Week 3, Class 5

Transforming and Creating Variables

Sean Westwood

Today

- Understanding conditional logic step-by-step
- Creating new variables with mutate()
- Variable recoding with case_when() and if_else()
- Calculating proportions with summarise()
- Extracting values with pull()
- Combining multiple transformation steps

| Why Must We Transform Data? | |
|-----------------------------|--|
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |

Real Political Science Data is Messy

Example: Survey responses about income:

- Raw: \$35,000, \$45000, 45k, forty-five thousand
- Analysis needs: 35000, 45000, 45000, 45000

Example: Education levels to numeric:

- Raw: high school, HS grad, 12 years, completed secondary
- Analysis needs: All coded as 12 (years of education)

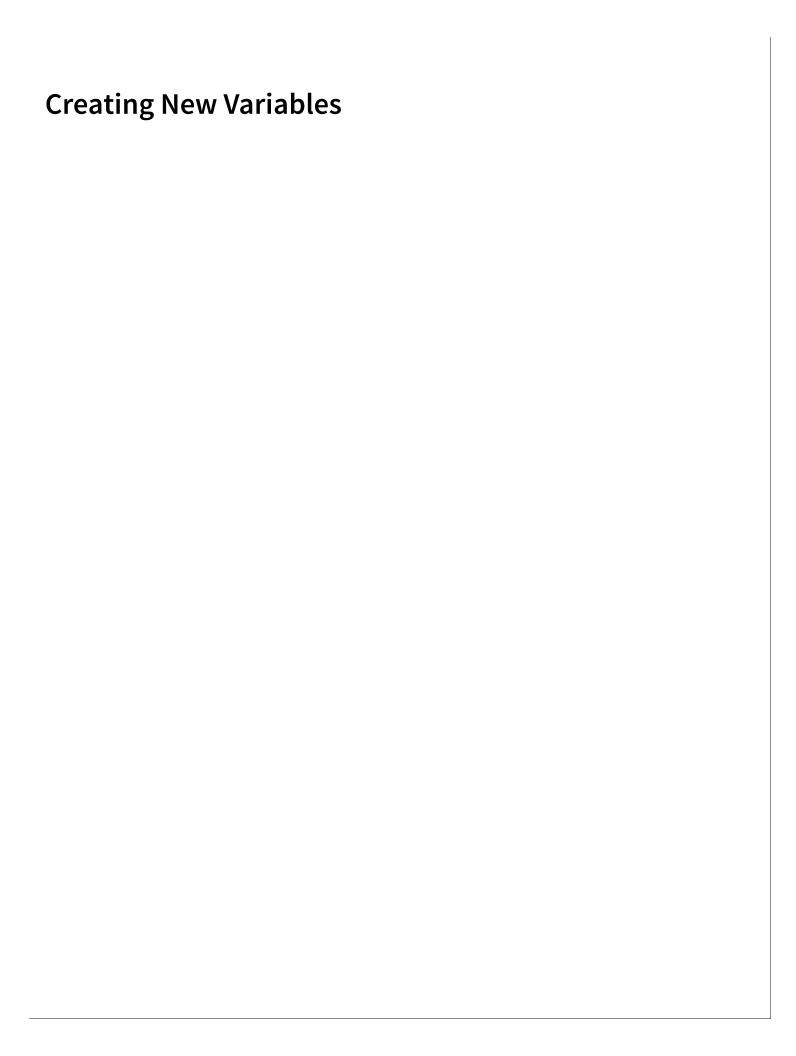
Example: Age data:

- Raw: Birth years like 1985, 1992, 1967
- Analysis needs: Age groups like "18-29", "30-44", "45-64"

The transformation process

- 1. Clean inconsistent formats
- 2. **Recode** into meaningful categories
- 3. Calculate new measures (ratios, percentages)

The majority of data analysis time is spent on data transformation!



Creating Variables with mutate()

The mutate() function adds new columns to your dataframe Basic syntax:

```
1 data %>%
2 mutate(
3    new_variable = calculation,
4    another_var = some_function(existing_var)
5  )
```

summarise() creates a new dataframe with new columns; mutate() alters the original dataframe

What mutate() does:

- 1. Takes your existing data
- 2. Adds new column(s) based on calculations
- 3. Keeps all original columns
- 4. Returns the expanded dataset

Mathematical Operations in mutate()

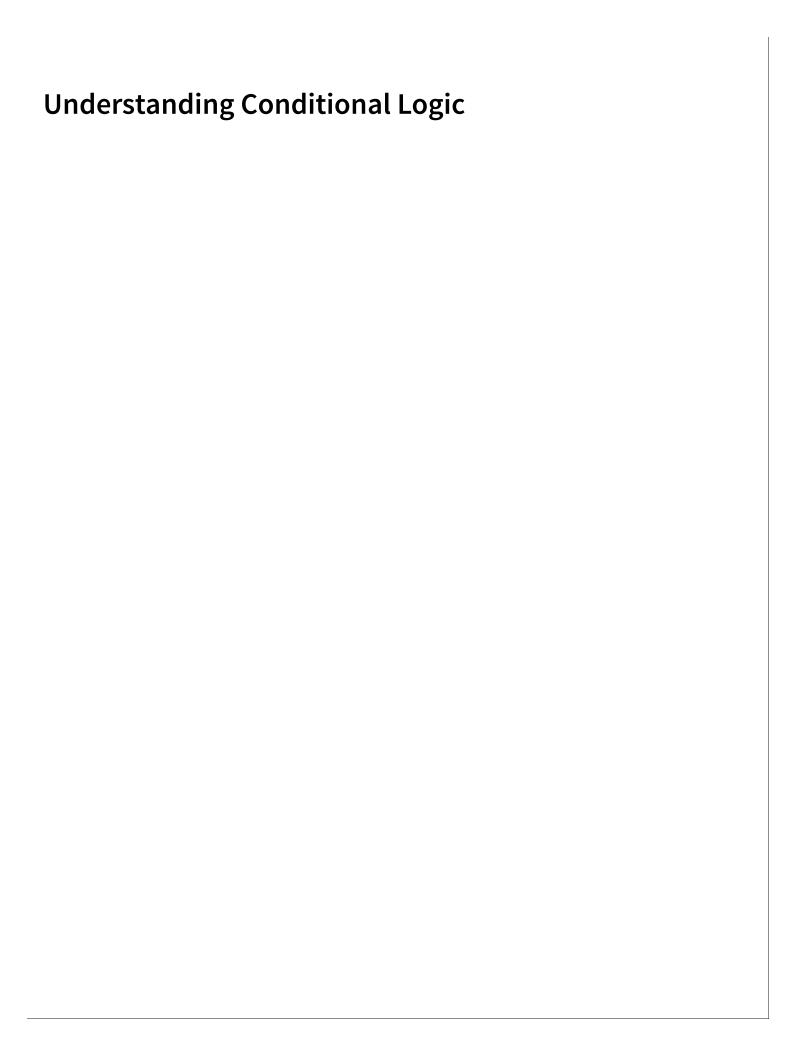
Basic math:

```
1 mutate(
2  total = votes_dem + votes_rep,  # Addition
3  margin = votes_dem - votes_rep,  # Subtraction
4  vote_share = votes_dem / total_votes,  # Division
5  doubled = approval_rating * 2  # Multiplication
6 )
```

This is row-wise! Each calculation will be independent for each row.

Useful calculations for political science:

```
1 mutate(
2  turnout_rate = votes_cast / eligible_voters,
3  spending_per_vote = total_spent / votes_received,
4  ideology_squared = ideology_score^2,
5  log_income = log(household_income)
6 )
```



What is Conditional Logic?

In order to create new variables that are useful for analysis, we often need to use conditional logic.

Conditional logic: Making decisions based on conditions **In everyday life**:

- If it's raining, then bring an umbrella
- If you're 21 or older, then you can drink alcohol
- If it is after 6:00 AM and before 10:00AM, then you can order breakfast at McDonalds

In R:

- If party == "Democrat", then assign "Liberal"
- If age >= 65, then assign "Senior"
- If approval_rating > 50, then assign "Majority Approves"

The Building Blocks: Comparison Operators

Comparison operators:

- == (equal to)
- != (not equal to)
- > (greater than)
- >= (greater than or equal to)
- < (less than)
- <= (less than or equal to)

The triad of 'equals'

- <- used to assign values to objects (make a variable 'get' a value)
- = Specify values when using a function (set an argumen to a value)
- == Test whether two values are equal (returning a logical result: TRUE, FALSE, or NA)

Why We Want TRUE/FALSE Outcomes

TRUE/FALSE is the foundation of data transformation because:

- All data decisions can ultimately be reduced down to "yes" or "no" questions
- We're essentially asking: "Does this observation meet our criteria?"
- These logical outcomes let us sort, categorize, and transform our data
- This binary logic is what allows us to create meaningful new variables from raw data with simple rules

Examples of TRUE/FALSE Outcomes

Equal to: ==

```
1 party_id == "Republican"
2 # TRUE if Republican, FALSE otherwise
```

Greater than or equal: >=

```
1 income >= 50000
2 # TRUE if 50,000 or more
```

Not equal to: !=

```
1 education != "College"
2 # TRUE if not College, FALSE if College
```

Less than: <

```
1 approval < 40
2 # TRUE if less than 40
```

Greater than: >

```
1 age > 65
2 # TRUE if older than 65
```

Less than or equal: <=

```
1 years_in_office <= 2
2 # TRUE if 2 years or fewer</pre>
```

Understanding TRUE and FALSE

Every comparison returns TRUE or FALSE:

```
# Examples with actual values
age <- 25
age >= 18  # Returns: TRUE
age >= 65  # Returns: FALSE

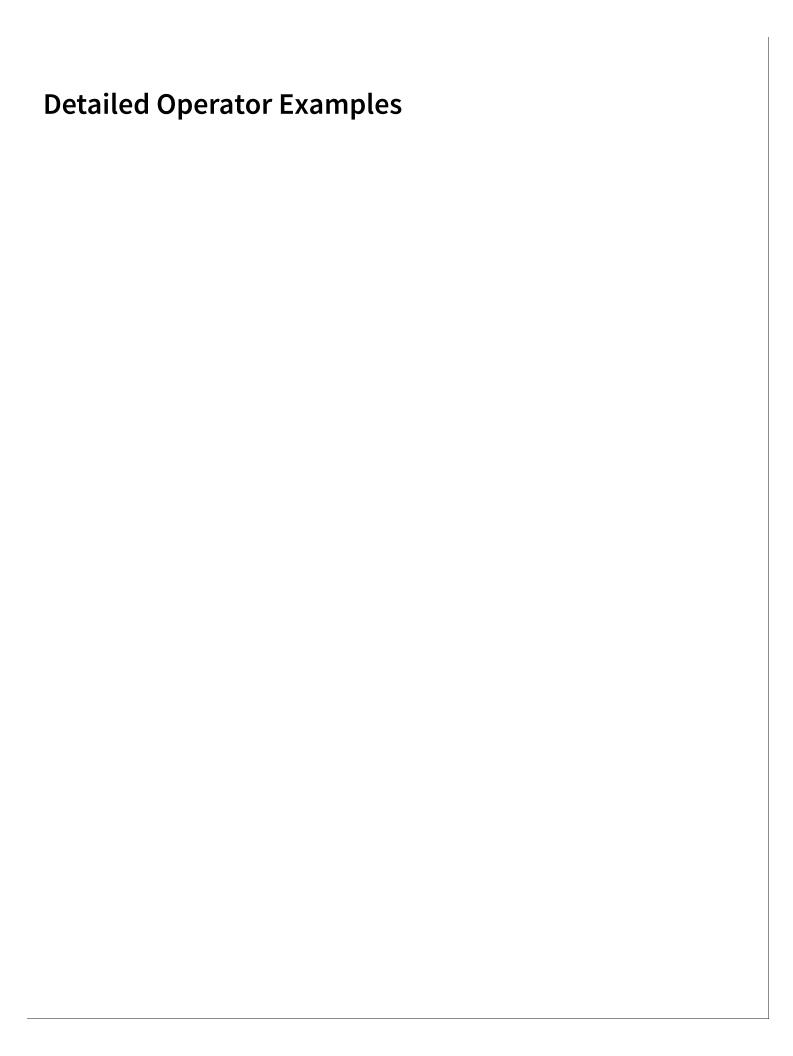
party <- "Democrat"
party == "Republican"  # Returns: FALSE
party != "Republican"  # Returns: TRUE</pre>
```

R treats TRUE as 1 and FALSE as 0:

We can apply mathimatical operators with TRUE/FALSE:

- sum(TRUE, FALSE, TRUE) = 2 (counts the TRUEs)
- mean(TRUE, FALSE, TRUE) = 0.67 (proportion that are TRUE)

Why mean() gives proportions: Since TRUE = 1 and FALSE = 0, the mean is the sum divided by the total count, which is exactly what a proportion is!



Example 1: Creating a New Variable with == (Equal To)

Question: How many survey respondents identify as Democrats?

Using == to Find Democrats

```
# Test each row: is party_id equal to "Democrat"?
survey_sample %>%
mutate(
    is_democrat = party_id == "Democrat"
)
```

What happened:

- Row 1: "Democrat" == "Democrat" → TRUE
- Row 2: "Republican" == "Democrat" → FALSE
- Row 3: "Democrat" == "Democrat" → TRUE
- And so on...

Counting with ==

```
survey_sample %>%
summarise(
total_respondents = n(),
democrats = sum(party_id == "Democrat"),
democrat_percentage = mean(party_id == "Democrat") * 100
)
```

```
# A tibble: 1 × 3
total_respondents democrats democrat_percentage
<int> <int> < dbl>
1 8 4 50
```

- sum(party_id == "Democrat") counts TRUE values
- mean(party_id == "Democrat") * 100 gives the proportion of Democrats

Example 2: The != (Not Equal To) Operator

Question: How many respondents are NOT Republicans?

```
# Count non-Republicans
survey_sample %>%
summarise(
non_republicans = sum(party_id != "Republican"),
non_republican_pct = mean(party_id != "Republican") * 100
)
```

Using >= to Find Seniors

```
# Create age sample data
age_sample <- tibble(
    voter_id = 1:10,
    age = c(22, 35, 67, 45, 72, 29, 81, 55, 19, 65)

# Show the data
age_sample</pre>
# Show the data
```

```
# Count seniors
age_sample %>%
summarise(
   total_voters = n(),
   seniors = sum(age >= 65),
   senior_percentage = mean(age >= 65) * 100

7 )
```

```
# A tibble: 1 × 3
total_voters seniors senior_percentage
<int> <int> <dbl>
1 10 4 40
```

Example 4: The <= (Less Than or Equal) Operator

Question: How many voters are young adults (30 or younger)?

```
# Count young adults
age_sample %>%
summarise(
young_adults = sum(age <= 30),
young_adult_pct = mean(age <= 30) * 100

)</pre>
```

```
# A tibble: 1 × 2
young_adults young_adult_pct
<int> <dbl>
1 3 30
```

Multiple Conditions

The world is complex - we often need to examine multiple variables simultaneously

Real-world questions require multiple conditions:

- Young AND Democrat voters (age < 30 AND party == "Democrat")
- High achievers (high IQ AND admitted to Dartmouth)
- At-risk populations (low income AND poor health)
- Swing voters (independent AND frequent voters)

Combining Conditions with AND (&) and OR (|)

AND (₺): Both conditions must be TRUE

1 age >= 18 & age < 65 # Working age adults

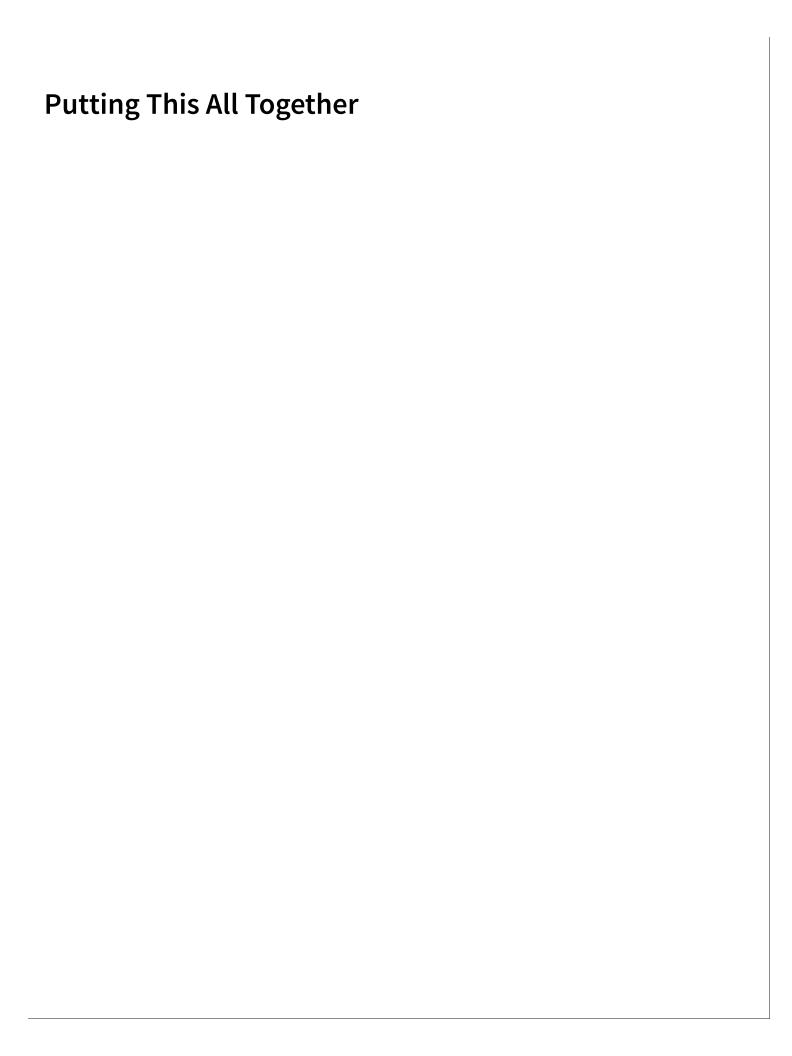
OR (|): Either condition can be TRUE

party == "Democrat" | party == "Republican" # Major party members

Example of multiple conditions with real data

```
1 # Complex conditions
2 age_sample %>%
3 mutate(
4 working_age = age >= 18 & age < 65,
5 voting_age = age >= 18
6 ) %>%
7 summarise(
8 working_age_count = sum(working_age),
9 voting_age_count = sum(voting_age)
10 )
```

```
# A tibble: 1 × 2
working_age_count voting_age_count
<int> <int> 
1 6 10
```



Logical Operations in mutate()

Creating TRUE/FALSE variables:

Why create logical variables?:

- Easy to count: sum(is_competitive)
- Easy to get percentages: mean(high_turnout)
- Clear for analysis: filter(major_party)

The if_else() Function

What if we want to create a new variable that is a categorical variable (e.g., "Senior", "Non-Senior")?

Instead of creating a variable that shows is someone is young or not, we can create a varaible that calls them "Young" or "Old".

For simple binary decisions:

```
1 mutate(
2 age_group = if_else(age >= 65, "Senior", "Non-Senior"),
3 result = if_else(vote_share > 0.5, "Won", "Lost"),
4 income_level = if_else(income >= 50000, "High", "Low")
5 )
```

Syntax: if_else(condition, value_if_true, value_if_false)
When to use if_else():

- Only two possible outcomes
- Simple, clear conditions where the labels are sensible

The case_when() Function

Real political data often has **multiple meaningful categories** that can't be captured with simple TRUE/FALSE:

Example 1: Education levels:

- Raw data: Years of education (8, 12, 14, 16, 18, 20)
- Need categories: "Less than HS", "High School", "Some College", "Bachelor's", "Graduate"

Example 2: Congressional ideology scores:

- Raw data: DW-NOMINATE scores (-0.8, -0.2, 0.1, 0.6, 0.9)
- Need categories: "Very Liberal", "Liberal", "Moderate", "Conservative", "Very Conservative"

case_when() handles multiple conditions in order:

Using case_when()

For multiple conditions and outcomes:

What if we wanted to make a new variable called "age_category" by checking each person's age and assigning them to one of four age groups based on which condition they meet first?

```
mutate(
    age_category = case_when(
    age < 30 ~ "Young",
    age < 50 ~ "Middle-aged",
    age < 70 ~ "Older",
    TRUE ~ "Senior" # catch-all for age >= 70
    )
}
```

Key rules for case_when():

- 1. Conditions are tested in order
- 2. First TRUE condition wins
- 3. Use ~ to separate condition from result
- 4. Always end with TRUE ~ "catch-all" for safety

Understanding the ~ Symbol

The tilde (∼) means "then":

- age < 30 ~ "Young" means "if age < 30, then assign 'Young'"
- income >= 100000 ~ "High" means "if income >= 100000, then assign 'High"

Think of it as an arrow:

- condition ∼ result
- "When this condition is TRUE, give this result"

Only used in case_when(), not in if_else()

case_when() vs if_else()

Use if_else() for 2 categories:

```
party_type = if_else(party == "Independent", "Independent", "Major Party")
```

Use case_when() for 3+ categories:

```
party_group = case_when(
party == "Democrat" ~ "Democrat",
party == "Republican" ~ "Republican",
party == "Independent" ~ "Independent",
TRUE ~ "Other"

6 )
```

Both create new columns, but case_when() is more flexible

Common case_when() Patterns

Age groups:

```
1 age_group = case_when(
2 age < 25 ~ "18-24",
3 age < 35 ~ "25-34",
4 age < 50 ~ "35-49",
5 age < 65 ~ "50-64",

TRUE ~ "65+"

7 )
```

Income brackets:

```
income_bracket = case_when(
income < 25000 ~ "Low",
income < 75000 ~ "Middle",

TRUE ~ "High"

)</pre>
```

Performance ratings:

```
1 rating = case_when(
2    score >= 90 ~ "Excellent",
3    score >= 80 ~ "Good",
4    score >= 70 ~ "Fair",
5    TRUE ~ "Poor"
6 )
```

| Examples: Summarize Without Mutate | | |
|------------------------------------|--|--|
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |

Example 1: Congressional Voting Patterns Analysis

Research Question: What are the basic characteristics of congressional voting behavior by party?

Data: DW-NOMINATE dataset with 1,000 members of Congress, including ideology scores, party affiliation, chamber, and state information.

Code: Results:

```
congress_data %>%
group_by(party_code) %>%
summarise(
count = n(),
mean_ideology = mean(nominate_dim1, na.rm = TRUE),
median_ideology = median(nominate_dim1, na.rm = TRUE),
percentage = n() / nrow(congress_data) * 100
)
```

AI Prompts for This Congressional Analysis

Prompt to Create Code:

```
"I have congressional data that look like this:
Rows: 1000Rows: 1,000
Columns: 5
$ member_id
14, 15, 16, 1...
               <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13,
$ party_code <chi
"Democrat", "Repub...</pre>
              <chr> "Democrat", "Republican", "Democrat",
Create tidyverse code to summarize by party: count of members, mean
```

Al Response:

I'll help you create a summary by party using group_by() and summarise(). Here's the code:

[code provided above]

Prompt to Explain Code:

"Explain this summarise code step by step. What does each function do and why do we use na.rm = TRUE for the mean and median calculations?"

Al Response:

Let me break this down:

- group_by(party_code): Creates separate groups for each party
- n(): Counts rows in each group
 mean(nominate_dim1, na.rm=TRUE): Calculates average ideology, ignoring missing values
- 4. na.rm=TRUE prevents errors if any ideology scores are missing (NA)
- 5. percentage calculation: n() divided by total dataset size * 100

Example 2: Election Turnout Demographics Analysis

Research Question: How does voter turnout vary across different demographic groups?

Data: Simulated voter turnout dataset with 1,000 respondents, including age groups, education levels, income brackets, and turnout status.

Code: Results:

```
voter_data %>%
group_by(education_level) %>%
summarise(
total_respondents = n(),
voters = sum(voted == "Yes"),
non_voters = sum(voted == "No"),
turnout_rate = mean(voted == "Yes"),
median_age = median(age)
)
```

```
# A tibble: 4 × 6
 education_level total_respondents voters non_voters
turnout_rate median_age
 <chr>
                             <int> <int>
<dbl>
1 College
0.718
2 Graduate
                               170
           52
0.806
3 High School
                                     229
                                                187
                               416
           53.5
0.550
4 Some College
                                      83
                                                 51
                               134
0.619
```

AI Prompts for Turnout Analysis

Prompt to Create Code:

```
"I have voter data tha tlook liek this

Rows: 1,000
Columns: 4
$ respondent_id <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16,...
$ education_level <chr> "High School", "Graduate", "College", "High School", "...

Using tidyverse Create a summary by education level showing: total respondents, number who voted, number who didn't vote, turnout rate, and median age."
```

Al Response:

```
I'll create a grouped summary that analyzes turnout by education
level:
[code provided above]

This uses:
    group_by() to analyze each education level
    sum(voted == "Yes") to count "Yes" responses
    mean(voted == "Yes") for turnout percentage
    median() for typical age in each group
```

Prompt to Explain Code:

"Explain how sum(voted == 'Yes') works to count voters and how mean(voted == 'Yes') calculates the turnout rate. Why does this work?"

Al Response:

```
Great question! This works because R treats TRUE/FALSE as 1/0:

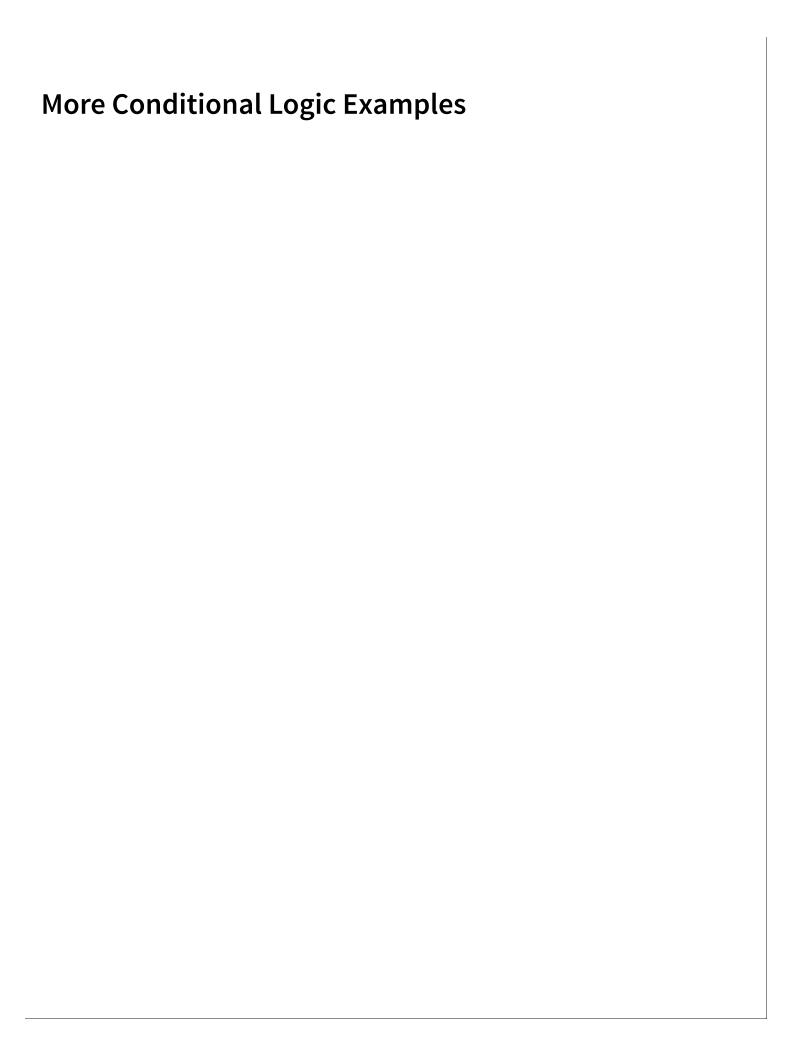
sum(voted == "Yes"):

- Creates TRUE/FALSE for each row
- sum() treats TRUE as 1, FALSE as 0

- Result: count of "Yes" responses

mean(voted == "Yes"):

- Same TRUE/FALSE creation
- mean() of 1s and 0s = proportion
- Result: percentage who voted
```



Example: Income Analysis with Multiple Operators

Research Question: How do income patterns vary in survey data?

Using Different Operators on Income Data

```
# A tibble: 12 \times 5
   person_id annual_income low_income high_income middle_income
                          <dbl> <lgl>
                                                <lgl>FALSE
        <int>
                                                                <lgl>
                          25000 TRUE
75000 FALSE
                                                                FALSE
TRUE
                           45000 FALSE
                                                 FALSE
                         120000 FALSE
                                                 TRUE
                           35000 FALSE
                                                 FALSE
                           95000 FALSE
                                                 FALSE
6
7
8
9
10
11
12
                           15000 TRUE
                                                 FALSE
                                                                FALSE
            8
9
10
                           65000 FALSE
                                                 FALSE
                                                                 TRUE
                         85000 FALSE
200000 FALSE
42000 FALSE
58000 FALSE
                                                FALSE
TRUE
FALSE
                                                                TRUE
                                                                FALSE
                                                                TRUE
                                                 FALSE
                                                                TRUE
```

Income Analysis Summary

```
income_demo %>%
summarise(
total_people = n(),
low_income_count = sum(annual_income <= 30000),
high_income_count = sum(annual_income >= 100000),
middle_income_count = sum(annual_income > 30000 & annual_income < 100000),
median_income = median(annual_income),
avg_income = mean(annual_income)
)</pre>
```

AI Prompts for Income Analysis

Prompt to Create Code:

"I have income data with person_id, annual_income, and state variables. I want to create logical variables that categorize people as: low income (≤30K), high income (≥100K), middle income (between 30K-100K), and identify anyone earning exactly the median. Then summarize the counts. The data looks like this:

person_id annual_income state

I want tidyverse code to do this."

Simulated AI Response:

I'll help you create logical variables and summarize the income data:

[code provided above]

Key techniques used:

- <= and >= operators for thresholds
- Combined conditions with & operator
- != operator for exclusions
- == operator for exact matches
- sum() counts TRUE values

Prompt to Explain Code:

"Explain how the logical operators work in this income analysis. Why do we use sum() to count TRUE values? What's the difference between & and | in logical conditions?"

Simulated AI Response:

Logical operators create TRUE/FALSE values:

- <= 30000 checks if income is at or below threshold
 >= 100000 checks if income meets high threshold
 & requires BOTH conditions to be true

- | requires at least ONE condition to be true

sum() with logical variables:

- TRUE = 1, FALSE = 0 in R
- sum(condition) counts how many TRUEs- Efficient way to count observations meeting criteria

Example 3: Campaign Spending Analysis

Research Question: What are the spending patterns across different types of political campaigns?

Data: Campaign finance dataset with 1,000 candidates, including office type, party, total raised, total spent, and campaign outcome.

Code: Results:

```
campaign_data %>%
group_by(office_type, party) %>%
summarise(
num_candidates = n(),
avg_raised = mean(total_raised),
median_spent = median(total_spent),
win_rate = mean(outcome == "Won"),
total_spending = sum(total_spent),
groups = "drop"
)
```

| # A tibble: 6 | × 7 | | | |
|----------------------------|-------------|----------------|-------------|--------------|
| office_type | party | num_candidates | avg_raised | median_spent |
| win_rate <chr></chr> | <chr></chr> | <int></int> | <dbl></dbl> | <dbl></dbl> |
| <dbl></dbl> | | | | |
| 1 Governor 0.472 | Democrat | 142 | 3126128. | 2852395. |
| 2 Governor 0.551 | Republican | 138 | 3229189. | 2855065. |
| 3 House | Democrat | 328 | 238780. | 223938. |
| 0.506 4 House | Republican | 299 | 243071. | 229582. |
| 0.482 5 Senate | Democrat | 40 | 1708165. | 1576734. |
| 0.425 6 Senate 0.547 | Republican | 53 | 1776919. | 1631585. |
| 0.347 | | | | |

AI Prompts for Campaign Analysis

Prompt to Create Code:

```
"I have campaign data with office_type, party, total_raised,
total_spent, and outcome variables. Group by office type and party to show: candidate count, average raised, median spent, win rate, and total spending. The data look like this: Rows: 12
Columns: 3
$ person_id <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12
$ annual_income <dbl> 25000, 75000, 45000, 120000, 35000, 95000,
15000, 65000,...
$ state ..."
"FL", "CA", "T..."
                         <chr> "CA", "TX", "NY", "CA", "FL", "TX", "NY",
I want tidyverse code to do this.
```

Simulated AI Response:

I'll create a grouped summary by office type and party. Here's the code:

[code provided above]

This groups by TWO variables simultaneously: - Each combination gets its own row

- Calculations done within each group
- .groups="drop" removes grouping after

Prompt to Explain Code:

"Walk me through this grouped summary. How does grouping by two variables work? Why might median be better than mean for spending analysis?"

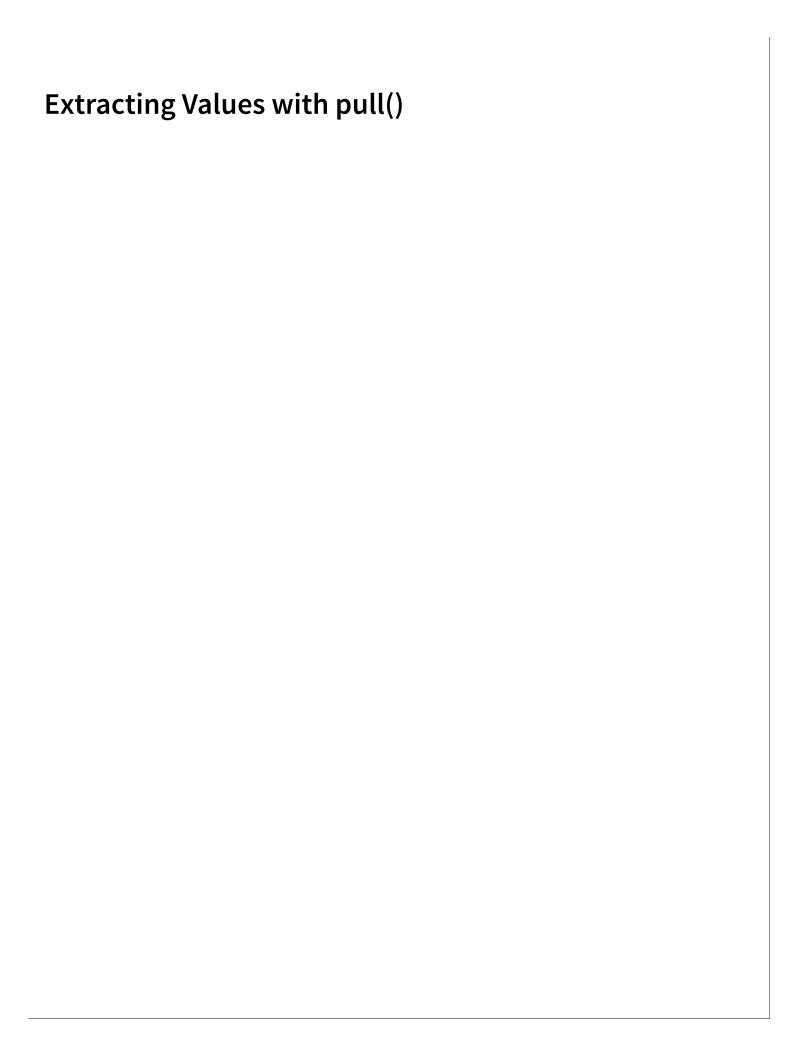
Simulated AI Response:

Multi-variable grouping creates subgroups:

- group_by(office_type, party) creates 6 groupsEach office type + party combination
- Statistics calculated within each subgroup

Median vs Mean for spending:

- Campaign spending often has outliers
- Median less affected by extremely high spenders
 Mean pulled up by billionaire self-funders
 Median shows "typical" candidate spending



When to Use pull()

Purpose: Extract a single column as a vector

Common use cases:

- Get a single statistic for further calculation
- Extract values for comparison
- Create variables from summary statistics

Basic syntax:

```
data %>%
pull(column_name)

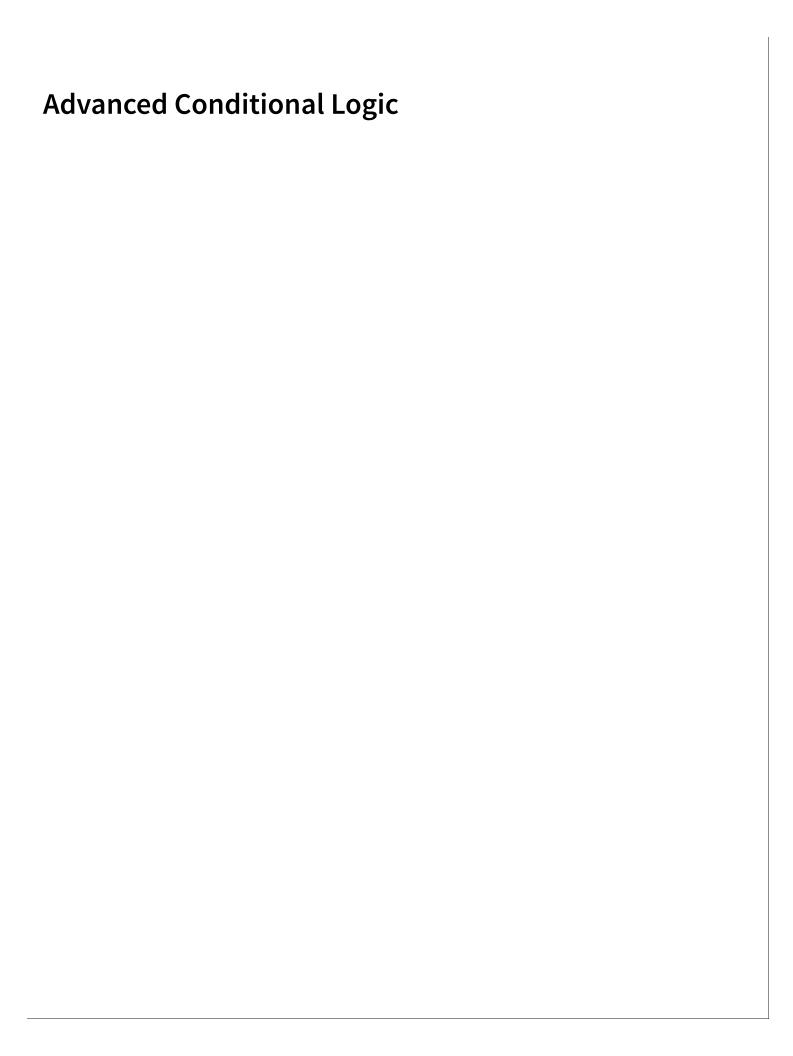
# Or extract calculated values
data %>%
summarise(mean_value = mean(variable)) %>%
pull(mean_value)
```

Example: Using pull() for Calculations

```
# Extract median income for threshold calculation
median_income <- survey_data %>%
summarise(median_inc = median(income)) %>%
pull(median_inc)
median_income
```

```
[1] 51198.99
```

```
# Use extracted value in new calculation
survey_data %>%
mutate(
   above_median = income > median_income,
   income_ratio = income / median_income
) %>%
select(respondent_id, income, above_median, income_ratio) %>%
slice_head(n = 5)
```



Example: Complex Campaign Classifications

Research Question: How can we classify campaigns by competitiveness and resource level?

```
# Create sample campaign data with more variables
campaign_complex <- tibble(
    candidate = paste("Candidate", 1:10),
    vote_share = c(0.52, 0.48, 0.67, 0.33, 0.51, 0.49, 0.78, 0.22, 0.55, 0.45),
    total_spent = c(150000, 145000, 300000, 100000, 180000, 175000, 500000, 80000, 200000, 190000),
    incumbent = c(TRUE, FALSE, TRUE, FALSE, TRUE, TRUE, FALSE, FALSE, TRUE)
}
campaign_complex</pre>
```

```
# A tibble: 10 × 4
  candidate vote_share total_spent incumbent
1 Candidate 1
                      0.52
                                  150000 TRUE
 2 Candidate 2
                       0.48
                                  145000 FALSE
3 Candidate 3
4 Candidate 4
                       0.67
0.33
                                  300000 TRUE
100000 FALSE
                       0.51
                                  180000 FALSE
5 Candidate 5
6 Candidate 6
                                  175000 TRUE
                       0.49
7 Candidate 7
                       0.78
                                  500000 TRUE
                                  80000 FALSE
8 Candidate 8
                       0.22
9 Candidate 9
                                 200000 FALSE
                       0.55
                                  190000 TRUE
10 Candidate 10
                       0.45
```

Creating Complex Classifications with case_when()

```
campaign_complex %>%
        # Competitiveness based on vote margin
        competitiveness = case_when(
abs(vote_share - 0.5) <= 0.02 ~ "Extremely Close", # Within 2%
abs(vote_share - 0.5) <= 0.05 ~ "Competitive", # Within 5%
abs(vote_share - 0.5) <= 0.10 ~ "Somewhat Safe", # Within 10%
                                                                                       # Within 10%
           TRUE ~ "Safe"
                                                                                       # More than 10%
        # Spending level
        # Spending level = case_when(
  total_spent >= 400000 ~ "High Budget",
  total_spent >= 200000 ~ "Medium Budget",
  total_spent >= 100000 ~ "Low Budget",
  TRUE ~ "Minimal Budget"
        # Advantage type using multiple conditions
        advantage = case_when(
  incumbent == TRUE & total_spent >= 200000 ~ "Incumbent + Money",
  incumbent == TRUE & total_spent < 200000 ~ "Incumbent Only",
  incumbent == FALSE & total_spent >= 200000 ~ "Money Only",
           TRUE ~ "Neither Advantage
     select(candidate, vote_share, competitiveness, spending_level, advantage)
# A tibble: 10 × 5
    {\tt candidate} \qquad {\tt vote\_share} \ {\tt competitiveness} \ {\tt spending\_level} \ {\tt advantage}
    <chr>
                             <dbl> <chr>
                                                             <chr>
 1 Candidate 1
                               0.52 Competitive
                                                              Low Budget
                                                                                    Incumbent Only
 2 Candidate 2
                               0.48 Competitive
                                                              Low Budget
                                                                                    Neither Advantage
 3 Candidate 3
                               0.67 Safe
                                                              Medium Budget
                                                                                    Incumbent + Money
                               0.33 Safe
4 Candidate 4
                                                              Low Budget
                                                                                    Neither Advantage
                               0.51 Extremely Close Low Budget
                                                                                    Neither Advantage
5 Candidate 5
6 Candidate 6
                               0.49 Extremely Close Low Budget
                                                                                    Incumbent Only
7 Candidate 7
                               0.78 Safe
                                                              High Budget
                                                                                    Incumbent + Money
                              0.22 Safe Minimal Budget No.
0.55 Somewhat Safe Medium Budget Money Only
Low Budget Incumbent Only
8 Candidate 8
                                                              Minimal Budget Neither Advantage
9 Candidate 9
10 Candidate 10
```

Understanding Complex Conditions

Breaking down the advantage classification:

- 1. incumbent == TRUE & total_spent >= 200000
 - Must be incumbent AND high spending
 - Both conditions must be TRUE
- 2. incumbent == TRUE & total_spent < 200000</pre>
 - Must be incumbent AND low spending
 - Checked only if #1 is FALSE
- 3.incumbent == FALSE & total_spent >= 200000
 - Must be challenger AND high spending
 - Checked only if #1 and #2 are FALSE
- 4. TRUE
 - Catches all remaining cases
 - Challenger with low spending

Example 4: Public Opinion Polling

Research Question: How do approval ratings vary across different demographic and geographic groups?

Code:

```
polling_data %>%
    group_by(region, party_id) %>%
summarise(
    sample_size = n(),
    avg_approval = mean(approval_rating),
    approval_range = max(approval_rating) - min(approval_rating),
    high_approval_pct = mean(approval_rating > 50),
    response_rate = mean(!is.na(approval_rating)),
    .groups = "drop"
)
```

```
# A tibble: 12 × 7
            party_id sample_size avg_approval approval_range
  region
high_approval_pct
<chr>
                             <int>
                                          <dbl>
                                                         <dbl>
            <chr>
<dbl>
 1 Midwest Democrat
                                           61.8
0.681
2 Midwest Independ...
                               88
                                           54.1
                                                            60
0.614
 3 Midwest Republic...
                                           45.1
                                                            60
0.453
4 Northeast Democrat
                                88
                                           57.4
0.625
5 Northeast Independ...
                                94
                                           52.7
                                                            60
0.5
6 Northeast Republic...
                                           44.8
                                                            60
0.423
```

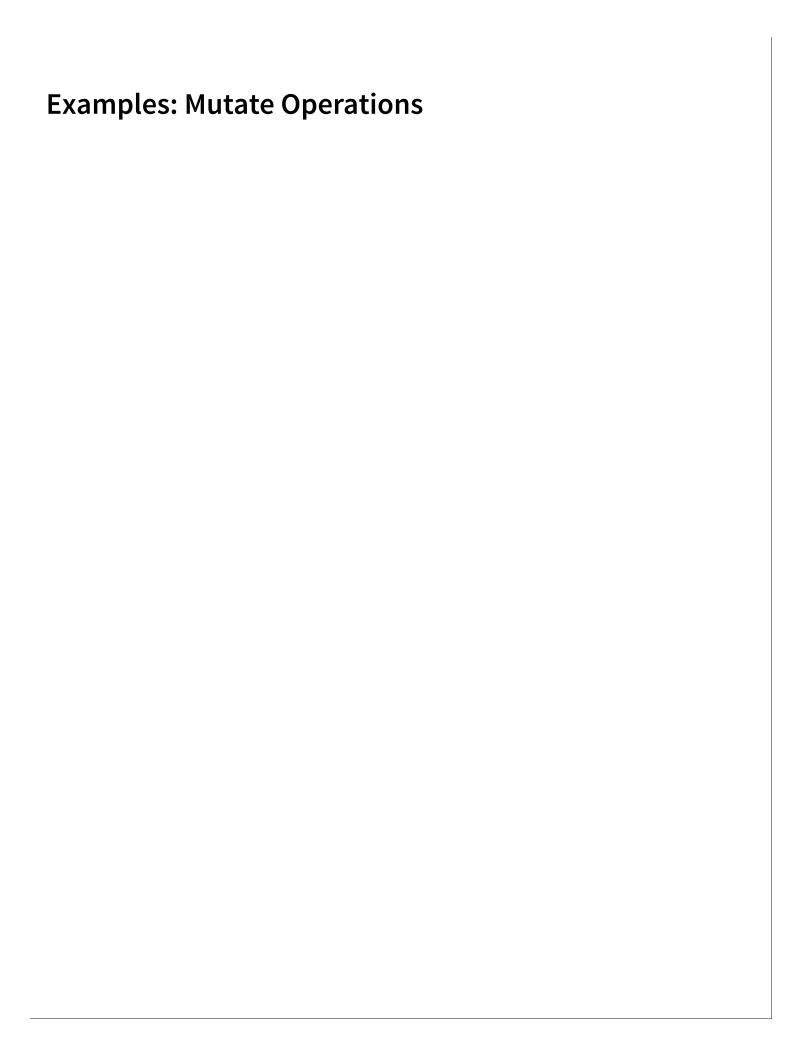
Al Prompt to Create:

"I have polling data with region, party_id, and approval_rating variables. Create a summary by region and party showing: sample size, average approval, approval range, percentage with high approval (>50), and response rate with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"<Paste the results from R.>

Explain how mean(approval_rating > 50) calculates the percentage
with high approval, and how mean(!is.na(approval_rating))
calculates response rate."



Example 5: Creating Campaign Efficiency Metrics

Research Question: How can we create meaningful efficiency metrics for campaign spending?

```
campaign_data %>%
mutate(
spending_efficiency = vote_share / (total_spent / 1000),
fundraising_ratio = total_raised / total_spent,
cost_per_vote = total_spent / total_votes,
budget_category = case_when(
total_spent < 100000 ~ "Low Budget",
total_spent < 500000 ~ "Medium Budget",
TRUE ~ "High Budget"
)
)
) >>%
select(candidate_name, spending_efficiency, fundraising_ratio,
cost_per_vote, budget_category) %>%
slice_head(n = 5)
```

```
# A tibble: 5 \times 5
 candidate_name spending_efficiency fundraising_ratio
cost_per_vote
                             0.00128
                                                  0.842
1 Anderson, John
2 Williams, Susan
                             0.00355
                                                   1.19
3 Davis, John
                             0.000134
                                                  0.867
50.5
4 Smith, David
                             0.00170
                                                  1.19
3.58
5 Miller, Bob
                             0.000334
                                                   1.42
17.7
# i 1 more variable: budget_category <chr>
```

Al Prompt to Create:

"I have campaign data with vote_share, total_spent, total_raised, and total_votes. Create new variables for: spending efficiency(vote share per \$1000 spent), fundraising ratio, cost per vote, and budget categories with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain each mutate calculation. Why divide total_spent by 1000 for efficiency? How does case_when work for creating budget categories?"

Example 6: Demographic Recoding and Analysis

Research Question: How can we create standardized demographic categories for analysis?

```
survey_data %>%
mutate(
    age_group = case_when(
    age < 30 ~ "18-29",
    age < 45 ~ "30-44",
    age < 65 ~ "45-64",
    TRUE ~ "65+"
    ),
    income_quintile = ntile(income, 5),
    education_years_std = (education_years - mean(education_years)
    high_earner = income > median(income),
    political_engagement = (political_interest + civic_participati
    ) %>%
    select(respondent_id, age_group, income_quintile, education_year
    high_earner, political_engagement) %>%
    slice_head(n = 5)
```

```
# A tibble: 5 × 6
 respondent_id age_group income_quintile education_years_std
high_earner
         <int> <chr>
                                                        <dbl>
<lgl>
              1 45-64
                                                        0.845
TRUE
              2 65+
                                                        0.583
FALSE
             3 30-44
                                                       -1.25
TRUE
             4 65+
                                                       -0.201
TRUE
              5 45-64
                                                       -0.724
TRUE
# i 1 more variable: political_engagement <dbl>
```

Al Prompt to Create:

"I have survey data with age, income, education_years, political_interest, and civic_participation. Create: age groups, income quintiles, standardized education, high earner indicator, and political engagement score with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain how ntile() creates quintiles, what standardization does to education_years, and how the political engagement score is calculated."

Example 7: Electoral Performance Metrics

Research Question: How can we measure electoral competitiveness and performance?

```
1 election_data %>%
2  mutate(
3  vote_margin = abs(dem_votes - rep_votes),
4  competitive_race = vote_margin < (total_votes * 0.05),
5  dem_vote_share = dem_votes / total_votes,
6  turnout_rate = total_votes / eligible_voters,
7  incumbent_advantage = if_else(incumbent_party == winner, 1, 0)
8  upset_victory = (dem_vote_share > 0.5 & dem_expected_share < 0
9  (dem_vote_share < 0.5 & dem_expected_share > 0.0
10 ) %>%
11  select(district, competitive_race, dem_vote_share, turnout_rate, incumbent_advantage, upset_victory) %>%
12  slice_head(n = 5)
```

```
# A tibble: 5 × 6
 district competitive_race dem_vote_share turnout_rate
incumbent_advantage
1 NV-02
                                     0.537
                                                  0.716
2 ID-03
           FALSE
                                     0.261
                                                  0.825
 VT-06
           TRUE
                                     0.518
                                                  0.579
4 IL-02
           FALSE
                                     0.301
                                                  0.733
5 MI-01
           FALSE
                                     0.707
                                                  0.612
# i 1 more variable: upset_victory <lgl>
```

Al Prompt to Create:

"I have election data with dem_votes, rep_votes, total_votes, eligible_voters, incumbent_party, winner, and expected vote shares. Create metrics for: vote margin, competitive races, vote shares, turnout, incumbent advantage, and upsets with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain the logic for identifying competitive races and upset victories. How do the logical operators & and | work in the upset_victory calculation?"

Example 8: Policy Position Scaling

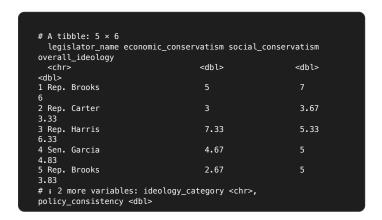
Research Question: How can we create meaningful scales from multiple policy positions?

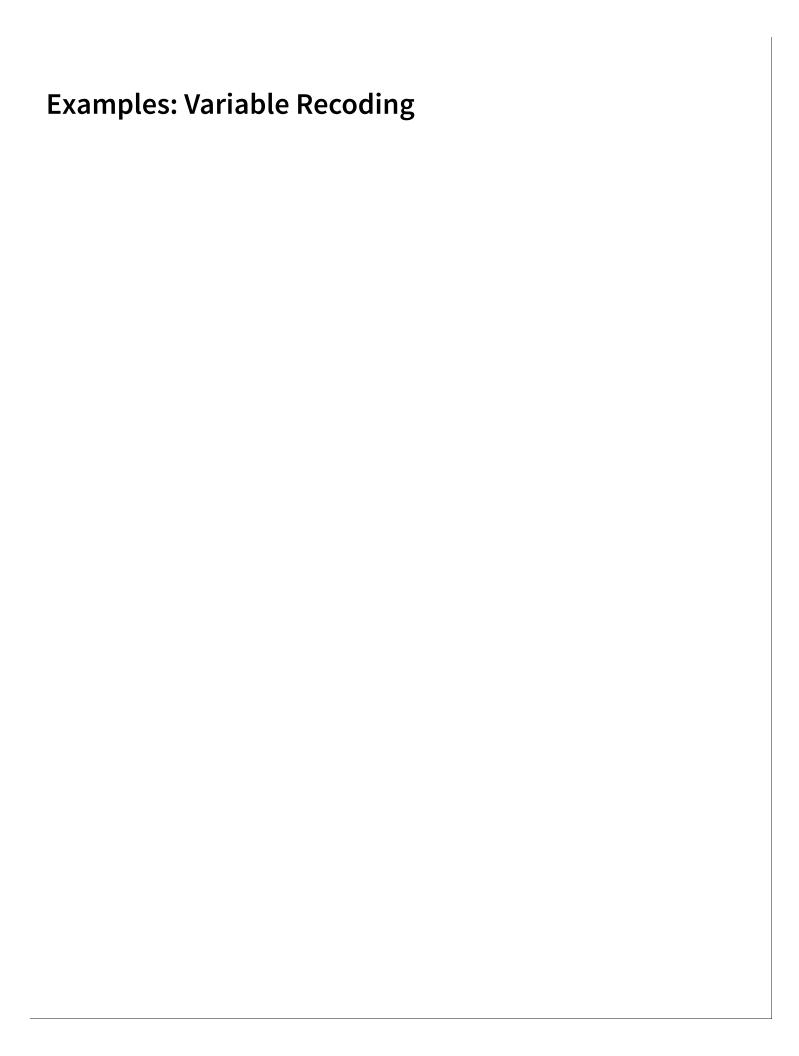
Al Prompt to Create:

"I have legislator data with multiple policy voting scores (1–10 scale). Create composite scores for economic and social conservatism, overall ideology, ideology categories, extreme position indicators, and policy consistency measures with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain how averaging creates composite scores and how policy_consistency measures the difference between economic and social positions. Why use abs()?"





Example 9: Income Bracket Recoding

Research Question: How can we create meaningful income categories from continuous income data?

Al Prompt to Create:

"I have household income and savings data. Create income brackets (6 categories), poverty status (3 categories), and economic security (combining income and savings thresholds). Show the distribution across all categories."

Al Prompt to Explain:

"Explain how case_when processes conditions in order and why TRUE is used as the final condition. How does the compound condition in economic_security work?"

```
# A tibble: 10 × 4
  {\tt income\_bracket} \quad {\tt poverty\_status} \ {\tt economic\_security}
   <chr>>
                     <chr>
                                     <chr>
                                                          <int>
                     Above Poverty Insecure
1 High Income
                                                            50
                                                            91
12
72
252
65
247
                     Above Poverty Secure
2 High Income
                     Below Poverty Insecure
3 Low Income
                     Near Poverty
4 Low Income
                                     Insecure
5 Lower Middle
                    Above Poverty Insecure
                    Near Poverty Insecure
Above Poverty Insecure
6 Lower Middle
 7 Middle Class
8 Upper Middle
                     Above Poverty Insecure
9 Upper Middle
                     Above Poverty Secure
10 Very High Income Above Poverty Insecure
```

Example 10: Geographic Region Consolidation

Research Question: How can we group states into meaningful political regions?

Al Prompt to Create:

"I have state-level data with state abbreviations and electoral_votes. Create political regions (5 main regions), identify battleground states, and categorize states by electoral vote count with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain how %in% works to check if a state belongs to a vector of states. Why is this more efficient than multiple | (OR) conditions?"

```
# A tibble: 10 × 4
  political_region battleground_state electoral_votes_category
  <chr>
1 Midwest
                   FALSE
                                       Medium State
80
2 Midwest
                   FALSE
                                       Small State
120
3 Midwest
                                       Medium State
                    TRUE
40
4 Mountain West
                    FALSE
                                       Medium State
5 Mountain West
                    FALSE
                                       Small State
6 Mountain West
                    TRUE
                                       Medium State
20
```

Example 11: Survey Response Cleaning

Research Question: How can we clean and standardize survey response data?

```
survey_clean_data %>%
mutate(
party_clean = case_when(
str_detect(party_id, regex("dem|liberal", ignore_case = TRUE
str_detect(party_id, regex("rep|conserv|gop", ignore_case =
str_detect(party_id, regex("ind|none|other", ignore_case = Tis.na(party_id) ~ "Unknown",
TRUE ~ "Other"
),
approval_numeric = case_when(
approval %in% c("Strongly Approve", "Approve") ~ 1,
approval %in% c("Strongly Disapprove", "Disapprove") ~ 0,
approval == "Neither" ~ 0.5,
TRUE ~ NA_real_
),
deducation_level = case_when(
education = 12 ~ "High School Graduate",
education == 12 ~ "High School Graduate",
education == 16 ~ "College Graduate",
education > 16 ~ "Post-Graduate",
TRUE ~ "Unknown"
)
TRUE ~ "Unknown"
)
%>%
count(party_clean, education_level) %>%
slice_head(n = 10)
```

```
# A tibble: 10 × 3
  party_clean education_level
   <chr>
               <chr>
                                       <int>
1 Democrat
               College Graduate
                                          21
2 Democrat
              High School Graduate
                                          24
3 Democrat
              Less than High School
4 Democrat
              Post-Graduate
                                         102
              Some College
                                          91
5 Democrat
6 Independent College Graduate
7 Independent High School Graduate
                                          21
8 Independent Less than High School
9 Independent Post-Graduate
10 Independent Some College
```

Al Prompt to Create:

"I have messy survey data with party_id (text), approval (categories), and education (years). Clean party affiliations using text detection, convert approval to numeric (1=approve, 0=disapprove, 0.5=neutral), and create education levels with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain how str_detect with regex works for party identification and why ignore_case = TRUE is important. How does NA_real_ differ from just NA?"

Example 12: Policy Position Recoding

Research Question: How can we simplify complex policy positions into analyzable categories?

```
voting_data %>%
mutate(
healthcare_position = case_when(
healthcare_vote == "Yes" & medicare_expansion == "Yes" ~ "Pr
healthcare_vote == "Yes" & medicare_expansion == "No" ~ "Mod
healthcare_vote == "No" & medicare_expansion == "No" ~ "Cons
healthcare_vote == "Abstain" | medicare_expansion == "Abstai
TRUE ~ "Mixed"
),
environmental_stance = case_when(
climate_action == "Yes" & clean_energy == "Yes" ~ "Pro-Envir
climate_action == "No" & clean_energy == "No" ~ "Anti-Enviro
TRUE ~ "Mixed Environment"
),
fiscal_ideology = case_when(
tax_cuts == "Yes" & spending_cuts == "Yes" ~ "Fiscal Conserv
tax_cuts == "No" & spending_cuts == "No" ~ "Fiscal Liberal",
tax_cuts == "Yes" & spending_cuts == "No" ~ "Tax Cutter",
tax_cuts == "No" & spending_cuts == "Yes" ~ "Spending Cutter
TRUE ~ "Fiscal Moderate"
)
) %>%
count(healthcare_position, environmental_stance, fiscal_ideology
slice_head(n = 10)
```

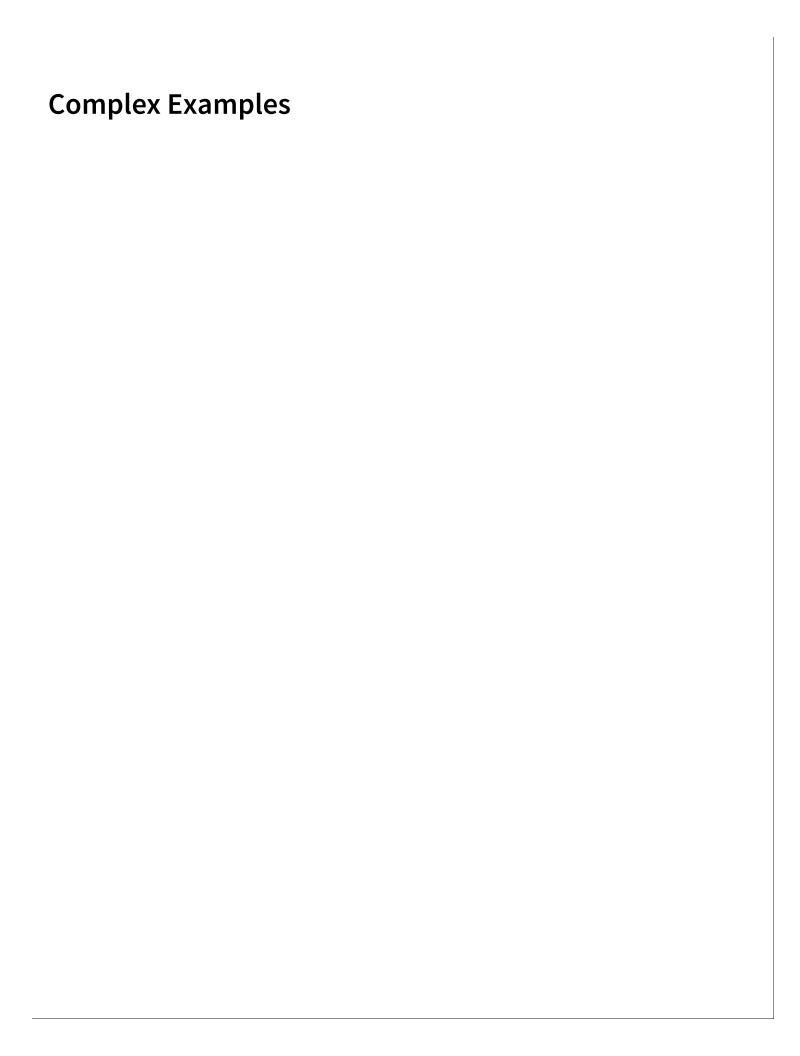
```
# A tibble: 10 × 4
  healthcare_position environmental_stance fiscal_ideology
  <chr>
<int>
                                           Fiscal Conservative
1 Conservative
                      Anti-Environment
                                            Fiscal Liberal
2 Conservative
                      Anti-Environment
3 Conservative
                       Anti-Environment
                                            Fiscal Moderate
4 Conservative
                      Anti-Environment
                                            Spending Cutter
5 Conservative
                      Mixed Environment
                                            Fiscal Conservative
                                            Fiscal Liberal
6 Conservative
                      Mixed Environment
```

Al Prompt to Create:

"I have legislator voting data with Yes/No/Abstain votes on healthcare, medicare_expansion, climate_action, clean_energy, tax_cuts, and spending_cuts. Create position categories combining related votes into meaningful ideological groups with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain how combining multiple vote conditions creates more nuanced categories than single votes. Why use 'Mixed' and 'Unclear' categories instead of forcing binary classifications?"



Complex Example 1: Comprehensive Electoral Analysis

Research Question: What factors predict electoral success across different types of races and regions?

```
1 electoral_analysis <- complex_election_data %>%
2  # Filter to competitive races only
3  filter(vote_margin < 0.1) %>%
4
5  # Create new variables
6  mutate(
7  spending_advantage = (candidate_spending - opponent_spending)
8  (candidate_spending + opponent_spending),
9  experience_level = case_when(
10  prior_office == "None" ~ "Newcomer",
11  prior_office %in% c("Local", "County") ~ "Local Experience",
12  prior_office %in% c("State", "Federal") ~ "Experienced",
13  TRUE ~ "Other"
14  ),
15  demographic_match = case_when(
16  candidate_race == district_majority_race ~ "Matches District
17  TRUE ~ "Different from District"
18  ),
19  upset_potential = if_else(
20  polling_deficit > 5 & spending_advantage > 0.2,
21  "High", "Low"
22  )
23  ) %>%
24
25  # Group analysis by region and office type
26  group_by(region, office_type, experience_level) %>%
27  summarise(
28  tech_account of the content of the
```

```
# A tibble: 8 × 7
 region
           office_type experience_level total_races
avg_spending_advantage
 <chr>
                       <chr>
                                              <int>
<dbl>
1 Midwest Governor
                       Newcomer
0.151
           Senate
                       Local Experience
2 West
-0.0585
3 South
           Governor
                       Newcomer
-0.132
4 Northeast Governor
                       Newcomer
0.150
5 Midwest
           Senate
                       Experienced
-0.214
6 West
           Governor
                       Experienced
-0.125
```

Al Prompt to Create:

"I have comprehensive election data with candidate info, spending, demographics, and results. Create an analysis that: 1) filters to competitive races, 2) creates spending advantage and experience metrics, 3) groups by region/office/experience, 4) calculates win rates and other success metrics with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain this complex analysis workflow. How do the filter, mutate, group_by, and summarise functions work together? Why use .groups = 'drop' and what does na.rm = TRUE do in the mean calculations?"

Complex Example 2: Public Opinion Trend Analysis

Research Question: How do public opinion trends vary across demographic groups?

```
opinion_trends <- complex_polling_data %>%
 # Create time-based variables
   year_month = floor_date(survey_date, "month"),
    approval_binary = if_else(approval_rating > 50, 1, 0),
   TRUE ~ "Other Issues'
   age_generation = case_when(
     age < 30 ~ "Gen Z",
age < 45 ~ "Millennial",
age < 65 ~ "Gen X",
     TRUE \sim "Boomer+"
 # Group by demographics and time
 group_by(year_month, age_generation, party_id) %>%
  summarise(
   sample_size = n(),
   avg_approval = mean(approval_rating),
   approval_rate = mean(approval_binary)
   economic_optimism = mean(economic_outlook == "Positive"),
```

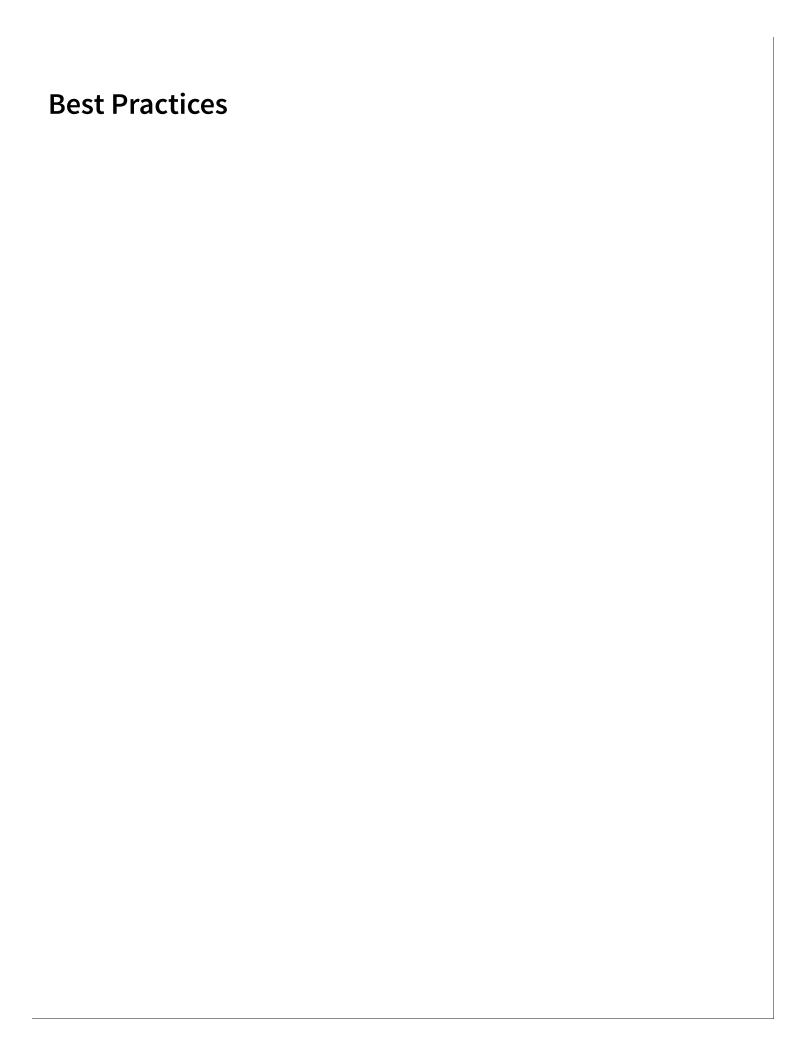
```
# A tibble: 8 × 9
 year_month age_generation party_id
                                       sample_size avg_approval
approval_rate
 <date>
            <chr>
                            <chr>
                                              <int>
                                                           <dbl>
<dbl>
1 2023-06-01 Boomer+
                            Democrat
                                                            52.3
0.583
2 2023-06-01 Boomer+
                            Independent
                                                            42.1
3 2023-06-01 Boomer+
                            Republican
                                                 10
                                                            53.9
4 2023-06-01 Gen X
                            Democrat
                                                            48.1
5 2023-06-01 Gen X
                            Independent
                                                            51.8
6 2023-06-01 Gen X
                            Republican
                                                            57.8
0.6
```

Al Prompt to Create:

"I have longitudinal polling data with dates, demographics, approval ratings, and issue priorities. Create an analysis that: 1) creates time and demographic categories, 2) groups by multiple variables for monthly trends, 3) calculates trend changes with lag functions with tidyverse. Here is what the data look like <insert glimpse>"

Al Prompt to Explain:

"Explain this complex analysis workflow. How do multiple group_by operations work? What does lag() do for trend calculation? How does the combination of filter, mutate, group_by, and summarise create this analysis?"



Writing Effective AI Prompts

Structure your requests clearly:

- 1. Data description: "I have a dataset with variables X, Y, Z..."
- 2. **Goal**: "I want to create/calculate/recode..."
- 3. **Specific requirements**: "Show me the code and explain each step" **Example prompt structure**:

"I have campaign finance data with variables: candidate_name, total_spent, total_raised, vote_share, and outcome. I want to create efficiency metrics including cost per vote and spending per percentage point of vote share with tidyverse. Here is what the data look like <insert glimpse>. Show me the mutate code and explain how each calculation works."

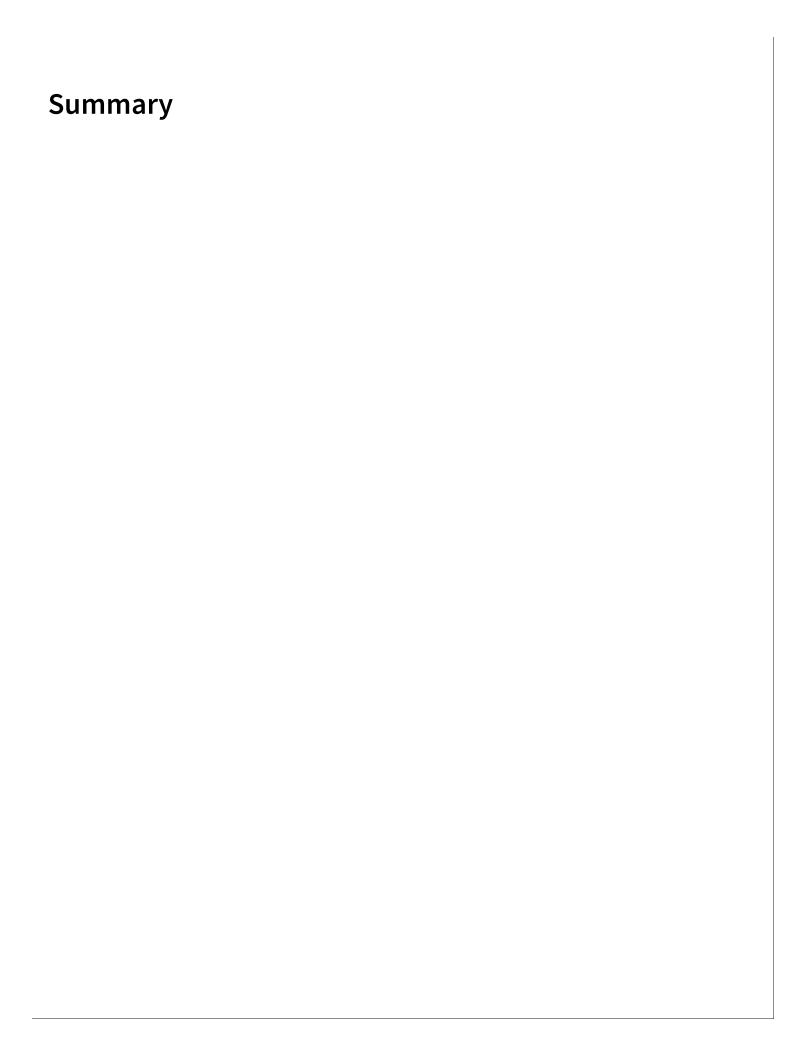
Common Mistakes to Avoid

Logical operator confusion:

- Use & for AND, | for OR in conditions
- Remember operator precedence with parentheses case_when() ordering:
 - Conditions are evaluated in order
 - More specific conditions should come first
 - Always include a catch-all with TRUE

Missing value handling:

- Use na. rm = TRUE in summary functions when needed
- Consider is na () for missing value conditions



Key Functions Learned

Variable creation:

- mutate(): Create new variables with calculations
- case_when(): Complex conditional recoding
- if_else(): Simple binary recoding

Advanced summarization:

- summarise() with proportions and percentages
- group_by() with multiple grouping variables
- pull() for extracting single values

Mathematical and logical operations:

- Arithmetic calculations and ratios
- Logical comparisons and boolean logic
- Text pattern matching with str_detect()



Speaker notes