## Subsetting Data in R

Data Wrangling in R

# Subsetting part 2

#### Data

# A tibble: 6 x 10

Let's continue to work with the Diamond dataset from the ggplot2 package of the tidyverse.

We will often use the glimpse() function of the dplyr package of the tidyverse to look at a rotated view of the data.

```
library(tidyverse)
head(diamonds)
```

```
color clarity depth table price
 carat cut
                          <dbl> <ord>
              <ord> <ord>
  0.23 Ideal
           F.
                   ST2
                           61.5
                                 55
                                      326
                                          3.95
2 0.21 Premium E
                   SI1
                           59.8
                                 61
                                      326
                                          3.89
3
              F.
                                 65
                                      327
  0.23 Good
                   VS1
                           56.9
                                          4.05
  0.29 Premium I
                   VS2
                           62.4
                                 58
                                      334
                                          4.2
4
5 0.31 Good
                   SI2
                           63.3
                                 58
                                      335 4.34
                                              4
6
                                  57
  0.24 Very Good J
                   VVS2
                           62.8
                                      336
                                          3.94
```

3

3

4

4

3

### Let's learn more about this data

We can use ?diamonds to get more informatin in the Help pane.

We might decide to rename some columns,

- ▶ x to be length
- y to be width
- z to be depth
- but first changing depth to be depth\_percentage

### Renaming Columns of a data frame or tibble

To rename columns in dplyr, you can use the rename function.

Notice the new name is listed first!

3 0.23 Good E VS1

56.9 65 327

### More Renaming

\$ table

\$ price

\$ length

\$ width

```
diamonds_2<- diamonds %>%
              rename(depth_percentage = depth,
                                                                                                                                          length = x,
                                                                                                                                                 width = y,
                                                                                                                                                 depth = z)
glimpse(diamonds 2)
Rows: 53,940
Columns: 10
                                                                                                                                          <dbl> 0.23, 0.21, 0.23, 0.29, 0.31, 0.24
$ carat
                                                                                                                                          <ord> Ideal, Premium, Good, Premium,
$ cut
                                                                                                                                          <ord> E, E, E, I, J, J, I, H, E, H, J, .
$ color
$ clarity
                                                                                                                                          <ord> SI2, SI1, VS1, VS2, SI2, VVS2, VVS
```

\$ depth\_percentage <dbl> 61.5, 59.8, 56.9, 62.4, 63.3, 62.8

<dbl> 55, 61, 65, 58, 58, 57, 57, 55, 63

<int> 326, 326, 327, 334, 335, 336, 336

<dbl> 3.95, 3.89, 4.05, 4.20, 4.34, 3.94

<dbl> 3.98, 3.84, 4.07, 4.23, 4.35, 3.96

### Take Care with Column Names

When you can, avoid spaces, special punctuation, or numbers in column names, as these require special treatment to refer to them.

See https://jhudatascience.org/intro\_to\_r/quotes\_vs\_backticks.html for more guidance.

```
diamonds %>% rename(depth percentage = depth) # this will condition diamonds %>% rename(depth_percentage = depth) # this will a
```

diamonds %>% rename(`depth percentage` = depth) # not record

### Unusual Column Names

It's best to avoid unusual column names where possible, as things get tricky later.

We just showed the use of 'backticks'. You may see people use quotes as well.



Other atypical column names are those with:

- spaces
- number without characters
- number starting the name
- other punctuation marks (besides "\_" or ":" and not at the beginning)



Rename tricky column names so that you don't have to deal with them later!

### Example

### glimpse(diamonds\_bad\_names)

# A tibble: 53,940 x 10

```
Rows: 53,940
Columns: 10
                         <dbl> 0.23, 0.21, 0.23, 0.29, 0.3
$ carat
$ cut
                         <ord> Ideal, Premium, Good, Premium
                         <ord> E, E, E, I, J, J, I, H, E, I
$ color
$ clarity
                         <ord> SI2, SI1, VS1, VS2, SI2, VVS
$ depth
                         <dbl> 61.5, 59.8, 56.9, 62.4, 63.3
                         <dbl> 55, 61, 65, 58, 58, 57, 57,
$ table
$ `Price(in US dollars)` <int> 326, 326, 327, 334, 335, 336
$ `Length (in mm)` <dbl> 3.95, 3.89, 4.05, 4.20, 4.34
$ `Width in mm`
                        <dbl> 3.98, 3.84, 4.07, 4.23, 4.39
$ `Depth percentage` <dbl> 2.43, 2.31, 2.31, 2.63, 2.79
diamonds_bad_names %>%
        rename(price = `Price(in US dollars)`)
```

# Renaming all columns of a data frame: dplyr

To rename all columns you use the rename\_with(). In this case we will use toupper() to make all letters upper case. Could also use tolower() function.

diamonds upper <- diamonds %>% rename with(toupper) head(diamonds upper, 2)

```
# A tibble: 2 \times 10
```

CARAT CUT COLOR CLARITY DEPTH TABLE PRICE X

<dbl> <ord> <ord> <dbl> <dbl> <int> <dbl> <dbl> <</pre> 1 0.23 Ideal E SI2 61.5 55

2 0.21 Premium E SI1

59.8 61 326 3.89 3.84

diamonds\_upper %>% rename\_with(tolower) %>% head(n = 2)

326 3.95 3.98

# A tibble: 2 x 10

carat cut color clarity depth table price

<dbl> <ord> <ord> <ord>

<dbl> <dbl> <dbl> <dbl> <dbl> <dbl> < 1 0.23 Ideal E SI2 61.5 55 326 3.95 3.98

59.8 61

326 3 89 3 84

2 0 21 Premium F ST1

### Janitor package

Rows: 53,940

```
#install.packages("janitor")
library(janitor)
clean_names(diamonds_bad_names) %>% glimpse()
```

```
Columns: 10
$ carat
                                                                                                               <dbl> 0.23, 0.21, 0.23, 0.29, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 0.31, 
                                                                                                               <ord> Ideal, Premium, Good, Premium,
$ cut
                                                                                                               <ord> E, E, E, I, J, J, I, H, E, H, .
$ color
$ clarity
                                                                                                               <ord> SI2, SI1, VS1, VS2, SI2, VVS2,
$ depth
                                                                                                                <dbl> 61.5, 59.8, 56.9, 62.4, 63.3, 6
$ table
                                                                                                                <dbl> 55, 61, 65, 58, 58, 57, 57, 55
$ price_in_us_dollars <int> 326, 326, 327, 334, 335, 336, 3
$ length_in_mm
                                                                                                               <dbl> 3.95, 3.89, 4.05, 4.20, 4.34, 3
$ width_in_mm
                                                                                                               <dbl> 3.98, 3.84, 4.07, 4.23, 4.35, 3
$ depth_percentage
                                                                                                              <dbl> 2.43, 2.31, 2.31, 2.63, 2.75, 3
```

## Subset based on a class

# The where() function can help select columns of a specific class

head(diamonds, 2)

# A tibble: 2 x 10

 $\verb"is.character"() and \verb"is.numeric"() are often the most helpful$ 

2 0.21 59.8 61 326 3.89 3.84 2.31

# The where() function can help select columns of a specific class

 $\verb"is.character"() and \verb"is.numeric"() are often the most helpful$ 

### distinct() function

To filter for distinct values from a variable, multiple variables, or an entire tibble you can use the distinct() function from the dplyr package. Similar to count, but without the number of times the value shows up.

#### distinct(diamonds, cut)

```
# A tibble: 5 x 1
  cut
  <ord>
1 Ideal
Premium
3 Good
```

4 Very Good 5 Fair

# Adding/Removing Columns

## Adding columns to a data frame: dplyr (tidyverse way) The mutate function in dplyr allows you to add or modify columns

```
of a data frame.
# General format - Not the code!
{data object to update} <- {data to use} %>%
                         mutate({new variable name} = {new variable name} = {new variable name}
```

1 US dollar = 1.32 Canadian dollars diamonds %>%

\$ depth

\$ table

```
mutate(price_canadian = price * 1.32) %>% glimpse()
Rows: 53,940
```

Columns: 11

\$ carat <dbl> 0.23, 0.21, 0.23, 0.29, 0.31, 0.24,

\$ cut <ord> Ideal, Premium, Good, Premium, Good</br>

\$ color <ord> E, E, E, I, J, J, I, H, E, H, J, J,

\$ clarity <ord> SI2, SI1, VS1, VS2, SI2, VVS2, VVS1

<dbl> 61.5, 59.8, 56.9, 62.4, 63.3, 62.8,

<dbl> 55, 61, 65, 58, 58, 57, 57, 55, 61,

# Use mutate to modify existing columns

The mutate function in dplyr allows you to add or modify columns of a data frame.

mutate(diamonds, price = price \* 1.32) %>% glimpse()

```
Rows: 53,940
Columns: 10
```

\$ x

<dbl> 3.95, 3.89, 4.05, 4.20, 4.34, 3.94, 3.95, 4

### remember to save your data

If you want to actually make the change you need to reassign the

```
data object.
```

diamonds <- diamonds %>% mutate(price = price \* 1.32) %>% y

### Removing columns of a data frame: dplyr

select(diamonds, - price) %>% glimpse()

The select function can remove a column with minus (-)

\$ clarity <ord> SI2, SI1, VS1, VS2, SI2, VVS2, VVS1, SI1, V

Or, you can simply select the columns you want to keep, ignoring the ones you want to remove.

### Removing columns in a data frame: dplyr

You can use c() to list the columns to remove.

Remove newcol and drat:

```
select(diamonds, -c("x", "y", "z")) %>% glimpse()
Rows: 53,940
Columns: 7
$ carat <dbl> 0.23, 0.21, 0.23, 0.29, 0.31, 0.24, 0.24,
$ cut <ord> Ideal, Premium, Good, Premium, Good, Very (
$ color <ord> E, E, E, I, J, J, I, H, E, H, J, J, F, J, I
$ clarity <ord> SI2, SI1, VS1, VS2, SI2, VVS2, VVS1, SI1, V
$ depth <dbl> 61.5, 59.8, 56.9, 62.4, 63.3, 62.8, 62.3, 63.3, 63.8, 63.3, 63.8, 63.3, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8, 63.8
$ table <dbl> 55, 61, 65, 58, 58, 57, 57, 55, 61, 61, 55
$ price <int> 326, 326, 327, 334, 335, 336, 336, 337, 33
```

# Ordering columns

The select function can reorder columns.

326 61.5 0.23 Ideal E

326 59.8 0.21 Premium E

2

```
head(diamonds, n = 2)
# A tibble: 2 x 10
 carat cut color clarity depth table price x
 <dbl> <ord> <ord> <dbl> <dbl> <int> <dbl> <dbl> <dbl> <
1 0.23 Ideal E SI2 61.5 55 326 3.95 3.98
2 0.21 Premium E SI1
                          59.8 61 326 3.89 3.84
diamonds %>% select(price, depth, carat, cut, color) %>% he
# A tibble: 2 x 5
 price depth carat cut color
 <int> <dbl> <ord> <ord>
```

2

The select function can reorder columns. Put price first, then select the rest of columns:

```
head(diamonds, n = 2)
# A tibble: 2 x 10
  carat cut color clarity depth table price x
  <dbl> <ord> <ord> <dbl> <dbl> <int> <dbl> <dbl> <dbl> <
1 0.23 Ideal E SI2 61.5 55 326 3.95 3.98
2 0.21 Premium E SI1
                           59.8 61 326 3.89 3.84
diamonds %>% select(price, everything()) %>% head(n = 2)
# A tibble: 2 x 10
 price carat cut color clarity depth table x
  <int> <dbl> <ord> <ord> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
```

326 0.23 Ideal E SI2 61.5 55 3.95 3.98

326 0.21 Premium E SI1 59.8 61 3.89 3.84

head(diamonds, n = 2)

Put price at the end ("remove, everything, then add back in"):

In addition to select we can also use the relocate() function of dplyr to rearrange the columns for more complicated moves.

For example, let say we just wanted price to be before carat.

```
head(diamonds, n = 2)
```

```
# A tibble: 2 x 10
```

carat cut color clarity depth table price x

```
<dbl> <ord>      <ord>      <dbl> <dbl> <int> <dbl> <dbl> <</pre>
```

```
1 0.23 Ideal E
               SI2 61.5 55 326 3.95 3.98
```

2 0.21 326 Premium E

```
diamonds %>% relocate( price, .before = cut) %>% head(n = 5
```

# A tibble: 2 x 10

carat price cut color clarity depth table x

<dbl> <int> <ord> <ord> <dbl> <dbl> <dbl> <dbl> <dbl> < 1 0.23 326 Ideal E SI2 61.5 55 3.95 3.98

SI1

59.8

61 3.89 3.84

```
Ordering the column names of a data frame: alphabetically
   Using the base R order() function.
```

```
order(colnames(diamonds))
```

[1] 6 9 10 diamonds %>% select(order(colnames(diamonds)))

# A tibble: 53,940 x 10	
carat clarity color cut	depth price table

Ε

Ε

Ε

Ι

J

J

Ι

Η

Ε

# A tibb.	Le: 53,94	40 x 10	
carat	clarity	color cut	
<dbl></dbl>	<ord></ord>	<ord> <ord></ord></ord>	

0.23 SI2

0.21 SI1

0.23 VS1

0.29 VS2

0.31 SI2

0.24 VVS2

0.24 VVS1

0.26 SI1

0.22 VS2

3

5

6

8

Ideal

Good

Good

Fair

Premium

Premium

Very Good

Very Good

Very Good

59.8

62.4

62.8

62.3

61.9

65.1

56.9

63.3

334

335

336

336

337

337

58

58

57

57

55

61

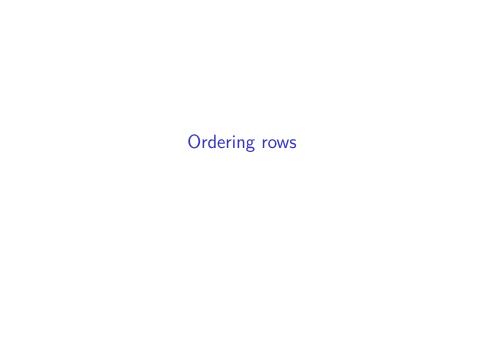
3.94

3.95

4.07

3.87

X



The arrange function can reorder rows By default, arrange orders in increasing order:

```
diamonds %>% arrange(cut)
```

... with 53,930 more rows

```
# A tibble: 53,940 x 10
               color clarity depth table price
   carat cut
   <dbl> <ord> <ord> <dbl> <int> <dbl> <dbl> <int> <dbl> <dbl> <
```

1	0.22 Fair	Ε	VS2	65.1	61	337	3.87	3.78
2	0.86 Fair	E	SI2	55.1	69	2757	6.45	6.33
3	0.96 Fair	F	SI2	66.3	62	2759	6.27	5.95

2	0.86	rair	Ł	S12	55.1	69	2/5/	6.45	6.3
3	0.96	Fair	F	SI2	66.3	62	2759	6.27	5.9
4	0.7	Fair	F	VS2	64.5	57	2762	5.57	5.5

3	0.96	Fair	F	SI2	66.3	62	2759	6.27	5.95
4	0.7	Fair	F	VS2	64.5	57	2762	5.57	5.53
5	0.7	Fair	F	VS2	65.3	55	2762	5.63	5.58
6	0.91	Fair	H	SI2	64.4	57	2763	6.11	6.09

3	0.96	Fair	F	SI2	66.3	62	2759	6.27	5.98
4	0.7	Fair	F	VS2	64.5	57	2762	5.57	5.53
5	0.7	Fair	F	VS2	65.3	55	2762	5.63	5.58
6	0.91	Fair	H	SI2	64.4	57	2763	6.11	6.09
7	0.91	Fair	H	SI2	65.7	60	2763	6.03	5.99
8	0 98	Fair	H	ST2	67 9	60	2777	6 05	5 97

3	0.96	Fair	F	SI2	66.3	62	2759	6.27	5.9
4	0.7	Fair	F	VS2	64.5	57	2762	5.57	5.5
5	0.7	Fair	F	VS2	65.3	55	2762	5.63	5.5
6	0.91	Fair	H	SI2	64.4	57	2763	6.11	6.0
7	0.91	Fair	H	SI2	65.7	60	2763	6.03	5.9
8	0.98	Fair	Н	SI2	67.9	60	2777	6.05	5.9

0.84 Fair G 55.1 67 2782 6.39 6.2 SI1

1.01 Fair E 10 I1 64.5 58 2788 6.29 6.21

Use the desc to arrange the rows in descending order:

```
diamonds %>% arrange(depth)
```

```
# A tibble: 53,940 x 10
               color clarity depth table price
   carat cut
   <dbl> <ord> <ord> <ord>
                              <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
         Fair
                     SI1
                               43
                                       59
                                           3634
                                                 6.32
                                                        6.27
               G
                     VS2
    1.09 Ideal J
                               43
                                       54
                                           4778 6.53
                                                        6.55
 3
         Fair
               G
                     VS2
                               44
                                       53
                                           4032
                                                 6.31
                                                        6.24
    1.43 Fair
                               50.8
                                           6727 7.73
                                                        7.25
              Т
                     VS1
                                       60
 5
    0.3
               F.
                     VVS2
                               51
                                       67
                                            945
                                                 4.67
                                                        4.62
        Fair
        Fair
                     SI1
                               52.2
                                       65
                                           1895
                                                 6.04
 6
    0.7
               D
                                                        5.99
 7
    0.37 Fair
               F
                     ΤF
                               52.3
                                       61
                                           1166
                                                 4.96
                                                        4.91
 8
    0.56 Fair
               Η
                     VS2
                               52.7
                                       70
                                           1293
                                                 5.71
                                                        5.57
    1.02 Fair
                     SI1
                               53
                                       63
                                           2856
                                                 6.84
                                                        6.77
10
    0.96 Fair
               Ε
                     SI2
                               53.1
                                       63
                                           2815
                                                 6.73
                                                        6.65
  ... with 53,930 more rows
```

Use the desc to arrange the rows in descending order:

```
diamonds %>% arrange(desc(price))
```

```
# A tibble: 53,940 x 10
                 color clarity depth table price
  carat cut
                                                   X
  <dbl> <ord>
                 <ord> <ord>
                               <dbl> <dbl> <int> <dbl> <
   2.29 Premium
                       VS2
                                60.8
                                       60 18823 8.5
                                63.5
                                       56 18818 7.9
        Very Good G
                       SI1
3
   1.51 Ideal
                       ΙF
                                61.7
                                       55 18806 7.37
                 G
                                62.5
                                       55 18804
   2.07 Ideal
                       SI2
                                                 8.2
5 2
        Very Good H
                       SI1
                                62.8
                                       57 18803
                                                 7.95
6
                                61.8
                                       59 18797
                                                 8.52
   2.29 Premium
                       SI1
   2.04 Premium
                       SI1
                                58.1
                                       60 18795
                                                 8.37
               Н
8
        Premium
               I
                       VS1
                                60.8
                                       59 18795
                                                 8.13
```

59 18791 7.57

8.29

54 18791

# ... with 53,930 more rows

VS2

SI2

62.3 62.6

F

G

1.71 Premium

2.15 Ideal

10

# A tibble: 53,940 x 10

<dbl> <ord>

5.01 Fair

carat cut

You can combine increasing and decreasing orderings:

```
arrange(diamonds, desc(carat), table)
```

<ord> <ord>

Ι1

```
2 4.5 Fair
                     Ι1
                             65.8
                                    58 18531 10.2
                                                 1
  4.13 Fair
            H I1
                             64.8
                                    61 17329 10
  4.01 Premium I I1
                             61
                                    61 15223 10.1
                                                 10
5
  4.01 Premium
                J I1
                             62.5
                                    62 15223 10.0
       Very Good I I1
                             63.3
                                    58 15984 10.0
6
7 3.67 Premium
                     Ι1
                             62.4
                                    56 16193 9.86
              I
8
   3.65 Fair
                Η
                     Ι1
                             67.1
                                    53 11668 9.53 9
   3.51 Premium
                .J
                     VS2
                             62.5
                                    59 18701 9.66
10
   3.5 Ideal
                Η
                     Ι1
                             62.8
                                    57 12587 9.65
# ... with 53,930 more rows
```

color clarity depth table price

65.5

X

<dbl> <dbl> <int> <dbl> <

59 18018 10.7 10

Lab

 ${\sf Link}\ {\sf to}\ {\sf Lab}$