

IE605: Engineering Statistics

Lecture 07

Manjesh K. Hanawal

Previous Lecture:

- ▶ Exponential Family of Distributions
- ▶ Population and Random Sampling
- ▶ Sample mean, variance and standard deviation
- ▶ Sampling from Normal distribution
- ▶ Student's t-distribution

This Lecture:

- ▶ F-distributions
- ▶ Convergence of RVs
- ▶ Consistency
- ▶ Order Statistics
- ▶ Generating Random Samples

F-distributions

We would be interested in variability of populations:

$$(X_1, X_2, \dots, X_n) \text{ are iid and } X_i \sim \mathcal{N}(\mu_X, \sigma_X^2) \quad \forall i$$

$$(Y_1, Y_2, \dots, Y_m) \text{ are iid and } Y_j \sim \mathcal{N}(\mu_Y, \sigma_Y^2) \quad \forall j$$

We would estimate $\frac{S_X^2}{S_Y^2}$. What is its distribution?

$$\frac{S_X^2/S_Y^2}{\sigma_X^2/\sigma_Y^2} = \frac{S_X^2/\sigma_X^2}{S_Y^2/\sigma_Y^2} = \frac{(n-1)S_X^2/(n-1)\sigma_X^2}{(m-1)S_Y^2/(m-1)\sigma_Y^2} = \frac{\chi_{n-1}^2/(n-1)}{\chi_{m-1}^2/(m-1)}$$

$\frac{S_X^2/S_Y^2}{\sigma_X^2/\sigma_Y^2}$ has F-distribution with (n-1) and (m-1) degree of freedom

F-distribution is named in the honor of Sir Ronald Fisher!

F-distributions contd...

PDF of F distribution with p and q degrees of freedom ($F_{p,q}$)

$$f_F(x) = \frac{\Gamma\left(\frac{p+q}{2}\right)}{\Gamma\left(\frac{p}{2}\right)\Gamma\left(\frac{q}{2}\right)} \left(\frac{p}{q}\right)^2 \frac{x^{p/2-1}}{[1 + (p/q)x]^{(p+q)/2}} \quad x > 0$$

$$F_{p,q} = \frac{U/p}{V/q} \text{ where } U \sim \chi_p^2, V \sim \chi_q^2, \text{ and independent}$$

Derivation of pdf of $F_{p,q}$

- ▶ $X = \frac{U/p}{V/q} = \frac{q}{p} \frac{U}{V}$ and $Y = V$
- ▶ As U, V are independent $f(U, V) = f(U)f(V)$
- ▶ Find joint distribution of (X, Y) by applying transformations
- ▶ Find marginal distribution of X

Properties of F -distribution

- ▶ **Claim 1:** If $X \sim F_{p,q}$, then $1/X \sim F_{q,p}$

$X = \frac{U/p}{V/p}$ where $U \sim \chi_p^2$, $V \sim \chi_q^2$ and are independent

$1/X = \frac{V/q}{U/p}$, hence $1/X \sim F_{q,p}$

- ▶ **Claim 2:** if $X \sim t_p$, then $X^2 \sim F_{1,p}$

$X = \frac{U}{\sqrt{V/p}}$, where $U \sim \mathcal{N}(0, 1)$, $V \sim \chi_p^2$ and are independent

$X^2 = U^2/(V/p) = \chi_1^2/(V/p) = (\chi_1^2/1)/(\chi_p^2/p) \sim F_{1,p}$

- ▶ **Claim 3:** if $X \sim F_{p,q}$, then $\frac{(p/q)X}{1+(p/q)X} \sim \text{beta}(p/2, q/2)$
(Exercise!)

Convergence of Sequence of RVs

What happens when the number of samples goes to infinity
(theoretical artifact)

Convergence in Probability: A sequence of RVs X_1, X_2, \dots , converge in probability to a random variable X if, for an $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0 \text{ or } \lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1$$

- ▶ In the definition X_1, X_2, \dots need not be i.i.d or independent
- ▶ Compactly written as $X_n \xrightarrow{P} X$ in probability.

- ▶ Suppose X_1, X_2, \dots are i.i.d. with common mean μ and variance $\sigma^2 > \infty$. From LLN, we know

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} \mu$$

- ▶ For an $\epsilon > 0$

$$\begin{aligned} P(|\bar{X}_n - \mu| \geq \epsilon) &\leq \frac{\mathbb{E}(|\bar{X}_n - \mu|^2)}{\epsilon^2} \\ &= \frac{\text{Var}(\bar{X}_n)}{\epsilon^2} = \frac{\sigma^2/n}{\epsilon} \rightarrow 0 \end{aligned}$$

- ▶ Sample mean converges to population mean!
- ▶ $\mathbb{E}(\bar{X}_n) = \mu$ (unbiased).

Consistency of Sample mean and Sample Variance

Consistency: A sample quantity is consistent if its sequence converges to a constant

▶ **Sample mean** is consistent: $\bar{X}_n \xrightarrow{P} \mu$ (by LLN)

▶ Is **sample variance** consistent?

$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$. We know $\mathbb{E}(S_n^2) = \sigma^2$

$$P(|S_n^2 - \sigma^2| \geq \epsilon) \leq \frac{\mathbb{E}((S_n^2 - \sigma^2)^2)}{\epsilon^2} = \frac{\text{Var}(S_n^2)}{\epsilon^2}$$

if $\text{Var}(S_n^2) \rightarrow 0$, then $S_n^2 \xrightarrow{P} \sigma^2$ (hence consistent)

▶ Is **sample standard deviation** consistent? (Exercise!)

Other Convergence types

Almost sure convergence: A sequence of RVs X_1, X_2, \dots convergence to X almost surely if $P\left(\lim_{n \rightarrow \infty} X_n = X\right) = 1$.

Denoted as $X_n \xrightarrow{a.s.} X$.

Convergence in distribution: A sequence of RVs X_1, X_2, \dots convergence to X in distribution if $\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$ for

all continuity points of F_X . Denoted as $X_n \xrightarrow{d} X$

$$X_n \xrightarrow{a.s.} X \implies X_n \xrightarrow{p} X \implies X_n \xrightarrow{d} X.$$

Order Statistics

Smallest, largest, middle observation of a random sample are useful

- ▶ Highest temperature in the last 50 years
- ▶ Lowest rainfall in the last 50 years
- ▶ median value of stock index in the last month

The order static of a random sample X_1, X_2, \dots, X_n are the sample value placed in the ascending order, denotes by $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ where $X_{(1)} \leq X_{(2)} \leq \dots, \leq X_{(n)}$

$$X_{(1)} = \min_{1 \leq i \leq n} X_i$$

$$X_{(2)} = \text{second smallest } X_i$$

\vdots

$$X_{(n)} = \max_{1 \leq i \leq n} X_i$$

Sample mean vs Sample Median

- ▶ Sample range: $X_{(n)} - X_{(1)}$
- ▶ Sample median:

$$M = \begin{cases} X_{((n+1)/2)} & \text{if } n \text{ is odd} \\ (X_{(n/2)} + X_{((n/2)+1)})/2 & \text{if } n \text{ is even} \end{cases}$$

- ▶ Example:
Random sample: 24, 89, 59, 34, 55, 81, 45, 93, 85, 50
Order statistic: 24, 34, 45, 50, 55, 59, 81, 85, 89, 93
Sample range: $93 - 24 = 69$
Sample mean: 61.5
Median: 57
- ▶ Median gives better indication of "typical" values than means!

Sample Percentile

For any $p \in [0, 1]$, the $(100p)$ th percentile is the observation such that approximately np of the observations are less than this observation and $n(1 - p)$ of the observations are greater.

- ▶ For $p = 0.5$, 50th percentile gives median
- ▶ For any $b \in \mathbb{R}_+$, define

$$\{b\} = \begin{cases} \lceil b \rceil & \text{if } \lceil b \rceil \leq b + 0.5 \\ \lfloor b \rfloor & \text{if } b - 0.5 < \lfloor b \rfloor \end{cases}$$

- ▶ $\frac{1}{2} < np < n - \frac{1}{2} \implies \frac{1}{2n} < p < 1 - \frac{1}{2n}$

Lower and Upper Quartile

$$(100p)\text{th sample percentile is } = \begin{cases} X_{(\{np\})} & \text{if } p < 0.5 \\ X_{(n+1-\{n(1-p)\})} & \text{if } p > 0.5 \end{cases}$$

Example 1: $n = 50, p = .35, np = 17.5, \{np\} = 18$. 35th sample percentile is $X_{(18)}$

Example 2: $n = 50, p = .65, n(1-p) = 17.5, \{n(1-p)\} = 18$
 $n + 1 - \{n(1-p)\} = 50 + 1 - 18 = 33$. 65th sample percentile is $X_{(33)}$

- ▶ For $p < 0.5$ and $p > 0.5$ sample percentiles exhibit symmetry
- ▶ if $(100p)$ th sample percentile is i th smallest observation, then $100(1-p)$ th sample percentile is the i th largest observation
- ▶ 25th sample percentile is called lower quartile
- ▶ 75th sample percentile is called upper quartile

Distribution of Order Statistics

Discrete Case:

Random sample X_1, X_2, \dots, X_n come from a discrete distributions with pmf $P_X(x_i) = p_i$, where $x_1 < x_2 < \dots$ are the possible realizations in ascending order. For any x_i , what is $P(X_{(j)} \leq x_i)$?

$$P_0 = 0$$

$$P(X \leq x_1) = P_1 = p_1$$

$$P(X \leq x_2) = P_2 = p_1 + p_2$$

$$\vdots$$

$$P(X \leq x_i) = P_i = p_1 + p_2 + \dots + p_i$$

$$\vdots$$

Discrete Case contd..

- ▶ Fix some x_i . Define $Y_j = \mathbb{1}_{\{X_j \leq x_i\}}$ for all $j = 1, 2, \dots, n$
- ▶ $P(Y_j = 1) = P_i$ for all $j = 1, 2, \dots, n$
- ▶ $Y = \sum_{j=1}^n Y_j$, $Y \in \{0, 1, 2, \dots, n\}$
- ▶ As X_i s are i.i.d, Y_j s are i.i.d. $Y_j \sim \text{Ber}(P_i)$.
- ▶ $Y \sim \text{Bin}(n, P_i)$. Y is sum of n $\text{Ber}(P_i)$ RVs
- ▶ $\{X_{(j)} \leq x_i\} = \{Y \geq j\}$. Hence $P(X_{(j)} \leq x_i) = P(Y \geq j)$
- ▶ $P(Y \geq j) = \sum_{k=j}^n \binom{n}{k} P_i^k (1 - P_i)^{n-k}$.

$$P(X_{(j)} \leq x_i) = \sum_{k=j}^n \binom{n}{k} P_i^k (1 - P_i)^{n-k}$$

$$\begin{aligned} P(X_{(j)} = x_i) &= P(X_{(j)} \leq x_i) - P(X_{(j)} \leq x_{i-1}) \\ &= \sum_{k=j}^n \binom{n}{k} \left(P_i^k (1 - P_i)^{n-k} - P_{i-1}^k (1 - P_{i-1})^{n-k} \right) \end{aligned}$$

Continuous case

Random sample X_1, X_2, \dots, X_n come from a population with pdf $f_X(x)$, and CDF $F_X(x)$. Let $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ denote the order statistics. Then, pdf of $X_{(j)}$ is

$$f_{X_{(j)}}(x) = \frac{n!}{(j-1)!(n-j)!} f_X(x) (F_X(x))^{j-1} (1 - F_X(x))^{n-j}$$

Joint pdf of $X_{(i)}$ and $X_{(j)}$ for $1 \leq i < j \leq n$ is

$$f_{X_{(i)}, X_{(j)}}(u, v) = \frac{n!}{(i-1)!(j-1-i)(n-j)!} \times \\ f_X(u) f_X(v) (F_X(u))^{i-1} (F_X(v) - F_X(u))^{j-1-i} (1 - F_X(v))^{n-j}$$