

Chapter 2 Categorical Data Analysis

Section 2.1: Rates

Categorical Variables

- Ordinal:** Categories with natural ordering (e.g., stone size: small/large).
- Nominal:** Categories without ordering (e.g., gender: male/female).

Rate Definition

- Rate measures the proportion or percentage for a category:

$$\text{rate}(X) = \frac{\text{Count of } X}{\text{Total Count}}$$

- Properties:
 - $0 \leq \text{rate}(X) \leq 1$ (as a fraction)
 - $0\% \leq \text{rate}(X) \leq 100\%$ (as a percentage)
 - Rates provide a normalized comparison for categorical data.

Example (Treatment Outcomes)

- Success rate:

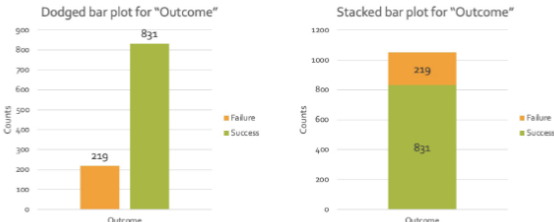
$$\text{rate}(\text{Success}) = \frac{831}{1050} = 0.791 \text{ (79.1\%)}$$

- Failure rate:

$$\text{rate}(\text{Failure}) = \frac{219}{1050} = 0.209 \text{ (20.9\%)}$$

Plots

- Dodged Bar Plot:** Side-by-side comparison of categories.
- Stacked Bar Plot:** Combined view showing percentages or proportions.



Example

Treatment\Outcome	Success	Failure	Row Total
X	542	158	700
Y	289	61	350
Total	831	219	1050

- What proportion of patients *were given treatment Y and had an unsuccessful outcome*?
 - Refers to the joint rate/proportion/percentage.

$$\text{rate}(\text{Unsuccessful and Y}) = \frac{61}{1050} = 0.0581$$

- What proportion of patients *given treatment Y had an unsuccessful outcome*?
 - Refers to the conditional rate/proportion/percentage.

$$\text{rate}(\text{Unsuccessful} \mid Y) = \frac{61}{350} = 0.174$$

Section 2.2: Association

Definition of Association

- Two variables are **associated** if the presence/absence of one variable changes the rate of another.
 - Positive association: $\text{rate}(A \mid B) > \text{rate}(A \mid NB)$.
 - Negative association: $\text{rate}(A \mid B) < \text{rate}(A \mid NB)$.
 - Note: *NB* means Not *B*

Positive association between A and B	Negative association between A and B
$\text{rate}(A \mid B) > \text{rate}(A \mid NB)$	$\text{rate}(A \mid B) < \text{rate}(A \mid NB)$
$\text{rate}(B \mid A) > \text{rate}(B \mid NA)$	$\text{rate}(B \mid A) < \text{rate}(B \mid NA)$
$\text{rate}(NA \mid NB) > \text{rate}(NA \mid B)$	$\text{rate}(NA \mid NB) < \text{rate}(NA \mid B)$
$\text{rate}(NB \mid NA) > \text{rate}(NB \mid A)$	$\text{rate}(NB \mid NA) < \text{rate}(NB \mid A)$

Example

Based on the previous data

Example

Treatment\Outcome	Success	Failure	Row Total
X	542	158	700
Y	289	61	350
Total	831	219	1050

- What proportion of patients *were given treatment Y and had an unsuccessful outcome?*
 - Refers to the joint rate/proportion/percentage.

$$rate(\text{Unsuccessful and Y}) = \frac{61}{1050} = 0.0581$$

- What proportion of patients *given treatment Y had an unsuccessful outcome?*
 - Refers to the conditional rate/proportion/percentage.

$$rate(\text{Unsuccessful} \mid Y) = \frac{61}{350} = 0.174$$

- Treatment X success rate:

$$rate(A \mid B) = rate(\text{Success} \mid X) = \frac{542}{700} = 0.774 \text{ (77.4\%)}$$

- Treatment Y success rate:

$$rate(A \mid \text{NB}) = rate(\text{Success} \mid Y) = \frac{289}{350} = 0.826 \text{ (82.6\%)}$$

- Since

$$rate(A \mid B) < rate(A \mid \text{NB})$$

- The success of treatment is **negatively associated with treatment X** because the presence of A when B is present is weaker compared to when B is absent, indicating fewer successful treatments under treatment X.
- Conversely, the success of treatment is **positively associated with treatment Y** because it shows more successful treatments compared to treatment X.

Section 2.3: Two Rules on Rates

Symmetry Rule

- $rate(A \mid B) > rate(A \mid \text{NB}) \iff rate(B \mid A) > rate(B \mid \text{NA})$
- $rate(A \mid B) < rate(A \mid \text{NB}) \iff rate(B \mid A) < rate(B \mid \text{NA})$
- $rate(A \mid B) = rate(A \mid \text{NB}) \iff rate(B \mid A) = rate(B \mid \text{NA})$

	B	Not B	Row Total
A	w	x	$w + x$
Not A	y	z	$y + z$
Column Total	$w + y$	$x + z$	$w + x + y + z$

Based on the first rule

$$\begin{aligned}
 rate(A \mid B) > rate(A \mid \text{NB}) &\iff rate(B \mid A) > rate(B \mid \text{NA}) \\
 \frac{w}{w+y} > \frac{x}{x+z} &\iff \frac{w}{w+x} > \frac{y}{y+z} \\
 w(x+z) > x(w+y) &\iff w(y+z) > y(w+x) \\
 wx + wz > xw + xy &\iff wy + wz > yw + yx \\
 wz > xy &\iff wz > xy
 \end{aligned}$$

Example

Based on the previous data

Example

Treatment\Outcome	Success	Failure	Row Total
X	542	158	700
Y	289	61	350
Total	831	219	1050

- What proportion of patients *were given treatment Y and had an unsuccessful outcome?*
 - Refers to the joint rate/proportion/percentage.

$$rate(\text{Unsuccessful and Y}) = \frac{61}{1050} = 0.0581$$

- What proportion of patients *given treatment Y had an unsuccessful outcome?*
 - Refers to the conditional rate/proportion/percentage.

$$rate(\text{Unsuccessful} \mid Y) = \frac{61}{350} = 0.174$$

$$rate(A \mid B) = rate(\text{Success} \mid X) = 0.774, rate(A \mid \text{NB}) = rate(\text{Success} \mid Y) = 0.826$$

Using symmetry:

$$rate(B \mid A) = rate(X \mid \text{Success}) = 0.652 < rate(B \mid \text{NA}) = rate(X \mid \text{Failure}) = 0.721$$

Basic Rule

- Overall rate (A) is will always lie between $\text{rate}(A \mid B)$ and $\text{rate}(A \mid \text{NB})$:

$$\text{rate}(A \mid B) \leq \text{rate}(A) \leq \text{rate}(A \mid \text{NB})$$

Consequence 1:

- The closer $\text{rate}(B)$ is to 100%, the closer $\text{rate}(A)$ is to $\text{rate}(A \mid B)$.

Consequence 2:

- If $\text{rate}(B) = 50\%$, then $\text{rate}(A) = \frac{1}{2} [\text{rate}(A \mid B) + \text{rate}(A \mid \text{NB})]$

Consequence 3:

- If $\text{rate}(A \mid B) = \text{rate}(A \mid \text{NB})$, then $\text{rate}(A) = \text{rate}(A \mid B) = \text{rate}(A \mid \text{NB})$

Example:

If $\text{rate}(A \mid B) = 75\%$, $\text{rate}(A \mid \text{NB}) = 55\%$, overall rate $\text{rate}(A) \in [55\%, 75\%]$.

Section 2.4: Simpson’s Paradox

Definition

- A phenomenon where a trend reverses when data is combined vs. when divided into subgroups.
- Simpson’s Paradox highlights the importance of subgroup analysis.

Example

Based on the previous data

Example

Treatment\Outcome	Success	Failure	Row Total
X	542	158	700
Y	289	61	350
Total	831	219	1050

- What proportion of patients *were given treatment Y and had an unsuccessful outcome*?
 - Refers to the joint rate/proportion/percentage.

$$\text{rate}(\text{Unsuccessful and Y}) = \frac{61}{1050} = 0.0581$$

- What proportion of patients *given treatment Y had an unsuccessful outcome*?
 - Refers to the conditional rate/proportion/percentage.

$$\text{rate}(\text{Unsuccessful} \mid Y) = \frac{61}{350} = 0.174$$

- Combined: $\text{rate}(\text{Success} \mid X) < \text{rate}(\text{Success} \mid Y)$.
- Split by stone size:
 - Large stones: $\text{rate}(\text{Success} \mid X) > \text{rate}(\text{Success} \mid Y)$.
 - Small stones: $\text{rate}(\text{Success} \mid X) > \text{rate}(\text{Success} \mid Y)$.

Section 2.5: Confounders

Definition

A confounder is a third variable associated with both the independent and dependent variables, affecting the observed relationship. Association does not imply causation; consider confounders.

Example

Based on the previous data

Example

Treatment\Outcome	Success	Failure	Row Total
X	542	158	700
Y	289	61	350
Total	831	219	1050

- What proportion of patients *were given treatment Y and had an unsuccessful outcome*?
 - Refers to the joint rate/proportion/percentage.

$$\text{rate}(\text{Unsuccessful and Y}) = \frac{61}{1050} = 0.0581$$

- What proportion of patients *given treatment Y had an unsuccessful outcome*?
 - Refers to the conditional rate/proportion/percentage.

$$\text{rate}(\text{Unsuccessful} \mid Y) = \frac{61}{350} = 0.174$$

- **Stone size** is a confounder:
 - Large stones are more likely treated with X ($\text{rate}(\text{Large} \mid \text{X}) = 75.1\%$).
 - Large stones have lower success rates overall ($\text{rate}(\text{Success} \mid \text{Large}) = 71.9\%$).

Addressing Confounders

1. **Slicing**: Analyze subgroups (e.g., by stone size).
2. **Randomization**: Randomly assign patients to treatments to equalize confounders.

