



instacart [Project Name]

Project Name: Instacart Grocery Basket Analysis

Date: January 17, 2024

Analyst Name: Sarah Boller

Contents:

Population Flow

Consistency checks

Wrangling steps

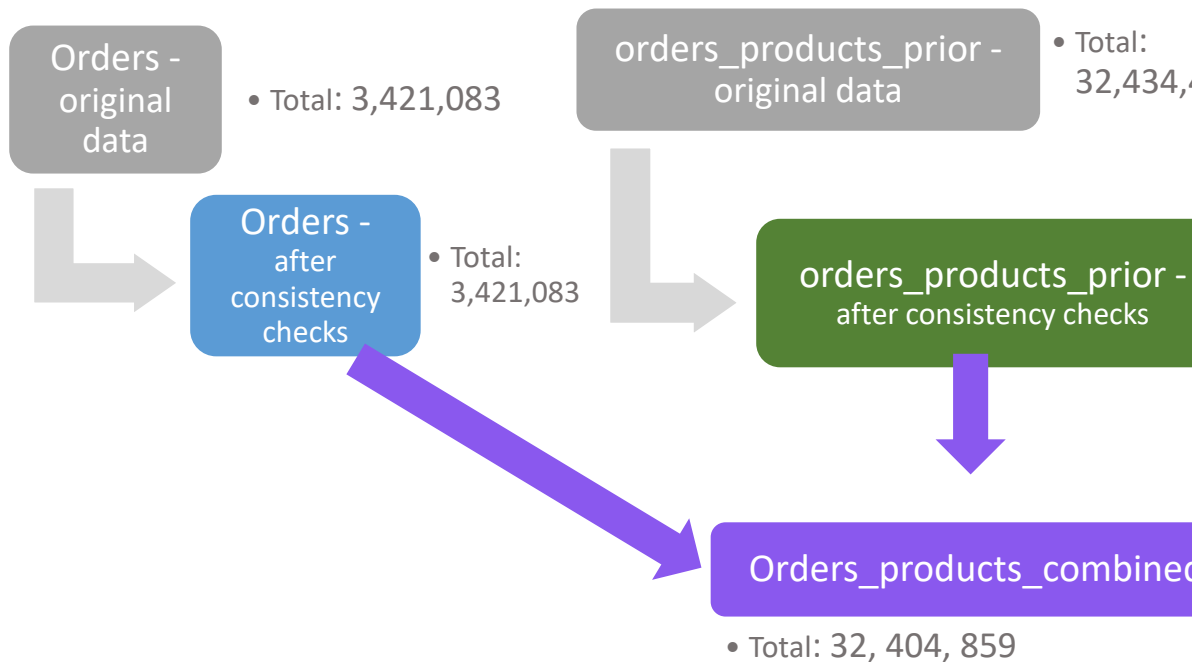
Column derivations

Visualizations

Recommendations



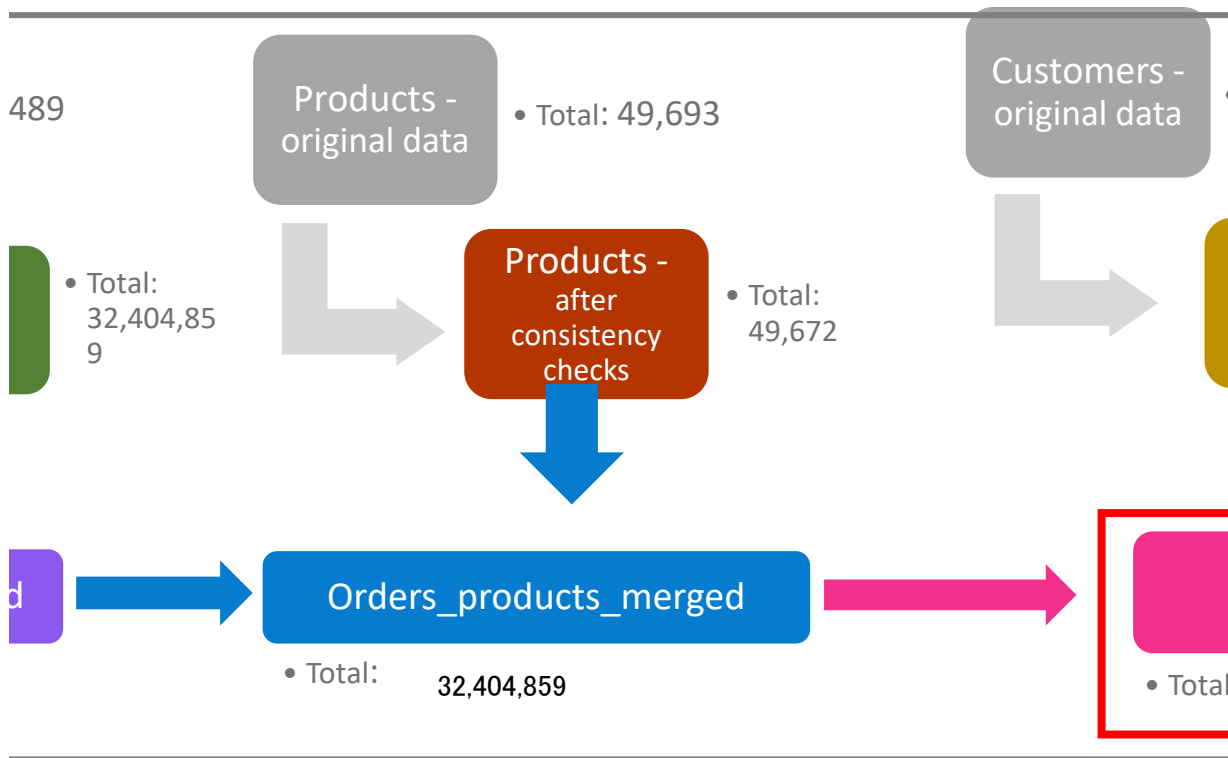
Population flow



1.) The grey boxes in the first row of the population flow represent the original data sets imported into Jupyter.

2.) The second row of boxes (coloured) represents the data sets **after** you have conducted these operations. This offers a visual overview of how the data flows.

3.) The third row, where also the arrows are coloured, represents the merge of the two datasets into the final dataset (in the red box). Keep in mind the final dataset is the result of the merge.



nal data sets as they were when you downloaded them. In the Total fields you ne

anipulated them, e.g., removed missing values and duplicates. In the Total fields
ows throughout the data consistency checks.

is you performed between the datasets. In the Total fields you need to add the c
aset should be without exclusions (based on the exclusion flag).

- Total: 206, 209

Customers -
after
consistency
checks

- Total:
206,209



Orders_products_all

l: 32,404,859

Exclusion flag

Condition: max_order <

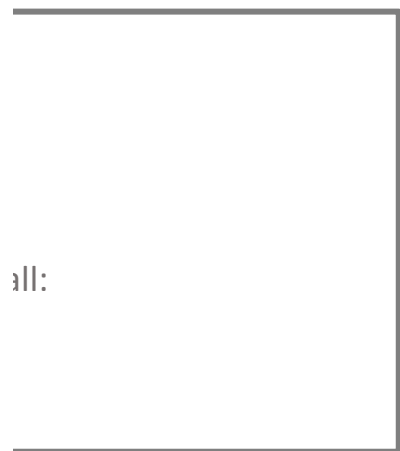
Observations to be removed:

Final total count of order_products_a

eed to add the count of the rows when you

you need to add the count of the rows after

ount of the rows in the merged datasets, so that





Consistency checks

[illegible]

[illegible]



Wrangling steps

Columns dropped	Columns renamed	Columns' type changed
eval_set		
	order_dow -> order_day_of_week	order_day_of_week
		user_id
	n_dependants -> n_dependents	
	Gender -> gender	
	STATE -> state	
	Age -> age	
First Name		
Surnam		
	order_hour_of_day -> time_of_order_hr	
add_to_cart_order		
aisle_id		

Comment/Reason
unnecessary to data
int -> string
miss spelled
Consistancy
Consistancy
Consistancy
PII
PII
clarity
irrelevant
irrelevant

Column derivations and aggregations

Dataset	New column	Column/s it was derived from
final with profiles	returning customer	days_since_prior_order
final with profiles	loyalty_flag	max_order
	mean_expenses	prices, user_id
final with profiles	Spender_flag	mean_expenses
	median_freq	days_since_prior_order
final with profiles	freq_flag	days_since_prior_order
final with profiles	income_flag	income
final with profiles	parent_flag	n_dependents
final with profiles	Age_Group	age
final with profiles	customer_profiles	age_group, income_flag, parent_fl

Conditions
if nan returning customer is False. If there is a number returning customer is True
if max_orders is more than 40, then set Loyal customer
else if max_orders is more than 10 and less than equal to 40, then set Regular customer
else max_orders is less than equal to 10, set New customer
mean prices per user_id
if mean_expenses is less than 10, then set Low spender
else if expenses is more than equal to 10, then set High spender
median of days_since_prior_order per user_id
if median_days_prior_purchase is more than 20, then set Non-frequent customer
else if median_days_prior_purchase is more than 10 and less than equal to 20, then set Regular customer
else median_days_prior_purchase is less than equal to 10, set Frequent customer
If income <= \$67,584, then low income
If income > \$67,584 & < \$127,912, then regular income
if income >= \$127,912, then high income
If n_dependents >= 1 then parent
Else, not a parent
If age <= 33, then Age_Group 18-33
If age > 33 & <= 49, then 34-49
If age > 49 & <= 65, then 50 to 65
Else, Over 65
Age_Group_income flag_ Parent_flag

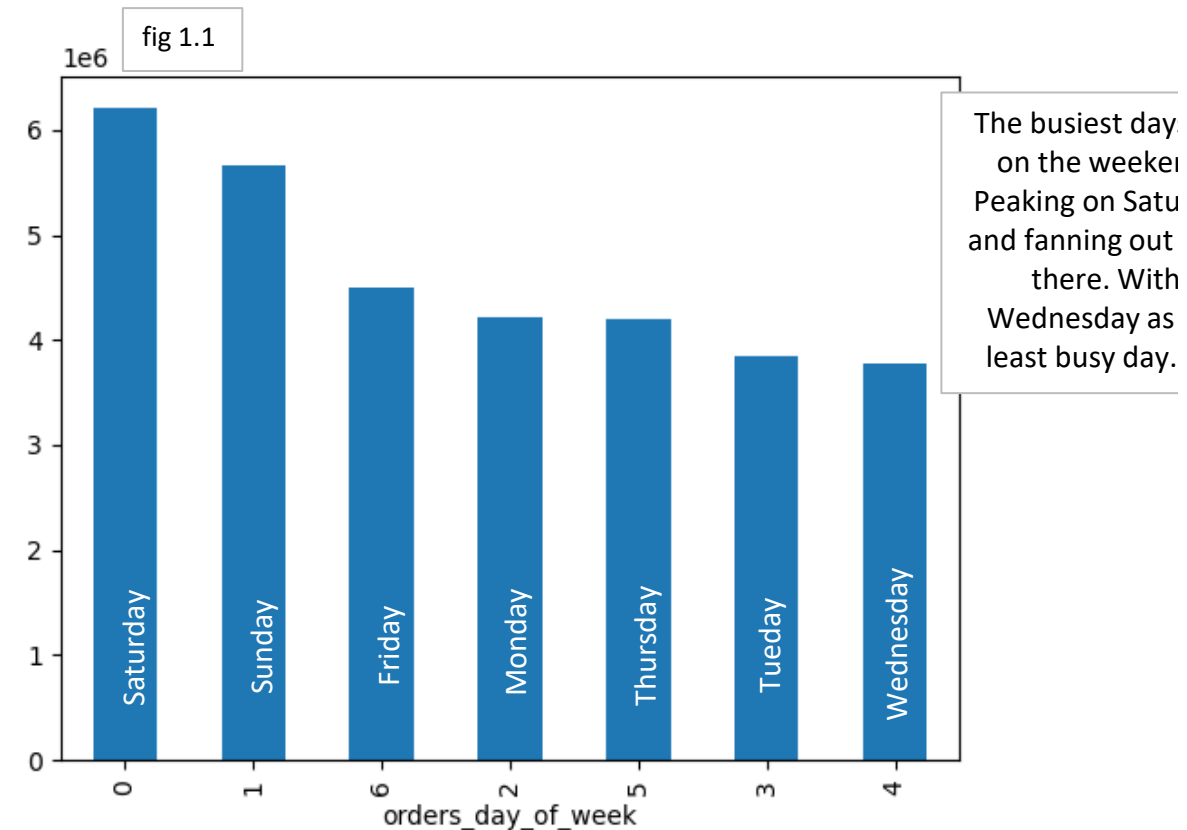
In this tab
them. NB:

Title page

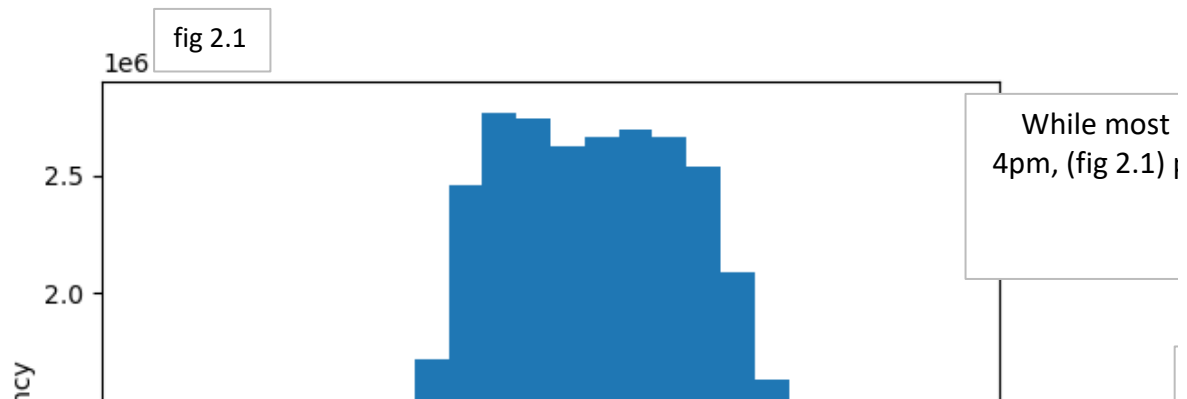
› you should add tables below that the frequencies of flags/label variables that you produced after deriving
: don't do this for continuous variables, only for flags.

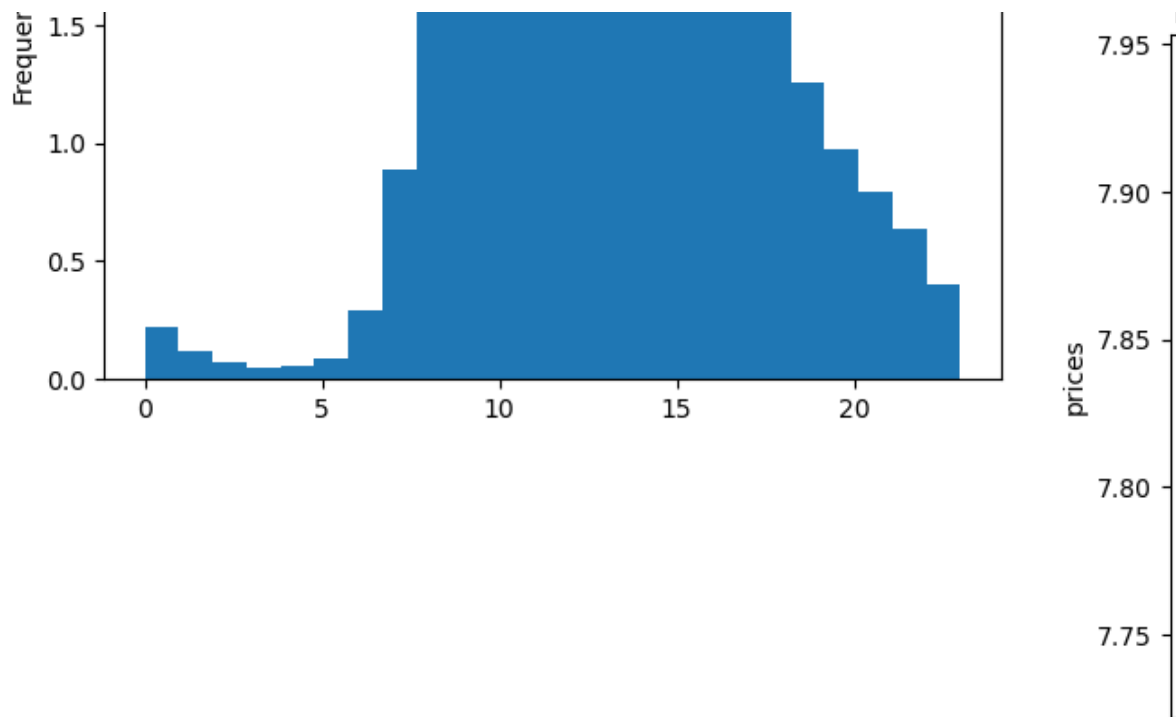
Visualisations

The sales team needs to know what the busiest days of the week and hours (orders) in order to schedule ads at times when there are fewer orders.

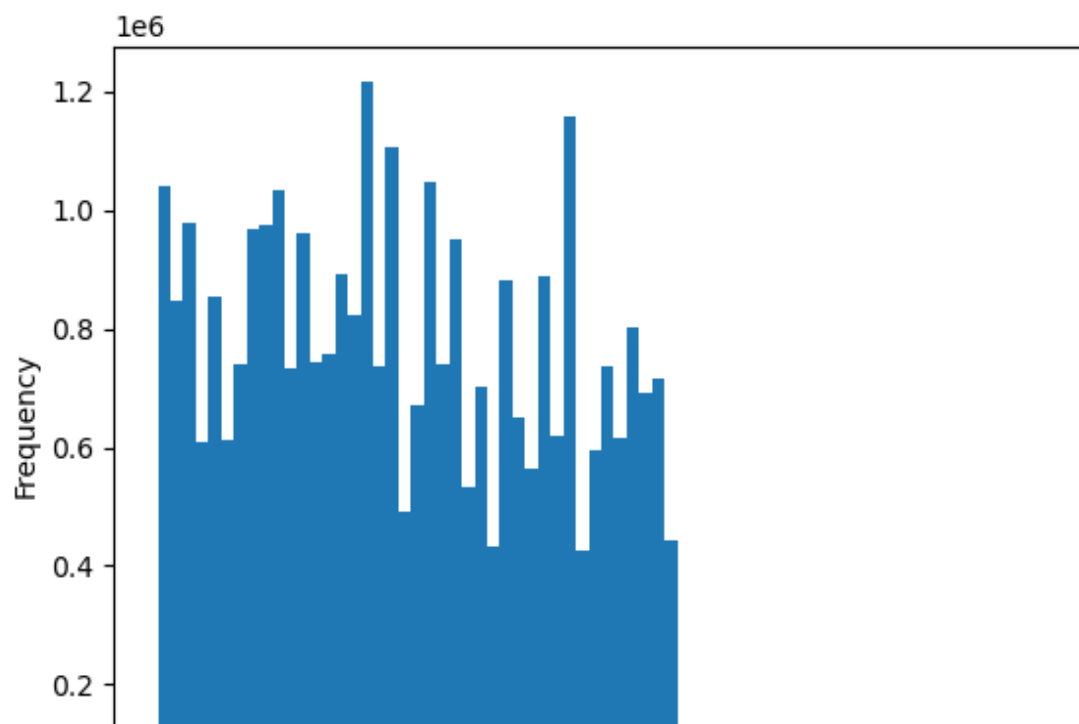


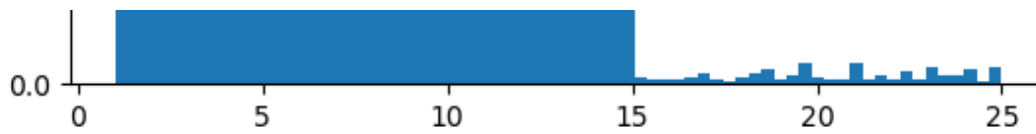
They also want to know whether there are particular times of the day when inform the type of products they advertise at these times.



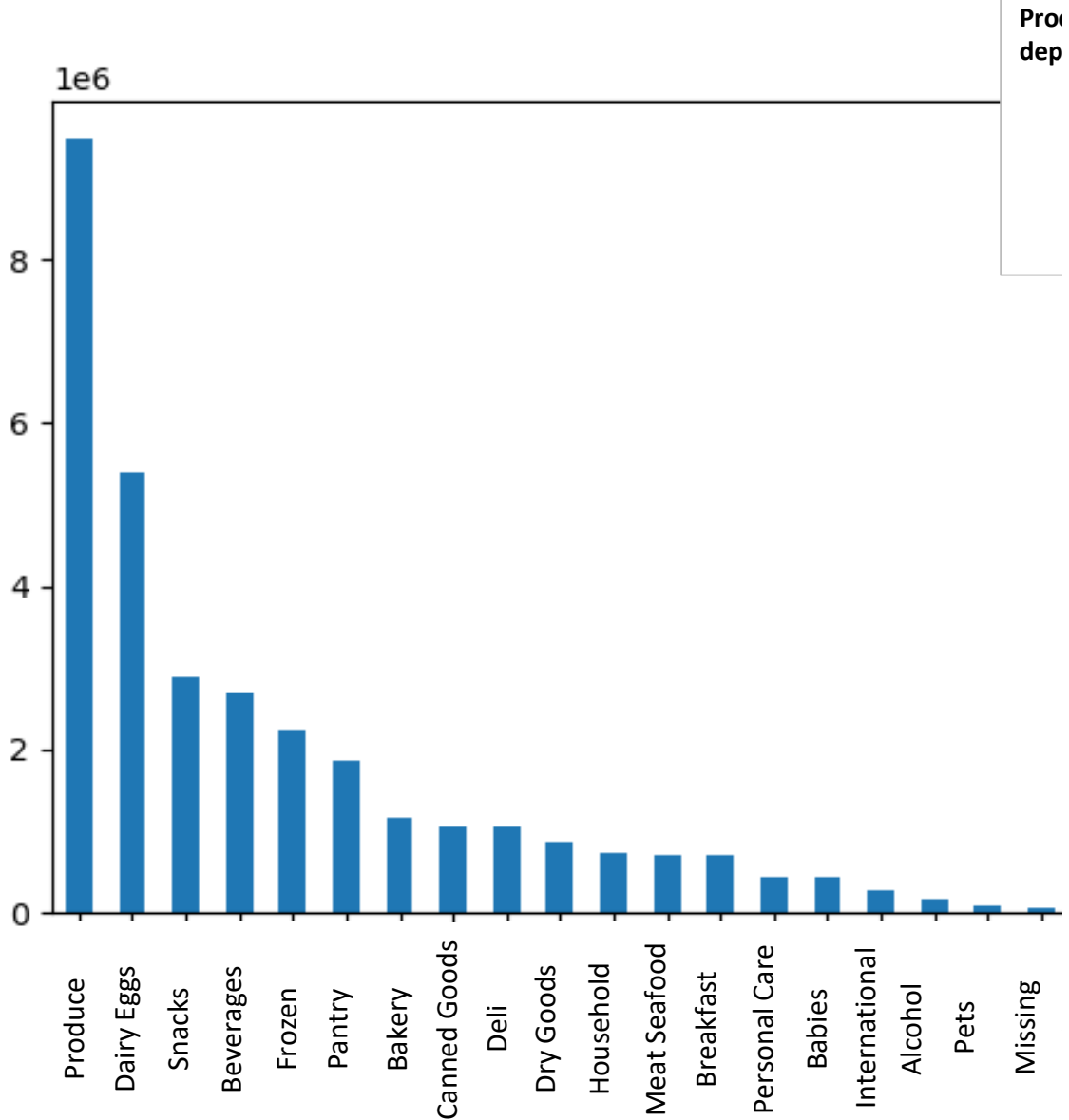


Instacart has a lot of products with different price tags. Marketing and sales w groupings to help direct their efforts

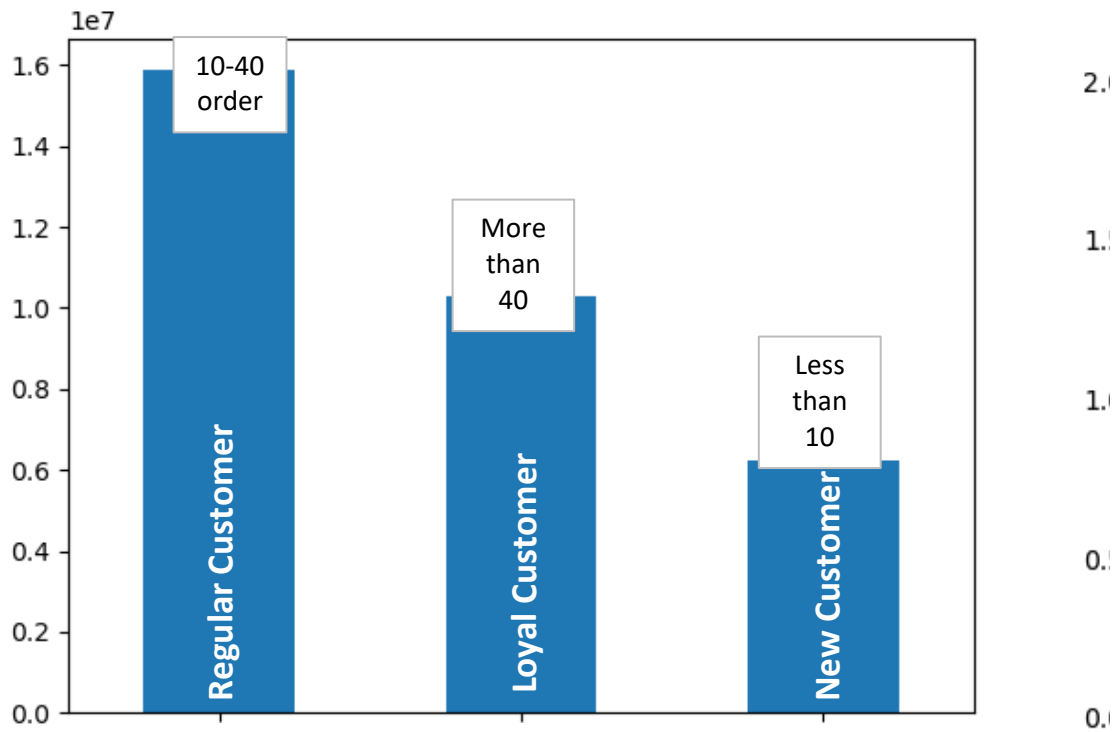




Are there certain types of products that are more popular than others? The more we know which departments have the highest frequency of product orders.



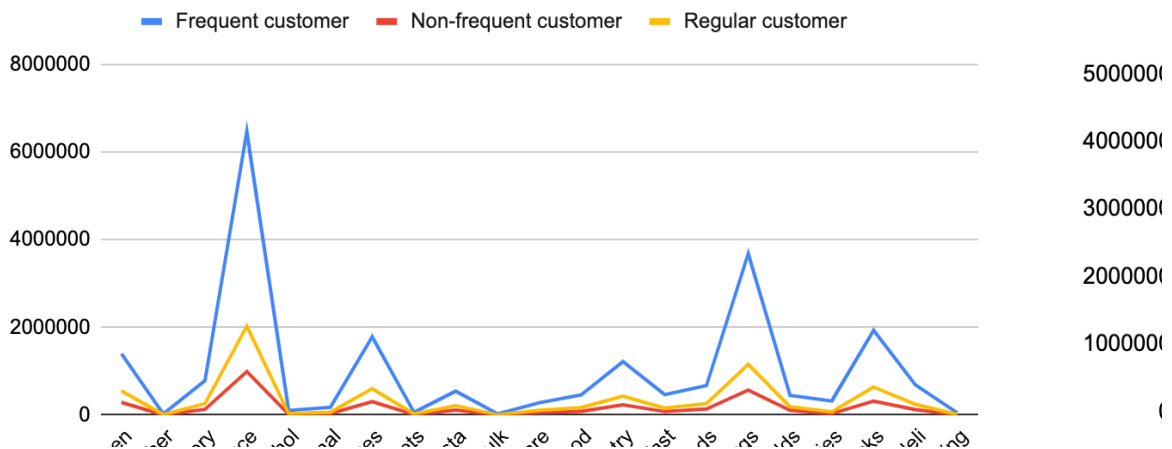
What's the distribution among users in regards to their brand loyalty (i.e., how many times they've purchased a product from a specific brand)?

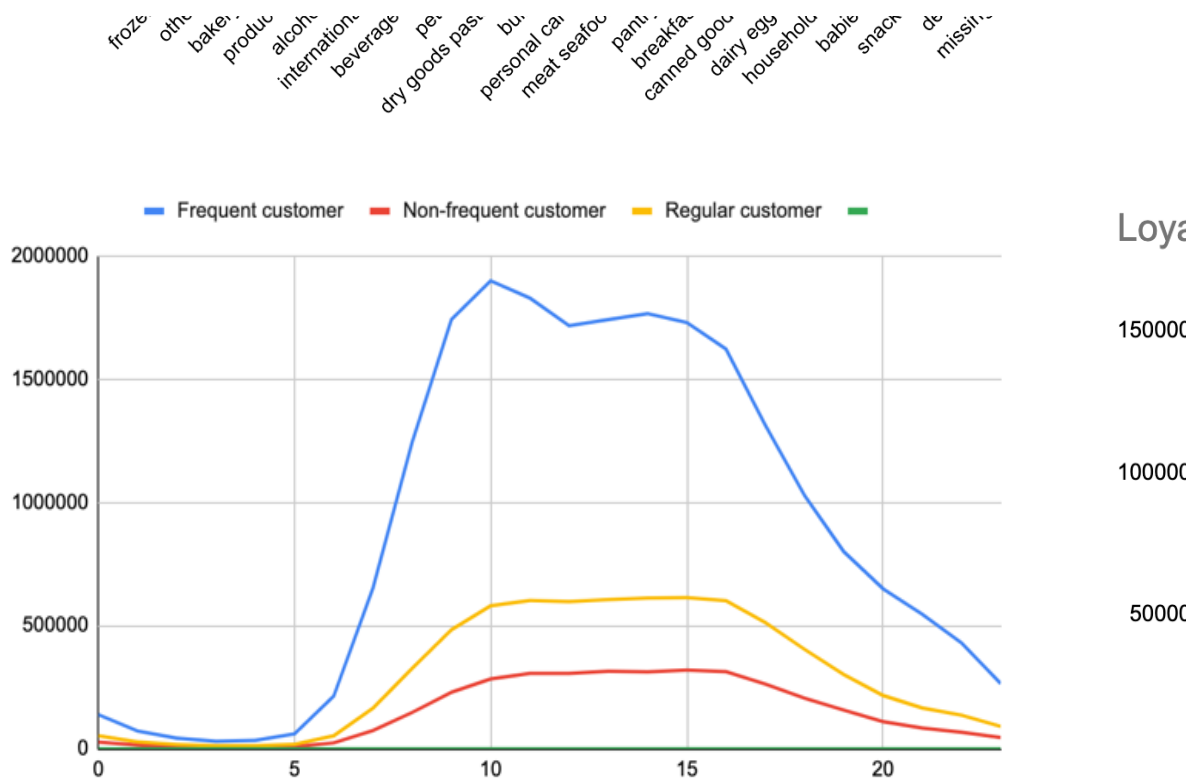


Are there differences in ordering habits based on customer's loyalty status?

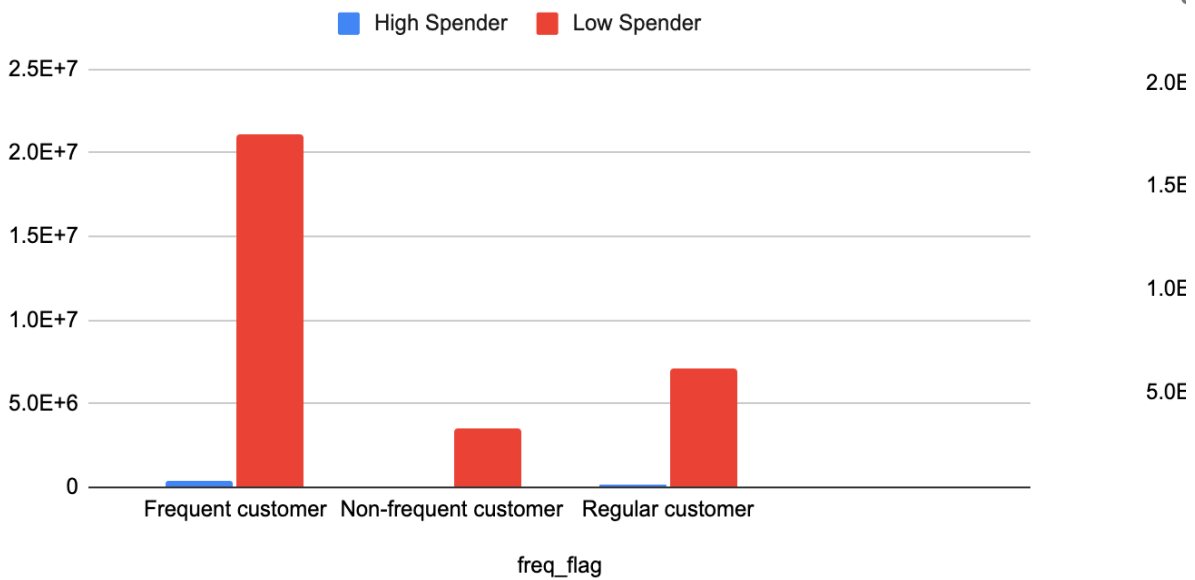
There is no outstanding difference in spending habits based on loyalty or frequency

Department Buying Habits Based on Frequency





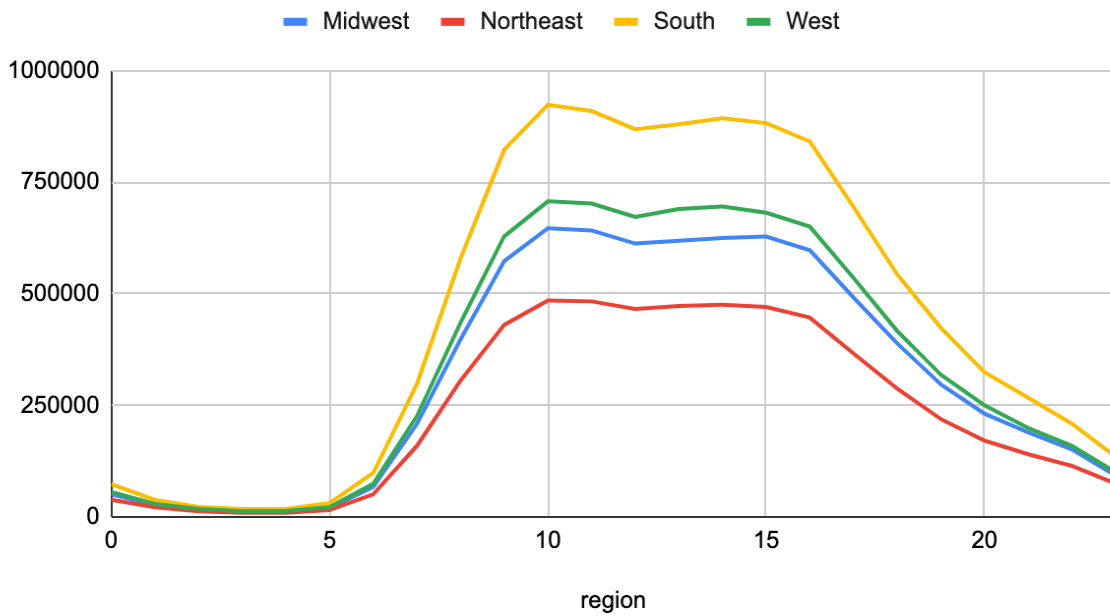
High Spender and Low Spender



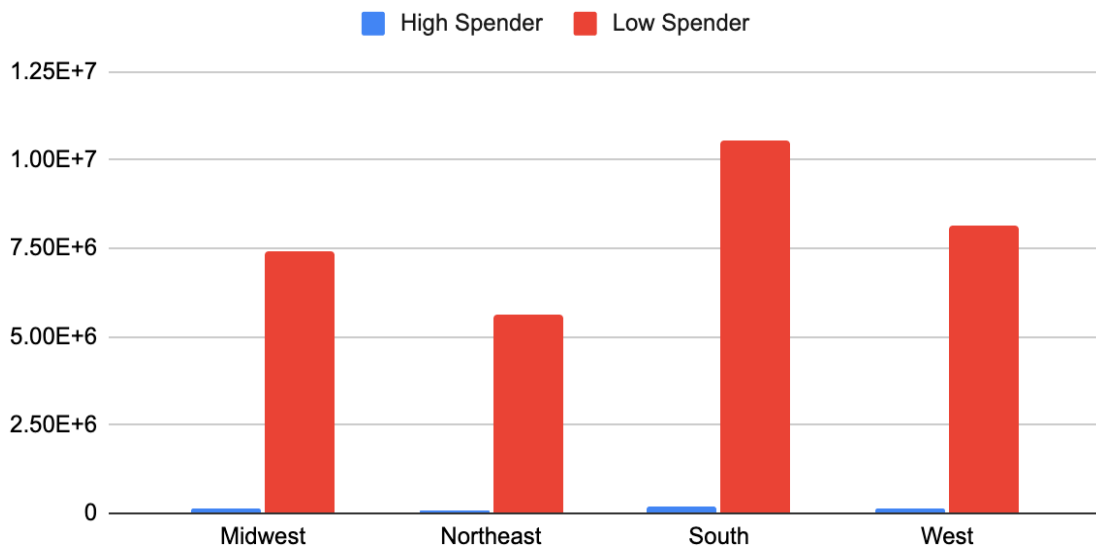
Are there differences in ordering habits based on a customer's region

There is no outstanding difference in spending habits based on region. They f

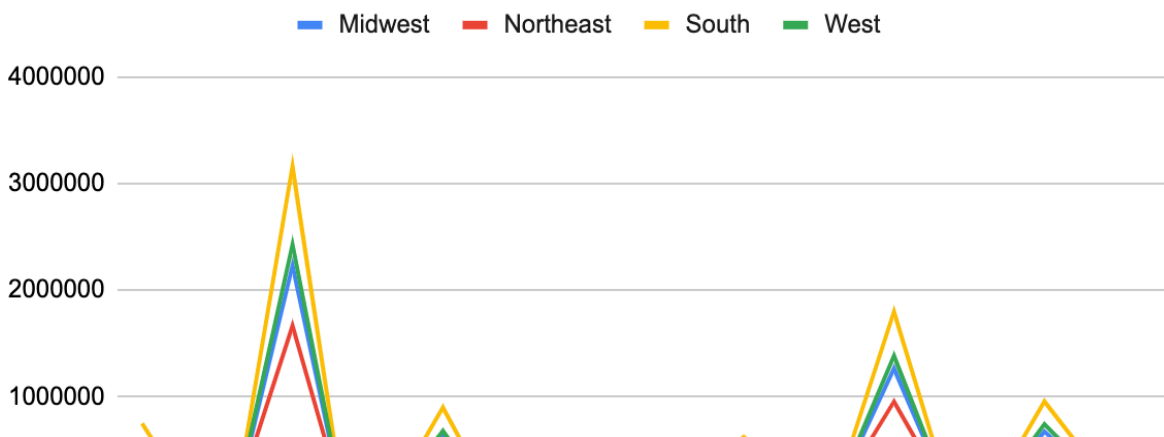
Midwest, Northeast, South and West

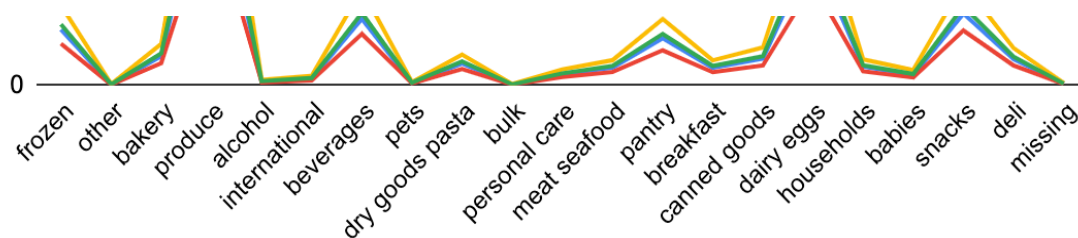


High Spender and Low Spender



Midwest, Northeast, South and West



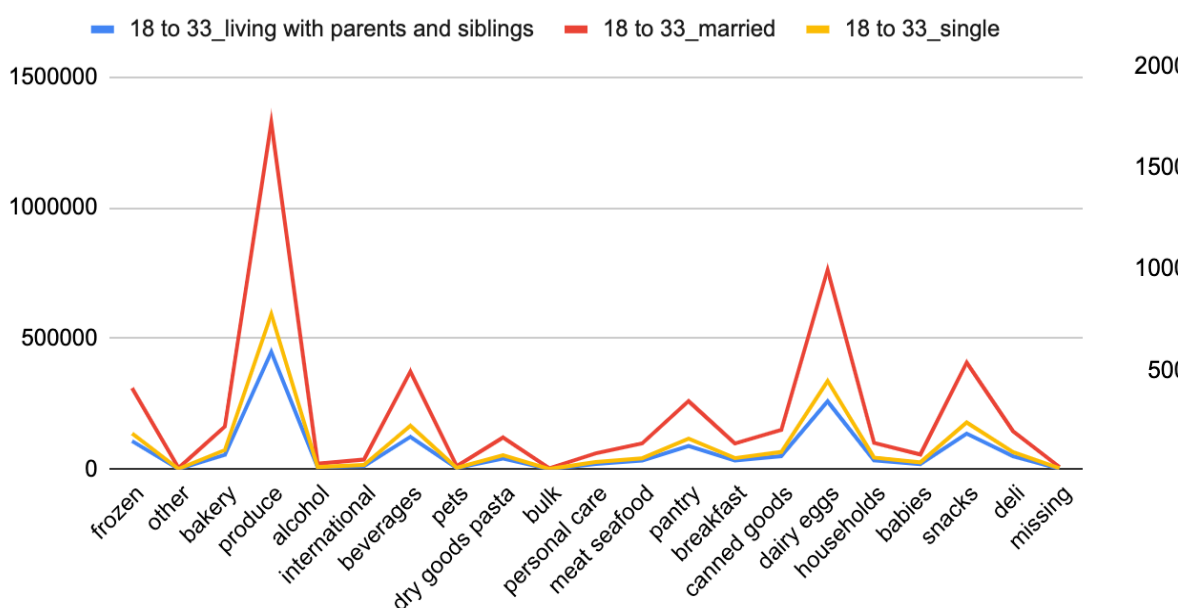


Is there a connection between age and family status in terms of ordering hab

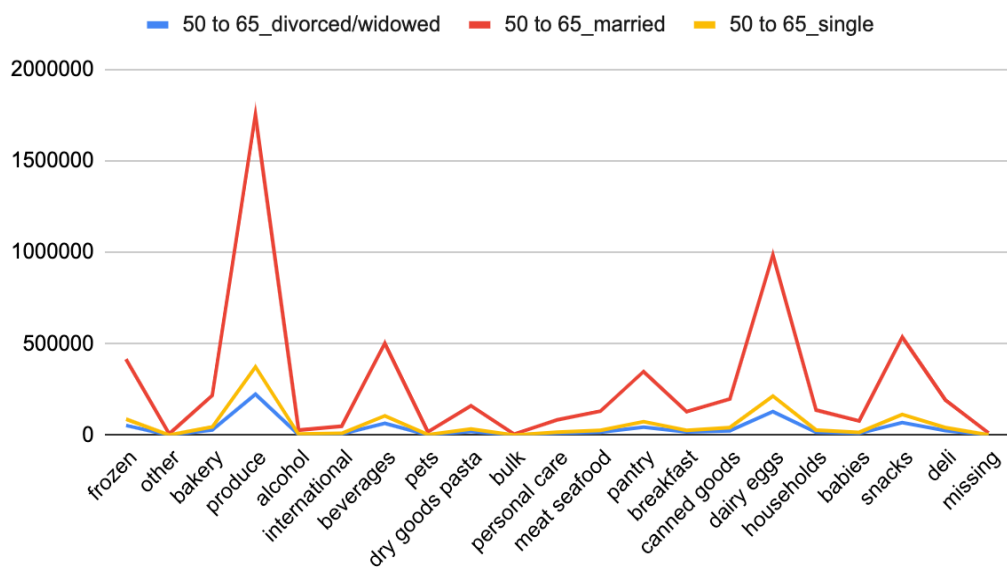
There is no outstanding difference in spending habits based on age and family statu

18 to 33 Family Status Department Buying Habits

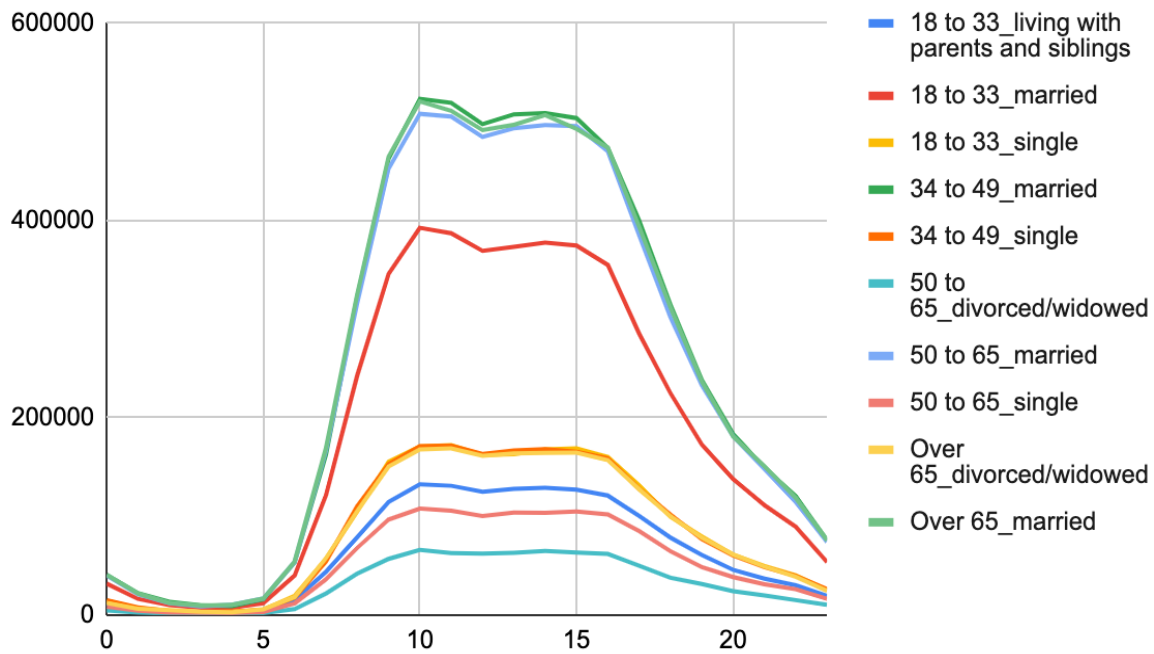
34



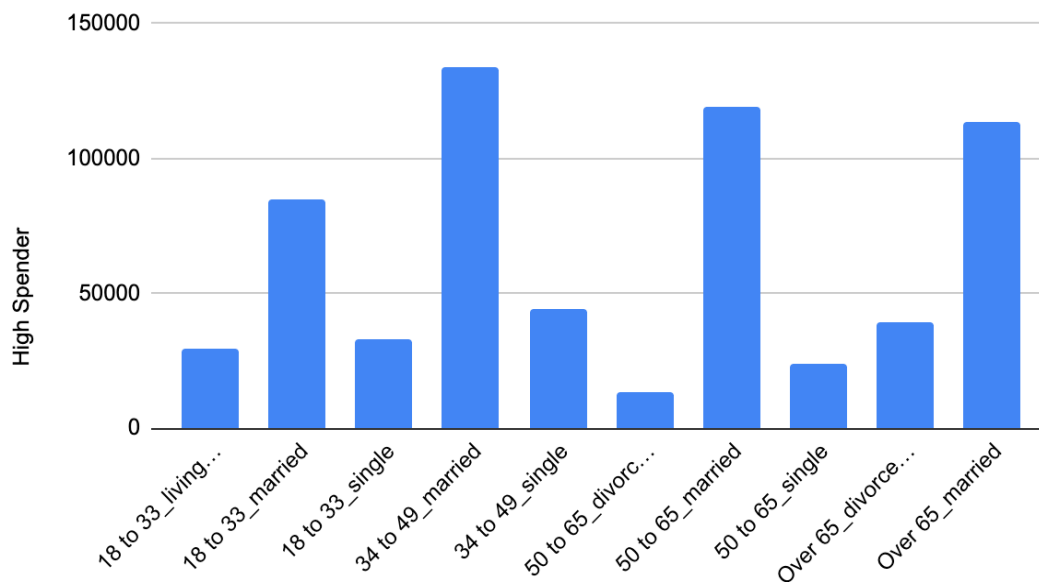
50 to 65 Family Status Department Buying Habits



Age, family status, by time



High Spender Age and Fam Status

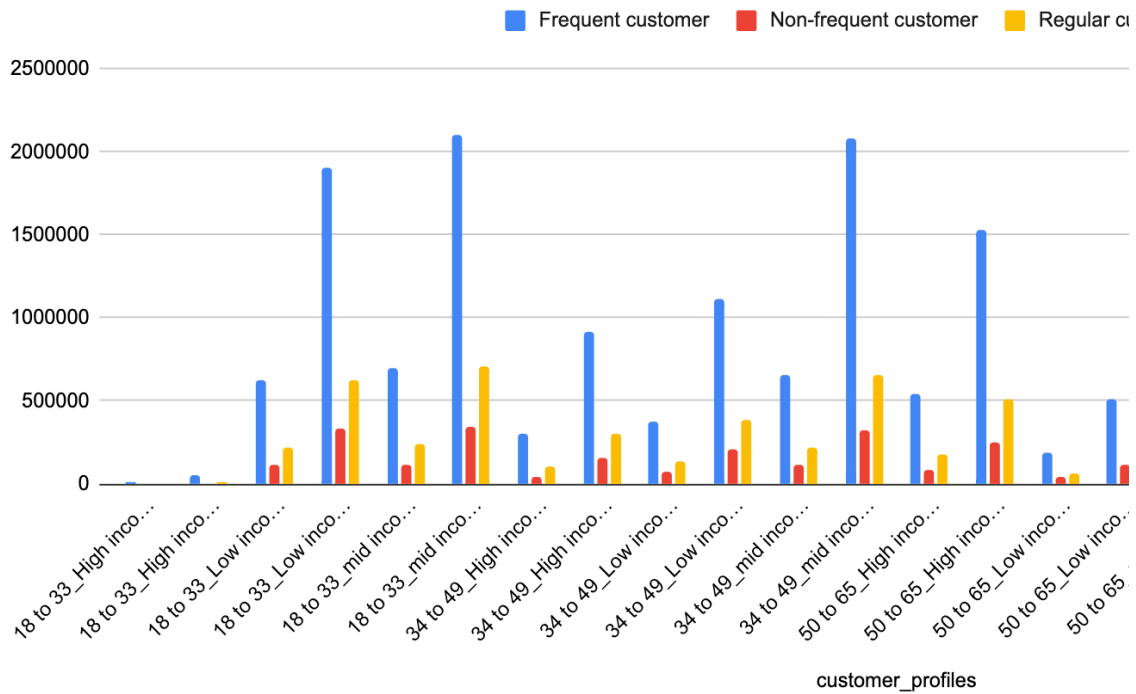


I think
the h
ma
tho

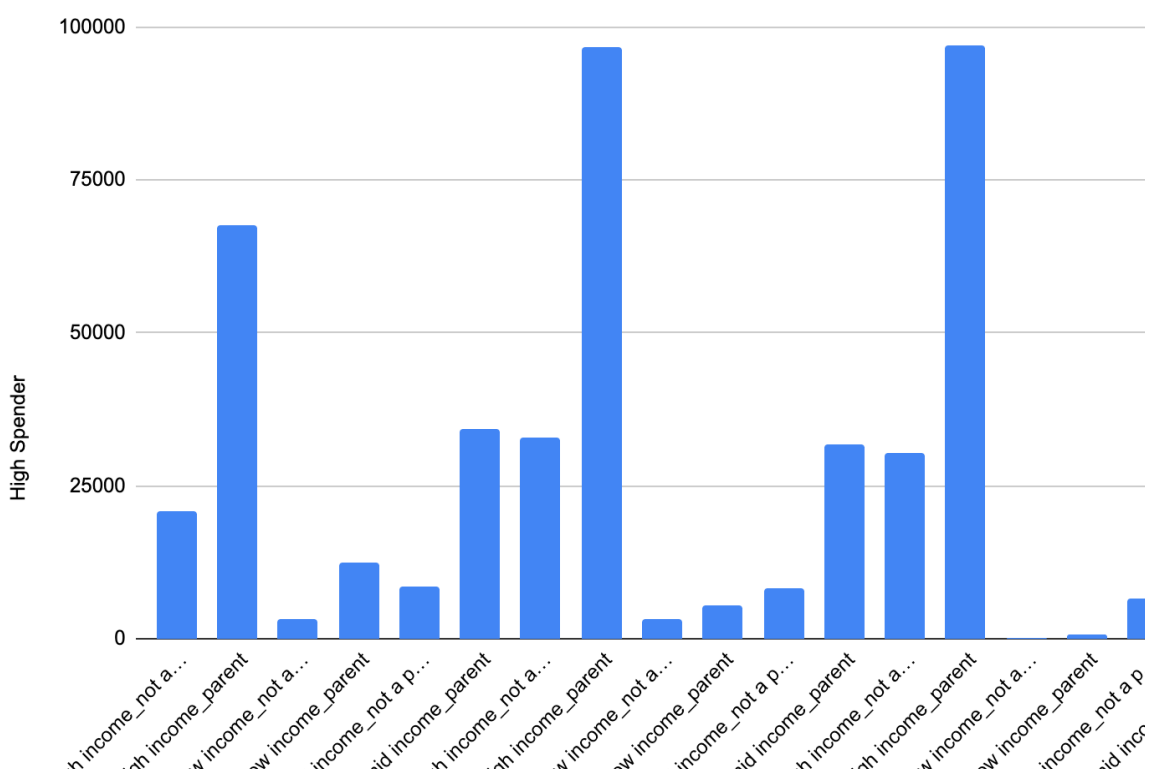
What different classifications does the demographic information suggest goods? Family status?

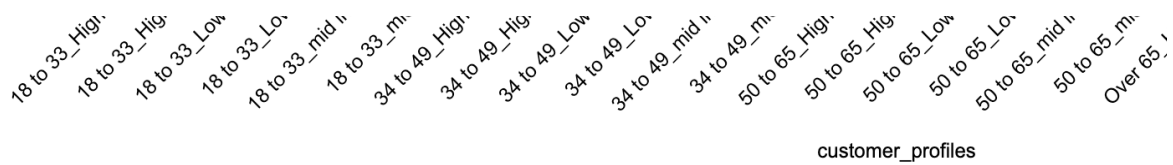
goods: Family status:

Frequent customer, Non-frequent customer and Regular customer

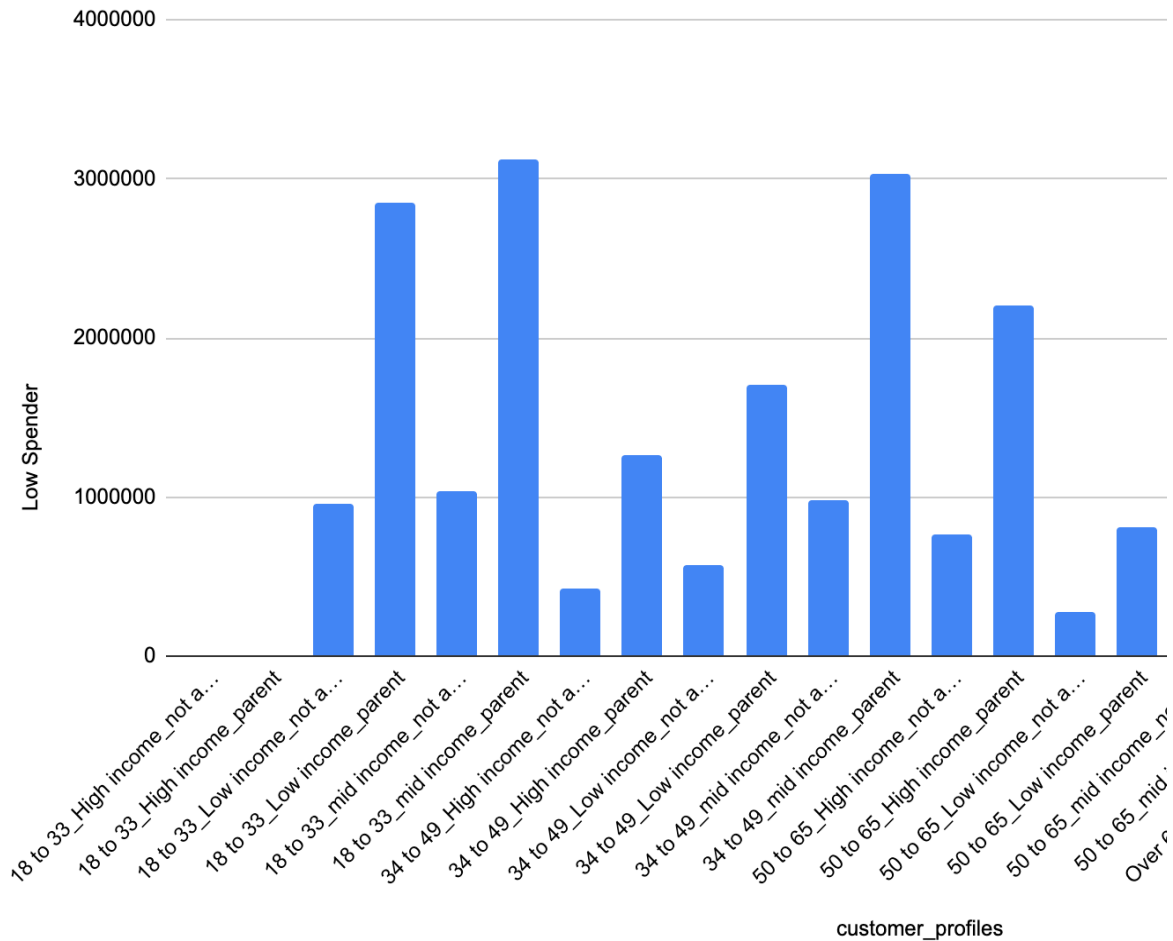


High Spender vs. customer_profiles





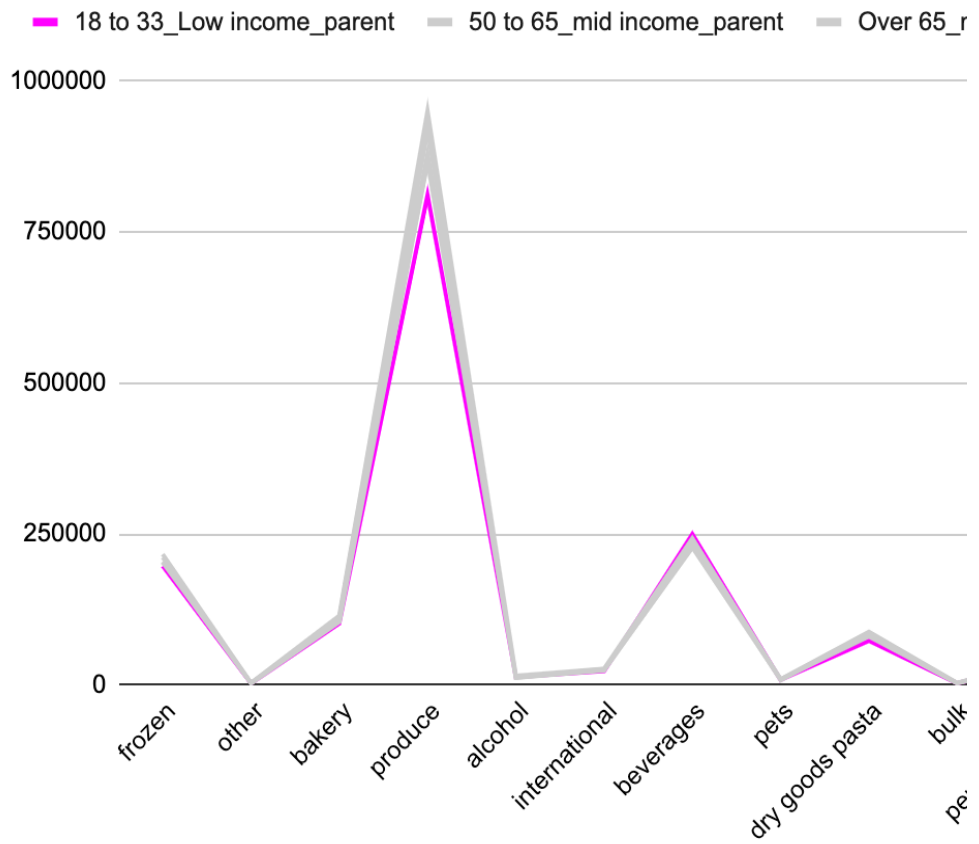
Low Spender vs. customer_profiles



What differences can you find in ordering habits of different customer profiles? Consider the frequency of orders, the products customers are ordering, and the amount spent.

One profile stands out is low income parents in the 18-33 age group. They tend to purchase more snacks, beverages, and breakfast items, over other customer profiles.

18 to 33 low income parent compared to other top cl



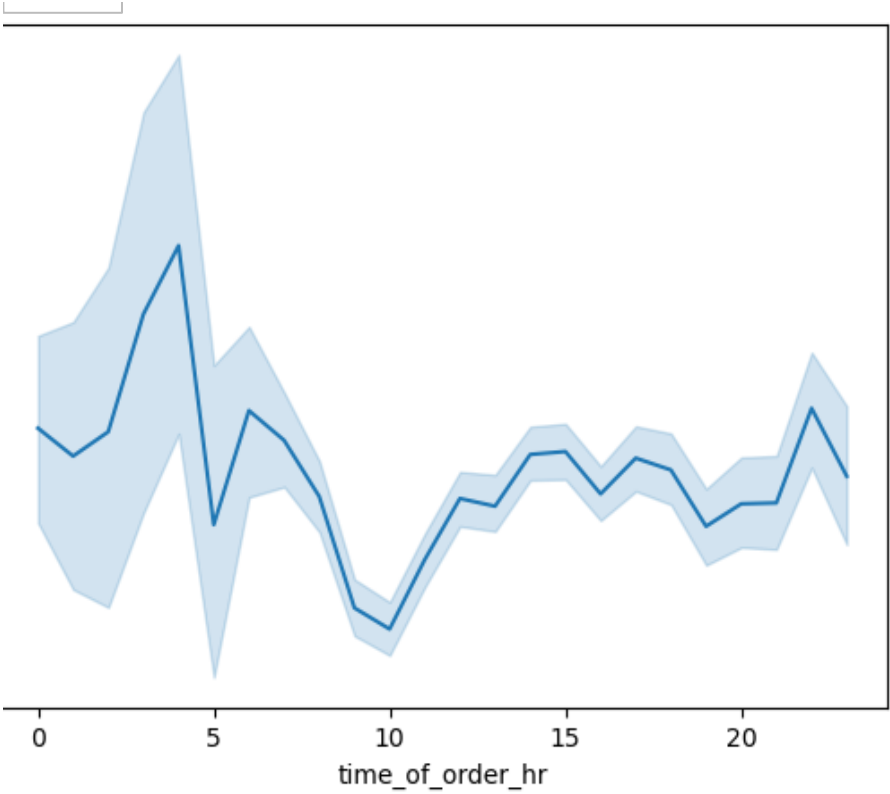
of the day are (i.e., the days and times with the most

s are
nd.
rday
from
,
the
(fig

people spend the most money, as this might

people are shopping between 9am and
people spend the most money between 2
and 4am (fig 2.2)

fig 2.2



want to use simpler price range

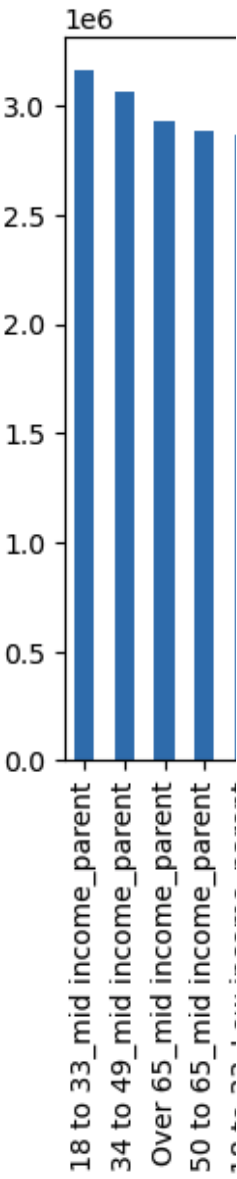
Price Range	Number of Products
mid-range (\$5-\$15)	22,086,764
low-range (less than \$5)	9,900,417
high-range (more than \$	714,678

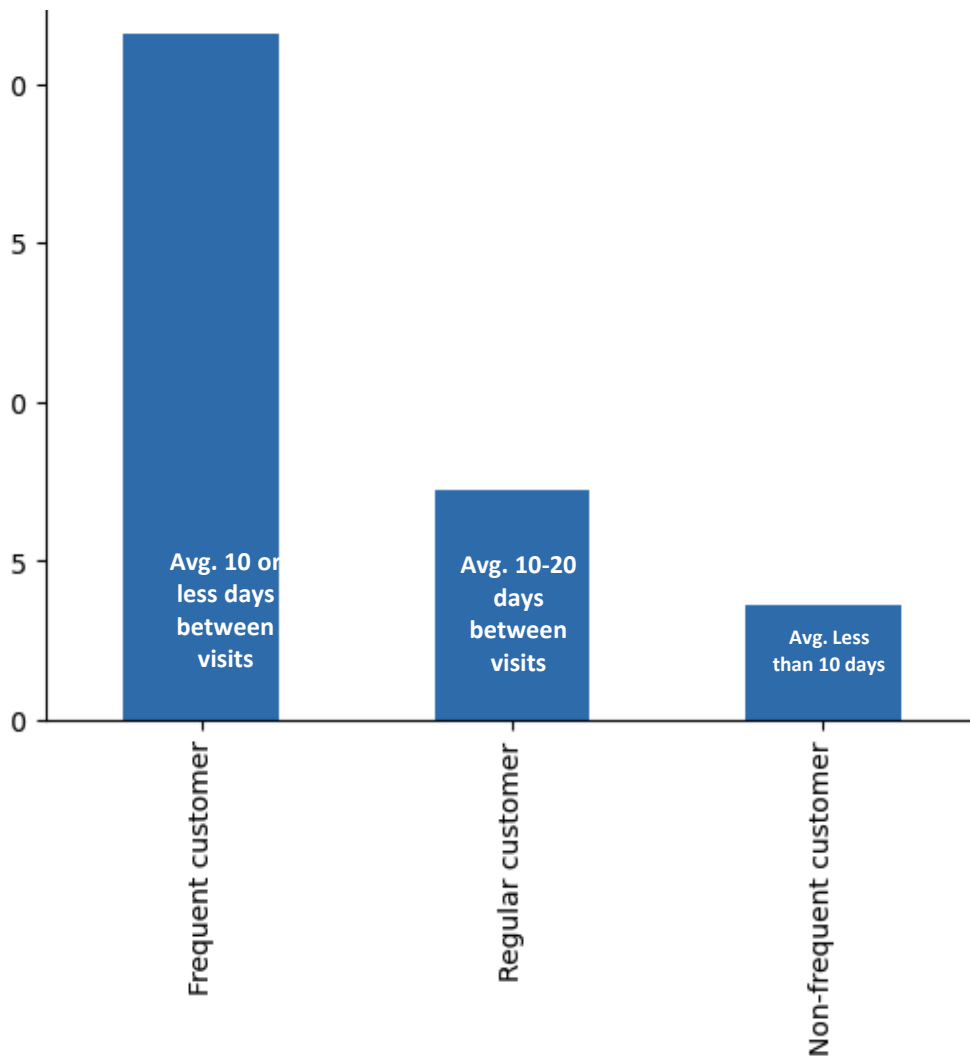
Marketing and sales teams want

Product is by far the most popular department, overall.



How often do they return to





rs?

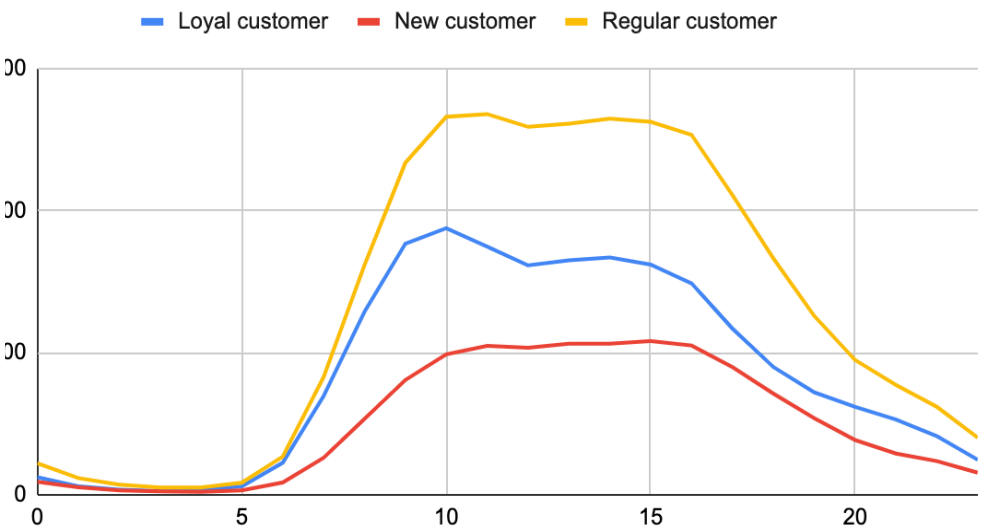
equency. They follow the same patterns as the data as a whole.

Ordering habits based on Loyalty Flag

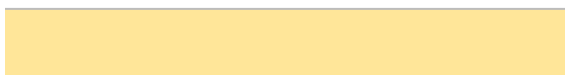
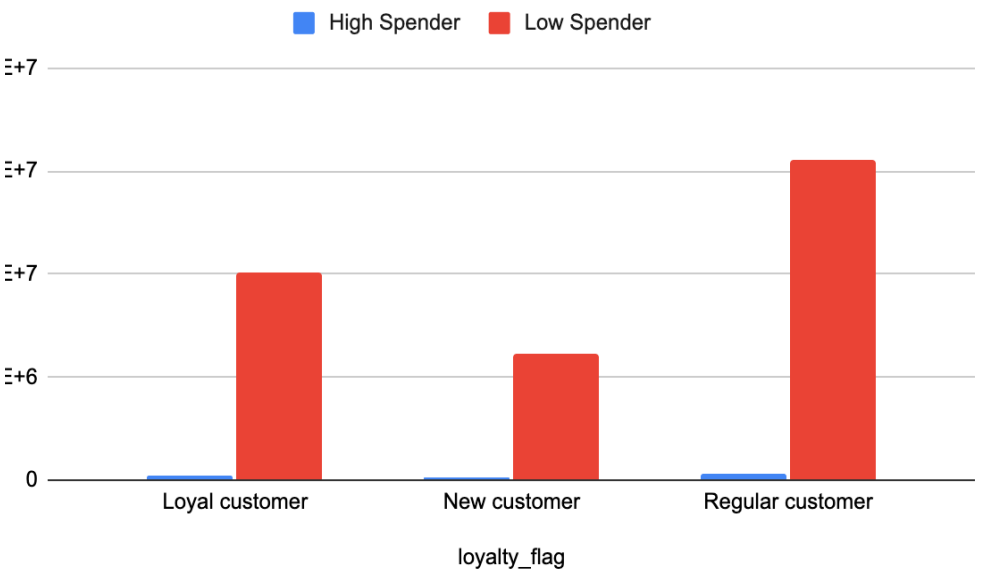




alty Flag Order Times



gh Spender and Low Spender

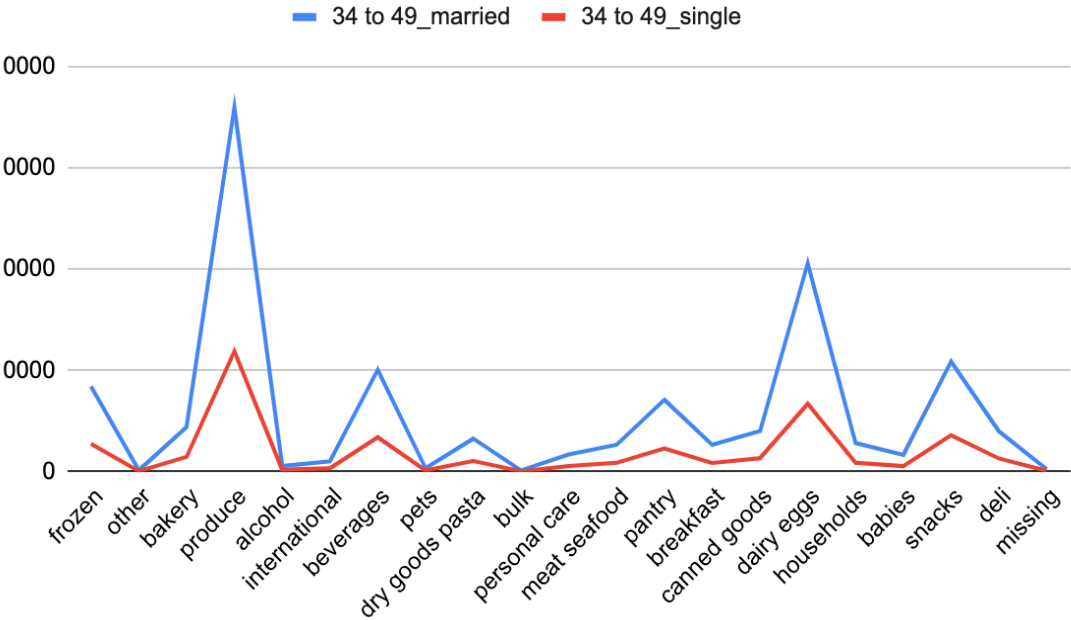


follow the same patterns as the data as a whole.

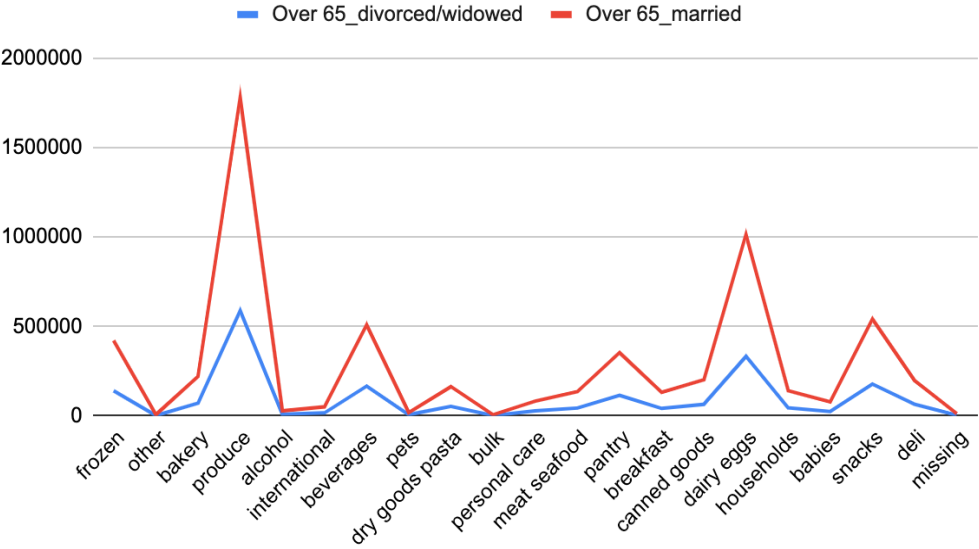
bits?

is. They follow the same patterns as the data as a whole.

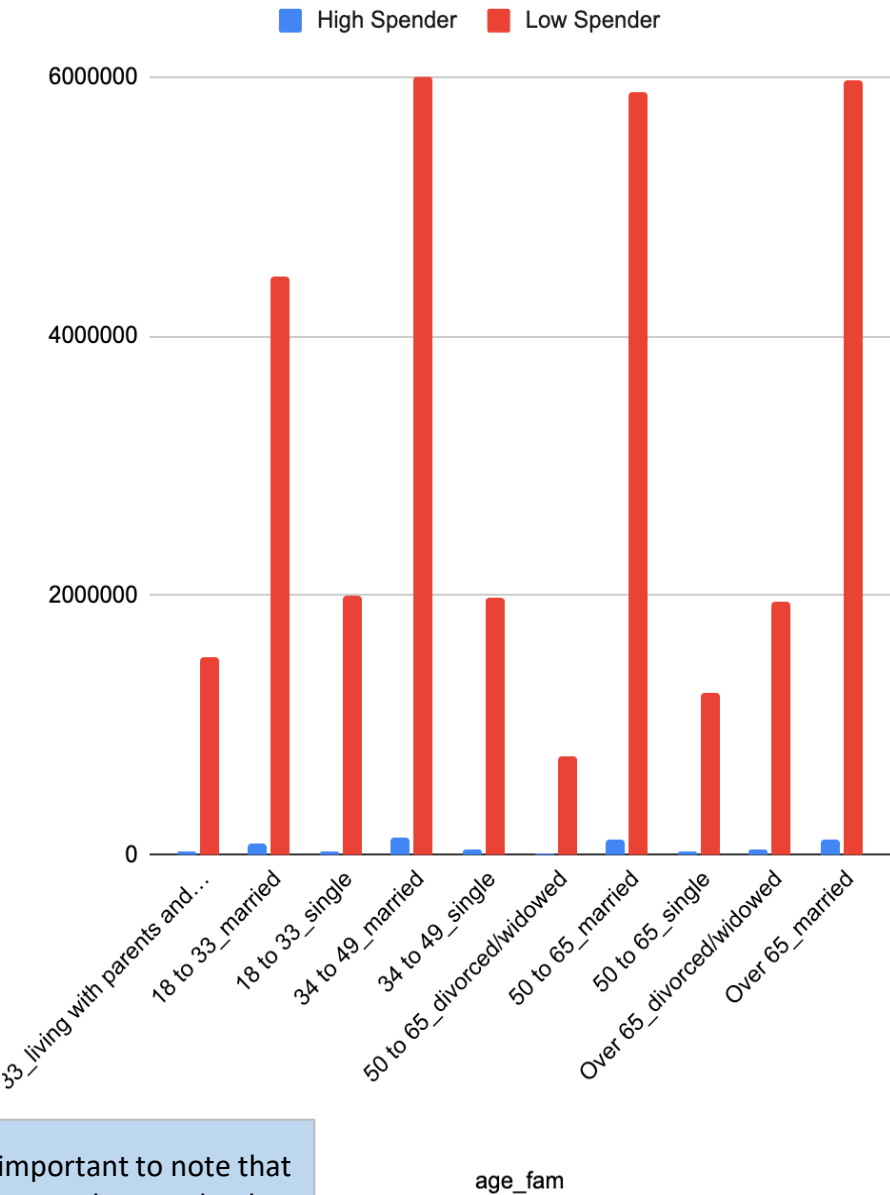
34 to 49 Family Status Department Buying Habits



Over 65_divorced/widowed and Over 65_married



High Spender and Low Spender

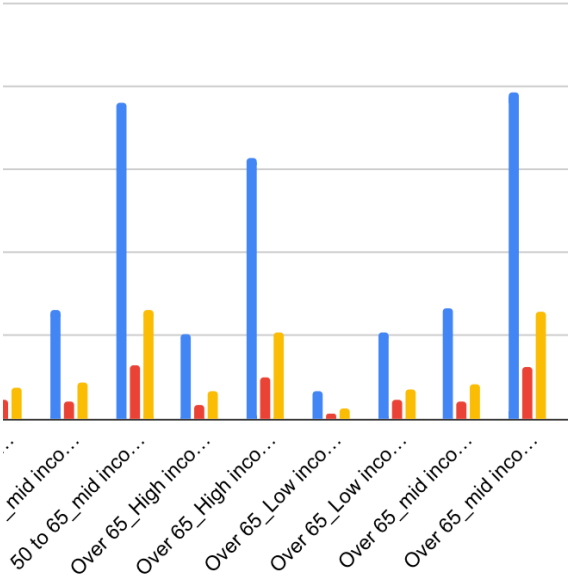


It is important to note that highest spenders tend to be married, with the highest of use being between 34 - 49 years of age.

? Age? Income? Certain types of

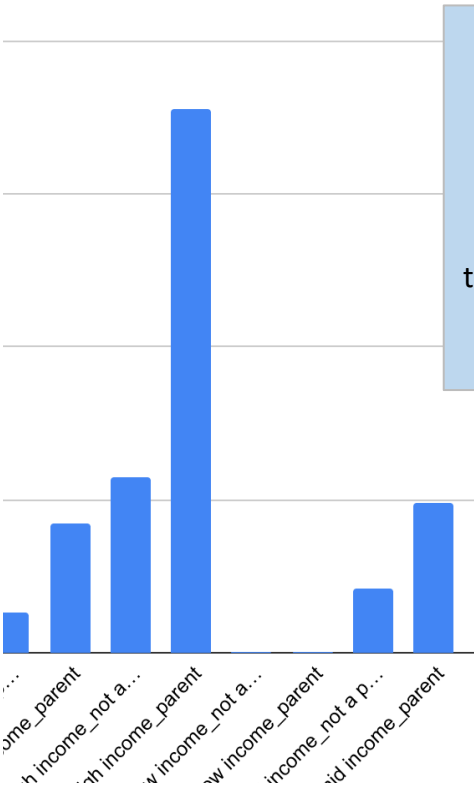


customer



The most frequent customers tend to be parents between 18 and 65 in the mid income range.

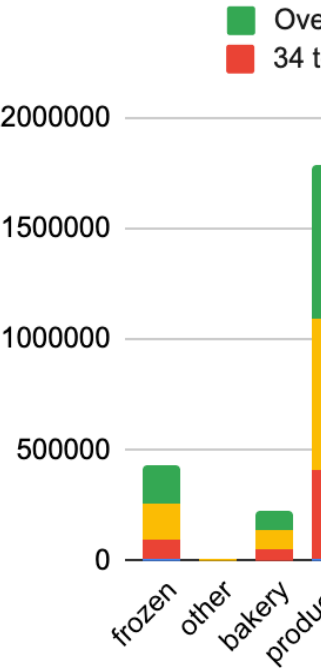
An important note is the higher the income the less frequent the customer, and 18-33 high income have some of the lowest frequencies, while thier mid to low income profiles have some of the



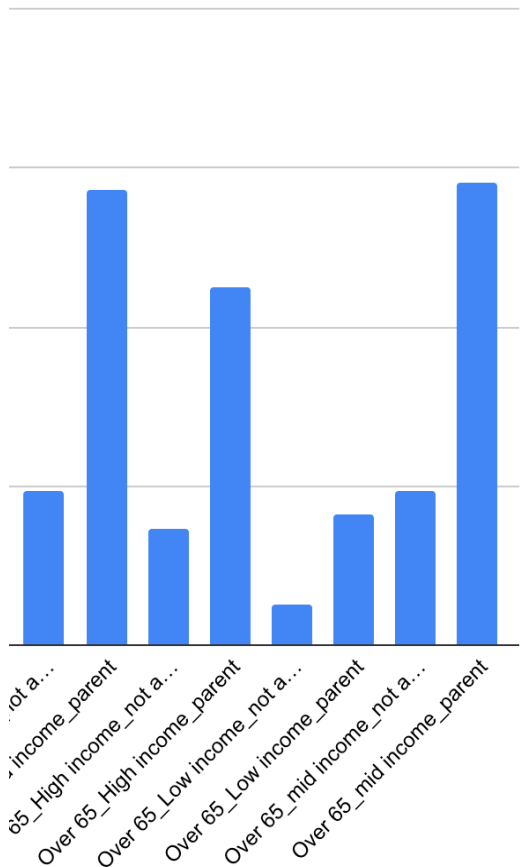
The highest spenders are parents in the high income range.

There are no departments that deviate from the general buying patterns of the data

Department Bu

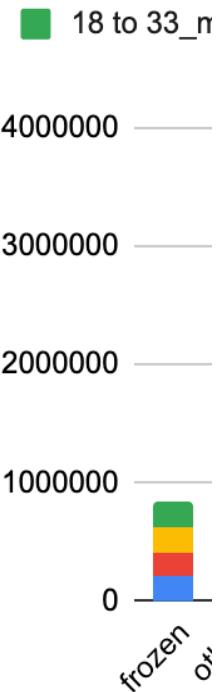


High
Over 65_High
Over 65_Low
Over 65_Low
Over 65_mid
Over 65_mid



The lowest spenders are parents in the middle income range, but they buy the most amount of items.

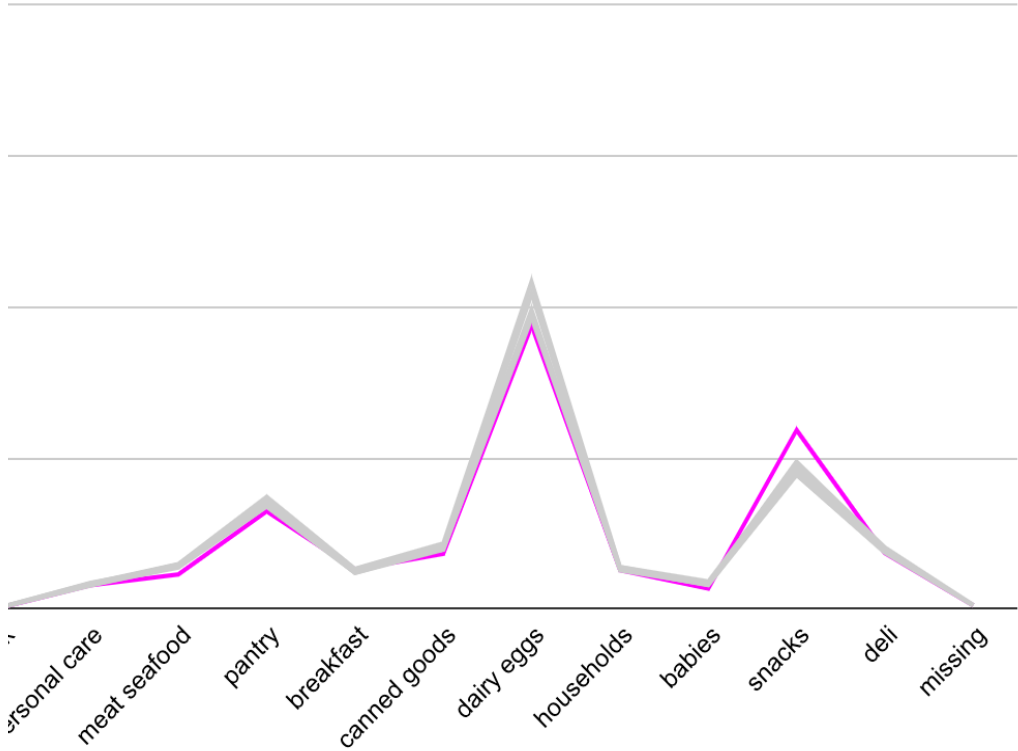
Buying Habits



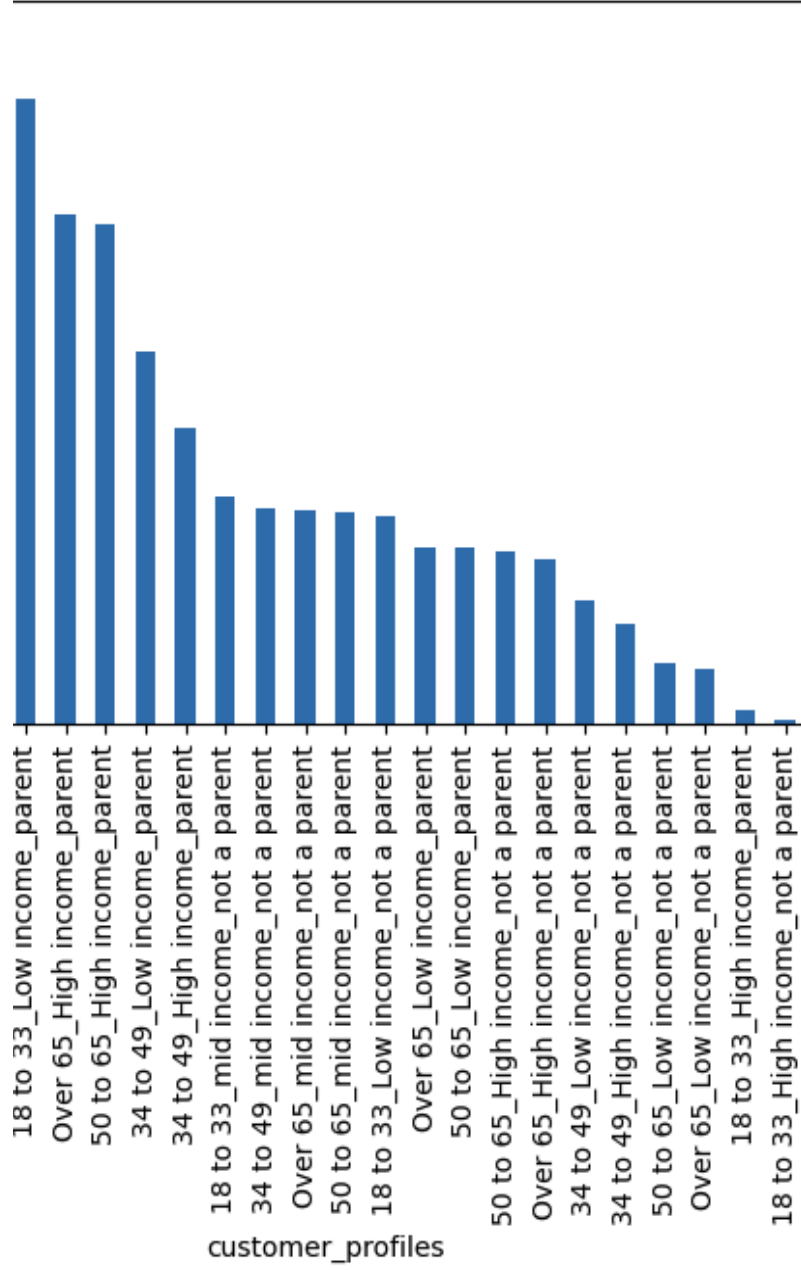
other profiles? Consider the price of orders, anything else you can think of.

Customer profiles

mid income_parent 34 to 49_mid income_parent 18 to 33_mid income_parent

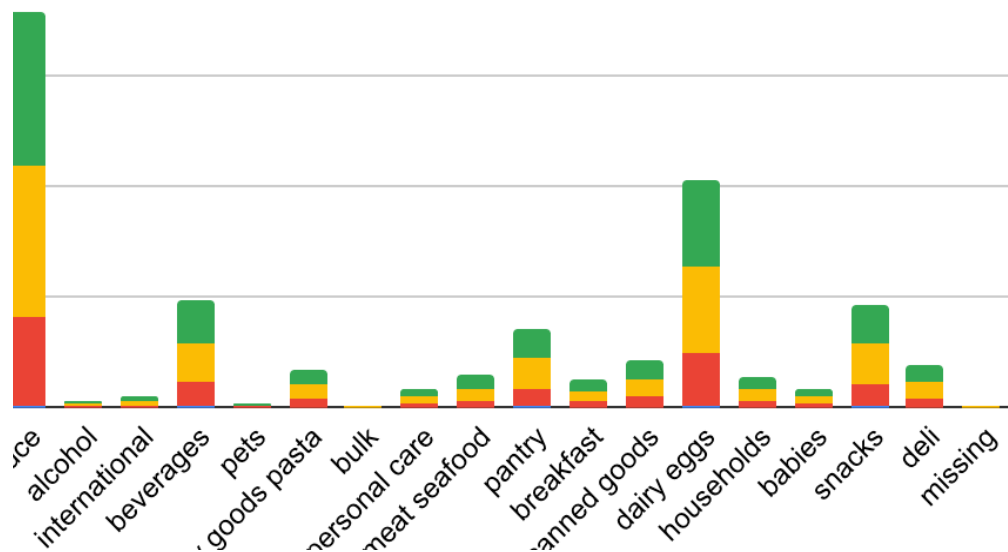


Title page



Spending Habits of Highest Spenders

65+ 50 to 65 18 to 33



dry - p n -

abits of Middle Income Range Parents

mid income_parent 34 to 49_mid income_parent Over 65_mid income_parent
50 to 65_mid income_parent

