MS2 Submission

Mariusz Derezinski-Choo, Nicholas Aldana, Robert Chen, Justin Du, Karoline Xiong

Data Procurement Plan:

We have found a dataset of 30,000 amazon items from kaggle, available at
https://www.kaggle.com/promptcloud/amazon-product-details. To procure the data, we are
writing a python script that will load the dataset using the pandas package, and validate each
row of the dataset to make sure it is complete. For example, by visually inspecting the dataset
we have noticed that some rows do not have a properly formatted price, or may be missing an
image url. We will validate each row against a set of conditions to ensure that the data entry is
complete, and delete all the rows that are incomplete. We will also have to dynamically a table
of storefront rows based on the sellers listed in the dataset. Each item in our dataset has
anywhere between 1-3 sellers listed, so we will keep track of all the sellers we encounter and
generate emails, passwords, and balances for each of them. Then, we can use sqlalchemy to
insert all of these rows into the database. A rough outline of the script is shown below.

```python
data_procure.py
1    import pandas as pd
2    import sqlite3
3    conn = sqlite3.connect('./database.db')
4
5    dataset = pd.read_csv("item_data.csv")
6    for index, row in dataset.iterrows():
7        if '$' not in str(row['Price']):
8            dataset.drop(index,inplace=True)
9        # TODO: add other validation conditions
10
11       # TODO: generate seller accounts based on row
12
13   for index, row in dataset.iterrows():
14       c.execute("INSERT into item values (?,?,?,?)",row['Uniq id'],row['Category'],row['Title'],row['Description']))
15       #TODO: insert rows into seller table
16
17       #TODO: insert rows into Listing table
18
19   conn.commit()
20
21
22   dataset.to_csv("cleaned_data.csv")
23
```

Changes from MS2:
- User table was split into Buyer and Seller tables. The schema for each remained largely the same except sellers do not have a photo.
- The selling table was renamed to Listing to provide more clarity.