

# **WEEK 2**

# **TIDYVERSE & MANIPULATING DATA**

DATA VISUALIZATION FOR SOCIAL SCIENTISTS

LECTURER: JEFFREY ZIEGLER, PHD

TEACHING FELLOW: SHEKHAR KEDIA

ASDS - TRINITY COLLEGE DUBLIN

SPRING 2026

# ROAD MAP FOR TODAY

## ■ Today:

- ▶ Work with rectangular data using `gapminder` dataset
- ▶ Manipulate data with core `dplyr` verbs
- ▶ Summarize amounts and proportions and plot them with `ggplot2`
  - Practice histograms, density plots, ridge plots, boxplots, violins, and jittered points
  - Emphasize choices (bin width, bandwidth, outlier rules) that change how uncertainty appears

## ■ By next week, please...

- ▶ Problem set #2

# KEY DPLYR VERBS

- **filter()**: extract rows (cases) that meet logical conditions
- **select()**: choose columns (variables) from a data frame
- **arrange()**: sort rows by one or more variables
- **mutate()**: create or transform columns
- **group\_by() + summarize()**: compute group-level summaries

## FILTER(): SUBSETTING ROWS

- Pattern: `filter(data, condition)` returns rows where condition is TRUE
- Example: `filter(gapminder, country == "Denmark")` returns all Denmark rows
- Uses logical operators such as `==`, `<`, `>`, `<=`, `>=`

# LOGICAL TESTS

- $x < y$ : less than
- $x > y$ : greater than
- $x == y$ : equal to (used for testing, not assignment)
- $x \leq y, x \geq y$ : less/greater than or equal to

## EX: STUDYING AFRICAN FARMER-LED IRRIGATION

- SAFI study looked at farming and irrigation methods in Tanzania and Mozambique
- Survey data was collected through interviews conducted between November 2016 and June 2017

| Variable         | Description   |
|------------------|---|
| village          | Village name  |
| interview_date   | Date of interview   |
| no_membrs        | How many members in the household?  |
| years_liv        | How many years have you been living in this village or neighboring village? |
| rooms            | How many rooms in the main house are used for sleeping?                     |
| memb_assoc       | Are you a member of an irrigation association?                              |
| affect_conflicts | Have you been affected by conflicts with other irrigators in the area?      |
| liv_count        | Number of livestock owned   |
| no_meals         | How many meals do people in your household normally eat in a day?           |

## EX: FILTERING WITH FILTER()

```
1 # Import dataset and load library
2 SAFI <- read.csv("https://raw.githubusercontent.com/ASDS-TCD/
   DataViz_2026/refs/heads/main/datasets/SAFI.csv")
3
4 # All households in village "God"
5 SAFI_god <- filter(SAFI, village == "God")
6
7 # Households with more than 6 members
8 large_households <- filter(SAFI, no_membrs > 6)
9
10 # Households in "God" or "Ruaca" with more than 4 members
11 god_ruaca_large <- filter(SAFI, village %in% c("God", "Ruaca") &
   no_membrs > 4)
```

## EX: SELECTING WITH SELECT()

```
1 # Keep only identification and location variables
2 id_loc <- SAFI |>
3   select(key_ID, village)
4
5 # Move village to the first position
6 village_first <- SAFI |>
7   select(village, everything())
8
9 # Drop columns related to survey timing
10 no_timing <- SAFI |>
11   select(-starts_with("interview_date"))
```



## Ex: ARRANGING DATA WITH ARRANGE()

```
1 # Households ordered by size (ascending)
2 SAFI_by_size <- SAFI |>
3   arrange(no_membrs)
4
5 # Households ordered by size (descending) within village
6 SAFI_by_village_size <- SAFI |>
7   arrange(village, desc(no_membrs))
```

## MUTATE(): CREATING NEW VARIABLES

- Pattern: `mutate(data, new_var = expression)` adds or modifies columns
- Examples: indicator for African countries, logged GDP per capita, classification of regions
- Often combined with `group_by()` and `summarize()` for richer summaries

## EX: ADDING NEW COLUMNS WITH MUTATE()

```
1 # Add a logical indicator for "large" households
2 SAFI <- SAFI |>
3   mutate(large_hh = no_membrs >= 6)
4
5 # Years of living in village in decades
6 SAFI <- SAFI |>
7   mutate(living_decades = years_liv / 10)
```

## EX: COMBINING MULTIPLE VERBS WITH PIPES (|>)

```
1 # Categorize household size
2 SAFI <- SAFI |>
3   mutate(
4     hh_size_cat = case_when(
5       no_membrs <= 3 ~
6         "small",
7       no_membrs >= 4 & no_membrs
8         <= 7 ~ "medium",
9       no_membrs > 7 ~
10        "large"
11     )
12 )
```

```
1 # Large households in Ruaca with
2   selected columns
3 ruaca_large_small_df <- SAFI |>
4   filter(village == "Ruaca", no_
5     membrs > 6) |>
6   mutate(
7     large_hh = TRUE,
8     members_per_room = no_membrs
9     / rooms
10  ) |>
11  select(village, no_membrs,
12    rooms, members_per_room,
13    large_hh)
```

## TWO-STEP WORKFLOW FOR PROPORTIONS

- Step 1: Summarize data with `dplyr`, often using `group_by()` and `summarize()`
- Step 2: Visualize summarized data with `ggplot2`
- Clarifies how proportions are computed before plotting them

## SUMMARIZING WITH GROUP\_BY() AND SUMMARIZE()

- Workflow: `data |> group_by(group_vars) |> summarize(stats)`
- Typical statistics: mean, median, min, max, count, etc.
- Used to compute quantities like median house members by village

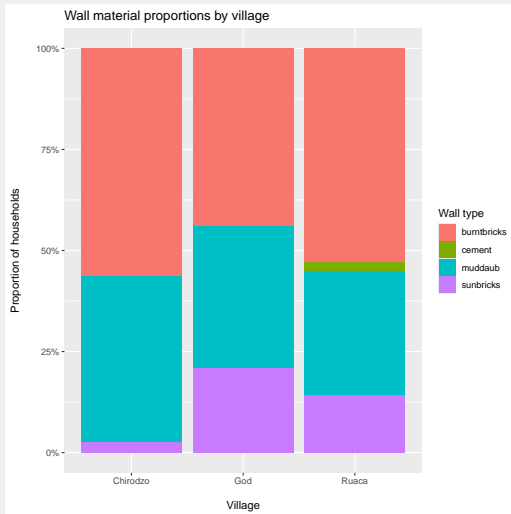
```
1 # Mean household size by village
2 village_hh_summary <- SAFI |>
3   group_by(village) |>
4   summarize(
5     mean_members = mean(no_membrs, na.rm = TRUE),
6     median_members = median(no_membrs, na.rm = TRUE),
7     n_households = n(),
8     .groups = "drop"
9   )
```

## EX: PROPORTIONS WITH GROUP\_BY() |> SUMMARIZE()

```
1 # Mean years in village by roof type
2 wall_living_summary <- SAFI |>
3   count(wall_type, village, name = "n") |>
4   group_by(village) |>
5   mutate(prop = n / sum(n))
```

# EX: PROPORTIONS WITH GROUP\_BY() |> SUMMARIZE()

```
1 ggplot(wall_living_summary,  
2       aes(x = village, y =  
3         prop, fill = wall_type)  
4       ) +  
5       geom_col(position = "fill")  
6       +  
7       scale_y_continuous(labels =  
8         scales::percent_format  
9         ()) +  
10      labs(  
11        x = "\nVillage",  
12        y = "Proportion of  
13          households\n",  
14        fill = "Wall type",  
15        title = "Wall material  
16          proportions by village"  
17      )
```





## FROM SUMMARIES TO GRAPHICS

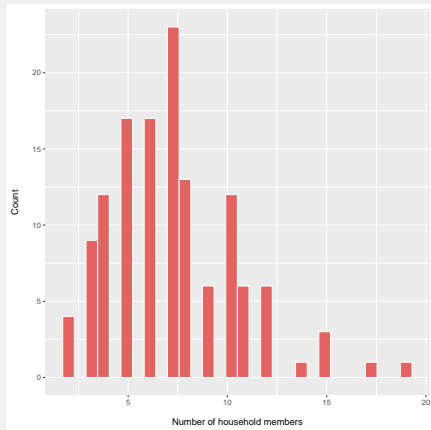
- Use summarized tables (e.g., counts or proportions by group) as ggplot inputs
- Choose appropriate geometries (bars, lines, etc.) to represent amounts and proportions
- Emphasize truthful, clear visual communication of aggregated data

# HISTOGRAMS AND BIN WIDTH

- Histograms show counts of observations within ranges (bins) of a variable such as `no_membrs`
- Bin width strongly affects the apparent shape of the distribution and must be chosen deliberately
- Options include `binwidth`, bar border color, fill color (e.g. `#E16462`), and `boundary` to align bins at whole numbers

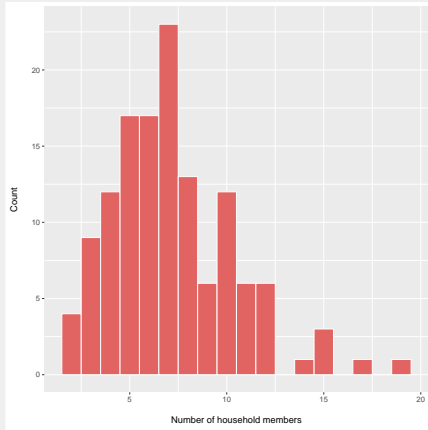
# EX: HISTOGRAM OF HOUSEHOLD SIZE

Default bin widths (=30)



```
1 ggplot(SAFI, aes(x = no_membrs)) +  
2   geom_histogram(color = "white", fill = "#  
   E16462") +  
3   labs(x = "\nNumber of household members",  
4        y = "Count\n")
```

bin=1



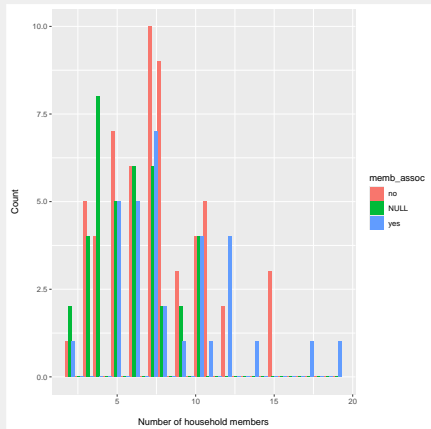
```
1 ggplot(SAFI, aes(x = no_membrs)) +  
2   geom_histogram(binwidth = 1, color = "  
   white", fill = "#E16462") +  
3   labs(x = "\nNumber of household members",  
4        y = "Count\n")
```

# HISTOGRAMS WITH GROUPS

- Additional variables like `memb_assoc` (member of irrigation association) can be mapped to `fill`
- Overlaid histograms with multiple fills can be hard to interpret when groups overlap
- Faceting (e.g. `facet_wrap(vars(village))`) separates groups into small multiples for clearer comparisons

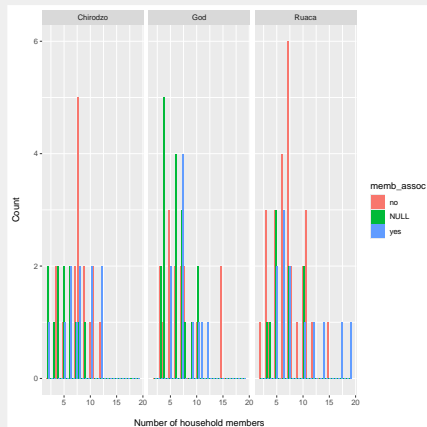
# EX: HISTOGRAMS WITH GROUPS

## Member association



```
1 ggplot(SAFI, aes(x = no_membrs, fill = memb_
  assoc)) +
2 geom_histogram(position="dodge") +
3 labs(x = "\nNumber of household members",
4 y = "Count\n")
```

## Member association by village



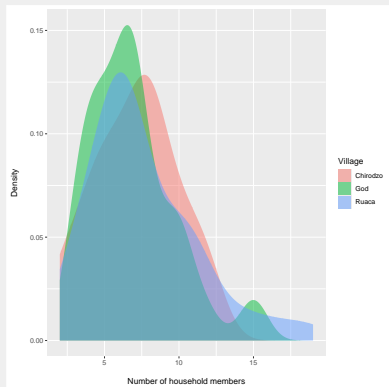
```
1 ggplot(SAFI, aes(x = no_membrs, fill = memb_
  assoc)) +
2 geom_histogram(position="dodge") +
3 facet_wrap(vars(village)) +
4 labs(x = "\nNumber of household members",
5 y = "Count\n")
```

## DENSITY AND RIDGE PLOTS

- Density plots provide a smoothed estimate of the distribution; bandwidth (bw) controls smoothing
- Aesthetics include fill (e.g. #E16462), border color (e.g. #9C3836), and line width
- Mapping fill = drv and using transparency (alpha) overlays group densities; ggribes can convert these to ridge plots with geom\_density\_ridges()

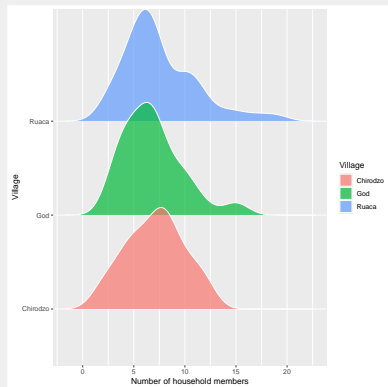
# EX: # OF HOUSEHOLD MEMBERS BY VILLAGE

## Density plots



```
1 ggplot(SAFI, aes(x = no_membrs, fill =  
  village)) +  
2 geom_density(alpha = 0.5, color = NA,  
  adjust = 1) +  
3 labs(x = "\nNumber of household members",  
4 y = "Density\n",  
5 fill = "Village")
```

## Ridge plots



```
1 ggplot(SAFI, aes(x = no_membrs, y = village,  
  fill = village)) +  
2 geom_density_ridges(alpha = 0.7, color = "  
  white", scale = 1.2) +  
3 labs(x = "\nNumber of household members",  
4 y = "Village\n",  
5 fill = "Village")
```

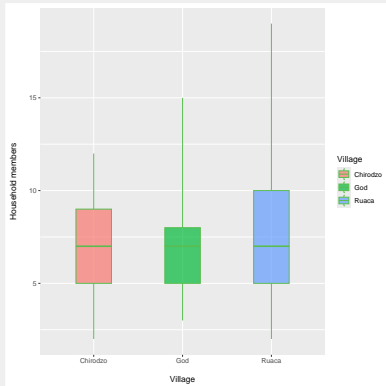
## BOXPLOTS, VIOLINS, AND DOTS

- Boxplots show quartiles, median, and outliers; `coef` controls the outlier rule (e.g.  $5 \times \text{IQR}$ )
- Visual settings include `fill` (e.g. `#E6AD3C`) and `color` (e.g. `#5ABD51`)
- Violin plots show mirrored density shapes; adding jittered points (e.g. `width 0.1`, `size 0.5`) reveals individual observations by category such as `village`



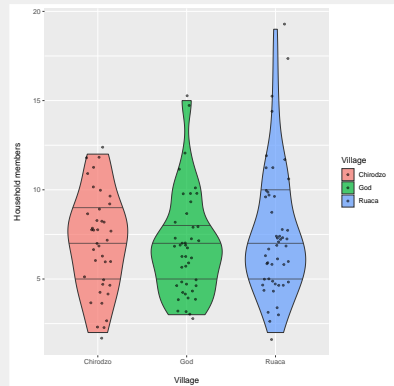
# Ex: # OF HOUSEHOLD MEMBERS BY VILLAGE

## Box plot



```
1 ggplot(SAFI, aes(x = village, y = no_membrs,  
  fill = village)) +  
2   geom_boxplot(alpha = 0.7, color = "#5ABD51",  
  coef = 3, width = 0.4) +  
3   labs(fill = "Village", x = "\nVillage", y =  
  "Household members\n")
```

## Violin plot



```
1 ggplot(SAFI, aes(x = village, y = no_membrs,  
  fill = village)) +  
2   geom_violin(alpha = 0.7, width = 0.8, draw  
  _quantiles = c(0.25, 0.5, 0.75)) +  
3   geom_jitter(width = 0.15, alpha = 0.6,  
  size = 0.8) +  
4   labs(fill = "Village", x = "\nVillage", y =  
  "Household members\n")
```

# CLASS BUSINESS

- Read required (and suggested) online materials
- Fork GitHub repository
- Problem set #2 is up on GitHub