

STK-IN4300
**Statistical Learning Methods in Data
Science**

OBLIG 1

Egil Furnes
Student: 693784

Problem 1

a)

I have chosen the dataset `HistData::PearsonLee` which is inspired by a similar dataset from Galton on heights of parents and their offspring. From the description of the dataset:

Wachsmuth et. al (2003) noticed that a loess smooth through Galton's data on heights of mid-parents and their offspring exhibited a slightly non-linear trend, and asked whether this might be due to Galton having pooled the heights of fathers and mothers and sons and daughters in constructing his tables and graphs.

This is an expansion on Galton's data to find if there are also gender-effects exhibited on the relation between parent-child height.

```
1 df <- HistData::PearsonLee %>% as_tibble()
2 df %>% summary()
```

	child	parent	frequency	gp
1	Min. :52.50	Min. :52.50	Min. : 0.250	fd:206
2	1st Qu.:62.50	1st Qu.:61.50	1st Qu.: 0.750	fs:179
3	Median :65.50	Median :65.50	Median : 2.500	md:185
4	Mean :66.05	Mean :65.01	Mean : 6.558	ms:176
5	3rd Qu.:69.50	3rd Qu.:68.50	3rd Qu.: 9.250	
6	Max. :79.50	Max. :75.50	Max. :46.500	
7				
8	par	chl		
9	Father:385	Daughter:391		
10	Mother:361	Son :355		

b)

In this task I design a few bad and good charts. First we look at the distribution of the categorical variable `par`, which is a binary variable with two levels, Father and Mother. To make it even 'better' for our readers we have plotted this distribution using polar coordinates, such that one must read the distribution by angle rather than height or length. Furthermore we use blue for Mother and red for Father which is probably against convention and can contribute to making the plot more difficult to read.

```
1 ggplot(df, aes(x = par, weight = frequency, fill = par)) +
2   geom_bar() +
3   coord_polar(theta = "y") +
4   labs(title = "Parent sex distribution",
5         x = NULL,
6         y = NULL) +
7   theme_minimal()
```

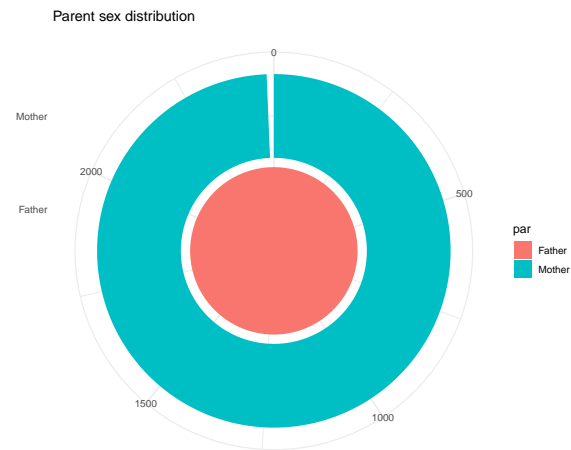


Figure 1: Parent sex distribution

Next we plot a distribution of the height of the children in the variable `child`, which is a continuous variable. It does make sense to plot a distribution of a continuous variable – but in our plot we use way too few bins in `geom_histogram` with 3, which only divides the height of the children into 3 categories. Furthermore, using `coord_flip` in this plot makes it have an unfortunate visual representation.

```
1 ggplot(df, aes(x = child, weight = frequency)) +
2   geom_histogram(bins = 3,
3                   fill = "#D2B48C",
4                   color = "black") +
5   ggtitle("Distribution of child height") +
6   coord_flip() +
7   theme_minimal()
```

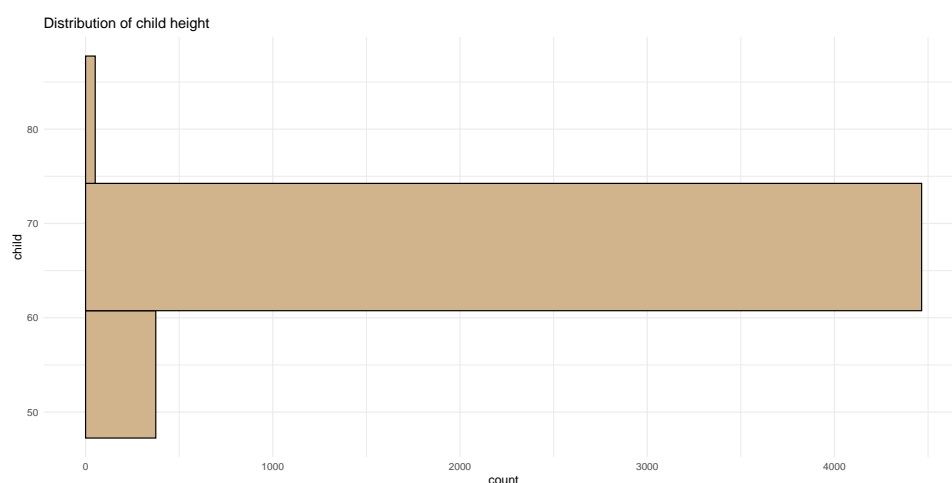


Figure 2: Distribution of child height

c)

In this task we are set to improve the plots or maybe even make them ‘good’. Regarding the distribution of children height in inches I have now removed the `coord_flip` command and furthermore increased the bins to show an approximated normal distribution in height.

```
1 ggplot(df, aes(x = child, weight = frequency)) +
2   geom_histogram(bins = 29,
3                 fill = "#D2B48C",
4                 color = "black") +
5   ggtitle("Distribution of child height") +
6   xlab("Height of children in inches") +
7   ylab("Count") +
8   theme_minimal()
```

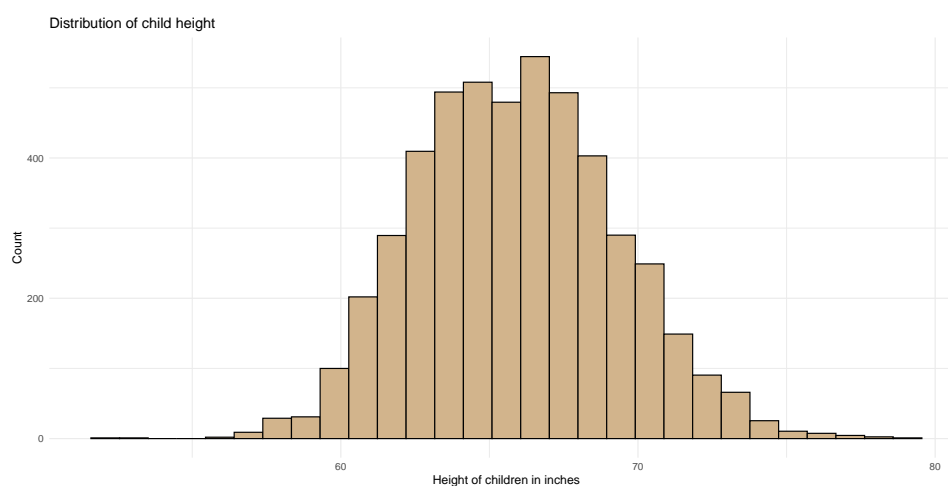


Figure 3: Caption

Regarding the distribution of parents, I found that just plotting the count for Father 2454 and Mother 2438 would be two boring bars. Therefore I added the distribution of the parent heights as density plots and split them by sex. As such we find some cool insights, such as a left-shift for mothers and right-shift for father indicating fathers on average are taller than mothers. Furthermore the distribution of mothers look ‘taller’ and for fathers ‘wider’.

```
1 ggplot(df, aes(
2   x = parent,
3   weight = frequency,
4   fill = par,
5   colour = par
6 )) +
7   geom_density(alpha = 0.35, adjust = 1.1) +
8   scale_fill_manual(values = c(Father = "#4C78A8", Mother = "#F56598")) +
9   scale_colour_manual(values = c(Father = "#4C78A8", Mother = "#F56598")) +
```

```

10 labs(
11   title = "Parent height distribution by sex",
12   x = "Parent height in inches",
13   y = "Density",
14   fill = NULL,
15   colour = NULL
16 ) +
17 theme_minimal()

```

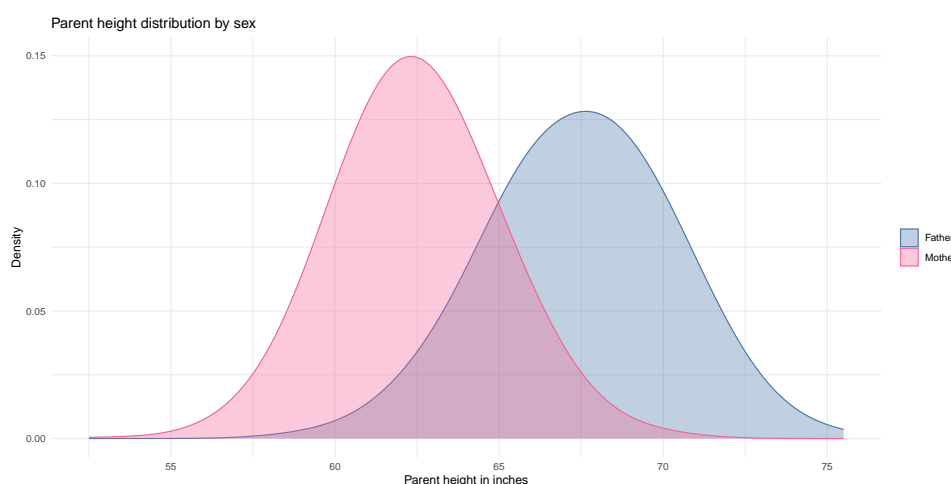


Figure 4: Parent height distribution in inches by sex

Problem 2

In this dataset we initially have 6 variables `child`, `parent`, `frequency`, `gp`, `par`, and `chl`.

However, looking at values for `gp` this is directly colinear with the pair of values from `par` and `chl`. As such I first remove the `gp` and rather just keep the two others. Furthermore the `frequency` variable is an aggregated count of families from methodology by Pearson and Lee, and as such is not relevant to explain the height of the children, its rather just a useful data marker. We therefore remove this one too.

```

1 > df %>% head(3)
2 # A tibble: 3 × 6
3   child parent frequency gp    par    chl
4   <dbl>  <dbl>      <dbl> <fct> <fct> <fct>
5 1  59.5   62.5         0.5 fs    Father Son
6 2  59.5   63.5         0.5 fs    Father Son
7 3  59.5   64.5         1  fs    Father Son

```

Now fitting the data to the model we just use a simple linear regression `lm()` model and try out a stepwise model selection using both ‘forward’ and ‘backward’ selection with the argument both in the function `stepAIC()` from the `MASS` package.

```

1 library(MASS)
2 df <- HistData::PearsonLee %>% as_tibble()
3 lm1 <- lm(child~parent+par+chl, data = df, weights = frequency)
4 lm2 <- stepAIC(lm1, direction = "both", trace = FALSE)

```

From the regression output we actually find that both the `lm1()` simple linear regression and the `lm2()` stepwise model selection are equal. This is explained by all the used variables in the model being significant and so at an 0 *** level.

```

1 > lm1
2
3 Call:
4 lm(formula = child ~ parent + par + chl, data = df, weights =
   frequency)
5
6 Coefficients:
7 (Intercept)      parent      parMother      chlSon
8      28.4825      0.5216      2.8293      4.6985
9
10 > lm2
11
12 Call:
13 lm(formula = child ~ parent + par + chl, data = df, weights =
   frequency)
14
15 Coefficients:
16 (Intercept)      parent      parMother      chlSon
17      28.4825      0.5216      2.8293      4.6985

```

Now looking at the regression output of the model we can assess how it is at explaining the height of children using the other variables. First off, we see that all input variables are statistically significant *** including the Intercept.

Although this doesn't necessarily mean that the variables explain the output well. Now looking at the explanatory power of this model, we read the R-statistic, either the Multiple R-squared or the Adjusted R-squared which are 0.5658 and 0.5641 respectively, both being higher than the sort of threshold at 0.3, and with a statistically significant p-value of $< 2.2e-16$.

```

1 > lm1 %>% summary()
2
3 Call:
4 lm(formula = child ~ parent + par + chl, data = df, weights =
   frequency)
5
6 Weighted Residuals:
7      Min       1Q   Median       3Q      Max
8 -14.0597  -4.5536   0.2846   4.6882  12.7085
9

```

```

10 Coefficients:
11      Estimate Std. Error t value Pr(>|t|)
12 (Intercept)  28.4825     2.2150   12.86  <2e-16 ***
13 parent       0.5216     0.0327   15.95  <2e-16 ***
14 parMother    2.8293     0.2356   12.01  <2e-16 ***
15 chlSon       4.6985     0.1708   27.50  <2e-16 ***
16 ---
17 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
18
19 Residual standard error: 5.917 on 742 degrees of freedom
20 Multiple R-squared:  0.5658,    Adjusted R-squared:  0.5641
21 F-statistic: 322.3 on 3 and 742 DF,  p-value: < 2.2e-16

```

In sum, we can probably say that the linear model appears to be sufficient in explaining this scenario.