Lecture 12: Data wrangling

Last time

- filter: choose certain rows
- summarize: calculate summary statistics
- group_by: group rows together
- mutate: create new columns

Data for today

- Data on professional baseball teams between 1871 and 2022
- 3015 rows and 48 columns
- Each row represents one year (season) for one team
- Variables include:
 - yearID: Year
 - franchID: Franchise
 - W: Wins
 - L: Losses

Data for today

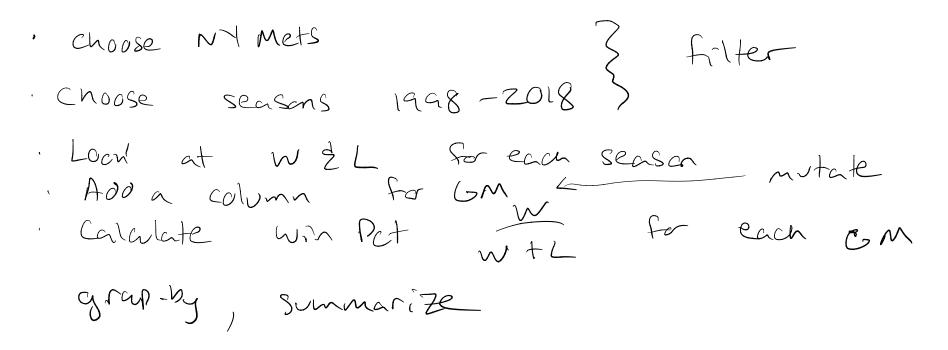
- Variables include:
 - yearID: Year
 - franchID: Franchise
 - W: Wins
 - L: Losses

We want to know: which NY Mets general manager performed best between 1998 - 2018

Making a plan

We want to know: which NY Mets general manager performed best between 1998 - 2018

Question: What steps could we take to answer this question?



Step 0: Make the columns more manageable

There are 48 columns in the initial data! Let's only focus on the ones we care about:

```
columns to Heep
                    select
    Teams(|>
      select(yearID, franchID, W, L)
     yearID franchID
                             T.
       1871
                  BNA
                            10
1
                       20
       1871
                  CNA
                       19
3
                  CFC
       1871
                       10
                           19
       1871
                  KEK
                           12
       1871
                  NNA
                       16
                           17
       1871
                  PNA
                       21
       1871
                  ROK
                        4
                            21
8
       1871
                            15
                  TRO
                       13
       1871
                  OLX
                       15
                           15
10
       1872
                  BLC
                       35
                           19
11
       1872
                  ECK
                           26
12
       1872
                  BRA
                            28
       1872
13
                             Q
                  BMD
                        39
```

Step 1: Focus on the Mets

```
1 Teams |>
2 select(yearID, FranchID, W, L) |>
3 ...(franchID == "NYM")
```

What function do I use to choose only the rows corresponding to the Mets?

Step 1: Focus on the Mets

NYM 100

83

83

83

82

71

NYM

NYM

NYM

NYM

NYM

62

79

79

73

79

91

8

9

10

11

12

13

1969

1970

1971

1972

1973

1974

```
Teams |>
      select(yearID, franchID, W, L) |>
      filter(franchID == "NYM")
   yearID franchID
                           \mathbf{L}
     1962
                     40 120
1
                MYM
     1963
                     51 111
                NYM
     1964
                NYM
                     53 109
4
     1965
                NYM
                     50 112
5
     1966
                     66 95
                NYM
                     61 101
6
     1967
                NYM
     1968
                     73
                         89
                NYM
```

Step 2: Focus on the Mets between 1998 and 2018

```
1 Teams |>
2 select(yearID, franchID, W, L) |>
3 filter(franchID == "NYM",
4 ...)
```

How do I specify the range of years I want?

Step 2: Focus on the Mets between 1998 and 2018

```
Teams >
      select(yearID, franchID, W, L) |>
      filter(franchID == "NYM",
             yearID >= 1998, yearID <= 2018)</pre>
   yearID franchID
                               requivalent to

FranchD == "WIM"
1
     1998
               NYM 88 74
               NYM 97 66
     1999
              NYM 94 68
    2000
                                Teas 10>=1998 & year 10 <= 2018
    2001
              NYM 82 80
    2002
              NYM 75 86
6
    2003
              NYM 66 95
    2004
              NYM 71 91
     2005
               NYM 83 79
     2006
9
               NYM 97 65
10
    2007
               NYM 88 74
11
    2008
               NYM 89 73
12
    2009
               NYM 70 92
13
     2010
               NYM 79 83
```

Step 3: Who was the GM?

- 1998 2003: Steve Phillips
- 2004: Jim Duquette
- 2005 2010: Omar Minaya
- 2011 2018: Sandy Alderson

How should we add this information to the data?

Step 3: Who was the GM?

```
yearID franchID
                   WL
                                qm
     1998
               NYM 88 74 Phillips
1
               NYM 97 66 Phillips
     1999
               NYM 94 68 Phillips
3
     2000
               NYM 82 80 Phillips
4
     2001
    2002
               NYM 75 86 Phillips
               NYM 66 95 Phillips
6
     2003
     2004
               NYM 71 91 Duquette
8
     2005
               NYM 83 79
                          Minaya
     2006
9
               NYM 97 65
                         Minaya
10
     2007
               NYM 88 74
                           Minaya
11
    2008
               NYM 89 73
                          Minaya
     2009
                          Minaya
12
               NYM 70 92
13
     2010
               NYM 79 83
                           Minaya
```

Step 4: Summarize performance

How do I calculate performance for each GM?

wpct 1 0.5019112

Step 4: Summarize performance

```
Teams |>
      select(yearID, franchID, W, L) |>
      filter(franchID == "NYM",
             yearID >= 1998, yearID <= 2018) |>
 4
      mutate(gm = case when(
        yearID <= 2003 ~ "Phillips",</pre>
     yearID == 2004 ~ "Duquette",
   yearID \leq 2010 \sim \text{"Minaya"}
      yearID <= 2018 ~ "Alderson"</pre>
10
    )) |>
11
   group by(gm) >
      summarize(wpct = sum(W)/sum(W + L))
12
# A tibble: 4 \times 2
      wpct
  qm
  <chr> <dbl>
1 Alderson 0.485
2 Duquette 0.438
3 Minaya 0.521
4 Phillips 0.517
```

Finally: arrange results

```
Teams |>
      select(yearID, franchID, W, L) |>
      filter(franchID == "NYM",
             yearID >= 1998, yearID <= 2018) |>
      mutate(qm = case when(
        yearID <= 2003 ~ "Phillips",</pre>
       yearID == 2004 ~ "Duguette",
     yearID \le 2010 \sim "Minaya",
       yearID <= 2018 ~ "Alderson"</pre>
10
      )) |>
      group by (qm) |>
      summarize(wpct = sum(W)/sum(W + L)) |>
12
      arrange(desc(wpct))
# A tibble: 4 \times 2
          wpct
                        serts the rows by one or more columns in the data
  gm
  <chr>
          <dbl>
1 Minaya
          0.521
                        arrange (upct)
2 Phillips 0.517
                                                     (larst to highest upct)
3 Alderson 0.485
4 Duquette 0.438
                                                    (highest to lowest wast
                        arrange (desc(wpct))
```

Class activity

https://sta279-

f23.github.io/class_activities/ca_lecture_12.html