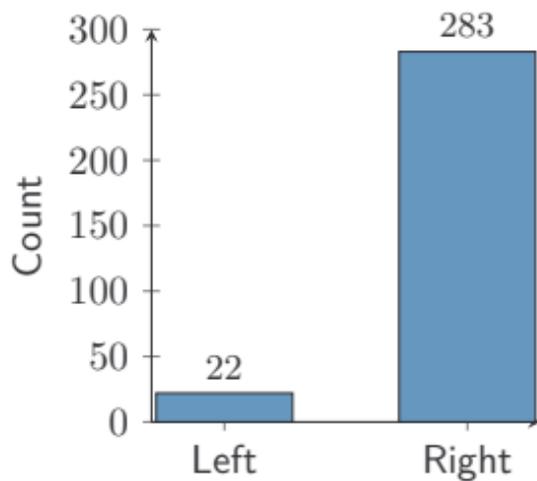


Chapter 6

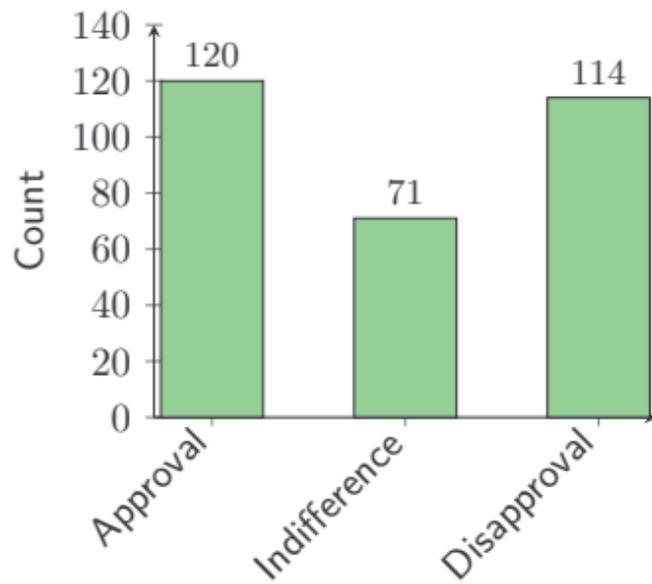
Two-Way Tables & Risk

Visualizing Categorical Data: Bar Charts

Handedness Distribution



Pineapple Preference Distribution



From Quantitative to Categorical Associations

Previously: We studied how two quantitative variables can be associated (e.g., scatterplots, correlation).

- Example: Height vs. weight, hours studied vs. exam score
- Tools: scatterplots, regression, correlation coefficient

Now: We try to answer similar questions when both variables are categorical.

- Are two categorical variables associated?
- How do we summarize and visualize their relationship?
- What tools can we use to compare groups?

Categorical Data

The following are examples of questions we might answer:

- Does a new vaccine actually **reduce the risk** of getting sick?
- Is a hiring process **fair across demographics**?
- Do left-handed people really prefer different foods than right-handed people?
- When a hospital reports better outcomes, is it truly better- or does it just treat **easier cases**?

 **Key Point:** To answer these questions, we need tools for analyzing two categorical variables **simultaneously**.

The Framework: Exposure and Outcome

In this chapter, we generally view our data as having two distinct roles. For every person (or observation) in our dataset, we track:

1. Exposure (X)

- Also called the **Explanatory Variable**.
- This defines the **groups** we are comparing. It does **not** need to be the row variable of the contingency table, but it frequently is.
- *Examples: Vaccination status, Handedness, Study Method.*

2. Outcome (Y)

- Also called the **Response Variable**.
- This is the **result** we are measuring.
- *Examples: Getting the flu, Pineapple preference, Passing the exam.*

 **Key Point:** Our goal is to determine if the likelihood of the **Outcome** depends on the **Exposure**.

The Core Question, Restated

Throughout this chapter, we will ask variations of this question:

Does the proportion with a given outcome differ between exposure groups?

If YES (proportions differ):

The variables are **associated**.

Exposure tells us something about outcome.

We want to quantify the difference.

If NO (proportions are the same):

The variables are **independent**.

Exposure tells us nothing about outcome.

Intended Learning Outcomes

- Construct and interpret two-way tables
- Calculate marginal distributions
- Calculate conditional distributions
- Distinguish between marginal and conditional distributions
- Calculate relative risk and odds ratio
- Interpret RR and OR in context
- Recognize Simpson's paradox
- Identify confounding variables

PART 1

Two-Way Tables & Distributions

Analyzing Relationships Between Categorical Variables

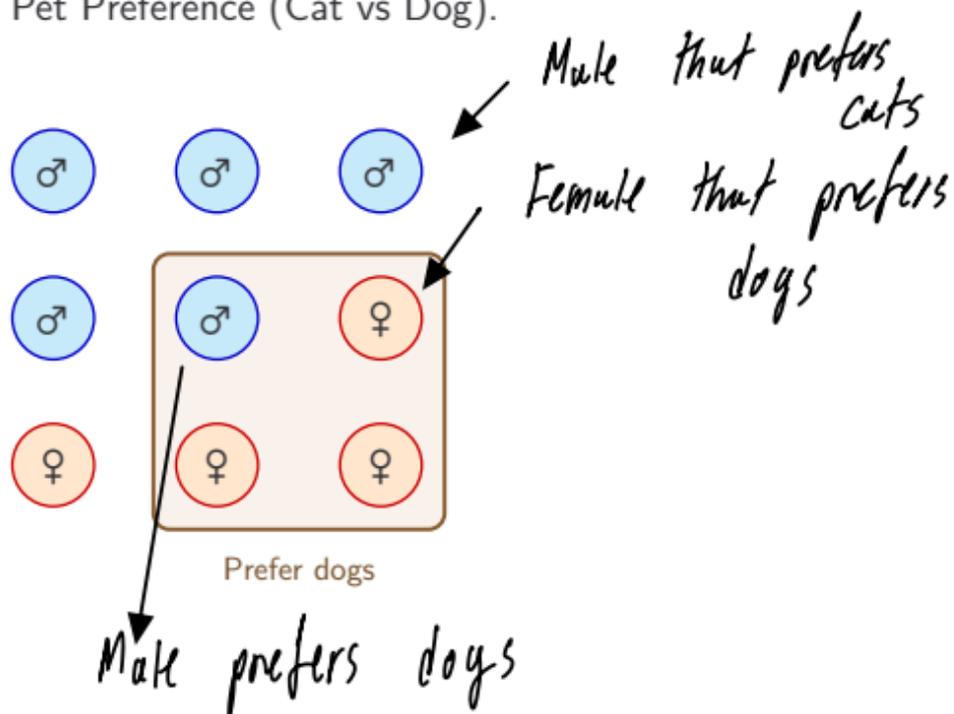
- So far, we have analyzed **one categorical variable** at a time using bar graphs and pie charts.
- Many real-world questions involve **two categorical variables**:
 1. Does **treatment type** affect **recovery status**?
 2. Is **political affiliation** related to **opinion on an issue**?
 3. Does **vaccination status** relate to **disease outcome**?

 **Key Point:** To study two categorical variables together, we use **two-way tables**.

Visualizing the Data

Example: Data collected on 9 individuals, tracking

- **Variables:** Gender (σ , ρ) and Pet Preference (Cat vs Dog).



Two-Way Contingency Tables

Two-Way Table

A **two-way table** (or contingency table) displays counts for each combination of values of two categorical variables.

- One variable defines the **rows**
- The other variable defines the **columns**
- Each cell contains a **count**

The diagram shows a 2x2 contingency table with the following structure:

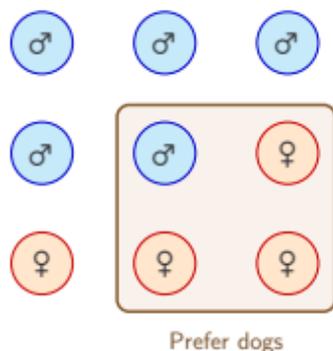
	Col 1	Col 2	Total
Row 1	count	count	Σ
Row 2	count	count	Σ
Total	Σ	Σ	n

Annotations and labels:

- Column variable:** A blue arrow points to the column headers.
- Row variable:** A blue arrow points to the row headers.
- Row totals:** A green arrow points to the Σ values in the rightmost column.
- Column totals:** A green arrow points to the Σ values in the bottom row.
- Handwritten note:** An arrow points from the top-right cell to the text: "# of observations in category 1 of row variable and category 1 of the column variable".

Example 6.0: Constructing a Table

Context: Recall our toy example with 9 individuals, tracking Gender (σ , ϕ) and Pet Preference (Cat vs Dog). Construct the two-way contingency table showing the relationship between Gender and Pet Preference.

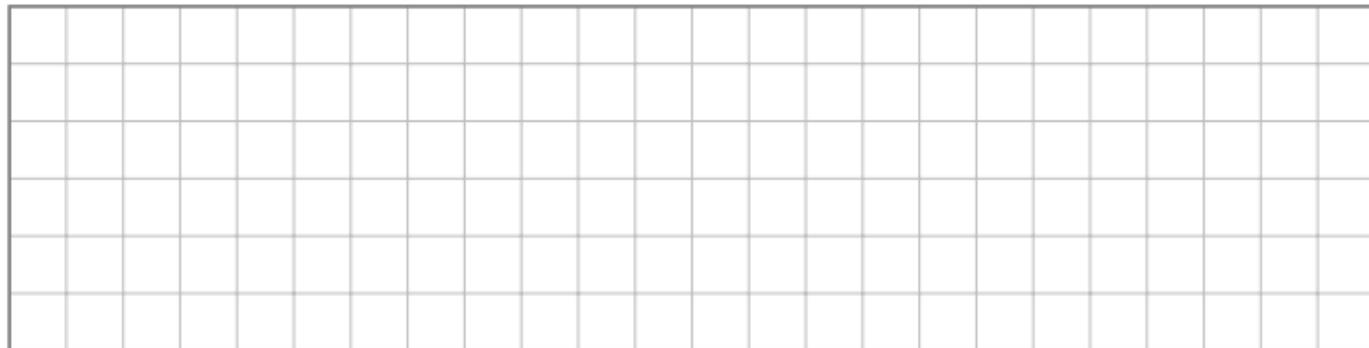


		Pet Preference	
		Dog	Cat
Gender	M	1	4
	F	3	1

Example 6.1: Pineapple on Pizza

Context: Last term, I asked DS1000A students about their handedness and pineapple-on-pizza preference. The results are summarized below:

Pineapple?	Left-handed	Right-handed	Total
Approval	11	109	120
Indifference	5	66	71
Disapproval	6	108	114
Total	22	283	305



Uncertainty: Why We Need Proportions

Key insight: When we select an individual from a population, we **don't know in advance** what their exposure or outcome will be.

- Pick a random student from this class: are they left-handed or right-handed?
- Pick a random patient from a hospital: did they receive Treatment A or B?
- Pick a random website visitor: did they click the "Buy" button?

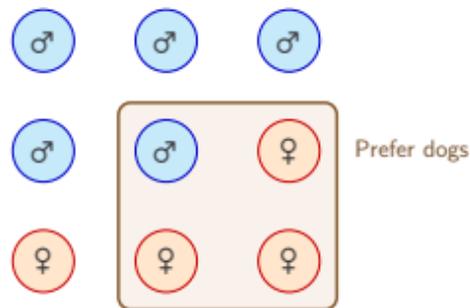
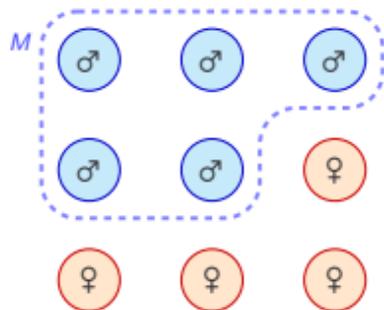
 **Key Point:** Because we face **uncertainty**, we describe patterns using **proportions**:
"What fraction of people in Group A have Outcome Y?"

These proportions tell us **how likely** different outcomes are within different groups. Comparing them tells us whether exposure and outcome are **associated**.

Marginal Distributions

Marginal Distribution

The **marginal distribution** of a variable describes the distribution of that variable alone, *ignoring the categories of the other variable*



a) Find marginal distⁿ of Gender

b) Pet preference

Gender	Probability	Pet	Probability
M	5/9	Dog	4/9
F	4/9	Cat	5/9

Marginal Distributions

Marginal Distribution

The **marginal distribution** of a variable describes the distribution of that variable alone, *ignoring the other variable*

		Pet Preference	
		Dog	Cat
Gender	M	1	4
	F	3	1

		Pet Preference	
		Dog	Cat
Gender	M	1	4
	F	3	1

a) Find marginal distⁿ of Gender

Gender	Probability
M	$(1+4)/9$
F	$(3+1)/9$

b) Marginal distⁿ for pet preference

Pet	Probability
Dog	$(1+3)/9$
Cat	$(4+1)/9$

Notation: Conditional Probability

Denominator

Before we calculate conditional distributions, let's introduce a standard piece of notation that makes it easier to write down what we mean.

The "Given" Bar (|)

In probability and statistics, the vertical bar | is read as "given."

$$P(A | B)$$

This translates to: "The probability of outcome A **given** that we already know condition B is true."

$P(\text{Task completion} | \text{Notif. on})$
numerator (i.e. subset of B that we should consider)
Denominator
population under consideration

- **Marginal:** $P(\text{Approval})$ means "What % of **everyone** approves?"
- **Conditional:** $P(\text{Approval} | \text{Left-Handed})$ means "If we look **only** at Left-Handed people, what % of **them** approve?"

Conditional Distributions

Conditional Distribution

The **conditional distribution** of a variable describes the distribution **within a specific group** defined by the other variable. There is a separate distribution for each value of the conditioning variable.

Pineapple?	Left-handed	Right-handed
Approval	11	109
Indifference	5	66
Disapproval	6	108
Total	22	283

Pineapple?	Left-handed	Right-handed	Total
Approval	11	109	120
Indifference	5	66	71
Disapproval	6	108	114
Total	22	283	305

Find the conditional distribution of pineapple preference *given* handedness.

Left handed people $P(\text{Outcome}|\text{Left})$:

- Approval: $11/22$
- Indifference: $5/22$
- Disapproval: $6/22$

Right handed people $P(\text{Outcome}|\text{Right})$:

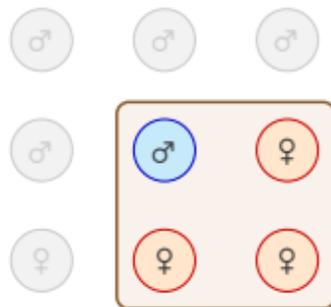
- Approval: $109/283$
- Indifference: $66/283$
- Disapproval: $108/283$

Visualizing Conditional Probability

Conditional means we **zoom in** on a specific group (e.g., dog enthusiasts) and ignore everyone else.

Example: Find the conditional distribution of gender among dog enthusiasts.

Gender given dog enthusiast dist ^b	
Gender	Prob [#]
M	1/4
F	3/4



→ Dog enthusiasts
The "faded" nodes no longer count.

Reading and Computing Conditional Proportions

 **Key Point:** $P(A | B) =$ “probability of A given B ” $= \frac{\# \text{ with both } A \text{ and } B}{\# \text{ with } B}$

The key: The denominator is **restricted** to only those with condition B .

Example Table:

	Flu	No Flu	Total
Vaccinated	26	3874	3900
Unvaccinated	70	3830	3900

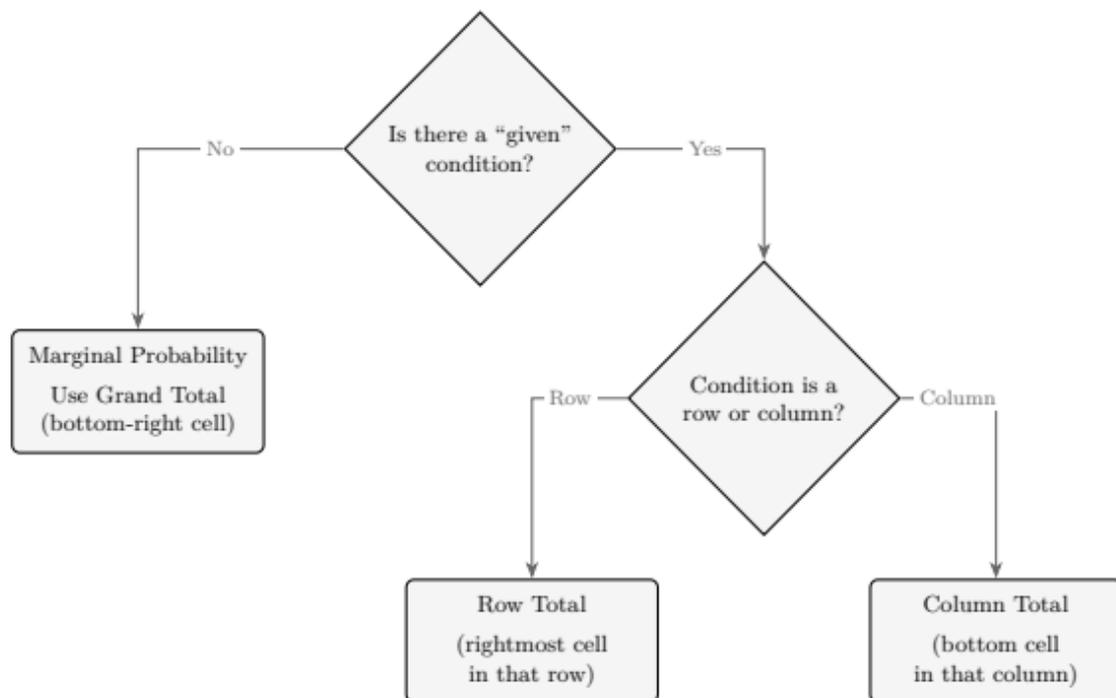
Computing:

$$P(\text{Flu} | \text{Vaccinated}) = \frac{26}{3900}$$

$$P(\text{Flu} | \text{Unvaccinated}) = \frac{70}{3900}$$

 **Caution:** The denominator changes depending on what comes after the “|” symbol.

Which Denominator to Use?



Example 6.1: Conditional Distribution of Handedness

Context: Using the DS1000A pineapple-on-pizza data:

Pineapple?	Left-handed	Right-handed	Total
Approval	11	109	120
Indifference	5	66	71
Disapproval	6	108	114
Total	22	283	305

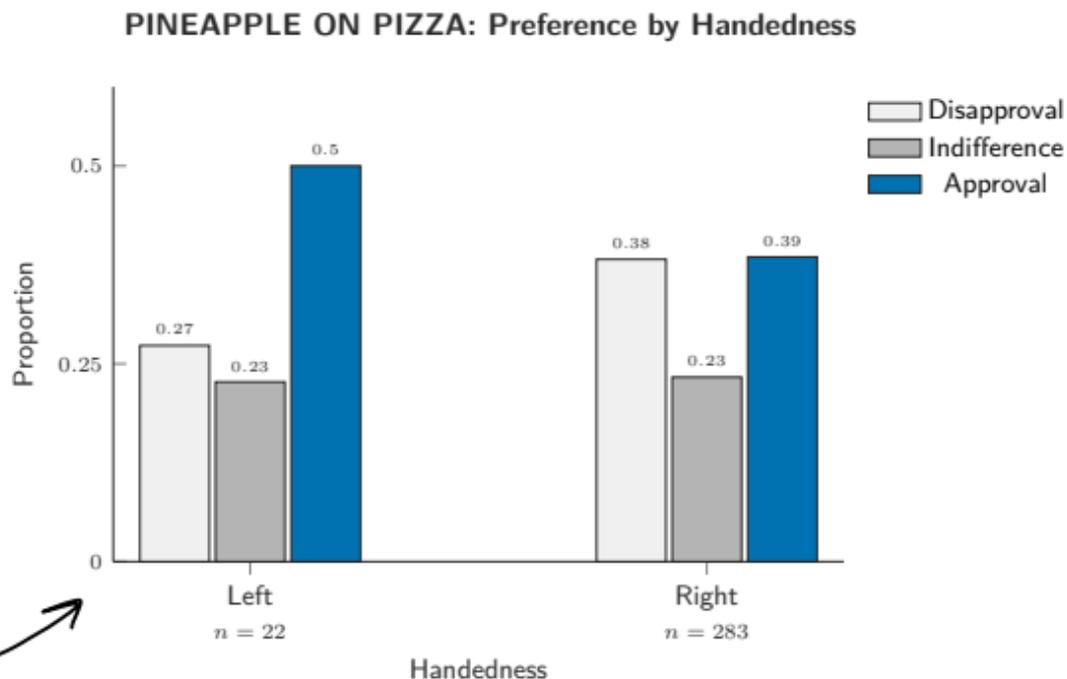
Find the conditional distribution of handedness **given** pineapple preference of **indifference**.

<u>Handedness</u>	<u>Probⁿ</u>
Right	$5/71 \doteq 0.07$
Left	$66/71 \doteq 0.93$

Visualizing Conditional Distributions: Grouped Bar Chart

Grouped Bar Chart

- Bars grouped by one variable
- Colors represent the other variable
- Easy to compare within groups



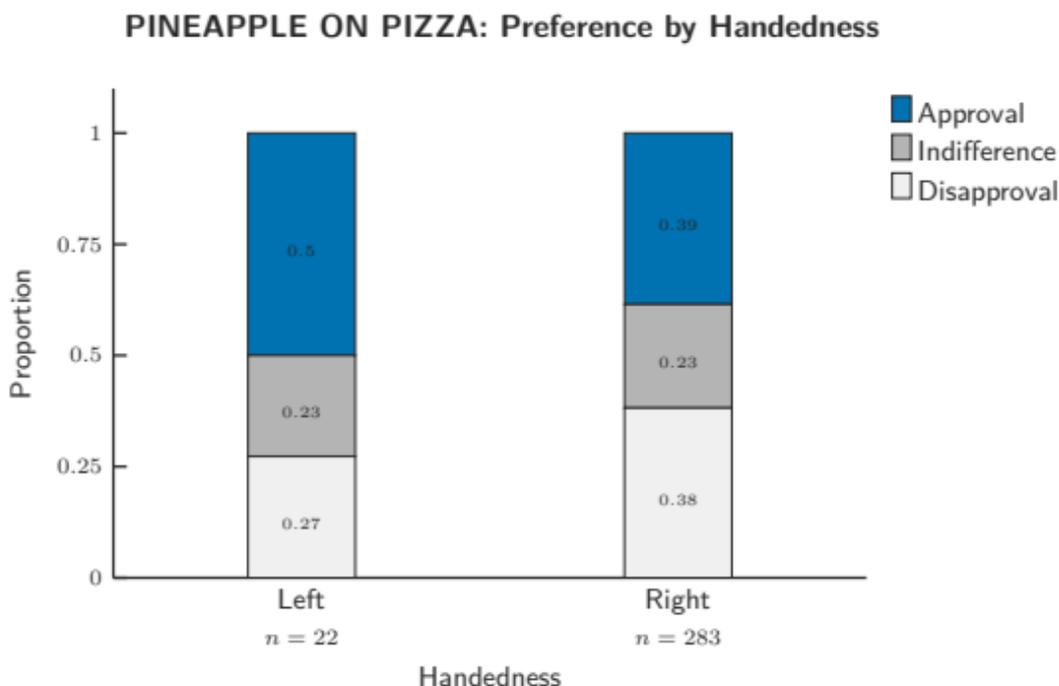
Conditional distⁿ Pineapple | left handed



Visualizing Conditional Distributions: Stacked Bar Chart

Segmented (Stacked) Bar Chart

- Each bar represents 100%
- Segments show proportions
- Easy to compare proportions

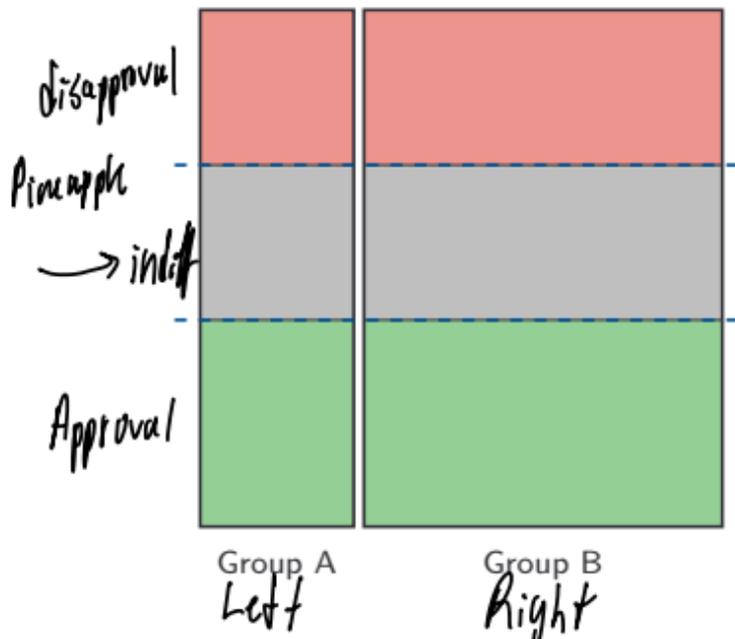


Visualizing Associations: Mosaic Plots

A **mosaic plot** shows cell counts as rectangular areas. Width = column proportion, Height = row proportion within column.

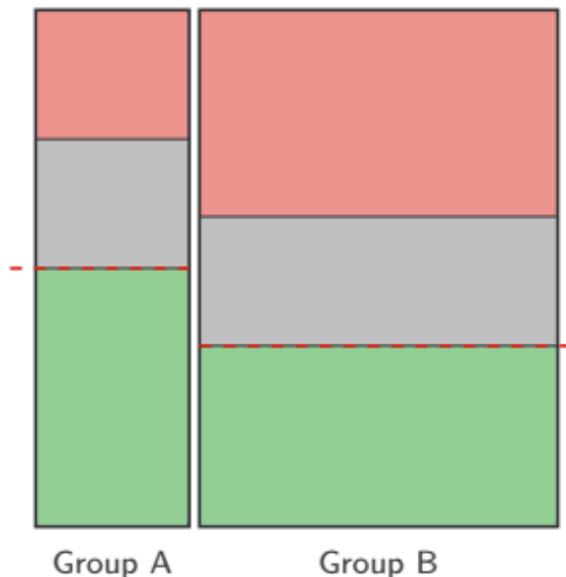
Independent Variables

(Equal row heights across columns)



Associated Variables

(Different row heights across columns)



Independence of Categorical Variables

Independence

Two categorical variables are **independent** if the conditional distribution of one variable is the *same regardless of the level* of the other variable.

Key Point: Checking for independence:

Compare conditional distributions across groups. If they are (approximately) equal, the variables may be independent.

What Would Independence Look Like?

Hypothetical: Suppose pineapple preference was truly independent of handedness. Then the conditional distributions would be (approximately) equal:

	Left-handed	Right-handed	Marginal
Approval	39.3%	39.3%	39.3%
Indifference	23.3%	23.3%	23.3%
Disapproval	37.4%	37.4%	37.4%

 **Key Point:** Under independence, every conditional distribution equals (approximately) the **marginal distribution**.

 **Caution:** Compare this with the **actual data**: left-handed approval was 50.0% vs. right-handed 38.5%. These are quite far apart, suggesting the variables are **not** independent.

Example 6.3: Instrument and Music Genre

Context: A survey of 120 students asked about their primary instrument and favorite music genre. Are instrument and favorite genre independent?

Instrument	Classical	Pop	Total
Piano	25	15	40
Guitar	10	30	40
Violin	20	20	40
Total	55	65	120

What should B be?

Conditioning on genre		B = "Pop"	
Instrument	Probability	Instrument	Probability
Piano	$25/55 \approx 0.45$	Piano	$15/65 \approx 0.23$
Guitar	$10/55 \approx 0.18$	Guitar	$30/65 \approx 0.46$
Violin	$20/55 \approx 0.36$	Violin	$20/65 \approx 0.31$

Genre and Instrument are dependent, they are not independent.

Example 6.4: Movie Genre and Snack Preference

Context: A movie theater surveyed 500 customers about their movie genre choice and snack preference.

Genre	Popcorn	Candy	Total
Action	84	56	140
Comedy	126	84	210
Drama	90	60	150
Total	300	200	500

Are genre and snack preference independent?

Condition on "snacks":

Genre	Popcorn		Genre	Candy	
	Count	Prob ^y		Count	Prob ^y
Action	84	$84/300 = 0.28$	Action	56	$56/200 = 0.28$
Comedy	126	$126/300 = 0.42$	Comedy	84	$84/200 = 0.42$
Drama	90	$90/300 = 0.30$	Drama	60	$60/200 = 0.30$

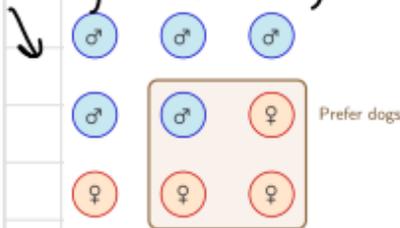
As Genre | Snack = "Popcorn" is equal to Genre | Snack = "Candy", these variables are independent.

Last Time



- Calculating a conditional prob^y or proportion

$$P(M) = \frac{5}{9} \quad P(F) = \frac{4}{9}$$



$$P(A|B) = \frac{\# \text{ of observations in categories A and B}}{\# \text{ of observations in category B}}$$

$$\left. \begin{aligned} P(M|C) &= \frac{\# \text{ men prefer cats}}{\# \text{ cat enthusiasts}} = \frac{4}{5} \\ P(F|C) &= \frac{\# \text{ women prefer cats}}{\# \text{ cat enthusiasts}} = \frac{1}{5} \end{aligned} \right\} \begin{array}{l} \text{Conditional} \\ \text{dist}^n \text{ of} \\ \text{gender given} \\ \text{cat enthusiasts} \end{array}$$

- Two categorical vars. X and Y are independent if conditional distⁿ of X given $Y=y$ is the same as the marginal distⁿ of X .

PART 2

Quantifying Associations

Beyond Independence

In Part 1, we asked a binary question: **Are these variables independent?**

- If $P(\text{Outcome}|\text{Group A}) \neq P(\text{Outcome}|\text{Group B})$, we say they are **associated**.
- But simply knowing an association *exists* is rarely enough for decision-making.

The Next Logical Questions:

1. **Direction:** Does the exposure increase or decrease the likelihood of the outcome?
2. **Strength:** How *large* is that effect? Is it negligible or substantial?

Quantifying the Association

To make decisions, we need to measure the **effect size** of the relationship.

Scenario A:

Treatment A: 1.0% recovery

Treatment B: 1.1% recovery

Associated? Yes.

Meaningful? Perhaps not.

Scenario B:

Treatment A: 1.0% recovery

Treatment B: 50.0% recovery

Associated? Yes.

Meaningful? Absolutely.

Risk Measures

Risk Measures (like Relative Risk and Odds Ratios) are statistics that quantify the **strength** of the association between two categorical variables.

Why “Risk”?

The terminology comes from **epidemiology** and **biostatistics**, where these tools were developed to study disease.

- In those fields, the “Success” outcome is often a negative event (death, infection, side effect).
- Therefore, the probability of the outcome is called the **Risk**.

Probability, Risk, and Odds

Key Terms

For a categorical variable with category c :

- **Percentage:** $\frac{\# \text{ in category } c}{\text{total } \#} \times 100\%$

- **Proportion/Probability/Risk:** $\frac{\# \text{ in category } c}{\text{total } \#}$

- **Odds:** $\frac{\# \text{ in category } c}{\# \text{ not in category } c} = \frac{\text{Probability}}{1 - \text{Probability}} =$

$$= \frac{p \cdot \# \text{ of outcomes}}{(1-p) \cdot \# \text{ of outcomes}}$$
$$= \frac{\# \text{ in cat } c}{\# \text{ not in category } c}$$

Converting Between Probability and Odds

$$\text{Odds of } A = \frac{P(A)}{1-P(A)}$$

 **Key Point:** You can convert back and forth between probability and odds:

Probability \rightarrow Odds

$$\text{Odds} = \frac{p}{1-p}$$

Example: If $p = 0.75$,

$$\begin{aligned} \text{Odds} &= \text{[input box]} \\ &= \frac{0.75}{1-0.75} \\ &= 0.75/0.25 = 3. \end{aligned}$$

Odds \rightarrow Probability

$$p = \frac{\text{Odds}}{1 + \text{Odds}}$$

Example: If Odds = 4,

$$\begin{aligned} p &= \text{[input box]} \\ &= \frac{4}{1+4} \\ &= \frac{4}{5} = 0.8 \end{aligned}$$

Example 6.5: Rare Outcomes

Context: In a population of 10,000 individuals, 6 have heterochromia (two different colored eyes). Find the percentage, proportion, probability, risk, and odds of heterochromia.

$$\text{Percentage} : 0.0006 \times 100\% = 0.06\%$$

$$\begin{aligned} \text{Proportion / Probability / Risk} : p &= \frac{\# \text{ heterochromia}}{\# \text{ observations}} \\ &= \frac{6}{10,000} = 0.0006 \end{aligned}$$

$$\text{Odds} : \frac{p}{1-p} = \frac{0.0006}{0.9994} = \frac{6}{9994} \approx 0.0006$$

When the probability is small, then odds \approx prob^{ty}.

Example 6.6: Conditional Probability

Context: A neighborhood survey on rollerblading and having children.

	Kids	No Kids	Total
Rollerblader	41	9	50
Not a rollerblader	3	9	12
Total	44	18	62

Suppose that David, Anastasia and Miles were survey participants and that:

- David is focusing on his career and does not have children.
- Anastasia has never rollerbladed.
- Miles is very tall.

Example 6.6: Neighborhood Survey (cont.)

Context: A neighborhood survey on rollerblading and having children.

David is focusing on his career and does not have children. What is the probability that David rollerblades? $R = \text{"Rollerblader"}$ $K^c = \text{"no kids"}$

	Kids	No Kids	Total
Rollerblader	41	9	50
Not a rollerblader	3	9	12
Total	44	18	62

• Find $P(R | K^c)$

= $\frac{\# \text{ obs that rollerblade and child free}}{\# \text{ obs that are child free}}$

$$= \frac{9}{9+9} = \frac{9}{18} = 0.5$$

Example 6.6: Neighborhood Survey (cont.)

R = "Rollerblader" K^c = "no kids"
 R^c = "not a rollerblader" K = "has kids"

Context: A neighborhood survey on rollerblading and having children.

	Kids	No Kids	Total
Rollerblader	41	9	50
Not a rollerblader	3	9	12
Total	44	18	62

Anastasia has never rollerbladed. What are the odds that Anastasia has children?

$$\begin{aligned} \text{Odds Anastasia has children} &= \frac{P(\text{Anastasia has children})}{P(\text{Anastasia has no children})} \\ &= \frac{P(K|R^c)}{P(K^c|R^c)} \\ P(K|R^c) &= \frac{3}{12} = 0.25 \\ \therefore \text{odds} &= \frac{0.25}{1-0.25} = \frac{0.25}{0.75} = \frac{1}{3} \approx 0.33 \end{aligned}$$

Example 6.6: Neighborhood Survey (cont.)

	Kids	No Kids	Total
Rollerblader	41	9	50
Not a rollerblader	3	9	12
Total	44	18	62

Context: A neighborhood survey on rollerblading and having children.

Miles is very tall. What is the risk of Miles rollerblading?

$$P(R) = \frac{50}{62} \doteq 0.81$$

PART 3

Comparing Groups: RR & OR

Relative Risk (RR)

The **relative risk** compares the risk (probability) in one group to the risk in another group (baseline):

$$\text{Relative Risk} = \frac{P(\text{Outcome} \mid \text{Exposed})}{P(\text{Outcome} \mid \text{Unexposed})}$$

Interpreting RR:

- RR = 1: Both groups have equal risk (no association)
- RR > 1: Exposed group has **higher** risk than baseline
- RR < 1: Exposed group has **lower** risk than baseline

Note: The baseline group is typically the **control** or **unexposed** group.

Visualizing Relative Risk: Two Towns

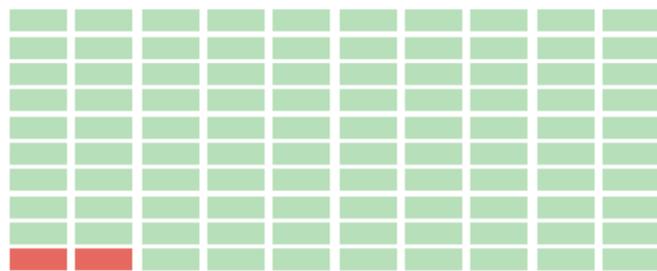
Example: Disease outbreak: 8% in Town A vs 2% in Town B \implies RR = 4.0

Town A (100 residents)



■ 8 sick (8%)

Town B (100 residents)



■ 2 sick (2%)

RR of disease in town A vs. town B is $\frac{0.08}{0.02} = 4.$

Example 6.7: Flu Vaccine Study

Context: A study comparing flu rates between vaccinated and control groups.

Group	Developed Flu	No Flu	Total
Vaccinated	26	3874	3900
Control	70	3830	3900
Total	96	7704	7800

$$RR = \frac{P(\text{Outcome} | \text{Exposed})}{P(\text{Outcome} | \text{Unexposed})}$$

Find the relative risk of developing flu for the vaccinated group compared to the control group

	Outcome	exposed	unexposed
$RR = \frac{P(\text{flu} \text{vaccinated})}{P(\text{flu} \text{control})}$		$= \frac{26}{3900}$	$\frac{70}{3900}$
		$= \frac{26}{70}$	$\approx 0,37$
Vaccinated individuals had only 37% of the risk of developing the flu compared to control group.			

Example 6.2: Pet Ownership (Relative Risk)

Context: Dog ownership by sex.

Sex	Cat	Dog	Total
Female	20	30	50
Male	40	40	80

Calculate RR of dog ownership for females compared to males.

Outcome	exposed	unexposed
$RR = \frac{P(D F)}{P(D M)} = \frac{30}{50} / \frac{40}{80} = 1.20$		
Interpretation: Women have a 20% higher risk of dog ownership than men.		
Interpretation: The risk of dog ownership for women is 1.2 that of men.		

Relative Risk: Terminology

All of the following are equivalent ways to ask for **relative risk** between two groups:

- Calculate the relative risk of dog ownership for females **compared to** males.
- Find the relative risk for dog ownership for females **relative to** males.
- Determine the RR of dog ownership, females **compared with** males.
- Estimate the relative risk for dog ownership: females **vs** males.
- Compute the relative risk for dog ownership in females **versus** males.

Odds Ratio

Odds Ratio (OR)

The **odds ratio** compares the odds in one group to the odds in another group:

$$\text{Odds Ratio} = \frac{\text{Odds in Exposed Group}}{\text{Odds in Unexposed Group}}$$

Handwritten annotations:
An arrow points from "Exposed Group" to the handwritten text "treatment".
An arrow points from "Unexposed Group" to the handwritten text "control" or "baseline".

Interpreting OR:

- OR = 1: Both groups have equal odds (no association)
- OR > 1: Exposed group has **higher** odds than baseline
- OR < 1: Exposed group has **lower** odds than baseline

Example 6.7: Flu Vaccine (Odds Ratio)

Context: Same flu vaccine study.

Group	Developed Flu	No Flu	Total
Vaccinated	26	3874	3900
Control	70	3830	3900

Calculate odds ratio of developing flu (vaccinated vs. control).

$$\text{Odds}(A) = \frac{p}{1-p}$$

$$\text{OR} = \frac{\text{odds of flu for vaccinated}}{\text{odds of flu for control}}$$

$$= \frac{26}{3874} / \frac{70}{3830} \doteq 0.36$$

$$= \frac{\# \text{ in } A}{\# \text{ outside of } A}$$

Interpretation: The odds of flu for vaccinated group are 0.36 that of the control group.

Last Time

- Risk, proportion and probability of category C:

$$P(C) := \frac{\# \text{ observations in category C}}{\# \text{ observations}}$$

- Odds of category C:

$$\text{Odds}(C) = \frac{\# \text{ observations in category C}}{\# \text{ observations outside of category C}} = \frac{(\# \text{ observations in category C} / \# \text{ observations})}{(\# \text{ observations outside of category C} / \# \text{ observations})} = \frac{P(C)}{1 - P(C)}$$

- Relative Risk = $\frac{P(\text{outcome} \mid \text{exposed group})}{P(\text{outcome} \mid \text{baseline (unexposed group)})}$

- Odds Ratio = $\frac{\text{Odds of outcome in exposed group}}{\text{Odds of outcome in baseline (unexposed group)}}$

Example 6.2: Pet Ownership (Odds Ratio)

Context: Dog ownership by sex.

Sex	Cat	Dog	Total
Female	20	30	50
Male	40	40	80
	60	70	

Calculate OR of dog ownership for females compared to males:

$$OR = \frac{\text{Odds}(F)}{\text{Odds}(M)} = \frac{30/20}{40/40} = 1.5$$

Interpretation: The odds of dog ownership are 50% higher for women than for men.

NOTE: it would be incorrect to interpret this as "women are 50% more likely to own a dog than men"

Odds and risk (probability) are different concepts. You can check that $RR = 1.20 \neq 1.50 = OR$.

Odds Ratio: Cross-Product Shortcut

For a 2×2 table, there is a shortcut to compute OR:

	Outcome Yes	Outcome No
Group 1 (Exposed)	<i>a</i>	<i>b</i>
Group 2 (Unexposed)	<i>c</i>	<i>d</i>

$$\text{OR} = \frac{a \times d}{b \times c}$$

Relative Risk vs. Odds Ratio

Relative Risk (RR)

- Compares **probabilities**
- Easier to interpret
- “X times more likely”
- Requires knowing group totals

Odds Ratio (OR)

- Compares **odds**
- Used in logistic regression
- “X times higher odds”
- Cross-product shortcut

PART 4

Simpson's Paradox

Example 6.8: Simpson's Paradox

Context: UC Berkeley Graduate Admissions (1973). Was there bias against women?

Status	# Men	# Women	Total
Admitted	1,198	557	1,755
Rejected	1,493	1,278	2,771
Total	2,691	1,835	4,526

Admission rates:

- Men: $1198 / 2691$
- Women: $557 / 1835$

Example 6.8: By Department

When we look at the data by department:

Dept	Number of Applicants		Admission Rate	
	Men	Women	Men	Women
A	825	108	62%	82%
B	560	25	63%	68%
C	325	593	37%	34%
D	417	375	33%	35%
E	191	393	28%	24%
F	373	341	6%	7%

In 4 of 6 departments, women had **higher** admission rates than men.

Understanding Simpson's Paradox

Simpson's Paradox

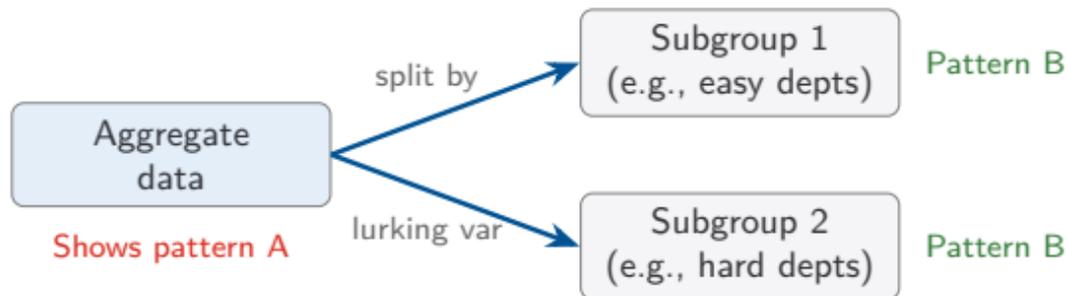
A trend that appears in aggregate data **reverses** or disappears when the data is split into subgroups.

What happened at Berkeley?

- Women applied more to competitive departments (C, D, E, F) with low admit rates
- Men applied more to less competitive departments (A, B) with high admit rates
- **Department choice** was a confounding variable!

Why Does Simpson's Paradox Happen?

The key mechanism: Groups have **unequal representation** across subgroups.



Unequal group sizes across subgroups create a misleading aggregate

Key Point: Simpson's Paradox occurs when a lurking variable is **unevenly distributed** across the groups being compared, and that lurking variable is also **related to the outcome**.

Example 6.10: Basketball Win Rates

Scenario: Two players, Brandon and Luca, play basketball in different settings.

Brandon's Record:

- Plays mostly against **children**
- Rarely plays against **peers**

Luca's Record:

- Plays mostly against **peers**
- Rarely plays against **children**

 **Key Point:** Let's look at their overall win rates and see if we can determine who is the better player.

Overall Win Rates

Player	Wins	Total Games	Win Rate
Brandon	78	100	78%
Luca	70	100	70%

First impression: Brandon has a higher overall win rate. However, what happens when we break down their performance by opponent type?

Breaking Down by Opponent Type

Games Against Children:

	Wins	Total
Brandon	75	80
Luca	10	10

- Brandon: $75/80 = 93.8\%$
- Luca: $10/10 = 100.0\%$

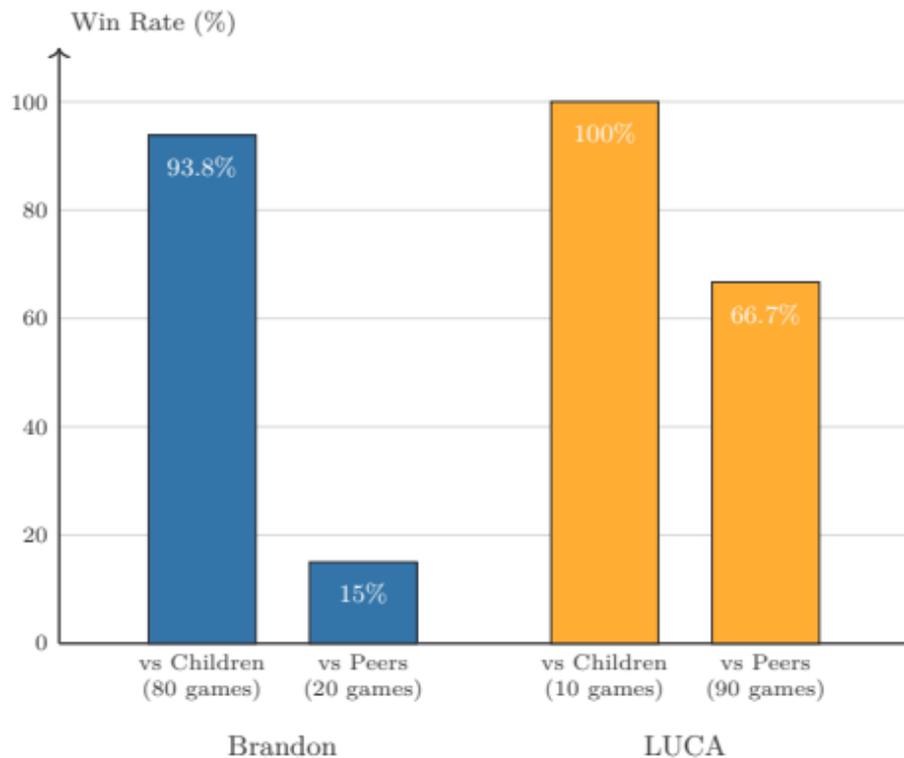
Games Against Peers:

	Wins	Total
Brandon	3	20
Luca	60	90

- Brandon: $3/20 = 15.0\%$
- Luca: $60/90 = 66.7\%$

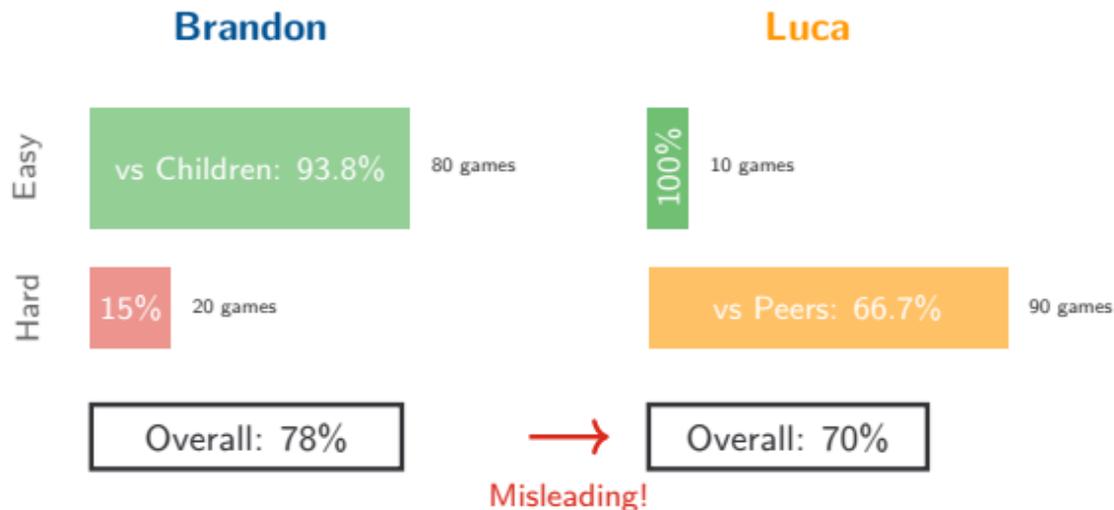
 **Key Point:** Luca has a higher win rate in **both** opponent categories.

Visualizing the Win Rates



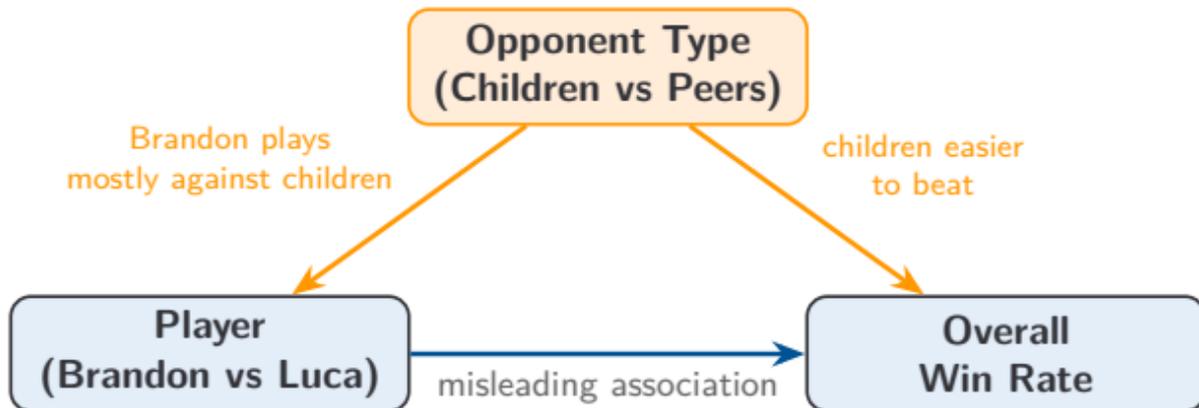
The Weighting Trap: Why Aggregates Mislead

Example: Two players with different opponent mixes.



Key Point: Bar width = sample size. Brandon's high average is inflated by playing mostly easy games.

The Hidden Variable: Opponent Type



Key Point: Opponent type is the confounding variable:

- Related to player (Brandon chooses children, Luca chooses peers)
- Related to outcome (children are easier opponents)
- Creates a misleading aggregate comparison

Simpson's Paradox in Basketball

Caution: The Paradox:

Brandon has a 78% overall win rate vs Luca's 70%

BUT Luca wins more often than Brandon against **both** children and peers.

Why this happens:

- Brandon plays **80%** of games against children (easy wins)
- Luca plays **90%** of games against peers (harder competition)
- The **mix of opponents** differs dramatically between players
- This creates a misleading aggregate picture

Example 6.9: By Stone Size

When we segment the data by stone size, we get:

Small Stones:

	A	B
Success	81	234
Total	87	270
Rate	93%	87%

Treatment
A is
superior

Large Stones:

	A	B
Success	192	55
Total	263	80
Rate	73%	69%

Treatment A
is superior.

Which is the preferred treatment for each stone size?

Treatment A was applied to 87 small stones	• Large stones are harder to treat
B was applied to 270 small stones	
Treatment A was applied to 263 large stones	• Treatment A has considerably more trials with large stones.
B was applied to 80 large stones	

Confounding Variables

Lurking Variable

A **lurking variable** is an unobserved variable that influences both the explanatory and response variables, potentially leading to a spurious association between them.

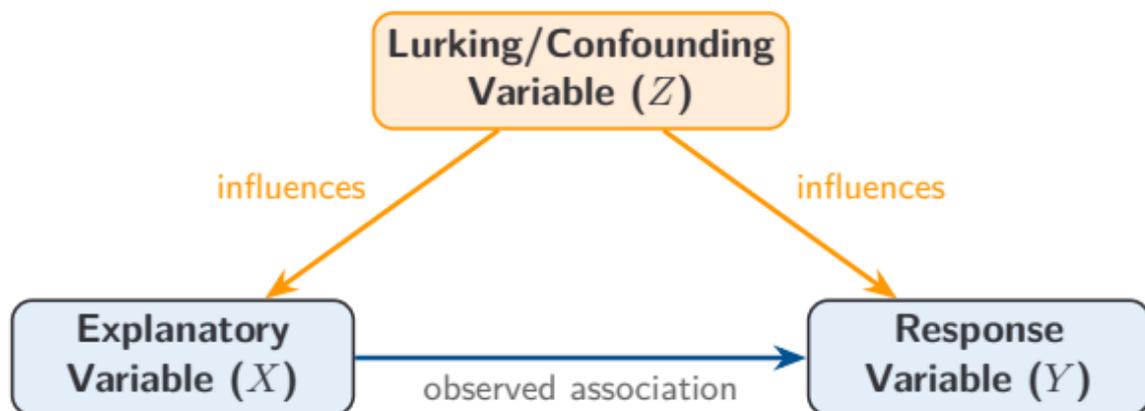
Confounding Variable

A **confounding variable** is a variable that influences both the explanatory and response variables, and is **measured** in the study.

↑ kidney stone size was first lurking, and then became confounding.

⚠ Caution: Confounding can cause associations to appear, disappear, or even reverse.

Confounding: A Causal Diagram



Key Point: A confounding variable Z creates a **spurious association** between X and Y because Z affects both.

Berkeley example: X = Gender, Y = Admission, Z = Department choice.

CHAPTER 6

Summary

Chapter 6 Summary

Two-Way Tables

- **Two-way tables:** counts for two categorical variables
- **Marginal:** totals (ignoring the other variable)
- **Conditional:** within each level of another variable
- **Independence:** conditional = marginal for all groups

Risk Measures

- **Relative risk:** ratio of probabilities
- **Odds ratio:** ratio of odds (cross-product: $\frac{ad}{bc}$)

Cautions

- **Simpson's paradox:** trends reverse when accounting for a third variable
- **Confounding:** affects both explanatory and response
- **Lurking:** unmeasured confounder

From Tables to Probability

 **Key Point:** Throughout this chapter, every proportion you computed from a contingency table was a **probability**.

- A **marginal proportion** $P(Y)$ is the probability of an outcome in the whole population
- A **conditional proportion** $P(Y | X)$ is the probability of an outcome *within a specific group*
- **Relative risk** compares two such probabilities: $\frac{P(Y | \text{Exposed})}{P(Y | \text{Unexposed})}$

Looking Ahead

In Chapter 12, we develop a formal framework for probability, with precise rules for combining and reasoning about uncertain events. The intuition you built here with contingency tables carries directly into that framework.

PRACTICE

Exercises

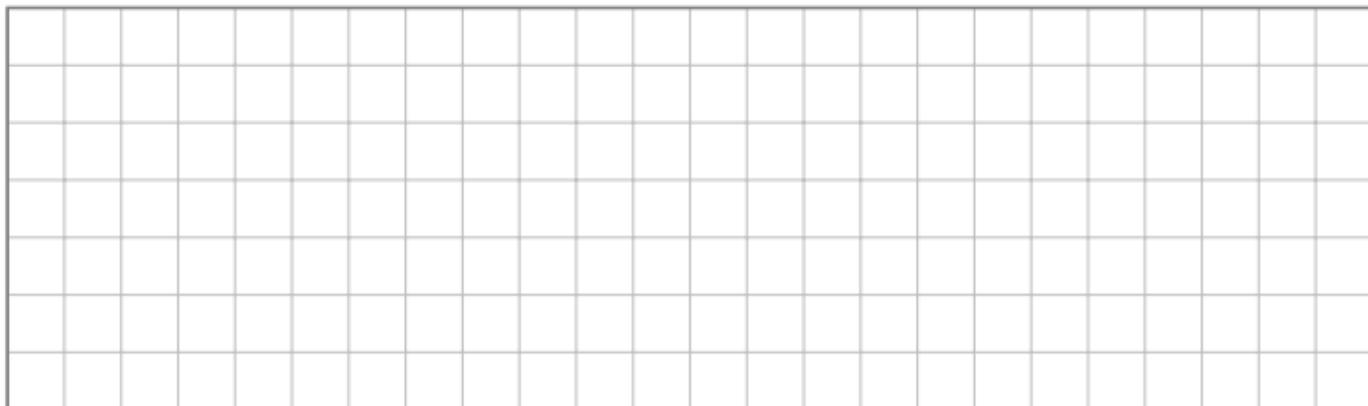
Example 6.12: Drug Side Effects

Context: A clinical trial with 500 participants testing a new drug.

	Side Effect	No SE	Total
Drug	36	214	250
Placebo	12	238	250
Total	48	452	500

Find:

1. Risk of side effect in each group
2. Relative risk (Drug vs. Placebo)
3. Odds ratio using cross-product
4. Interpret RR and OR in context



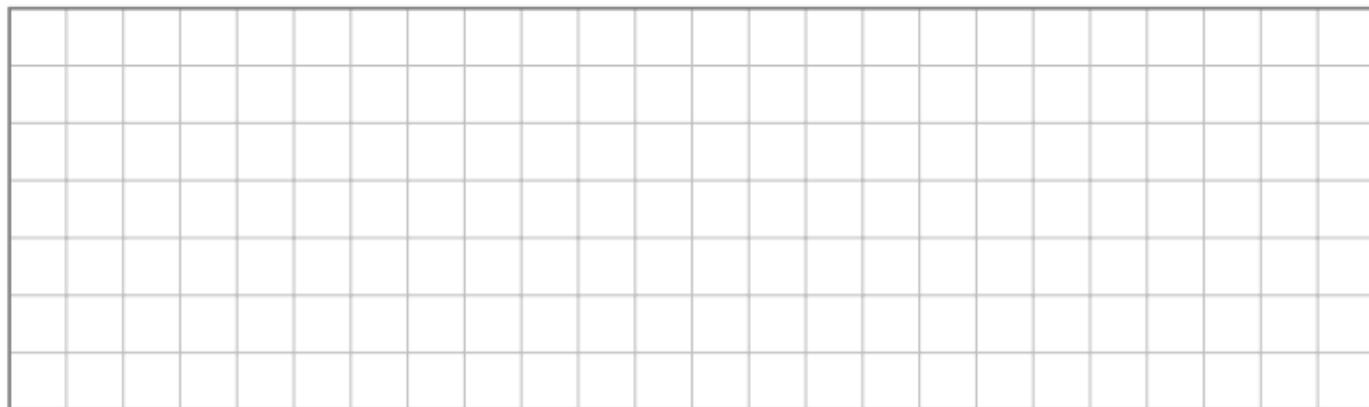
Example 6.13: Study Methods and Exam Outcomes

Context: A university tracked 300 students by study method and exam result.

	Pass	Fail	Total
Online study	90	60	150
In-person study	120	30	150
Total	210	90	300

Find:

1. Marginal distribution of exam result
2. Conditional distribution for each study method
3. RR and OR of passing (in-person vs. online)
4. Is in-person associated with better outcomes?



Example 6.14: Working Backwards

Context: In a study of 600 patients (300 per group), the relative risk of infection for the treatment group compared to the control group is $RR = 0.40$. The risk of infection in the **control group** is 0.25.

Find:

1. Risk of infection in treatment group
2. Number of infections in each group
3. Reconstruct the full two-way table
4. Odds ratio of infection (treatment vs. control)

3)

Group	Infection status		Sum
	Infected	Not infected	
Treatment	30	270	300
Control	75	225	300

1)	$RR = 0.40 = \frac{P(\text{infection} \text{treatment})}{P(\text{infection} \text{control})} = \frac{P(\text{infection} \text{treatment})}{0.25}$
	$\Rightarrow P(\text{infection} \text{treatment}) = 0.40 \cdot 0.25 = 0.10$
2)	$\# \text{ infections in treatment} = P(\text{infection} \text{treatment}) \# \text{ observations in treatment}$
	$= 0.10 \cdot 300 = 30$
	$\# \text{ infections} = P(\text{infection} \text{control}) \# \text{ observations in control}$
	$= 0.25 \cdot 300 = 75$
4)	$OR = (30/270) / (75/225) = 0.33$