

Probability & Statistics Review

Dave Goldsman

Updated 8/21/12

1 Getting Started — The Gambler's Ruin

Each time a gambler plays, he wins \$1 with probability p and loses \$1 with probability $1 - p = q$. Each play is independent. Suppose he starts with \$ i . Find the probability that his fortune will hit \$ N (i.e., he breaks the bank) before it hits \$0 (i.e., he is ruined).

Let X_n denote his fortune at time n . It turns out that X_1, X_2, \dots is a *Markov chain* — a stochastic process where the next state depends only the current state (more on this later).

To get things going, we'll use a common trick — a so-called *one-step analysis*. Let

$$\begin{aligned} P_i &\equiv \Pr(\text{Eventually hit } \$N | X_0 = i) \\ &= \Pr(\text{Event. hit } N | X_1 = i + 1 \text{ and } X_0 = i) \Pr(X_1 = i + 1 | X_0 = i) \\ &\quad + \Pr(\text{Event. hit } N | X_1 = i - 1 \text{ and } X_0 = i) \Pr(X_1 = i - 1 | X_0 = i) \\ &= \Pr(\text{Event. hit } N | X_1 = i + 1)p + \Pr(\text{Event. hit } N | X_1 = i - 1)q \\ &= pP_{i+1} + qP_{i-1}, \quad i = 1, 2, \dots, N - 1. \end{aligned}$$

Since $p + q = 1$, we have

$$pP_i + qP_i = pP_{i+1} + qP_{i-1}$$

iff

$$p(P_{i+1} - P_i) = q(P_i - P_{i-1})$$

iff

$$P_{i+1} - P_i = \frac{q}{p}(P_i - P_{i-1}), \quad i = 1, 2, \dots, N - 1.$$

Since $P_0 = 0$, we have

$$\begin{aligned} P_2 - P_1 &= \frac{q}{p}P_1 \\ P_3 - P_2 &= \frac{q}{p}(P_2 - P_1) = \left(\frac{q}{p}\right)^2 P_1 \\ &\vdots \\ P_i - P_{i-1} &= \frac{q}{p}(P_{i-1} - P_{i-2}) = \left(\frac{q}{p}\right)^{i-1} P_1. \end{aligned}$$

Summing up the LHS terms and the RHS terms,

$$\sum_{j=2}^i (P_j - P_{j-1}) = P_i - P_1 = \sum_{j=1}^{i-1} \left(\frac{q}{p}\right)^j P_1.$$

This implies that

$$P_i = P_1 \sum_{j=0}^{i-1} \left(\frac{q}{p}\right)^j = \begin{cases} \frac{1-(q/p)^i}{1-(q/p)} P_1 & \text{if } q \neq p \text{ } (p \neq 1/2) \\ iP_1 & \text{if } q = p \text{ } (p = 1/2) \end{cases}.$$

In particular, note that

$$1 = P_N = \begin{cases} \frac{1-(q/p)^N}{1-(q/p)} P_1 & \text{if } p \neq 1/2 \\ NP_1 & \text{if } p = 1/2 \end{cases}.$$

Thus,

$$P_1 = \begin{cases} \frac{1-(q/p)^N}{1-(q/p)} & \text{if } p \neq 1/2 \\ 1/N & \text{if } p = 1/2 \end{cases},$$

so that

$$P_i = \begin{cases} \frac{1-(q/p)^i}{1-(q/p)^N} & \text{if } p \neq 1/2 \\ i/N & \text{if } p = 1/2 \end{cases}. \quad \diamond$$

By the way, as $N \rightarrow \infty$,

$$P_i \rightarrow \begin{cases} 1 - (q/p)^i & \text{if } p > 1/2 \\ 0 & \text{if } p \leq 1/2 \end{cases}. \quad \diamond$$

Example: A guy can somehow win any blackjack hand w.p. 0.6. If he wins, he fortune increases by \$100; a loss costs him \$100. Suppose he starts out with \$500, and that he'll quit playing as soon as his fortune hits \$0 or \$1500. What's the probability that he'll eventually hit \$1500?

$$P_5 = \frac{1 - (0.4/0.6)^5}{1 - (0.4/0.6)^{15}} = 0.870. \quad \diamond$$

2 Probability Review

2.1 Preliminaries

To start with, I'll assume that you know the following things:

- The set of all possible outcomes of an experiment is the *sample space* Ω .
- Any subset E of a sample space Ω is an *event*.
- The set of all possible events is denoted by \mathcal{F} , which is called a *sigma field* of Ω .

For example, if $\Omega = \{H, T\}$, then $\mathcal{F} = \{\emptyset, \{H\}, \{T\}, \{H, T\}\}$.

A sigma field must satisfy the following:

1. $A \in \mathcal{F}$ implies the complement $\bar{A} \in \mathcal{F}$.
 2. $A_1, A_2, \dots \in \mathcal{F}$ implies $\cup_{j=1}^{\infty} A_j \in \mathcal{F}$.
- The *probability* function $P(\cdot)$ must satisfy 3 axioms:
 1. For any event $E \in \mathcal{F}$, we must have $0 \leq P(E) \leq 1$
 2. $P(\Omega) = 1$
 3. For any *disjoint* sequence of events E_1, E_2, \dots (i.e., $E_i \cap E_j = \emptyset$ if $i \neq j$), we have $P(\cup_{i=1}^{\infty} E_i) = \sum_{i=1}^{\infty} P(E_i)$.
 - A *probability space* is the triple (Ω, \mathcal{F}, P) .

2.2 Conditional Probability and Independence

Definition: If $P(B) > 0$, then $P(A|B) \equiv P(A \cap B)/P(B)$ is the *conditional probability* of A given B .

Example: Toss a fair die. Let $A = \{1, 2, 3\}$ and $B = \{3, 4, 5, 6\}$. Then

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{1/6}{4/6} = 1/4. \quad \diamond$$

Definition: If $P(A \cap B) = P(A)P(B)$, then A and B are *independent* events.

Theorem: If A and B are independent, then $P(A|B) = P(A)$.

Proof: Easy. \diamond

Example: Toss two dice. Let $A =$ “Sum is 7” and $B =$ “First die is 4”. Then $P(A) = 1/6$, $P(B) = 1/6$, and $P(A \cap B) = P((4, 3)) = 1/36 = P(A)P(B)$; so A and B are independent. \diamond

2.3 Random Variables

Definition: A *random variable* (RV) X is a function from Ω to the real line \mathbb{R} , i.e., $X : \Omega \rightarrow \mathbb{R}$.

Example: Let X be the sum of two dice rolls. Then $X((4, 6)) = 10$. In addition,

$$P(X = x) = \begin{cases} 1/36 & \text{if } x = 2 \\ 2/36 & \text{if } x = 3 \\ \vdots & \\ 1/36 & \text{if } x = 12 \\ 0 & \text{otherwise} \end{cases} \quad \diamond$$

Definition: If the number of possible values of a RV X is finite or countably infinite, then X is a *discrete* RV. Its *probability mass function* (pmf) is $f(x) \equiv P(X = x)$. Note that $\sum_x f(x) = 1$.

Example: Flip 2 coins. Let X be the number of heads.

$$f(x) = \begin{cases} 1/4 & \text{if } x = 0 \text{ or } 2 \\ 1/2 & \text{if } x = 1 \\ 0 & \text{otherwise} \end{cases} \quad \diamond$$

Examples: Here are some well-known discrete RV's that you should review: Bernoulli(p), Binomial(n, p), Geometric(p), Negative Binomial, Poisson(λ), etc.

Definition: A *continuous* RV is one with probability zero at every individual point. A RV is continuous if there exists a *probability density function* (pdf) $f(x)$ such that $P(X \in A) = \int_A f(x) dx$ for every set A . Note that $\int_{\mathbb{R}} f(x) dx = 1$.

Example: Pick a random number between 3 and 7. Then

$$f(x) = \begin{cases} 1/4 & \text{if } 3 \leq x \leq 7 \\ 0 & \text{otherwise} \end{cases} \quad \diamond$$

Examples: Here are some well-known continuous RV's that you should review: $\text{Uniform}(a, b)$, $\text{Exponential}(\lambda)$, $\text{Normal}(\mu, \sigma^2)$, etc.

Definition: For any RV X (discrete or continuous), the *cumulative distribution function* (cdf) is

$$F(x) \equiv P(X \leq x) = \begin{cases} \sum_{y \leq x} f(y) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^x f(y) dy & \text{if } X \text{ is continuous} \end{cases}$$

For convenience, we'll henceforth write $F(x) = \int_{-\infty}^x dF(y)$ to denote *both* the discrete and continuous cases. Note that $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow \infty} F(x) = 1$.

Example: Flip 2 coins. Let X be the number of heads.

$$F(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1/4 & \text{if } 0 \leq x < 1 \\ 3/4 & \text{if } 1 \leq x < 2 \\ 1 & \text{if } x \geq 2 \end{cases} \quad \diamond$$

Example: Suppose $X \sim \text{Exp}(\lambda)$ (i.e., X has the exponential distribution with parameter λ). Then $f(x) = \lambda e^{-\lambda x}$, $x \geq 0$, and the cdf is $F(x) = 1 - e^{-\lambda x}$, $x \geq 0$. \diamond

2.4 Expectation

Definition: The *expected value* (or *mean*) of a RV X is

$$\mathbb{E}[X] \equiv \int_{\mathbb{R}} x dF(x) = \begin{cases} \sum_x x P(X = x) & \text{if } X \text{ is discrete} \\ \int_{\mathbb{R}} x f(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

Example: Suppose that $X \sim \text{Bernoulli}(p)$. Then

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

and we have $\mathbb{E}[X] = \sum_x x f(x) = p$. \diamond

Example: Suppose that $X \sim \text{Uniform}(a, b)$. Then

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$

and we have $\mathbb{E}[X] = \int_{\mathbb{R}} x f(x) dx = (a + b)/2$. \diamond

“Definition” (the “law of the unconscious statistician”): Suppose that $g(X)$ is some function of the RV X . Then $E[g(X)] \equiv \int_{\mathbb{R}} g(x) dF(x)$.

Example: Suppose X is the following discrete RV:

x	2	3	4
$f(x)$	0.3	0.6	0.1

Then $E[X^3] = \sum_x x^3 f(x) = 8(0.3) + 27(0.6) + 64(0.1) = 25$. \diamond

Example: Suppose $X \sim U(0, 2)$. Then $E[X^n] = \int_{\mathbb{R}} x^n f(x) dx = 2^n/(n+1)$. \diamond

Definitions: $E[X^n]$ is the n th *moment* of X . $E[(X - E[X])^n]$ is the n th *central moment* of X . $\text{Var}(X) \equiv E[(X - E[X])^2]$ is the *variance* of X .

Theorem: $\text{Var}(X) = E[X^2] - (E[X])^2$.

Proof: Easy. \diamond

Example: Suppose $X \sim \text{Bern}(p)$. Recall that $E[X] = p$. Further,

$$E[X^2] = \sum_x x^2 f(x) = 0^2(1-p) + 1^2 p = p$$

and

$$\text{Var}(X) = E[X^2] - (E[X])^2 = p - p^2 = p(1-p). \quad \diamond$$

Example: Suppose $X \sim U(0, 2)$. By previous examples, $E[X] = 1$ and $E[X^2] = 4/3$. So $\text{Var}(X) = E[X^2] - (E[X])^2 = 1/3$. \diamond

Theorem: $E[aX + b] = aE[X] + b$ and $\text{Var}(aX + b) = a^2 \text{Var}(X)$.

Proof: Easy. \diamond

Definition: The *moment generating function* (mgf) of X is $M_X(t) \equiv E[e^{tX}]$. For now, we'll assume that this expectation is finite in a neighborhood of $t = 0$.

Example: Suppose $X \sim \text{Bern}(p)$. $M_X(t) = \sum_x e^{tx} f(x) = pe^t + q$. \diamond

Example: Suppose $X \sim \text{Exp}(\lambda)$. $M_X(t) = \int_0^\infty e^{tx} \lambda e^{-\lambda x} dx = \lambda/(\lambda - t)$, $t < \lambda$. \diamond

Theorem (why we call them moment generating functions): Assuming that the mgf exists in a neighborhood around $t = 0$,

$$\mathbb{E}[X^k] = \frac{d^k}{dt^k} M_X(t)|_{t=0}, \quad k = 1, 2, \dots$$

“Proof” We’ll just do the first moment. (The others are similar.)

$$\frac{d}{dt} M_X(t) = \frac{d}{dt} \mathbb{E}[e^{tX}] \stackrel{“=”}{=} \mathbb{E} \left[\frac{d}{dt} e^{tX} \right] \stackrel{“=”}{=} \mathbb{E}[X e^{tX}].$$

If you believe the above steps, then

$$\frac{d}{dt} M_X(t)|_{t=0} = \mathbb{E}[X]. \quad \diamond$$

Theorem: If X and Y have the same mgf, then they have the *same distribution*. Note that there are problems if the mgf doesn’t exist around $t = 0$.

Bonus Definition (which is sometimes useful): The *probability generating function* (pgf) of a random variable X is $P_X(s) \equiv \mathbb{E}[s^X]$. It can be shown that

$$\frac{d^k}{ds^k} P_X(s)|_{s=1} = \mathbb{E}[X(X-1)\cdots(X-k+1)], \quad k = 1, 2, \dots$$

2.5 Functions of a RV

Problem: Suppose we have a RV X with p.d.f./p.m.f. $f(x)$. Let $Y = h(X)$. Find $g(y)$, the p.d.f./p.m.f. of Y .

Discrete case: If X is discrete, then Y will be discrete, in which case

$$g(y) = \Pr(Y = y) = \Pr[h(X) = y] = \Pr\{x : h(x) = y\} = \sum_{x:h(x)=y} f(x).$$

Example: Let X denote the number of H ’s from two coin tosses. We want the p.m.f. for $Y = X^2 - X$.

x	0	1	2
$f(x)$	1/4	1/2	1/4
$y = x^2 - x$	0	0	2

This implies that $g(0) = \Pr(Y = 0) = \Pr(X = 0 \text{ or } 1) = 3/4$ and $g(2) = \Pr(Y = 2) = 1/4$. In other words,

$$g(y) = \begin{cases} 3/4 & \text{if } y = 0 \\ 1/4 & \text{if } y = 2 \\ 0 & \text{otherwise} \end{cases} \quad \diamond$$

Continuous Case: We'll assume that if X is continuous, then so is Y . The usual method is to first compute the c.d.f. of Y ,

$$G(y) = \Pr(Y \leq y) = \Pr[h(X) \leq y] = \int_{\{x: h(x) \leq y\}} f(x) dx,$$

and then take the derivative, $g(y) = G'(y)$.

Example: Suppose X has p.d.f. $f(x) = |x|$, $-1 \leq x \leq 1$. Find the p.d.f. of $Y = X^2$. First of all, the c.d.f. of Y is

$$\begin{aligned} G(y) &= \Pr(Y \leq y) = \Pr(X^2 \leq y) = \Pr(-\sqrt{y} \leq X \leq \sqrt{y}) \\ &= \int_{-\sqrt{y}}^{\sqrt{y}} |x| dx = y, \quad 0 < y < 1. \end{aligned}$$

Thus, the p.d.f. of Y is $g(y) = G'(y) = 1$, $0 < y < 1$, indicating that $Y \sim \text{Unif}(0, 1)$. \diamond

Example: Here is a great result sometimes called the Inverse Transform Theorem. Suppose X is a continuous random variable having c.d.f. $F(x)$. Then, amazingly, $F(X) \sim \text{Unif}(0, 1)$.

Proof: Let $Y = F(X)$. Then the c.d.f. of Y is

$$\Pr(Y \leq y) = \Pr(F(X) \leq y) = \Pr(X \leq F^{-1}(y)) = F(F^{-1}(y)) = y,$$

which is the c.d.f. of the $\text{Unif}(0, 1)$. \diamond

Here is a *more-direct* method for dealing with functions of RV's...

Theorem: Suppose that X has p.d.f. $f(x)$, $a \leq x \leq b$. Let $Y = h(X)$ be a *monotone* function (either increasing or decreasing) of X . Then the p.d.f. of Y is

$$g(y) = f(h^{-1}(y)) \left| \frac{d}{dy} h^{-1}(y) \right|, \quad h(a) \leq y \leq h(b) \quad (\text{or } h(b) \leq y \leq h(a)).$$

Remarks: (i) Warning: You can only use this method of $h(x)$ if *monotone*! (The p.d.f. $f(x)$ doesn't have to be monotone.) (ii) Think of the inverse function $h^{-1}(y) = x$, and

the quantity in the $|\cdot|$ as the Jacobian of the transformation.

Example: Suppose that $f(x) = 3x^2$, $0 \leq x \leq 1$. Find the p.d.f. of $Y = h(X) = X^2$. Note that $f(x)$ is only defined on the domain $0 \leq x \leq 1$; and on this range, $h(x)$ is monotone increasing — so it's OK to use the wonderful theorem.

First, we have $x = h^{-1}(y) = \pm\sqrt{y} = \sqrt{y}$ (since we're only concerned with positive x 's). The theorem then implies that

$$\begin{aligned} g(y) &= f(\sqrt{y}) \left| \frac{d}{dy} \sqrt{y} \right|, \quad h(0) \leq y \leq h(1) \\ &= 3y \times \frac{1}{2\sqrt{y}} = \frac{3}{2}\sqrt{y}, \quad 0 \leq y \leq 1. \quad \diamond \end{aligned}$$

Remark: We can also look at functions of ≥ 2 RV's, but this takes more work. See any probability text for more info on this important topic.

2.6 Jointly Distributed RV's

Definition: The *joint cdf* of X and Y is $F(x, y) \equiv P(X \leq x, Y \leq y)$, for all x, y .

Remark: The *marginal cdf* of X is $F_X(x) = F(x, \infty)$. (We use the X subscript to remind us that it's just the cdf of X all by itself.) Similarly, the marginal cdf of Y is $F_Y(y) = F(\infty, y)$.

Definition: If X and Y are discrete, then the *joint pmf* of X and Y is $f(x, y) \equiv P(X = x, Y = y)$.

Remark: The *marginal pmf* of X is

$$f_X(x) = P(X = x) = \sum_y f(x, y).$$

The marginal pmf of Y is

$$f_Y(y) = P(Y = y) = \sum_x f(x, y).$$

Example: Suppose the following table gives the joint pmf of X and Y , along with the accompanying marginals.

	$X = 2$	$X = 3$	$X = 4$	$f_Y(y)$
$Y = 4$	0.3	0.2	0.1	0.6
$Y = 6$	0.1	0.2	0.1	0.4
$f_X(x)$	0.4	0.4	0.2	1

Definition: If X and Y are continuous, then the *joint pdf* of X and Y is $f(x, y) \equiv \frac{\partial^2}{\partial x \partial y} F(x, y)$.

Remark: The *marginal pdf* of X is

$$f_X(x) = \int_{\mathbb{R}} f(x, y) dy.$$

The *marginal pdf* of Y is

$$f_Y(y) = \int_{\mathbb{R}} f(x, y) dx.$$

Example: This example shows that you have to be careful about “funny” limits when computing marginals. Suppose the joint pdf is

$$f(x, y) = \frac{21}{4}x^2y, \quad x^2 \leq y \leq 1.$$

Then the marginal pdf's are:

$$f_X(x) = \int_{\mathbb{R}} f(x, y) dy = \int_{x^2}^1 \frac{21}{4}x^2y dy = \frac{21}{8}x^2(1 - x^4), \quad -1 \leq x \leq 1$$

and

$$f_Y(y) = \int_{\mathbb{R}} f(x, y) dx = \int_{-\sqrt{y}}^{\sqrt{y}} \frac{21}{4}x^2y dx = \frac{7}{2}y^{5/2}, \quad 0 \leq y \leq 1. \quad \diamond$$

2.7 Independent RV's

Definition: X and Y are *independent* RV's if $f(x, y) = f_X(x)f_Y(y)$ for all x, y .

Examples: If $f(x, y) = cxy$ for $0 \leq x \leq 2, 0 \leq y \leq 3$, then X and Y are independent. If $f(x, y) = \frac{21}{4}x^2y$ for $x^2 \leq y \leq 1$, then X and Y are *not* independent. If $f(x, y) = c/(x+y)$ for $1 \leq x \leq 2, 1 \leq y \leq 3$, then X and Y are *not* independent. \diamond

Definition: If $f_X(x) > 0$, then $f(y|x) \equiv f(x, y)/f_X(x)$ is the *conditional pdf* (or *pmf*) of Y given $X = x$.

Example: Suppose $f(x, y) = \frac{21}{4}x^2y$ for $x^2 \leq y \leq 1$. By a previous example, we find that

$$f(y|x) = \frac{f(x, y)}{f_X(x)} = \frac{\frac{21}{4}x^2y}{\frac{21}{8}x^2(1 - x^4)} = \frac{2y}{1 - x^4}, \quad x^2 \leq y \leq 1. \quad \diamond$$

“Definition”: Suppose that $h(X, Y)$ is some function of the RV's X and Y . Then

$$\mathbb{E}[h(X, Y)] = \begin{cases} \sum_x \sum_y h(x, y) f(x, y) & \text{if } (X, Y) \text{ is discrete} \\ \int_{\mathbb{R}} \int_{\mathbb{R}} h(x, y) f(x, y) dx dy & \text{if } (X, Y) \text{ is continuous} \end{cases}$$

Example/Theorem: Whether or not X and Y are independent, we have $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$. In fact, if X_1, X_2, \dots are RV's, then $\mathbb{E}[\sum_i X_i] = \sum_i \mathbb{E}[X_i]$.

Theorem: If X and Y are *independent*, then $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ and $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$.

Proof: Easy algebra. \diamond

Theorem: Suppose that X_1, \dots, X_n are *independent* RV's. If $Y = \sum_{i=1}^n X_i$, then

$$M_Y(t) = \mathbb{E}[e^{tY}] = \mathbb{E}[e^{t \sum X_i}] = \prod_{i=1}^n \mathbb{E}[e^{tX_i}] = \prod_{i=1}^n M_{X_i}(t).$$

Definition: X_1, \dots, X_n form a *random sample* from $f(x)$ is

1. X_1, \dots, X_n are independent, and
2. Each X_i has the same pdf (or pmf) $f(x)$.

Notation: $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x)$. (The term “iid” reads *independent and identically distributed*)

Corollary: X_1, \dots, X_n iid implies that $M_Y(t) = [M_{X_i}(t)]^n$.

Example: Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bern}(p)$. Then $M_{\sum_i X_i}(t) = (pe^t + q)^n$. It turns out that this is the mgf for the $\text{Bin}(n, p)$ distribution. Thus, by a previous theorem, we have $\sum_{i=1}^n X_i \sim \text{Bin}(n, p)$. \diamond

Example: If $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x)$ and $\bar{X} \equiv \sum_{i=1}^n X_i/n$, then $\mathbb{E}[\bar{X}] = \mathbb{E}[X_i]$ and $\text{Var}(\bar{X}) = \text{Var}(X_i)/n$. Thus, the variance *decreases*. \diamond

2.8 Covariance and Correlation

Definition: The *covariance* between X and Y is $\text{Cov}(X, Y) \equiv \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$. Note that $\text{Var}(X) = \text{Cov}(X, X)$.

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

Proof: Easy. \diamond

Theorem: If X and Y are independent RV's, then $\text{Cov}(X, Y) = 0$.

Proof: Since X and Y are independent, we have $E[XY] = E[X]E[Y]$. \diamond

Remark: $\text{Cov}(X, Y) = 0$ does *not* imply that X and Y are independent!

Example: Suppose $X \sim \text{Unif}(-1, 1)$ and $Y = X^2$. Then X and Y are clearly dependent. However,

$$\text{Cov}(X, Y) = E[X^3] - E[X]E[X^2] = E[X^3] = \int_{-1}^1 \frac{x^3}{2} dx = 0. \quad \diamond$$

Theorem: If a and b are constants, the $\text{Cov}(aX, bY) = ab\text{Cov}(X, Y)$.

Definition: The *correlation* between X and Y is

$$\rho \equiv \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}.$$

Theorem: $-1 \leq \rho \leq 1$.

Proof: Follows from the Cauchy-Schwarz inequality. \diamond

Remark: If $\rho \approx 1$, we say that X and Y have “high positive” correlation. If $\rho \approx 0$, X and Y have “low” correlation. If $\rho \approx -1$, there is “high negative” correlation.

Example: Suppose that X is the average yards per carry gained by a University of Georgia fullback and Y is his IQ. Further suppose that the joint pmf $f(x, y)$ is given in the following table.

	$X = 2$	$X = 3$	$X = 4$	$f_Y(y)$
$Y = 40$	0.00	0.20	0.10	0.3
$Y = 50$	0.15	0.10	0.05	0.3
$Y = 60$	0.30	0.00	0.10	0.4
$f_X(x)$	0.45	0.30	0.25	1

Then we have $E[X] = 2.8$, $\text{Var}(X) = 0.66$, $E[Y] = 51$, $\text{Var}(Y) = 69$, $E[XY] = \sum_x \sum_y xyf(x, y) = 140$, and

$$\rho = \frac{E[XY] - E[X]E[Y]}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = -0.415. \quad \diamond$$

Theorem: $\text{Var}(\sum_{i=1}^n a_i X_i) = \sum_{i=1}^n a_i^2 \text{Var}(X_i) + 2 \sum \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$.

Corollary: $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$.

Corollary: $\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y)$.

2.9 Some Fun Distributions

First, some discrete distributions...

2.9.1 Bernoulli

$X \sim \text{Bernoulli}(p)$.

$$f(x) = \begin{cases} p & \text{if } x = 0 \\ 1 - p & \text{if } x = 1 \end{cases}$$

$E[X] = p$, $\text{Var}(X) = p(1 - p)$, $M_X(t) = pe^t + q$.

If X_1, X_2, \dots, X_n are i.i.d. $\text{Bern}(p)$, we say that they form a series of *Bernoulli(p) trials*.

2.9.2 Binomial

$X \sim \text{Binomial}(n, p)$.

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

$E[X] = np$, $\text{Var}(X) = np(1 - p)$, $M_X(t) = (pe^t + q)^n$. If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bern}(p)$, then $\sum_{i=1}^n X_i \sim \text{Bin}(n, p)$.

2.9.3 Geometric

$X \sim \text{Geom}(p)$ is the number of $\text{Bern}(p)$ trials until a success occurs. For example, “FFFS” implies that $X = 4$.

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

$E[X] = 1/p$, $\text{Var}(X) = q/p^2$.

2.9.4 Negative Binomial

$X \sim \text{NegBin}(r, p)$ is the sum of r i.i.d. $\text{Geom}(p)$ RV's, i.e., the time until the r th success occurs. For example, “FFFSSFS” implies that $\text{NegBin}(3, p) = 7$.

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

$$\mathbb{E}[X] = r/p, \quad \text{Var}(X) = qr/p^2.$$

2.9.5 Poisson

A *counting process* $N(t)$ tallies the number of “arrivals” observed in $[0, t]$. A *Poisson process* is a counting process satisfying the following.

- i. Arrivals occur one-at-a-time.
- ii. Independent increments, i.e., the numbers of arrivals in disjoint time intervals are independent.
- iii. Stationary increments, i.e., the distribution of the number of arrivals only depends on the length of the time interval under observation.

$X \sim \text{Pois}(\lambda)$ is the number of arrivals that a Poisson processes experiences in one time unit, i.e., $N(1)$.

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

$$\mathbb{E}[X] = \lambda = \text{Var}(X).$$

Now, some continuous distributions...

2.9.6 Uniform

$X \sim \text{Unif}(a, b)$.

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

$$\mathbb{E}[X] = \frac{a+b}{2}, \quad \text{Var}(X) = \frac{(b-a)^2}{12}, \quad M_X(t) = \frac{e^{tb} - e^{ta}}{t}.$$

2.9.7 Exponential

$X \sim \text{Exp}(\lambda)$.

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

$$\mathbb{E}[X] = 1/\lambda, \quad \text{Var}(X) = 1/\lambda^2, \quad M_X(t) = \frac{\lambda}{\lambda - t}, \quad t < \lambda.$$

Theorem. The exponential distribution has the *memoryless property*, i.e., for $s, t > 0$,

$$\Pr(X > s + t | X > s) = \Pr(X > t).$$

By the way, the $\text{Exp}(\lambda)$ is the only continuous distribution with this property.

Example: Suppose that a light bulb has a lifetime that is exponential with mean 1000 hours. Suppose it has already survived 500 hours. Then the probability that it makes it to 2000 is

$$\Pr(X > 2000 | X > 500) = \Pr(X > 1500) = e^{-\lambda t} = e^{-1500/1000}. \quad \diamond$$

2.9.8 Gamma

$X \sim \text{Gamma}(\alpha, \lambda)$.

$$f(x) = \frac{\lambda^\alpha x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)}, \quad x \geq 0,$$

where the gamma function is

$$\Gamma(z) \equiv \int_0^\infty t^{z-1} e^{-t} dt.$$

$\mathbb{E}[X] = \alpha/\lambda$, $\text{Var}(X) = \alpha/\lambda^2$, $M_X(t) = (\frac{\lambda}{\lambda - t})^\alpha$. If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Exp}(\lambda)$, then $Y \equiv \sum_{i=1}^n X_i \sim \text{Gamma}(n, \lambda)$. The $\text{Gamma}(n, \lambda)$ is also called the $\text{Erlang}_n(\lambda)$. It has c.d.f.

$$F_Y(y) = 1 - e^{-\lambda y} \sum_{j=0}^{n-1} \frac{(\lambda y)^j}{j!}, \quad y \geq 0.$$

2.9.9 Normal

$X \sim \text{Nor}(\mu, \sigma^2)$.

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right], \quad x \in \mathbb{R}.$$

$$\mathbb{E}[X] = \mu, \quad \text{Var}(X) = \sigma^2, \quad M_X(t) = \exp\{\mu t + \frac{1}{2}\sigma^2 t^2\}.$$

Theorem (Additive Property of Normals): Suppose that X_1, X_2, \dots, X_n are *independent* with $X_i \sim \text{Nor}(\mu_i, \sigma_i^2)$, $i = 1, 2, \dots, n$. Then

$$Y = \sum_{i=1}^n a_i X_i + b \sim \text{Nor}\left(\sum_{i=1}^n a_i \mu_i + b, \sum_{i=1}^n a_i^2 \sigma_i^2\right).$$

Proof: Use m.g.f.'s.

$$\begin{aligned} M_Y(t) &= \mathbb{E}(e^{tY}) = \mathbb{E}\left(\exp\left\{t\left(\sum_{i=1}^n a_i X_i + b\right)\right\}\right) \\ &= e^{tb} \mathbb{E}\left(\exp\left\{\sum_{i=1}^n (a_i t) X_i\right\}\right) \\ &= e^{tb} \prod_{i=1}^n \mathbb{E}\left(e^{(a_i t) X_i}\right) \quad (\text{by independence}) \\ &= e^{tb} \prod_{i=1}^n M_{X_i}(a_i t) \\ &= e^{tb} \prod_{i=1}^n \exp\left\{\mu_i(a_i t) + \frac{1}{2}\sigma_i^2(a_i t)^2\right\} \\ &= \exp\left\{\left(\sum_{i=1}^n \mu_i a_i + b\right)t + \frac{1}{2}\left(\sum_{i=1}^n a_i^2 \sigma_i^2\right)t^2\right\}. \quad \diamond \end{aligned}$$

Example: Suppose $X \sim \text{Nor}(3, 4)$, $Y \sim \text{Nor}(4, 6)$, and X and Y are independent. Then

$$2X - 3Y + 1 \sim \text{Nor}(2\mathbb{E}[X] - 3\mathbb{E}[Y] + 1, 4\text{Var}(X) + 9\text{Var}(Y)) \sim \text{Nor}(-5, 70). \quad \diamond$$

Corollary: If $X \sim \text{Nor}(\mu, \sigma^2)$, then $aX + b \sim \text{Nor}(a\mu + b, a^2\sigma^2)$.

Corollary: If $X \sim \text{Nor}(\mu, \sigma^2)$, then $Z \equiv \frac{X-\mu}{\sigma} \sim \text{Nor}(0, 1)$, the standard normal distribution.

Notation: The standard normal's p.d.f. is $\phi(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$, and the c.d.f. is $\Phi(x)$, which is usually tabled. For example, $\Phi(1.96) \doteq 0.975$.

2.10 A First Look at Some Limit Theorems

Corollary (of theorem on linear combinations of normals from previous subsection): If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Nor}(\mu, \sigma^2)$, then the sample mean

$$\bar{X} \equiv \frac{1}{n} \sum_{i=1}^n X_i \sim \text{Nor}(\mu, \sigma^2/n).$$

This is a special case of the *Law of Large Numbers*, which says that \bar{X} approximates μ well as n becomes large.

Markov's Inequality: If X is a non-negative RV, then for all $\epsilon > 0$, we have

$$\Pr(X \geq \epsilon) \leq \mathbb{E}[X]/\epsilon.$$

Proof: Since X is non-negative,

$$\mathbb{E}[X] = \int_0^\infty x f(x) dx \geq \int_\epsilon^\infty x f(x) dx \geq \epsilon \int_\epsilon^\infty f(x) dx = \epsilon \Pr(X \geq \epsilon). \quad \diamond$$

Chebychev's Inequality: For any RV X and for all $\epsilon > 0$, we have

$$\Pr(|X - \mathbb{E}[X]| \geq \epsilon) \leq \text{Var}(X)/\epsilon^2.$$

Proof: Uses Markov's Inequality; see any probability test. \diamond

Bonus Generalization: $\Pr(|X| \geq \epsilon) \leq \mathbb{E}[|X|^r]/\epsilon^r$.

Remark: These inequalities are usually pretty crude!

Example: Suppose that $X \sim \text{Unif}(0, 1)$. Then the probability that X deviates from its mean by at least $1/4$ is exactly

$$\begin{aligned} \Pr\left(\left|X - \frac{1}{2}\right| \geq \frac{1}{4}\right) &= 1 - \Pr\left(\left|X - \frac{1}{2}\right| < \frac{1}{4}\right) = 1 - \Pr\left(-\frac{1}{4} < X - \frac{1}{2} < \frac{1}{4}\right) \\ &= 1 - \Pr\left(\frac{1}{4} < X < \frac{3}{4}\right) = \frac{1}{2}. \end{aligned}$$

Meanwhile, by Chebychev (with $\mathbb{E}[X] = 1/2$, $\text{Var}(X) = 1/12$, and $\epsilon = 1/4$), we have

$$\Pr\left(\left|X - \frac{1}{2}\right| \geq \frac{1}{4}\right) \leq \frac{\text{Var}(X)}{\epsilon^2} = \frac{16}{12} = \frac{4}{3},$$

which is a very crude upper bound indeed! \diamond

Definition: The sequence of random variables Y_1, Y_2, \dots with respective c.d.f.'s $F_{Y_1}(y), F_{Y_2}(y), \dots$ *converges in distribution* to the random variable Y having c.d.f. $F_Y(y)$ if $\lim_{n \rightarrow \infty} F_{Y_n}(y) = F_Y(y)$ for all y belonging to the continuity set of Y (i.e., the set of all points y at which $F_Y(y)$ is continuous). Notation: $Y_n \xrightarrow{\mathcal{D}} Y$. (Also sometimes called *convergence in law* or *weak convergence*.)

Idea: If $Y_n \xrightarrow{\mathcal{D}} Y$, then you would expect to be able to approximate the distribution of Y_n by the limiting distribution of Y , at least for large enough n .

Central Limit Theorem: If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} f(x)$ with mean μ and variance σ^2 , then

$$Z_n \equiv \frac{\sum_{i=1}^n X_i - n\mu}{\sqrt{n}\sigma} = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{\mathcal{D}} \text{Nor}(0, 1),$$

where \bar{X}_n is the sample mean. Thus, the c.d.f. of Z_n approaches that of the standard normal as n increases. The CLT usually works pretty well if the pdf/pmf is fairly symmetric and $n \geq 15$.

Example: Suppose that $X_1, X_2, \dots, X_{100} \stackrel{\text{iid}}{\sim} \text{Exp}(1)$. Then

$$\begin{aligned} \Pr\left(90 \leq \sum_{i=1}^{100} X_i \leq 110\right) &= \Pr\left(\frac{90 - 100}{\sqrt{100}} \leq Z_{100} \leq \frac{110 - 100}{\sqrt{100}}\right) \\ &= \Pr(-1 \leq Z_{100} \leq 1) \approx \Pr(-1 \leq \text{Nor}(0, 1) \leq 1) = 2\Phi(1) - 1 \approx 0.683. \quad \diamond \end{aligned}$$

Definition: The sequence of random variables Y_1, Y_2, \dots is said to *converge in probability* to Y (often a constant) if for all $\epsilon > 0$, $\Pr(|Y_n - Y| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$. Notation: $Y_n \xrightarrow{\mathcal{P}} Y$.

Theorem: $Y_n \xrightarrow{\mathcal{P}} Y$ implies $Y_n \xrightarrow{\mathcal{D}} Y$. In other words, convergence in probability is a bit stronger than convergence in distribution.

Weak Law of Large Numbers: If X_1, X_2, \dots are i.i.d. with mean μ , then $\bar{X}_n \xrightarrow{\mathcal{P}} \mu$. Why is this called the *weak* LLN? Simply because there's a stronger one coming up later.

Continuous Mapping Theorem: If $Y_n \xrightarrow{\mathcal{P}} Y$ and $g(\cdot)$ is a nice, continuous function, then $g(Y_n) \xrightarrow{\mathcal{P}} g(Y)$. The CMT is often useful for characterizing the convergence of nasty functions of the Y_i 's.

Definition: The sequence of random variables Y_1, Y_2, \dots is said to *converge in r th mean* to Y (often a constant) if $\mathbb{E}[|Y_n - Y|^r] \rightarrow 0$ as $n \rightarrow \infty$. Notation: $Y_n \xrightarrow{r} Y$.

Theorem: $Y_n \xrightarrow{r} Y$ implies $Y_n \xrightarrow{\mathcal{P}} Y$. In other words, convergence in r th mean is a bit stronger than convergence in probability.

Proof: Follows immediately from bonus version of Chebychev. \diamond

Definition: The sequence of random variables Y_1, Y_2, \dots *converges almost surely* (or *with probability one*) to Y if $\Pr(Y_n \text{ converges to } Y) = 1$ as $n \rightarrow \infty$. Notation: $Y_n \xrightarrow{a.s.} Y$.

Theorem: $Y_n \xrightarrow{a.s.} Y$ implies $Y_n \xrightarrow{\mathcal{P}} Y$. In other words, convergence almost surely is a bit stronger than convergence in probability.

Strong Law of Large Numbers: If X_1, X_2, \dots are i.i.d. with mean μ , then $\bar{X}_n \xrightarrow{a.s.} \mu$. It turns out that the SLLN implies the WLLN.

How do almost sure and r th mean convergence relate to each other?

Dominated Convergence Theorem: If $Y_n \xrightarrow{a.s.} Y$ and there exists a random variable W such that $\Pr(|Y_n| \leq W) = 1$ for every n , then $E[Y_n] \rightarrow E[Y]$.