# Statistical Analysis Presentation

A STUDY OF ALCOHOL CONSUMPTION IN ENGLAND

# Introduction

**Background**

Alcohol is a significant public health issue, cited as "the leading cause of ill-health, disability, and death" for people in England aged 15-49 (Department of Health and Social Care, 2021).

Alcohol related hospital admissions in England have been rising year on year. 62% of admissions in 2018/19 were men, and admissions rise with age until 55-64 before reducing (NHS Digital, 2020).

Drinking amongst young women in the UK has been increasing, whilst conditions such as liver disease have also been increasing for the same group (Plant, 2008).

Secondary problems such as public order offences, accidents and injuries caused by binge drinking are more likely to be felt by men (Miller et al., 2005), whilst "higher percentages of heavier-typical quantity drinking are found in the younger age groups" (Chaiyasong et al., 2018).

**Purpose of the presentation**

This presentation uses data from Health Survey for England, 2011* to explore relationships and correlations between attributes such as alcohol consumption, gender, region, height and weight.

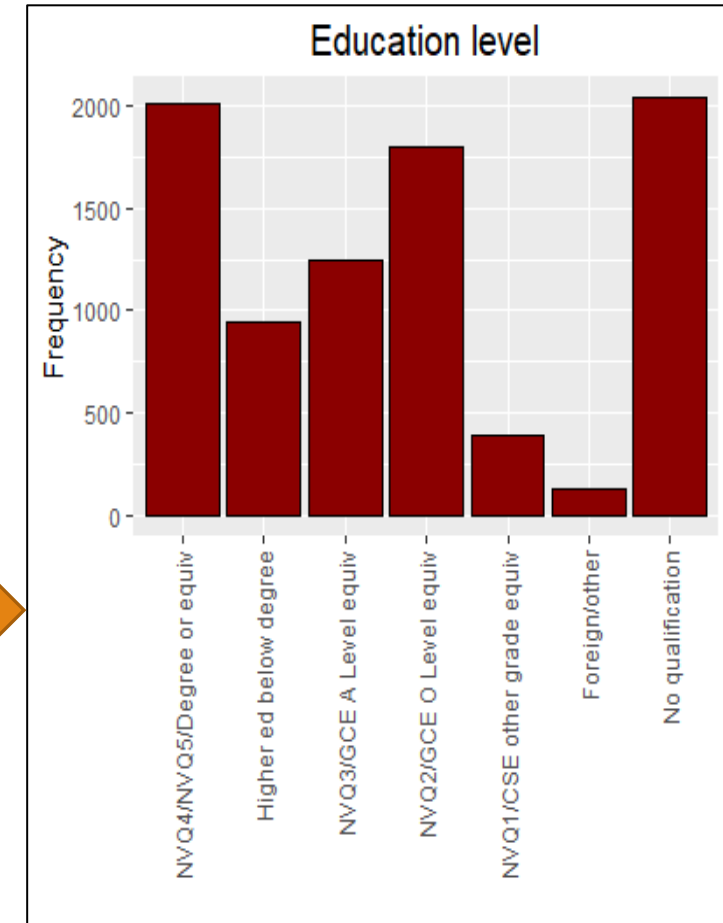T-test, Chi-Square and Kruskal-Wallis tests for inferential statistics

Shapiro Wilk and Kolmogorov-Smirnov tests for distribution normality

R was used to perform the analysis because it is able to quickly perform statistical analysis and produce data visualisations.

*https://www.my-course.co.uk/pluginfile.php/891744/mod_assign/intro/HSE%202011.sav

# Descriptive statistics

| Question | Answer | Observation |
|---|---|---|
| How many people are included in the sample? | 10,617 | Very large sample |
| What is the percentage of people who drink alcohol? | 63.2% | Figure includes N/A at 19.6% which is predominantly children under 18. If N/A is excluded it increases to 78.7% |
| What is the percentage of women in the sample? | 54.3% | Similar to actual 50.9% females in the UK in 2011 (Clark, 2023) |
| What is the highest educational level? | No qualification (2,037) | 'No qualification' had most responses. Highest 'actual' education level was NVQ4/NVQ5/Degree or equivalent: |
| What is the percentage of divorced people? | 5.6% | Figure includes N/A at 18.9% If N/A is excluded it increases to 6.9% |
| What is the percentage of separated people? | 2.1% | Figure includes N/A at 18.9% If N/A is excluded it increases to 2.6% |

# Descriptive statistics

| | Household size | BMI | Age at last birthday |
|---|---|---|---|
| Mean | 2.9 | 25.9 | 41.6 |
| Median | 3 | 25.6 | 42 |
| Mode | 2 | Multi modal *<br>See table | 42 & 64 |
| Minimum | 1 | 8.34011 | 0 |
| Maximum | 10 | 65.27721 | 100 |
| Range | 1 - 10 | 8.34011 - 65.27721 | 0 - 100 |
| Standard deviation | 1.37 | 6.14 | 23.83 |

| 180 Modes for BMI | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 13.76706 | 20.52755 | 22.44604 | 23.768 | 24.63076 | 25.3505 | 26.75475 | 27.91995 | 30.10691 |
| 15.07419 | 20.64521 | 22.47292 | 23.82012 | 24.66207 | 25.41852 | 26.79123 | 27.93738 | 30.15232 |
| 15.10565 | 20.96107 | 22.53362 | 23.82994 | 24.6755 | 25.44888 | 26.89171 | 28.03387 | 30.34498 |
| 15.11716 | 21.10538 | 22.76147 | 23.94858 | 24.72154 | 25.69331 | 26.97778 | 28.06514 | 30.55718 |
| 15.39402 | 21.16579 | 22.7878 | 23.97406 | 24.73146 | 25.77486 | 27.00408 | 28.11204 | 30.69307 |
| 15.58987 | 21.36725 | 22.82347 | 24.04789 | 24.7586 | 25.99036 | 27.02979 | 28.11526 | 30.81796 |
| 16.54064 | 21.54066 | 22.89376 | 24.0497 | 24.79339 | 26.03177 | 27.10957 | 28.12389 | 30.98742 |
| 16.56805 | 21.61622 | 22.90585 | 24.05754 | 24.80159 | 26.09568 | 27.11678 | 28.21979 | 31.04269 |
| 17.2607 | 21.62143 | 23.03436 | 24.11305 | 24.8826 | 26.10882 | 27.13884 | 28.27788 | 31.6068 |
| 17.62971 | 21.68367 | 23.08026 | 24.12092 | 24.96797 | 26.18647 | 27.19212 | 28.42288 | 31.93521 |
| 18.15462 | 21.72132 | 23.08843 | 24.27915 | 24.97874 | 26.23748 | 27.32885 | 28.46231 | 32.00386 |
| 18.38062 | 21.7333 | 23.18113 | 24.28338 | 25.02917 | 26.26407 | 27.37818 | 28.62897 | 32.2211 |
| 18.40747 | 21.74837 | 23.2767 | 24.29688 | 25.06672 | 26.2674 | 27.41375 | 28.65216 | 32.23528 |
| 19.2635 | 21.80385 | 23.28977 | 24.31359 | 25.10867 | 26.2701 | 27.42296 | 28.91054 | 32.28888 |
| 19.56781 | 22.04789 | 23.36088 | 24.35262 | 25.18414 | 26.34628 | 27.44727 | 28.9566 | 32.68683 |
| 19.7607 | 22.12727 | 23.39094 | 24.3598 | 25.22811 | 26.43323 | 27.49056 | 29.22573 | 32.88863 |
| 19.85798 | 22.13931 | 23.51504 | 24.46706 | 25.30832 | 26.50954 | 27.76621 | 29.88813 | 33.08129 |
| 20.01842 | 22.18928 | 23.56401 | 24.49239 | 25.32541 | 26.55283 | 27.82247 | 29.89543 | 33.3317 |
| 20.21217 | 22.21588 | 23.5782 | 24.5957 | 25.34064 | 26.65846 | 27.85467 | 29.93344 | 34.44429 |
| 20.43035 | 22.36915 | 23.73614 | 24.63003 | 25.34732 | 26.7128 | 27.91761 | 30.0971 | 36.84641 |

*BMI has 180 modes with just two instances each so could be considered no mode rather than multi-modal

# Inferential statistics – which gender drinks more alcohol

Chi-square test, because both variables are categorical, to see if there is a difference between whether males and females drink alcohol, α=0.05.

$H_0$: There is no difference between the percentage of men saying that they drink alcohol and women saying that they drink alcohol.

$H_a$: There is a difference between the percentage of men saying that they drink alcohol and women saying that they drink alcohol.
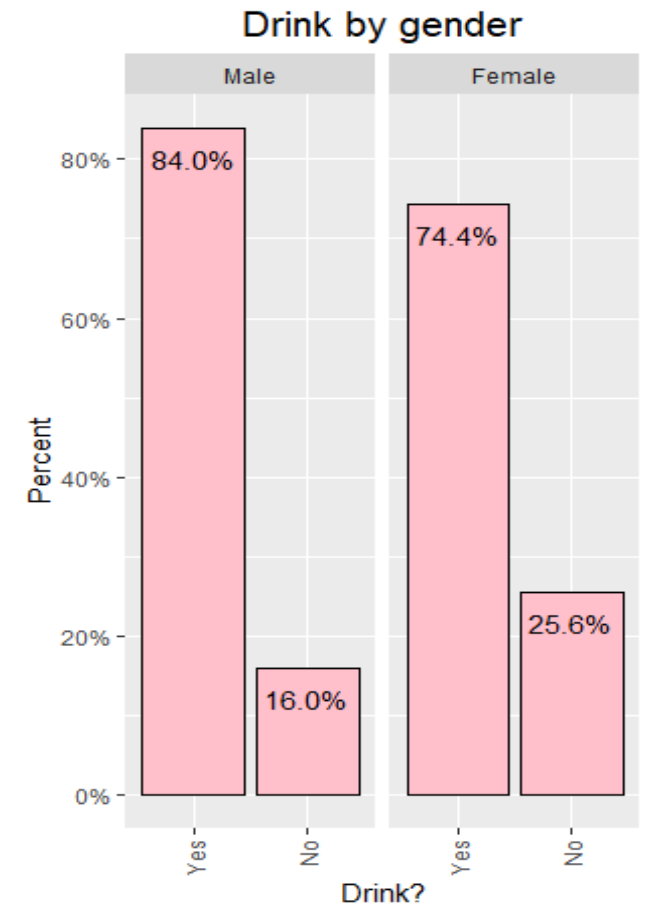
| Gender | Does drink | Does not drink | Test value and p value |
|--------|-----------|----------------|------------------------|
| Male | 3,172 (84.0%) | 605 (16.0%) | Chi-square = 114.15 |
| Female | 3,540 (74.4%) | 1,217 (25.6%) | p-value < 2.2e-16 |

p-value < 0.05
We reject the null hypothesis
There is a difference between the percentage of men saying that they drink alcohol and women saying that they drink alcohol
A higher percentage of men drink than women



Drink by gender

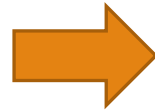# Inferential statistics – which gender drinks more alcohol

Test to see if there is a difference between the weekly quantity of alcohol consumed between males and females at 0.05 significance level ($\alpha$=0.05).

$H_0$ : There is no difference in the average weekly quantity of alcohol consumed between men and women.
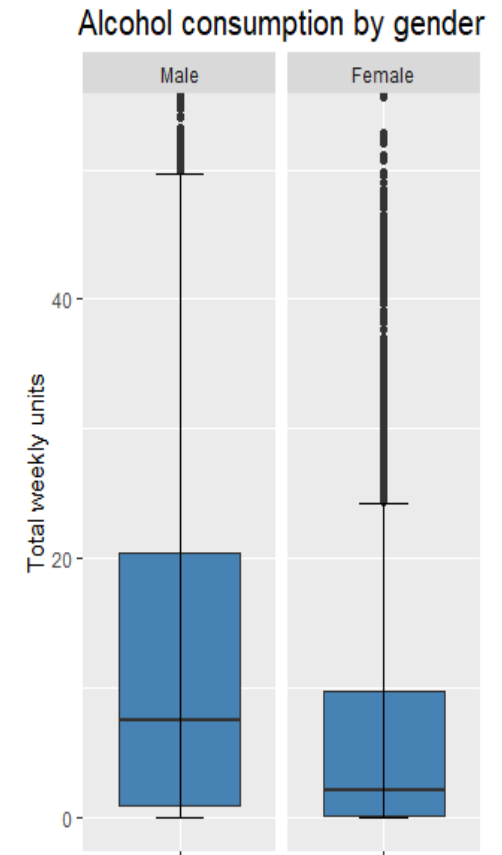
$H_a$ : There is a difference in the average weekly quantity of alcohol consumed between men and women.

A Kolmogorov-Smirnov test was run plus a histogram and Q-Q plot examined of the total weekly units data and the distribution was not normal. However, due to the large sample size a t-test was still appropriate.

```
        welch Two Sample t-test

data:  totalwu by as_factor(Sex)
t = 16.602, df = 5885.1, p-value < 2.2e-16
alternative hypothesis: true difference in means
 between group Male and group Female is not equal
 to 0
95 percent confidence interval:
 6.631599 8.407449
sample estimates:
  mean in group Male mean in group Female
        15.090413              7.570889
```

p-value < 0.05
We reject the null hypothesis
There is a difference in the average weekly quantity of alcohol consumed between men and women
The mean for men is 15.1 units
The mean for women is 7.6 units
The difference between the means at 95% confidence interval is 6.6 – 8.4 units


Alcohol consumption by gender

# Inferential statistics – which region drinks the most alcohol

Chi-square test to see if there is a difference in the percentage of people who consume alcohol by region, $\alpha = 0.05$.

$H_0$: There is no difference in the percentage of people who consume alcohol by region.

$H_a$: There is a difference in the percentage of people who consume alcohol by region.

| Region | Does drink | Does not drink | Test value and p value |
|---|---|---|---|
| North East | 576 (81.0%) | 135 (19.0%) | Chi-square = 98.53 p-value < 2.2e-16 |
| North West | 833 (75.5%) | 270 (24.5%) | |
| Yorkshire and The Humber | 686 (77.3%) | 201 (22.7%) | |
| East Midlands | 624 (82.1%) | 136 (17.9%) | |
| West Midlands | 686 (76.8%) | 207 (23.2%) | |
| East of England | 763 (81.6%) | 172 (18.4%) | |
| London | 674 (68.9%) | 304 (31.1%) | |
| South East | 1,130 (81.6%) | 255 (18.4%) | |
| South West | 740 (83.9%) | 142 (16.1%) | |

p-value < 0.05
We reject the null hypothesis

There is a difference in the percentage of people who consume alcohol by region

South West consumes most alcohol

# Inferential statistics –
# which region drinks the most alcohol

Test to see if there is a difference between the average weekly units of alcohol consumed between regions at 0.05 significance level ($\alpha=0.05$).

$H_0$ : There is no difference in the average weekly quantity of alcohol consumed between regions.

$H_a$ : There is a difference in the average weekly quantity of alcohol consumed between region.

Kruskal-Wallis test to compare multiple regions (groups) at the same time
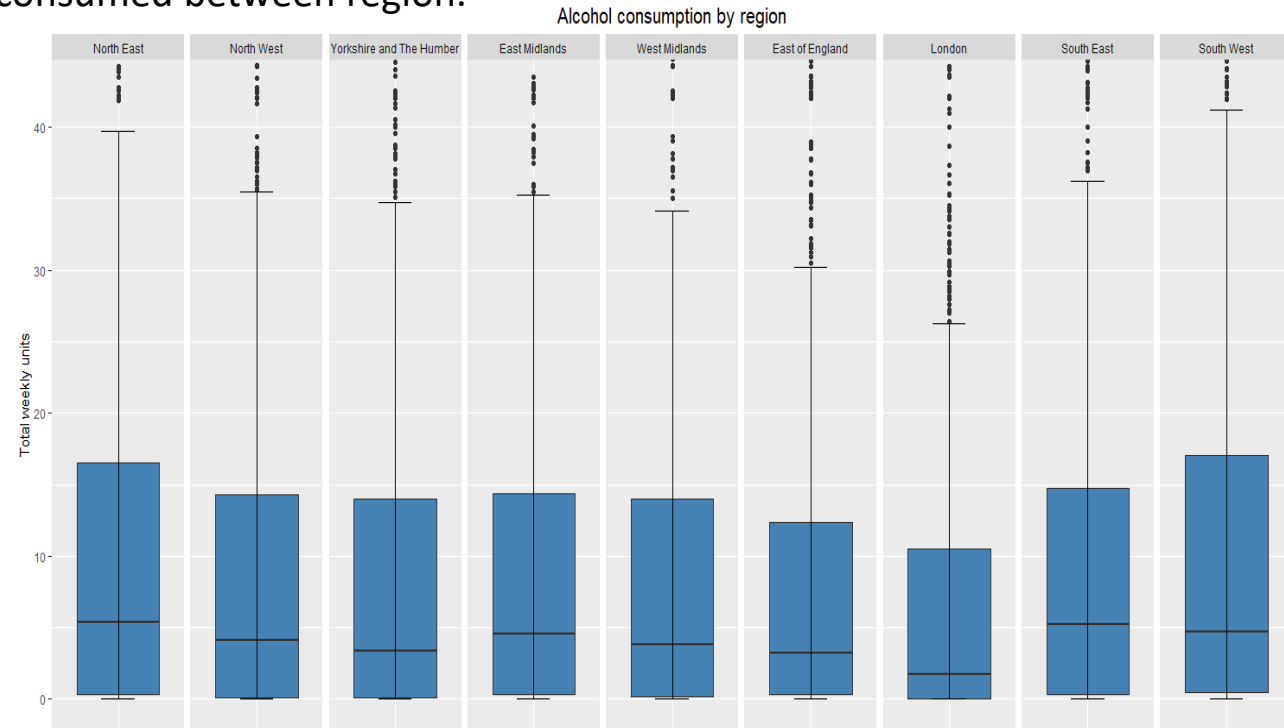
```
        Kruskal-Wallis rank sum test

data:  totalwu by gor1
Kruskal-Wallis chi-squared = 81.722, df = 8, p-value = 2.199e-14
```

p-value < 0.05
We reject the null hypothesis
There is a difference in the average weekly quantity of alcohol consumed between regions



Alcohol consumption by region

# Is there is a statistical difference in height between men and women?

Subset of Adults created where Age > 17.

Shapiro Wilk test to check for normal distributions in height of men and women.

```
        Shapiro-Wilk normality test

data:  Adults$htval[as_factor(Adults$Sex) == "Male"]
W = 0.99919, p-value = 0.1674
```
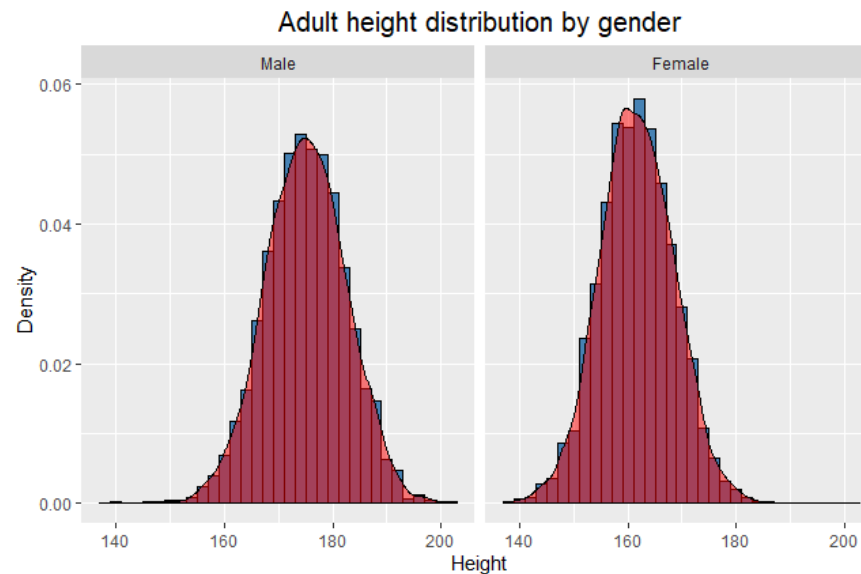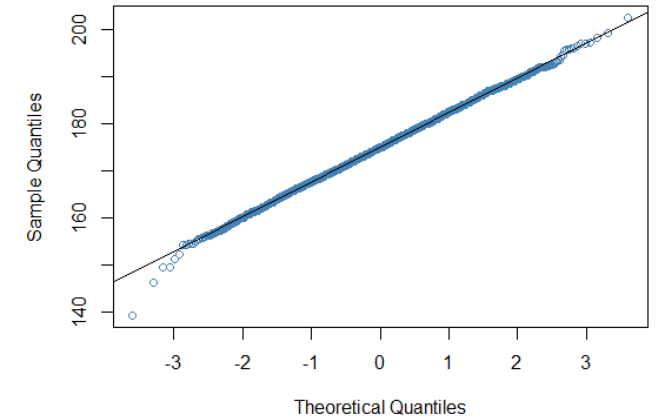
```
        Shapiro-Wilk normality test

data:  Adults$htval[as_factor(Adults$Sex) == "Female"]
W = 0.99961, p-value = 0.6573
```

With p-values > 0.05 both distributions are normal. The density plots and Q-Q plots also confirm normality.
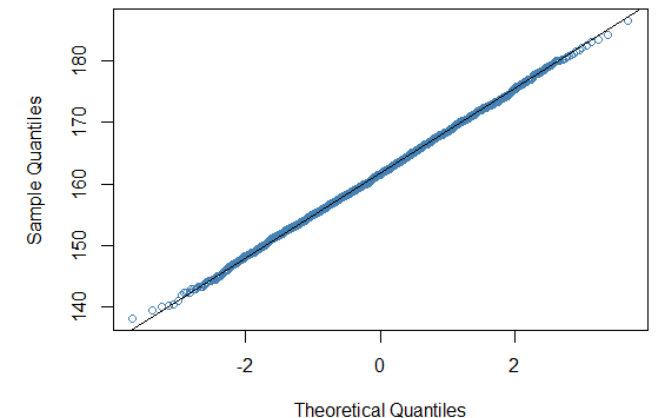
A t-test was therefore be used.



Adult height distribution by gender



Q-Q plot to test for normal distribution of height - Male



Q-Q plot to test for normal distribution of height - Female

# Is there is a statistical difference in height between men and women?

$H_0$ : There is no difference in height between men and women.

$H_a$ : There is a difference in height between men and women.

```
        Welch Two Sample t-test

data:  htval by as_factor(Sex)
t = 77.707, df = 6444.4, p-value < 2.2e-16
alternative hypothesis: true difference in means between group Male and group Female is not equal to 0
95 percent confidence interval:
 12.98930 13.66162
sample estimates:
  mean in group Male mean in group Female
          175.0300              161.7045
```
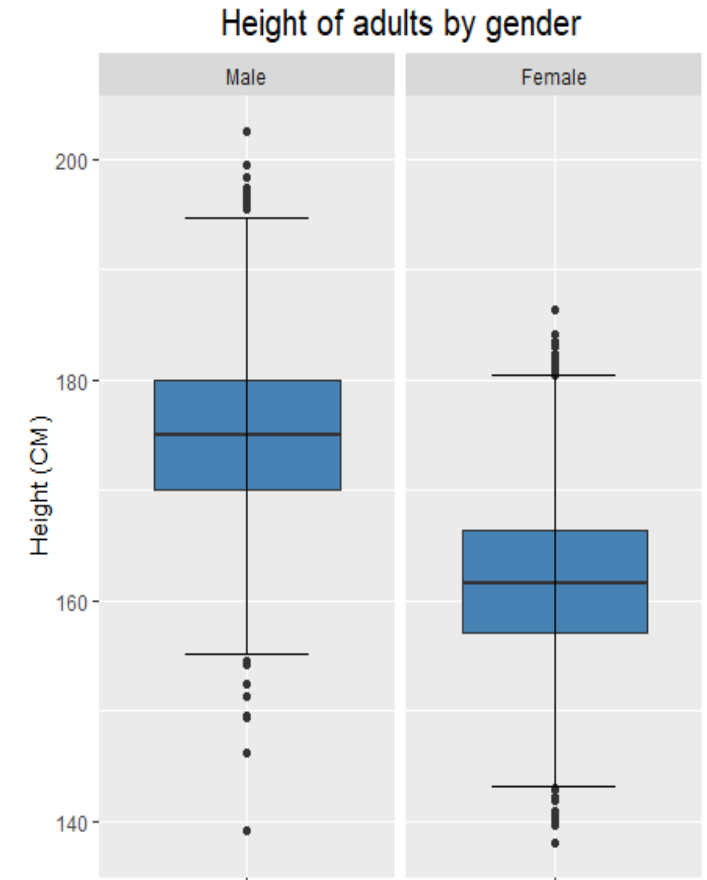
p-value < 0.05

We reject the null hypothesis

There is a difference in height between men and women – men are taller

The mean for men in the sample is 175.0cm

The mean for women in the sample is 161.7cm

The difference between the means at 95% confidence interval is 13.0 – 13.7cm


Height of adults by gender

# Is there is a statistical difference in weight between men and women?

Shapiro Wilk test to check for normal distributions in weight of men and women.

```
        Shapiro-Wilk normality test

data:  Adults$wtval[as_factor(Adults$Sex) == "Male"]
W = 0.95707, p-value < 2.2e-16
```
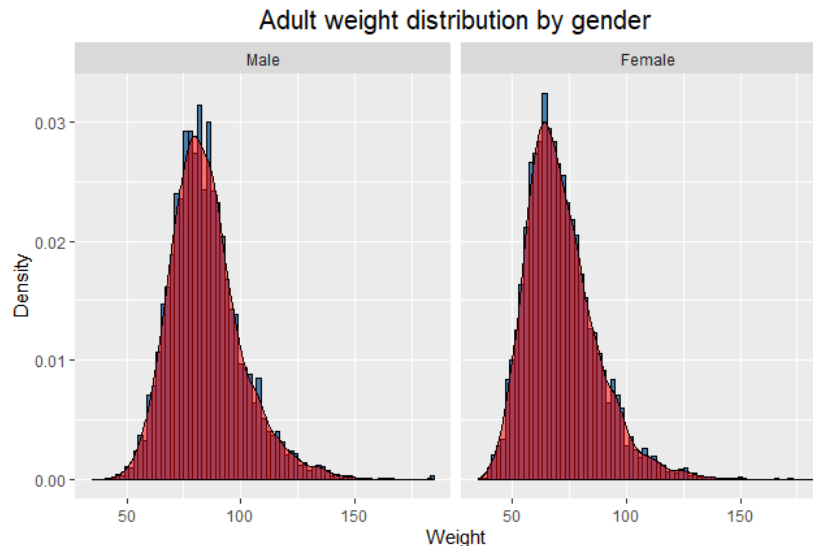
```
        Shapiro-Wilk normality test

data:  Adults$wtval[as_factor(Adults$Sex) == "Female"]
W = 0.94493, p-value < 2.2e-16
```

With p-values < 0.05, neither distribution was normal. The density plots were slightly skewed and the Q-Q plots deviated from normality.

With a large sample size the central limit theorem allows the use of a t-test even when the distribution is not normal.

A t-test was therefore used in favour of the non-parametric Mann-Whitney U test.



Q-Q plot to test for normal distribution of weight - Male



Q-Q plot to test for normal distribution of weight - Female



Adult weight distribution by gender

# Is there is a statistical difference in weight between men and women?

$H_0$ : There is no difference in weight between men and women.

$H_a$ : There is a difference in weight between men and women.

```
        Welch Two Sample t-test

data:  wtval by as_factor(Sex)
t = 34.808, df = 6738.3, p-value < 2.2e-16
alternative hypothesis: true difference in means between group Male and group Female is not equal to 0
95 percent confidence interval:
 12.57371 14.07447
sample estimates:
  mean in group Male mean in group Female
          84.89677              71.57267
```
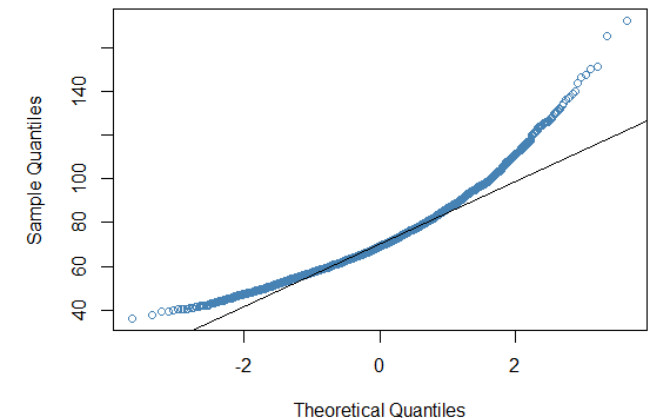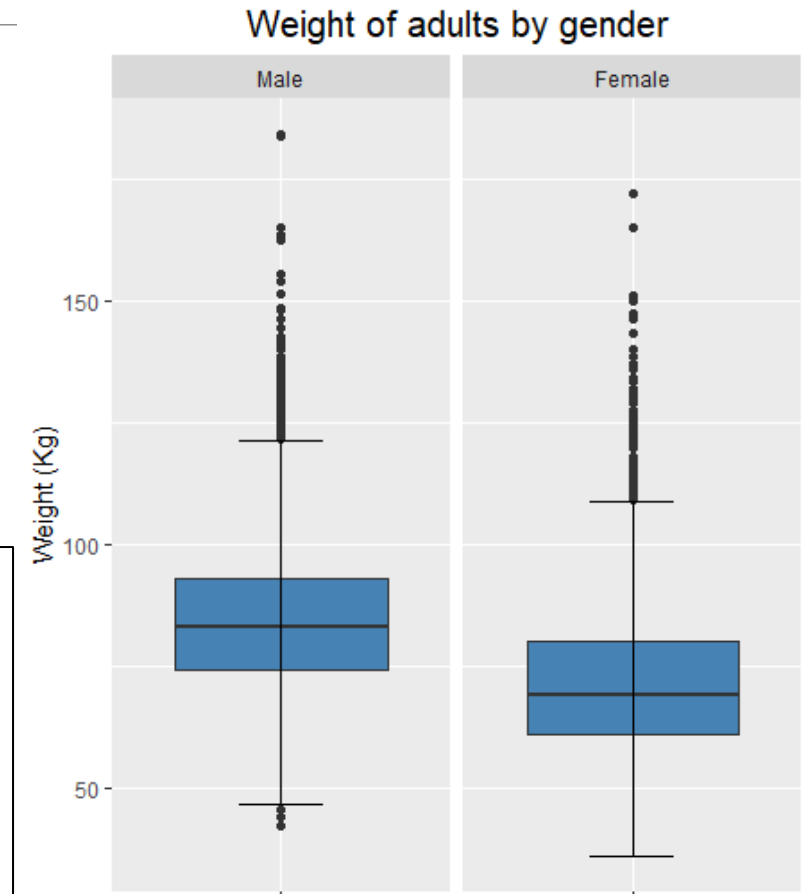
p-value < 0.05

We reject the null hypothesis

There is a difference in height between men and women – men are heavier

The mean for men in the sample is 84.9kg

The mean for women in the sample is 71.6kg

The difference between the means at 95% confidence interval is 12.6 – 14.1kg



Weight of adults by gender

# What is the correlation between whether a person drinks nowadays, total household income, age at last birthday and gender?

Kolmogorov-Smirnov tests to check for normal distributions for the four variables

```
        Asymptotic one-sample Kolmogorov-Smirnov test

data:  HSE_2011$dnnow
D = 0.84134, p-value < 2.2e-16
alternative hypothesis: two-sided
```

```
        Asymptotic one-sample Kolmogorov-Smirnov test

data:  HSE_2011$totinc
D = 0.99448, p-value < 2.2e-16
alternative hypothesis: two-sided
```
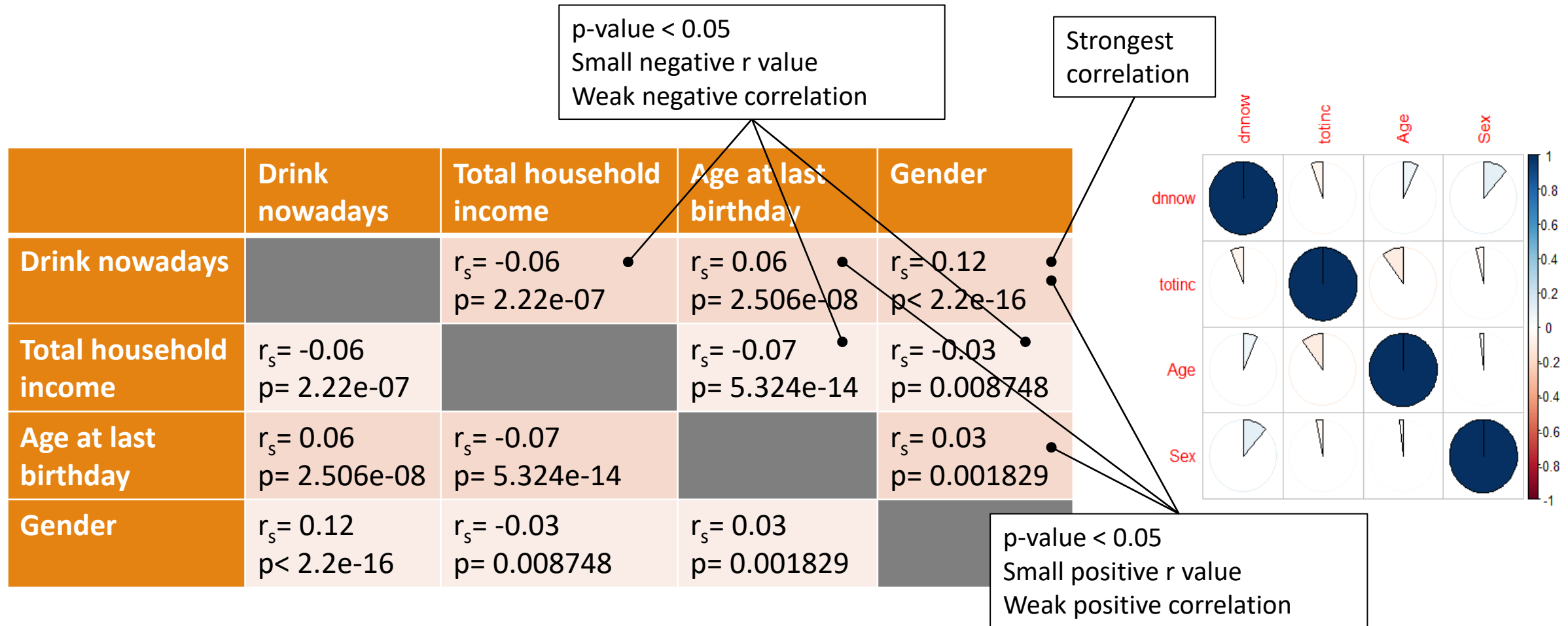
```
        Asymptotic one-sample Kolmogorov-Smirnov test

data:  HSE_2011$Age
D = 0.95683, p-value < 2.2e-16
alternative hypothesis: two-sided
```

```
        Asymptotic one-sample Kolmogorov-Smirnov test

data:  HSE_2011$Sex
D = 0.84134, p-value < 2.2e-16
alternative hypothesis: two-sided
```

Histograms and Q-Q plots confirmed non-normal distributions so Spearman's rank was used

# What is the correlation between whether a person drinks nowadays, total household income, age at last birthday and gender?

p-value < 0.05
Small negative r value
Weak negative correlation

Strongest correlation

| | Drink nowadays | Total household income | Age at last birthday | Gender |
|---|---|---|---|---|
| **Drink nowadays** | | $r_s$= -0.06 <br> p= 2.22e-07 | $r_s$= 0.06 <br> p= 2.506e-08 | $r_s$=0.12 <br> p< 2.2e-16 |
| **Total household income** | $r_s$= -0.06 <br> p= 2.22e-07 | | $r_s$= -0.07 <br> p= 5.324e-14 | $r_s$= -0.03 <br> p= 0.008748 |
| **Age at last birthday** | $r_s$= 0.06 <br> p= 2.506e-08 | $r_s$= -0.07 <br> p= 5.324e-14 | | $r_s$= 0.03 <br> p= 0.001829 |
| **Gender** | $r_s$= 0.12 <br> p< 2.2e-16 | $r_s$= -0.03 <br> p= 0.008748 | $r_s$= 0.03 <br> p= 0.001829 | |

p-value < 0.05
Small positive r value
Weak positive correlation

# Discussion and recommendations

The analysis found that men drink more than women, both when measured as number of people who drink, and average weekly units. This is consistent with the findings in NHS Digital (2020). Also consistent with the same study is that there is regional variation in levels of drinking. Men were also found on average to be taller and heavier than women.

Chaiyasong et al. (2018) found that high frequency drinking is higher in older age groups, and heavier drinking in middle and older age groups. This study found a weak positive correlation between drinking and age, but by measuring only the drink now variable rather than total weekly units it couldn't differentiate heavier drinking.

The study found that a high number of people are higher education educated. Alcohol consumption amongst students is increasing (Davoren et al., 2016), so it is recommended to start a targeted campaign to increase awareness of the risks of alcohol throughout universities in England, as well as a wider public health campaign focusing on the secondary effects of excessive alcohol consumption such as public disorder as well as primary health risks.

# References

Chaiyasong, S. et al. (2018) Drinking patterns vary by gender, age and country-level income: Cross-country analysis of the International Alcohol Control Study. *Drug and Alcohol Review*, 37(S2): S53-S62. https://doi.org/10.1111/dar.12820

Clark, D. (2023) Population of the United Kingdom from 1953 to 2020, by gender. Available from: https://www.statista.com/statistics/281240/population-of-the-united-kingdom-uk-by-gender/ [Accessed 23 March 2023].

Davoren, M.P., Demant, J., Shiely, F. & Perry, I.J. (2016) Alcohol consumption among university students in Ireland and the United Kingdom from 2002 to 2014: a systematic review. *BMC public health*, 16(1): 1-13. https://doi.org/10.1186/s12889-016-2843-1

Department of Health and Social Care (2021) Delivering better oral health: an evidence-based toolkit for prevention. Chapter 12: Alcohol. Available from: https://www.gov.uk/government/publications/delivering-better-oral-health-an-evidence-based-toolkit-for-prevention/chapter-12-alcohol [Accessed 23 March 2023].

Miller, P., Plant, M. & Plant, M. (2005) Spreading out or concentrating weekly consumption: alcohol problems and other consequences within a UK population sample. *Alcohol and Alcoholism*, 40(5): 461-468. https://doi.org/10.1093/alcalc/agh169

NHS Digital (2020) Statistics on Alcohol, England 2020. Available from: https://digital.nhs.uk/data-and-information/publications/statistical/statistics-on-alcohol/2020/part-1 [Accessed 23 March 2023].

Plant, M.L. (2008) The role of alcohol in women's lives: A review of issues and responses. *Journal of Substance Use*, 13(3): 155-191. https://doi.org/10.1080/14659890802040880

# Appendix A: R code used to produce statistics and visualisations

```r
1   library(haven)
2   library(dplyr)
3   library(ggplot2)
4   library(corrplot)
5
6   HSE_2011 <- read_sav("HSE 2011.sav")
7   View(HSE_2011)
8
9   # How many people are included in the sample?
10  # pserial is the serial number of the individual, so count unique individuals
11  n_distinct(HSE_2011$pserial)
12
13  # What is the percentage of people who drink alcohol?
14  # Create a frequency table of whether people drink alcohol
15  table(
16    as_factor(HSE_2011$dnnow)
17  )
18
19  # Make a proportional table that can be turned into percentages
20  prop.table(
21    table(
22      as_factor(HSE_2011$dnnow)
23    )
24  )
25
26  # Repeat including the NA values
27  table(
28    as_factor(HSE_2011$dnnow)
29    ,useNA="ifany"
30  )
31
32  prop.table(
33    table(
34      as_factor(HSE_2011$dnnow)
35      ,useNA="ifany"
36    )
37  )
38
```

```r
39  # Plot percentage of people who drink alcohol
40  ggplot(HSE_2011, aes(as_factor(dnnow))) +
41    geom_bar(aes(y=after_stat(count/sum(count))), colour='black' ,fill='pink') +
42    labs(title='Percentage of people who drink alcohol', x='Drink alcohol?', y='Percent') +
43    geom_text(aes( label = scales::percent(after_stat(count/sum(count))),
44              y=after_stat(count/sum(count))), stat= 'count', vjust = 2) +
45    scale_y_continuous(labels=scales::percent) +
46    theme_grey() +
47    theme(plot.title = element_text(hjust = 0.5, size = 15)) +
48    theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5))
49
50
51  # Plot Male/female split
52  filter(HSE_2011, (dnnow > 0)) %>%
53    ggplot(aes(as_factor(dnnow))) +
54    geom_bar(aes(y=after_stat(prop), group=1), colour='black', fill='pink') +
55    labs(title='Percentage people who drink alcohol, by gender', x='Drink alcohol?', y='Percent') +
56    geom_text(aes(label = scales::percent(after_stat(prop)),
57              y=after_stat(prop),group=1), stat= 'count', vjust = 2) +
58    scale_y_continuous(labels=scales::percent) +
59    facet_grid(~as_factor(Sex)) +
60    theme_grey() +
61    theme(plot.title = element_text(hjust = 0.5, size = 15)) +
62    theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5))
63
64
65  # What is the percentage of women in the sample?
66  # Create a frequency table of whether people drink alcohol
67  table(
68    as_factor(HSE_2011$Sex)
69  )
70
71
72  # Make a proportional table that can be turned into percentages
73  prop.table(
74    table(
75      as_factor(HSE_2011$Sex)
76    )
77  )*100
78
```

```r
79   # Plot all people
80   ggplot(HSE_2011, aes(as_factor(Sex))) +
81     geom_bar(aes(y=after_stat(prop), group=1), colour='black', fill='pink') +
82     labs(title='Gender breakdown', x='Gender', y='Percent') +
83     geom_text(aes(label = scales::percent(after_stat(prop)),
84                   y=after_stat(prop),group=1), stat= 'count', vjust = 2) +
85     scale_y_continuous(labels=scales::percent) +
86     theme_grey() +
87     theme(plot.title = element_text(hjust = 0.5, size = 15)) +
88     theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5))
89
90   # Differentiate adults
91   age_labels <- c(
92     'FALSE' = 'Under 18',
93     'TRUE' = 'Adult'
94   )
95
96   # What is the highest educational level?
97   # Discover the attributes of topqual3
98   attributes(HSE_2011$topqual3)
99
100  # Now find the range
101  range(HSE_2011$topqual3, na.rm=TRUE)
102
103  table(as_factor(HSE_2011$topqual3))
104
105  filter(HSE_2011, topqual3 > 0) %>% # to remove NA values
106    ggplot(aes(as_factor(topqual3))) +
107    geom_bar(colour='black' ,fill='dark red') +
108    labs(title='Qualifications', x='', y='Frequency') +
109    theme_grey() +
110    theme(plot.title = element_text(hjust = 0.5, size = 15)) +
111    theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5))
112
113  # What is percentage of divorced and separated people?
114  # Plot all people
115  filter(HSE_2011, (marstatc > 0)) %>% # to remove NA values
116    ggplot(aes(as_factor(marstatc))) +
117    geom_bar(aes(y=after_stat(prop), group=1), colour='black', fill='pink') +
118    labs(title='Marital status', x='', y='Percent') +
119    geom_text(aes(label = scales::percent(after_stat(prop)),
120                  y=after_stat(prop),group=1), stat= 'count', vjust = -0.5) +
121    scale_y_continuous(labels=scales::percent) +
122    theme_grey() +
123    theme(plot.title = element_text(hjust = 0.5, size = 15)) +
124    theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5))
125
126  attributes(HSE_2011$marstatc)
127
128  table(
129    as_factor(HSE_2011$marstatc)
130    ,useNA="ifany"
131  )
132
133  prop.table(
134    table(
135      as_factor(HSE_2011$marstatc)
136      ,useNA="ifany"
137    )
138  )*100
139
140
141  # Find the mean, median, mode, minimum, maximum, range and standard deviation of household size,
142  # BMI and age at last birthday.
143  # First using the total sample
144  mean(HSE_2011$HHSize)
145  median(HSE_2011$HHSize)
146
147  # Modes function from https://stackoverflow.com/questions/2547402/how-to-find-the-statistical-mode
148  # Uses a function that calculated Mode, enhanced to remove NA values plus another that calculates
149  # multimodes. Combined here
150  Modes <- function(x, na.rm = FALSE) {
151    if(na.rm){
152      x = x[!is.na(x)]
153    }
154    ux <- unique(x)
155    tab <- tabulate(match(x, ux))
156    ux[tab == max(tab)]
157  }
158
159  Modes(HSE_2011$HHSize)
160  sort(table(HSE_2011$HHSize),decreasing=T)
161  min(HSE_2011$HHSize)
162  max(HSE_2011$HHSize)
163  range(HSE_2011$HHSize)
164  sd(HSE_2011$HHSize)
165  par(mar=c(5,5,4,4))
166  hist(HSE_2011$HHSize,main="Household size",xlab="Number of people",col="lightblue")
167
```

```r
168  ggplot(HSE_2011,
169         aes(x=HHSize)) +
170    geom_histogram(colour='black' ,fill='light blue', binwidth=1) +
171    labs(title='Household size', x='Number of people', y='Frequency') +
172    scale_x_continuous(breaks=1:10) +
173    theme(plot.title = element_text(hjust = 0.5, size = 15))
174
175  SummaryDataHHS <- data.frame(Attribute = c('Mean', 'Median', 'Mode'),
176                               Value = c(mean(HSE_2011$HHSize),
177                                         median(HSE_2011$HHSize),
178                                         Mode(HSE_2011$HHSize)
179                               ))
180
181
182  ggplot(HSE_2011,
183         aes(HHSize)) +
184    geom_histogram(colour='black' ,fill='light blue', binwidth=1) +
185    geom_vline(data=SummaryDataHHS,
186               aes(xintercept=Value,
187                   col=Attribute), size=1) +
188    labs(title='Household size', x='Number of people', y='Frequency') +
189    scale_x_continuous(breaks=1:10) +
190    theme(plot.title = element_text(hjust = 0.5, size = 15))
191
192  SummaryDataHHS <- data.frame(Attribute = c('Mean', 'Mean -1 SD', 'Mean +1 SD'),
193                               Value = c(mean(HSE_2011$HHSize),
194                                         mean(HSE_2011$HHSize)-sd(HSE_2011$HHSize),
195                                         mean(HSE_2011$HHSize)+sd(HSE_2011$HHSize)
196                               ))
197
198
199  ggplot(HSE_2011,
200         aes(HHSize)) +
201    geom_histogram(colour='black' ,fill='light blue', binwidth=1) +
202    geom_vline(data=SummaryDataHHS,
203               aes(xintercept=Value,
204                   col=Attribute), size=1) +
205    labs(title='Household size', x='Number of people', y='Frequency') +
206    scale_x_continuous(breaks=1:10) +
207    theme(plot.title = element_text(hjust = 0.5, size = 15))
208
```

```r
209  # Boxplot of household size
210  ggplot(HSE_2011,
211         aes(HHSize)) +
212    geom_boxplot(colour='black', fill='light blue') +
213    stat_boxplot(geom ='errorbar', width=0.3) +
214    coord_flip() +
215    labs(title='Household size', x='Number of people') +
216    scale_x_continuous(breaks=1:10) +
217    scale_y_discrete() +
218    theme(plot.title = element_text(hjust = 0.5, size = 15))
219
220
221  # Need to remove NA values
222  mean(HSE_2011$bmival, na.rm=TRUE)
223  median(HSE_2011$bmival, na.rm=TRUE)
224  Modes(HSE_2011$bmival, na.rm=TRUE)
225  head(sort(table(HSE_2011$bmival),decreasing=T), n=200)
226  # We can see multiple BMI values with two counts
227
228  HSE_2011 %>% count(bmival) %>% arrange(desc(n))
229
230  min(HSE_2011$bmival, na.rm=TRUE)
231  max(HSE_2011$bmival, na.rm=TRUE)
232  range(HSE_2011$bmival, na.rm=TRUE)
233  sd(HSE_2011$bmival, na.rm=TRUE)
234  par(mar=c(5,5,4,4))
235  hist(HSE_2011$bmival,breaks=50,main="BMI value - 50 breaks",xlab="BMI",col="lightblue")
236
237  SummaryDataBMI <- data.frame(Attribute = c('Mean', 'Median'),
238                               Value = c(mean(HSE_2011$bmival, na.rm=TRUE),
239                                         median(HSE_2011$bmival, na.rm=TRUE)
240                               ))
241
242
243  # Histogram of BMI
244  ggplot(HSE_2011,
245         aes(bmival)) +
246    geom_histogram(colour='black' ,fill='light blue', bins=50, na.rm = TRUE) +
247    geom_vline(data=SummaryDataBMI,
248               aes(xintercept=Value,
249                   col=Attribute), size=1) +
250    labs(title='BMI', x='BMI', y='Count') +
251    scale_x_continuous(breaks = seq(from = 10, to = 100, by = 10)) +
252    theme(plot.title = element_text(hjust = 0.5, size = 15))
253
```

```r
254  SummaryDataBMI <- data.frame(Attribute = c('Mean', 'Mean -1 SD', 'Mean +1 SD'),
255                               Value = c(mean(rounded_bmi_1dp, na.rm=TRUE),
256                                         mean(rounded_bmi_1dp, na.rm=TRUE)-sd(HSE_2011$bmival,
257                                                                              na.rm=TRUE),
258                                         mean(rounded_bmi_1dp, na.rm=TRUE)+sd(HSE_2011$bmival,
259                                                                              na.rm=TRUE)
260                                ))
261  # Histogram of BMI
262  ggplot(HSE_2011,
263         aes(bmival)) +
264    geom_histogram(colour='black' ,fill='light blue', bins=50, na.rm = TRUE) +
265    geom_vline(data=SummaryDataBMI,
266               aes(xintercept=Value,
267                   col=Attribute), size=1) +
268    labs(title='BMI - total sample', x='BMI', y='Frequency') +
269    scale_x_continuous(breaks = seq(from = 10, to = 100, by = 10)) +
270    theme(plot.title = element_text(hjust = 0.5, size = 15))
271
272  # Boxplot of BMI by gender
273  ggplot(HSE_2011,
274         aes(bmival)) +
275    geom_boxplot(colour='black', fill='light blue', na.rm = TRUE) +
276    stat_boxplot(geom ='errorbar', width=0.3) +
277    coord_flip() +
278    labs(title='BMI value', x='') +
279    scale_x_continuous(breaks=seq(0, 70, by=10)) +
280    scale_y_discrete() +
281    facet_grid(~as_factor(Sex)) +
282    theme(plot.title = element_text(hjust = 0.5, size = 15))
283
284  mean(HSE_2011$Age)
285  median(HSE_2011$Age)
286  Modes(HSE_2011$Age)
287  # check with sorted table
288  sort(table(HSE_2011$Age),decreasing=T)
289  min(HSE_2011$Age)
290  max(HSE_2011$Age)
291  range(HSE_2011$Age)
292  sd(HSE_2011$Age)
293
294  par(mar=c(5,5,4,4))
295  hist(HSE_2011$Age,breaks=25,main="Age - 25 breaks",xlab="Age",col="lightblue")
296
297  SummaryDataAge <- data.frame(Attribute = c('Mean', 'Median', 'Mode 1', 'Mode 2'),
298                               Value = c(mean(HSE_2011$Age, na.rm=TRUE),
299                                         median(HSE_2011$Age, na.rm=TRUE),
300                                         Modes(HSE_2011$Age, na.rm=TRUE)
301                                ))
302
303  ggplot(HSE_2011,
304         aes(Age)) +
305    geom_histogram(colour='black' ,fill='light blue', bins=101, na.rm = TRUE) +
306    geom_vline(data=SummaryDataAge,
307               aes(xintercept=Value,
308                   col=Attribute), size=1.25) +
309    labs(title='Age', x='Age', y='Frequency') +
310    scale_x_continuous(breaks = seq(from = 0, to = 100, by = 10)) +
311    theme(plot.title = element_text(hjust = 0.5, size = 15))
312
313  SummaryDataAge <- data.frame(Attribute = c('Mean', 'Mean -1 SD', 'Mean +1 SD'),
314                               Value = c(mean(HSE_2011$Age, na.rm=TRUE),
315                                         mean(HSE_2011$Age, na.rm=TRUE)-sd(HSE_2011$Age, na.rm=TRUE),
316                                         mean(HSE_2011$Age, na.rm=TRUE)+sd(HSE_2011$Age, na.rm=TRUE)
317                                ))
318
319  ggplot(HSE_2011,
320         aes(Age)) +
321    geom_histogram(colour='black' ,fill='light blue', bins=101, na.rm = TRUE) +
322    geom_vline(data=SummaryDataAge,
323               aes(xintercept=Value,
324                   col=Attribute), size=1.25) +
325    labs(title='Age', x='Age', y='Frequency') +
326    scale_x_continuous(breaks = seq(from = 0, to = 100, by = 10)) +
327    theme(plot.title = element_text(hjust = 0.5, size = 15))
328
329
330
331  # Inferential Statistics.
332  # Run a significance test to find out which gender drinks more alcohol.
333  # First using total weekly units as the measure of 'more alcohol'
334  # Check for normality using Kolmogorov-Smirnov test because the sample size is large
335  ks.test(HSE_2011$totalwu, 'pnorm')
336
337  hist(HSE_2011$totalwu,
338       breaks=20,
339       col='steelblue',
340       main='Histogram to test for normal distribution of alcohol consumption',
341       xlab='Units of alcohol per week')
342
```

```r
343  qqnorm(HSE_2011$totalwu,
344        col='steelblue',
345        main='Q-Q plot to test for normal distribution of alcohol consumption')
346  qqline(HSE_2011$totalwu)
347
348  # Mann-Whitney U test to test for significance since the data is not normally distributed
349  wilcox.test(totalwu ~ Sex, data=HSE_2011)
350  t.test(totalwu ~ as_factor(Sex), data=HSE_2011)
351
352  # Next using ddnow as the measure of 'more alcohol'
353  # drink nowadays and gender are both binary values so use a contingency table and perform a
354  # chi-squared test.
355  table(as_factor(HSE_2011$dnnow), as_factor(HSE_2011$Sex))
356  chisq.test(table(HSE_2011$dnnow, HSE_2011$Sex))
357
358  filter(HSE_2011, (dnnow > 0)) %>% # to remove NA values
359    ggplot(aes(as_factor(dnnow))) +
360    geom_bar(aes(y=after_stat(prop), group=1), colour='black', fill='pink') +
361    labs(title='Drink by gender', x='Drink?', y='Percent') +
362    geom_text(aes(label = scales::percent(after_stat(prop)),
363              y=after_stat(prop),group=1), stat= 'count', vjust = 2) +
364    scale_y_continuous(labels=scales::percent) +
365    facet_grid(~as_factor(Sex)) +
366    theme_grey() +
367    theme(plot.title = element_text(hjust = 0.5, size = 15)) +
368    theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5))
369
370
371  # Boxplot focusing on Q2, median and Q3
372  ggplot(HSE_2011, aes(x = '', y=totalwu)) +
373    geom_boxplot(fill='steel blue') +
374    stat_boxplot(geom ='errorbar', width=0.3) +
375    labs(title = "Alcohol consumption by gender", y = 'Total weekly units', x='') +
376    facet_grid(~as_factor(Sex)) +
377    coord_cartesian(ylim = quantile(HSE_2011$totalwu, na.rm=TRUE, c(0.1, 0.97))) +
378    theme_grey() +
379    theme(plot.title = element_text(hjust = 0.5, size = 15))
380
381  mean(HSE_2011$totalwu[as_factor(HSE_2011$Sex)=='Male'], na.rm=TRUE)
382  mean(HSE_2011$totalwu[as_factor(HSE_2011$Sex)=='Female'], na.rm=TRUE)
383  median(HSE_2011$totalwu[as_factor(HSE_2011$Sex)=='Male'], na.rm=TRUE)
384  median(HSE_2011$totalwu[as_factor(HSE_2011$Sex)=='Female'], na.rm=TRUE)
385
386  # Run a significance test to find out which region drinks the most alcohol.
387  # Use the Krushal-Wallis test to compare the alcohol consumption across the nine government regions:
388  kruskal.test(totalwu ~ gor1, data=HSE_2011)

389  |
390  table(as_factor(HSE_2011$gor1), as_factor(HSE_2011$dnnow))
391
392  chisq.test(table(HSE_2011$gor1, HSE_2011$dnnow))
393
394  filter(HSE_2011, (dnnow > 0)) %>% # to remove NA values
395    ggplot(aes(as_factor(dnnow))) +
396    geom_bar(aes(y=after_stat(prop), group=1), colour='black', fill='pink') +
397    labs(title='Drink by region', x='Drink?', y='Percent') +
398    geom_text(aes(label = scales::percent(after_stat(prop)),
399              y=after_stat(prop),group=1), stat= 'count', vjust = 2) +
400    scale_y_continuous(labels=scales::percent) +
401    facet_grid(~as_factor(gor1)) +
402    theme_grey() +
403    theme(plot.title = element_text(hjust = 0.5, size = 15)) +
404    theme(axis.text.x=element_text(angle=90,hjust=1,vjust=0.5))
405
406
407  # Boxplot of alcohol by region focusing on Q1, median and Q3
408  ggplot(HSE_2011, aes(x = '', y=totalwu)) +
409    geom_boxplot(fill='steel blue') +
410    stat_boxplot(geom ='errorbar', width=0.3) +
411    labs(title = "Alcohol consumption by region", y = 'Total weekly units', x='') +
412    facet_grid(~as_factor(gor1)) +
413    coord_cartesian(ylim = quantile(Adults$totalwu, na.rm=TRUE, c(0, 0.95))) +
414    theme_grey() +
415    theme(plot.title = element_text(hjust = 0.5, size = 15))
416
417
418  # Investigate whether there is a statistical difference between men and women on
419  # the following variables:
420  #    valid height.
421  # Create a subset of adults
422  Adults <- subset(HSE_2011,Age > 17)
423  # Check to ensure only Adults are included:
424  range(Adults$Age)
425  # Now count the number of adults in the sample
426  n_distinct(Adults$pserial)
427
428  # Test for normal distribution
429  shapiro.test(Adults$htval[as_factor(Adults$Sex)=='Male'])
430  shapiro.test(Adults$htval[as_factor(Adults$Sex)=='Female'])
431
```

```r
432  ggplot(Adults,
433         aes(x=htval)) +
434    geom_histogram(aes(y=..density..), binwidth=2,
435                   colour='black', fill='steel blue') +
436    geom_density(alpha=.5, fill='red') +
437    facet_grid(~as_factor(Sex)) +
438    labs(title = "Adult height distribution by gender", y = 'Density', x='Height') +
439    theme_grey() +
440    theme(plot.title = element_text(hjust = 0.5, size = 15))
441
442  qqnorm(Adults$htval[as_factor(Adults$Sex)=='Male'],
443         col='steelblue',
444         main='Q-Q plot to test for normal distribution of height - Male')
445  qqline(Adults$htval[as_factor(Adults$Sex)=='Male'])
446
447  qqnorm(Adults$htval[as_factor(Adults$Sex)=='Female'],
448         col='steelblue',
449         main='Q-Q plot to test for normal distribution of height - Female')
450  qqline(Adults$htval[as_factor(Adults$Sex)=='Female'])
451
452  # t-test to test for significance
453  t.test(htval ~ as_factor(Sex), data=Adults)
454
455  ggplot(Adults, aes(x ='', y=htval)) +
456    geom_boxplot(fill='steel blue') +
457    stat_boxplot(geom ='errorbar', width=0.5) +
458    labs(title = "Height of adults by gender", y = 'Height (CM)', x='') +
459    facet_grid(~as_factor(Sex)) +
460    theme_grey() +
461    theme(plot.title = element_text(hjust = 0.5, size = 15))
462
463  #   Valid height.
464  # Test for normal distribution
465  shapiro.test(Adults$wtval[as_factor(Adults$Sex)=='Male'])
466  shapiro.test(Adults$wtval[as_factor(Adults$Sex)=='Female'])
467
468  ggplot(Adults,
469         aes(x=wtval)) +
470    geom_histogram(aes(y=..density..), binwidth=2,
471                   colour='black', fill='steel blue') +
472    geom_density(alpha=.5, fill='red') +
473    facet_grid(~as_factor(Sex)) +
474    labs(title = "Adult weight distribution by gender", y = 'Density', x='Weight') +
475    theme_grey() +
476    theme(plot.title = element_text(hjust = 0.5, size = 15))
477
```

```r
478  qqnorm(Adults$wtval[as_factor(Adults$Sex)=='Male'],
479         col='steelblue',
480         main='Q-Q plot to test for normal distribution of weight - Male')
481  qqline(Adults$wtval[as_factor(Adults$Sex)=='Male'])
482
483  qqnorm(Adults$wtval[as_factor(Adults$Sex)=='Female'],
484         col='steelblue',
485         main='Q-Q plot to test for normal distribution of weight - Female')
486  qqline(Adults$wtval[as_factor(Adults$Sex)=='Female'])
487
488  # Even though the distribution is not normal, we will still use the t-test
489  # because the sample size is large
490  t.test(wtval ~ as_factor(Sex), data=Adults)
491  #wilcox.test(wtval ~ Sex, data=Adults)
492
493  ggplot(Adults, aes(x = '', y=wtval)) +
494    geom_boxplot(fill='steel blue') +
495    stat_boxplot(geom ='errorbar', width=0.3) +
496    labs(title = "Weight of adults by gender", y = 'Weight (Kg)', x='') +
497    facet_grid(~as_factor(Sex)) +
498    theme_grey() +
499    theme(plot.title = element_text(hjust = 0.5, size = 15))
500
501  # What is the correlation between whether a person drinks nowadays, total
502  # household income, age at last birthday and gender?
503  # ks.test(HSE_2011$Sex, 'pnorm')
504
505  # drink nowadays is binary and total household income is ratio
506  # Kolmogorov-Smirnov because sample size is too large for Shapiro Wilk test
507  ks.test(HSE_2011$dnnow, 'pnorm')
508  ks.test(HSE_2011$totinc, 'pnorm')
509  ks.test(HSE_2011$Age, 'pnorm')
510  ks.test(HSE_2011$Sex, 'pnorm')
511
512  hist(HSE_2011$dnnow,
513       breaks=100,
514       col='steelblue',
515       main='Histogram to test for normal distribution of drink nowadays',
516       xlab='Units of alcohol per week')
517
518  hist(HSE_2011$totinc,
519       breaks=100,
520       col='steelblue',
521       main='Histogram to test for normal distribution of total household income',
522       xlab='Units of alcohol per week')
523
```

```r
524  hist(HSE_2011$Age,
525       breaks=100,
526       col='steelblue',
527       main='Histogram to test for normal distribution of age',
528       xlab='Units of alcohol per week')
529
530  hist(HSE_2011$Sex,
531       breaks=100,
532       col='steelblue',
533       main='Histogram to test for normal distribution of gender',
534       xlab='Units of alcohol per week')
535
536  qqnorm(HSE_2011$dnnow,
537         col='steelblue',
538         main='Q-Q plot to test for normal distribution of drink nowadays')
539  qqline(HSE_2011$dnnow)
540
541  qqnorm(HSE_2011$totinc,
542         col='steelblue',
543         main='Q-Q plot to test for normal distribution of total household income')
544  qqline(HSE_2011$totinc)
545
546  qqnorm(HSE_2011$Age,
547         col='steelblue',
548         main='Q-Q plot to test for normal distribution of age')
549  qqline(HSE_2011$Age)
550
551  qqnorm(HSE_2011$Sex,
552         col='steelblue',
553         main='Q-Q plot to test for normal distribution of gender')
554  qqline(HSE_2011$Sex)
555
556  # Not normal distribution so use spearman's rank
557  cor.test(HSE_2011$dnnow, HSE_2011$totinc, method='spearman', use = 'complete.obs')
558  cor.test(HSE_2011$dnnow, HSE_2011$Age, method='spearman', use = 'complete.obs')
559  cor.test(HSE_2011$dnnow, HSE_2011$Sex, method='spearman', use = 'complete.obs')
560
561  cor.test(HSE_2011$totinc, HSE_2011$Age, method='spearman', use = 'complete.obs')
562  cor.test(HSE_2011$totinc, HSE_2011$Sex, method='spearman', use = 'complete.obs')
563
564  cor.test(HSE_2011$Age, HSE_2011$Sex, method='spearman', use = 'complete.obs')
565
566  cordata <- HSE_2011[,c('dnnow', 'totinc', 'Age', 'Sex')]
567  corrplot(cor(cordata, method='spearman', use = 'complete.obs'), method='pie')
568
569
```