

# Generative AI for Data Analysis

Writing code, method descriptions, and important considerations

Greta Linse, Sara Mannheimer, Ava Yazdian, Sally Slipher

1

---

---

---

---

---

---

---

# Montana State University Library

Guidance on responsible use of generative AI in academic work

2

---

---

---

---

---

---

---

# Using Generative AI:

Writing code and more for data analysis

3

---

---

---

---

---

---

---

## Example use cases

- Translating code from one language to another
- Writing code given a set of instructions/steps/pseudo code
- Create code for a desired plot by describing it in words or by giving it an example

4

---

---

---

---

---

---

---

---

## Translating code from one language to another

- PROMPT: Given the following code translate to R.
- Demonstration:
  - Translate NHANES Sample Code for Logistic Regression that is only available in SAS-callable SUDAAN version 11
- FOLLOW UP PROMPT 1: What will the output look like?
- FOLLOW UP PROMPT 2: Write a results paragraph using these results.

5

---

---

---

---

---

---

---

---

## SAS code

```
*****;
**Linear trends analysis by age group using Logistic Regression**
*****;

proc logist data = Rx FILETYPE = SAS DESIGN = WR;
  * Next statement: PSUs nested within SUDAAN accounts for the design effects*;
  NEST SOMVSTRA SOMVPSU / MISSUNIT;
  * Weight statement: specify appropriate weight, accounts for the unequal probability of sampling and non-response*;
  WEIGHT wintdy;
  * Class statement: specify categorical variable(s). Reflevel statement can be included to choose reference category
  for the categorical variables. By default SUDAAN uses the highest category*;
  class one / noREQ;
  * Subpop statement: specify the subpopulation of interest (the inclusion criteria)*;
  subpoppx inAnalysis;
  * Model statement: specifies dependent variable and independent variable(s)*;
  MODEL_ anyOpioid* ageCat;
  * Output statement: outputs the results to a file*;
  output / betas= default filename=fig3_logistReg_betas_all replace;
  * Test statement: produces statistics and P values for the Satterthwaite adjusted CHI square (satadjchi),
  the Satterthwaite adjusted F (satadjf), and Satterthwaite adjusted degrees of freedom (printed by default).
  If this statement is omitted, the nominal degrees of freedom,
  the WALDF and the p-value corresponding to the WALDF and WALDP will be produced*;
  test waldf satadjf satadjchi;
run;
```

6

---

---

---

---

---

---

---

---

## 8

## 9

## 3

## Resulting R Code

```

# R code for a split-plot experiment design
# This code generates a balanced split-plot design, simulates data, fits a mixed-effects model, and visualizes the results.

# 1. Design Generation
set.seed(1234)
n_reps = 4
n_blocks = 2
n_treatments = 3
n_main_effects = 2
n_sub_effects = 2

# Generate design matrix
design = generate_split_plot_design(
  n_reps = n_reps,
  n_blocks = n_blocks,
  n_treatments = n_treatments,
  n_main_effects = n_main_effects,
  n_sub_effects = n_sub_effects
)

# 2. Data Simulation
simulate_data = function(design) {
  # Simulate data for each treatment and block
  data = data.frame()
  for (block in 1:n_blocks) {
    for (treatment in 1:n_treatments) {
      data = rbind(data, simulate_data_block(
        block = block,
        treatment = treatment,
        design = design
      ))
    }
  }
  return(data)
}

# 3. Model Fitting
fit_model = function(data) {
  # Fit a mixed-effects model using lme4 and DHARMa
  model = lmer(yield ~ treatment + (1|block), data = data)
  return(model)
}

# 4. Visualization
visualize_results = function(model) {
  # Visualize the results of the model fit
  plot_results(model)
}

# Run the entire process
design = generate_split_plot_design(
  n_reps = 4,
  n_blocks = 2,
  n_treatments = 3,
  n_main_effects = 2,
  n_sub_effects = 2
)
data = simulate_data(design)
model = fit_model(data)
visualize_results(model)

```

10

## Knowledge required

- Need to understand how to install packages
  - What is lme4 and DHARMa? Are they actually popular libraries?
- The suggested code uses functions with messages and errors. Seems pretty complicated!
- Have to read the code carefully because there are some optional code chunks that are included.
- Realized that the treatment was included as a random effect instead of a fixed effect.
- Follow up Prompt 1: What is MainPlot?
- Follow up Prompt 2: Add a fixed treatment effect
- Follow up Prompt 3: Give me a concrete example
- Follow up Prompt 4: What will the output look like?

11

## Output from second example

**Crop Yield Split-Plot Experimental Design**

This example document illustrates:

- Hierarchical design with multiple randomization levels
- Simulated crop yield data
- Split-plot plus main effect
- Fixed treatment effect
- Visualization
- Statistical modeling
- Model diagnostics

The experimental factors are:

- MainPlot: Irrigation Methods
- Subplot: Fertilizer Types
- Subsubplot: Crop Varieties
- Treatment: Pest Control Methods

Would you like me to elaborate on any specific aspect?

Block	Treatment	Main Effect	Subplot	Subsubplot	Yield
1	1	1	1	1	1.0
1	1	1	1	2	1.5
1	1	1	1	3	2.0
1	1	1	2	1	1.5
1	1	1	2	2	2.0
1	1	1	2	3	2.5
1	1	1	3	1	2.0
1	1	1	3	2	2.5
1	1	1	3	3	3.0
1	1	2	1	1	1.5
1	1	2	1	2	2.0
1	1	2	1	3	2.5
1	1	2	2	1	2.0
1	1	2	2	2	2.5
1	1	2	2	3	3.0
1	1	2	3	1	2.5
1	1	2	3	2	3.0
1	1	2	3	3	3.5
1	2	1	1	1	2.0
1	2	1	1	2	2.5
1	2	1	1	3	3.0
1	2	1	2	1	2.5
1	2	1	2	2	3.0
1	2	1	2	3	3.5
1	2	1	3	1	3.0
1	2	1	3	2	3.5
1	2	1	3	3	4.0
1	2	2	1	1	2.5
1	2	2	1	2	3.0
1	2	2	1	3	3.5
1	2	2	2	1	3.0
1	2	2	2	2	3.5
1	2	2	2	3	4.0
1	2	2	3	1	3.5
1	2	2	3	2	4.0
1	2	2	3	3	4.5
1	3	1	1	1	3.0
1	3	1	1	2	3.5
1	3	1	1	3	4.0
1	3	1	2	1	3.5
1	3	1	2	2	4.0
1	3	1	2	3	4.5
1	3	1	3	1	4.0
1	3	1	3	2	4.5
1	3	1	3	3	5.0
1	3	2	1	1	3.5
1	3	2	1	2	4.0
1	3	2	1	3	4.5
1	3	2	2	1	4.0
1	3	2	2	2	4.5
1	3	2	2	3	5.0
1	3	2	3	1	4.5
1	3	2	3	2	5.0
1	3	2	3	3	5.5
2	1	1	1	1	1.0
2	1	1	1	2	1.5
2	1	1	1	3	2.0
2	1	1	2	1	1.5
2	1	1	2	2	2.0
2	1	1	2	3	2.5
2	1	1	3	1	2.0
2	1	1	3	2	2.5
2	1	1	3	3	3.0
2	1	2	1	1	1.5
2	1	2	1	2	2.0
2	1	2	1	3	2.5
2	1	2	2	1	2.0
2	1	2	2	2	2.5
2	1	2	2	3	3.0
2	1	2	3	1	2.5
2	1	2	3	2	3.0
2	1	2	3	3	3.5
2	2	1	1	1	2.0
2	2	1	1	2	2.5
2	2	1	1	3	3.0
2	2	1	2	1	2.5
2	2	1	2	2	3.0
2	2	1	2	3	3.5
2	2	1	3	1	3.0
2	2	1	3	2	3.5
2	2	1	3	3	4.0
2	2	2	1	1	2.5
2	2	2	1	2	3.0
2	2	2	1	3	3.5
2	2	2	2	1	3.0
2	2	2	2	2	3.5
2	2	2	2	3	4.0
2	2	2	3	1	3.5
2	2	2	3	2	4.0
2	2	2	3	3	4.5
2	3	1	1	1	3.0
2	3	1	1	2	3.5
2	3	1	1	3	4.0
2	3	1	2	1	3.5
2	3	1	2	2	4.0
2	3	1	2	3	4.5
2	3	1	3	1	4.0
2	3	1	3	2	4.5
2	3	1	3	3	5.0
2	3	2	1	1	3.5
2	3	2	1	2	4.0
2	3	2	1	3	4.5
2	3	2	2	1	4.0
2	3	2	2	2	4.5
2	3	2	2	3	5.0
2	3	2	3	1	4.5
2	3	2	3	2	5.0
2	3	2	3	3	5.5

12