

# Review of Statistics

Ming-Ching Luoh

2022.2.

Estimation of the Population Mean

Hypothesis Tests Concerning the Population Mean

Confidence Intervals for the Population Mean

Comparing Means from Different Populations

Scatterplots and Sample Correlation

# Estimation of the Population Mean

- One natural way to estimate the **population mean**,  $\mu_Y$ , is simply to compute the **sample average**  $\bar{Y}$  from a sample of  $n$  i.i.d. observations. This can also be motivated by **law of large numbers**.

## Estimators (估計式) and Their Properties

- The sample average  $\bar{Y}$  is a natural way to estimate  $\mu_Y$ , but,  $\bar{Y}$  is **not** the only way. For example, the first observation  $Y_1$  **can be** another **estimator** of  $\mu_Y$ .
- **What** makes one estimator “**better**” than another? What are desirable characteristics of the sampling distribution of an estimator?

- In general, we want an estimator that gets as **close** as possible to the unknown true value, at least in some average sense.
- In other words, we want the **sampling distribution** of an estimator to be as **tightly** centered around the unknown value as possible.
- This leads to **three** specific desirable characteristics of an estimator: **unbiasedness, consistency, and efficiency.**

### Three desirable characteristics of an estimator.

Let  $\hat{\mu}_Y$  denote some estimator of  $\mu_Y$ ,

- Unbiasedness:  $E(\hat{\mu}_Y) = \mu_Y$ .
- Consistency:  $\hat{\mu}_Y \xrightarrow{P} \mu_Y$ .
- Efficiency.

Let  $\tilde{\mu}_Y$  be another estimator of  $\mu_Y$ , and suppose **both**  $\hat{\mu}_Y$  and  $\tilde{\mu}_Y$  are **unbiased**. Then  $\hat{\mu}_Y$  is said to be **more efficient** than  $\tilde{\mu}_Y$  if  $\text{Var}(\hat{\mu}_Y) < \text{Var}(\tilde{\mu}_Y)$ .

## Properties of $\bar{Y}$

- It can be shown that  $E(\bar{Y}) = \mu_Y$  and  $\bar{Y} \xrightarrow{P} \mu_Y$  (from law of large numbers),  $\bar{Y}$  is both unbiased and consistent.
- But, is  $\bar{Y}$  **efficient**?

## Examples of alternative estimators.

*Example 1:* The first observation  $Y_1$ ?

Since  $E(Y_1) = \mu_Y$ ,  $Y_1$  is an unbiased estimator of  $\mu_Y$ .

But,

$$\text{Var}(Y_1) = \sigma_Y^2 \geq \text{Var}(\bar{Y}) = \frac{\sigma_Y^2}{n},$$

if  $n \geq 2$ ,  $\bar{Y}$  is more **efficient** than  $Y_1$ .



*Example 2:*

$$\tilde{Y} = \frac{1}{n} \left( \frac{1}{2} Y_1 + \frac{3}{2} Y_2 + \cdots + \frac{1}{2} Y_{n-1} + \frac{3}{2} Y_n \right),$$

where  $n$  is assumed to be an even number.

The mean of  $\tilde{Y}$  is  $\mu_Y$  and its variance is

$$\text{Var}(\tilde{Y}) = \frac{1.25\sigma_Y^2}{n} > \text{Var}(\bar{Y})$$

Thus  $\tilde{Y}$  is unbiased and, because  $\text{Var}(\tilde{Y}) \rightarrow 0$  as  $n \rightarrow \infty$ ,  $\tilde{Y}$  is consistent.

However,  $\bar{Y}$  is more efficient than  $\tilde{Y}$ .

- In fact,  $\bar{Y}$  is the most efficient estimator of  $\mu_Y$  among all unbiased estimators that are **weighted averages** of  $Y_1, \dots, Y_n$ . (Weighted average implies that the estimators are all unbiased.)
- Said differently,  $\bar{Y}$  is the Best Linear Unbiased Estimator (BLUE).
- It is the most efficient (**best**) estimator among all estimators that are **unbiased** and are **linear** function of  $Y_1, \dots, Y_n$ .

$\bar{Y}$  is the least squares estimator of  $\mu_Y$ .

- The sample average  $\bar{Y}$  provides the best fit to the data in the sense that the average squared differences between the observation and  $\bar{Y}$  are the **smallest** of all possible estimators.
- The solution to the problem of minimizing

$$\sum_{i=1}^n (Y_i - m)^2$$

is  $\hat{m} = \bar{Y}$ , which is called the **least squares estimator**.

# Hypothesis Tests Concerning the Population Mean

## Null and Alternative Hypotheses

The **hypothesis testing** problem: **make** a provisional **decision, based on the evidence** at hand, whether a null hypothesis is true, or instead that some alternative hypothesis is true.

$$H_0: E(Y) = \mu_{Y,o} \text{ v.s. } H_1: E(Y) > \mu_{Y,o} \text{ (1 - sided, } > \text{)}$$

$$H_0: E(Y) = \mu_{Y,o} \text{ v.s. } H_1: E(Y) < \mu_{Y,o} \text{ (1 - sided, } < \text{)}$$

$$H_0: E(Y) = \mu_{Y,o} \text{ v.s. } H_1: E(Y) \neq \mu_{Y,o} \text{ (2 - sided)}$$

- If the null hypothesis is “accepted,” this does **not** mean that it is **true**. It is accepted **tentatively** with the recognition that it might be **rejected** later based on additional data.
- The  $p$ -value is the **probability** of drawing a statistic (e.g.  $\bar{Y}$ ) at least as **adverse** to the null as the value actually computed with your data, assuming that the null hypothesis is true.
- For the case of population mean, the  $p$ -value is the probability of drawing  $\bar{Y}$  at least **as far** in the tails of its distribution under the null hypothesis as the sample average you **actually computed**.

## Calculating the $p$ -value based on $\bar{Y}$ :

$$p\text{-value} = \Pr_{H_0} (|\bar{Y} - \mu_{Y,o}| > |\bar{Y}^{act} - \mu_{Y,o}|),$$

where  $\bar{Y}^{act}$  is the value of  $\bar{Y}$  actually observed.

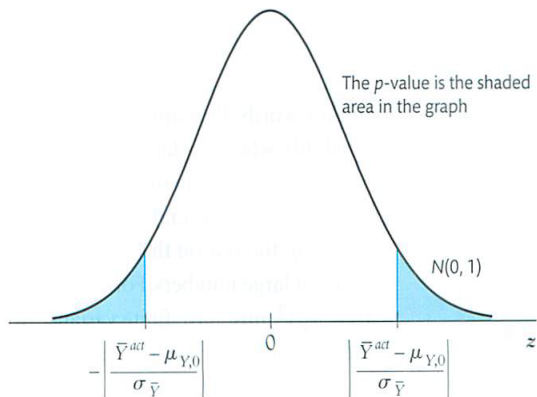
- To compute the  $p$ -value, we need to know the **sampling distribution** of  $\bar{Y}$  under the null hypothesis.
- If  $n$  is large,  $\bar{Y}$  is well approximated by a normal distribution.

$$\begin{aligned}
 p\text{-value} &= \Pr_{H_0} \left( |\bar{Y} - \mu_{Y,o}| > |\bar{Y}^{act} - \mu_{Y,o}| \right) \\
 &= \Pr_{H_0} \left( \left| \frac{\bar{Y} - \mu_{Y,o}}{\sigma_Y/\sqrt{n}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,o}}{\sigma_Y/\sqrt{n}} \right| \right) \\
 &= \Pr_{H_0} \left( \left| \frac{\bar{Y} - \mu_{Y,o}}{\sigma_{\bar{Y}}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,o}}{\sigma_{\bar{Y}}} \right| \right) \\
 &\cong \text{probability under left + right } N(0, 1) \text{ tails}
 \end{aligned}$$

where  $\sigma_{\bar{Y}}$  denotes the standard deviation of the distribution of  $\bar{Y}$ .

**FIGURE 3.1** Calculating a  $p$ -value

The  $p$ -value is the probability of drawing a value of  $\bar{Y}$  that differs from  $\mu_{Y,0}$  by at least as much as  $\bar{Y}^{act}$ . In large samples,  $\bar{Y}$  is distributed  $N(\mu_{Y,0}, \sigma_{\bar{Y}}^2)$  under the null hypothesis, so  $(\bar{Y} - \mu_{Y,0}) / \sigma_{\bar{Y}}$  is distributed  $N(0, 1)$ . Thus the  $p$ -value is the shaded standard normal tail probability outside  $\pm |(\bar{Y}^{act} - \mu_{Y,0}) / \sigma_{\bar{Y}}|$ .





## The Sample Variance, Sample Standard Deviation, and Standard Error

- In practice,  $\sigma_{\bar{Y}}$  is **unknown** and needs to be estimated.
- Estimator of the variance of Y:

$$s_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

Fact: If  $(Y_1, \dots, Y_n)$  are i.i.d. and  $E(Y^4) < \infty$ , then

$$s_Y^2 \xrightarrow{p} \sigma_Y^2$$

- Why does the law of large numbers apply?  
Because  $s_Y^2$  is a **sample average**.
- Technical note: we assume  $E(Y^4) < \infty$  because here the average is not of  $Y_i$ , but of its square.

Prove that  $s_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \xrightarrow{p} \sigma_Y^2$ .

$$\begin{aligned} \text{First, } & (Y_i - \bar{Y})^2 \\ &= [(Y_i - \mu_Y) - (\bar{Y} - \mu_Y)]^2 \\ &= (Y_i - \mu_Y)^2 - 2(Y_i - \mu_Y)(\bar{Y} - \mu_Y) + (\bar{Y} - \mu_Y)^2 \end{aligned}$$

$$\begin{aligned} s_Y^2 &= \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n (Y_i - \mu_Y)^2 - \frac{2}{n-1} \sum_{i=1}^n (Y_i - \mu_Y)(\bar{Y} - \mu_Y) \\ &\quad + \frac{1}{n-1} \sum_{i=1}^n (\bar{Y} - \mu_Y)^2 \\ &= \frac{n}{n-1} \left[ \frac{1}{n} \sum_{i=1}^n (Y_i - \mu_Y)^2 \right] - \frac{n}{n-1} (\bar{Y} - \mu_Y)^2 \end{aligned}$$

For the first term,

- Define  $W_i = (Y_i - \mu_Y)^2$ , then  $E(W_i) = \sigma_Y^2$ , and  $W_i, \dots, W_n$  are i.i.d.
- $E(W_i^2) = E[(Y_i - \mu_Y)^4] < \infty$  because  $E(Y_i^4) < \infty$ .
- Thus  $W_i, \dots, W_n$  are i.i.d. and  $\text{Var}(W_i) < \infty$ , so  $\bar{W} \xrightarrow{P} E(W_i) = \sigma_Y^2$ , and  $\frac{n}{n-1} \rightarrow 1$ .
- Therefore,  $\frac{n}{n-1} \left[ \frac{1}{n} \sum_{i=1}^n (Y_i - \mu_Y)^2 \right] = \frac{n}{n-1} \bar{W} \xrightarrow{P} \sigma_Y^2$ .

For the second term, because  $\bar{Y} \xrightarrow{P} \mu_Y$ ,  $(\bar{Y} - \mu_Y)^2 \xrightarrow{P} 0$ .

Therefore,  $s_Y^2 \xrightarrow{P} \sigma_Y^2$ .

## Computing the $p$ -value with estimated $\sigma_Y^2$ :

$$\begin{aligned}
 p\text{-value} &= \Pr_{H_0} (|\bar{Y} - \mu_{Y,o}| > |\bar{Y}^{act} - \mu_{Y,o}|) \\
 &= \Pr_{H_0} \left( \left| \frac{\bar{Y} - \mu_{Y,o}}{\sigma_Y/\sqrt{n}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,o}}{\sigma_Y/\sqrt{n}} \right| \right) \\
 &\cong \Pr_{H_0} \left( \left| \frac{\bar{Y} - \mu_{Y,o}}{s_Y/\sqrt{n}} \right| > \left| \frac{\bar{Y}^{act} - \mu_{Y,o}}{s_Y/\sqrt{n}} \right| \right) \text{ (large } n) \\
 &= \Pr_{H_0} (|t| > |t^{act}|) \\
 &\cong \text{probability under normal tails (large } n)
 \end{aligned}$$

where  $t = \frac{\bar{Y} - \mu_{Y,o}}{s_Y/\sqrt{n}}$  is the  $t$ -statistic or  **$t$ -ratio**.

## The p-value and the significance level

- Type I error: the null hypothesis (無罪) is **rejected** when in fact it is **true**. (誤判)
- Type II error: the null hypothesis (無罪) is **not rejected** when in fact it is **false**. (縱放)
- The prespecified probability of type I error is the **significance level** of the test.
- With a prespecified significance level (e.g. 5%):
  - reject if  $|t| > 1.96$ .
  - equivalently: reject if  $p \leq 0.05$ .

- The probability that the test actually **incorrectly** rejects the null hypothesis when it is true is the **size** of the test.
- The probability that the test correctly rejects the null hypothesis when the alternative is true is the **power** of the test.

## Digression: The Student $t$ -distribution

If  $Y$  is distributed  $N(\mu_Y, \sigma_Y^2)$ , then the  $t$ -statistic has the Student  $t$ -distribution (tabulated in back of all stats books)

Some comments:

- For  $n > 30$ , the  $t$ -distribution and  $N(0, 1)$  are very close.
- The assumption that  $Y$  is distributed  $N(\mu_Y, \sigma_Y^2)$  is rarely plausible in practice (income? number of children?)
- The  $t$ -distribution is an historical artifact from days when sample sizes were very small.
- In this class, we won't use the  $t$  distribution - we rely solely on the large- $n$  approximation given by the Central Limit Theorem.



# Confidence Intervals for the Population Mean

- Because of random sampling error, it is impossible to learn the **exact** value of the population mean of  $Y$  using only the information in a sample.
- It is possible to use data from a random sample to **construct** a set of values that **contains** the true population mean  $\mu_Y$  with a certain **prespecified probability**.

- A 95% **confidence interval** for  $\mu_Y$  is an interval that **contains the** true value of  $Y$  in **95% of repeated samples**.
- *Digression:* **What is random here?**  
the confidence interval— it will differ from one sample to the next; the population parameter,  $\mu_Y$ , is not random.

A 95% confidence interval can always be constructed as the set of values of  $\mu_Y$  not rejected by a hypothesis test with a 5% significance level.

$$\begin{aligned} & \left\{ \mu_Y \mid \left| \frac{\bar{Y} - \mu_Y}{s_Y / \sqrt{n}} \right| \leq 1.96 \right\} \\ &= \left\{ \mu_Y \mid -1.96 \leq \frac{\bar{Y} - \mu_Y}{s_Y / \sqrt{n}} \leq 1.96 \right\} \\ &= \left\{ \mu_Y \mid -1.96 \frac{s_Y}{\sqrt{n}} \leq \bar{Y} - \mu_Y \leq 1.96 \frac{s_Y}{\sqrt{n}} \right\} \\ &= \left\{ \mu_Y \in \left( \bar{Y} - 1.96 \frac{s_Y}{\sqrt{n}}, \bar{Y} + 1.96 \frac{s_Y}{\sqrt{n}} \right) \right\} \end{aligned}$$

**Summary:** From the assumptions of:

- (1) simple random sampling of a population, that is,  $\{Y_i, i = 1, \dots, n\}$  are i.i.d.
- (2)  $0 < E(Y^4) < \infty$ .

we developed, for large samples (large  $n$ ):

- Theory of estimation (sampling distribution of  $\bar{Y}$ )
- Theory of hypothesis testing (large- $n$  distribution of  $t$ -statistic and computation of the  $p$ -value).
- Theory of confidence intervals (constructed by inverting test statistic).

Are assumptions (1) & (2) plausible in practice? Yes

# Comparing Means from Different Populations

Let  $\mu_w$  be the mean hourly earning in the population of women recently graduated from college and let  $\mu_m$  be population mean for recently graduated men. Consider the null hypothesis that earnings for these two populations differ by certain amount  $d$ , then

$$H_0: \mu_m - \mu_w = d \text{ v.s. } H_1: \mu_m - \mu_w \neq d.$$

Since  $\bar{Y}_m \sim N(\mu_m, \frac{\sigma_m^2}{n_m})$  and  $\bar{Y}_w \sim N(\mu_w, \frac{\sigma_w^2}{n_w})$ , then

$$\bar{Y}_m - \bar{Y}_w \sim N(\mu_m - \mu_w, \frac{\sigma_m^2}{n_m} + \frac{\sigma_w^2}{n_w})$$

Replace population variances by sample variances, we have the standard error ( $SE$ )

$$SE(\bar{Y}_m - \bar{Y}_w) = \sqrt{\frac{s_m^2}{n_m} + \frac{s_w^2}{n_w}}$$

and the  $t$ -statistic is

$$t = \frac{\bar{Y}_m - \bar{Y}_w - d}{SE(\bar{Y}_m - \bar{Y}_w)}$$

If both  $n_m$  and  $n_w$  are large, the  $t$ -statistic has a standard normal distribution.

**TABLE 3.1** Trends in Hourly Earnings in the United States of Working College Graduates, Ages 25–34, 1992 to 2012, in 2012 Dollars

Year	Men			Women			Difference, Men vs. Women		
	$\bar{Y}_m$	$s_m$	$n_m$	$\bar{Y}_w$	$s_w$	$n_w$	$\bar{Y}_m - \bar{Y}_w$	$SE(\bar{Y}_m - \bar{Y}_w)$	95% Confidence Interval for $d$
1992	24.83	10.85	1594	21.39	8.39	1368	3.44**	0.35	2.75–4.14
1996	23.97	10.79	1380	20.26	8.48	1230	3.71**	0.38	2.97–4.46
2000	26.55	12.38	1303	22.13	9.98	1181	4.42**	0.45	3.54–5.30
2004	26.80	12.81	1894	22.43	9.99	1735	4.37**	0.38	3.63–5.12
2008	26.63	12.57	1839	22.26	10.30	1871	4.36**	0.38	3.62–5.10
2012	25.30	12.09	2004	21.50	9.99	1951	3.80**	0.35	3.11–4.49

These estimates are computed using data on all full-time workers ages 25–34 surveyed in the Current Population Survey conducted in March of the next year (for example, the data for 2012 were collected in March 2013). The difference is significantly different from zero at the \*\*1% significance level.

## Another Example.

**TABLE 3.1** Differences in Household Income According to Childhood Socioeconomic Circumstances, Grouped by Level of Highest Qualification

Qualification	Father's NS-SEC = Higher			Father's NS-SEC = Routine			Difference, Higher vs. Routine			
	$Y_h$	$s_h$	$n_h$	$Y_r$	$s_r$	$n_r$	$Y_h - Y_r$	$SE(Y_h - Y_r)$	95% Confidence Interval for $d$	
None	£2,223.13	£2,115.12	1129	£1,842.98	£1,487.29	6383	£380.15	£65.64	£251.38	£508.93
GCSE/O-Level	£2,837.18	£1,819.73	1962	£2,596.93	£1,738.47	4042	£240.25	£49.35	£143.49	£337.00
A-Level	£3,045.99	£2,451.81	1216	£2,745.70	£1,912.50	1169	£300.30	£89.85	£124.11	£476.49
Undergraduate degree or more	£3,690.51	£2,743.55	4359	£3,370.96	£2,443.58	2505	£319.55	£64.11	£193.86	£445.23
All categories	£3,215.71	£2,497.73	8666	£2405.45	£1,886.86	14099	£810.25	£31.18	£749.13	£871.38

Source: Understanding Society.



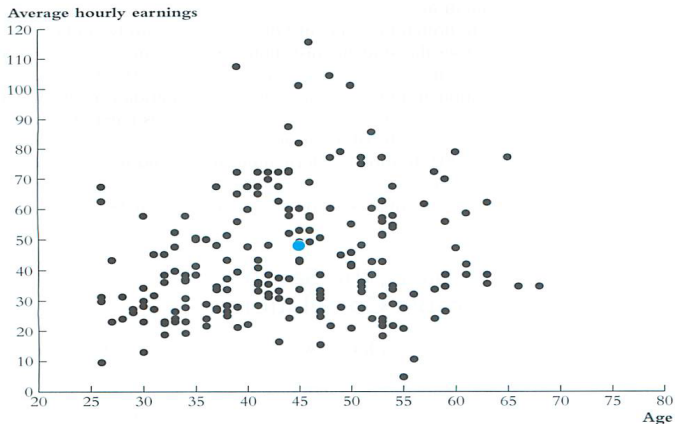
# Scatterplots, the Sample Covariance, and the Sample Correlation

Three ways to summarize the relationship between two variables

- scatterplot,
- sample covariance,
- sample correlation coefficient.

# Scatterplots

**FIGURE 3.2** Scatterplot of Average Hourly Earnings vs. Age



Each point in the plot represents the age and average earnings of one of the 200 workers in the sample. The highlighted dot corresponds to a 45-year-old worker who earns \$49.15 per hour. The data are for computer and information systems managers from the March 2016 CPS.

## Sample Covariance and Correlation

- The population covariance and correlation can be estimated by the **sample covariance** and **sample correlation**.
- The **sample covariance** is

$$s_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

- The **sample correlation** is

$$r_{XY} = \frac{s_{XY}}{s_X s_Y}, |r_{XY}| \leq 1$$

- It can be shown that under the assumptions that  $(X_i, Y_i)$  are i.i.d. and that  $X_i$  and  $Y_i$  have finite **fourth** moments,

$$s_Y^2 \xrightarrow{p} \sigma_Y^2$$

$$s_{XY} \xrightarrow{p} \sigma_{XY}$$

$$r_{XY} \xrightarrow{p} \text{Corr}(X, Y)$$

**Prove that**  $\underline{s_{XY}} \xrightarrow{p} \sigma_{XY}$ .

$$\begin{aligned}
 & s_{XY} \\
 = & \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \\
 = & \frac{1}{n-1} \sum_{i=1}^n [(X_i - \mu_X) - (\bar{X} - \mu_X)] [(Y_i - \mu_Y) - (\bar{Y} - \mu_Y)] \\
 = & \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu_X)(Y_i - \mu_Y) - \frac{1}{n-1} \sum_{i=1}^n (\bar{X} - \mu_X)(Y_i - \mu_Y) \\
 & - \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu_X)(\bar{Y} - \mu_Y) + \frac{1}{n-1} \sum_{i=1}^n (\bar{X} - \mu_X)(\bar{Y} - \mu_Y)
 \end{aligned}$$

- Use the fact that  $\sum_{i=1}^n (Y_i - \mu_Y) = n(\bar{Y} - \mu_Y)$ ,  $\sum_{i=1}^n (X_i - \mu_X) = n(\bar{X} - \mu_X)$  and collect terms, we have

$$s_{XY} = \left(\frac{n}{n-1}\right) \frac{1}{n} \sum_{i=1}^n (X_i - \mu_X)(Y_i - \mu_Y) - \left(\frac{n}{n-1}\right) (\bar{X} - \mu_X)(\bar{Y} - \mu_Y)$$

- It is easy to see that the second term converges in probability to zero because  $\bar{X} \xrightarrow{p} \mu_X$  and  $\bar{Y} \xrightarrow{p} \mu_Y$  so  $(\bar{X} - \mu_X)(\bar{Y} - \mu_Y) \xrightarrow{p} 0$  by Slutsky's theorem.

- By the definition of covariance, we have  $E((X_i - \mu_X)(Y_i - \mu_Y)) = \sigma_{XY}$ . To apply the law of large numbers on the first term, we need to have

$$\text{Var}((X_i - \mu_X)(Y_i - \mu_Y)) < \infty$$

which is satisfied since

$$\begin{aligned} & \text{Var}((X_i - \mu_X)(Y_i - \mu_Y)) \\ &= E((X_i - \mu_X)^2(Y_i - \mu_Y)^2) \\ &\leq \sqrt{E(X_i - \mu_X)^4 E(Y_i - \mu_Y)^4} < \infty \end{aligned}$$

The second inequality follows by applying the Cauchy-Schwartz inequality, and the last inequality follows because of the finite fourth moments for  $(X_i, Y_i)$ .

- *The Cauchy-Schwartz inequality is*

$$|E(XY)| \leq \sqrt{E(X^2)E(Y^2)}$$

- Applying the law of large numbers, we have

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n (X_i - \mu_X)(Y_i - \mu_Y) \\ & \xrightarrow{p} E((X - \mu_X)(Y - \mu_Y)) = \sigma_{XY} \end{aligned}$$

- Also,  $\frac{n}{n-1} \rightarrow 1$ , therefore

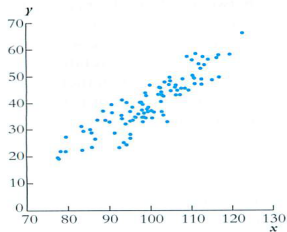
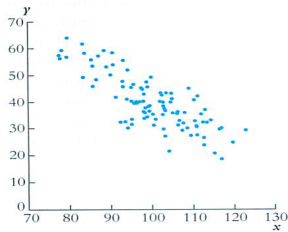
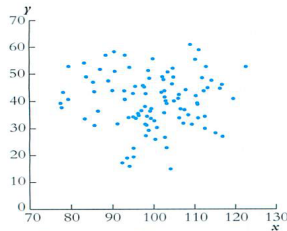
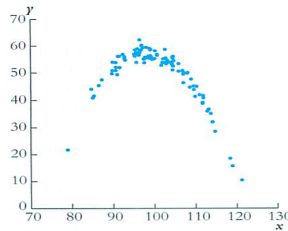
$$s_{XY} \xrightarrow{p} \sigma_{XY}$$



# Scatterplots and Sample Correlation

**FIGURE 3.3** Scatterplots for Four Hypothetical Data Sets

The scatterplots in Figures 3.3a and 3.3b show strong linear relationships between  $X$  and  $Y$ . In Figure 3.3c,  $X$  is independent of  $Y$  and the two variables are uncorrelated. In Figure 3.3d, the two variables also are uncorrelated even though they are related nonlinearly.

**(a)** Correlation = +0.9**(b)** Correlation = -0.8**(c)** Correlation = 0.0**(d)** Correlation = 0.0 (quadratic)