

## Lecture 36: Dec 7

### Last time

- MGF cont.
- Covariance and Correlation

### Today

- Course evaluations (13/38)
- Final exam format
  - Final exam will be take home
  - Open book, open note, not open internet
  - Final exam will be released on Friday (12/09/2022) right after class
  - Final exam due 23:59 pm on Friday 12/16/2022.
  - Scan and submit your exam via email with a single pdf file
  - Send your email to both your instructor and your TA.
  - Submitted exams should be human-readable to receive non-zero scores.
- Random Samples
- Convergence
- Central Limit Theorem

### Random Samples

**Definition** 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 and identically distributed (iid) random variables with the same pdf or pmf  $f(x)$ .

If  $X_1, \dots, X_n$  are iid, then their joint pdf or pmf is

$$f(x_1, \dots, x_n) = f(x_1)f(x_2) \dots f(x_n) = \prod_{j=1}^n f(x_j)$$

**Statistics** Let  $X_1, \dots, X_n$  be a random sample and let  $T(x_1, \dots, x_n)$  be a function defined on  $\mathbb{R}^n$ . Then the random variable  $Y = T(X_1, \dots, X_n)$  is called a *statistic*. The probability distribution of  $Y$  is called the *sampling distribution* of  $Y$ .

Note:  $T$  is only a function of  $(x_1, \dots, x_n)$ , no parameters.

## Examples

$$\begin{aligned}\text{sample mean} \quad \bar{X} &= \frac{1}{n} \sum_{j=1}^n X_j \\ \text{sample variance} \quad S^2 &= \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})^2 \\ \text{sample standard deviation} \quad S &= \sqrt{S^2} \\ \text{minimum} \quad X_{(1)} &= \min_{1 \leq i \leq n} X_i\end{aligned}$$

**Properties** Let  $x_1, \dots, x_n$  be  $n$  numbers and define

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j, \quad s^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2$$

Then

$$\begin{aligned}\min_a \sum_{j=1}^n (x_j - a)^2 &= \sum_{j=1}^n (x_j - \bar{x})^2 \\ (n-1)s^2 &= \sum_{j=1}^n (x_j - \bar{x})^2 = \sum_{j=1}^n x_j^2 - n\bar{x}^2\end{aligned}$$

**Residuals** Lemma: Let  $X_1, \dots, X_n$  be a random sample from a population with mean  $\mu$  and variance  $\sigma^2$ . Define the residuals  $R_i = X_i - \bar{X}$ . Then

$$\begin{aligned}E(R_i) &= 0, \quad \text{Var}(R_i) = \frac{n-1}{n}\sigma^2 \\ \text{Cov}(R_i, \bar{X}) &= 0, \quad \text{Cov}(R_i, R_j) = -\sigma^2/n \text{ if } i \neq j\end{aligned}$$

**Theorem** Let  $X_1, \dots, X_n$  be a random sample from a population with mgf  $M_X(t)$ . Then the mgf of the sample mean is

$$M_{\bar{X}}(t) = [M_X(t/n)]^n$$

## Convergence

**Convergence in Probability** A sequence of random variables  $X_1, \dots, X_n$  *converges in probability* to a random variable  $X$ , denoted

$$X_n \xrightarrow{p} X$$

if for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr(|X_n - X| < \epsilon) = 1$$

or equivalently

$$\lim_{n \rightarrow \infty} \Pr(|X_n - X| > \epsilon) = 0$$

In other words,  $X_n$  is more and more likely to be close to  $X$ , or less and less likely to be far from  $X$ .

**Example** Let  $X_n = X + \epsilon_n$ , where  $\epsilon_n \sim N(0, 1/n)$  and  $X$  is an arbitrary random variable. Then, as  $n \rightarrow \infty$ ,

$$X_n \xrightarrow{p} X$$

**Weak law of large numbers (WLLN)** Let  $Y_1, \dots, Y_n$  be iid with common mean  $\mu$  and variance  $\sigma^2$ . Then, as  $n \rightarrow \infty$ ,

$$\bar{Y}_n = \frac{1}{n} \sum_{j=1}^n Y_j \xrightarrow{p} \mu$$

*Proof:*

The proof is quite simple, being a straightforward application of Chebychev's Inequality. We have, for every  $\epsilon > 0$ ,

$$\Pr(|\bar{Y}_n - \mu| \geq \epsilon) = \Pr(|\bar{Y}_n - \mu|^2 \geq \epsilon^2) \leq \frac{E(\bar{Y}_n - \mu)^2}{\epsilon^2} = \frac{Var(\bar{Y}_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty$$

**Convergence in Distribution** A sequence of random variables  $X_1, \dots, X_n$  converges in distribution to a random variable  $X$ , denoted

$$X_n \xrightarrow{d} X$$

if

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$$

This is also called *convergence in law* or *weak convergence*. In other words, the distribution of  $X_n$  is closer and closer to the distribution of  $X$ .

**Relation between "in distribution" and "in probability"** Theorem:

1. Convergence in probability implies convergence in distribution:

$$X_n \xrightarrow{p} X \Rightarrow X_n \xrightarrow{d} X$$

2. Suppose  $X_n \xrightarrow{d} X$  where  $X$  has a degenerate distribution, i.e.  $\Pr\{X = a\} = 1$  for some  $a \in \mathbb{R}$ . Then,

$$X_n \xrightarrow{d} a \Rightarrow X_n \xrightarrow{p} a$$

**Convergence in Distribution via Convergence of Mgf's** Theorem: Suppose the mgf  $M_n(t)$  of  $Y_n$  exists for  $|t| < h$ , and the mgf  $M(t)$  of  $Y$  exists for  $|t| < h_1 < h$ . Then,

$$Y_n \xrightarrow{d} Y \iff \lim_{n \rightarrow \infty} M_n(t) = M(t), \quad |t| < h_1$$

**Example** Let  $X_\lambda \sim \text{Poisson}(\lambda)$ . Then, as  $\lambda \rightarrow \infty$ ,

$$\frac{X_\lambda - \lambda}{\lambda} \xrightarrow{p} 0$$

$$\frac{X_\lambda - \lambda}{\sqrt{\lambda}} \xrightarrow{d} N(0, 1)$$

**Central Limit Theorem** Let  $X_1, X_2, \dots, X_n$  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  $h > 0$ ). Let  $EX_i = \mu$  and  $Var(X_i) = \sigma^2 > 0$ . (Both  $\mu$  and  $\sigma^2$  are finite since the mgf exists) Define  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . 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 limiting standard normal distribution, in other words,  $\sqrt{n}(\bar{X}_n - \mu)/\sigma \xrightarrow{d} N(0, 1)$

*Proof:*

Define  $Y_i = (X_i - \mu)/\sigma$ , and let  $M_Y(t)$  denote the common mgf of  $Y_i$ s, which exists for  $|t| < \sigma h$  and  $M_Y(t) = M_{\frac{1}{\sigma}X_i - \mu/\sigma}(t) = e^{-\frac{t}{\sigma} \mu} M_X(\frac{t}{\sigma})$ . Since

$$\frac{\sqrt{n}(\bar{X}_n)}{\sigma} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i,$$

we have,

$$\begin{aligned} M_{\sqrt{n}(\bar{X}_n - \mu)/\sigma}(t) &= M_{\sum_{i=1}^n Y_i/\sqrt{n}}(t) \\ &= M_{\sum_{i=1}^n Y_i}(t/\sqrt{n}) \\ &= [M_Y(t/\sqrt{n})]^n. \end{aligned}$$

We now expand  $M_Y(t/\sqrt{n})$  in a Taylor series (power series) around 0.

$$M_Y\left(\frac{t}{\sqrt{n}}\right) = \sum_{k=0}^{\infty} M_Y^{(k)}(0) \frac{(t/\sqrt{n})^k}{k!},$$

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

$$M_Y\left(\frac{t}{\sqrt{n}}\right) = 1 + \frac{(t/\sqrt{n})^2}{2!} + R_Y\left(\frac{t}{\sqrt{n}}\right),$$

where  $R_Y$  is the remainder term in the Taylor expansion such that

$$\lim_{n \rightarrow \infty} \frac{R_Y(t/\sqrt{n})}{(t/\sqrt{n})^2} = 0.$$

Therefore, for any fixed  $t$ , we can write

$$\begin{aligned} \lim_{n \rightarrow \infty} \left[ M_Y\left(\frac{t}{\sqrt{n}}\right) \right]^n &= \lim_{n \rightarrow \infty} \left[ 1 + \frac{(t/\sqrt{n})^2}{2!} + R_Y\left(\frac{t}{\sqrt{n}}\right) \right]^n \\ &= \lim_{n \rightarrow \infty} \left[ 1 + \frac{1}{n} \left( \frac{t^2}{2} + nR_Y\left(\frac{t}{\sqrt{n}}\right) \right) \right]^n \\ &= e^{t^2/2} \end{aligned}$$