SQL + Python Retail Order Analysis

Performed Extract, Transform and Load(ETL)

```
# Install nbconvert if not already installed (uncomment if needed)
# !pip install nbconvert

# Convert and export notebook to PDF
!jupyter nbconvert --to pdf 'SQL_Python_Retail_Order_Analysis.ipynb'
```

1.Extraction

Extracted the data by using kaggle api

```
##pip install kaggle
```

```
→ Collecting kaggle

      Downloading kaggle-1.6.17.tar.gz (82 kB)
      Preparing metadata (setup.py) ... done
    Requirement already satisfied: six>=1.10 in /Users/itspawanrajput/Library/Python/3.11/lib/python/site-packages (from kag
    Requirement already satisfied: certifi>=2023.7.22 in /opt/homebrew/lib/python3.11/site-packages (from kaggle) (2023.7.22
    Requirement already satisfied: python-dateutil in /Users/itspawanrajput/Library/Python/3.11/lib/python/site-packages (fr
    Requirement already satisfied: requests in /opt/homebrew/lib/python3.11/site-packages (from kaggle) (2.31.0)
    Requirement already satisfied: tqdm in /opt/homebrew/lib/python3.11/site-packages (from kaggle) (4.66.1)
    Collecting python-slugify (from kaggle)
      Downloading python_slugify-8.0.4-py2.py3-none-any.whl.metadata (8.5 kB)
    Requirement already satisfied: urllib3 in /opt/homebrew/lib/python3.11/site-packages (from kaggle) (2.0.7)
    Collecting bleach (from kaggle)
      Downloading bleach-6.2.0-py3-none-any.whl.metadata (30 kB)
    Collecting webencodings (from bleach->kaggle)
      Using cached webencodings-0.5.1-py2.py3-none-any.whl.metadata (2.1 kB)
    Collecting text-unidecode>=1.3 (from python-slugify->kaggle)
      Downloading text_unidecode-1.3-py2.py3-none-any.whl.metadata (2.4 kB)
    Requirement already satisfied: charset-normalizer<4,>=2 in /opt/homebrew/lib/python3.11/site-packages (from requests->ka
    Requirement already satisfied: idna<4,>=2.5 in /opt/homebrew/lib/python3.11/site-packages (from requests->kaggle) (3.10)
    Downloading bleach-6.2.0-py3-none-any.whl (163 kB)
    Downloading python_slugify-8.0.4-py2.py3-none-any.whl (10 kB)
    Downloading text_unidecode-1.3-py2.py3-none-any.whl (78 kB)
    Using cached webencodings-0.5.1-py2.py3-none-any.whl (11 kB)
    Building wheels for collected packages: kaggle
      Building wheel for kaggle (setup.py) ... done
Created wheel for kaggle: filename=kaggle-1.6.17-py3-none-any.whl size=105786 sha256=49f01b5487f522139ca4e3e21f98e62d7
      Stored in directory: /Users/itspawanrajput/Library/Caches/pip/wheels/ff/55/fb/b27a466be754d2a06ffe0e37b248d844f090a63b
    Successfully built kaggle
    Installing collected packages: webencodings, text-unidecode, python-slugify, bleach, kaggle
    Successfully installed bleach-6.2.0 kaggle-1.6.17 python-slugify-8.0.4 text-unidecode-1.3 webencodings-0.5.1
    [notice] A new release of pip is available: 24.3.1 -> 25.0.1
    [notice] To update, run: python3.11 -m pip install --upgrade pip
Note: you may need to restart the kernel to use updated packages.
```

!kaggle datasets download ankitbansal06/retail-orders -f orders.csv

2. Transform

#reading the file and handling the missing values and null values import pandas as pd

_ _		order_id	order_date	ship_mode	segment	country	city	state	postal_code	region	category	sub_category	prod
	0	1	2023-03-01	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Bookcases	F 10
	1	2	2023-08-15	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	F 10
	2	3	2023-01-10	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	C 10
	3	4	2022-06-18	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Furniture	Tables	F 10
	4	5	2022-07-13	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Office Supplies	Storage	C 10

df.info()

df.head(5)

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
Data columns (total 19 columns):

df['profit']=df['sale_price']-df['cost_price']

#	Column	Non-Null Count	Dtype
0	order_id	9994 non-null	int64
1	order_date	9994 non-null	object
2	ship_mode	9988 non-null	object
3	segment	9994 non-null	object
4	country	9994 non-null	object
5	city	9994 non-null	object
6	state	9994 non-null	object
7	postal_code	9994 non-null	int64
8	region	9994 non-null	object
9	category	9994 non-null	object
10	sub_category	9994 non-null	object
11	product_id	9994 non-null	object
12	cost_price	9994 non-null	int64
13	list_price	9994 non-null	int64
14	quantity	9994 non-null	int64
15	discount_percent	9994 non-null	int64
16	discount	9994 non-null	float64
17	sale_price	9994 non-null	float64
18	profit	9994 non-null	float64
dtype	es: float64(3), in	t64(6) , object(10))
memoı	ry usage: 1.4+ MB		

df.info()

```
<<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
Data columns (total 19 columns):
# Column Non-Null Count
```

#	Column	Non-Null Count	Dtype
0	order_id	9994 non-null	int64
1	order_date	9994 non-null	datetime64[ns]
2	ship_mode	9988 non-null	object

```
segment
                                       object
                       9994 non-null
                       9994 non-null
4
    country
                                       object
                       9994 non-null
                                       object
5
    city
    state
                       9994 non-null
                                       object
    postal_code
                       9994 non-null
                                       int64
    region
                       9994 non-null
                                       object
    category
                       9994 non-null
                                       object
10 sub_category
                       9994 non-null
                                       object
                       9994 non-null
    product_id
                                       object
11
                       9994 non-null
                                       int64
    cost_price
    list_price
                       9994 non-null
13
                                       int64
14
                       9994 non-null
                                       int64
    quantity
                       9994 non-null
15
    discount_percent
                                       int64
16
    discount
                       9994 non-null
                                       float64
                       9994 non-null
17
    sale_price
                                       float64
18 profit
                       9994 non-null
                                       float64
dtypes: datetime64[ns](1), float64(3), int64(6), object(9)
memory usage: 1.4+ MB
```

df.head()

→		order_id	order_date	ship_mode	segment	country	city	state	postal_code	region	category	sub_category	prod
	0	1	2023-03-01	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Bookcases	F 10
	1	2	2023-08-15	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	F 10
	2	3	2023-01-10	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	C 10
	3	4	2022-06-18	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Furniture	Tables	F 10
	4	5	2022-07-13	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Office Supplies	Storage	C 10

df.drop(columns=['list_price','cost_price','discount_percent'],inplace=True) #dropping the columns
df.head()

_		order_id	order_date	ship_mode	segment	country	city	state	postal_code	region	category	sub_category	prod
	0	1	2023-03-01	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Bookcases	F 10
	1	2	2023-08-15	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	F 10
	2	3	2023-01-10	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	C 10

df.head()

_		order_id	order_date	ship_mode	segment	country	city	state	postal_code	region	category	sub_category	prod
	0	1	2023-03-01	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Bookcases	F 10
	1	2	2023-08-15	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	F 10
	2	3	2023-01-10	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	C 10

3. Load Data

Load data into SQL

```
#loading the data into the database
import sqlite3
conn = sqlite3.connect('orders.db')
df.to_sql('orders',conn,if_exists='replace',index=False)
conn.close()

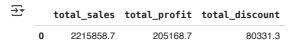
#read the data from the database
import sqlite3
```

conn = sqlite3.connect('orders.db')
data = pd.read_sql('SELECT * FROM orders',conn)
data.head()

Ť	order_id	order_date	ship_mode	segment	country	city	state	postal_code	region	category	sub_category	prod
0	1	2023-03- 01T00:00:00	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Bookcases	F 10
1	2	2023-08- 15T00:00:00	Second Class	Consumer	United States	Henderson	Kentucky	42420	South	Furniture	Chairs	F 10
2	3	2023-01- 10T00:00:00	Second Class	Corporate	United States	Los Angeles	California	90036	West	Office Supplies	Labels	C 10
3	4	2022-06- 18T00:00:00	Standard Class	Consumer	United States	Fort Lauderdale	Florida	33311	South	Furniture	Tables	F 10
4	5	2022-07-	Standard	Consumer	United	Fort	Florida	33311	South	Office	Storage	C

#total sales, profit, discount

data = pd.read_sql('SELECT SUM(sale_price) as total_sales, SUM(profit) as total_profit, SUM(discount) as total_discount FROM
data.head()



find top 10 highest revenew generating products

data = pd.read_sql('SELECT product_id, SUM(sale_price) as total_sales FROM orders GROUP BY product_id ORDER BY total_sales [
data.head()

} ▼		product_id	total_sales
	0	TEC-CO-10004722	59514.0
	1	OFF-BI-10003527	26525.3
	2	TEC-MA-10002412	21734.4
	3	FUR-CH-10002024	21096.2
	4	OFF-BI-10001359	19090.2

find top 10 highest revenew generating cities

data = pd.read_sql('SELECT city, SUM(sale_price) as total_sales FROM orders GROUP BY city ORDER BY total_sales DESC LIMIT 10
data.head()

_ →		city	total_sales
	0	New York City	247205.7
	1	Los Angeles	169758.4
	2	Seattle	115358.7
	3	San Francisco	108890.1
	4	Philadelphia	105258.3

find top 5 highest selling products in each region

data = pd.read_sql('SELECT region, product_id, SUM(sale_price) as total_sales FROM orders GROUP BY region, product_id ORDER
data.groupby('region').head(5)
data.head()

Z *		region	product_id	total_sales
	0	Central	TEC-CO-10004722	16975.0
	1	Central	TEC-MA-10000822	13770.0
	2	Central	OFF-BI-10001120	11056.5
	3	Central	OFF-BI-10000545	10132.7

for each category find the month with the highest sales

result = data.groupby('category').head(1)

4 Central OFF-BI-10004995

8416.1

```
print("\nHighest sales month for each category:")
print(result)
    Highest sales month for each category:
               category month total_sales
              Furniture
                                    71649.5
                           ดล
                                    77959.5
        Office Supplies
    12
                            02
    24
              Technology
                           10
                                   103021.1
# Install nbconvert if not already installed (uncomment if needed)
# !pip install nbconvert
# Convert and export notebook to PDF
!jupyter nbconvert --to pdf 'SQL_Python_Retail_Order_Analysis.ipynb'

    usage: jupyter [-h] [--version] [--config-dir] [--data-dir] [--runtime-dir]

                     --paths] [--json] [--debug]
                    [subcommand]
    Jupyter: Interactive Computing
    positional arguments:
                     the subcommand to launch
      subcommand
    options:
      -h, --help
                      show this help message and exit
      --version
                      show the versions of core jupyter packages and exit
      --config-dir
                      show Jupyter config dir
                      show Jupyter data dir
      --data-dir
      --runtime-dir
                      show Jupyter runtime dir
      --paths
                      show all Jupyter paths. Add -- json for machine-readable
                      format.
                      output paths as machine-readable json
      --ison
      --debug
                      output debug information about paths
    Available subcommands: kernel kernelspec migrate run troubleshoot
    Jupyter command `jupyter-nbconvert` not found.
```

pip install nbconvert

🚌 Requirement already satisfied: nbconvert in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/python3.1 Requirement already satisfied: beautifulsoup4 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pyth Requirement already satisfied: bleach!=5.0.0 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pytho Requirement already satisfied: defusedxml in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/python3. Requirement already satisfied: jinja2>=3.0 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/python3 Requirement already satisfied: jupyter-core>=4.7 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/p Requirement already satisfied: jupyterlab-pygments in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib Requirement already satisfied: markupsafe>=2.0 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pyt Requirement already satisfied: mistune<4,>=2.0.3 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/p Requirement already satisfied: nbclient>=0.5.0 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pyt Requirement already satisfied: nbformat>=5.7 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pytho Requirement already satisfied: packaging in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/python3.1 Requirement already satisfied: pandocfilters>=1.4.1 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/li Requirement already satisfied: pygments>=2.4.1 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pyt Requirement already satisfied: traitlets>=5.1 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pyth Requirement already satisfied: webencodings in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/python Requirement already satisfied: tinycss2<1.5,>=1.1.0 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/li Requirement already satisfied: platformdirs>=2.5 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/p Requirement already satisfied: jupyter-client>=6.1.12 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/ Requirement already satisfied: fastjsonschema>=2.15 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/li Requirement already satisfied: jsonschema>=2.6 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pyt Requirement already satisfied: soupsieve>1.2 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pytho Requirement already satisfied: typing-extensions>=4.0.0 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.ven Requirement already satisfied: attrs>=22.2.0 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pytho Requirement already satisfied: jsonschema-specifications>=2023.03.6 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngi Requirement already satisfied: referencing>=0.28.4 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib Requirement already satisfied: rpds-py>=0.7.1 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/pyth Requirement already satisfied: python-dateutil>=2.8.2 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/ Requirement already satisfied: pyzmq>=23.0 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/python3 Requirement already satisfied: tornado>=6.2 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/python Requirement already satisfied: six>=1.5 in /Users/itspawanrajput/Desktop/Goal/RedditDataEngineering/.venv/lib/python3.13 Note: you may need to restart the kernel to use updated packages.

Start coding or generate with AI.