

## statistical inference:

### Recap:

Population:  $X \sim f_{\theta}(x)$

$\theta$  is unknown

$x_1, x_2, \dots, x_n$  (random sample)

- use data to learn about  $\theta$

(Point) estimator:  $\hat{\theta}$

a statistic used to estimate  $\theta$

### Desirable properties of $\hat{\theta}$ :

- Unbiased
- Small variance
- Consistent
- Asymptotically normal

## Descriptive statistics:

- sample mean  $\bar{X}$  estimates population mean  $\mu$ .

# Some descriptive statistics and what they

Mean:  $\rightarrow$  measure of center of the dist.

Have a random sample  $x_1, \dots, x_n$  from the population of  $X$ .

Population mean:

$$\mu = E[X]$$

Sample mean:

$$\bar{X} = \frac{\sum x_i}{n}$$

Properties of  $\bar{X}$ : natural estimator of  $\mu$  [may be possible to find better estimator]

- $\bar{X}$  is ~~a~~ consistent.  $\bar{X}$  is unbiased.  $\text{Var}[\bar{X}] = \frac{1}{n} \text{Var}(X)$
- LLN:  $\bar{X}$  is ~~a~~ consistent.  $E[\bar{X}] = \mu$  for all  $n \rightarrow \infty$
- CLT:  $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$  if  $n$  is large (asymptotic normality)  
→ Have distribution of the data if the parent distn. is not symmetric.
- Greatly affected by outliers

Ex: (CPU data):  $\bar{X} = ?$   $\bar{X} = \frac{70 + 36 + \dots + 19}{30} = 48.22$  (?)

using R

$SE[\bar{X}] = ?$  — will see this later.

## Median:

Another measure of center of a dist.

→ more wifte is  
sol'n be dist.  
if discrete.

Population median: A value  $M$  such that

$$\{ P(X > M) \leq 0.5, \text{ and } P(X < M) \leq 0.5. \} \quad \textcircled{X}$$

Essentially  $M$  is the ~~middle~~ middle value — it divides the probability distribution in two ~~equal~~ halves.

$M$  for a Continuous distribution:

- ~~F(x)~~ F( $M$ ) ~~is~~  $\frac{1}{2}$
- CDF  $F(x) = P[x \leq x] = P[x < x]$  is const. and  $\uparrow$  fn.  $\forall x$ .
- $P[X > M] = 1 - F(M) \leq 0.5$

$$F(M) \leq 0.5$$

$$P[X > M] = 1 - F(M) \leq 0.5$$

$$\Rightarrow (F(M) \geq 0.5)$$

$$F(M) = 0.5$$

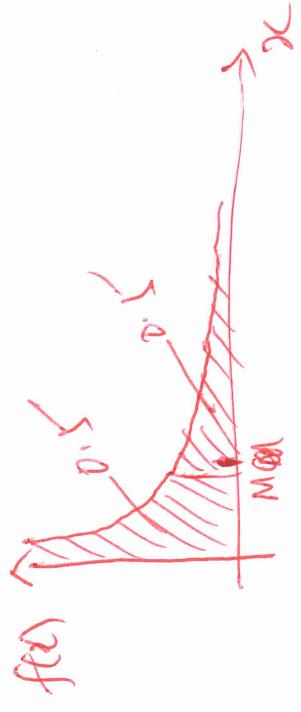
$$\Rightarrow$$



Solve this equation to  
get  $M$  — unique.

To force alignment, revise  $\hat{x}$  w.  
 $M$  is the smallest  $x$  such that  $F(x) \geq 0.5$ : —  $\hat{x}$

**Ex:** Suppose  $X \sim \text{Exponential}(\lambda)$ . Recall its cdf,  
 $F(x) = 1 - e^{-\lambda x}$  for  $x > 0$ . What is  $M$ ?



Solve for  $M$ :

$$F(M) = 1 - e^{-\lambda M} = 0.5$$

$$\Rightarrow e^{-\lambda M} = 0.5$$

$$\Rightarrow -\lambda M = \log(0.5)$$

$$\Rightarrow M = -\frac{1}{\lambda} \log(0.5)$$

## $M$ for a discrete distribution:

Problem 1:  $F(M) = 0.5$  may have a whole interval of roots.

- Median not unique
- Take the mid-point of the interval as the median.

Problem 2:  $F(M) = 0.5$  may not have any root. Take the smallest  $x$  with  $F(x) \geq 0.5$  as the median.

**Ex:** Look at Figure 8.4 and find the median.

From textbook:

## SIMPLE DESCRIPTIVE STATISTICS

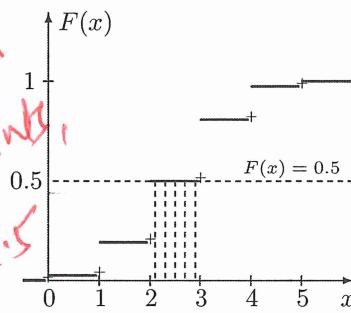
231

Any value  
you want & to  
is valid median.

Alt 1:

Arg of  
the endpoints  
i.e.,  $2+3=2.5$   
is the median.

(a) Binomial ( $n=5, p=0.5$ )  
many roots



(b) Binomial ( $n=5, p=0.3$ )  
no roots

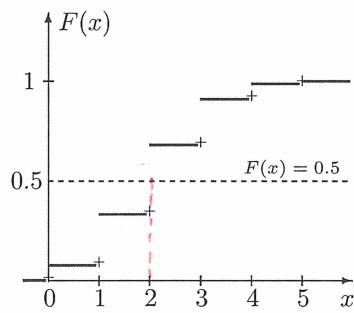


Figure 8.4 Computing medians of discrete distributions.

### Computing sample medians

Alt 2: New def.  
M=2.

A sample is always discrete, it consists of a finite number of observations. Then, computing a sample median is similar to the case of discrete distributions.

In simple random sampling, all observations are equally likely, and thus, equal probabilities on each side of a median translate into an equal number of observations.

Again, there are two cases, depending on sample size  $n$ .

#### Sample median

If  $n$  is odd, the  $\left(\frac{n+1}{2}\right)$ -th smallest observation is a median.

If  $n$  is even, any number between the  $\left(\frac{n}{2}\right)$ -th smallest and the  $\left(\frac{n+2}{2}\right)$ -th smallest observations is a median.

**Example 8.11 (MEDIAN CPU TIME).** Let's compute the median of  $n = 30$  CPU times from the data on p. 225.

## Sample median

- midpoint (or middle value) in the sample.

→ book:

$$\hat{M} = \begin{cases} \text{middle obs.} & , n \text{ is odd} \\ \text{average of the two middle obs.} & , n \text{ is even} \end{cases}$$

average of the  
two middle  
obs.

- can also use new def. of median:

- estimate  $\hat{M}$ .
- estimate  $M_{\text{var}}[\hat{M}]$  → will use bootstrap!
- difficult to estimate  $M_{\text{var}}[\hat{M}]$  for this purpose.

# 15      # 16.

139

n = 30.

Ex: GPU data:

9, 15, 19,

42, 43,

46,

139

sorted data:

14 obs.

$$\hat{M} = \frac{42+43}{2} \quad 14 \text{ obs.}$$

= 42.5!

$X$  = income of a randomly selected Dallas resident  $\rightarrow$  typically right-skewed.

Symmetric dist:

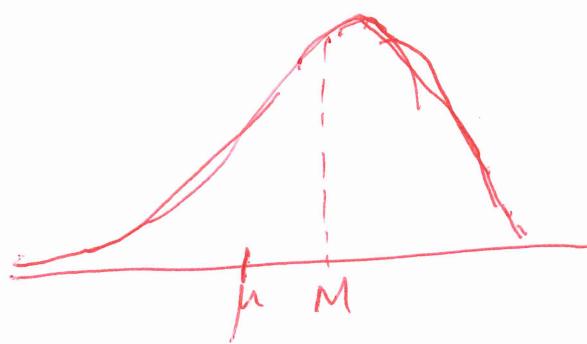


Two tails are equal

- makes sense to use mean ~~as center~~ for center of the dist.
- use  $\bar{X}$  to estimate  $\mu$ .

provided  $\mu$  exists.

left-skewed dist:

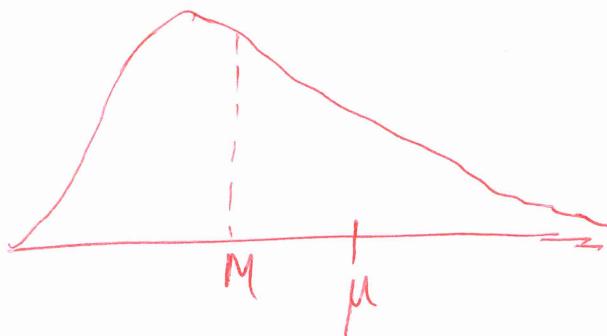


left tail is longer

- Mean ( $\mu$ ) is not a good measure of center of dist.
- don't use  $\bar{X}$ .

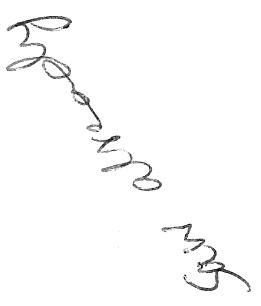
Right-skewed dist:

↑  
same situation with right-skewed dist.  
long right tail



# Characteristic shapes of a distribution

Symmetric:

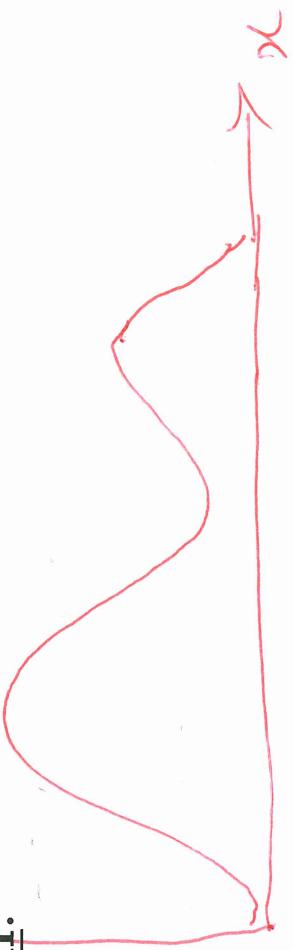


Right-skewed:

Left-skewed:

## Bimodal or multi-modal:

Several bumps in  
the dist.  
(mixture distr.)



Which measure of center to use — mean or median?:

When the  
dist. is  
symmetric.

# Descriptive statistics and what they estimate (continued)

*(continues from previous slide)*

p-quantile of a population: A value  $q_p$  such that

$$P(X < q_p) \leq p, \text{ and } P(X > q_p) \leq 1 - p.$$

Essentially  $X$  has  $p$  probability on the left of  $q_p$ .

p-quantile of a sample: A number  $\hat{q}_p$  that exceeds *at most*  $100p\%$  of the sample (i.e.,  $np$  observations) and is exceeded by *at most*  $100(1-p)\%$  of the sample (i.e.,  $n(1-p)$  observations).

- $\hat{q}_p$  estimates  $q_p$
- 0.5-quantile =  $M$

Population quartiles:  $(Q_1, Q_2, Q_3) = (q_{0.25}, q_{0.50}, q_{0.75})$

— they divide the distribution in four equal parts.

- Sample quartiles: *Simply divide the data in 4 equal parts using sample counterpart of Q<sub>1</sub>, Q<sub>2</sub>, Q<sub>3</sub>, max.*
- 5-number summary:
  - *If a sample data set: (min, Q<sub>1</sub>, Q<sub>2</sub>, Q<sub>3</sub>, max).*