

Python for Clinical Study Reports and Submission

Yilong Zhang

Nan Xiao

Table of contents

Welcome	4
1 Introduction	5
2 Polars	6
2.1 Polars	6
2.2 I/O	6
2.3 Filtering	6
2.4 Deriving	7
2.5 Grouping	8
2.6 Joining	9
2.7 Pivoting	9
I Tables, Listings, and Figures	11
3 Disposition of Participants Table	12
3.1 Overview	12
3.2 Step 1: Load Data	12
3.3 Step 2: Count Total Participants	13
3.4 Step 3: Count Completed Participants	14
3.5 Step 4: Count Discontinued Participants	15
3.6 Step 5: Break Down Discontinuation Reasons	16
3.7 Step 6: Combine All Results	17
3.8 Step 7: Generate Publication-Ready Output	18
4 Study Population Table	19
4.1 Overview	19
4.2 Step 1: Load Data	19
4.3 Step 2: Calculate Treatment Group Totals	20
4.4 Step 3: Define Helper Function	20
4.5 Step 4: Count Each Population	21
4.5.1 All Randomized Participants	21
4.5.2 Intent-to-Treat Population	21
4.5.3 Efficacy Population	22
4.5.4 Safety Population	22

4.6	Step 5: Combine All Populations	23
4.7	Step 6: Calculate Percentages	23
4.8	Step 7: Format Display Values	24
4.9	Step 8: Create Final Table	25
4.10	Step 9: Generate Publication-Ready Output	26
5	Baseline Characteristics Table	27
5.1	Overview	27
5.2	Step 1: Load Data	27
5.3	Step 2: Calculate Summary Statistics	28
5.3.1	Continuous Variables (Age)	28
5.3.2	Categorical Variables (Sex, Race)	29
5.4	Step 3: Format Results	30
5.4.1	Format Age Statistics	30
5.4.2	Format Categorical Statistics	30
5.5	Step 4: Create Table Structure	31
5.6	Step 5: Generate Publication-Ready Output	33
6	Adverse Events Summary Table	35
6.1	Overview	35
6.2	Step 1: Load Data	35
6.3	Step 2: Filter Safety Population	36
6.4	Step 3: Define AE Categories	37
6.5	Step 4: Combine and Calculate Percentages	39
6.6	Step 5: Format for Display	40
6.7	Step 6: Create Final Table Structure	41
6.8	Step 7: Generate Publication-Ready Output	42
7	Specific Adverse Events Table	45
7.1	Setup	45
7.2	Prepare AE Summary Data	45
7.3	Create RTF Output	46
8	ANCOVA Efficacy Analysis	48
8.1	Setup	48
8.2	Prepare Analysis Data	48
8.3	Create Tables for RTF Output	50
8.4	Create RTF Document	51
9	Summary	53
	References	54

Welcome

Welcome to Python for Clinical Study Reports and Submission. Clinical study reports (CSR) are crucial components in clinical trial development. A CSR is an “integrated” full scientific report of an individual clinical trials.

The [ICH E3: Structure and Content of Clinical Study Reports](#) offers comprehensive instructions to sponsors on the creation of a CSR. This book is a clear and straightforward guide on using Python to streamline the process of preparing CSRs. Additionally, it provides detailed guidance on the submission process to regulatory agencies. Whether you are a beginner or an experienced developer, this book is an indispensable asset in your clinical reporting toolkit.

This is a work-in-progress draft.

1 Introduction

This is a book created from Markdown and executable code.

See Knuth (1984) for additional discussion of literate programming.

```
import polars as pl
from datetime import datetime

df = pl.DataFrame(
    {
        "integer": [1, 2, 3],
        "date": [
            datetime(2025, 1, 1),
            datetime(2025, 1, 2),
            datetime(2025, 1, 3),
        ],
        "float": [4.0, 5.0, 6.0],
        "string": ["a", "b", "c"],
    }
)

df
```

integer	date	float	string
i64	datetime[s]	f64	str
1	2025-01-01 00:00:00	4.0	"a"
2	2025-01-02 00:00:00	5.0	"b"
3	2025-01-03 00:00:00	6.0	"c"

2 Polars

2.1 Polars

We use [Polars](#) in this book. The syntax aligns with tidyverse-style pipelines. The sections below cover basic examples. For deeper dives, see the [Polars user guide](#) or the book [Python Polars: The Definitive Guide](#).

2.2 I/O

Polars supports multiple data formats for input and output (see the [I/O guide](#)). For clinical development, we recommend the `.parquet` format because tools in Python, R, and Julia can read and write it without conversion. The example below loads subject-level ADSL data with Polars.

```
import polars as pl
adsl = pl.read_parquet("data/adsl.parquet")
adsl = adsl.select("STUDYID", "USUBJID", "TRT01A", "AGE", "SEX") # select columns
adsl.head()
```

STUDYID	USUBJID	TRT01A	AGE	SEX
str	str	str	f64	str
"CDISCPILLOT01"	"01-701-1015"	"Placebo"	63.0	"Female"
"CDISCPILLOT01"	"01-701-1023"	"Placebo"	64.0	"Male"
"CDISCPILLOT01"	"01-701-1028"	"Xanomeline High Dose"	71.0	"Male"
"CDISCPILLOT01"	"01-701-1033"	"Xanomeline Low Dose"	74.0	"Male"
"CDISCPILLOT01"	"01-701-1034"	"Xanomeline High Dose"	77.0	"Female"

2.3 Filtering

Filtering in Polars uses the `.filter()` method with column expressions. Below are examples applied to the ADSL data.

```
# Filter female subjects
adsl.filter(pl.col("SEX") == "Female").head()
```

STUDYID	USUBJID	TRT01A	AGE	SEX
str	str	str	f64	str
"CDISCPILLOT01"	"01-701-1015"	"Placebo"	63.0	"Female"
"CDISCPILLOT01"	"01-701-1034"	"Xanomeline High Dose"	77.0	"Female"
"CDISCPILLOT01"	"01-701-1047"	"Placebo"	85.0	"Female"
"CDISCPILLOT01"	"01-701-1111"	"Xanomeline Low Dose"	81.0	"Female"
"CDISCPILLOT01"	"01-701-1133"	"Xanomeline High Dose"	81.0	"Female"

```
# Filter subjects with Age >= 65
adsl.filter(pl.col("AGE") >= 65).head()
```

STUDYID	USUBJID	TRT01A	AGE	SEX
str	str	str	f64	str
"CDISCPILLOT01"	"01-701-1028"	"Xanomeline High Dose"	71.0	"Male"
"CDISCPILLOT01"	"01-701-1033"	"Xanomeline Low Dose"	74.0	"Male"
"CDISCPILLOT01"	"01-701-1034"	"Xanomeline High Dose"	77.0	"Female"
"CDISCPILLOT01"	"01-701-1047"	"Placebo"	85.0	"Female"
"CDISCPILLOT01"	"01-701-1097"	"Xanomeline Low Dose"	68.0	"Male"

2.4 Deriving

Deriving new variables is common in clinical data analysis for creating age groups, BMI categories, or treatment flags. Polars uses `.with_columns()` to add new columns while keeping existing ones.

```
# Create age groups
adsl.with_columns([
    pl.when(pl.col("AGE") < 65)
        .then(pl.lit("<65"))
        .otherwise(pl.lit(">=65"))
        .alias("AGECAT")
]).head()
```

STUDYID	USUBJID	TRT01A	AGE	SEX	AGECAT
str	str	str	f64	str	str
"CDISCPILLOT01"	"01-701-1015"	"Placebo"	63.0	"Female"	"<65"
"CDISCPILLOT01"	"01-701-1023"	"Placebo"	64.0	"Male"	"<65"
"CDISCPILLOT01"	"01-701-1028"	"Xanomeline High Dose"	71.0	"Male"	">=65"
"CDISCPILLOT01"	"01-701-1033"	"Xanomeline Low Dose"	74.0	"Male"	">=65"
"CDISCPILLOT01"	"01-701-1034"	"Xanomeline High Dose"	77.0	"Female"	">=65"

2.5 Grouping

Grouping operations are fundamental for creating summary statistics in clinical reports. Polars uses `group_by()` followed by aggregation functions to compute counts, means, and other statistics by categorical variables like treatment groups.

The `.count()` method provides a quick way to get subject counts by group.

```
# Count by treatment group
adsl.group_by("TRT01A").count().sort("TRT01A")
```

```
/tmp/ipykernel_7375/3512064118.py:2: DeprecationWarning: `GroupBy.count` was renamed; use `G
adsl.group_by("TRT01A").count().sort("TRT01A")
```

TRT01A	count
str	u32
"Placebo"	86
"Xanomeline High Dose"	84
"Xanomeline Low Dose"	84

You can also use `.agg()` with multiple aggregation functions:

```
# Age statistics by treatment group
adsl.group_by("TRT01A").agg([
    pl.col("AGE").mean().round(1).alias("mean_age"),
    pl.col("AGE").std().round(2).alias("sd_age")
]).sort("TRT01A")
```


TRT01A	mean_age	sd_age
str	f64	f64
"Placebo"	75.2	8.59
"Xanomeline High Dose"	74.4	7.89
"Xanomeline Low Dose"	75.7	8.29

2.6 Joining

Joining datasets is essential for combining subject-level data (ADSL) with event-level data (e.g. ADAE, ADLB). Polars supports various join types including inner, left, and full joins.

Here is a toy example that splits ADSL and joins it back by `USUBJID`.

```
# Create a simple demographics subset
demo = adsl.select("USUBJID", "AGE", "SEX").head(3)

# Create treatment info subset
trt = adsl.select("USUBJID", "TRT01A").head(3)

# Left join to combine datasets
demo.join(trt, on="USUBJID", how="left")
```

USUBJID	AGE	SEX	TRT01A
str	f64	str	str
"01-701-1015"	63.0	"Female"	"Placebo"
"01-701-1023"	64.0	"Male"	"Placebo"
"01-701-1028"	71.0	"Male"	"Xanomeline High Dose"

2.7 Pivoting

Pivoting transforms data from long to wide format, commonly needed for creating tables. Use `.pivot()` to reshape grouped data into columns.

```
# Create summary by treatment and sex
(
    adsl
    .group_by(["TRT01A", "SEX"])
    .agg(pl.len().alias("n"))
)
```

```

    .pivot(
      values="n",
      index="SEX",
      on="TRT01A"
    )
)

```

SEX str	Xanomeline Low Dose u32	Xanomeline High Dose u32	Placebo u32
"Male"	34	44	33
"Female"	50	40	53

Part I

Tables, Listings, and Figures

3 Disposition of Participants Table

3.1 Overview

Clinical trials need to track how participants flow through a study from enrollment to completion. Following [ICH E3 guidance](#), regulatory submissions require a disposition table in Section 10.1 that summarizes:

- **Enrolled:** Total participants who entered the study
- **Completed:** Participants who finished the study protocol
- **Discontinued:** Participants who left early and their reasons

This tutorial shows you how to create a regulatory-compliant disposition table using Python's `rtflite` package.

```
import polars as pl # Manipulate data
import rtflite as rtf # Reporting in RTF format
```

3.2 Step 1: Load Data

We start by loading the Subject-level Analysis Dataset (ADSL), which contains all participant information needed for our disposition table.

The ADSL dataset stores participant-level information including treatment assignments and study completion status. We're using the `parquet` format for data storage.

```
adsl = pl.read_parquet("data/adsl.parquet")
```

Let's examine the key variables we'll use to build our disposition table:

- **USUBJID:** Unique identifier for each participant
- **TRT01P:** Treatment name (text)
- **TRT01PN:** Treatment group (numeric code)
- **DISCONFL:** Flag indicating if participant discontinued (Y/N)
- **DCREASCD:** Specific reason for discontinuation

```
adsl.select(["USUBJID", "TRT01P", "TRT01PN", "DISCONFL", "DCREASCD"])
```

USUBJID	TRT01P	TRT01PN	DISCONFL	DCREASCD
str	str	i64	str	str
"01-701-1015"	"Placebo"	0	" "	"Completed"
"01-701-1023"	"Placebo"	0	"Y"	"Adverse Event"
"01-701-1028"	"Xanomeline High Dose"	81	" "	"Completed"
"01-701-1033"	"Xanomeline Low Dose"	54	"Y"	"Sponsor Decision"
"01-701-1034"	"Xanomeline High Dose"	81	" "	"Completed"
...
"01-718-1254"	"Xanomeline Low Dose"	54	" "	"Completed"
"01-718-1328"	"Xanomeline High Dose"	81	"Y"	"Withdrew Consent"
"01-718-1355"	"Placebo"	0	" "	"Completed"
"01-718-1371"	"Xanomeline High Dose"	81	"Y"	"Adverse Event"
"01-718-1427"	"Xanomeline High Dose"	81	"Y"	"Lack of Efficacy"

3.3 Step 2: Count Total Participants

First, we count how many participants were enrolled in each treatment group.

We group participants by treatment arm and count them using `.group_by()` and `.agg()`. The `.pivot()` operation reshapes our data from long format (rows for each treatment) to wide format (columns for each treatment), which matches the standard disposition table layout.

```
n_rand = (
    adsl
    .group_by("TRT01PN")
    .agg(n = pl.len())
    .with_columns([
        pl.lit("Participants in population").alias("row"),
        pl.lit(None, dtype=pl.Float64).alias("pct") # Placeholder for percentage (not applicable)
    ])
    .pivot(
        index="row",
        on="TRT01PN",
        values=["n", "pct"],
        sort_columns=True
    )
)
```

n_rand

row	n_0	n_54	n_81	pct_0	pct_54	pct_81
str	u32	u32	u32	f64	f64	f64
"Participants in population"	86	84	84	null	null	null

3.4 Step 3: Count Completed Participants

Next, we identify participants who successfully completed the study and calculate what percentage they represent of each treatment group.

We filter for participants where DCREASCD == "Completed", then calculate both counts and percentages. The .join() operation brings in the total count for each treatment group so we can compute percentages.

```
n_complete = (
    adsl
    .filter(pl.col("DCREASCD") == "Completed")
    .group_by("TRT01PN")
    .agg(n = pl.len())
    .join(
        adsl.group_by("TRT01PN").agg(total = pl.len()),
        on="TRT01PN"
    )
    .with_columns([
        pl.lit("Completed").alias("row"),
        (100.0 * pl.col("n") / pl.col("total")).round(1).alias("pct")
    ])
    .pivot(
        index="row",
        on="TRT01PN",
        values=["n", "pct"],
        sort_columns=True
    )
)

n_complete
```

row	n_0	n_54	n_81	pct_0	pct_54	pct_81
str	u32	u32	u32	f64	f64	f64
"Completed"	58	25	27	67.4	29.8	32.1

3.5 Step 4: Count Discontinued Participants

Now we count participants who left the study early, regardless of their specific reason.

We filter for participants where the discontinuation flag `DISCONFL == "Y"`, then follow the same pattern of counting and calculating percentages within each treatment group.

```
n_disc = (
    adsl
    .filter(pl.col("DISCONFL") == "Y")
    .group_by("TRT01PN")
    .agg(n = pl.len())
    .join(
        adsl.group_by("TRT01PN").agg(total = pl.len()),
        on="TRT01PN"
    )
    .with_columns([
        pl.lit("Discontinued").alias("row"),
        (100.0 * pl.col("n") / pl.col("total")).round(1).alias("pct")
    ])
    .pivot(
        index="row",
        on="TRT01PN",
        values=["n", "pct"],
        sort_columns=True
    )
)

n_disc
```

row	n_0	n_54	n_81	pct_0	pct_54	pct_81
str	u32	u32	u32	f64	f64	f64
"Discontinued"	28	59	57	32.6	70.2	67.9

3.6 Step 5: Break Down Discontinuation Reasons

For regulatory reporting, we need to show the specific reasons why participants discontinued.

We filter out completed participants, then group by both treatment and discontinuation reason. The indentation (four spaces) in the row labels helps show these are subcategories under “Discontinued”. We also use `.fill_null(0)` to handle cases where certain discontinuation reasons don’t occur in all treatment groups.

```
n_reason = (
    adsl
    .filter(pl.col("DCREASCD") != "Completed")
    .group_by(["TRT01PN", "DCREASCD"])
    .agg(n = pl.len())
    .join(
        adsl.group_by("TRT01PN").agg(total = pl.len()),
        on="TRT01PN"
    )
    .with_columns([
        pl.concat_str([pl.lit("    "), pl.col("DCREASCD")]).alias("row"),
        (100.0 * pl.col("n") / pl.col("total")).round(1).alias("pct")
    ])
    .pivot(
        index="row",
        on="TRT01PN",
        values=["n", "pct"],
        sort_columns=True
    )
    .with_columns([
        pl.col(["n_0", "n_54", "n_81"]).fill_null(0),
        pl.col(["pct_0", "pct_54", "pct_81"]).fill_null(0.0)
    ])
    .sort("row")
)

n_reason
```

row	n_0	n_54	n_81	pct_0	pct_54	pct_81
str	u32	u32	u32	f64	f64	f64
" Adverse Event"	8	44	40	9.3	52.4	47.6
" Death"	2	1	0	2.3	1.2	0.0

row	n_0	n_54	n_81	pct_0	pct_54	pct_81
str	u32	u32	u32	f64	f64	f64
" I/E Not Met"	1	0	2	1.2	0.0	2.4
" Lack of Efficacy"	3	0	1	3.5	0.0	1.2
" Lost to Follow-up"	1	1	0	1.2	1.2	0.0
" Physician Decision"	1	0	2	1.2	0.0	2.4
" Protocol Violation"	1	1	1	1.2	1.2	1.2
" Sponsor Decision"	2	2	3	2.3	2.4	3.6
" Withdrew Consent"	9	10	8	10.5	11.9	9.5

3.7 Step 6: Combine All Results

Now we stack all our individual summaries together to create the complete disposition table.

Using `pl.concat()`, we combine the enrollment counts, completion counts, discontinuation counts, and detailed discontinuation reasons into a single table that flows logically from top to bottom.

```
tbl_disp = pl.concat([
    n_rand,
    n_complete,
    n_disc,
    n_reason
])

tbl_disp
```

row	n_0	n_54	n_81	pct_0	pct_54	pct_81
str	u32	u32	u32	f64	f64	f64
"Participants in population"	86	84	84	null	null	null
"Completed"	58	25	27	67.4	29.8	32.1
"Discontinued"	28	59	57	32.6	70.2	67.9
" Adverse Event"	8	44	40	9.3	52.4	47.6
" Death"	2	1	0	2.3	1.2	0.0
...
" Lost to Follow-up"	1	1	0	1.2	1.2	0.0
" Physician Decision"	1	0	2	1.2	0.0	2.4
" Protocol Violation"	1	1	1	1.2	1.2	1.2
" Sponsor Decision"	2	2	3	2.3	2.4	3.6

row	n_0	n_54	n_81	pct_0	pct_54	pct_81
str	u32	u32	u32	f64	f64	f64
" Withdrew Consent"	9	10	8	10.5	11.9	9.5

3.8 Step 7: Generate Publication-Ready Output

Finally, we format our table in RTF format using the `rtflite` package.

The `RTFDocument` class handles the complex formatting required for clinical reports, including proper column headers, borders, and spacing. The resulting RTF file can be directly included in regulatory submissions or converted to PDF for review.

```
doc_disp = rtf.RTFDocument(
    df=tbl_disp.select("row", "n_0", "pct_0", "n_54", "pct_54", "n_81", "pct_81"),
    rtf_title=rtf.RTFTitle(text=["Disposition of Participants"]),
    rtf_column_header=[
        rtf.RTFColumnHeader(
            text=["", "Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"],
            col_rel_width=[3] + [2] * 3,
            text_justification=["l"] + ["c"] * 3,
        ),
        rtf.RTFColumnHeader(
            text=["", "n", "(%)", "n", "(%)", "n", "(%)"],
            col_rel_width=[3] + [1] * 6,
            text_justification=["l"] + ["c"] * 6,
            border_top=[""] + ["single"] * 6,
            border_left=["single"] + ["single", ""] * 3
        )
    ],
    rtf_body=rtf.RTFBody(
        col_rel_width=[3] + [1] * 6,
        text_justification=["l"] + ["c"] * 6,
        border_left=["single"] + ["single", ""] * 3
    ),
    rtf_source=rtf.RTFSource(text=["Source: ADSL dataset"]) # Required source attribution
)

doc_disp.write_rtf("rtf/tlf_disposition.rtf") # Save as RTF for submission
```

rtf/tlf_disposition.rtf

4 Study Population Table

4.1 Overview

Clinical trials define multiple analysis populations based on different inclusion criteria. Following [ICH E3 guidance](#), regulatory submissions must clearly document the number of participants in each analysis population to support the validity of statistical analyses.

The key analysis populations typically include:

- **All Randomized:** Total participants who entered the study
- **Intent-to-Treat (ITT):** Participants included in the primary efficacy analysis
- **Efficacy Population:** Participants who meet specific criteria for efficacy evaluation
- **Safety Population:** Participants who received at least one dose of study treatment

This tutorial shows you how to create a population summary table using Python's `rtflite` package.

```
import polars as pl # Data manipulation
import rtflite as rtf # RTF reporting
```

4.2 Step 1: Load Data

We start by loading the Subject-level Analysis Dataset (ADSL), which contains population flags for each participant.

```
adsl = pl.read_parquet("data/adsl.parquet")
```

Let's examine the key population flag variables we'll use:

- **USUBJID:** Unique participant identifier
- **TRT01P:** Planned treatment group
- **ITTFL:** Intent-to-treat population flag (Y/N)
- **EFFFL:** Efficacy population flag (Y/N)
- **SAFFL:** Safety population flag (Y/N)

```
adsl.select(["USUBJID", "TRT01P", "ITTFL", "EFFFL", "SAFFL"])
```

USUBJID	TRT01P	ITTFL	EFFFL	SAFFL
str	str	str	str	str
"01-701-1015"	"Placebo"	"Y"	"Y"	"Y"
"01-701-1023"	"Placebo"	"Y"	"Y"	"Y"
"01-701-1028"	"Xanomeline High Dose"	"Y"	"Y"	"Y"
"01-701-1033"	"Xanomeline Low Dose"	"Y"	"Y"	"Y"
"01-701-1034"	"Xanomeline High Dose"	"Y"	"Y"	"Y"
...
"01-718-1254"	"Xanomeline Low Dose"	"Y"	"Y"	"Y"
"01-718-1328"	"Xanomeline High Dose"	"Y"	"Y"	"Y"
"01-718-1355"	"Placebo"	"Y"	"Y"	"Y"
"01-718-1371"	"Xanomeline High Dose"	"Y"	"Y"	"Y"
"01-718-1427"	"Xanomeline High Dose"	"Y"	"Y"	"Y"

4.3 Step 2: Calculate Treatment Group Totals

First, we calculate the total number of randomized participants in each treatment group, which will serve as the denominator for percentage calculations.

```
totals = adsl.group_by("TRT01P").agg(
    total = pl.len()
)

totals
```

TRT01P	total
str	u32
"Xanomeline High Dose"	84
"Placebo"	86
"Xanomeline Low Dose"	84

4.4 Step 3: Define Helper Function

We create a reusable function to count participants by treatment group for any population subset.

```
def count_by_treatment(data, population_name):
    """Count participants by treatment group and add population label"""
    return data.group_by("TRT01P").agg(
        n = pl.len()
    ).with_columns(
        population = pl.lit(population_name)
    )
```

4.5 Step 4: Count Each Population

Now we calculate participant counts for each analysis population.

4.5.1 All Randomized Participants

```
pop_all = count_by_treatment(
    data=adsl,
    population_name="Participants in population"
)

pop_all
```

TRT01P	n	population
str	u32	str
"Placebo"	86	"Participants in population"
"Xanomeline High Dose"	84	"Participants in population"
"Xanomeline Low Dose"	84	"Participants in population"

4.5.2 Intent-to-Treat Population

```
adsl_itt = adsl.filter(pl.col("ITTFL") == "Y")
pop_itt = count_by_treatment(
    data=adsl_itt,
    population_name="Participants included in ITT population"
)

pop_itt
```

TRT01P	n	population
str	u32	str
"Xanomeline High Dose"	84	"Participants included in ITT p...
"Placebo"	86	"Participants included in ITT p...
"Xanomeline Low Dose"	84	"Participants included in ITT p...

4.5.3 Efficacy Population

```
adsl_eff = adsl.filter(pl.col("EFFFL") == "Y")
pop_eff = count_by_treatment(
    data=adsl_eff,
    population_name="Participants included in efficacy population"
)

pop_eff
```

TRT01P	n	population
str	u32	str
"Xanomeline Low Dose"	81	"Participants included in effic...
"Xanomeline High Dose"	74	"Participants included in effic...
"Placebo"	79	"Participants included in effic...

4.5.4 Safety Population

```
adsl_saf = adsl.filter(pl.col("SAFFL") == "Y")
pop_saf = count_by_treatment(
    data=adsl_saf,
    population_name="Participants included in safety population"
)

pop_saf
```

TRT01P	n	population
str	u32	str
"Xanomeline Low Dose"	84	"Participants included in safet...

TRT01P	n	population
str	u32	str
"Xanomeline High Dose"	84	"Participants included in safet...
"Placebo"	86	"Participants included in safet...

4.6 Step 5: Combine All Populations

We stack all population counts together into a single dataset.

```
all_populations = pl.concat([
    pop_all,
    pop_itt,
    pop_eff,
    pop_saf
])

all_populations
```

TRT01P	n	population
str	u32	str
"Placebo"	86	"Participants in population"
"Xanomeline High Dose"	84	"Participants in population"
"Xanomeline Low Dose"	84	"Participants in population"
"Xanomeline High Dose"	84	"Participants included in ITT p...
"Placebo"	86	"Participants included in ITT p...
...
"Xanomeline High Dose"	74	"Participants included in effic...
"Placebo"	79	"Participants included in effic...
"Xanomeline Low Dose"	84	"Participants included in safet...
"Xanomeline High Dose"	84	"Participants included in safet...
"Placebo"	86	"Participants included in safet...

4.7 Step 6: Calculate Percentages

We join with the total counts and calculate what percentage each population represents of the total randomized participants.

```
stats_with_pct = all_populations.join(
  totals,
  on="TRT01P"
).with_columns(
  pct = (100.0 * pl.col("n") / pl.col("total")).round(1)
)

stats_with_pct
```

TRT01P	n	population	total	pct
str	u32	str	u32	f64
"Placebo"	86	"Participants in population"	86	100.0
"Xanomeline High Dose"	84	"Participants in population"	84	100.0
"Xanomeline Low Dose"	84	"Participants in population"	84	100.0
"Xanomeline High Dose"	84	"Participants included in ITT p..."	84	100.0
"Placebo"	86	"Participants included in ITT p..."	86	100.0
...
"Xanomeline High Dose"	74	"Participants included in effic..."	84	88.1
"Placebo"	79	"Participants included in effic..."	86	91.9
"Xanomeline Low Dose"	84	"Participants included in safet..."	84	100.0
"Xanomeline High Dose"	84	"Participants included in safet..."	84	100.0
"Placebo"	86	"Participants included in safet..."	86	100.0

4.8 Step 7: Format Display Values

For the final table, we format the display text. The total randomized count shows just “N”, while subset populations show “N (%)”.

```
formatted_stats = stats_with_pct.with_columns(
  display = pl.when(pl.col("population") == "Participants in population")
    .then(pl.col("n").cast(str))
    .otherwise(
      pl.concat_str([
        pl.col("n").cast(str),
        pl.lit(" (%)"),
        pl.col("pct").round(1).cast(str),
        pl.lit("%")
      ])
    )
)
```



```
)
```

```
formatted_stats
```

TRT01P str	n u32	population str	total u32	pct f64	display str
"Placebo"	86	"Participants in population"	86	100.0	"86"
"Xanomeline High Dose"	84	"Participants in population"	84	100.0	"84"
"Xanomeline Low Dose"	84	"Participants in population"	84	100.0	"84"
"Xanomeline High Dose"	84	"Participants included in ITT p...	84	100.0	"84 (100.0)"
"Placebo"	86	"Participants included in ITT p...	86	100.0	"86 (100.0)"
...
"Xanomeline High Dose"	74	"Participants included in effic...	84	88.1	"74 (88.1)"
"Placebo"	79	"Participants included in effic...	86	91.9	"79 (91.9)"
"Xanomeline Low Dose"	84	"Participants included in safet...	84	100.0	"84 (100.0)"
"Xanomeline High Dose"	84	"Participants included in safet...	84	100.0	"84 (100.0)"
"Placebo"	86	"Participants included in safet...	86	100.0	"86 (100.0)"

4.9 Step 8: Create Final Table

We reshape the data from long format (rows for each treatment-population combination) to wide format (columns for each treatment group).

```
df_overview = formatted_stats.pivot(  
    values="display",  
    index="population",  
    on="TRT01P",  
    maintain_order=True  
) .select(  
    ["population", "Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"]  
)
```

```
df_overview
```

population str	Placebo str	Xanomeline Low Dose str	Xanomeline High Dose str
"Participants in population"	"86"	"84"	"84"
"Participants included in ITT p...	"86 (100.0)"	"84 (100.0)"	"84 (100.0)"

population str	Placebo str	Xanomeline Low Dose str	Xanomeline High Dose str
"Participants included in effic...	"79 (91.9)"	"81 (96.4)"	"74 (88.1)"
"Participants included in safet...	"86 (100.0)"	"84 (100.0)"	"84 (100.0)"

4.10 Step 9: Generate Publication-Ready Output

Finally, we format the population table for regulatory submission using the `rtflite` package.

```
doc_overview = rtf.RTFDocument(
    df=df_overview,
    rtf_title=rtf.RTFTitle(
        text=["Analysis Population", "All Participants Randomized"]
    ),
    rtf_column_header=rtf.RTFColumnHeader(
        text=["", "Placebo\n\n (%)", "Xanomeline Low Dose\n\n (%)", "Xanomeline High Dose\n\n (%)"],
        col_rel_width=[4, 2, 2, 2],
        text_justification=["l", "c", "c", "c"],
    ),
    rtf_body=rtf.RTFBody(
        col_rel_width=[4, 2, 2, 2],
        text_justification=["l", "c", "c", "c"],
    ),
    rtf_source=rtf.RTFSource(text=["Source: ADSL dataset"])
)

doc_overview.write_rtf("rtf/tlf_population.rtf")
```

rtf/tlf_population.rtf

5 Baseline Characteristics Table

5.1 Overview

Baseline characteristics tables summarize the demographic and clinical characteristics of study participants at enrollment. Following [ICH E3 guidance](#), these tables are essential for understanding the study population and assessing comparability between treatment groups.

This tutorial shows you how to create a baseline characteristics table using Python's `rtflite` package.

```
import polars as pl # Data manipulation
import rtflite as rtf # RTF reporting
```

5.2 Step 1: Load Data

We start by loading the Subject-level Analysis Dataset (ADSL) and filtering to the safety population.

```
adsl = (
    pl.read_parquet("data/adsl.parquet")
    .select(["USUBJID", "TRT01P", "AGE", "SEX", "RACE"])
)

adsl
```

USUBJID	TRT01P	AGE	SEX	RACE
str	str	f64	str	str
"01-701-1015"	"Placebo"	63.0	"Female"	"White"
"01-701-1023"	"Placebo"	64.0	"Male"	"White"
"01-701-1028"	"Xanomeline High Dose"	71.0	"Male"	"White"
"01-701-1033"	"Xanomeline Low Dose"	74.0	"Male"	"White"
"01-701-1034"	"Xanomeline High Dose"	77.0	"Female"	"White"
...

USUBJID	TRT01P	AGE	SEX	RACE
str	str	f64	str	str
"01-718-1254"	"Xanomeline Low Dose"	78.0	"Male"	"White"
"01-718-1328"	"Xanomeline High Dose"	86.0	"Male"	"White"
"01-718-1355"	"Placebo"	79.0	"Male"	"White"
"01-718-1371"	"Xanomeline High Dose"	69.0	"Female"	"White"
"01-718-1427"	"Xanomeline High Dose"	74.0	"Female"	"Black Or African American"

5.3 Step 2: Calculate Summary Statistics

We'll create separate functions to handle continuous and categorical variables.

5.3.1 Continuous Variables (Age)

For continuous variables, we calculate mean (SD) and median [min, max].

```
def summarize_continuous(df, var):
    """Calculate summary statistics for continuous variables"""
    return df.groupby("TRT01P").agg([
        pl.col(var).mean().round(1).alias("mean"),
        pl.col(var).std().round(2).alias("sd"),
        pl.col(var).median().alias("median"),
        pl.col(var).min().alias("min"),
        pl.col(var).max().alias("max"),
        pl.len().alias("n")
    ])

age_stats = summarize_continuous(adsl, "AGE")
age_stats
```

TRT01P	mean	sd	median	min	max	n
str	f64	f64	f64	f64	f64	u32
"Xanomeline Low Dose"	75.7	8.29	77.5	51.0	88.0	84
"Xanomeline High Dose"	74.4	7.89	76.0	56.0	88.0	84
"Placebo"	75.2	8.59	76.0	52.0	89.0	86

5.3.2 Categorical Variables (Sex, Race)

For categorical variables, we calculate counts and percentages.

```
def summarize_categorical(df, var):
    """Calculate counts and percentages for categorical variables"""
    # Get counts by treatment and category
    counts = df.groupby(["TRT01P", var]).len()

    # Get treatment totals for percentage calculations
    totals = df.groupby("TRT01P").len().rename({"len": "total"})

    # Calculate percentages
    result = counts.join(totals, on="TRT01P").with_columns([
        (100.0 * pl.col("len") / pl.col("total")).round(1).alias("pct")
    ])

    return result

sex_stats = summarize_categorical(adsl, "SEX")
sex_stats
```

TRT01P	SEX	len	total	pct
str	str	u32	u32	f64
"Xanomeline High Dose"	"Female"	40	84	47.6
"Xanomeline Low Dose"	"Female"	50	84	59.5
"Xanomeline Low Dose"	"Male"	34	84	40.5
"Xanomeline High Dose"	"Male"	44	84	52.4
"Placebo"	"Male"	33	86	38.4
"Placebo"	"Female"	53	86	61.6

```
race_stats = summarize_categorical(adsl, "RACE")
race_stats
```

TRT01P	RACE	len	total	pct
str	str	u32	u32	f64
"Placebo"	"Black Or African American"	8	86	9.3
"Xanomeline High Dose"	"American Indian Or Alaska Nati..."	1	84	1.2
"Xanomeline Low Dose"	"Black Or African American"	6	84	7.1

TRT01P str	RACE str	len u32	total u32	pct f64
"Xanomeline Low Dose"	"White"	78	84	92.9
"Xanomeline High Dose"	"White"	74	84	88.1
"Xanomeline High Dose"	"Black Or African American"	9	84	10.7
"Placebo"	"White"	78	86	90.7

5.4 Step 3: Format Results

Now we format the statistics into the standard baseline table format.

5.4.1 Format Age Statistics

```
# Format age as "Mean (SD)" and "Median [Min, Max]"
age_formatted = age_stats.with_columns([
    pl.format("{} ({})", pl.col("mean"), pl.col("sd")).alias("mean_sd"),
    pl.format("{} [{}, {}]", pl.col("median"), pl.col("min"), pl.col("max")).alias("median_range")
]).select(["TRT01P", "mean_sd", "median_range"])

age_formatted
```

TRT01P str	mean_sd str	median_range str
"Xanomeline Low Dose"	"75.7 (8.29)"	"77.5 [51.0, 88.0]"
"Xanomeline High Dose"	"74.4 (7.89)"	"76.0 [56.0, 88.0]"
"Placebo"	"75.2 (8.59)"	"76.0 [52.0, 89.0]"

5.4.2 Format Categorical Statistics

```
# Format categorical as "n (%)"
sex_formatted = sex_stats.with_columns(
    pl.format("{} ({}%)", pl.col("len"), pl.col("pct")).alias("n_pct")
).select(["TRT01P", "SEX", "n_pct"])

race_formatted = race_stats.with_columns(
```

```
pl.format("{} ({}%)", pl.col("len"), pl.col("pct")).alias("n_pct")
).select(["TRT01P", "RACE", "n_pct"])
```

```
sex_formatted
```

TRT01P	SEX	n_pct
str	str	str
"Xanomeline High Dose"	"Female"	"40 (47.6%)"
"Xanomeline Low Dose"	"Female"	"50 (59.5%)"
"Xanomeline Low Dose"	"Male"	"34 (40.5%)"
"Xanomeline High Dose"	"Male"	"44 (52.4%)"
"Placebo"	"Male"	"33 (38.4%)"
"Placebo"	"Female"	"53 (61.6%)"

5.5 Step 4: Create Table Structure

We'll build the table row by row following the standard baseline table format.

```
# Helper function to get value for a treatment group
def get_value(df, treatment):
    """Get value for a specific treatment group or return default"""
    result = df.filter(pl.col("TRT01P") == treatment)
    return result[result.columns[-1]][0] if result.height > 0 else "0 (0.0%)"

# Build the baseline table structure
table_rows = []

# Age section
table_rows.append(["Age (years)", "", "", ""])

# Age Mean (SD) row
age_mean_row = ["  Mean (SD)"] + [
    get_value(age_formatted.select(["TRT01P", "mean_sd"]), trt).replace("0 (0.0%)", "")
    for trt in ["Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"]
]
table_rows.append(age_mean_row)

# Age Median [Min, Max] row
age_median_row = ["  Median [Min, Max]"] + [
```

```

        get_value(age_formatted.select(["TRT01P", "median_range"]), trt).replace("0 (0.0%)", "")
        for trt in ["Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"]
    ]
    table_rows.append(age_median_row)

# Sex section
table_rows.append(["Sex", "", "", ""])

for sex_cat in ["Female", "Male"]:
    sex_data = sex_formatted.filter(pl.col("SEX") == sex_cat)
    sex_row = [f" {sex_cat}"] + [
        get_value(sex_data, trt)
        for trt in ["Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"]
    ]
    table_rows.append(sex_row)

# Race section
table_rows.append(["Race", "", "", ""])

for race_cat in ["White", "Black Or African American", "American Indian Or Alaska Native"]:
    race_data = race_formatted.filter(pl.col("RACE") == race_cat)
    race_row = [f" {race_cat}"] + [
        get_value(race_data, trt)
        for trt in ["Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"]
    ]
    table_rows.append(race_row)

# Create DataFrame from table rows
baseline_table = pl.DataFrame(
    table_rows,
    schema=["Characteristic", "Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"],
    orient="row"
)

baseline_table

```

Characteristic str	Placebo str	Xanomeline Low Dose str	Xanomeline High Dose str
"Age (years)"	"	"	"
" Mean (SD)"	"75.2 (8.59)"	"75.7 (8.29)"	"74.4 (7.89)"
" Median [Min, Max]"	"76.0 [52.0, 89.0]"	"77.5 [51.0, 88.0]"	"76.0 [56.0, 88.0]"

Characteristic str	Placebo str	Xanomeline Low Dose str	Xanomeline High Dose str
"Sex"	" "	" "	" "
" Female"	"53 (61.6%)"	"50 (59.5%)"	"40 (47.6%)"
" Male"	"33 (38.4%)"	"34 (40.5%)"	"44 (52.4%)"
"Race"	" "	" "	" "
" White"	"78 (90.7%)"	"78 (92.9%)"	"74 (88.1%)"
" Black Or African American"	"8 (9.3%)"	"6 (7.1%)"	"9 (10.7%)"
" American Indian Or Alaska Na..."	"0 (0.0%)"	"0 (0.0%)"	"1 (1.2%)"

5.6 Step 5: Generate Publication-Ready Output

Finally, we format the baseline table for regulatory submission using the `rtflite` package.

```
# Get treatment group sizes for column headers
treatment_n = adsl.group_by("TRT01P").len().sort("TRT01P")
n_placebo = treatment_n.filter(pl.col("TRT01P") == "Placebo")["len"][0]
n_low = treatment_n.filter(pl.col("TRT01P") == "Xanomeline Low Dose")["len"][0]
n_high = treatment_n.filter(pl.col("TRT01P") == "Xanomeline High Dose")["len"][0]

doc_baseline = rtf.RTFDocument(
    df=baseline_table,
    rtf_title=rtf.RTFTitle(
        text=[
            "Baseline Characteristics of Participants",
            "(All Participants Randomized)"
        ]
    ),
    rtf_column_header=rtf.RTFColumnHeader(
        text=[
            "Characteristic",
            f"Placebo\n(N={n_placebo})",
            f"Xanomeline Low Dose\n(N={n_low})",
            f"Xanomeline High Dose\n(N={n_high})"
        ],
        text_justification=["l", "c", "c", "c"],
        col_rel_width=[3, 2, 2, 2]
    ),
    rtf_body=rtf.RTFBody(
        text_justification=["l", "c", "c", "c"],
```

```
        col_rel_width=[3, 2, 2, 2]
    ),
    rtf_source=rtf.RTFSource(text=["Source: ADSL dataset"])
)

doc_baseline.write_rtf("rtf/tlf_baseline.rtf") # Save as RTF for submission
```

rtf/tlf_baseline.rtf

6 Adverse Events Summary Table

6.1 Overview

Adverse events (AE) summary tables are critical safety assessments required in clinical study reports. Following [ICH E3 guidance](#), these tables summarize the overall safety profile by showing the number and percentage of participants experiencing various categories of adverse events across treatment groups.

Key categories typically include:

- **Any adverse event:** Total participants with at least one AE
- **Drug-related events:** Events potentially related to study treatment
- **Serious adverse events:** Events meeting regulatory criteria for seriousness
- **Deaths:** Fatal outcomes
- **Discontinuations:** Participants who stopped treatment due to AEs

This tutorial shows you how to create an AE summary table using Python's `rtflite` package.

```
import polars as pl
import rtflite as rtf
```

6.2 Step 1: Load Data

We need two datasets for AE analysis: the subject-level dataset (ADSL) and the adverse events dataset (ADAE).

```
# Load datasets
adsl = pl.read_parquet("data/adsl.parquet")
adae = pl.read_parquet("data/adae.parquet")

# Display key variables from ADSL
adsl.select(["USUBJID", "TRT01A", "SAFFL"]).head()
```

USUBJID str	TRT01A str	SAFFL str
"01-701-1015"	"Placebo"	"Y"
"01-701-1023"	"Placebo"	"Y"
"01-701-1028"	"Xanomeline High Dose"	"Y"
"01-701-1033"	"Xanomeline Low Dose"	"Y"
"01-701-1034"	"Xanomeline High Dose"	"Y"

```
# Display key variables from ADAE
adae.select(["USUBJID", "AEREL", "AESER", "AEOUT", "AEACN"]).head()
```

USUBJID str	AEREL str	AESER str	AEOUT str	AEACN str
"01-701-1015"	"PROBABLE"	"N"	"NOT RECOVERED/NOT RESOLVED"	" "
"01-701-1015"	"PROBABLE"	"N"	"NOT RECOVERED/NOT RESOLVED"	" "
"01-701-1015"	"REMOTE"	"N"	"RECOVERED/RESOLVED"	" "
"01-701-1023"	"POSSIBLE"	"N"	"NOT RECOVERED/NOT RESOLVED"	" "
"01-701-1023"	"PROBABLE"	"N"	"NOT RECOVERED/NOT RESOLVED"	" "

Key ADAE variables used in this analysis:

- **USUBJID**: Unique subject identifier to link with ADSL
- **AEREL**: Relationship of adverse event to study drug (e.g., "RELATED", "POSSIBLE", "PROBABLE", "DEFINITE", "NOT RELATED")
- **AESER**: Serious adverse event flag ("Y" = serious, "N" = not serious)
- **AEOUT**: Outcome of adverse event (e.g., "RECOVERED", "RECOVERING", "NOT RECOVERED", "FATAL")
- **AEACN**: Action taken with study treatment (e.g., "DOSE NOT CHANGED", "DRUG WITHDRAWN", "DOSE REDUCED")

6.3 Step 2: Filter Safety Population

For safety analyses, we focus on participants who received at least one dose of study treatment.

```
# Filter to safety population
adsl_safety = adsl.filter(pl.col("SAFFL") == "Y").select(["USUBJID", "TRT01A"])

# Get treatment counts for denominators
```

```

pop_counts = adsl_safety.group_by("TRT01A").agg(
    N = pl.len()
).sort("TRT01A")

# Preserve the treatment level order for downstream joins
treatment_levels = pop_counts.select(["TRT01A"])

# Safety population by treatment
pop_counts

```

TRT01A	N
str	u32
"Placebo"	86
"Xanomeline High Dose"	84
"Xanomeline Low Dose"	84

```

# Join treatment information to AE data
adae_safety = adae.join(adsl_safety, on="USUBJID")

# Total AE records in safety population
adae_safety.height

```

1191

6.4 Step 3: Define AE Categories

We'll calculate participant counts for standard AE categories used in regulatory submissions.

```

def count_participants(df, condition=None):
    """
    Count unique participants meeting a condition

    Args:
        df: DataFrame with adverse events
        condition: polars expression for filtering (None = count all)

    Returns:
        DataFrame with counts by treatment
    """

```

```

"""
if condition is not None:
    df = df.filter(condition)

counts = df.group_by("TRT01A").agg(
    n = pl.col("USUBJID").n_unique()
)

return treatment_levels.join(counts, on="TRT01A", how="left").with_columns(
    pl.col("n").fill_null(0)
)

# Calculate each category
categories = []

# 1. Participants in population (no filtering)
pop_row = pop_counts.with_columns(
    category = pl.lit("Participants in population")
).rename({"N": "n"})
categories.append(pop_row)

# 2. With any adverse event
any_ae = count_participants(adae_safety).with_columns(
    category = pl.lit("With any adverse event")
)
categories.append(any_ae)

# 3. With drug-related adverse event
drug_related = count_participants(
    adae_safety,
    pl.col("AEREL").is_in(["POSSIBLE", "PROBABLE", "DEFINITE", "RELATED"])
).with_columns(
    category = pl.lit("With drug-related adverse event")
)
categories.append(drug_related)

# 4. With serious adverse event
serious = count_participants(
    adae_safety,
    pl.col("AESER") == "Y"
).with_columns(
    category = pl.lit("With serious adverse event")
)

```

```

)
categories.append(serious)

# 5. With serious drug-related adverse event
serious_drug_related = count_participants(
    adae_safety,
    (pl.col("AESER") == "Y") &
    pl.col("AEREL").is_in(["POSSIBLE", "PROBABLE", "DEFINITE", "RELATED"])
).with_columns(
    category = pl.lit("With serious drug-related adverse event")
)
categories.append(serious_drug_related)

# 6. Who died
deaths = count_participants(
    adae_safety,
    pl.col("AEOUT") == "FATAL"
).with_columns(
    category = pl.lit("Who died")
)
categories.append(deaths)

# 7. Discontinued due to adverse event
discontinued = count_participants(
    adae_safety,
    pl.col("AEACN") == "DRUG WITHDRAWN"
).with_columns(
    category = pl.lit("Discontinued due to adverse event")
)
categories.append(discontinued)

```

6.5 Step 4: Combine and Calculate Percentages

Now we combine all categories and calculate percentages based on the safety population.

```

# Combine all categories
ae_summary = pl.concat(categories, how="diagonal")

# Add population totals and calculate percentages
ae_summary = ae_summary.join(

```

```

pop_counts.select(["TRT01A", "N"]),
on="TRT01A",
how="left"
).with_columns([
    # Fill missing counts with 0
    pl.col("n").fill_null(0),
    # Calculate percentage
    pl.when(pl.col("category") == "Participants in population")
        .then(None) # No percentage for population row
        .otherwise((100.0 * pl.col("n") / pl.col("N")).round(1))
        .alias("pct")
])

ae_summary.sort(["category", "TRT01A"])

```

TRT01A	n	category	N	pct
str	u32	str	u32	f64
"Placebo"	0	"Discontinued due to adverse ev...	86	0.0
"Xanomeline High Dose"	0	"Discontinued due to adverse ev...	84	0.0
"Xanomeline Low Dose"	0	"Discontinued due to adverse ev...	84	0.0
"Placebo"	86	"Participants in population"	86	null
"Xanomeline High Dose"	84	"Participants in population"	84	null
...
"Xanomeline High Dose"	2	"With serious adverse event"	84	2.4
"Xanomeline Low Dose"	1	"With serious adverse event"	84	1.2
"Placebo"	0	"With serious drug-related adve...	86	0.0
"Xanomeline High Dose"	1	"With serious drug-related adve...	84	1.2
"Xanomeline Low Dose"	1	"With serious drug-related adve...	84	1.2

6.6 Step 5: Format for Display

We'll format the counts and percentages for the final table display.

```

# Format display values
ae_formatted = ae_summary.with_columns([
    # Show counts as strings, including zeros
    pl.col("n").cast(str).alias("n_display"),
    # Format percentages with parentheses; blank out population row
    pl.when(pl.col("category") == "Participants in population")

```



```

        .then(pl.lit(""))
        .otherwise(
            pl.format("{} ", pl.col("pct").fill_null(0).round(1).cast(str))
        )
        .alias("pct_display")
    ])

ae_formatted.select(["category", "TRT01A", "n_display", "pct_display"])

```

category str	TRT01A str	n_display str	pct_display str
"Participants in population"	"Placebo"	"86"	" "
"Participants in population"	"Xanomeline High Dose"	"84"	" "
"Participants in population"	"Xanomeline Low Dose"	"84"	" "
"With any adverse event"	"Placebo"	"69"	"(80.2)"
"With any adverse event"	"Xanomeline High Dose"	"79"	"(94.0)"
...
"Who died"	"Xanomeline High Dose"	"0"	"(0.0)"
"Who died"	"Xanomeline Low Dose"	"1"	"(1.2)"
"Discontinued due to adverse ev...	"Placebo"	"0"	"(0.0)"
"Discontinued due to adverse ev...	"Xanomeline High Dose"	"0"	"(0.0)"
"Discontinued due to adverse ev...	"Xanomeline Low Dose"	"0"	"(0.0)"

6.7 Step 6: Create Final Table Structure

We reshape the data to create the final table with treatments as columns.

```

# Define category order for consistent display
category_order = [
    "Participants in population",
    "With any adverse event",
    "With drug-related adverse event",
    "With serious adverse event",
    "With serious drug-related adverse event",
    "Who died",
    "Discontinued due to adverse event"
]

# Pivot to wide format

```

```

ae_wide = ae_formatted.pivot(
    values=["n_display", "pct_display"],
    index="category",
    on="TRT01A",
    maintain_order=True
)

# Reorder columns for each treatment group
treatments = ["Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"]
column_order = ["category"]
for trt in treatments:
    column_order.extend([f"n_display_{trt}", f"pct_display_{trt}"])

# Create final table with proper column order
final_table = ae_wide.select(column_order).sort(
    pl.col("category").cast(pl.Enum(category_order))
)

final_table

```

category str	n_display_Placebo str	pct_display_Placebo str	n_display_Xanomeline Low str
"Participants in population"	"86"	"	"84"
"With any adverse event"	"69"	"(80.2)"	"77"
"With drug-related adverse even..."	"44"	"(51.2)"	"73"
"With serious adverse event"	"0"	"(0.0)"	"1"
"With serious drug-related adve..."	"0"	"(0.0)"	"1"
"Who died"	"2"	"(2.3)"	"1"
"Discontinued due to adverse ev..."	"0"	"(0.0)"	"0"

6.8 Step 7: Generate Publication-Ready Output

Finally, we format the AE summary table for regulatory submission using the `rtflite` package.

```

# Get population sizes for column headers
n_placebo = pop_counts.filter(pl.col("TRT01A") == "Placebo")["N"][0]
n_low = pop_counts.filter(pl.col("TRT01A") == "Xanomeline Low Dose")["N"][0]
n_high = pop_counts.filter(pl.col("TRT01A") == "Xanomeline High Dose")["N"][0]

```

```

doc_ae_summary = rtf.RTFDocument(
    df=final_table.rename({"category": ""}),
    rtf_title=rtf.RTFTitle(
        text=[
            "Analysis of Adverse Event Summary",
            "(Safety Analysis Population)"
        ]
    ),
    rtf_column_header=[
        rtf.RTFColumnHeader(
            text = [
                "",
                "Placebo",
                "Xanomeline Low Dose",
                "Xanomeline High Dose"
            ],
            col_rel_width=[4, 2, 2, 2],
            text_justification=["l", "c", "c", "c"],
        ),
        rtf.RTFColumnHeader(
            text=[
                "",          # Empty for first column
                "n", "(%)",  # Placebo columns
                "n", "(%)",  # Low Dose columns
                "n", "(%)"   # High Dose columns
            ],
            col_rel_width=[4] + [1] * 6,
            text_justification=["l"] + ["c"] * 6,
            border_left = ["single"] + ["single", ""] * 3,
            border_top = [""] + ["single"] * 6
        )
    ],
    rtf_body=rtf.RTFBody(
        col_rel_width=[4] + [1] * 6,
        text_justification=["l"] + ["c"] * 6,
        border_left = ["single"] + ["single", ""] * 3
    ),
    rtf_footnote=rtf.RTFFootnote(
        text=[
            "Every subject is counted a single time for each applicable row and column."
        ]
    ),

```

```
    rtf_source=rtf.RTFSource(  
        text=["Source: ADSL and ADAE datasets"]  
    )  
)  
  
doc_ae_summary.write_rtf("rtf/tlf_ae_summary.rtf")
```

rtf/tlf_ae_summary.rtf

7 Specific Adverse Events Table

This article demonstrates how to create a specific adverse events table by System Organ Class and Preferred Term.

7.1 Setup

```
import polars as pl
import rtflite as rtf

adsl = pl.read_parquet("data/adsl.parquet")
adae = pl.read_parquet("data/adae.parquet")
treatments = ["Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"]
```

7.2 Prepare AE Summary Data

```
# Get safety population counts and AE data
adsl_safety = adsl.filter(pl.col("SAFFL") == "Y").select(["USUBJID", "TRT01A"])
adae_safety = adae.join(adsl_safety, on="USUBJID", how="inner")
pop_counts = adsl_safety.group_by("TRT01A").agg(N=pl.len()).sort("TRT01A")

# Calculate AE counts by SOC and term
ae_counts = (
    adae_safety.with_columns(pl.col("AEDECOD").str.to_titlecase())
    .group_by(["TRT01A", "AEBODSYS", "AEDECOD"])
    .agg(n=pl.col("USUBJID").n_unique())
    .sort(["AEBODSYS", "AEDECOD", "TRT01A"])
)

# Build table rows
table_data = [
    ["Participants in population"] + [str(pop_counts.filter(pl.col("TRT01A") == t)["N"][0])]
```

```

[""] * 4 # Blank row
]

# Add SOC and AE term rows
for soc in ae_counts["AEBODSYS"].unique().sort():
    table_data.append([soc] + [""] * 3)
    soc_data = ae_counts.filter(pl.col("AEBODSYS") == soc)

    for ae in soc_data["AEDECOD"].unique().sort():
        row = [f" {ae}"]
        for trt in treatments:
            count = soc_data.filter((pl.col("AEDECOD") == ae) & (pl.col("TRT01A") == trt))
            row.append(str(count["n"][0]) if count.height > 0 else "0")
        table_data.append(row)

df_ae_specific = pl.DataFrame(table_data, schema=[""] + treatments, orient="row")

```

7.3 Create RTF Output

```

doc_ae_specific = rtf.RTFDocument(
    df=df_ae_specific,
    rtf_title=rtf.RTFTitle(text=["Specific Adverse Events", "(Safety Analysis Population)"]),
    rtf_column_header=rtf.RTFColumnHeader(
        text=["", "Placebo\n\n", "Xanomeline Low Dose\n\n", "Xanomeline High Dose\n\n"],
        col_rel_width=[4, 1, 1, 1],
        text_justification=["l", "c", "c", "c"],
    ),
    rtf_body=rtf.RTFBody(
        col_rel_width=[4, 1, 1, 1],
        text_justification=["l", "c", "c", "c"],
        text_font_style=lambda df, i, j: "bold" if j == 0 and " " not in str(df[i, j]) else "",
    ),
    rtf_footnote=rtf.RTFFootnote(text=["Number of participants with specific adverse events."]),
    rtf_source=rtf.RTFSource(text=["Source: ADSL and ADAE datasets"])
)

doc_ae_specific.write_rtf("rtf/tlf_ae_specific.rtf")

```

rtf/tlf_ae_specific.rtf

```
PosixPath('pdf/tlf_ae_specific.pdf')
```

8 ANCOVA Efficacy Analysis

This article demonstrates how to create an ANCOVA efficacy table for glucose levels at Week 24 with LOCF imputation.

8.1 Setup

```
import polars as pl
import rtflite as rtf
import pandas as pd
import numpy as np
import statsmodels.formula.api as smf
from scipy import stats as scipy_stats
from importlib.resources import files

adsl = pl.read_parquet("data/adsl.parquet")
adlbc = pl.read_parquet("data/adlbc.parquet")
treatments = ["Placebo", "Xanomeline Low Dose", "Xanomeline High Dose"]
```

8.2 Prepare Analysis Data

```
# Clean data types and filter for efficacy population
adlbc_clean = adlbc.with_columns(
    [pl.col(c).cast(str).str.strip_chars() for c in ["USUBJID", "PARAMCD", "AVISIT", "TRTP"]]
)
adsl_eff = adsl.filter(pl.col("EFFFL") == "Y").select(["USUBJID"])
adlbc_eff = adlbc_clean.join(adsl_eff, on="USUBJID", how="inner")

# Apply LOCF for glucose data up to Week 24
gluc_data = (
    adlbc_eff.filter((pl.col("PARAMCD") == "GLUC") & (pl.col("AVISITN") <= 24))
    .sort(["USUBJID", "AVISITN"])
```



```

.group_by("USUBJID")
.agg([
    pl.col("TRTP").first(),
    pl.col("BASE").first(),
    pl.col("AVAL").filter(pl.col("AVISITN") == 0).first().alias("Baseline"),
    pl.col("AVAL").last().alias("Week 24")
])
.filter(pl.col("Baseline").is_not_null() & pl.col("Week 24").is_not_null())
.with_columns((pl.col("Week 24") - pl.col("Baseline")).alias("CHG"))
)

# Calculate descriptive statistics
desc_stats = []
for trt in treatments:
    trt_data = gluc_data.filter(pl.col("TRTP") == trt)
    baseline_full = adlbc_eff.filter(
        (pl.col("PARAMCD") == "GLUC") & (pl.col("AVISIT") == "Baseline") & (pl.col("TRTP") ==
    )

    desc_stats.append({
        "Treatment": trt,
        "N_Baseline": baseline_full.height,
        "Baseline_Mean": baseline_full["AVAL"].mean() if baseline_full.height > 0 else np.nan,
        "Baseline_SD": baseline_full["AVAL"].std() if baseline_full.height > 0 else np.nan,
        "N_Week24": trt_data.height,
        "Week24_Mean": trt_data["Week 24"].mean() if trt_data.height > 0 else np.nan,
        "Week24_SD": trt_data["Week 24"].std() if trt_data.height > 0 else np.nan,
        "N_Change": trt_data.height,
        "Change_Mean": trt_data["CHG"].mean() if trt_data.height > 0 else np.nan,
        "Change_SD": trt_data["CHG"].std() if trt_data.height > 0 else np.nan
    })

# Perform ANCOVA
ancova_df = gluc_data.to_pandas()
ancova_df["TRTP"] = pd.Categorical(ancova_df["TRTP"], categories=treatments)
model = smf.ols("CHG ~ TRTP + BASE", data=ancova_df).fit()

# Calculate LS means and confidence intervals
base_mean = ancova_df["BASE"].mean()
var_cov = model.cov_params()
ls_means = []

```

```

for i, trt in enumerate(treatments):
    x_pred = np.array([1, int(i==1), int(i==2), base_mean])
    ls_mean = model.predict(pd.DataFrame({"TRTP": [trt], "BASE": [base_mean]}))[0]
    se_pred = np.sqrt(x_pred @ var_cov @ x_pred.T)

    ls_means.append({
        "Treatment": trt,
        "LS_Mean": ls_mean,
        "CI_Lower": ls_mean - 1.96 * se_pred,
        "CI_Upper": ls_mean + 1.96 * se_pred
    })

```

8.3 Create Tables for RTF Output

```

# Table 1: Descriptive Statistics
tbl1_data = [
    [
        s["Treatment"],
        str(s["N_Baseline"]),
        f"{s['Baseline_Mean']:.1f} ({s['Baseline_SD']:.2f})",
        str(s["N_Week24"]),
        f"{s['Week24_Mean']:.1f} ({s['Week24_SD']:.2f})",
        str(s["N_Change"]),
        f"{s['Change_Mean']:.1f} ({s['Change_SD']:.2f})",
        f"{ls['LS_Mean']:.2f} ({ls['CI_Lower']:.2f}, {ls['CI_Upper']:.2f})"
    ]
    for s, ls in zip(desc_stats, ls_means)
]

tbl1 = pl.DataFrame(tbl1_data, orient="row", schema=[
    "Treatment", "N_Base", "Mean_SD_Base", "N_Wk24", "Mean_SD_Wk24",
    "N_Chg", "Mean_SD_Chg", "LS_Mean_CI"
])

# Table 2: Pairwise Comparisons
tbl2_data = []
for comp_name, trt_name in [
    ("Xanomeline Low Dose - Placebo", "TRTP[T.Xanomeline Low Dose]"),
    ("Xanomeline High Dose - Placebo", "TRTP[T.Xanomeline High Dose]"),
]:
    coef = model.params[trt_name]
    se = model.bse[trt_name]

```

```

t_stat = coef / se
p_value = 2 * (1 - scipy_stats.t.cdf(abs(t_stat), model.df_resid))

tbl2_data.append([
    comp_name,
    f"{coef:.2f} ({coef - 1.96*se:.2f}, {coef + 1.96*se:.2f})",
    f"{p_value:.3f}"
])

tbl2 = pl.DataFrame(tbl2_data, orient="row", schema=["Comparison", "Diff_CI", "P_Value"])

```

8.4 Create RTF Document

```

# Create RTF document with two sections
doc_ancova = rtf.RTFDocument(
    df=[tbl1, tbl2],
    rtf_title=rtf.RTFTitle(text=[
        "ANCOVA of Change from Baseline Glucose (mmol/L) at Week 24", "LOCF",
        "Efficacy Analysis Population"
    ]),
    rtf_column_header=[
        [rtf.RTFColumnHeader(text=["", "Baseline", "Week 24", "Change from Baseline"],
            col_rel_width=[3, 2, 2, 4], text_justification=["l", "c", "c", "c"],
            border_bottom="single"),
        rtf.RTFColumnHeader(text=["Treatment", "N", "Mean (SD)", "N", "Mean (SD)", "N",
            "Mean (SD)", "LS Mean (95% CI){^a}"],
            col_rel_width=[3, 0.7, 1.3, 0.7, 1.3, 0.7, 1.3, 2],
            text_justification=["l"] + ["c"] * 7, border_bottom="single")],
        [rtf.RTFColumnHeader(text=["Pairwise Comparison", "Difference in LS Mean (95% CI){^a}"],
            col_rel_width=[5, 4, 2], text_justification=["l", "c", "c"])]
    ],
    rtf_body=[
        rtf.RTFBody(col_rel_width=[3, 0.7, 1.3, 0.7, 1.3, 0.7, 1.3, 2],
            text_justification=["l"] + ["c"] * 7),
        rtf.RTFBody(col_rel_width=[5, 4, 2], text_justification=["l", "c", "c"])
    ],
    rtf_footnote=rtf.RTFFootnote(text=[
        "{^a}Based on an ANCOVA model after adjusting baseline value. LOCF approach is used t
        "ANCOVA = Analysis of Covariance, LOCF = Last Observation Carried Forward",
        "CI = Confidence Interval, LS = Least Squares, SD = Standard Deviation"
    ])

```

```
    ]),  
    rtf_source=rtf.RTFSource(text=["Source: ADLBC dataset"])  
)  
  
doc_ancova.write_rtf("rtf/tlf_efficacy_ancova.rtf")
```

rtf/tlf_efficacy_ancova.rtf

PosixPath('pdf/tlf_efficacy_ancova.pdf')

9 Summary

In summary, this book has no content whatsoever.

1 + 1

2

References

Knuth, Donald E. 1984. “Literate Programming.” *Comput. J.* 27 (2): 97–111.