

# Lecture Assignment 9

Taiki Yamashita

2024-05-02

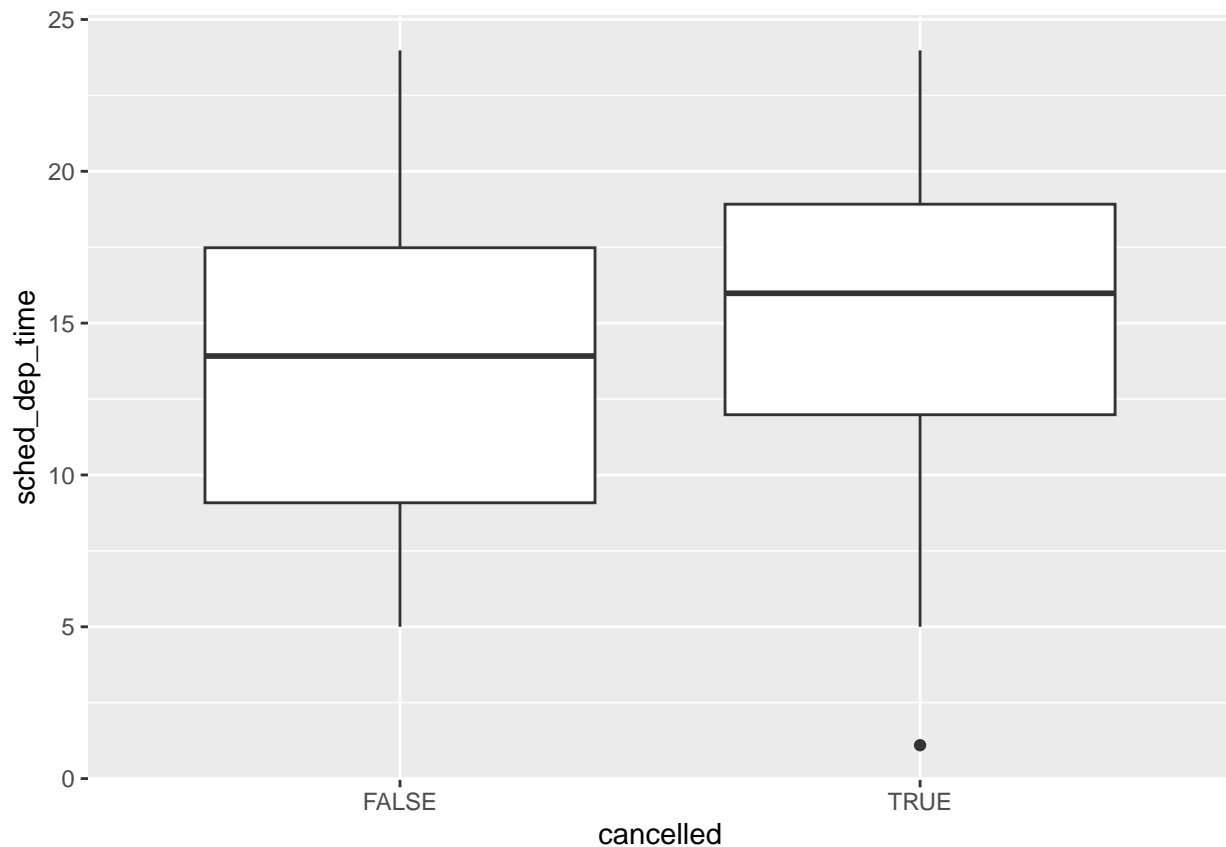
```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.5
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2    3.5.0      v tibble    3.2.1
## v lubridate  1.9.3      v tidyr     1.3.1
## v purrr      1.0.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
##
## Attaching package: 'ggstance'
##
##
## The following objects are masked from 'package:ggplot2':
##
##   geom_errorbarh, GeomErrorbarh
```

## 7.5.1.1 1)

Use what you've learned to improve the visualisation of the departure times of cancelled vs. non-cancelled flights.

What we can do is instead of using a freqplot is that we can now use a box plot!

```
nycflights13::flights %>%
  mutate(
    cancelled = is.na(dep_time),
    sched_hour = sched_dep_time %/% 100,
    sched_min = sched_dep_time %% 100,
    sched_dep_time = sched_hour + sched_min / 60
  ) %>%
  ggplot() +
  geom_boxplot(mapping = aes(y = sched_dep_time, x = cancelled))
```

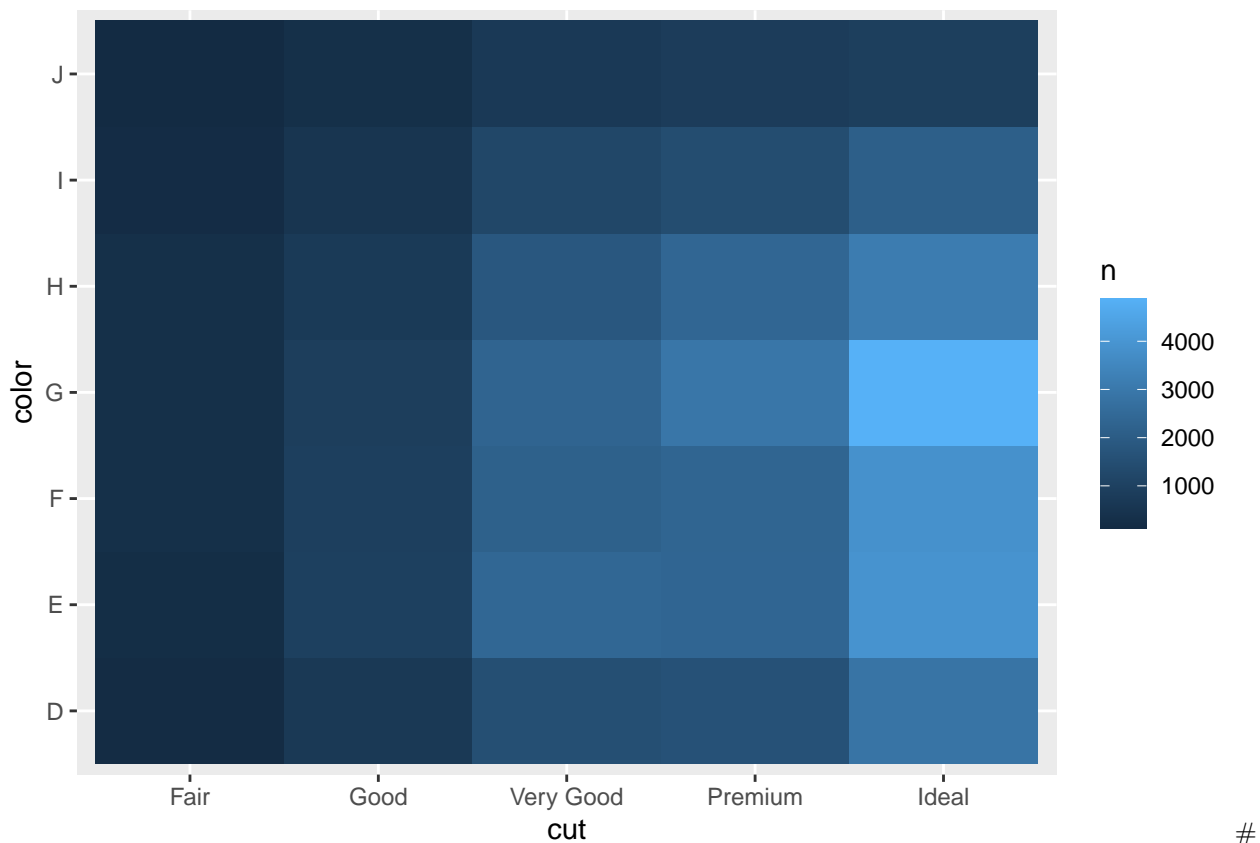


### 7.5.2.1 3)

Why is it slightly better to use `aes(x = color, y = cut)` rather than `aes(x = cut, y = color)` in the example above?

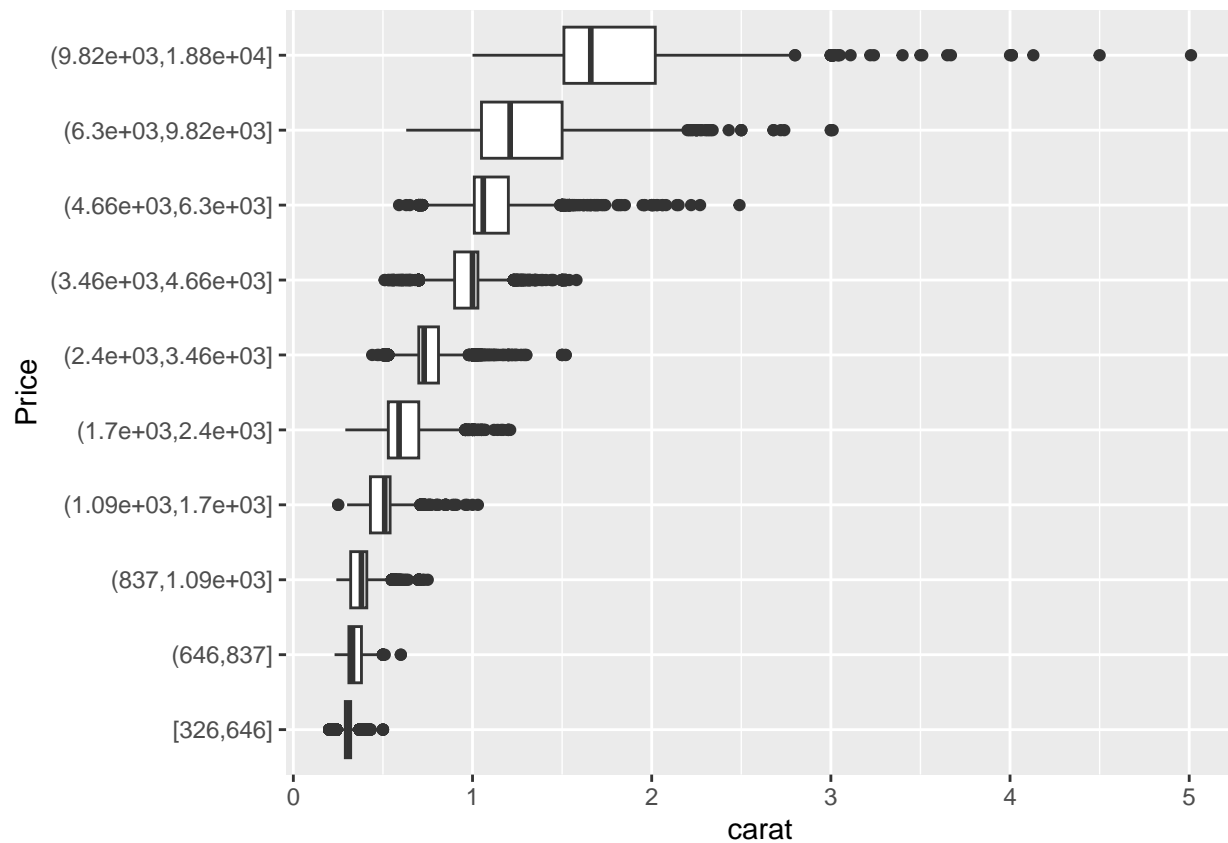
Usually it's better to use the categorical variable with a larger number of categories or the longer labels on the y axis. But, switching the order will not result in overlapping labels. Labels should be horizontal because it is easier to read.

```
diamonds %>%
  count(color, cut) %>%
  ggplot(mapping = aes(y = color, x = cut)) +
  geom_tile(mapping = aes(fill = n))
```



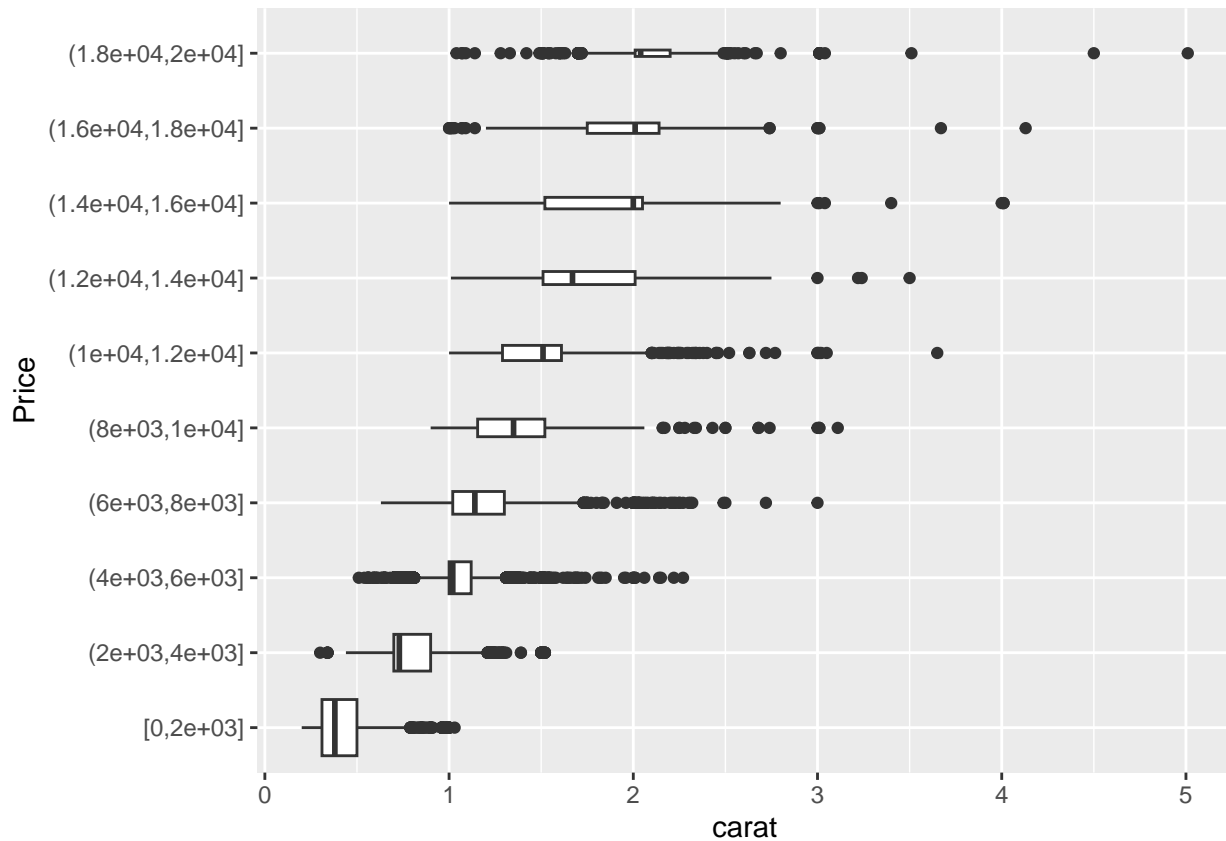
7.5.3.1 2) # Visualise the distribution of carat, partitioned by price. A graph of a box plot with 10 bins an equal number of observations. The width is determined by the number of observations.

```
ggplot(diamonds, aes(x = cut_number(price, 10), y = carat)) +
  geom_boxplot() +
  coord_flip() +
  xlab("Price")
```



Another visualization would be a box plot with 10 equal-width bins of 2000 dollars. `boundary = 0` is what ensures that the first bin is 0 to 2000 dollars.

```
ggplot(diamonds, aes(x = cut_width(price, 2000, boundary = 0), y = carat)) +
  geom_boxplot(varwidth = TRUE) +
  coord_flip() +
  xlab("Price")
```

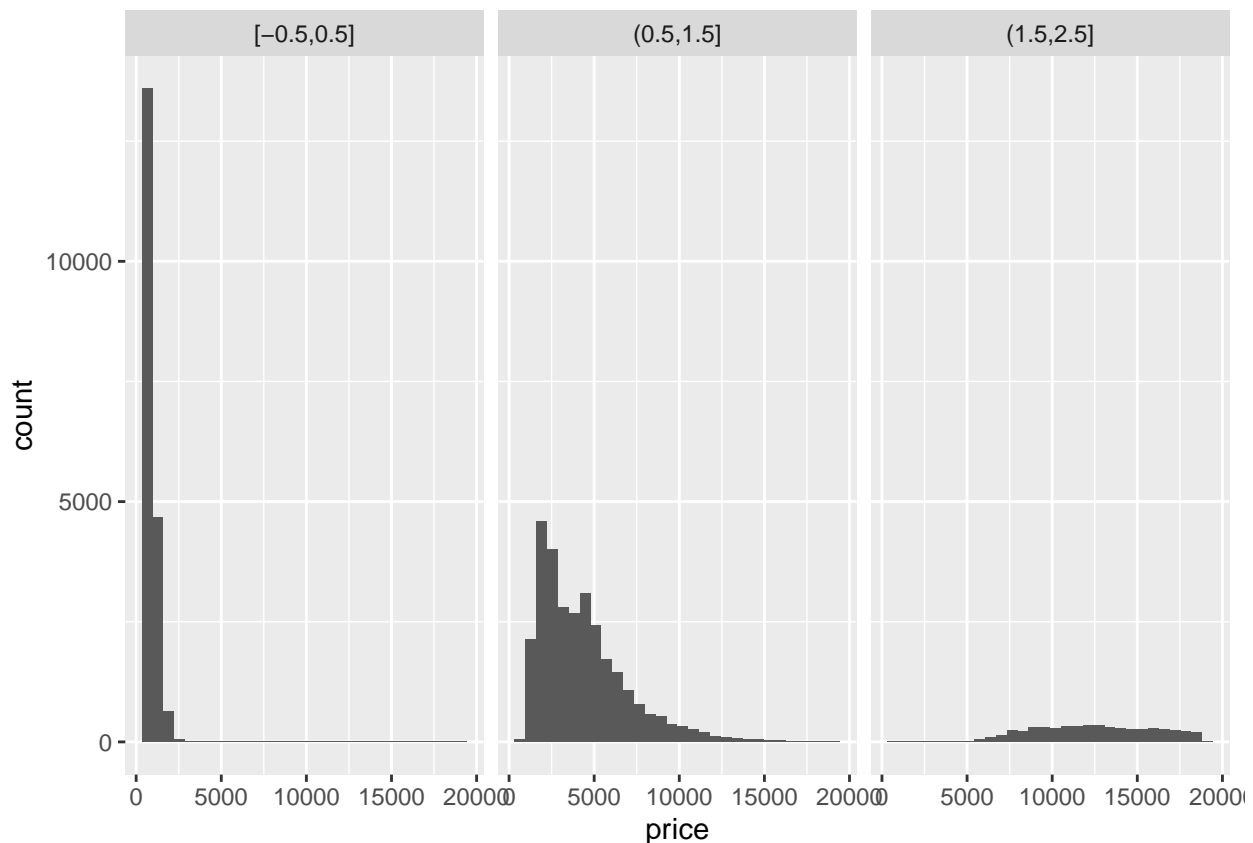


### 7.5.3.1 3)

How does the price distribution of very large diamonds compare to small diamonds? Is it as you expect, or does it surprise you?

```
diamonds %>%
  filter(between(carat, 0, 2.5)) %>%
  mutate(carat = cut_width(carat, 1)) %>%
  ggplot(aes(price)) +
  geom_histogram() +
  facet_wrap(~ carat)
```

## 'stat\_bin()' using 'bins = 30'. Pick better value with 'binwidth'.

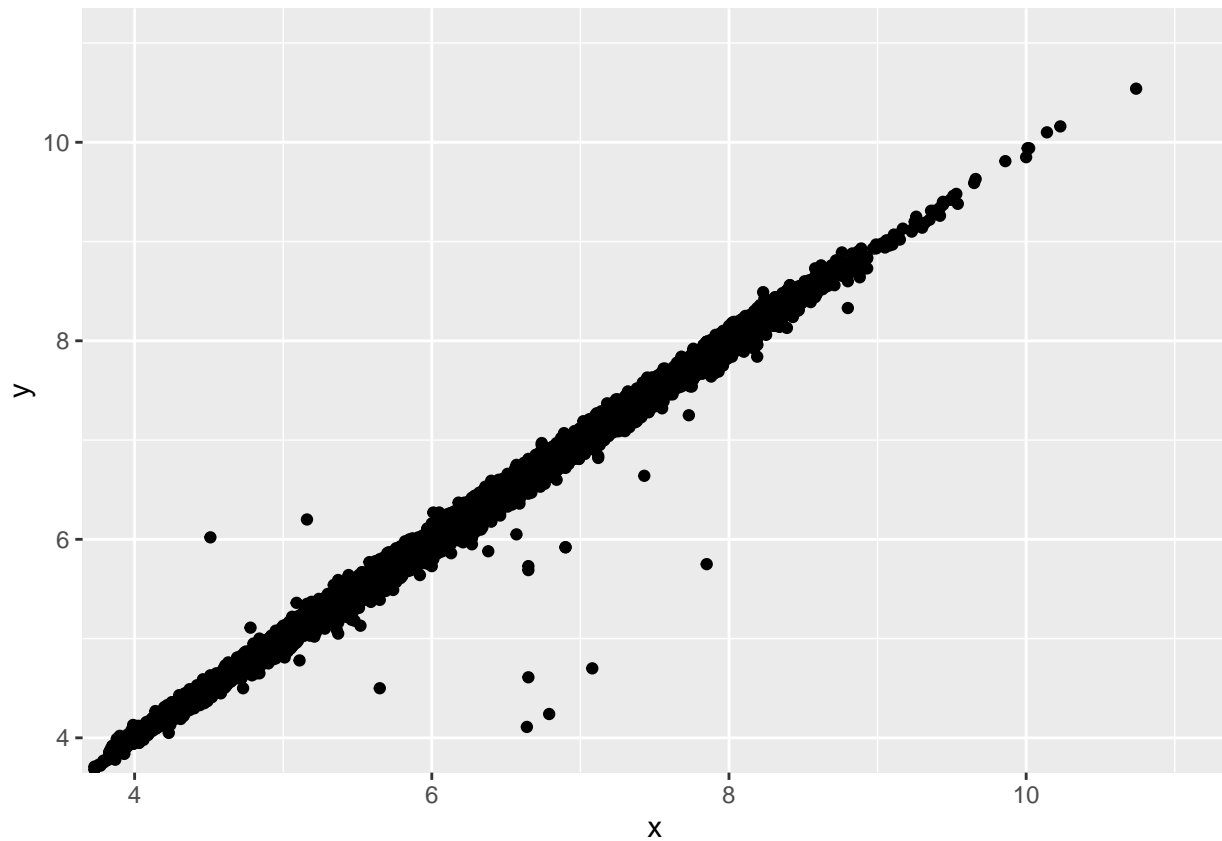


This was very surprising as I was expecting very little variance. It seems as though big diamonds can cost anything between 5000 and 18000. Whereas small ones have very little variance.

### 7.5.3.1 5)

Two dimensional plots reveal outliers that are not visible in one dimensional plots. For example, some points in the plot below have an unusual combination of x and y values, which makes the points outliers even though their x and y values appear normal when examined separately.

```
ggplot(data = diamonds) +
  geom_point(mapping = aes(x=x, y=y)) +
  coord_cartesian(xlim = c(4, 11), ylim = c(4, 11))
```



There is a strong relationship between  $x$  and  $y$ . The outliers are not as extreme in either  $x$  or  $y$ . A binned plot would not reveal these outliers, and may lead us to conclude that the largest value of  $x$  was an outlier even though it appears to fit the bivariate pattern well.