

AIT580 Final Project

Anil Kumar Murebina

05/12/2022

```
library(tidyverse)
library(readr)
library(ggplot2)
library(ggdist)
library(scales)

# Import data
udemy_dataset = read_csv("https://raw.githubusercontent.com/victorfogos/Text-multiclassification/-/PySpark/main/udemy_dataset.csv")

# Clean and transform data
udemy_dataset = subset(udemy_dataset, select = c(1,4,14))
colnames(udemy_dataset)

## [1] "course_id" "course_title" "is_paid"
## [4] "price" "num_subscribers" "num_reviews"
## [7] "num_lectures" "level" "content_duration"
## [10] "published_timestamp" "subject"

udemy_dataset$price[udemy_dataset$price == "free"] = 0
udemy_dataset$price = as.numeric(as.character(udemy_dataset$price))
na_values = udemy_dataset[is.na(udemy_dataset)] > 0
udemy_dataset = subset(udemy_dataset, course_id!=na_values$course_id )
sum(is.na(udemy_dataset))

## [1] 0

#set working directory
setwd("Users\\anil\\Desktop\\First Semester\\AIT580\\Final Project")
# Write data set to a new csv file
write_csv(udemy_dataset, "udemy_dataset.csv", row.names = FALSE)

# Import the final clean data set that has been done by Python
udemy_dataset = read_csv("clean_udemy_dataset.csv")
udemy_dataset = subset(udemy_dataset, select = c(12:16))
sum(is.na(udemy_dataset))

## [1] 0

summary(udemy_dataset)

## course_id course_title is_paid price
## Min. : 8324 Length:3677 Mode :logical Min. : 9.88
## 1st Qu.: 497098 Class :character FALSE:329 1st Qu.: 29.88
## Median : 688992 Mode :character TRUE :3367 Median : 45.88
## Mean : 675965 Mean : 156.3 Mean : 49.12
## 3rd Qu.: 963509 3rd Qu.: 67.8 3rd Qu.: 46.99
## Max. : 1282064 3rd Qu.: 27445.8 Max. : 209.88
## num_subscribers num_reviews num_lectures level
## Min. : 0 Min. : 0.8 Min. : 4.88 Length:3677
## 1st Qu.: 131 1st Qu.: 4.8 1st Qu.: 15.09 Class :character
## Median : 912 Median : 29.8 Median : 25.99 Mode :character
## Mean : 3198 Mean : 156.3 Mean : 49.12
## 3rd Qu.: 2547 3rd Qu.: 67.8 3rd Qu.: 46.99
## Max. : 268923 Max. : 27445.8 Max. : 279.09
## content_duration published_timestamp subject
## Min. : 6.3333 Min. : 2021-07-09 05:43:31 Length:3677
## 1st Qu.: 5.0000 1st Qu.: 2025-03-16 18:36:17 Class :character
## Median : 5.0000 Median : 2025-01-27 18:10:28 Mode :character
## Mean : 4.0006 Mean : 2025-11-26 14:49:11
## 3rd Qu.: 4.5000 3rd Qu.: 2026-18-29 21:56:38
## Max. : 78.5000 Max. : 2027-07-06 21:40:30
## category
## Min. : 1.0000
## 1st Qu.:1.0000
## Median :2.0000
## Mean :2.402
## 3rd Qu.:4.0000
## Max. :4.0000

# Write data set to a new csv file
write_csv(udemy_dataset, "clean_udemy_dataset.csv", row.names = FALSE)

# Average price of the courses for all the categories?
Business_Finance = subset(udemy_dataset, subject == "Business Finance")
Graphic_Design = subset(udemy_dataset, subject == "Graphic Design")
Musical_Instruments = subset(udemy_dataset, subject == "Musical Instruments")
Web_Development = subset(udemy_dataset, subject == "Web Development")

summary(Business_Finance$price)

## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.00 29.00 45.00 68.59 95.00 280.00

summary(Graphic_Design$price)

## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.00 29.00 39.00 57.83 89.00 280.00

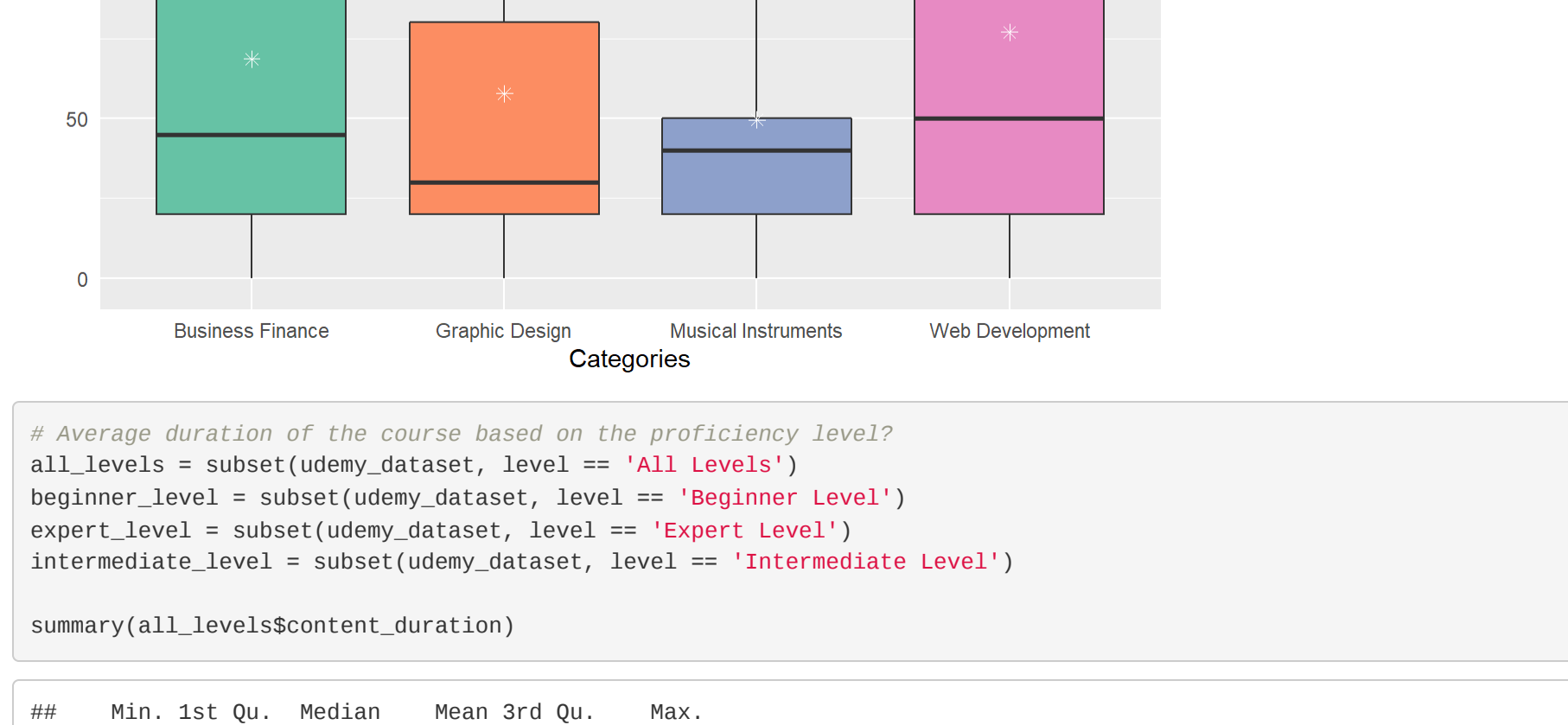
summary(Musical_Instruments$price)

## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.00 29.00 40.00 49.56 59.00 280.00

summary(Web_Development$price)

## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.00 29.00 50.00 77.05 115.00 280.00

ggplot(udemy_dataset, aes(x = subject, y = price, fill = subject)) +
  geom_boxplot() +
  stat_summary(fun = "mean", geom = "point", shape = 8, size = 2, color = "white") +
  labs(x="Categories", y = "Price", title = "Statistics of Price based on Category") +
  scale_fill_brewer(palette = "Set2") +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5),
        axis.ticks = element_blank())
```



Average duration of the course based on the proficiency level?

```
all_levels = subset(udemy_dataset, level == "All Levels")
beginner_level = subset(udemy_dataset, level == "Beginner Level")
expert_level = subset(udemy_dataset, level == "Expert Level")
intermediate_level = subset(udemy_dataset, level == "Intermediate Level")

summary(all_levels$content_duration)
```

```
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.1333 1.5000 2.5000 4.8722 5.5000 76.5000
```

```
summary(beginner_level$content_duration)
```

```
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.450 1.000 2.000 3.091 3.500 78.500
```

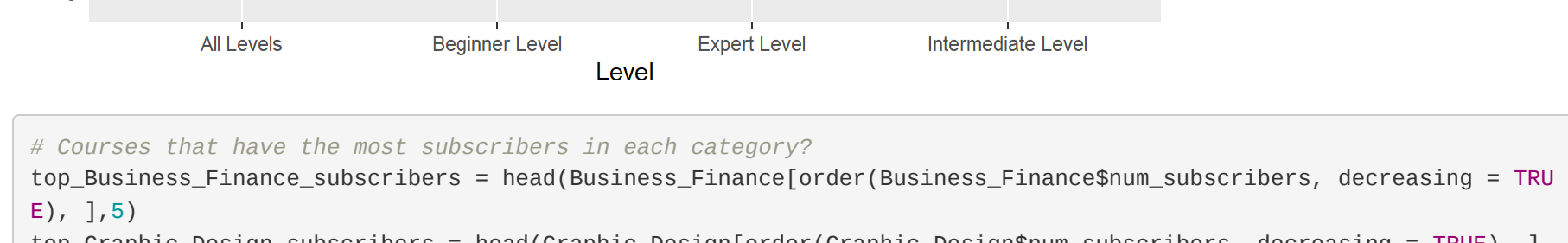
```
summary(expert_level$content_duration)
```

```
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.5167 1.0000 2.0000 2.9055 3.8750 12.5000
```

```
summary(intermediate_level$content_duration)
```

```
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.500 1.000 2.000 2.735 4.500 31.500
```

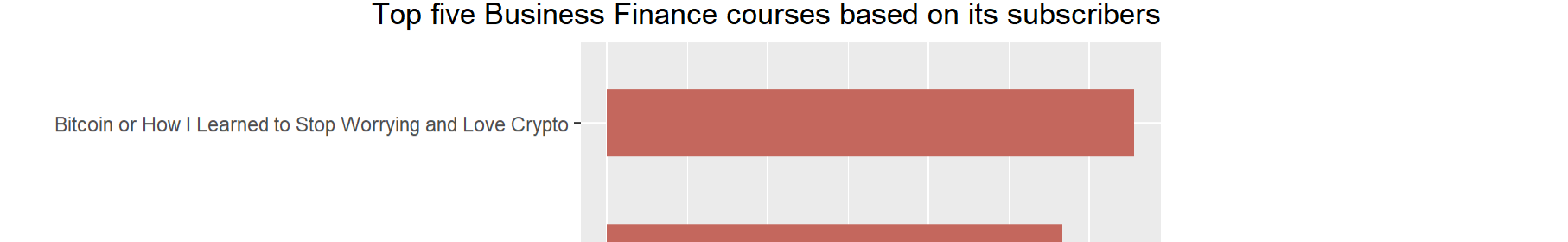
```
ggplot(udemy_dataset, aes(x = level, y = content_duration, fill = level)) +
  geom_boxplot() +
  stat_summary(fun = "mean", geom = "point", shape = 8, size = 2, color = "white") +
  labs(x="level", y = "Duration", title = "Statistics of duration based on level") +
  scale_fill_brewer(palette = "Set2") +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```



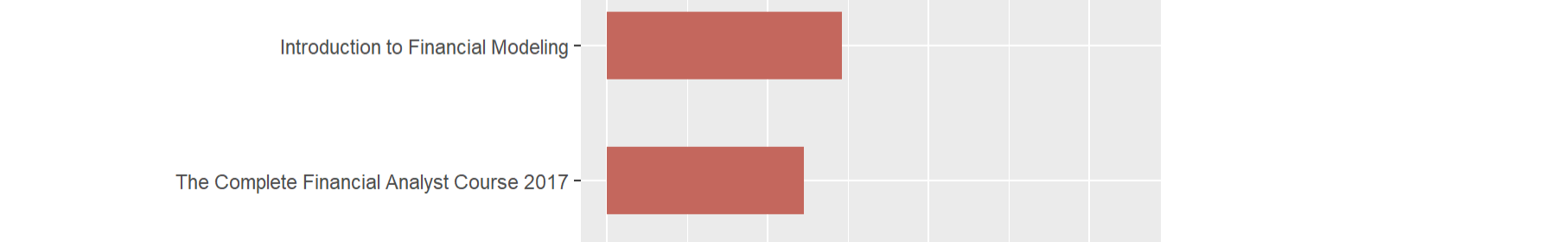
Courses that have the most subscribers in each category?

```
top_Business_Finance_subscribers = head(Business_Finance[order(Business_Finance$num_subscribers, decreasing = TRUE), ], 5)
top_Graphic_Design_subscribers = head(Graphic_Design[order(Graphic_Design$num_subscribers, decreasing = TRUE), ], 5)
top_Musical_Instruments_subscribers = head(Musical_Instruments[order(Musical_Instruments$num_subscribers, decreasing = TRUE), ], 5)
top_Web_Development_subscribers = head(Web_Development[order(Web_Development$num_subscribers, decreasing = TRUE), ], 5)
```

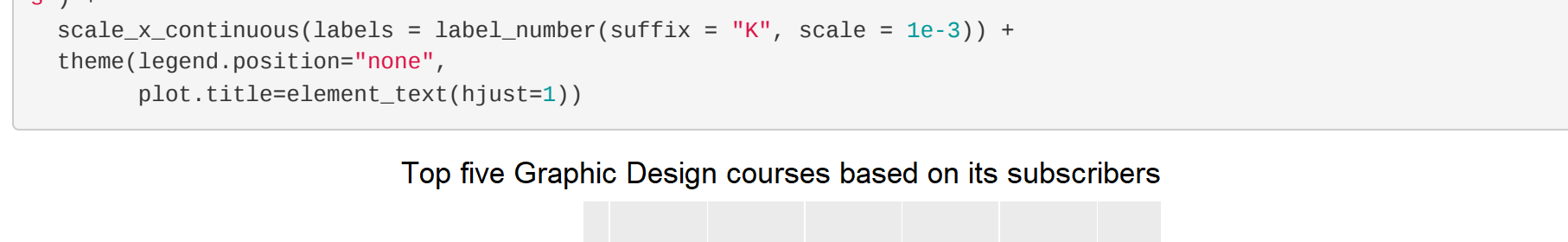
```
ggplot(top_Business_Finance_subscribers, aes(x = num_subscribers, y = reorder(course_title, num_subscribers))) +
  geom_bar(stat = "identity", width = 0.5, fill = "#C67870") +
  labs(x="Number of subscribers", y = "Courses", title = "Top five Business Finance courses based on its subscribers") +
  scale_x_continuous(labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```



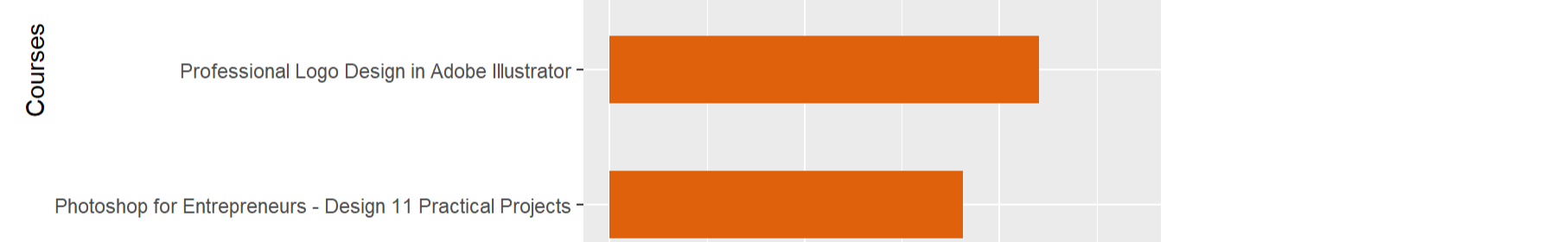
```
ggplot(top_Graphic_Design_subscribers, aes(x = num_subscribers, y = reorder(course_title, num_subscribers))) +
  geom_bar(stat = "identity", width = 0.5, fill = "#00A0A0") +
  labs(x="Number of subscribers", y = "Courses", title = "Top five Graphic Design courses based on its subscriber") +
  scale_x_continuous(labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```



```
ggplot(top_Musical_Instruments_subscribers, aes(x = num_subscribers, y = reorder(course_title, num_subscribers))) +
  geom_bar(stat = "identity", width = 0.5, fill = "#4682B4") +
  labs(x="Number of subscribers", y = "Courses", title = "Top five Musical Instruments courses based on its subscribers") +
  scale_x_continuous(labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```



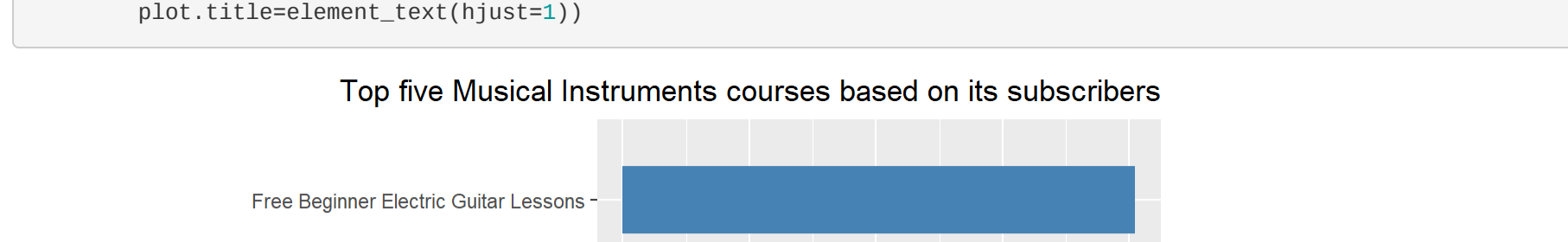
```
ggplot(top_Web_Development_subscribers, aes(x = num_subscribers, y = reorder(course_title, num_subscribers))) +
  geom_bar(stat = "identity", width = 0.5, fill = "#4682B4") +
  labs(x="Number of subscribers", y = "Courses", title = "Top five Web Development courses based on its subscribers") +
  scale_x_continuous(labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```



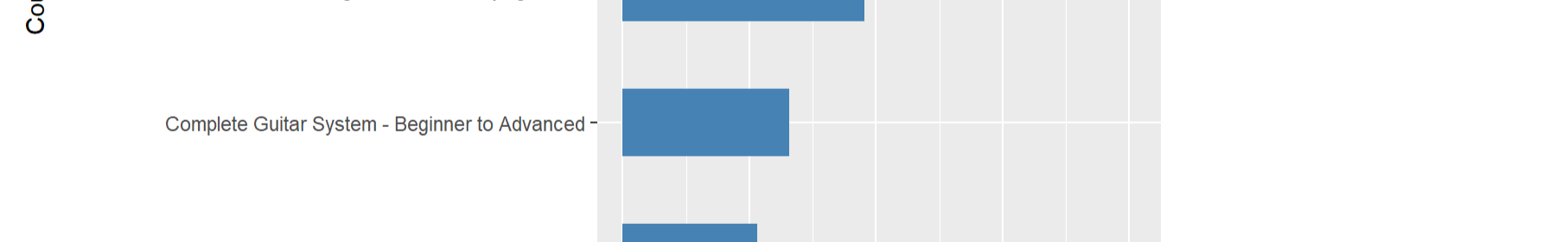
Courses that have the most reviews in each category?

```
top_Business_Finance_reviews = head(Business_Finance[order(Business_Finance$num_reviews, decreasing = TRUE), ], 5)
top_Graphic_Design_reviews = head(Graphic_Design[order(Graphic_Design$num_reviews, decreasing = TRUE), ], 5)
top_Musical_Instruments_reviews = head(Musical_Instruments[order(Musical_Instruments$num_reviews, decreasing = TRUE), ], 5)
top_Web_Development_reviews = head(Web_Development[order(Web_Development$num_reviews, decreasing = TRUE), ], 5)
```

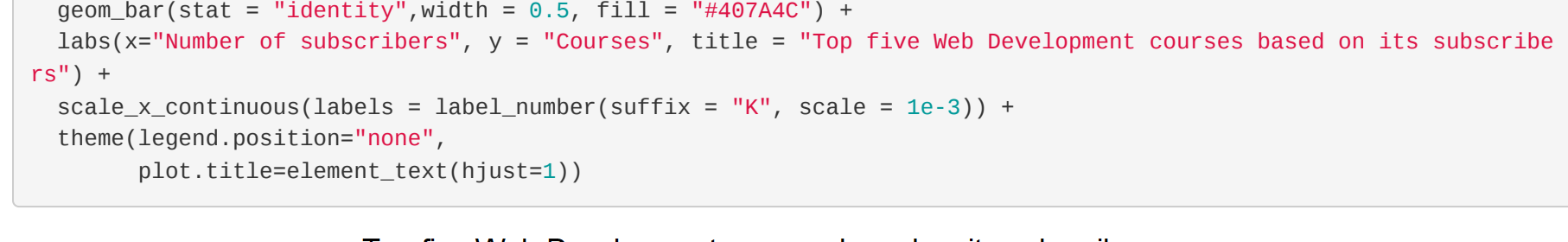
```
ggplot(top_Business_Finance_reviews, aes(x = num_reviews, y = reorder(course_title, num_reviews))) +
  geom_bar(stat = "identity", width = 0.5, fill = "#C67870") +
  labs(x="Number of reviews", y = "Courses", title = "Top five Business Finance courses based on its reviews") +
  scale_x_continuous(labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```



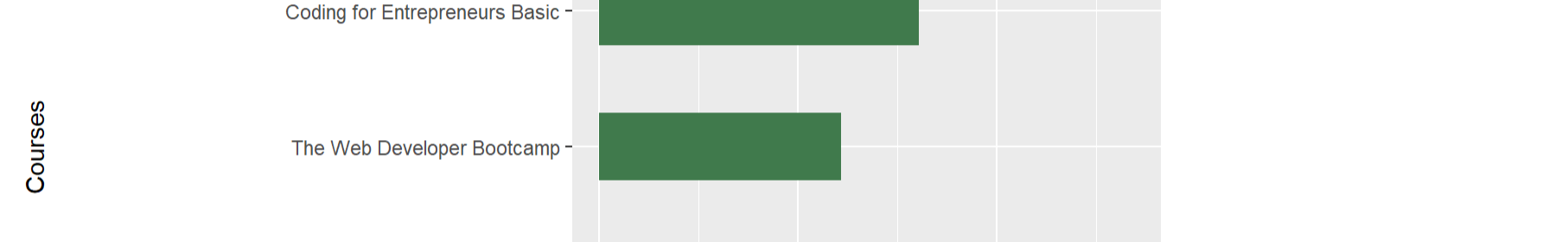
```
ggplot(top_Graphic_Design_reviews, aes(x = num_reviews, y = reorder(course_title, num_reviews))) +
  geom_bar(stat = "identity", width = 0.5, fill = "#00A0A0") +
  labs(x="Number of reviews", y = "Courses", title = "Top five Graphic Design courses based on its reviews") +
  scale_x_continuous(labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```



```
ggplot(top_Musical_Instruments_reviews, aes(x = num_reviews, y = reorder(course_title, num_reviews))) +
  geom_bar(stat = "identity", width = 0.5, fill = "#4682B4") +
  labs(x="Number of reviews", y = "Courses", title = "Top five Musical Instruments courses based on its reviews") +
  scale_x_continuous(labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```

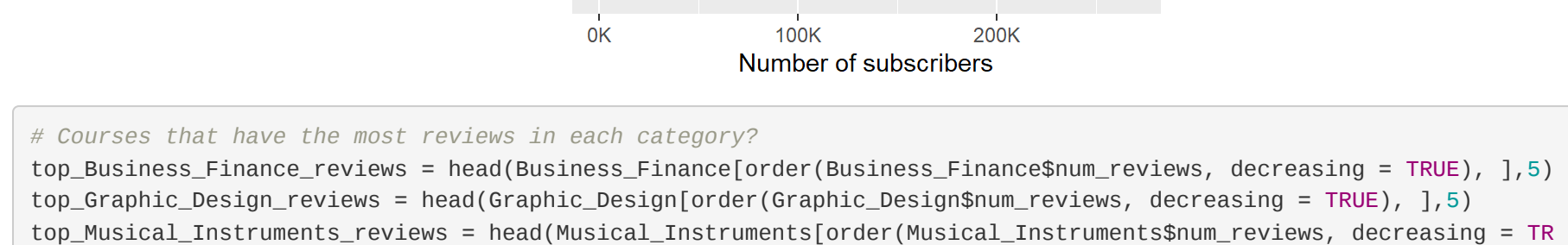


```
ggplot(top_Web_Development_reviews, aes(x = num_reviews, y = reorder(course_title, num_reviews))) +
  geom_bar(stat = "identity", width = 0.5, fill = "#4682B4") +
  labs(x="Number of reviews", y = "Courses", title = "Top five Web Development courses based on its reviews") +
  scale_x_continuous(labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5))
```



Level of proficiency count for all the courses in each category

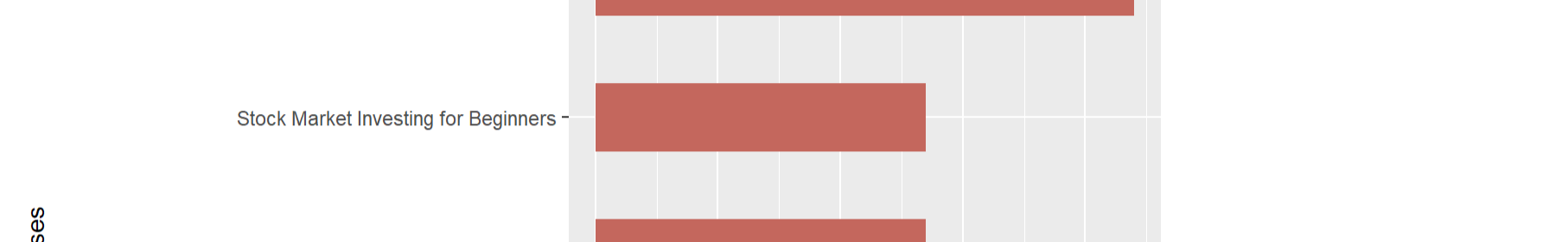
```
ggplot(udemy_dataset, aes(x = subject, fill = level)) +
  geom_bar(position = position_dodge(), width = 0.8) +
  labs(x="Categories", y = "Number of levels", title = "Number of levels in all categories") +
  scale_fill_brewer(palette = "Set2", direction = -1) +
  theme(legend.position = "bottom",
        legend.position = "bottom",
        plot.title=element_text(hjust=0.5),
        axis.ticks = element_blank())
```



Most popular courses and when they are published

```
popular_courses = head(udemy_dataset[order(udemy_dataset$num_subscribers, decreasing = TRUE), ], 5)
```

```
ggplot(popular_courses, aes(x = num_subscribers, y = reorder(course_title, num_subscribers), fill = course_title)) +
  geom_bar(stat = "identity", width = 0.5, fill = "#800080") +
  labs(x="Number of subscribers", y = "Courses", title = "Most Popular courses") +
  scale_x_continuous(breaks = seq(0, 30000, 5000), labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5),
        axis.ticks = element_blank())
```



Most engaging courses

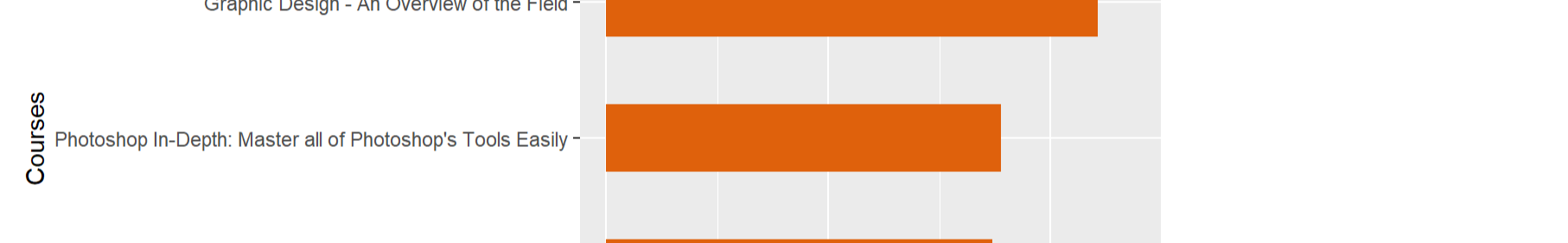
```
engaging_courses = head(udemy_dataset[order(udemy_dataset$num_reviews, decreasing = TRUE), ], 10)
```

```
ggplot(engaging_courses, aes(x = num_reviews, y = reorder(course_title, num_reviews), fill = course_title)) +
  geom_bar(stat = "identity", fill = "#00A0A0") +
  labs(x="Number of Reviews", y = "Courses", title = "Most Engaging courses") +
  scale_x_continuous(breaks = seq(0, 30000, 5000), labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5),
        axis.ticks = element_blank())
```



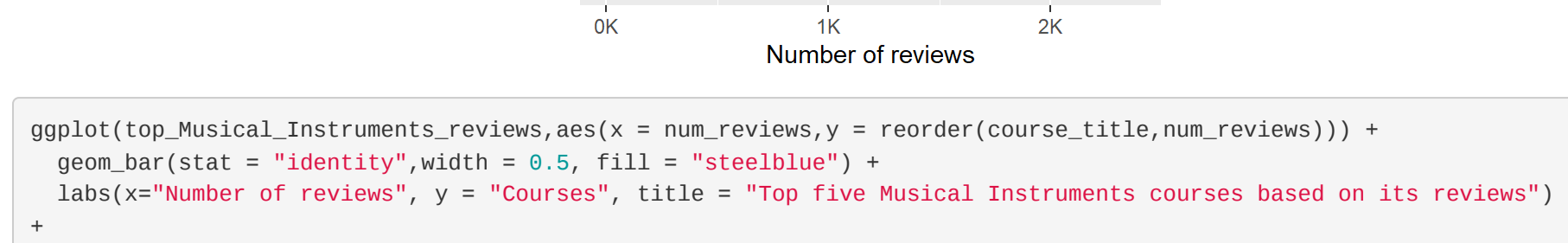
Number of courses per category

```
ggplot(udemy_dataset, aes(x = subject, y = frequency(subject), fill = subject)) +
  geom_bar(stat = "identity", width = 0.5, fill = "#800080") +
  labs(x="level", y = "Number of subscribers") +
  scale_x_continuous(breaks = seq(0, 30000, 5000), labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5),
        axis.ticks = element_blank())
```



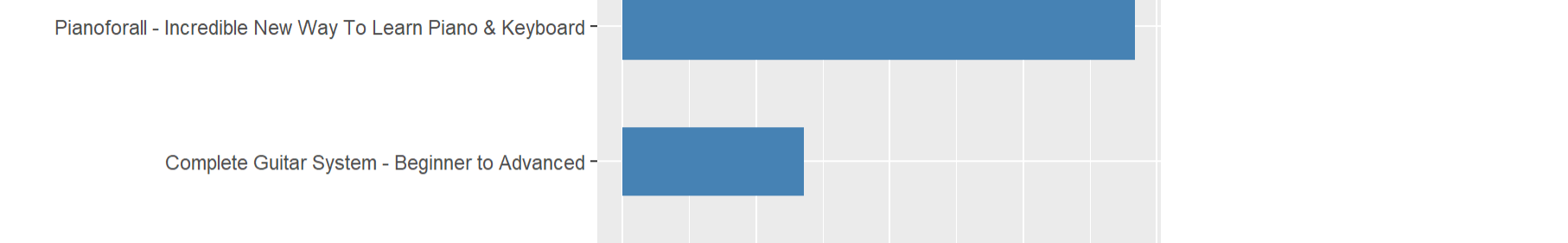
Number of subscribers per level

```
subscribers_per_level = ggplot(udemy_dataset, aes(x = reorder(level, num_subscribers), y = num_subscribers, fill = level)) +
  geom_bar(stat = "identity") +
  scale_fill_brewer(palette = "Set2") +
  labs(x="level", y = "Number of subscribers") +
  scale_x_continuous(breaks = seq(0, 30000, 5000), labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5),
        axis.ticks = element_blank())
```



Number of courses per level

```
courses_per_level = ggplot(udemy_dataset, aes(x = reorder(level, num_subscribers), y = frequency(course_id), fill = level)) +
  geom_bar(stat = "identity") +
  scale_fill_brewer(palette = "Set2") +
  labs(x="level", y = "Number of subscribers") +
  scale_x_continuous(breaks = seq(0, 30000, 5000), labels = label_number(suffix = "K", scale = 1e-3)) +
  theme(legend.position="none",
        plot.title=element_text(hjust=0.5),
        axis.ticks = element_blank())
```



Number of courses and their subscribers per level

```
grid.arrange(courses_per_level, subscribers_per_level, top="Number of courses and their subscribers per level")
```

