# CWRU DSCI351-351M-451: Week04a(CWRU, Pitt, UCF, UTRGV)

Profs: R. H. French, L. S. Bruckman, P. Leu, K. Davis, S. Cirlos

TAs: W. Oltjen, K. Hernandez, M. Li, M. Li, D. Colvin

13 September, 2022

## Contents

## 4.1.2.1   Some simple settings suggestions for windows or mac

- In Windows File Manager, make sure (show file extensions) is checked
- On a Mac, no Insert Key, so Shift Insert doesn't work
    - Use two-finger tap on trackpad, to get popup menu
    - Then choose paste

## 4.1.2.2   An example of EDA with pipes

- Library in the packages we will use

```
library(ggplot2) # load package

? ggplot2

library(dplyr)  #load package

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
? dplyr
```

## 4.1.2.2.1   Diamonds dataset

- Do load the Diamonds dataset and do some quick EDA

```r
data(diamonds)  # load dataset that comes with ggplot

head(diamonds) # Return the First or Last Part of an Object
```

```
## # A tibble: 6 x 10
##   carat cut       color clarity depth table price     x     y     z
##   <dbl> <ord>     <ord> <ord>   <dbl> <dbl> <int> <dbl> <dbl> <dbl>
## 1  0.23 Ideal     E     SI2      61.5    55   326  3.95  3.98  2.43
## 2  0.21 Premium   E     SI1      59.8    61   326  3.89  3.84  2.31
## 3  0.23 Good      E     VS1      56.9    65   327  4.05  4.07  2.31
## 4  0.29 Premium   I     VS2      62.4    58   334  4.2   4.23  2.63
## 5  0.31 Good      J     SI2      63.3    58   335  4.34  4.35  2.75
## 6  0.24 Very Good J     VVS2     62.8    57   336  3.94  3.96  2.48
```

```r
str(diamonds) # Compactly Display the Structure of an Arbitrary R Object
```

```
## tibble [53,940 x 10] (S3: tbl_df/tbl/data.frame)
##  $ carat  : num [1:53940] 0.23 0.21 0.23 0.29 0.31 0.24 0.24 0.26 0.22 0.23 ...
##  $ cut    : Ord.factor w/ 5 levels "Fair"<"Good"<..: 5 4 2 4 2 3 3 3 1 3 ...
##  $ color  : Ord.factor w/ 7 levels "D"<"E"<"F"<"G"<..: 2 2 2 6 7 7 6 5 2 5 ...
##  $ clarity: Ord.factor w/ 8 levels "I1"<"SI2"<"SI1"<..: 2 3 5 4 2 6 7 3 4 5 ...
##  $ depth  : num [1:53940] 61.5 59.8 56.9 62.4 63.3 62.8 62.3 61.9 65.1 59.4 ...
##  $ table  : num [1:53940] 55 61 65 58 58 57 57 55 61 61 ...
##  $ price  : int [1:53940] 326 326 327 334 335 336 336 337 337 338 ...
##  $ x      : num [1:53940] 3.95 3.89 4.05 4.2 4.34 3.94 3.95 4.07 3.87 4 ...
##  $ y      : num [1:53940] 3.98 3.84 4.07 4.23 4.35 3.96 3.98 4.11 3.78 4.05 ...
##  $ z      : num [1:53940] 2.43 2.31 2.31 2.63 2.75 2.48 2.47 2.53 2.49 2.39 ...
```

```r
glimpse(diamonds)  # the tidyverse version of str
```

```
## Rows: 53,940
## Columns: 10
## $ carat   <dbl> 0.23, 0.21, 0.23, 0.29, 0.31, 0.24, 0.24, 0.26, 0.22, 0.23, 0.~
## $ cut     <ord> Ideal, Premium, Good, Premium, Good, Very Good, Very Good, Ver~
## $ color   <ord> E, E, E, I, J, J, I, H, E, H, J, J, F, J, E, E, I, J, J, J, I,~
## $ clarity <ord> SI2, SI1, VS1, VS2, SI2, VVS2, VVS1, SI1, VS2, VS1, SI1, VS1, ~
## $ depth   <dbl> 61.5, 59.8, 56.9, 62.4, 63.3, 62.8, 62.3, 61.9, 65.1, 59.4, 64~
## $ table   <dbl> 55, 61, 65, 58, 58, 57, 57, 55, 61, 61, 55, 56, 61, 54, 62, 58~
## $ price   <int> 326, 326, 327, 334, 335, 336, 336, 337, 337, 338, 339, 340, 34~
## $ x       <dbl> 3.95, 3.89, 4.05, 4.20, 4.34, 3.94, 3.95, 4.07, 3.87, 4.00, 4.~
## $ y       <dbl> 3.98, 3.84, 4.07, 4.23, 4.35, 3.96, 3.98, 4.11, 3.78, 4.05, 4.~
## $ z       <dbl> 2.43, 2.31, 2.31, 2.63, 2.75, 2.48, 2.47, 2.53, 2.49, 2.39, 2.~
```

```r
summary(diamonds) #  produce result summaries of the results of
```

```
##      carat               cut          color        clarity          depth
##  Min.   :0.2000   Fair     : 1610   D: 6775   SI1    :13065   Min.   :43.00
##  1st Qu.:0.4000   Good     : 4906   E: 9797   VS2    :12258   1st Qu.:61.00
##  Median :0.7000   Very Good:12082   F: 9542   SI2    : 9194   Median :61.80
##  Mean   :0.7979   Premium  :13791   G:11292   VS1    : 8171   Mean   :61.75
##  3rd Qu.:1.0400   Ideal    :21551   H: 8304   VVS2   : 5066   3rd Qu.:62.50
##  Max.   :5.0100                     I: 5422   VVS1   : 3655   Max.   :79.00
##                                     J: 2808   (Other): 2531
##      table           price             x                y
##  Min.   :43.00   Min.   :  326   Min.   : 0.000   Min.   : 0.000
##  1st Qu.:56.00   1st Qu.:  950   1st Qu.: 4.710   1st Qu.: 4.720
```
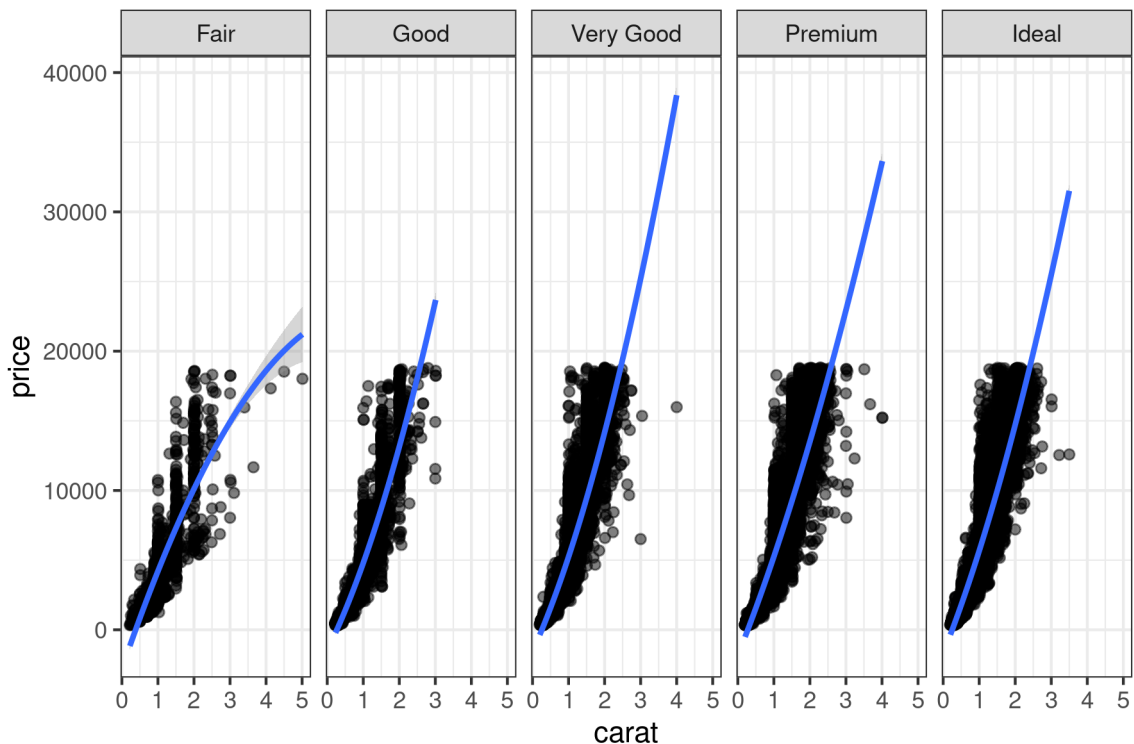
```
##  Median :57.00    Median : 2401    Median : 5.700    Median : 5.710
##  Mean   :57.46    Mean   : 3933    Mean   : 5.731    Mean   : 5.735
##  3rd Qu.:59.00    3rd Qu.: 5324    3rd Qu.: 6.540    3rd Qu.: 6.540
##  Max.   :95.00    Max.   :18823    Max.   :10.740    Max.   :58.900
##
##         z
##  Min.   : 0.000
##  1st Qu.: 2.910
##  Median : 3.530
##  Mean   : 3.539
##  3rd Qu.: 4.040
##  Max.   :31.800
##
```
```
# various model fitting functions
```

#### 4.1.2.2.2   Now some compact ggplot2 EDA code

- Use the pipe **%>%** operator

    - Which simply passes the output of the left operator
    - As the first argument to the right operator

```
diamonds %>%
  ggplot(aes(x = carat, y = price)) +  # aes is aesthetic mapping
  geom_point(alpha = 0.5) + # each data point as a point
  facet_grid( ~ cut) + # facet the scatter plot on cut, color or clarity
  stat_smooth(method = lm, formula = y ~ poly(x, 2)) + # fit a 2nd order lin. model
  theme_bw()
```



With this simple visualization,

- We can quickly see that price increases with carat size,
    - The relationship is nonlinear,
- There are some outliers,
    - And the relationship does not depend too heavily on cut.

Now lets use GGally and GGpairs packages

- These packages are extensions to GGplot 2

```r
library(GGally) # a ggplot2 extention; gallery of plot templates
```

```
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg   ggplot2
```

```r
? GGally  # This doesn't work for the GGally package
```
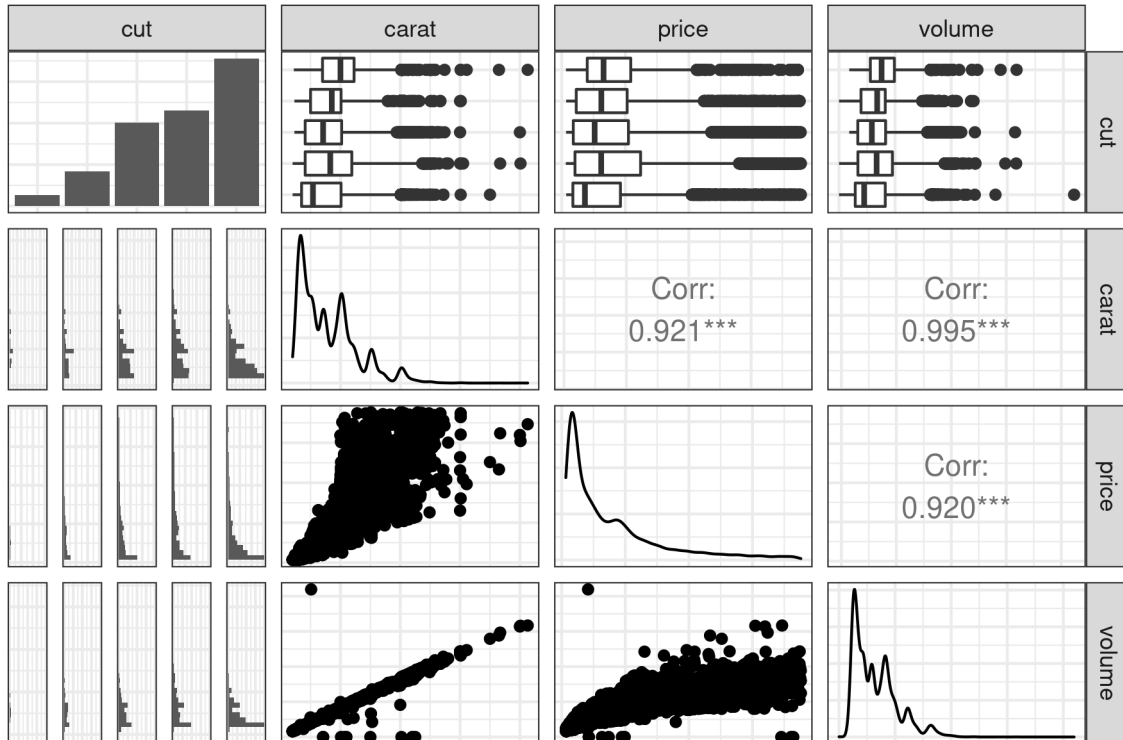
```r
?? GGally # So try this
```

```r
diamonds %>%
  mutate(volume = x * y * z) %>%     # in the pipe calculate the volume
  select(cut, carat, price, volume) %>%
  sample_frac(0.5, replace = TRUE) %>%
  ggpairs(axisLabels = "none") +
  theme_bw()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



The Tidyverse functions

- mutate, select, and sample_frac (verbs!)

- Are part of the dplyr data manipulation library
- And are loaded when you load the tidyverse metapackage

### 4.1.2.3 References

- Deep Ganguli The Grammar of DSCI