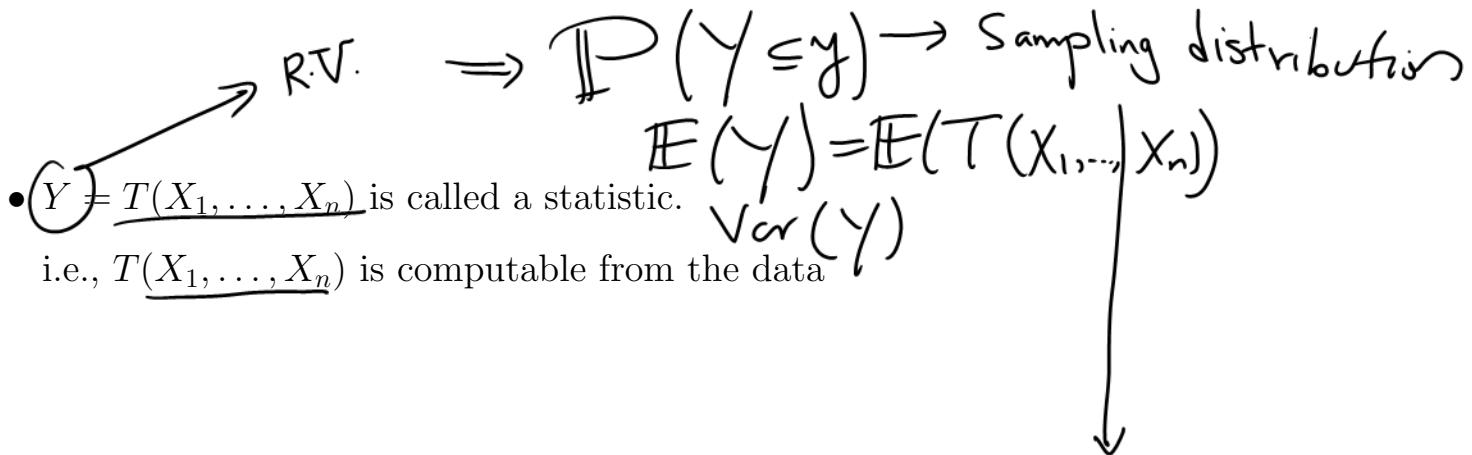# Random samples and iid variables

Definitions

- *Definition:* If $X_1, \ldots, X_n$ are **independent identically distributed** (iid) with $X_i \sim f_X(x_i)$, then we call $X_1, \ldots, X_n$ a **random sample** from the population $f_X(x)$.

R.V. $\implies \mathbb{P}(Y \leq y) \to$ Sampling distribution

$\mathbb{E}(Y) = \mathbb{E}(T(X_1, \ldots, X_n))$

$\text{Var}(Y)$

- $Y = T(X_1, \ldots, X_n)$ is called a statistic.

  i.e., $T(X_1, \ldots, X_n)$ is computable from the data

- The distribution of a statistic $Y$ is sometime called the **sampling distribution** of the statistic.

- Examples

  1. sample mean: $\bar{X}_n = \sum_{i=1}^n X_i / n$

  2. sample variance:

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

  $X_{(1)}$  3. minimum: $\min\{X_1, \ldots, X_n\}$

  $X_{(n)}$  4. maximum: $\max\{X_1, \ldots, X_n\}$

186

# Random samples and iid variables

Distribution of $\bar{X}_n$

Let $X_1, \ldots, X_n$ be a random sample from $f_X(x)$ with $\mu = \mathrm{E}X_i$ and $\sigma^2 = \mathrm{Var}(X_i)$

*i. i. d*

**Important Results for** $\bar{X}_n$: If $X_1, \ldots, X_n$ is a sample random with $\mu = \underline{\mathrm{E}X_i}$ and $\sigma^2 = \underline{\mathrm{Var}(X_i)}$, then

1. $\boxed{\mathrm{E}\bar{X}_n = \mu}$

2. $\mathrm{Var}(\bar{X}_n) = \dfrac{\sigma^2}{n}$

$$M_{\bar{X}_n}(t) = \left[M_{X_1}\left(t/n\right)\right]^n \quad \text{If } X_1 \cdots X_n \text{ are } iid$$

MGF approach can sometimes apply for determining the exact distribution of $\bar{X}_n$

(MGF of *i.i.d* $X_i$'s)

$$M_{\bar{X}_n}(t) = \mathrm{E}e^{t\bar{X}_n} = \mathrm{E}e^{n^{-1}t(X_1+\cdots+X_n)} = \mathrm{E}\prod_{i=1}^{n} e^{n^{-1}tX_i} = \prod_{i=1}^{n} \mathrm{E}e^{n^{-1}tX_i} = [M_{X_1}(t/n)]^n$$

def of MGF     def of $\bar{X}$     $e^{a+b} = e^a \cdot e^b$

$\mathrm{E}\left[e^{t/n X_1} e^{t/n X_2} \cdots e^{t/n X_n}\right]$ $X_i$ are ind.

$\mathrm{E}\left[e^{t/n X_1}\right] \cdots \mathrm{E}\left[e^{t/n X_n}\right]$ $X_i$ are identically dis.

$\left(\mathrm{E}\left[e^{t/n X_1}\right]\right)^n = \left(M_{X_1}(t/n)\right)^n$

**Examples**

1. Suppose $X_1, \ldots, X_n$ are iid Gamma$(\alpha, \beta)$

$$M_{\bar{X}_n}(t) = \left[M_{X_1}\left(t/n\right)\right]^n$$

$X_i \sim$ Gamma$(\alpha, \beta)$

$M_X(t) = (1-\beta t)^{-\alpha}$

$\left[M_{X_1}(t/n)\right]$

$$= \left[\left(1-\beta t/n\right)^{-\alpha}\right]^n = \left(1-\dfrac{\beta t}{n}\right)^{-n\alpha} \Rightarrow \bar{X}_n \sim \text{Gamma}\left(n\alpha, \dfrac{\beta}{n}\right)$$

2. Suppose $X_1, \ldots, X_n$ are iid $N(\mu, \sigma^2)$

$$M_{\bar{X}_n}(t) = \left[M_{X_1}\left(t/n\right)\right]^n$$

$X_i \sim N(\mu, \sigma^2)$

$M_X(t) = e^{\mu t + \sigma^2 \frac{t^2}{2}}$

$$\left[e^{\frac{\mu t}{n} + \frac{\sigma^2 t^2}{2n^2}}\right]^n = e^{\mu t + \left(\frac{\sigma^2}{n}\right)\frac{t^2}{2}}$$

$$\Rightarrow \bar{X}_n \sim N\left(\mu, \dfrac{\sigma^2}{n}\right)$$

# Random samples and iid variables

Distribution of $S^2$ (Sample Variance $S^2$)

Let $X_1, \ldots, X_n$ be a *i.i.d* random sample from $f_X(x)$ with $\mu = EX_i$ and $\sigma^2 = \text{Var}(X_i)$

- The exact sampling distribution of sample variance

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

  is difficult to obtain in general

- However, if $X_1, \ldots, X_n$ are iid $N(\mu, \sigma^2)$ then the sampling distribution of $S^2$ can be found (later... after scaling, the distribution is chi-square with $n-1$ degrees of freedom)

- **Result:** For random samples with $\mu = EX_i$ and $\sigma^2 = \text{Var}(X_i)$,

$$ES^2 = \sigma^2 = \text{Var}(X_i)$$

def of $S^2$

$EX_i^2 = \text{Var}X_i + (EX_i)^2$

$\text{Var}(\bar{X}_n) = \dfrac{\sigma^2}{n}$ (⊛)

$\text{Var}(\bar{X}_n) = E[(\bar{X}_n)^2] - (E(\bar{X}_n))^2$

$\dfrac{\sigma^2}{n} = E[(\bar{X}_n)^2] - \mu^2$

$\Rightarrow E[(\bar{X}_n)^2] = \mu^2 + \dfrac{\sigma^2}{n}$

Proof:
$$E[S^2] = E\left[\left(\frac{1}{n-1}\right)\left[\sum_{i=1}^{n} X_i^2 - n(\bar{X}_n)^2\right]\right]$$

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

$$= \frac{1}{n-1}\sum_{i=1}^{n}(\sigma^2 + \mu^2) \qquad -n\left(\mu^2 + \frac{\sigma^2}{n}\right)$$

$$= \frac{1}{n-1}\left[n\sigma^2 + n\mu^2 - n\mu^2 - \sigma^2\right] = \frac{\sigma^2[n-1]}{n-1} = \sigma^2$$

$$\Rightarrow \boxed{E[S^2] = \sigma^2}$$

188

# Random samples and iid variables

## Distribution of Maximum and Minimum

$F_X(x) = P(X \leq x)$ where $X$ has the same distribution of $X_1, \ldots, X_n$

i.i.d

Let $X_1, \ldots, X_n$ be a random sample with common cdf $F_{X_1}(x) = P(X_1 \leq x)$

Let $X_{(n)} = \max\{X_1, \ldots, X_n\}$ and $X_{(1)} = \min\{X_1, \ldots, X_n\}$

**Important results:**

1. $F_{X_{(n)}}(x) = P(X_{(n)} \leq x) = [F_{X_1}(x)]^n$ for $x \in \mathbb{R}$

   def of CDF

2. $F_{X_{(1)}}(x) = P(X_{(1)} \leq x) = 1 - [1 - F_{X_1}(x)]^n$, for $x \in \mathbb{R}$

   def of CDF

3. If the population cdf $F_{X_1}(x) = P(X_1 \leq x)$ is continuous with pdf $f_{X_1}(x) = \frac{dF_{X_1}(x)}{dx}$, then $X_{(n)}$ and $X_{(1)}$ both have pdfs given by

$$f_{X_{(n)}}(x) = nf_{X_1}(x)[F_{X_1}(x)]^{n-1}, \qquad f_{X_{(1)}}(x) = nf_{X_1}(x)[1 - F_{X_1}(x)]^{n-1}$$

$f_{X_{(n)}}(x) \overset{def}{=} \frac{d}{dx} F_{X_{(n)}}(x) = \frac{d}{dx}\left[F_{X_1}(x)\right]^n = n\left(\frac{d}{dx}F_{X_1}(x) = f_{X_1}(x)\right)\left(F_{X_1}(x)\right)^{n-1} = nf_{X_1}(x)\left[F_{X_1}(x)\right]^{n-1}$

*Proofs:* (These are proofs that are useful to remember.)

$F_{X_{(n)}}(x) \overset{def}{=} P(X_{(n)} \leq x) = P\left(\max(X_1, \ldots, X_n) \leq x\right)$

$X_i$'s are independent $\overset{def\ of\ X_{(n)}}{=} P(X_1 \leq x, X_2 \leq x, \ldots, X_n \leq x)$

$X_i$'s are identically dist. $= P(X_1 \leq x)P(X_2 \leq x) \cdots P(X_n \leq x)$

$= [P(X_1 \leq x)]^n = [F_{X_1}(x)]^n$

$F_{X_{(1)}}(x) \overset{def\ of\ CDF}{=} P(X_{(1)} \leq x) \overset{P(A^c) = 1 - P(A)}{=} 1 - P(X_{(1)} > x)$

$= 1 - P(X_1 > x, X_2 > x, \ldots, X_n > x)$

$= 1 - P(X_1 > x)P(X_2 > x) \cdots P(X_n > x)$

$= 1 - [P(X_1 > x)]^n = 1 - [1 - P(X_1 \leq x)]^n$

$F_{X_{(1)}}(x) = 1 - [1 - F_{X_1}(x)]^n$

$= 1 - [1 - F_{X_1}(x)]^n$

# Random samples and iid variables

Order statistics

- *Definition:* The **order statistics** for a sample $X_1, \ldots, X_n$ are the values in ascending order denoted as

$$X_{(1)} \leq X_{(2)} \leq \cdots \leq X_{(n)}$$

- Primarily interested in iid $X_1, \ldots, X_n$ having a continuous distribution

- For random samples we may be interested in

  1. the distribution of a single order statistic $X_{(i)}$

  2. the distribution of two or more order statistics $(X_{(i)}, X_{(j)})$

  3. function of two or more order statistics

     e.g., range $R = X_{(n)} - X_{(1)}$

- order statistics are a type of (discontinuous) transformation of $X_1, \ldots, X_n$

# Random samples and iid variables

Distribution of $k$th order statistic

**Result 1:** If $X_1, \ldots, X_n$ are a $\overset{i.i.d}{\underline{\text{random sample}}}$ with common cdf $F_{X_1}(x)$, then the cdf of the $\underline{k\text{th order statistic}}$ (given some $k = 1, \ldots, n$) is given by

$$F_{X_{(k)}}(x) = P(X_{(k)} \leq x) = P(\text{at least } k \; X_i\text{'s} \leq x) = \sum_{j=k}^{n} \binom{n}{j} [F_{X_1}(x)]^j [1 - F_{X_1}(x)]^{n-j}$$

*Proof:* Let $Y = \#$ of $X_i$'s which are less than or equal to $x$

$$\Rightarrow Y \sim Bin(n, F_{X_1}(x))$$

$$P(Y \geq k) \underset{Y \sim Bin(n, F_{X_1}(x))}{=\!=\!=} \sum_{j=k}^{n} P(Y=j) = \sum_{j=k}^{n} \binom{n}{j} \left(F_{X_1}(x)\right)^j \left(1 - F_{X_1}(x)\right)^{n-j}$$

**Result 2** (pdf in continuous case): If $X_1, \ldots, X_n$ are a random sample with common continuous cdf $F_{X_1}(x)$ and pdf $f_{X_1}(x)$, the pdf of the $k$th order statistic is

$$f_{X_{(k)}}(x) = \frac{dF_{X_{(k)}}(x)}{dx} = \frac{n!}{(k-1)!(n-k)!} f_{X_1}(x) [F_{X_1}(x)]^{k-1} [1 - F_{X_1}(x)]^{n-k}$$

- Heuristic argument for the form of the pdf $f_{X_{(k)}}(x)$:

  $k-1$ observations $\leq x$; 1 observation in $(x, x+dx)$; $n-k$ observations $> x$

- A formal proof uses derivative of cdf + algebra (see next slide)

Note: in the discrete case, the pmf of $X_{(k)}$ is obtained as

$$f_{X_{(k)}}(x) = P(X_{(k)} = x) = P(X_{(k)} \leq x) - P(X_{(k)} < x) = F_{X_{(k)}}(x) - \lim_{y \uparrow x} F_{X_{(k)}}(y)$$

191