

Expectation

- Recall. Expectation for univariate random variable. LNp. 4-11 ~ 16
LNp. 5-10 ~ 13
- Theorem. For random variables $\mathbf{X}=(X_1, \dots, X_n)$ with joint pmf $p_{\mathbf{X}}$ /pdf $f_{\mathbf{X}}$, the *expectation* of a univariate random variable Y , where

$$Y=g(X_1, \dots, X_n), g:\mathbb{R}^n \rightarrow \mathbb{R}^1$$

is $E(Y) \equiv \sum_{y \in \mathcal{Y}} y p_Y(y)$ (1)

no need to calculate $P_Y(y)$ $\equiv \sum_{\mathbf{x}=(x_1, \dots, x_n) \in \mathcal{X}} g(x_1, \dots, x_n) p_{\mathbf{X}}(x_1, \dots, x_n)$ (2)
 $\equiv E[g(X_1, \dots, X_n)]$ i.e. $\sum |g(\mathbf{x})| P_{\mathbf{X}}(\mathbf{x}) < \infty$

if X_1, \dots, X_n are discrete and the sum converges absolutely, or

$E(Y) \equiv \int_{-\infty}^{\infty} y f_Y(y) dy$ (3)

no need to calculate $f_Y(y)$ $\equiv \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} g(x_1, \dots, x_n) f_{\mathbf{X}}(x_1, \dots, x_n) dx_1 \dots dx_n$ (4)
 $\equiv E[g(X_1, \dots, X_n)]$ i.e. $\int |g(\mathbf{x})| f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} < \infty$

if Y and X_1, \dots, X_n are continuous and the integrals converges absolutely

Proof. Like the univariate case. LNp. 4-13
LNp. 5-11

eg. $Y = \begin{cases} 1, & X \in A \subset \mathbb{R}^n \\ 0, & \text{o.w.} \end{cases}$

➤ Q: What if Y is discrete and X_1, \dots, X_n are continuous?

➤ Notation.

- Shorthand notation. Combine (1) and (3) by writing

$$E(Y) = \int_{-\infty}^{\infty} y dF_Y(y) = \begin{cases} \sum_{y \in \mathcal{Y}} y p_Y(y), & \text{for discrete case,} \\ \int_{-\infty}^{\infty} y f_Y(y) dy, & \text{for continuous case,} \end{cases}$$

and combine (2) and (4) by writing Note: $\frac{dF_Y(y)}{dy} = f_Y(y) \Rightarrow dF_Y(y) = f_Y(y) dy$

$$E[g(\mathbf{X})] = \int_{\mathbb{R}^n} g(\mathbf{x}) dF_{\mathbf{X}}(\mathbf{x}) \quad \text{joint cdf.}$$

$$= \begin{cases} \sum_{\mathbf{x} \in \mathcal{X}} g(\mathbf{x}) p_{\mathbf{X}}(\mathbf{x}), & \text{for discrete case.} \\ \int_{\mathbb{R}^n} g(\mathbf{x}) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}, & \text{for continuous case.} \end{cases}$$

- Riemann-Stieltjes Integral. For example, for non-negative g ,

$$\int_a^b g(x) dF(x) = \lim \sum_{i=1}^n g(x_i) [F(x_i) - F(x_{i-1})].$$

where the limit is taken over all $a=x_0 < x_1 < \dots < x_n=b$ as $n \rightarrow \infty$ and $\max_{i=1, \dots, n} (x_i - x_{i-1}) \rightarrow 0$.

[Recall. The integral of g over $(a, b]$ is defined as

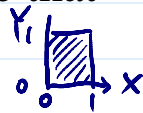
$$\int_a^b g(x) dx = \lim \sum_{i=1}^n g(x_i) (x_i - x_{i-1}).]$$

➤ Note.

- $g(X_1, \dots, X_n) = X_i \Rightarrow E[g(X_1, \dots, X_n)] = E(X_i) \equiv \mu_{X_i}$.
- $g(X_1, \dots, X_n) = (X_i - \mu_{X_i})^2 \Rightarrow E[g(X_1, \dots, X_n)] = \text{Var}(X_i) \equiv \sigma_{X_i}^2$.

➤ Example (Distance between two points). Suppose that

$g(x,y)$ X, Y are i.i.d. \sim Uniform(0, 1).



Let $D = |X - Y|$. Find $E(D)$.

▪ The joint pdf of (X, Y) is

Note: not necessary to derive the pdf of D

$$f(x, y) = \begin{cases} 1, & 0 \leq x \leq 1, 0 \leq y \leq 1, \\ 0, & \text{otherwise.} \end{cases}$$

$$\begin{aligned} E(D) &= \int_0^1 \int_0^1 |x - y| \, dy \, dx = \int_0^1 \left[\int_0^x (x - y) \, dy + \int_x^1 (y - x) \, dy \right] dx \\ &= \int_0^1 \left[-\frac{1}{2}(y - x)^2 \Big|_{y=0}^x + \frac{1}{2}(y - x)^2 \Big|_x^1 \right] dx \\ &= \int_0^1 \frac{1}{2} [x^2 + (1 - x)^2] \, dx = \frac{1}{6} [x^3 - (1 - x)^3] \Big|_{x=0}^1 = \frac{1}{3}. \end{aligned}$$

• Theorem (Mean of Sum). For r.v.'s X_1, \dots, X_n and constants $-\infty < a_0, a_1, \dots, a_n < \infty$,

$$E(a_0 + a_1 X_1 + \dots + a_n X_n) = a_0 + a_1 E(X_1) + \dots + a_n E(X_n).$$

Proof. $E(a_0 + a_1 X_1 + \dots + a_n X_n)$

$$\begin{aligned} &= \int_{\mathbb{R}^n} (a_0 + a_1 x_1 + \dots + a_n x_n) \, dF_{\mathbf{X}}(\mathbf{x}) \\ &= \int_{\mathbb{R}^n} a_0 \, dF_{\mathbf{X}}(\mathbf{x}) + a_1 \int_{\mathbb{R}^n} x_1 \, dF_{\mathbf{X}}(\mathbf{x}) + \dots + a_n \int_{\mathbb{R}^n} x_n \, dF_{\mathbf{X}}(\mathbf{x}) \\ &= a_0 + a_1 E(X_1) + \dots + a_n E(X_n). \end{aligned}$$

➤ Corollary. Suppose that $\mu = E(X_1) = \dots = E(X_n)$. Let

$$\bar{X}_n = \frac{X_1 + \dots + X_n}{n}, \quad \leftarrow a_0 = 0, a_1 = a_2 = \dots = a_n = \frac{1}{n}$$

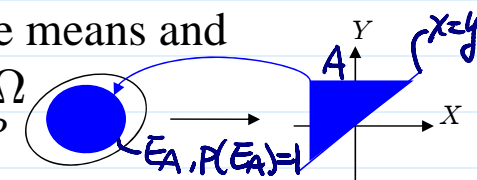
then, $E(\bar{X}_n) = \mu$.

➤ Corollary. If X and Y are r.v.'s with finite means and

eg. X : 出生時身高
 Y : 10歲時身高

$P(X \leq Y) = 1$, $X \leq Y$ with prob. one, almost surely

then $E(X) \leq E(Y)$.



Proof. First, if Z is a random variable with finite mean and

$$P(Z \geq 0) = 1, \Rightarrow Z \geq 0 \text{ with prob. one, almost surely} \Rightarrow P(Z < 0) = 0.$$

then $E(Z) = \int_0^{\infty} z \, dF_Z(z) \geq 0$.
For the general case, let $Z = Y - X$, then $Z \geq 0$ with probability one, and therefore, $0 \leq E(Z) = E(Y - X) = E(Y) - E(X)$.

➤ Corollary. If $P(a < X \leq b) = 1$ for some constants a, b , then

$$P(X - a \geq 0) = 1 \Rightarrow a \leq E(X) \leq b.$$

• Theorem. If two random vectors $\mathbf{X} (\in \mathbb{R}^m)$ and $\mathbf{Y} (\in \mathbb{R}^n)$ are independent (i.e., $F_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = F_{\mathbf{X}}(\mathbf{x}) \times F_{\mathbf{Y}}(\mathbf{y})$, or

$$f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = f_{\mathbf{X}}(\mathbf{x}) \times f_{\mathbf{Y}}(\mathbf{y}), \text{ or } p_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = p_{\mathbf{X}}(\mathbf{x}) \times p_{\mathbf{Y}}(\mathbf{y}),$$

then for $g: \mathbb{R}^m \rightarrow \mathbb{R}$ and $h: \mathbb{R}^n \rightarrow \mathbb{R}$,

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

Note: $(g(\mathbf{X}), h(\mathbf{Y}))$ are independent (LN p. 6-20 Example)

can be relaxed to uncorrelated

Proof. We only prove it for the continuous case:

h: concave

$$\begin{aligned}
 E[g(\mathbf{X})h(\mathbf{Y})] &= \int_{\mathbb{R}^m} \int_{\mathbb{R}^n} g(\mathbf{x})h(\mathbf{y})f_{\mathbf{X},\mathbf{Y}}(\mathbf{x},\mathbf{y}) dyd\mathbf{x} \\
 &\stackrel{\text{indep}}{=} \int_{\mathbb{R}^m} \int_{\mathbb{R}^n} g(\mathbf{x})h(\mathbf{y})f_{\mathbf{X}}(\mathbf{x})f_{\mathbf{Y}}(\mathbf{y}) dyd\mathbf{x} \\
 &= \int_{\mathbb{R}^m} g(\mathbf{x})f_{\mathbf{X}}(\mathbf{x}) \left[\int_{\mathbb{R}^n} h(\mathbf{y})f_{\mathbf{Y}}(\mathbf{y}) dy \right] d\mathbf{x} \quad \text{This is a constant for } \mathbf{x}. \\
 &= \left[\int_{\mathbb{R}^m} g(\mathbf{x})f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \right] \left[\int_{\mathbb{R}^n} h(\mathbf{y})f_{\mathbf{Y}}(\mathbf{y}) dy \right] \\
 &= E[g(\mathbf{X})]E[h(\mathbf{Y})].
 \end{aligned}$$

➤ Corollary. For 2 independent r.v.'s X and Y , $E(XY)=E(X)E(Y)$.

proof: $g(x)=x, h(y)=y$

➤ **Q:** For independent r.v.'s X and Y , $E(X/Y)=E(X)/E(Y)$?

$$E\left(\frac{X}{Y}\right) = E\left[X \cdot \frac{1}{Y}\right] \stackrel{\text{indep}}{=} E(X) \cdot E\left(\frac{1}{Y}\right) \neq E(X) \cdot \frac{1}{E(Y)}$$

∵ X & $\frac{1}{Y}$ indep by an example

➤ Note. $E[h(Y)] \neq h(E(Y))$ in general, e.g., $E(1/Y) \neq 1/E(Y)$. *LNp.6-20*

• Covariance and Correlation between 2 random variables

Recall mean variance

Definition. Suppose that X and Y are two random variables with finite means μ_X, μ_Y and variances σ_X^2, σ_Y^2 , respectively.

1. Let $g(x, y) = (x - \mu_X)(y - \mu_Y)$, then

$$Cov(X, Y) \equiv E[g(X, Y)] = E[(X - \mu_X)(Y - \mu_Y)]$$

is called the *covariance* between X and Y , denoted by σ_{XY} .

can be calculate from the marginal distribution of X & Y

2. The *correlation* (coefficient) between X and Y is defined as ^{p. 7-6}

$$Cor(X, Y) = \sigma_{XY} / (\sigma_X \sigma_Y)$$

standard deviation (LNp.4-14, 5-12)

and denoted by ρ_{XY} .

3. X and Y are called *uncorrelated* if $\rho_{XY}=0$. *i.e. $Cov(X, Y)=0$*

■ A special case of covariance: $Cov(X, X) = Var(X)$.

$$= E[(X - \mu_X)(X - \mu_X)] = E[(X - \mu_X)^2]$$

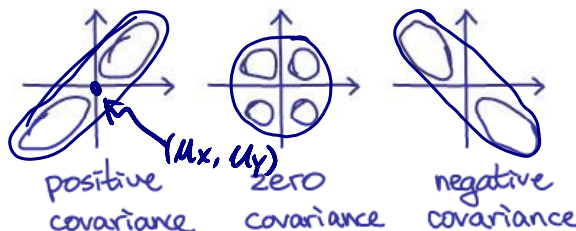
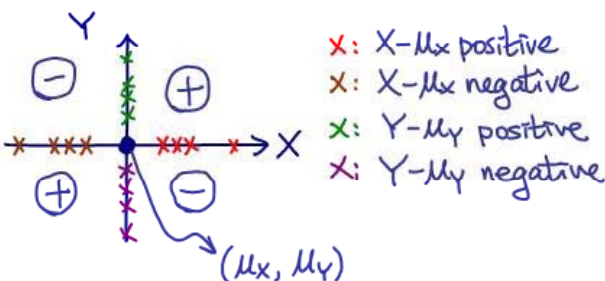
➤ Intuitive explanation of covariance and correlation

■ Covariance is a measure of the joint variability of X and Y , or their degree of **association**. *whether $Y \uparrow$ (or $Y \downarrow$) when $X \uparrow$. e.g. X : height. Y : weight.*

■ Covariance is the average value of the product of the deviation of X from its mean and the deviation of Y from its mean. *drawback: covariance depends on the units/scales of X & Y , e.g. height: m \rightarrow cm, 10^2 larger*

definition

■ Positive Covariance and Negative Covariance



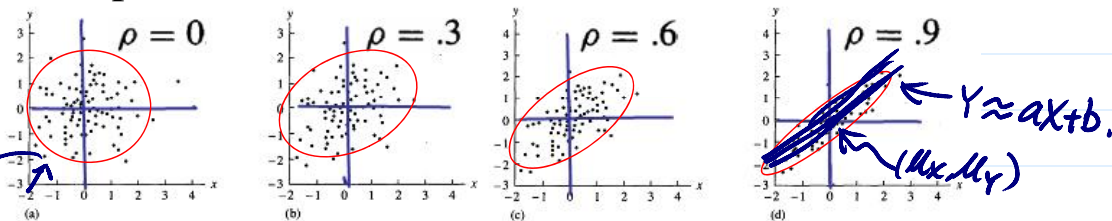
Correlation Coefficient is unit free.

Correlation coefficient measures the strength of the linear relationship between X and Y .

why? check definition

1/3

joint pdf.



Theorem. $Cov(X, Y) = E(XY) - \mu_X \mu_Y$. \leftarrow c.f. $cov(x, x) = Var(X) = E(X^2) - \mu_X^2$

Proof. $Cov(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$
 $= E(XY - \mu_X Y - \mu_Y X + \mu_X \mu_Y)$ (constants)
 $= E(XY) - \mu_X E(Y) - \mu_Y E(X) + \mu_X \mu_Y$
 $= E(XY) - \mu_X \mu_Y - \mu_Y \mu_X + \mu_X \mu_Y$

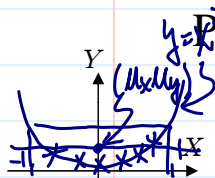
Corollary. If X and Y are independent, then $Cov(X, Y) = 0$, i.e. X and Y are uncorrelated.

corollary Lnp. 7-5

Proof. When X, Y are independent, $E(XY) = E(X)E(Y) = \mu_X \mu_Y$.

uncorrelated is a weaker property than indep.

However, the converse statement is not necessarily true (e.g., let $X \sim \text{Uniform}(-1, 1)$ and $Y = X^2$, then $Cov(X, Y) = 0$, but X and Y are not independent). $E(XY) = \int_{-1}^1 x^3 dx = 0$



Corollary. $\rho_{XY} = E \left[\left(\frac{X - \mu_X}{\sigma_X} \right) \left(\frac{Y - \mu_Y}{\sigma_Y} \right) \right]$

standardization (標準化, Lnp. 5-27)

Proof. By definition.

Example. If $(X_1, \dots, X_m) \sim \text{Multinomial}(n, m, p_1, \dots, p_m)$, then

$Cov(X_i, X_j) = -np_i p_j$, for $1 \leq i \neq j \leq m$.

Because $(X_1, X_2, X_3 + \dots + X_m) \sim \text{Multinomial}(n, 3, p_1, p_2, p_3 + \dots + p_m)$, and

$X_3 + \dots + X_m = n - X_1 - X_2$,

$p_3 + \dots + p_m = 1 - p_1 - p_2$,

we have

$0 \leq X_1 + X_2 \leq n$

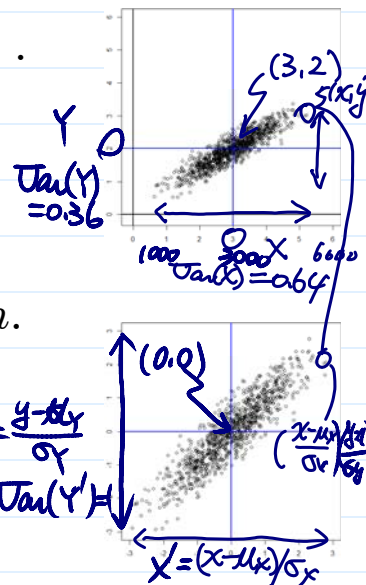
$E(X_1 X_2) = \sum x_1 x_2 \binom{n}{x_1, x_2, n-x_1-x_2} p_1^{x_1} p_2^{x_2} (1-p_1-p_2)^{n-x_1-x_2}$

$0 \leq x_1 + x_2 \leq n$
 $x_1 \neq 0$
 $x_2 \neq 0$

$= \sum x_1 x_2 \frac{n!}{x_1! x_2! (n-x_1-x_2)!} p_1^{x_1} p_2^{x_2} (1-p_1-p_2)^{n-x_1-x_2}$

$= n(n-1)p_1 p_2 \left[\sum \frac{(n-2)!}{(x_1-1)!(x_2-1)!(n-x_1-x_2)!} p_1^{x_1-1} p_2^{x_2-1} (1-p_1-p_2)^{n-x_1-x_2} \right]$

$= n(n-1)p_1 p_2$. joint pmf of multinomial $(n-2, 3, p_1, p_2, 1-p_1-p_2)$



Q: Why cor unit free?

- WLOG, we can get $E(X_i X_j) = n(n-1)p_i p_j$, for $i \neq j$.

Therefore,
$$\begin{aligned} Cov(X_i, X_j) &= E(X_i X_j) - E(X_i)E(X_j) \\ &= n(n-1)p_i p_j - (np_i)(np_j) = -np_i p_j. \end{aligned}$$

- And, for $i \neq j$,

$$Cor(X_i, X_j) = \frac{-np_i p_j}{\sqrt{np_i(1-p_i)}\sqrt{np_j(1-p_j)}} = \ominus \sqrt{\frac{p_i p_j}{(1-p_i)(1-p_j)}}.$$

why negative?

• Expectations for Sums of Random Variables

➤ Notation. In the following, let X_1, \dots, X_n and Y_1, \dots, Y_m be r.v.'s and $-\infty < a_0, a_1, \dots, a_n, b_0, b_1, \dots, b_m < \infty$ are constants.

➤ Recall. $E(a_0 + a_1 X_1 + \dots + a_n X_n) = a_0 + a_1 E(X_1) + \dots + a_n E(X_n)$.

➤ Theorem (covariance of two sums).

$$\begin{aligned} Cov(a_0 + a_1 X_1 + \dots + a_n X_n, b_0 + b_1 Y_1 + \dots + b_m Y_m) \\ = \sum_{i=1}^n \sum_{j=1}^m a_i b_j Cov(X_i, Y_j). = [a_1 \dots a_n] \begin{bmatrix} Cov(X_1, Y_1) \\ \vdots \\ Cov(X_1, Y_m) \end{bmatrix} \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} \end{aligned}$$

a0, b0 are gone

Proof. Let $S = a_0 + a_1 X_1 + \dots + a_n X_n$ and $T = b_0 + b_1 Y_1 + \dots + b_m Y_m$,

then
$$\begin{aligned} S - E(S) &= \sum_{i=1}^n a_i (X_i - \mu_{X_i}), \\ T - E(T) &= \sum_{j=1}^m b_j (Y_j - \mu_{Y_j}), \\ [S - E(S)][T - E(T)] &= \sum_{i=1}^n \sum_{j=1}^m a_i b_j (X_i - \mu_{X_i})(Y_j - \mu_{Y_j}). \end{aligned}$$

Y1...Ym Cov(Xi, Yj)

Therefore, $Cov(S, T) = E \{ [S - E(S)][T - E(T)] \}$

$$\begin{aligned} [a_1, \dots, a_n] \begin{bmatrix} Cov(X_i, X_j) \\ \vdots \\ Cov(X_i, X_n) \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} &= \sum_{i=1}^n \sum_{j=1}^m a_i b_j E[(X_i - \mu_{X_i})(Y_j - \mu_{Y_j})] \\ &= \sum_{i=1}^n \sum_{j=1}^m a_i b_j Cov(X_i, Y_j). \end{aligned}$$

symmetric nxn

➤ Theorem (variance of sum). $Cov(a_0 + a_1 X_1 + \dots + a_n X_n, a_0 + a_1 X_1 + \dots + a_n X_n)$

$$\begin{aligned} Var(a_0 + a_1 X_1 + \dots + a_n X_n) &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j Cov(X_i, X_j) \\ &= \sum_{i=1}^n a_i^2 Var(X_i) + 2 \sum_{1 \leq i < j \leq n} a_i a_j Cov(X_i, X_j). \end{aligned}$$

Proof. $Cov(X_i, X_i) = Var(X_i)$. Why? ① $X_i \approx X_i, Var(X_i + X_i) = 4Var(X_i)$
② $X_i \approx -X_i, Var(X_i - X_i) = 0 = 2Var(X_i) - 2Cov(X_i, X_i)$

■ Corollary. If X_1, \dots, X_n are uncorrelated, then $Cov(X_i, X_j) = 0, \forall i, j$.

mean of sum of n.v. Wp. 7-3

c.f. $Var(\underbrace{a_0}_{0} + \underbrace{a_1 X_1}_{Y_1} + \dots + \underbrace{a_n X_n}_{Y_n}) = \sum_{i=1}^n a_i^2 Var(X_i)$

square

■ Corollary. If X_1, \dots, X_n are uncorrelated and $Var(X_i) = \sigma^2 < \infty$ (var exists).

then $Var(\bar{X}_n) = \sigma^2/n \approx 0$ when $n \rightarrow \infty$ i.e. $\bar{X}_n \approx C_n$ when n is large enough!

■ Corollary. Suppose that X_1, \dots, X_n are uncorrelated and have same mean μ and variance σ^2 . Let

not assume i.i.d

c.f. $E(\bar{X}_n) = \mu$ c.f. definition of variance of X

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{n-1}$$

then $E(S^2) = \sigma^2$

Proof.

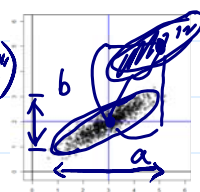
$$\begin{aligned} (n-1)S^2 &= \sum_{i=1}^n (X_i - \bar{X}_n)^2 \\ &= \sum_{i=1}^n [(X_i - \mu) - (\bar{X}_n - \mu)]^2 \\ &= \left[\sum_{i=1}^n (X_i - \mu)^2 \right] + \left[\sum_{i=1}^n (\bar{X}_n - \mu)^2 \right] \\ &\quad - 2(\bar{X}_n - \mu) \left[\sum_{i=1}^n (X_i - \mu) \right] = n(\bar{X}_n - \mu) \\ &= \left[\sum_{i=1}^n (X_i - \mu)^2 \right] + n(\bar{X}_n - \mu)^2 - 2n(\bar{X}_n - \mu)^2 \\ &= \left[\sum_{i=1}^n (X_i - \mu)^2 \right] - n(\bar{X}_n - \mu)^2. \end{aligned}$$

Note: $E(\bar{X}_n) = \mu$
 $Var(\bar{X}_n) = \frac{\sigma^2}{n}$

Therefore

$$\begin{aligned} (n-1)E(S^2) &= \left\{ \sum_{i=1}^n \underbrace{E[(X_i - \mu)^2]}_{Var(X_i)} \right\} - nE[(\bar{X}_n - \mu)^2] \\ &= n\sigma^2 - nVar(\bar{X}_n) = (n-1)\sigma^2. \end{aligned}$$

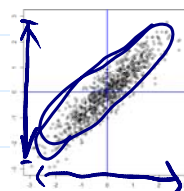
■ Note. The previous three corollaries also hold if X_1, \dots, X_n are independent (∴ "indep" implies "uncorrelated")



➤ Theorem (ρ of linear transformation).

$$\begin{aligned} Cor(a_0 + a_1X_1, b_0 + b_1Y_1) &= \text{sign}(a_1b_1) \times Cor(X_1, Y_1), \\ \text{and } |Cor(a_0 + a_1X, b_0 + b_1Y)| &= |Cor(X, Y)|, \end{aligned}$$

i.e., $|\rho_{XY}|$ is invariant under location and scale changes.
 ↪ why? check corollary in LN p1 7-8.



Proof. Let $S = a_0 + a_1X_1$ and $T = b_0 + b_1Y_1$, then

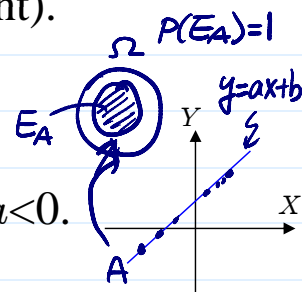
$$\begin{aligned} Cov(S, T) &= Cov(a_0 + a_1X_1, b_0 + b_1Y_1) = a_1b_1Cov(X_1, Y_1), \\ Var(S) &= a_1^2Var(X_1), \quad \text{and} \quad Var(T) = b_1^2Var(Y_1). \end{aligned}$$

Therefore,

$$\rho_{ST} = \frac{Cov(S, T)}{\sigma_S\sigma_T} = \frac{a_1b_1Cov(X_1, Y_1)}{|a_1||b_1|\sigma_X\sigma_Y} = \frac{a_1b_1}{|a_1b_1|} \rho_{XY} = \text{sign}(a_1b_1) \rho_{XY}.$$

➤ Theorem (some properties of correlation coefficient).

- (1) $-1 \leq \rho_{XY} \leq 1$. ($\Leftrightarrow |Cov(X, Y)| \leq \sigma_X\sigma_Y$)
- (2) $\rho_{XY} = \pm 1$ if and only if $P(Y = aX + b) = 1$.
 ↪ $Y = aX + b$ almost surely.
- (3) Furthermore, $\rho_{XY} = 1$, if $a > 0$ and $\rho_{XY} = -1$, if $a < 0$.



Proof of (1). $0 \leq Var\left(\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right)$ constants

$$\begin{aligned} &= Var\left(\frac{X}{\sigma_X}\right) + Var\left(\frac{Y}{\sigma_Y}\right) + 2Cov\left(\frac{X}{\sigma_X}, \frac{Y}{\sigma_Y}\right) \\ &= \frac{Var(X)}{\sigma_X^2} + \frac{Var(Y)}{\sigma_Y^2} + 2\frac{Cov(X, Y)}{\sigma_X\sigma_Y} \\ &= 1 + 1 + 2\rho_{XY} \Rightarrow \rho_{XY} \geq -1. \end{aligned}$$

Similarly,

$$0 \leq Var\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right) = 1 + 1 - 2\rho_{XY} \Rightarrow \rho_{XY} \leq 1.$$

Proof of (2) and (3). We see from the proof of (1),

$$\rho_{XY} = 1 \Leftrightarrow \text{Var} \left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y} \right) = 0.$$

$$\Leftrightarrow P \left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y} = c \right) = 1,$$

where c is a constant.

$$\Leftrightarrow P \left(Y = \frac{\sigma_Y}{\sigma_X} X + c\sigma_Y \right) = 1.$$

Note: $\text{Var}(Z) = 0$
if and only if
 $P(Z = a \text{ constant}) = 1$
i.e. $Z = a \text{ constant}$
almost surely

Similarly, $\rho_{XY} = -1 \Leftrightarrow P \left(Y = -\frac{\sigma_Y}{\sigma_X} X + c\sigma_Y \right) = 1.$

• **Q:** How to use expectations to (roughly) characterize random variables X_1, \dots, X_n ?

➤ $g(X_1, \dots, X_n) = X_i \Rightarrow E[g(\mathbf{X})] = \mu_{X_i}$: mean of X_i . *g: 1st order polynomial of X_1, \dots, X_n*

➤ $g(X_1, \dots, X_n) = (X_i - \mu_{X_i})^2 \Rightarrow E[g(\mathbf{X})] = \sigma_{X_i}^2$: variance of X_i .

➤ $g(X_1, \dots, X_n) = (X_i - \mu_{X_i})(X_j - \mu_{X_j})$ for $i \neq j$

$\Rightarrow E[g(\mathbf{X})] = \sigma_{X_i X_j}$: covariance of X_i and X_j .

➤ $g(X_1, \dots, X_n) = [(X_i - \mu_{X_i})/\sigma_{X_i}][(X_j - \mu_{X_j})/\sigma_{X_j}]$ for $i \neq j$

$\Rightarrow E[g(\mathbf{X})] = \rho_{X_i X_j}$: correlation coefficient of X_i and X_j . *g: 2nd order polynomial of X_1, \dots, X_n*

➤ Notes. $\mu_{X_i}, \sigma_{X_i}^2, \sigma_{X_i X_j}, \rho_{X_i X_j}$ are constants, not r.v.'s.

❖ **Reading:** textbook, Sec 7.1, 7.2, 7.4

Conditional Expectation ← conditional distribution LNp. 6-44 ~ 51

• Recall. $p_{Y|X}(y|x)$ or $f_{Y|X}(y|x)$ is a pmf/pdf for y, z

• Definition. The conditional expectation of $h(\mathbf{Y})$ given $\mathbf{X} = \mathbf{x}$, where

$h: \mathbb{R}^m \rightarrow \mathbb{R}^1$, is $E_{Y|X} \left(E(h(\mathbf{Y}) | \mathbf{X} = \mathbf{x}) \right) = \sum_{y \in \mathcal{Y}} h(y) p_{Y|X}(y|x)$,

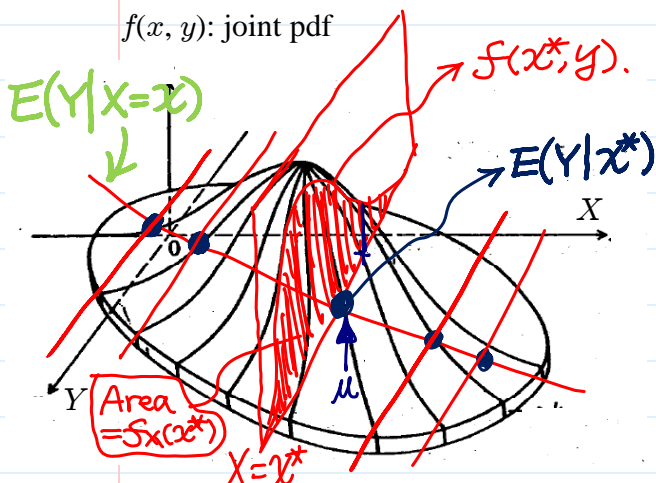
in the discrete case, or,

$$E_{Y|X} \left(E(h(\mathbf{Y}) | \mathbf{X} = \mathbf{x}) \right) = \int_{\mathbb{R}^m} h(y) f_{Y|X}(y|x) dy,$$

in the continuous case, provided that the sum or integral $= \int_{\mathcal{Y}} h(y) f_{Y|X}(y|x) dy$ converges absolutely.

i.e. exists

$f(x, y)$: joint pdf



➤ $f(x, y)$: a joint pdf.

➤ Fix x^* , is $f(x^*, y)$ a pdf of y ? i.e.,

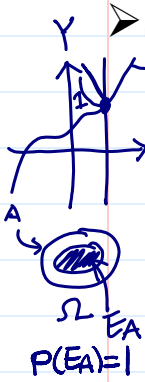
$$\int_{-\infty}^{\infty} f(x^*, y) dy = f_X(x^*) \stackrel{?}{=} 1$$

➤ $f_{Y|X}(y|x^*) = f(x^*, y) / f_X(x^*)$ is a pdf of y since

$$\frac{\int_{-\infty}^{\infty} f(x^*, y) dy}{f_X(x^*)} = 1.$$

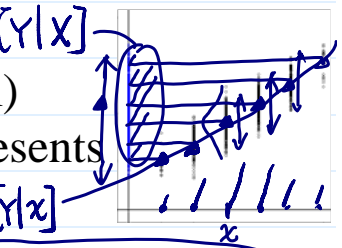
➤ $E(Y|x^*)$: mean of $f_{Y|X}(y|x^*)$.

➤ Do it for any $x = x^*$, and get a function of $x \Rightarrow E(Y|x)$



Some Notes.

- $E(h(Y)|X=x)$ is a function of x and is free of Y . *It's a constant line for x.*
- If X and Y are independent, then $E(h(Y)|X=x) = E[h(Y)]$. *$E_{Y|x}$ E_Y*
- $E[h(X)|X=x] = h(x)$. *$\because f_{Y|x}(y|x) = f_Y(y)$
 $P_{Y|x}(y|x) = P_Y(y)$*
- Let $g(x) = E[h(Y)|X=x]$, where $g: \mathbb{R}^n \rightarrow \mathbb{R}$, then we write $E(h(Y)|X)$ when x (a fixed value) replaced by X (a r.v.) in g .
- Notice that $g(X)$ is a random variable.



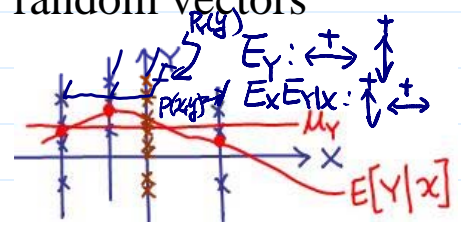
Q: What's the source of their randomness?

- Example. X =age (unit=year), Y =height (unit=cm)
- $Y|X=x$: a random variable (unit=cm) that represents the height distribution of people with age= x .
- $E(Y|X=x)$: a function maps from age (year) to average height (cm) of people with age= x . It is not a random variable.
- $E(Y|X)$: a random variable because it is a function of age, where age is treated as random. Notice that the unit of $E(Y|X)$ is "cm". $P(E(Y|X) = E(Y|x)) = P(X=x)$
- $Var(Y|X=x)$ and $Var(Y|X)$ can be similarly defined.
- $E(Y)$: average height of all people
- $Var(Y)$: variation of height of all people

c.f.

Theorem (Law of Total Expectation). For two random vectors X and Y ,

$$E_X\{E_{Y|X}[h(Y)|X]\} = E_Y[h(Y)].$$



In particular, let $h(Y) = Y_i$, we have

$$E_X[E_{Y|X}(Y_i|X)] = E_Y(Y_i) = E_{X,Y}[Y_i]$$

use the example given in LNp. 7-15 to realize what the terms mean.

Proof. (only prove it for the continuous case)

$$\begin{aligned} E_X\{E_{Y|X}[h(Y)|X]\} &= \int_{\mathbb{R}^n} E_{Y|X}(h(Y)|x) f_X(x) dx \\ &= \int_{\mathbb{R}^n} \left[\int_{\mathbb{R}^m} h(y) f_{Y|X}(y|x) dy \right] f_X(x) dx \\ &= \int_{\mathbb{R}^m} \int_{\mathbb{R}^n} h(y) \frac{f_{X,Y}(x,y)}{f_X(x)} f_X(x) dx dy \\ &= \int_{\mathbb{R}^m} h(y) \left[\int_{\mathbb{R}^n} f_{X,Y}(x,y) dx \right] dy = E_{X,Y}(h(Y)) \\ &= \int_{\mathbb{R}^m} h(y) f_Y(y) dy = E_Y[h(Y)]. \end{aligned}$$

Example. If a sample of n balls is drawn without replacement from a box containing R red balls, W white balls, and $N-R-W$ blue balls. Let

- X = # of red balls in the sample, *generalization of hypergeometric (similar to the binomial to multinomial)*
- Y = # of white balls in the sample, *generalization from binomial to multinomial*

then, the joint pmf of (X, Y) is

$$p_{X,Y}(x, y) = \frac{\binom{R}{x} \binom{W}{y} \binom{N-R-W}{n-x-y}}{\binom{N}{n}}$$

Find $E(Y)$.

Sol. Because $Y|X=x \sim \text{hypergeometric}(n-x, N-R, W)$,

$$g(x) \equiv E(Y|X=x) = (n-x)[W/(N-R)].$$

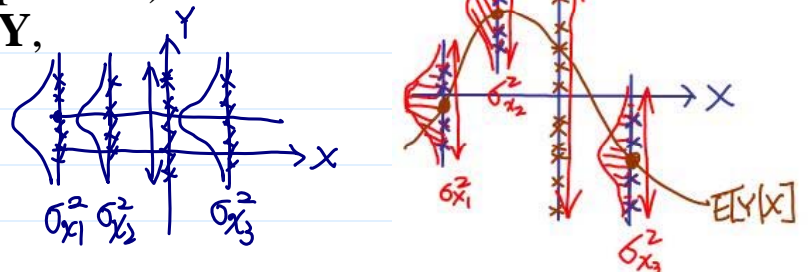
Because $X \sim \text{hypergeometric}(n, N, R) \Rightarrow E(X) = n(R/N)$, and then $E(Y) = E_X[E_{Y|X}(Y|X)] = E_X[g(X)]$

$$\begin{aligned} &= E_X \left[(n-X) \frac{W}{N-R} \right] = \frac{W}{N-R} [n - E_X(X)] \\ &= \frac{W}{N-R} \left(n - n \frac{R}{N} \right) = n \frac{W}{N}. \end{aligned}$$

Note that $Y \sim \text{hypergeometric}(n, N, W) \Rightarrow E(Y) = n(W/N)$.

The concept leads to the "Analysis of Variance (ANOVA)" Theorem (Variance Decomposition). For i s statistics, two random vectors \mathbf{X} and \mathbf{Y} ,

$$\begin{aligned} \text{Var}_{\mathbf{Y}}(Y_i) &= \text{Var}_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})] \\ &\quad + E_{\mathbf{X}}[\text{Var}_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]. \end{aligned}$$



Proof. $\text{Var}_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{x}) = E_{\mathbf{Y}|\mathbf{X}}\{[Y_i - E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{x})]^2|\mathbf{x}\}$

Note $\text{Var}(Z) = E(Z^2) - [E(Z)]^2$

$$\begin{aligned} &\cong E_{\mathbf{Y}|\mathbf{X}}(Y_i^2|\mathbf{x}) - [E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{x})]^2, \\ \text{and, } E_{\mathbf{X}}[\text{Var}_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})] &= E_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i^2|\mathbf{X})] - E_{\mathbf{X}}\{[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]^2\}. \end{aligned}$$

Also, $\text{Var}_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})] \equiv g(x)$

$$\cong E_{\mathbf{X}}\{[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]^2\} - \{E_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]\}^2.$$

Now, $\text{Var}_{\mathbf{Y}}(Y_i) = E_{\mathbf{Y}}(Y_i^2) - [E_{\mathbf{Y}}(Y_i)]^2$

Law of Total Expectation

$$\begin{aligned} &\ominus E_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i^2|\mathbf{X})] - \{E_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]\}^2 \\ &= E_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i^2|\mathbf{X})] \ominus E_{\mathbf{X}}\{[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]^2\} \\ &\quad \oplus E_{\mathbf{X}}\{[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]^2\} - \{E_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]\}^2 \\ &= E_{\mathbf{X}}[\text{Var}_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})] + \text{Var}_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]. \end{aligned}$$

➤ Corollary.

- $\text{Var}_{\mathbf{Y}}(Y_i) \geq E_{\mathbf{X}}[\text{Var}_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]$ and the equality holds if and only if $E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X}) = E_{\mathbf{Y}}(Y_i)$ with probability one.

$\text{Var}[E(Y|X)] = 0 \Rightarrow E(Y|X)$ is a constant over x & $E[E(Y|X)] = \mu_Y$.

- $\text{Var}_{\mathbf{Y}}(Y_i) \geq \text{Var}_{\mathbf{X}}[E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})]$ and the equality hold if and only if $\text{Var}_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X}) = 0$ ($\Rightarrow Y_i = E_{\mathbf{Y}|\mathbf{X}}(Y_i|\mathbf{X})$) with probability one.

$\leftarrow E[\text{Var}(Y|X)] = 0$.

Moment Generating Function

• Definition (Moment and Central Moment). If a random variable X has a cdf F_X , then

$$\mu_k \equiv E(X^k) = \int_{-\infty}^{\infty} x^k dF_X(x), \quad k = 1, 2, 3, \dots,$$

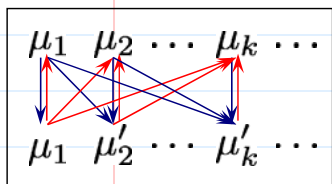
are called the k^{th} moments of X provided that the integral converges absolutely, and

$$\mu'_k \equiv E[(X - \mu_X)^k] = \int_{-\infty}^{\infty} (x - \mu_X)^k dF_X(x), \quad k = 1, 2, 3, \dots,$$

are called k^{th} moment about the mean μ_X or central moment of X provided that the integral converges absolutely.

► Some Notes.

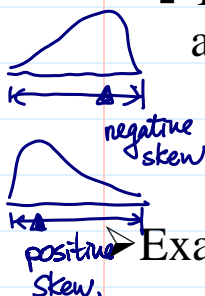
$$\begin{aligned} \mu'_k &= E[(X - \mu_X)^k] = E\left[\sum_{i=0}^k \binom{k}{i} (-\mu_X)^{k-i} X^i\right] \\ &= \sum_{i=0}^k \binom{k}{i} (-\mu_X)^{k-i} E(X^i) = \sum_{i=0}^k \binom{k}{i} (-\mu_X)^{k-i} \mu_i. \end{aligned}$$



$$\begin{aligned} \text{and, } \mu_k &= E(X^k) = E\{[(X - \mu_X) + \mu_X]^k\} \\ &= \sum_{i=0}^k \binom{k}{i} (\mu_X)^{k-i} E[(X - \mu_X)^i] \\ &= \sum_{i=0}^k \binom{k}{i} (\mu_X)^{k-i} \mu'_i. \end{aligned}$$

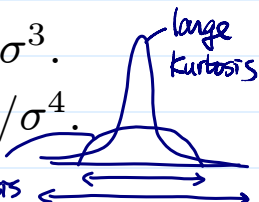
In particular, $E(X) = \mu_X = \mu_1$, and $Var(X) = \sigma_X^2 = \mu_2 - \mu_1^2 = \mu'_2$.

• The (central) moments give a lot of useful information about the distribution, e.g., in addition to mean and variance,



▫ Skewness (a measure of the asymmetry): μ'_3/σ^3 .

▫ Kurtosis (a measure of the "peakedness"): μ'_4/σ^4 .



Example (Uniform). If $X \sim \text{Uniform}(0, 1)$, then

$$\mu_k = \int_0^1 x^k dx = \frac{1}{k+1},$$

therefore, $\mu_X = \mu_1 = 1/2$, and,

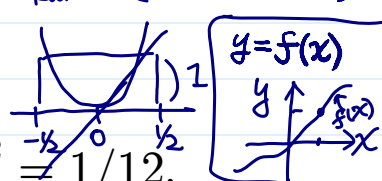
$$\sigma_X^2 = \mu_2 - \mu_1^2 = 1/3 - (1/2)^2 = 1/12.$$

how to characterize a distribution?

- ① pdf/pmf
- ② cdf
- ③ mgf

And, $\mu'_k = \int_0^1 (x - 1/2)^k dx = \frac{1}{k+1} [(1/2)^{k+1} - (-1/2)^{k+1}]$

$$= \begin{cases} 0, & k \text{ is odd,} \\ \frac{1}{(k+1)2^k}, & k \text{ is even.} \end{cases}$$



• Definition (Moment Generating Function). If X is a random variable with the cdf F_X , then

$$M_X(t) = E(e^{tX}) = \int_{-\infty}^{\infty} e^{tx} dF_X(x)$$

is called the *moment generating function* (mgf) of X provided that the integral converges absolutely in some non-degenerate interval of t .

$f(x) = \sum_{k=0}^{\infty} a_k \cdot x^k$

Taylor expansion $\frac{f^{(k)}(0)}{k!}$

Laplace Transformation $E(e^{tx})$

$f(x) = \begin{cases} \text{pdf} \\ \text{pmf} \end{cases}$

➤ Some Notes.

- The mgf is a function of the variable t .
- The mgf may only exist for some particular values of t .

➤ Example.

→ i.e. not for all $t \in \mathbb{R}$.

- If X is a discrete r.v. taking on values x_i with probability p_i , $i=1, 2, 3, \dots$, then $M_X(t) = \sum_{i=1}^{\infty} e^{tx_i} p_i$.

- If $X \sim \text{Poisson}(\lambda)$, then for $-\infty < t < \infty$,

$$M_X(t) = \sum_{x=0}^{\infty} \left(e^{tx} \times \frac{e^{-\lambda} \lambda^x}{x!} \right)$$

$$= e^{-\lambda} \left(e^{\lambda e^t} \right) \sum_{x=0}^{\infty} \frac{e^{-(\lambda e^t)} (\lambda e^t)^x}{x!} = e^{-\lambda} e^{\lambda e^t} = e^{\lambda(e^t - 1)}$$

pmf of Poisson(λe^t)

- If $X \sim \text{Exponential}(\lambda)$, then for $t < \lambda$,

$$M_X(t) = \int_0^{\infty} e^{tx} \times \lambda e^{-\lambda x} dx$$

$$= \lambda \left(\frac{1}{\lambda - t} \right) \int_0^{\infty} (\lambda - t) e^{-(\lambda - t)x} dx = \frac{\lambda}{\lambda - t}$$

and $M_X(t)$ does not exist for $t \geq \lambda$. *pdf of exponential ($\lambda - t$)*

- A list of some mgfs (**exercise**)

this must be > 0

- If $X \sim \text{Binomial}(n, p)$, *use binomial expansion (LNp. 4-20)*

$$M_X(t) = (1 - p + pe^t)^n, \text{ for } t < -\log(1 - p).$$

- If $X \sim \text{Negative Binomial}(r, p)$, *use negative binomial expansion (LNp. 4-25)*

$$M_X(t) = \left[\frac{pe^t}{1 - (1-p)e^t} \right]^r, \text{ for } t < -\log(1 - p).$$

- If $X \sim \text{Uniform}(\alpha, \beta)$, $M_X(t) = \frac{e^{\beta t} - e^{\alpha t}}{t(\beta - \alpha)}$.

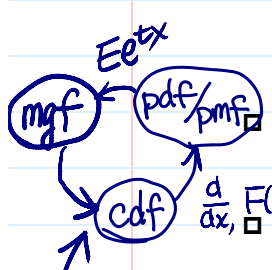
- If $X \sim \text{Gamma}(\alpha, \lambda)$,

$$M_X(t) = \left(\frac{\lambda}{\lambda - t} \right)^\alpha, \text{ for } t < \lambda.$$

*use $e^u = \sum_{k=0}^{\infty} \frac{u^k}{k!}$
& use STO*

- If $X \sim \text{Beta}(\alpha, \beta)$, $M_X(t) = 1 + \sum_{k=1}^{\infty} \left(\prod_{r=0}^{k-1} \frac{\alpha + r}{\alpha + \beta + r} \right) \frac{t^k}{k!}$

- If $X \sim \text{Normal}(\mu, \sigma^2)$, $M_X(t) = e^{\mu t + \frac{\sigma^2 t^2}{2}}$.



• Theorem (Uniqueness Theorem). Suppose that the mgfs $M_X(t)$ and $M_Y(t)$ of random variables X and Y exist for all $|t| < h$ for some $h > 0$.

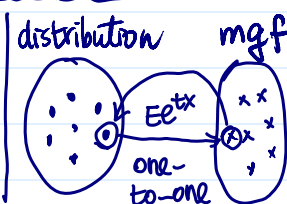
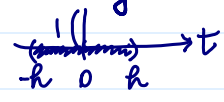
If *does not mean $X=Y$*

$$M_X(t) = M_Y(t),$$

s.i.e. an open interval containing zero

for $|t| < h$, then

$$F_X(z) = F_Y(z)$$



for all $z \in \mathbb{R}$, where F_X and F_Y are the cdfs of X and Y , respectively.

Proof. Skipped (by the uniqueness theorem of Laplace transform.) *Fourier transformation*

➤ Application of the uniqueness theorem

- When a moment generating function exists, there is a unique distribution corresponding to that mgf.
- This allows us to use mgfs to find distributions of transformed random variables in some cases.
- This technique is most commonly used for linear combinations of independent random variables

why? check the thms in Lnp. 724~25

➤ Example. If $M_X(t) = p_1 e^{a_1 t} + \dots + p_k e^{a_k t}$, where $p_1 + \dots + p_k = 1$, then X is a discrete r.v. and its pmf is

$$p_X(x) = \begin{cases} p_i, & \text{for } x = a_i, i = 1, \dots, k, \\ 0, & \text{otherwise.} \end{cases}$$

• Theorem (Moments and MGF). If $M_X(t)$ exist for $|t| < h$ for some $h > 0$, then

$$M_X(0) = 1,$$

This explains why it's called moment generating function

and,

$$M_X^{(k)}(0) = \mu_k, \quad k = 1, 2, 3, \dots$$

Proof. First,

$$M_X(0) = \int_{-\infty}^{\infty} e^{0 \cdot x} dF_X(x) = \int_{-\infty}^{\infty} 1 dF_X(x) = 1.$$

$\int_{-\infty}^{\infty} F(x) \Big|_{-\infty}^{\infty}$

$$\begin{aligned} M_X'(0) &= \left. \frac{d}{dt} M_X(t) \right|_{t=0} = \left[\left. \frac{d}{dt} \int_{-\infty}^{\infty} e^{tx} dF_X(x) \right|_{t=0} \right] \\ &= \int_{-\infty}^{\infty} \left(\left. \frac{d}{dt} e^{tx} \right|_{t=0} \right) dF_X(x) = \int_{-\infty}^{\infty} (x e^{tx} \Big|_{t=0}) dF_X(x) \\ &= \int_{-\infty}^{\infty} x \cdot 1 dF_X(x) = E(X) = \mu_1. \end{aligned}$$

... = ...

$$\begin{aligned} M_X^{(k)}(0) &= \left. \frac{d^k}{dt^k} M_X(t) \right|_{t=0} = \left[\left. \frac{d^k}{dt^k} \int_{-\infty}^{\infty} e^{tx} dF_X(x) \right|_{t=0} \right] \\ &= \int_{-\infty}^{\infty} \left(\left. \frac{d^k}{dt^k} e^{tx} \right|_{t=0} \right) dF_X(x) = \int_{-\infty}^{\infty} (x^k e^{tx} \Big|_{t=0}) dF_X(x) \\ &= \int_{-\infty}^{\infty} x^k \cdot 1 dF_X(x) = E(X^k) = \mu_k. \end{aligned}$$

➤ Example. If $X \sim \text{Exponential}(\lambda)$, then $M_X(t) = \frac{\lambda}{\lambda - t}$.

Because

$$M_X^{(k)}(t) = \frac{k! \lambda}{(\lambda - t)^{k+1}},$$

we get

$$\mu_k = M_X^{(k)}(0) = \frac{k!}{\lambda^k}.$$

We can use kth moments to calculate mean, variance, skewness, kurtosis, central moments, ...

• Theorem (MGF for linear transformation). For constants a and b ,

$$M_{a+bX}(t) = e^{at} M_X(bt).$$

can be use to identify the dist. of $a+bX$ from the dist. of X .

Proof. $M_{a+bX}(t) = E[e^{t(a+bX)}] = e^{at} E[e^{(bt)X}] = e^{at} M_X(bt).$

• Theorem (MGF for sum of independent r.v.'s). If X_1, \dots, X_n are independent each with mgfs $M_1(t), \dots, M_n(t)$, respectively, then the mgf of $S = X_1 + \dots + X_n$ is

$$\mathbb{R}^1 \rightarrow \mathbb{R}^1: M_S(t) = M_1(t) \times \dots \times M_n(t). \quad \text{--- (*)}$$

Proof. $M_S(t) = E(e^{tS}) = E[e^{t(X_1 + \dots + X_n)}]$
 $= E(e^{tX_1} \times \dots \times e^{tX_n}) \stackrel{\text{independent}}{=} E(e^{tX_1}) \times \dots \times E(e^{tX_n})$
 $= M_1(t) \times \dots \times M_n(t).$

Note: Geo. is a special case of negative binomial $r=1$

Example. If X_1, \dots, X_n are i.i.d. \sim Geometric(p), then $S = X_1 + \dots + X_n \sim$ Negative Binomial(n, p).
 use convolution (LN p. 6-25) to prove it

Proof. $M_S(t) = M_{X_1}(t) \times \dots \times M_{X_n}(t)$
 $= \frac{pe^t}{1-(1-p)e^t} \times \dots \times \frac{pe^t}{1-(1-p)e^t} = \left[\frac{pe^t}{1-(1-p)e^t} \right]^n.$

c.f. convolution approach LN p. 6-27

Example. If X_1, \dots, X_n are independent and $X_i \sim$ Normal(μ_i, σ_i^2), for $i=1, \dots, n$.

Let $S = a_0 + a_1X_1 + \dots + a_nX_n$, then

$$S \sim \text{Normal}(a_0 + a_1\mu_1 + \dots + a_n\mu_n, a_1^2\sigma_1^2 + \dots + a_n^2\sigma_n^2).$$

Proof. $M_S(t) = e^{a_0t} \times \prod_{i=1}^n e^{\mu_i(a_it) + \frac{\sigma_i^2(a_it)^2}{2}}$
 $= e^{(a_0 + a_1\mu_1 + \dots + a_n\mu_n)t + \frac{(a_1^2\sigma_1^2 + \dots + a_n^2\sigma_n^2)t^2}{2}}$

• Definition (Joint Moment Generating Function). For random variables X_1, \dots, X_n , their joint mgf is defined as

$$\mathbb{R}^n \rightarrow \mathbb{R}^1: M_{X_1, \dots, X_n}(t_1, \dots, t_n) = E(e^{t_1X_1 + \dots + t_nX_n})$$

provided that the expectation exists.

different t's \leftrightarrow (*) in LN p. 7-25.

Example. If $X_1, \dots, X_m \sim$ Multinomial(n, m, p_1, \dots, p_m),

The multinomial expansion in LN p. 6-13

$$M_{X_1, \dots, X_m}(t_1, \dots, t_m) = \sum_{x_1 + \dots + x_m = n} e^{t_1x_1 + \dots + t_mx_m} \binom{n}{x_1, \dots, x_m} p_1^{x_1} \dots p_m^{x_m}$$

$$= \sum_{x_1 + \dots + x_m = n} \binom{n}{x_1, \dots, x_m} (p_1 e^{t_1})^{x_1} \dots (p_m e^{t_m})^{x_m}$$

$$= (p_1 e^{t_1} + \dots + p_m e^{t_m})^n.$$

relationship between joint mgf & marginal mgf

• Some Properties of Joint mgf

➤ $M_{X_1}(t) = M_{X_1, X_2, \dots, X_n}(t, 0, \dots, 0).$

➤ uniqueness theorem

compare it with the (*) in LN p. 7-25

Note: same property holds for cdf, pdf/pmf

X_1, \dots, X_n are independent if and only if $M_{X_1, \dots, X_n}(t_1, \dots, t_n) = M_{X_1}(t_1) \times \dots \times M_{X_n}(t_n).$

$\frac{\partial^{k_1 + \dots + k_n}}{\partial t_1^{k_1} \dots \partial t_n^{k_n}} M_{X_1, \dots, X_n}(0, \dots, 0) = E(X_1^{k_1} \times \dots \times X_n^{k_n}).$