# GenAI to Enhance Your Statistical Programming

Phil Bowsher, Posit; Boston, USA
Cole Arendt, Posit, Boston, USA

## Abstract

This paper will cover examples of using generative artificial intelligence (GenAI) with an emphasis on statistical programming tasks. The information below offers insight and information on using AI in the clinical domain. The following examples are focused on GenAI tools available in an IDE that can help users be a more productive coder in creating new content.

Many organizations are looking to GenAI to help in accelerating the learning of open source languages often while transitioning from other languages. Andy Nicholls (GSK) discusses this topic in his Phuse talk "*Heading to base camp on our way up the open-source mountain…*":
https://phuse.s3.eu-central-1.amazonaws.com/Archive/2023/Connect/EU/Birmingham/PRE_TT09.pdf

Brittany Long (IQVIA) discussed at Phuse how AI can support programming tasks and outlined the following benefits:
- Encourage you to think of new options
- Review your code
- Chain of thought prompting

https://phuse.s3.eu-central-1.amazonaws.com/Archive/2023/Connect/EU/Birmingham/POS_PP14.pdf

The slides for this paper are here:
https://colorado.posit.co/rsc/ai-open-source/

Much of the information below is adapted from the paper "AI Exploration and Innovation for the Clinical Data Scientist" and slides. The following information is focused on information for program leaders and managers, while this paper is for end-users:
https://phuse.s3.eu-central-1.amazonaws.com/Archive/2024/Connect/US/Bethesda/PAP_ET08.pdf
https://phuse.s3.eu-central-1.amazonaws.com/Archive/2024/Connect/US/Bethesda/PRE_ET08.pdf

## Introduction

Many statistical programmers are migrating from commercial software to open source. Below will provide examples for using various GenAI tools with an emphasis on IDE programming aids. Each option below has benefits and challenges. Programmers will also have personal preferences, so it is recommended to experiment to find the right mix.

## GenAI to Enhance Your Code-First Data Science

*It is our belief that software development happens in an IDE*, and the focus of this paper will be on tools that incorporate GenAI in an IDE like RStudio. This concept is often referred to as **Code-First Data Science**: Building your analyses with code as the workflows can be reproduced, adapted, and scaled to solve similar problems in the future.

While there are popular public GenAI LLMs that have public interfaces like **chat.openai.com**, this paper will focus on tools incorporated in an IDE.These tools bring the power of public LLMs within the coding environment.

Many of the examples below will highlight Github Copilot. "Copilot" is becoming a popular term in GenAI. For example, a popular tool used by programmers is Obsidian, a writing app that adapts to the way you think. There is a tool called "Copilot" for Obsidian that is not connected to Github Copilot, but rather a ChatGPT interface inside Obsidian. Other popular tools like JetBrains IDE do have GitHub Copilot support.

The are 5 popular ways of using GenAI for programming in an IDE:

I. GitHub Copilot
   A. RStudio
   B. GitHub Copilot in VS Code
II. ChatGPT or other public LLM via chattr or other R package
III. cursor.sh or other AI IDE
IV. Local Large Language Models via an interface or IDE

**Sample Tasks to Try in GenAI:**
Below we will cover various tools for using GenAI. Here are examples to try using in each tool:

- Try "Write R code that generates an Adverse Summary Table with GT R package"
- Try "Write R code that uses the adae pharmaverse adam dataset with Age greater than 18. Then use the gt R package to generate an Adverse Summary Table with GT R package. Also create a ggplot2 graphic on AESEV count. The first table should show the total number of patients with at least one adverse event. The ggplot2 graphic should show a bar chart with total counts by AGE. Add thorough documentation to your code."

Andreas Handel has great advice on using GenAI for programming. It is recommended to read the sections on **Good prompting and Iterating** here:
https://andreashandel.github.io/MADAcourse/content/module-ai/ai-write-code-r/ai-write-code-r.html

Andreas Handel has some great examples of AI tools for fixing code with examples here:
https://andreashandel.github.io/MADAcourse/content/module-ai/ai-analysis-r/ai-analysis-r.html

The examples point to code examples that users can test here:
https://github.com/andreashandel/MADAcourse/tree/main/code

We encourage users to test out the examples in the various methods outlined below in an IDE.

## I. GitHub Copilot
**GitHub Copilot Information**
GitHub Copilot is an AI pair programmer tool that is available in many IDEs. The tool provides autocomplete-style suggestions while programming. GitHub Copilot is not free and requires a subscription and an active GitHub Copilot account.

### A. Copilot in RStudio
Please see the presentation below from Posit Conf 2023 that introduced Github Copilot in RStudio:
https://colorado.posit.co/rsc/rstudio-copilot/

The talk highlights the differences between autocomplete and Github Copilot:

## Autocomplete vs Copilot

```
1   # Take the mean of the MPG column in mtcars dataset
2
3   mea|
```

| | | |
|---|---|---|
| mean | {base} | mean(x, ...) |
| mean.Date | {base} | **Arithmetic Mean** |
| mean.default | {base} | Generic function for the (trimmed) arithmetic mean. |
| mean.difftime | {base} | Press F1 for additional help |
| mean.POSIXct | {base} | |
| mean.POSIXlt | {base} | |

```
1   # Take the mean of the MPG column in mtcars dataset
2
3   mean(mtcars$mpg)
```

Here is an example of Copilot generated Code:

```
1   # calculate the average fuel efficiency of     Context
    cars grouped by cyl
2
3   library(dplyr)
4
5   mtcars %>%                                      Copilot generated
        group_by(cyl) %>%                          'ghost text'
        summarise(avg_mpg = mean(mpg))
```

**Ghost text doesn't "exist" in document until accepted**

It helps in GenAI to provide context of the programming task. The description below of the objective will show an example. To help achieve better results, programmers will iterate on specific tasks. The talk "A hacker's guide to open source LLMs" provides additional information:

https://www.youtube.com/watch?v=sYliwvml9Es

An example of context is below at the top of the code:

## Simple: Just get the hint words

Break it down into component parts, let's work with the hint words!

```
1   # create a json_url function that accepts a date argument
2   # the date argument should be in the format of "YYYY-MM-DD" but
3   # will need to use stringr to replace the "-" with "/" to match the url format
4   # using glue, create a url that matches
5   # "https://keyword-client-prod.red.aws.wapo.pub/levels/2023/08/09.json"
6   # but replace the 2023/08/09 with the cleaned date
7
8   json_url <- function(date = Sys.Date()){
      date_clean <- stringr::str_replace_all(date, "-", "/")

      url <- glue::glue("https://keyword-client-prod.red.aws.wapo.pub/levels/{date_clean}.json")

      return(url)
    }
```

**GitHub Copilot Setup in RStudio**

There are various versions of the RStudio IDE such as RStudio desktop, RStudio Server, and Posit Cloud, all of which support GitHub Copilot. The downloads for RStudio are here:
https://posit.co/downloads/

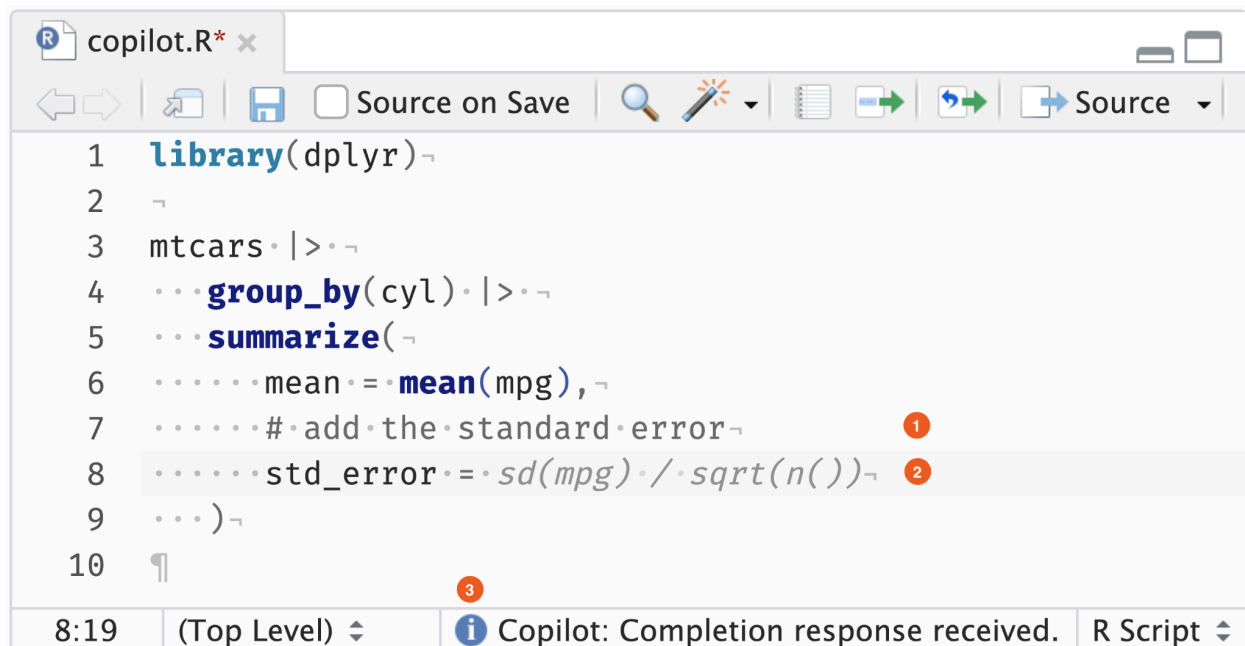Posit Cloud is an online platform for learning R. The article below shows how to use Github Copilot in Posit Cloud:
https://posit.co/blog/github-copilot-on-posit-cloud/

Once a user has a Github Copilot account and subscription, they can configure the connection in RStudio using the following direction:
https://docs.posit.co/ide/user/ide/guide/tools/copilot.html

The link shows the following steps:

- Navigate to Tools > Global Options > Copilot.
- Check the box to "Enable GitHub Copilot".
- Download and install the Copilot Agent components.
- Click the "Sign In" button.
- In the "GitHub Copilot: Sign in" dialog, copy the Verification Code.
- GitHub Copilot: Sign in
- Navigate to or click on the link to https://github.com/login/device, paste the Verification Code and click "Continue".
- GitHub will request the necessary permissions for GitHub Copilot. To approve these permissions, click "Authorize GitHub Copilot Plugin".
- After the permissions have been approved, your RStudio IDE will indicate the currently signed in user.
- Close the Global Options dialogue, open a source file (.R, .py, .qmd, etc) and begin coding with Copilot!

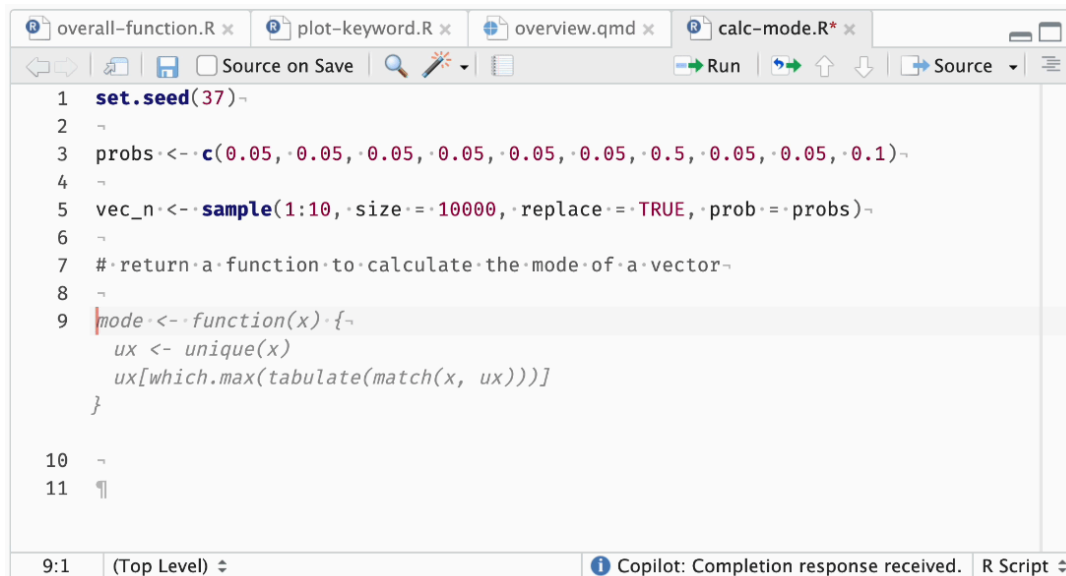Here is an example of Github Copilot in RStudio:



The numbers correspond to:

1. A simple but specific comment providing additional context to Copilot.
2. Copilot's code suggestion, shown in light gray "ghost text".
3. The Copilot status bar, which indicates whether RStudio is waiting on a response to be generated, a completion response has been received, or no completions are available.

**RStudio Copilot Example 1:**

Below is a video showing an example of returning a function to calculate the mode of a vector where the function and output is generated from Github Copilot

**https://colorado.posit.co/rsc/New_Innovations_OpenSource_Drug_Development/images/copilot-gif.gif**



**RStudio Copilot Example 2:**

Below is a video showing an example of creating a Shiny app with Copilot. The gray code is generated by Github Copilot.

https://www.youtube.com/watch?v=40dTx5pEtyQ&t=7s

### B. Copilot in VS Code

There are various versions of VS Code. VS Code for a desktop machine can be downloaded here:
**https://code.visualstudio.com/**

It is common to use Github Copilot in VS Code with the GitHub Copilot Extension:
https://marketplace.visualstudio.com/items?itemName=GitHub.copilot

Additional information about Github Copilot is below::
**https://code.visualstudio.com/blogs/2023/03/30/vscode-copilot**
**https://code.visualstudio.com/docs/editor/github-copilot**
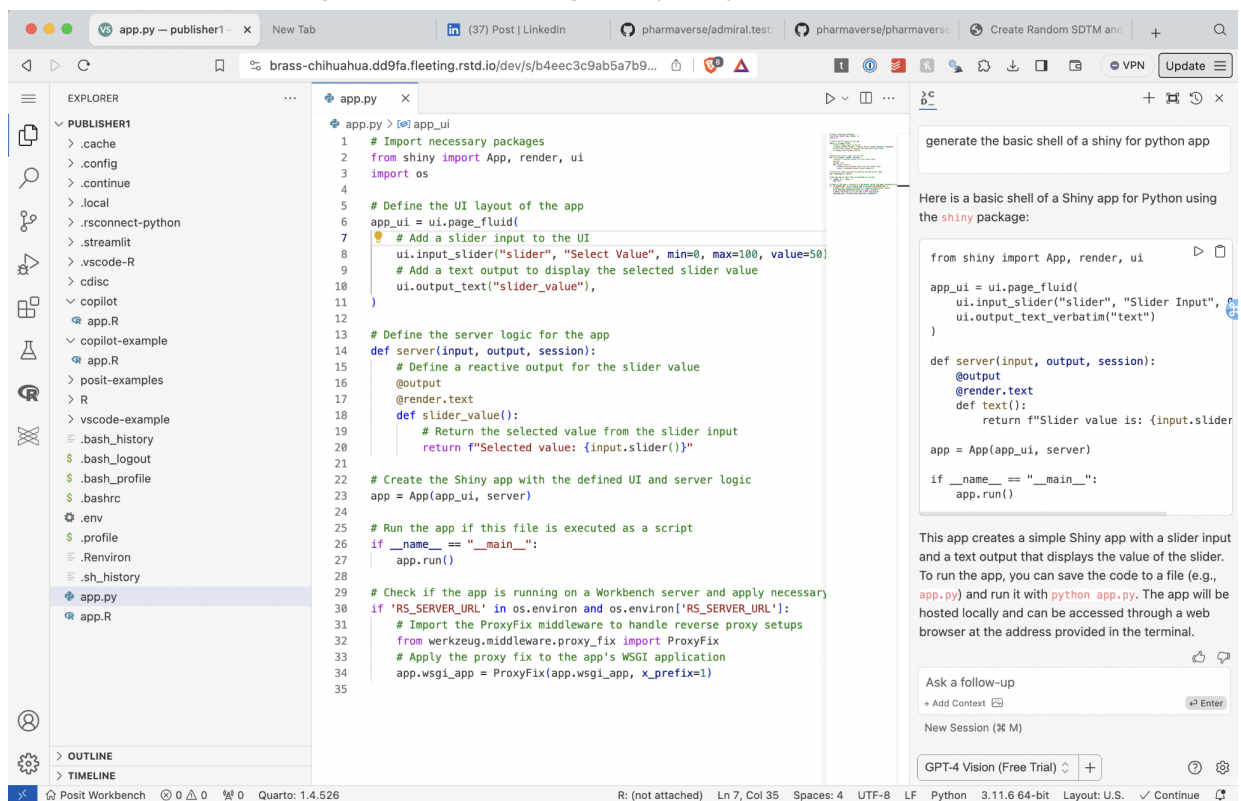**https://docs.github.com/en/copilot/using-github-copilot/getting-started-with-github-copilot**

Devin Pastoor provided the R in Pharma workshop Quarto All the Things! During the workshop, he uses Github Copilot to generate regular expressions within a clinical context as seen here:
https://youtu.be/k-dQ36sx4Rk?t=8276

## VS Code Copilot Example 1:

Below is a screenshot showing an example of creating a Shiny for Python app with Copilot in VS Code:



Other Python examples to test can be found here:
https://github.com/sol-eng/python-examples

## II.    ChatGPT

ChatGPT (Chat Generative Pre-Training Transformer) provides access to a public LLM via an interface. Below we will show tools that bring ChatGPT with an IDE. There are many extensions for using ChatGPT in VS Code. These can be found by doing a search in VS Code.

## ChatGPT in RStudio with chattr R Package

Below we will feature how to access ChatGPT in RStudio. ChatGPT provides a REST API endpoint that can be accessed with a secret key as explained here:
https://platform.openai.com/account/api-keys

The chattr package will discover the secret key via the Environment Variable called OPENAI_API_KEY. The usethis R package can be used to set the variable with the code below: usethis::edit_r_environ()

Create a OPENAI_API_KEY and set it equal to the OpenAI key. Save the file and restart R. Function chattr_test() will confirm that the connection works:
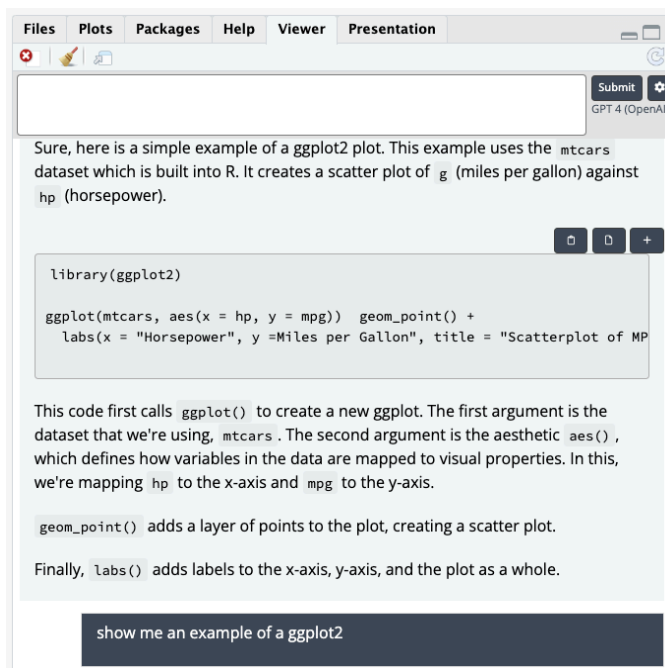
```
> chattr::chattr_test()

— Testing chattr
● Provider: Open AI - Chat Completions
● Path/URL: https://api.openai.com/v1/chat/completions
● Model: gpt-4
✔ Connection with OpenAI cofirmed
|--Prompt: Hi!
|--Response: Hi! I see that you have specific requirements for the code. Could you
please provide more details on what you would like the code to do?
```
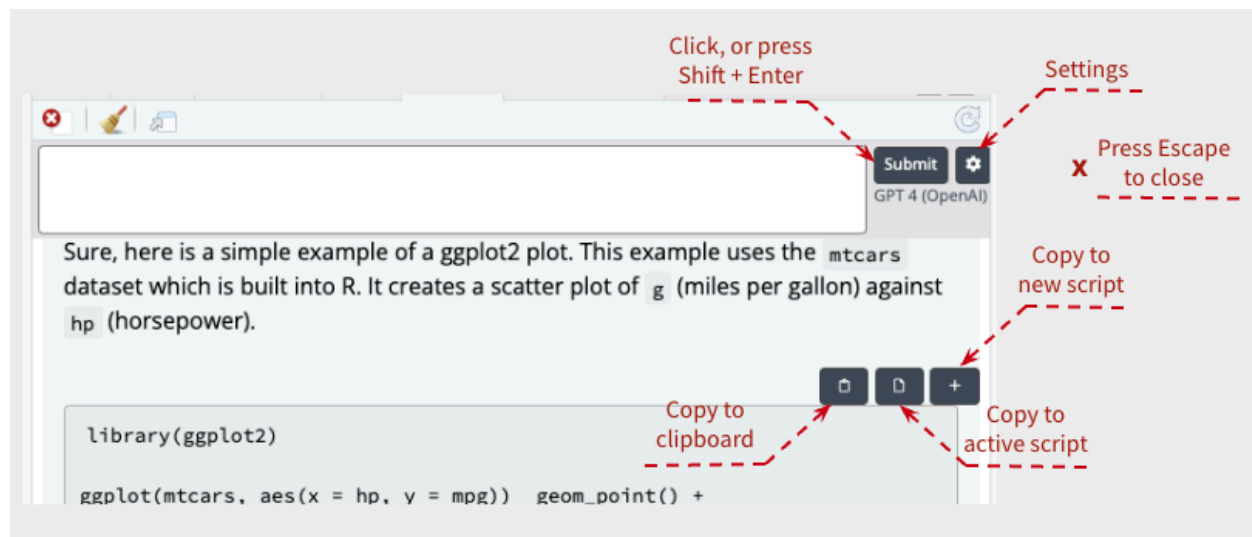
## chattr Example 1:

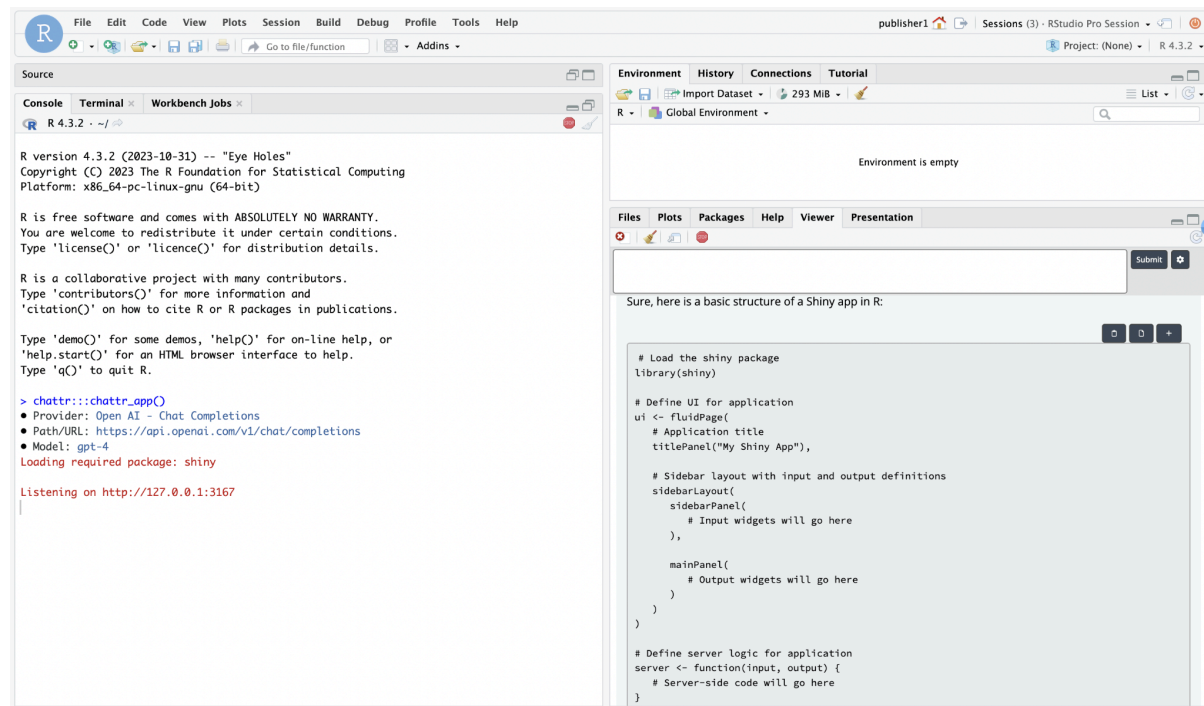Below is an example of generating ggplot2 code with ChatGPT via chattr.
https://blogs.rstudio.com/ai/posts/2024-04-04-chat-with-llms-using-chattr/

## chattr Example 2:

Below is an example of generating a Shiny app with ChatGPT via chattr.
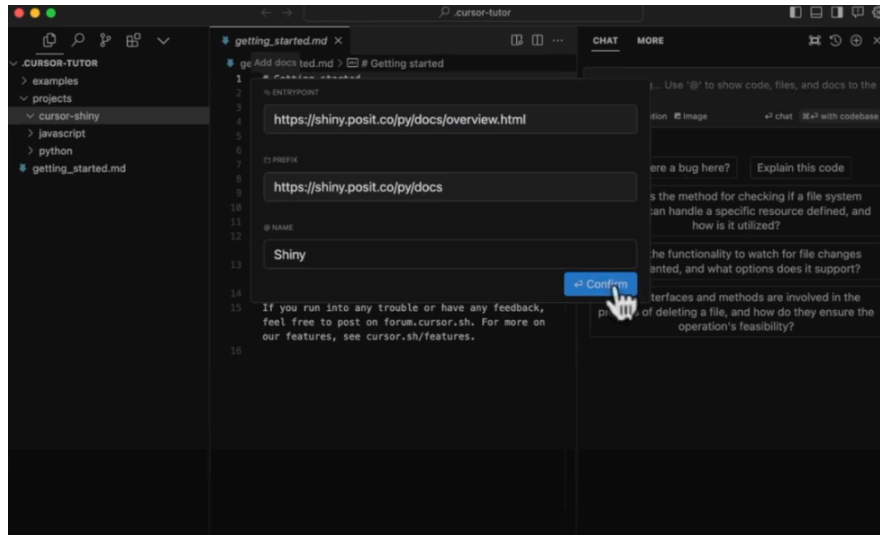


### III. Example with cursor.sh

cursor.sh is an IDE built for pair-programming. It supports **RAG** (Retrieval Augmented Generation). **RAG** was introduced by Meta in 2020 here: https://arxiv.org/abs/2005.11401v4

cursor.sh can be supplemented with current documentation and additional context via RAG for enhancing the programming experience as seen below.

**Cursor Example 1:**

Below is a video of using Cursor.sh with RAG. The video shows the user asking for Shiny for Python examples yet the LLM was trained before it existed and references another tool. Then the user loads the documentation in for Shiny for Python and is able to get an example!

https://www.youtube.com/watch?v=megC3ulM1Wo&t=14s



### IV. Local Models with chattr

Sharon Machlis covers "5 easy ways to run an LLM locally" that discusses various ways to deploy a LLM on a local system: https://www.infoworld.com/article/3705035/5-easy-ways-to-run-an-llm-locally.html

Below is an example diagram showing the internal nature of local LLMs:

OpenAI

Firewall

My Network

The Local Approach

Local LLM

My Computer

chattr will work with local models on a local machine (laptop, server etc.) with LlamaGPTJ-chat. The local models supported as April 2024 are:

- GPT-J (ggml and gpt4all models)
- LLaMA (ggml Vicuna models from Meta)
- Mosaic Pretrained Transformers (MPT)

.

**chattr Example 2:**
Below is an example showing chattr working with a local LLM on the same server as Posit Workbench:

Sure! Here's a boilerplate code for a Shiny app in R that includes a pretty graph using the ggplot2 package:

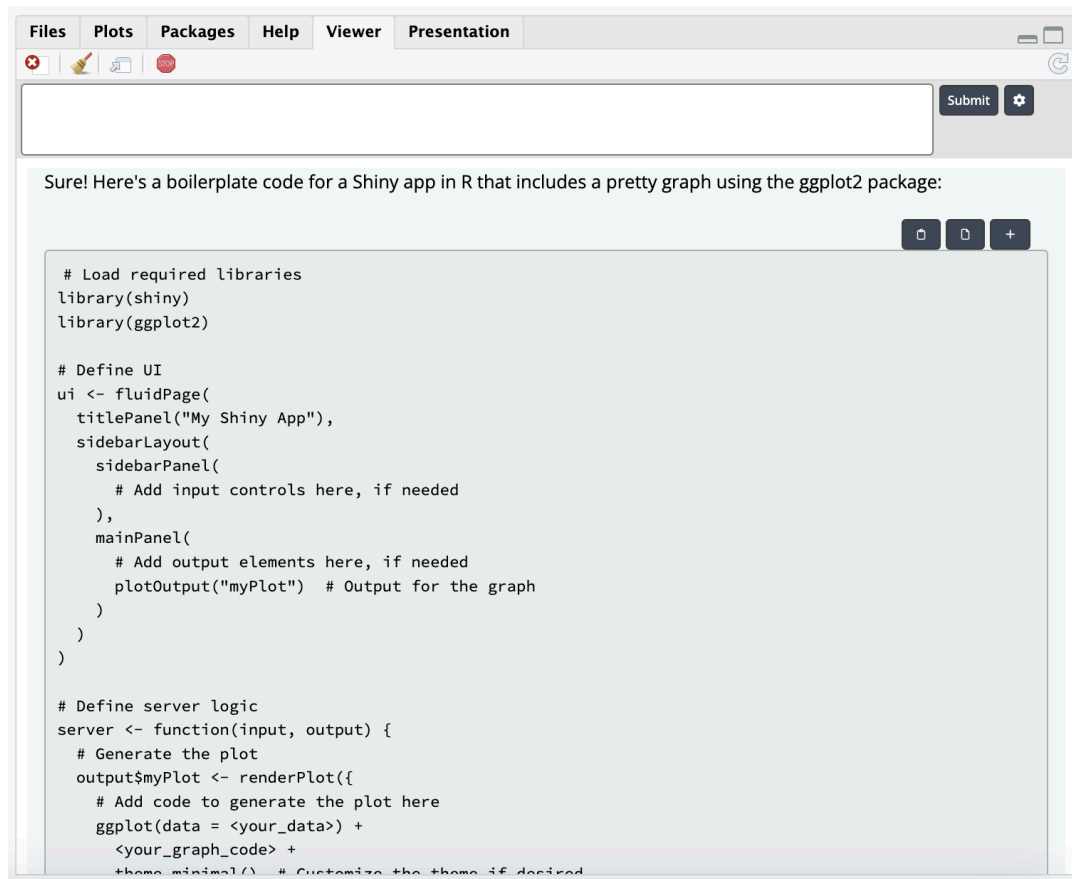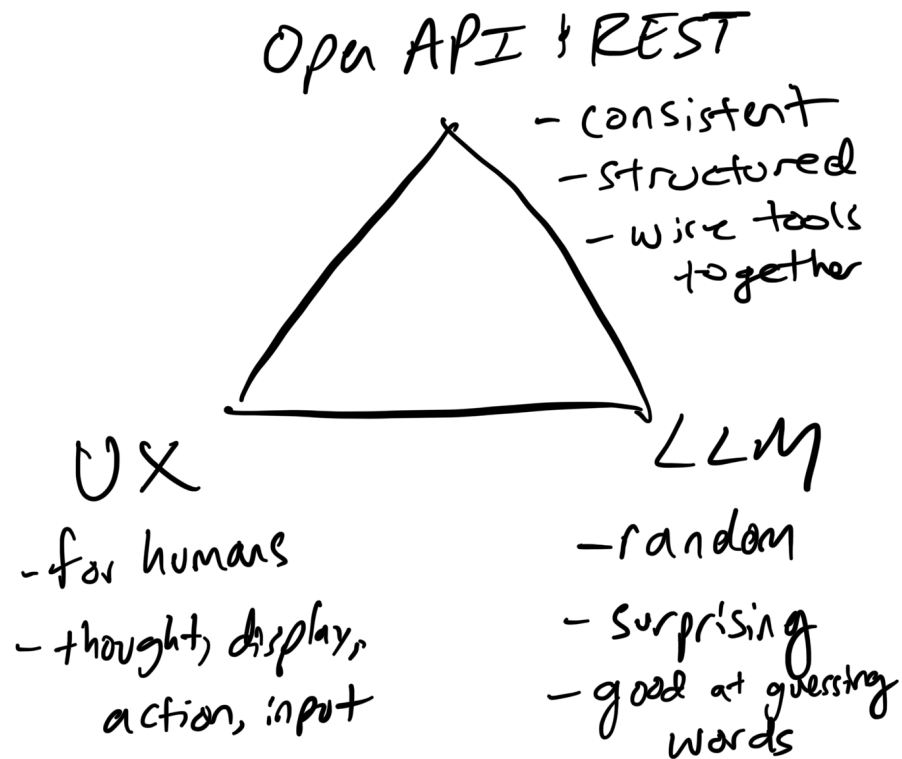```r
# Load required libraries
library(shiny)
library(ggplot2)

# Define UI
ui <- fluidPage(
  titlePanel("My Shiny App"),
  sidebarLayout(
    sidebarPanel(
      # Add input controls here, if needed
    ),
    mainPanel(
      # Add output elements here, if needed
      plotOutput("myPlot")  # Output for the graph
    )
  )
)

# Define server logic
server <- function(input, output) {
  # Generate the plot
  output$myPlot <- renderPlot({
    # Add code to generate the plot here
    ggplot(data = <your_data>) +
      <your_graph_code> +
      theme_minimal()  # Customize the theme if desired
```

While we believe that software development and data science happens in an IDE, there are many valid reasons to host an interface for non-coding professionals. The focus here is on optimizing tools such as:

- Computing (OpenAPI & Rest)
- User Experience (UX)
- Large Language Models (LLM)

Opa API & REST
- consistent
- structured
- wire tools together

UX
- for humans
- thoughts, displays, action, input

LLM
- random
- surprising
- good at guessing words

If an organization has an internal local model (LLM), API (plumber), Shiny or other interface (Streamlit etc.) can be used to share it with non-programming professionals via an OpenAI type of interface. With Shiny, there are various options for surfacing the local LLM such as Shiny via R, Shiny with Reticulate, or Shiny for Python. Below is an example using Shiny for Python.
https://github.com/wch/chatstream/blob/main/examples/doc_query/app.py

## Conclusion
This paper was to highlight practical examples of how statistical programmers can use GenAI for learning open source for statistical programming tasks. As the use of AI increases by clinical statistical programmers, it is important organizations understand the tools to provide to users. We discussed common tools for using GenAI as well as how to access local models. The information outlined in this paper highlights the quickly changing space of GenAI and examples of how it can be used today by statistical programmers to learn open source.

## CONTACT INFORMATION
Your comments and questions are valued and encouraged. Contact the author at:
Phil Bowsher
phil@posit.co