

iFood CRM Case Study

Presenters:

JANN Ang

MAGGIE Ding

AISHWARYA Panse

QUEENIE Chao

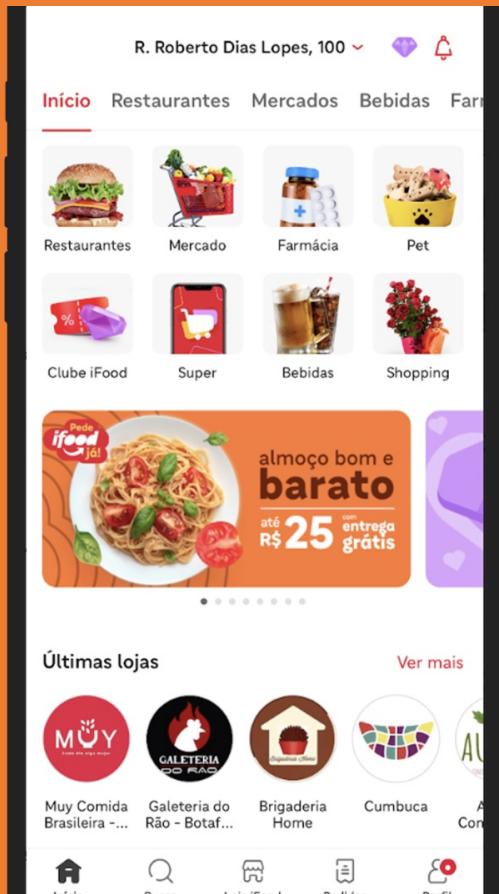
FARIN Fukunaga



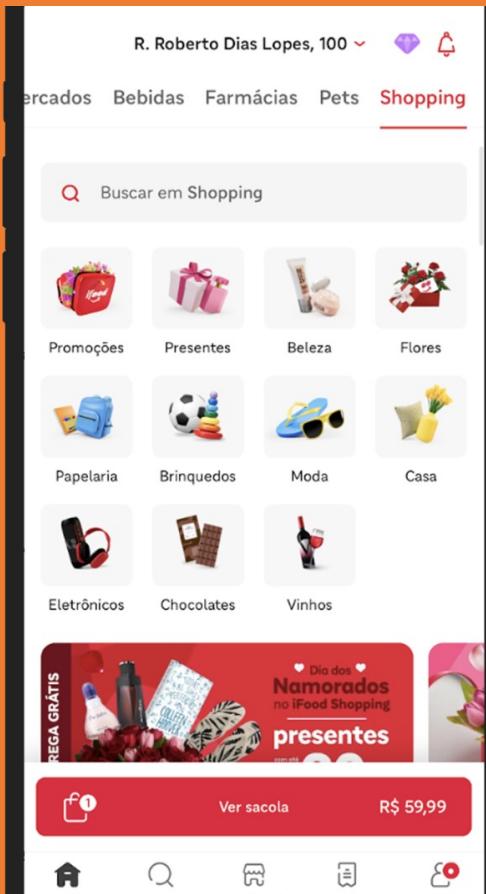
•
•
•
Pedi
chegou.

ifood

iFood's home screen



iFood has a variety of categories for different products



"Who wants a coupon?" is the title of the app's promotions page



Description

Data's origin

Case study competition “iFood Brain team data challenge” provided by the iFood

With this sample dataset, we will analyze meta information on the customer and on iFood campaign interactions with their customers

Our challenge is understanding the data, finding business opportunities & insights, and analyze different segment of customer's spending



Key Objective

Explore the data – Provide insights, and define cause and effect. Provide a better understanding of the characteristic features of respondents

Describe customer segmentation based on customers' behaviors

Create 2 predictive models which allow the company to maximize the profit of their next marketing campaign



Data Dictionary

Products from 5 major categories



Further divided into:



3 sales channels:



Feature	Description
ID	Unique identifier for each customer
Year_Birth	Customer's year of birth
DtCustomer	Date of customer's enrollment with the company
Education	Customer's level of education
Marital	Customer's marital status
Kidhome	Number of small children in customer's household
Teenhome	Number of teenagers in customer's household
Income	Customer's yearly household income
MntFishProducts	Amount spent on fish products in the last 2 years
MntMeatProducts	Amount spent on meat products in the last 2 years
MntFruits	Amount spent on fruit products in the last 2 years
MntSweetProducts	Amount spent on sweet products in the last 2 years
MntWines	Amount spent on wine products in the last 2 years
MntGoldProds	Amount spent on gold products in the last 2 year
NumDealsPurchases	Number of purchases made with discount
NumCatalogPurchases	Number of purchases made using catalogue
NumStorePurchases	Number of purchases made directly in stores
NumWebPurchases	Number of purchases made through company's website
NumWebVisitsMonth	Number of visits to company's website in the last month
Recency	Number of days since the last purchase

Data Cleaning

```
(select *,  
case  
when Income = '' then null  
else Income  
end as Income_clean  
from ifood_data) as sub1)
```

1. Missing values: Some columns may have missing values represented as NULL using CASE statement
2. Dummy variables on the amount spent (wine, meat, fruit, etc): <= \$10 : 0 and >\$10 : 1

Well-documented 5

Well-maintained 5

Clean data 5

Challenges

- We can only select a few datasets that are significant/interesting enough to analyze it

Useful Data to have...

The region of Brazil these people came from

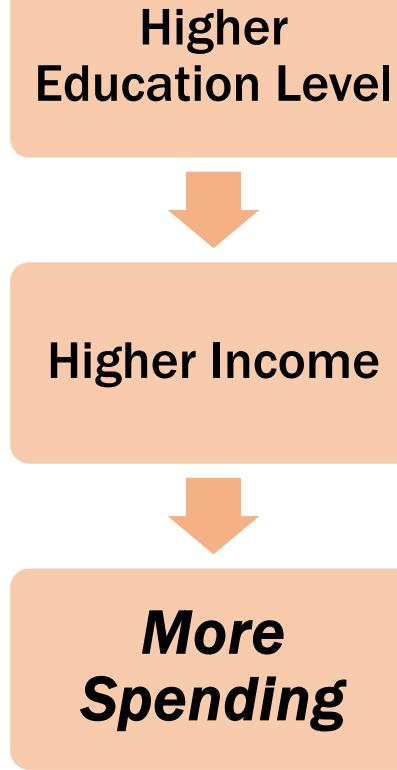
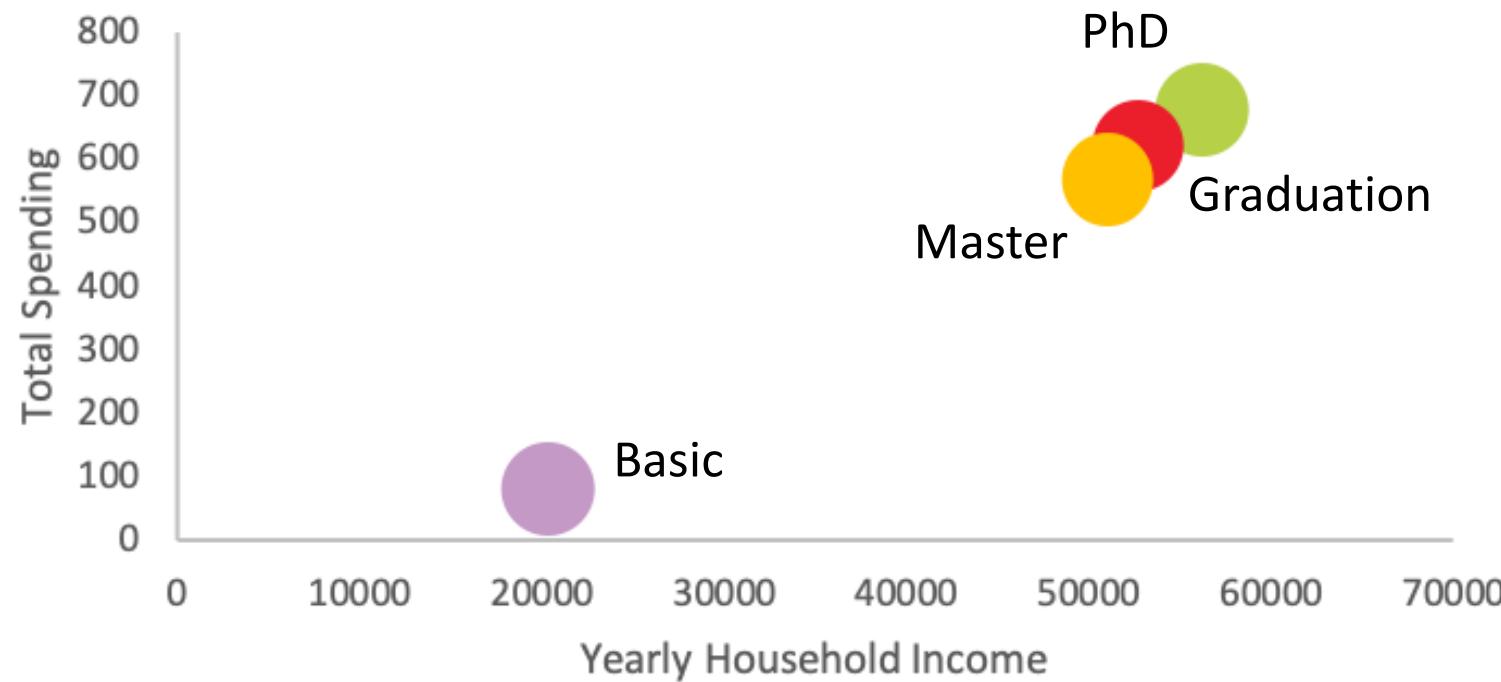
What exact products customer brought

What promotions customers responded to



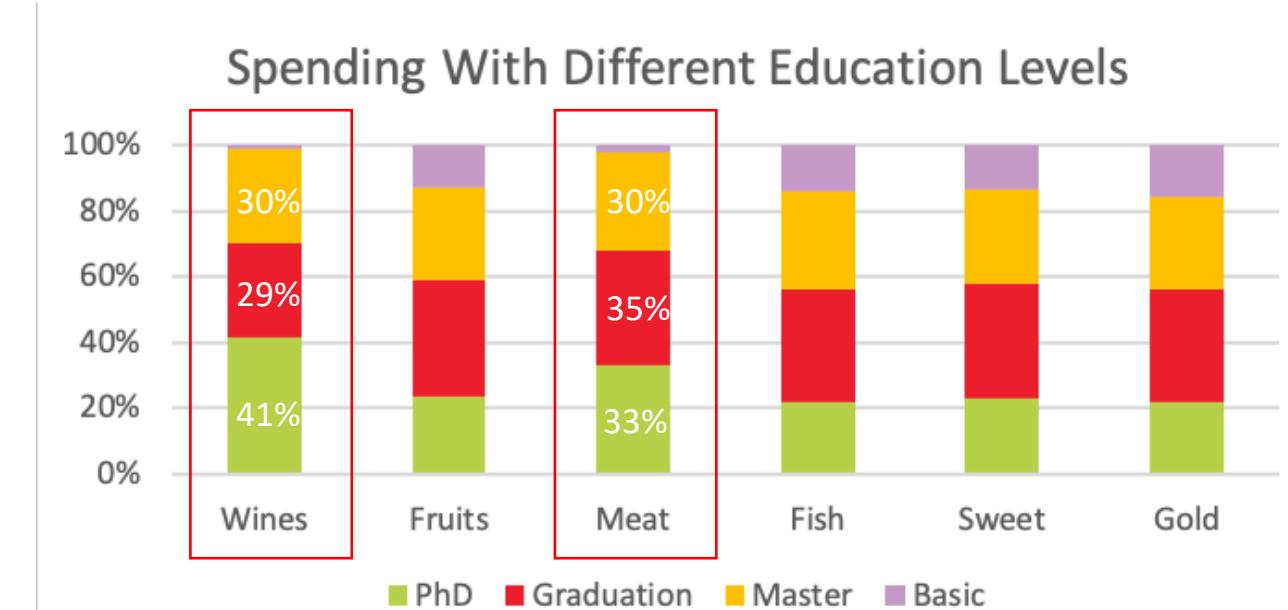
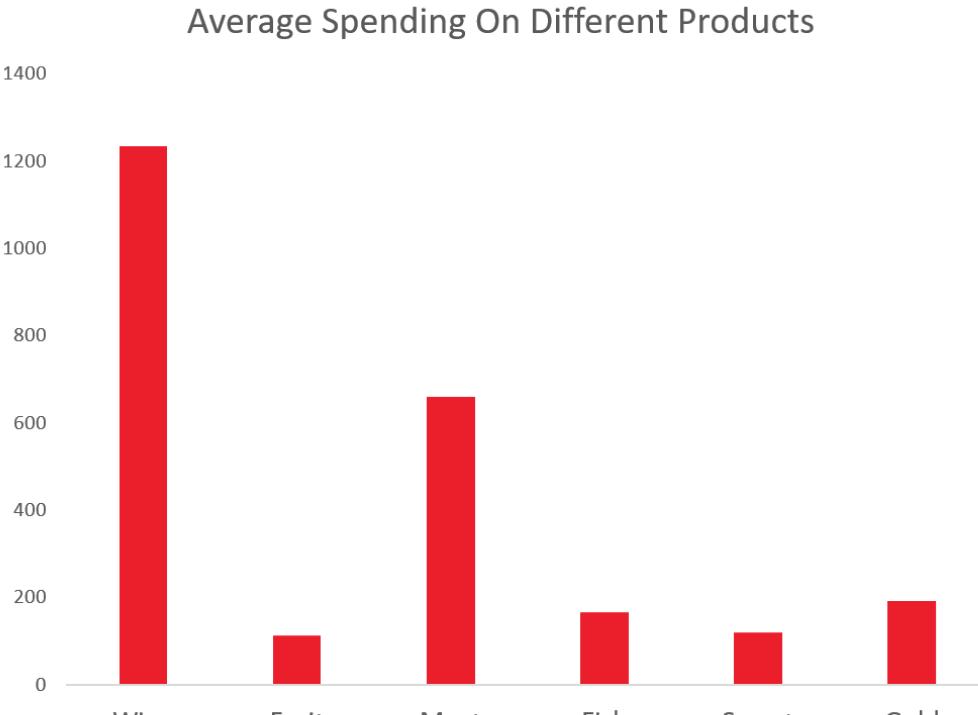
**Does a person's
education level affect
buying decisions?**

How Education Level Can Influence Grocery Spending?



How Education Level Can Influence Grocery Spending?

Customers with higher education levels are more willing to spend more on higher-priced groceries.



Master, Single
PhD, Couple
Without Kids/
Teens ...



Assumption:
Customers with higher education levels and households without kids or teens are more willing to spend more on higher-priced groceries.

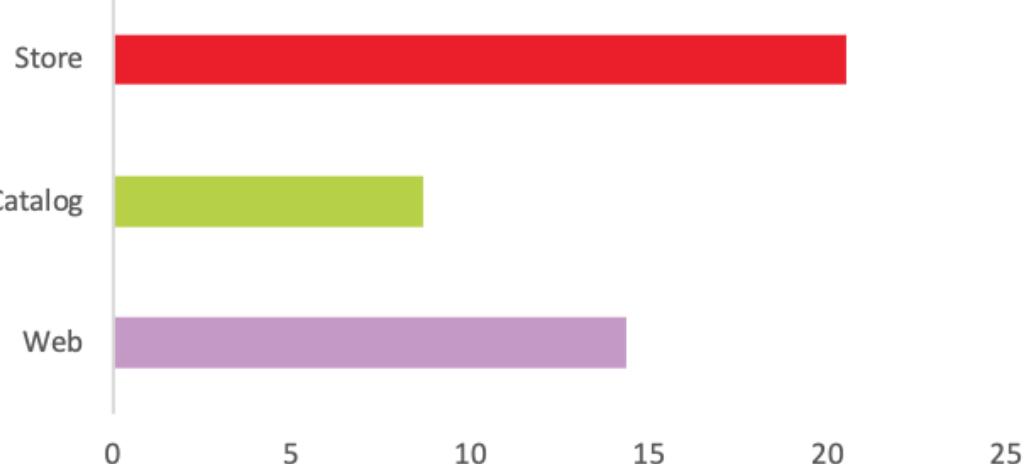
High Income, Consumer Hedonism...

How Education Level Can Influence Channel Choice?

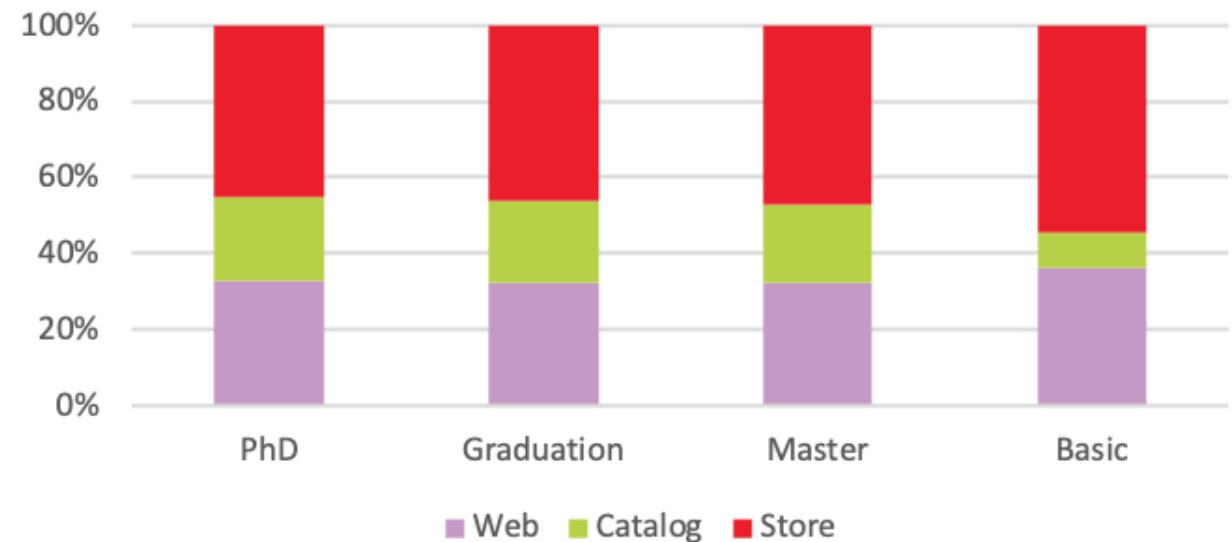
Customers make more purchases directly in stores.

There's no significant correlation between
Education Level and Channel Choice.

Average Purchase On Different Channels



Purchases With Different Education Levels



Assumption:
Households with kids or teens are
more willing to make purchases online.



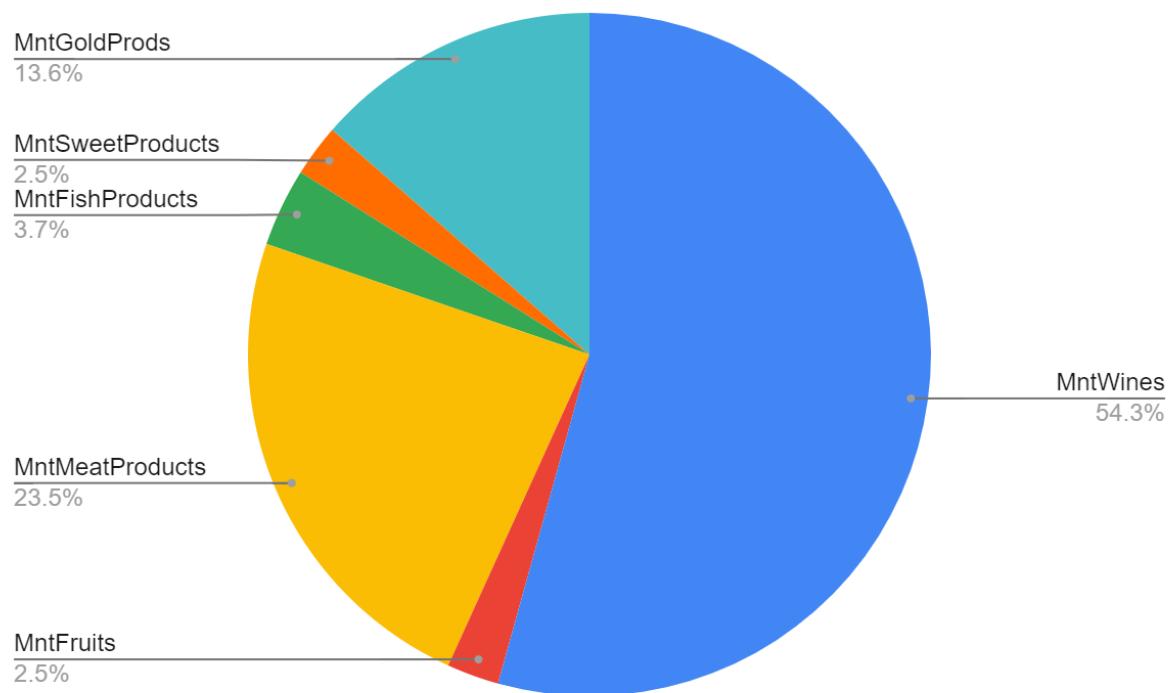


**Does having a family
affect buying decisions?**

Amount Spent On Various Categories By Non-Parents vs Parents

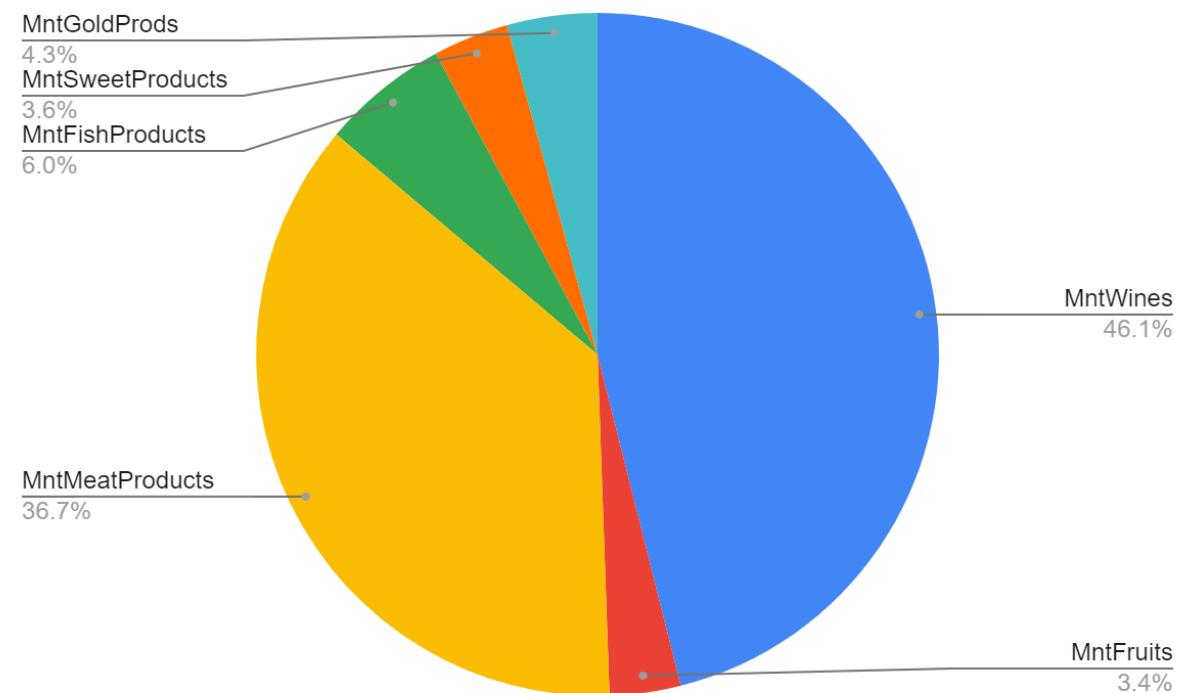
MntWines	MntFruits	MntMeat Products	MntFish Products	MntSweet Products	MntGold Prods
44	2	19	3	2	11

Parents
Count: 419



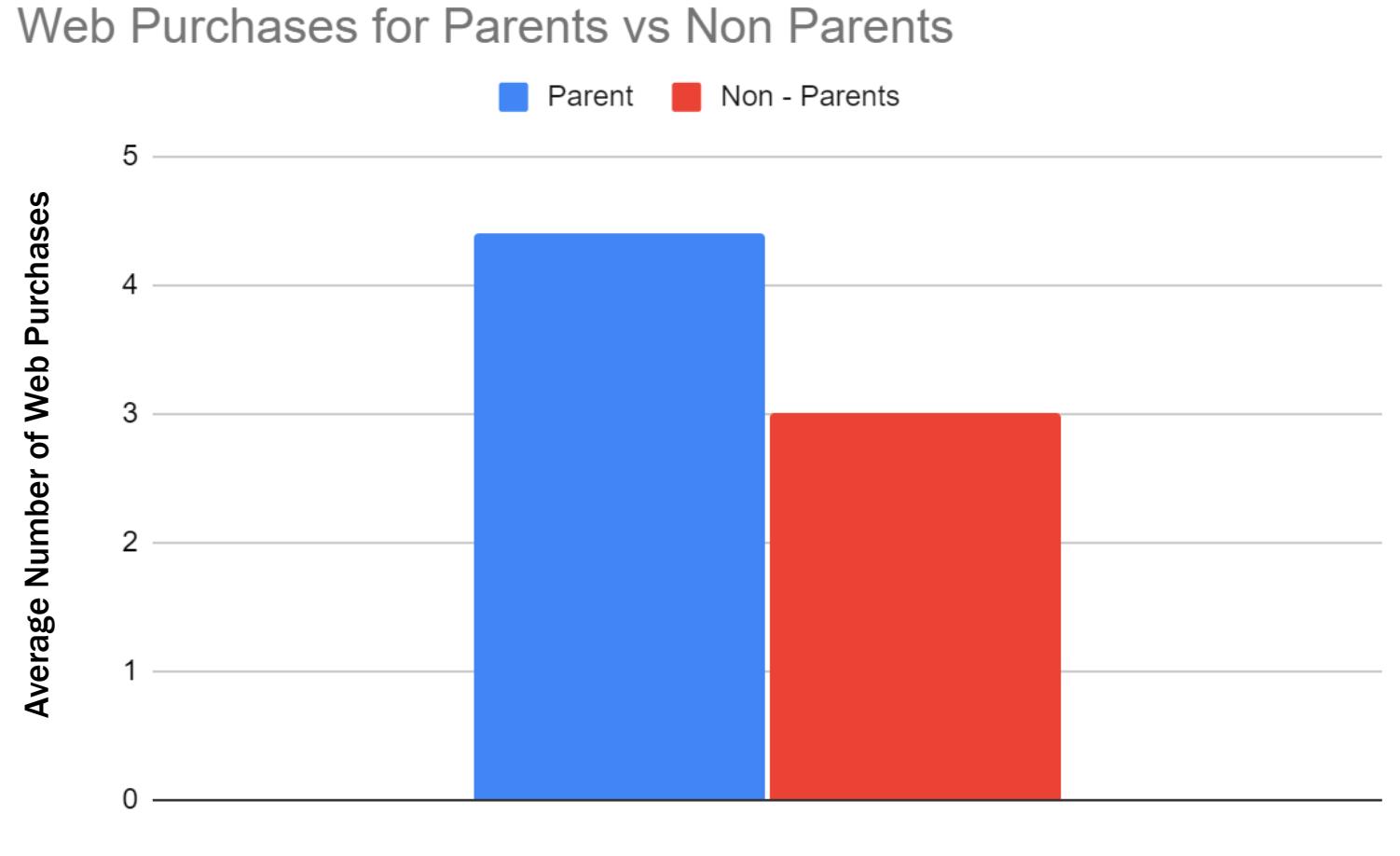
MntWine s	MntFruit s	MntMea tProduct s	MntFish Products	MntSweet Products	MntGol dProds
452	33	360	59	35	42

Non-Parents
Count: 633



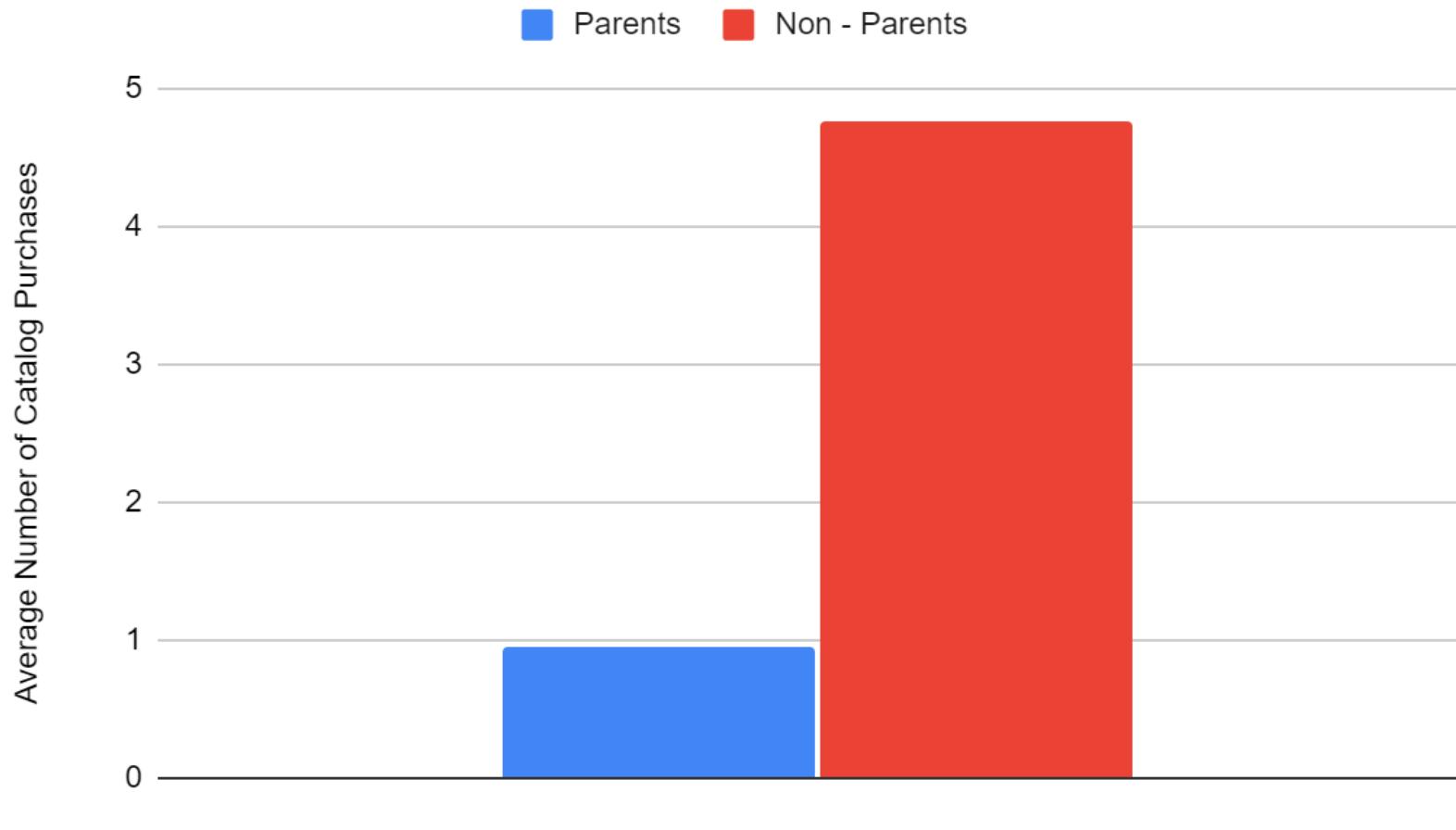
But Where Do They Shop From?

We tried to analyze if having a family also affects where one buys their products from.



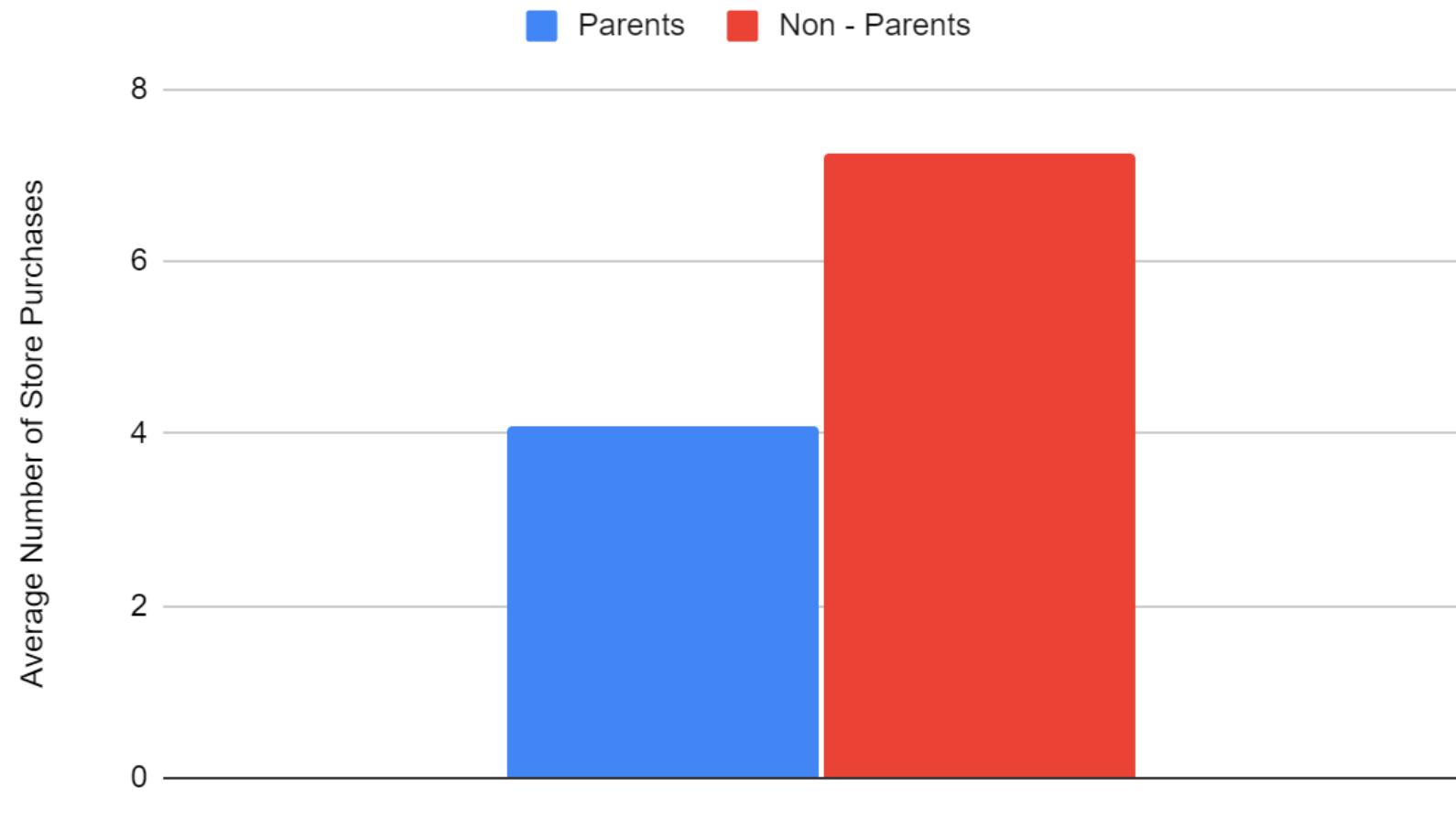
But Where Do They Shop From?

Catalog Purchases for Parents vs Non Parents



But Where Do They Shop From?

Store Purchases for Parents vs Non Parents



Final Thought

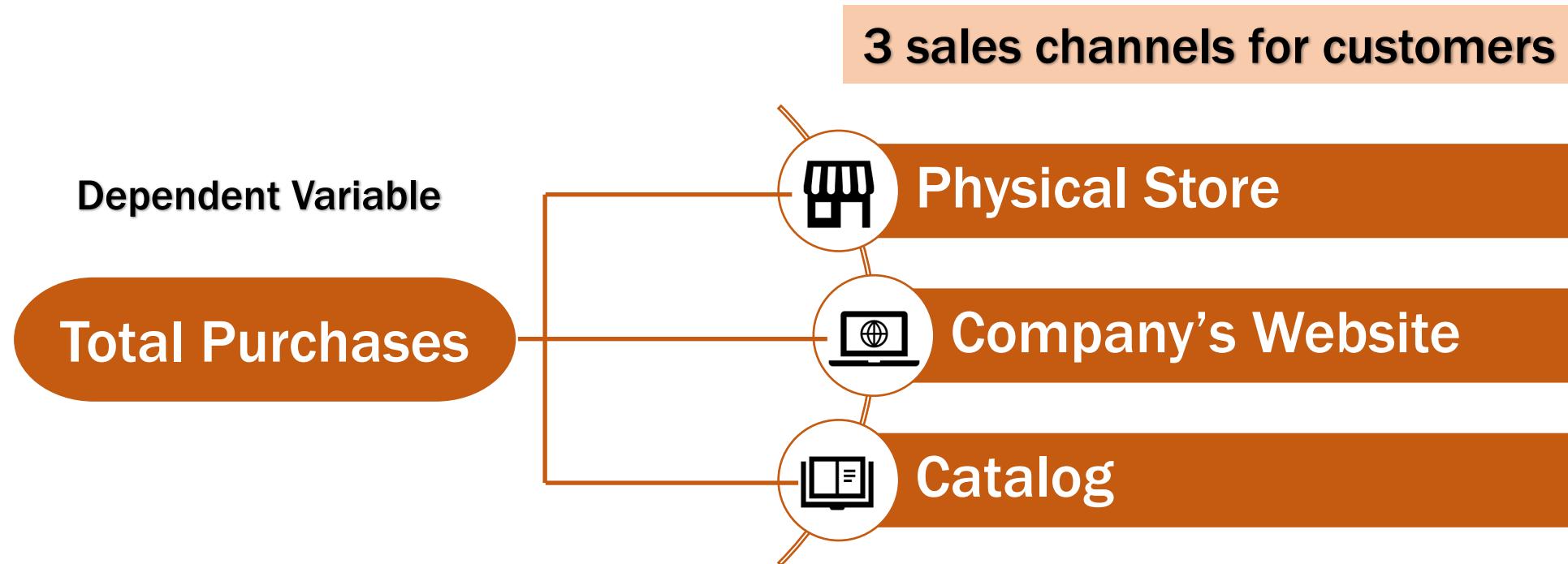
Having a family does not necessarily change what the client purchases, but it does affect which channel they choose to purchase from.



A yellow circular icon resembling a bomb or a bombshell, positioned on the left side of the slide. It has a yellow circle with a slightly darker yellow outline. From the top of the circle, five short, curved yellow lines radiate outwards, suggesting motion or an explosion.

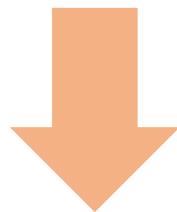
Predictive Models

Predictive Model A: No. of Purchases vs Sales Channel



Predictive Model A: No. of Purchases vs Sales Channel

Independent Variables:
30 variables



Stepwise Regression:

Stepwise regression performs selection based on statistical criteria to add or remove variables from the model until the best-fitting model is achieved.

Income	Response
Kidhome	Age
Teenhome	Customer_Days
Recency	marital_Divorced
MntWines	marital_Married
MntFruits	marital_Single
MntMeatProducts	marital_Together
MntFishProducts	marital_Widow
MntSweetProducts	education_2n Cycle
MntGoldProds	education_Basic
NumDealsPurchases	education_Graduation
NumWebVisitsMonth	education_Master
Complain	education_PhD
Z_CostContact	MntTotal
Z_Revenue	MntRegularProds

Predictive Model A: Stepwise Regression: 30 > 14 Variables

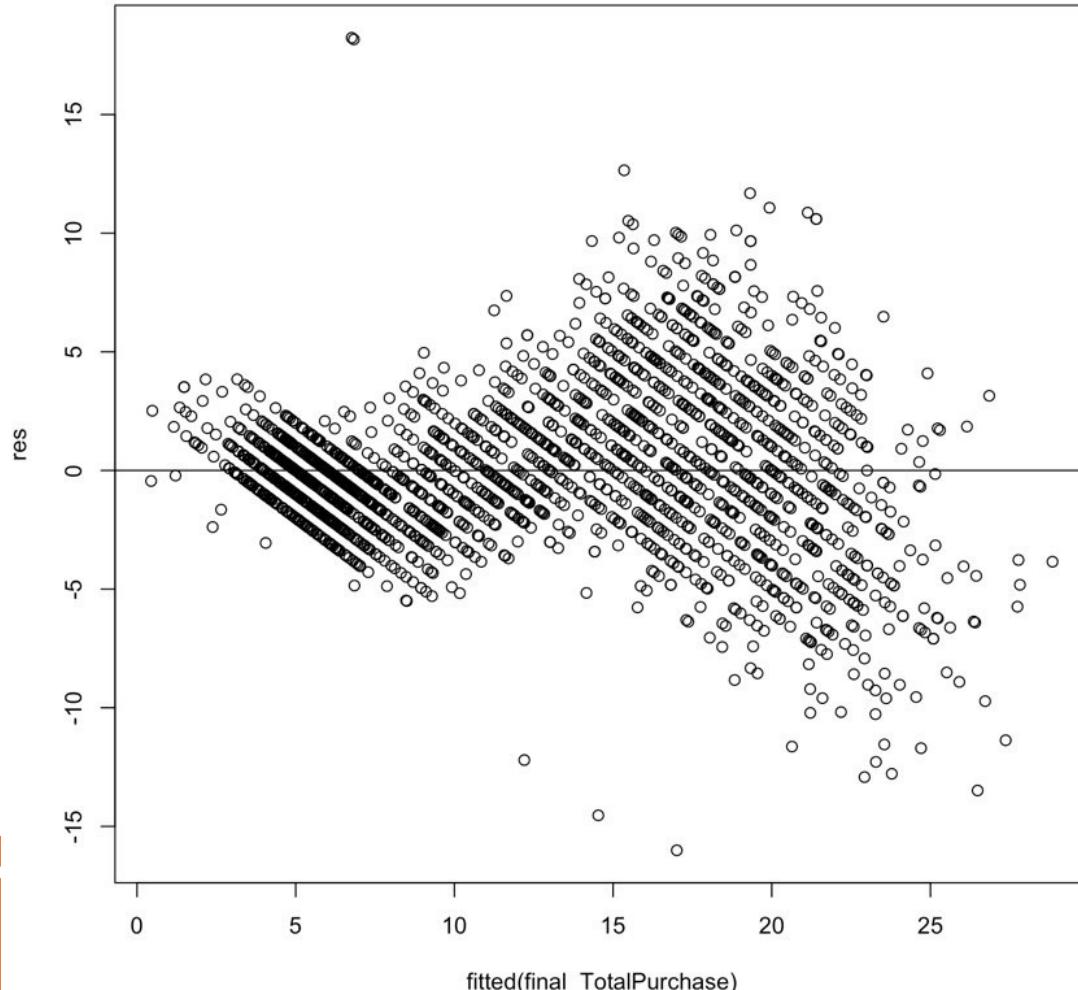
Coefficients:					
	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-1.0070	1.0450	-0.9640	0.3352	
Income	0.0001	0.0000	17.3720	< 2e-16	***
Kidhome	-2.3730	0.1721	-13.7870	< 2e-16	***
Teenhome	-0.4282	0.1654	-2.5890	0.0097	**
MntWines	0.0062	0.0003	17.9290	< 2e-16	***
MntFruits	0.0094	0.0025	3.7700	0.0002	***
MntMeatProducts	0.0014	0.0006	2.5380	0.0112	*
MntFishProducts	0.0072	0.0019	3.7950	0.0002	***
MntSweetProducts	0.0118	0.0024	4.9060	0.0000	***
MntGoldProds	0.0135	0.0017	8.1340	0.0000	***
NumDealsPurchases	0.8311	0.0457	18.2020	< 2e-16	***
Response	-0.8131	0.2138	-3.8030	0.0001	***
Customer_Days	0.0014	0.0004	3.5760	0.0004	***
education_Graduation	-0.2332	0.1637	-1.4250	0.1544	
education_Master	-0.3446	0.2162	-1.5940	0.1110	

- **Income, Kidhome, MntWines, NumDealsPurchases** are strongly significant
- **Education** is not statistically significant to the model

Legend on the significance level:

0 : ***
 0.001: **
 0.01: *
 0.05
 0.1

Predictive Model A: Residual Plot & Result



Residual Standard Error: 3.353

R-squared: 0.7833

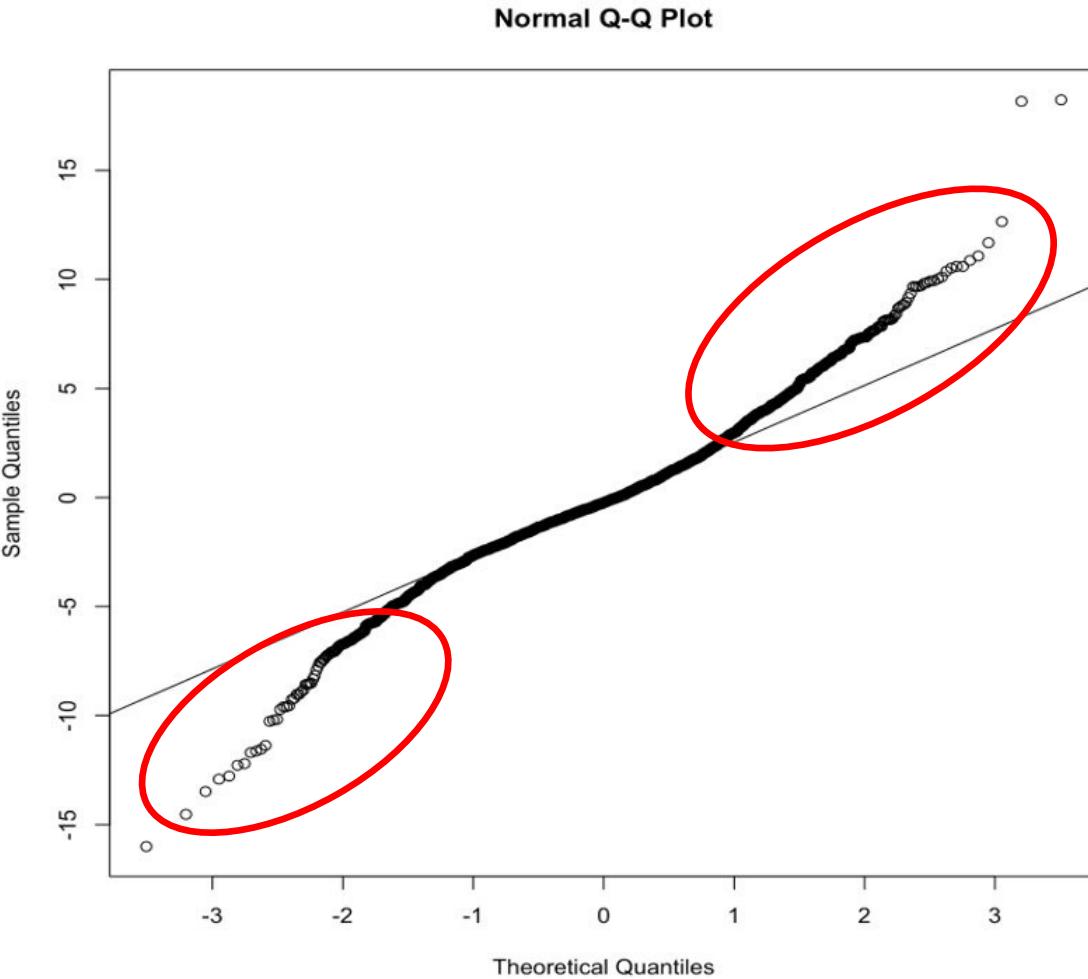
Adjusted R-squared: 0.782

F-statistics: 565.6

P-value: $< 2 \times 10^{-16}$

- ✓ 78.33% of the variability in the dependent variable is accounted for by the independent variables.
- ✓ The p-value is less than 0.0000 indicating that the model is statistically significant.

Predictive Model A: Normal QQ Plot



A normal Q-Q plot is a graphical tool used to assess whether a dataset follows a normal distribution.

- ✓ The largest values are larger than would be expected
- ✓ The smallest values are smaller than would be expected

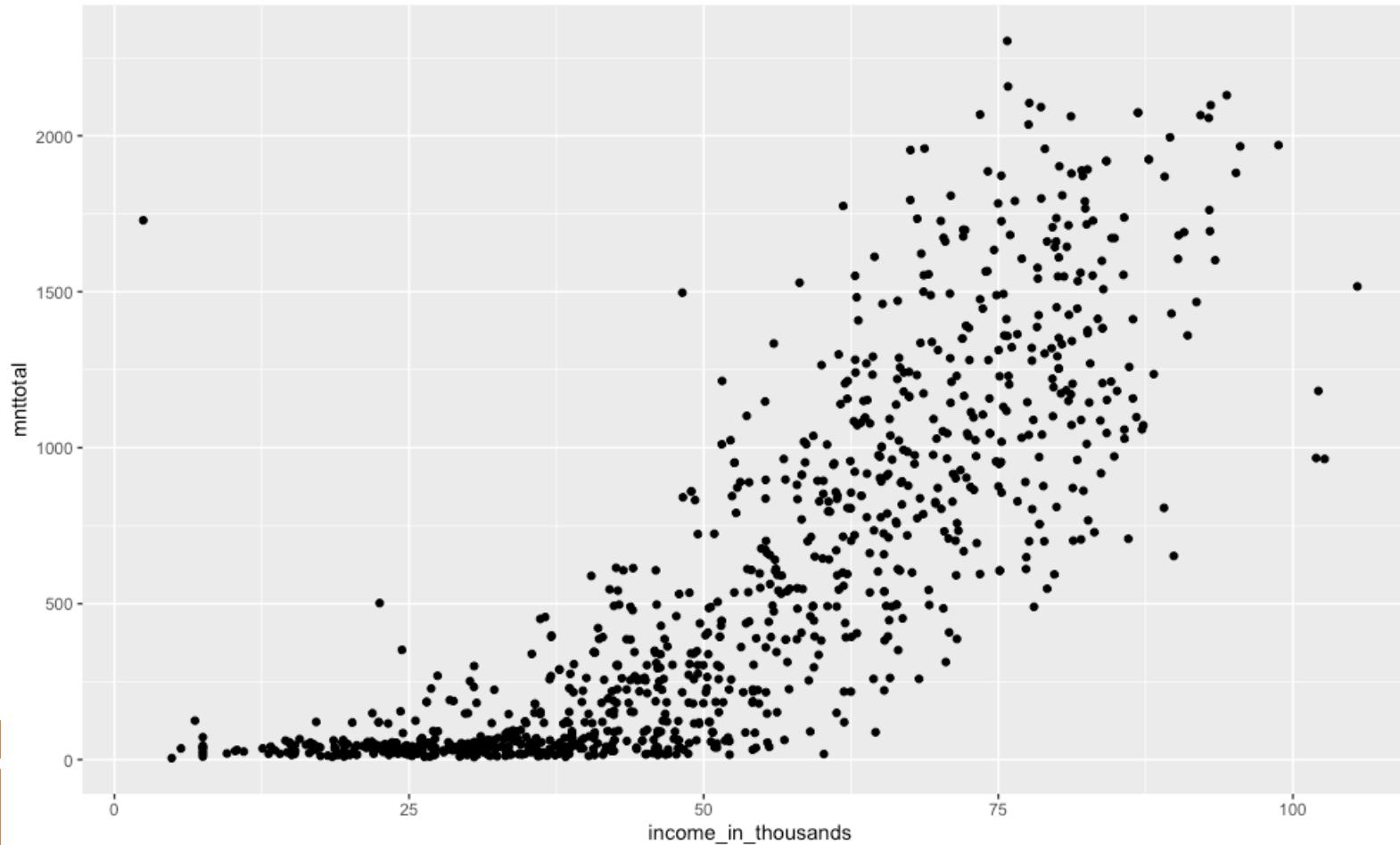
Predictive Model A: Amount Spent on Wines vs No. of Purchases

Coefficients:					
	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-1.0070	1.0450	-0.9640	0.3352	
MntWines	0.0062	0.0003	17.9290	< 2e-16	***

	# of Purchases below median	# of Purchases above median
Average spent on Wines	\$63.5415	\$ 553.0055

Customers spend more on wines, the more purchases they make!!

Predictive Model B: Income vs Total Spend (mnttotal)



Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-611.057	27.731	-22.04	<2e-16 ***
income_in_thousands	22.709	0.495	45.88	<2e-16 ***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

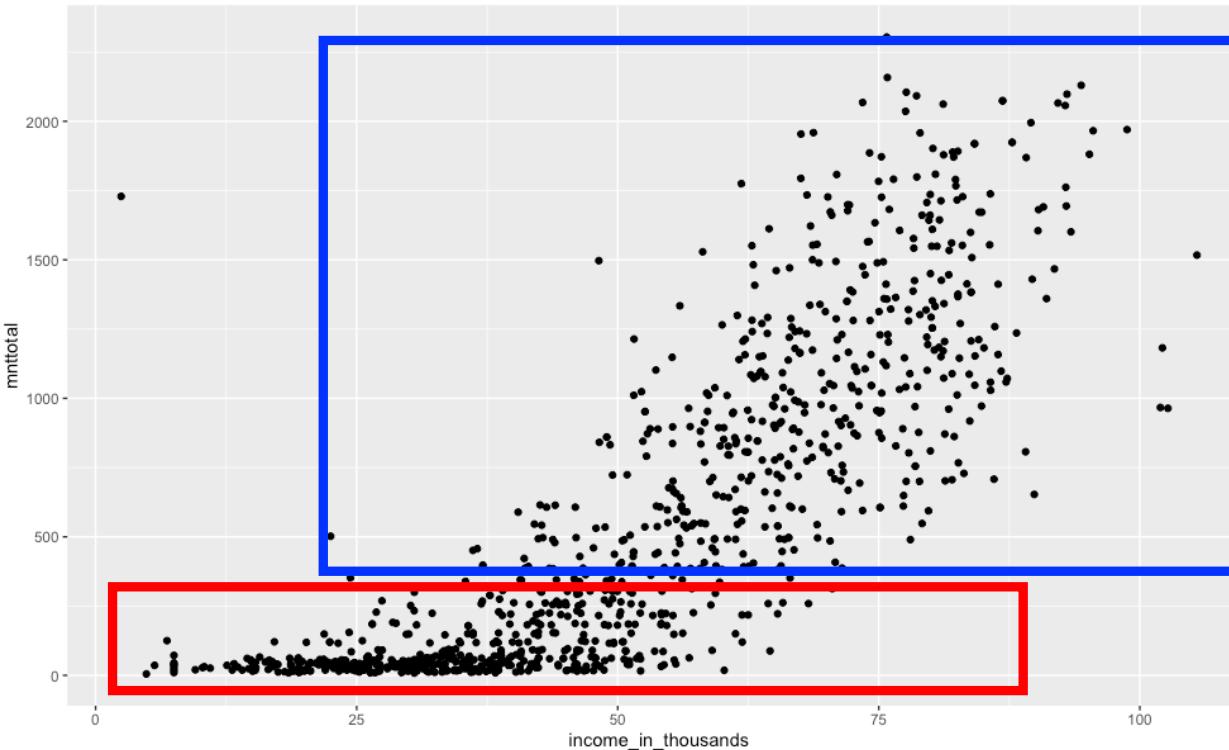
Residual standard error: 326.6 on 998 degrees of freedom

Multiple R-squared: 0.6783, Adjusted R-squared: 0.678

F-statistic: 2104 on 1 and 998 DF, p-value: < 2.2e-16

Predictive Model B: Two Issues

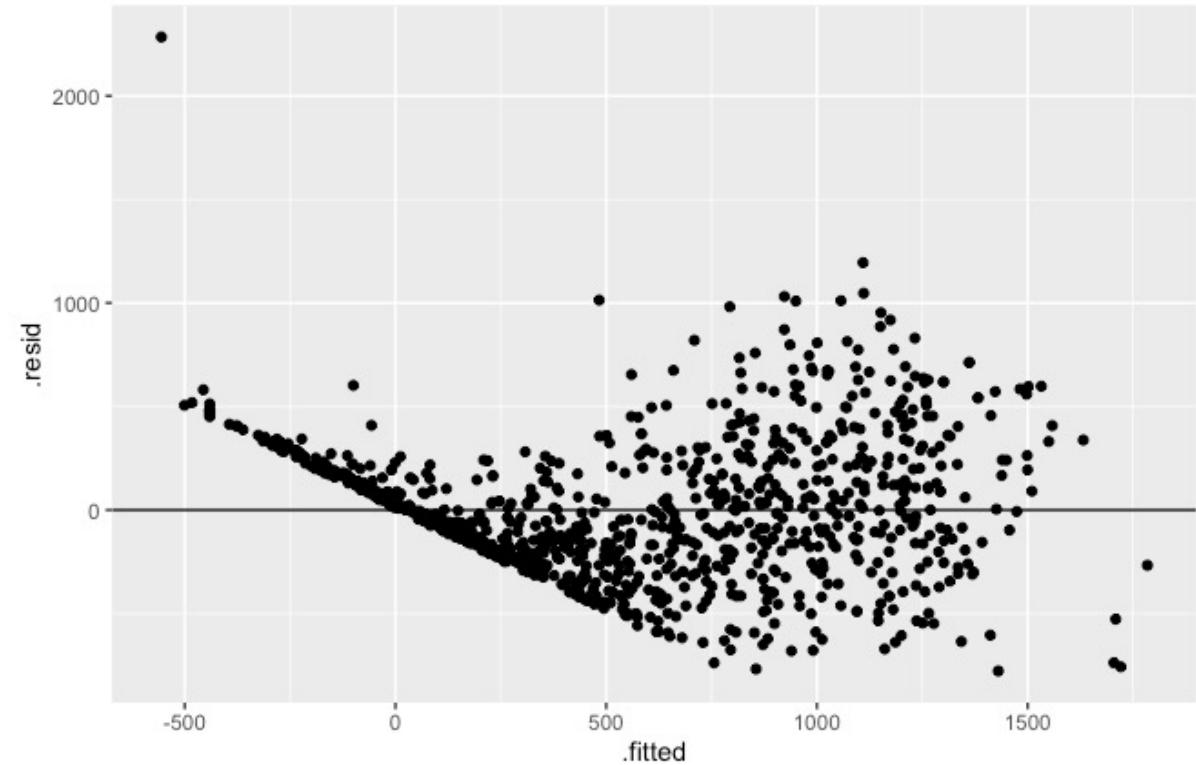
Two groups that behaved differently:



Make it rain

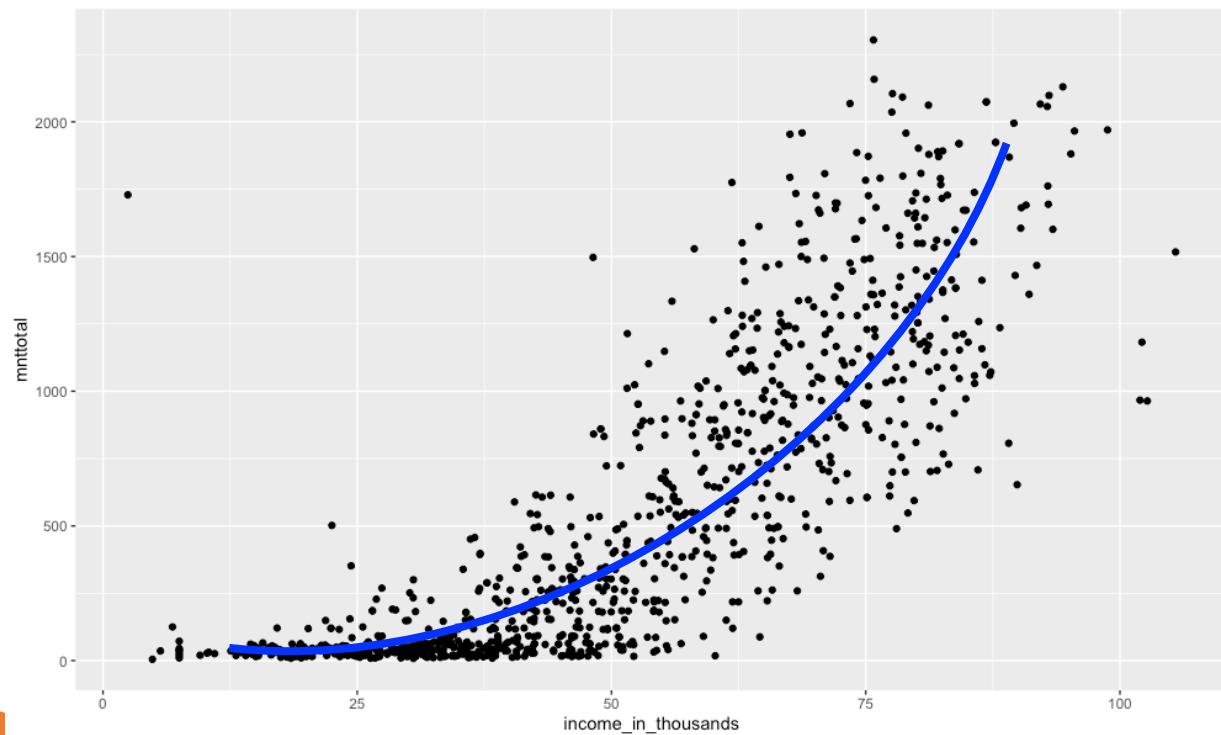
No money, no groceries

Heteroskedasticity (residuals are not normally distributed):

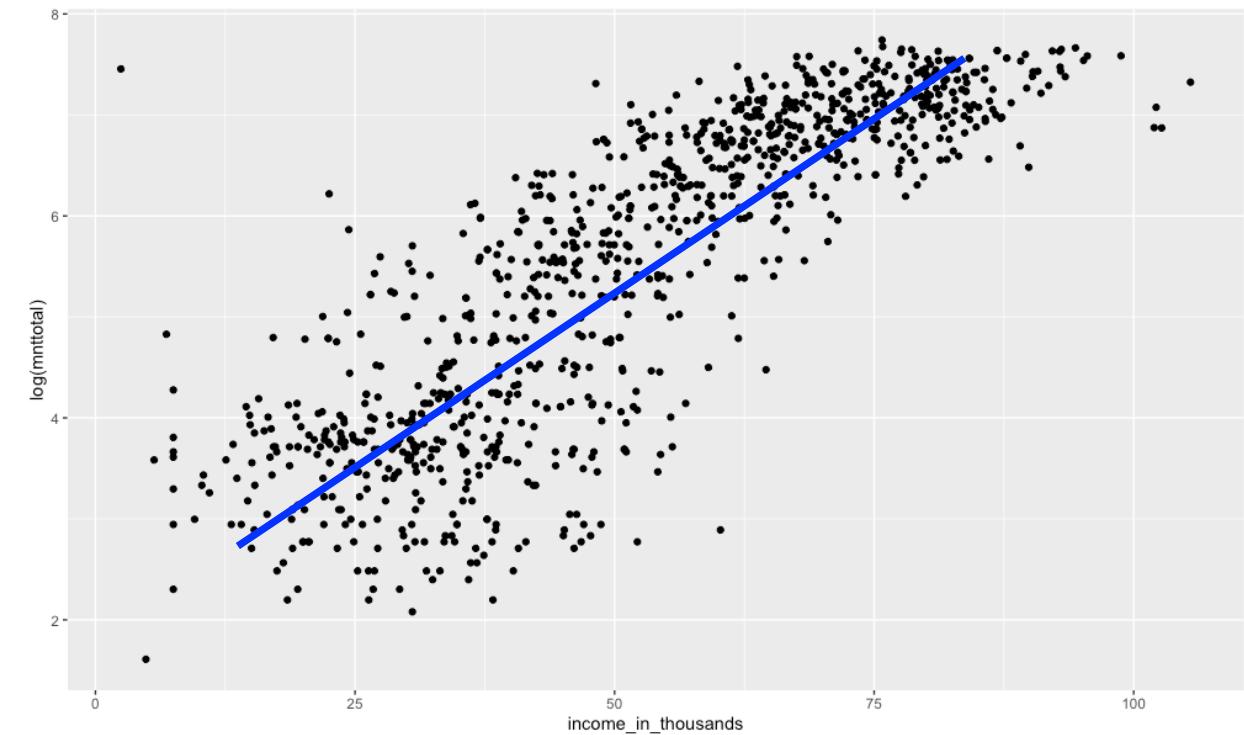


Predictive Model B: Solve for Heteroskedasticity

Raw (mnttotal):



Log (mnttotal):



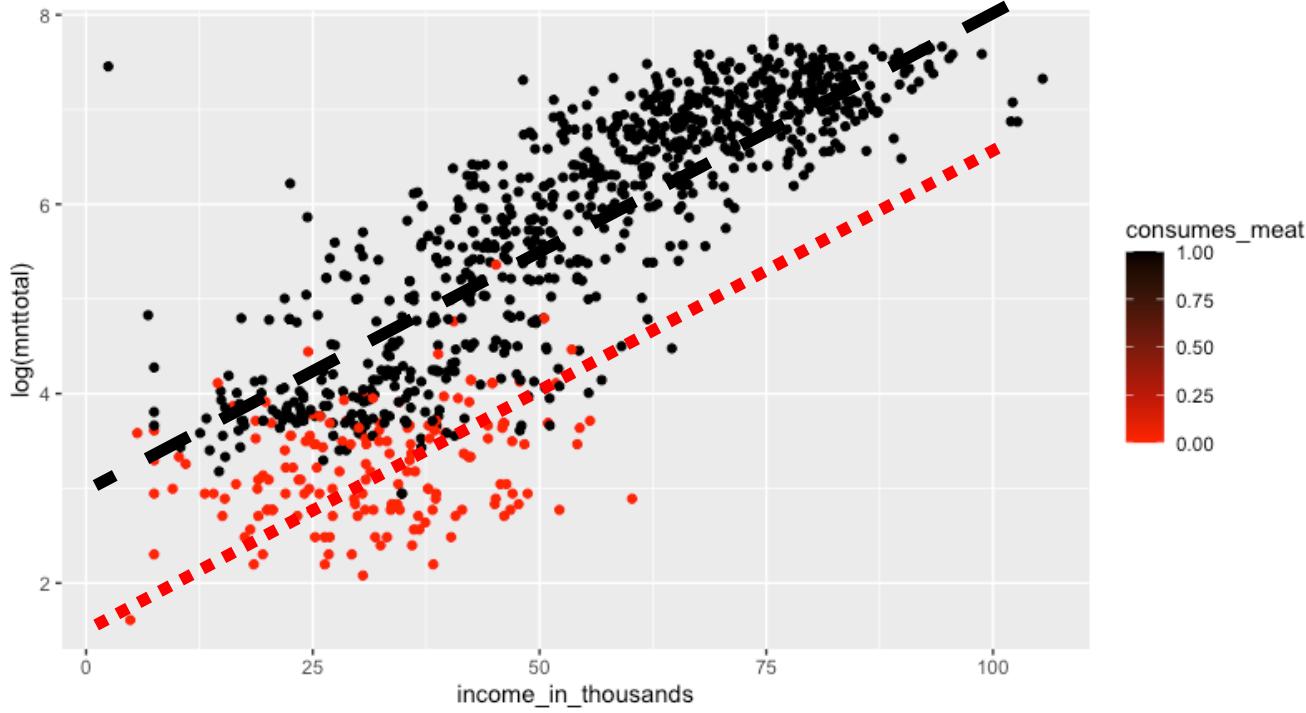
Predictive Model B: Solve for Differing Behavior

Guess this variable:



Predictive Model B: Final Model

Log(mnttotal) vs. Income and Consumes Meat



Regression Output

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.611101	0.063181	25.50	<2e-16 ***
income_in_thousands	0.051203	0.001127	45.43	<2e-16 ***
consumes_meat	1.464287	0.063341	23.12	<2e-16 ***

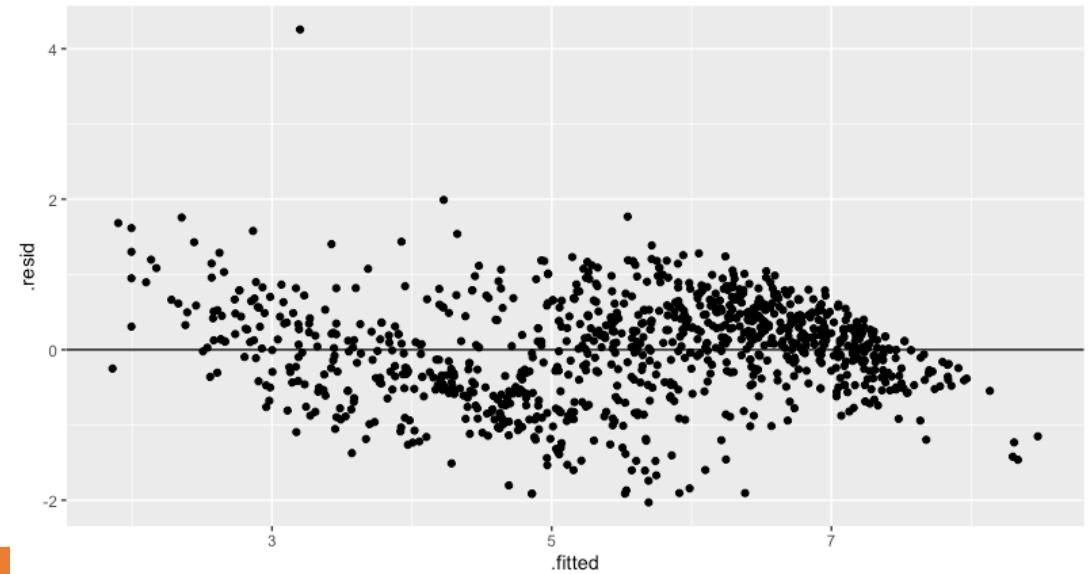
Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 0.67 on 997 degrees of freedom

Multiple R-squared: 0.8126, Adjusted R-squared: 0.8122

F-statistic: 2161 on 2 and 997 DF, p-value: < 2.2e-16

Fitted vs. Residuals



Conclusion: Business Applications

- Acquire High-Income Customers >> More Purchases and LTV
- Meat and Wine consumption = Proxy for Income
 - Align campaigns to social and cultural norms that influence meat/wine consumption
- Family status and education level are not as predictive as spending power
- Stakeholders should target customers based on income
 - Rather than what they purchase or where they purchase



Questions?



Appendix A

```
-- Clean Data
#Remove the row with Income blank
#Remove the row with Marital Status "Absurd" and "Yolo"
#Change the marital status of the customer with "Alone" to "Single"
#Add a column of "Have_Kid_Or_Not", described as whether the customer has kids or teens at home
#Change the education level of the customer with "2n Cycle" to "Master"
create table ifood_data_cleaned_1
as
select *,
case
when Education = "2n Cycle" then "Master"
else Education
end as Education_clean
from
(select *,
case
when Kidhome = 0 and Teenhome = 0 then "No"
else "Yes"
end as Have_Kid_Or_Not
from
(select *,
case
when Marital_Status = "Alone" then "Single"
else Marital_Status
end as Marital_Status_clean
from
(select *,
case
when Income = " then null
else Income
end as Income_clean
from ifood_data) as sub1) as sub2) as sub3
where Income_clean is not null and Marital_Status <> "YOLO" and Marital_Status <> "Absurd";
```

```
#Check if there're any duplicate IDs
select count(*) as total_count, count(distinct ID) as
count_remove_duplicate
from ifood_data_cleaned_1;
```

Appendix A



-- Explore the influence of Education Levels on Grocery Spending

#Overall picture

```
select Education_clean, avg_income, wines_spending + fruits_spending + meat_spending +
fish_spending + sweet_product + gold_spending as avg_total,
      wines_spending, fruits_spending, meat_spending, fish_spending, sweet_product,
      gold_spending
from
(select Education_clean, avg(Income) as avg_income,
      avg(MntWines) as wines_spending, avg(MntFruits) as fruits_spending,
      avg(MntMeatProducts) as meat_spending, avg(MntFishProducts) as fish_spending,
      avg(MntSweetProducts) as sweet_product, avg(MntGoldProds) as gold_spending
from ifood_data_cleaned_1
group by Education_clean
order by avg_income desc) as sub
order by avg_total desc;
```

#Customer Profile for different products

```
##Wine (Top1 product having the most spending amount)
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,
      avg(MntWines) as wines_spending
from ifood_data_cleaned_1
group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not
order by wines_spending desc
limit 5;
##Meat (Top2 product having the most spending amount)
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,
      avg(MntMeatProducts) as meat_spending
from ifood_data_cleaned_1
group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not
order by meat_spending desc
limit 5;
```

##Gold

```
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,
      avg(MntGoldProds) as gold_spending
from ifood_data_cleaned_1
group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not
order by gold_spending desc
limit 5;
```

##Fish

```
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,
      avg(MntFishProducts) as fish_spending
from ifood_data_cleaned_1
group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not
order by fish_spending desc
limit 5;
```

##Sweet

```
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,
      avg(MntSweetProducts) as sweet_spending
from ifood_data_cleaned_1
group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not
order by sweet_spending desc
limit 5;
```

##Fruit

```
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,
      avg(MntFruits) as fruit_spending
from ifood_data_cleaned_1
group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not
order by fruit_spending desc
limit 5;
```

Appendix A



```
-- Explore the influence of Education Levels on choosing Channel Purchase
```

```
#Overall picture
```

```
select Education_clean, web, catalog, store,  
      web + catalog + store as total, deal  
  from  
(select Education_clean, avg(NumWebPurchases) as web, avg(NumCatalogPurchases) as catalog,  
      avg(NumStorePurchases) as store,  
      avg(NumDealsPurchases) as deal  
  from ifood_data_cleaned_1  
  group by Education_clean) as sub  
  group by Education_clean  
order by total desc;
```

```
#Customer Profile for different channels
```

```
##Web
```

```
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,  
      avg(NumWebPurchases) as web  
  from ifood_data_cleaned_1  
  group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not  
order by web desc  
limit 5;
```

```
##Store
```

```
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,  
      avg(NumStorePurchases) as store  
  from ifood_data_cleaned_1  
  group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not  
order by store desc  
limit 5;
```

```
##Catalog
```

```
select Education_clean, Marital_Status_clean, Have_Kid_Or_Not,  
      avg(NumCatalogPurchases) as catalog  
  from ifood_data_cleaned_1  
  group by Education_clean, Marital_Status_clean, Have_Kid_Or_Not  
order by catalog desc  
limit 5;
```

Appendix B

#count of parents

```
select count(ID) as Parents
from ml_project1_data
where kidhome > 0 and teenhome > 0;
```

#count of non-parents

```
select count(ID) as Not_Parents
from ml_project1_data
where kidhome = 0 and teenhome = 0;
```

#amnt spent for parents for each category

```
select ID as Parents_ID, MntWines, MntFruits,
MntMeatProducts, MntFishProducts, MntSweetProducts,
MntGoldProds from ml_project1_data
```

where ID in (select ID from ml_project1_data

where kidhome > 0 and teenhome > 0)

order by ID;

#amnt spent for non-parents for each category

```
select ID as Non_Parents_ID, MntWines, MntFruits,
MntMeatProducts, MntFishProducts, MntSweetProducts,
MntGoldProds from ml_project1_data
```

where ID in (select ID from ml_project1_data

where kidhome = 0 and teenhome = 0)

order by ID;

Appendix B

#WEB - Parents

```
select ID as Parents, NumWebPurchases from  
ml_project1_data  
  
where ID in (select ID from ml_project1_data  
  
where kidhome > 0 and teenhome > 0 )  
  
order by ID;
```

###WEB - Non - Parents

```
select ID as Non_Parents, NumWebPurchases from  
ml_project1_data  
  
where ID in (select ID from ml_project1_data  
  
where kidhome = 0 and teenhome = 0 )  
  
order by ID;
```

#Catalog - Parents

```
select ID as Parents, NumCatalogPurchases from  
ml_project1_data  
  
where ID in (select ID from ml_project1_data  
  
where kidhome > 0 and teenhome > 0 )  
  
order by ID;
```

#Catalog - Non -Parents

```
select ID as Non_Parents, NumCatalogPurchases from  
ml_project1_data  
  
where ID in (select ID from ml_project1_data  
  
where kidhome = 0 and teenhome = 0 )  
  
order by ID;
```

#Store Purchases - Parents

```
select ID as Parents, NumStorePurchases, Income from  
ml_project1_data  
  
where ID in (select ID from ml_project1_data  
  
where kidhome > 0 and teenhome > 0 )  
  
order by ID;
```

#Store Purchases - Non - Parents

```
select ID as Non_Parents, NumStorePurchases from  
ml_project1_data  
  
where ID in (select ID from ml_project1_data  
  
where kidhome = 0 and teenhome = 0 )  
  
order by ID;
```

Appendix C

```
##R code for performing stepwise regression##
ifood <- read.csv("/Users/chaomanling/Desktop/MSBA/Summer/Data Management and
SQL/iFood_df.csv", header=T)
head(ifood)

###Total number of purchase #R=0.7833
ifood$sum_purchase <-
ifood$NumStorePurchases + ifood$NumWebPurchases + ifood$NumCatalogPurchases
x4 <- colnames(ifood[, -c(12:14, 16:20, 39, 40)])
formula_str4 <- paste("sum_purchase~", paste(x4, collapse = " + "))
TotalPurchase <- lm(formula_str4, data = ifood)
final_TotalPurchase <- step(TotalPurchase, direction = "both")
summary(final_TotalPurchase)

#residula plot
res <- resid(final_TotalPurchase)
plot(fitted(final_TotalPurchase), res)
abline(0,0)

#create Q-Q plot for residuals
qqnorm(res)
#add a straight diagonal line to the plot
qqline(res)
```

Appendix C

```
##My SQL Code##
ALTER TABLE ifood ADD sum_purchase INT;

UPDATE ifood
SET sum_purchase = NumWebPurchases + NumCatalogPurchases + NumStorePurchases;

#Compute the median of total purchase
SELECT AVG(sum_purchase) AS Median
FROM (
    SELECT sum_purchase, @rowindex:="@rowindex + 1 AS rowindex,
    @total_rows AS total_rows
    FROM ifood, (SELECT @rowindex := 0, @total_rows := (SELECT COUNT(*)
    FROM ifood)) r
    ORDER BY sum_purchase
) AS d
WHERE rowindex IN (FLOOR(total_rows / 2), CEIL(total_rows / 2));
```

Appendix C

```
#Compute the average amount spend on wines that purchases below the median
SELECT AVG(MntWines) AS avg_wines_lower
FROM ifood
WHERE sum_purchase <= (
    SELECT AVG(sum_purchase) AS Median
    FROM (
        SELECT sum_purchase, @rowindex:="@rowindex + 1 AS rowindex,
        @total_rows AS total_rows
        FROM ifood, (SELECT @rowindex := 0, @total_rows := (SELECT COUNT(*)
        FROM ifood)) r
        ORDER BY sum_purchase
    ) AS d
    WHERE rowindex IN (FLOOR(total_rows / 2), CEIL(total_rows / 2))
);
#Compute the average amount spend on wines that purchases above the median
SELECT AVG(MntWines) AS avg_wines_upper
FROM ifood
WHERE sum_purchase > (
    SELECT AVG(sum_purchase) AS Median
    FROM (
        SELECT sum_purchase, @rowindex:="@rowindex + 1 AS rowindex,
        @total_rows AS total_rows
        FROM ifood, (SELECT @rowindex := 0, @total_rows := (SELECT COUNT(*)
        FROM ifood)) r
        ORDER BY sum_purchase
    ) AS d
    WHERE rowindex IN (FLOOR(total_rows / 2), CEIL(total_rows / 2))
);
```

Appendix D

SQL script to clean data for regression analysis:

```
#DROP TABLE ifood_clean;
```

```
CREATE TABLE ifood_clean AS (
SELECT
    income
    , income / 1000 AS income_in_thousands
    , kidhome
    , teenhome
    , recency
    , mntwines
    , CASE
        WHEN mntwines > 10 THEN 1
        ELSE 0
    END AS consumes_wine
    , mntfruits
    , CASE
        WHEN mntfruits > 10 THEN 1
        ELSE 0
    END AS consumes_fruits
    , mntmeatproducts
    , CASE
        WHEN mntmeatproducts > 10 THEN 1
        ELSE 0
    END AS consumes_meat
    , mntfishproducts
    , CASE
        WHEN mntfishproducts > 10 THEN 1
        ELSE 0
    END AS consumes_fish
    , mntsweetproducts
```

```
, CASE
    WHEN mntsweetproducts > 10 THEN 1
    ELSE 0
END AS consumes_sweets
, mntgoldprods
, CASE
    WHEN mntgoldprods > 10 THEN 1
    ELSE 0
END AS consumes_gold
, NumDealsPurchases AS num_deals_purchases
, NumWebPurchases AS num_web_purchases
, NumCatalogPurchases AS num_catalog_purchases
, NumStorePurchases AS num_store_purchases
, NumWebVisitsMonth AS num_web_visits_month
, AcceptedCmp1 AS accepted_cmp1
, AcceptedCmp1 AS accepted_cmp2
, AcceptedCmp1 AS accepted_cmp3
, AcceptedCmp1 AS accepted_cmp4
, AcceptedCmp1 AS accepted_cmp5
, complain , Z_CostContact AS z_cost_contact
, z_revenue
, response
, age
, customer_days
, marital_divorced
, marital_single
, marital_together
, marital_widow
, 'education_2n Cycle' AS education_2n_cycle
, education_basic
, education_graduation
, education_master
, education_phd
, mnntotal
, mnntregularprods
, AcceptedCmpOverall AS accepted_cmp_overall
FROM ifood_df_full
)
```

Appendix D



R code to create linear models and plots:

```
#install.packages("tidyverse")
#install.packages("SciViews")
library(ggplot2)
library("SciViews")
```

```
#regression, no log no meat
lm0 = lm(mnttotal~income_in_thousands,data=ifood_clean)
summary(lm0)
res = resid(lm0)
plot(fitted(lm0), res)
ggplot(data=lm0,aes(x=.fitted,y=.resid)) + geom_point() +
geom_hline(yintercept = 0)
```

```
#regression, no log on mnttotal
lm1 =
lm(mnttotal~income_in_thousands+consumes_meat,data=ifoo
d_clean)
summary(lm1)
res = resid(lm1)
plot(fitted(lm1), res)
```

```
#regression, log on mnttotal
lmlog = lm(log(mnttotal)~income_in_thousands+consumes_meat,data=ifood_clean)
summary(lmlog)
res = resid(lmlog)
plot(fitted(lmlog), res)
ggplot(data=lmlog,aes(x=.fitted,y=.resid)) + geom_point() + geom_hline(yintercept = 0)
```

```
#no log on income
plot_no_color = ggplot(data=ifood_clean,
aes(x=income_in_thousands,y=mnttotal))+geom_point()
plot_no_color
plot = ggplot(data=ifood_clean,
aes(x=income_in_thousands,y=mnttotal,color=consumes_meat))+geom_point()
plot +geom_point(size=0.1)+scale_color_gradient(low = "red", high = "black")
```

```
#log income
plot_no_color2 = ggplot(data=ifood_clean,
aes(x=income_in_thousands,y=log(mnttotal)))+geom_point()
plot_no_color2
plot2 = ggplot(data=ifood_clean,
aes(x=income_in_thousands,y=log(mnttotal),color=consumes_meat))
plot2 +geom_point()+scale_color_gradient(low = "red", high = "black")
```

```
#hist of meat consumption
ggplot(data=ifood_clean, aes(x=mntmeatproducts,color=income))+geom_histogram()
```