

Lab 9

Juan D Astudillo

11:59PM April 14, 2019

“data wrangling / munging / carpentry” with dplyr.

First load dplyr, tidyr, magrittr and lubridate in one line.

```
pacman::p_load(dplyr, tidyr, magrittr, lubridate)
```

Load the `storms` dataset from the `dplyr` package and investigate it using `str` and `summary` and `head`. Which two columns should be converted to type factor? Do so below using the `mutate` and the overwrite pipe operator `%<>%`. Verify.

```
data("storms")
str(storms)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 10010 obs. of 13 variables:
## $ name      : chr "Amy" "Amy" "Amy" "Amy" ...
## $ year      : num 1975 1975 1975 1975 1975 ...
## $ month     : num 6 6 6 6 6 6 6 6 6 6 ...
## $ day       : int 27 27 27 27 28 28 28 28 29 29 ...
## $ hour      : num 0 6 12 18 0 6 12 18 0 6 ...
## $ lat       : num 27.5 28.5 29.5 30.5 31.5 32.4 33.3 34 34.4 34 ...
## $ long      : num -79 -79 -79 -79 -78.8 -78.7 -78 -77 -75.8 -74.8 ...
## $ status    : chr "tropical depression" "tropical depression" "tropical depression" "tropical dep
## $ category  : Ord.factor w/ 7 levels "-1"<"0"<"1"<"2"<...: 1 1 1 1 1 1 1 1 2 2 ...
## $ wind      : int 25 25 25 25 25 25 25 30 35 40 ...
## $ pressure  : int 1013 1013 1013 1013 1012 1012 1011 1006 1004 1002 ...
## $ ts_diameter: num NA NA NA NA NA NA NA NA NA NA ...
## $ hu_diameter: num NA NA NA NA NA NA NA NA NA NA ...
```

```
summary(storms)
```

```
##      name      year      month      day
## Length:10010   Min.   :1975   Min.   : 1.000   Min.   : 1.00
## Class :character 1st Qu.:1990   1st Qu.: 8.000   1st Qu.: 8.00
## Mode :character  Median :1999   Median : 9.000   Median :16.00
##                Mean   :1998   Mean   : 8.779   Mean   :15.86
##                3rd Qu.:2006   3rd Qu.: 9.000   3rd Qu.:24.00
##                Max.   :2015   Max.   :12.000   Max.   :31.00
##
##      hour      lat      long      status
## Min.   : 0.000   Min.   : 7.20   Min.   : -109.30   Length:10010
## 1st Qu.: 6.000   1st Qu.:17.50   1st Qu.: -80.70   Class :character
## Median :12.000   Median :24.40   Median : -64.50   Mode  :character
## Mean   : 9.114   Mean   :24.76   Mean   : -64.23
## 3rd Qu.:18.000   3rd Qu.:31.30   3rd Qu.: -48.60
## Max.   :23.000   Max.   :51.90   Max.   : -6.00
##
## category      wind      pressure      ts_diameter
## -1:2545   Min.   : 10.00   Min.   : 882.0   Min.   : 0.00
## 0 :4373   1st Qu.: 30.00   1st Qu.: 985.0   1st Qu.: 69.05
```

```
## 1 :1685 Median : 45.00 Median : 999.0 Median : 138.09
## 2 : 628 Mean : 53.49 Mean : 992.1 Mean : 166.76
## 3 : 363 3rd Qu.: 65.00 3rd Qu.:1006.0 3rd Qu.: 241.66
## 4 : 348 Max. :160.00 Max. :1022.0 Max. :1001.18
## 5 : 68 NA's :6528
## hu_diameter
## Min. : 0.00
## 1st Qu.: 0.00
## Median : 0.00
## Mean : 21.41
## 3rd Qu.: 28.77
## Max. :345.23
## NA's :6528
```

```
head(storms)
```

```
## # A tibble: 6 x 13
##   name   year month   day hour   lat   long status category wind pressure
##   <chr> <dbl> <dbl> <int> <dbl> <dbl> <dbl> <chr>   <ord>    <int>    <int>
## 1 Amy   1975     6    27     0  27.5 -79   tropi~ -1      25     1013
## 2 Amy   1975     6    27     6  28.5 -79   tropi~ -1      25     1013
## 3 Amy   1975     6    27    12  29.5 -79   tropi~ -1      25     1013
## 4 Amy   1975     6    27    18  30.5 -79   tropi~ -1      25     1013
## 5 Amy   1975     6    28     0  31.5 -78.8 tropi~ -1      25     1012
## 6 Amy   1975     6    28     6  32.4 -78.7 tropi~ -1      25     1012
## # ... with 2 more variables: ts_diameter <dbl>, hu_diameter <dbl>
```

```
storms %<>%
  mutate(name = factor(name), status = factor(status))
str(storms)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 10010 obs. of 13 variables:
## $ name : Factor w/ 198 levels "AL011993","AL012000",...: 44 44 44 44 44 44 44 44 44 44 ...
## $ year : num 1975 1975 1975 1975 1975 1975 ...
## $ month : num 6 6 6 6 6 6 6 6 6 6 ...
## $ day : int 27 27 27 27 28 28 28 28 29 29 ...
## $ hour : num 0 6 12 18 0 6 12 18 0 6 ...
## $ lat : num 27.5 28.5 29.5 30.5 31.5 32.4 33.3 34 34.4 34 ...
## $ long : num -79 -79 -79 -79 -78.8 -78.7 -78 -77 -75.8 -74.8 ...
## $ status : Factor w/ 3 levels "hurricane","tropical depression",...: 2 2 2 2 2 2 2 2 3 3 ...
## $ category : Ord.factor w/ 7 levels "-1"<"0"<"1"<"2"<...: 1 1 1 1 1 1 1 1 2 2 ...
## $ wind : int 25 25 25 25 25 25 25 30 35 40 ...
## $ pressure : int 1013 1013 1013 1013 1012 1012 1011 1006 1004 1002 ...
## $ ts_diameter: num NA NA NA NA NA NA NA NA NA NA ...
## $ hu_diameter: num NA NA NA NA NA NA NA NA NA NA ...
```

Reorder the columns so name is first, status is second, category is third and the rest are the same. Verify.

```
storms %<>%
  select(name, status, category, everything())
storms
```

```
## # A tibble: 10,010 x 13
##   name status category year month day hour lat long wind pressure
##   <fct> <fct> <ord>    <dbl> <dbl> <int> <dbl> <dbl> <dbl> <int>    <int>
## 1 Amy   tropi~ -1      1975     6    27     0  27.5 -79     25     1013
## 2 Amy   tropi~ -1      1975     6    27     6  28.5 -79     25     1013
```

```
## 3 Amy tropi~ -1 1975 6 27 12 29.5 -79 25 1013
## 4 Amy tropi~ -1 1975 6 27 18 30.5 -79 25 1013
## 5 Amy tropi~ -1 1975 6 28 0 31.5 -78.8 25 1012
## 6 Amy tropi~ -1 1975 6 28 6 32.4 -78.7 25 1012
## 7 Amy tropi~ -1 1975 6 28 12 33.3 -78 25 1011
## 8 Amy tropi~ -1 1975 6 28 18 34 -77 30 1006
## 9 Amy tropi~ 0 1975 6 29 0 34.4 -75.8 35 1004
## 10 Amy tropi~ 0 1975 6 29 6 34 -74.8 40 1002
## # ... with 10,000 more rows, and 2 more variables: ts_diameter <dbl>,
## # hu_diameter <dbl>
```

Sort the dataframe by year (most recent first) then category of the storm (most severe first). Verify.

```
storms %<>%
  arrange(desc(year), desc(category))
storms
```

```
## # A tibble: 10,010 x 13
##   name status category year month day hour lat long wind pressure
##   <fct> <fct> <ord>   <dbl> <dbl> <int> <dbl> <dbl> <dbl> <int>   <int>
## 1 Joaq~ hurri~ 4      2015 10 1 12 23.1 -73.7 115 942
## 2 Joaq~ hurri~ 4      2015 10 1 18 23 -74.2 115 936
## 3 Joaq~ hurri~ 4      2015 10 2 0 22.9 -74.4 120 931
## 4 Joaq~ hurri~ 4      2015 10 2 6 23 -74.7 120 935
## 5 Joaq~ hurri~ 4      2015 10 2 12 23.4 -74.8 115 937
## 6 Joaq~ hurri~ 4      2015 10 3 0 24.3 -74.3 115 943
## 7 Joaq~ hurri~ 4      2015 10 3 6 24.8 -73.6 120 945
## 8 Joaq~ hurri~ 4      2015 10 3 12 25.4 -72.6 135 934
## 9 Joaq~ hurri~ 4      2015 10 3 18 26.3 -71 130 934
## 10 Joaq~ hurri~ 4      2015 10 4 0 27.4 -69.5 115 941
## # ... with 10,000 more rows, and 2 more variables: ts_diameter <dbl>,
## # hu_diameter <dbl>
```

Create a new feature wind_speed_per_unit_pressure.

```
storms %<>%
  mutate(wind_speed_per_unit_pressure = wind / pressure)
storms
```

```
## # A tibble: 10,010 x 14
##   name status category year month day hour lat long wind pressure
##   <fct> <fct> <ord>   <dbl> <dbl> <int> <dbl> <dbl> <dbl> <int>   <int>
## 1 Joaq~ hurri~ 4      2015 10 1 12 23.1 -73.7 115 942
## 2 Joaq~ hurri~ 4      2015 10 1 18 23 -74.2 115 936
## 3 Joaq~ hurri~ 4      2015 10 2 0 22.9 -74.4 120 931
## 4 Joaq~ hurri~ 4      2015 10 2 6 23 -74.7 120 935
## 5 Joaq~ hurri~ 4      2015 10 2 12 23.4 -74.8 115 937
## 6 Joaq~ hurri~ 4      2015 10 3 0 24.3 -74.3 115 943
## 7 Joaq~ hurri~ 4      2015 10 3 6 24.8 -73.6 120 945
## 8 Joaq~ hurri~ 4      2015 10 3 12 25.4 -72.6 135 934
## 9 Joaq~ hurri~ 4      2015 10 3 18 26.3 -71 130 934
## 10 Joaq~ hurri~ 4      2015 10 4 0 27.4 -69.5 115 941
## # ... with 10,000 more rows, and 3 more variables: ts_diameter <dbl>,
## # hu_diameter <dbl>, wind_speed_per_unit_pressure <dbl>
```

Create a new feature: average_diameter which averages the two diameters.

```
storms %<>%
  mutate(average_diameter = (ts_diameter + hu_diameter) / 2)
storms

## # A tibble: 10,010 x 15
##   name status category year month day hour lat long wind pressure
##   <fct> <fct> <ord>    <dbl> <dbl> <int> <dbl> <dbl> <dbl> <int>    <int>
## 1 Joaq~ hurri~ 4      2015 10    1    12 23.1 -73.7 115     942
## 2 Joaq~ hurri~ 4      2015 10    1    18 23   -74.2 115     936
## 3 Joaq~ hurri~ 4      2015 10    2     0 22.9 -74.4 120     931
## 4 Joaq~ hurri~ 4      2015 10    2     6 23   -74.7 120     935
## 5 Joaq~ hurri~ 4      2015 10    2    12 23.4 -74.8 115     937
## 6 Joaq~ hurri~ 4      2015 10    3     0 24.3 -74.3 115     943
## 7 Joaq~ hurri~ 4      2015 10    3     6 24.8 -73.6 120     945
## 8 Joaq~ hurri~ 4      2015 10    3    12 25.4 -72.6 135     934
## 9 Joaq~ hurri~ 4      2015 10    3    18 26.3 -71   130     934
## 10 Joaq~ hurri~ 4      2015 10    4     0 27.4 -69.5 115     941
## # ... with 10,000 more rows, and 4 more variables: ts_diameter <dbl>,
## #   hu_diameter <dbl>, wind_speed_per_unit_pressure <dbl>,
## #   average_diameter <dbl>
```

Calculate the distance from each storm observation to Miami in a new variable `distance_to_miami`.

```
MIAMI_COORDS = c(25.7617, -80.1918)
RAD_EARTH = 3963

compute_globe_distance = function(destination, origin){
  destination_radians = destination*pi/180
  origin_radians = origin*pi/180
  lat_change = destination_radians[1] - origin_radians[1]
  long_change = destination_radians[2] - origin_radians[2]
  a = (sin(lat_change/2))^2 + cos(origin_radians[1]) * cos(destination_radians[1]) * (sin(long_change/2))^2
  return(RAD_EARTH * 2 * asin(sqrt(a)))
}

storms %<>%
  rowwise() %>%
  mutate(distance_to_miami = compute_globe_distance(MIAMI_COORDS, c(lat, long))) %>%
  select(lat, long, distance_to_miami, everything())
storms
```

```
## Source: local data frame [10,010 x 16]
## Groups: <by row>
##
## # A tibble: 10,010 x 16
##   lat long distance_to_mia~ name status category year month day
##   <dbl> <dbl>          <dbl> <fct> <fct> <ord>    <dbl> <dbl> <int>
## 1 23.1 -73.7         448. Joaq~ hurri~ 4      2015 10    1
## 2 23   -74.2         423. Joaq~ hurri~ 4      2015 10    1
## 3 22.9 -74.4         415. Joaq~ hurri~ 4      2015 10    2
## 4 23   -74.7         395. Joaq~ hurri~ 4      2015 10    2
## 5 23.4 -74.8         376. Joaq~ hurri~ 4      2015 10    2
## 6 24.3 -74.3         383. Joaq~ hurri~ 4      2015 10    3
## 7 24.8 -73.6         418. Joaq~ hurri~ 4      2015 10    3
## 8 25.4 -72.6         474. Joaq~ hurri~ 4      2015 10    3
## 9 26.3 -71          572. Joaq~ hurri~ 4      2015 10    3
```

```
## 10 27.4 -69.5          671. Joaq~ hurri~ 4          2015    10    4
## # ... with 10,000 more rows, and 7 more variables: hour <dbl>, wind <int>,
## #   pressure <int>, ts_diameter <dbl>, hu_diameter <dbl>,
## #   wind_speed_per_unit_pressure <dbl>, average_diameter <dbl>
```

At home: convert year, month, day, hour into the variable `timestamp` using the `lubridate` package.

```
storms1 = storms

storms1 %<>%
  unite(timestamp, year, month, day, hour, sep = " - ") %<>%
  mutate(timestamp = ymd_h(timestamp))
```

At home: using the `lubridate` package, create new variables `day_of_week` which is a factor with levels “Sunday”, “Monday”, ... “Saturday” and `week_of_year` which is integer 1, 2, ..., 52.

```
storms1 %<>%
  mutate(day_of_week = wday(timestamp, label = TRUE)) %<>%
  mutate(week_of_the_year = week(timestamp))
```

Create a new data frame `serious_storms` which are category 3 and above hurricanes.

```
serious_storms = storms %>%
  filter(category >= 3)
serious_storms
```

```
## Source: local data frame [779 x 16]
## Groups: <by row>
##
## # A tibble: 779 x 16
##   lat long distance_to_mia~ name status category year month day
##   <dbl> <dbl>          <dbl> <fct> <fct> <ord>    <dbl> <dbl> <int>
## 1 23.1 -73.7          448. Joaq~ hurri~ 4          2015    10    1
## 2 23   -74.2          423. Joaq~ hurri~ 4          2015    10    1
## 3 22.9 -74.4          415. Joaq~ hurri~ 4          2015    10    2
## 4 23   -74.7          395. Joaq~ hurri~ 4          2015    10    2
## 5 23.4 -74.8          376. Joaq~ hurri~ 4          2015    10    2
## 6 24.3 -74.3          383. Joaq~ hurri~ 4          2015    10    3
## 7 24.8 -73.6          418. Joaq~ hurri~ 4          2015    10    3
## 8 25.4 -72.6          474. Joaq~ hurri~ 4          2015    10    3
## 9 26.3 -71           572. Joaq~ hurri~ 4          2015    10    3
## 10 27.4 -69.5          671. Joaq~ hurri~ 4          2015    10    4
## # ... with 769 more rows, and 7 more variables: hour <dbl>, wind <int>,
## #   pressure <int>, ts_diameter <dbl>, hu_diameter <dbl>,
## #   wind_speed_per_unit_pressure <dbl>, average_diameter <dbl>
```

In `serious_storms`, merge the variables `lat` and `long` together into `lat_long` with values `lat / long` as a string.

```
serious_storms %<>%
  unite(lat_long, lat, long, sep = " / ")
serious_storms
```

```
## # A tibble: 779 x 15
##   lat_long distance_to_mia~ name status category year month day hour
##   <chr>          <dbl> <fct> <fct> <ord>    <dbl> <dbl> <int> <dbl>
## 1 23.1 / ~          448. Joaq~ hurri~ 4          2015    10    1    12
## 2 23 / -7~          423. Joaq~ hurri~ 4          2015    10    1    18
```

```
## 3 22.9 / ~ 415. Joaq~ hurri~ 4 2015 10 2 0
## 4 23 / -7~ 395. Joaq~ hurri~ 4 2015 10 2 6
## 5 23.4 / ~ 376. Joaq~ hurri~ 4 2015 10 2 12
## 6 24.3 / ~ 383. Joaq~ hurri~ 4 2015 10 3 0
## 7 24.8 / ~ 418. Joaq~ hurri~ 4 2015 10 3 6
## 8 25.4 / ~ 474. Joaq~ hurri~ 4 2015 10 3 12
## 9 26.3 / ~ 572. Joaq~ hurri~ 4 2015 10 3 18
## 10 27.4 / ~ 671. Joaq~ hurri~ 4 2015 10 4 0
## # ... with 769 more rows, and 6 more variables: wind <int>,
## #   pressure <int>, ts_diameter <dbl>, hu_diameter <dbl>,
## #   wind_speed_per_unit_pressure <dbl>, average_diameter <dbl>
```

Back to the main dataframe `storms`, create a new feature `decile_windspeed` by binning wind speed into 10 bins.

```
storms %<>%
  mutate(decile_windspeed = factor(ntile(wind, 10)))
storms
```

```
## Source: local data frame [10,010 x 17]
## Groups: <by row>
##
## # A tibble: 10,010 x 17
##   lat long distance_to_mia~ name status category year month day
##   <dbl> <dbl>          <dbl> <fct> <fct> <ord>   <dbl> <dbl> <int>
## 1 23.1 -73.7          448. Joaq~ hurri~ 4 2015 10 1
## 2 23 -74.2           423. Joaq~ hurri~ 4 2015 10 1
## 3 22.9 -74.4          415. Joaq~ hurri~ 4 2015 10 2
## 4 23 -74.7           395. Joaq~ hurri~ 4 2015 10 2
## 5 23.4 -74.8          376. Joaq~ hurri~ 4 2015 10 2
## 6 24.3 -74.3          383. Joaq~ hurri~ 4 2015 10 3
## 7 24.8 -73.6          418. Joaq~ hurri~ 4 2015 10 3
## 8 25.4 -72.6          474. Joaq~ hurri~ 4 2015 10 3
## 9 26.3 -71           572. Joaq~ hurri~ 4 2015 10 3
## 10 27.4 -69.5          671. Joaq~ hurri~ 4 2015 10 4
## # ... with 10,000 more rows, and 8 more variables: hour <dbl>, wind <int>,
## #   pressure <int>, ts_diameter <dbl>, hu_diameter <dbl>,
## #   wind_speed_per_unit_pressure <dbl>, average_diameter <dbl>,
## #   decile_windspeed <fct>
```

Let's summarize some data. Find the strongest storm by wind speed per year.

```
storms %>%
  group_by(year) %>%
  summarize(max_wind_speed = max(wind))
```

```
## Warning: Grouping rowwise data frame strips rowwise nature
```

```
## # A tibble: 41 x 2
##   year max_wind_speed
##   <dbl>          <dbl>
## 1 1975           100
## 2 1976           105
## 3 1977           150
## 4 1978            80
## 5 1979           150
## 6 1980            90
```

```
## 7 1981      115
## 8 1982      115
## 9 1983      100
## 10 1984     115
## # ... with 31 more rows
```

For each status, find the average category, wind speed, pressure and diameters (do not allow the average to be NA).

```
storms %>%
  group_by(status) %>%
  summarise(avg_category = mean(as.numeric(as.character(category))), avg_wind_speed = mean(wind), avg_p
```

```
## Warning: Grouping rowwise data frame strips rowwise nature
```

```
## # A tibble: 3 x 6
##   status avg_category avg_wind_speed avg_pressure avg_ts_diameter
##   <fct>      <dbl>      <dbl>      <dbl>      <dbl>
## 1 hurri~    1.86        86.0        969.        288.
## 2 tropi~   -1         27.3       1008.         0
## 3 tropi~   0.000229    45.8        999.        160.
## # ... with 1 more variable: avg_hu_diameter <dbl>
```

For each named storm, find its maximum category, wind speed, pressure and diameters (do not allow the max to be NA) and the number of readings (i.e. observations).

```
storms %>%
  group_by(name) %>%
  summarise(max_category = max(as.numeric(as.character(category))), max_wind_speed = max(wind), max_pre
```

```
## Warning: Grouping rowwise data frame strips rowwise nature
```

```
## # A tibble: 198 x 6
##   name max_category max_wind_speed max_pressure max_ts_diameter
##   <fct>      <dbl>      <dbl>      <dbl>      <dbl>
## 1 AL01~    -1         30        1003      -Inf
## 2 AL01~    -1         25        1010      -Inf
## 3 AL02~    -1         30        1009      -Inf
## 4 AL02~    -1         30        1017      -Inf
## 5 AL02~    -1         30        1006      -Inf
## 6 AL02~    -1         30        1010      -Inf
## 7 AL02~    -1         25        1012      -Inf
## 8 AL02~    -1         30        1010      -Inf
## 9 AL02~     0         45        1008       69.0
## 10 AL03~     0         40        1015      -Inf
## # ... with 188 more rows, and 1 more variable: max_hu_diameter <dbl>
```

For each category, find its average wind speed, pressure and diameters (do not allow the max to be NA).

```
storms %>%
  group_by(category) %>%
  summarise(avg_wind_speed = mean(wind), avg_pressure = mean(pressure), avg_ts_diameter = mean(ts_diamete
```

```
## Warning: Grouping rowwise data frame strips rowwise nature
```

```
## # A tibble: 7 x 5
##   category avg_wind_speed avg_pressure avg_ts_diameter avg_hu_diameter
##   <ord>      <dbl>      <dbl>      <dbl>      <dbl>
## 1 -1         27.3       1008.         0         0
```

```
## 2 0          45.8          999.          160.          0
## 3 1          70.9          982.          278.          57.3
## 4 2          89.4          967.          282.          78.8
## 5 3         105.          954.          307.          91.4
## 6 4         122.          940.          315.         102.
## 7 5         145.          916.          317.         120.
```

for each named storm, find its duration in hours.

```
storms %>%
  group_by(name, status = 'tropical storm' ) %>%
  summarize(dur_hours = sum(hour))
```

```
## Warning: Grouping rowwise data frame strips rowwise nature
```

```
## # A tibble: 198 x 3
## # Groups:   name [198]
##   name      status      dur_hours
##   <fct>    <chr>      <dbl>
## 1 AL011993 tropical storm      72
## 2 AL012000 tropical storm      36
## 3 AL021992 tropical storm      48
## 4 AL021994 tropical storm      56
## 5 AL021999 tropical storm      28
## 6 AL022000 tropical storm     108
## 7 AL022001 tropical storm      54
## 8 AL022003 tropical storm      36
## 9 AL022006 tropical storm      42
## 10 AL031987 tropical storm     288
## # ... with 188 more rows
```

```
storms

## Source: local data frame [10,010 x 17]
## Groups: <by row>
##
## # A tibble: 10,010 x 17
##   lat long distance_to_mia~ name status category year month day
##   <dbl> <dbl>      <dbl> <fct> <fct> <ord>   <dbl> <dbl> <int>
## 1 23.1 -73.7      448. Joaq~ hurri~ 4      2015    10    1
## 2 23   -74.2      423. Joaq~ hurri~ 4      2015    10    1
## 3 22.9 -74.4      415. Joaq~ hurri~ 4      2015    10    2
## 4 23   -74.7      395. Joaq~ hurri~ 4      2015    10    2
## 5 23.4 -74.8      376. Joaq~ hurri~ 4      2015    10    2
## 6 24.3 -74.3      383. Joaq~ hurri~ 4      2015    10    3
## 7 24.8 -73.6      418. Joaq~ hurri~ 4      2015    10    3
## 8 25.4 -72.6      474. Joaq~ hurri~ 4      2015    10    3
## 9 26.3 -71       572. Joaq~ hurri~ 4      2015    10    3
## 10 27.4 -69.5      671. Joaq~ hurri~ 4      2015    10    4
## # ... with 10,000 more rows, and 8 more variables: hour <dbl>, wind <int>,
## #   pressure <int>, ts_diameter <dbl>, hu_diameter <dbl>,
## #   wind_speed_per_unit_pressure <dbl>, average_diameter <dbl>,
## #   decile_windspeed <fct>
```

For each named storm, find the distance from its starting position to ending position in kilometers.

Now we want to transition to building real design matrices for prediction. We want to predict the following: given the first three readings of a storm, can you predict its maximum wind speed? Identify the y and identify

which features you need x_1, \dots, x_p and build that matrix with `dplyr` functions. This is not easy, but it is what it's all about. Feel free to “featurize” (as Dana Chandler spoke about) as creatively as you would like. You aren't going to overfit if you only build a few features relative to the total 198 storms.

#TO-DO

Interactions in linear models

Load the Boston Housing Data from package `MASS` and use `str` and `summary` to remind yourself of the features and their types and then use `?MASS::Boston` to read an English description of the features.

```
data(Boston, package = "MASS")
str(Boston)

## 'data.frame':  506 obs. of  14 variables:
## $ crim   : num  0.00632 0.02731 0.02729 0.03237 0.06905 ...
## $ zn     : num  18 0 0 0 0 12.5 12.5 12.5 12.5 ...
## $ indus  : num  2.31 7.07 7.07 2.18 2.18 2.18 7.87 7.87 7.87 7.87 ...
## $ chas   : int   0 0 0 0 0 0 0 0 0 0 ...
## $ nox    : num  0.538 0.469 0.469 0.458 0.458 0.458 0.524 0.524 0.524 0.524 ...
## $ rm     : num  6.58 6.42 7.18 7 7.15 ...
## $ age    : num  65.2 78.9 61.1 45.8 54.2 58.7 66.6 96.1 100 85.9 ...
## $ dis    : num  4.09 4.97 4.97 6.06 6.06 ...
## $ rad    : int   1 2 2 3 3 3 5 5 5 5 ...
## $ tax    : num  296 242 242 222 222 222 311 311 311 311 ...
## $ ptratio: num  15.3 17.8 17.8 18.7 18.7 18.7 15.2 15.2 15.2 15.2 ...
## $ black  : num  397 397 393 395 397 ...
## $ lstat  : num  4.98 9.14 4.03 2.94 5.33 ...
## $ medv   : num  24 21.6 34.7 33.4 36.2 28.7 22.9 27.1 16.5 18.9 ...

summary(Boston)
```

```
##      crim              zn              indus              chas
## Min.   : 0.00632   Min.    : 0.00   Min.    : 0.46   Min.    :0.00000
## 1st Qu.: 0.08204   1st Qu.: 0.00   1st Qu.: 5.19   1st Qu.:0.00000
## Median : 0.25651   Median : 0.00   Median : 9.69   Median :0.00000
## Mean   : 3.61352   Mean    :11.36   Mean    :11.14   Mean    :0.06917
## 3rd Qu.: 3.67708   3rd Qu.:12.50   3rd Qu.:18.10   3rd Qu.:0.00000
## Max.   :88.97620   Max.    :100.00   Max.    :27.74   Max.    :1.00000
##      nox              rm              age              dis
## Min.   :0.3850   Min.    :3.561   Min.    : 2.90   Min.    : 1.130
## 1st Qu.:0.4490   1st Qu.:5.886   1st Qu.:45.02   1st Qu.: 2.100
## Median :0.5380   Median :6.208   Median :77.50   Median : 3.207
## Mean   :0.5547   Mean    :6.285   Mean    :68.57   Mean    : 3.795
## 3rd Qu.:0.6240   3rd Qu.:6.623   3rd Qu.:94.08   3rd Qu.: 5.188
## Max.   :0.8710   Max.    :8.780   Max.    :100.00   Max.    :12.127
##      rad              tax              ptratio              black
## Min.   : 1.000   Min.    :187.0   Min.    :12.60   Min.    : 0.32
## 1st Qu.: 4.000   1st Qu.:279.0   1st Qu.:17.40   1st Qu.:375.38
## Median : 5.000   Median :330.0   Median :19.05   Median :391.44
## Mean   : 9.549   Mean    :408.2   Mean    :18.46   Mean    :356.67
## 3rd Qu.:24.000   3rd Qu.:666.0   3rd Qu.:20.20   3rd Qu.:396.23
## Max.   :24.000   Max.    :711.0   Max.    :22.00   Max.    :396.90
##      lstat              medv
## Min.   : 1.73   Min.    : 5.00
```

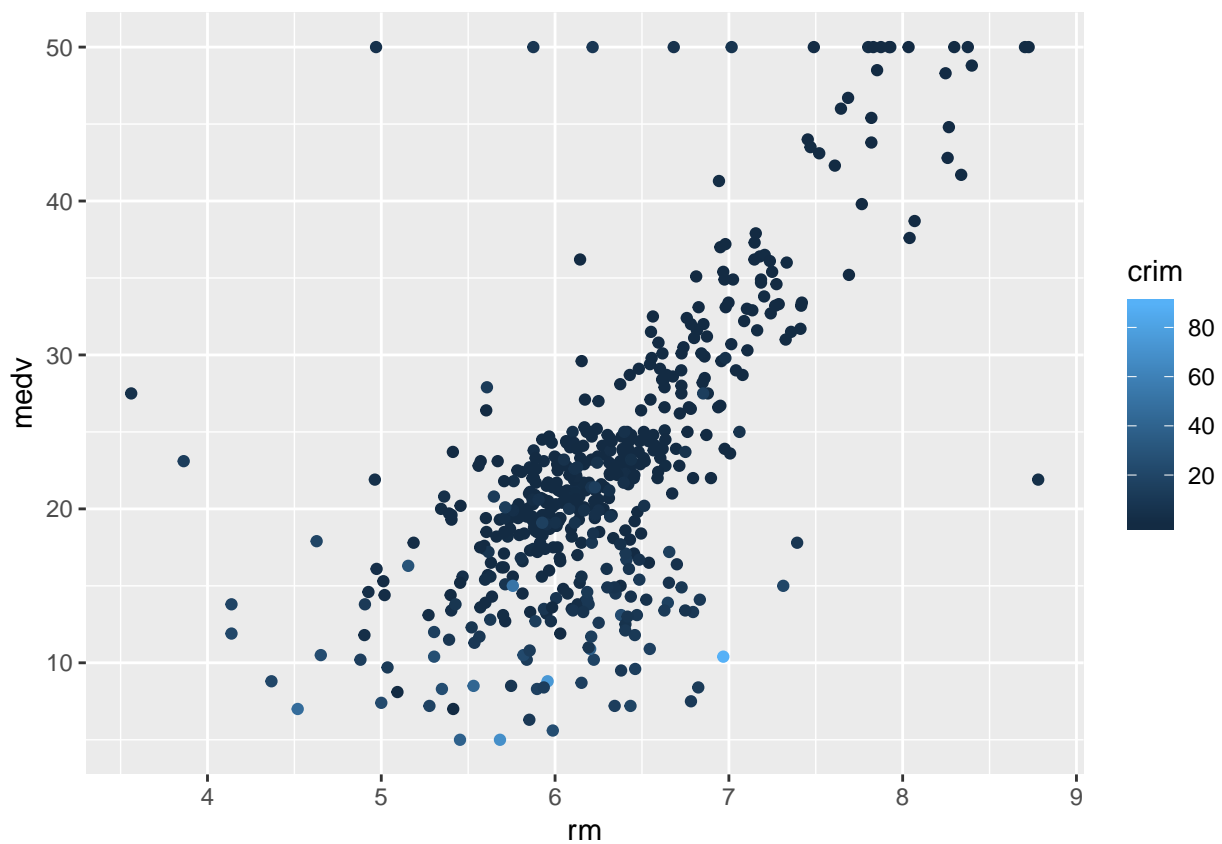
```
## 1st Qu.: 6.95    1st Qu.:17.02
## Median :11.36   Median :21.20
## Mean   :12.65   Mean   :22.53
## 3rd Qu.:16.95   3rd Qu.:25.00
## Max.    :37.97   Max.    :50.00
```

```
?MASS::Boston
```

```
## starting httpd help server ... done
```

Using your knowledge of the modeling problem, try to guess which features are interacting. Confirm using plots in ggplot that illustrate three (or more) features.

```
pacman::p_load(ggplot2)
base = ggplot(Boston, aes(x = rm, y = medv))
base + geom_point(aes(col = crim))
```



Once an interaction has been located, confirm the “non-linear linear” model with the interaction term does better than just the vanilla linear model.

```
mod = lm(medv ~ rm * crim, Boston)
coef(mod)
```

```
## (Intercept)      rm      crim  rm:crim
## -37.257338    9.651470    1.462943   -0.287657
```

```
mod_vanilla = lm(medv ~ rm + crim, Boston)
coef(mod_vanilla)
```

```
## (Intercept)      rm      crim
```

```
## -29.2447195    8.3910682   -0.2649133
```

```
summary(mod_vanilla)$r.squared
```

```
## [1] 0.5419592
```

```
summary(mod_vanilla)$sigma
```

```
## [1] 6.236844
```

```
summary(mod)$r.squared
```

```
## [1] 0.5814763
```

```
summary(mod)$sigma
```

```
## [1] 5.967672
```

Repeat this procedure for another interaction with two different features (not used in the previous interaction you found) and verify.

```
mod = lm(medv ~ rm * zn, Boston)
coef(mod)
```

```
## (Intercept)          rm          zn          rm:zn
## -26.9934476    7.7661501  -0.4697937    0.0791624
```

```
mod_vanilla = lm(medv ~ rm + zn, Boston)
summary(mod_vanilla)$r.squared
```

```
## [1] 0.5063381
```

```
summary(mod_vanilla)$sigma
```

```
## [1] 6.474818
```

```
summary(mod)$r.squared
```

```
## [1] 0.5223732
```

```
summary(mod)$sigma
```

```
## [1] 6.375133
```

Fit a model using all possible first-order interactions. Verify it is “better” than the linear model. Do you think you overfit? Why or why not?

```
mod2 = lm(medv ~ .*. , Boston)
summary(mod2)
```

```
##
```

```
## Call:
```

```
## lm(formula = medv ~ . * ., data = Boston)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -7.9374 -1.5344 -0.1068   1.2973  17.8500
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.579e+02  6.800e+01  -2.323 0.020683 *
## crim        -1.707e+01  6.554e+00  -2.605 0.009526 **
```

## zn	-7.529e-02	4.580e-01	-0.164	0.869508	
## indus	-2.819e+00	1.696e+00	-1.663	0.097111	.
## chas	4.451e+01	1.952e+01	2.280	0.023123	*
## nox	2.006e+01	7.516e+01	0.267	0.789717	
## rm	2.527e+01	5.699e+00	4.435	1.18e-05	***
## age	1.263e+00	2.728e-01	4.630	4.90e-06	***
## dis	-1.698e+00	4.604e+00	-0.369	0.712395	
## rad	1.861e+00	2.464e+00	0.755	0.450532	
## tax	3.670e-02	1.440e-01	0.255	0.798978	
## ptratio	2.725e+00	2.850e+00	0.956	0.339567	
## black	9.942e-02	7.468e-02	1.331	0.183833	
## lstat	1.656e+00	8.533e-01	1.940	0.053032	.
## crim:zn	4.144e-01	1.804e-01	2.297	0.022128	*
## crim:indus	-4.693e-02	4.480e-01	-0.105	0.916621	
## crim:chas	2.428e+00	5.710e-01	4.251	2.63e-05	***
## crim:nox	-1.108e+00	9.285e-01	-1.193	0.233425	
## crim:rm	2.163e-01	4.907e-02	4.409	1.33e-05	***
## crim:age	-3.083e-03	3.781e-03	-0.815	0.415315	
## crim:dis	-1.903e-01	1.060e-01	-1.795	0.073307	.
## crim:rad	-6.584e-01	5.815e-01	-1.132	0.258198	
## crim:tax	3.479e-02	4.287e-02	0.812	0.417453	
## crim:ptratio	4.915e-01	3.328e-01	1.477	0.140476	
## crim:black	-4.612e-04	1.793e-04	-2.572	0.010451	*
## crim:lstat	2.964e-02	6.544e-03	4.530	7.72e-06	***
## zn:indus	-6.731e-04	4.651e-03	-0.145	0.885000	
## zn:chas	-5.230e-02	6.450e-02	-0.811	0.417900	
## zn:nox	1.998e-03	4.721e-01	0.004	0.996625	
## zn:rm	-7.286e-04	2.602e-02	-0.028	0.977672	
## zn:age	-1.249e-06	8.514e-04	-0.001	0.998830	
## zn:dis	1.097e-02	7.550e-03	1.452	0.147121	
## zn:rad	-3.200e-03	6.975e-03	-0.459	0.646591	
## zn:tax	3.937e-04	1.783e-04	2.209	0.027744	*
## zn:ptratio	-4.578e-03	7.015e-03	-0.653	0.514325	
## zn:black	1.159e-04	7.599e-04	0.153	0.878841	
## zn:lstat	-1.064e-02	4.662e-03	-2.281	0.023040	*
## indus:chas	-3.672e-01	3.780e-01	-0.971	0.331881	
## indus:nox	3.138e+00	1.449e+00	2.166	0.030855	*
## indus:rm	3.301e-01	1.327e-01	2.488	0.013257	*
## indus:age	-4.865e-04	3.659e-03	-0.133	0.894284	
## indus:dis	-4.486e-02	6.312e-02	-0.711	0.477645	
## indus:rad	-2.089e-02	5.020e-02	-0.416	0.677560	
## indus:tax	3.129e-04	6.034e-04	0.519	0.604322	
## indus:ptratio	-6.011e-02	3.783e-02	-1.589	0.112820	
## indus:black	1.122e-03	2.034e-03	0.552	0.581464	
## indus:lstat	5.063e-03	1.523e-02	0.332	0.739789	
## chas:nox	-3.272e+01	1.243e+01	-2.631	0.008820	**
## chas:rm	-5.384e+00	1.150e+00	-4.681	3.87e-06	***
## chas:age	3.040e-02	5.840e-02	0.521	0.602982	
## chas:dis	9.022e-01	1.334e+00	0.676	0.499143	
## chas:rad	-7.773e-01	5.707e-01	-1.362	0.173907	
## chas:tax	4.627e-02	3.645e-02	1.270	0.204930	
## chas:ptratio	-6.145e-01	6.914e-01	-0.889	0.374604	
## chas:black	2.500e-02	1.567e-02	1.595	0.111423	
## chas:lstat	-2.980e-01	1.845e-01	-1.615	0.107008	

```
## nox:rm      5.990e+00  5.468e+00   1.095 0.273952
## nox:age     -7.273e-01  2.340e-01  -3.108 0.002012 **
## nox:dis      5.694e+00  3.723e+00   1.529 0.126969
## nox:rad     -1.994e-01  1.897e+00  -0.105 0.916360
## nox:tax     -2.793e-02  1.312e-01  -0.213 0.831559
## nox:ptratio -3.669e+00  3.096e+00  -1.185 0.236648
## nox:black   -1.854e-02  3.615e-02  -0.513 0.608298
## nox:lstat    1.119e+00  6.511e-01   1.719 0.086304 .
## rm:age     -6.277e-02  2.203e-02  -2.849 0.004606 **
## rm:dis      3.190e-01  3.295e-01   0.968 0.333516
## rm:rad     -8.422e-02  1.527e-01  -0.552 0.581565
## rm:tax     -2.242e-02  9.910e-03  -2.262 0.024216 *
## rm:ptratio  -4.880e-01  2.172e-01  -2.247 0.025189 *
## rm:black   -4.528e-03  3.351e-03  -1.351 0.177386
## rm:lstat   -2.968e-01  4.316e-02  -6.878 2.24e-11 ***
## age:dis    -1.678e-02  8.882e-03  -1.889 0.059589 .
## age:rad     1.442e-02  4.212e-03   3.423 0.000682 ***
## age:tax    -3.403e-04  2.187e-04  -1.556 0.120437
## age:ptratio -7.520e-03  6.793e-03  -1.107 0.268946
## age:black  -7.029e-04  2.136e-04  -3.291 0.001083 **
## age:lstat  -6.023e-03  1.936e-03  -3.111 0.001991 **
## dis:rad    -5.580e-02  7.075e-02  -0.789 0.430678
## dis:tax    -3.882e-03  2.496e-03  -1.555 0.120623
## dis:ptratio -4.786e-02  9.983e-02  -0.479 0.631920
## dis:black  -5.194e-03  5.541e-03  -0.937 0.349116
## dis:lstat   1.350e-01  4.866e-02   2.775 0.005774 **
## rad:tax     3.131e-05  1.446e-03   0.022 0.982729
## rad:ptratio -4.379e-02  8.392e-02  -0.522 0.602121
## rad:black  -4.362e-04  2.518e-03  -0.173 0.862561
## rad:lstat  -2.529e-02  1.816e-02  -1.392 0.164530
## tax:ptratio  7.854e-03  2.504e-03   3.137 0.001830 **
## tax:black  -4.785e-07  1.999e-04  -0.002 0.998091
## tax:lstat  -1.403e-03  1.208e-03  -1.162 0.245940
## ptratio:black 1.203e-03  3.361e-03   0.358 0.720508
## ptratio:lstat 3.901e-03  2.985e-02   0.131 0.896068
## black:lstat -6.118e-04  4.157e-04  -1.472 0.141837
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.852 on 414 degrees of freedom
## Multiple R-squared:  0.9212, Adjusted R-squared:  0.9039
## F-statistic: 53.18 on 91 and 414 DF,  p-value: < 2.2e-16

mod_vanilla = lm(medv ~ rm + zn + crim, Boston)
summary(mod_vanilla)$r.squared

## [1] 0.5558362

summary(mod_vanilla)$sigma

## [1] 6.147755

summary(mod2)$r.squared

## [1] 0.9211876
```

```
summary(mod2)$sigma
```

```
## [1] 2.851634
```

```
## mod2 is better , we are not overfitting since our n is 506 and our p = 91
```

CV

Use 5-fold CV to estimate the generalization error of the model with all interactions.

```
pacman::p_load(mlr)
```

```
library(mlr)
```

```
modeling_task = makeRegrTask(data = Boston, target = "medv") #instantiate the task
```

```
algorithm = makeLearner("regr.lm") #instantiate the OLS learner algorithm on the diamonds dataset and s
```

```
validation = makeResampleDesc("CV", iters = 5) #instantiate the 5-fold CV
```

```
resample(algorithm, modeling_task, validation)
```

```
## Resampling: cross-validation
```

```
## Measures:           mse
```

```
## [Resample] iter 1:   18.8164429
```

```
## [Resample] iter 2:   20.6391585
```

```
## [Resample] iter 3:   29.8384484
```

```
## [Resample] iter 4:   29.1942087
```

```
## [Resample] iter 5:   23.3788547
```

```
##
```

```
## Aggregated Result: mse.test.mean=24.3734226
```

```
##
```

```
## Resample Result
```

```
## Task: Boston
```

```
## Learner: regr.lm
```

```
## Aggr perf: mse.test.mean=24.3734226
```

```
## Runtime: 0.050518
```