Chapter 5

Estimation, testing hypothesis and regression analysis

Professor Jung Jin Lee Soongsil University, Korea New Uzbekistan University, Uzbekistan

Chapter 5 Estimation, testing hypothesis and regression analysis

- 5.1 Sampling distribution and estimation 5.1.1 Sampling distribution of sample means 5.1.2 Estimation of a population mean 5.2 Testing hypothesis for a population mean 5.3 Testing hypothesis for two populations means 5.4 Testing hypothesis for several population means: Analysis of variance 5.5 Regression analysis 5.5.1 Correlation analysis 5.5.2 Simple linear regression
 - 5.5.3 Multiple linear regression

- Inferential statistics is used to find out characteristics of unknown populations.
- Characteristics of a population usually refers to the population mean, variance, etc.,
- Characteristic values of a population are called parameters.
- The parameters are estimated using sample characteristics, such as a sample mean and sample variance.



Central Limit Theorem(CLT)



Central limit theorem

If a population has an infinite elements with a mean μ and variance σ^2 , then, if the sample size is large enough, the distribution of all possible sample means is an approximately normal distribution $N(\mu, \frac{\sigma^2}{n})$. We can summarize specifically the central limit theorem as follows.

1) The average of all possible sample means, $\mu_{\overline{x}}$, is equal to the population mean μ .

(i.e., $\mu_{\overline{X}} = \mu$)

2) The variance of all possible sample means, $\sigma_{\overline{x}}^2$, is the population variance divided by n.

 $(i.e., \sigma_{\overline{X}}^2 = \frac{\sigma^2}{n})$

3) The distribution of all possible sample means is approximately a normal distribution.

The above facts can be briefly written as $\overline{X} \sim N(\mu, \frac{\sigma^2}{n})$.

Estimation

- A. Point estimation of population mean
 - An observed value of the sample mean is a point estimate of the population mean.

B. Interval estimation of population mean

^{IIII} 100(1-α)% Confidence Interval for Population Mean μ --- Population is normal and population variance σ^2 is known

$$\left[\overline{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \ , \ \overline{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right]$$



^{IIII} 100(1- α)% Confidence Interval for Population Mean μ ---- Population is normal and population variance σ^2 is unknown

$$\left[\overline{X} - t_{n-1;\,\alpha/2} \cdot \frac{S}{\sqrt{n}} , \overline{X} + t_{n-1;\,\alpha/2} \cdot \frac{S}{\sqrt{n}}\right]$$

n is the sample size and S is the sample standard deviation.



Example 5.1.2 Suppose we do not know the population variance in Example 4.4.2. If the sample size is 25 and the sample standard deviation is 5 (unit: 10,000 KRW), estimate the mean of the starting salary of college graduates at the 95% confidence level.

Answer

Since we do not know the population variance, we should use the t distribution for interval estimation of the population mean. Since $t_{n-1: \alpha/2} = t_{25-1: 0.05/2} = t_{25-1: 0.025} = 2.0639$, the 95% confidence interval of the population mean is as follows.

$$\begin{bmatrix} \overline{X} - t_{n-1: \alpha/2} \frac{S}{\sqrt{n}}, \overline{X} + t_{n-1: \alpha/2} \frac{S}{\sqrt{n}} \end{bmatrix}$$

$$\Leftrightarrow \begin{bmatrix} 275 - 2.0639(5/5), 275 + 2.0639(5/5) \end{bmatrix}$$

$$\Leftrightarrow \begin{bmatrix} 272.9361, 277.0639 \end{bmatrix}$$

Note that the smaller the sample size, the wider the interval width.

- Examples of testing hypothesis for a population mean.
 - Capacity of a cookie bag is indicated as 200g. Will there be enough cookies in the indicated capacity?
 - At a light bulb factory, a newly developed light bulb advertises a longer bulb life than the past. Is this propaganda reliable?
 - In this year's academic test, students said that there will be an average English score of 5 points higher than last year. How can you investigate if this is true?

[Example 5.2.1] At a light bulb factory, the average life expectancy of a light bulb made by a conventional method is known to be 1,500 hours, and the standard deviation is 200 hours. The company introduced a new production method with an average life expectancy of 1,600 hours. To confirm this argument, 30 samples were taken, and the sample mean was 1555 hours. Does the new type of light bulb have a life of 1600 hours?

<Answer>

• Make two assumptions about the different arguments for the population mean μ .

 $H_0: \mu = 1500, \quad H_1: \mu = 1600$

• H_0 is a null hypothesis and H_1 is an alternative hypothesis

- Unless there is a significant reason, keep the null hypothesis
- ⇒ 'conservative decision making'
- Testing hypothesis is based on sampling distribution of \overline{X} .
- Select a critical value C based on sampling distribution
- ⇒ Decision rule:

'If \overline{X} < C, then accept H_0 , else reject H_0 '



Table 7.1.1 Two types of errors in testing hypothesis

	Act	tual
	H_0 is true	H_1 is true
Decision: H_0 is true H_1 is true	Correct Type 1 Error	Type 2 Error Correct

• If we set the significance level is 5%, C can be calculated by finding the percentile of $N(1500, \frac{200^2}{30})$

C = 1500 + 1.645
$$\sqrt{\frac{200^2}{30}}$$
 = 1560.06

- Decision rule: 'If \overline{X} < 1560.06, then accept H_0 , else reject H_0 '
- Since $\overline{X} = 1555$, it is less than 1560.06, accept H_0 .

- p-value is the probability of a type 1 error when the observed sample mean value is considered as the critical value for decision
 - ⇒ p-value indicates where the observed sample mean is located among all possible sample means

Table 5.2.2 Testing hypothesis for a population mean - unknown σ case

Decision Rule

Type of Hypothesis

- 1) $H_0: \mu = \mu_0$ $H_1: \mu > \mu_0$ If $rac{\overline{X} - \mu_0}{rac{S}{\sqrt{n}}} > t_{n-1: \ lpha}$, then reject H_0
- 2) $H_0: \mu = \mu_0$ $H_1: \mu < \mu_0$ If $rac{\overline{X}-\mu_0}{\frac{S}{\sqrt{n}}} < -t_{n-1: \ lpha}$, then reject H_0
- 3) $egin{array}{ll} H_0:\mu=\mu_0 \ H_1:\mu
 eq\mu_0 \end{array}$ If $\left|rac{\overline{X}-\mu_0}{rac{S}{\sqrt{n}}}
 ight|>t_{n-1;\;lpha/2}$, then reject H_0

Note: Assume that the population is a normal distribution. The H_0 of 1) can be written as $H_0:\mu\leq \mu_0$, 2) as $H_0:\mu\geq \mu_0$

Example 5.2.2 The weight of a bag of cookies is supposed to be 250 grams. Suppose the weight of all bags of cookies follows a normal distribution. In the survey of 16 random samples of bags, the sample mean was 253 grams, and the sample standard deviation was 10 grams. Test the hypothesis whether the weight of the bag of cookies is 250g or larger using $\alpha = 1\%$ and find the *p*-value. Use ^reStatU_J to test the hypothesis above.

Answer

Since the population standard deviation is unknown and the sample

$$ext{'If} \; rac{\overline{X}-\mu_0}{rac{S}{\sqrt{n}}} > t_{n-1:\;lpha}, ext{ then reject } H_0 ext{ else accept } H_0'$$

 $ext{'If } \; rac{253-250}{rac{10}{\sqrt{16}}} > t_{16:\;0.01}, ext{ then reject } H_0 ext{ else accept } H_0'$



Since the value of test statistic is $\frac{253-250}{\sqrt{16}} = 1.2$, and $t_{15:\ 0.01} = 2.602$, we accept H_0 . Note that the decision rule can be written as follows.

'If
$$\overline{X} > 250 + 2.602 rac{10}{\sqrt{16}}, ext{ then reject } H_0 ext{ else accept } H_0'$$

Test statistic and sampling distribution

$$rac{\overline{x}_1-\overline{x}_2)-D_0}{\sqrt{rac{s_p^2}{p}+rac{s_p^2}{p}}} \qquad ext{where } s_p^2 = rac{(n_1-1)s_1^2+(n_2-1)s_2^2}{n_1+n_2-2}$$

 n_1

 n_2

$$\sim t_{n_1+n_2-2}$$

Table 5.3.1 Testing hypothesis of two populations means

Type of Hypothesis	Decision Rule
1) $egin{array}{ll} H_0: \mu_1-\mu_2=D_0\ H_1: \mu_1-\mu_2>D_0 \end{array}$	If $rac{(\overline{x}_1-\overline{x}_2)-D_0}{\sqrt{rac{s_p^2}{n_1}+rac{s_p^2}{n_2}}}>t_{n_1+n_2-2;\;lpha}$, then reject H_0 , else accept H_0
2) $egin{array}{ll} H_0: \mu_1-\mu_2=D_0\ H_1: \mu_1-\mu_2 < D_0 \end{array}$	If $rac{(\overline{x}_1-\overline{x}_2)-D_0}{\sqrt{rac{s_p^2}{n_1}+rac{s_p^2}{n_2}}} < -t_{n_1+n_2-2;\;lpha}$, then reject H_0 , else accept H_0
3) $H_0: \mu_1 - \mu_2 = D_0$ $H_1: \mu_1 - \mu_2 eq D_0$	If $\left rac{(ar{x}_1-ar{x}_2)-D_0}{\sqrt{rac{s_p^2}{n_1}+rac{s_p^2}{n_2}}} ight >t_{n_1+n_2-2;\;lpha/2},$ then reject H_0 , else accept H_0

Example 5.3.2 (Monthly wages by male and female) Random samples of 10 male and female college graduates this year showed their monthly wages as follows. (Unit 10,000 KRW)

Male 272 255 278 282 296 312 356 296 302 312 Female 276 280 369 285 303 317 290 250 313 307 Ex ⇔ DataScience ⇔ WageByGender.csv.

Using ^reStat₁, answer the following questions.

- 1) If population variances are assumed to be the same, test the hypothesis at the 5% significance level of whether the average monthly wage for males and females is the same.
- 2) If population variances are assumed to be different, test the hypothesis at the 5% significance level of whether the average monthly wage for males and females is the same.

H₀: µ₁ - µ₂ = 0.00 , H₁: µ₁ - µ₂ ≠ 0.00



	-							
File	E	EX080103	_WageB	yGender.csv				
Analysis Var by Group								
2: 1	2: Income I: Gender							
(Selected data: Raw Data) (or Paired Var)								
Selec	tedvar	2 Dy V1,						
	Gender	Income	V3	V4				
1	Μ	272						
2	Μ	255						
3	Μ	278						
4	Μ	282						
5	Μ	296						
6	Μ	312						
7	Μ	356						
8	Μ	296						
9	М	302						
10	Μ	312						
11	F	276						
12	F	280						
13	F	369						
14	F	285						
15	F	303						
16	F	317						
17	F	290						
18	F	250						
19	F	313						
20	F	307						

p-value = 0.8302

(Group Gender) Income Confidence Interval Graph



	240 260 280 300 320 340 Income	9 360 380					
(Group	Gender) Income Testing Hypothesis: Two Population N He: $\mu_1 - \mu_2 = D$, $H_1: \mu_1 - \mu_2 \neq D$, $D = 0.00$ [TestStat] = ($R_1 - R_2 - D$) / (pooledStd * $\sqrt{(1/n_1+1/n_2)}$ ~ t(18) Distribution	Testing Hypothesis: Two Population Means	Analysis Var	Income	Group Name	Gender	
0.45 -		Statistics	Observation	Mean	Std Dev	std err	Population Mean 95% Confidence Interval
0.35 -		1 (F)	10	299.000	31.742	10.038	(276.293, 321.707)
0.30 -		2 (M)	10	296.100	27.739	8.772	(276.257, 315.943)
0.25 -		Total	20	297.550	29.051	6.496	(283.954, 311.146)
0.20 -		Missing Observations	0				
0.15 -		Hypothesis	Variance Assumption	$\sigma_1^2 = \sigma_2^2$			
0.05 -	0.025 -4 -3 -2 -1 0 1 2 3 4	H ₀ : µ ₁ - µ ₂ = D	D	[TestStat]	t value	p-value	μ ₁ -μ ₂ 95% Confidence Interval
	Reject Ho -> -2.101 <- Accept Ho -> 2.101 <- Reject Ho [TestStat] = 0.218	H ₁ : µ ₁ - µ ₂ ≠ D	0.00	Difference of Sample	0.218	0.8302	(-25.106, 30.906)

Means

- Examples to compare means of several populations.
 - Are average hours of library usage for each grade the same?
 - Are yields of three different rice seeds equal?
 - In a chemical reaction, are response rates the same at four different temperatures?
 - Are average monthly wages of college graduates the same in three different cities?
- A factor is a variable that distinguishes populations, such as grade or rice.

[Example 5.4.1] To compare the English proficiency of each grade at a university, samples were randomly selected from each grade to take the same English test.

Grade	English Proficiency Score	Average
1	81 75 69 90 72 83	$\overline{y}_{1.} = 78.3$
2	65 80 73 79 81 69	$\overline{y}_{2.} = 74.5$
3	72 67 62 76 80	$\overline{y}_{3.} = 71.4$
4	89 94 79 88	$\overline{y}_{4\cdot}$ = 87.5

- 1) Using [[]eStat], draw a dot graph of exam scores for each grade and compare average.
- 2) We want to test a hypothesis whether the average scores of each grade are the same or not. Write a null hypothesis and an alternative hypothesis.
- 3) Apply the analysis of variances to test the hypothesis in question 2).
- 4) Use **"eStat_** to check the results of the ANOVA test.

<Answer of Example 5.4.1>

File		Ex911EnglishScoreByGrade.csv									
Anal	ysis Var			b	y Group						
2: 5	Score			• 1:	Grade						
(Se	elected dat	ta:	Raw Data)	(Select	: up to two g	roups)					
Selec	tedVar	V	2 by V1,								
	Grade		Score	V3	V4	٧					
1		1	81								
2		1	75								
3		1	69								
4		1	90								
5		1	72								
6		1	83								
7		2	65								
8		2	80								
9		2	73								
10		2	79								
11		2	81								
12		2	69								
13		3	72								
14		3	67								
15		3	62								
16		3	76								
17		3	80								
18		4	89								
19		4	94								
20		4	79								
21		4	88								



Confidence Interval Graph Histogram

$H_o: \mu_1 = \mu_2 = \dots = \mu_k$	H_l : At least one pair of	of means is different			
Significance Level $\alpha = \odot$ 5% \bigcirc 1% Confidence Level \odot 95% \bigcirc 99					
ANOVA F test Stan	dardized Residual Plot	Kruskal-Wallis Test			



<Answer of Example 5.4.1>

- 2) Null hypothesis $H_o: \mu_1 = \mu_2 = \mu_3 = \mu_4$ Alternative hpothesis $H_1:$ at least one pair of μ_i is not the same
- 3) Between sum of squares (SSB) or Treatment sum of squares (SSTr)

SSTr = $6(78.3 - \bar{y}_{..})^2 + 6(74.5 - \bar{y}_{..})^2 + 5(71.4 - \bar{y}_{..})^2 + 4(87.5 - \bar{y}_{..})^2 = 643.633$ \Rightarrow If SSTr is close to zero, all sample means for four grades are similar.

Within sum of squares (SSW) or Error sum of squares (SSE)

$$SSE = (81 - \bar{y}_{1.})^{2} + (75 - \bar{y}_{1.})^{2} + \dots + (83 - \bar{y}_{1.})^{2} + (65 - \bar{y}_{2.})^{2} + (80 - \bar{y}_{2.})^{2} + \dots + (69 - \bar{y}_{2.})^{2} + (72 - \bar{y}_{3.})^{2} + (67 - \bar{y}_{3.})^{2} + \dots + (80 - \bar{y}_{3.})^{2} + (89 - \bar{y}_{4.})^{2} + (94 - \bar{y}_{4.})^{2} + \dots + (88 - \bar{y}_{4.})^{2} = 839.033$$

<Answer of Example 9.1.1>

$$F_{0} = \frac{\frac{3317}{(4-1)}}{\frac{55E}{(21-4)}} = \frac{Treatment Mean Square (MSTr)}{Error Mean Square (MSE)} \sim F_{3,17}$$

$$F_0 = \frac{\frac{643.633}{3}}{\frac{839.033}{17}} = 4.347$$
 $F_{3,17;0.05} = 3.20$

• Hence Reject
$$H_o: \mu_1 = \mu_2 = \mu_3 = \mu_4$$

ANOVA Table

Factor	Sum of Squares	Degree of freedom	Mean Squares	F value
Treatment Error	<u>SSTr</u> = 643.633 <u>SSE</u> = 839.033	4-1 21-4	<u>MSTr</u> = 643.633/3 <u>MSE</u> = 839.033/17	<u>Fo</u> = 4.347
Total	<u>SST</u> =1482.666	20		

<Answer of Example 5.4.1>



Statistics	Analysis Var	Score	G	roup lame	Gra	de						
Group Variable (Grade)	Observation	Mean	Mean Std		std	err	Populati 95% Co Inte	on Mean nfidence erval	955	Population Variance % Confidence Interval		
1 (Group 1)	6	78.333		7.789		3.180	(70.159	86.507)	(23	.638, 364.929)		
2 (Group 2)	6	74.500		6.565		2.680	(67.610	81.390)	(16	.793, 259.260)		
3 (Group 3)	5	71.400		7.127		3.187	(62.550	80.250)	(18	.235, 419.472)		
4 (Group 4)	4	87.500		6.245	3.122 (7		(77.563, 97.437)		3.122 (77.563, 97.437)		(12	.516, 542.181)
Total	21	77.333		8.610	1.879		(73.414, 81.253)		(43	.391, 154.593)		
Missing Observations	0											
Analysis of Variance												
Factor	Sum of Squares	deg o freedo	of m	Mean S	quares	F	value	p valu	е			
Treatment	643.6	33	3	ĩ	214.544		4.347	0.	0191			
Error	839.0	33	17		49.355							
Total	1482.6	67	20									

<Answer of Example 5.4.1>

Testing Hypothesis ANOVA

[Hypothesis] $H_o: \mu_1 = \mu_2 = \dots = \mu_k$

 H_1 : At least one pair of means is different

[Test Type] F test (ANOVA)

Significance Level $\alpha = \bigcirc 5\% \bigcirc 1\%$

[Sample Data] Input either sample data using BSV or sample statistics at the next boxes

Sample 1	81	75	69	90	72	83					
Sample 2	65	80	73	79	81	69					
Sample 3	72	67	62	76	80						
Sample 4	89	94	79	88							
[Sample Sta	atisti	cs]									
$n_1 = [$	(6		n ₂ =	=	6	$n_3 =$	5] /	1 ₄ =	4
$\bar{x}_I = ($	78	.33		$\bar{x}_{2} =$	=	74.50	$\bar{x}_3 =$	71.40]	$\bar{c}_4 =$	87.50
$s_I^2 = [$	60	.67		s_2^{2}	=	43.10	$s_3^2 =$	50.80]	${a_4}^2 =$	39.00

Execute

Menu

Table 9.1.2 Notation of one-way ANOVA											
Factor	Obser	Observed values of sample Average									
Level 1	Y ₁₁	Y_{12}		$Y_{1n_{1}}$	$\overline{Y}_{1.}$						
Level 2	Y_{21}	Y_{22}		Y_{2n_2}	\overline{Y}_2 .						
Level k	Y_{k1}	Y_{k2}		$Y_{k\!n_k}$	\overline{Y}_{k} .						

ANOVA Model

$$\begin{aligned} \mathcal{X}_{ij} &= \mu_i + \varepsilon_{ij} \\ &= \mu + \alpha_i + \varepsilon_{ii}, \ i = 1, 2, \dots, k; j = 1, 2, \dots, n_i \end{aligned}$$

Hypothesis

 $H_0: \alpha_1 = \alpha_2 = \cdots = \alpha_k = 0$ $H_1:$ At least one pair of α_i is not equal to 0

	Table 9.1	.3 Analysis	of variance table of one-	way ANOVA
Factor	Sum of Squares	Degree of freedom	Mean Squares	F value
Treatment	SSTr	k-1	MSTr = SSTr / (k-1)) $F_0 = MSTr/MSE$
Error	SSE	n-k	MSE = SSE / (n - k)	(;)
Total	<u>SST</u>	n-1	$(n = \sum_{i=1}^{k} n_i)$	
 Total S 	um of Squar	es	SST =	$\sum_{i=1}^{k} \sum_{i=1}^{n_i} (Y_{ii} - \overline{Y}_{})^2$

SSTr = $\sum_{i=1}^{k} \sum_{j=1}^{n_i} (\overline{Y}_{i \cdot} - \overline{Y}_{\cdot \cdot})^2$

SSE = $\sum_{i=1}^{k} \sum_{j=1}^{n_i} (Y_{ij} - \overline{Y}_{i})^2$

- **Treatment Sum of Squares**
- **Error Sum of Squares**
- SST = SSTr + SSE
- If $F_0 > F_{k-1,n-k;\alpha}$, then reject H_0

Correlation analysis

- Population with (μ_X, μ_Y) and (σ_X^2, σ_Y^2)
- Random Sample $(X_1, Y_1), (X_2, Y_2), ..., (X_n, Y_n)$
- Population Covariance $\sigma_{XY} = Cov(X, Y) = E(X_i \mu_X)(Y_i \mu_Y)$
- Sample Covariance
- Population Correlation $\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$

$$S_{XY} = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \overline{X})(Y_i - \overline{$$

• Sample Correlation $r = \frac{S_{XY}}{S_X S_Y} = \frac{\sum_{i=1}^n (X_i - \overline{X})(Y_i - \overline{Y})}{\sqrt{\sum_{i=1}^n (X_i - \overline{X})^2 \sum_{i=1}^n (Y_i - \overline{Y})^2}}$

Y)



Characteristics of ρ

1) ρ has a value between -1 and +1.

- closer to +1 ⇒ strong positive linear relation
 - closer to -1 ⇒ strong negative linear relation.
 - closer to 0 ⇒ weak linear relation
- 2) If all values of X and Y are located on a straight line, ρ is either +1 or -1.
- 3) ρ is only a measure of linear relationship between two variables.
 - if $\rho = 0$, there is no linear relationship between the two variables, but there may be a different relationship

Correlation analysis





Correlation analysis

 \Box Testing population correlation coefficient ρ

Null Hypothesis: $H_0: \rho = 0$ Test Statistic: $t_0 = \sqrt{n-2} \frac{r}{\sqrt{1-r^2}} \sim t_{n-2}$

Rejection Region of H_0 :

1) $H_1: \rho < 0$ Reject H_0 if $t_0 < -t_{n-2; \alpha}$ 2) $H_1: \rho > 0$ Reject H_0 if $t_0 > t_{n-2; \alpha}$ 3) $H_1: \rho \neq 0$ Reject H_0 if $|t_0| > t_{n-2; \alpha/2}$



[Example 5.5.1] Based on the survey of advertising costs and sales for 10 companies that make the same product, we obtained the following data. Test the hypothesis that the population correlation coefficient is zero with the significance level 0.05.

Company	1	2	3	4	5	6	7	8	9	10	
Advertise (X)	4	6	6	8	8	9	9	10	12	12	
Sales (Y)	39	42	45	47	50	50	52	55	57	60	

Correlation analysis

<Answer of Example 5.5.1>

File	ſ	EX120101_SalesByAdvertise.csv							
Y Va	r			b	y X Var				
2: 5	Sales		~	1:	Advertise	е			
(S	elected da	ta: Raw Data)	(M	Iultiple Sele	ction)			
Selec	tedVar	/2 by V1	,						
	Advertis	e Sales	V3		V4	VS			
1	4	39							
2	6	42							
3	6	45							
4	8	47							
5	8	50							
6	9	50							
7	9	52							
8	10	55							
9	12	57							
10	12	60							
11									



Correlation analysis	id	X	Y	X^2	<i>Y</i> ²	XY
	1	4	39	16	1521	156
<answer 5.5.1="" example="" of=""></answer>	2	6	42	36	1764	252
$\mathbf{SXX} = \sum_{i=1}^{n} (X_i - \overline{X})^2$	3	6	45	36	2025	270
$= \sum_{i=1}^{n} X_i^2 - n \overline{X}^2$	4	8	47	64	2209	376
$= 766 - 10 \times 8.4^2 = 60.4$	5	8	50	64	2500	400
$SYY = \sum_{i=1}^{n} (Y_i - \overline{Y})^2$	6	9	50	81	2500	450
$-\sum_{i=1}^{n} \langle i \rangle^2 - n \overline{V}^2$	7	9	52	81	2704	468
$- \sum_{i=1}^{n-1} \sum_{i=1}^{n-1$	8	10	55	100	3025	550
$= 25097 - 10 \times 49.7^{-} = 396.1$	9	12	57	144	3249	684
$SXY = \sum_{i=1}^{n} (X_i - X)(Y_i - Y)$	10	12	60	144	3600	720
$= \sum_{i=1}^{n} X_i Y_i - n X Y$	Sum	84	497	766	25097	4326
$= 4326 - 10 \times 8.4 \times 49.7 = 151.2$	Mean	8.4	49.7			

$$S_{XY} = \frac{1}{n-1} SXY = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \overline{X})(Y_i - \overline{Y}) = \frac{151.2}{10-1}$$

$$r = \frac{S_{XY}}{S_X S_Y} = \frac{\sum_{i=1}^n (X_i - \overline{X})(Y_i - \overline{Y})}{\sqrt{\sum_{i=1}^n (X_i - \overline{X})^2 \sum_{i=1}^n (Y_i - \overline{Y})^2}} = \frac{SXY}{\sqrt{SXX SXY}} = \frac{151.2}{\sqrt{60.4 \times 396.1}} = 0.978$$

Correlation analysis

<Answer of Example 5.5.1> $t_0 = \sqrt{n-2} \frac{r}{\sqrt{1-r^2}} = \sqrt{10-2} \frac{0.978}{\sqrt{1-0.978^2}} = 13.26$ $t_{10-2;0.025} = 2.306$ Hence $H_0: \rho = 0$ is rejected

Regression Analysis				
Regression	y =	28.672 +	2.503 x	
Correlation Coefficient	r = 0.978	$H_0: \rho = 0$ $H_1: \rho \neq 0$	t value = 13.117	p value < 0.0001
Coefficient of Determination	$r^2 = 0.956$			
Standard Error	s = 1.483			



- Regression analysis is a statistical method
 - a mathematical model of relationships between variables,
 - estimates model using measured values of the variables,
 - uses estimated model to describe the relationship between variables, or to apply it to the analysis such as forecasting.
- Mathematical model ⇒ regression equation
- A variable affected by other variables is called a dependent variable. ⇒ response variable
- Variables that affect dependent variable are called independent variables.
 explanatory variable

Simple regression analysis

- Population regression model $Y_i = \alpha + \beta X_i + \epsilon_i$, i = 1, 2, ..., nEstimated regression equation $\widehat{Y}_i = a + b X_i$ Residuals $e_i = Y_i - \widehat{Y}_i$
- Method of Least Squares

A method of estimating regression coefficients so that total sum of the squared errors occurring in each observation is minimized.

Find α and β which minimize $\sum_{i=1}^{n} \epsilon_i^2 = \sum_{i=1}^{n} (Y_i - \alpha - \beta X_i)^2$

• Least square estimator of α and β

$$b = \frac{\sum_{i=1}^{n} (X_i - \overline{X})(Y_i - \overline{Y})}{\sum_{i=1}^{n} (X_i - \overline{X})^2}$$
$$a = \overline{Y} - b \overline{X}$$

5.5 Regression analysis Simple regression analysis

- Goodness of Fit for Regression Line
- Residual standard error s is a measure of the extent to which observations are scattered around the estimated line.

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \widehat{Y}_i)^2$$

The residual standard error s is defined as the square root of s^2 .

SST = SSE + SSR

$$SST = \sum_{i=1}^{n} (Y_i - \overline{Y})^2 \qquad df \quad n-1$$

$$SSE = \sum_{i=1}^{n} (Y_i - \widehat{Y}_i)^2 \qquad df \quad n-2 \qquad R^2 = \frac{SSR}{SST}$$

$$SSR = \sum_{i=1}^{n} (\widehat{Y}_i - \overline{Y})^2 \qquad df \quad 1$$

Simple regression analysis

 \Box Inference for the parameter β

- Point estimate: $b = \frac{\sum_{i=1}^{n} (X_i \overline{X})(Y_i \overline{Y})}{\sum_{i=1}^{n} (X_i \overline{X})^2} \sim N(\beta, \frac{\sigma^2}{\sum_{i=1}^{n} (X_i \overline{X})^2})$
- Standard error of estimate *b*: $SE(b) = \frac{s}{\sqrt{\sum_{i=1}^{n} (X_i \bar{X})^2}}$
- Confidence interval of β : $b \pm t_{n-2; \alpha/2} \times SE(b)$
- Testing hypothesis:

Null hypothesis:
$$H_0: \beta = \beta_0$$
Test statistic: $t = \frac{b - \beta_0}{SE(b)}$

1) $H_1: \beta < \beta_0$ Reject H_0 if $t < -t_{n-2; \alpha}$ 2) $H_1: \beta > \beta_0$ Reject H_0 if $t > t_{n-2; \alpha}$ 3) $H_1: \beta \neq \beta_0$ Reject H_0 if $|t| > t_{n-2; \alpha/2}$

Simple regression analysis

- **\Box** Inference for the parameter α
- Point estimate: $a = \overline{Y} b \overline{X} \sim N(\alpha, (\frac{1}{n} + \frac{\overline{X}^2}{\sum_{i=1}^n (X_i \overline{X})^2})\sigma^2)$
- Standard error of estimate *a*: $SE(a) = s \sqrt{\frac{1}{n} + \frac{\overline{X}^2}{\sum_{i=1}^n (X_i \overline{X})^2}}$
- Confidence interval of β : $a \pm t_{n-2; \alpha/2} \times SE(a)$
- Testing hypothesis:

Null hypothesis: $H_0: \alpha = \alpha_0$ Test statistic: $t = \frac{a - \alpha_0}{SE(a)}$

1) $H_1: \alpha < \alpha_0$ Reject H_0 if $t < -t_{n-2; \alpha}$ 2) $H_1: \alpha > \alpha_0$ Reject H_0 if $t > t_{n-2; \alpha}$ 3) $H_1: \alpha \neq \alpha_0$ Reject H_0 if $|t| > t_{n-2; \alpha/2}$

5.5 Regression analysis Simple regression analysis

 $\Box \text{ Inference for the average value } \mu_{Y|x} = \alpha + \beta X_0$

- Point estimate: $\widehat{Y}_0 = a + bX_0$
- Standard error of estimate \widehat{Y}_0 : $SE(\widehat{Y}_0) = s \sqrt{\frac{1}{n} + \frac{(X_0 \overline{X})^2}{\sum_{i=1}^n (X_i \overline{X})^2}}$
- Confidence interval of $\mu_{Y|x}$:

$$\widehat{Y}_0 \pm t_{n-2; \alpha/2} \times SE(\widehat{Y}_0)$$

Table 5.5.2 Analysis of variance table for simple linear regression

Source	Sum of squares	Degrees of freedom	Mean Squares	F value
Regression	SSR	1	$MSR = \frac{SSR}{1}$	$F_0 = \frac{MSR}{MSE}$
Error	SSE	n-2	$MSE = \frac{SSE}{n-2}$	
Total	SST	n-1		

Simple regression analysis

[Example 5.5.2] In [Example 5.5.1], find the least squares estimate of the slope and intercept if the sales amount is a dependent variable and the advertising cost is an independent variable.

- Predict amount of sales when you have spent on advertising by 10.
- Calculate the value of the residual standard error and the coefficient of determination in the data on advertising costs and sales
- Prepare an ANOVA table and test it using the 5% significance level



Simple regression analysis

<Answer of Example 5.5.2>

$$b = \frac{\sum_{i=1}^{n} (X_i - \overline{X})(Y_i - \overline{Y})}{\sum_{i=1}^{n} (X_i - \overline{X})^2} = \frac{151.2}{60.4} = 2.503$$

$$a = \overline{Y} - b\overline{X} = 49.7 - 2.503 \times 8.4 = 28.672$$

• Forecasting \widehat{Y}_i = 28.672 + 2.503 X_i 28.671 + 2.503 × 10 = 53.705

•
$$s^2 = \frac{1}{n-2} \sum_{i=1}^{n} (Y_i - \widehat{Y}_i)^2$$

= $\frac{17.622}{(10-2)} = 2.203$
• $R^2 = \frac{SSR}{SST} = \frac{378.429}{396.1} = 0.956$

Simple regression analysis

	X _i	Y _i	\widehat{Y}_i	$\frac{\text{SST}}{\sum (Y_i - \overline{Y})}$)2	$\frac{SSR}{\sum(\mathbf{\hat{Y}}_i - \mathbf{\hat{Y}})}$	$\frac{SSE}{\sum (Y_i - \widehat{Y}_i)}$	$)^2$			
1	4	39	38.639	114.49		122.346	0.130				
2	6	42	43.645	59.29		36.663	2.706				
3	6	45	43.645	22.09		36.663	1.836				
4	8	47	48.651	7.29		1.100	2.726				
5	8	50	48.651	0.09		1.100	1.820				
6	9	50	51.154	0.09		2.114	1.332				
7	9	52	51.154	5.29		2.114	0.716				
8	10	55	53.657	28.09		15.658	1.804				
9	12	57	58.663	53.29		80.335	2.766				
10	12	60	58.663	106.09		80.335	1.788				
Sum	84	497	496.522	396.1		378.429	17.622				
verage	8.4	49.7	[AN	[ANOVA]							
			Fa	ctor	tor Sum of Squares		deg of freedom	Mean Squares			
			Regr	ession	378.501		378.501		1	378.501	
			E	rror		17.599	8	2.200			

Total

A

p value

< 0.0001

F value

9

396.100

172.052

Simple regression analysis

- 1) Inference for β
- b = 2.50333

$$SE(b) = \frac{s}{\sqrt{\sum_{i=1}^{n} (X_i - \overline{X})^2}} = \frac{1.484}{60.4} = 0.1908$$

- Confidence interval of β : $b \pm t_{n-2; \alpha/2} \times SE(b)$ 2.5033 \pm 3.833 \times 0.1908 \Leftrightarrow (1.7720,3.2346)
- Test statistic for H_0 : $\beta = 0$ H_1 : $\beta \neq 0$

Reject
$$H_0$$
 if $|t| > t_{n-2; \alpha/2}$
 $t = \frac{b - \beta_0}{SE(b)} = \frac{2.5033 - 0}{0.1908} = 13.22$
Since $t_{8; 0.025} = 3.833$, H_0 is rejected.

Simple regression analysis

2) Inference for α

• a = 29.672

$$SE(a) = s \sqrt{\frac{1}{n} + \frac{\overline{X}^2}{\sum_{i=1}^n (X_i - \overline{X})^2}} = 1.484 \sqrt{\frac{1}{10} + \frac{8.4^2}{60.4}} = 1.670$$

• Test statistic for
$$H_0: \alpha = 0$$
 $H_1: \alpha \neq 0$
Reject H_0 if $|t| > t_{n-2; \alpha/2}$
 $t = t = \frac{a - \alpha_0}{SE(a)} = \frac{29.672 - 0}{1.670} = 17.1657$
Since $t_{8; 0.025} = 3.833$, H_0 is rejected.

3) Confidence interval of $\mu_{Y|x}$: $\widehat{Y}_0 \pm t_{n-2; \alpha/2} \times SE(\widehat{Y}_0)$ if x = 8, \widehat{Y}_0 = 49.699, \Rightarrow 49.699 \pm 3.833 \times 0.475

Simple regression analysis

12 13 62 62 60 -60 58 58 -56 56 54 54 52 52 -Sales 50 -50 48 -48 46 -46 44 -44 42 -42 40 -40 38 -38 y = (28.67)+(2.50)x Advertise $r = 0.98 r^2 = 0.96$

Sales(y) : Advertise(x) Scatter Plot

Parameter	Estimated Value	std err	t value	p value
Intercept	28.672	1.670	17.166	< 0.0001
Slope	2.503	0.191	13.117	< 0.0001

*****Simple regression analysis

Residual analysis

Standardized Residual vs Forecasting Plot





Standardized Residual Q-Q Plot



Population regression model

$$Y_i = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon_i$$
, $i = 1, 2, \dots, n$

 $Y = X \boldsymbol{\beta} + \boldsymbol{\epsilon}$

$$\boldsymbol{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \boldsymbol{X} = \begin{bmatrix} 1 & X_{11} & X_{12} & X_{1k} \\ 1 & X_{21} & X_{22} & X_{2k} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & X_{n2} & X_{nk} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

Multiple regression analysis

Least squares method

A method of estimating regression coefficients so that total sum of the squared errors occurring in each observation is minimized.

Find α and β which minimize $\sum_{i=1}^{n} \epsilon_i^2 = \epsilon' \epsilon = (Y - X \beta)' (Y - X \beta)$

Least Square Estimator of α and β
 b = (X'X)⁻¹(X'Y)

• Residuals $e_i = Y_i - \hat{Y}_i = Y_i - b_0 + b_1 X_{i1} + b_2 X_{i2} + \cdots + b_k X_{ik}$ Residual standard error *s*

$$s = \sqrt{\frac{1}{n-k-1}\sum_{i=1}^{n}(Y_i - \widehat{Y}_i)^2}$$

Multiple regression analysis

Analysis of Variance for Multiple Linear Regression

Source	Sum of Squares	Degrees of Freedom	Mean Squares	F value
Regression Error	SSR SSE	<i>k</i> n – <i>k</i> – 1	MSR=SSR / <i>k</i> MSE=SSE/ (n - <i>k</i> - 1)	$F_0 = \frac{MSR}{MSE}$
Total	SST	n - 1		

- $H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0$ $H_1:$ At least one of k number of $\beta'_i s$ is not equal to 0
- Reject H_0 if $F_0 > F_{k,n-k-1;\alpha}$

Multiple regression analysis

 \Box Inference for the parameter β_i

- Point estimate: b_i
- Standard error of estimate b_i : $SE(b_i) = \sqrt{c_{ii}} s$
- Confidence interval of b_i : $b_i \pm t_{n-k-1; \alpha/2} \times SE(b_i)$
- Testing hypothesis: Null hypothesis: H₀: Test statistic: t =

$$H_0: \beta_i = \beta_{i0}$$
$$t = \frac{b_i - \beta_{i0}}{SE(b_i)}$$

1) $H_1: \beta_i < \beta_{i0}$ Reject H_0 if $t < -t_{n-k-1; \alpha}$ **2)** $H_1: \beta_i > \beta_{i0}$ Reject H_0 if $t > t_{n-k-1; \alpha}$ **3)** $H_1: \beta_i \neq \beta_{i0}$ Reject H_0 if $|t| > t_{n-k-1; \alpha/2}$

Multiple regression analysis

[Example 5.5.3] When logging trees in forest areas, it is necessary to investigate the amount of timber in those areas. Since it is difficult to measure the volume of a tree directly, we can think of ways to estimate the volume using the diameter and height of a tree that is relatively easy to measure. Draw a scatter plot matrix of this data and consider a regression model for this problem. Diameter(cm) Height(m) Volume()

21.0	21.33	0.291	
21.8	19.81	0.291	
22.3	19.20	0.288	
26.6	21.94	0.464	
27.1	24.68	0.532	
27.4	25.29	0.557	
27.9	20.11	0.441	
27.9	22.86	0.515	
29.7	21.03	0.603	
32.7	22.55	0.628	
32.7	25.90	0.956	
33.7	26.21	0.775	
34.7	21.64	0.727	
35.0	19.50	0.704	
40.6	21.94	1.084	

Multiple regression analysis

<Answer of Example 5.5.3>



Regression Analysis					
Regression y =	(-1.024)	+ (0.037) X ₁	+ (0.024) X ₂		
Multiple Correlation Coeff	0.961	Coefficient of Determination	0.924	Standard Error	0.069
Parameter	Estimated Value	std err	t value	p value	95% Confidence Interval
0	-1.024	0.188	-5.458	0.0001	(-1.358 ,-0.689)
1 Diameter	0.037	0.003	10.590	< 0.0001	(0.031 ,0.043)
₂ Height	0.024	0.008	2.844	0.0148	(0.009 ,0.038)
[ANOVA]					
Factor	Sum of Squares	deg of freedom	Mean Squares	F value	p value
Regression	0.7058	2	0.3529	73.1191	< 0.0001
Error	0.0579	12	0.0048		
Total	0.7638	14			

Multiple regression analysis

[Example 5.5.4]

Standardized Residual vs Forecasting Plot





Summary

- Sampling distribution and estimation:
 - Central limit theorem
 - Point and interval estimation for a population mean
- Testing hypothesis for a population mean:
 - Type 1 and 2 error, significance level, p-value
 - Residual analysis
- Testing hypothesis for two population means
- Testing hypothesis for several population means
- Regression analysis:
 - Correlation analysis
 - Simple linear regression and multiple linear regression



Thank you !!!