

# Repeat Buyers Prediction

**CECS 550 : Pattern Recognition  
Spring 2023**

**Group :- 7**

<b>Anthony Martinez</b>	<b>012500759</b>
<b>Diksha Patil</b>	<b>030849765</b>
<b>Pratik Jadhav</b>	<b>030880471</b>
<b>Pavan More</b>	<b>030873555</b>
<b>Sudarshan Powar</b>	<b>030845787</b>

**Instructor: Prof Mahshid Fardadi**

**Teaching Assistant: Rahul Deo Vishwakarma**



CALIFORNIA STATE UNIVERSITY  
**LONG BEACH**



# Agenda

- Background and Goals
- Data Interpretation and Visualization
- Feature Engineering and ranking
- Prediction models
- Performance
- Recommendations
- Conclusion



# Background and Goals

1. A shop runs big promotions on “Double 11” - the biggest online shopping event, in order to attract a large number of new buyers.
2. Unfortunately, many of the attracted buyers are one-time deal hunters, and these promotions may barely have long-lasting impact on sales.
3. To reduce the promotion cost and enhance the return on investment (ROI), they want to identify who can be converted into repeated buyers.

## Goals

1. Predict the probability of the given user becoming a repeat buyer of the given merchant in the future
2. To find the most important factor to predict repeat buyers



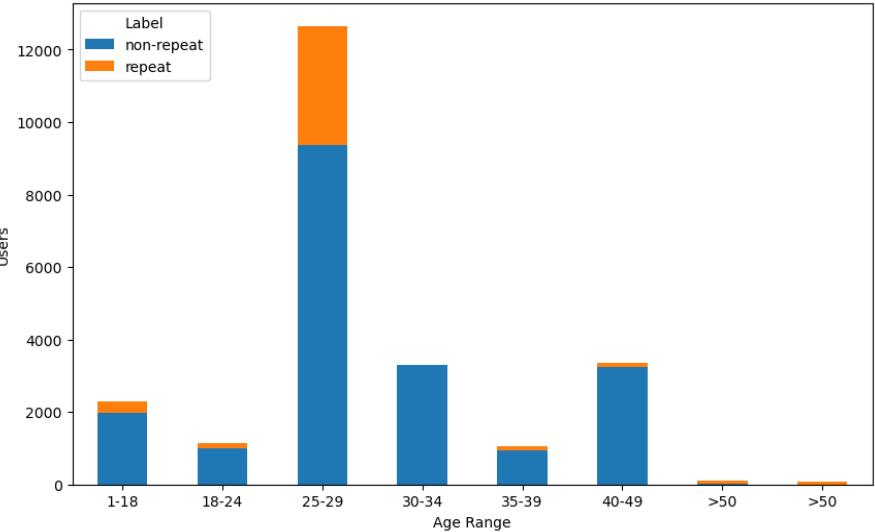
# Dataset

- The merged dataset from the given data profiles : user\_info, user\_log, train\_data

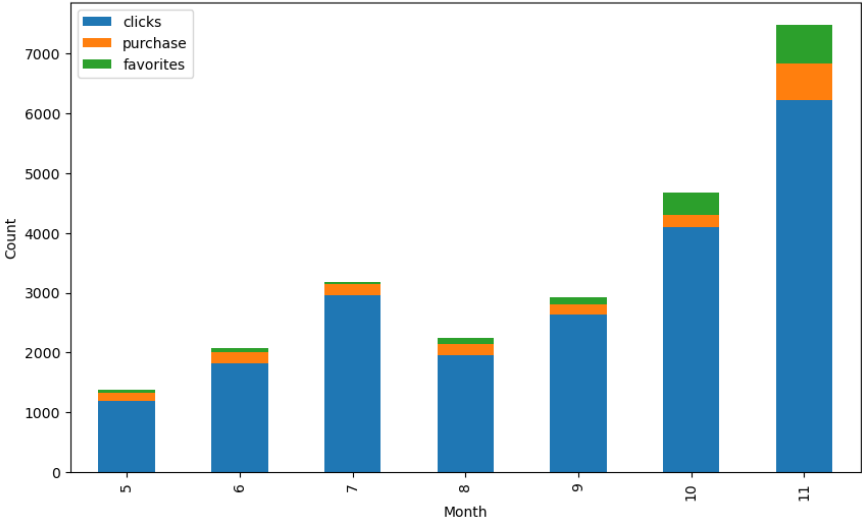
	user_id	merchant_id	label	item_id	cat_id	seller_id	brand_id	time_stamp	action_type	age_range	gender
<b>8291</b>	18306	4733	0	1014866	1397	127	6434	701	0	3	2
<b>19281</b>	99459	2032	0	162933	1389	4963	1991	1002	0	6	0
<b>2893</b>	18306	1710	1	66013	1577	3877	3213	1005	0	3	2
<b>22182</b>	238467	4966	0	935412	1349	1629	2292	821	3	6	1
<b>10805</b>	149634	742	0	236773	1505	416	4014	812	0	6	0

# Data Visualization

Number of Users by Age Range and Label

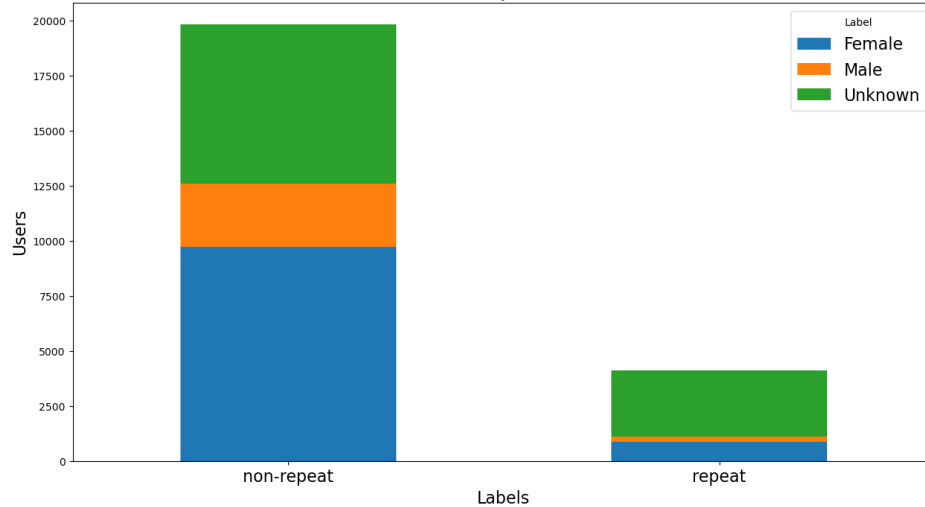


Counts of actions by month and type

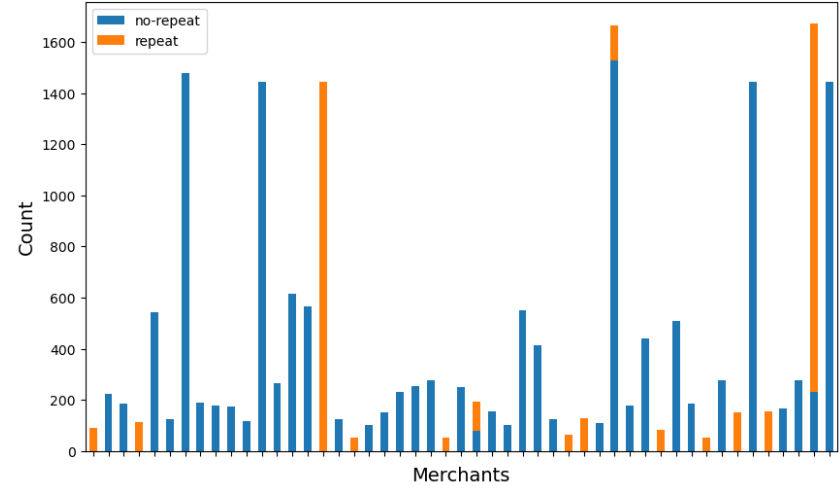


# Data Visualization

Number of Users by Gender and Label

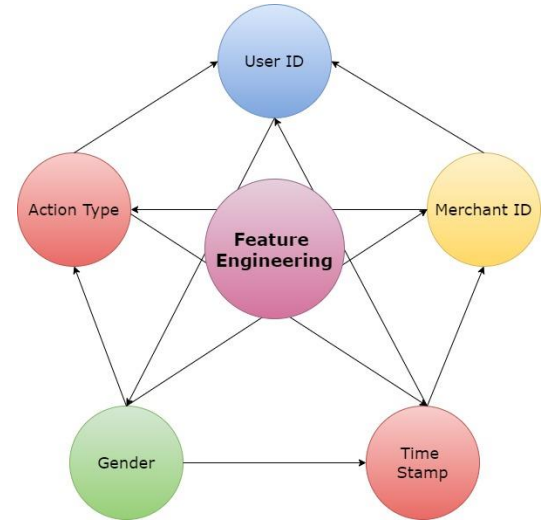


Count of actions by Merchants and label



# Feature Engineering

- Average User Age for each Category
- Purchase Average Time
- Purchase Ratio
- Purchase Frequency
- Average User Age for each Merchant
- Ratio of add-to-cart actions to clicks for each Merchant
- Number of distinct brands a user has interacted with for each Merchant
- Number of distinct categories a user has interacted with for each Merchant



# Feature Ranking

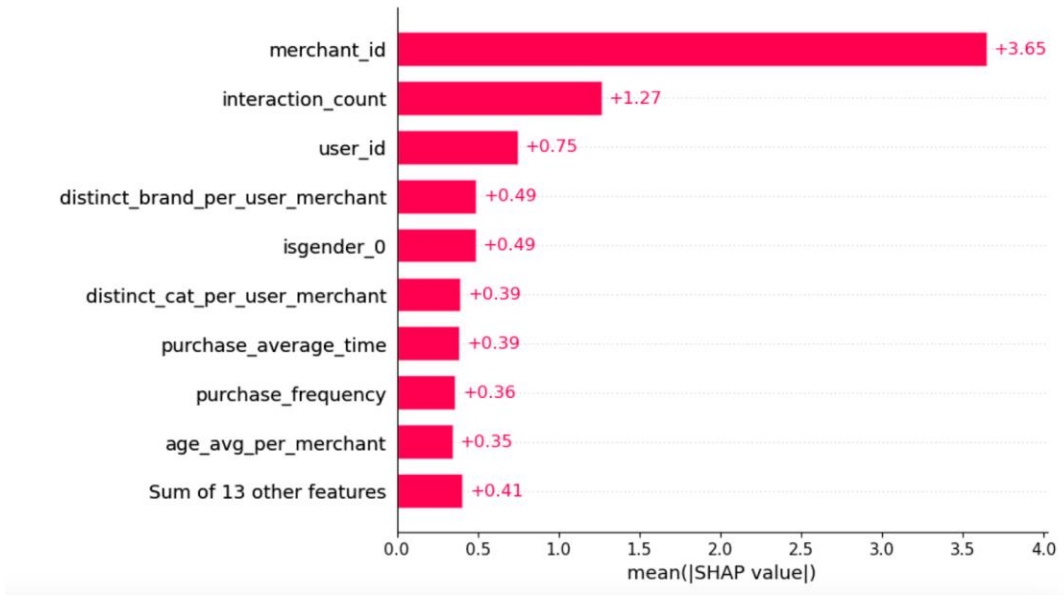
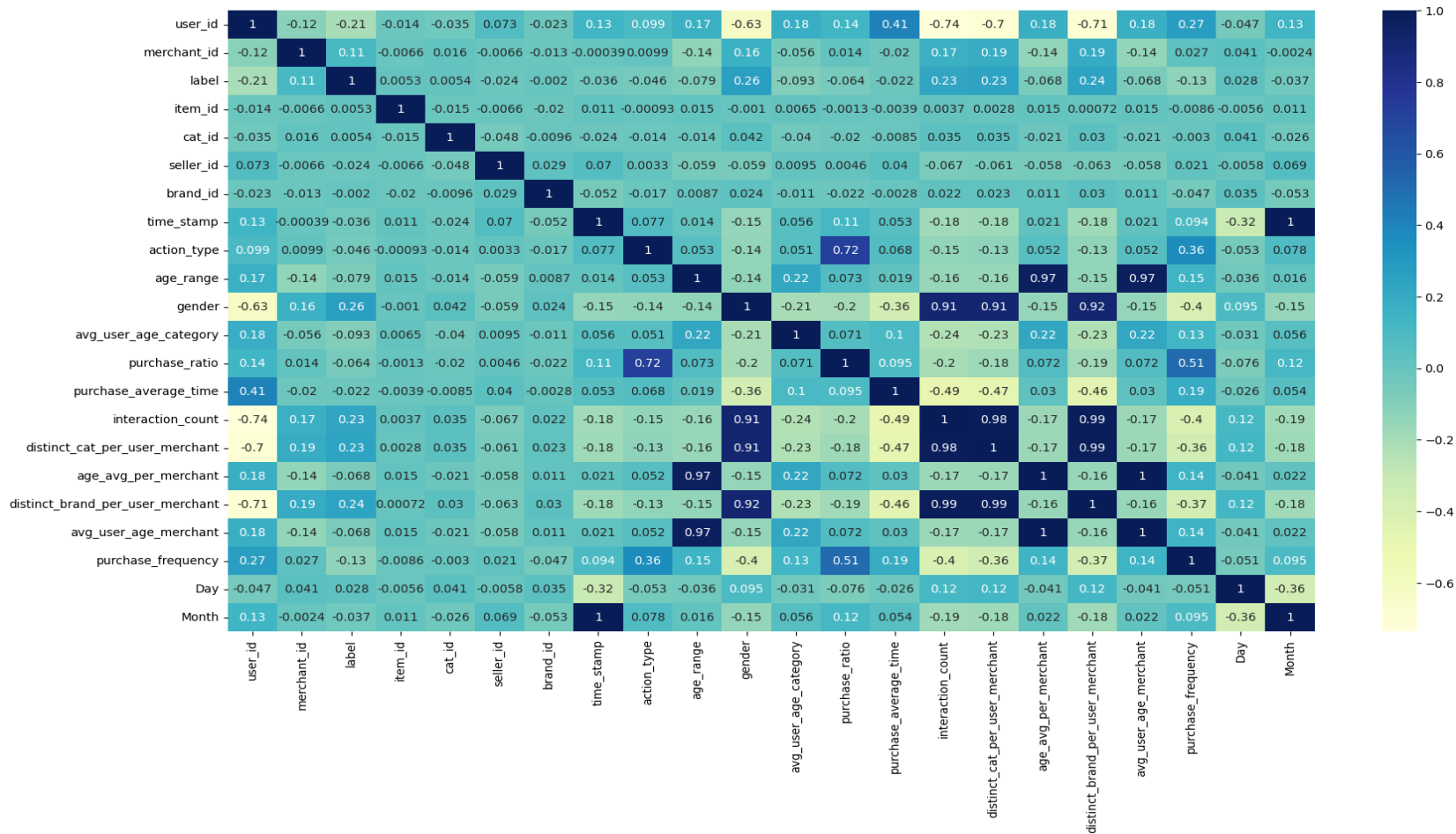


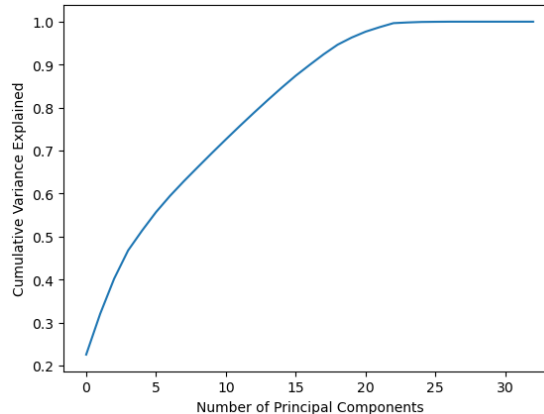
Figure 1: Important features based on SHAP





# Principal Component Analysis (PCA)

- We performed principal component analysis (PCA) on the data to reduce its dimensionality.
- We determined the optimal number of components by selecting the smallest number that captured at least 95% of the variance in the original data. 20 components in our case.
- We then used the PCA-transformed data to train the Bayesian Gaussian classifier . We found that this approach resulted in a much better accuracy of 75% from 63%.



# Prediction Models

## Splitting of Data

- Training Set : 80% of data
- Validation Set : 20% of Training Set
- Testing Set : 20% of data

## Trained on the following models

- Bayes classifier
- Random Forest
- KNN
- Neural Networks

# Hyperparameter Tuning

- Hyperparameter tuning involves selecting the optimal values for model parameters that are set before training to improve the model's performance. Such as loss functions, learning rate, activation functions and optimizers.
- By tuning hyperparameters, we can improve a model's performance on unseen data and avoid overfitting or underfitting.
- A validation set is used to evaluate the model during hyperparameter tuning, which helps in selecting the best set of hyperparameters that generalize well to new data.



## Bayes classifier:

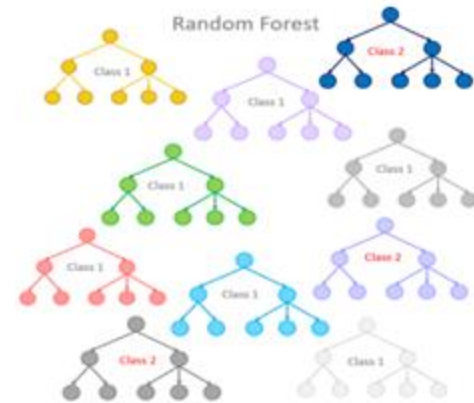
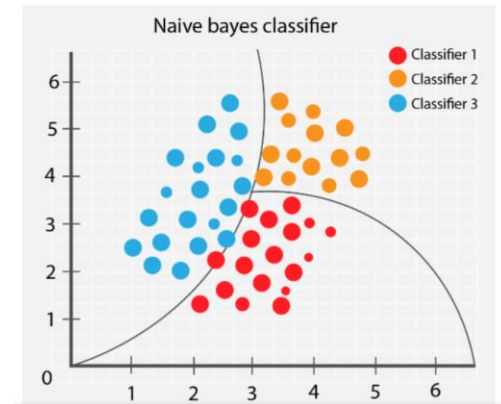
- Bayes' theorem is a fundamental principle in probability theory that describes the probability of an event based on prior knowledge or information.

Model	Accuracy	F1 Score
Bayes	0.75	0.36

## Random forest:

- Different classifiers overfit the data in a different way, and through voting those differences are **averaged out**.

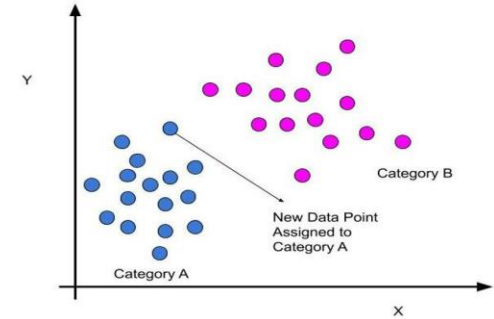
Model	Accuracy	F1 Score
RFC	0.86	0.58



## KNN:

- KNN with  $k=1$  and distance metric 'manhattan' performed best.
- F1-score was low due to the data being unbalanced.
- After **upsampling** on training data the metrics improved.

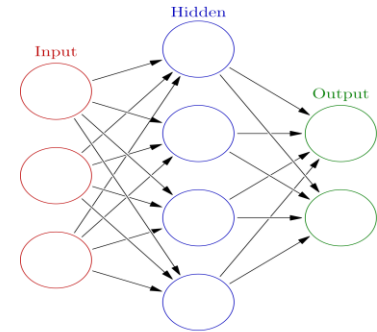
Model	Accuracy	F1 Score
KNN	0.84	0.56



## Neural Network

- The optimal model has 4 dense layers with 'relu' activation and an output layer with 'sigmoid' activation function.
- Binary Cross Entropy loss showed the best accuracy with 'adam' optimizer.
- **Dropout** layers were used for regularization.

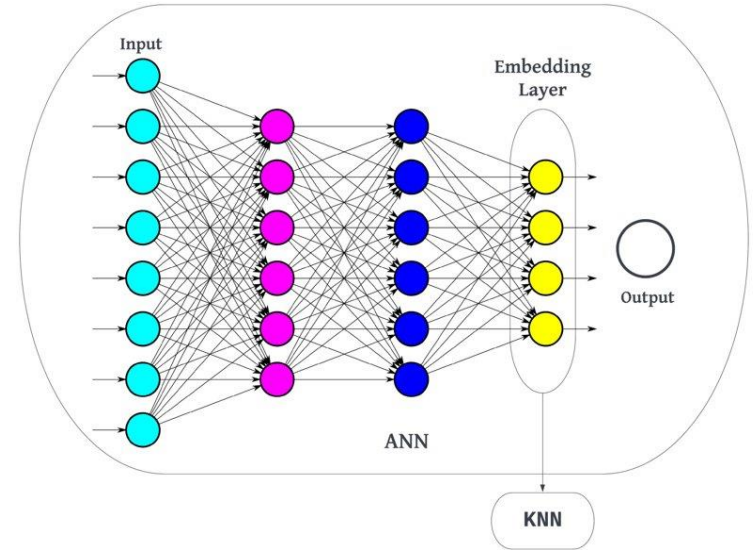
Model	Accuracy	F1 Score
ANN	0.94	0.85



## ANN combined with KNN

- The ANN model was combined with KNN.
- ANN was trained on the data and the penultimate layer(embedding) was then fed to the KNN.
- ANNs and KNN have complementary strengths and weaknesses. ANNs are good at learning complex non-linear patterns in the data, while KNN is good at capturing the local structure of the data.

Model	Accuracy	F1 Score
ANN + KNN	0.96	0.90



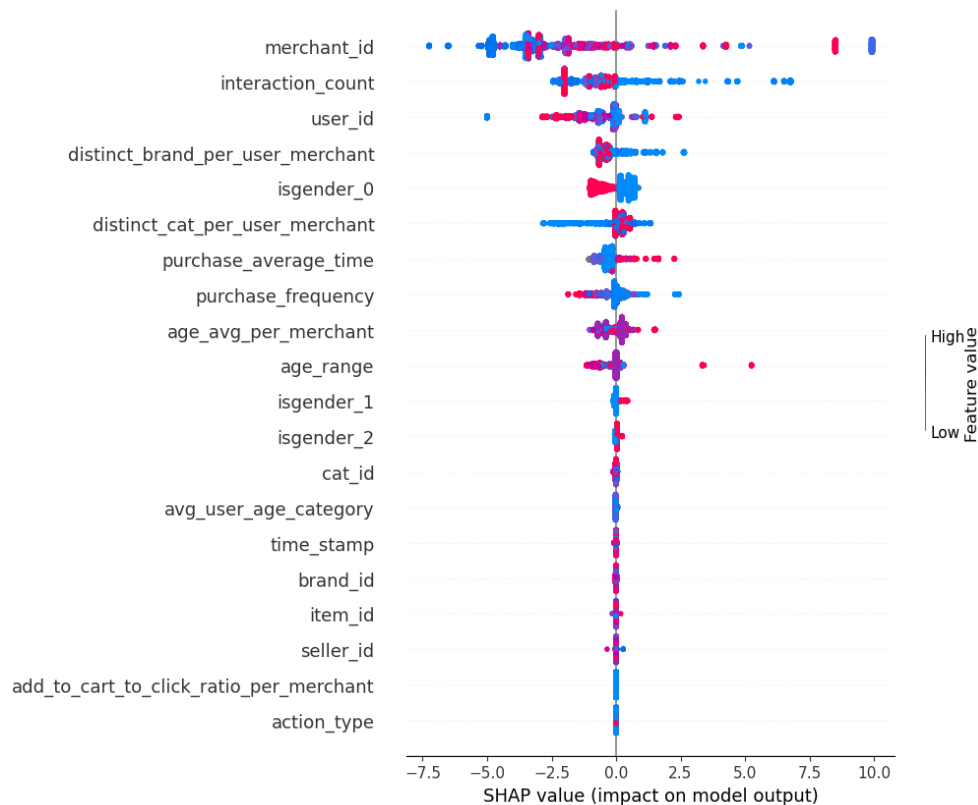
# Models Performance

Model	Accuracy	F1 Score
Bayes	0.75	0.36
RFC	0.86	0.68
KNN	0.84	0.56
ANN	0.94	0.85
ANN + KNN	<b>0.96</b>	<b>0.90</b>





# Results

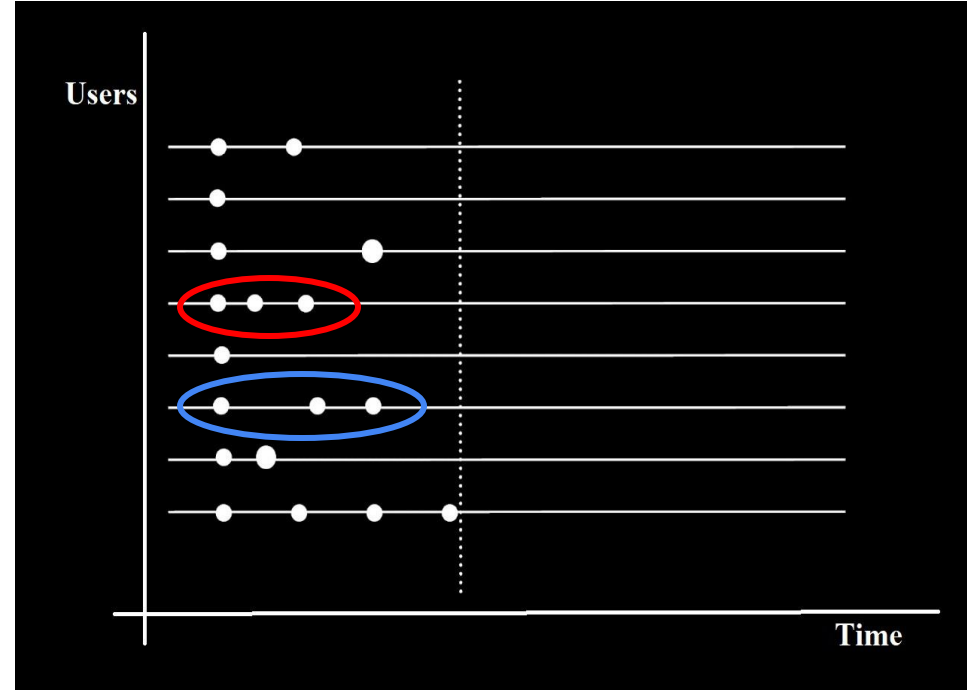
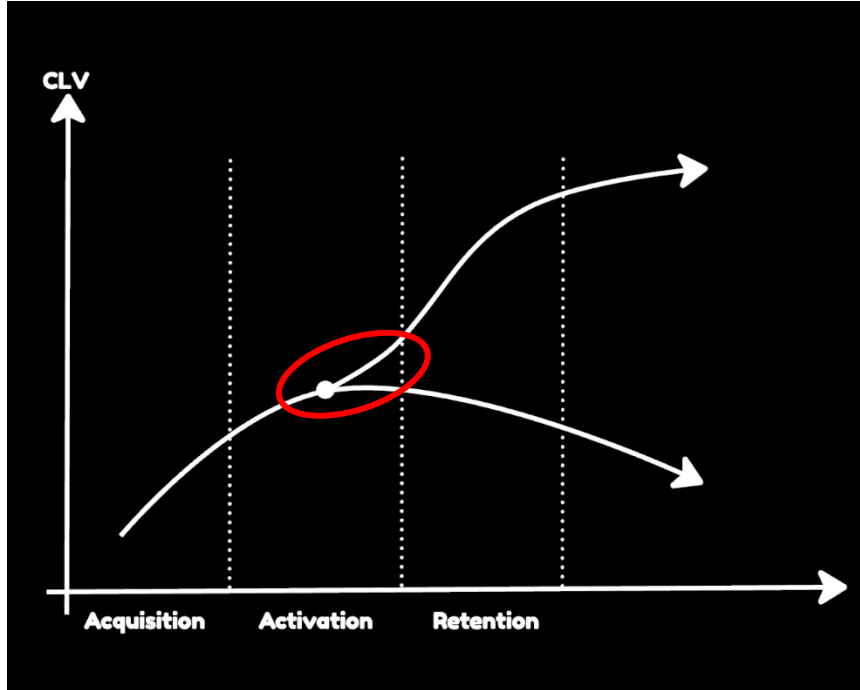


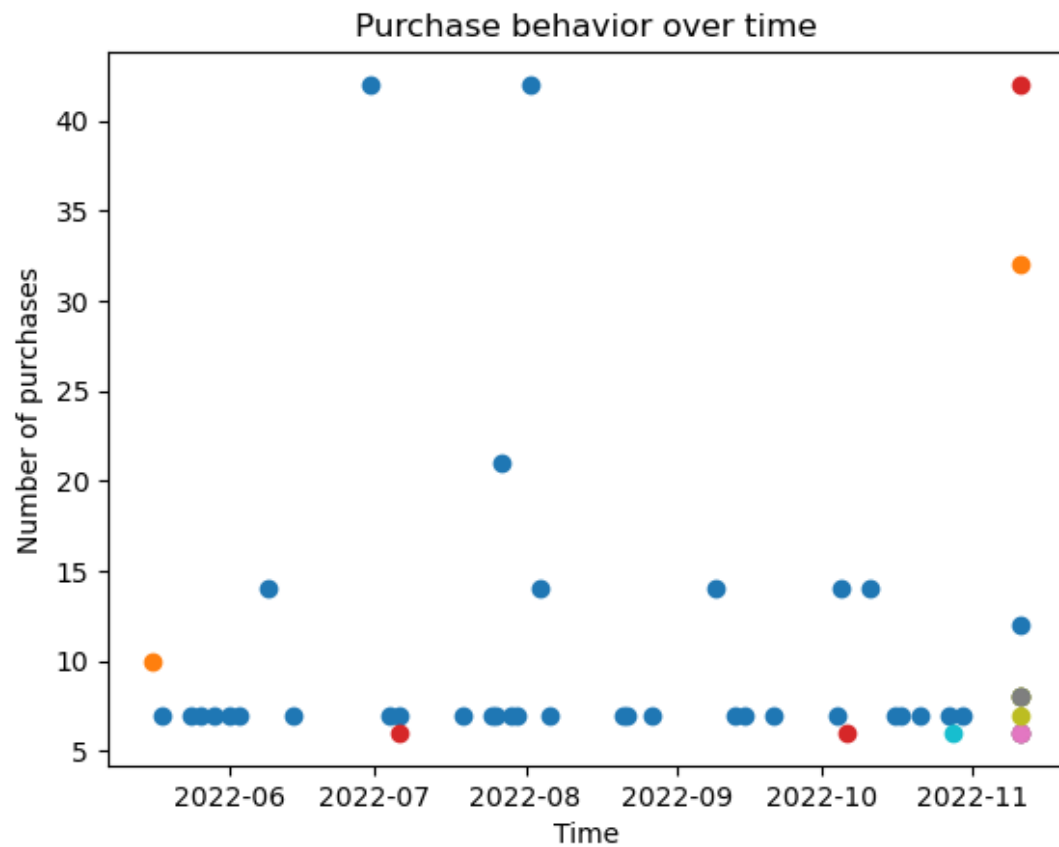
# Recommendations

- **Improve user engagement**
  - businesses can consider offering personalized recommendations based on the user's browsing and purchase history.
- **Segments users by gender and age range**
  - businesses can tailor their offerings and marketing strategies to specific groups
- **Increase purchase frequency**
  - businesses can offer promotions or incentives to encourage users to make repeat purchases.
- **Personalize marketing strategies**
  - businesses can personalize their marketing strategies to target specific users with products or promotions that are most likely to appeal to them.



# Recommendations





# Conclusion

- The methods predict whether a user would be a repeat customer or not by training different models on a dataset containing user interactions with merchants.
- The model ANN combined KNN achieved the highest accuracy and F1-score.
- The most important features for predicting repeat customers were found to be interaction\_count, merchant\_id, and gender.

Overall, the study provides valuable insights for businesses looking to improve customer retention.

