# QUIZ

Fill in the blanks below:

For a large random sample of size $n$ from a distribution with mean $\mu$ and standard deviation $\sigma$, sample mean has an approximately ……………. distribution with mean …… and standard deviation ……...

# QUIZ

Fill in the blanks below:

For a large random sample of size $n$ from a distribution with mean $\mu$ and standard deviation $\sigma$, sample mean has an approximately …… Normal …… distribution with mean ……$\mu$…… and standard deviation ……$\sigma/\sqrt{n}$……

$$\left( \text{Say } \sigma^2 \text{ is variance of population.} \right.$$

Variance of $\bar{X}$ is $\dfrac{\sigma^2}{n}$.

$$\left. \text{std. dev. } '' \quad '' \quad '' \quad \sqrt{\dfrac{\sigma^2}{n}} = \dfrac{\sigma}{\sqrt{n}} \right)$$

# Outline

➢ We covered sampling distribution of means and CLT (8.1-8.4).

➢ Sampling distribution of variance and t-distribution (8.5, 8.6): We will cover them when discussing estimation problems.

Today

➢ We will continue with some more examples of CLT.

➢ Next we will discuss estimation problems: Chapter 9

# Another Example on Central Limit Theorem

Example: $X_i$: customer spending with $\mu = 80, \sigma = 40$.

Approximate the probability that the average spending of $100$ customers is $10\%$ or more below average. Use norm.cdf(-2)=0.023.

Solution:

Example: $X_i$: customer spending with $\mu = 80, \sigma = 40$.

Approximate the probability that the average spending of $100$ customers is $10\,\%$ or more below average. Use norm.cdf(-2)=0.023.

Solution:

$$P(\bar{X} \leq 72) = ?$$

# Another Example on Central Limit Theorem

Example: $X_i$: customer spending with $\mu = 80, \sigma = 40$.

Approximate the probability that the average spending of $100$ customers is $10\,\%$ or more below average. Use norm.cdf(-2)=0.023.

Solution:

$$P\left(\bar{X} \leq 72\right) = ?$$

$$P\left(\frac{\bar{X} - \mu_{\bar{X}}}{\sigma_{\bar{X}}} \leq \frac{72 - 80}{\frac{40}{\sqrt{100}}}\right)$$

$$2$$

# Another Example on Central Limit Theorem

Example: $X_i$: customer spending with $\mu = 80, \sigma = 40$.

Approximate the probability that the average spending of $100$ customers is $10\,\%$ or more below average. Use norm.cdf(-2)=0.023.

Solution:

$$P\left(\overline{X} \leq 72\right) = ?$$

$$P\left(\frac{\overline{X} - \mu_{\overline{x}}}{\sigma_{\overline{x}}} \leq \frac{72-80}{\frac{40}{\sqrt{100}}}\right) = P\left(z \leq \frac{-8 \times 10}{40}\right)$$

$$z$$

$$= P(z \leq -2)$$

$$= 0.023$$

Revisit:

Linear combinations/transformations of normal independent variables:

If $X_1, X_2$ are independent and each is normally distributed

then $Y = a_1 X_1 + a_2 X_2 + b$ has a normal distribution.

Its mean is $a_1 \mu_1 + a_2 \mu_2 + b$ and its variance is $a_1^2 \sigma_1^2 + a_2^2 \sigma_2^2$.

Apply it to the difference of two sample means:

If independent samples of size $n_1$, $n_2$ are drawn at random from two populations, with means $\mu_1$, $\mu_2$ and variances $\sigma_1^2$, $\sigma_2^2$, respectively, then the sampling distribution of the differences of means, $\overline{X}_1 - \overline{X}_2$ is approximately normally distributed with mean $\mu_1 - \mu_2$ and variance $\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}$.

**Case Study 8.2:** **Paint Drying Time**: Two independent experiments are run in which two different types of paint are compared. Eighteen specimens are painted using type $A$, and the drying time, in hours, is recorded for each. The same is done with type $B$. The population standard deviations are both known to be 1.0.

Assuming that the mean drying time is equal for the two types of paint, find $P(\bar{X}_A - \bar{X}_B > 1.0)$, where $\bar{X}_A$ and $\bar{X}_B$ are average drying times for samples of size $n_A = n_B = 18$. Use norm.cdf(3)=0.9987.

Solution:

# Example on CLT and Difference Between Two Means

**Case Study 8.2:** **Paint Drying Time**: Two independent experiments are run in which two different types of paint are compared. Eighteen specimens are painted using type $A$, and the drying time, in hours, is recorded for each. The same is done with type $B$. The population standard deviations are both known to be 1.0.

Assuming that the mean drying time is equal for the two types of paint, find $P(\bar{X}_A - \bar{X}_B > 1.0)$, where $\bar{X}_A$ and $\bar{X}_B$ are average drying times for samples of size $n_A = n_B = 18$. Use norm.cdf(3)=0.9987.

Solution:

$$\text{Let} \quad Y = \bar{X}_A - \bar{X}_B$$

$$\mu_Y = 0$$

$$\sigma^2_Y = \frac{1}{18} + \frac{1}{18} = \frac{1}{9} \implies \sigma_Y = \frac{1}{3}$$

**Case Study 8.2:** **Paint Drying Time**: Two independent experiments are run in which two different types of paint are compared. Eighteen specimens are painted using type $A$, and the drying time, in hours, is recorded for each. The same is done with type $B$. The population standard deviations are both known to be 1.0.

Assuming that the mean drying time is equal for the two types of paint, find $P(\bar{X}_A - \bar{X}_B > 1.0)$, where $\bar{X}_A$ and $\bar{X}_B$ are average drying times for samples of size $n_A = n_B = 18$. Use norm.cdf(3)=0.9987.

Solution:

Let $Y = \bar{X}_A - \bar{X}_B$

Assuming we can apply CLT on $\bar{X}_A, \bar{X}_B$
$\Rightarrow Y$ has normal distribution.

$\mu_Y = 0$

$$P(Y > 1.0) = P\left(\frac{Y - \mu_Y}{\sigma_Y} > \frac{1.0 - 0}{\frac{1}{3}}\right)$$

$$\sigma^2_Y = \frac{1}{18} + \frac{1}{18} = \frac{1}{9} \Rightarrow \sigma_Y = \frac{1}{3}$$

**Case Study 8.2:** **Paint Drying Time**: Two independent experiments are run in which two different types of paint are compared. Eighteen specimens are painted using type $A$, and the drying time, in hours, is recorded for each. The same is done with type $B$. The population standard deviations are both known to be 1.0.

Assuming that the mean drying time is equal for the two types of paint, find $P(\bar{X}_A - \bar{X}_B > 1.0)$, where $\bar{X}_A$ and $\bar{X}_B$ are average drying times for samples of size $n_A = n_B = 18$. Use norm.cdf(3)=0.9987.

Solution:

$$\text{Let } Y = \bar{X}_A - \bar{X}_B$$

$$\mu_Y = 0$$

$$\sigma^2_Y = \frac{1}{18} + \frac{1}{18} = \frac{1}{9} \implies \sigma_Y = \frac{1}{3}$$

Assuming we can apply CLT on $\bar{X}_A, \bar{X}_B$
$\implies Y$ has normal distribution.

$$P(Y > 1.0) = P\left(\frac{Y - \mu_Y}{\sigma_Y} > \frac{1.0 - 0}{\frac{1}{3}}\right)$$

$$= P(Z > 3)$$

$$= 1 - P(Z \leq 3) = 1 - 0.9987 = 0.0013$$

# Outline

➢ We covered sampling distribution of means and CLT (8.1-8.4).

➢ Sampling distribution of variance and t-distribution (8.5, 8.6): We will discuss them when relevant in estimation problems.

Today

➢ We will continue with some more examples of CLT.

➢ Next we will discuss estimation problems: Chapter 9

Probability:

Distribution $\longrightarrow$ Samples

Given distribution find probabilities of data/events.

Ex: $X$ has distribution $Binomial(X; n, p)$

Probability of $x = 3$ successes in $n = 10$ trials with $p = 0.7$?

Probability:

Distribution $\longrightarrow$ Samples

Given distribution find probabilities of data/events.

Ex: $X$ has distribution $Binomial(X; n, p)$

Probability of $x = 3$ successes in $n = 10$ trials with $p = 0.7$?

Statistics:

Sample $\longrightarrow$ Distribution

Given data find parameters/properties of distribution.

Ex: We observed $X_1 = 0, X_2 = 1, \cdots, X_{10} = 0$

What is the distribution parameter $p$, i.e. probability of heads?

# Point Estimation

**Point Estimate:**

Find single "good estimate" of a quantity of interest of a distribution/population using statistics.

**Statistics:**

Any function of the sample: average, max, max-min, ⋯

# Point Estimation

Formally, $X_1, \cdots, X_n$ iid data points from some distribution. An estimate of parameter $\theta$ of the distribution is:
$$\widehat{\Theta} = r(X_1, \cdots, X_n), \quad \text{for some appropriate function } r.$$

Note 1:  A single value of $\widehat{\Theta}$ is denoted with $\hat{\theta}$.

Note 2:  $\theta$ is considered fixed, unknown quantity. $\widehat{\Theta}$ is a random variable.

# Point Estimation

Ex: Estimate $\theta = \mu = \sum_{x} x f(x)$ of an unknown distribution.

Say, true (unknown) value $\theta = 2.5$.

Sample 4 data points $X_1, X_2, X_3, X_4$, say $3, 6, 5, -2$.

Ex: Estimate $\theta = \mu = \sum_x x\, f(x)$ of an unknown distribution.

Say, true (unknown) value $\theta = 2.5$.

Sample $4$ data points $X_1, X_2, X_3, X_4$, say $3, 6, 5, -2$.

We can try to estimate $\theta$ with any function of $X_1, \cdots, X_4$:

$\widehat{\Theta}$:  $\dfrac{X_1 + \cdots + X_n}{n}$  $\dfrac{min(X_1, \cdots, X_n) + max(X_1, \cdots, X_n)}{2}$  $X_1 \cdot X_n$

$\widehat{\theta}$ values:

Ex: Estimate $\theta = \mu = \sum_x x \, f(x)$ of an unknown distribution.

Say, true (unknown) value $\theta = 2.5$.

Sample 4 data points $X_1, X_2, X_3, X_4$, say 3,6,5, $-2$.

We can try to estimate $\theta$ with any function of $X_1, \cdots, X_4$:

$\widehat{\Theta}$:

$$\frac{X_1 + \cdots + X_n}{n} \qquad \frac{min(X_1, \cdots, X_n) + max(X_1, \cdots, X_n)}{2} \qquad X_1 \cdot X_n$$

$\widehat{\theta}$ values:

3

2

-6

A desired property of an estimator:

A statistic $\hat{\Theta}$ is said to be an **unbiased estimator** of the parameter $\theta$ if

$$\mu_{\hat{\Theta}} = E(\hat{\Theta}) = \theta.$$

A desired property of an estimator:

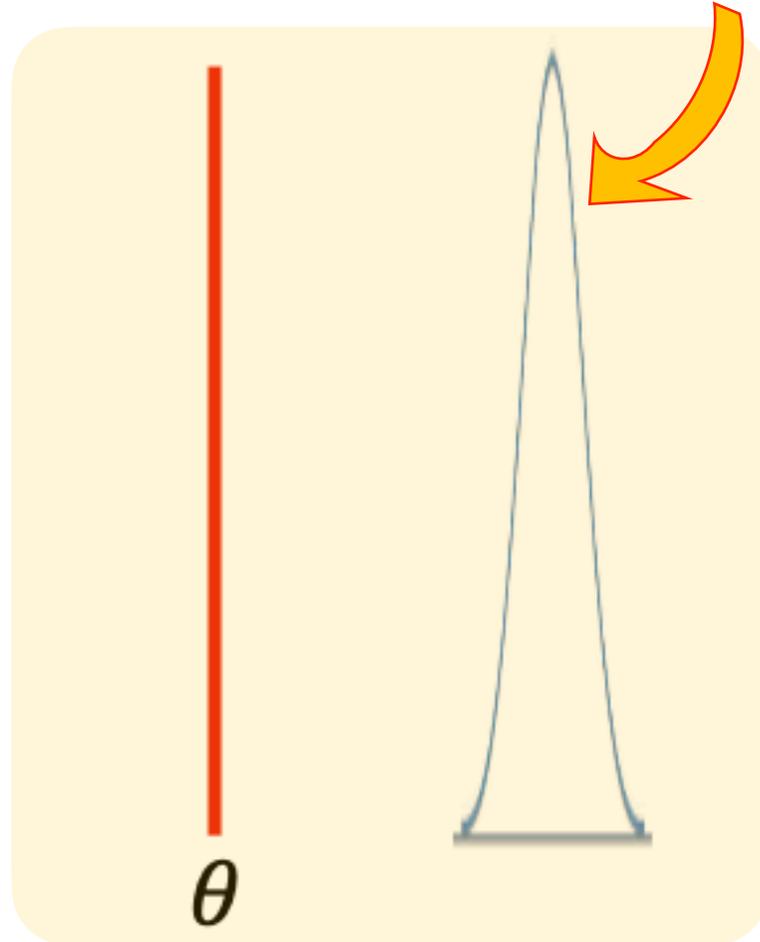A statistic $\hat{\Theta}$ is said to be an **unbiased estimator** of the parameter $\theta$ if

$$\mu_{\hat{\Theta}} = E(\hat{\Theta}) = \theta.$$

Multiple unbiased estimators:

If we consider all possible unbiased estimators of some parameter $\theta$, the one with the smallest variance is called the **most efficient estimator** of $\theta$.
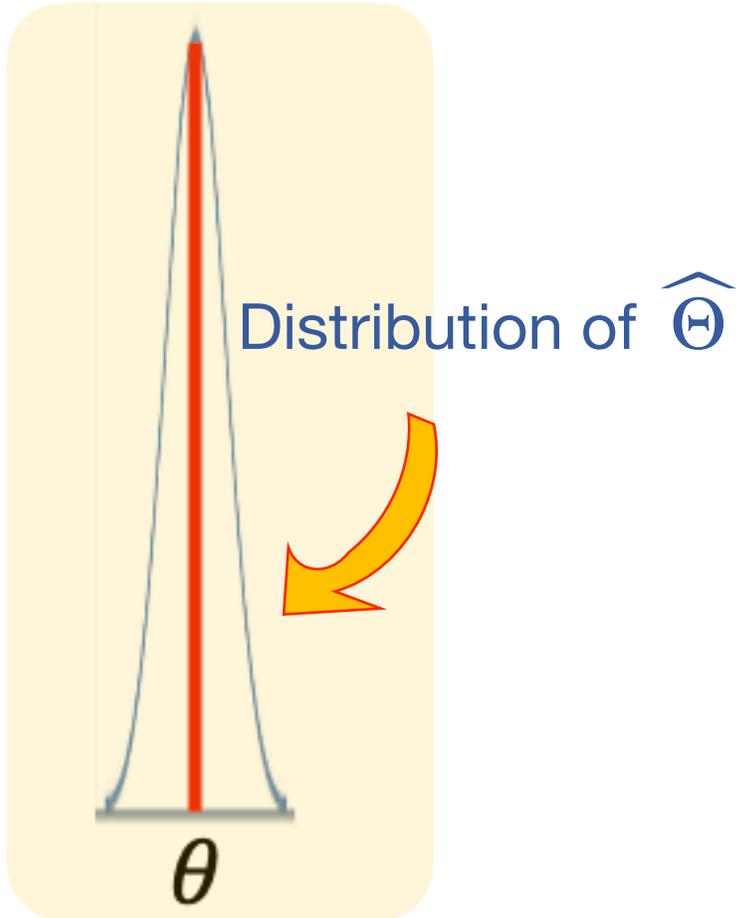
# Point Estimation

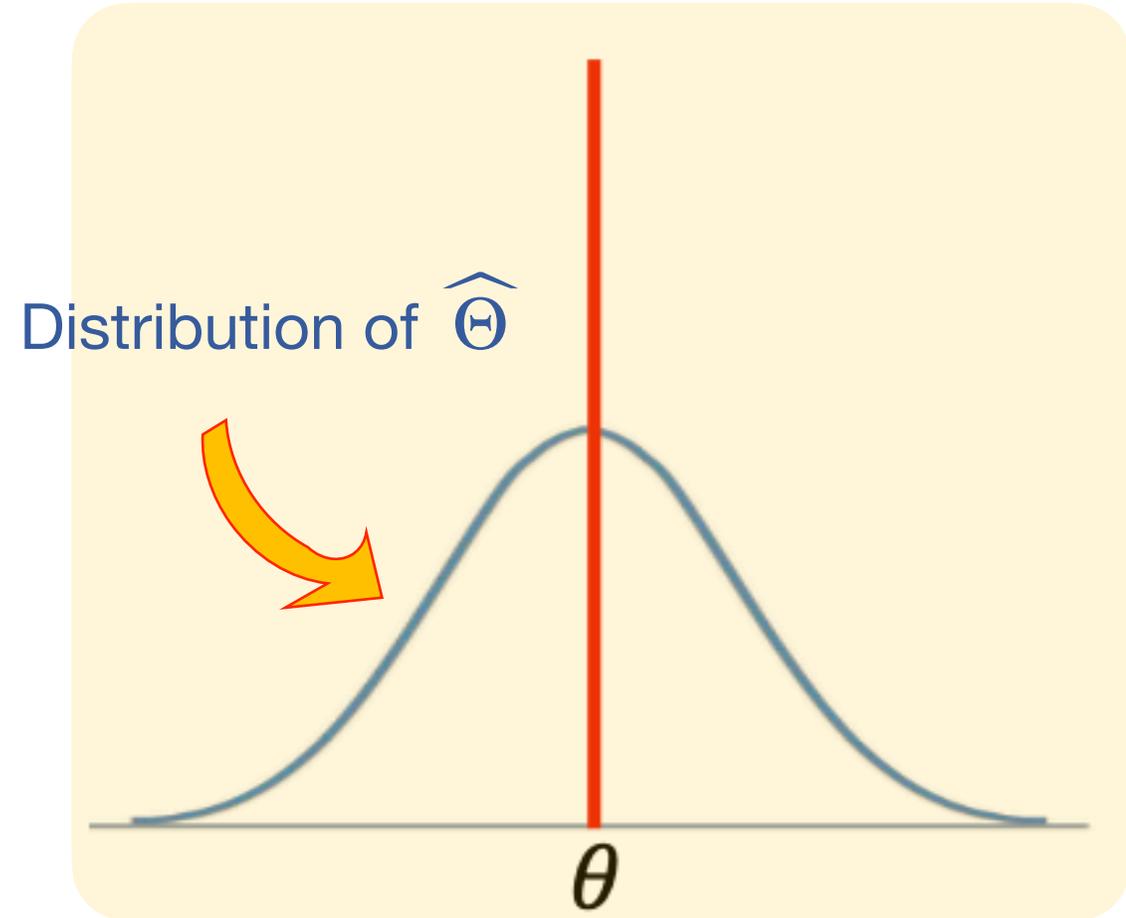Distribution of $\widehat{\Theta}$



$\theta$

Biased

# Point Estimation



Distribution of $\widehat{\Theta}$

$\theta$

$E(\hat{\theta})$

Biased

# Point Estimation



Distribution of $\widehat{\Theta}$

$\theta$

Unbiased
Small variance

Distribution of $\widehat{\Theta}$

$\theta$

Unbiased
Large variance

# Point Estimation



Distribution of $\widehat{\Theta}$

also
$E(\hat{\theta})$

$\theta$

Unbiased
Small variance

Distribution of $\widehat{\Theta}$

also
$E(\hat{\theta})$
$\theta$

Unbiased
Large variance

# Point Estimation



Distribution of $\widehat{\Theta}$

also
$E(\hat{\theta})$

$\theta$

This is
more efficient
than this

Distribution of $\widehat{\Theta}$

also
$E(\hat{\theta})$   $\theta$

Unbiased
Small variance

Unbiased
Large variance

Example: Observe $n$ coin flips $X_1, \cdots, X_n \sim Bernoulli(p)$.

True value of $p$ unknown. We want to estimate it.

Possible estimator: Sample mean $\overline{X} = \dfrac{1}{n} \sum X_i$

Is it biased? $E(\overline{X}) \overset{?}{=} p$

Example: Observe $n$ coin flips $X_1, \cdots, X_n \sim Bernoulli(p)$.

True value of $p$ unknown. We want to estimate it.

Possible estimator: Sample mean $\overline{X} = \dfrac{1}{n} \sum X_i$

Is it biased?

$$E\left(\overline{X}\right) \overset{?}{=} p$$

$$E\left(\overline{X}\right) = \frac{1}{n} \sum E(x_i) = \frac{n \cdot p}{n} = p$$

$$\Rightarrow \text{Sample mean is an unbiased estimator of } p.$$

Example: Observe $n$ coin flips $X_1, \cdots, X_n \sim Bernoulli(p)$.

True value of $p$ unknown. We want to estimate it.

Possible estimator: Sample mean $\overline{X} = \dfrac{1}{n} \sum X_i$

Is it biased?   $E(\overline{x}) \overset{?}{=} p$

$$E(\overline{x}) = \frac{1}{n} \sum E(x_i) = \frac{n \cdot p}{n} = p$$

$\Rightarrow$ Sample mean is an unbiased estimator of $p$.

(In general, sample mean is an unbiased estimator of $\mu$, since $E(\overline{X}) = \mu$, for any distribution.)