

E-commerce Website

CSE4/560 Data Models and Query Language Semester Project
Milestone-2 Report

Aman Prakash
Engineering and Applied Sciences
M.S in Data Science
University at buffalo
email address

Serath Chandra Nutakki
Engineering and Applied Sciences
M.S in Data Science
University at buffalo
serathch@buffalo.edu

Prashant Upadhyay
Engineering and Applied Sciences
M.S in Data Science
University at buffalo
pupadhya@buffalo.edu

Abstract— In the present environment of business shopping manually at the store is a hassle, especially with multiple products, discounts, and variations. Considering all the parameters we would like to build an e-commerce website to serve the customer with efficiency, time management and provide insights for which a database is required.

I. INTRODUCTION

E-commerce is quickly becoming a common business model. Web sites with the capability for conducting commercial transactions over the internet are being implemented by an increasing number of companies. Shopping via the internet may be described as a simple norm.

In this existing modern world, there is a need for modern shopping which can be satisfied by online shopping through an e-commerce website mainly due to its simplicity and easiest way to communicate and transition between buyer and seller.

The main objective of our project is to implement an e-commerce store website where multiple products like clothes, electronics, and many more can access through the internet from our own comfort by connecting to a database.

II. IMPLEMENTATION

We are going to implement an e-commerce website, where we want to store multiple users, products, discount coupons, etc. and for that, we will need a database. An e-commerce website has multiple products, and the end-user is not going to buy any product without having a view of the item, so we need to insert images of the product into our database. There will be multiple users buying the same items and this creates data redundancy which will hamper the data processing.

Taking the above points in mind we need a database instead of an excel file because:

Data Purpose

What type of data are we collecting? Excel files are great for text and numeric values in relatively low volume, but databases can not only handle numeric and text values but can easily handle other types of information, such as documents and images. Moreover, in our project, we will store multiple images which are not feasible for the database to capture.

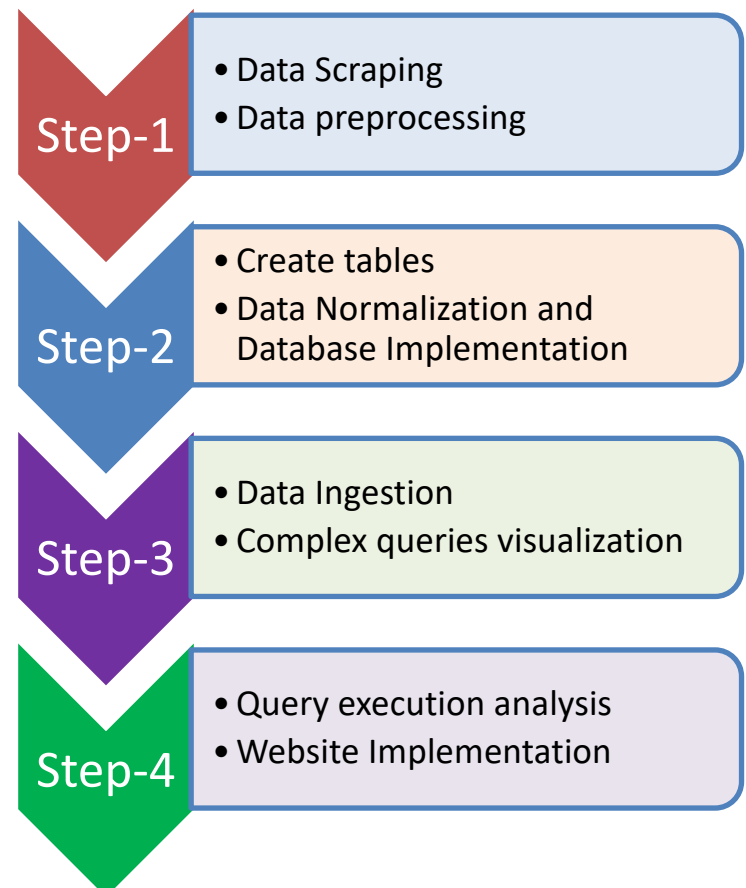
Data Volume

Spreadsheets have record limitations whereas databases do not. Spreadsheets require a large amount of hard-drive space for data storage than databases. When a spreadsheet has many fields, the spreadsheet can be hard to read, and finding specific data can be challenging. Relational databases on the other hand uses querying tools to overcome these issues.

Redundancy

The database structure avoids data redundancy. Since the data in different database tables are linked, there is little or no duplication of source data. On our project, one user can have multiple purchases of single items and of different items with the same quantity or different, so this information of duplicate items will make the excel sheet redundant.

Design process:



Data Scraping:

The process of extracting the data from the website into a well-maintained spreadsheet that can be used as our data for further implementation is called Data Scraping.

It allows us to scrap specific data like product details, inputs, user details, order details, categories, and many more. This way we can export large product data in CSV for our desired output.

In this project, we extracted our essential data from the Amazon website through various sources and toolkits.

Data preprocessing:

Our sophisticated databases and websites cannot be implemented on poor data because they can yield insignificant and inefficient results. So, in order to avoid that we need to perfectly clean our data with special techniques, and this process is known as data preprocessing.

In this process, we have taken care of all the missing values product details, and order details. Filling the unknown prices with average mean values.

The sample representation of our dataset:

1	orderId	userid	email	firstName	lastName	address1	address2	zipcode	city	state	country	phone	productid	productna	price	description	stock	image	categoryid	categoryn	discountid	discountperce
2	20889	954	vanda_tenVanda	Tender	7	Shelford/Ct		96335	Kimberley	BC	Canada	250-605-3	7038	Anico Wel	19.78	Male sure	26	https://m	39	Toys & Ga	NULL	NULL
3	20891	280	elmer@noElmer	Reddon	53	Euclid Ave		64624	Fortes	NSW	Australia	02-3075-4	7190	Funko Adv	97.86	Male sure	83	https://m	618	Home & Ki	NULL	NULL
4	20892	628	theresa.daiTheresa	Davirro	863	E Macdonell Rd		64089	Kamloops	BC	Canada	250-940-5	1745	Barbie Be	34.29	Male sure	36	https://m	39	Toys & Ga	NULL	NULL
5	20894	1686	916-770-7	Kerry	Theodoros	6915	W Main St	57231	Sacramento	Sacramento	United Sta	95807	3279	Littlest Pet Shop Cat	1	Male sure	17	https://m	34	Toys & Ga	228	2
6	70811	683	lythe_caslythe	Czapinski	4	4th St		16755	Calgary	AB	Canada	403-775-4	6688	Hekkat Ra	55.35	Genuine R	66	https://m	40	Toys & Ga	NULL	NULL
7	70813	1570	410-937-4	Irene	Eroman	2853	S Central Expy	21534	Glen Burni	Anne Arvon	United Sta	21061	5596	Funko Pop	82.91	Male sure	56	https://m	3	Toys & Ga	384	64
8	70820	1065	sol@gmailSol	Cowser	6448	Tillard St		98041	Conidrou	South Vi	United Ki	01412-528	1991	Prezmann	51.23	Male sure	79	https://m	62	Toys & Ga	348	81
9	70823	383	alya@hotAlya	Lekhou	186	Geary Blvd #923		14483	Trughey	H NSW	Australia	02-5385-3	7708	Hect Wlhed	10.24	Male sure	31	https://m	190	Toys & Ga	NULL	NULL
10	70825	251	johanna@Johanna	Saffer	750	Lancaster Ave		77009	Campsie	NSW	Australia	02-9470-1	4326	Philips Hee	96.36	Product Da	56	https://m	255	Toys & Ga	NULL	NULL

Create tables:

After getting our data we generated tables for different features by specific scripts that we have written.

- Users – table contains all the user-specific data (email, userID, password, first and last name, zip code)
- Products – table created to store all the product details (productid, name, price, description, image, etc.)
- Orders – table created to store order details (ordered, userid, productid, amount, etc)
- Discount – table created to store discount details (discount Id, percentage, description, etc)
- Categories – contains all the product categories' specific details (name, category id, etc)
- Cart – To store the cart details that we have added (userid, product id, etc)
- Support – Table to store the support related details (supportId, user id, body, etc)

Each table was created and defined as follows:

```
try:
cur.execute('''CREATE TABLE users
(userid INTEGER PRIMARY KEY,
password TEXT,
email TEXT,
firstName TEXT,
lastName TEXT,
address1 TEXT,
address2 TEXT,
zipcode TEXT,
city TEXT,
state TEXT,
country TEXT,
phone TEXT
)''')
```

```
cur.execute('''CREATE TABLE products
(productid INTEGER PRIMARY KEY,
name TEXT,
price TEXT,
description TEXT,
image TEXT,
stock INTEGER,
categoryid INTEGER,
FOREIGN KEY(categoryid) REFERENCES categories(categoryid)
)''')

cur.execute('''CREATE TABLE cart
(userid INTEGER,
productid INTEGER,
FOREIGN KEY(userid) REFERENCES users(userid),
FOREIGN KEY(productid) REFERENCES products(productid)
)''')

cur.execute('''CREATE TABLE IF NOT EXISTS discount (
discountid integer NOT NULL,
description text NULL DEFAULT NULL,
discountPercentage integer NULL DEFAULT NULL,
PRIMARY KEY (discountid))''')

cur.execute('''CREATE TABLE IF NOT EXISTS orders (
orderid integer NOT NULL,
userid integer NULL DEFAULT NULL,
productid integer NULL DEFAULT NULL,
discountid integer NULL DEFAULT NULL,
PRIMARY KEY (orderid),
```

Data Normalization:

Currently, in our dataset, we have 1,50,000 records of the products but we observed some undesirable characteristics features and data redundancies. Therefore, this creates a need to normalize the data by forming a relationship linkage between smaller tables that are been derived from the larger tables.

Removing all the partial dependencies and nonprime attributes and transforming the data to BCNF(Boyce- Codd Normal Form) which are having working functional dependencies.

Therefore, we transformed the large unstructured data by normalizing it to a structured data form manually.

Database Implementation:

We have normalized our data and we created all the tables. So, we implement this database and establish the connection.

```
import psycopg2

def conn_create():
    conn = psycopg2.connect(
        host="localhost",
        database="ecom",
        user="postgres",
        port="8888",
        password="abc12345")

    cur = conn.cursor()

    return conn, cur
```

```
import psycopg2
from connect_db import conn_create

conn, cur = conn_create()

try:
    cur.execute('''CREATE TABLE users
        (userId INTEGER PRIMARY KEY,
        password TEXT,
        email TEXT,
        firstName TEXT,
        lastName TEXT,
        address1 TEXT,
        address2 TEXT,
        zipcode TEXT,
        city TEXT,
        state TEXT,
        country TEXT,
        phone TEXT
        )''')

    cur.execute('''CREATE TABLE categories
        (categoryId INTEGER PRIMARY KEY,
        name TEXT
        )''')

except (Exception, psycopg2.DatabaseError) as error:
    print(error)
finally:
    if conn is not None:
        conn.close()
    print('Tables Created Successfully.')
    print('Database connection closed.')
```

Above is the code snippet of the sample execution of the table created. By compiling the create_database.py it automates the process of creating the tables by establishing a connection with the database.

Data Ingestion:

The next step after creating the tables and database implementation is data ingestion. We implemented programs and scripts to import the data into our database.

Below is the sample code snippet for data ingestion of products data:

```
products_ingestion.py 1 X
C:\Users> serat > AppData > Local > Temp > Temp1_Archive (7).zip > ingestion_scripts > products_ingestion.py > ...

2 import csv
3 import random
4
5
6 def category_id_fetcher(data, category):
7     for i, item in enumerate(data):
8         if category in item:
9             # print('Index of 1 is {}'.format(i, item.index(1)))
10            return data[i][0]
11
12
13 conn, cur = conn_create()
14 cur = conn.cursor()
15
16 cur.execute('SELECT * FROM categories')
17 data = cur.fetchall()
18 print(data)
19 file = open('webscrap_data/amazon.csv')
20
21 print(type(file))
22
23 csvreader = csv.reader(file)
24 header = []
25 header = next(csvreader)
26 print(header)
27
28 insert_val = []
29 count = 1
30 for row in csvreader:
31     if row[1] != '' and row[4] != '':
32         insert_val.append((
33             "productid": count,
34             "name": row[1],
35             "price": row[7].replace('$', '') if row[7] == '' else str(random.randint(1, 100)) + '.' + str(random.randint(10, 99)),
36             "description": row[10] if row[10] != '' else "Product Description Yet To Be Updated By Seller.",
37             "stock": row[8] if row[8] != '' else random.randint(0, 100),
```

In this way, we import all the data to the database to the respective schema.

The dependencies of each table are as follows:

User:

Type	Name	Restriction
Sequence	public.id_seq	auto
Function	nextval('id_seq::regclass)	auto
Primary Key	public.users_pkey	auto
Foreign Key	public.cart.kart_userid_fkey	normal
Foreign Key	public.orders.orders_userid_fkey	normal
Foreign Key	public.support.support_userid_fkey	normal
Rule	_RETURN ON public.order_details	normal

Support:

Type	Name	Restriction
Foreign Key	public.support.support_userid_fkey	auto
Primary Key	public.support_pkey	auto

Products:

Type	Name	Restriction
Foreign Key	public.products.products_categoryid_fkey	auto
Primary Key	public.products_pkey	auto
Foreign Key	public.cart.kart_productid_fkey	normal
Foreign Key	public.orders.orders_productid_fkey	normal
Rule	_RETURN ON public.order_details	normal

Order:

Type	Name	Restriction
Foreign Key	public.orders.orders_discountid_fkey	auto
Foreign Key	public.orders.orders_productid_fkey	auto
Foreign Key	public.orders.orders_userid_fkey	auto
Primary Key	public.orders_pkey	auto
Rule	_RETURN ON public_order_details	normal

Discount:

Type	Name	Restriction
Primary Key	public.discount_pkey	auto
Foreign Key	public.orders.orders_discountid_fkey	normal
Rule	_RETURN ON public_order_details	normal

Categories:

Type	Name	Restriction
Primary Key	public.categories_pkey	auto
Foreign Key	public.products.products_categoryid_fkey	normal
Rule	_RETURN ON public_order_details	normal

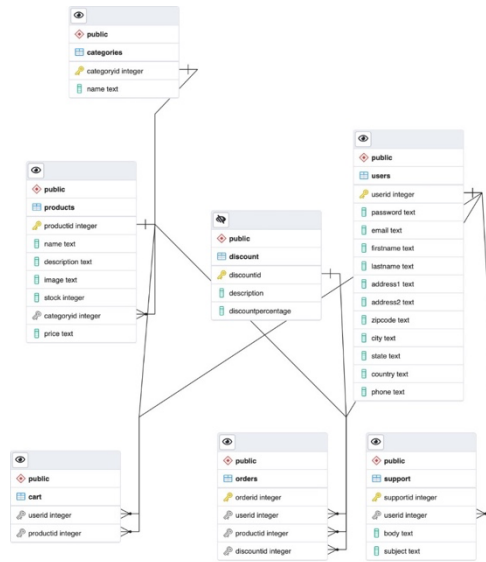
Cart:

Type	Name	Restriction
Foreign Key	public.cart.kart_productid_fkey	auto
Foreign Key	public.cart.kart_userid_fkey	auto

III. Issues Faced

We have not downloaded the data from any private repository rather we web scrapped the entire data from Amazon. We have used BeautifulSoup python package to extract the data from Amazon. It took us a lot of time to web scrape the data and clean it before implementing the database and running the query.

IV. ER Diagram



The ER diagram between all tables can be visualized

V. Query Implementation

Below are some of the queries that we implemented to get the overall scenario of the database and the relationship between all the tables.

We need to check the number of values in the dataset by connecting all the tables to each other and for that we implemented the below query which will give the total count of the number of cells in the database.

```
SELECT * FROM orders o
FULL JOIN users u ON o.userid = u.userid
FULL JOIN products p ON o.productid = p.productid
FULL JOIN discount d ON o.discountid = d.discountid
FULL JOIN categories c ON c.categoryid = p.categoryid ;
```

Output: It took 510ms to load 150006 rows.

	orderid integer	userid integer	productid integer	discountid integer	userid integer	password text	email text	firstname text	lastname text
1	20989	954	7038	[null]	954	81d9b52b04d:20036db8313ed055	vanda_tentler@tentler.org	Vanda	Tentler
2	20991	280	7198	[null]	280	81d9b52b04d:20036db8313ed055	elmer@hotmail.com	Elmer	Redlon
3	20992	629	1745	[null]	629	81d9b52b04d:20036db8313ed055	theresa.davino@aol.com	Theresa	Davino
4	20994	1666	3279	228	1666	81d9b52b04d:20036db8313ed055	916-770-7448	Kerry	Czaplinski
5	70811	683	6688	[null]	683	81d9b52b04d:20036db8313ed055	blythe_czapinski@czaplinski.com	Blythe	Czaplinski
6	70813	1570	5596	384	1570	81d9b52b04d:20036db8313ed055	410-937-4543	Ilene	Eroman
7	70820	1065	1991	348	1065	81d9b52b04d:20036db8313ed055	sol@gmail.com	Sol	Cowser
8	70823	393	7708	[null]	393	81d9b52b04d:20036db8313ed055	alysa@hotmail.com	Alysa	Lehoux
9	70825	250	4326	[null]	250	81d9b52b04d:20036db8313ed055	johanna@yahoo.com	Johanna	Saffer
10	70832	165	4148	[null]	165	81d9b52b04d:20036db8313ed055	hector.barnas@barnas.com.au	Hector	Barnas
11	70834	1970	3779	851	1970	81d9b52b04d:20036db8313ed055	703-874-4248	Malcolm	Tromblay

As our dataset is an online ecommerce website dataset. We wanted to have a look of all the customers with their orders and their details. For such we implemented view for that because view doesn't store data permanently and it also encapsulate the name of the table.

Below is the code for view in sql.

```
CREATE VIEW order_details AS
SELECT
o.orderid as orderid,
u.userid as userid,
u.email as email,
u.firstName as firstName,
u.lastName as lastname,
u.address1 as address1,
u.address2 as address2,
u.zipcode as zipcode,
u.city as city,
u.state as state,
u.country as country,
u.phone as phone,
p.productId as productId,
p.name as productname,
p.price as price,
p.description as description,
p.stock as stock,
p.image as image,
p.categoryid as categoryid,
c.name as categoryname,
d.discountid as discountid,
d.discountpercentage as discountpercentage
FROM orders o
FULL JOIN users u ON o.userid = u.userid
FULL JOIN products p ON o.productId = p.productId
FULL JOIN discount d ON o.discountid = d.discountid
FULL JOIN categories c ON c.categoryid = p.categoryid;
```

For running the above view query we have to write the below code:

```
SELECT * FROM order_details;
```

Output: It took 450 ms to load the database with the order details for all the users with their personal details.

orderid	userid	email	firstname	lastname	address1	address2	zipcode	city
1	20989	vanda_vantier@entier.org	Vanda	Tentier	7 Sheffield Ct		96335	Kimberley
2	20991	elmer@hotmail.com	Elmer	Redlon	53 Euclid Ave		64624	Forbes
3	20992	thera.davemo@aol.com	Theresa	Davemo	863 E Midwood Rd		64909	Kamloops
4	20994	916-710-7148	Kerry	Theodorov	6916 W Main St		97231	Sacramento
5	70811	683_blythe_czapinski@icazapinski.com	Blythe	Czapinski	4 S 4th St		16755	Calgary
6	70813	1570_410-937-4543	Hene	Enoman	2833 S Central Expy		21534	Glen Burnie
7	70820	1065_sot@gmail.com	Sol	Cowser	6448 Tillard St		98141	Coniborough and
8	70823	393_alysa@hotmail.com	Alysa	Lehoux	186 Geary Blvd #923		14483	Trumpley Hall
9	70825	250_johanna@yahoo.com	Johanna	Saffler	790 Lancaster Ave		77009	Campsie
10	70832	165_hector.barras@barras.com.au	Hector	Barras	62 J St #450		89585	Comberbar
11	70834	1970_703-874-4248	Malcolm	Tromblay	747 Leona Blvd		88080	Annandale

orderid	count
1	77

Output: It took 88ms to load 1 rows with all the orders of the user.

As we are going through only the users data and their orders, now let's also go through the country stats and check the country with the maximum sale and order.

```
CREATE VIEW get_country_stats AS
SELECT
country,
SUM(CAST(price AS DOUBLE PRECISION)) as total_revenue_country,
COUNT(*) AS order_count,
MODE() WITHIN GROUP (ORDER BY productname) as most_sold_product,
MODE() WITHIN GROUP (ORDER BY productId) as most_sold_product_id,
MODE() WITHIN GROUP (ORDER BY price) as most_sold_product_price,
MODE() WITHIN GROUP (ORDER BY description) as most_sold_product_details,
MODE() WITHIN GROUP (ORDER BY stock) as most_sold_product_stock,
MODE() WITHIN GROUP (ORDER BY image) as most_sold_product_image,
MODE() WITHIN GROUP (ORDER BY categoryname) as most_sold_category,
MODE() WITHIN GROUP (ORDER BY categoryid) as most_sold_category_id
FROM order_details
WHERE price > 0
GROUP BY country;
```

To run the above code, we have used the below query:

```
SELECT * from get_country_stats;
```

country	total_revenue_country	order_count	most_sold_product	most_sold_product_id	most_sold_product_price	most_sold_product_details
Australia	1910358.9099999194	37245	Barbie Dolphin Magic Sharkel Fun Friends	241	70.44	Product Description Yet To Be Updated
Canada	1906472.6699999075	37277	American Party Supplies, Multicolor	6627	88.89	Product Description Yet To Be Updated
United Kingdom	1895039.8099999158	36814	K2 Skate Alexia 84 Boa Inline Skate	863	93.91	Product Description Yet To Be Updated
United States	1897719.7499999402	36998	Btwswm NFL Pool Noodles (Pack of 3)	4314	34.84	Product Description Yet To Be Updated

Output: It took 1 sec 829ms to load 4 rows with the country stats.

VI. Modules used

Now we wanted to see the stored orders of the user in the database and for that we have used stored procedure. Stored procedure has multiple advantages like better performance, higher productivity, scalability, easy to use and secure .Below is the code for stored procedure.

```
Stored Procedure 1 -
DROP FUNCTION order_count_per_user(text);
DROP TYPE order_details_type;
CREATE TYPE order_details_type AS
(
orderid bigint
);
CREATE OR REPLACE FUNCTION order_count_per_user(email_input text)
RETURNS SETOF order_details_type
AS
$func$
BEGIN
RETURN QUERY
SELECT count(*) as order_count FROM order_details where order_details.email=email_input;
END
$func$
LANGUAGE plpgsql;
```

To run the above code, we have used the below query:

```
SELECT * FROM
order_count_per_user('amanprakash9597@gmail.com');
```

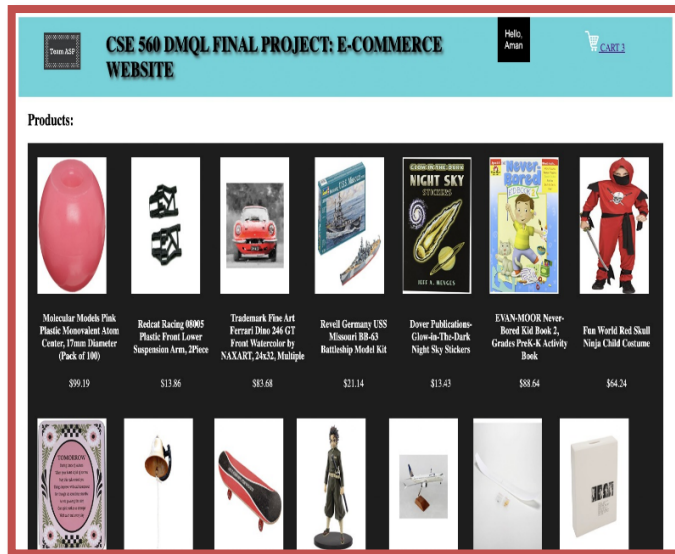
1. autopep8==1.6.0
2. click==8.1.3
3. Flask==2.1.2
4. importlib-metadata==4.11.3
5. itsdangerous==2.1.2
6. Jinja2==3.1.2
7. MarkupSafe==2.1.1
8. psycpg2==2.9.3
9. pycodestyle==2.8.0
10. toml==0.10.2
11. Werkzeug==2.1.2
12. zipp==3.8.0

VII. Website Design

Our website is designed for anyone who wants to shop a wide range of products, and different categories by login into our website with their credentials.

UI:

The homepage of our website linked to our database consists of multiple products that can be accessed by the user.



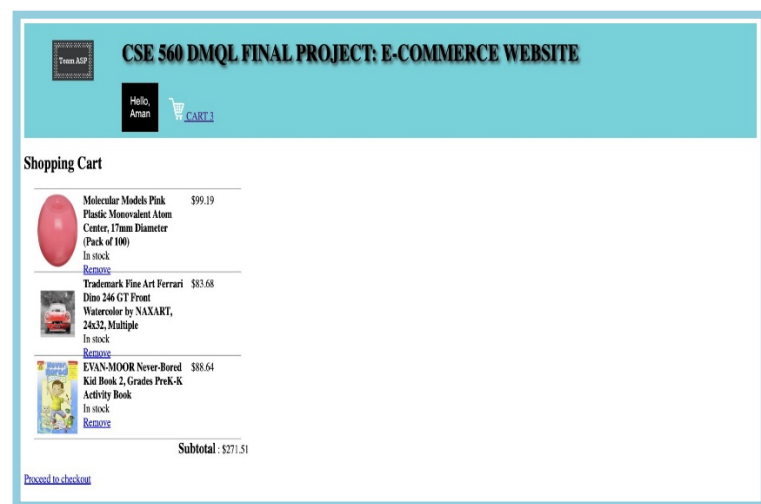
Products by Category: The user can view the products for each category and choose accordingly.



Product Details: This page consists of details of the product selected



Order page: This page consists of the items the user wanted to order and present on the card. It also displays the total amount of the products.



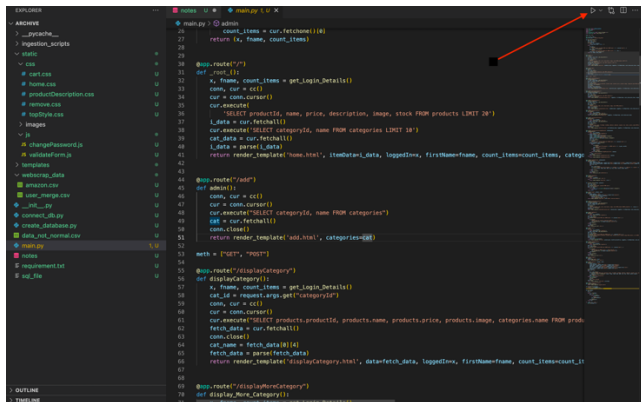
Register Page: The user can register himself with his credentials and login the next time he visits the website

Email:	<input type="text"/>
Password:	<input type="password"/>
Confirm Password:	<input type="password"/>
First Name:	<input type="text"/>
Last Name:	<input type="text"/>
Address Line 1:	<input type="text"/>
Address Line 2:	<input type="text"/>
Zipcode:	<input type="text"/>
City:	<input type="text"/>
State:	<input type="text"/>
Country:	<input type="text"/>
Phone Number:	<input type="text"/>
Register:	<input type="button" value="Register"/>
Login:	<input type="button" value="Login"/>

Edit Profile: The user can edit his/her profile using the edit profile page and it will store the latest details of the user in the database.

Way to run the project:

1. Open VS code and open the project directory.
2. Click on main.py and click on the top left to run the code.



3. A terminal will pop in the below just run the below command to install all the requirements.

pip install -r requirement.txt

or

pip3 install -r requirement.txt

4. After installing all the requirements, if you used pip just write in the terminal the below code:

python main.py

If you used pip3 while installing all the requirements, then use the below code in the terminal:

python3 main.py

5. Click ctrl + left mouse on the link (http://127.0.0.1:5000/) to redirect to the website.

Team Members Contribution:

We have successfully completed the project with equal contribution from each teammate.

Team Member	Project + Report	Contribution
Aman Prakash	<ul style="list-style-type: none"> • Data scrapping • Data preprocessing • Create tables • Data Normalization and Database Implementation • Data Ingestion • Complex queries visualization • Query execution analysis • Website Implementation 	<ul style="list-style-type: none"> • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3%
Serath Chandra Nutakki	<ul style="list-style-type: none"> • Data scrapping • Data preprocessing • Create tables • Data Normalization and Database Implementation • Data Ingestion • Complex queries visualization • Query execution analysis • Website Implementation 	<ul style="list-style-type: none"> • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3%
Prashant Upadhyay	<ul style="list-style-type: none"> • Data scrapping • Data preprocessing • Create tables • Data Normalization and Database Implementation • Data Ingestion • Complex queries visualization • Query execution analysis • Website Implementation 	<ul style="list-style-type: none"> • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3% • 33.3%

For convenient access of files please projects github repo by accessing the below link:-

https://github.com/amanprak/cse560_final_project