

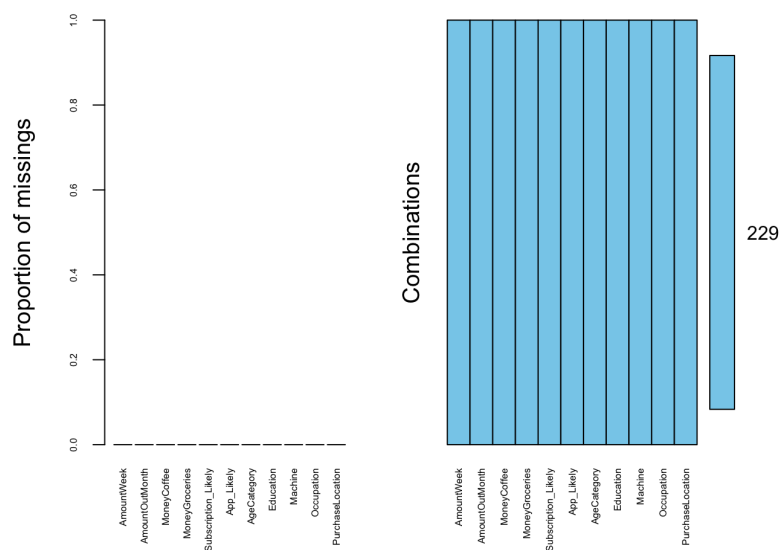
Thesis Data Analysis

13 March, 2021

Steps data analysis

- Univariate descriptions - categorical variables
 - Data table
 - Graphs
- Univariate descriptions - numerical variables
 - Summary
 - Confidence intervals
 - Graphs
- Boxplots - numerical
- Joint distribution tables
- Outliers
- Parametric testing
- Relationships & correlations
 - Residual plots
- Regressions
- Data problems

Introduction



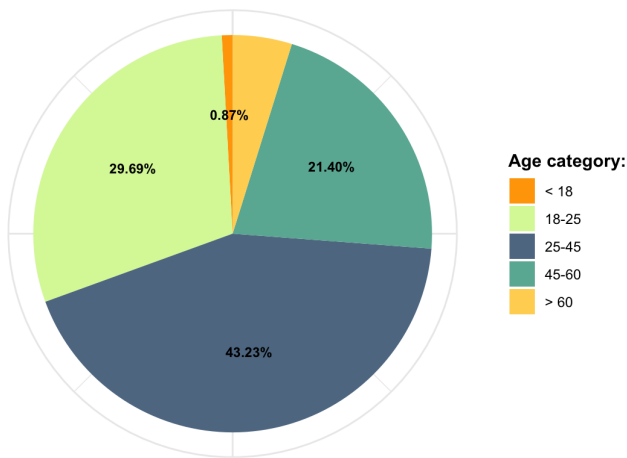
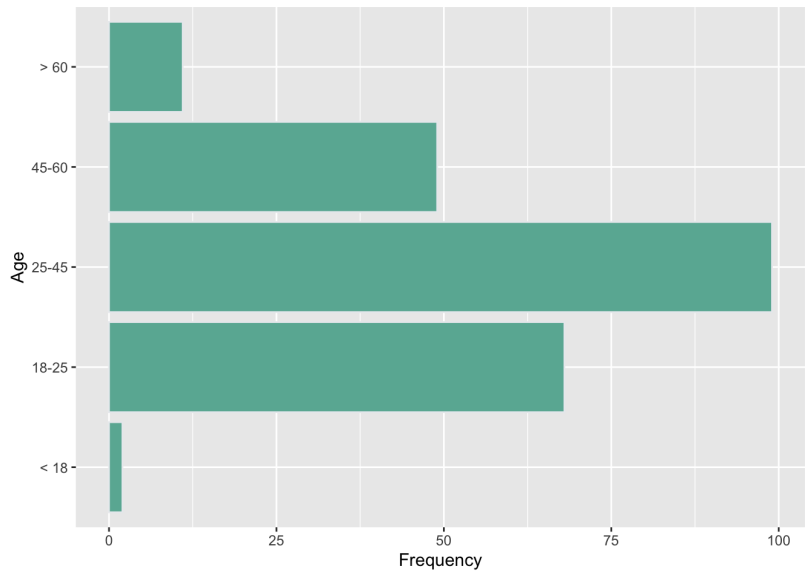
The variables included in the data set are:

Field	Description
AmountWeek	How many cups of coffee do you typically consume weekly?
AmountOutMonth	How frequently do you drink out-of-home per month on average?
MoneyCoffee	How much money on average do you estimate you spend on coffee per month?
MoneyGroceries	How much on average do you spend on general groceries per month?
Machine	How do you brew your coffee at home?
Brand change	How often do you switch between coffee brands?
Purchase location	Where do you usually purchase your coffee?
Supermarket_Positive_Reasons	When you purchase coffee from the supermarket what are your main reasons for doing so?
Supermarket_Negative_Reasons	What would be reasons why you would not purchase coffee from the supermarket?
Criteria_Type_Coffee	What are your main criteria's or evaluation points for choosing the type of coffee?
KnowledgeCoffee	How would you describe your knowledge level regarding coffee in general?
Purchase_Price	I believe that the ____ is important to my decision on which coffee to purchase.
Purchase_Sustainability	I believe that the ____ is important to my decision on which coffee to purchase.
Purchase_Sustainability	I believe that the ____ is important to my decision on which coffee to purchase.
Purchase_Fairtrade	I believe that the ____ is important to my decision on which coffee to purchase.
Purchase_Packaging	I believe that the ____ is important to my decision on which coffee to purchase.
Frequency_Specialty	How often do you drink specialty coffee?
Subscription_Likely	How likely are you to have an online subscription for (specialty) coffee?
Subscription_Not_Likely	What is the number one reasons why you would be hesitant?
App_Likely	How likely are you to value and use an app for your online subscription?
Gender	What is your gender?
AgeCategory	What is your age category?
Occupation	What is your occupational status?
Education	What level of education have you completed?
Home	How would you describe the place you currently live in?

Univariate descriptions - Categorical variables

Age category

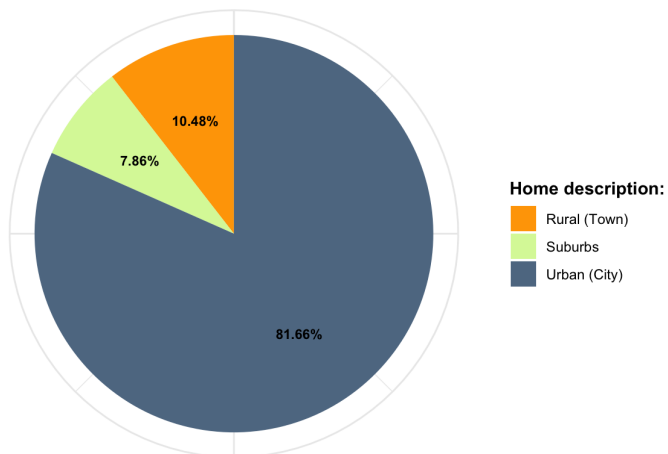
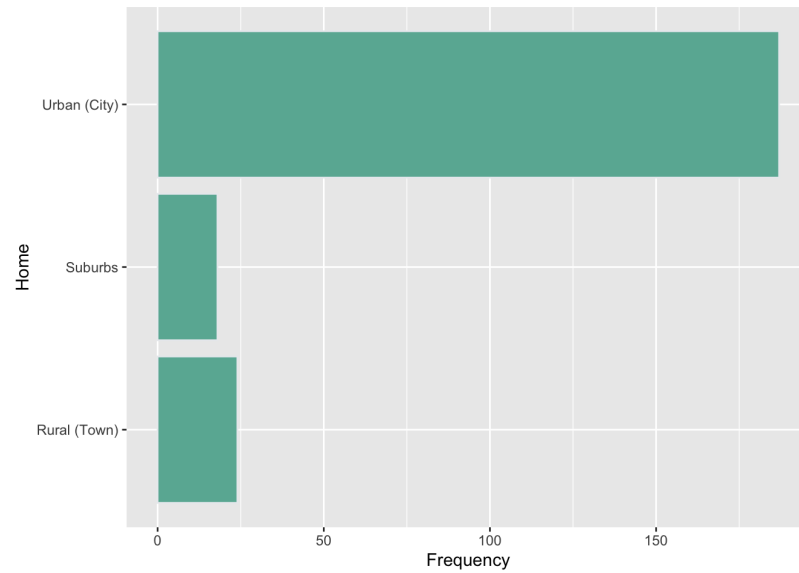
Age Category	Absolute	Relative
< 18	2	0.87%
18-25	68	29.69%
25-45	99	43.23%
45-60	49	21.40%
> 60	11	4.80%



Home

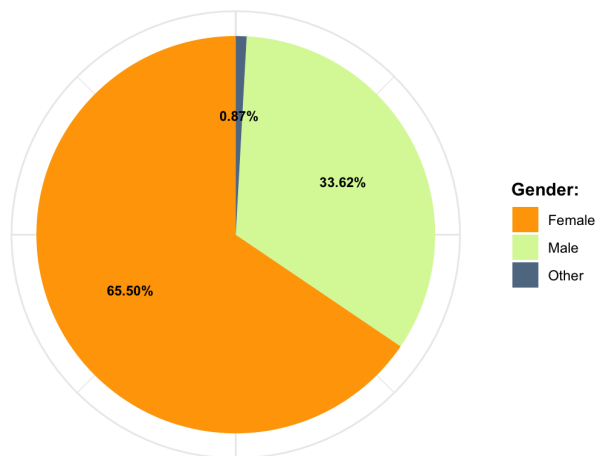
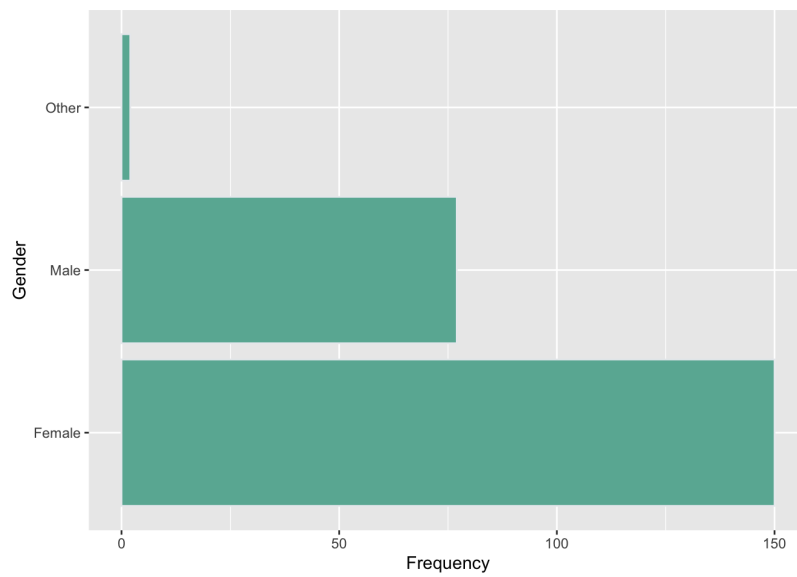


Home	Absolute	Relative
Rural (Town)	24	10.48%
Suburbs	18	7.86%
Urban (City)	187	81.66%



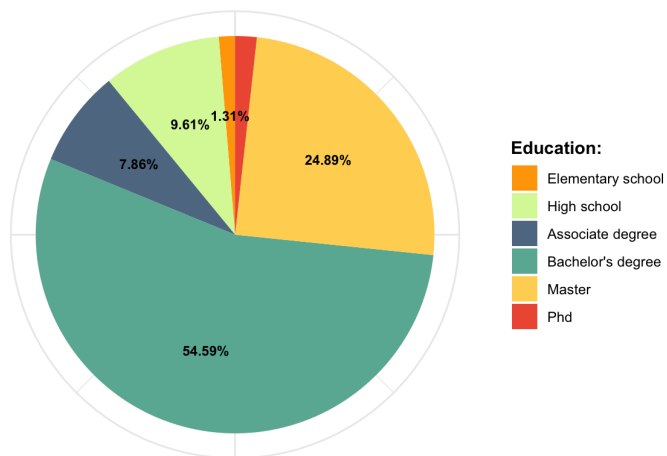
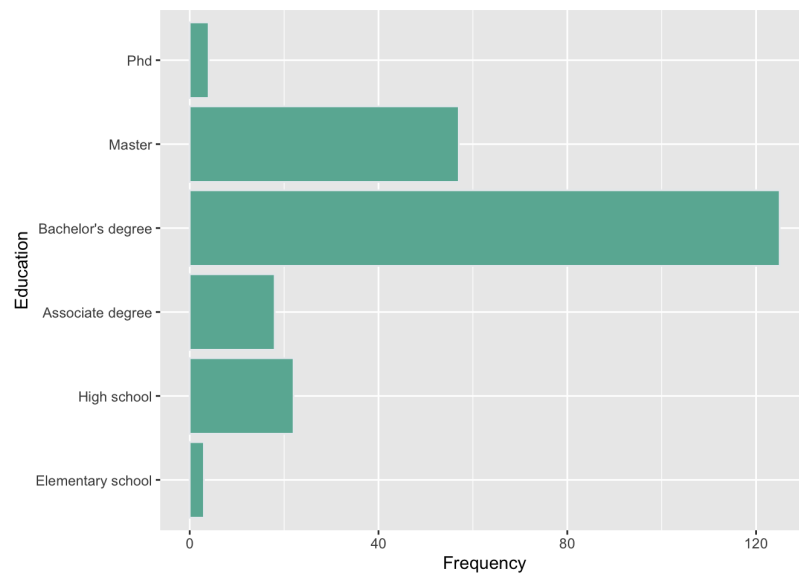
Gender

Gender	Absolute	Relative
Female	150	65.50%
Male	77	33.62%
Other	2	0.87%



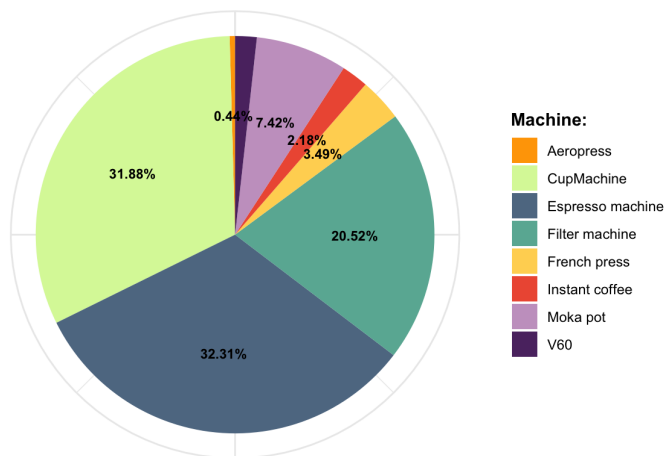
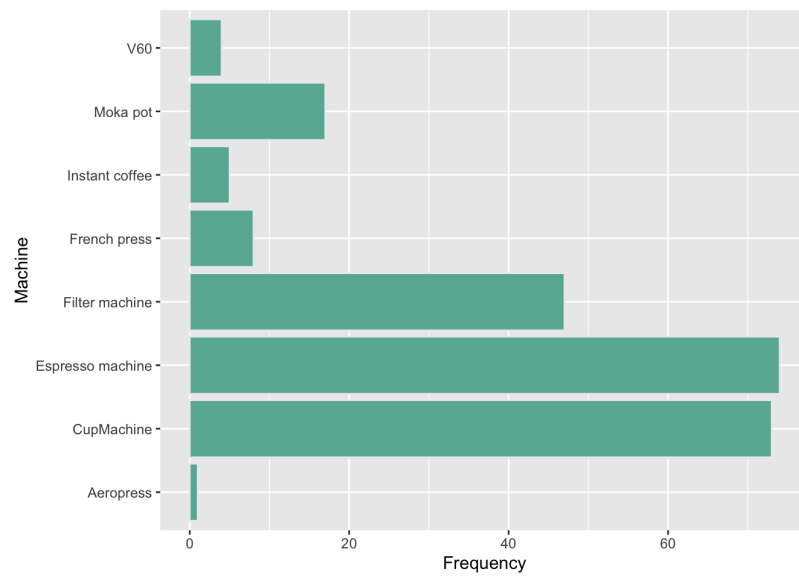
Education

Education	Absolute	Relative
Elementary school	3	1.31%
High school	22	9.61%
Associate degree	18	7.86%
Bachelor's degree	125	54.59%
Master	57	24.89%
Phd	4	1.75%



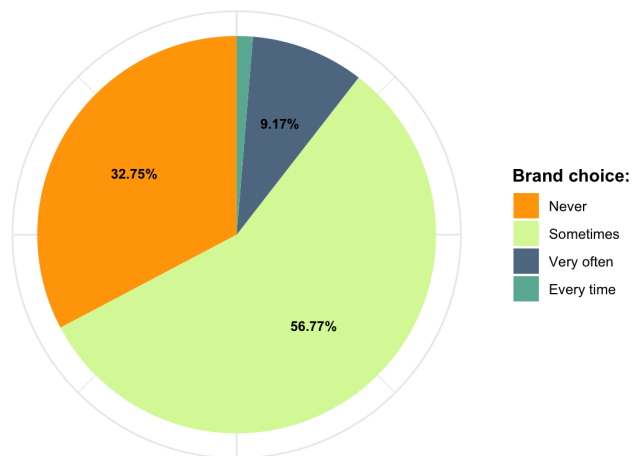
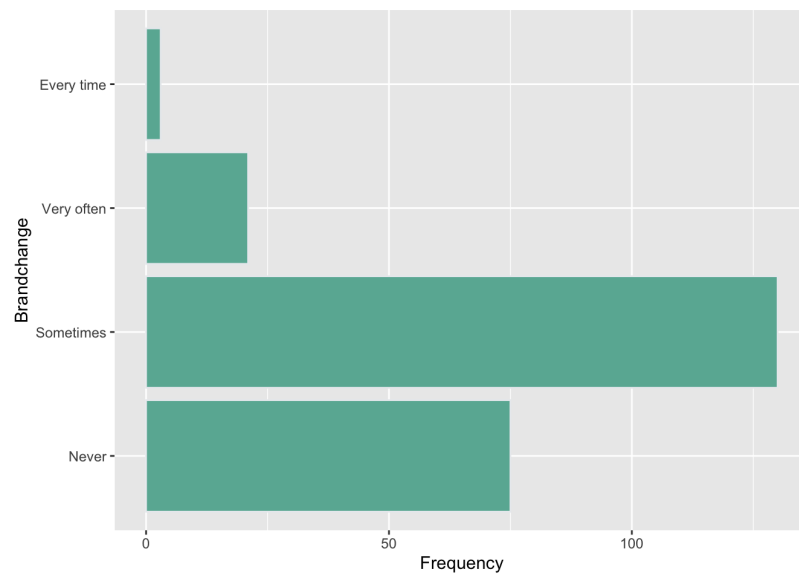
Machine

Machine	Absolute	Relative
Aeropress	1	0.44%
CupMachine	73	31.88%
Espresso machine	74	32.31%
Filter machine	47	20.52%
French press	8	3.49%
Instant coffee	5	2.18%
Moka pot	17	7.42%
V60	4	1.75%



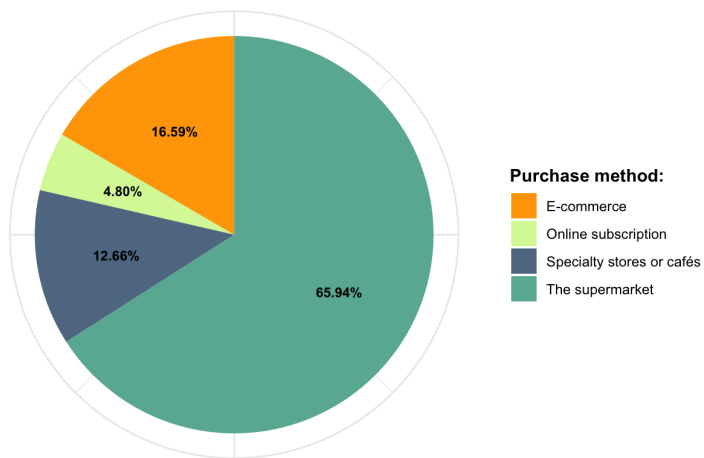
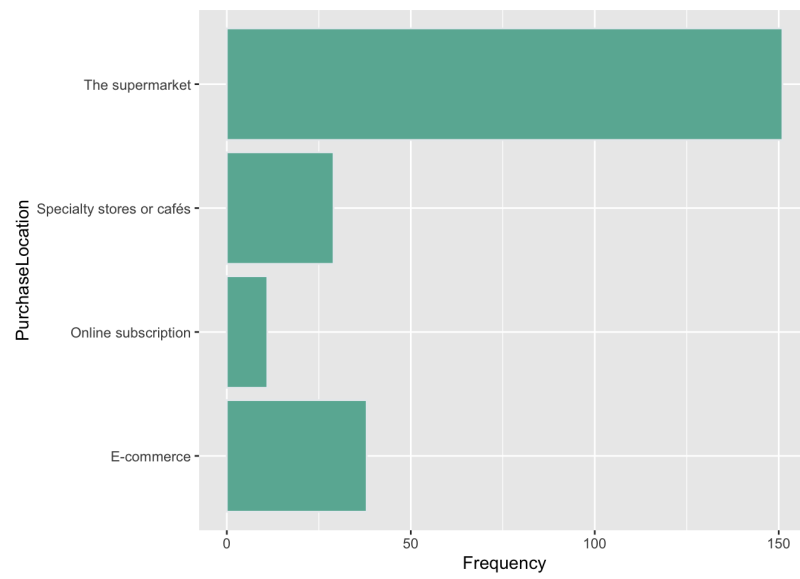
Brand choose

Brand choice	Absolute	Relative
Never	75	32.75%
Sometimes	130	56.77%
Very often	21	9.17%
Every time	3	1.31%



Purchase Method

Purchase Method	Absolute	Relative
E-commerce	38	16.59%
Online subscription	11	4.80%
Specialty stores or cafés	29	12.66%
The supermarket	151	65.94%

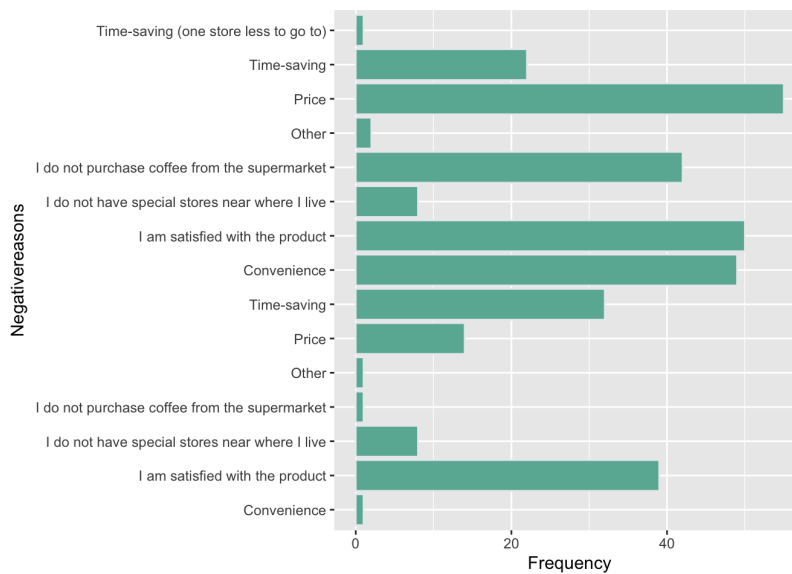


Multiple option answers:

Reasons buying from the supermarket

N
1 325

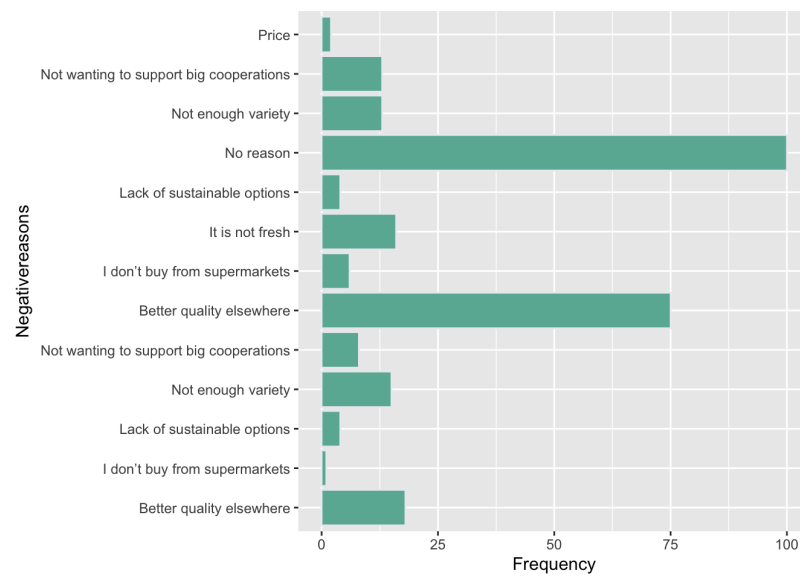
Reason	Frequency
Convenience	1
I am satisfied with the product	39
I do not have special stores near where I live	8
I do not purchase coffee from the supermarket	1
Other	1
Price	14
Time-saving	32
Convenience	49
I am satisfied with the product	50
I do not have special stores near where I live	8
I do not purchase coffee from the supermarket	42
Other	2
Price	55
Time-saving	22
Time-saving (one store less to go to)	1



Reasons for not buying from the supermarket

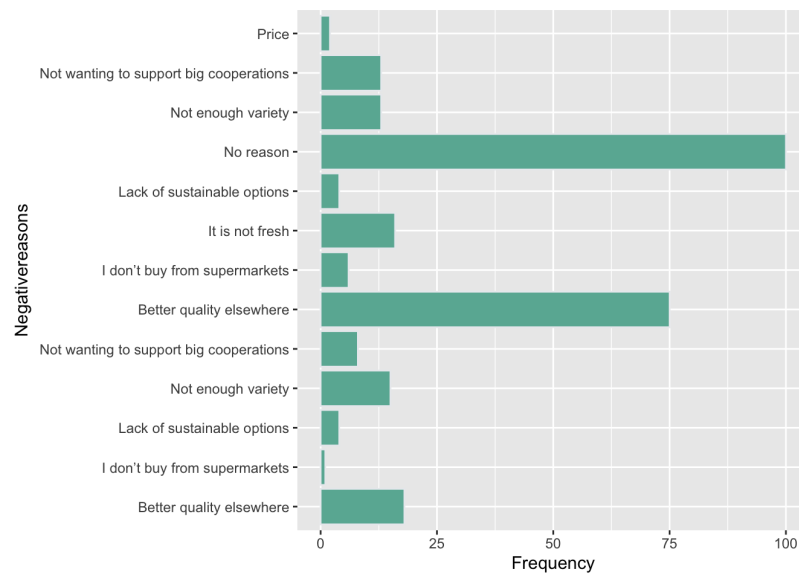
Reason	Frequency
Better quality elsewhere	18
I don't buy from supermarkets	1

Lack of sustainable options	4
Not enough variety	15
Not wanting to support big cooperations	8
Better quality elsewhere	75
I don't buy from supermarkets	6
It is not fresh	16
Lack of sustainable options	4
No reason	100
Not enough variety	13
Not wanting to support big cooperations	13
Price	2



Criteria for choosing the type of coffee

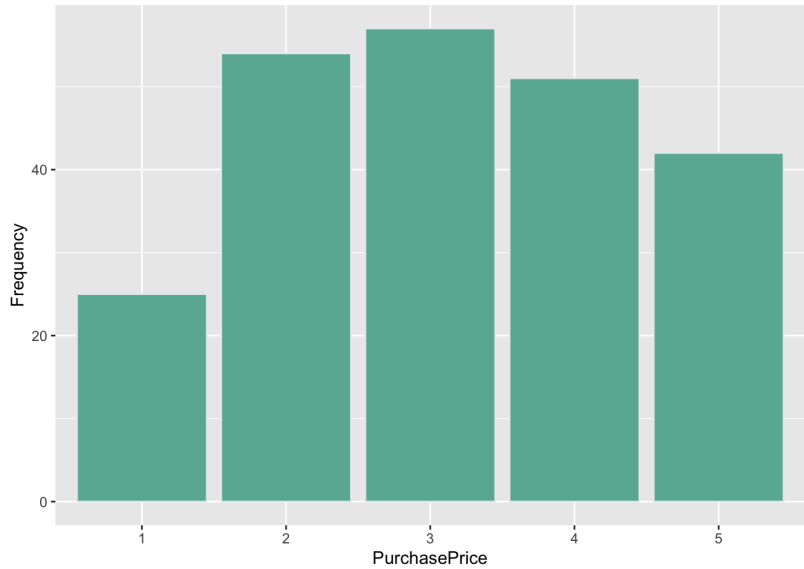
Reason	Frequency
Better quality elsewhere	18
I don't buy from supermarkets	1
Lack of sustainable options	4
Not enough variety	15
Not wanting to support big cooperations	8
Better quality elsewhere	75
I don't buy from supermarkets	6
It is not fresh	16
Lack of sustainable options	4
No reason	100
Not enough variety	13
Not wanting to support big cooperations	13
Price	2



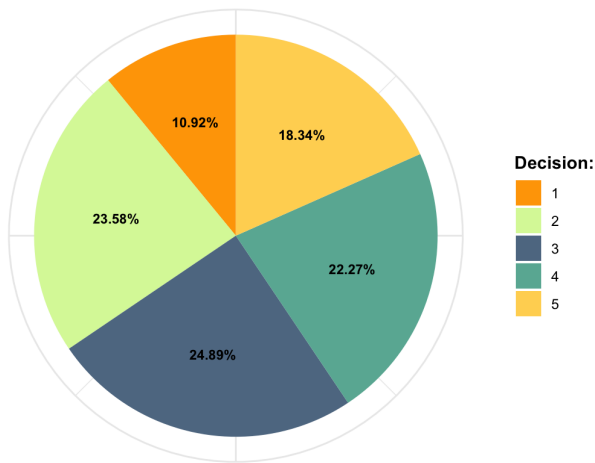
Purchase decisions 1-5

Price

Purchase decision - price	Absolute	Relative
1	25	10.92%
2	54	23.58%
3	57	24.89%
4	51	22.27%
5	42	18.34%



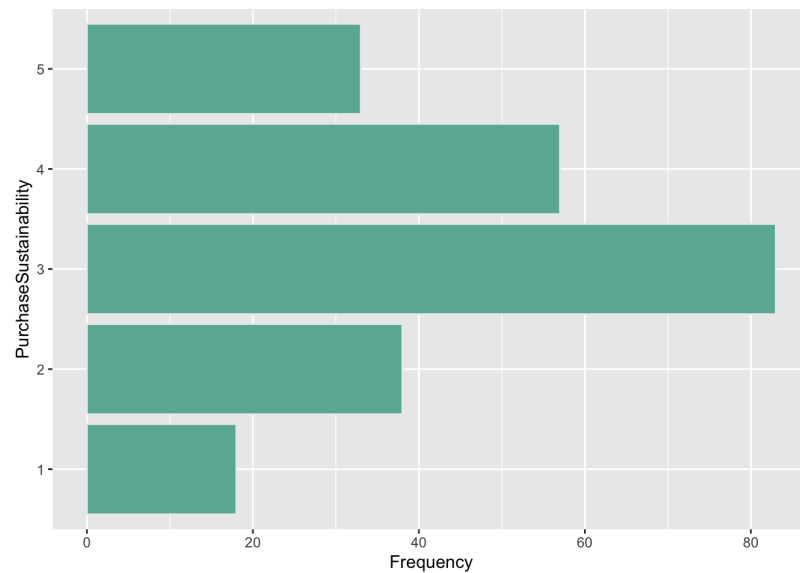
Price



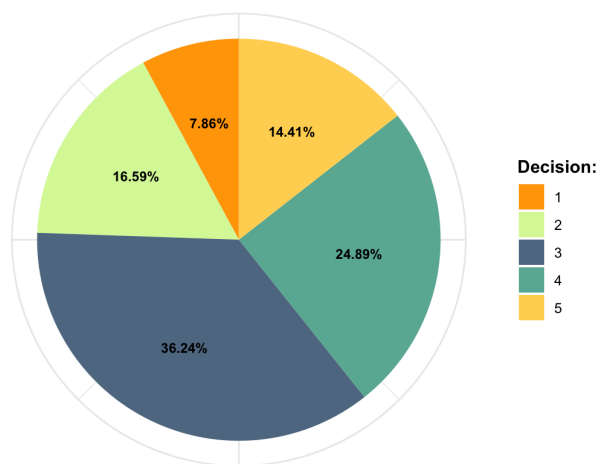
Sustainability



Purchase decision - sustainability	Absolute	Relative
1	18	7.86%
2	38	16.59%
3	83	36.24%
4	57	24.89%
5	33	14.41%



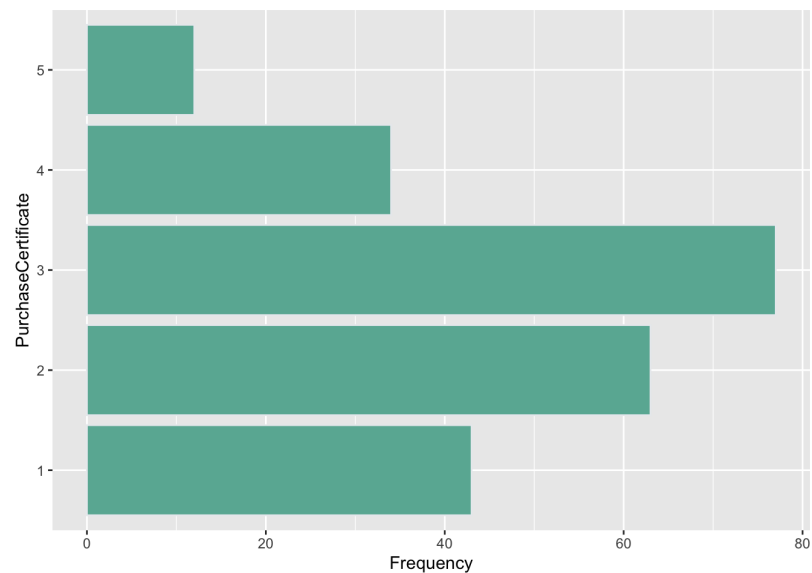
Sustainability



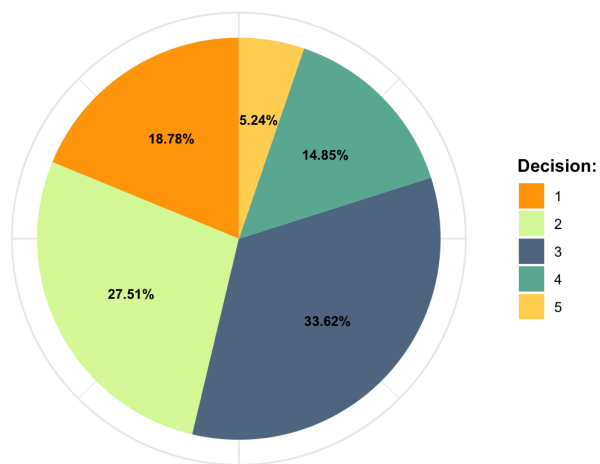
Certificates

Purchase decision - certificate	Absolute	Relative
1	43	18.78%
2	63	27.51%
3	77	33.62%

4	34	14.85%
5	12	5.24%

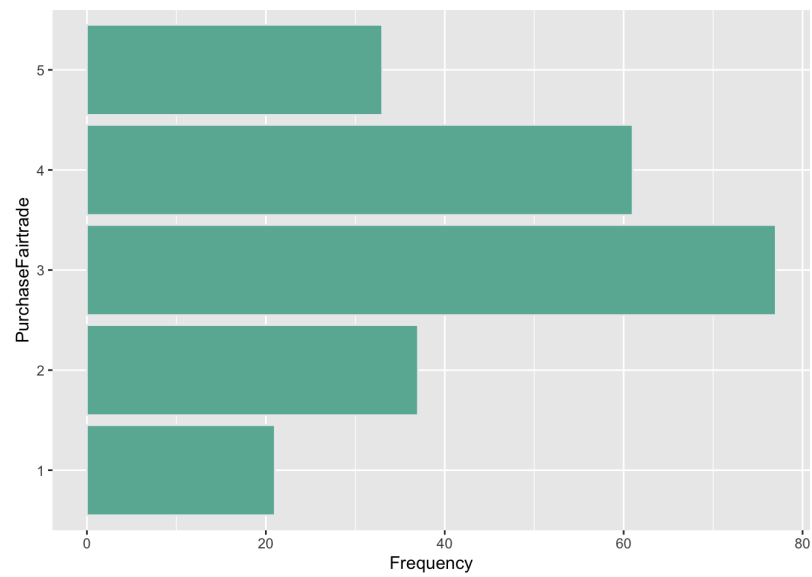


Certificate

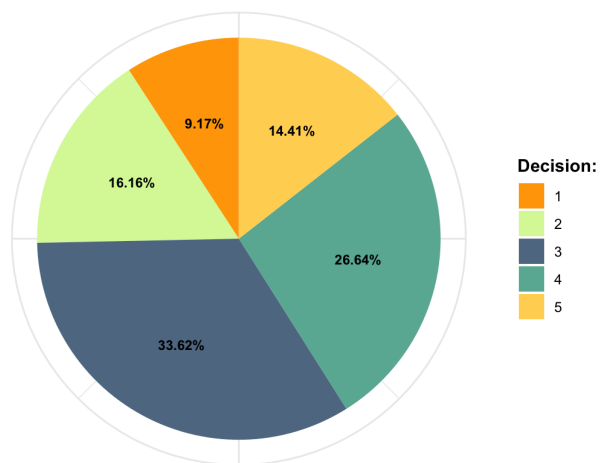


Fairtrade

Purchase decision - fairtrade	Absolute	Relative
1	21	9.17%
2	37	16.16%
3	77	33.62%
4	61	26.64%
5	33	14.41%

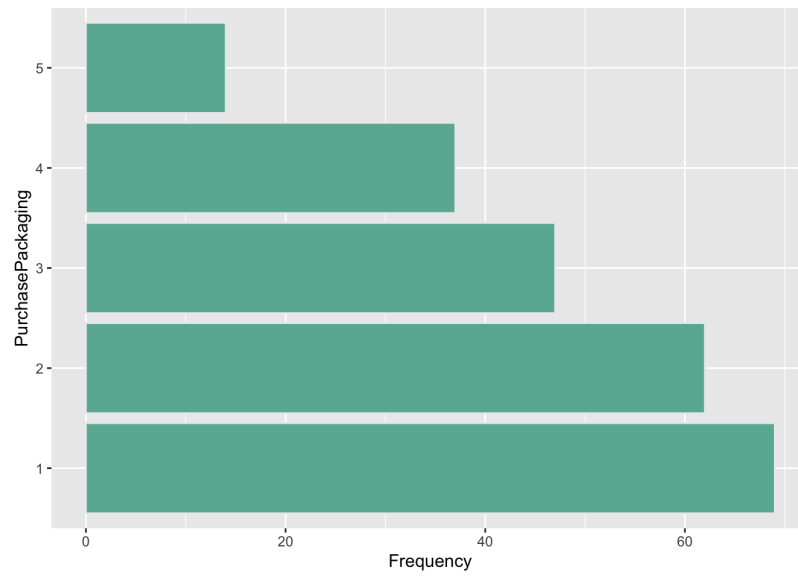


Fair trade

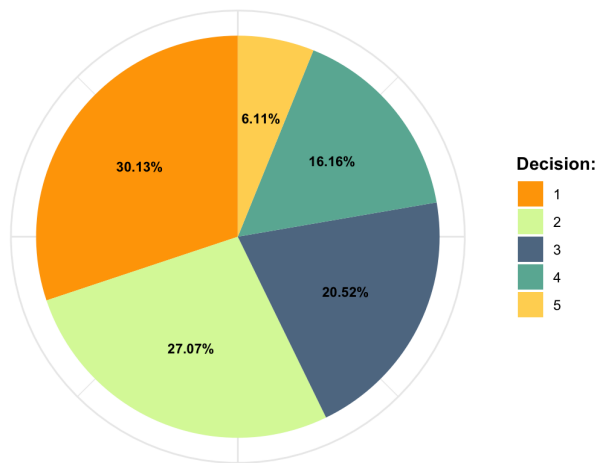


Packaging

Purchase decision - packaging	Absolute	Relative
1	69	30.13%
2	62	27.07%
3	47	20.52%
4	37	16.16%
5	14	6.11%

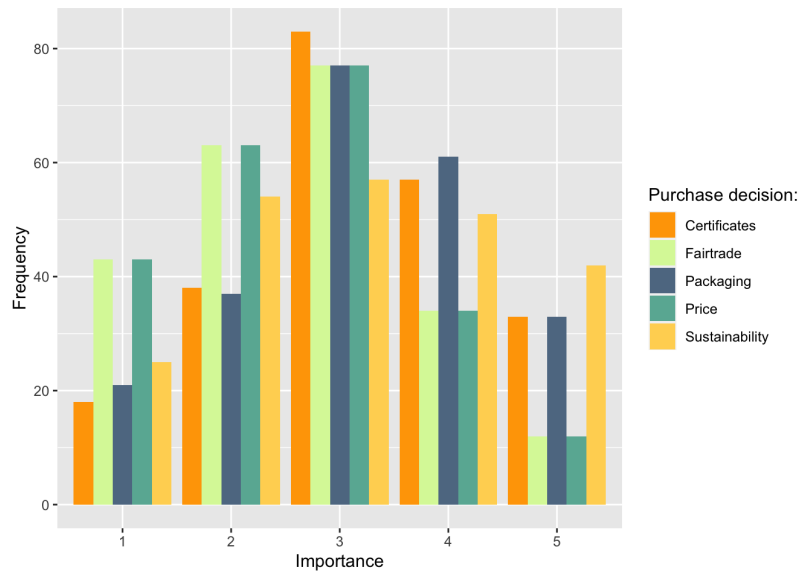


Packaging

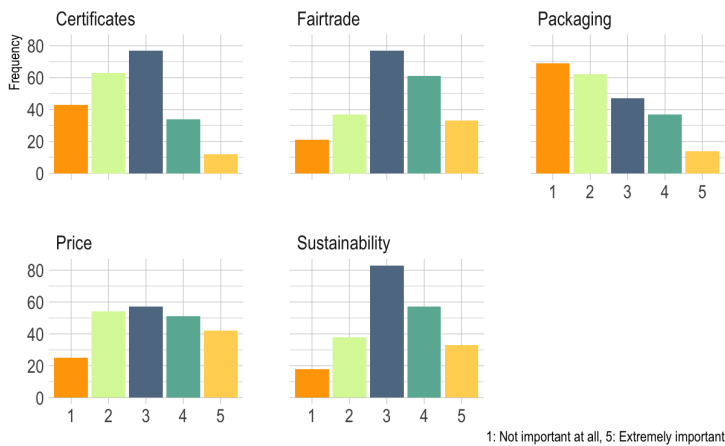


Combined data

Importance	Price	Sustainability	Certificates	Fairtrade	Packaging
1	43	25	18	43	21
2	63	54	38	63	37
3	77	57	83	77	77
4	34	51	57	34	61
5	12	42	33	12	33



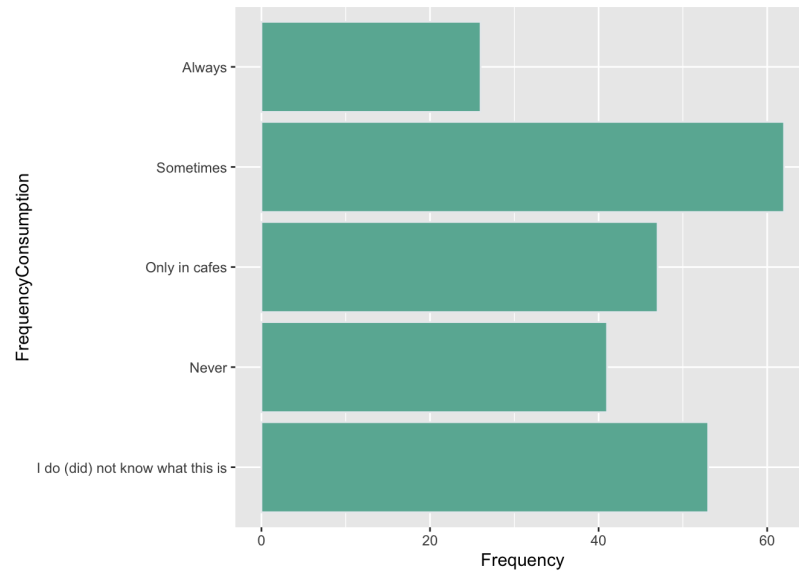
Importance of several factors on the purchasing decision.



Frequency specialty coffee consumption

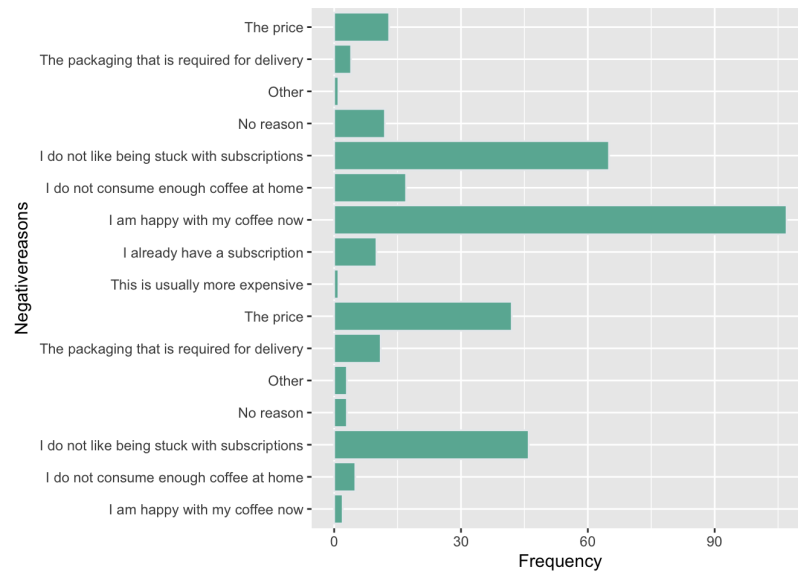
Frequency coffee consumption	Absolute	Relative
------------------------------	----------	----------

Never	41	17.90%
Only in cafes	47	20.52%
Sometimes	62	27.07%
Always	26	11.35%



Reasons for not being likely to set up a subscription

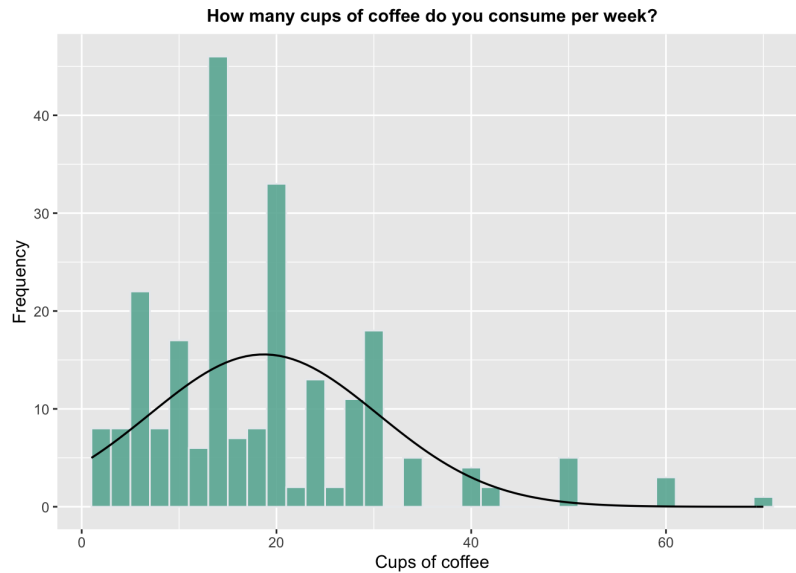
Reason	Frequency
I am happy with my coffee now	2
I do not consume enough coffee at home	5
I do not like being stuck with subscriptions	46
No reason	3
Other	3
The packaging that is required for delivery	11
The price	42
This is usually more expensive	1
I already have a subscription	10
I am happy with my coffee now	107
I do not consume enough coffee at home	17
I do not like being stuck with subscriptions	65
No reason	12
Other	1
The packaging that is required for delivery	4
The price	13



Univariate descriptions - Numerical variables

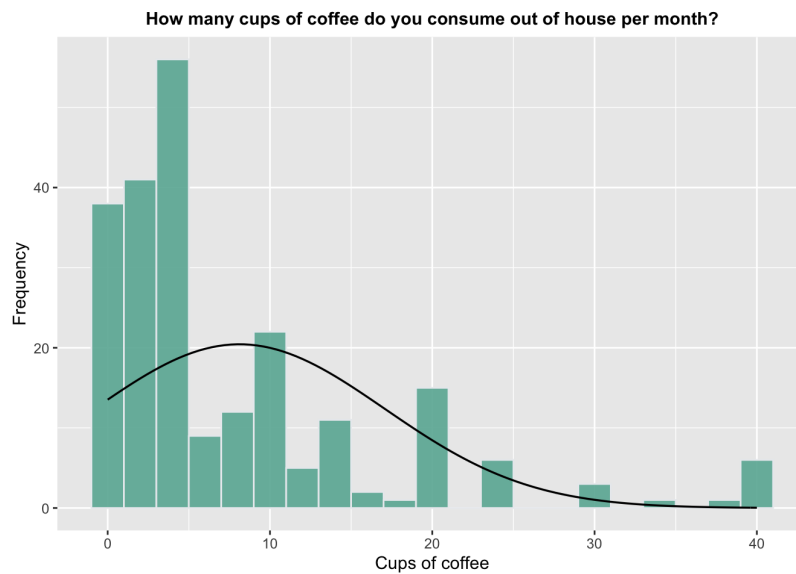
Amount coffe consumed weekly

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.0	10.0	15.0	18.7	25.0	70.0



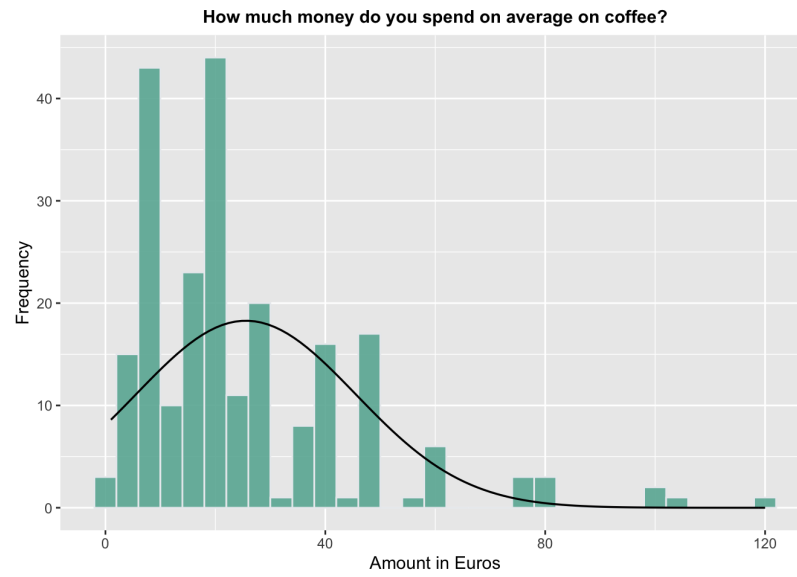
Amount per month out of house

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.000	2.000	5.000	8.122	10.000	40.000



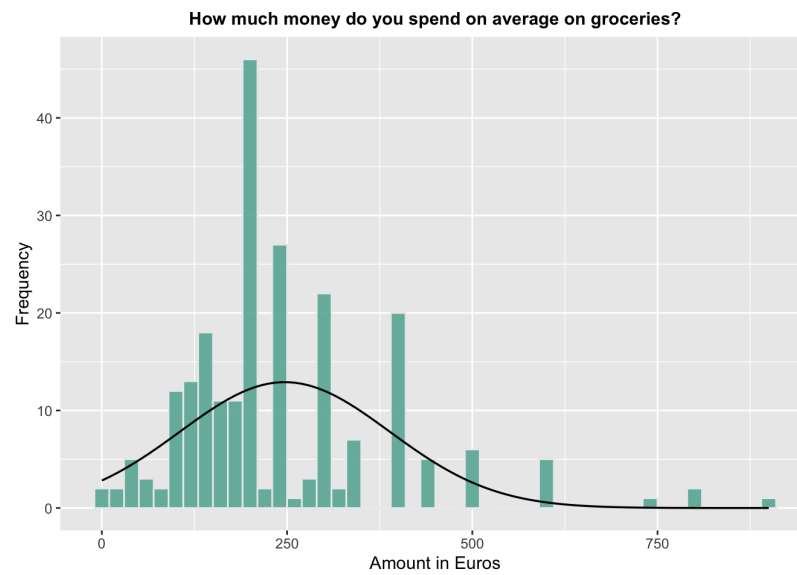
Money coffee

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.00	10.00	20.00	25.55	35.00	120.00



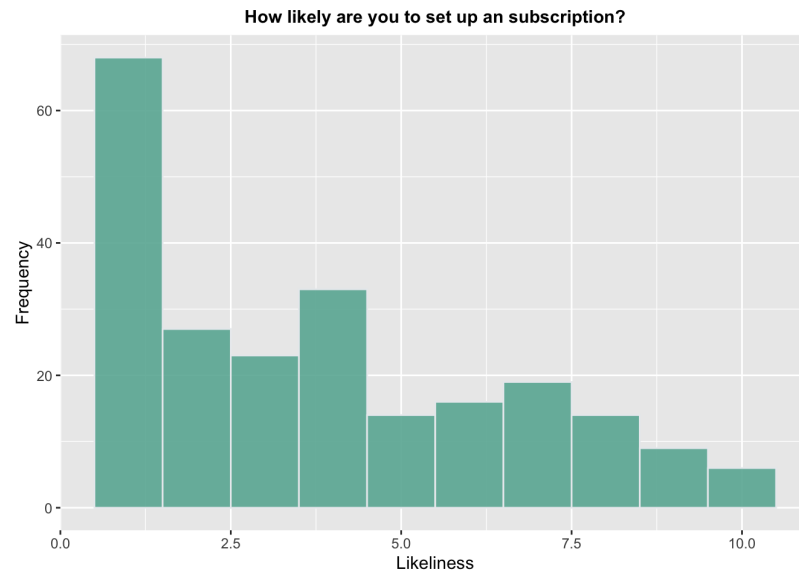
Money groceries

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.0	160.0	200.0	246.9	300.0	900.0



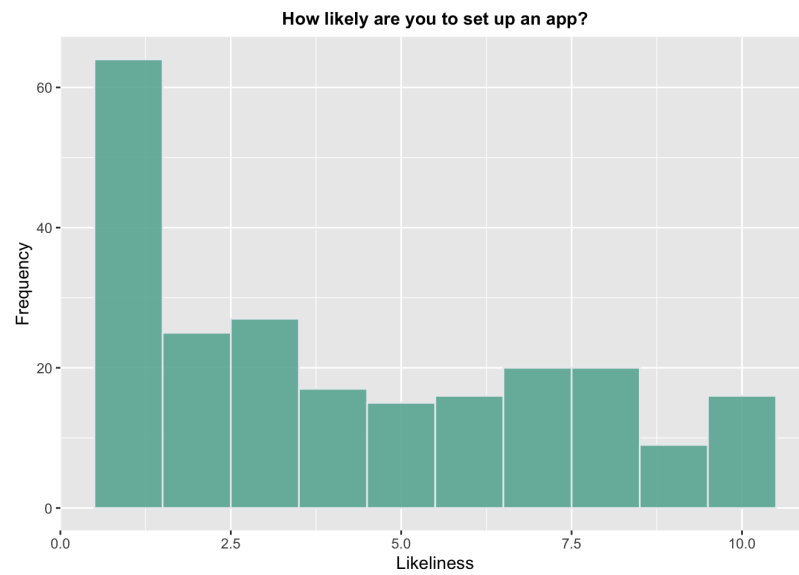
Subscription likely

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.000	1.000	3.000	3.821	6.000	10.000

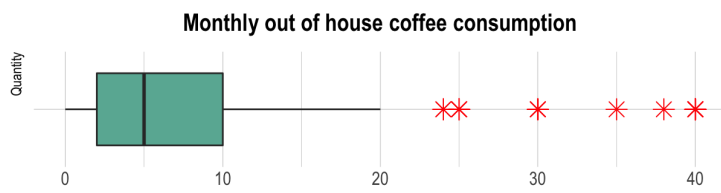
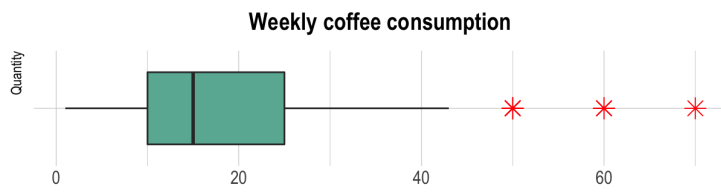


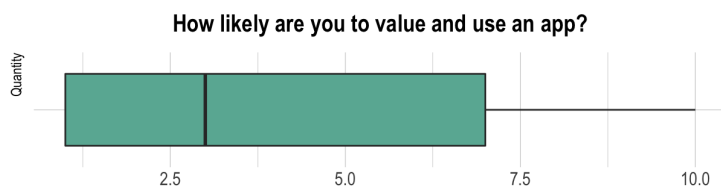
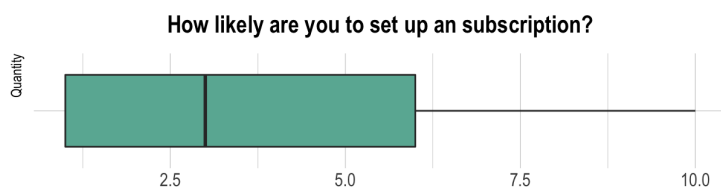
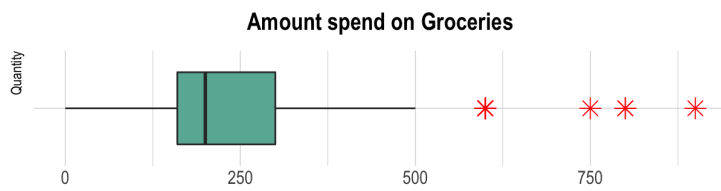
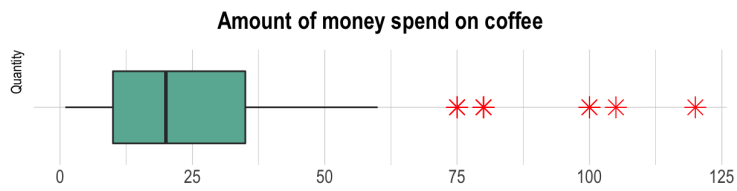
App likely

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.000	1.000	3.000	4.258	7.000	10.000



Boxplots





Parametric testing

H₀ <- There is no association between the two variables.

H_a <- There is a association.

Age - Amount coffee drank

Pearson's Chi-squared test

data: AmountWeek and AgeCategory

X-squared = 236.36, df = 136, p-value = 0.0000002125

Pearson's Chi-squared test with simulated p-value (based on 500 replicates)

data: AmountWeek and AgeCategory

X-squared = 236.36, df = NA, p-value = 0.005988

Education - Amount coffee drank

Pearson's Chi-squared test

data: AmountWeek and Education

X-squared = 225.64, df = 170, p-value = 0.002762

Pearson's Chi-squared test with simulated p-value (based on 500 replicates)

data: AmountWeek and Education

X-squared = 225.64, df = NA, p-value = 0.07385

Gender - Amount coffee drank

Pearson's Chi-squared test

data: AmountWeek and Gender

X-squared = 68.66, df = 68, p-value = 0.4548

Pearson's Chi-squared test with simulated p-value (based on 500 replicates)

data: AmountWeek and Gender

X-squared = 68.66, df = NA, p-value = 0.3473

Home - Amount coffee drank

Pearson's Chi-squared test

data: AmountWeek and Home

X-squared = 65.386, df = 68, p-value = 0.5674

Pearson's Chi-squared test with simulated p-value (based on 500 replicates)

```
data: AmountWeek and Home
X-squared = 65.386, df = NA, p-value = 0.6128
```

App - Age

Pearson's Chi-squared test

```
data: App_Likely and AgeCategory
X-squared = 56.56, df = 36, p-value = 0.01585
```

Pearson's Chi-squared test with simulated p-value (based on 500 replicates)

```
data: App_Likely and AgeCategory
X-squared = 56.56, df = NA, p-value = 0.02395
```

Coffee knowledge - Age

Pearson's Chi-squared test

```
data: KnowledgeCoffee and AgeCategory
X-squared = 152.13, df = 36, p-value = 0.0000000000000003177
```

Pearson's Chi-squared test with simulated p-value (based on 500 replicates)

```
data: KnowledgeCoffee and AgeCategory
X-squared = 152.13, df = NA, p-value = 0.001996
```

Coffee knowledge - Purchase location

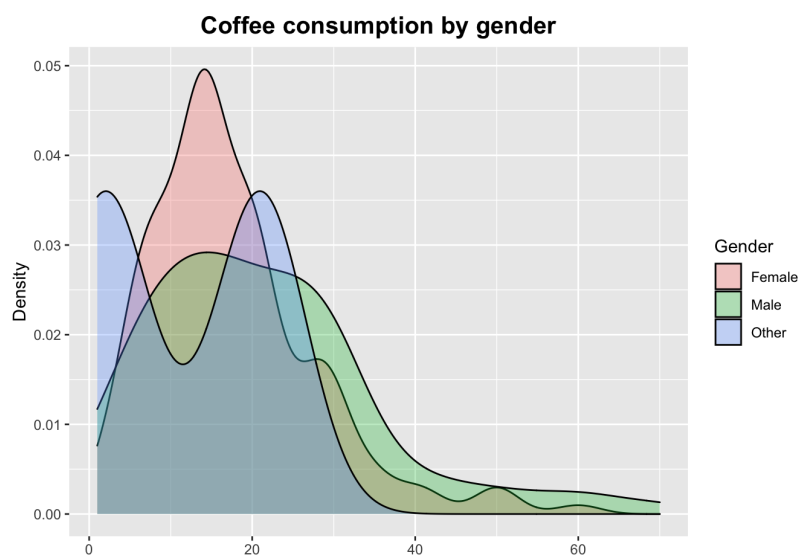
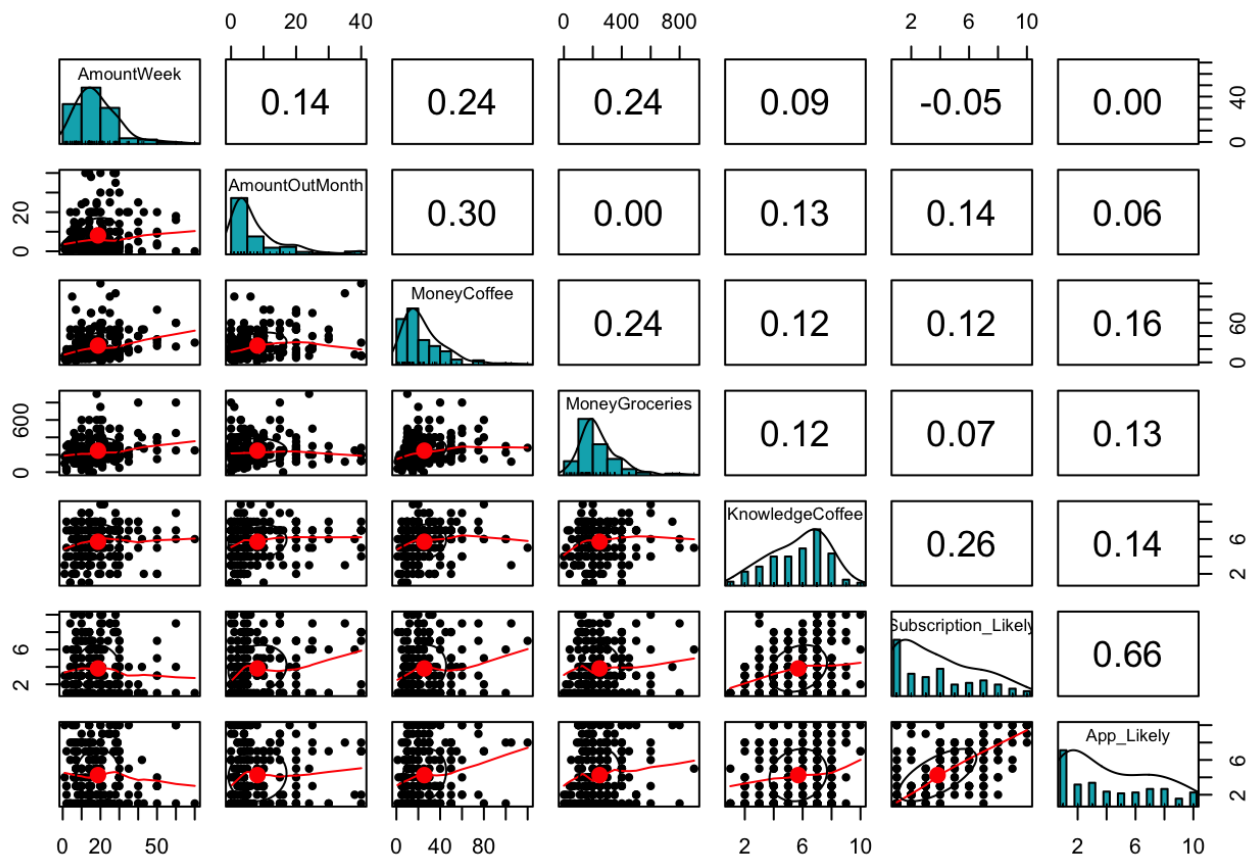
Pearson's Chi-squared test

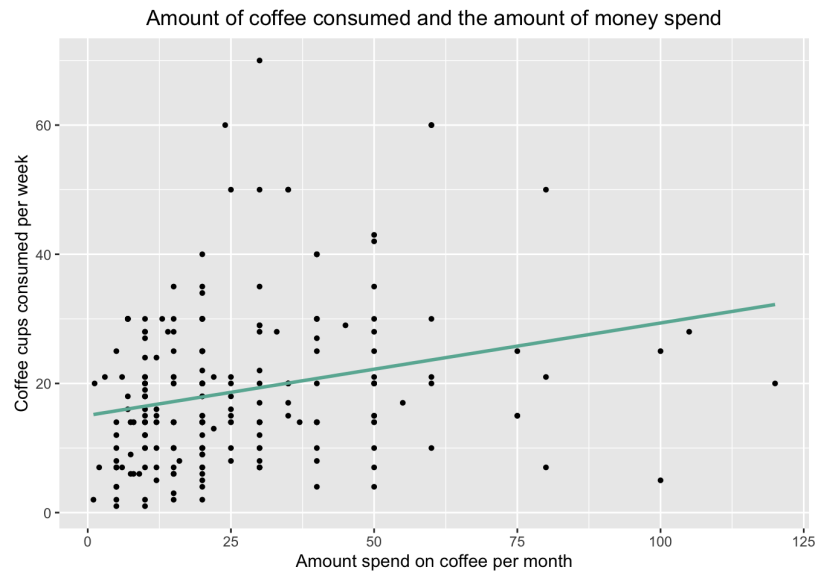
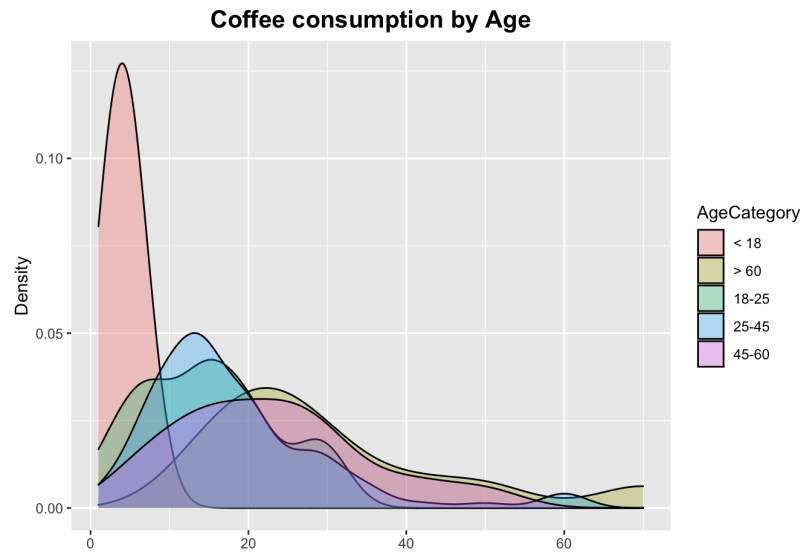
```
data: KnowledgeCoffee and PurchaseLocation
X-squared = 35.471, df = 27, p-value = 0.1273
```

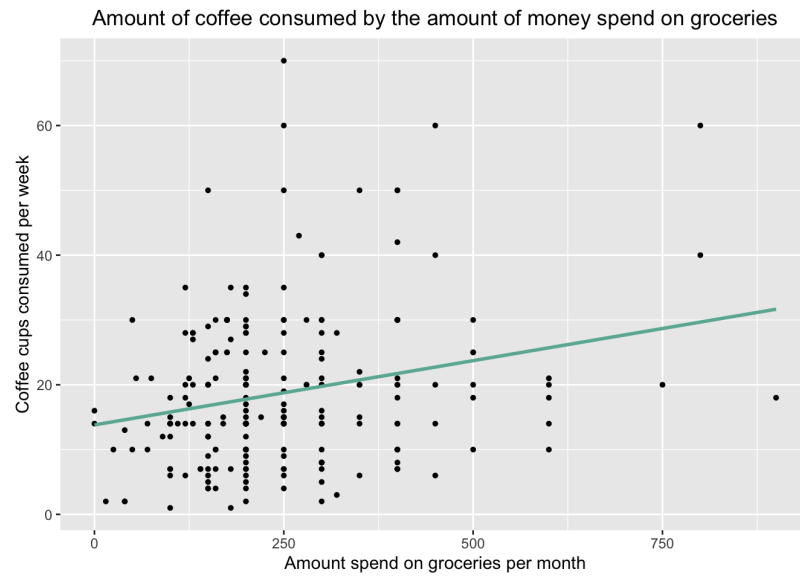
Pearson's Chi-squared test with simulated p-value (based on 500 replicates)

```
data: KnowledgeCoffee and PurchaseLocation
X-squared = 35.471, df = NA, p-value = 0.1118
```

Relationships







Regressions

```
Call:
lm(formula = Subscription_Likely ~ KnowledgeCoffee)

Residuals:
    Min       1Q   Median       3Q      Max
-3.963 -2.236 -0.272  2.074  7.110

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   1.85403    0.52260   3.548  0.000472 ***
KnowledgeCoffee 0.34542    0.08665   3.986  0.0000904 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.606 on 227 degrees of freedom
Multiple R-squared:  0.06543,    Adjusted R-squared:  0.06131
F-statistic: 15.89 on 1 and 227 DF,  p-value: 0.00009044
```

Incl categorical variables as dummies

Cooks distance -> outliers

Data problems