# Final Project Document

**Project Topic**

## Alibaba Digital Marketing



**Team Details**

## Team 5

## Prathamesh Verlekar

## Yash Pandya

# Overview

**Alibaba Group Holding Limited** (also known as **Alibaba Group** and **Alibaba.com**) is a Chinese multinational technology company specializing in e-commerce, retail, Internet, and technology. On 19th September 2014, Alibaba's market value was US$231 billion. It is one of the top 10 most valuable and is the 59th biggest public company in the world by the Global 2000 list.

As of 2020, Alibaba has the 6th highest global brand valuation. Alibaba is the world's largest retailer and e-commerce company, is on the list of largest Interest companies and artificial intelligence companies, is one of the biggest venture capital firms, and one of the biggest investment corporations in the world. Its online sales and profits surpassed all US retailers (including Walmart, Amazon, and eBay) combined since 2015.

# Goals

➢ Finding Recency, Frequency, and Monetary Value to find the company's best customers by using certain measures.
➢ To find the customers lost by calculating the churn rate and the customer's lifetime value.
➢ Build a recommendation system for suggesting products to the customers that they might also like.
➢ To Crawl, Index, and Rank data using the Search engine for better customer experience.
➢ Build a dashboard for business owners to better understand their sales and decide future marketing strategies.

➤ Create a web application for better user experience.

# Dataset

➤ We will be using datasets available on Kaggle along with data available through different open sources and websites.
➤ Kaggle Dataset Link:https://www.kaggle.com/AppleEcomerceInfo/ecommerce-information?select=products.txt
➤ http://yongfeng.me/dataset/

## Data Sample

```
                                                           product_id    product_name    descriptions
carts_has_products.txt                                     1200    Macbook Pro (2017)      The ultimate pro notebook. MacBook Pro features faster processors ;upgraded memory;the Apple
category.txt                                               1300    Macbook Air (2015)      MacBook Air lasts up to an incredible 12 hours between charges So from your morning coffee
option.txt                                                 1400    Iphone X       The iPhone X display is so immersive the device itself disappears into the experience.
order.txt                                                  1500    Iphone 7       Great connectivity of this device includes Bluetooth 4.2 version with A2DP
orders_has_products.txt                                    1600    Iphone 8       iPhone 8 introduces a glass design. The glass back enables easy wireless charging.
orders_paid_creditcard.txt                                 1700    Ipad Air       4th gen The iPad Air is unbelievably thin and light. And yet it is so much more powerful and capable
orders_placed_user.txt                                     1800    Ipad Mini 3th gen       3th gen Everything you love about iPad ◊ the beautiful screen and fast
product_sold_vendor.txt                                    1900    ESC8000 G3      G3 High-density GPU server with hybrid computing power. ASUS-patented Adaptable Topology design.
products.txt                                               2000    ESC8000 G4      G4 High performance ASUS 2U server with hybrid-storage design and high power-efficiency
products_belong_category.txt                               2100    XPS 13 - 5080   Thinner and more powerful than ever the Dell XPS reinforces its lofty standing with an 8th Gen Intel
products_has_options.txt                                   2200    XPS 15 - 5070   Ultra-thin and distinctly refined the stylish Dell Inspiron gives definitive elegance to a powerful
shoppingcart.txt                                           2300    Monoprice Ultra Slim Series High Speed HDMI Cable       The Monoprice Ultra Slim Active High Speed HDMI Cable series
user.txt                                                   2400    Monoprice Ultra Slim Series High Speed HDMI Cable - 4K  Monoprice Commercial Cable supports the following HDMI feat
user_has_creditcard.txt                                    2500    Avantree HT3189 Wireless Headphones     Avantree HT3189 Wireless Headphones for TV Watching & PC Gaming with Bluetoo
vendor.txt                                                 2600    COWIN E7 PRO   Active Noise Cancelling Headphone Bluetooth Headphones with Microphone Hi-Fi Deep Bass Wireless Hea
```

# Personas

**Who –** Alibaba's technical teams, stakeholders, and sponsors for deciding their future marketing strategies.

**What** – Build a recommendation system for the company's customers and dashboards for stakeholders to analyze their sales and other insights.

**When** – It will be completed within 2 weeks of a given timeline.

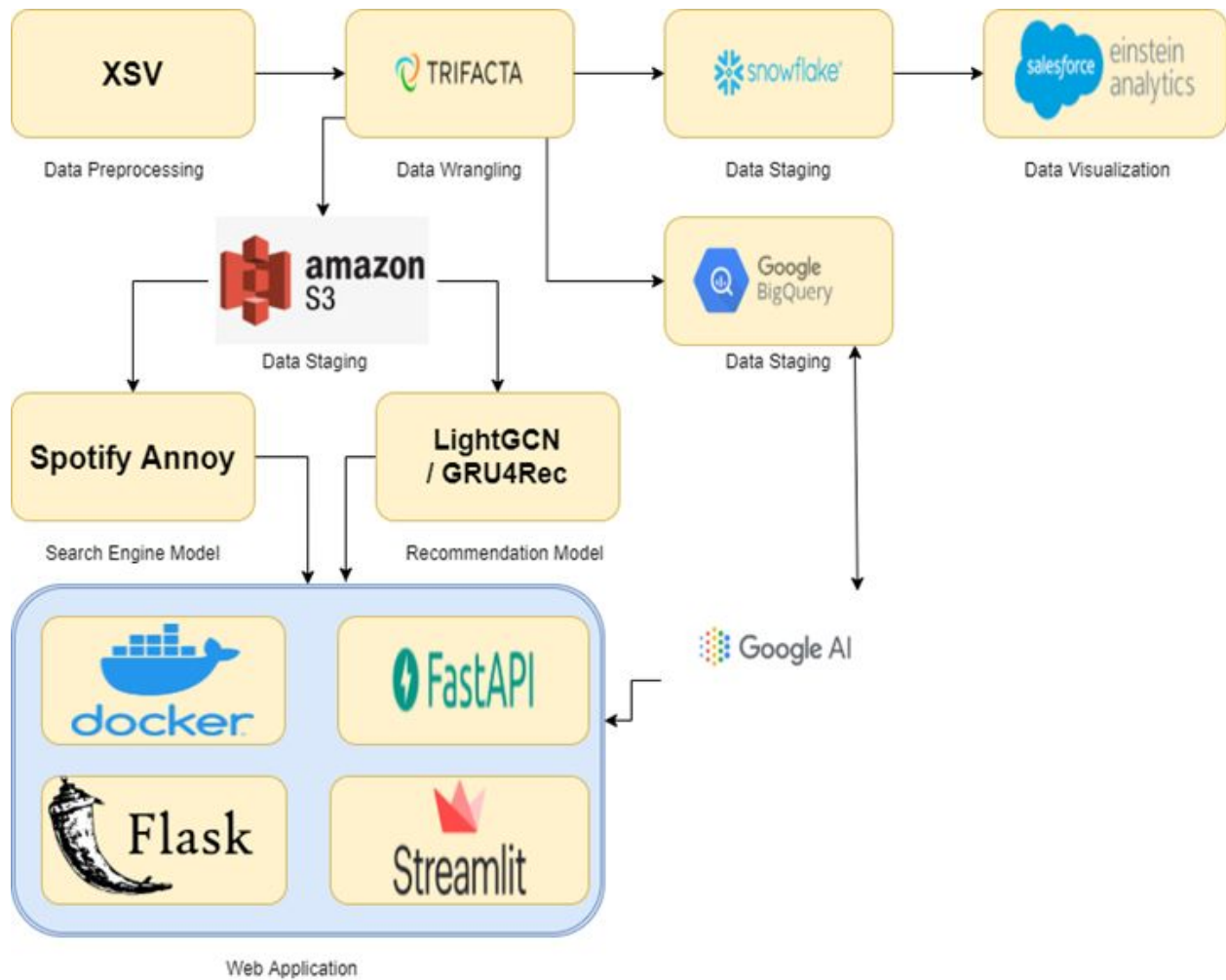**Where** – We will be working at our remote locations and then collaborate on our work.

**Why** – To get data insights that will help the stakeholders decide their marketing strategies and to improve the customer experience by building a recommendation system that will recommend the customers product that may also like.

**How** – Using the tools and technologies learned in the course.

- **XSV**
- **Python**
- **Trifacta**
- **Salesforce Einstein Analytics**
- **Streamlit**
- **Flask**

➢ **JMeter**

# Workflow



| Stage | Tool |
|---|---|
| Data Preprocessing | XSV |
| Data Wrangling | TRIFACTA |
| Data Staging | snowflake |
| Data Visualization | salesforce einstein analytics |
| Data Staging | amazon S3 |
| Data Staging | Google BigQuery |
| Search Engine Model | Spotify Annoy |
| Recommendation Model | LightGCN / GRU4Rec |
| Web Application | docker, FastAPI, Flask, Streamlit |

## Use Cases

➢ **Making user shopping experience more efficient.**
➢ **Know the product trend to make future decisions based on it to increase profitability.**
➢ **Help stakeholders to make important decisions and changes.**

## Web scraping

Web Scraping (also termed Screen Scraping, Web Data Extraction, Web Harvesting, etc.) is a technique employed to extract large amounts of data from websites whereby the data is extracted and saved to a local file in your computer or to a database in table (spreadsheet) format.

## Code Sample:

```python
#For web scrapping we will import urlopen and BeautifulSoup
from urllib.request import urlopen
from bs4 import BeautifulSoup
```

```python
#Link of website which we will be scrapping
url = "https://www.alibaba.com/"
html = urlopen(url)
```

```python
#Creating a Beautiful Soup object
soup = BeautifulSoup(html, 'lxml')
type(soup)

bs4.BeautifulSoup
```

```python
#Getting the title of the website which we are scrapping
title = soup.title
print(title)
```

```python
#Printing the text on the webpage
text = soup.get_text()
#print(text)
```

```python
soup.find_all('a')
```

# Website:



# Trifacta (Data Wrangling)

Trifacta's data wrangling software allows you to prepare & visualize complex data in no time.

# Previewing Data_1.csv

| # Quantity | 🕐 InvoiceDate | ## UnitPrice | # CustomerID |
|---|---|---|---|
| 6 | 12/1/2010 8:26 | 2.55 | 17850 |
| 6 | 12/1/2010 8:26 | 3.39 | 17850 |
| 8 | 12/1/2010 8:26 | 2.75 | 17850 |
| 6 | 12/1/2010 8:26 | 3.39 | 17850 |
| 6 | 12/1/2010 8:26 | 3.39 | 17850 |
| 2 | 12/1/2010 8:26 | 7.65 | 17850 |
| 6 | 12/1/2010 8:26 | 4.25 | 17850 |
| 6 | 12/1/2010 8:28 | 1.85 | 17850 |
| 6 | 12/1/2010 8:28 | 1.85 | 17850 |
| 32 | 12/1/2010 8:34 | 1.69 | 13047 |
| 6 | 12/1/2010 8:34 | 2.1 | 13047 |
| 6 | 12/1/2010 8:34 | 2.1 | 13047 |
| 8 | 12/1/2010 8:34 | 3.75 | 13047 |
| 6 | 12/1/2010 8:34 | 1.65 | 13047 |
| 6 | 12/1/2010 8:34 | 4.25 | 13047 |
| 3 | 12/1/2010 8:34 | 4.95 | 13047 |
| 2 | 12/1/2010 8:34 | 9.95 | 13047 |

## Data_1

**Edit Recipe**  **Add**  ...

Recipe    Data

Steps Preview

1  **Rename** StockCode to 'ProductId'

2  **Delete rows where** ISMISSING([Description])

3  **Delete rows where** ISMISSING([CustomerID])

4  **Create** column1 from IF(Country == 'United Kingdom', IF(UnitPrice > 5, 'United States', IF(UnitPrice > 2, 'India', 'United Kingdom')), Country)

5  **Delete** Country

6  **Rename** column1 to 'Country'

7  **Create** Age from RANDBETWEEN(10, 80)

8  **Move** InvoiceNo before Age

9  **Move** CustomerID before InvoiceNo

10  **Move** ProductId after CustomerID

11  **Change** InvoiceNo **type to** String

12  **Change** ProductId **type to** String

Dataset
Data_1.csv

Recipe
Data_1

Output
Data_1

# RFM and Customer Segmentation

**RFM stands for three dimensions:**

- Recency – How recently did the customer purchase?
- Frequency – How often do they purchase?
- Monetary Value – How much do they spend?

```
In [71]:  ▶  df_cleaned.head(10)
```

Out[71]:

| ge | Description | Quantity | InvoiceDate | UnitPrice | Country | QuantityCanceled | TotalPrice | min_recency | max_recency | frequency | monetary_value | RFMScore |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 25 | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 2010-12-01 08:26:00 | 2.55 | 17 | 0 | 15.30 | 372.0 | 373.0 | 34 | 5327.79 | 411 |
| 68 | WHITE METAL LANTERN | 6 | 2010-12-01 08:26:00 | 3.39 | 17 | 0 | 20.34 | 372.0 | 373.0 | 34 | 5327.79 | 411 |
| 49 | CREAM CUPID HEARTS COAT HANGER | 8 | 2010-12-01 08:26:00 | 2.75 | 17 | 0 | 22.00 | 372.0 | 373.0 | 34 | 5327.79 | 411 |
| 69 | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 2010-12-01 08:26:00 | 3.39 | 17 | 0 | 20.34 | 372.0 | 373.0 | 34 | 5327.79 | 411 |
| 34 | RED WOOLLY HOTTIE WHITE HEART. | 6 | 2010-12-01 08:26:00 | 3.39 | 17 | 0 | 20.34 | 372.0 | 373.0 | 34 | 5327.79 | 411 |
| 60 | SET 7 BABUSHKA NESTING BOXES | 2 | 2010-12-01 08:26:00 | 7.65 | 37 | 0 | 15.30 | 372.0 | 373.0 | 34 | 5327.79 | 411 |
| 23 | GLASS STAR FROSTED T-LIGHT HOLDER | 6 | 2010-12-01 08:26:00 | 4.25 | 17 | 0 | 25.50 | 372.0 | 373.0 | 34 | 5327.79 | 411 |

# SnowFlake

Snowflake is a cloud-based Data Warehouse solution provided as a Saas (Software-as-a-Service) with full support for ANSI SQL. It also has a unique architecture that enables users to just create tables and start querying data with very little administration or DBA activities needed.

# Salesforce Einstein Analytics Dashboard

Salesforce.com, inc. is an American cloud-based software company headquartered in San Francisco, California. It provides customer relationship management (CRM) service and also sells a complementary suite of enterprise applications focused on customer service, marketing automation, analytics, and application development.

- Total Revenue generated , Total Quantity sold, Top country that use the website

- Which Country has purchased most items

- Category wise items ordered and Cancelled

- How many items after placing in cart where purchased

- Number of orders based on weekdays

- Which has platform(Phone and Web) used for Most

- Quantity sold based on segments (Loyal, potential Loyalist, Cannot lose them, new Customer and Lost customers)

- Tablewise segmentation

# ALIBABA WEBSITE INSIGHTS

| TOTAL REVENUE GENERATED | TOTAL QUANTITY SOLD | TOP COUNTRY THAT USE WEBSITE |
|---|---|---|
| 5.9M | 3.6M | CHINA |

## Category

Search for values...

Show Selected (0)

| | |
|---|---|
| Accessories | 121.9k |
| Apparel | 690.6k |
| Books | 50.0k |
| Design | 2.3k |
| Electronics | 423.6k |
| Gift & Food | 176.2k |
| Health | 164.5k |
| Home & Decor | 1.7x |
| Kitchen & Dining | 225.8k |

## Quantity Sold Country Wise

1.5M

63k

## Category Wise Quantity ordered and Cancelled

Measure
- Sum of Quantity
- Sum of QuantityCanceled

Sum of Quantity / Sum of QuantityCancelled

2M, 1.5M, 1M, 500k, 0 ... 100k, 80k, 60k, 40k, 20k

Accessories 122k, Apparel 691k, Books 50k, Design 2.3k, Electronics 424k, Gift & Food 176k, Health 164k, Home & Decor 1.7M, Kitchen & Dining 226k

Category

## Quantity Sold Category Wise

Category
- Accessories
- Apparel
- Books
- Design
- Electronics
- Gift & Food
- Health
- Home & Decor
- Kitchen & Dining

Home & Decor (1.7M)
Apparel (691k)
Kitchen & Dining (226k)
Gift & Food (176k)
Health (164k)
Electronics (424k)
Accessories (122k)
Books (50k)

## Analysis of Website

Analysis
- AddToCart
- ClickProduct
- Purchased

44.76%

22.14%
33.1%
44.76%

## Sum Of Orders on Weekdays

Sum of Number Of Orders

80k, 60k, 40k, 20k, 0

81k, 24k, 45k, 56k, 14k, 10k, 15k

Average Line

Friday, Monday, Saturday, Sunday, Thrusday, Tuesday, Wednesday

## Quantity Sold Segment Wise

Sum of Quantity: 3.6M

| Cannot Loose Them | Lost Customers | Loyal | New Customer | Potential Loyalist |
|---|---|---|---|---|
| 730k | 123k | 431k | 252k | 2M |

## Platform Used

Platform Used
- Phone
- Web

55.32%

44.68%

| # | Category | Description | Segments | ProductId | Monetary | Quantity | Recency | Frequency | Country | TotalPrice | Age | InvoiceDate | UnitPrice |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Electronics | WHITE HANGING HEART T-LIGHT HOLDER | New Customer | 85123A | 5,328 | 6 | 373 | 34 | China | 15 | 25 | 12/1/2019 8:26 | 2.6 |
| 2 | Electronics | WHITE METAL LANTERN | New Customer | 71053 | 5,328 | 6 | 373 | 34 | China | 20 | 68 | 12/1/2019 8:26 | 3.4 |
| 3 | Home & Decor | CREAM CUPID HEARTS COAT HANGER | New Customer | 84406B | 5,328 | 8 | 373 | 34 | China | 22 | 49 | 12/1/2019 8:26 | 2.8 |
| 4 | Kitchen & Dining | KNITTED UNION FLAG HOT WATER BOTTLE | New Customer | 84029G | 5,328 | 6 | 373 | 34 | China | 20 | 69 | 12/1/2019 8:26 | 3.4 |
| 5 | Home & Decor | RED WOOLLY HOTTIE WHITE HEART. | New Customer | 84029E | 5,328 | 6 | 373 | 34 | China | 20 | 34 | 12/1/2019 8:26 | 3.4 |
| 6 | Home & Decor | SET 7 BABUSHKA NESTING BOXES | New Customer | 22752 | 5,328 | 2 | 373 | 34 | China | 15 | 60 | 12/1/2019 8:26 | 7.7 |
| 7 | Home & Decor | GLASS STAR FROSTED T-LIGHT HOLDER | New Customer | 21730 | 5,328 | 6 | 373 | 34 | China | 26 | 23 | 12/1/2019 8:26 | 4.3 |
| 8 | Apparel | HAND WARMER UNION JACK | New Customer | 22633 | 5,328 | 6 | 373 | 34 | China | 11 | 42 | 12/1/2019 8:28 | 1.9 |
| 9 | Apparel | HAND WARMER RED POLKA DOT | New Customer | 22632 | 5,328 | 6 | 373 | 34 | China | 11 | 29 | 12/1/2019 8:28 | 1.9 |
| 10 | Apparel | HAND WARMER RED POLKA DOT | New Customer | 22632 | 5,328 | 6 | 373 | 34 | China | 11 | 30 | 12/1/2019 9:01 | 1.9 |
| 11 | Apparel | HAND WARMER UNION JACK | New Customer | 22633 | 5,328 | 6 | 373 | 34 | China | 11 | 44 | 12/1/2019 9:01 | 1.9 |
| 12 | Electronics | WHITE HANGING HEART T-LIGHT HOLDER | New Customer | 85123A | 5,328 | 6 | 373 | 34 | China | 15 | 80 | 12/1/2019 9:02 | 2.6 |
| 13 | Electronics | WHITE METAL LANTERN | New Customer | 71053 | 5,328 | 6 | 373 | 34 | China | 20 | 44 | 12/1/2019 9:02 | 3.4 |
| 14 | Home & Decor | CREAM CUPID HEARTS COAT HANGER | New Customer | 84406B | 5,328 | 8 | 373 | 34 | China | 22 | 64 | 12/1/2019 9:02 | 2.8 |

# Recommendation Model

## LightGCN

LightGCN is a simplified design of GCN to make it more concise and appropriate for a recommendation.



**Light Graph Convolution (LGC)**

In LightGCN, we adopt the simple weighted sum aggregator and abandon the use of feature transformation and nonlinear activation.

```
In [12]:  ▶| topk_scores = model.recommend_k_items(test, top_k=TOP_K, remove_seen=True)

            topk_scores.head()
```

Out[12]:

|   | userID | itemID | prediction |
|---|--------|--------|------------|
| 0 | 12349 | 23245 | 8.174675 |
| 1 | 12349 | 22839 | 7.465651 |
| 2 | 12349 | 22423 | 7.303855 |
| 3 | 12349 | 23284 | 6.950679 |
| 4 | 12349 | 22507 | 6.786175 |

```
In [13]:  ▶| eval_map = map_at_k(test, topk_scores, k=TOP_K)
            eval_ndcg = ndcg_at_k(test, topk_scores, k=TOP_K)
            eval_precision = precision_at_k(test, topk_scores, k=TOP_K)
            eval_recall = recall_at_k(test, topk_scores, k=TOP_K)

            print("MAP:\t%f" % eval_map,
                  "NDCG:\t%f" % eval_ndcg,
                  "Precision@K:\t%f" % eval_precision,
                  "Recall@K:\t%f" % eval_recall, sep='\n')
```

```
MAP:      0.041018
NDCG:     0.124642
Precision@K:    0.088660
Recall@K:       0.090790
```

## Surprise Singular Value Decomposition (SVD)

SVD introduces two new scalar variables: the user biases bu and item biases bi. The user biases are supposed to capture the tendency of some users to rate items higher (or lower) than the average. The same goes for items: some items are usually rated higher than some others.

```
In [41]:  ▶| svd = surprise.SVD(random_state=0, n_factors=200, n_epochs=30, verbose=True)

          with Timer() as train_time:
              svd.fit(train_set)

          print("Took {} seconds for training.".format(train_time.interval))
```

```
Processing epoch 0
Processing epoch 1
Processing epoch 2
Processing epoch 3
Processing epoch 4
Processing epoch 5
Processing epoch 6
Processing epoch 7
Processing epoch 8
Processing epoch 9
Processing epoch 10
Processing epoch 11
Processing epoch 12
Processing epoch 13
Processing epoch 14
```

```
print("RMSE:\t\t%f" % eval_rmse,
      "MAE:\t\t%f" % eval_mae,
      "rsquared:\t%f" % eval_rsquared,
      "exp var:\t%f" % eval_exp_var, sep='\n')

print('----')

print("MAP:\t%f" % eval_map,
      "NDCG:\t%f" % eval_ndcg,
      "Precision@K:\t%f" % eval_precision,
      "Recall@K:\t%f" % eval_recall, sep='\n')
```

```
RMSE:           2.451093
MAE:            2.000027
rsquared:      -1.992368
exp var:       -0.000013
----
MAP:    0.002902
NDCG:   0.016522
Precision@K:    0.013356
Recall@K:       0.008077
```

# Search Engine

## Search By an Artistic Style

One typical business case is an eCommerce website that allows to search of a poster based on the example uploaded by the user. A user would usually expect to get results that are similar in terms of artistic style. A real search system would typically mix style similarity with other similarity scores such as image subject or category (landscape, still life, etc.) that can also be obtained using deep neural networks.

```python
def search_by_style(reference_image, max_results=10):
    v0 = image_style_embeddings[reference_image]
    distances = {}
    for k,v in image_style_embeddings.items():
        d = sc.spatial.distance.cosine(v0, v)
        distances[k] = d

    sorted_neighbors = sorted(distances.items(), key=lambda x: x[1], reverse=False)

    f, ax = plt.subplots(1, max_results, figsize=(16, 8))
    for i, img in enumerate(sorted_neighbors[:max_results]):
        ax[i].imshow(images[img[0]])
        ax[i].set_axis_off()

    plt.show()
```

```python
search_by_style('23200.jpg')
```



# Streamlit Application

Streamlit is an open-source app framework for Machine Learning and Data Science teams. Create beautiful data apps in hours, not weeks. All in pure Python.

https://finalproject-streamlit-app.herokuapp.com/

# SEARCH FOR SIMILAR PRODUCTS



Choose a product

Products:

| 1528.jpg | ▾ |
|---|---|

## Product Selected



## Similar Products



1528.jpg



1529.jpg



1530.jpg



1531.jpg

# Heroku Deployment

Heroku is a cloud platform as a service supporting several programming languages. One of the first cloud platforms, Heroku has been in development since June 2007, when it supported only the Ruby programming language, but now supports Java, Node.js, Scala, Clojure, Python, PHP, and Go.



# FastAPI

**FastAPI** is a modern, fast (high-performance), a web framework for building APIs with Python 3.6+ based on standard Python type hints. The key features are: ... One of the

fastest Python frameworks available. Fast to code: Increase the speed to develop features by about 200% to 300%.

# JMeter

Apache **JMeter** is a **testing** tool used for analyzing and measuring the performance of different software services and products. It is a pure Java open source software used for **testing** Web Application or FTP application. It is used to execute performance **testing**, load **testing,** and functional **testing** of web applications.

We used Jmeter to test our three recommendation models and based on the throughput rate we were able to decide that LightGCN was much better and faster as compared to other models so we used LightGCN for building our recommendation web application.

# Flash Application Deployed on AWS

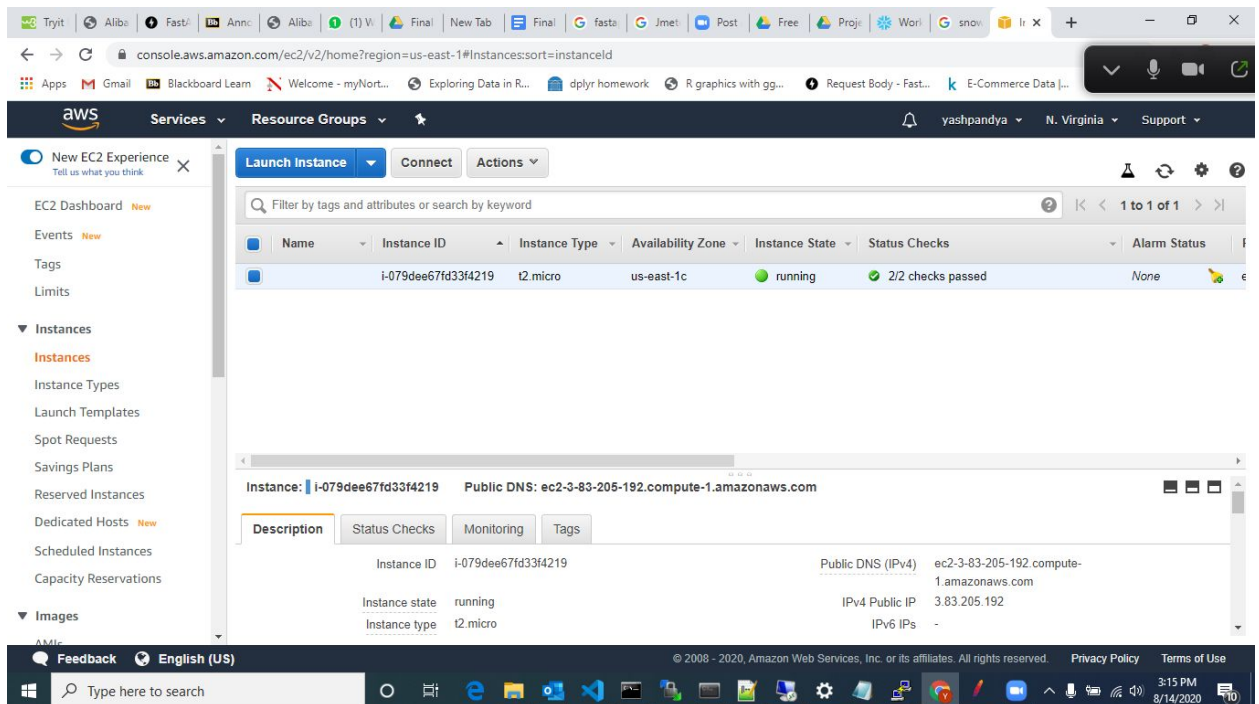http://ec2-3-83-205-192.compute-1.amazonaws.com:5000/home.html

**Flask** is a micro web framework written in Python. It is classified as a microframework because it does not require particular tools or libraries. It has no database abstraction layer, form validation, or any other components where pre-existing third-party libraries provide common functions.

**AWS** provides on-demand access to scalable web and **application** servers, storage, databases, content delivery, cache, search, and other **application services** that make it easier to build and run apps that deliver a great customer experience.
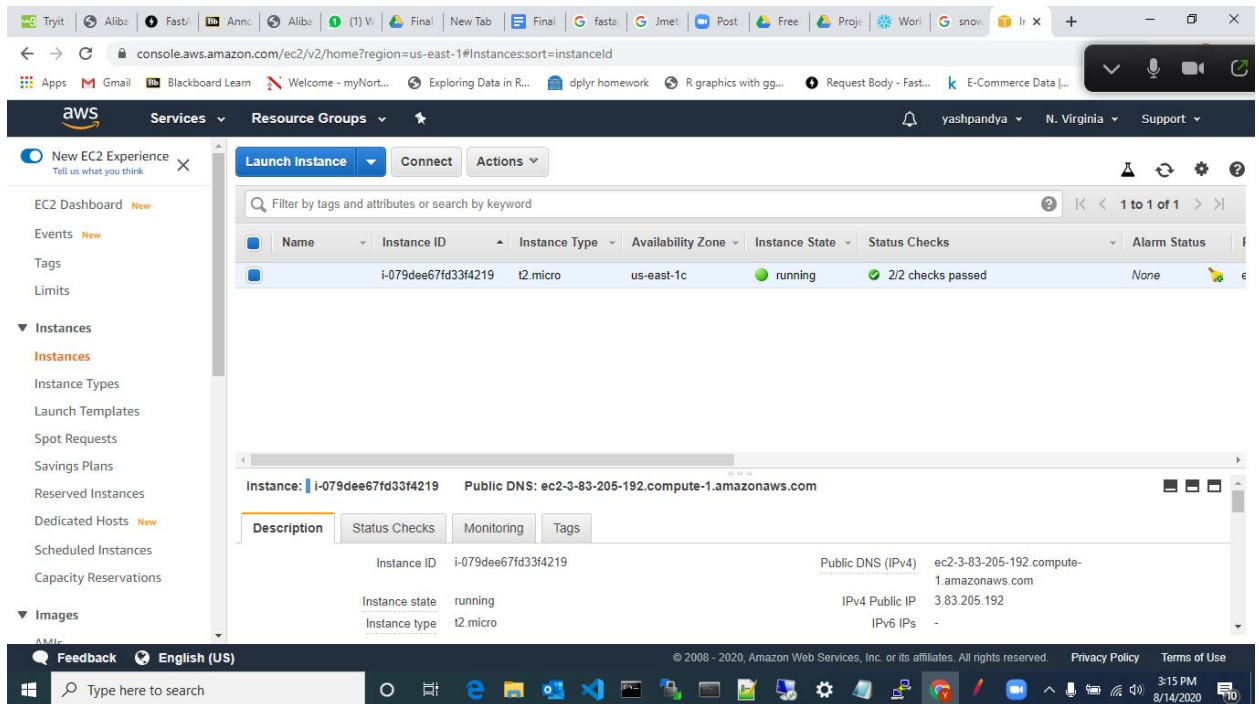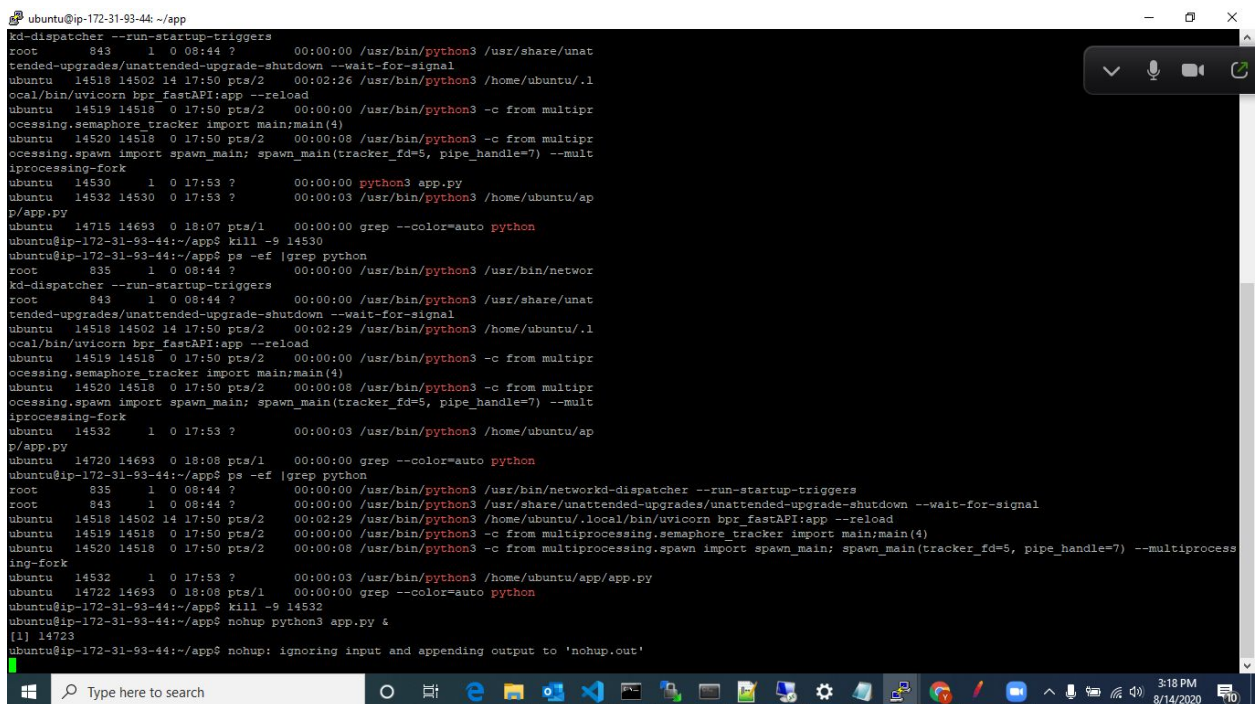

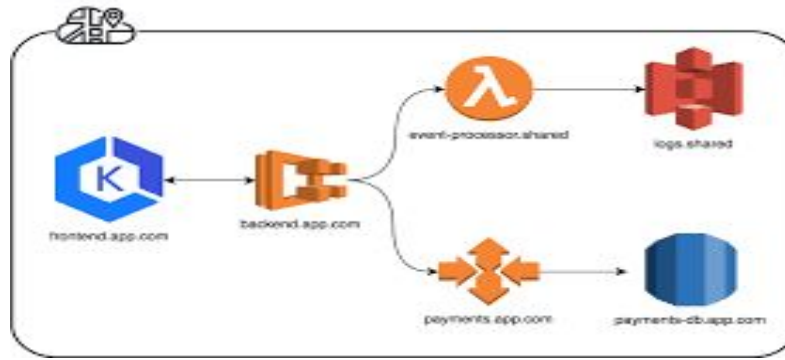
Amazon EC2

1) AWS Instance created on AWS account



2) File Transferred using WinSCP

3) Used putty for deployment

# Conclusion

We were able to implement various tools we have learned in this course to help support Alibaba's Website analyze their sales and increase their profit.

# References

1. https://github.com/microsoft/recommenders
2. https://fastapi.tiangolo.com/tutorial/body/
3. https://github.com/microsoft/recommenders/blob/master/examples/00_quick_start/sequential_recsys_amazondataset.ipynb
4. https://github.com/microsoft/recommenders/blob/master/examples/02_model_collaborative_filtering/lightgcn_deep_dive.ipynb
5. https://github.com/microsoft/recommenders/tree/master/tests
6. https://github.com/davidefiocco/streamlit-fastapi-model-serving