

# DSX Data Scientist Coding Assessment

## [DSX Data Scientist Coding Assessment](#)

[Introduction](#)

[Assessment Format](#)

[Questions](#)

[What to Submit](#)

[Timeline](#)

[Evaluation Process](#)

[Posit Cloud](#)

[Required R Packages](#)

[R Version](#)

[Learning Resources](#)

[R Package Development](#)

[Pharmaverse & CDISC](#)

[Data Manipulation & Visualization](#)

[Testing & Best Practices](#)

[Additional Resources](#)

[R Package Development](#)

[Question 1: Descriptive Statistics](#)

[Pharmaverse](#)

[Question 2: SDTM DS Domain Creation using {sdm.oak}](#)

[Question 3: ADaM ADSL Dataset Creation](#)

[Question 4: TLG - Adverse Events Reporting](#)

[Python Coding assessment](#)

[Question 5: Clinical Data API \(FastAPI\)](#)

[Question 6: GenAI Clinical Data Assistant \(LLM & LangChain\)](#)

[FAQ](#)

# Introduction

Welcome to the R Package Development, Pharmaverse Expertise, and Python Coding Assessment. This assessment is designed to evaluate your skills and proficiency in the following areas:

- **R Package Development:** Creating R packages that are well-structured, documented, and maintainable.
- **Pharmaverse Ecosystem:** Working with open-source packages designed for clinical trial data standards, including SDTM (Study Data Tabulation Model) and ADaM (Analysis Data Model).
- **Data Manipulation & Derivations:** Utilizing modern R tools like {dplyr}, {tidyr}, and Pharmaverse packages such as {admiral}, {sdm.oak}, and {gtsummary}.
- **Clinical Reporting:** Generating Tables, Listings, and Graphs (TLGs) for regulatory submissions
- **Python:** Applying Python for data science
- **Clean Code Practices:** Writing reproducible, efficient, and well-documented R code
- **Problem-Solving & Adaptability:** Demonstrating the ability to quickly learn new tools and apply them effectively to solve challenges.
- **Efficient use of AI:** Leverage modern AI tools to accelerate learning new concepts and complete coding exercises.

The assessment includes practical coding exercises focused on clinical trial reporting in the pharmaceutical industry and core software development. It is designed to evaluate your technical expertise, critical thinking, and approach to software design, and candidates are allowed to use AI-enabled coding assistants.

## Assessment Format

### Questions

This coding assessment consists of **6 questions**:

- **Question 1:** R package development
- **Questions 2-4:** SDTM/ADaM and Table creation using open source R packages from Pharmaverse
- **Questions 5-6:** Python questions

Candidates are **not required** to complete all questions and may choose which ones to answer. However, answering more or all of the questions is recommended as it will increase the likelihood of advancing to the subsequent stage in the interview process.

## What to Submit

### 1. GitHub Repository Link

- Create a public GitHub repository for your solutions
- Repository should be well-organized with clear folder structure (one for each question)
- Include a comprehensive README.md explaining the repo structure, contents of each folder that could be helpful for the reviewer.

### 2. Code Files

- Each question should have a dedicated folder and its own file (e.g., `question_1.R`, `question_2.R`, etc.)
- For R packages, follow standard package structure (DESCRIPTION, NAMESPACE, R/, man/)
- Include comments explaining key logic, especially for derivations

### 3. Video Explanation (2 minutes)

- Record a brief screen-share walkthrough of your repo
- Explain your approach, design decisions, and key challenges overcome
- Discuss what you learned from the exercise

## Timeline

- **Submission Deadline:** 2 weeks
- **Submission Format:** GitHub repo link + video link submitted as part of the Phenom interview
- **Interview Scheduling:** Selected candidates will be contacted for follow-up discussion

## Evaluation Process

Your submission will be evaluated on:

- **Code Quality:** Clarity, efficiency, and adherence to R best practices
- **Correctness:** Does the code produce the expected output?
- **Documentation:** Are functions and code well-documented?
- **Problem-Solving:** How creative and thorough are your solutions?
- **Communication:** Can you explain your approach clearly?

## Posit Cloud

Candidates can use a posit cloud free plan or any other plan as mentioned in the Posit website (<https://posit.cloud/plans>) to work on R related questions. Once logged in, create a new workspace to clone their Github repository to get started.

## Required R Packages

The necessary packages, along with all dependencies, can be installed from CRAN using the provided command. Please note this is only an example; additional R packages will be needed to complete the assessment.

None

```
install.packages(c("admiral", "sdm.oak", "gt", "ggplot2"))
```

## R Version

R 4.2.0 and above.

## Learning Resources

Each question has question specific learning resources and a hint to answer the questions. Also, here is the consolidated list of all learning resources

### R Package Development

- **R Packages (2nd Edition):** <https://r-pkgs.org/>
  - Chapters: Package structure, Documentation (Roxysen2), Testing, Namespace
- **Writing R Extensions:** <https://cran.r-project.org/doc/manuals/R-exts.pdf>
- **usethis Package Documentation:** <https://usethis.r-lib.org/>

### Pharmaverse & CDISC

- **Pharmaverse Examples:** <https://pharmaverse.github.io/examples/>
- **{admiral} Documentation:** <https://pharmaverse.github.io/admiral/>
- **{sdm.oak} Documentation:** <https://pharmaverse.github.io/sdm.oak/>
- **CDISC Standards:**
  - SDTM IG: <https://www.cdisc.org/standards/foundational/sdmig>
  - ADaM IG: <https://www.cdisc.org/standards/foundational/adam>
- **Coursera:** [Hands On Clinical Reporting Using R](#)
- **Pharmaverse site:** <https://pharmaverse.org/>

### Data Manipulation & Visualization

- **dplyr Documentation:** <https://dplyr.tidyverse.org/>
- **tidyr Documentation:** <https://tidyr.tidyverse.org/>
- **ggplot2 Documentation:** <https://ggplot2.tidyverse.org/>
- **gt (Great Tables):** <https://gt.rstudio.com/>

## Testing & Best Practices

- **testthat Package:** <https://testthat.r-lib.org/>
- **R Style Guide:** <https://style.tidyverse.org/>
- **Advanced R (OOP & Functional Programming):** <https://adv-r.hadley.nz/>

## Additional Resources

- **Pharmaverse Blog:** Articles on Pharmaverse packages and clinical trial data programming <https://pharmaverse.github.io/blog/>
- **YouTube:** Search for "admiral R package tutorial" or "SDTM programming in R using {sdm.oak} package" to view the videos in "R in Pharma" channel. There are many more videos in R in Pharma or R consortium youtube channels.
- **RStudio Community:** <https://community.rstudio.com/> (Q&A forum)

## R Package Development

Reference Book: R Packages by Hadley Wickham and Jennifer Bryan <https://r-pkgs.org/>

### Question 1: Descriptive Statistics

Objective:

Create a well-structured R package called `descriptiveStats` that implements descriptive statistics functions. This question tests your understanding of R package structure, documentation, and best practices.

None

```
### Question 1: R Package Development - Descriptive Statistics
```

```
#### Requirements
```

```
1. Package Structure
```

- Create proper DESCRIPTION file with metadata (Title, Version, Authors, License, Description)
- Create NAMESPACE file with appropriate exports
- Organize R code in the `R/` folder

```
2. Functions to Implement
```

```
Create the following functions with full Roxygen2 documentation:
```

- `calc\_mean(x)` - Calculate arithmetic mean
- `calc\_median(x)` - Calculate median

- ``calc_mode(x)`` - Calculate mode (handle ties and no mode cases)
- ``calc_q1(x)`` - Calculate first quartile (Q1)
- ``calc_q3(x)`` - Calculate third quartile (Q3)
- ``calc_iqr(x)`` - Calculate Interquartile Range ( $IQR = Q3 - Q1$ )

3. **\*\*Documentation\*\***

- Use Roxygen2 format for all functions (``#' @param``, ``#' @return``, ``#' @examples``, ``#' @export``)
- Include meaningful examples for each function
- Document expected behavior for edge cases (e.g., empty vectors, NA values)

4. **\*\*Quality Standards\*\***

- Handle edge cases gracefully (empty vectors, NA values, single values)
- Provide informative error messages for invalid inputs
- Functions should work with numeric vectors of any length

#### Deliverables

- Complete R package in ``question_1/descriptive_stats/``
- Documentation and Test cases for the functions
- README explaining the package

#### Example Usage

```

```r
# Install and load
devtools::install("question_1/descriptive_stats")
library(descriptiveStats)

# Example data
data <- c(1, 2, 2, 3, 4, 5, 5, 5, 6, 10)

# Use functions
calc_mean(data)      # 3.3
calc_median(data)    # 4.5
calc_mode(data)      # 5
calc_q1(data)        # 2.5
calc_q3(data)        # 5.5
calc_iqr(data)       # 3

```

# Pharmaverse

Learning Resources:

Coursera: [Hands On Clinical Reporting Using R](#)

Pharmaverse Examples: <https://pharmaverse.github.io/examples/>

This section has three questions (questions 2, 3 and 4)

## Question 2: SDTM DS Domain Creation using {sdm.oak}

### Objective

Create an SDTM Disposition (DS) domain dataset from raw clinical trial data using the {sdm.oak}.

{sdm.oak} learning Resources:

<https://pharmaverse.github.io/examples/> (SDTM section)

Slides and Training Videos <https://pharmaverse.github.io/rinpharma-SDTM-workshop/>

Package Documentation - <https://pharmaverse.github.io/sdm.oak/>

CDISC SDTM Implementation Guide - Refer to the DS domain in the [SDTMIG v3.4](#) in the CDISC Website. Refer to the PDF file named SDTMIG v3.4 in Files and Links tab,

Question:

Develop an R program to create the DS domain using the below

**Input raw data:** pharmaverseraw::ds\_raw

**Study controlled terminology:** The `study_ct` required to solve this exercise can be downloaded from [Github](#). If the Github link is not accessible, you can follow the instructions in [Pharmaverse Running the example](#) page and any of the [examples](#) in the SDTM section can provide the `study_ct` object. If this doesn't work, create the required file using the instructions following the below code

```
None
study_ct <-
data.frame(
  stringsAsFactors = FALSE,
  codelist_code = c("C66727", "C66727",
    "C66727", "C66727", "C66727", "C66727",
    "C66727", "C66727"),
```

```

        term_code = c("C41331", "C25250",
"C28554", "C48226", "C48227", "C48250", "C142185", "C49628",
        "C49632", "C49634"),
        term_value = c("ADVERSE EVENT",
        "COMPLETED", "DEATH", "LACK OF EFFICACY", "LOST TO
FOLLOW-UP",
        "PHYSICIAN DECISION", "PROTOCOL VIOLATION",
        "SCREEN FAILURE", "STUDY TERMINATED BY SPONSOR",
        "WITHDRAWAL BY SUBJECT"),
        collected_value = c("Adverse Event",
        "Complete", "Dead", "Lack of Efficacy", "Lost To
Follow-Up",
        "Physician Decision", "Protocol Violation",
        "Trial Screen Failure", "Study Terminated By Sponsor",
        "Withdrawal by Subject"),
        term_preferred_term = c("AE", "Completed", "Died",
        NA, NA, NA, "Violation",
        "Failure to Meet Inclusion/Exclusion
Criteria", NA, "Dropout"),
        term_synonyms = c("ADVERSE EVENT",
        "COMPLETE", "Death", NA, NA, NA, NA, NA, NA,
        "Discontinued Participation")
)

```

**SDTM Programming details:** This is the mock up eCRF that represents the data in `pharmaverse::ds_raw` and also has the programming details. Please refer to the 'General Notes' in this pdf file.

[https://github.com/pharmaverse/pharmaverse/blob/main/vignettes/articles/aCRFs/Subject\\_Disposition\\_aCRF.pdf](https://github.com/pharmaverse/pharmaverse/blob/main/vignettes/articles/aCRFs/Subject_Disposition_aCRF.pdf)

Expected Result:

Error free program with good documentation that will create the DS domain with the following variables STUDYID, DOMAIN, USUBJID, DSSEQ, DSTERM, DSDECOD, DSCAT, VISITNUM, VISIT, DSDTC, DSSTDTC, DSSTDY

HINT:

This example is very similar to the AE example in the Pharmaverse Examples.

Deliverable:

- SDTM creation script: `question_2_sdtm/02_create_ds_domain.R`



## Question 3: ADaM ADSL Dataset Creation

### Objective

Create an ADSL (Subject Level) dataset using SDTM source data, the {admiral} family of packages, and tidyverse tools.

### {admiral} learning Resources

<https://pharmaverse.github.io/examples/> (ADaM section)

Package Documentation - <https://pharmaverse.github.io/admiral/>

### Question

Develop an R program to create the ADSL using the input SDTM data, the {admiral} family of packages, and tidyverse tools as explained in the [Pharmaverse examples - ADSL](#) or in {admiral} [documentation](#). Adjust the logic and derive additional variables as mentioned below.

The DM domain is used as the basis of the ADSL. Start by assigning `pharmaversesdtm::dm` to adsl object as explained in [this](#) section of the ADSL article

Derive the variables mentioned below:

- AGEGR9 & AGEGR9N: age grouping into the following categories: "<18", "18 - 50", ">50"
- TRTSDTM: Treatment start date-time (using the first exposure record for each participant and imputing missing hours and minutes but not seconds)
- ITTFL: "Y"/"N" flag identifying patients who have been randomized, that is, where ARM is populated in pharmaversesdtm::dm domain.
- ABNSBPFL: "Y"/"N" flag identifying all patients with a supine systolic blood pressure <100 or >=140 mmHg.
- LSTALVDT: last known alive date using any vital signs visit date, any adverse event start date, any disposition record and any exposure record.
- CARPOPFL: "Y"/NA flag identifying patients with a cardiac adverse event.

The below section provides detailed specification to derive the above mentioned variables.

- **AGEGR9:** Age grouping of Analysis Age [DM.AGE]. Categories are "<18", "18 - 50", ">50".
- **AGEGR9N:** Numeric Age grouping of Analysis Age [DM.AGE]. Categories are "<18", "18 - 50", ">50". Numeric groupings are 1, 2, 3.
- **TRTSDTM/TRTSTMF:** Set to datetime of patient's first exposure observation Start Date/Time of Treatment [EX.EXSTDTC] converted to numeric datetime when sorted in date/time order. Derivation only includes observations where the patient received a valid

dose (see NOTE) and datepart of Start Date/Time of Treatment [EX.EXSTDTC] is complete. If time is missing, ie. not collected, then impute completely missing time with 00:00:00, partially missing time with 00 for missing hours, 00 for missing minutes, 00 for missing seconds. **If only seconds are missing then do not populate the imputation flag.**

**NOTE:** A valid dose is defined as (Dose per Administration [EX.EXDOSE] greater than 0 or (Dose per Administration [EX.EXDOSE] equal to 0 and Name of Actual Treatment [EX.EXTRT] contains 'PLACEBO' )).

- **ITTF:** Set to "Y" if [DM.ARM] not equal to missing Else set to "N"
- **ABNSBPFL:** Set to "Y" if patient has an observation where [VS.VSTESTCD] = "SYSBP" and [VS.VSSTRESU] is "mmHg" and [VS.VSSTRESN] is greater than or equal to 140 or less than 100. Else set to "N".
- **LSTALVDT:** Set to the last date patient has documented clinical data to show him/her alive, converted to numeric date, using the following dates:
  - (1) last complete date of vital assessment with a valid test result ([VS.VSSTRESN] and [VS.VSSTRESC] not both missing) and datepart of [VS.VSDTC] not missing.
  - (2) last complete onset date of AEs (datepart of Start Date/Time of Adverse Event [AE.AESTDTC]).
  - (3) last complete disposition date (datepart of Start Date/Time of Disposition Event [DS.DSSTDTC]).
  - (4) last date of treatment administration where patient received a valid dose (datepart of Datetime of Last Exposure to Treatment [ADSL.TRTEDTM]).Set to max of (Vitals complete, AE onset complete, disposition complete, treatment complete).
- **CARPOPF:** Set to "Y" if patient has an observation where uppercase of [AE.AESOC] = "CARDIAC DISORDERS". Else set to missing.

**Input datasets:** pharmaversesdtm::dm, pharmaversesdtm::vs, pharmaversesdtm::ex, pharmaversesdtm::ds, pharmaversesdtm::ae

### Expected Result

Error free program with good documentation that will create the ADSL dataset with all the requested variables in the question. These should be derived using {admiral} functions where possible.

### Hint

This additional variable derivation is very similar to the ADSL example in the Pharmaverse Examples: ADSL or in {admiral} documentation -

<https://pharmaverse.github.io/admiral/cran-release/articles/adsl.html>

### Deliverable

- ADSL creation script: [question\\_3\\_adam/create\\_adsl.R](#)

## Question 4: TLG - Adverse Events Reporting

### Objective

Create Tables, Listings, and Graphs (TLGs) for adverse events summary using the ADAE dataset and {gtsummary}. This tests your ability to create regulatory-compliant clinical reports.

Input : `pharmaverseadam::adae` and `pharmaverseadam::ads1`

### Learning Resources

Ggplot2 Documentation: <https://ggplot2.tidyverse.org/index.html>

FDA TLG Catalogue: <https://pharmaverse.github.io/cardinal/quarto/index-catalog.html>

Pharmaverse Examples - TLG section

### Questions

#### 1. Summary Table using {gtsummary} - HINT - [FDA Table 10](#)

Create a summary table of treatment-emergent adverse events (TEAEs).

- Treatment-emergent AE records will have TRTEMFL == "Y" in `pharmaverseadam::adae`
- Rows: AETERM or AESOC
- Columns: Treatment groups (ACTARM)
- Cell values: Count (n) and percentage (%)
- Include total row with all subjects
- Sort by descending frequency

Sample Output -

Primary System Organ Class Reported Term for the Adverse Event	Placebo N = 86 <sup>1</sup>	Xanomeline High Dose N = 72 <sup>1</sup>	Xanomeline Low Dose N = 96 <sup>1</sup>
Treatment Emergent AEs	65 (76%)	68 (94%)	84 (88%)
CARDIAC DISORDERS	12 (14%)	14 (19%)	14 (15%)
ATRIAL FIBRILLATION	1 (1.2%)	2 (2.8%)	2 (2.1%)
ATRIAL FLUTTER	0 (0%)	1 (1.4%)	1 (1.0%)
ATRIAL HYPERTROPHY	1 (1.2%)	0 (0%)	0 (0%)
ATRIOVENTRICULAR BLOCK FIRST DEGREE	1 (1.2%)	0 (0%)	1 (1.0%)
ATRIOVENTRICULAR BLOCK SECOND DEGREE	1 (1.2%)	0 (0%)	0 (0%)
BRADYCARDIA	1 (1.2%)	0 (0%)	0 (0%)

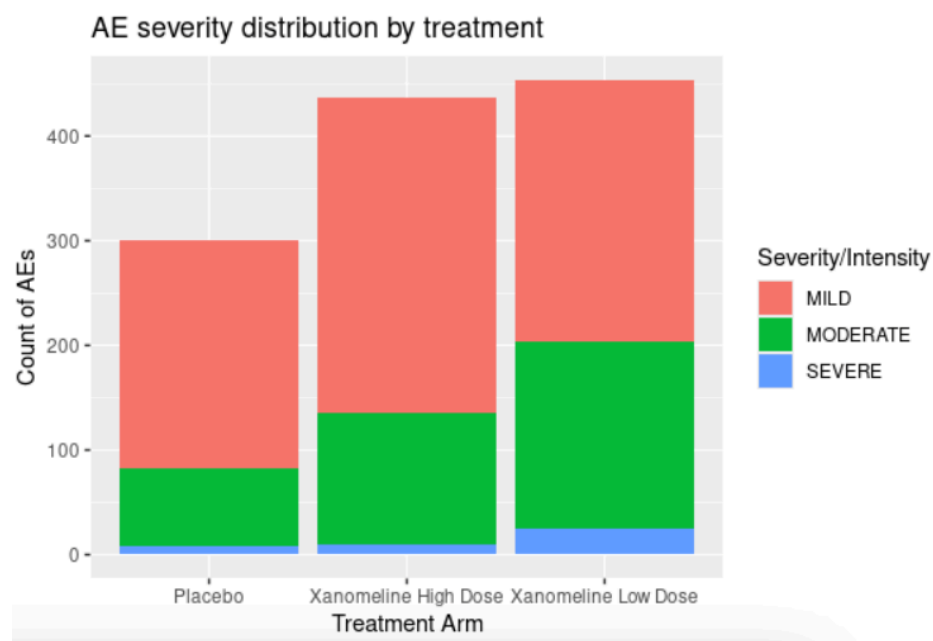
Primary System Organ Class Reported Term for the Adverse Event	Placebo N = 86 <sup>1</sup>	Xanomeline High Dose N = 72 <sup>1</sup>	Xanomeline Low Dose N = 96 <sup>1</sup>
Treatment Emergent AEs	65 (76%)	68 (94%)	84 (88%)
GENERAL DISORDERS AND ADMINISTRATION SITE CONDITIONS	21 (24%)	36 (50%)	51 (53%)
APPLICATION SITE PRURITUS	6 (7.0%)	21 (29%)	23 (24%)
APPLICATION SITE ERYTHEMA	3 (3.5%)	14 (19%)	13 (14%)
APPLICATION SITE DERMATITIS	5 (5.8%)	7 (9.7%)	9 (9.4%)
APPLICATION SITE IRRITATION	3 (3.5%)	9 (13%)	9 (9.4%)
APPLICATION SITE VESICLES	1 (1.2%)	5 (6.9%)	5 (5.2%)
FATIGUE	1 (1.2%)	5 (6.9%)	5 (5.2%)

Output format: html file

## 2. Visualizations using {ggplot2}

- **Plot 1:** AE severity distribution by treatment (bar chart or heatmap). AE Severity is captured in the AESEV variable in `pharmaverseadam::adae` dataset.

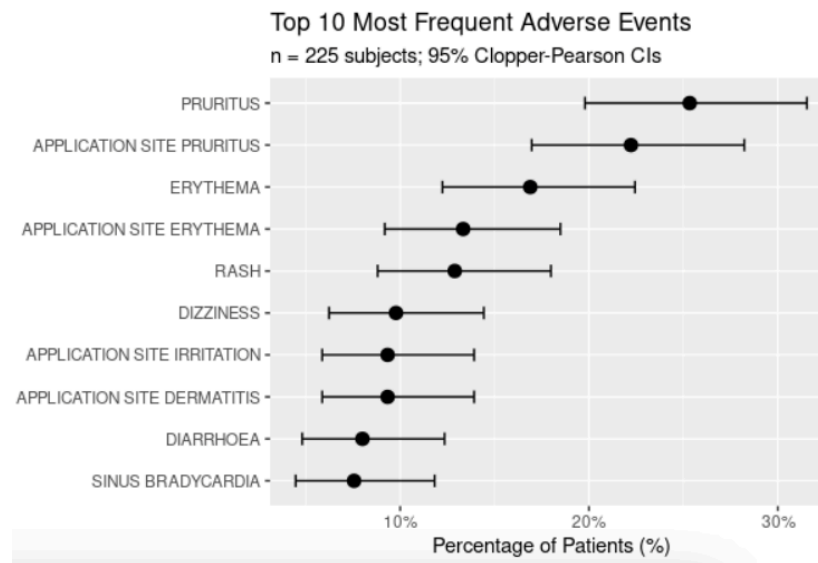
**Sample Output:**



Output format: png file

- **Plot 2:** Top 10 most frequent AEs (with 95% CI for incidence rates). AEs are captured in the **AETERM** variable in **pharmaverseadam::adae** dataset.

### Sample Output:



Output format: png file

### 3. Listing using {gtsummary}

- Using **pharmaverseadam::adae** create a detailed listing of all AEs with:
  - Subject ID, Treatment, AE Term, Severity, Relationship to drug, Start/End dates
  - Filtered for treatment-emergent events
  - Sorted by subject and event date

### Sample Output:

Listing of Treatment-Emergent Adverse Events by Subject  
Excluding Screen Failure Patients

Unique Subject Identifier	Description of Actual Arm	Reported Term for the Adverse Event	Severity/Intensity	Causality	Start Date/Time of Adverse Event	End Date/Time of Adverse
01-701-1015	Placebo	APPLICATION SITE ERYTHEMA	MILD	PROBABLE	2014-01-03	NA
		APPLICATION SITE PRURITUS	MILD	PROBABLE	2014-01-03	NA
		DIARRHOEA	MILD	REMOTE	2014-01-09	2014-01-11
01-701-1023	Placebo	ERYTHEMA	MODERATE	PROBABLE	2012-08-07	NA
			MILD	POSSIBLE	2012-08-07	2012-08-30
			MILD	POSSIBLE	2012-08-07	2012-08-30
01-701-1028	Xanomeline High Dose	ATRIOVENTRICULAR BLOCK SECOND DEGREE	MILD	POSSIBLE	2012-08-24	NA
		APPLICATION SITE ERYTHEMA	MILD	POSSIBLE	2013-07-21	NA
		APPLICATION SITE PRURITUS	MILD	PROBABLE	2013-08-08	NA
01-701-1034	Xanomeline High Dose	APPLICATION SITE PRURITUS	MILD	PROBABLE	2014-08-27	NA
01-701-1047	Placebo	PATITIS	MILD	POSSIBLE	2014-11-02	NA
		HIATUS HERNIA	MODERATE	NONE	2013-02-12	2013-02-12
			MODERATE	NONE	2013-02-12	2013-02-12
		UPPER RESPIRATORY TRACT INFECTION	MILD	NONE	2013-03-06	NA
		BUNDLE BRANCH BLOCK LEFT	MILD	NONE	2013-03-10	NA
01-701-1097	Xanomeline Low Dose	ERYTHEMA	MILD	POSSIBLE	2014-01-03	NA
		PRURITUS GENERALISED	MODERATE	POSSIBLE	2014-02-20	2014-02-22
		APPLICATION SITE VESICLES	MILD	PROBABLE	2014-02-20	NA
		APPLICATION SITE PRURITUS	MODERATE	POSSIBLE	2014-02-21	NA
			MILD	POSSIBLE	2014-02-21	NA
		PRURITUS GENERALISED	MODERATE	POSSIBLE	2014-03-21	2014-03-21
			MODERATE	POSSIBLE	2014-03-31	2014-03-31
			MODERATE	POSSIBLE	2014-04-19	2014-04-20
		PHARYNGOLARYNGEAL PAIN	MILD	NONE	2014-04-19	2014-04-22

Output Format: HTML file

## Deliverables

- Script to create summary table:  
`question_4_tlg/01_create_ae_summary_table.R`
- Script to create visualizations: `question_4_tlg/02_create_visualizations.R`
- Script to create Listings: `question_4_tlg/03_create_listings.R`
- Output files:
  - `ae_summary_table.html` (or .docx/.pdf)
  - Two PNG files
  - `ae_listings.html`

## Python Coding assessment

### Question 5: Clinical Data API (FastAPI)

#### Objective

Create a RESTful API using FastAPI that serves clinical trial data, performs dynamic cohort analysis, and calculates patient risk scores. This tests your ability to build interactive backends for data-driven applications.

#### Learning Resources

- **FastAPI Request Body:** <https://fastapi.tiangolo.com/tutorial/body/>
- **Pandas Filtering:** [https://pandas.pydata.org/docs/user\\_guide/indexing.html](https://pandas.pydata.org/docs/user_guide/indexing.html)

#### Input

- **File:** `adae.csv`, exported from your previous R work `pharmaversesdtm::ae`).

#### Requirements

Develop a Python script using FastAPI that exposes the following endpoints:

1. **GET /**
  - Returns a JSON welcome message: `{"message": "Clinical Trial Data API is running"}`.
2. **POST /ae-query (Dynamic Filtering)**
  - **Goal:** Allow users to filter data dynamically based on a JSON payload.
  - **Input Body:** A JSON object with *optional* filters.  
JSON

None

```
{
  "severity": ["MILD", "MODERATE"],
  "treatment_arm": "Placebo"
}
```

- **Logic:**
  - Accept a list of severities (e.g., **AESEV**) and/or a specific treatment arm (**ACTARM**).
  - Filter the dataset to match *all* provided criteria.
  - If a field is missing or null in the request, ignore that filter (return all records for that dimension).
- **Output:** JSON returning the count of matching records and a list of unique **USUBJIDs** in that cohort.
- 3. **GET /subject-risk/{subject\_id} (Calculation Logic)**
  - **Goal:** Calculate a "Safety Risk Score" for a specific patient.
  - **Logic:**
    - Filter AEs for the given **subject\_id**.
    - Calculate a **weighted score** based on **AESEV**:
      - **MILD**: 1 point
      - **MODERATE**: 3 points
      - **SEVERE**: 5 points
    - Sum the points to get a total **risk\_score**.
    - Assign a **risk\_category**:
      - **Low**: Score < 5
      - **Medium**: 5 <= Score < 15
      - **High**: Score >= 15
  - **Output:** JSON

None

```
{
  "subject_id": "01-701-1015",
  "risk_score": 8,
  "risk_category": "Medium"
}
```

- **Error Handling:** Return a 404 HTTP exception if the `subject_id` does not exist.

## Deliverables

- **Code:** Python scripts
- **Documentation:** A brief instruction in your README on how to run the API locally (e.g., `uvicorn main:app --reload`).

## Question 6: GenAI Clinical Data Assistant (LLM & LangChain)

### Objective

Develop a Generative AI Assistant that translates natural language questions into structured Pandas queries. The goal is to test your ability to use LLMs (e.g., OpenAI via LangChain) to dynamically map user intent to the correct dataset variable without hard-coding rules.

### Scenario

A clinical safety reviewer wants to ask free-text questions about the AE dataset. They don't know the column names. Your Agent must "understand" the dataset schema and route the question to the correct variable. For example,

- If they ask about "severity" or "intensity" → Map to **AESEV**.
- If they ask about a specific condition (e.g., "Headache") → Map to **AETERM**.
- If they ask about a body system (e.g., "Cardiac", "Skin") → Map to **AESOC**.

### Input

- **File:** `adae.csv` (`pharmaversesdtm:ae`).
- **API Key:** You may use your own OpenAI API key or any other solution. If you do not have one, you may mock the LLM response in your code, but the logic flow (Prompt -> Parse -> Execute) must be complete.

### Requirements

1. **Schema Definition:** Understand the data and define a dictionary or string in your code describing the relevant columns (**AESEV**, **AETERM**, **AESOC**, etc) to the LLM.
2. **LLM Implementation:**
  - Create a function or class `ClinicalTrialDataAgent`.
  - Use an LLM to parse a user's question into a **Structured JSON Output** containing:
    - **target\_column:** The column to filter.
    - **filter\_value:** The value to search for (extracted from the question).
3. **Execution:**



- Write a function that takes the LLM's output and applies the actual Pandas filter to the `ae` dataframe.
- Return the count of unique subjects (`USUBJID`) and a list of the matching IDs.

## Deliverables

- **Code** of the solution developed
- **Test Script:** A simple block of code that runs 3 example queries (of your choice) and prints the results. An example question is, `Give me the subjects who had Adverse events of Moderate severity`

## FAQ & Troubleshooting

Q: Do I have to use the recommended packages like `{sdtm.oak}` for Question 2,3 & 4?

A: Yes, it's recommended as it aligns with Pharmaverse best practices. Demonstrating familiarity with `{sdtm.oak}`, `{admiral}` and `{gtsummary}` is a plus.

Q: What if I'm not familiar with CDISC standards?

A: Start by reviewing the CDISC resources listed above. The standards are well-documented, and the Pharmaverse examples provide excellent practical illustrations.

Q: Should I test my package installation before submission?

A: Definitely! Before submitting, ensure:

- `devtools::check()` passes (for packages)
- Your code runs from scratch without errors
- Documentation builds correctly with `roxygen2::roxygenise()`
- All deliverables are in the correct locations