



DSA3101 ASSIGNMENT 2

6 OCTOBER 2020

MEET THE TEAM!



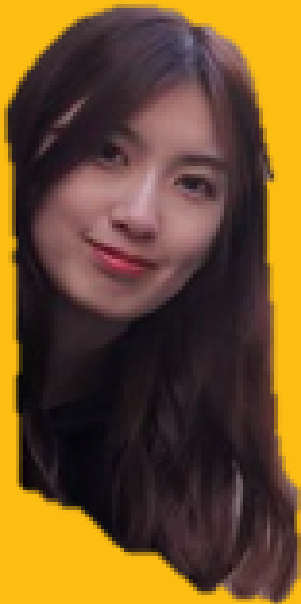
CHAN YU HANG



CHUA CHENG LING



CHUA HUAREN



CAI AN QI




JOASH CHIN




CLARENCE ONG


ASSIGNMENT 2



CUSTOMER
SEGMENTATION



PRINCIPAL
COMPONENT
ANALYSIS



COLLABORATIVE
FILTERING

DATA IMPUTATION

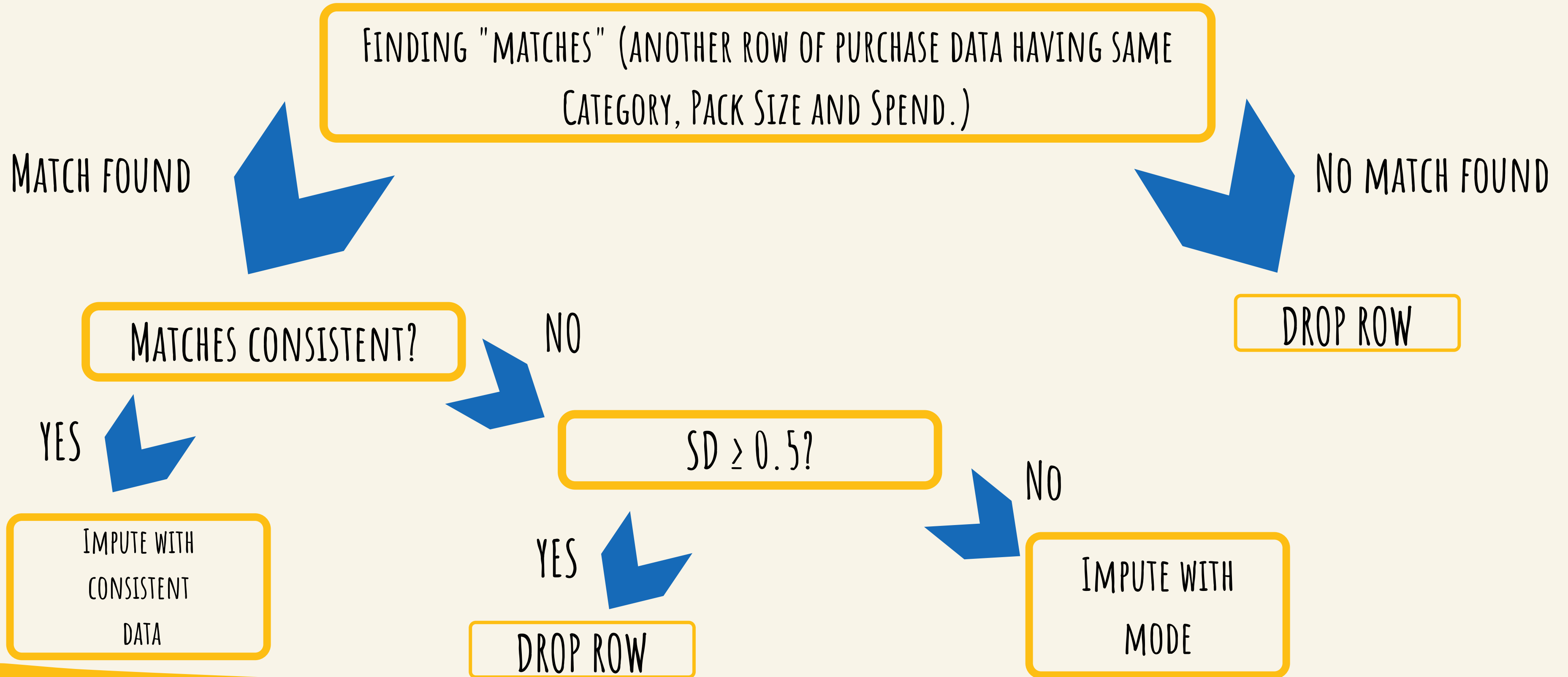
PREVIOUSLY IN ASSIGNMENT 1, THERE WERE 49,539 INSTANCES OF MISSING/MIS-RECORDED DATA. OUR TEAM TRIED TO USE KNN IMPUTATION, HOWEVER WE WERE UNABLE TO DO SO AND CHOSE TO REMOVE THE DATA INSTEAD AS IT WAS ONLY 3.8% OF THE DATASET.



DATA IMPUTATION

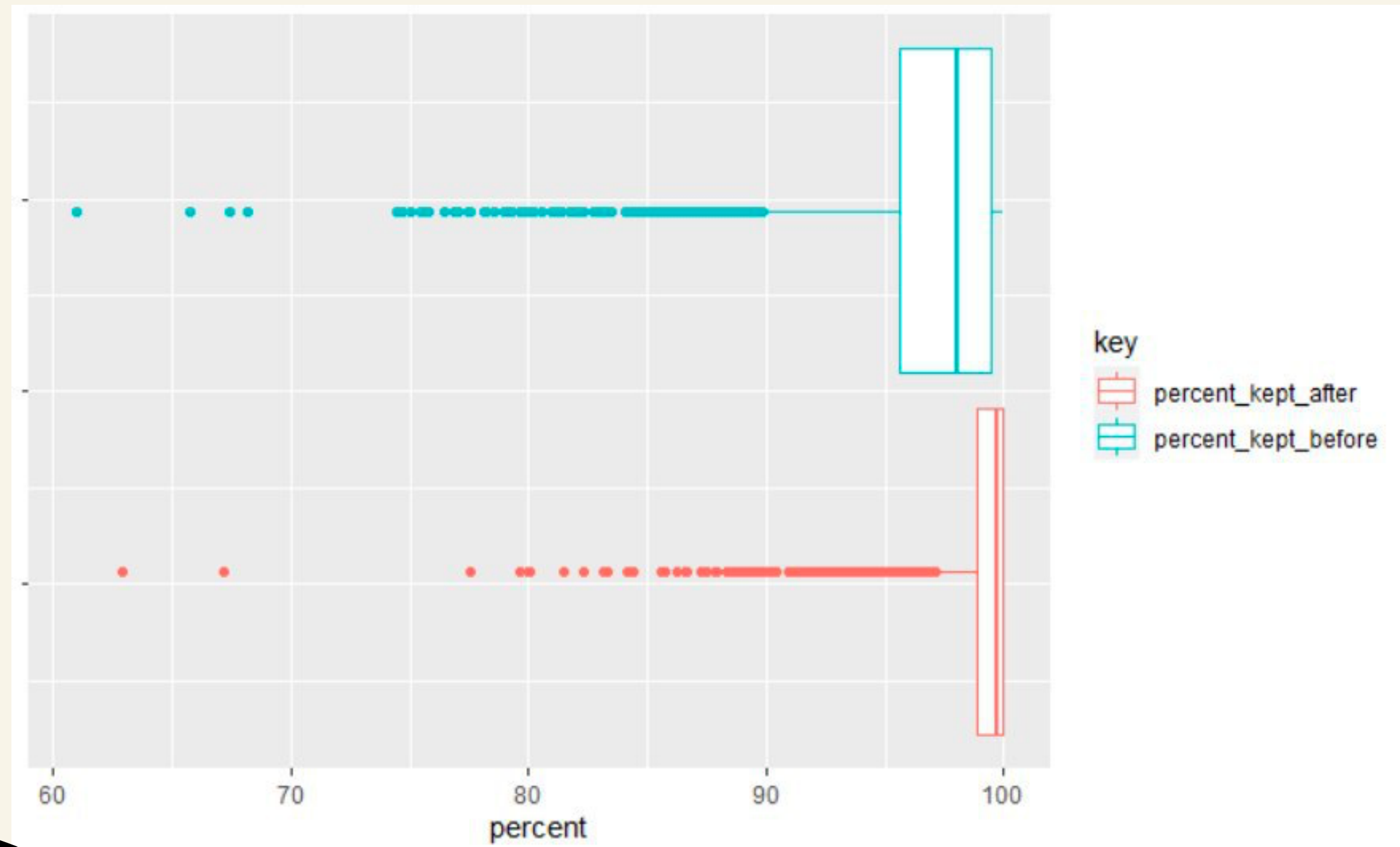
FOR ASSIGNMENT 2, WE DECIDED TO DO DATA IMPUTATION BECAUSE
REMOVING ALL INCORRECT DATA WOULD DROP AROUND 207 CUSTOMERS
(LOSING MORE THAN 10% OF THEIR PURCHASE DATA), WHICH COULD AFFECT
THE ACCURACY OF SUGGESTIONS ESPECIALLY IF THERE IS SOME SIGNIFICANCE
TO THE MISSING DATA

CONSERVATIVE APPROACH TOWARDS DATA IMPUTATION



EFFECTIVENESS?

THE MEAN INCREASE IN USABLE
ROWS IS 9.1%!



CUSTOMERS' LOCATION

NORTH

SOUTH

“CENTRAL”

EAST COAST



CENTRAL

CSD

Rice

....

Belacan

CUSTOMER A



....



CUSTOMER B



....



CUSTOMER SEGMENTATION - PRINCIPAL COMPONENT ANALYSIS

- LOW CORRELATION BETWEEN THE VARIOUS COMPONENTS IN THE INPUT VECTORS, MAKING IT DIFFICULT TO USE PCA TO SIGNIFICANTLY REDUCE THE NUMBER OF DIMENSIONS USED TO REPRESENT EACH CUSTOMER EFFECTIVELY

Top Absolute Correlations

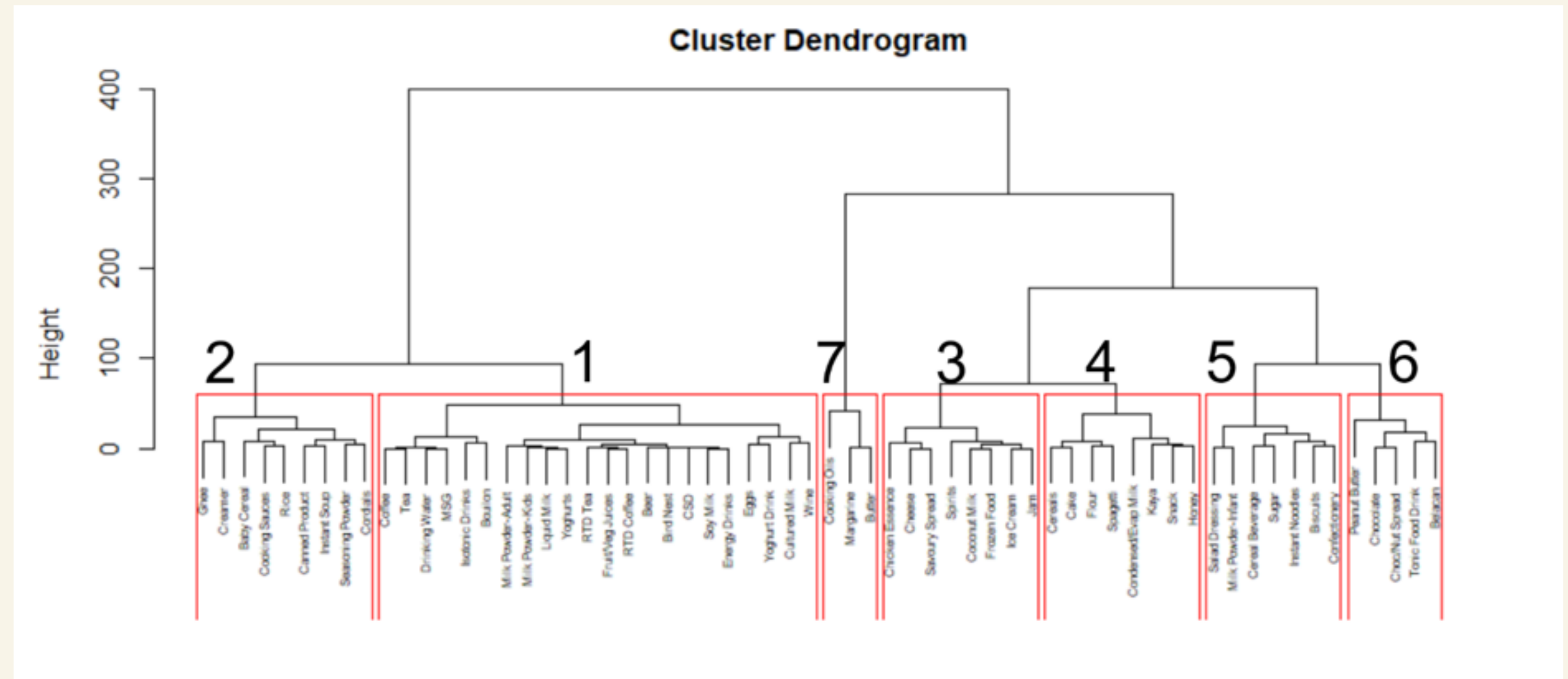
Category	Category	
Flour	Sugar	0.406735
Beer	Isotonic Drinks	0.404076
Condensed/Evap Milk	Sugar	0.399785
Cooking Oils	Rice	0.391161
Isotonic Drinks	RTD Tea	0.37523
Biscuits	Liquid Milk	0.36821
CSD	Isotonic Drinks	0.36497
Biscuits	Tonic Food Drink	0.35215
	Cooking Oils	0.35134
	Instant Food	0.348424
Beer	RTD Tea	0.345429
Cheese	Yoghurts	0.344703
Cooking Oils	Tonic Food Drink	0.343524
Yoghurt Drink	Yoghurts	0.326345
Rice	Sugar	0.325497
Salad Dressing	Spagetti	0.324520
Cheese	Instant Soup	0.323259
Beer	CSD	0.322699
Cooking Oils	Cooking Sauces	0.322407
Eggs	Sugar	0.322243

“ HIERARCHICAL CLUSTERING ”

- IDEA: SEGMENT THE DIFFERENT CATEGORIES INTO 7 LARGER OVERARCHING CATEGORIES BASED ON THE CALORIES PER 100G OF EACH OF THE 62 ORIGINAL CATEGORIES
- AGGLOMERATIVE MANNER



CLUSTERS WITH A SMALLER LABEL
NUMBER VALUE REPRESENT
CATEGORIES WITH LOWER VALUES
OF CALORIES PER 100G



PRINCIPAL COMPONENT ANALYSIS

USING $K = 2$ (PC1 AND PC2), OUR TEAM HAS IDENTIFIED THAT:

1. PC1 IS HAS ITEMS OF RELATIVELY HIGH CALORIES LIKE CHOCOLATE, PEANUT BUTTER AND ICECREAM WHICH CAN BE CONSIDERED AS COMFORT FOOD.
2. PC2 CAN BE REPRESENTED BY ITEMS OF LOW CALORIES LIKE SOY MILK, FRUITS AND VEGETABLE, AND ITEMS THAT ARE RICH IN CARBOHYDRATES, LIKE BISCUITS, INSTANT NOODLES AND CEREAL BEVERAGES

“

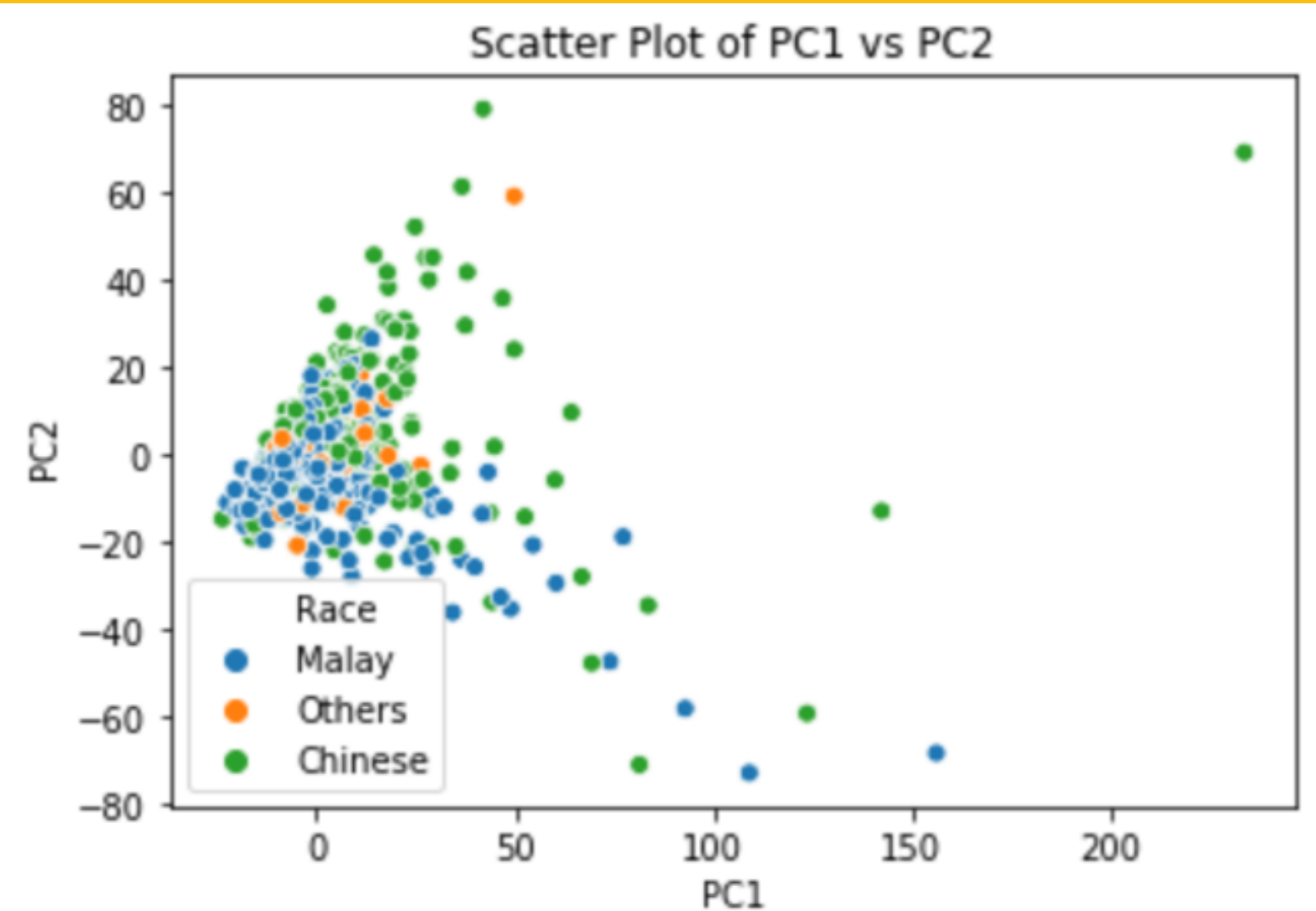


	PC1	PC2
rank		
1		0.827071
2		
3		
4		
5		0.791783
6	0.711918	
7	0.741586	

”

VERY IMPORTANT FINDING!

THERE SEEMS TO BE A RELATIONSHIP BETWEEN RACE AND HIS/HER PURCHASING HABITS!



FOR MALAYS, THERE IS A
NEGATIVE LINEAR
RELATIONSHIP BETWEEN
THE 2 PC

FOR CHINESE, THERE ARE
VARYING PURCHASE
PATTERNS!

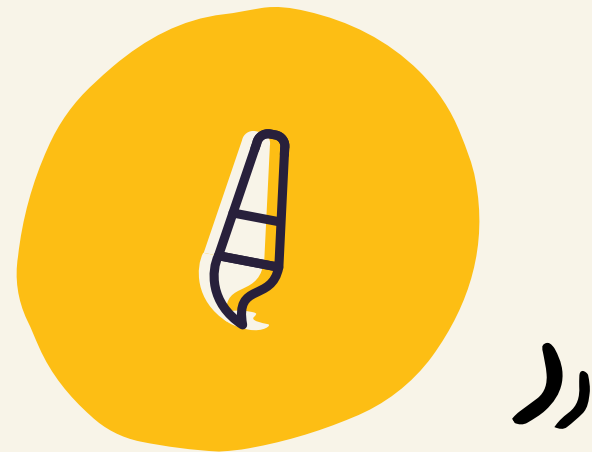


IDEA!

SINCE WE HAVE ESTABLISHED A STRONG RELATIONSHIP BETWEEN A
CUSTOMER'S RACE AND THEIR PURCHASING HABIT, LET'S APPLY IT
TO OUR RECOMMENDATION ENGINE!

PROBLEM WITH THE DATA

WE ARE ONLY PROVIDED WITH THE 62 CATEGORIES. USING THESE CATEGORIES WILL NOT PROVIDE ANY MEANINGFUL INSIGHT.



SOLUTION

CREATE ARBITRARILY UNIQUE ITEMS CATEGORIZED BY THE SPEND/VOLUME RATIO. WE GENERATED A LIST OF 36,237 ITEMS.



WEB APPLICATION



HISTORICAL
PURCHASES




DISCOUNTS
AND PROMOS



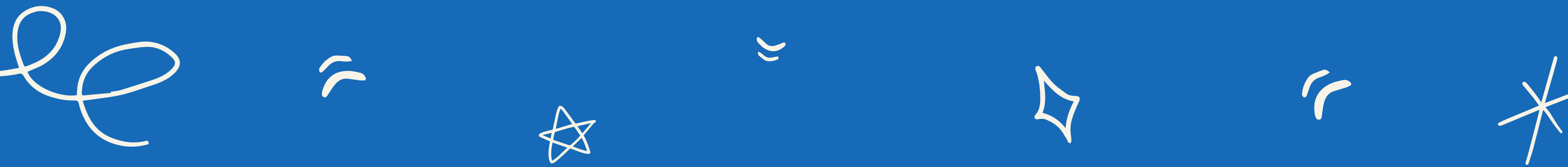
HISTORICAL
VIEWS



TRENDING
PURCHASES



COLLABORATIVE
FILTERING



NAIVE MODELS

GLOBAL MEAN TO
PREDICT EVERY
ITEM'S RATINGS

RMSE = 1.59

PANEL ID'S MEAN TO
PREDICT FOR THE
CORRESPONDING PANEL
ID'S ITEMS' RATINGS

RMSE = 1.56

ITEM'S MEAN TO
PREDICT FOR THE
CORRESPONDING
ITEM'S RATINGS

RMSE = 1.55

RMSE FOR COLLABORATIVE FILTERING = 1.52

CONCLUSION

GIVEN SLIGHTLY MORE INFORMATION (EITHER KNOWING THE
PANEL/PRODUCT MEAN) CAN OBTAIN AN EVEN LOWER RMSE

SUMMARY



DATA IMPUTATION

SINCE RECO ENGINE IS MOSTLY ABOUT CUSTOMER SATISFACTION, WE DO NOT WANT TO RISK LOSING THE 10% OF CUSTOMERS WHICH COULD BE COSTLY



PRINCIPAL COMPONENT ANALYSIS

RACE HAS A STRONG CORRELATION ON THE PURCHASING HABITS OF THE CUSTOMER.



RECOMMENDATION ENGINE

- DEPLOYABLE AND PROVIDES KEY INSIGHTS
- EMA USED SO MODEL LEARNS WITH TIME
- COLLAB FILTERING MODEL (MATRIX FACTORISATION) OUTPERFORMS NAIVE MODELS



THANK YOU!