



# AMAZON PROJECT

Team “E”

Raúl Alonzo

Victor González

Roberto Pérez

# Background

This dataset consists of reviews of fine foods from amazon. The data span a period of more than 10 years, including all ~500,000 reviews up to October 2012. Reviews include product and user information, ratings, and a plain text review. It also includes reviews from all other Amazon categories.

Data includes:

Reviews from Oct 1999 - Oct 2012

568,454 reviews

256,059 users

74,258 products

260 users with > 50 reviews

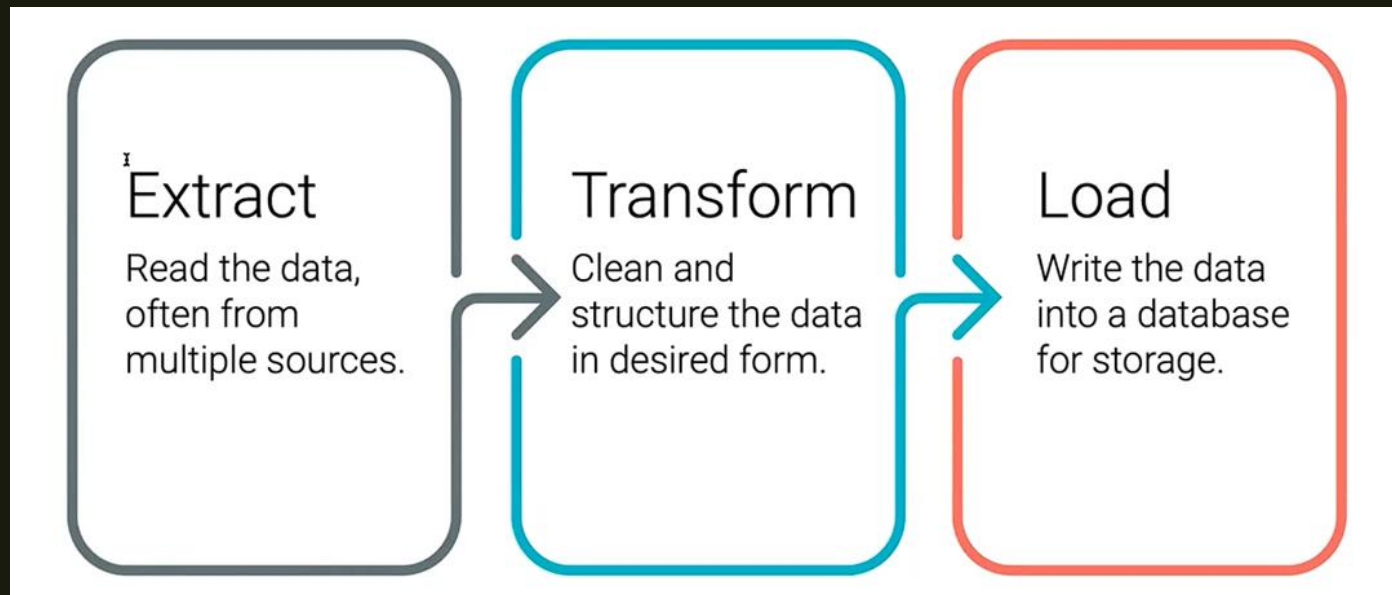
wordcloud

# Project

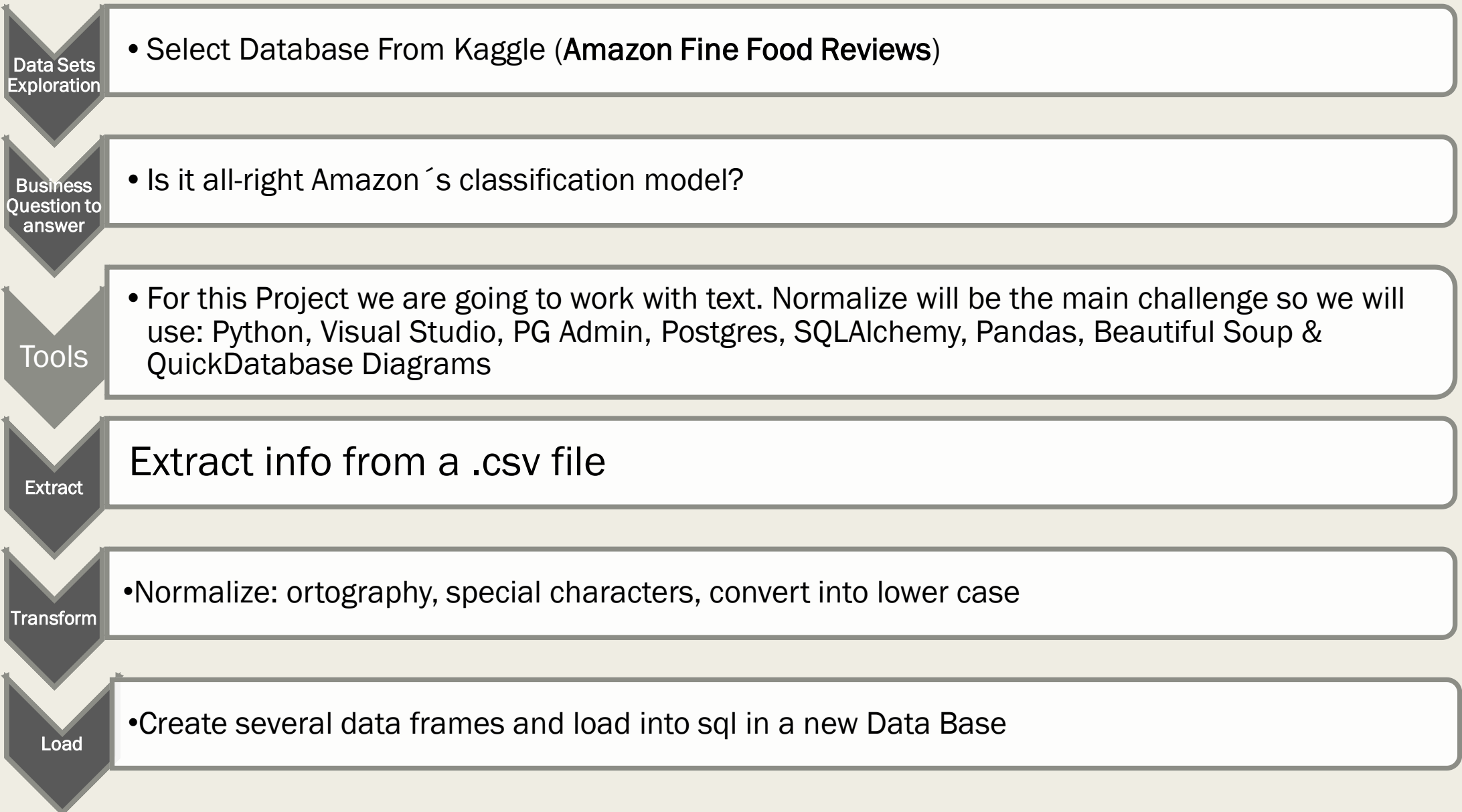
Our Project consist on reviewing Amazon ´s satisfaction method. Actually Amazon does a posotive/negative qualification and we believe ther must be a more wide range n order to take actions and get a Better feeling about customer service.

.

# PROCESS



# Workflow



# Extract data From a Kaggle CSV file

< Amazon\_rev.csv (269.3 MB)

#	Id	ProductId	UserId	ProfileName	HelpfulnessNum...	HelpfulnessDen...	Score	Time	Summary
		72005 unique values	243414 unique values	208273 unique values			positive 84% negative 16%		271924 unique values
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	positive	1303862400	Good Quality Dog Food
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	negative	1346976000	Not as Advertised
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	positive	1219017600	"Delight" says it all
3	4	B000UA0QIQ	A395B0RC6FGVXV	Karl	3	3	negative	1307923200	Cough Medicine
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	positive	1350777600	Great taffy
5	6	B006K2ZZ7K	ADT0SRK1MG0EU	Twoapennything	0	0	positive	1342051200	Nice Taffy
6	7	B006K2ZZ7K	A1SP2KVKFXRU1	David C. Sullivan	0	0	positive	1340150400	Great! Just as good as the expensive brands!
7	8	B006K2ZZ7K	A3JRGQVEQN31IQ	Pamela G. Williams	0	0	positive	1336003200	Wonderful, tasty taffy
8	9	B000E7L2R4	A1MZY09TZK0BBI	R. James	1	1	positive	1322006400	Yay Barley
9	10	B00171APVA	A21BT40VZCCYT4	Carol A. Reed	0	0	positive	1351209600	Healthy Dog Food
10	11	B0001PB9FE	A3HDK070W0QNK4	Canadian Fan	1	1	positive	1107820800	The Best Hot Sauce in the World
11	12	B0009XLVG0	A2725IB4YY9JEB	A Poeng "SparkyGoHome"	4	4	positive	1282867200	My cats LOVE this "diet" food better than their regular food
12	13	B0009XLVG0	A327PCT23YH90	LT	1	1	negative	1339545600	My Cats Are Not Fans of the New Food
13	14	B001GVISJM	A18ECVX2RJ7HUE	willie "roadie"	2	2	positive	1288915200	fresh and greasy!
14	15	B001GVISJM	A2MUGFV2TDQ47K	Lynrie "Oh HELL no"	4	5	positive	1268352000	Strawberry Twizzlers - Yummy
15	16	B001GVISJM	A1CZX3CP8IKQIJ	Brian A. Lee	4	5	positive	1262044800	Lots of twizzlers, just what you expect.

Data set column titles

Data set detail information

# Transform Data

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 from bs4 import BeautifulSoup as bs
4 import re
```

```
In [2]: 1 # List of unpleasant words to check
2 unpleasant = ['lousy','disappointed','discouraged','ashamed','powerless','diminished','guilty','dissatisfied',
3 'miserable','detestable','repugnant','despicable','disgusting','abominable','terrible','in despair',
4 'sulky','bad','upset','doubtful','uncertain','indecisive','perplexed','embarrassed','hesitant',
5 'shy','stupefied','disillusioned','unbelieving','skeptical','distrustful','misgiving','lost',
6 'unsure','uneasy','pessimistic','tense','incapable','alone','paralyzed','fatigued','useless',
7 'inferior','vulnerable','empty','forced','hesitant','despair','frustrated','distressed','woeful',
8 'pathetic','tragic','in a stew','dominated','irritated','enraged','hostile','insulting',
9 'annoyed','upset','hateful','offensive','bitten','aggressive','resentful','inflamed','provoked',
10 'incensed','infuriated','cross','worked up','boiling','fuming','fearful','terrified','suspicious',
11 'anxious','alarmed','panic','nervous','scared','worried','frightened','timid','shaky','restless',
12 'doubtful','threatened','cowardly','quaking','wary','crushed','tormented','deprived','pained',
13 'tortured','dejected','rejected','injured','offended','afflicted','aching','victimized',
14 'heartbroken','agonized','appalled','humiliated','wronged','alienated','tearful','sorrowful',
15 'pained','grief','anguish','desolate','desperate','pessimistic','unhappy','lonely','grieved',
16 'mournful','dismayed','insensitive','dull','nonchalant','neutral','reserved','weary','bored',
17 'preoccupied','cold','disinterested','lifeless','never'];
```

```
In [3]: 1 # List of pleasant words to check
2 pleasant = ['understanding','confident','reliable','easy','amazed','free','sympathetic','interested','satisfied',
3 'receptive','accepting','kind','great','joyous','lucky','fortunate','delighted','overjoyed','gleeful','thankful',
4 'festive','ecstatic','glad','cheerful','elated','jubilant','playful','courageous','energetic','liberated','opti',
5 'impulsive','free','animated','spirited','thrilled','wonderful','calm','peaceful','at ease','comfortable','pleas',
6 'clever','surprised','content','quiet','certain','relaxed','serene','reassured','loving','considerate','affecti',
7 'tender','devoted','attracted','passionate','admiration','warm','touched','close','comforted','loved','concerned',
8 'intrigued','absorbed','inquisitive','engrossed','curious','drawn toward','eager','keen','earnest','intent','ins',
9 'enthusiastic','bold','brave','daring',
10 'optimistic','impulsive','free','sure','certain','rebellious','unique','dynamic','tenacious','hardy','secure',
11 'confident','challenged','love'];
```

```
In [6]: 1 # Read csv
2 path = 'Resources/Amazon_rev.csv'
```

```
In [7]: 1 amazon_foods_df = pd.read_csv(path)
```

```
In [8]: 1 # Drop first column because it creates junk info
2 amazon_foods_df = amazon_foods_df.drop(columns=["Unnamed: 0"])
```

```
In [9]: 1 # Copy of dataframe
2 amazon_foods_modified_df = amazon_foods_df
```

```
In [10]: 1 # Add a regular expression to look for all the symbols
2 regex = "[.\\\"\\',!@#%$^&*()_\\-+=?:;|/!]"
```

```
In [11]: 1 # Text replacement for all rules in the regex
2 text_replace = ''
```

```
In [12]: 1 # Create first column with text in lower case and no special characters
2 amazon_foods_modified_df['Lowercase Text'] = ''
```

```
In [15]: 1 # Split data frame to run demo version
2 amazon_foods_modified_df = amazon_foods_df.iloc[:100]
```

```
In [16]: 1 %%time
2 # Populate new column column
3 amazon_foods_modified_df['Lowercase Text'] = amazon_foods_df.apply(
4     lambda row: (
5         re.sub(
6             regex,
7             text_replace,
8             row[9].lower()
9         )
10     ),
11     axis=1)
```

Wall time: 17.5 s

C:\Users\victo\anaconda3\lib\site-packages\ipykernel\_launcher.py:10: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)  
# Remove the CWD from sys.path while we load stuff.

```
In [17]: 1 # Create second column and clear br html tags
2 amazon_foods_modified_df['Lowercase Text Clean'] = ''
```

# Transform Data

```
In [18]: 1 # Define function to extract with beautiful soup br and just return text property
2 def extract_br_tags(soup):
3     for e in soup.findAll('br'):
4         e.extract()
5     if soup.find('p'):
6         return soup.find('p').text
7     elif soup.find('span'):
8         return soup.find('span').text
9     elif soup.find('a'):
10        return soup.find('a').text
11    return soup.find('body').text
```

```
In [19]: 1 # Define function to count positive and negative words based on pleasant and unpleasant Lists
2 def add_positive_negative(text, is_positive):
3     text_list = text.split()
4     count = 0
5     for word in text_list:
6         if is_positive:
7             for pleasant_word in pleasant:
8                 if word == pleasant_word:
9                     count+=1
10        else:
11            for unpleasant_word in unpleasant:
12                if word == unpleasant_word:
13                    count+=1
14    return count
15
16
```

```
In [20]: 1 %%time
2 amazon_foods_modified_df['Lowercase Text Clean'] = amazon_foods_modified_df.apply(
3     lambda row: (
4         extract_br_tags(bs(row[10], 'xml'))
5     ),
6     axis=1)
```

Wall time: 94.9 ms

```
In [21]: 1 # Create column for positive count
2 amazon_foods_modified_df['Positive Count'] = ''
```

C:\Users\victo\anaconda3\lib\site-packages\ipykernel\_launcher.py:2: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
In [22]: 1 %%time
2 amazon_foods_modified_df['Positive Count'] = amazon_foods_modified_df.apply(
3     lambda row: (
4         add_positive_negative(row[11], True)
5     ),
6     axis=1)
```

Wall time: 36.6 ms

C:\Users\victo\anaconda3\lib\site-packages\ipykernel\_launcher.py:5: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
In [23]: 1 # Create column for negative count
2 amazon_foods_modified_df['Negative Count'] = ''
```

C:\Users\victo\anaconda3\lib\site-packages\ipykernel\_launcher.py:2: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
In [24]: 1 %%time
2 amazon_foods_modified_df['Negative Count'] = amazon_foods_modified_df.apply(
3     lambda row: (
4         add_positive_negative(row[11], False)
5     ),
6     axis=1)
```

Wall time: 46 ms

# Transform Data

```
In [25]: 1 amazon_foods_modified_df
```

```
Out[25]:
```

Userid	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text	Lowercase Text	Lowercase Text Clean	Positive Count	Negative Count
AUHU8GW	delmartian	1	1	positive	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...	i have bought several of the vitality canned d...	i have bought several of the vitality canned d...	0	0
SZCVE5NK	dll pa	0	0	negative	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...	product arrived labeled as jumbo salted peanut...	product arrived labeled as jumbo salted peanut...	1	0
WJIXXAIN	Natalia Corres "Natalia Corres"	1	1	positive	1219017600	"Delight" says it all	This is a confection that has been around a fe...	this is a confection that has been around a fe...	this is a confection that has been around a fe...	0	0
3C6FGVXV	Karl	3	3	negative	1307923200	Cough Medicine	If you are looking for the secret ingredient I...	if you are looking for the secret ingredient I...	if you are looking for the secret ingredient I...	0	0
2LF8GW1T	Michael D. Bigham "M. Wassir"	0	0	positive	1350777600	Great taffy	Great taffy at a great price. There was a wid...	great taffy at a great price there was a wide...	great taffy at a great price there was a wide...	2	0
...	...	...	...	...	...	...	...	...	...	...	...
v6MXO9B0	pionex1796	0	0	positive	1326412800	Loved these Tartlets	What a nice alternative to an apple pie. Love ...	what a nice alternative to an apple pie love t...	what a nice alternative to an apple pie love t...	3	0
521O626G	Rachel Westendorf	0	0	positive	1308700800	The best	I like Creme Brulee. I loved that these were s...	i like creme brulee i loved that these were so...	i like creme brulee i loved that these were so...	3	0
5UIR6F8KB	Adam E. Smith	2	2	positive	1295308800	Wasting Vinegar on a Cucumber is a Shame!	I first bought pickled asparagus at an Amish m...	i first bought pickled asparagus at an amish m...	i first bought pickled asparagus at an amish m...	1	0

```
In [26]: 1 # Drop columns not needed for modified csv
2 amazon_foods_modified_df = amazon_foods_modified_df.drop(columns=['Lowercase Text'])
```

```
In [27]: 1 amazon_foods_modified_df.to_csv('Amazon_rev_modified.csv')
```

```
In [ ]: 1
```



# Load Data

```
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3
```

```
CREATE DATABASE CONNECTION

In [129]: 1 connection_string = "postgres:B00tc@mp@localhost:5432/Amazon"

In [130]: 1 # confirm tables
          2 engine = create_engine(f"postgresql+psycopg2://{connection_string}")

LOAD DATAFRAME INTO DATABASE

In [131]: 1 data_df.to_sql(name="data", con=engine, if_exists="append", index=False)

CREATE & LOAD 'USER' DATAFRAME INTO DATABASE

In [132]: 1 # Create DF for table 'user'
          2 user_df_table = data_df[['user_id', 'profile_name']]

In [133]: 1 user_df_table
          ...

In [134]: 1 # Identify Duplicated
          2 user_df_table = user_df_table.assign(repeat=lambda r:r['user_id'].duplicated())

In [135]: 1 # Query
          2 user_df_table = user_df_table.query('repeat == False')

In [136]: 1 # Drop repeated
          2 user_df_table = user_df_table.drop(columns=['repeat'])

In [137]: 1 # Print
          2 user_df_table
          ...

In [138]: 1 # Send data to SQL
          2 user_df_table.to_sql(name='user', con=engine, if_exists="append", index=False)
```

```
CREATE & LOAD 'PRODUCT' DATAFRAME INTO DATABASE

In [139]: 1 # Create DF for table 'product'
          2 product_df_table = data_df[['product_id']]

In [140]: 1 # Print
          2 product_df_table
          ...

In [141]: 1 # Identify Duplicated
          2 product_df_table = product_df_table.assign(repeat=lambda r:r['product_id'].duplicated())

In [142]: 1 # Query
          2 product_df_table = product_df_table.query('repeat == False')

In [143]: 1 # Drop repeated
          2 product_df_table = product_df_table.drop(columns=['repeat'])

In [144]: 1 product_df_table
          ...

In [145]: 1 # Send data to SQL
          2 product_df_table.to_sql(name='product', con=engine, if_exists="append", index=False)

CREATE & LOAD 'REVIEW' DATAFRAME INTO DATABASE

In [146]: 1 # Create DF for table 'review'
          2 review_df_table = data_df[['review_id', 'score', 'summary', 'text']]

In [147]: 1 # Print
          2 review_df_table
          ...

In [148]: 1 # Identify Duplicated
          2 review_df_table = review_df_table.assign(repeat=lambda r:r['review_id'].duplicated())

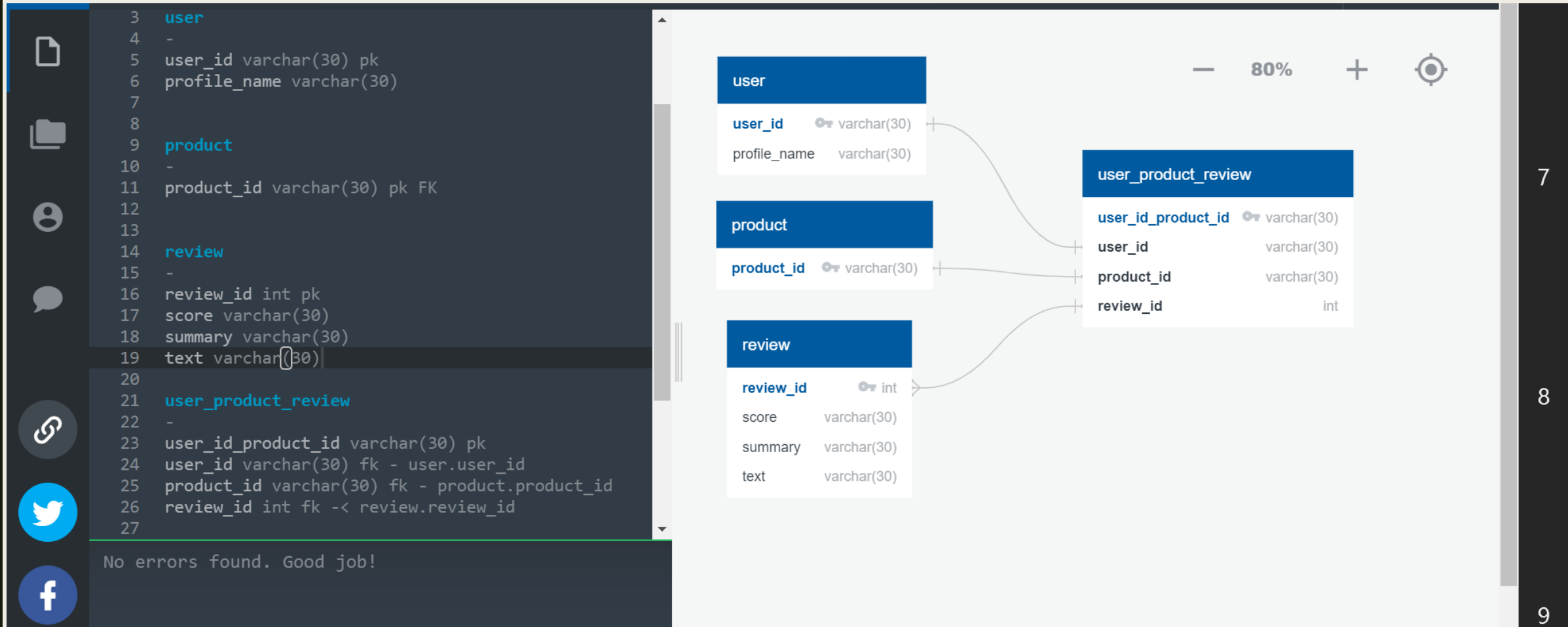
In [149]: 1 # Query
          2 review_df_table = review_df_table.query('repeat == False')

In [150]: 1 # Drop repeated
          2 review_df_table = review_df_table.drop(columns=['repeat'])

In [151]: 1 # Print
          2 review_df_table
          ...

In [152]: 1 # Send data to SQL
          2 review_df_table.to_sql(name='review', con=engine, if_exists="append", index=False)
```

# Relational Diagram



# Amazon new Data Base

**pgAdmin** File Object Tools Help

Browser

- Servers (1)
  - PostgreSQL 11
    - Databases (6)
      - Amazon
        - Casts
        - Catalogs
        - Event Triggers
        - Extensions
        - Foreign Data Wrappers
        - Languages
        - Schemas (1)
          - public
            - Collations
            - Domains
            - FTS Configurations
            - FTS Dictionaries
            - FTS Parsers
            - FTS Templates
            - Foreign Tables
            - Functions
            - Materialized Views
            - Procedures
            - Sequences
            - Tables (4)
              - product
              - review
              - user
              - user\_product\_review
            - Trigger Functions
            - Types
            - Views
      - customer\_db
      - emoji
      - employee\_db

Dashboard Properties SQL Statistics Dependencies Dependents Amazon/postgres@PostgreSQL 11 \*

Amazon/postgres@PostgreSQL 11

Query Editor Query History

```
1 -- Exported from QuickDBD: https://www.quickdatabasediagrams.com/
2 -- Link to schema: https://app.quickdatabasediagrams.com/#/d/JuWDih
3 -- NOTE! If you have used non-SQL datatypes in your design, you will have to change these here.
4
5 -- develop model
6
7 CREATE TABLE "user" (
8     "user_id" varchar(30) NOT NULL,
9     "profile_name" varchar(30) NOT NULL,
10    CONSTRAINT "pk_user" PRIMARY KEY (
11        "user_id"
12    )
13 );
14
15 CREATE TABLE "product" (
16     "product_id" varchar(30) NOT NULL,
17    CONSTRAINT "pk_product" PRIMARY KEY (
18        "product_id"
19    )
20 );
21
22 CREATE TABLE "review" (
23     "review_id" int NOT NULL,
24     "score" varchar(30) NOT NULL,
25     "summary" varchar(30) NOT NULL,
26     "text" varchar(30) NOT NULL,
```

Data Output Explain Messages Notifications

ALTER TABLE

Query returned successfully in 74 msec.