

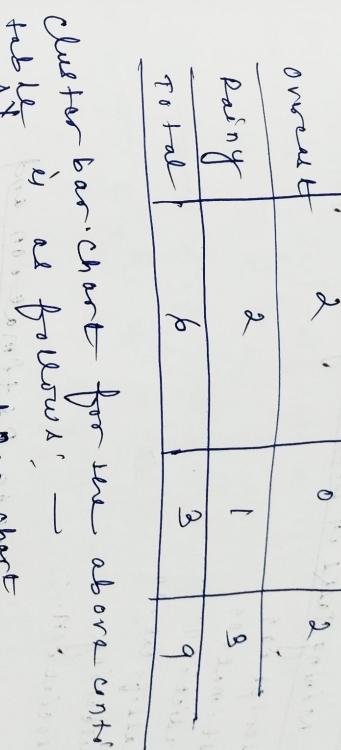
$$\begin{aligned}
 v_6' &= \frac{58-58}{19.47} = -0.308 \\
 v_7' &= \frac{58-58}{19.47} = -0.102 \\
 v_8' &= \frac{60-58}{19.47} = -0.102 \\
 v_9' &= \frac{63-58}{19.47} = 0.256
 \end{aligned}
 \quad
 \begin{aligned}
 v_{10}' &= \frac{70-58}{19.47} = 0.616 \\
 v_{11}' &= \frac{70-58}{19.47} = 0.616 \\
 v_{12}' &= \frac{110-58}{19.47} = 2.670
 \end{aligned}$$

Contingency Table

→ A contingency table (or cross-tabulation) is a matrix that displays the frequency distribution of variables. It is especially useful when both variables are categorical.

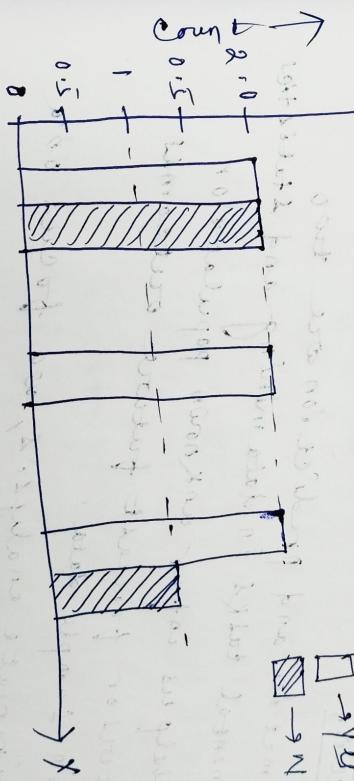
Example

weather	play	no play	Total
sunny	2	2	4
overcast	2	0	2
Rainy	1	3	4
Total	5	5	10



Cluster bar chart for the above contingency table is as follows:

clustered Bar-Chart



HATIC
AN



Chapter-5

Univariate Statistical Analysis

Data Mining Tasks in Discovering Knowledge in Data

- & Descriptive
- & Estimation
- & Prediction
- & Classification
- & Clustering
- & Association.

Statistical Approaches to Estimation and Prediction

- * Estimation and prediction are two fundamental tasks in data mining and statistics.
- * They help us infer unknown population parameters or forecast future outcomes based on sample data.
- * In univariate analysis, we focus on a single variable at a time.

I) Estimation

- * Estimation involves using sample data to approximate population parameters.
- * Point-estimation provides a single best guess of a parameter.

Example: Sample mean \bar{x} as an estimation of population mean μ .

$$\therefore \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- * Interval Estimation (Confidence Intervals) provide a range of plausible values for the parameter μ .
Example: 95% confidence interval for mean is $\bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$

- * Captures uncertainty in estimation.
- * Properly written $\bar{x} = \text{sample mean}$ (σ = sample standard deviation if σ unknown)

$$\bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} = \text{sample size}$$

$$Z_{\alpha/2} = \text{critical value from standard normal distribution}$$

$$\alpha/2 = 0.025$$

$$\text{For } 95\%, Z_{0.025} = 1.96.$$

- * Captures uncertainty in estimation.

a) Estimator Properties

- 1) Unbiasedness: Expected value equals true parameter.
- 2) Consistency: The estimator converges to the true value as $n \rightarrow \infty$.
- 3) Efficiency: Minimum variance among unbiased estimators.

II - Prediction

- a) regression - Based prediction:
 - & regression - Based prediction to forecast future samples data to fit with sample data
 - & it uses prediction of univariate analysis, regression models
 - & in univariate linear regression often use simple linear model
 - $y = \beta_0 + \beta_1 x$
 - & prediction intervals given to fit uncertainty to mean to confidence interval to fit to different forecasts, given account for to different intervals between actual observation and random estimation error and random estimation error
 - & prediction error by the difference between observed and predicted values
 - & prediction error is measured using mean squared error (MSE)
 - $$\text{example: } \begin{bmatrix} \text{actual} = [150, 17] \\ \text{predicted} = [2, 16, 17, 8] \end{bmatrix} \text{ MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$\text{RMSE} = \sqrt{\text{MSE}} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Application of Estimation and Prediction

In Data Mining

- a) Estimation: customer average spending
- & Estimation: defect rate in manufacturing
- & Prediction: forecasting sales, predicting loan default
- & Both tasks under class of classification models (Predicting categories) and regression models (Predicting continuous values)

Confidence Interval Estimation of the Mean

When we collect a sample from a population, we usually compute a point estimate (like the sample mean \bar{x}) to approximate the population mean μ .

But a single "number down't tell us how reliable the estimate is."

Confidence Intervals (C.I.) provide a range of possible values for the population parameter incorporating the uncertainty due to sampling.

C.I. sampling introduces variability as different samples are having different estimates.

A C.I. expresses this variability by giving a range rather than a single point.

Example: instead of saying "average height = 165cm" we say "average height is between 162 - 168cm with 95% confidence".

C.I. is more informative than point estimates as a point estimate alone doesn't show precision but C.I. shows both the estimates and the margin of error.

& margin of error indicates the precision where as wider C.I. indicates low precision.

C.I. helps to determine whether a population parameter is significantly different from a benchmark.

Example

Example
→ A 95% CI for mean exam score is (72, 78)
→ If a 95% CI for mean is 70, we can be
95% confident that the population mean
and the passing threshold is exceeded.
and the passing threshold is exceeded.

confidence interval

- * confidence intervals and hypothesis testing are closely related.
- * a hypothesized value (e.g. $\mu = 0$) lies outside of a hypothesized value (e.g. $\mu = 0$) → the null hypothesis would be rejected.
- * if the null hypothesis were true, it would be rejected at the same confidence level but at the same confidence level.

Practical Applications

- a) Medical studies: CI for drug effectiveness shows we're likely to get improvement.
- b) Business Analytics: CI for average customer spending helps in forecasting revenue.
- c) Predictive modelling: CI as part of predictions quantifies reliability.

* Key Insights

If confidence interval (e.g. 95%) means that repeated sampling many times, 95% of the intervals would contain the true mean. If sample size is large, CI is narrower.

$$CI = \bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

Steps in constructing a CI

1. Identify sample statistics: \bar{x} , s , n
2. Choose confidence level (commonly 90%, 95% or 99%)
3. find critical value

 - a) Z_{α/2} for large samples or known σ
 - b) t_{df=n-1} for small samples or unknown σ.

* A CI provides a range of possible values for a population mean based on sample data.

* It reflects both the point estimate (sample mean \bar{x}) and the uncertainty due to sampling variability.

* The confidence level (e.g. 95%) indicates the proportion of intervals that would contain the true mean if repeated sampling were done.

General formulae

- * When population standard deviation is known → use Z-distribution
- * $CI = \bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$
- * When σ is unknown → use t-distribution with $n-1$ degrees of freedom

$$CI = \bar{x} \pm t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}}$$

4. compute margin of error(ME):

$$ME = \text{critical value} \times \frac{1}{\sqrt{n}}$$

5. Construct Interval

$$CI = (\bar{x} - ME, \bar{x} + ME)$$

Example

Suppose

→ sample mean, $\bar{x} = 50$

→ sample standard deviation, $s = 10$

→ sample size, $n = 25$

→ confidence level = 95%.

Find out CI.

$$\text{Ans: } df = n - 1 = 25 - 1 = 24$$

$$\alpha = 1 - 0.95 \rightarrow \alpha = 0.05$$

$$\alpha/2 = 0.05/2 = 0.025$$

so critical value = $t_{\alpha/2, n-1}$ or $t_{0.025, 24}$

$$= t_{0.025, 24} = 2.064$$

$$= 2.064 \times 10 = 20.64$$

∴ CI = $(\bar{x} - ME, \bar{x} + ME)$

$$= (50 - 20.64, 50 + 20.64)$$

$$= (29.36, 79.64)$$

∴ CI = $(29.36, 79.64)$

Now we are using a z-table, then for

a two-tailed confidence interval, divide

α by 2 ($\alpha/2 = 0.025$ for 95% confidence).

Look up the z-value in the z-table to find

the corresponding critical value which is

the corresponding critical value which is typically 1.96 for a 95% confidence interval.

How to Reduce the Margin of Error

Margin of Error(ME) = Critical value $\times \frac{1}{\sqrt{n}}$

→ so the margin of error we can reduce by changing the sample size as follows:-

- 1) $t_{\alpha/2}$, which depends on the confidence level and the sample size.

- 2) use sample standard deviation, which is a characteristic of the data, and may not be changed.

- 3) n , the sample size.

- Thus we may decrease our margin of error in two ways, which are as follows:-

$$= 2.064 \times \frac{10}{\sqrt{25}} \\ = 2.064 \times 2 = 4.128$$

$$= (29.36, 79.64)$$

$$= \frac{2t_{\alpha/2}}{\sqrt{n}} \sqrt{n}$$

$$= (4.128, 4.128)$$

$$= 2.064 \times \frac{10}{\sqrt{n}}$$

1) By decreasing the confidence level, i.e., which reduced the value of $Z_{\alpha/2}$, and therefore reduced NE. But it is not recommended.

2) By increasing the sample size, the margin recommended is not increasing. It is only way to decrease the margin of error while maintaining the level of confidence.

Confidence Interval Estimation of Proportion

Shows that 46% of 8333 customers had changed so that an estimate of the company's total population proportion π of all customers who churn is

Example
Let $n = 200$, $x = 60$, confidence = 95%

$$p = \frac{x}{n} = \frac{60}{200} = 0.30$$

$$\sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{0.30 \times (1-0.30)}{200}} \approx 0.032$$

$\alpha = 0.05$ since 95% confidence $\alpha/2 = 0.025$, $2 \times Z_{\alpha/2} = 1.96$

$$\text{so } C.I. = p \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} = 0.30 \pm 1.96 \times 0.032$$

The confidence interval for the population proportion is given

$$C.I. = p \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

$$\text{where } p = \frac{x}{n} = \frac{70 - 45}{3333} = 0.1449$$

$$= \frac{n}{x} = \frac{45}{3333} = 0.0135$$

$Z_{\alpha/2} \rightarrow$ critical value from the standard normal distribution (e.g. 1.96 for 95% confidence
 $\alpha = 0.05$
 $\alpha/2 = 0.025$)

$$\text{C.I.} = 0.30 \pm 1.96 \sqrt{\frac{0.30 \times 0.70}{200}}$$

Interpretation

We are 95% confident that the true population proportion lies between 23.7% and 36.3%.

- Wider intervals = more uncertainty
- narrower intervals = more precision

- Increasing sample size reduces uncertainty

Hypothesis Testing for the Mean

→ Hypothesis testing is a procedure which claims about the value of a population parameter (such as mean or standard deviation) may be considered using the evidence from the sample.

→ Two competing statements on hypotheses are crafted about the parameter value which are as follows:

		Reality	
		H_0 true: Defendant did not commit crime	H_0 false: Defendant committed crime
Jury's Decision	Reject H_0 : finding defendant guilty	Type I error	Correct decision
	Do not reject H_0 : finding defendant not guilty	Correct decision	Type II error

- The two possible conclusions are:
 - a) reject H_0 and b) do not reject H_0 .

Example

- H_0 : Defendant is innocent
 H_a : Defendant is guilty

- The following table illustrates the four possible outcomes of the criminal trial with respect to the jury's decision and what is true in reality.

- a) The null hypothesis H_0 is tested to any hypothesis we wish to state and denoted by H_a . It is the statement of no hypothesis. H_a represents an alternative hypothesis.
- b) The research hypothesis H_a represents a research claim about the value of the parameter.

- The rejection of H_0 leads to the acceptable alternative hypothesis, denoted by H_a .
- The alternative hypothesis H_a usually represents one guess to be an answer or one theory to be tested and thus it's specification is crucial.

→ The probability of a Type I error is denoted as α . (Want to denote β .)
 → The probability of a Type II error is denoted as β . (Want to denote α .)
 of a Type II sample size increase in
 → for a constant β , α decrease.

→ α is associated with an increase in
 β , and vice versa.

→ In statistical analysis α is usually fixed at some small level of such as 0.05 and called the significance level.

* A test of any statistical hypothesis where the alternative is two-sided such as $H_0: \theta = \theta_0$,
 $H_a: \theta \neq \theta_0$
 is called a two-tailed test.
 where θ_0 represents the hypothesized value of θ .

	H_0 is true	H_0 is false
H_0 is rejected	correct decision	Type I error
H_0 is not rejected	Type II error	correct decision
Reject H_0		

*
 → Point Estimation:
 Estimation: → Estimation is the process of using sample data to infer the numerical value of an unknown population parameter such as mean or proportion.

Forms of Hypothesis testing

One Tailed Test

* A test of any statistical hypothesis where the alternative is one-sided such as $H_0: \theta = \theta_0$, $H_a: \theta > \theta_0$.

$H_a: \theta > \theta_0 \rightarrow$ alternative $< \theta_0$.

or perhaps $H_0: \theta = \theta_0$,

$H_a: \theta < \theta_0$,

is called one-tailed test

Sampling Error:

* A sampling error is a statistical error that occurs when a sample does not represent the entire population.

* Sampling error refers to the difference between estimated obtained from a sample and the true population value.
 e.g. sampling error for mean = $|\bar{x} - \mu|$

Confidence Interval

Confidence interval is a confidence interval as unknown & on statistics is used to estimate an unknown a range of values used to population statistical parameters, such as a mean.

A confidence interval consists of an interval population parameter consists of a point estimator produced by a confidence level of number produced by a confidence level to other width an associated probability that the interval specifies fitting the parameter.

contains the parameter.

- * the general form of confidence intervals
- C.I = point estimate ± margin of error

Margins of Error
Estimate of the precision of the interval estimation

* Smaller margin of error indicates greater precision.

The duration customer is normally distributed with mean 20 minutes and standard deviation 5 minutes. for the 95% sample size, construct a 95% confidence interval for the service calls

$$\begin{aligned} b) n &= 100 \\ \text{so } t_{\alpha/2, n-1} &= t_{0.025, 99} = 1.984 \\ \text{so C.I.} &= 20 \pm 1.984 \times \frac{5}{\sqrt{100}} = 20 \pm 0.992 \\ \text{C.I.} &= (20 - 0.992, 20 + 0.992) \\ &= (19.008, 20.992) \\ c) n &= 400, t_{\alpha/2, n-1} = t_{0.025, 399} = 1.966 \\ \text{so C.I.} &= 20 \pm 1.966 \times \frac{5}{\sqrt{400}} = 20 \pm 0.495 \\ \text{C.I.} &= (20 - 0.495, 20 + 0.495) \\ &= (19.5085, 20.4915) \end{aligned}$$

to calculate and interpret margin of error for the previous question for $n=25$ for the previous question for $n=25$ $\text{Margin of error} = t_{\alpha/2, n-1} \times \left(\frac{s}{\sqrt{n}} \right)$

$$\begin{aligned} \text{Margin of error} &= 2.064 \times \frac{5}{\sqrt{25}} = 2.064 \\ \text{Interpretation:} &\text{mean will} \\ &\text{be located no more than 2.064 minutes from the sample mean} \end{aligned}$$

cumulative probability of 95% mean = 20 minutes
confidence level = $1 - \frac{95\%}{100\%} = 0.95$

$$\alpha/2 = 0.025 \rightarrow \text{complementary probability} = 1 - \alpha/2 = 1 - 0.025 = 0.975$$

$$a) n = 25 \\ \text{Degree of freedom} = 25 - 1 = 24, s = 5, m = 20$$

$$\begin{aligned} \text{t-value at } 0.025 \text{ for degree of freedom 24} \\ i.e. t_{\alpha/2, n-1} = 2.064 \\ \text{so confidence interval} &= 20 \pm 2.064 \times \frac{5}{\sqrt{25}} \\ &= 20 \pm 2.064 \\ &i.e. (17.936, 22.064) \text{ minutes} \end{aligned}$$

Forms of Hypothesis testing

→ common treatment of hypothesis testing
for the mean is to restrict the hypothesis to the following three forms:

i) Left-tailed test,

$H_0: \mu \leq \mu_0$ versus $H_a: \mu > \mu_0$

↳ always take alternative hypothesis to know the difference.

ii) Right-tailed test

$H_0: \mu \geq \mu_0$ versus $H_a: \mu < \mu_0$

iii) Two-tailed test

$H_0: \mu = \mu_0$ versus $H_a: \mu \neq \mu_0$

where μ_0 represent hypothesized value of μ ,
when H_0 represent hypothesis (more than 30)
when the sample size is large or fine
& when the sample size is large distributed, the
population is normally distributed; the

$$t\text{-statistic} = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

follows a t-distribution, with $(n-1)$ degrees of freedom.

& data — no. of standard error above/below the
hypothesized mean μ_0 , that are
sample mean \bar{x} consider.

$$\frac{\bar{x} - \mu_0}{s/\sqrt{n}} - \text{standard error}$$

When the value of t-data is extreme, will indicate a conflict between the null hypothesis and the observed data.

& if the data represents empirical evidence where as the null hypothesis represents merely a claim, such conflicts are resolved in favour of the data.

so that when t-data is extremely, the null hypothesis H_0 is rejected. This extremes & is measured using the p-value.

& the p-value is the probability of observing a sample statistic (such as \bar{x} or t-data) at least as extreme as the statistic actually observed if the null hypothesis is true.

& as the p-value represents a probability, it's value must always fall below 0.1.

Table: How to calculate p-value

<u>Form of Hypothesis Test</u>	<u>p-value</u>
i) Left-tailed test	$P(t < t\text{-data})$
ii) Right-tailed test	$P(t > t\text{-data})$

$$P(t > t\text{-data})$$

$$= 2 \times P(t > t\text{-data})$$

$$= 2 \times P(t > 0, \text{then } t\text{-data})$$

$$= 2 \cdot P(t > t\text{-data})$$

$$= 2 \cdot (1 - P(t \leq t\text{-data}))$$

* A small p-value indicates conflict between the null hypothesis and the data and the null p-value is small if we consider the p-value to be small if it is less than α .

Reject H_0 if the p-value $< \alpha$

$$\text{or } H_0: \text{mean} = 2.4, H_1: \text{mean} \neq 2.4 \\ \bar{x} = 1.607, s = 1.892, n = 28, \bar{x}_0 = 2.4 \\ \alpha = 0.05$$

Perform hypothesis testing

$$H_0: \text{mean} = 2.4, H_1: \text{mean} \neq 2.4$$

so it's two tailed test

$$\text{so } t_{\text{data}} = \frac{\bar{x} - \bar{x}_0}{s/\sqrt{n}} \\ = \frac{1.607 - 2.4}{1.892/\sqrt{28}} \\ = -\frac{0.793}{0.359}$$

$$= -2.2178$$

t_{table}

$$\text{so } p\text{-value} = 2 \times P(t < |t_{\text{data}}|) \\ = \text{since } t \text{ table for } df = 27$$

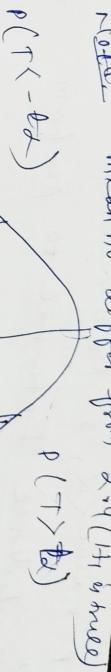
and $|t_{\text{data}}| > t_{\text{table}}$
take average of $2.7 + 2.5$

$$= 2 \times 0.0175 \alpha = 0.035$$

$$\text{so } \alpha = 0.05$$

~~reject~~ since p-value $< \alpha$, so we can't reject null hypothesis

so we can't reject null hypothesis and always same is to reject null hypothesis and accept alternative hypothesis



Note: mean no differ from 2.4 (H_1 is true)

$$p(T < -2.2178) = p(T > 2.2178) \\ p(T < -2.2178)$$

Q1. A random sample of 100 recorded deaths in the United States during the past 7 years showed an average life span of 71.6 years. Assuming a population standard deviation of 8.9 years does this seem to indicate that the mean life span today is greater than 70 years? Use $\alpha = 0.05$. $\bar{x} = 71.8$ $s = 8.9$ $n = 100$

$$H_1: \mu > 70 \quad H_0: \mu \leq 70 \quad \alpha = 0.05$$

$$\text{so } t_{\text{data}} = \frac{\bar{x} - \bar{x}_0}{s/\sqrt{n}} = \frac{71.8 - 70}{8.9/\sqrt{100}} = 1.8$$

$$= \frac{8.9/10}{0.89} = 1.8$$

$$\text{so } p\text{-value} = p(Z > 1.8) \approx 0.12359$$

so $p\text{-value} = p(Z > Z_{\text{data}}) \approx 0.12359$
but for right tailed need $Z > Z_{\text{data}}$
 $\text{so } p\text{-value} = p(Z > Z_{\text{data}}) = 1 - 0.876 = 0.12359$

Since p -value $\leq \alpha$, so H_0 is rejected.

~~skip for students~~
confidence intervals to perform
~~skip for students~~
Hypothesis is True

<u>p-value</u>	<u>Strength of Evidence Against H_0</u>
p -value ≤ 0.001	Extremely strong evidence
$0.001 < p$ -value ≤ 0.01	Very strong evidence
$0.01 < p$ -value ≤ 0.05	Solid evidence
$0.05 < p$ -value ≤ 0.10	Mild evidence
$0.10 < p$ -value ≤ 0.15	slight evidence
$0.15 < p$ -value	No evidence

<u>Confidence Level $100(1-\alpha)\%$</u>	<u>Level of Significance</u>
90%	0.10
95%	0.05
99%	0.01

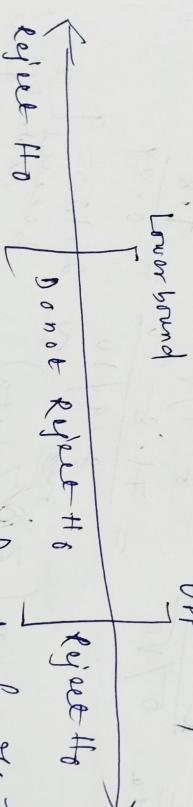
- (Lower bound, upper bound) = $(0.875, 2.339)$
- $\alpha = 0.05$
- For any no. of possible values of μ
- a) 0.5 b) 1.0 c) 2.4
- ① $H_0: \mu = 0.5$
 $H_1: \mu \neq 0.5$
- use α to find C.I., then see if the given μ values are given in C.I. or not.
 If yes accept H_0 , otherwise reject.
- ② $H_0: \mu = 1.0$
 $H_1: \mu \neq 1.0$

use α to find C.I., then see if the given μ values are given in C.I. or not.
 If yes accept H_0 , otherwise reject.

$$H_0: \mu = 2.4$$

$$H_1: \mu \neq 2.4$$

use α to find C.I., then see if the given μ values are given in C.I. or not.
 If yes accept H_0 , otherwise reject.



+ if a certain hypothesized value for μ falls outside the confidence level with confidence level $100(1-\alpha)\%$, then we two-tailed hypothesis

test with level of significance α will reject H_0 for μ at value $\mu \neq \mu_0$.

* if the hypothesized value for μ falls inside the confidence interval with confidence level $100(1-\alpha)\%$ then the two-tailed hypothesis test with level of significance α will not reject H_0 for that value of μ .

Table: Conclusion for three hypothesis tests using confidence level position in relation to CI.

H_0	H_0 vs H_a : $\pi \neq \pi_0$	Conclusion	concluding
$H_0: \pi = 0.5$	outside	reject H_0	
$H_0: \pi = 1.0$	inside	do not reject H_0	
$H_0: \pi = 2.4$	inside	reject H_0	
	at boundary 0.875 to 2.339		

Hypothesis Testing for the proportion

Hypothesis testing may also be performed about the population proportion π . The test statistic is

$$Z_{\text{data}} = \frac{\hat{p} - \pi_0}{\sqrt{\pi_0(1-\pi_0)/n}}$$

where π_0 is the hypothesized value of π and \hat{p} is the sample proportion.

$$\hat{p} = \frac{\text{number of successes } n}{n}$$

Table: Hypotheses and p-values for hypotheses tests about π

Hypothesis with $\alpha = 0.05$	p-value
Left-tailed test: $H_0: \pi \geq \pi_0$ vs. $H_a: \pi < \pi_0$	$P(Z < Z_{\text{data}})$
Right-tailed test: $H_0: \pi \leq \pi_0$ vs. $H_a: \pi > \pi_0$	$P(Z > Z_{\text{data}})$

Example

$$\text{Suppose } 3333 \text{ customers have charged } 483 \text{ at } 3333. \text{ Sample size } n = 3333. \text{ Hypothesis: } H_0: \pi = \frac{x}{n} = \frac{483}{3333} = 0.449$$

Suppose we would like to test using level of significance $\alpha = 0.10$. Whether π differs from 0.15. The hypothesis are: $H_0: \pi = 0.15$ versus $H_a: \pi \neq 0.15$.

$$\text{The test statistic is: } Z_{\text{data}} = \frac{\hat{p} - \pi_0}{\sqrt{\pi_0(1-\pi_0)/n}} = \frac{0.1449 - 0.15}{\sqrt{0.15(1-0.15)/3333}} = -0.8246$$

As $Z_{\text{data}} < 0$, we $p\text{-value} = 2 \times P(Z < Z_{\text{data}})$

$$= 2 \times P(Z < -0.8246)$$

$$= 2 \times 0.2048$$

$$= 0.4096$$

$$p\text{-value} = 0.4096 > \alpha (= 0.10)$$

fail to reject H_0

$$\frac{\text{By } 95\% \text{ confidence level}}{\alpha = 1 - \frac{95}{100} = 1 - 0.95 = 0.05}$$

$$\text{So critical value} = Z_{\alpha/2} = Z_{0.025} \approx 1.96$$

$$\text{Margin of Error (ME)} = \text{critical value} \times \sqrt{\frac{p(1-p)}{n}}$$

$$= 1.960 \times \sqrt{\frac{0.10 \times (1-0.10)}{1000}}$$

$$= 0.0588$$

$$\text{So confidence interval (CI)} \\ = p \pm ME = 0.10 \pm 0.0588 = (0.0412, 0.1588)$$

c) similarly for 99% confidence

Since population proportion is given 0.10
 Since population proportion is given 0.10
 since sample size > 30 , so population, which
 is 1000, so 2-distribution will be used
 since a) As confidence = 90%
 $\alpha = 1 - 0.90 = 0.10$, $\alpha/2 = 0.10 = 0.05$
 $\alpha/2 = 1 - \frac{90}{100} = 1 - 0.9 = 0.10$, $\alpha/2 = 0.10 = 0.05$
 so critical value = $Z_{\alpha/2} = Z_{0.05} \approx 1.645$

so Margin of Error (ME) = critical value $\times \sqrt{\frac{p(1-p)}{n}}$

$$= 1.645 \times \sqrt{\frac{0.10(1-0.10)}{1000}}$$

$$= 0.0494$$

so the confidence Interval (CI)

$$= \text{critical value } p \pm ME$$

$$= 0.10 \pm 0.0494 = (0.0506, 0.1494)$$

because after any test
 because after any test
 because after any test
 because after any test

- provide the hypothesis. State the meaning of the rejection rule.
- What is the rejection and interpretation? What is the conclusion?

Ans-a) The hypotheses are

- Null hypothesis is H_0 : The population mean μ is less than or equal to the hypothesized value.
- Alternative hypothesis is H_a : The population mean μ is greater than the hypothesized value.

mean \bar{x} is greater than the hypothesized mean $\mu_0 = 45.0$.
 \bar{x} is a right-tailed test.

b) Rejection Rule

We will reject the null hypothesis if H_0 , if the p-value for the test statistic is less than the significance level α . For our problem, we would reject the null hypothesis if the p-value $< \alpha = 0.05$.

Since sample size and sample standard deviation given

$$c) \text{So } t\text{-stat} = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

$$\bar{x} = 55$$

$$\begin{cases} n = 50 \\ s = 2.5 \\ \alpha = 0.05 \end{cases}$$

The first confidence interval we calculated for population mean turned out to be 15 minutes. Since our confidence level is 95%, our significance level α for these tests is calculated as follows:

$$\alpha = (1 - 95/100) = 0.05$$

For right-tailed test, we calculate the critical value $t_{\alpha/2}$ which is $t_{0.025}$ with 49 degrees of freedom.

$$= \frac{55 - 50}{2.5/\sqrt{50}} = 2$$

So for right-tailed test, p-value = $P(t > 2)$

$$= P(t > 2) \text{ with } 49$$

degree of freedom

$= 0.02412 \rightarrow$ (critical value of t statistic for one-tailed hypothesis)

Since p-value ($\alpha = 0.05$), so the null-hypothesis is rejected.

Ans-b)

Refer to the first confidence interval you calculated for the population mean duration for customer to wait whether this population is at least 15 minutes. To test whether this population mean differs from the following values giving level of significance $\alpha = 0.05$

a. 15 minutes b. 20 minutes c. 25 minutes

i. We can determine a confidence interval to perform hypothesis testing as it is complete to perform a two-tailed test. Our hypothesis for this type of test would be formed as $H_0: \mu = 15$, $H_a: \mu \neq 15$. Our confidence level is 95%, our significance level α for these tests is since our confidence level α for these tests is.

$$\alpha = (1 - 95/100) = 0.05$$

For right-tailed test, we calculate the critical value $t_{\alpha/2}$ which is $t_{0.025}$ with 49 degrees of freedom.

$$c) \bar{x} = 20 \pm 2.064 = (17.936, 22.064) \text{ minutes}$$

a. 15 minutes
 Since 15 minutes lies outside of (17.936, 22.064) in confidence interval, we would reject the null hypothesis as the test conclusion that the population mean duration is not 15 minutes.

b. 20 minutes
 Since 20 minutes lies within the confidence interval, we would not reject the null-hypothesis.

Hypothesis for null test and conclude that we have enough evidence to indicate that the population mean duration is not 20 minutes.

25 minutes vs outside (less than). Since 25 minutes is our wrong reject, we conclude no confidence interval, we don't have confidence for null test and conclude the null hypothesis for mean duration is not true.

that we're population that uses 15 minutes for proportion

→ hypothesis for proportion
25% in a sample of 1000 customers is 240. Then we can say the company's population proportion of chumbers when they compare their proportion, but whether the population proportion of chumbers is less than 25%, using a level of significance $\alpha = 0.01$

$$\text{test} - p = \frac{n}{n} = \frac{\text{respondent}}{\text{sample size}} = \frac{240}{1000} = 0.24$$

in the alternative hypothesis: Null hypothesis: population proportion $p > 0.25$

Alternative Hypothesis: H_a : population proportion $p < 0.25$.

→ it is one-tailed test to proportion. Since, it is given population proportion $p_{\text{chamber}} = 0.25$ and will be used to calculate Z-value.

$$\rightarrow \text{so } Z\text{-value} = \frac{p - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.24 - 0.25}{\sqrt{\frac{0.25 \times 0.75}{1000}}} = -0.01$$

$$= \frac{-0.01}{0.01369} = -0.7305$$

for left-tailed test
 $p\text{-value} = P(Z < -0.7305) = 0.232542$

Since the $p\text{-value} > \alpha = 0.01$, we fail to reject the null hypothesis concluding that there is no evidence to disprove it. In other words, there is no evidence to disprove that the population proportion is indeed less than the population proportion chamber will remain less than 25%

Chapter - 6 Statistical

Multivariate statistical methods

→ so far we have discussed inference methods for one variable at a time. Inference methods are also interested in relationships between two variables or multi-variate relationship between two variables which are analyzed.

Between one target variable we have two predictors or variables we have two

→ In discriminate analysis we wish to test independent samples and check if there for significance between two samples for proportion of the two training mean of proportion, we illustrate how

→ In this chapter, we partitioned our data into training and test data set for

the data is partitioned into training

and test and test data set for

data set and test data set for

Type of Multivariate Analysis

For a continuous variable, we use the two-sample t-test for the difference of means.

To test for the differences in proportions for a multinomial variable, we use the test for the homogeneity of proportions.

Two-sample T-test for Differences in Means

To test for the differences in population means, we use the following test statistic:

$$t\text{-stat} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{(\sigma_1^2/n_1) + (\sigma_2^2/n_2)}}$$

which follows an approximate t-distribution with degrees of freedom the smaller of $n_1 - 1$ and $n_2 - 1$, whenever either both populations are normally distributed or both sample size large.

Example

We partitioned the churn data set into a training set of 2529 records and a test set of 804 records. We would like to test the null hypothesis of the proportion by testing whether the population mean number of customer service calls differ between the two data sets. The summary statistics is

Dataset	Sample Mean	Sample SD	Sampling
Training set	$\bar{x}_1 = 1.5714$	$s_1 = 1.3126$	$n_1 = 2529$
Test set	$\bar{x}_2 = 1.5361$	$s_2 = 1.3251$	$n_2 = 804$

So we hypothesised are
 $H_0: \pi_1 = \pi_2$ Vs. $H_a: \pi_1 \neq \pi_2$

The test statistic is

$$t\text{data} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{(x_1^2/n_1) + (x_2^2/n_2)}{((1.3125)^2/2529) + ((1.325)^2/804)}}$$

$$= \frac{1.5714 - 1.536}{\sqrt{(1.3125)^2/2529 + (1.325)^2/804}} \\ = 0.6595$$

Since the given test is two-tailed
so p-value = $\exp(-|t|)$ where $t > 0$

$$= \exp(-0.6595)$$

= 0.5098

Since the p-value is greater than 0.05, there is no evidence
that the mean no. of customers in the training
data differs between the training data set and test data set.
and we test data set is valid at least, we partitioned it into
at least, we partitioned it into

for example, our partitioned resulted in
 $x_1 = 707$ of $n_1 = 2529$ customers in the training
data belonging to the voice mail plan
while $x_2 = 115$ if $n_2 = 804$ customers in
the test set belonging to that $\pi_1 = \frac{x_1}{n_1}$
 $= \frac{707}{2529} = 0.2795$

$$\hat{p}_2 = \frac{x_2}{n_2} = \frac{215}{804} = 0.2674$$

$$\hat{p}_{\text{pool}} = \frac{x_1+x_2}{n_1+n_2} = \frac{707+215}{2529+804} = 0.2766$$

The hypothesis are:-

$H_0: \pi_1 = \pi_2$ Vs. $H_a: \pi_1 \neq \pi_2$

The test statistic is

$$Z\text{data} = \frac{\sqrt{\hat{p}_{\text{pool}} \cdot (1-\hat{p}_{\text{pool}})((1/n_1) + (1/n_2))}}{\hat{p}_1 - \hat{p}_2}$$

Two Sample Z-test for Difference in Proportions

→ So we could turn to the two-sample
Z-test for the difference in proportions
→ The test statistic is

$$= \frac{0.2766 - 0.2674}{\sqrt{0.2766 \cdot 0.7234 ((1/2529) + (1/804))}} = 0.6736$$

The p-value is
 $p\text{-value} = 2 \times P(Z > |Z_{\text{data}}|)$ (for two-tailed)

$$= 2 \times P(Z > 0.6736) \\ = 0.5006.$$

Since p-value $> \alpha$, we fail to reject the null hypothesis.
 So there is no evidence that the proportion of voice mail plan members different between the training and the test data set.
 For this variable, test partition is valid.

Test for the Homogeneity of Proportions

Statistical data is a collection of binomial data to K+2 categories. For example, suppose a multinomial variable marital status has the values married, single and other. Suppose we have a training set of 1000 people with the frequencies of 250 people with the following distribution in Table 1.

Dataset	MARRIED	SINGLE	OTHER	Total
Training set	410	340	250	1000
Test set	95	85	70	250
Total	505	425	320	1250

Z_{data})

To determine whether significant differences exist between the proportions at the two data sets, we could then test for the homogeneity of proportions.

The hypothesis are:

H₀: P_{married, training} = P_{married, test},

P_{singer, training} = P_{singer, test},

Other, training = Other, test

As we see one of the claims in H₀ is wrong.

To determine whether there are significant differences between the training frequencies reported from the training and test data sets, we compare the observed frequencies with the expected frequencies that we would expect if the were true. For example, to find the expected frequency for the training set, we first calculate proportions of married and people overall proportion and we (i) find the overall proportion of married people in both the training and test sets, $\frac{505}{1250}$ and (ii) we multiply this overall proportion by the number of people in the training set, 1000, giving us the expected proportion of married people in the training set. To be people in the training set. To be

Expected frequency = $\frac{\text{Row total} \times \text{Column total}}{\text{Grand total}}$

$$\text{So expected frequency}_{\text{married, training}} = \frac{1000 \times 505}{1250} = 404.$$

As we see overall proportion of C_1 because the total proportions are equal.

Applying this formula to each cell in the table gives us the table of expected frequencies as follows:

Table: Expected frequency

Dataset	Mammal	Conifer	Other	Total
Training set	404	340	256	1000
Test set	101	85	64	250
Total	505	425	320	1250

→ the observed frequencies (O) and the expected frequencies (E) are compared using a test statistic

from the χ^2 -square distribution

$$\chi^2_{\text{data}} = \sum \frac{(O - E)^2}{E}$$

cell

$$\frac{\text{Observed frequency}}{\text{Expected frequency}} = \frac{O/E}{(O-E)/E} = \frac{O/E}{1 - O/E}$$

Mammal, training

$$\frac{404}{425} = \frac{404/425}{(425-404)/425} = 0.904$$

Mammal, test

$$\frac{101}{85} = \frac{101/(85-101)}{101} = 0.32$$

Conifer, training

$$\frac{340}{320} = \frac{(340-340)/320}{320} = 0$$

Conifer, test

$$\frac{85}{64} = \frac{(85-85)/64}{64} = 0$$

Other, training

$$\frac{256}{320} = \frac{(256-256)/320}{320} = 0$$

Other, test

$$\frac{64}{1250} = \frac{(64-64)/1250}{1250} = 0$$

$$\chi^2_{\text{data}} = \frac{(404-425)^2}{425} + \dots + \frac{(64-64)^2}{1250} = 1.15$$

p-value = value one to the right of

$$\chi^2_{\text{data}} \text{ under the } \chi^2 \text{ curve with degrees of freedom equal to } (\text{number of rows} - 1)(\text{number of columns} - 1) = 2$$

$$\text{So p-value} = P(\chi^2 > \chi^2_{\text{data}}) = P(\chi^2 > 1.15) = 0.5627$$

With df = 2

Since p-value is and is very large

so there is no evidence that the observed frequencies represent proportion that are significantly different from the training and test data sets. In other words, for this variable, the proportion

values:

Chi-square Test for Goodness of Fit of Mammal Data

→ Use chi-square goodness-of-fit for one categorical variable to see if its distribution

category match an expected distribution

→ Use chi-square test for Homogeneity when

you have one categorical variable but

you have two or more distributions for

two or more populations testing if homogenous

or not variable are the same (homogeneous)

in each group.

Example
Suppose a multinomial variable marital status takes the values married, single and other and that we suppose that 40% of the population we know that 40% of the population are married, 35% are single and 25% are married.

We are taking a sample and would

like to determine if the sample is drawn from the population. We could then test

the χ^2 (chi-square) goodness of fit test

\rightarrow The hypotheses for this χ^2 goodness of fit test would be as follows:

H_0 : $p_{\text{married}} = 0.40$, $p_{\text{single}} = 0.35$, $p_{\text{other}} = 0.25$

H_a : At least one of the proportions in H_0 is wrong

\rightarrow Our sample of size $n = 100$, yields the following observed frequencies represented by the letter O

$O_{\text{married}} = 36$, $O_{\text{single}} = 35$, $O_{\text{other}} = 29$

\rightarrow To determine whether these counts represent proportions that are significantly different from those expected in H_0 , we compare the observed frequencies with the expected frequencies that we would expect

\rightarrow If H_0 were true, then we would expect 40% of our sample to be married, that is the expected frequency for married is

$$E_{\text{married}} = n \cdot p_{\text{married}} = 100 \cdot 0.40 = 40$$

Similarly,

$$E_{\text{single}} = 100 \cdot p_{\text{single}} = 100 \times 0.35 = 35$$

$$E_{\text{other}} = 100 \cdot p_{\text{other}} = 100 \times 0.25 = 25$$

These frequencies are compared using the test statistic:

$$\chi^2_{\text{data}} = \sum \frac{(O - E)^2}{E}$$

\rightarrow Again, larger differences between the observed and expected frequencies, and thus a large value for χ^2_{data} , will lead to a small p-value and a rejection of the null hypothesis.

\rightarrow The test statistic is calculated as

etc

Marital Status	Observed frequency	Expected frequency	$(O_E - E)^2 / E$
Married	36	40	$\frac{(36-40)^2}{40} = 0.4$
Single	35	35	$\frac{(35-35)^2}{35} = 0$
Other	29	25	$\frac{(29-25)^2}{25} = 0.16$

\rightarrow The p-value is the area to the right of χ^2_{data} under the χ^2 -curve with $K-1$ degrees of freedom, where K is the no. of categories (here $K = 3$)

$$p\text{-value} = P(\chi^2 > \chi^2_{\text{data}}) = P(\chi^2 > 1.04) = 0.5945$$

To see the p-value is very large, so there is no evidence that the observed difference is significant to represent proportions that differ significantly from those in the null hypothesis. In other words, our sample is representative of the population.

Chapter - 7

Preparing to Model the Data

→ supervised and unsupervised learning are two main types of machine learning / data mining methods. In supervised learning, the model is trained with labeled data where each input has a corresponding output. On the other hand, unsupervised learning involves training the model with unlabeled data which helps to uncover patterns, structures or relationships within the data without predefined outputs.

Based	Supervised Learning	Unsupervised Learning
1) Definition	Supervised learning algorithm trains them to find patterns in data, where every data that has no preceding input has a corresponding output.	Unsupervised learning is to discover hidden patterns in data that has no preceding labels.

2) Goal	The goal of supervised learning is to discover rules or identify hidden patterns, structures based on input and relationships features.
3) I/P Data	Unlabeled input data is raw and unlabeled.

4) Human supervision	Supervised learning algorithm needs continuous human supervision to train the model.
5) Tasks	Regression, classification and clustering as association and dimensionality reduction.

6) Complexity	Supervised methods are computationally simple.
7) Algorithms	Linear Regression, K-Means clustering, DBSCAN, Decision tree, Random Forest, Naive Bayes, SVM.

8) Assumptions	Supervised methods are accurate if input features are independent and continuous, custom or segmentation.
9) Applications	Image classification, Anomaly detection, sentiment analysis, recommendation systems.

Statistical Methodology and Data Mining

Methodology

① Data Type	Work with structured data using probabilistic reasoning (mean, regression) for insights, divide data mining in massive, often messy data (clustering, classification) to find patterns algorithms (clustering, classification) forming the core of many prediction, w/ statistical techniques
	Work mainly with quantitative data
② Visual	Derive insights often through prediction.
	Based on probability theory, distribution.
④ Domain Knowledge	Exploratory focused on discovering hidden patterns.
	Relies on mathematics, heuristics (rules of thumb) play an important role.
⑤ Data Collection	Focuses on data collection and cleaning.
	Emphasis on working with existing data, not collecting it.
⑥ Hypotheses	No "a priori" hypotheses, <u>at priori</u> hypothesis
	No "a priori" hypotheses

Cross Validation

- cross validation is a technique used to check how well a machine learning model performs on unseen data while preventing overfitting.
- It works by:
 - * splitting the dataset into several parts.
 - * training the model off some part and testing it on the remaining part.
 - * Repeating this resampling process multiple times by choosing different parts of the dataset.
 - * Average the results from each validation step to get the final performance.

Types of cross validation

① Holdout validation

- in holdout validation method typically 50% data is used for training & 50% for testing making it simple and quick to apply.
- The major drawback of this method is that only 50% data is used for training the model which may make it important patterns in the other half which leads to high bias.
- ② Leave One Out Cross Validation

- on the entire dataset sweep →
 - * this method the model is trained on the entire dataset except for one data point which is used for testing. This process is repeated for each data point & are used for training → the data points are used for training resulting in low bias.

* Testing on a single data point can cause high variance especially if the point is an outlier.

It can be very time consuming for large datasets as it requires one iteration per data point.

(3) Stratified Cross Validation

→ It is a technique that ensures each fold of validation process has the same class distribution as the full dataset.

→ This is useful for imbalanced datasets where

→ one is under-represented.

→ Some classes are under-represented in a dataset divided into k-folds.

→ The dataset is divided into each fold.

→ Class proportions remain the same for every

→ In each iteration, one fold is used for testing and the remaining folds for training.

→ This process is repeated K times so that each fold is used once as the test set.

→ This process is repeated K times so that each fold is used once as the test set.

→ This process is repeated K times so that each fold is used once as the test set.

(4) K-fold Cross Validation

→ K-fold cross validation splits the dataset into K equal-sized folds.

→ The model is trained on K-1 folds and tested on the remaining fold.

→ This process is repeated K-times each time using a different fold for testing.

Note: → It is always suggested that the value of K - should be 10 as the lower value of K

takes towards validation and higher values of K leads to LOOCV method.

Dataset	Training	Test	Training
100% k=5	Training	Test	Training

Training	Test	Training
Training	Test	Training

Test	Training
K-fold Cross Validation	

K-fold Cross Validation

→ And iteration → Let 20% of data is used for testing and 80% is used for training

→ Let iteration → Let 20% of data is used for testing and 80% is used for training

→ This process continues until each fold has been used once as the test set.

→ This process continues until each fold has been used once as the test set.

Table → Suggested hypothesis is to test for relatedness between target variables.

Table → Suggested hypothesis is to test for relatedness between target variables.

Type of Target Variable Hypothesis Test

Continuous Two sample t-test for two sample difference in means

Two sample z-test for two sample difference in proportions

Multinomial Test for homogeneity of proportions

Testing set

Training set / test set or validation test

Training set:

- this is the actual dataset from which a model trains i.e. the model sees and learns from this data to predict the outcome or to make the right decisions.
- most of the training data is collected from several resources and then processed and organized to provide proper performance of the model.

Testing set

- type of training data highly determines the ability of the model to generalize and diversity of the training data, the better will be the performance of the model.

Testing set

- the dataset is independent of the training set but has a somewhat similar type of probability distribution of class labels and is used as a benchmark to evaluate the model, usually only after the training of the model is complete.

- testing set is usually a properly organized dataset having all kinds of data for scenarios that the model would probably be facing when used in the real world.

- often we hold out a testing set which is not used as a testing set which is not considered as a good practice - ~~because it's not used~~

Validation set

- the validation set is used to fine-tune the hyper-parameters of the model and is considered apart of the training of the model.
- the model is held this data for evaluation but does not lead from this data, providing an objective unbiased evaluation of the model.

Overfitting and Underfitting

- * Machine learning models aim to perform well on both training data and new unseen data, and is considered "good".
- i) to learn patterns differently from the training data.
- ii) to generalize well to new unseen data.
- iii) to avoid memorizing the training data (overfitting) or failing to capture relevant patterns (underfitting).
- to evaluate how well a model learns and generalizes, we monitor its performance on the training data and a separate validation dataset which is often measured by its accuracy or prediction error.
- However, achieving this balance can be challenging.
- two common issues that affect a model's performance and generalization ability are overfitting and underfitting.
- underfitting are major contributions to poor performance in machine learning models.

Bias and Variance in Machine Learning

→ Bias and variance are two key sources of error in machine learning models that directly impact their performance and generalization ability.

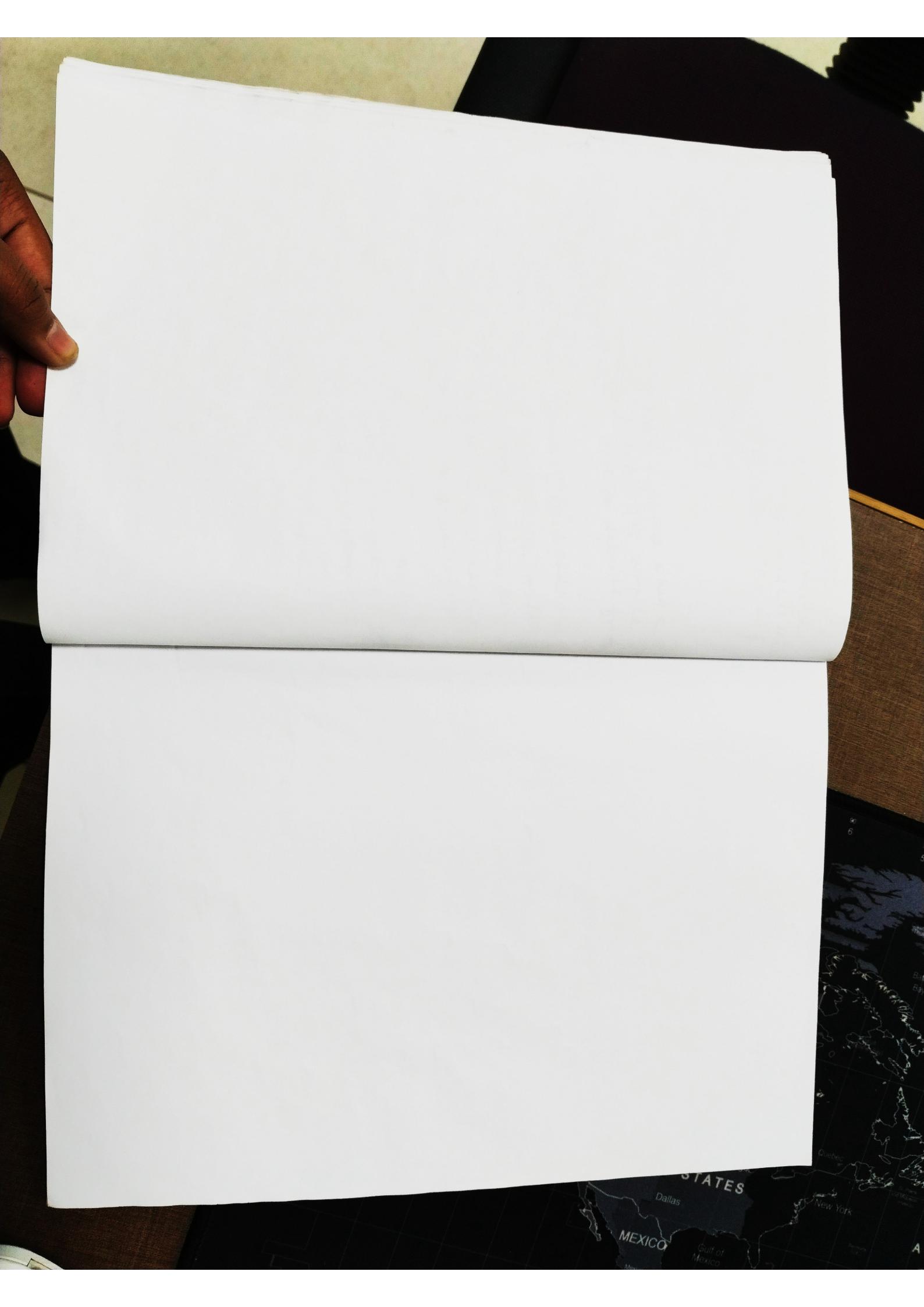
→ Bias:

• Bias is an error that happens when a machine learning model is too simple and does not learn enough details from the data.

• E.g.: It is like assuming all birds say only the same and flew, so the model fails to recognize big birds like ostriches or penguins that can't fly and get biased with predictions.

• These assumptions make the model easier to build but may prevent it from capturing the underlying complexities of the data.

• High bias typically



Resampling

→ Resampling is a technique used to adjust the class distribution of an imbalanced dataset. The class distribution of machine learning models which helps to prevent "machine learning models from being biased towards the majority class." → This is typically achieved through over-sampling (increasing the minority class) or under-sampling (decreasing the majority class).

→ General equation for resampling is

$$\text{Rare} + \kappa = p(\text{Record} + \kappa)$$

Where, $\kappa \rightarrow$ required no. of resampled records.

$p \rightarrow$ desired proportion of rare values in the balanced dataset.

records \rightarrow no. of records in the unbalanced data set.

rare \rightarrow current no. of rare target values.

Note: → The test dataset should never be balanced.

challenges with oversampling → aim to increase the no. of

oversampling minority class. → examples instead minority class 2.

→ The associated challenges are:

① Oversampling can lead to random oversampling, the minority class can lead to the model overfitting. The instances can lead to the model overfitted instances class, as it memorizes the repeated patterns rather than learning general patterns.

(i) Introduction to noise

- Advanced techniques like synthetic minority over sampling technique generate synthetic examples via interpolation
- The data is inherently noisy or the class overlap, these methods may introduce noisy or unrealistic samples that do not accurately represent the real world minority class distribution, potentially degrading performance.

(ii) Increased computation cost

- Creating a larger training dataset increases memory usage and the computational time required to train the model.

(iii) False confidence

- Models trained on synthetically balanced data might appear to have high performance metric during validation but can fail when applied to real world imbalanced scenario as the synthetic data did not accurately capture the true data distribution.

Challenges with undersampling

- Undersampling aims to reduce the number of samples in the minority class, the associated challenges are:

- 1) Loss of Information: The most significant draw back of undersampling is the potential

Loss of valuable information

- The loss of valuable and informative data from the majority class
- Randomly removing instances may discard crucial patterns that are important for the model to learn the decision boundary, leading to underfitting.

(i) Risk of Bias

- If the removed minority samples were critical or represented unique sub patterns, the resulting model may be biased and fail to generalize well to unseen data.

(ii) Poor Minority Class Recognition

- While undersampling helps balance the classes, reducing the overall dataset size might lead the model with insufficient data (especially if the original minority class was not small) to learn meaningless patterns in the minority class effectively.

(iii) Difficulty in Generalization

- Undersampling might create an artificial distribution that is too different from real-world conditions, hindering the model's ability to make accurate predictions in practice.

UNITED STATES
Chicago
Dallas
Toronto
Quebec

Example

over sampling - 100000 non fraudulent records

1000 - fraudulent records

We want to make 25% of the records
fraudulent with over sampling.

$$Rate + \kappa = P(\text{Fraudulent})$$

$$\Rightarrow 1000 + \kappa = 0.25 \times (100000 + \kappa)$$

$$\Rightarrow 1000 + \kappa = 25000 + 0.25\kappa$$

$$\Rightarrow 0.75\kappa = 24000$$

$$\Rightarrow \kappa = 24000 / 0.75 = 32000$$

Total rare records after over sampling

$$= 1000 + 32000 = 33000$$

(and)

Chapter - 8

Simple Linear Regression

→ We consider the modelling between variable

→ the dependent or independent variable

→ When there is only one independent variable

→ we linear regression model, the model is

generally termed as a simple linear regression

model.

→ When there are more than one independent

variable in the model, then the linear model is

termed as the multiple linear regression model.

→ Consider a simple linear regression model

$$y = \cancel{f(x)} b_0 + b_1 x + \epsilon$$

where y → dependent or study variable i.e

x → independent or exploratory variable

b_0 & b_1 → parameters of the model

regression coefficients

ϵ → error term of the regression line

↳ y - intercept of the regression line

b_1 → slope of the regression line

ϵ → unobservable error component

↳ it accounts for the failure of data

fitting to lie on a straight line and

represents the difference between

the true and observed realization

of y)

→ the goal of linear regression is to find a straight line that minimizes the error (the difference) between the observed points and the predicted values.
and one predicted value.
→ this line helps us predict the dependent variable for new, unseen data.

variable for new, unseen data.

$$y = b_0 + b_1 x$$

observed value for x_i

predicted value for x_i

$$\left\{ \begin{array}{l} \text{Random error} \\ \epsilon_i \\ \text{slope} = \tan \beta = b_1 \\ b_0 \end{array} \right\}$$

figure: Linear Regression

Equation of the Best-fit Line

→ for simple linear regression (with one independent variable), see best-fit line

represented by the equation

$$y = b_0 + b_1 x$$

Minimizing the Error - The Least Square

Method

→ To find the best fit line, we use a method called least square estimation -

the idea behind this method is to minimize the sum of squared differences between

$$y_i = b_0 + b_1 x_i + \epsilon_i \quad \text{for } i=1, 2, \dots, n \quad (2)$$

→ the least squares line is the line that minimizes the population sum of squared error (SSE),

$$\text{SSE} = \sum_{i=1}^n \epsilon_i^2 \quad (3)$$

where $\epsilon_i = \text{actual observed value} - \text{predicted value}$

$$\text{SSE} = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 \quad (4)$$

→ we may find values of b_0 and b_1 that minimizes $\sum_{i=1}^n \epsilon_i^2$ by differentiating eq (4) w.r.t. to b_0 and b_1 and setting the result equal to zero.

→ the partial derivative of eq (4) w.r.t to b_0 and b_1 are respectively -

$$\frac{\partial \text{SSE}}{\partial b_0} = \sum_{i=1}^n 2(y_i - b_0 - b_1 x_i) \cdot (-1) = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) \quad (5)$$

$$\frac{dSLEP}{db_1} = \sum_{i=1}^n 2 \cdot (y_i - b_0 - b_1 x_i) \cdot \frac{d}{db_1} (-b_1 x_i)$$

$$= -2 \sum_{i=1}^n x_i (y_i - b_0 - b_1 x_i) \quad \textcircled{G}$$

Solving eqⁿ(5) & (6) \rightarrow zero, we have

$$-\sum_{i=1}^n (y_i - b_0 - b_1 x_i) = 0$$

$$-\sum_{i=1}^n x_i (y_i - b_0 - b_1 x_i) = 0$$

$$\Rightarrow \left\{ \begin{array}{l} \sum_{i=1}^n (y_i - b_0 - b_1 x_i) = 0 \\ \sum_{i=1}^n x_i (y_i - b_0 - b_1 x_i) = 0 \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} \sum_{i=1}^n y_i - \sum_{i=1}^n b_0 - \sum_{i=1}^n b_1 x_i = 0 \\ \sum_{i=1}^n x_i (y_i - b_0 - b_1 x_i) = 0 \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} \sum_{i=1}^n y_i - n b_0 - b_1 \sum_{i=1}^n x_i = 0 \\ \sum_{i=1}^n x_i (y_i - n b_0 - b_1 \sum_{i=1}^n x_i) = 0 \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} \frac{\sum_{i=1}^n y_i}{n} - \frac{b_1 \sum_{i=1}^n x_i}{n} = \sum_{i=1}^n x_i y_i \\ \frac{\sum_{i=1}^n y_i}{n} - \frac{b_1 \sum_{i=1}^n x_i}{n} - \frac{b_1 \sum_{i=1}^n x_i}{n} + b_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \end{array} \right.$$

$$\Rightarrow \frac{\sum_{i=1}^n x_i y_i}{n} - b_1 \frac{(\sum_{i=1}^n x_i)^2}{n} + b_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i$$

$$\Rightarrow \left\{ \begin{array}{l} \cancel{n} b_0 + b_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \quad \textcircled{H} \\ b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \quad \textcircled{I} \end{array} \right.$$

$$\Rightarrow b_1 = \frac{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n}{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n}$$

Solving eqⁿ(7) and eqⁿ(8) for b_0 and b_1 , we have

$$b_0 = \sum_{i=1}^n y_i - b_1 \sum_{i=1}^n x_i \quad (\text{from eq}(7))$$

$$\Rightarrow b_0 = \frac{\sum_{i=1}^n y_i}{n} - \frac{b_1 \sum_{i=1}^n x_i}{n}$$

$$\Rightarrow b_0 = \frac{\sum_{i=1}^n y_i - b_1 \sum_{i=1}^n x_i}{n} \quad \textcircled{J}$$

From eqⁿ(8), we have b_1 putting b_0 value from eqⁿ(J) in eqⁿ(8)

Numerical problem on Simple Linear Regression

Consider the following set of points $\{(1, -1), (1, 1), (3, 2)\}$

a) Find the least square regression line for the given data points.

Ans:- We are the following equations to find the regression line coefficients.

$$b_0 = \bar{y} - b_1 \bar{x}, \quad b_1 = \frac{\sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)/n}{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n}$$

x	y	xy	x^2
-1	-1	1	1
1	1	1	1
3	2	6	9
$\sum x = 2$	$\sum y = 2$	$\sum xy = 9$	$\sum x^2 = 14$

Now, we can calculate the value of b_0 and b_1 (regression coefficient)

$$b_1 = \frac{\sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)/n}{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n}$$

$$= \frac{9 - (2 \times 2)/3}{(4 - (2)^2/3} = \frac{23}{28} = 0.605$$

$$b_0 = \bar{y} - b_1 \bar{x} = \frac{2}{3} - \frac{23}{28} \times \frac{2}{3} = 0.263$$

$$\text{So equation of regression line is } y = 0.263 + 0.605x \quad (\text{Ans})$$

Record No.	x	y	xy	x^2
1	43	99	4257	1849
2	21	65	1365	491
3	25	79	1975	625
4	42	75	3150	1764
5	57	87	4959	3249
6	59	81	4799	3481

Ans:-

Record No.	x	y	xy	x^2
1	43	99	4257	1849
2	21	65	1365	491
3	25	79	1975	625
4	42	75	3150	1764
5	57	87	4959	3249
6	59	81	4799	3481

Record No.	x	y	xy	x^2
1	43	99	4257	1849
2	21	65	1365	491
3	25	79	1975	625
4	42	75	3150	1764
5	57	87	4959	3249
6	59	81	4799	3481

Pind out the regression line.

$$b_1 = \frac{\sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)/n}{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n}$$

$$= \frac{11409 - (247)^2/6}{(247)^2/6 - ((247) \times 486)/6} = 0.385$$

$$b_0 = \bar{y} - b_1 \bar{x} = \frac{486}{6} - 0.385 \times \frac{247}{6} = 65.15$$

$$\text{So equation of regression line is } y = b_0 + b_1 x = 65.15 + 0.385x \quad (\text{Ans})$$

Calculation of SSE , SSR , SST , r^2

Sum of Squared Error (SSE)

- * It is the sum of the squared differences between each observed and the predicted value.

$$\text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

y_i → actual value

\hat{y}_i → predicted value

Sum of Squared Regression (SSR)

- * It is the sum of the squared differences between the predicted value and the mean of the dependent variable.

* It can be defined as

$$\text{SSR} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

\bar{y} → mean of dependent variable.

Sum of Squared Total (SST)

- * It is the sum of the squared differences between the observed dependent variables and the overall mean.

* It can be defined as

$$\text{SST} = \sum_{i=1}^n (y_i - \bar{y})^2$$

y_i → actual value for i^{th} data point
 \hat{y}_i → predicted value for i^{th} data point of variable y

Coefficient of Determination (r^2)

- * It explains how much a dependent variable varies when the independent variable changes.

* It can be defined as

$$r^2 = \frac{\text{SSR}}{\text{SST}}$$

Numerical problem on SSE , SSR and SST

Subject	$x = \text{Time}$	$y = \text{Distance}$	Predicted scores
1	2	10	16
2	2	11	16
3	3	12	12
4	4	13	14
5	4	14	14
6	5	15	16
7	6	20	18
8	7	18	20
9	8	22	22
10	9	25	24

Subject	$x = \text{Time}$	$y = \text{Distance}$	Predicted scores	Error in prediction, $(\hat{y} - y)^2$
1	2	10	10	0
2	2	11	16	25
3	3	12	12	0
4	4	13	14	1
5	4	14	14	0
6	5	15	16	1
7	6	20	18	4
8	7	18	20	4
9	8	22	22	0
10	9	25	24	1

$$\text{SSE} = \sum_{i=1}^n (y_i - \bar{y})^2 = 12$$

or

Subject	$x = \text{Time}$	$y = \text{Distance}$	Predicted	$(y - \bar{y})(\hat{y} - \bar{y})$	$(\hat{y} - \bar{y})^2$
1	2	10	10	-6	36
2	2	11	10	-6	36
3	3	12	12	-4	16
4	4	13	14	-3	9
5	4	14	14	-2	4
6	5	15	16	0	0
7	6	20	18	2	4
8	7	18	20	4	16
9	8	22	22	0	36
10	9	25	24	8	64

\bar{y}	$\text{SSE} = \sum_{i=1}^n (y - \bar{y})^2 = 216$
12	$= 12 + 216 = 228$

$$r^2 = \frac{\text{SST}}{\text{SST} + \text{SSE}} = \frac{12}{228} = 0.0526$$

or, SST can also be calculated as

$$\boxed{\text{SST} = \text{SST} + \text{SSE}}$$

Subject	$x = \text{Time}$	$y = \text{Distance}$	\bar{y}	$(y - \bar{y})$	$(\hat{y} - \bar{y})$	$(y - \bar{y})(\hat{y} - \bar{y})$	$(\hat{y} - \bar{y})^2$
1	2	10	10	-8	36		
2	2	11	10	-5	25		
3	3	12	10	-4	16		
4	4	13	10	-3	9		
5	4	14	10	-2	4		
6	5	15	10	-1	1		
7	6	20	10	4	16		
8	7	18	10	2	4		
9	8	22	10	6	36		
10	9	25	10	9	81		

$$\bar{y} = 16, \quad \text{SST} = \sum_{i=1}^n (y - \bar{y})^2 = 228$$

