⊙ Theorem (best *constant* predictor under MSE). { a constant function of $x$ } minimum

$G_1$　$\underline{E_{X,Y}\ (Y-\underline{c})^2} = \underline{E_Y(Y-\underline{c})^2} \geq E_Y[Y-\underline{E_Y(Y)}]^2 = \underline{Var_Y(Y)}$

The equality holds if and only if $c=E_Y(Y)$. — only need to know $\mu_Y$

Proof. $R(Y)=(Y-c)^2$: a function of $Y$

$E_{X,Y}[R(Y)] = E_Y[R(Y)]$ ｜ Thm in LNp.5-19

‖

$E_Y E_{X|Y}[R(Y)|Y]$, $R(Y)$ — LNp.8-22

| 12/11 |

$\dfrac{E_Y(Y-\underline{c})^2}{}$

$= \dfrac{Var_Y(Y)+(\mu_Y-\underline{c})^2}{}$

$\geq Var_Y(Y)$ — cf. LNp.8-24

LNp.8-24 — minimum

⊙ Theorem (best predictor under MSE).

$G_3$　$E_{X,Y}[Y-\underline{g(X)}]^2 \geq E_{X,Y}[Y-\underline{E_{Y|X}(Y|X)}]^2 = \underline{E_X[Var_{Y|X}(Y|X)]}$

The equality holds if and only if $g(x)=E_{Y|X}(Y|x)$. —(*)　cf.(LNp.8-21)

a function of $X$ only　　a function of $X$ only

Proof. $E_{X,Y}[Y-g(X)]^2$

$= E_{X,Y}\{[\underline{Y}-E_{Y|X}(Y|\underline{X})]+[E_{Y|X}(Y|\underline{X})-g(\underline{X})]\}^2$

$= E_{X,Y}[Y-E_{Y|X}(Y|X)]^2 + E_X[E_{Y|X}(Y|X)-g(X)]^2$

$\quad + 2\cdot E_{X,Y}\{[Y-E_{Y|X}(Y|X)][E_{Y|X}(Y|X)-g(X)]\}$　=0

last "=" $= E_{X,Y}[Y-E_{Y|X}(Y|X)]^2 + E_X[E_{Y|X}(Y|X)-g(X)]^2$　=0 iff $g(X)=E_{Y|X}(Y|X)$

$\geq E_{X,Y}[Y-E_{Y|X}(Y|X)]^2$

where the last "=" comes from

$E_{X,Y}\{[\underline{Y}-E_{Y|X}(Y|\underline{X})][E_{Y|X}(Y|\underline{X})-g(\underline{X})]\}$ — $R(X,Y)$

$= E_X E_{Y|X}\{[\underline{Y}-E_{Y|X}(Y|X)][E_{Y|X}(Y|X)-g(X)]\,|\,\underline{X}\}$

By the law of total expectation (LNp.8-22)

$E_{XY}[R(X,Y)]=E_X E_{Y|X}[R(X,Y)|X]$　— this is a constant when conditioned on $X$

$= E_X\{[E_{Y|X}(Y|X)-g(X)]\ E_{Y|X}[Y-E_{Y|X}(Y|X)|\underline{X}]\} = 0.$　= $E_{Y|X}(Y|X)-E_{Y|X}(Y|X)$

important concept: mean is best predictor under MSE

Furthermore, (for (*) in LNp.8-27)

$E_{X,Y}[\underline{Y}-E_{Y|X}(Y|\underline{X})]^2$

$= E_X E_{Y|X}\{[Y-E_{Y|X}(Y|X)]^2|\underline{X}\} = E_X[Var_{Y|X}(Y|X)]$

➤ Some notes for the best predictor in $G_3$

cf. best in $G_1$ $E_Y(Y)$

● $E_{Y|X}(Y|x)$ is the best predictor of $Y$ based on $X$, in the sense of mean square prediction error — intuition ← check the graph in LNp.8-20

cf. ● Its calculation requires to know the joint distribution of $X$ and $Y$, or at least $E_{Y|X}(Y|x)$

　▪ $E_{Y|X}(Y|x)$ is called the regression function of $Y$ on $X$ 回錄

**⊙ Theorem (best _linear_ predictor under MSE).** $\boxed{-1 \le \rho_{XY} \le 1}$

$G_2$

$$E_{X,Y}[Y-(a+bX)]^2 \ge E_{X,Y}\left\{Y-\left[\mu_Y+\rho_{XY}\frac{\sigma_Y}{\sigma_X}(X-\mu_X)\right]\right\}^2$$

$$\boxed{minimum} = \sigma_Y^2(1-\rho_{XY}^2)$$

The equality holds if and only if $a=\mu_Y-b\mu_X$ and $b=\rho_{XY}\sigma_Y/\sigma_X$ $\boxed{unit=?}$

Proof. $E_{X,Y}(Y-a-bX)^2 - h(X,Y)\equiv Z$

$\boxed{\because Var(Z)=E(Z^2)-[E(Z)]^2}$

$= Var_{X,Y}(Y-a-bX)+[E_{X,Y}(Y-a-bX)]^2$

$= Var_{X,Y}(Y-bX)+(\mu_Y-a-b\mu_X)^2$

$\ge Var_{X,Y}(Y-bX)$  ($\Rightarrow$ setting $a=\mu_Y-b\mu_X$)

$\boxed{Thm in LNp.8-13}$

$= \sigma_Y^2+b^2\sigma_X^2-2b\sigma_{XY}$

$= \sigma_X^2\left(b^2-2b\frac{\sigma_{XY}}{\sigma_X^2}+\frac{\sigma_{XY}^2}{\sigma_X^4}\right)+\sigma_Y^2-\frac{\sigma_{XY}^2}{\sigma_X^2}$

$= \sigma_X^2\left(b-\frac{\sigma_{XY}}{\sigma_X^2}\right)^2+\sigma_Y^2(1-\rho_{XY}^2)$

$\ge \sigma_Y^2(1-\rho_{XY}^2)$   ($\Rightarrow$ setting $b=\frac{\sigma_{XY}}{\sigma_X^2}=\frac{\sigma_{XY}}{\sigma_X\sigma_Y}\times\frac{\sigma_Y}{\sigma_X}=\rho_{XY}\frac{\sigma_Y}{\sigma_X}$)

---

➤ Some notes for the best _linear_ predictor in $G_2$

$\boxed{best in G_3}$ ■ $E_{Y|X}(Y|x)=\mu_Y+(\rho_{XY}\sigma_Y/\sigma_X)(x-\mu_X)$ if $(X,Y)$ is distributed as bivariate normal. $\boxed{linear regression analysis}$ $\boxed{best in G_2}$

■ Its calculation requires to know the means, variances, and covariance of $X$ and $Y$. cf. $\boxed{best in G_3, G_1}$ $\boxed{Which one require more information?}$

$\boxed{MSE=0 if \rho_{XY}=\pm1, MSE=\sigma_Y^2 if \rho_{XY}=0}$ ■ $\sigma_Y^2(1-\rho_{XY}^2)$ is small if $\rho_{XY}$ is close to +1 or −1, and large if $\rho_{XY}$ is close to 0. $\boxed{intuition?}$ $\boxed{check the correlation plots in LNp.8-9}$ $\boxed{more information better predictor}$

• A comparison of these minimum MSEs

$\boxed{\because G_1 \cap G_2 \cap G_3}$ ➢ $\min_{a,b}E_{X,Y}[Y-(a+bX)]^2 \le \min_c E_{X,Y}(Y-c)^2$ and the equality holds if and only if $\rho_{XY}=0$.

➢ $\min_g E_{X,Y}[Y-g(X)]^2 \le \min_{a,b}E_{X,Y}[Y-(a+bX)]^2$ and the equality holds if and only if $E_{Y|X}(Y|x)=\mu_Y+(\rho_{XY}\sigma_Y/\sigma_X)(x-\mu_X)$.

❖ **Reading**: textbook, Sec 7.6

# Moment Generating Function

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

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

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

$$\mu'_k \equiv E[(X - \underbrace{\mu_X}_{\text{a constant}})^k] = \int_{-\infty}^{\infty} (x - \mu_X)^k \, dF_X(x), \quad k = 2, 3, \ldots,$$

are called $k^{\text{th}}$ *moment about the mean* $\mu_X$ or *central moment* of $X$ provided that the integral converges absolutely.
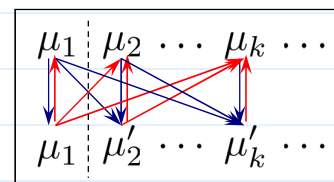
➤ Some notes.

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

  $= \sum_{i=0}^{k} \binom{k}{i} (-\mu_X)^{k-i} E(X^i) = \sum_{i=0}^{k} \binom{k}{i} (-\mu_X)^{k-i} \mu_i.$

- $\mu_k = E(X^k) = E\{[(X - \mu_X) + \mu_X]^k\}$

  $= \sum_{i=0}^{k} \binom{k}{i} (\mu_X)^{k-i} E[(X - \mu_X)^i]$

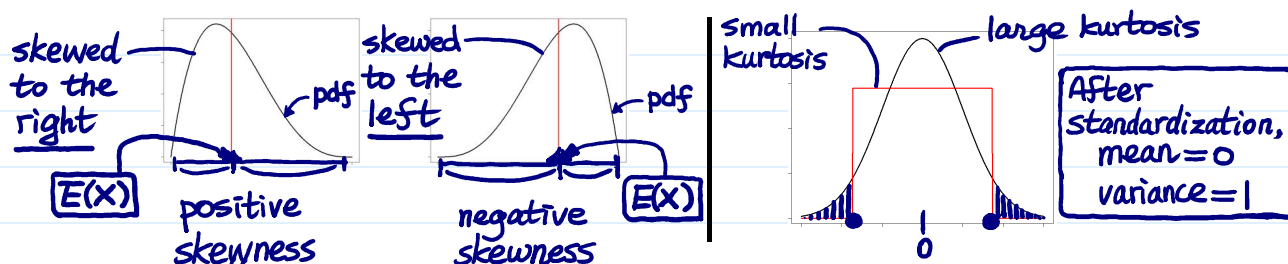  $= \sum_{i=0}^{k} \binom{k}{i} (\mu_X)^{k-i} \mu'_i.$   $\boxed{\mu'_0 = 1 \\ \mu'_1 = 0}$   $\boxed{\mu_0 = E(X^0) = 1}$

- In particular,

  $$E(X) = \mu_X = \mu_1, \quad \text{and,}$$
  $$Var(X) = \sigma_X^2 = \mu'_2 = \mu_2 - \mu_1^2. = E(X^2) - [E(X)]^2$$

  $$\begin{array}{ccc} \mu_1 & \mu_2 & \cdots & \mu_k & \cdots \\ \mu_1 & \mu'_2 & \cdots & \mu'_k & \cdots \end{array}$$

---

$\boxed{\text{Recall.} \\ \text{mean, var,} \\ \text{cov, cor}}$ → ■ The (central) moments give a lot of useful information about the distribution in addition to mean and variance, e.g.,

$\boxed{\text{defined by} \\ \text{expectation}}$

- ▫ Skewness (a measure of the asymmetry): $\mu'_3 / \sigma^3 := E\left(\frac{X - \mu}{\sigma}\right)^3$

- ▫ Kurtosis (a measure of the "heavy tails"): $\mu'_4 / \sigma^4 := E\left(\frac{X - \mu}{\sigma}\right)^4$



➤ Example (Uniform). If $X \sim \text{Uniform}(0, 1)$, then
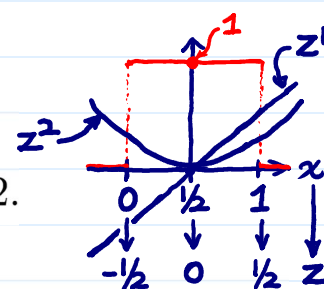
$$\mu_k = \int_0^1 x^k \, dx = \frac{1}{k+1},$$

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

$\mu'_2 \to \sigma_X^2 = \mu_2 - \mu_1^2 = 1/3 - (1/2)^2 = 1/12.$

And, $\mu'_k = \int_0^1 (x - 1/2)^k \, dx = \int_{-1/2}^{1/2} z^k \, dz$

$\boxed{\text{skewness} = 0 \\ \text{kurtosis} = 1.8}$

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

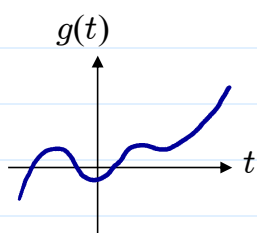- Recall. How to <u>characterize a distribution</u>?

  (1) <u>pdf/pmf</u>, (2) <u>cdf</u>, (3) <u>mgf</u>

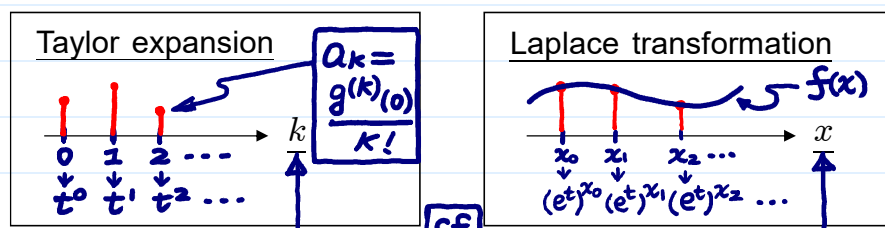- Definition (Moment Generating Function). If $X$ is a <u>random</u> <u>variable</u> with the <u>cdf $F_X$</u>, then

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

$$\color{blue}{= \int_{-\infty}^{\infty} e^{tx}\, \underline{f_X(x)dx} \text{ for } \underline{\text{continuous case}}}$$

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

$$g(t) \color{blue}{\boxed{=}} \sum_{k=0}^{\infty} a_k\, t^{\underline{k}} \qquad g(t) \color{blue}{\boxed{=}} \int_{\underline{\mathbb{R}}} \underline{f(x)}\, (\underline{e^t})^{\underline{x}}\, dx$$

<span style="color:blue">Q: how to express a function?</span>

$g(t)$



Taylor expansion

<span style="color:blue">$a_k = \dfrac{g^{(k)}(0)}{k!}$</span>

Laplace transformation

<span style="color:blue">cf.</span>

- ➤ Some <u>Notes</u>.
  - ▪ The <u>mgf</u> is a <u>function</u> of the <u>variable $t$</u>.
  - ▪ The <u>mgf</u> may <u>only exist</u> for <u>some particular values of $t$</u>. <span style="color:blue">← i.e., <u>not all</u> $t \in \mathbb{R}$</span>
  - ▪ <u>$M_X(t)$</u> always <u>exists</u> at $t=0$ and $M_X(0)=\underline{1}$ <span style="color:blue">← Thm in LNp.8-36</span>

- ➤ <u>Example</u>.
  - ▪ If <u>X</u> is a <u>discrete</u> r.v. taking on values <u>$x_i$</u>'s with probability <u>$p_i$</u>'s, $i=1, 2, 3, \ldots$, then

$$M_X(t) = E(\underline{e^{tX}}) = \sum_{i=1}^{\infty} e^{\underline{tx_i}}\underline{p_i}.$$

  - ▪ If $X \sim$ <u>Poisson($\lambda$)</u>, then for $-\infty<t<\infty$,

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

$$= e^{-\lambda}\left(e^{\lambda e^t}\right)\sum_{x=0}^{\infty}\underbrace{\frac{e^{-(\lambda e^t)}(\lambda e^t)^x}{x!}}_{\color{blue}{\text{pmf of Poisson}(\lambda e^t)}} = e^{-\lambda}e^{\lambda e^t} = \underline{e^{\lambda(e^t-1)}}.$$

  - ▪ If $X \sim$ <u>exponential($\lambda$)</u>, then for $t<\lambda$,

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

$$= \lambda\left(\tfrac{1}{\lambda-t}\right)\int_0^{\infty}\underbrace{(\lambda-t)\,e^{-(\lambda-t)x}}_{\color{blue}{\text{pdf of exponential}(\lambda-t)}}dx = \frac{\lambda}{\lambda-t},$$

<span style="color:blue">This must be $>0$</span>

<span style="color:red">STO</span>

  and $M_X(t)$ does <u>not exist</u> for $t \geq \lambda$.

  - ▪ A list of <u>some mgfs</u> (<span style="color:red">exercise</span>)
    - ▫ If $X \sim$ <u>binomial($n, p$)</u>, <span style="color:blue">use binomial expansion (LNp.5-23)</span>

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

**$r=1$, geometric distribution**

- If $X \sim$ negative binomial$(r, p)$, [use negative binomial expansion (LNp.5-29)]

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

**$\alpha=1$, exponential distribution**

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

$E(e^{tx})$

mgf — pdf/pmf — cdf $\frac{d}{dx}F_x$ — $F(x)-F(x-)$

- If $X \sim$ gamma$(\alpha, \lambda)$, [use STO]　[use $e^u = \sum_{k=0}^{\infty} \frac{u^k}{k!}$ & STO]

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

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

- If $X \sim$ normal$(\mu, \sigma^2)$, $M_X(\underline{t}) \doteq e^{\mu \underline{t}+(\sigma^2/2)\underline{t}^2}$. [use STO]

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

$$\underline{M_X(t) = M_Y(t)},$$

for $\underline{|t|<h}$, then

[i.e. an open interval containing zero]　[distribution $\xrightarrow[\text{Laplace transf.}]{}$ mgf]

[does not mean $X=Y$] ← $\underline{F_X(z) = F_Y(z)}$ ← [X,Y have same dist.]

　　 $E(e^{tx})$ one-to-one

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

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

---

➤ Application of the uniqueness theorem

- ▪ When a mgf exists for all $|t|<h$ for some $h>0$, 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.

[Find dist. of $X_1+\cdots+X_n$ (Check the Thms in LNp.8-38)]

- ▪ This technique is most commonly used for *linear combinations* of independent random variables $\underline{X_1, \ldots, X_n}$

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

[by uniqueness Thm & mgf in LNp.8-34]

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

- • Theorem (Moments and MGF). If $\underline{M_X(t)}$ exists for $|t|<h$ for some $h>0$, then

[can take derivative at $t=0$]

$$M_X(0)=\underline{1},$$

[related to the coefficients in the Taylor expansion of $M_X(t)$] → [Know all moments ⇒ know dist.]

and,　[kth derivative] →

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

[This explains why it's called *moment generating* function.]

Proof. First, $M_X(\underline{0}) = \int_{-\infty}^{\infty} e^{\underline{0 \cdot x}} \, d\underline{F_X(x)} = \underline{\int_{-\infty}^{\infty} 1 \, dF_X(x)} = \underline{1}.$

$\qquad = \underline{\underline{F_X(x)}}\Big|_{-\infty}^{\infty} =$

$\underline{M_X{'}(\underline{0})} = \frac{d}{dt} \underline{M_X(t)}\Big|_{\underline{t=0}} = \left[ \frac{d}{dt} \int_{-\infty}^{\infty} \underline{e^{tx}} \, dF_X(x) \right]\Big|_{\underline{t=0}}$

$\qquad = \underline{\int_{-\infty}^{\infty}} \left( \frac{d}{dt} \underline{e^{\underline{tx}}}\Big|_{\underline{t=0}} \right) dF_X(x) = \int_{-\infty}^{\infty} \left( \underline{xe^{tx}}\big|_{\underline{t=0}} \right) dF_X(x)$

$\qquad = \int_{-\infty}^{\infty} \underline{x \cdot 1} \, d\underline{F_X(x)} = \underline{E_X(X)} = \underline{\mu_1}.$

$\qquad \cdots = \cdots$

$\underline{M_X^{(k)}(\underline{0})} = \frac{d^k}{dt^k} \underline{M_X(t)}\Big|_{\underline{t=0}} = \left[ \frac{d^k}{dt^k} \int_{-\infty}^{\infty} \underline{e^{tx}} \, dF_X(x) \right]\Big|_{\underline{t=0}}$

$\qquad = \underline{\int_{-\infty}^{\infty}} \left( \frac{d^k}{dt^k} \underline{e^{\underline{tx}}}\Big|_{\underline{t=0}} \right) dF_X(x) = \int_{-\infty}^{\infty} \left( \underline{x^k e^{tx}}\big|_{\underline{t=0}} \right) dF_X(x)$

$\qquad = \int_{-\infty}^{\infty} \underline{x^k \cdot 1} \, d\underline{F_X(x)} = \underline{E_X(X^k)} = \underline{\mu_k}.$

➢ Example. If $\underline{X} \sim \underline{\text{exponential}(\lambda)}$, then $\underline{M_X(t) = \frac{\lambda}{\lambda - \underline{t}}}.$ ← LNp.8-34

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

Then, can use kth moments to obtain mean, variance, skewness, kurtosis, ⋯, kth central moments,⋯

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