

Assignment 2

ETC1010_5510

Mohammed Faizan

Monday, May 17 2021

```
library(naniar)
library(broom)
library(ggmap)
library(knitr)
library(lubridate)
library(timeDate)
library(tsibble)
library(here)
library(readr)
library(tidyverse)
library(ggResidpanel)
library(gridExtra)
library(kableExtra)
```

```
tree_data0 <- read_csv("Data/Assignment_data.csv")
```

Part I

Question 1: Rename the variables *Date Planted* and *Year Planted* to *Dateplanted* and *Yearplanted* using the *rename()* function. Make sure *Dateplanted* is defined as a date variable. Then extract from the variable *Dateplanted* the year and store it in a new variable called *Year*. Display the first 6 rows of the data frame. (5pts)

```
tree_data <- as_tibble(tree_data0) %>%
  rename(Dateplanted=c("Date Planted"),
         Yearplanted=c("Year Planted")) %>%
  mutate(Dateplanted = dmy(Dateplanted)) %>%
  mutate(Year = year(Dateplanted))

tree_data %>%
  head() %>%
  kable(caption = "Tree Data: All Variables") %>%
  kable_styling(latex_options = c("scale_down", "hold_position"))
```

Table 1: Tree Data: All Variables

CoM ID	Common Name	Scientific Name	Genus	Family	Diameter Breast Height	Yearplanted	Dateplanted	Age Description	Useful Life Expectancy	Useful Life Expectancy Value	Precinct	Located in	UploadDate	CoordinateLocation	Latitude	Longitude	Easting	Northing	Year
1057605	White Poplar	Populus alba	Populus	Salicaceae	NA	1900	2000-01-01	NA	NA	NA	NA	Park	30/7/20	(-37.790803752047804, 144.9229711420543)	-37.79089	144.9228	317080.6	581535.3	2000
1028440	London Plane	Platanus x acerifolia	Platanus	Platanaceae	62	1900	2000-01-02	Mature	6-10 years (>50% canopy)	10	NA	Park	30/7/20	(-37.84544711087132, 144.9790135884796)	-37.84545	144.9790	322184.5	580940.8	2000
1058665	Small-leaved Linden	Tilia cordata	Tilia	Malvaceae	19	2000	2000-05-29	Semi-Mature	31-60 years	60	NA	Street	30/7/20	(-37.789897713116666, 144.96853113075328)	-37.78989	144.9685	321148.5	581453.1	2000
1026352	Variegated Elm	Ulmus minor	Ulmus	Ulmaceae	26	1900	2000-01-02	Semi-Mature	21-30 years	30	NA	Street	30/7/20	(-37.811003176311319, 144.98350174882896)	-37.81100	144.9836	322560.1	581317.2	2000
1038440	Canary Island Pine	Pinus canariensis	Pinus	Pinaceae	91	1900	2000-01-02	Mature	31-60 years	60	NA	Park	30/7/20	(-37.78240938542976, 144.96107078451748)	-37.78241	144.9602	320173.4	581036.7	2000
1015128	London Plane	Platanus x acerifolia	Platanus	Platanaceae	99	1900	2000-01-02	Mature	11-20 years	20	NA	Street	30/7/20	(-37.79916133896208, 144.9484631137445)	-37.79912	144.9488	319560.1	581440.1	2000

```
tree_data %>%
  select(c(1:7)) %>%
  head() %>%
  kable(caption = "Tree Data") %>%
  kable_styling(latex_options = c("hold_position"))
```

Table 2: Tree Data

CoM ID	Common Name	Scientific Name	Genus	Family	Diameter Breast Height	Yearplanted
1057605	White Poplar	Populus alba	Populus	Salicaceae	NA	1900
1028440	London Plane	Platanus x acerifolia	Platanus	Platanaceae	62	1900
1058665	Small-leaved Linden	Tilia cordata	Tilia	Malvaceae	19	2000
1026352	Variegated Elm	Ulmus minor	Ulmus	Ulmaceae	26	1900
1038440	Canary Island Pine	Pinus canariensis	Pinus	Pinaceae	91	1900
1015128	London Plane	Platanus x acerifolia	Platanus	Platanaceae	99	1900

```
tree_data %>%
  select(c(8:13)) %>%
  head() %>%
  kable(caption = "Tree Data") %>%
  kable_styling(latex_options = c("hold_position"))
```

Table 3: Tree Data

Dateplanted	Age Description	Useful Life Expectancy	Useful Life Expectancy Value	Precinct	Located in
2000-01-01	NA	NA	NA	NA	Park
2000-01-02	Mature	6-10 years (>50% canopy)	10	NA	Park
2000-05-29	Semi-Mature	31-60 years	60	NA	Street
2000-01-02	Semi-Mature	21-30 years	30	NA	Street
2000-01-02	Mature	31-60 years	60	NA	Park
2000-01-02	Mature	11-20 years	20	NA	Street

```
tree_data %>%
  select(c(14:20)) %>%
  head() %>%
  kable(caption = "Tree Data") %>%
  kable_styling(latex_options = c("hold_position"))
```

Table 4: Tree Data

UploadDate	CoordinateLocation	Latitude	Longitude	Easting	Northing	Year
30/9/20	(-37.790893782047654, 144.92257141425543)	-37.79089	144.9226	317080.6	5815353	2000
30/9/20	(-37.84544751087332, 144.97904358884796)	-37.84545	144.9790	322184.5	5809408	2000
30/9/20	(-37.798902715216066, 144.96852132025538)	-37.79890	144.9685	321146.2	5814553	2000
30/9/20	(-37.81160317631319, 144.98355174888286)	-37.81160	144.9836	322500.1	5813172	2000
30/9/20	(-37.78240938533976, 144.96019784543748)	-37.78241	144.9602	320373.4	5816367	2000
30/9/20	(-37.79941634390208, 144.94984511347445)	-37.79942	144.9498	319503.0	5814460	2000

Question 2: Have you noticed any differences between the variables *Year* and *Yearplanted*? Why is that? Demonstrate your claims using R code. Fix the problem if there is one (Hint: Use *ifelse* inside a mutate function to fix the problem and store the data in *tree_data_clean*). After this question, please use the data in *tree_data_clean* to proceed. (3pts)

The corresponding values in the variables *Year* and *Yearplanted* are different. The newly created variable *Year* contains the year 2000 in all observations but one (1977). Correct value for the year of tree plantation is present in *Yearplanted*.

This difference is because the original variable, *Date Planted* in *tree_data0* has the *Date Planted* as “2/1/00” where the dmy() interprets the year “00” as 2000 and hence for both 1900 and 2000. The year 1977 is mapped correctly owing to the fact that 2077 has not yet arrived. These claims can be seen below.

```
tree_data %>%
  select(`CoM ID`, Yearplanted, Dateplanted, Year) %>%
  filter(`CoM ID` %in% c("1028440", "1058665", "1060068")) %>%
  kable(caption = "Mismatching Years") %>%
  kable_styling(latex_options = c("hold_position"))
```

Table 5: Mismatching Years

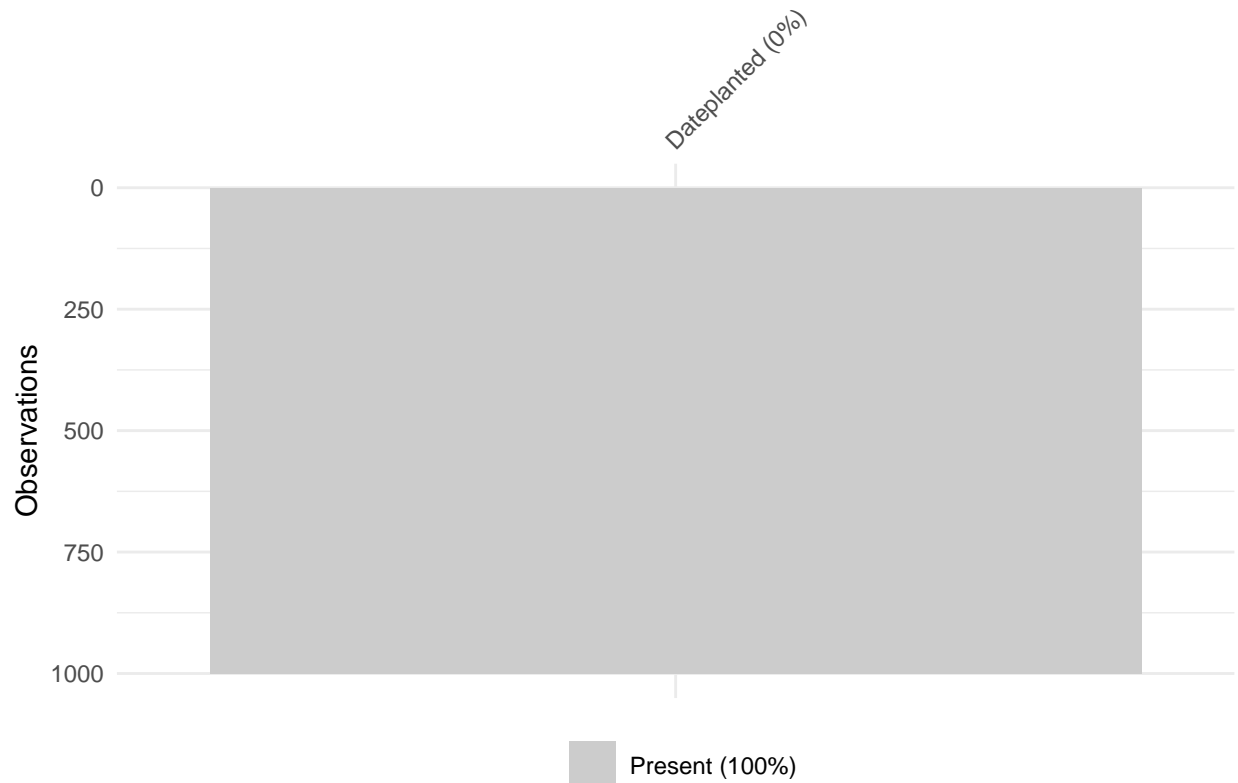
CoM ID	Yearplanted	Dateplanted	Year
1028440	1900	2000-01-02	2000
1058665	2000	2000-05-29	2000
1060068	1977	1977-07-07	1977

```
tree_data_clean <- tree_data %>%
  mutate(Dateplanted = str_replace(as.character(Dateplanted),
                                   "2000", as.character(Yearplanted))) %>%
  mutate(Year = Yearplanted) %>%
  mutate(ymd(Dateplanted))
```

Question 3: Investigate graphically the missing values in the variable *Dateplanted* for the last 1000 rows of the data set. What do you observe? (max 30 words) (2pts)

There is no missing values in the last 1000 observations of the variable *Dateplanted*.

```
tree_data_singlevariable <- tree_data_clean %>%  
  tail(1000) %>%  
  select(`Dateplanted`)  
visdat::vis_miss(tree_data_singlevariable)
```



Question 4: What is the proportion of missing values in each variable in the tree data set? Display the results in descending order of the proportion. (2pts)

The missingness in the variables of the tree data set is listed below in descending order of proportion. The first 8 rows below have some missing values. All other variables do not have any missing values.

```
miss_var_summary(tree_data_clean)
```

```
## # A tibble: 21 x 3
```

```
##      variable                n_miss pct_miss
##      <chr>                  <int>   <dbl>
## 1 Precinct                  6828 100
## 2 Diameter Breast Height    1454  21.3
## 3 Age Description            1454  21.3
## 4 Useful Life Expectency     1454  21.3
## 5 Useful Life Expectency Value 1454  21.3
## 6 Dateplanted                2    0.0293
## 7 ymd(Dateplanted)           2    0.0293
## 8 Common Name                1    0.0146
## 9 Located in                 1    0.0146
## 10 CoM ID                    0     0
## # ... with 11 more rows
```

Question 5: How many observations have a missing value in the variable *Dateplanted*? Identify the rows and display the information in those rows. Remove all the rows in the data set of which the variable *Dateplanted* has a missing value recorded and store the data in *tree_data_clean1*. Display the first 4 rows of *tree_data_clean1*. Use R inline code to complete the sentence below. (6pts)

There are 2 observations with missing values in Dateplanted variable. These observations are displayed below.

```
# n_missingobs <- as.data.frame(tree_data_clean$Dateplanted) %>% miss_case_table() %>% filter(n_miss_i
# tree_data_clean%>%filter(is.na(`Dateplanted`)) %>% count() #2 observations with NAs

tree_data_clean %>%
  filter(is.na(`Dateplanted`)) %>%
kable(caption = "Observations with missing Dateplanted") %>%
  kable_styling(latex_options = c("scale_down", "hold_position"))
```

Table 6: Observations with missing Dateplanted

CoM ID	Common Name	Scientific Name	Genus	Family	Diameter Breast Height	Yearplanted	Dateplanted	Age Description	Useful Life Expectency	Useful Life Expectency Value	Precinct	Located in	UploadDate	CoordinateLocation	Latitude	Longitude	Easting	Northing	Year	(ymd(Dateplanted))
1021157	Cyprus Plane	Platanus orientalis	Platanus	Platanaceae	22	1900	NA	Semi-Mature	21-30 years	80	NA	Street	30/9/20	[-37.81669289363358, 144.9691758246063]	-37.81669	144.9691	480060.3	5812825	1900	NA
1023092	London Plane	Platanus x acerifolia	Platanus	Platanaceae	29	1900	NA	Semi-Mature	6-30 years (>50% canopy)	80	NA	Street	30/9/20	[-37.80195671211686, 144.96871413201405]	-37.80196	144.9687	478994.6	5814256	1900	NA

```
tree_data_clean %>%
  filter(is.na(`Dateplanted`)) %>%
  select(c(1:7)) %>%
kable(caption = "Observations with missing Dateplanted") %>%
  kable_styling(latex_options = c("hold_position"))
```

Table 7: Observations with missing Dateplanted

CoM ID	Common Name	Scientific Name	Genus	Family	Diameter Breast Height	Yearplanted
1021155	Cyprus Plane	Platanus orientalis	Platanus	Platanaceae	22	1900
1023092	London Plane	Platanus x acerifolia	Platanus	Platanaceae	29	1900

```
tree_data_clean %>%
  filter(is.na(`Dateplanted`)) %>%
  select(c(8:13)) %>%
  kable(caption = "Observations with missing Dateplanted") %>%
  kable_styling(latex_options = c("hold_position"))
```

Table 8: Observations with missing Dateplanted

Dateplanted	Age Description	Useful Life Expectancy	Useful Life Expectancy Value	Precinct	Located in
NA	Semi-Mature	21-30 years	30	NA	Street
NA	Semi-Mature	6-10 years (>50% canopy)	10	NA	Street

```
tree_data_clean %>%
  filter(is.na(`Dateplanted`)) %>%
  select(c(14:20)) %>%
  kable(caption = "Observations with missing Dateplanted") %>%
  kable_styling(latex_options = c("hold_position"))
```

Table 9: Observations with missing Dateplanted

UploadDate	CoordinateLocation	Latitude	Longitude	Easting	Northing	Year
30/9/20	(-37.81605293763187, 144.95571785206653)	-37.81605	144.9557	320060.5	5812625	1900
30/9/20	(-37.80199573215045, 144.96671419201465)	-37.80200	144.9667	320994.6	5814206	1900

```
tree_data_clean1 <- tree_data_clean %>%
  filter(!is.na(`Dateplanted`))

head(tree_data_clean1, 4) %>%
  kable(caption = "tree_data_clean1") %>%
  kable_styling(latex_options = c("scale_down", "hold_position"))
```

\begin{table}[!h]

\caption{tree_data_clean1}

CoM ID	Common Name	Scientific Name	Genus	Family	Diameter Breast Height	Yearplanted	Dateplanted	Age Description	Useful Life Expectancy	Useful Life Expectancy Value	Precinct	Located in	UploadDate	CoordinateLocation	Latitude	Longitude	Easting	Northing	Year	(null)Dateplanted
1057605	White Poplar	Populus alba	Populus	Salicaceae	NA	1900	1900-01-01	NA	NA	NA	NA	Peak	30/9/20	(-37.798003952049054, 144.9229114429543)	-37.79809	144.9229	317680.6	5819365	1900	1900-01-01
1028440	London Plane	Platanus x acerifolia	Platanus	Platanaceae	62	1900	1900-01-02	Mature	6-10 years (>50% canopy)	10	NA	Peak	30/9/20	(-37.8464114957232, 144.9706448889196)	-37.84645	144.9706	322194.2	5809608	1900	1900-01-02
1058665	Small-leaved Linden	Tilia cordata	Tilia	Malvaceae	19	2000	1900-05-20	Semi-Mature	21-30 years	30	NA	Street	30/9/20	(-37.798003952049054, 144.96671419201465)	-37.79809	144.9667	321136.2	5819365	2000	2000-05-20
1026352	Variegated Elm	Ulmus minor	Ulmus	Ulmaceae	26	1900	1900-01-02	Semi-Mature	21-30 years	30	NA	Street	30/9/20	(-37.81605293763187, 144.96671419201465)	-37.81609	144.9667	322500.1	5811172	1900	1900-01-02

\end{table}

```
head(tree_data_clean1, 4) %>%
  select(c(1:7)) %>%
  kable(caption = "tree_data_clean1") %>%
  kable_styling(latex_options = c("hold_position"))
```

\begin{table}[!h]

\caption{tree_data_clean1}

CoM ID	Common Name	Scientific Name	Genus	Family	Diameter Breast Height	Yearplanted
1057605	White Poplar	Populus alba	Populus	Salicaceae	NA	1900
1028440	London Plane	Platanus x acerifolia	Platanus	Platanaceae	62	1900
1058665	Small-leaved Linden	Tilia cordata	Tilia	Malvaceae	19	2000
1026352	Variegated Elm	Ulmus minor	Ulmus	Ulmaceae	26	1900

\end{table}

```
head(tree_data_clean1, 4) %>%
  select(c(8:13)) %>%
  kable(caption = "tree_data_clean1") %>%
  kable_styling(latex_options = c("hold_position"))
```

\begin{table}[!h]

\caption{tree_data_clean1}

Dateplanted	Age Description	Useful Life Expectency	Useful Life Expectency Value	Precinct	Located in
1900-01-01	NA	NA	NA	NA	Park
1900-01-02	Mature	6-10 years (>50% canopy)	10	NA	Park
2000-05-29	Semi-Mature	31-60 years	60	NA	Street
1900-01-02	Semi-Mature	21-30 years	30	NA	Street

\end{table}

```
head(tree_data_clean1, 4) %>%
  select(c(14:20)) %>%
  kable(caption = "tree_data_clean1") %>%
  kable_styling(latex_options = c("hold_position"))
```

\begin{table}[!h]

\caption{tree_data_clean1}

UploadDate	CoordinateLocation	Latitude	Longitude	Easting	Northing	Year
30/9/20	(-37.790893782047654, 144.92257141425543)	-37.79089	144.9226	317080.6	5815353	1900
30/9/20	(-37.84544751087332, 144.97904358884796)	-37.84545	144.9790	322184.5	5809408	1900
30/9/20	(-37.798902715216066, 144.96852132025538)	-37.79890	144.9685	321146.2	5814553	2000
30/9/20	(-37.81160317631319, 144.98355174888286)	-37.81160	144.9836	322500.1	5813172	1900

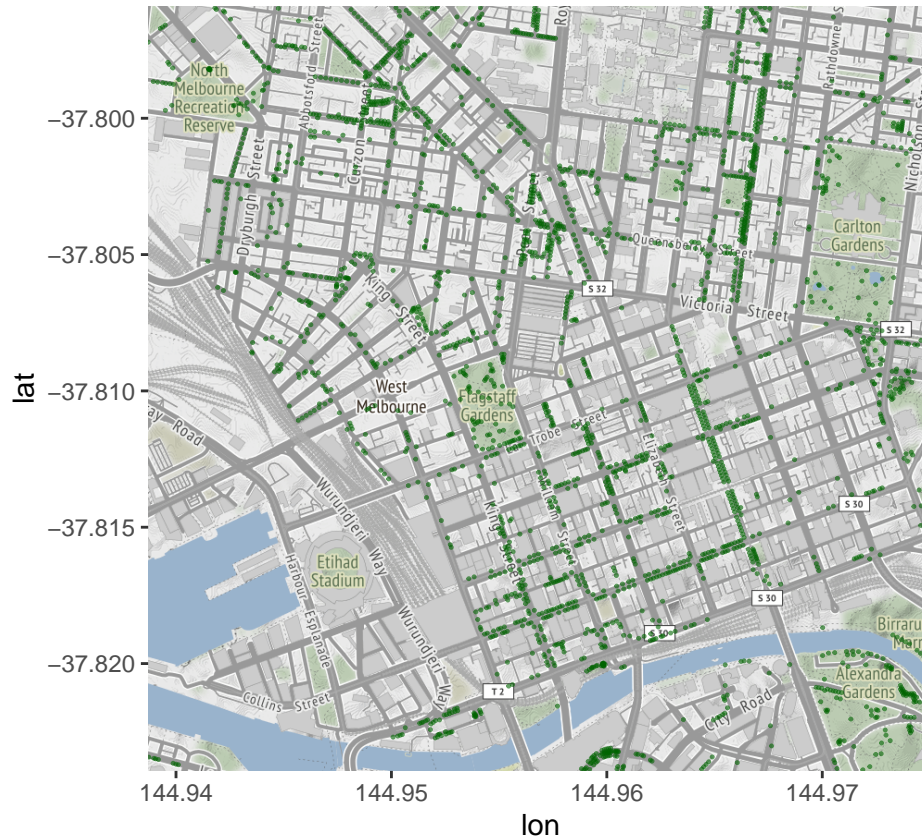
\end{table}

The number of rows in the cleaned data set are 6826 and the number of columns are 21.

Question 6: Create a map with the tree locations in the data set. (2pts)

```
# We have created the map below for you
melb_map <- read_rds(here::here("Data/melb-map.rds"))

# Here you just need to add the location for each tree into the map.
ggmap(melb_map) +
  geom_point(data = tree_data_clean1,
    aes(x = Longitude,
        y = Latitude),
    colour = "#006400",
    alpha = 0.6,
    size = 0.2)
```

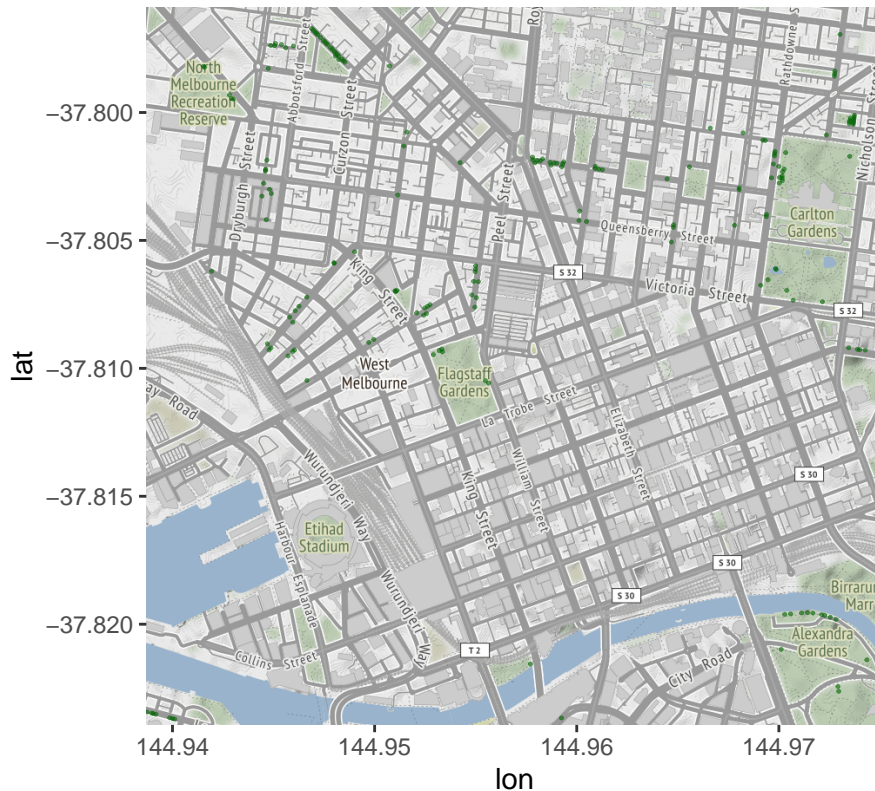


Question 7: Create another map and draw trees in the *Genus* groups of Eucalyptus, Macadamia, Prunus, Acacia, and Quercus. Use the “Dark2” color palette and display the legend at the bottom of the plot. (8pts)

```
selected_group <- tree_data_clean1 %>%
  filter(Genus %in% c("Eucalyptus", "Macadamia", "Prunus", "Acacia", "Quercus"))
```

```
ggmap(melb_map) +
  geom_point(data = selected_group,
    aes(x = Longitude,
        y = Latitude),
    colour = "#006400",
    alpha = 0.6,
    size = 0.2)+
  theme(legend.position = "bottom") +
  scale_color_brewer(palette = "Dark2") +
  labs(title = "Map of trees belonging to the selected genus group")
```


Map of trees belonging to the selected genus group



Question 8: Filter the data *tree_data_clean1* so that only the variables *Year*, *Located in*, and *Common Name* are displayed. Arrange the data set by *Year* in descending order and display the first 4 lines. Call this new data set *tree_data_clean_filter*. Then answer the following question using inline R code: When (*Year*), where (*Located in*) and what tree (*Common Name*) was the first tree planted in Melbourne according to this data set? (8pts)

```
tree_data_clean_filter <- tree_data_clean1 %>%
  select(Year, `Located in`, `Common Name`) %>%
  arrange(desc(Year))
head(tree_data_clean_filter, 4) %>%
  kable(caption = "Selected Variables of Tree Data") %>%
  kable_styling(latex_options = "hold_position")
```

```
tree_data_clean_filter_rename <- tree_data_clean1 %>%
  filter(`Dateplanted` == min(tree_data_clean1$Dateplanted)) %>%
  rename(location = `Located in`, common_name = `Common Name`)
```

The first tree was planted in 1900 at a Street and the tree name is London Plane. '

Table 10: Selected Variables of Tree Data

Year	Located in	Common Name
2000	Street	Small-leaved Linden
2000	Street	Spotted Gum
2000	Street	Drooping sheoak
2000	Park	Kanooka

2062 trees were planted on 1900-01-01, the first day when plantations began. These are located at both locations. The above answer is based on the last row of the above dataset.

Question 9: How many trees were planted in parks and how many in streets? Tabulate the results (only for locations in parks and streets) using the function *kable()* from the *kableExtra* R package. (3pts)

```
tree_data_clean1 %>%
  filter(`Located in` %in% c("Park", "Street")) %>%
  group_by(`Located in`) %>%
  summarise(`number of trees` = n()) %>%
  kableExtra::kable(caption = "Tree Count by Location") %>%
  kable_styling(latex_options = "hold_position")
```

Table 11: Tree Count by Location

Located in	number of trees
Park	2737
Street	4088

Question 10: How many trees are there in each of the Family groups in the data set *tree_data_clean1* (display the first 5 lines of the results in descending order)? (2pt)

```
tree_data_clean1 %>%
  group_by(`Family`) %>%
  summarise(`number of trees` = n()) %>%
  arrange(desc(`number of trees`)) %>%
  head(5) %>%
  kable(caption = "Tree Count by Family") %>%
  kable_styling(latex_options = "hold_position")
```

Table 12: Tree Count by Family

Family	number of trees
Myrtaceae	2102
Platanaceae	1512
Ulmaceae	1125
Fabaceae	327
Fagaceae	254

Question 11: Create a markdown table displaying the number of trees planted in each year (use variable *Yearplanted*) with common names Ironbark, Olive, Plum, Oak, and Elm (Hint: Use `kable()` from the `gridExtra` R package). What is the oldest most abundant tree in this group? (8pts)

Elm is the oldest most abundant tree in this group.

```
tree_data_clean1 %>%
  filter(`Common Name` %in% c("Ironbark", "Olive", "Plum", "Oak", "Elm")) %>%
  group_by(`Yearplanted`, `Common Name`) %>%
  summarise(`number of trees` = n()) %>%
  arrange(`Yearplanted`, desc(`number of trees`)) %>%
  kableExtra::kable(caption = "Tree Count by Year") %>%
  kable_styling(latex_options = "hold_position")
```

Table 13: Tree Count by Year

Yearplanted	Common Name	number of trees
1900	Elm	179
1900	Ironbark	29
1900	Olive	17
1900	Oak	4
2000	Ironbark	23
2000	Elm	18
2000	Oak	9

Question 12: Select the trees with diameters (Diameter Breast Height) greater than 40 cm and smaller 100 cm and comment on where the trees are located (streets or parks). (max 25 words) (3pts)

We see that, for the diameters 41 to 56, there are more trees planted on the streets than in parks. Larger trees are prevalent more in parks and their number reduces with diameter.

```

large_trees_data <- tree_data_clean1 %>%
  filter(`Diameter Breast Height` %in% c(41:99)) %>%
  group_by(`Located in`, `Diameter Breast Height`) %>%
  summarise(`number of trees` = n()) %>%
  ungroup() %>%
  pivot_wider(names_from = `Located in`,
              values_from = `number of trees`)

large_trees_data %>%
  kableExtra::kable(caption = "Tree Count by Diameter Breast Height and Location")

```

```

tree_data_clean1 %>%
  filter(`Diameter Breast Height` %in% c(41:99)) %>%
  group_by(`Located in`, `Diameter Breast Height`) %>%
  summarise(`number of trees` = n()) %>%
  ungroup() %>%
  ggplot() +
  geom_col(mapping = aes(x = `Diameter Breast Height`,
                        y = `number of trees`,
                        fill = `Located in`,
                        alpha = 0.5),
           position = "identity") +
  labs(title = "Comparing Tree Locations: Parks and Streets")

```

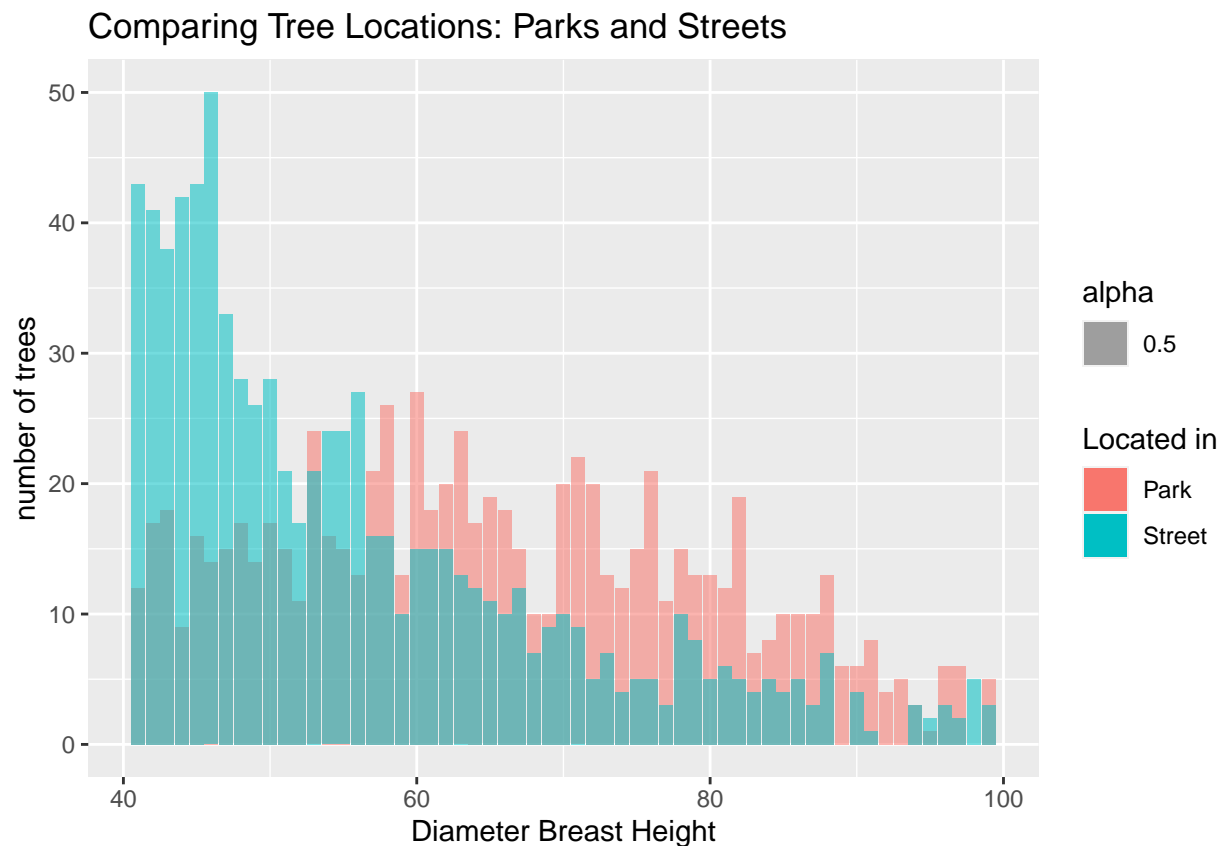


Table 14: Tree Count by Diameter Breast Height and Location

Diameter Breast Height	Park	Street
41	12	43
42	17	41
43	18	38
44	9	42
45	16	43
46	14	50
47	15	33
48	17	28
49	14	26
50	17	28
51	15	21
52	11	17
53	24	21
54	16	24
55	15	24
56	13	27
57	21	16
58	26	16
59	13	10
60	27	15
61	18	15
62	20	15
63	24	13
64	17	12
65	19	11
66	18	10
67	15	12
68	10	7
69	10	9
70	20	10
71	22	9
72	20	5
73	13	7
74	12	4
75	15	5
76	21	5
77	11	3
78	15	10
79	13	8
80	13	5
81	12	6
82	19	5
83	7	4
84	8	5
85	10	4
86	10	5
87	10	3
88	13	7
89	6	NA
90	6	4
91	8	1
1392	4	NA
93	5	NA
94	3	3
95	1	2

Question 13: Plot the trees within the diameter range that you have selected in Question 12, which are located in parks and streets on a map using 2 different colours to differentiate their locations (streets or parks). (6pts)

```
large_trees_data_parks <- tree_data_clean1 %>%
  filter(`Diameter Breast Height` %in% c(41:99))
```

```
ggmap(melb_map) +
  geom_point(data = large_trees_data_parks ,
            aes(x = Longitude,
                y = Latitude,
                colour = `Located in`),
            alpha = 0.6,
            size = 1.0) +
  theme(legend.position = "bottom") +
  scale_color_brewer(palette = "Dark2") +
  labs(title = "Map of Large Trees")
```



Question 14: Create a time series plot (using `geom_line`) that displays the total number of trees planted per year in the data set `tree_data_clean1` that belong to the Families: Myrtaceae, Arecaceae, and Ulmaceae. What do you observe from the plot?
(6pts)

We see that the number of trees that were planted decreases from 1900 to 2000. More trees belonging to Myrtaceae were planted with one tree uniquely planted in 1977.

```
Fig_data <- tree_data_clean1 %>%
  filter(`Family` %in% c("Myrtaceae", "Arecaceae", "Ulmaceae")) %>%
  group_by(`Yearplanted`, `Family`) %>%
  summarise(`number of trees` = n()) %>%
  arrange(desc(`number of trees`))
```

```
Fig_data %>%
  ggplot() +
  geom_line(mapping = aes(x = `Yearplanted`, y = `number of trees`, colour = `Family`)) +
  geom_point(mapping = aes(x = `Yearplanted`, y = `number of trees`, colour = `Family`))+
  theme(legend.position = "bottom") +
  theme_bw() +
  labs(title = "Year Planted vs Number of Trees")
```



Part 2: Simulation Exercise

Question 15: Create a data frame called *simulation_data* that contains 2 variables with names *response* and *covariate*. Generate the variables according to the following model:

$response = 3.5 \times covariate + \epsilon$ where *covariate* is a variable that takes values $0, 1, 2, \dots, 100$ and ϵ is generated according to a Normal distribution (Hint: Use the function *rnorm()* to generate *epsilon*.)
(3pts)

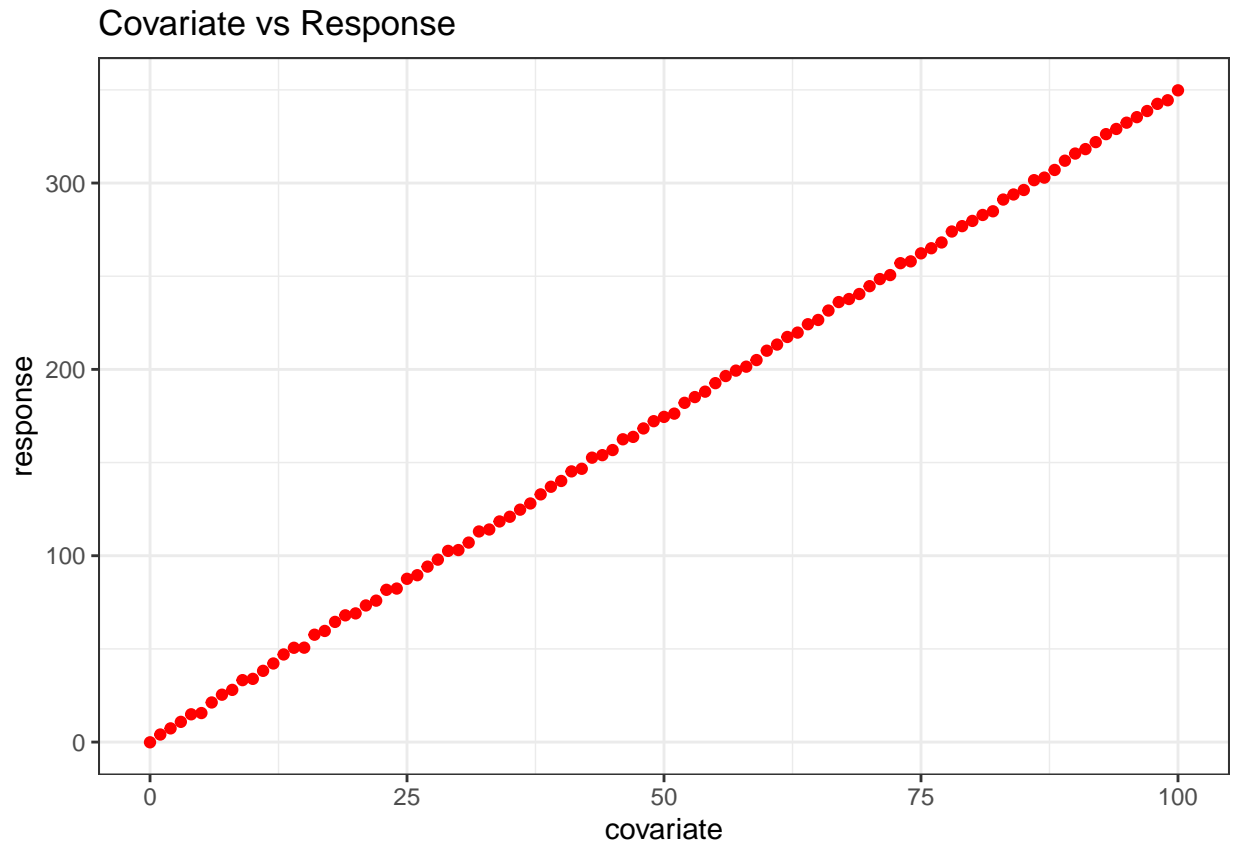
```
set.seed(2021)

simulation_data <- tibble(covariate = 0:100) %>%
  mutate(response = 3.5 * covariate + rnorm(101, 0, 1))
```

Question 16: Display graphically the relationship between the variables *response* and *covariate* (1pt) using a point plot. Which kind of relationship do you observe? (2pts)

We observe a linear relationship where the response variable increases with the covariate.

```
simulation_data %>%
  ggplot() +
  geom_point(mapping = aes(x = `covariate`,
                           y = `response`),
             colour = "red") +
  theme_bw() +
  labs(title = "Covariate vs Response")
```

Question 17: Fit a linear model between the variables *response* and *covariate* that you generate in Question 15 and display the model summary. (2pts)

```
simulation_data_lm <- lm(response~covariate, data=simulation_data)
summary(simulation_data_lm)
```

```
##
## Call:
## lm(formula = response ~ covariate, data = simulation_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.07431 -0.71466  0.05844  0.64196  2.25176
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.135896   0.199948    0.68   0.498
## covariate    3.493775   0.003455 1011.35 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 1.012 on 99 degrees of freedom
## Multiple R-squared:  0.9999, Adjusted R-squared:  0.9999
## F-statistic: 1.023e+06 on 1 and 99 DF,  p-value: < 2.2e-16
```

Question 18: What are the values for the intercept and the slope in the estimated model in Question 17 (Hint: Use the function `coef()`)? How do these values compare with the values in the simulation model? (max 50 words) (2pts)

```
#coef(summary(simulation_data_lm))
slope_intercept <- tidy(summary(simulation_data_lm)) %>%
  select(term, estimate)
```

The generated model has a slope of 3.49 and an intercept of 0.14

The simulation data was generated from the equation, $response = 3.5 \times covariate + \epsilon$ where ϵ is an error factor. The generated linear model is of the form $response = 3.4937754 \times covariate + 0.1358957$.

The value $3.49 \sim 3.5$ is the slope of the linear equation and the intercept of the model is 0.14. The fitted model differs from the simulation data in ϵ , which is centered around zero. The intercept of the model is close to zero.

```
#coef(summary(simulation_data_lm))
slope_intercept %>%
  kable(caption = "Slope and Intercept")%>%
  kable_styling(latex_options = "hold_position")
```

Table 15: Slope and Intercept

term	estimate
(Intercept)	0.1358957
covariate	3.4937754

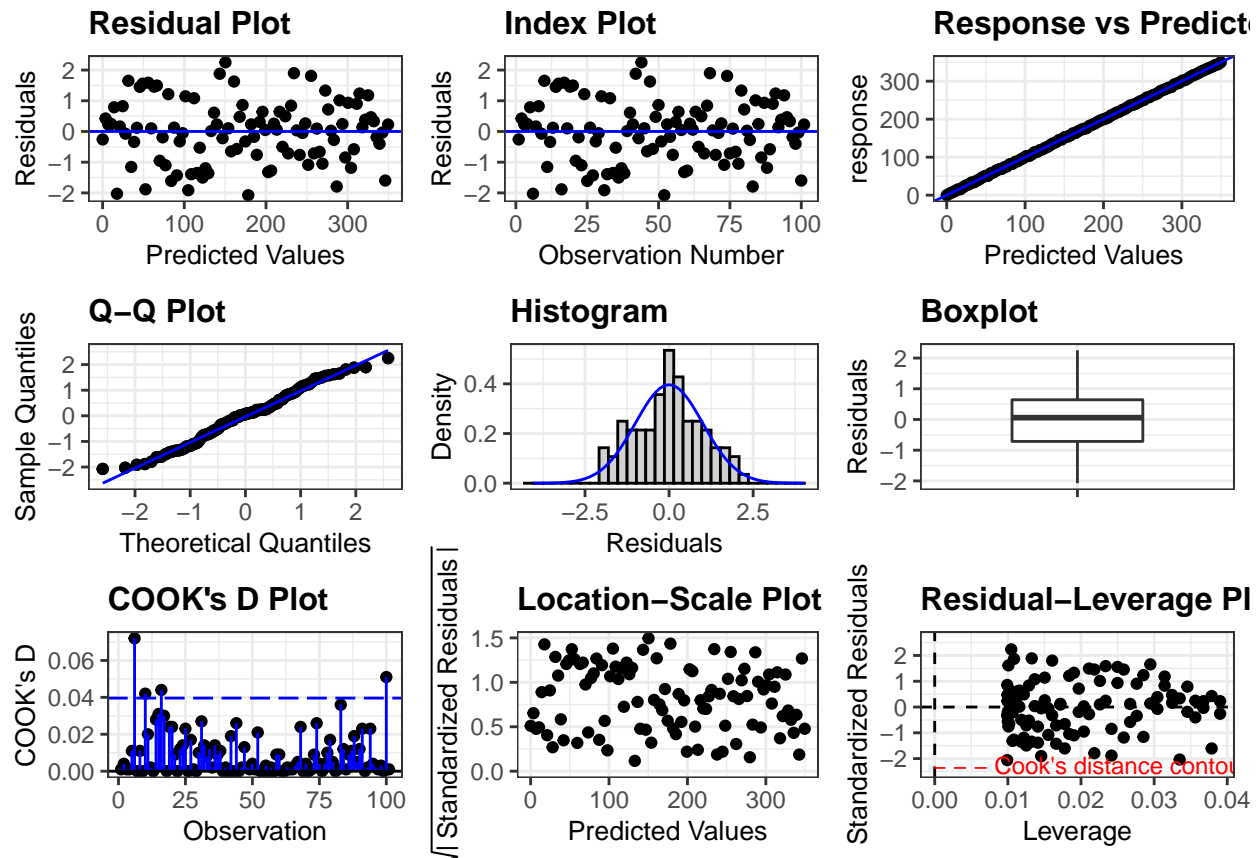
Question 19: Create a figure to display the diagnostic plots of the linear model that you fit in Question 17. Comment on the diagnostic plots (max 50 words). Is this a good/bad model and why? (max 30 words) (4pts)

- The Residual plot is a scatter plot of predicted values vs residuals. Residual is the difference between actual values and the predicted values. For a good model, the residual ~ 0 . The residual plot for a model having randomly dispersed points suggests that the model is good.
- The Index plot is similar to residual plot whereas it plots the observation number on the x axis.
- the Response vs Predicted plot is a scatter plot. A good model will have points aligned in straight line fashion showing that, predicted values \sim response.

- The plots in the second row show the distribution of the residuals. A good model has a normal distribution of residuals centered around 0.

The plots below show the goodness of fit of the model representing the simulation data. The residual plot has points scattered indefinitely, the response vs predicted plot is a straight line (slope = 1, response ~ predicted), showing that it is a well fitted model. The residuals lie within (-1,1) with a median of 0 suggesting goodness of the model.

```
resid_panel(simulation_data_lm, plots = "all")
```



```
augment(simulation_data_lm)
```

Question 20: Report R2, Radjusted, AIC, and BIC. Is this a good/bad model? Please explain your answer. (max 30 words) (2pts)

The model generated for the simulation data is a good model.

We know that the R squared value is a measure for the goodness of fit of a linear model which has values in the range [0, 1]. A good model has an R squared close to 1. The generated model has an r squared and r squared adjusted value of 0.9999 and, hence is a good model. The model with lowest AIC and BIC is a good model. For this model, the AIC and BIC are comparable and have low values. However, we do not have any other model for comparison and hence this model a good model as suggested by the r squared values.

```
glance(simulation_data_lm) %>%
  select(r.squared, adj.r.squared, AIC, BIC) %>%
  kable(caption = "Measures of Goodness of Fit")%>%
  kable_styling(latex_options = "hold_position")
```

Table 16: Measures of Goodness of Fit

r.squared	adj.r.squared	AIC	BIC
0.9999032	0.9999022	293.0547	300.9001