# RETAIL CLICKSTREAM ANALYSIS AND PREDICTION



**PROJECT REPORT for CS-GY 6513 BIG DATA**

*at*

*Submitted by*

**NIHARIKA KRISHNAN (nk2982)**

**VAIBHAV SINGH (vs2410)**

# CONTENTS

# INTRODUCTION

Over the last few years, e-commerce has become an indispensable part of the global retail framework. Like many other industries, the retail landscape has undergone a substantial transformation following the advent of the internet, and thanks to the ongoing digitalization of modern life, consumers from virtually every country now profit from the perks of online transactions. According to statistics, in the year 2021, almost 2 Billion people purchased goods online worldwide resulting in USD 4.9 Trillion dollars in revenue. The key drivers of success over the next decade will be centered on building a deep understanding of and connection to the empowered consumer, and the only way to understand consumer behavior is to measure and analyze. Hence, retailers work towards gathering data from different areas that have an impact on their online store and use it to understand the trends and the shift in consumers' behavior to make data-driven decisions that will drive more online sales.
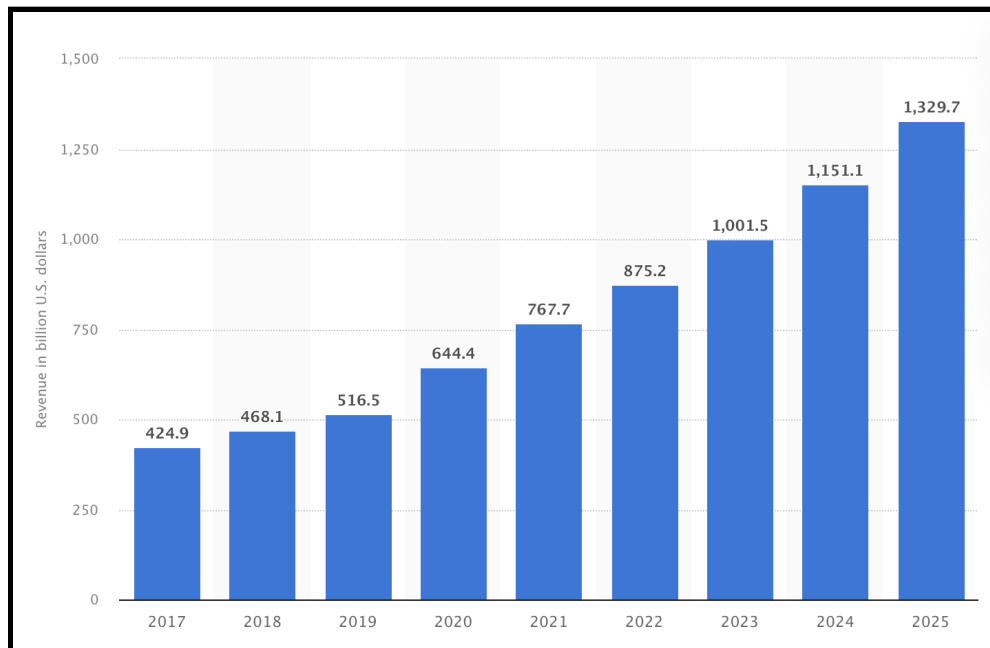


Figure-1: Retail e-commerce revenue in the US from 2017 to 2025 *(in billion U.S. dollars)*

Clickstream Data is the user's digital footprint left on a specific website during a browsing session. Clickstream analytics is the process of collecting, analyzing, and reporting aggregate data regarding pages a user visits on a website and the order in which they visit them. The path a user takes through a website is called the clickstream.

# ARCHITECTURE

The data layer is currently a single CSV file, and this can be scaled with the help of any current big data storage solution like Amazon S3 or MongoDB. To perform big data analytics we have chosen Pyspark which is the Python API for Apache Spark, an open-source, distributed computing framework with a set of libraries for real-time, large-scale data processing.

For visualization, we have used Matplotlib and Plotly which is an interactive, open-source, and browser-based graphing library for Python. Finally, we used Streamlit, which is an open-source Python library that makes it easy to create and share beautiful, custom web apps for machine learning and data science. This is used in predictive analytics to simulate a real-time deployment.
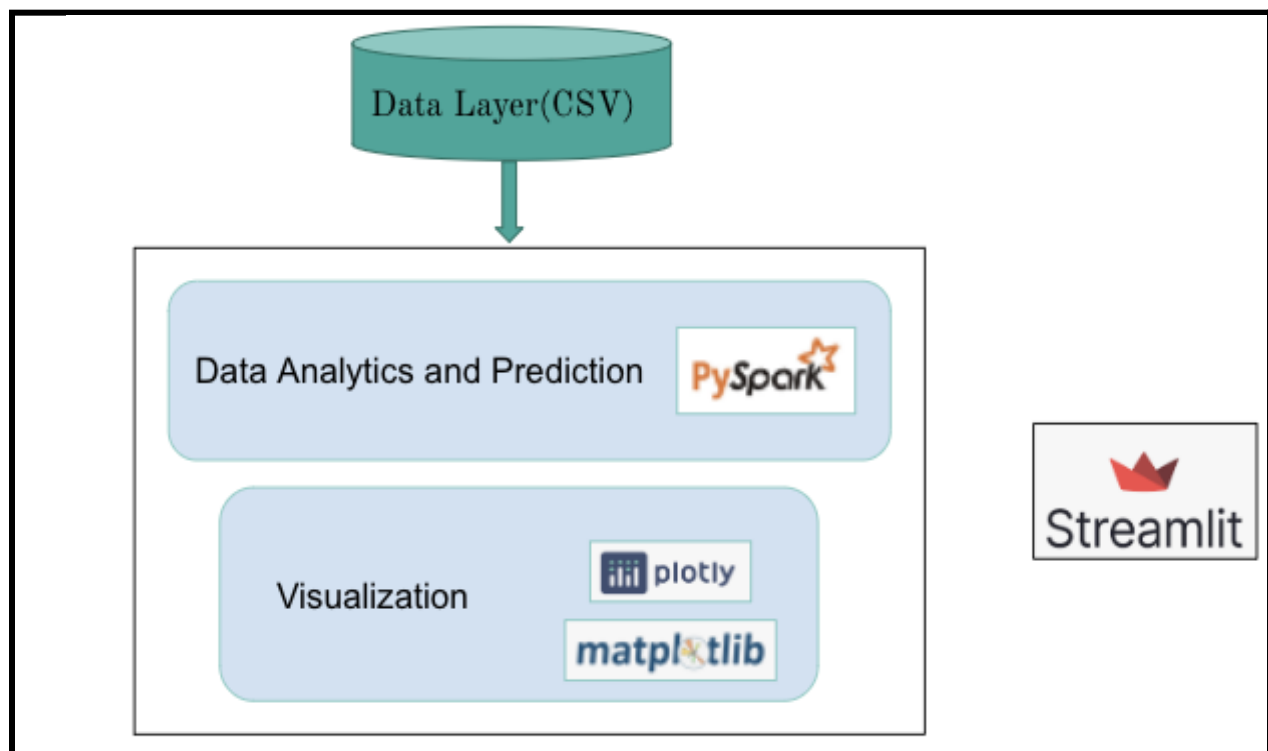


Figure-2: Architecture

# DATASET

Dataset used for this Big Data project is clickstream user behavior data for October 2019 available from a large multi-category online store.

**Dataset Source:** It was collected by the Open CDP project

**Dataset Size:** 5.6 GB (Oct 2019)

If data is collected for one year, this would be an estimated 60GB+, making centralized computing futile in both offline and real-time inference. Hence, use of Big Data technologies is imperative to make analytics computationally fast leading to valuable insights and making critical decisions.

**Dataset Overview:**



**Dataset without Pre-Processing:**

# DATA PROCESSING AND IMPUTATION

Data preprocessing and imputation are some of the most important steps in the analytics pipeline that help in providing clean, usable and efficient data. We extract informative features like Category and Product from Category_Code, Date and Time from Timestamp and handle Null values in our category column by imputing them with brand values.

**<u>Extract Category & Product from Category Code:</u>**

```python
@udf
def extract_category(category, brand):
    newlist = str(category).split('.')
    if newlist[0]=="empty":
      if brand == "empty":
        return "unknown"
      return brand
    return newlist[0]

@udf
def extract_product(category, brand):
    newlist = str(category).split('.')
    if newlist[-1] == "empty":
      if brand == "empty":
        return "unknown"
      return brand
    return newlist[-1]

df_category_product_extracted = preprocessed_df.select("*", extract_category("category_code", "brand"),
 extract_product("category_code", "brand"))
df_category_product_extracted = df_category_product_extracted.withColumnRenamed("extract_category(category_code,
brand)","category").withColumnRenamed("extract_product(category_code, brand)", "product")
df_category_product_extracted = df_category_product_extracted.drop("category_code")
df_category_product_extracted.limit(5).toPandas()
```

**<u>Extract Date & Hour from Timestamp:</u>**

```python
df_time = df_category_product_extracted.withColumn('Date', split(preprocessed_df['event_time'], '
').getItem(0)).withColumn('Time', split(preprocessed_df['event_time'], ' ').getItem(1))
df_time = df_time.withColumn('Day', split(df_time['Date'], '-').getItem(2)).withColumn('Hour',
split(df_time['Time'], ':').getItem(0))
df_time = df_time.drop("Date")
df_time.limit(5).toPandas()
```

**Processed Dataset**

| event_time | event_type | product_id | category_id | brand | price | user_id | user_session | category | product | Time | Day | Hour |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2019-10-01 00:00:00 UTC | view | 44600062 | 2103807459595387724 | shiseido | 35.79 | 541312140 | 72d76fde-8bb3-4e00-8c23-a032dfed738c | None | None | 00:00:00 | 01 | 00 |
| 2019-10-01 00:00:00 UTC | view | 3900821 | 2053013552326770905 | aqua | 33.20 | 554748717 | 9333dfbd-b87a-4708-9857-6336556b0fcc | appliances | water_heater | 00:00:00 | 01 | 00 |
| 2019-10-01 00:00:01 UTC | view | 17200506 | 2053013559792632471 | None | 543.10 | 519107250 | 566511c2-e2e3-422b-b695-cf8e6e792ca8 | furniture | sofa | 00:00:01 | 01 | 00 |
| 2019-10-01 00:00:01 UTC | view | 1307067 | 2053013558920217191 | lenovo | 251.74 | 550050854 | 7c90fc70-0e80-4590-96f3-13c02c18c713 | computers | notebook | 00:00:01 | 01 | 00 |
| 2019-10-01 00:00:04 UTC | view | 1004237 | 2053013555631882655 | apple | 1081.98 | 535871217 | c6bd7419-2748-4c56-95b4-8cec9ff8b80d | electronics | smartphone | 00:00:04 | 01 | 00 |

# **OBJECTIVES**

The goal of the project is to use this processed clickstream data and big data techniques:

- To Analyze key performance indicators and find insights to improve revenue

- Provide insights for personalized digital marketing

- Measure marketing efforts of retailers, and optimize the overall user experience

Such analysis will help us gain insights and make critical business decisions. By using historical clickstream data, retailers can get an understanding of user behavior, what products the user is likely to purchase, etc. Hence, the objectives we've explored in this project are:

- Category Analysis
  - Determine best performing categories on the e-commerce site based on purchases
  - Find Brands that generate the highest traction in these best performing categories
- Effect of Adding to Cart
  - Correlation between impact Adding to Cart → Purchase: Cart Conversion Ratio
  - Evaluate Cart Abandonment Rate across categories and brands
- Effect of day-time on purchase trends
  - Analyze the Purchase trends across the month
  - Determine E-Commerce Prime Time
- Build a Real-Time classification model that predicts a purchase using clickstream features

Hence, in this project, we aim to use Big Data techniques on handcrafted features to find insights, correlations and key actions. We also work towards building a predictive modeling Machine Learning pipeline to boost purchases of items in real-time.

# TRACKING A USER'S JOURNEY

Querying for a specific user session, we can track a user's behavior through our website. According to the Figure below, we can infer that the user is interested in buying Apple smartphones. The user browses through various products with different prices, some viewing multiple times and eventually adding to cart and purchasing them.



Given below is the eCommerce Conversion Funnel that represents the user journey across 3 different event_types.



Figure 3: eCommerce Conversion Funnel

# OBJECTIVE 1 - CATEGORY ANALYSIS

We grouped data by category to analyze the performance of products and different brands. This can be the first step in attempting to gauge the user's interest and understand what actionable steps can the organization take to increase its sales.

**OBJECTIVE 1.a Determining the best performing categories based on purchases.**

**Techniques:** GroupBy, Filters, UDF, Count, Bar Plots

**Insights:** It is observed that Electronics, Appliances, and Computers are the most Browsed (both Viewed and Carted) and Purchased Categories. Moreover, Electronics is the highest revenue-generating category. We observed an interesting trend in Furniture, which is although highly browsed but has a low purchase rate. This reflects a general notion that people prefer in-person viewing of furniture-based items for comfort and first-hand experience.



Figure 4: Top Brands in Top Performing Categories



Figure 5: Top Brands in Top Performing Categories

**Actions:** Since electronics is the highest browsed as well as highest purchased category, allocation of required resources, and data should be done to ensure its performance. Onboarding of new vendors can be done to keep the revenue incoming. On the other hand, we need to analyze the root cause of conversion sales in apparel/furniture. Enhancing customer experience by introducing new 3D viewing techniques can be adopted.

**OBJECTIVE 1.b: Determining the best performing product in each of the best performing categories:**

```
df_cb_purchased_count = df_cat_purchased.groupBy("category","brand").count().orderBy(desc("count"))
window = Window.partitionBy(df_cb_purchased_count['category']).orderBy(df_cb_purchased_count['count'].desc())
ranked =  df_cb_purchased_count.select("*", rank().over(window).alias("rank")).filter(col('rank')≤5)
ranked.orderBy(desc("count")).show(5)
```

**Insights**: We went one step further in top-performing categories to analyze top-performing brands as shown in Fig-5. We observed that Samsung has a very strong brand presence across Electronics and Appliances with 159k purchases amounting to nearly 38% of all electronic purchases.

**Actions:** From these insights, we can take some concrete steps like onboarding more products/vendors under these categories and brands. Identifying the Brand's top user group for personalized targeted marketing can incentivize regular customers to increase their purchases. Tracking brand reviews, and performing sentiment analysis for brands that aren't getting good purchase conversion could potentially be one of the few measures to be taken.



Figure 6:  Top Brands in Top Performing Categories

# OBJECTIVE 2 - EFFECT OF ADDING TO CART

Here we analyzed the correlation between the impact of adding to cart → purchase: cart conversion ratio and evaluated cart abandonment rate across categories and brands.

**Objective 2.a: Correlation between impact Adding to Cart → Purchase: Cart Conversion Ratio**

**Insights:** It's observed that there is an 80% Cart Conversion Ratio (Cart: 926k, Purchase: 742k). We also observe from Fig-7.1 & 7.2 that smartphones are the most added to the cart as well as the highest purchased and have a 63% conversion rate. On the contrary, clocks have a 90% conversion ratio and we believe it is because there are fewer options to choose from, hence the specs are pre-determined leading to a purchase.



Figure 7.1:  Cart vs Purchase count



Figure 7.2:  Cart vs Purchase count(zoomed in)

**Actions:** Based on these insights, we can identify key reasons why products already in the cart don't get purchased. The reasons might be due to better deals, return policy, etc from other sites. It's useful to perform A/B testing when introducing a new product or feature in the products.

## Objective 2.b: Evaluate Cart Abandonment Rate across categories and brands

**Insights:** Further drilling down, we calculated cart abandonment rate for categories as well as for brands and observed that Construction and appliances have the highest cart abandonment rate. And although electronics was highly purchased and browsed, it also has a high Cart Abandonment Ratio of 37%. A similar analysis is made on brands and we see that Oppo, Huawei, Xiaomi, Samsung, and Apple have 46%, 44%, 43.5%, 43.2%, and 31% CAR respectively, clearly showing that the Apple brand is more trustworthy with a significantly lower CAR.

```python
#Creating a UDF for calculating Cart Abandonment Rate
@udf(returnType=FloatType())
def cart_miss_rate_udf(cart_count, purchase_count):
  rate = (1 - purchase_count/cart_count)*100
  return rate

df_cart_miss_rate = cart_purchase.select("*", cart_miss_rate_udf("cart_count", "purchase_count"))
df_cart_miss_rate = df_cart_miss_rate.withColumnRenamed("cart_miss_rate_udf(cart_count, purchase_count)",
"cart_miss_rate")
df_cart_miss_rate = df_cart_miss_rate.filter("cart_miss_rate>0")
df_cart_miss_rate = df_cart_miss_rate.filter("cart_count>5000")
df_cart_miss_rate = df_cart_miss_rate.orderBy(desc("cart_miss_rate")).limit(10)
```
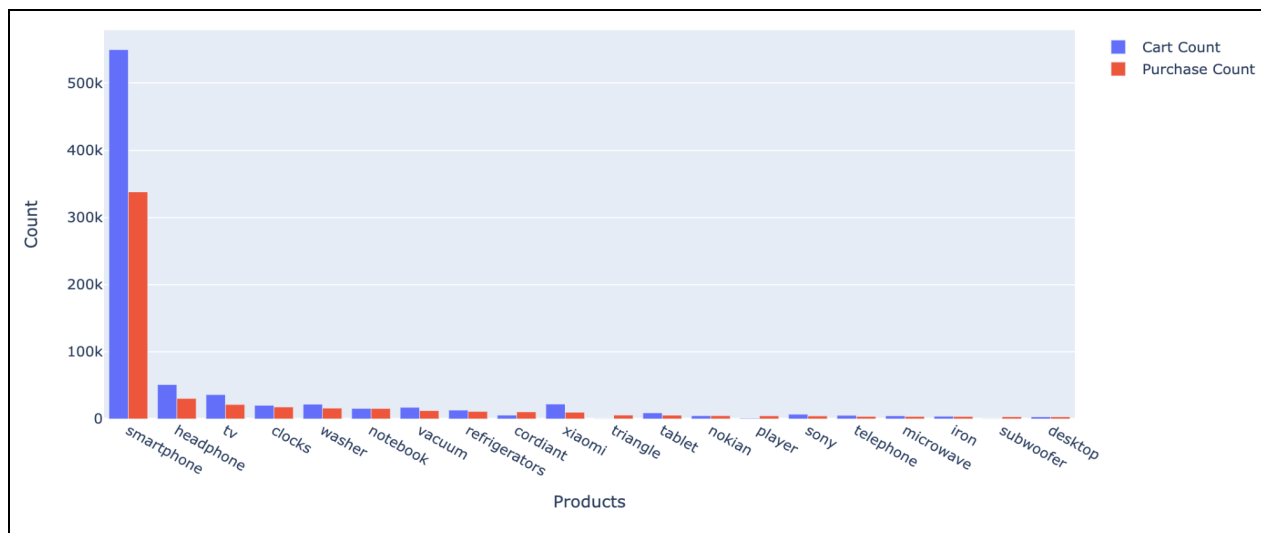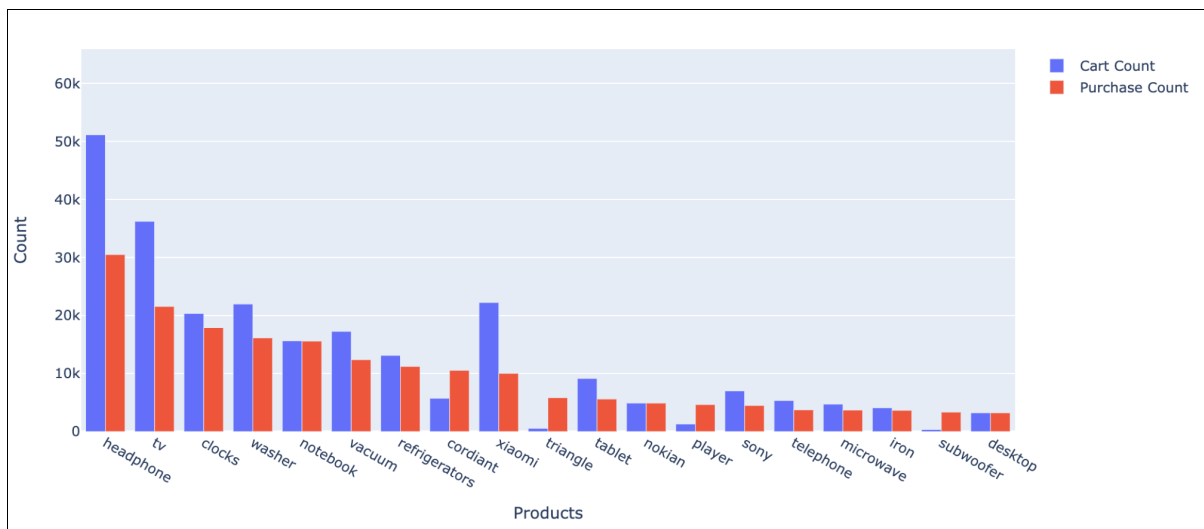
**Actions:** To curb this high cart abandonment, we need to thoroughly analyze the return options, shipping costs, and payment methods to understand why and what exactly is causing such items to go unpurchased. Some viable solutions can be partnering with vendors/brands to offer deals and discounts on products that are not converted. Finding specific user groups that purchase these products and provide personalized marketing to customers from these user groups who abandon the same
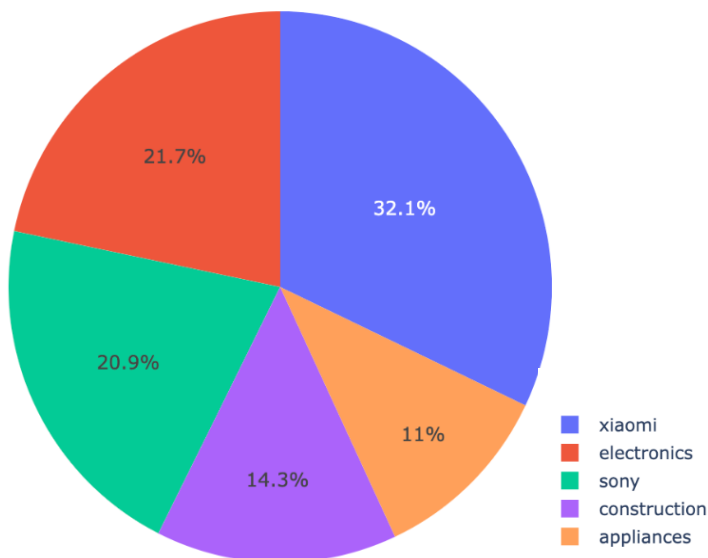


Fig 8.1 - Cart Abandonment Rate by Category



Fig-8.2 - Cart Abandonment Rate by Brands

# OBJECTIVE 3 - EFFECT OF DAY-TIME WITH PURCHASE

**Objective 3.a: Finding e-commerce Prime Time**

```python
df_purchase = df[df['event_type'] == "purchase"]
df_purchase_date_count = df_purchase.groupby("Day").count()
count = [val[0] for val in df_purchase_date_count.select('count').collect()]
date = [val.Day for val in df_purchase_date_count.select('Day').collect()]
plt.figure(figsize=(20, 8))
plt.bar(date, count)
```

**Insights:** We continue our analysis of the purchasing pattern of users in October. The bar plot in Fig-9, depicts that the sale of items is highest around the middle of the month which can be attributed to a couple of reasons like generally salaries get credited at the beginning of the month and people are figuring out expenses like rent, loans, etc.

**Actions:** Offering mid-month sale/discount from day 11 until 17 would act as a catalyst to increase sales. Further, this analysis can be scaled across all the months to identify peak traffic thereby giving users lucrative offers. It would be tremendously helpful to do such an analysis throughout the year to gauge a similar pattern in a monthly fashion.



Fig: 9 - Purchase Trends in October.

**Objective 3.b: Finding e-commerce Prime Time**

**Insights:** We also focussed on analyzing shopping patterns per day to understand peak shopping time as shown in Fig 10. Almost 2M users have already accessed the e-commerce site by 11:00. We observed that it's generally around evening that users are actively browsing. There is also a high purchase in the morning denoted by the green bar, and it's significantly low through the tail end.

**Actions:** Time-bound flash sale/discount would be an actionable item to coerce users to complete their purchase journey.



Fig-10 -Peak Traffic during different hours of the day.

# OBJECTIVE 4: Build A Real-Time Machine Learning Classification Model To Predict Purchase

For Objective 4, we aimed to build a binary classification model that can predict whether there would be a purchase in a given user session. Through our analysis, we see that there is a high conversion ratio once the user adds a product to the cart. Features like br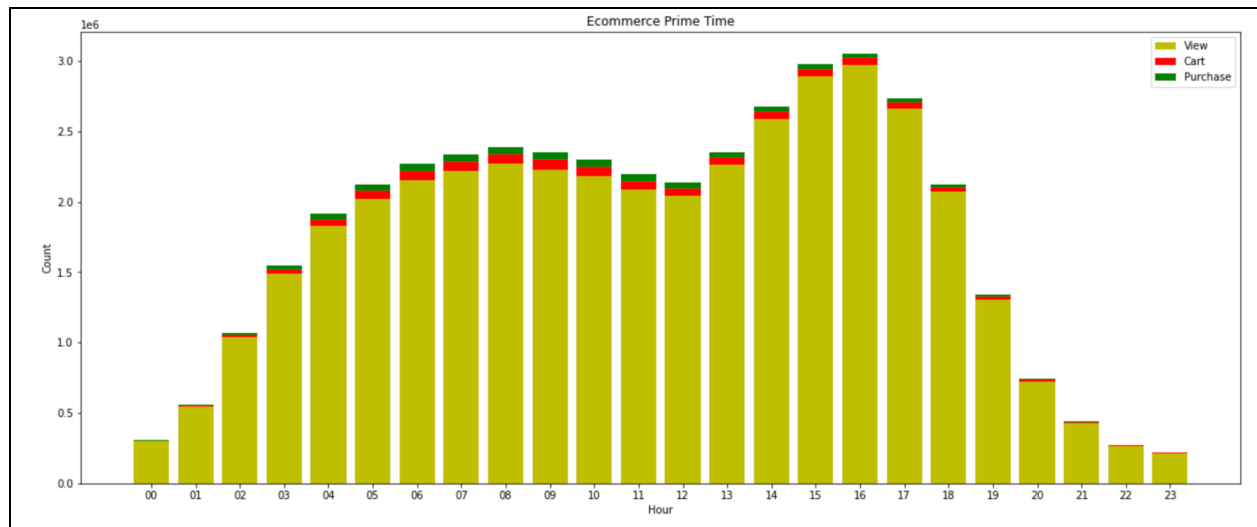and, price, and how many times carted, have a direct impact on the purchasing trend of a customer. Through our predictive analysis pipeline, based on ordinal features like brand, category, product_type, and numerical features like price, week of the month, and activity count, we get the prediction result.
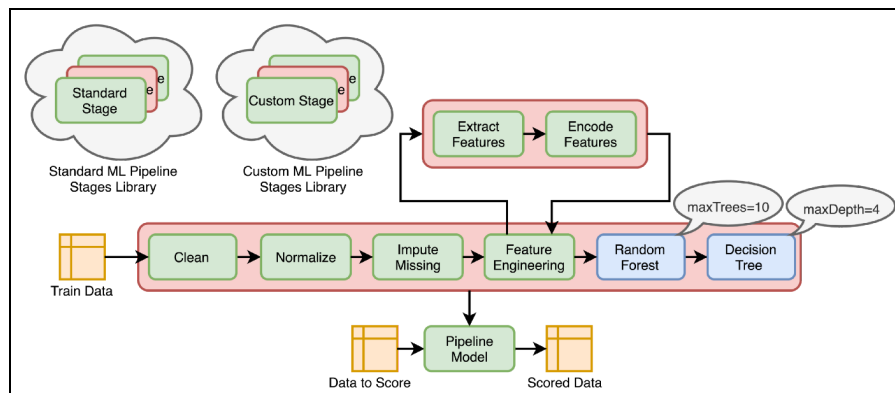


Fig-11 -ML pipeline for predictive modeling

PySpark ML was used for building the pipeline. We empirically tested a few classification models available in Pyspark like Logistic Regression, and XGBoost, and selected the random forest model which gave the highest accuracy of 78.42%. To simulate the inference pipeline deployment in the real world, we built an interactive web app using streamlit for demonstration.

# DEMO USING STREAMLIT

Streamlit is a free, open-source, all-python framework that enables data scientists to quickly build interactive dashboards and machine learning web apps. For our use case, we've used Streamlit to demo a mockup of a real-time inference engine that can be plugged into an existing retail analytics solution that captures clickstream data.

# CONCLUSION

Clickstream analysis is used by leading Retailers to make important business decisions using Big Data techniques. It also requires specific skills and resources necessary to capture, collect and analyze this information. As we saw through our project there are many benefits of such a detailed analysis such as optimizing user routes customers take to reach a page or to make a purchase. Deeper insights into consumer behavior like how visitors get to the website; what they do once there; how long they stay on a page; the number of page visits visitors make; and the number of unique and repeat visitors can help in the development of quantifiable metrics to boost revenue generation. Such analysis can also help in running narrowly-targeted marketing campaigns. In short, we observed that

*"Tracking customer behavior in an online store is instrumental to offering a personalized customer experience and selling the right products in the right way"*

# FUTURE SCOPE

When solving any problems for Big Data, the first and most critical step is to actually understand what is in the data, which is difficult because of the size of the data. When we first started looking at the data, it was difficult to get a full picture quickly. Even a simple count or filter on such huge data using pandas failed. This helped in appreciating big data technologies in a practical scenario. There was also a paradigm shift in our understanding and implementation from centralized and functional programming to distributed and parallel computing. This project not only provided us with hands-on experience on a practical problem but also made us technically skilled to tackle such problems in real-world scenarios. This work can be extended to incorporate more components like the integration of databases and message queues, to ingest dynamic data in real-time. We also limited ourselves to pyspark, but modern state-of-the-art solutions like Dask, Fueger, and Ray could also be explored to enhance performance.

Github Link - https://github.com/niharikakrishnan/Retail-Clickstream-Analysis-and-Prediction