

# **Task 4: Data Munging**

## **Here is your task**

### **Part 1: Get the data**

First, you need to get your hands on the relevant data. The shipping department has been kind enough to provide you with a repository containing all of their spreadsheets, as well as a copy of the sqlite database. First, fork and clone the repository at: <https://github.com/theforage/forage-walmart-task-4>

### **Part 2: Populate the database**

Your task is to insert all of the data contained in the provided spreadsheets into the SQLite database. You will write a Python script which:

- Reads each row from the spreadsheets.
- Extracts the relevant data.
- Munges it into a format that fits the database schema.
- Inserts the data into the database.

Spreadsheet 0 is self contained and can simply be inserted into the database, but spreadsheets 1 and 2 are dependent on one another. Spreadsheet 1 contains a single product per row, you will need to combine each row based on its shipping identifier, determine the quantity of goods in the shipment, and add a new row to the database for each product in the shipment. The origin and destination for each shipment in spreadsheet 1 are contained in spreadsheet 2. You may assume that all the given data is valid - product names are always spelled the same way, quantities are positive, etc.

## Solution :

```
1 - import pandas as pd
2 - import sqlite3
3 - # Step 1: Connect to SQLite database
4 - DB_FILE = 'shipping.db'
5 - conn = sqlite3.connect(DB_FILE)
6 - cursor = conn.cursor()
7 - # Step 2: Read all spreadsheets
8 - spreadsheet0 = pd.read_csv('spreadsheet0.csv')
9 - spreadsheet1 = pd.read_csv('spreadsheet1.csv')
10 - spreadsheet2 = pd.read_csv('spreadsheet2.csv')
11 - # Step 3: Process Spreadsheet 0 (self-contained)
12 - for _, row in spreadsheet0.iterrows():
13 -     cursor.execute('''
14 -         INSERT OR IGNORE INTO Shipments
15 -             (ShipmentID, ProductName, Quantity, Origin, Destination)
16 -             VALUES (?, ?, ?, ?, ?)
17 - ''', (
18 -     row['ShippingID'],
19 -     row['ProductName'],
20 -     row['Quantity'],
21 -     row['Origin'],
22 -     row['Destination']
23 - ))
24 - # Step 4: Process Spreadsheets 1 & 2 (dependent)
25 - merged_shipments = pd.merge(
26 -     spreadsheet1,
27 -     spreadsheet2,
28 -     on='ShippingID',
29 -     suffixes=('_product', '_location'))
30 -     row['ProductName'],
31 -     row['Quantity'],
32 -     row['Origin'],
33 -     row['Destination']
34 - )
35 - # Step 4: Process Spreadsheets 1 & 2 (dependent)
36 - merged_shipments = pd.merge(
37 -     spreadsheet1,
38 -     spreadsheet2,
39 -     on='ShippingID',
40 -     how='left'
41 - )
42 - for _, row in merged_shipments.iterrows():
43 -     cursor.execute('''
44 -         INSERT OR IGNORE INTO Shipments
45 -             (ShipmentID, ProductName, Quantity, Origin, Destination)
46 -             VALUES (?, ?, ?, ?, ?)
47 - ''', (
48 -     row['ShippingID'],
49 -     row['ProductName'],
50 -     row['Quantity'],
51 -     row['Origin'],
52 -     row['Destination']
53 - ))
54 - # Step 5: Commit and close
55 - conn.commit()
56 - conn.close()
57 - print('All spreadsheets have been successfully loaded into the database!')
```