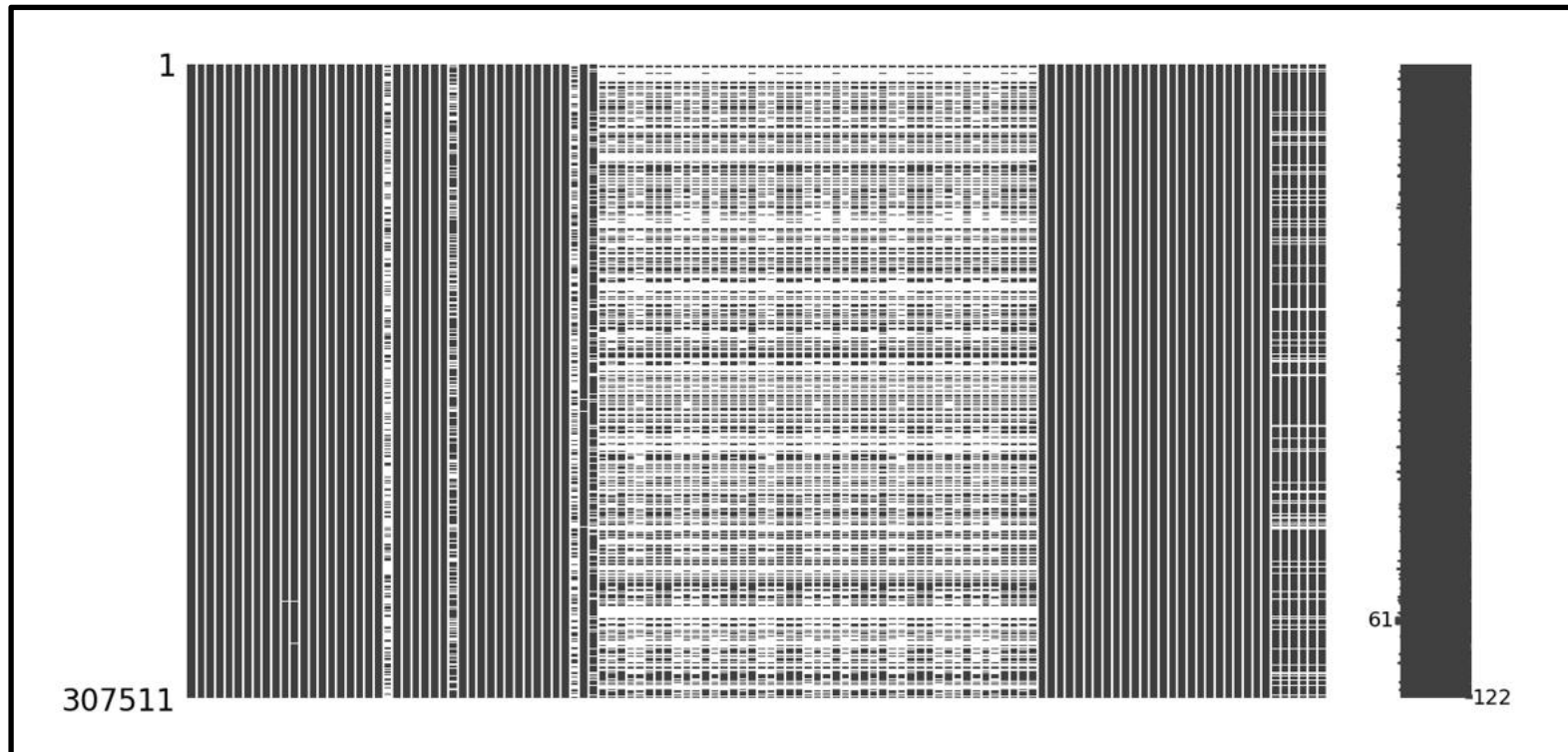


**Graph1:-** from what I can see, this image appears to be a **missing values heatmap (null-value matrix)**, likely generated using **Python's seaborn or missingno library** for a dataset.



### .Graph Overview

- Each vertical column represents one feature (column) in your dataset.
- Each horizontal line represents one data record (row) — looks like there are 307,511 rows (shown bottom left).

### Color meaning:

- Black (or dark areas) → Values present (non-null).
- White (or empty stripes) → Missing (NaN) values.

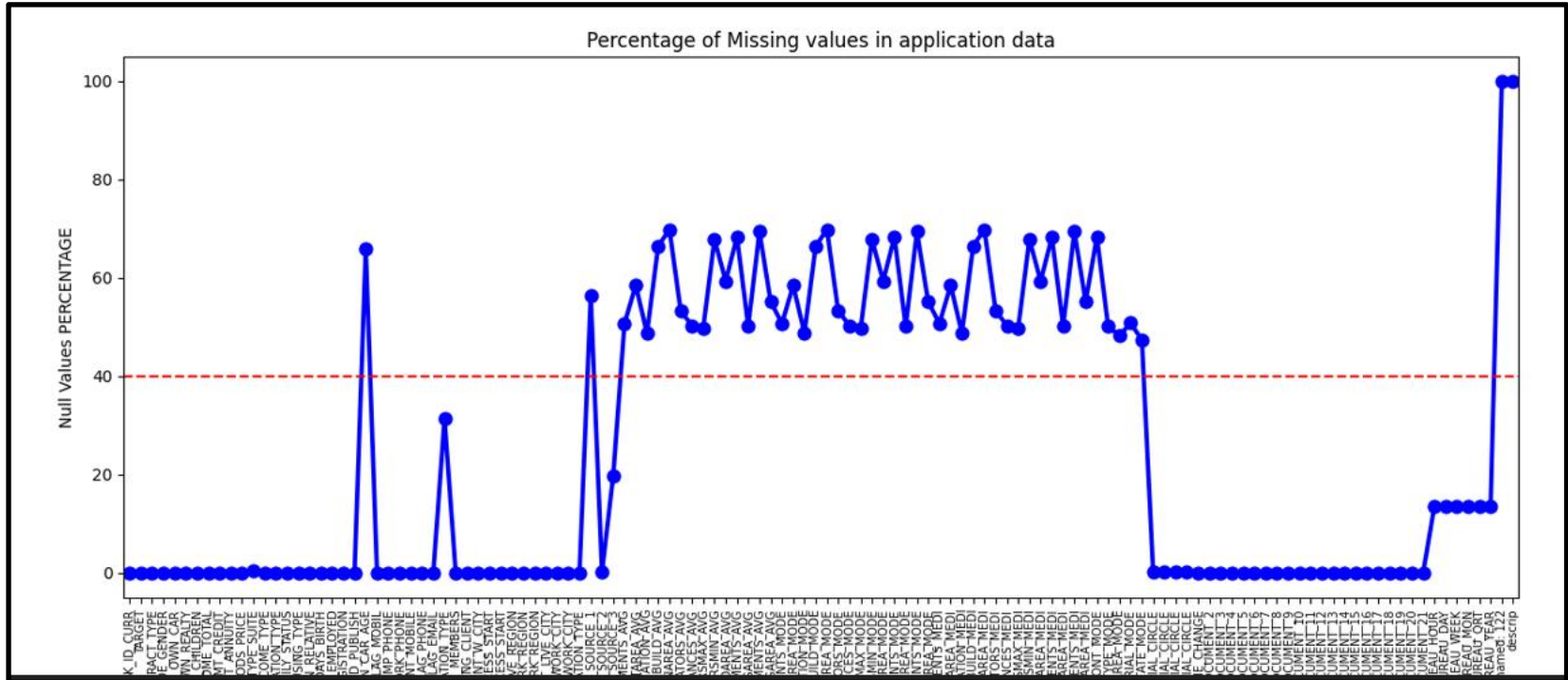
## Interpretation

- The **left side** of the graph has many solid black columns — meaning **those features have no missing values**.
- The **middle region** has a mix of black and white — these columns have **some missing data**.
- The **right side** also has mostly black columns, suggesting **complete data again**.
- The **color bar on the right** (with values like 122, 61, etc.) likely indicates **the range of missingness** or some **data index reference** depending on how the heatmap was plotted.

This graph visually helps you:

- Spot which features (columns) have **missing data**.
- Estimate **how much** data is missing per feature.
- Identify **patterns of missingness** — for example, if certain rows have multiple missing features together.

**Graph2:-** This graph represents the **percentage of missing (null) values** in each feature (column) of an application dataset.



## 1. Title

**“Percentage of Missing values in application data” —**

It indicates that this visualization is used to assess data quality by showing how much data is missing per column.

## 2. X-Axis (Horizontal)

It shows feature/column names from your dataset.

The labels are rotated and dense, suggesting there are many columns.

### 3. Y-Axis (Vertical)

- It represents the **percentage of missing values** for each column.
- The scale goes from 0% to 100%.

### 4. Blue Line and Dots

- Each blue **dot** corresponds to a single column's percentage of missing data.
- The **blue line** connects them, showing trends across columns.
- Columns with higher points (close to 100%) have more missing data.

### 5. Red Dashed Line

- This line appears around **40%** on the Y-axis.
- It likely represents a **threshold** — columns with more than 40% missing values may be considered for **removal or special treatment** during data cleaning.

### 6. Observations

- Many features have **0% missing values** (flat sections at the bottom).
- A cluster of features in the middle range have **50–70% missing values**, which is significant.
- A few extreme columns (like the last one) have **100% missing values** — these are completely empty.
- Some isolated spikes (around 30%, 65%, etc.) indicate inconsistent data quality across fields.

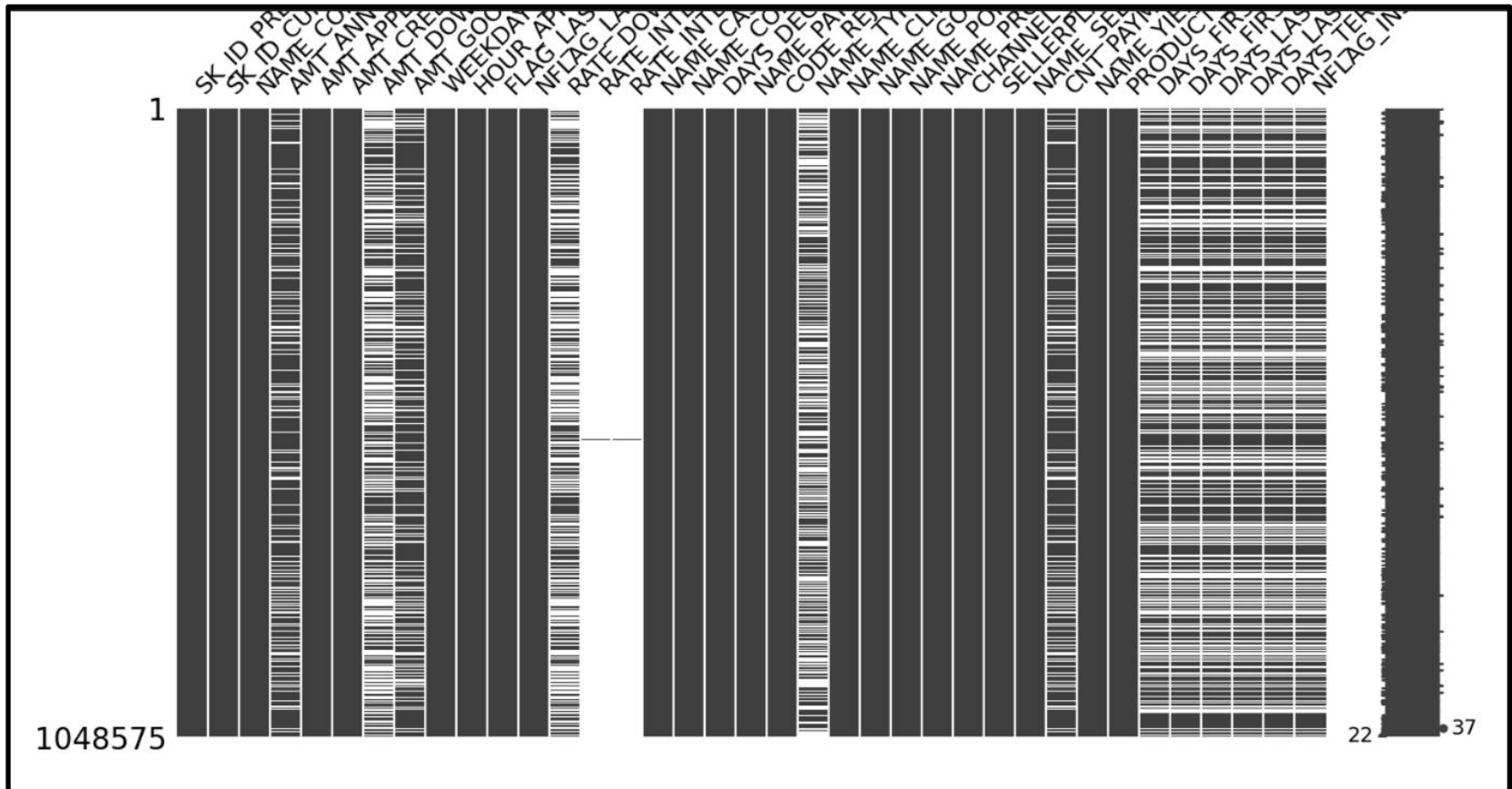
### 7. Insights / Implications

- The dataset has **uneven data completeness**.
- Features above the 40% red line may **reduce model reliability** if not handled properly.

You could:

- Drop columns with **>40% missing** (or 100% missing).
- Impute values for partially missing ones.
- Investigate why some fields are heavily incomplete (data collection issues, optional fields, etc.).

**Graph3:-** The image you uploaded is a **missing data matrix plot**, most likely created using the missingno or seaborn library in Python (for example, msno.matrix(df)).



What the graph shows

- Each vertical bar represents a column (feature) in your dataset.
- Each horizontal line represents a row (data record).
- The white lines represent missing (null) values.
- The black (or dark grey) areas represent non-missing (valid) values.

## Key observations

### Columns on the x-axis:

The names like SK\_ID\_PREV, NAME\_CONTRACT\_TYPE, AMT\_ANNUITY, etc., are the feature names from your dataset.

### Rows on the y-axis:

The number 1048575 on the left side indicates that your dataset has around **1,048,575 rows** (over 1 million records).

### Missing values pattern:

- Columns that are **entirely black** have **no missing values**.
- Columns that have **white stripes or gaps** have **missing entries**.
- Columns on the **right side** appear to have more missing data (lots of white), suggesting they might be optional or derived fields.

### Correlation between missing values:

The continuous blocks of missing data suggest that certain features might be missing **together**, meaning they could be related.

## Interpretation / Insights

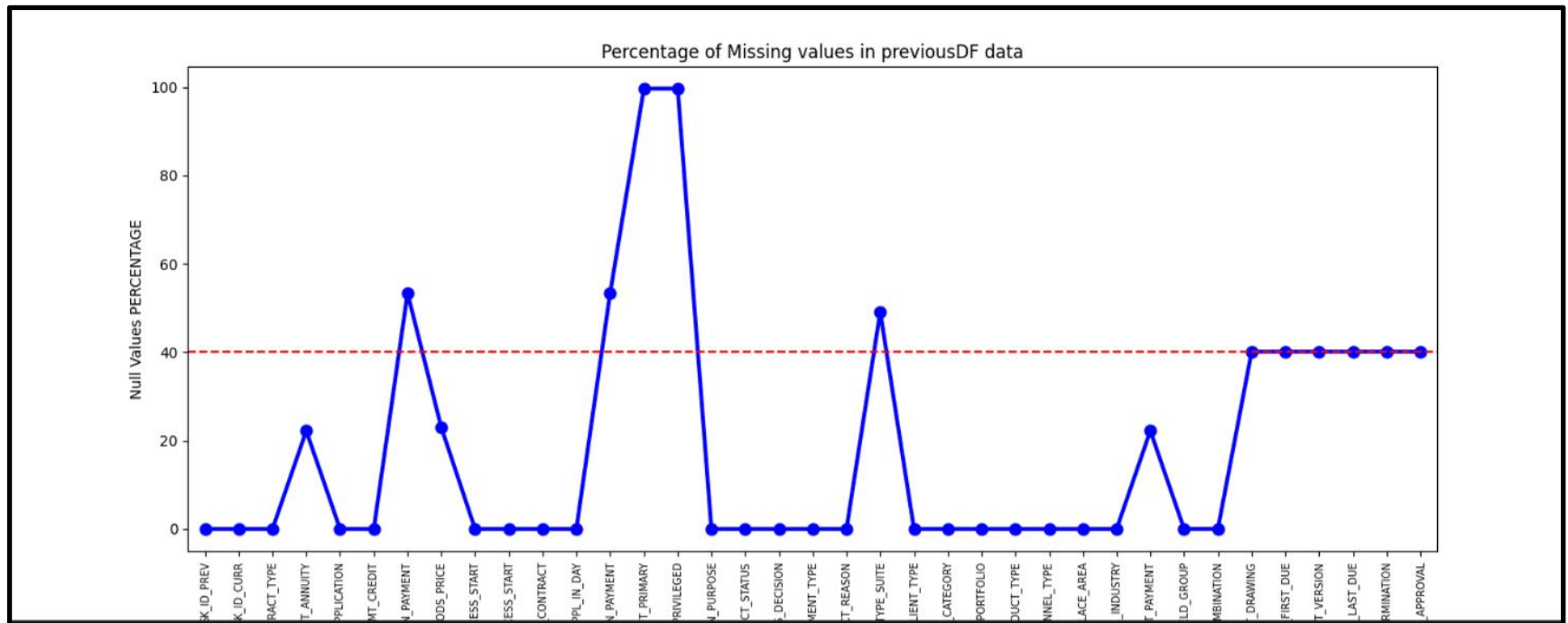
- You can identify **which columns are most complete or incomplete**.

This helps you decide:

- Whether to **impute**, **drop**, or **ignore** certain features.
- If **entire rows** with too many missing values should be removed.
- It also gives you a quick idea of **data quality** and **readiness** for modeling.

This matrix visualization gives an **at-a-glance overview of missing values** across all features and observations in your dataset.

**Graph4:-** The graph shows the **percentage of missing (null) values** for each column in a DataFrame named previousDF.



X-axis: Column names of the dataset (e.g., SK\_ID\_PREV, NAME\_CONTRACT\_TYPE, AMT\_ANNUITY, etc.)

Y-axis: Percentage of missing values in each column

Blue line: Represents the percentage of null values per column

Red dashed line: Represents a reference threshold (at 40%) — used to highlight which columns have too many missing values.

## Interpretation

### Columns with Very High Missing Values (~100%)

A few columns have nearly **100% missing values** — meaning almost no data is available in those features. Examples (from the visible spikes): DAYS\_FIRST\_DRAWING, DAYS\_FIRST\_DUE, or similar columns. Such columns often need to be **dropped** since they don't contribute useful information.

### Columns with Moderate Missing Values (40–60%)

Some columns (e.g., around AMT\_PAYMENT, NAME\_TYPE\_SUITE) have **40–60% missing data**. These may still hold partial information; you could **impute** (fill) them or **analyze their importance** before deciding.

### Columns with Low or No Missing Values (0–20%)

Most columns (majority near the baseline) have **very few or no missing values**, which is good. These columns are reliable for analysis and modeling.

### Red Threshold Line (40%)

The red dashed line marks a **cutoff** point — typically, analysts decide to **drop columns above this threshold** since they have too many missing values to be useful.

## Insights

- The dataset has a **mixed data quality** — most columns are fine, but a few are severely incomplete. For preprocessing:
- Columns above 40% missing could be **dropped or treated separately**.
- Remaining missing values could be handled using **mean/median imputation, forward fill, or model-based imputation** depending on data type.

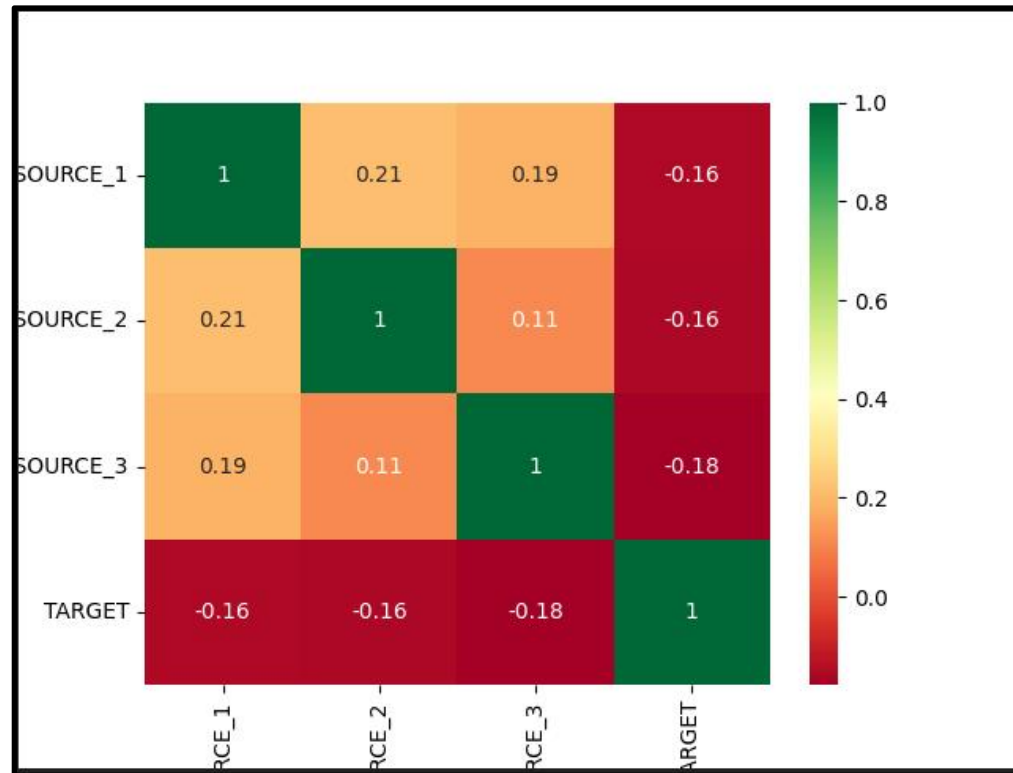
## Summary

Category	Missing % Range	Suggested Action
Low	0–20%	Keep and fill missing values
Moderate	20–40%	Consider imputation or partial removal
High	>40%	Drop or analyze necessity before keeping



**Graph5:-** This is a **correlation heatmap**, which visually represents the correlation coefficients between variables.

Let's interpret it step by step



What It Shows

The variables are:

- SOURCE\_1
- SOURCE\_2
- SOURCE\_3
- TARGET

Each cell shows the **correlation coefficient (r)** between two variables, ranging from **-1 to +1**:

**+1** → Perfect positive correlation (move together)

**-1** → Perfect negative correlation (move in opposite directions)

**0** → No correlation

The color scale:

**Green** → Positive correlation

**Red** → Negative correlation

**Yellow/Orange** → Weak or moderate correlation

### Observations

Relationship	Correlation	Interpretation
SOURCE_1 vs SOURCE_2	0.21	Weak positive correlation
SOURCE_1 vs SOURCE_3	0.19	Weak positive correlation
SOURCE_2 vs SOURCE_3	0.11	Very weak positive correlation
SOURCE_1 vs TARGET	-0.16	Weak negative correlation
SOURCE_2 vs TARGET	-0.16	Weak negative correlation
SOURCE_3 vs TARGET	-0.18	Weak negative correlation

### Discussion

**Sources are weakly correlated with each other**, meaning they capture somewhat different information — which is good if they are input features for a model (less redundancy).

**All sources have weak negative correlations with the target.**

- This suggests that as values of SOURCE\_1, SOURCE\_2, or SOURCE\_3 increase, the TARGET tends to slightly decrease.

- However, since correlations are small (around -0.16 to -0.18), these relationships are **not strong**.

**None of the correlations are high**, so **multicollinearity is low**, and the model can safely include all sources without risk of feature redundancy.

## Summary

**Low inter-feature correlation** → Good for modeling.

**Weak negative correlation with target** → Features may have limited direct predictive power.

Further analysis (e.g., feature importance or regression) is needed to confirm their contribution.