**Task:** A4 Data Preparation
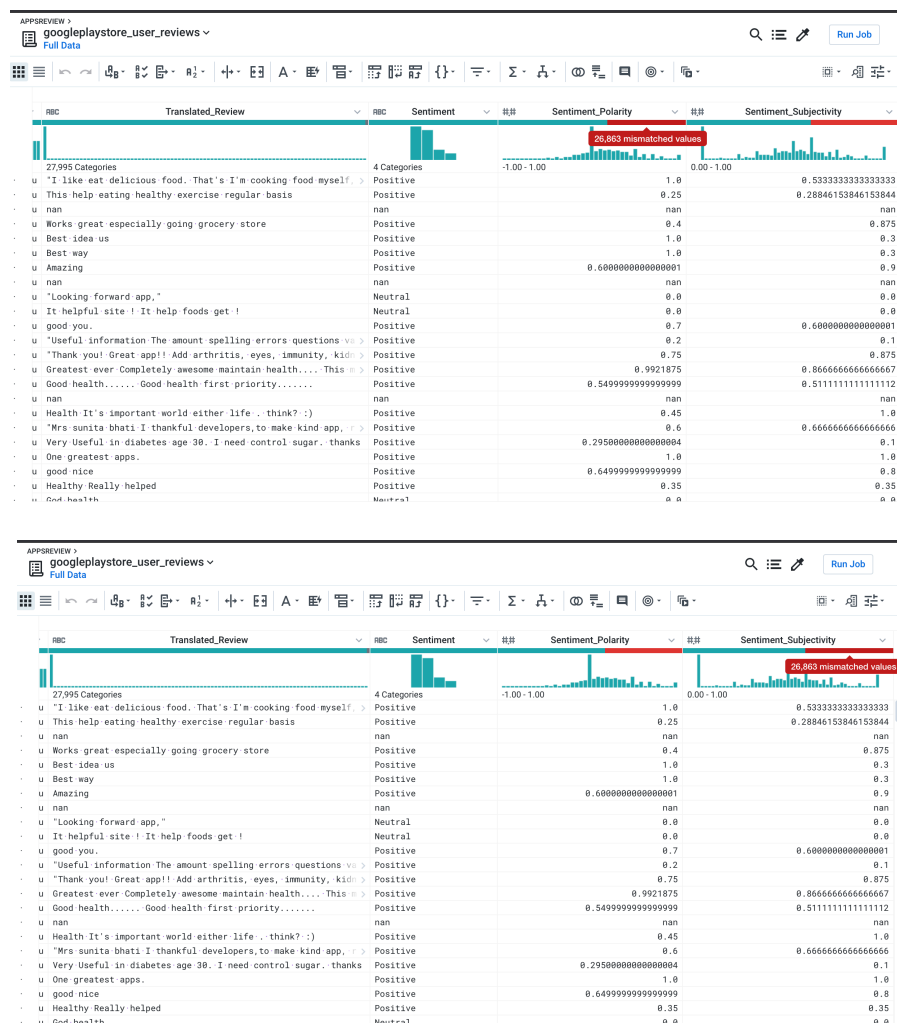**Name:** Edgardo Andres Panza Penso
**Date:** 10/Jan/2020

In the git repository you can find:

- Folder with datasets that i use for testing TRIFACTA
- LinksTask4.txt file with the links for the data cleaning training
- Readme.md
- A4_DataPreparation.pdf

https://github.com/edpape007/DataPreparationExercise

**TRIFACTA EXERCISE**

- First flow:

## Suggestions

### Delete rows

with mismatched values in Sentiment_Polarity

Edit  Add

### Keep rows

with mismatched values in Sentiment_Polarity

### Create a new column

flag mismatched values in Sentiment_Polarity

### Set

mismatched values in NULL()

mismatched values in 0

---

**New Step**  Recipe

☐  ⋯  ⚙

1  **Delete rows where**
ISMISMATCHED(Sentiment_Polarity,
['Float'])

2  **Delete rows where**
ISMISMATCHED(Sentiment_Subjectivity,
['Float'])

---

AppsReview › googleplaystore_user_reviews
**Job 112667**
Finished Today at 7:19 PM

**Download results**  ⋯

Overview    Output Destinations    Profile    Dependencies

### Completed stages

✓ **Transform**
Completed Today at 7:17 PM, started Today at 7:16 PM • Ran for 1 min

Environment    **Spark**

View steps and dependencies

✓ **Profile**
Completed Today at 7:19 PM, started Today at 7:17 PM • Ran for 1 min

● > 99% valid values    ● 0% mismatching values    ● < 1% missing values

View profile

✓ **Publish**
Completed Today at 7:17 PM, started Today at 7:17 PM • Ran for <1 sec

Activity

📄 googleplaystore_user_reviews.csv    ✓ Completed

View all

### Job summary

| | |
|---|---|
| Job ID | 112667 |
| Job status | Completed |
| Flow | AppsReview |
| Output | googleplaystore_user_reviews |

### Execution summary

| | |
|---|---|
| Job type | Manual |
| User | Edgardo Panza |
| Start time | January 10th 2020, 7:16 pm |
| Finish time | January 10th 2020, 7:19 pm |
| Last update | January 10th 2020, 7:19 pm |
| Duration | 3 minutes |

- Second Flow
  - In this one I identify a pattern and create new column.
  - Delete the rows with missing values
  - Format the date column

🔍 ☰ ✎   Run Job

| keywords ⌄ | ABC original_language ⌄ | ABC original_title ⌄ |
|---|---|---|
| | 37 Categories | 4,801 Categories |
| ""culture·clash""}, {""id"": 296 › | en | Avatar |
| ocean""}, {""id"": 726, ""name"" › | en | Pirates·of·the·Caribbean:·At·World's·End |
| spy""}, {""id"": 818, ""name"" › | en | Spectre |
| dc·comics""}, {""id"": 853, ""n › | en | The·Dark·Knight·Rises |
| based·on·novel""}, {""id"": 839 › | en | John·Carter |
| dual·identity""}, {""id"": 1452 › | en | Spider-Man·3 |
| hostage""}, {""id"": 2343, ""n › | en | Tangled |
| ""marvel·comic""}, {""id"": 9663 › | en | Avengers:·Age·of·Ultron |

New Step    Recipe    ✕

☐ ...                                    ⚙

1  Set homepage to IFMISSING($col, NULL())

2  Create column1 from MATCHES([genres], `{start}("\[{""{lower}{2}"": {digit}{2}, ""{lower}{4}"": ""{upper}{lower}+""{any}+\]"){end}`)

3  Change date format of release_date to M/d/yyyy

---

▶ Movies › tmdb_5000_movies
**Job 112670**
Finished Today at 7:30 PM

**Download results**   ...

**Overview**   Output Destinations   Profile   Dependencies

## Completed stages

### ✅ Transform
Completed Today at 7:29 PM, started Today at 7:28 PM • Ran for 1 min

Environment   Spark

View steps and dependencies

### ✅ Profile
Completed Today at 7:30 PM, started Today at 7:29 PM • Ran for 1 min

● 96% valid values   ● 0.1% mismatching values   ● 4% missing values

View profile

### ✅ Publish
Completed Today at 7:29 PM, started Today at 7:29 PM • Ran for <1 sec

Activity

| 📄 tmdb_5000_movies.csv | ✅ Completed |
|---|---|

View all

## Job summary

| | |
|---|---|
| Job ID | 112670 |
| Job status | Completed |
| Flow | Movies |
| Output | tmdb_5000_movies |

## Execution summary

| | |
|---|---|
| Job type | Manual |
| User | Edgardo Panza |
| Start time | January 10th 2020, 7:28 pm |
| Finish time | January 10th 2020, 7:30 pm |
| Last update | January 10th 2020, 7:30 pm |
| Duration | 2 minutes |

---

▶ Movies › tmdb_5000_movies
**Job 112670**
Finished Today at 7:30 PM

**Download results**   ...

Overview   **Output Destinations**   Profile   Dependencies

| Name | Location | Status |
|---|---|---|
| 📄 tmdb_5000_movies.csv | ⬇ s3://trifacta-saas-prod/trifacta-saas/17988/18201/queryResults/edgardo.panza@gmail.com/.trifacta/87… | ✅ Completed |