\infty} P(X = x|Y = y)P(Y = y) && \left( \begin{array}{l} \text{definition of} \\ \text{conditional probability} \end{array} \right) \\ &= \sum_{y=x}^{\infty} \left[ \binom{y}{x} p^x (1-p)^{y-x} \right] \left[ \frac{e^{-\lambda} \lambda^y}{y!} \right], && \left( \begin{array}{l} \text{conditional probability} \\ \text{is 0 if } y < x \end{array} \right) \end{aligned}$$

since  $X|Y = y$  is  $\text{binomial}(y, p)$  and  $Y$  is  $\text{Poisson}(\lambda)$ . If we now simplify this last expression, canceling what we can and multiplying by  $\lambda^x/\lambda^x$ , we get

$$\begin{aligned} P(X = x) &= \frac{(\lambda p)^x e^{-\lambda}}{x!} \sum_{y=x}^{\infty} \frac{((1-p)\lambda)^{y-x}}{(y-x)!} \\ &= \frac{(\lambda p)^x e^{-\lambda}}{x!} \sum_{t=0}^{\infty} \frac{((1-p)\lambda)^t}{t!} (t = y - x) \\ &= \frac{(\lambda p)^x e^{-\lambda}}{x!} e^{(1-p)\lambda} && \left( \begin{array}{l} \text{sum is a kernel for} \\ \text{a Poisson distribution} \end{array} \right) \\ &= \frac{(\lambda p)^x}{x!} e^{-\lambda p}, \end{aligned}$$

so  $X \sim \text{Poisson}(\lambda p)$ . Thus, any marginal inference on  $X$  is with respect to a Poisson( $\lambda p$ ) distribution, with  $Y$  playing no part at all. Introducing  $Y$  in the hierarchy was mainly to aid our understanding of the model. There was an added bonus in that the parameter of the distribution of  $X$  is the product of two parameters, each relatively simple to understand.

The answer to the original question is now easy to compute:

$$EX = \lambda p,$$

so, on the average,  $\lambda p$  eggs will survive. If we were interested only in this mean and did not need the distribution, we could have used properties of conditional expectations.

**Theorem 4.4.3** *If  $X$  and  $Y$  are any two random variables, then*

$$(4.4.1) \quad EX = E(E(X|Y)), \quad E[g(x)] = E[E[g(x)|Y]]$$

*provided that the expectations exist.*

**Proof:** Let  $f(x, y)$  denote the joint pdf of  $X$  and  $Y$ . By definition, we have

$$(4.4.2) \quad EX = \int \int xf(x, y) dx dy = \int \left[ \int xf(x|y) dx \right] f_Y(y) dy,$$

where  $f(x|y)$  and  $f_Y(y)$  are the conditional pdf of  $X$  given  $Y = y$  and the marginal pdf of  $Y$ , respectively. But now notice that the inner integral in (4.4.2) is the conditional expectation  $E(X|y)$ , and we have

$$EX = \int E(X|y) f_Y(y) dy = E(E(X|Y)),$$

as desired. Replace integrals by sums to prove the discrete case.  $\square$

We can now easily compute the expected number of survivors in Example 4.4.1. From Theorem 4.4.3 we have

$$\begin{aligned} EX &= E(E(X|Y)) \\ &= E(pY) \quad (\text{since } X|Y \sim \text{binomial}(Y, p)) \\ &= p\lambda. \quad (\text{since } Y \sim \text{Poisson}(\lambda)) \end{aligned}$$

**Example 4.4.5 (Generalization of Example 4.4.1)** Consider a generalization of Example 4.4.1, where instead of one mother insect there are a large number of mothers and one mother is chosen at random. We are still interested in knowing the average number of survivors, but it is no longer clear that the number of eggs laid follows the same Poisson distribution for each mother. The following three-stage hierarchy may be more appropriate. Let  $X$  = number of survivors in a litter; then

$$X|Y \sim \text{binomial}(Y, p),$$

$$Y|\Lambda \sim \text{Poisson}(\Lambda),$$

$$\Lambda \sim \text{exponential}(\beta),$$

where the last stage of the hierarchy accounts for the variability across different mothers.

The mean of  $X$  can easily be calculated as

$$\begin{aligned}
 E(X) &= E(E(X|Y)) \\
 &= E(pY) && (\text{as before}) \\
 &= E(E(pY|\Lambda)) \\
 &= E(p\Lambda) \\
 &= p\beta, && (\text{exponential expectation})
 \end{aligned}$$

completing the calculation. ||

Note that this three-stage model can also be thought of as a two-stage hierarchy by combining the last two stages. If  $Y|\Lambda \sim \text{Poisson}(\Lambda)$  and  $\Lambda \sim \text{exponential}(\beta)$ , then

$$\begin{aligned}
 P(Y = y) &= P(Y = y, 0 < \Lambda < \infty) \\
 &= \int_0^\infty f(y, \lambda) d\lambda \\
 &= \int_0^\infty f(y|\lambda)f(\lambda) d\lambda \\
 &= \int_0^\infty \left[ \frac{e^{-\lambda}\lambda^y}{y!} \right] \frac{1}{\beta} e^{-\lambda/\beta} d\lambda \\
 &= \frac{1}{\beta y!} \int_0^\infty \lambda^y e^{-\lambda(1+\beta^{-1})} d\lambda && \begin{pmatrix} \text{gamma} \\ \text{pdf kernel} \end{pmatrix} \\
 &= \frac{1}{\beta y!} \Gamma(y+1) \left( \frac{1}{1+\beta^{-1}} \right)^{y+1} \\
 &= \frac{1}{(1+\beta)} \left( \frac{1}{1+\beta^{-1}} \right)^y.
 \end{aligned}$$

This expression for the pmf of  $Y$  is the form (3.2.10) of the negative binomial pmf. Therefore, our three-stage hierarchy in Example 4.4.5 is equivalent to the two-stage hierarchy

$$X|Y \sim \text{binomial}(Y, p),$$

$$Y \sim \text{negative binomial} \left( p = \frac{1}{1+\beta}, r = 1 \right).$$

**Theorem 4.4.7 (Conditional variance identity)** *For any two random variables  $X$  and  $Y$ ,*

$$(4.4.4) \quad \text{Var } X = E(\text{Var}(X|Y)) + \text{Var}(E(X|Y)),$$

*provided that the expectations exist.*

**Proof:** By definition, we have

$$\text{Var } X = E([X - EX]^2) = E([X - E(X|Y) + E(X|Y) - EX]^2),$$

where in the last step we have added and subtracted  $E(X|Y)$ . Expanding the square in this last expectation now gives

$$(4.4.5) \quad \begin{aligned} \text{Var } X &= E([X - E(X|Y)]^2) + E([E(X|Y) - EX]^2) \\ &\quad + 2E([X - E(X|Y)][E(X|Y) - EX]). \end{aligned}$$

The last term in this expression is equal to 0, however, which can easily be seen by iterating the expectation:

$$(4.4.6) \quad E([X - E(X|Y)][E(X|Y) - EX]) = E(E\{[X - E(X|Y)][E(X|Y) - EX]|Y\}).$$

In the conditional distribution  $X|Y$ ,  $X$  is the random variable. So in the expression

$$E\{[X - E(X|Y)][E(X|Y) - EX]|Y\},$$

$E(X|Y)$  and  $EX$  are constants. Thus,

$$\begin{aligned} E\{[X - E(X|Y)][E(X|Y) - EX]|Y\} &= (E(X|Y) - EX)(E\{[X - E(X|Y)]|Y\}) \\ &= (E(X|Y) - EX)(E(X|Y) - E(X|Y)) \\ &= (E(X|Y) - EX)(0) \\ &= 0. \end{aligned}$$

Thus, from (4.4.6), we have that  $E((X - E(X|Y))(E(X|Y) - EX)) = E(0) = 0$ . Referring back to equation (4.4.5), we see that

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

and

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

establishing (4.4.4). □

**Example 4.4.6 (Beta-binomial hierarchy)** One generalization of the binomial distribution is to allow the success probability to vary according to a distribution. A standard model for this situation is

$$X|P \sim \text{binomial}(P), \quad i = 1, \dots, n,$$

$$P \sim \text{beta}(\alpha, \beta).$$

By iterating the expectation, we calculate the mean of  $X$  as

$$EX = E[E(X|P)] = E[nP] = n \frac{\alpha}{\alpha + \beta}. \quad \|$$

**Example 4.4.8 (Continuation of Example 4.4.6)** To calculate the variance of  $X$ , we have from (4.4.4),

$$\text{Var } X = \text{Var}(E(X|P)) + E(\text{Var}(X|P)).$$

Now  $E(X|P) = nP$ , and since  $P \sim \text{beta}(\alpha, \beta)$ ,

$$\text{Var}(E(X|P)) = \text{Var}(nP) = n^2 \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}.$$

Also, since  $X|P$  is binomial( $n, P$ ),  $\text{Var}(X|P) = nP(1 - P)$ . We then have

$$E[\text{Var}(X|P)] = nE[P(1 - P)] = n \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 p(1 - p)p^{\alpha-1}(1 - p)^{\beta-1} dp.$$

Notice that the integrand is the kernel of another beta pdf (with parameters  $\alpha + 1$  and  $\beta + 1$ ) so

$$E(\text{Var}(X|P)) = n \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \left[ \frac{\Gamma(\alpha + 1)\Gamma(\beta + 1)}{\Gamma(\alpha + \beta + 2)} \right] = n \frac{\alpha\beta}{(\alpha + \beta)(\alpha + \beta + 1)}.$$

Adding together the two pieces and simplifying, we get

$$\text{Var } X = n \frac{\alpha\beta(\alpha + \beta + n)}{(\alpha + \beta)^2(\alpha + \beta + 1)}. \quad \|$$

~~Oct. 26th 2021~~

## 4.5 Covariance and Correlation

Throughout this section we will frequently be referring to the mean and variance of  $X$  and the mean and variance of  $Y$ . For these we will use the notation  $\text{E}X = \mu_X$ ,  $\text{E}Y = \mu_Y$ ,  $\text{Var } X = \sigma_X^2$ , and  $\text{Var } Y = \sigma_Y^2$ . We will assume throughout that  $0 < \sigma_X^2 < \infty$  and  $0 < \sigma_Y^2 < \infty$ .

**Definition 4.5.1** The covariance of  $X$  and  $Y$  is the number defined by

$$\text{Cov}(X, Y) = \text{E}((X - \mu_X)(Y - \mu_Y)).$$

**Definition 4.5.2** The correlation of  $X$  and  $Y$  is the number defined by

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}.$$

The value  $\rho_{XY}$  is also called the correlation coefficient.

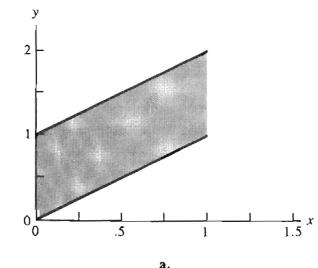
**Theorem 4.5.3** For any random variables  $X$  and  $Y$ ,

$$\text{Cov}(X, Y) = \text{E}XY - \mu_X \mu_Y.$$

**Proof:**  $\text{Cov}(X, Y) = \text{E}(X - \mu_X)(Y - \mu_Y)$   
 $= \text{E}(XY - \mu_X Y - \mu_Y X + \mu_X \mu_Y)$  (expanding the product)  
 $= EXY - \mu_X EY - \mu_Y EX + \mu_X \mu_Y$  ( $\mu_X$  and  $\mu_Y$  are constants)  
 $= EXY - \mu_X \mu_Y - \mu_Y \mu_X + \mu_X \mu_Y$   
 $= EXY - \mu_X \mu_Y.$  □

**Example 4.5.4 (Correlation-I)** Let the joint pdf of  $(X, Y)$  be  $f(x, y) = 1$ ,  $0 < x < 1$ ,  $x < y < x+1$ . See Figure 4.5.1. The marginal distribution of  $X$  is uniform( $0, 1$ ) so  $\mu_X = \frac{1}{2}$  and  $\sigma_X^2 = \frac{1}{12}$ . The marginal pdf of  $Y$  is  $f_Y(y) = y$ ,  $0 < y < 1$ , and  $f_Y(y) = 2 - y$ ,  $1 \leq y < 2$ , with  $\mu_Y = 1$  and  $\sigma_Y^2 = \frac{1}{6}$ . We also have

$$\begin{aligned} \text{EXY} &= \int_0^1 \int_x^{x+1} xy \, dy \, dx = \int_0^1 \frac{1}{2} xy^2 \Big|_x^{x+1} \, dx \\ &= \int_0^1 \left( x^2 + \frac{1}{2}x \right) \, dx = \frac{7}{12}. \end{aligned}$$



Using Theorem 4.5.3, we have  $\text{Cov}(X, Y) = \frac{7}{12} - (\frac{1}{2})(1) = \frac{1}{12}$ . The correlation is

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{1/12}{\sqrt{1/12} \sqrt{1/6}} = \frac{1}{\sqrt{2}}.$$

**Theorem 4.5.5** If  $X$  and  $Y$  are independent random variables, then  $\text{Cov}(X, Y) = 0$  and  $\rho_{XY} = 0$ .

**Proof:** Since  $X$  and  $Y$  are independent, from Theorem 4.2.10 we have  $\text{E}XY = (\text{E}X)(\text{E}Y)$ . Thus

$$\text{Cov}(X, Y) = \text{E}XY - (\text{E}X)(\text{E}Y) = (\text{E}X)(\text{E}Y) - (\text{E}X)(\text{E}Y) = 0$$

and

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{0}{\sigma_X \sigma_Y} = 0. \quad \square$$

Thus, the values  $\text{Cov}(X, Y) = \rho_{XY} = 0$  in some sense indicate that there is no relationship between  $X$  and  $Y$ . It is important to note, however, that Theorem 4.5.5 does *not* say that if  $\text{Cov}(X, Y) = 0$ , then  $X$  and  $Y$  are independent. For example, if  $X \sim f(x - \theta)$ , symmetric around 0 with  $\text{E}X = \theta$ , and  $Y$  is the indicator function  $Y = I(|X - \theta| < 2)$ , then  $X$  and  $Y$  are obviously not independent. However,

$$\text{E}(XY) = \int_{-\infty}^{\infty} xI(|x - \theta| < 2)f(x - \theta) dx = \int_{-2}^2 (t + \theta)f(t) dt = \theta \int_{-2}^2 f(t) dt = \text{E}X\text{E}Y,$$

where we used the fact that, by symmetry,  $\int_{-2}^2 tf(t) dt = 0$ . So it is easy to find uncorrelated, dependent random variables.

**Theorem 4.5.6** If  $X$  and  $Y$  are any two random variables and  $a$  and  $b$  are any two constants, then

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

If  $X$  and  $Y$  are independent random variables, then

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

**Proof:** The mean of  $aX + bY$  is  $\text{E}(aX + bY) = a\text{E}X + b\text{E}Y = a\mu_X + b\mu_Y$ . Thus

$$\begin{aligned} \text{Var}(aX + bY) &= \text{E}((aX + bY) - (a\mu_X + b\mu_Y))^2 \\ &= \text{E}(a(X - \mu_X) + b(Y - \mu_Y))^2 \\ &= \text{E}(a^2(X - \mu_X)^2 + b^2(Y - \mu_Y)^2 + 2ab(X - \mu_X)(Y - \mu_Y)) \\ &= a^2 \text{E}(X - \mu_X)^2 + b^2 \text{E}(Y - \mu_Y)^2 + 2ab \text{E}(X - \mu_X)(Y - \mu_Y) \\ &= a^2 \text{Var } X + b^2 \text{Var } Y + 2ab \text{Cov}(X, Y). \end{aligned}$$

If  $X$  and  $Y$  are independent, then, from Theorem 4.5.5,  $\text{Cov}(X, Y) = 0$  and the second equality is immediate from the first.  $\square$

**Theorem 4.5.7** For any random variables  $X$  and  $Y$ ,

- $-1 \leq \rho_{XY} \leq 1$ .
- $|\rho_{XY}| = 1$  if and only if there exist numbers  $a \neq 0$  and  $b$  such that  $P(Y = aX + b) = 1$ . If  $\rho_{XY} = 1$ , then  $a > 0$ , and if  $\rho_{XY} = -1$ , then  $a < 0$ .

**Proof:** Consider the function  $h(t)$  defined by

$$h(t) = E((X - \mu_X)t + (Y - \mu_Y))^2.$$

Expanding this expression, we obtain

$$\begin{aligned} h(t) &= t^2 E(X - \mu_X)^2 + 2t E(X - \mu_X)(Y - \mu_Y) + E(Y - \mu_Y)^2 \\ &= t^2 \sigma_X^2 + 2t \text{Cov}(X, Y) + \sigma_Y^2. \end{aligned}$$

This quadratic function of  $t$  is greater than or equal to 0 for all values of  $t$  since it is the expected value of a nonnegative random variable. Thus, this quadratic function can have at most one real root and thus must have a nonpositive discriminant. That is,

$$(2\text{Cov}(X, Y))^2 - 4\sigma_X^2\sigma_Y^2 \leq 0.$$

This is equivalent to

$$-\sigma_X\sigma_Y \leq \text{Cov}(X, Y) \leq \sigma_X\sigma_Y.$$

Dividing by  $\sigma_X\sigma_Y$  yields

$$-1 \leq \frac{\text{Cov}(X, Y)}{\sigma_X\sigma_Y} = \rho_{XY} \leq 1.$$

Also,  $|\rho_{XY}| = 1$  if and only if the discriminant is equal to 0. That is,  $|\rho_{XY}| = 1$  if and only if  $h(t)$  has a single root. But since  $((X - \mu_X)t + (Y - \mu_Y))^2 \geq 0$ , the expected value  $h(t) = E((X - \mu_X)t + (Y - \mu_Y))^2 = 0$  if and only if

$$P([(X - \mu_X)t + (Y - \mu_Y)]^2 = 0) = 1.$$

This is equivalent to

$$P((X - \mu_X)t + (Y - \mu_Y) = 0) = 1.$$

This is  $P(Y = aX + b) = 1$  with  $a = -t$  and  $b = \mu_Xt + \mu_Y$ , where  $t$  is the root of  $h(t)$ . Using the quadratic formula, we see that this root is  $t = -\text{Cov}(X, Y)/\sigma_X^2$ . Thus  $a = -t$  has the same sign as  $\rho_{XY}$ , proving the final assertion.  $\square$

If there is a line  $y = ax + b$ , with  $a \neq 0$ , such that the values of  $(X, Y)$  have a high probability of being near this line, then the correlation between  $X$  and  $Y$  will be near 1 or -1. But if no such line exists, the correlation will be near 0. This is an intuitive notion of the linear relationship that is being measured by correlation. This idea will be illustrated further in the next two examples.

**Example 4.5.8 (Correlation-II)** This example is similar to Example 4.5.4, but we develop it differently to illustrate other model building and computational techniques. Let  $X$  have a uniform(0, 1) distribution and  $Z$  have a uniform( $0, \frac{1}{10}$ ) distribution. Suppose  $X$  and  $Z$  are independent. Let  $Y = X + Z$  and consider the random vector  $(X, Y)$ . The joint distribution of  $(X, Y)$  can be derived from the joint distribution of  $(X, Z)$  using the techniques of Section 4.3. The joint pdf of  $(X, Y)$  is

$$f(x, y) = 10, \quad 0 < x < 1, \quad x < y < x + \frac{1}{10}.$$

Rather than using the formal techniques of Section 4.3, we can justify this as follows. Given  $X = x$ ,  $Y = x + Z$ . The conditional distribution of  $Z$  given  $X = x$  is just uniform( $0, \frac{1}{10}$ ) since  $X$  and  $Z$  are independent. Thus  $x$  serves as a location parameter in the conditional distribution of  $Y$  given  $X = x$ , and this conditional distribution is just uniform( $x, x + \frac{1}{10}$ ). Multiplying this conditional pdf by the marginal pdf of  $X$  (uniform(0, 1)) yields the joint pdf above. This representation of  $Y = X + Z$  makes the computation of the covariance and correlation easy. The expected values of  $X$  and  $Y$  are  $EX = \frac{1}{2}$  and  $EY = E(X + Z) = EX + EZ = \frac{1}{2} + \frac{1}{20} = \frac{11}{20}$ , giving

$$\begin{aligned} \text{Cov}(X, Y) &= EXY - (EX)(EY) \\ &= EX(X + Z) - (EX)(E(X + Z)) \\ &= EX^2 + EXZ - (EX)^2 - (EX)(EZ) \\ &= EX^2 - (EX)^2 + (EX)(EZ) - (EX)(EZ) \quad (\text{independence of } X \text{ and } Z) \\ &= \sigma_X^2 = \frac{1}{12}. \end{aligned}$$

From Theorem 4.5.6, the variance of  $Y$  is  $\sigma_Y^2 = \text{Var}(X + Z) = \text{Var } X + \text{Var } Z = \frac{1}{12} + \frac{1}{1200}$ . Thus

$$\rho_{XY} = \frac{\frac{1}{12}}{\sqrt{\frac{1}{12} \sqrt{\frac{1}{12} + \frac{1}{1200}}}} = \sqrt{\frac{100}{101}}.$$

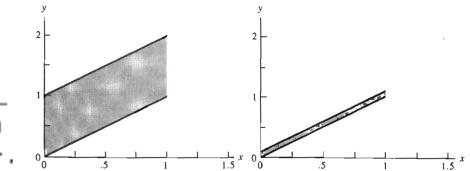


Figure 4.5.1. (a) Region where  $f(x, y) > 0$  for Example 4.5.4; (b) region where  $f(x, y) > 0$  for Example 4.5.8

This is much larger than the value of  $\rho_{XY} = 1/\sqrt{2}$  obtained in Example 4.5.4. The sets on which  $f(x, y)$  is positive for Example 4.5.4 and this example are illustrated in Figure 4.5.1. (Recall that this set is called the support of a distribution.) In each case,  $(X, Y)$  is a random point from the set. In both cases there is a linearly increasing relationship between  $X$  and  $Y$ , but the relationship is much stronger in Figure 4.5.1b. Another way to see this is by noting that in this example, the conditional distribution of  $Y$  given  $X = x$  is uniform( $x, x + \frac{1}{10}$ ). In Example 4.5.4, the conditional distribution of  $Y$  given  $X = x$  is uniform( $x, x + 1$ ). The knowledge that  $X = x$  gives us much more information about the value of  $Y$  in this model than in the one in Example 4.5.4. Hence the correlation is nearer to 1 in this example. ||

The next example illustrates that there may be a strong relationship between  $X$  and  $Y$ , but if the relationship is not linear, the correlation may be small.

**Example 4.5.9 (Correlation–III)** In this example, let  $X$  have a uniform( $-1, 1$ ) distribution and let  $Z$  have a uniform( $0, \frac{1}{10}$ ) distribution. Let  $X$  and  $Z$  be independent. Let  $Y = X^2 + Z$  and consider the random vector  $(X, Y)$ . As in Example 4.5.8, given  $X = x$ ,  $Y = x^2 + Z$  and the conditional distribution of  $Y$  given  $X = x$  is uniform( $x^2, x^2 + \frac{1}{10}$ ). The joint pdf of  $X$  and  $Y$ , the product of this conditional pdf and the marginal pdf of  $X$ , is thus

$$f(x, y) = 5, \quad -1 < x < 1, \quad x^2 < y < x^2 + \frac{1}{10}.$$

The set on which  $f(x, y) > 0$  is illustrated in Figure 4.5.2. There is a strong relationship between  $X$  and  $Y$ , as indicated by the conditional distribution of  $Y$  given  $X = x$ . But the relationship is not linear. The possible values of  $(X, Y)$  cluster around a parabola rather than a straight line. The correlation does not measure this non-linear relationship. In fact,  $\rho_{XY} = 0$ . Since  $X$  has a uniform( $-1, 1$ ) distribution,  $EX = EX^3 = 0$ , and since  $X$  and  $Z$  are independent,  $EXZ = (EX)(EZ)$ . Thus,

$$\begin{aligned}\text{Cov}(X, Y) &= E(X(X^2 + Z)) - (EX)(E(X^2 + Z)) \\ &= EX^3 + EXZ - 0E(X^2 + Z) \\ &= 0 + (EX)(EZ) = 0(EZ) = 0,\end{aligned}$$

and  $\rho_{XY} = \text{Cov}(X, Y)/(\sigma_X \sigma_Y) = 0$ .

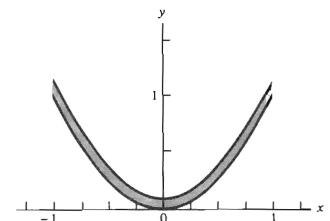


Figure 4.5.2. Region where  $f(x, y) > 0$  for Example 4.5.9

Oct. 28th 2021

NO.

Date

Ex: Let  $X_1, \dots, X_n$  be independent random variables all with  $E(X_i) = \mu$ ,  $\text{Var}(X_i) = \sigma^2$ . Let  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  be the sample mean. The random variable  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$

is called the sample variance, calculate  $\text{Var}(\bar{X})$  (b)  $E(S^2)$

$$(a) \text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) \stackrel{\text{indep}}{=} \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$
$$E(\bar{X}) = \mu$$

$$(b) \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2$$
$$= \sum_{i=1}^n (X_i - \mu)^2 + \sum_{i=1}^n (\bar{X} - \mu)^2 - 2 \sum_{i=1}^n (X_i - \mu)(\bar{X} - \mu)$$
$$\Rightarrow E\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) = \sum_{i=1}^n E((X_i - \mu)^2) + \sum_{i=1}^n E(\bar{X} - \mu)^2 - 2 \sum_{i=1}^n E(X_i - \mu)(\bar{X} - \mu)$$
$$= \sum_{i=1}^n (\sigma^2 + \frac{\sigma^2}{n} - 2\text{cov}(X_i, \bar{X})) \quad \text{cov}(X_i, \bar{X})$$
$$= \sum_{i=1}^n (\sigma^2 + \frac{\sigma^2}{n} - 2 \frac{\sigma^2}{n}) \quad \text{cov}(X_i, \frac{1}{n} \sum_{j=1}^n X_j)$$
$$= (n-1) \sigma^2 \quad = \frac{1}{n} \sum_{i=1}^{n-1} \text{cov}(X_i, X_j)$$
$$\Rightarrow E(S^2) = \frac{1}{n-1} E\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) = \sigma^2 \quad = \frac{\sigma^2}{n}$$

Ex A group of  $N$  people throw their hats into the center of the room. The hats are mixed, and each person randomly selects one

Let  $X = \#$  of people who select their own hat

Compute  $E(X)$ ,  $\text{Var}(X)$

- Number the people from 1 to  $N$ . Let:  $X_i = \begin{cases} 1 & \text{if person } i \text{ choose own hat} \\ 0 & \text{otherwise} \end{cases}$

$$\text{So } X = X_1 + X_2 + \dots + X_N \quad P(X_i = 1) = \frac{1}{N} \Rightarrow E(X_i) = 1 \cdot \frac{1}{N} + 0 \cdot (1 - \frac{1}{N}) = \frac{1}{N}$$

$$\Rightarrow E(X) = E\left(\sum_{i=1}^N X_i\right) = \sum_{i=1}^N E(X_i) = N \cdot \frac{1}{N} = 1$$

$$[\text{Var}\left(\sum_{i=1}^N X_i\right) = \sum_{i=1}^N \text{Var}(X_i) + 2 \sum_{1 < i < j} \text{cov}(X_i, X_j)]$$

NO.....

Date .....

$$X_i X_j = \begin{cases} 1 & \text{if persons } i \& j \text{ both choose their own hat} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{So } P(X_i=1, X_j=1) = P(X_j=1 | X_i=1) P(X_i=1) = \frac{1}{N(N-1)}$$

$$\Rightarrow E(X_i, X_j) = \frac{1}{N(N-1)}$$

$$\Rightarrow \text{Cov}(X_i, X_j) = E(X_i X_j) - E(X_i) E(X_j) = \frac{1}{N(N-1)} - \frac{1}{N^2} = \frac{1}{N^2(N-1)}$$

$$\text{var}(X) = N \frac{N-1}{N^2} + 2 \left(\frac{N}{2}\right) \frac{1}{N^2(N-1)} = \frac{N-1}{N} + \frac{1}{N} = 1$$

$$N \cdot \text{Var}(\text{Bernoulli}(\frac{1}{N}))$$

- Let  $N$  be a r.v. (discrete) with values in the natural numbers  $\{1, 2, \dots\}$  &  $E(N) < +\infty$ . Let  $X_1, X_2, \dots$  be a sequence of i.i.d. random variables (independent & identically distributed) that are independent of  $N$ . Also, assume that  $E(X_i) = E(X_1) < +\infty$ . Define  $Y = X_1 + X_2 + \dots + X_N$ .

What is  $E(Y)$ ?  $\text{Var}(Y)$ ? (assume  $\text{var}(N) < +\infty$ ,  $\text{var}(X_i) < +\infty$ )

$$E(Y) = E(E(Y|N)) = E(E(\sum_{i=1}^N X_i | N=n)) = E(N \cdot E(X_1)) = E(X_1) E(N)$$

$$E(Y|N) = N E(X_1) \Rightarrow \text{var}(E(Y|N)) = \text{var}(N E(X_1)) = (E(X_1))^2 \text{var}(N)$$

$$\text{var}(Y|N) = N \text{var}(X_1) \Rightarrow E(\text{Var}(Y|N)) = \text{Var}(X_1) E(N)$$

$$\Rightarrow \text{Var}(Y) = E(\text{Var}(Y|N)) + \text{Var}(E(Y|N)) = \cancel{\text{Var}(E(Y|N))}$$

$$= E(N) \text{Var}(X_1) + (E(X_1))^2 \text{Var}(N)$$

Def (bivariate Normal distribution):

Let  $M_x, M_y \in \mathbb{R}$ ,  $\sigma_x^2, \sigma_y^2 > 0$ , and  $-1 < \rho < 1$ . The bivariate normal pdf

With means  $M_x$  &  $M_y$ , variances  $\sigma_x^2$  &  $\sigma_y^2$ , and covariance  $\rho$ , is the

bivariate pdf given by

$$f(x,y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)}\left(\left(\frac{x-\mu_x}{\sigma_x}\right)^2 - 2\rho\left(\frac{x-\mu_x}{\sigma_x}\right)\left(\frac{y-\mu_y}{\sigma_y}\right) + \left(\frac{y-\mu_y}{\sigma_y}\right)^2\right)\right]$$

Properties

- ① Marginal distribution of  $X$  is Normal ( $\mu_x, \sigma_x^2$ )
- ② Marginal distribution of  $Y$  is Normal ( $\mu_y, \sigma_y^2$ )
- ③ The correlation between  $X$  &  $Y$  is  $\rho_{xy} = \rho$
- ④  $Y|X=x$  is also Normal ( $\mu_y + \rho(\sigma_y/\sigma_x)(x - \mu_x), \sigma_y^2(1-\rho^2)$ )
- ⑤ For any constants  $a, b$ , the distribution of  $aX+bY$  is Normal with mean  $= a\mu_x + b\mu_y$ . and var  $= a^2\sigma_x^2 + b^2\sigma_y^2 + 2ab\rho\sigma_x\sigma_y$

Proof ②  $f_X(x) = \int_R f(x,y) dy = \int_{-\infty}^{+\infty} \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)}(w^2 - 2\rho w z + z^2)\right) dz$

Define  $z = \frac{y-\mu_y}{\sigma_y}$

$$= \frac{\exp\left(-\frac{w^2}{2(1-\rho^2)}\right)}{2\pi\sigma_x\sqrt{1-\rho^2}} \int_{-\infty}^{+\infty} \exp\left(-\frac{1}{2(1-\rho^2)}(z^2 - 2\rho w z + \rho^2 w^2 - \rho^2 w^2)\right) dz$$

$dy = \sigma_y dz$   $w = \frac{y-\mu_y}{\sigma_y}$

$$= \frac{e^{-w^2/2}}{2\pi\sigma_x\sqrt{1-\rho^2}} e^{\frac{\rho^2 w^2}{2(1-\rho^2)}} \int_{-\infty}^{+\infty} \exp\left(-\frac{1}{2(1-\rho^2)}(z - \rho w)^2\right) dz$$

↓ pdf with normal  
with mean  $\rho w$   
variance  $1-\rho^2$

$$= \frac{e^{-w^2/2}}{2\pi\sigma_x\sqrt{1-\rho^2}} \sqrt{2\pi} \sqrt{1-\rho^2}$$
$$= \frac{1}{\sqrt{2\pi\sigma_x}} e^{-\frac{(x-\mu_x)^2}{2\sigma_x^2}}, x \in R$$

② is the same

$$f_{Y|X}(y|x) = \frac{\frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)}\left(\left(\frac{x-\mu_x}{\sigma_x}\right)^2 - 2\rho\left(\frac{x-\mu_x}{\sigma_x}\right)\left(\frac{y-\mu_y}{\sigma_y}\right) + \left(\frac{y-\mu_y}{\sigma_y}\right)^2\right]\right]}{\frac{1}{\sqrt{2\pi\sigma_x}} \exp\left(-\frac{(x-\mu_x)^2}{2\sigma_x^2}\right)}$$
$$= \frac{1}{\sqrt{2\pi\sigma_y\sqrt{1-\rho^2}}} \exp\left(-\frac{1}{2(1-\rho^2)}\left(\left(\frac{x-\mu_x}{\sigma_x}\right)^2 - (1-\rho^2)\left(\frac{x-\mu_x}{\sigma_x}\right)^2 + \left(\frac{y-\mu_y}{\sigma_y}\right)^2 - 2\rho\left(\frac{x-\mu_x}{\sigma_x}\right)\left(\frac{y-\mu_y}{\sigma_y}\right)\right)\right)$$

NO. \_\_\_\_\_ Date \_\_\_\_\_

$$= \frac{1}{\sqrt{2\pi\sigma_y\sqrt{1-\rho^2}}} \exp\left(-\frac{1}{2(1-\rho^2)}\left(\rho^2\left(\frac{x-\mu_x}{\sigma_x}\right)^2 - 2\rho\left(\frac{x-\mu_x}{\sigma_x}\right)\left(\frac{y-\mu_y}{\sigma_y}\right) + \left(\frac{y-\mu_y}{\sigma_y}\right)^2\right)\right)$$
$$= \frac{1}{\sqrt{2\pi\sigma_y\sqrt{1-\rho^2}}} \exp\left(-\frac{1}{2\sigma_y^2(1-\rho^2)}\left((y-\mu_y) - \rho\frac{\sigma_y}{\sigma_x}(x-\mu_x)\right)^2\right)$$
$$\Rightarrow Y|X \sim \text{Normal}(\mu_y + \rho\frac{\sigma_y}{\sigma_x}(x-\mu_x), \sigma_y^2(1-\rho^2))$$

Nov. 2nd 2021

**Definition 4.5.10** Let  $-\infty < \mu_X < \infty$ ,  $-\infty < \mu_Y < \infty$ ,  $0 < \sigma_X$ ,  $0 < \sigma_Y$ , and  $-1 < \rho < 1$  be five real numbers. The *bivariate normal pdf with means  $\mu_X$  and  $\mu_Y$ , variances  $\sigma_X^2$  and  $\sigma_Y^2$ , and correlation  $\rho$*  is the bivariate pdf given by

$$f(x, y) = \left( 2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2} \right)^{-1} \times \exp \left( -\frac{1}{2(1-\rho^2)} \left( \left( \frac{x-\mu_X}{\sigma_X} \right)^2 - 2\rho \left( \frac{x-\mu_X}{\sigma_X} \right) \left( \frac{y-\mu_Y}{\sigma_Y} \right) + \left( \frac{y-\mu_Y}{\sigma_Y} \right)^2 \right) \right)$$

for  $-\infty < x < \infty$  and  $-\infty < y < \infty$ .

- a. The marginal distribution of  $X$  is  $n(\mu_X, \sigma_X^2)$ .
- b. The marginal distribution of  $Y$  is  $n(\mu_Y, \sigma_Y^2)$ .
- c. The correlation between  $X$  and  $Y$  is  $\rho_{XY} = \rho$ .
- d. For any constants  $a$  and  $b$ , the distribution of  $aX + bY$  is  $n(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\rho\sigma_X\sigma_Y)$ .

We will leave the verification of properties (a), (b), and (d) as exercises (Exercise 4.45). Assuming (a) and (b) are true, we will prove (c). We have by definition

$$\begin{aligned} \rho_{XY} &= \frac{\text{Cov}(X, Y)}{\sigma_X\sigma_Y} \\ &= \frac{E(X - \mu_X)(Y - \mu_Y)}{\sigma_X\sigma_Y} \\ &= E\left(\frac{X - \mu_X}{\sigma_X}\right)\left(\frac{Y - \mu_Y}{\sigma_Y}\right) \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left( \frac{x - \mu_X}{\sigma_X} \right) \left( \frac{y - \mu_Y}{\sigma_Y} \right) f(x, y) dx dy. \end{aligned}$$

Make the change of variable

$$s = \left( \frac{x - \mu_X}{\sigma_X} \right) \left( \frac{y - \mu_Y}{\sigma_Y} \right) \quad \text{and} \quad t = \left( \frac{x - \mu_X}{\sigma_X} \right).$$

Then  $x = \sigma_X t + \mu_X$ ,  $y = (\sigma_Y s/t) + \mu_Y$ , and the Jacobian of the transformation is  $J = \sigma_X\sigma_Y/t$ . With this change of variable, we obtain

$$\begin{aligned} \rho_{XY} &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} sf\left(\sigma_X t + \mu_X, \frac{\sigma_Y s}{t} + \mu_Y\right) \left| \frac{\sigma_X\sigma_Y}{t} \right| ds dt \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} s \left( 2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2} \right)^{-1} \\ &\quad \times \exp \left( -\frac{1}{2(1-\rho^2)} \left( t^2 - 2\rho s + \left( \frac{s}{t} \right)^2 \right) \right) \frac{\sigma_X\sigma_Y}{|t|} ds dt. \end{aligned}$$

Noting that  $|t| = \sqrt{t^2}$  and  $t^2 - 2\rho s + \left(\frac{s}{t}\right)^2 = \left(\frac{s-\rho t^2}{t}\right)^2 + (1-\rho^2)t^2$ , we can rewrite this as

$$\rho_{XY} = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) \left[ \int_{-\infty}^{\infty} \frac{s}{\sqrt{2\pi} \sqrt{(1-\rho^2)t^2}} \exp\left(-\frac{(s-\rho t^2)^2}{2(1-\rho^2)t^2}\right) ds \right] dt.$$

The inner integral is  $ES$ , where  $S$  is a normal random variable with  $ES = \rho t^2$  and  $\text{Var } S = (1-\rho^2)t^2$ . Thus the inner integral is  $\rho t^2$ . Hence we have

$$\rho_{XY} = \int_{-\infty}^{\infty} \frac{\rho t^2}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt.$$

But this integral is  $\rho ET^2$ , where  $T$  is a  $n(0, 1)$  random variable. Hence  $ET^2 = 1$  and  $\rho_{XY} = \rho$ .

All the conditional distributions of  $Y$  given  $X = x$  and of  $X$  given  $Y = y$  are also normal distributions. Using the joint and marginal pdfs given above, it is straightforward to verify that the conditional distribution of  $Y$  given  $X = x$  is

$$n(\mu_Y + \rho(\sigma_Y/\sigma_X)(x - \mu_X), \sigma_Y^2(1 - \rho^2)).$$

As  $\rho$  converges to 1 or  $-1$ , the conditional variance  $\sigma_Y^2(1 - \rho^2)$  converges to 0. Thus, the conditional distribution of  $Y$  given  $X = x$  becomes more concentrated about the point  $\mu_Y + \rho(\sigma_Y/\sigma_X)(x - \mu_X)$ , and the joint probability distribution of  $(X, Y)$  becomes more concentrated about the line  $y = \mu_Y + \rho(\sigma_Y/\sigma_X)(x - \mu_X)$ . This illustrates again the point made earlier that a correlation near 1 or  $-1$  means that there is a line  $y = ax + b$  about which the values of  $(X, Y)$  cluster with high probability.

Note one important fact: All of the normal marginal and conditional pdfs are derived from the starting point of bivariate normality. The derivation does not go in the opposite direction. That is, marginal normality does not imply joint normality. See Exercise 4.47 for an illustration of this.

proof of (d):

$$\begin{aligned} M_{ax+by}(t) &= E(e^{t(ax+by)}) = E(e^{taX+tbY}) \\ &= E(E(e^{taX+tbY}|X)) \\ &= E(e^{taX} \cdot E(e^{tbY}|X)) \\ &= E(e^{taX} \cdot e^{tb(\mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(X - \mu_X) + \frac{t^2 b^2 (1 - \rho^2) \sigma_Y^2}{2})}) \\ &= e^{tb\mu_Y - tb \rho \frac{\sigma_Y}{\sigma_X} \mu_X + \frac{t^2 b^2 (1 - \rho^2) \sigma_Y^2}{2}} E(e^{ta + \frac{tb \rho \sigma_Y}{\sigma_X} X}) \\ &= \exp\left[tb\mu_Y - \frac{tb \rho \sigma_Y}{\sigma_X} \mu_X + \frac{t^2 b^2 (1 - \rho^2) \sigma_Y^2}{2} + ta\mu_X + \frac{tb \rho \sigma_Y}{\sigma_X} \mu_X + \frac{t^2 a^2 \sigma_X^2}{2} + t^2 ab \rho \sigma_X \sigma_Y + \frac{t^2 b^2 \rho^2 \sigma_Y^2}{2}\right] \\ &= \exp\left[t(a\mu_X + b\mu_Y) + \frac{t^2 (a^2 \sigma_X^2 + b^2 \sigma_Y^2 + 2ab \sigma_X \sigma_Y)}{2}\right] \end{aligned}$$

$$\Rightarrow aX + bY \sim N(a\mu_X + b\mu_Y, a^2 \sigma_X^2 + b^2 \sigma_Y^2 + 2ab \sigma_X \sigma_Y)$$

## 4.6 Multivariate Distributions

The random vector  $\mathbf{X} = (X_1, \dots, X_n)$  has a sample space that is a subset of  $\mathbb{R}^n$ . If  $(X_1, \dots, X_n)$  is a discrete random vector (the sample space is countable), then the *joint pmf of  $(X_1, \dots, X_n)$*  is the function defined by  $f(\mathbf{x}) = f(x_1, \dots, x_n) = P(X_1 = x_1, \dots, X_n = x_n)$  for each  $(x_1, \dots, x_n) \in \mathbb{R}^n$ . Then for any  $A \subset \mathbb{R}^n$ ,

$$(4.6.1) \quad P(\mathbf{X} \in A) = \sum_{\mathbf{x} \in A} f(\mathbf{x}).$$

If  $(X_1, \dots, X_n)$  is a continuous random vector, the *joint pdf of  $(X_1, \dots, X_n)$*  is a function  $f(x_1, \dots, x_n)$  that satisfies

$$(4.6.2) \quad P(\mathbf{X} \in A) = \int \cdots \int_A f(\mathbf{x}) d\mathbf{x} = \int \cdots \int_A f(x_1, \dots, x_n) dx_1 \cdots dx_n.$$

Let  $g(\mathbf{x}) = g(x_1, \dots, x_n)$  be a real-valued function defined on the sample space of  $\mathbf{X}$ . Then  $g(\mathbf{X})$  is a random variable and the *expected value of  $g(\mathbf{X})$*  is

$$(4.6.3) \quad \text{E}g(\mathbf{X}) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} \quad \text{and} \quad \text{E}g(\mathbf{X}) = \sum_{\mathbf{x} \in \mathbb{R}^n} g(\mathbf{x}) f(\mathbf{x})$$

in the continuous and discrete cases, respectively. These and other definitions are analogous to the bivariate definitions except that now the integrals or sums are over the appropriate subset of  $\mathbb{R}^n$  rather than  $\mathbb{R}^2$ .

The *marginal pdf or pmf* of any subset of the coordinates of  $(X_1, \dots, X_n)$  can be computed by integrating or summing the joint pdf or pmf over all possible values of the other coordinates. Thus, for example, the marginal distribution of  $(X_1, \dots, X_k)$ , the first  $k$  coordinates of  $(X_1, \dots, X_n)$ , is given by the pdf or pmf

$$(4.6.4) \quad f(x_1, \dots, x_k) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_{k+1} \cdots dx_n$$

or

$$(4.6.5) \quad f(x_1, \dots, x_k) = \sum_{(x_{k+1}, \dots, x_n) \in \mathbb{R}^{n-k}} f(x_1, \dots, x_n)$$

for every  $(x_1, \dots, x_k) \in \mathbb{R}^k$ . The *conditional pdf or pmf* of a subset of the coordinates of  $(X_1, \dots, X_n)$  given the values of the remaining coordinates is obtained by dividing the joint pdf or pmf by the marginal pdf or pmf of the remaining coordinates. Thus, for example, if  $f(x_1, \dots, x_k) > 0$ , the conditional pdf or pmf of  $(X_{k+1}, \dots, X_n)$  given  $X_1 = x_1, \dots, X_k = x_k$  is the function of  $(x_{k+1}, \dots, x_n)$  defined by

$$(4.6.6) \quad f(x_{k+1}, \dots, x_n | x_1, \dots, x_k) = \frac{f(x_1, \dots, x_n)}{f(x_1, \dots, x_k)}.$$

These ideas are illustrated in the following example.

**Example 4.6.1 (Multivariate pdfs)** Let  $n = 4$  and

$$f(x_1, x_2, x_3, x_4) = \begin{cases} \frac{3}{4}(x_1^2 + x_2^2 + x_3^2 + x_4^2) & 0 < x_i < 1, i = 1, 2, 3, 4 \\ 0 & \text{otherwise.} \end{cases}$$

This nonnegative function is the joint pdf of a random vector  $(X_1, X_2, X_3, X_4)$  and it can be verified that

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3, x_4) dx_1 dx_2 dx_3 dx_4 \\ &= \int_0^1 \int_0^1 \int_0^1 \int_0^1 \frac{3}{4}(x_1^2 + x_2^2 + x_3^2 + x_4^2) dx_1 dx_2 dx_3 dx_4 \\ &= 1. \end{aligned}$$

This joint pdf can be used to compute probabilities such as

$$\begin{aligned} P\left(X_1 < \frac{1}{2}, X_2 < \frac{3}{4}, X_4 > \frac{1}{2}\right) \\ = \int_{\frac{1}{2}}^1 \int_0^1 \int_0^{\frac{3}{4}} \int_0^{\frac{1}{2}} \frac{3}{4}(x_1^2 + x_2^2 + x_3^2 + x_4^2) dx_1 dx_2 dx_3 dx_4. \end{aligned}$$

Note how the limits of integration restrict the integration to those values of  $(x_1, x_2, x_3, x_4)$  that are in the event in question and for which  $f(x_1, x_2, x_3, x_4) > 0$ . Each of the four terms,  $\frac{3}{4}x_1^2$ ,  $\frac{3}{4}x_2^2$ , etc., can be integrated separately and the results summed. For example,

$$\int_{\frac{1}{2}}^1 \int_0^1 \int_0^{\frac{3}{4}} \int_0^{\frac{1}{2}} \frac{3}{4}x_1^2 dx_1 dx_2 dx_3 dx_4 = \frac{3}{256}.$$

The other three integrals are  $\frac{7}{1024}$ ,  $\frac{3}{64}$ , and  $\frac{21}{256}$ . Thus

$$P\left(X_1 < \frac{1}{2}, X_2 < \frac{3}{4}, X_4 > \frac{1}{2}\right) = \frac{3}{256} + \frac{7}{1024} + \frac{3}{64} + \frac{21}{256} = \frac{151}{1024}.$$

Using (4.6.4), we can obtain the marginal pdf of  $(X_1, X_2)$  by integrating out the variables  $x_3$  and  $x_4$  to obtain

$$\begin{aligned} f(x_1, x_2) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x_1, x_2, x_3, x_4) dx_3 dx_4 \\ &= \int_0^1 \int_0^1 \frac{3}{4}(x_1^2 + x_2^2 + x_3^2 + x_4^2) dx_3 dx_4 = \frac{3}{4}(x_1^2 + x_2^2) + \frac{1}{2} \end{aligned}$$

for  $0 < x_1 < 1$  and  $0 < x_2 < 1$ . Any probability or expected value that involves only  $X_1$  and  $X_2$  can be computed using this marginal pdf. For example,

$$\begin{aligned} EX_1 X_2 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 x_2 f(x_1, x_2) dx_1 dx_2 \\ &= \int_0^1 \int_0^1 x_1 x_2 \left( \frac{3}{4}(x_1^2 + x_2^2) + \frac{1}{2} \right) dx_1 dx_2 \\ &= \int_0^1 \int_0^1 \left( \frac{3}{4}x_1^3 x_2 + \frac{3}{4}x_1 x_2^3 + \frac{1}{2}x_1 x_2 \right) dx_1 dx_2 \\ &= \int_0^1 \left( \frac{3}{16}x_2 + \frac{3}{8}x_2^3 + \frac{1}{4}x_2 \right) dx_2 = \frac{3}{32} + \frac{3}{32} + \frac{1}{8} = \frac{5}{16}. \end{aligned}$$

For any  $(x_1, x_2)$  with  $0 < x_1 < 1$  and  $0 < x_2 < 1$ ,  $f(x_1, x_2) > 0$  and the conditional pdf of  $(X_3, X_4)$  given  $X_1 = x_1$  and  $X_2 = x_2$  can be found using (4.6.6). For any such  $(x_1, x_2)$ ,  $f(x_1, x_2, x_3, x_4) > 0$  if  $0 < x_3 < 1$  and  $0 < x_4 < 1$ , and for these values of  $(x_3, x_4)$ , the conditional pdf is

$$\begin{aligned} f(x_3, x_4 | x_1, x_2) &= \frac{f(x_1, x_2, x_3, x_4)}{f(x_1, x_2)} \\ &= \frac{\frac{3}{4}(x_1^2 + x_2^2 + x_3^2 + x_4^2)}{\frac{3}{4}(x_1^2 + x_2^2) + \frac{1}{2}} \\ &= \frac{x_1^2 + x_2^2 + x_3^2 + x_4^2}{x_1^2 + x_2^2 + \frac{2}{3}}. \end{aligned}$$

For example, the conditional pdf of  $(X_3, X_4)$  given  $X_1 = \frac{1}{3}$  and  $X_2 = \frac{2}{3}$  is

$$f\left(x_3, x_4 \mid x_1 = \frac{1}{3}, x_2 = \frac{2}{3}\right) = \frac{\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 + x_3^2 + x_4^2}{\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 + \frac{2}{3}} = \frac{5}{11} + \frac{9}{11}x_3^2 + \frac{9}{11}x_4^2.$$

This can be used to compute

$$\begin{aligned} P\left(X_3 > \frac{3}{4}, X_4 < \frac{1}{2} \mid X_1 = \frac{1}{3}, X_2 = \frac{2}{3}\right) &= \int_0^{\frac{1}{2}} \int_{\frac{3}{4}}^1 \left( \frac{5}{11} + \frac{9}{11}x_3^2 + \frac{9}{11}x_4^2 \right) dx_3 dx_4 \\ &= \int_0^{\frac{1}{2}} \left( \frac{5}{44} + \frac{111}{704} + \frac{9}{44}x_4^2 \right) dx_4 \\ &= \frac{5}{88} + \frac{111}{1408} + \frac{3}{352} = \frac{203}{1408}. \end{aligned}$$

**Definition 4.6.2** Let  $n$  and  $m$  be positive integers and let  $p_1, \dots, p_n$  be numbers satisfying  $0 \leq p_i \leq 1$ ,  $i = 1, \dots, n$ , and  $\sum_{i=1}^n p_i = 1$ . Then the random vector  $(X_1, \dots, X_n)$  has a *multinomial distribution with  $m$  trials and cell probabilities  $p_1, \dots, p_n$*  if the joint pmf of  $(X_1, \dots, X_n)$  is

$$f(x_1, \dots, x_n) = \frac{m!}{x_1! \cdots x_n!} p_1^{x_1} \cdots p_n^{x_n} = m! \prod_{i=1}^n \frac{p_i^{x_i}}{x_i!}$$

on the set of  $(x_1, \dots, x_n)$  such that each  $x_i$  is a nonnegative integer and  $\sum_{i=1}^n x_i = m$ .

**Theorem 4.6.4 (Multinomial Theorem)** Let  $m$  and  $n$  be positive integers. Let  $\mathcal{A}$  be the set of vectors  $\mathbf{x} = (x_1, \dots, x_n)$  such that each  $x_i$  is a nonnegative integer and  $\sum_{i=1}^n x_i = m$ . Then, for any real numbers  $p_1, \dots, p_n$ ,

$$(p_1 + \cdots + p_n)^m = \sum_{\mathbf{x} \in \mathcal{A}} \frac{m!}{x_1! \cdots x_n!} p_1^{x_1} \cdots p_n^{x_n}.$$

Theorem 4.6.4 shows that a multinomial pmf sums to 1. The set  $\mathcal{A}$  is the set of points with positive probability in Definition 4.6.2. The sum of the pmf over all those points is, by Theorem 4.6.4,  $(p_1 + \cdots + p_n)^m = 1^m = 1$ .

Now we consider some marginal and conditional distributions for the multinomial model. Consider a single coordinate  $X_i$ . If the occurrence of the  $i$ th outcome is labeled a “success” and anything else is labeled a “failure,” then  $X_i$  is the count of the number of successes in  $m$  independent trials where the probability of a success is  $p_i$  on each trial. Thus  $X_i$  should have a  $\text{binomial}(m, p_i)$  distribution. To verify this the marginal distribution of  $X_i$  should be computed using (4.6.5). For example, consider the marginal pmf of  $X_n$ . For a fixed value of  $x_n \in \{0, 1, \dots, n\}$ , to compute the marginal pmf  $f(x_n)$ , we must sum over all possible values of  $(x_1, \dots, x_{n-1})$ . That is, we must sum over all  $(x_1, \dots, x_{n-1})$  such that the  $x_i$ s are all nonnegative integers and  $\sum_{i=1}^{n-1} x_i = m - x_n$ . Denote this set by  $\mathcal{B}$ . Then

$$\begin{aligned} f(x_n) &= \sum_{(x_1, \dots, x_{n-1}) \in \mathcal{B}} \frac{m!}{x_1! \cdots x_n!} (p_1)^{x_1} \cdots (p_n)^{x_n} \\ &= \sum_{(x_1, \dots, x_{n-1}) \in \mathcal{B}} \frac{m!}{x_1! \cdots x_n!} p_1^{x_1} \cdots p_n^{x_n} \frac{(m - x_n)!}{(m - x_n)!} \frac{(1 - p_n)^{m - x_n}}{(1 - p_n)^{m - x_n}} \\ &= \frac{m!}{x_n! (m - x_n)!} p_n^{x_n} (1 - p_n)^{m - x_n} \\ &\quad \times \sum_{(x_1, \dots, x_{n-1}) \in \mathcal{B}} \frac{(m - x_n)!}{x_1! \cdots x_{n-1}!} \left( \frac{p_1}{1 - p_n} \right)^{x_1} \cdots \left( \frac{p_{n-1}}{1 - p_n} \right)^{x_{n-1}}. \end{aligned}$$

= |

But using the facts that  $x_1 + \dots + x_{n-1} = m - x_n$  and  $p_1 + \dots + p_{n-1} = 1 - p_n$  and Theorem 4.6.4, we see that the last summation is 1. Hence the marginal distribution of  $X_n$  is binomial( $m, p_n$ ). Similar arguments show that each of the other coordinates is marginally binomially distributed.

Given that  $X_n = x_n$ , there must have been  $m - x_n$  trials that resulted in one of the first  $n - 1$  outcomes. The vector  $(X_1, \dots, X_{n-1})$  counts the number of these  $m - x_n$  trials that are of each type. Thus it seems that given  $X_n = x_n$ ,  $(X_1, \dots, X_{n-1})$  might have a multinomial distribution. This is true. From (4.6.6), the conditional pmf of  $(X_1, \dots, X_{n-1})$  given  $X_n = x_n$  is

$$\begin{aligned} f(x_1, \dots, x_{n-1} | x_n) &= \frac{f(x_1, \dots, x_n)}{f(x_n)} \\ &= \frac{\frac{m!}{x_1! \dots x_n!} (p_1)^{x_1} \dots (p_n)^{x_n}}{\frac{m!}{x_n!(m-x_n)!} (p_n)^{x_n} (1-p_n)^{m-x_n}} \\ &= \frac{(m-x_n)!}{x_1! \dots x_{n-1}!} \left(\frac{p_1}{1-p_n}\right)^{x_1} \dots \left(\frac{p_{n-1}}{1-p_n}\right)^{x_{n-1}}. \end{aligned}$$

This is the pmf of a multinomial distribution with  $m - x_n$  trials and cell probabilities  $p_1/(1-p_n), \dots, p_{n-1}/(1-p_n)$ . In fact, the conditional distribution of any subset of the coordinates of  $(X_1, \dots, X_n)$  given the values of the rest of the coordinates is a multinomial distribution.

We see from the conditional distributions that the coordinates of the vector  $(X_1, \dots, X_n)$  are related. In particular, there must be some negative correlation. It turns out that all of the pairwise covariances are negative and are given by (Exercise 4.39)

$$\text{Cov}(X_i, X_j) = E[(X_i - p_i)(X_j - p_j)] = -mp_i p_j.$$

**Definition 4.6.5** Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be random vectors with joint pdf or pmf  $f(\mathbf{x}_1, \dots, \mathbf{x}_n)$ . Let  $f_{\mathbf{X}_i}(\mathbf{x}_i)$  denote the marginal pdf or pmf of  $\mathbf{X}_i$ . Then  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are called mutually independent random vectors if, for every  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n) = f_{\mathbf{X}_1}(\mathbf{x}_1) \dots f_{\mathbf{X}_n}(\mathbf{x}_n) = \prod_{i=1}^n f_{\mathbf{X}_i}(\mathbf{x}_i).$$

If the  $X_i$ s are all one-dimensional, then  $X_1, \dots, X_n$  are called mutually independent random variables.

**Theorem 4.6.6 (Generalization of Theorem 4.2.10)** Let  $X_1, \dots, X_n$  be mutually independent random variables. Let  $g_1, \dots, g_n$  be real-valued functions such that  $g_i(x_i)$  is a function only of  $x_i$ ,  $i = 1, \dots, n$ . Then

$$E(g_1(X_1) \dots g_n(X_n)) = (Eg_1(X_1)) \dots (Eg_n(X_n)).$$

**Theorem 4.6.7 (Generalization of Theorem 4.2.12)** Let  $X_1, \dots, X_n$  be mutually independent random variables with mgfs  $M_{X_1}(t), \dots, M_{X_n}(t)$ . Let  $Z = X_1 + \dots + X_n$ . Then the mgf of  $Z$  is

$$M_Z(t) = M_{X_1}(t) \cdots M_{X_n}(t).$$

In particular, if  $X_1, \dots, X_n$  all have the same distribution with mgf  $M_X(t)$ , then

$$M_Z(t) = (M_X(t))^n.$$

**Example 4.6.8 (Mgf of a sum of gamma variables)** Suppose  $X_1, \dots, X_n$  are mutually independent random variables, and the distribution of  $X_i$  is  $\text{gamma}(\alpha_i, \beta)$ . From Example 2.3.8, the mgf of a  $\text{gamma}(\alpha, \beta)$  distribution is  $M(t) = (1 - \beta t)^{-\alpha}$ . Thus, if  $Z = X_1 + \dots + X_n$ , the mgf of  $Z$  is

$$M_Z(t) = M_{X_1}(t) \cdots M_{X_n}(t) = (1 - \beta t)^{-\alpha_1} \cdots (1 - \beta t)^{-\alpha_n} = (1 - \beta t)^{-(\alpha_1 + \dots + \alpha_n)}.$$

This is the mgf of a  $\text{gamma}(\alpha_1 + \dots + \alpha_n, \beta)$  distribution. Thus, the sum of independent gamma random variables that have a common scale parameter  $\beta$  also has a gamma distribution. ||

**Corollary 4.6.9** Let  $X_1, \dots, X_n$  be mutually independent random variables with mgfs  $M_{X_1}(t), \dots, M_{X_n}(t)$ . Let  $a_1, \dots, a_n$  and  $b_1, \dots, b_n$  be fixed constants. Let  $Z = (a_1 X_1 + b_1) + \dots + (a_n X_n + b_n)$ . Then the mgf of  $Z$  is

$$M_Z(t) = (e^{t(\sum b_i)}) M_{X_1}(a_1 t) \cdots M_{X_n}(a_n t).$$

**Proof:** From the definition, the mgf of  $Z$  is

$$\begin{aligned} M_Z(t) &= E e^{tZ} \\ &= E e^{t \sum (a_i X_i + b_i)} \\ &= (e^{t(\sum b_i)}) E(e^{t a_1 X_1} \cdots e^{t a_n X_n}) \quad \left( \begin{array}{l} \text{properties of exponentials} \\ \text{and expectations} \end{array} \right) \\ &= (e^{t(\sum b_i)}) M_{X_1}(a_1 t) \cdots M_{X_n}(a_n t), \quad (\text{Theorem 4.6.6}) \end{aligned}$$

as was to be shown. □

**Corollary 4.6.10** Let  $X_1, \dots, X_n$  be mutually independent random variables with  $X_i \sim n(\mu_i, \sigma_i^2)$ . Let  $a_1, \dots, a_n$  and  $b_1, \dots, b_n$  be fixed constants. Then

$$Z = \sum_{i=1}^n (a_i X_i + b_i) \sim n \left( \sum_{i=1}^n (a_i \mu_i + b_i), \sum_{i=1}^n a_i^2 \sigma_i^2 \right).$$

**Proof:** Recall that the mgf of a  $n(\mu, \sigma^2)$  random variable is  $M(t) = e^{\mu t + \sigma^2 t^2/2}$ . Substituting into the expression in Corollary 4.6.9 yields

$$\begin{aligned} M_Z(t) &= (e^{t(\sum b_i)}) e^{\mu_1 a_1 t + \sigma_1^2 a_1^2 t^2/2} \cdots e^{\mu_n a_n t + \sigma_n^2 a_n^2 t^2/2} \\ &= e^{((\sum (a_i \mu_i + b_i)) t + (\sum a_i^2 \sigma_i^2) t^2/2)}, \end{aligned}$$

the mgf of the indicated normal distribution. □

**Theorem 4.6.11 (Generalization of Lemma 4.2.7)** Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be random vectors. Then  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are mutually independent random vectors if and only if there exist functions  $g_i(\mathbf{x}_i)$ ,  $i = 1, \dots, n$ , such that the joint pdf or pmf of  $(\mathbf{X}_1, \dots, \mathbf{X}_n)$  can be written as

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n) = g_1(\mathbf{x}_1) \cdots g_n(\mathbf{x}_n).$$

**Theorem 4.6.12 (Generalization of Theorem 4.3.5)** Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be independent random vectors. Let  $g_i(\mathbf{x}_i)$  be a function only of  $\mathbf{x}_i$ ,  $i = 1, \dots, n$ . Then the random variables  $U_i = g_i(\mathbf{X}_i)$ ,  $i = 1, \dots, n$ , are mutually independent.

Let  $(X_1, \dots, X_n)$  be a random vector with pdf  $f_{\mathbf{X}}(x_1, \dots, x_n)$ . Let  $\mathcal{A} = \{\mathbf{x} : f_{\mathbf{X}}(\mathbf{x}) > 0\}$ . Consider a new random vector  $(U_1, \dots, U_n)$ , defined by  $U_1 = g_1(X_1, \dots, X_n)$ ,  $U_2 = g_2(X_1, \dots, X_n), \dots, U_n = g_n(X_1, \dots, X_n)$ . Suppose that  $A_0, A_1, \dots, A_k$  form a partition of  $\mathcal{A}$  with these properties. The set  $A_0$ , which may be empty, satisfies  $P((X_1, \dots, X_n) \in A_0) = 0$ . The transformation  $(U_1, \dots, U_n) = (g_1(\mathbf{X}), \dots, g_n(\mathbf{X}))$  is a one-to-one transformation from  $A_i$  onto  $\mathcal{B}$  for each  $i = 1, 2, \dots, k$ . Then for each  $i$ , the inverse functions from  $\mathcal{B}$  to  $A_i$  can be found. Denote the  $i$ th inverse by  $x_1 = h_{1i}(u_1, \dots, u_n)$ ,  $x_2 = h_{2i}(u_1, \dots, u_n), \dots, x_n = h_{ni}(u_1, \dots, u_n)$ . This  $i$ th inverse gives, for  $(u_1, \dots, u_n) \in \mathcal{B}$ , the unique  $(x_1, \dots, x_n) \in A_i$  such that  $(u_1, \dots, u_n) = (g_1(x_1, \dots, x_n), \dots, g_n(x_1, \dots, x_n))$ . Let  $J_i$  denote the Jacobian computed from the  $i$ th inverse. That is,

$$J_i = \begin{vmatrix} \frac{\partial x_1}{\partial u_1} & \frac{\partial x_1}{\partial u_2} & \cdots & \frac{\partial x_1}{\partial u_n} \\ \frac{\partial x_2}{\partial u_1} & \frac{\partial x_2}{\partial u_2} & \cdots & \frac{\partial x_2}{\partial u_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_n}{\partial u_1} & \frac{\partial x_n}{\partial u_2} & \cdots & \frac{\partial x_n}{\partial u_n} \end{vmatrix} = \begin{vmatrix} \frac{\partial h_{1i}(\mathbf{u})}{\partial u_1} & \frac{\partial h_{1i}(\mathbf{u})}{\partial u_2} & \cdots & \frac{\partial h_{1i}(\mathbf{u})}{\partial u_n} \\ \frac{\partial h_{2i}(\mathbf{u})}{\partial u_1} & \frac{\partial h_{2i}(\mathbf{u})}{\partial u_2} & \cdots & \frac{\partial h_{2i}(\mathbf{u})}{\partial u_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_{ni}(\mathbf{u})}{\partial u_1} & \frac{\partial h_{ni}(\mathbf{u})}{\partial u_2} & \cdots & \frac{\partial h_{ni}(\mathbf{u})}{\partial u_n} \end{vmatrix},$$

the determinant of an  $n \times n$  matrix. Assuming that these Jacobians do not vanish identically on  $\mathcal{B}$ , we have the following representation of the joint pdf,  $f_{\mathbf{U}}(u_1, \dots, u_n)$ , for  $\mathbf{u} \in \mathcal{B}$ :

$$(4.6.7) \quad f_{\mathbf{U}}(u_1, \dots, u_n) = \sum_{i=1}^k f_{\mathbf{X}}(h_{1i}(u_1, \dots, u_n), \dots, h_{ni}(u_1, \dots, u_n)) |J_i|.$$

**Example 4.6.13 (Multivariate change of variables)** Let  $(X_1, X_2, X_3, X_4)$  have joint pdf

$$f_{\mathbf{X}}(x_1, x_2, x_3, x_4) = 24e^{-x_1-x_2-x_3-x_4}, \quad 0 < x_1 < x_2 < x_3 < x_4 < \infty.$$

Consider the transformation

$$U_1 = X_1, \quad U_2 = X_2 - X_1, \quad U_3 = X_3 - X_2, \quad U_4 = X_4 - X_3.$$

This transformation maps the set  $\mathcal{A}$  onto the set  $\mathcal{B} = \{\mathbf{u}: 0 < u_i < \infty, i = 1, 2, 3, 4\}$ . The transformation is one-to-one, so  $k = 1$ , and the inverse is

$$X_1 = U_1, \quad X_2 = U_1 + U_2, \quad X_3 = U_1 + U_2 + U_3, \quad X_4 = U_1 + U_2 + U_3 + U_4.$$

The Jacobian of the inverse is

$$J = \begin{vmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{vmatrix} = 1.$$

Since the matrix is triangular, the determinant is equal to the product of the diagonal elements. Thus, from (4.6.7) we obtain

$$\begin{aligned} f_{\mathbf{U}}(u_1, \dots, u_4) &= 24e^{-u_1-(u_1+u_2)-(u_1+u_2+u_3)-(u_1+u_2+u_3+u_4)} \\ &= 24e^{-4u_1-3u_2-2u_3-u_4} \end{aligned}$$

on  $\mathcal{B}$ . From this the marginal pdfs of  $U_1, U_2, U_3$ , and  $U_4$  can be calculated. It turns out that  $f_U(u_i) = (5-i)e^{-(5-i)u_i}, 0 < u_i$ ; that is,  $U_i \sim \text{exponential}(1/(5-i))$ . From Theorem 4.6.11 we see that  $U_1, U_2, U_3$ , and  $U_4$  are mutually independent random variables. ||

## 4.7 Inequalities

**Lemma 4.7.1** Let  $a$  and  $b$  be any positive numbers, and let  $p$  and  $q$  be any positive numbers (necessarily greater than 1) satisfying

$$(4.7.1) \quad \frac{1}{p} + \frac{1}{q} = 1. \quad (p+q=pq)$$

Then

$$(4.7.2) \quad \frac{1}{p}a^p + \frac{1}{q}b^q \geq ab$$

with equality if and only if  $a^p = b^q$ .

**Proof:** Fix  $b$ , and consider the function

$$g(a) = \frac{1}{p}a^p + \frac{1}{q}b^q - ab.$$

To minimize  $g(a)$ , differentiate and set equal to 0:

$$\frac{d}{da}g(a) = 0 \Rightarrow a^{p-1} - b = 0 \Rightarrow b = a^{p-1}.$$

A check of the second derivative will establish that this is indeed a minimum. The value of the function at the minimum is

$$\begin{aligned} \frac{1}{p}a^p + \frac{1}{q}(a^{p-1})^q - aa^{p-1} &= \frac{1}{p}a^p + \frac{1}{q}a^p - a^p \quad \left( \begin{array}{l} (p-1)q = p \text{ follows} \\ \text{from (4.7.1)} \end{array} \right) \\ &= 0. \quad (\text{again from (4.7.1)}) \end{aligned}$$

Hence the minimum is 0 and (4.7.2) is established. Since the minimum is unique (why?), equality holds only if  $a^{p-1} = b$ , which is equivalent to  $a^p = b^q$ , again from (4.7.1).  $\square$

**Theorem 4.7.2 (Hölder's Inequality)** *Let  $X$  and  $Y$  be any two random variables, and let  $p$  and  $q$  satisfy (4.7.1). Then*

$$\|XY\|_1 \leq \|X\|_p \|Y\|_q$$

$$(4.7.3) \quad |EXY| \leq E|XY| \leq (E|X|^p)^{1/p} (E|Y|^q)^{1/q}.$$

**Proof:** The first inequality follows from  $-|XY| \leq XY \leq |XY|$  and Theorem 2.2.5. To prove the second inequality, define

$$a = \frac{|X|}{(E|X|^p)^{1/p}} \quad \text{and} \quad b = \frac{|Y|}{(E|Y|^q)^{1/q}}.$$

Applying Lemma 4.7.1, we get

$$\frac{1}{p} \frac{|X|^p}{E|X|^p} + \frac{1}{q} \frac{|Y|^q}{E|Y|^q} \geq \frac{|XY|}{(E|X|^p)^{1/p} (E|Y|^q)^{1/q}}.$$

Now take expectations of both sides. The expectation of the left-hand side is 1, and rearrangement gives (4.7.3).  $\square$

**Theorem 4.7.3 (Cauchy–Schwarz Inequality)** *For any two random variables  $X$  and  $Y$ ,*

$$(4.7.4) \quad |EXY| \leq E|XY| \leq (E|X|^2)^{1/2} (E|Y|^2)^{1/2}.$$

**Example 4.7.4 (Covariance inequality)** If  $X$  and  $Y$  have means  $\mu_X$  and  $\mu_Y$  and variances  $\sigma_X^2$  and  $\sigma_Y^2$ , respectively, we can apply the Cauchy–Schwarz Inequality to get

$$E|(X - \mu_X)(Y - \mu_Y)| \leq \{E(X - \mu_X)^2\}^{1/2} \{E(Y - \mu_Y)^2\}^{1/2}.$$

Squaring both sides and using statistical notation, we have

$$(\text{Cov}(X, Y))^2 \leq \sigma_X^2 \sigma_Y^2.$$

Recalling the definition of the correlation coefficient,  $\rho$ , we have proved that  $0 \leq \rho^2 \leq 1$ . Furthermore, the condition for equality in Lemma 4.7.1 still carries over, and equality is attained here only if  $X - \mu_X = c(Y - \mu_Y)$ , for some constant  $c$ . That is, the correlation is  $\pm 1$  if and only if  $X$  and  $Y$  are linearly related. Compare the ease of this proof to the one used in Theorem 4.5.7, before we had the Cauchy–Schwarz Inequality.  $\parallel$

Some other special cases of Hölder's Inequality are often useful. If we set  $Y \equiv 1$  in (4.7.3), we get

$$(4.7.5) \quad E|X| \leq \{E(|X|^p)\}^{1/p}, \quad 1 < p < \infty.$$

For  $1 < r < p$ , if we replace  $|X|$  by  $|X|^r$  in (4.7.5), we obtain

$$E|X|^r \leq \{E(|X|^{pr})\}^{1/p}.$$

Now write  $s = pr$  (note that  $s > r$ ) and rearrange terms to get

$$(4.7.6) \quad \{E|X|^r\}^{1/r} \leq \{E|X|^s\}^{1/s}, \quad 1 < r < s < \infty, \quad ||X||_r \leq ||X||_s \quad (r < s)$$

which is known as Liapounov's Inequality.

**Theorem 4.7.5 (Minkowski's Inequality)** Let  $X$  and  $Y$  be any two random variables. Then for  $1 \leq p < \infty$ ,

$$(4.7.7) \quad [E|X+Y|^p]^{1/p} \leq [E|X|^p]^{1/p} + [E|Y|^p]^{1/p}.$$

**Proof:** Write

$$\begin{aligned} E|X+Y|^p &= E(|X+Y||X+Y|^{p-1}) \\ (4.7.8) \quad &\leq E(|X||X+Y|^{p-1}) + E(|Y||X+Y|^{p-1}), \end{aligned}$$

where we have used the fact that  $|X+Y| \leq |X| + |Y|$  (the *triangle inequality*; see Exercise 4.64). Now apply Hölder's Inequality to each expectation on the right-hand side of (4.7.8) to get

$$\begin{aligned} E(|X+Y|^p) &\leq [E(|X|^p)]^{1/p} [E|X+Y|^{q(p-1)}]^{1/q} \\ &\quad + [E(|Y|^p)]^{1/p} [E|X+Y|^{q(p-1)}]^{1/q}, \end{aligned}$$

where  $q$  satisfies  $1/p + 1/q = 1$ . Now divide through by  $[E(|X+Y|^{q(p-1)})]^{1/q}$ . Noting that  $q(p-1) = p$  and  $1 - 1/q = 1/p$ , we obtain (4.7.7).  $\square$

Nov. 18th. 2021

**Theorem 4.7.7 (Jensen's Inequality)** For any random variable  $X$ , if  $g(x)$  is a convex function, then

$$Eg(X) \geq g(EX).$$

Equality holds if and only if, for every line  $a + bx$  that is tangent to  $g(x)$  at  $x = EX$ ,  $P(g(X) = a + bX) = 1$ .

**Proof:** To establish the inequality, let  $l(x)$  be a tangent line to  $g(x)$  at the point  $g(EX)$ . (Recall that  $EX$  is a constant.) Write  $l(x) = a + bx$  for some  $a$  and  $b$ . The situation is illustrated in Figure 4.7.2.

Now, by the convexity of  $g$  we have  $g(x) \geq a + bx$ . Since expectations preserve inequalities,

$$Eg(X) \geq E(a + bX)$$

$$\begin{aligned} &= a + bEX && \left( \begin{array}{l} \text{linearity of expectation,} \\ \text{Theorem 2.2.5} \end{array} \right) \\ &= l(EX) && (\text{definition of } l(x)) \\ &= g(EX), && (l \text{ is tangent at } EX) \end{aligned}$$

as was to be shown.

If  $g(x)$  is linear, equality follows from properties of expectations (Theorem 2.2.5). For the “only if” part see Exercise 4.62.  $\square$

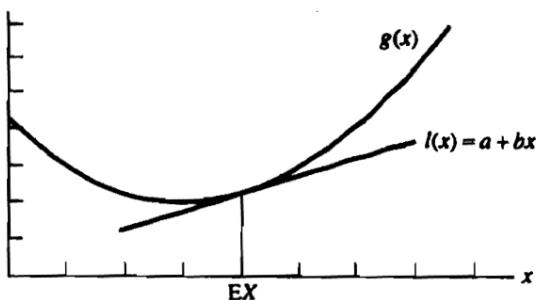


Figure 4.7.2. Graphical illustration of Jensen's Inequality

**Example 4.7.8 (An inequality for means)** Jensen's Inequality can be used to prove an inequality between three different kinds of means. If  $a_1, \dots, a_n$  are positive numbers, define

$$a_A = \frac{1}{n}(a_1 + a_2 + \dots + a_n), \quad (\text{arithmetic mean})$$

$$a_G = [a_1 a_2 \cdots a_n]^{1/n}, \quad (\text{geometric mean})$$

$$a_H = \frac{1}{\frac{1}{a_1} + \frac{1}{a_2} + \dots + \frac{1}{a_n}}. \quad (\text{harmonic mean})$$

An inequality relating these means is

$$a_H \leq a_G \leq a_A.$$

To apply Jensen's Inequality, let  $X$  be a random variable with range  $a_1, \dots, a_n$  and  $P(X = a_i) = 1/n, i = 1, \dots, n$ . Since  $\log x$  is a concave function, Jensen's Inequality shows that  $E(\log X) \leq \log(EX)$ ; hence,

$$\log a_G = \frac{1}{n} \sum_{i=1}^n \log a_i = E(\log X) \leq \log(EX) = \log \left( \frac{1}{n} \sum_{i=1}^n a_i \right) = \log a_A,$$

so  $a_G \leq a_A$ . Now again use the fact that  $\log x$  is concave to get

$$\log \frac{1}{a_H} = \log \left( \frac{1}{n} \sum_{i=1}^n \frac{1}{a_i} \right) = \log E \frac{1}{X} \geq E \left( \log \frac{1}{X} \right) = -E(\log X).$$

Since  $E(\log X) = \log a_G$ , it then follows that  $\log(1/a_H) \geq \log(1/a_G)$ , or  $a_G \geq a_H$ .  $\parallel$

The next inequality merely exploits the definition of covariance, but sometimes proves to be useful. If  $X$  is a random variable with finite mean  $\mu$  and  $g(x)$  is a nondecreasing function, then

$$E(g(X)(X - \mu)) \geq 0,$$

since

$$\begin{aligned} E(g(X)(X - \mu)) &= E(g(X)(X - \mu)I_{(-\infty, 0)}(X - \mu)) + E(g(X)(X - \mu)I_{[0, \infty)}(X - \mu)) \\ &\geq E(g(\mu)(X - \mu)I_{(-\infty, 0)}(X - \mu)) \\ &\quad + E(g(\mu)(X - \mu)I_{[0, \infty)}(X - \mu)) \quad (\text{since } g \text{ is nondecreasing}) \\ &= g(\mu)E(X - \mu) \\ &= 0. \end{aligned}$$

**Theorem 4.7.9 (Covariance Inequality)** Let  $X$  be any random variable and  $g(x)$  and  $h(x)$  any functions such that  $Eg(X)$ ,  $Eh(X)$ , and  $E(g(X)h(X))$  exist.

a. If  $g(x)$  is a nondecreasing function and  $h(x)$  is a nonincreasing function, then

$$E(g(X)h(X)) \leq (Eg(X))(Eh(X)).$$

b. If  $g(x)$  and  $h(x)$  are either both nondecreasing or both nonincreasing, then

$$E(g(X)h(X)) \geq (Eg(X))(Eh(X)).$$

The intuition behind the inequality is easy. In case (a) there is negative correlation between  $g$  and  $h$ , while in case (b) there is positive correlation. The inequalities merely reflect this fact. The usefulness of the Covariance Inequality is that it allows us to bound an expectation without using higher-order moments.

# Properties of a Random Sample

**Definition 5.1.1** The random variables  $X_1, \dots, X_n$  are called a random sample of size  $n$  from the population  $f(x)$  if  $X_1, \dots, X_n$  are mutually independent random variables and the marginal pdf or pmf of each  $X_i$  is the same function  $f(x)$ . Alternatively,  $X_1, \dots, X_n$  are called independent and identically distributed random variables with pdf or pmf  $f(x)$ . This is commonly abbreviated to iid random variables.

From Definition 4.6.5, the joint pdf or pmf of  $X_1, \dots, X_n$  is given by

$$(5.1.1) \quad f(x_1, \dots, x_n) = f(x_1)f(x_2) \cdots f(x_n) = \prod_{i=1}^n f(x_i).$$

**Example 5.1.2 (Sample pdf-exponential)** Let  $X_1, \dots, X_n$  be a random sample from an exponential( $\beta$ ) population. Specifically,  $X_1, \dots, X_n$  might correspond to the times until failure (measured in years) for  $n$  identical circuit boards that are put on test and used until they fail. The joint pdf of the sample is

$$f(x_1, \dots, x_n | \beta) = \prod_{i=1}^n f(x_i | \beta) = \prod_{i=1}^n \frac{1}{\beta} e^{-x_i/\beta} = \frac{1}{\beta^n} e^{-(x_1 + \dots + x_n)/\beta}.$$

This pdf can be used to answer questions about the sample. For example, what is the probability that all the boards last more than 2 years? We can compute

$$\begin{aligned} & P(X_1 > 2, \dots, X_n > 2) \\ &= \int_2^\infty \cdots \int_2^\infty \prod_{i=1}^n \frac{1}{\beta} e^{-x_i/\beta} dx_1 \cdots dx_n \\ &= e^{-2/\beta} \int_2^\infty \cdots \int_2^\infty \prod_{i=2}^n \frac{1}{\beta} e^{-x_i/\beta} dx_2 \cdots dx_n \quad (\text{integrate out } x_1) \\ &\vdots \qquad \qquad \qquad (\text{integrate out the remaining } x_i\text{'s successively}) \\ &= (e^{-2/\beta})^n \\ &= e^{-2n/\beta}. \end{aligned}$$

Or:  $P(X_1 > 2, \dots, X_n > 2)$

$$\begin{aligned} &= P(X_1 > 2) \cdots P(X_n > 2) \quad (\text{independence}) \\ &= [P(X_1 > 2)]^n \quad (\text{identical distributions}) \\ &= (e^{-2/\beta})^n \quad (\text{exponential calculation}) \\ &= e^{-2n/\beta}. \end{aligned}$$

Nov. 23rd 2021

**Definition 5.2.1** Let  $X_1, \dots, X_n$  be a random sample of size  $n$  from a population and let  $T(x_1, \dots, x_n)$  be a real-valued or vector-valued function whose domain includes the sample space of  $(X_1, \dots, X_n)$ . Then the random variable or random vector  $Y = T(X_1, \dots, X_n)$  is called a statistic. The probability distribution of a statistic  $Y$  is called the sampling distribution of  $Y$ .

**Definition 5.2.2** The sample mean is the arithmetic average of the values in a random sample. It is usually denoted by

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i.$$

**Definition 5.2.3** The sample variance is the statistic defined by

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

The sample standard deviation is the statistic defined by  $S = \sqrt{S^2}$ .

**Theorem 5.2.4** Let  $x_1, \dots, x_n$  be any numbers and  $\bar{x} = (x_1 + \dots + x_n)/n$ . Then

- $\min_a \sum_{i=1}^n (x_i - a)^2 = \sum_{i=1}^n (x_i - \bar{x})^2$ ,
- $(n-1)s^2 = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2$ .

**Proof:** To prove part (a), add and subtract  $\bar{x}$  to get

$$\begin{aligned} \sum_{i=1}^n (x_i - a)^2 &= \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - a)^2 \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + 2 \sum_{i=1}^n (x_i - \bar{x})(\bar{x} - a) + \sum_{i=1}^n (\bar{x} - a)^2 \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (\bar{x} - a)^2. \end{aligned} \quad (\text{cross term is 0})$$

It is now clear that the right-hand side is minimized at  $a = \bar{x}$ . (Notice the similarity to Example 2.2.6 and Exercise 4.13.)

To prove part (b), take  $a = 0$  in the above. □

**Theorem 5.2.6** Let  $X_1, \dots, X_n$  be a random sample from a population with mean  $\mu$  and variance  $\sigma^2 < \infty$ . Then

- $E\bar{X} = \mu$ ,
- $\text{Var } \bar{X} = \frac{\sigma^2}{n}$ ,
- $E S^2 = \sigma^2$ .

**Proof:** To prove (a), let  $g(X_i) = X_i/n$ , so  $Eg(X_i) = \mu/n$ . Then, by Lemma 5.2.5,

$$E\bar{X} = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} E\left(\sum_{i=1}^n X_i\right) = \frac{1}{n} nEX_1 = \mu.$$

Similarly for (b), we have

$$\text{Var } \bar{X} = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} n \text{Var } X_1 = \frac{\sigma^2}{n}.$$

For the sample variance, using Theorem 5.2.4, we have

$$\begin{aligned} ES^2 &= E\left(\frac{1}{n-1} \left[ \sum_{i=1}^n X_i^2 - n\bar{X}^2 \right] \right) \\ &= \frac{1}{n-1} (nEX_1^2 - nE\bar{X}^2) \\ &= \frac{1}{n-1} \left( n(\sigma^2 + \mu^2) - n \left( \frac{\sigma^2}{n} + \mu^2 \right) \right) = \sigma^2, \end{aligned}$$

establishing part (c) and proving the theorem.  $\square$

**Theorem 5.2.9** If  $X$  and  $Y$  are independent continuous random variables with pdfs  $f_X(x)$  and  $f_Y(y)$ , then the pdf of  $Z = X + Y$  is

$$(5.2.3) \quad f_Z(z) = \int_{-\infty}^{\infty} f_X(w)f_Y(z-w) dw.$$

**Proof:** Let  $W = X$ . The Jacobian of the transformation from  $(X, Y)$  to  $(Z, W)$  is 1. So using (4.3.2), we obtain the joint pdf of  $(Z, W)$  as

$$f_{Z,W}(z, w) = f_{X,Y}(w, z-w) = f_X(w)f_Y(z-w).$$

Integrating out  $w$ , we obtain the marginal pdf of  $Z$  as given in (5.2.3).  $\square$

**Example 5.2.10 (Sum of Cauchy random variables)** As an example of a situation where the mgf technique fails, consider sampling from a Cauchy distribution. We will eventually derive the distribution of  $\bar{Z}$ , the mean of  $Z_1, \dots, Z_n$ , iid Cauchy( $0, 1$ ) observations. We start, however, with the distribution of the sum of two independent Cauchy random variables and apply formula (5.2.3).

Let  $U$  and  $V$  be independent Cauchy random variables,  $U \sim \text{Cauchy}(0, \sigma)$  and  $V \sim \text{Cauchy}(0, \tau)$ ; that is,

$$f_U(u) = \frac{1}{\pi\sigma} \frac{1}{1 + (u/\sigma)^2}, \quad f_V(v) = \frac{1}{\pi\tau} \frac{1}{1 + (v/\tau)^2}, \quad \begin{array}{l} -\infty < u < \infty, \\ -\infty < v < \infty. \end{array}$$

Based on formula (5.2.3), the pdf of  $Z = U + V$  is given by

$$(5.2.4) \quad f_Z(z) = \int_{-\infty}^{\infty} \frac{1}{\pi\sigma} \frac{1}{1 + (w/\sigma)^2} \frac{1}{\pi\tau} \frac{1}{1 + ((z-w)/\tau)^2} dw, \quad -\infty < z < \infty.$$

This integral is somewhat involved but can be solved by a partial fraction decomposition and some careful antidifferentiation (see Exercise 5.7). The result is

$$(5.2.5) \quad f_Z(z) = \frac{1}{\pi(\sigma + \tau)} \frac{1}{1 + (z/(\sigma + \tau))^2}, \quad -\infty < z < \infty. \quad Z \sim \text{Cauchy}(0, \sigma + \tau)$$

Thus, the sum of two independent Cauchy random variables is again a Cauchy, with the scale parameters adding. It therefore follows that if  $Z_1, \dots, Z_n$  are iid  $\text{Cauchy}(0, 1)$  random variables, then  $\sum Z_i$  is  $\text{Cauchy}(0, n)$  and also  $\bar{Z}$  is  $\text{Cauchy}(0, 1)$ ! The sample mean has the same distribution as the individual observations. (See Example A.0.5 in Appendix A for a computer algebra version of this calculation.) ||

## 5.4 Order Statistics

**Definition 5.4.1** The order statistics of a random sample  $X_1, \dots, X_n$  are the sample values placed in ascending order. They are denoted by  $X_{(1)}, \dots, X_{(n)}$ .

The sample range,  $R = X_{(n)} - X_{(1)}$ , is the distance between the smallest and largest observations. It is a measure of the dispersion in the sample and should reflect the dispersion in the population.

The sample median, which we will denote by  $M$ , is a number such that approximately one-half of the observations are less than  $M$  and one-half are greater. In terms of the order statistics,  $M$  is defined by

$$(5.4.1) \quad M = \begin{cases} X_{((n+1)/2)} & \text{if } n \text{ is odd} \\ (X_{(n/2)} + X_{(n/2+1)}) / 2 & \text{if } n \text{ is even.} \end{cases}$$

mid-range  $V = \frac{X_{(1)} + X_{(n)}}{2}$

**Theorem 5.4.4** Let  $X_{(1)}, \dots, X_{(n)}$  denote the order statistics of a random sample,  $X_1, \dots, X_n$ , from a continuous population with cdf  $F_X(x)$  and pdf  $f_X(x)$ . Then the pdf of  $X_{(j)}$  is

$$(5.4.4) \quad f_{X_{(j)}}(x) = \frac{n!}{(j-1)!(n-j)!} f_X(x) [F_X(x)]^{j-1} [1 - F_X(x)]^{n-j}.$$

**Proof:** We first find the cdf of  $X_{(j)}$  and then differentiate it to obtain the pdf. As in Theorem 5.4.3, let  $Y$  be a random variable that counts the number of  $X_1, \dots, X_n$  less than or equal to  $x$ . Then, defining a “success” as the event  $\{X_j \leq x\}$ , we see that  $Y \sim \text{binomial}(n, F_X(x))$ . (Note that we can write  $P_i = F_X(x_i)$  in Theorem 5.4.3. Also, although  $X_1, \dots, X_n$  are continuous random variables, the counting variable  $Y$  is discrete.) Thus,

$$F_{X_{(j)}}(x) = P(Y \geq j) = \sum_{k=j}^n \binom{n}{k} [F_X(x)]^k [1 - F_X(x)]^{n-k},$$

and the pdf of  $X_{(j)}$  is

$$\begin{aligned}
 f_{X_{(j)}}(x) &= \frac{d}{dx} F_{X_{(j)}}(x) \\
 &= \sum_{k=j}^n \binom{n}{k} \left( k [F_X(x)]^{k-1} [1 - F_X(x)]^{n-k} f_X(x) \right. \\
 &\quad \left. - (n-k) [F_X(x)]^k [1 - F_X(x)]^{n-k-1} f_X(x) \right) \quad (\text{chain rule}) \\
 &= \binom{n}{j} j f_X(x) [F_X(x)]^{j-1} [1 - F_X(x)]^{n-j} \\
 &\quad + \sum_{k=j+1}^n \binom{n}{k} k [F_X(x)]^{k-1} [1 - F_X(x)]^{n-k} f_X(x) \\
 &\quad - \sum_{k=j}^{n-1} \binom{n}{k} (n-k) [F_X(x)]^k [1 - F_X(x)]^{n-k-1} f_X(x) \quad \left( \begin{array}{l} k = n \text{ term} \\ \text{is } 0 \end{array} \right) \\
 &= \frac{n!}{(j-1)!(n-j)!} f_X(x) [F_X(x)]^{j-1} [1 - F_X(x)]^{n-j} \\
 (5.4.5) \quad &\quad + \sum_{k=j}^{n-1} \binom{n}{k+1} (k+1) [F_X(x)]^k [1 - F_X(x)]^{n-k-1} f_X(x) \quad \left( \begin{array}{l} \text{change} \\ \text{dummy} \\ \text{variable} \end{array} \right) \\
 &\quad - \sum_{k=j}^{n-1} \binom{n}{k} (n-k) [F_X(x)]^k [1 - F_X(x)]^{n-k-1} f_X(x).
 \end{aligned}$$

Noting that

$$(5.4.6) \quad \binom{n}{k+1} (k+1) = \frac{n!}{k!(n-k-1)!} = \binom{n}{k} (n-k),$$

we see that the last two sums in (5.4.5) cancel. Thus, the pdf  $f_{X_{(j)}}(x)$  is given by the expression in (5.4.4).  $\square$

**Example 5.4.5 (Uniform order statistic pdf)** Let  $X_1, \dots, X_n$  be iid uniform(0, 1), so  $f_X(x) = 1$  for  $x \in (0, 1)$  and  $F_X(x) = x$  for  $x \in (0, 1)$ . Using (5.4.4), we see that the pdf of the  $j$ th order statistic is

$$\begin{aligned}
 f_{X_{(j)}}(x) &= \frac{n!}{(j-1)!(n-j)!} x^{j-1} (1-x)^{n-j} \quad \text{for } x \in (0, 1) \\
 &= \frac{\Gamma(n+1)}{\Gamma(j)\Gamma(n-j+1)} x^{j-1} (1-x)^{(n-j+1)-1}.
 \end{aligned}$$

Thus, the  $j$ th order statistic from a uniform(0, 1) sample has a beta( $j, n-j+1$ ) distribution. From this we can deduce that

$$\text{EX}_{(j)} = \frac{j}{n+1} \quad \text{and} \quad \text{Var } X_{(j)} = \frac{j(n-j+1)}{(n+1)^2(n+2)}. \quad \|$$

**Theorem 5.4.6** Let  $X_{(1)}, \dots, X_{(n)}$  denote the order statistics of a random sample,  $X_1, \dots, X_n$ , from a continuous population with cdf  $F_X(x)$  and pdf  $f_X(x)$ . Then the joint pdf of  $X_{(i)}$  and  $X_{(j)}$ ,  $1 \leq i < j \leq n$ , is

$$(5.4.7) \quad f_{X_{(i)}, X_{(j)}}(u, v) = \frac{n!}{(i-1)!(j-1-i)!(n-j)!} f_X(u) f_X(v) [F_X(u)]^{i-1} \\ \times [F_X(v) - F_X(u)]^{j-1-i} [1 - F_X(v)]^{n-j}$$

for  $-\infty < u < v < \infty$ .

The joint pdf of three or more order statistics could be derived using similar but even more involved arguments. Perhaps the other most useful pdf is  $f_{X_{(1)}, \dots, X_{(n)}}(x_1, \dots, x_n)$ , the joint pdf of all the order statistics, which is given by

$$f_{X_{(1)}, \dots, X_{(n)}}(x_1, \dots, x_n) = \begin{cases} n! f_X(x_1) \cdots f_X(x_n) & -\infty < x_1 < \cdots < x_n < \infty \\ 0 & \text{otherwise.} \end{cases}$$

**Example 5.4.7 (Distribution of the midrange and range)** Let  $X_1, \dots, X_n$  be iid uniform( $0, a$ ) and let  $X_{(1)}, \dots, X_{(n)}$  denote the order statistics. The range was earlier defined as  $R = X_{(n)} - X_{(1)}$ . The *midrange*, a measure of location like the sample median or the sample mean, is defined by  $V = (X_{(1)} + X_{(n)})/2$ . We will derive the joint pdf of  $R$  and  $V$  from the joint pdf of  $X_{(1)}$  and  $X_{(n)}$ . From (5.4.7) we have that

$$f_{X_{(1)}, X_{(n)}}(x_1, x_n) = \frac{n(n-1)}{a^2} \left( \frac{x_n}{a} - \frac{x_1}{a} \right)^{n-2} \\ = \frac{n(n-1)(x_n - x_1)^{n-2}}{a^n}, \quad 0 < x_1 < x_n < a.$$

Solving for  $X_{(1)}$  and  $X_{(n)}$ , we obtain  $X_{(1)} = V - R/2$  and  $X_{(n)} = V + R/2$ . The Jacobian for this transformation is  $-1$ . The transformation from  $(X_{(1)}, X_{(n)})$  to  $(R, V)$  maps  $\{(x_1, x_n) : 0 < x_1 < x_n < a\}$  onto the set  $\{(r, v) : 0 < r < a, r/2 < v < a - r/2\}$ . To see this, note that obviously  $0 < r < a$  and for a fixed value of  $r$ ,  $v$  ranges from  $r/2$  (corresponding to  $x_1 = 0, x_n = r$ ) to  $a - r/2$  (corresponding to  $x_1 = a - r, x_n = a$ ). Thus, the joint pdf of  $(R, V)$  is

$$f_{R, V}(r, v) = \frac{n(n-1)r^{n-2}}{a^n}, \quad 0 < r < a, \quad r/2 < v < a - r/2.$$

The marginal pdf of  $R$  is thus

$$(5.4.8) \quad f_R(r) = \int_{r/2}^{a-r/2} \frac{n(n-1)r^{n-2}}{a^n} dv \\ = \frac{n(n-1)r^{n-2}(a-r)}{a^n}, \quad 0 < r < a.$$

If  $a = 1$ , we see that  $r$  has a beta( $n - 1, 2$ ) distribution. Or, for arbitrary  $a$ , it is easy to deduce from (5.4.8) that  $R/a$  has a beta distribution. Note that the constant  $a$  is a scale parameter.

The set where  $f_{R,V}(r, v) > 0$  is shown in Figure 5.4.1, where we see that the range of integration of  $r$  depends on whether  $v > a/2$  or  $v \leq a/2$ . Thus, the marginal pdf of  $V$  is given by

$$f_V(v) = \int_0^{2v} \frac{n(n-1)r^{n-2}}{a^n} dr = \frac{n(2v)^{n-1}}{a^n}, \quad 0 < v \leq a/2,$$

and

$$f_V(v) = \int_0^{2(a-v)} \frac{n(n-1)r^{n-2}}{a^n} dr = \frac{n[2(a-v)]^{n-1}}{a^n}, \quad a/2 < v \leq a.$$

This pdf is symmetric about  $a/2$  and has a peak at  $a/2$ . ||

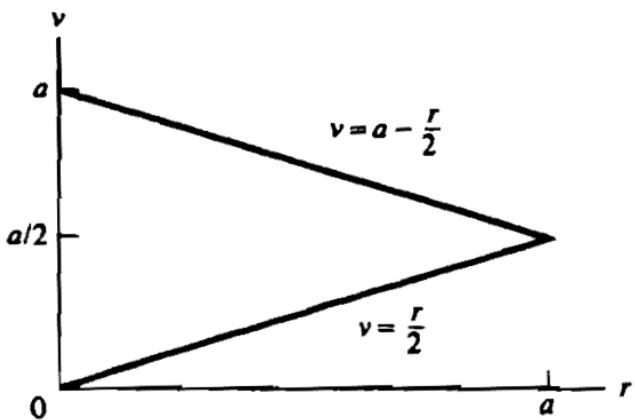


Figure 5.4.1. Region on which  $f_{R,V}(r, v) > 0$  for Example 5.4.7

Nov. 30th 2021

## 5.5 Convergence Concepts

**Definition 5.5.1** A sequence of random variables,  $X_1, X_2, \dots$ , converges in probability to a random variable  $X$  if, for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0 \quad \text{or, equivalently,} \quad \lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1.$$

**Theorem 5.5.2 (Weak Law of Large Numbers)** Let  $X_1, X_2, \dots$  be iid random variables with  $EX_i = \mu$  and  $\text{Var } X_i = \sigma^2 < \infty$ . Define  $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$ . Then, for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < \epsilon) = 1;$$

that is,  $\bar{X}_n$  converges in probability to  $\mu$ .  $\bar{X}_n \xrightarrow{P} \mu$

**Proof:** The proof is quite simple, being a straightforward application of Chebychev's Inequality. We have, for every  $\epsilon > 0$ ,

$$P(|\bar{X}_n - \mu| \geq \epsilon) = P((\bar{X}_n - \mu)^2 \geq \epsilon^2) \leq \frac{E(\bar{X}_n - \mu)^2}{\epsilon^2} = \frac{\text{Var } \bar{X}_n}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}.$$

Hence,  $P(|\bar{X}_n - \mu| < \epsilon) = 1 - P(|\bar{X}_n - \mu| \geq \epsilon) \geq 1 - \sigma^2/(n\epsilon^2) \rightarrow 1$ , as  $n \rightarrow \infty$ .  $\square$

**Example 5.5.3 (Consistency of  $S^2$ )** Suppose we have a sequence  $X_1, X_2, \dots$  of iid random variables with  $EX_i = \mu$  and  $\text{Var } X_i = \sigma^2 < \infty$ . If we define

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2,$$

can we prove a WLLN for  $S_n^2$ ? Using Chebychev's Inequality, we have

$$P(|S_n^2 - \sigma^2| \geq \epsilon) \leq \frac{E(S_n^2 - \sigma^2)^2}{\epsilon^2} = \frac{\text{Var } S_n^2}{\epsilon^2}$$

and thus, a sufficient condition that  $S_n^2$  converges in probability to  $\sigma^2$  is that  $\text{Var } S_n^2 \rightarrow 0$  as  $n \rightarrow \infty$ .  $\parallel$

**Theorem 5.5.4** Suppose that  $X_1, X_2, \dots$  converges in probability to a random variable  $X$  and that  $h$  is a continuous function. Then  $h(X_1), h(X_2), \dots$  converges in probability to  $h(X)$ .

**Example 5.5.5 (Consistency of  $S$ )** If  $S_n^2$  is a consistent estimator of  $\sigma^2$ , then by Theorem 5.5.4, the sample standard deviation  $S_n = \sqrt{S_n^2} = h(S_n^2)$  is a consistent estimator of  $\sigma$ . Note that  $S_n$  is, in fact, a biased estimator of  $\sigma$  (see Exercise 5.11), but the bias disappears asymptotically.  $\parallel$

**Definition 5.5.6** A sequence of random variables,  $X_1, X_2, \dots$ , converges almost surely to a random variable  $X$  if, for every  $\epsilon > 0$ ,

$$P(\lim_{n \rightarrow \infty} |X_n - X| < \epsilon) = 1.$$

**Example 5.5.7 (Almost sure convergence)** Let the sample space  $S$  be the closed interval  $[0, 1]$  with the uniform probability distribution. Define random variables  $X_n(s) = s + s^n$  and  $X(s) = s$ . For every  $s \in [0, 1]$ ,  $s^n \rightarrow 0$  as  $n \rightarrow \infty$  and  $X_n(s) \rightarrow s = X(s)$ . However,  $X_n(1) = 2$  for every  $n$  so  $X_n(1)$  does not converge to  $1 = X(1)$ . But since the convergence occurs on the set  $[0, 1]$  and  $P([0, 1]) = 1$ ,  $X_n$  converges to  $X$  almost surely.  $\parallel$

**Example 5.5.8 (Convergence in probability, not almost surely)** In this example we describe a sequence that converges in probability, but not almost surely. Again, let the sample space  $S$  be the closed interval  $[0, 1]$  with the uniform probability distribution. Define the sequence  $X_1, X_2, \dots$  as follows:

$$X_1(s) = s + I_{[0,1]}(s), \quad X_2(s) = s + I_{[0,\frac{1}{2}]}(s), \quad X_3(s) = s + I_{[\frac{1}{2},1]}(s),$$

$$X_4(s) = s + I_{[0,\frac{1}{3}]}(s), \quad X_5(s) = s + I_{[\frac{1}{3},\frac{2}{3}]}(s), \quad X_6(s) = s + I_{[\frac{2}{3},1]}(s),$$

etc. Let  $X(s) = s$ . It is straightforward to see that  $X_n$  converges to  $X$  in probability. As  $n \rightarrow \infty$ ,  $P(|X_n - X| \geq \epsilon)$  is equal to the probability of an interval of  $s$  values whose length is going to 0. However,  $X_n$  does not converge to  $X$  almost surely. Indeed, there is no value of  $s \in S$  for which  $X_n(s) \rightarrow s = X(s)$ . For every  $s$ , the value  $X_n(s)$  alternates between the values  $s$  and  $s + 1$  infinitely often. For example, if  $s = \frac{3}{8}$ ,  $X_1(s) = 1\frac{3}{8}$ ,  $X_2(s) = 1\frac{3}{8}$ ,  $X_3(s) = \frac{3}{8}$ ,  $X_4(s) = \frac{3}{8}$ ,  $X_5(s) = 1\frac{3}{8}$ ,  $X_6(s) = \frac{3}{8}$ , etc. No pointwise convergence occurs for this sequence. ||

**Theorem 5.5.9 (Strong Law of Large Numbers)** *Let  $X_1, X_2, \dots$  be iid random variables with  $\mathbb{E}X_i = \mu$  and  $\text{Var } X_i = \sigma^2 < \infty$ , and define  $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$ . Then, for every  $\epsilon > 0$ ,*

$$P(\lim_{n \rightarrow \infty} |\bar{X}_n - \mu| < \epsilon) = 1;$$

that is,  $\bar{X}_n$  converges almost surely to  $\mu$ .

**Definition 5.5.10** A sequence of random variables,  $X_1, X_2, \dots$ , converges in distribution to a random variable  $X$  if

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$$

$\xrightarrow{\text{d}}$   
 $\Rightarrow$

at all points  $x$  where  $F_X(x)$  is continuous.

**Example 5.5.11 (Maximum of uniforms)** If  $X_1, X_2, \dots$  are iid uniform(0, 1) and  $X_{(n)} = \max_{1 \leq i \leq n} X_i$ , let us examine if (and to where)  $X_{(n)}$  converges in distribution.

As  $n \rightarrow \infty$ , we expect  $X_{(n)}$  to get close to 1 and, as  $X_{(n)}$  must necessarily be less than 1, we have for any  $\epsilon > 0$ ,

$$\begin{aligned} P(|X_{(n)} - 1| \geq \epsilon) &= P(X_{(n)} \geq 1 + \epsilon) + P(X_{(n)} \leq 1 - \epsilon) \\ &= 0 + P(X_{(n)} \leq 1 - \epsilon). \end{aligned}$$

Next using the fact that we have an *iid* sample, we can write

$$P(X_{(n)} \leq 1 - \epsilon) = P(X_i \leq 1 - \epsilon, i = 1, \dots, n) = (1 - \epsilon)^n,$$

which goes to 0. So we have proved that  $X_{(n)}$  converges to 1 in probability. However, if we take  $\epsilon = t/n$ , we then have

$$P(X_{(n)} \leq 1 - t/n) = (1 - t/n)^n \rightarrow e^{-t},$$

which, upon rearranging, yields

$$P(n(1 - X_{(n)}) \leq t) \rightarrow 1 - e^{-t};$$

that is, the random variable  $n(1 - X_{(n)})$  converges in distribution to an exponential(1) random variable. ||

**Theorem 5.5.14 (Central Limit Theorem)** Let  $X_1, X_2, \dots$  be a sequence of iid random variables whose mgfs exist in a neighborhood of 0 (that is,  $M_{X_i}(t)$  exists for  $|t| < h$ , for some positive  $h$ ). Let  $EX_i = \mu$  and  $\text{Var } X_i = \sigma^2 > 0$ . (Both  $\mu$  and  $\sigma^2$  are finite since the mgf exists.) Define  $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$ . Let  $G_n(x)$  denote the cdf of  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$ . Then, for any  $x$ ,  $-\infty < x < \infty$ ,

$$\lim_{n \rightarrow \infty} G_n(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy;$$

that is,  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$  has a limit in a standard normal distribution.

**Proof of Theorem 5.5.14:** We will show that, for  $|t| < h$ , the mgf of  $\sqrt{n}(\bar{X}_n - \mu)/\sigma$  converges to  $e^{t^2/2}$ , the mgf of a  $n(0, 1)$  random variable.

Define  $Y_i = (X_i - \mu)/\sigma$ , and let  $M_Y(t)$  denote the common mgf of the  $Y_i$ s, which exists for  $|t| < \sigma h$  and is given by Theorem 2.3.15. Since

$$(5.5.1) \quad \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i,$$

we have, from the properties of mgfs (see Theorems 2.3.15 and 4.6.7),

$$(5.5.2) \quad \begin{aligned} M_{\sqrt{n}(\bar{X}_n - \mu)/\sigma}(t) &= M_{\sum_{i=1}^n Y_i / \sqrt{n}}(t) \\ &= M_{\sum_{i=1}^n Y_i} \left( \frac{t}{\sqrt{n}} \right) && (\text{Theorem 2.3.15}) \\ &= \left( M_Y \left( \frac{t}{\sqrt{n}} \right) \right)^n. && (\text{Theorem 4.6.7}) \end{aligned}$$

We now expand  $M_Y(t/\sqrt{n})$  in a Taylor series (power series) around 0. (See Definition 5.5.20.) We have

$$(5.5.3) \quad M_Y \left( \frac{t}{\sqrt{n}} \right) = \sum_{k=0}^{\infty} M_Y^{(k)}(0) \frac{(t/\sqrt{n})^k}{k!},$$

where  $M_Y^{(k)}(0) = (d^k/dt^k) M_Y(t)|_{t=0}$ . Since the mgfs exist for  $|t| < h$ , the power series expansion is valid if  $t < \sqrt{n}\sigma h$ .

Using the facts that  $M_Y^{(0)} = 1$ ,  $M_Y^{(1)} = 0$ , and  $M_Y^{(2)} = 1$  (by construction, the mean and variance of  $Y$  are 0 and 1), we have

$$(5.5.4) \quad M_Y \left( \frac{t}{\sqrt{n}} \right) = 1 + \frac{(t/\sqrt{n})^2}{2!} + R_Y \left( \frac{t}{\sqrt{n}} \right),$$

where  $R_Y$  is the remainder term in the Taylor expansion,

$$R_Y \left( \frac{t}{\sqrt{n}} \right) = \sum_{k=3}^{\infty} M_Y^{(k)}(0) \frac{(t/\sqrt{n})^k}{k!}.$$

An application of Taylor's Theorem (Theorem 5.5.21) shows that, for fixed  $t \neq 0$ , we have

$$\lim_{n \rightarrow \infty} \frac{R_Y(t/\sqrt{n})}{(t/\sqrt{n})^2} = 0.$$

Since  $t$  is fixed, we also have

$$(5.5.5) \quad \lim_{n \rightarrow \infty} \frac{R_Y(t/\sqrt{n})}{(1/\sqrt{n})^2} = \lim_{n \rightarrow \infty} n R_Y\left(\frac{t}{\sqrt{n}}\right) = 0,$$

and (5.5.5) is also true at  $t = 0$  since  $R_Y(0/\sqrt{n}) = 0$ . Thus, for any fixed  $t$ , we can write

$$(5.5.6) \quad \begin{aligned} \lim_{n \rightarrow \infty} \left( M_Y\left(\frac{t}{\sqrt{n}}\right) \right)^n &= \lim_{n \rightarrow \infty} \left[ 1 + \frac{(t/\sqrt{n})^2}{2!} + R_Y\left(\frac{t}{\sqrt{n}}\right) \right]^n \\ &= \lim_{n \rightarrow \infty} \left[ 1 + \frac{1}{n} \left( \frac{t^2}{2} + n R_Y\left(\frac{t}{\sqrt{n}}\right) \right) \right]^n \\ &= e^{t^2/2} \end{aligned}$$

by an application of Lemma 2.3.14, where we set  $a_n = (t^2/2) + n R_Y(t/\sqrt{n})$ . (Note that (5.5.5) implies that  $a_n \rightarrow t^2/2$  as  $n \rightarrow \infty$ .) Since  $e^{t^2/2}$  is the mgf of the  $n(0, 1)$  distribution, the theorem is proved.  $\square$

**Theorem 5.5.17 (Slutsky's Theorem)** *If  $X_n \rightarrow X$  in distribution and  $Y_n \rightarrow a$ , a constant, in probability, then*

- a.  $Y_n X_n \rightarrow aX$  in distribution.
- b.  $X_n + Y_n \rightarrow X + a$  in distribution.

**Example 5.5.18 (Normal approximation with estimated variance)** Suppose that

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightarrow n(0, 1),$$

but the value of  $\sigma$  is unknown. We have seen in Example 5.5.3 that, if  $\lim_{n \rightarrow \infty} \text{Var } S_n^2 = 0$ , then  $S_n^2 \rightarrow \sigma^2$  in probability. By Exercise 5.32,  $\sigma/S_n \rightarrow 1$  in probability. Hence, Slutsky's Theorem tells us

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} = \frac{\sigma}{S_n} \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightarrow n(0, 1). \quad \parallel$$

### 5.5.4 The Delta Method

The previous section gives conditions under which a standardized random variable has a limit normal distribution. There are many times, however, when we are not specifically interested in the distribution of the random variable itself, but rather some function of the random variable.

**Theorem 5.5.24 (Delta Method)** Let  $Y_n$  be a sequence of random variables that satisfies  $\sqrt{n}(Y_n - \theta) \rightarrow n(0, \sigma^2)$  in distribution. For a given function  $g$  and a specific value of  $\theta$ , suppose that  $g'(\theta)$  exists and is not 0. Then

$$(5.5.10) \quad \sqrt{n}[g(Y_n) - g(\theta)] \rightarrow n(0, \sigma^2[g'(\theta)]^2) \text{ in distribution.}$$

**Proof:** The Taylor expansion of  $g(Y_n)$  around  $Y_n = \theta$  is

$$(5.5.11) \quad g(Y_n) = g(\theta) + g'(\theta)(Y_n - \theta) + \text{Remainder},$$

where the remainder  $\rightarrow 0$  as  $Y_n \rightarrow \theta$ . Since  $Y_n \rightarrow \theta$  in probability it follows that the remainder  $\rightarrow 0$  in probability. By applying Slutsky's Theorem (Theorem 5.5.17) to

$$\sqrt{n}[g(Y_n) - g(\theta)] = g'(\theta)\sqrt{n}(Y_n - \theta),$$

the result now follows. See Exercise 5.43 for details. □

Ex.  $X_1, X_2 \stackrel{\text{i.i.d.}}{\sim} E(X_i) = \mu, \text{Var}(X_i) = \sigma^2 > 0$   
 $\sqrt{n}(\bar{X} - \mu) \Rightarrow N(0, \sigma^2)$

Define  $g(x) = \frac{1}{x}$

Delta method  $\sqrt{n}\left(\frac{1}{\bar{X}} - \frac{1}{\mu}\right) \Rightarrow N(0, \sigma^2 \cdot \frac{1}{\mu^4}) ; \mu \neq 0$

$\left[ \begin{array}{l} S_n \xrightarrow{P} \sigma^2 \\ \bar{X}_n \xrightarrow{P} \mu \end{array} \right] \frac{\sqrt{n}\left(\frac{1}{\bar{X}} - \frac{1}{\mu}\right)}{S_n \cdot \left(\frac{1}{\bar{X}}\right)^2} \Rightarrow N(0, 1) ; \mu \neq 0$