# Stat 652 Project

*Vincent Chiu*

*2019-10-16*

## Loading Libraries

```
library(tidyverse)
```

```
## -- Attaching packages ------------------------------------------------------------------ tidyverse 1
```

```
## v ggplot2 3.2.1      v purrr   0.3.3
## v tibble  2.1.3      v dplyr   0.8.3
## v tidyr   1.0.0      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.4.0
```

```
## -- Conflicts --------------------------------------------------------------------------- tidyverse_conflict
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

## Loading the data

```
library(nycflights13)
library(Hmisc)
```

```
## Loading required package: lattice
```

```
## Loading required package: survival
```

```
## Loading required package: Formula
```

```
##
## Attaching package: 'Hmisc'
```

```
## The following objects are masked from 'package:dplyr':
##
##     src, summarize
```

```
## The following objects are masked from 'package:base':
##
##     format.pval, units
```

```
set.seed(42)
original_data <- read_csv("fltrain.csv.gz")
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   carrier = col_character(),
##   tailnum = col_character(),
##   origin = col_character(),
##   dest = col_character(),
##   time_hour = col_datetime(format = ""),
##   name = col_character(),
##   dst = col_character(),
```

```
##   tzone = col_character(),
##   type = col_character(),
##   manufacturer = col_character(),
##   model = col_character(),
##   engine = col_character()
## )

## See spec(...) for full column specifications.
DF <- original_data
```

## turning all columns with datatype characters to factors.

```
DF[sapply(DF, is.character)] <- lapply(DF[sapply(DF, is.character)],
                                       as.factor)
DF$flight <- as.factor(DF$flight)
str(DF)

## Classes 'spec_tbl_df', 'tbl_df', 'tbl' and 'data.frame': 200000 obs. of  43 variables:
##  $ year.x        : num  2013 2013 2013 2013 2013 ...
##  $ month         : num  11 10 12 11 10 11 9 12 11 3 ...
##  $ day           : num  7 30 18 20 21 7 29 21 7 31 ...
##  $ dep_time      : num  600 1252 1723 2029 1620 ...
##  $ sched_dep_time: num  600 1250 1715 2030 1625 ...
##  $ dep_delay     : num  0 2 8 -1 -5 -8 -10 -4 0 -8 ...
##  $ arr_time      : num  826 1356 2008 2141 1818 ...
##  $ sched_arr_time: num  825 1400 2020 2205 1831 ...
##  $ arr_delay     : num  1 -4 -12 -24 -13 -18 -10 -16 4 -11 ...
##  $ carrier       : Factor w/ 16 levels "9E","AA","AS",..: 15 2 5 15 5 4 6 6 1 13 ...
##  $ flight        : Factor w/ 3672 levels "1","2","3","4",..: 1525 147 1400 2343 1860 24 3083 3351 207
##  $ tailnum       : Factor w/ 3957 levels "D942DN","N0EGMQ",..: 1437 1226 836 565 756 2459 204 2890 64
##  $ origin        : Factor w/ 3 levels "EWR","JFK","LGA": 3 2 3 1 3 1 1 3 2 3 ...
##  $ dest          : Factor w/ 104 levels "ABQ","ACK","ALB",..: 5 12 54 55 33 54 59 59 27 29 ...
##  $ air_time      : num  123 44 133 107 90 136 110 118 101 47 ...
##  $ distance      : num  762 187 950 711 502 937 725 738 589 214 ...
##  $ hour          : num  6 12 17 20 16 9 15 15 16 17 ...
##  $ minute        : num  0 50 15 30 25 0 29 30 50 0 ...
##  $ time_hour     : POSIXct, format: "2013-11-07 11:00:00" "2013-10-30 16:00:00" ...
##  $ temp          : num  63 59 34 37 63 ...
##  $ dewp          : num  55.9 46.9 17.1 18 41 ...
##  $ humid         : num  77.8 64.2 49.5 45.6 44.5 ...
##  $ wind_dir      : num  210 240 270 20 160 240 180 190 320 140 ...
##  $ wind_speed    : num  13.81 9.21 17.26 5.75 13.81 ...
##  $ wind_gust     : num  NA NA 21.9 NA NA ...
##  $ precip        : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ pressure      : num  1011 1025 1020 1036 1017 ...
##  $ visib         : num  10 10 10 10 10 10 10 10 10 10 ...
##  $ name          : Factor w/ 100 levels "Akron Canton Regional Airport",..: 37 31 67 17 26 67 32 32
##  $ lat           : num  33.6 42.4 28.4 41.8 42.2 ...
##  $ lon           : num  -84.4 -71 -81.3 -87.8 -83.4 ...
##  $ alt           : num  1026 19 96 620 645 ...
##  $ tz            : num  -5 -5 -5 -6 -5 -5 -6 -6 -5 -5 ...
##  $ dst           : Factor w/ 2 levels "A","N": 1 1 1 1 1 1 1 1 1 1 ...
```

```
##  $ tzone        : Factor w/ 7 levels "America/Anchorage",..: 5 5 5 2 5 5 2 2 5 5 ...
##  $ year.y       : num   2001 NA 2002 2006 1992 ...
##  $ type         : Factor w/ 3 levels "Fixed wing multi engine",..: 1 NA 1 1 1 1 1 1 1 1 ...
##  $ manufacturer : Factor w/ 35 levels "AGUSTA SPA","AIRBUS",..: 10 NA 2 10 3 2 18 11 11 3 ...
##  $ model        : Factor w/ 126 levels "150","172E","172M",..: 37 NA 80 37 84 88 106 98 99 79 ...
##  $ engines      : num   2 NA 2 2 2 2 2 2 2 2 ...
##  $ seats        : num   140 NA 145 140 182 200 55 80 95 179 ...
##  $ speed        : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ engine       : Factor w/ 6 levels "4 Cycle","Reciprocating",..: 3 NA 3 3 4 3 3 3 3 3 ...
##  - attr(*, "spec")=
##   .. cols(
##   ..   year.x = col_double(),
##   ..   month = col_double(),
##   ..   day = col_double(),
##   ..   dep_time = col_double(),
##   ..   sched_dep_time = col_double(),
##   ..   dep_delay = col_double(),
##   ..   arr_time = col_double(),
##   ..   sched_arr_time = col_double(),
##   ..   arr_delay = col_double(),
##   ..   carrier = col_character(),
##   ..   flight = col_double(),
##   ..   tailnum = col_character(),
##   ..   origin = col_character(),
##   ..   dest = col_character(),
##   ..   air_time = col_double(),
##   ..   distance = col_double(),
##   ..   hour = col_double(),
##   ..   minute = col_double(),
##   ..   time_hour = col_datetime(format = ""),
##   ..   temp = col_double(),
##   ..   dewp = col_double(),
##   ..   humid = col_double(),
##   ..   wind_dir = col_double(),
##   ..   wind_speed = col_double(),
##   ..   wind_gust = col_double(),
##   ..   precip = col_double(),
##   ..   pressure = col_double(),
##   ..   visib = col_double(),
##   ..   name = col_character(),
##   ..   lat = col_double(),
##   ..   lon = col_double(),
##   ..   alt = col_double(),
##   ..   tz = col_double(),
##   ..   dst = col_character(),
##   ..   tzone = col_character(),
##   ..   year.y = col_double(),
##   ..   type = col_character(),
##   ..   manufacturer = col_character(),
##   ..   model = col_character(),
##   ..   engines = col_double(),
##   ..   seats = col_double(),
##   ..   speed = col_double(),
##   ..   engine = col_character()
```

```
##   .. )
```

# Methods

## Preprocessing

Data preprocessing steps include the following: - Dropping columns that contain data from after the planes' departure which may leak information about the response variable dep_delay. - Dropping columns with too many NAs. - Impute NAs for the remaining columns. - Scaling the data to work well with methods like lasso regression.

## - Dropping columns that contain data from after the planes' departure which may leak information about the response variable dep_delay.

dropping the columns "dep_time", "arr_time", "air_time", "arr_delay", because that leaks the response variable. dropping column "year.x" because all the values are 2013 dropping tailnum because it produces too many dummy variable columns for one hot encoding.

```r
library(lubridate)
```

```
##
## Attaching package: 'lubridate'

## The following object is masked from 'package:base':
##
##     date
```

```r
DF$sched_arr_time_posix <- as.POSIXct(str_pad(as.character(DF$sched_arr_time), 4, pad="0"),format="%H%M
DF$sched_arr_time_hour <- hour(DF$sched_arr_time_posix)
DF$sched_arr_time_minute <- minute(DF$sched_arr_time_posix)

#num minute is number of minutes since start of day for scheduled arrival time
DF$sched_arr_time_num_minute <- 60*DF$sched_arr_time_hour + DF$sched_arr_time_minute

DF$sched_dep_time_posix <- as.POSIXct(str_pad(as.character(DF$sched_dep_time),4 , pad="0"),format="%H%M
DF$sched_dep_time_hour <- hour(DF$sched_dep_time_posix)
DF$sched_dep_time_minute <- minute(DF$sched_dep_time_posix)
#num minute is number of minutes since start of day for scheduled depival time
DF$sched_dep_time_num_minute <- 60*DF$sched_dep_time_hour + DF$sched_dep_time_minute
```

```r
select(original_data, time_hour, sched_dep_time, sched_arr_time, tz, tzone)
```

```
## # A tibble: 200,000 x 5
##    time_hour           sched_dep_time sched_arr_time    tz tzone
##    <dttm>                       <dbl>          <dbl> <dbl> <chr>
## 1 2013-11-07 11:00:00            600            825    -5 America/New_York
## 2 2013-10-30 16:00:00           1250           1400    -5 America/New_York
## 3 2013-12-18 22:00:00           1715           2020    -5 America/New_York
## 4 2013-11-21 01:00:00           2030           2205    -6 America/Chicago
## 5 2013-10-21 20:00:00           1625           1831    -5 America/New_York
## 6 2013-11-07 14:00:00            900           1157    -5 America/New_York
## 7 2013-09-29 19:00:00           1529           1649    -6 America/Chicago
## 8 2013-12-21 20:00:00           1530           1710    -6 America/Chicago
```

```
##  9 2013-11-07 21:00:00                1650              1906      -5 America/New_York
## 10 2013-03-31 21:00:00                1700              1821      -5 America/New_York
## # ... with 199,990 more rows
```

```r
select(DF, sched_arr_time, sched_arr_time_hour)
```

```
## # A tibble: 200,000 x 2
##    sched_arr_time sched_arr_time_hour
##             <dbl>               <int>
##  1            825                   8
##  2           1400                  14
##  3           2020                  20
##  4           2205                  22
##  5           1831                  18
##  6           1157                  11
##  7           1649                  16
##  8           1710                  17
##  9           1906                  19
## 10           1821                  18
## # ... with 199,990 more rows
```

```r
DF$sched_air_time <- DF$sched_arr_time_posix - DF$sched_dep_time_posix
drops <- c('sched_arr_time_posix', 'sched_arr_time_hour', 'sched_dep_time_posix', 'sched_dep_time_hour'
DF <- DF[ , !(names(DF) %in% drops)]
```

```r
drops <- c("dep_time", "arr_time", "air_time", "arr_delay", "year.x", 'tailnum')
DF <- DF[ , !(names(DF) %in% drops)]
```

```r
DF
```

```
## # A tibble: 200,000 x 37
##    month   day dep_delay carrier flight origin dest  distance  temp  dewp humid
##    <dbl> <dbl>     <dbl> <fct>   <fct>  <fct>  <fct>    <dbl> <dbl> <dbl> <dbl>
##  1    11     7         0 WN      1716   LGA    ATL        762  63.0  55.9  77.8
##  2    10    30         2 AA      178    JFK    BOS        187  59    46.9  64.2
##  3    12    18         8 DL      1585   LGA    MCO        950  34.0  17.1  49.5
##  4    11    20        -1 WN      3494   EWR    MDW        711  37.0  18.0  45.6
##  5    10    21        -5 DL      2231   LGA    DTW        502  63.0  41    44.5
##  6    11     7        -8 B6      27     EWR    MCO        937  64.4  55.4  77.3
##  7     9    29       -10 EV      4580   EWR    MKE        725  69.1  53.1  56.7
##  8    12    21        -4 EV      5207   LGA    MKE        738  57.9  46.0  64.5
##  9    11     7         0 9E      2910   JFK    CVG        589  53.6  48.2  81.9
## 10     3    31        -8 US      2183   LGA    DCA        214  51.1  36.0  56.0
## # ... with 199,990 more rows, and 26 more variables: wind_dir <dbl>,
## #   wind_speed <dbl>, wind_gust <dbl>, precip <dbl>, pressure <dbl>,
## #   visib <dbl>, name <fct>, lat <dbl>, lon <dbl>, alt <dbl>, tz <dbl>,
## #   dst <fct>, tzone <fct>, year.y <dbl>, type <fct>, manufacturer <fct>,
## #   model <fct>, engines <dbl>, seats <dbl>, speed <dbl>, engine <fct>,
## #   sched_arr_time_minute <int>, sched_arr_time_num_minute <dbl>,
## #   sched_dep_time_minute <int>, sched_dep_time_num_minute <dbl>,
## #   sched_air_time <drtn>
```

```r
## Remove columns with more than 50% NA
DF <- DF[, -which(colMeans(is.na(DF)) > 0.5)]
```

```r
DF$sched_air_time <- as.numeric(DF$sched_air_time)
library(imputeMissings)
```

```
## 
## Attaching package: 'imputeMissings'

## The following object is masked from 'package:Hmisc':
## 
##     impute

## The following object is masked from 'package:dplyr':
## 
##     compute
```

```r
impute_model <- imputeMissings::compute(DF, method="median/mode")
impute_model
```

```
## $month
## [1] 7
## 
## $day
## [1] 16
## 
## $dep_delay
## [1] -2
## 
## $carrier
## [1] "UA"
## 
## $flight
## [1] "15"
## 
## $origin
## [1] "EWR"
## 
## $dest
## [1] "ATL"
## 
## $distance
## [1] 872
## 
## $temp
## [1] 57.2
## 
## $dewp
## [1] 42.8
## 
## $humid
## [1] 57.69
## 
## $wind_dir
## [1] 220
## 
## $wind_speed
## [1] 10.35702
## 
## $precip
## [1] 0
## 
```

```
## $pressure
## [1] 1017.5
##
## $visib
## [1] 10
##
## $name
## [1] "Hartsfield Jackson Atlanta Intl"
##
## $lat
## [1] 36.09775
##
## $lon
## [1] -83.35339
##
## $alt
## [1] 433
##
## $tz
## [1] -5
##
## $dst
## [1] "A"
##
## $tzone
## [1] "America/New_York"
##
## $year.y
## [1] 2002
##
## $type
## [1] "Fixed wing multi engine"
##
## $manufacturer
## [1] "BOEING"
##
## $model
## [1] "A320-232"
##
## $engines
## [1] 2
##
## $seats
## [1] 149
##
## $engine
## [1] "Turbo-fan"
##
## $sched_arr_time_minute
## [1] 30
##
## $sched_arr_time_num_minute
## [1] 957
##
```

```
## $sched_dep_time_minute
## [1] 29
##
## $sched_dep_time_num_minute
## [1] 839
##
## $sched_air_time
## [1] 139
```

```
DF <- impute(DF, object=impute_model, , flag=TRUE)
```

```
numeric_only_df <- dplyr::select_if(DF, is.numeric)
library(corrplot)
```

```
## corrplot 0.84 loaded
```

```
corrplot(cor(numeric_only_df), type = 'lower')
```



## try features scaling

```
dep_delay_vec <- DF$dep_delay
DF$dep_delay <- NULL
head(DF)
```

```
##   month day carrier flight origin dest distance  temp  dewp humid wind_dir
## 1    11   7      WN   1716    LGA  ATL      762 62.96 55.94 77.83      210
## 2    10  30      AA    178    JFK  BOS      187 59.00 46.94 64.22      240
```

```
## 3    12   18        DL   1585   LGA   MCO       950 33.98 17.06 49.51      270
## 4    11   20        WN   3494   EWR   MDW       711 37.04 17.96 45.58       20
## 5    10   21        DL   2231   LGA   DTW       502 62.96 41.00 44.47      160
## 6    11    7        B6     27   EWR   MCO       937 64.40 55.40 77.29      240
##    wind_speed precip pressure visib                              name      lat
## 1   13.80936      0   1011.0    10     Hartsfield Jackson Atlanta Intl 33.63672
## 2    9.20624      0   1024.9    10 General Edward Lawrence Logan Intl 42.36435
## 3   17.26170      0   1019.8    10                        Orlando Intl 28.42939
## 4    5.75390      0   1035.6    10                 Chicago Midway Intl 41.78597
## 5   13.80936      0   1016.9    10              Detroit Metro Wayne Co 42.21244
## 6   16.11092      0   1017.5    10                        Orlando Intl 28.42939
##        lon  alt tz dst          tzone year.y                       type
## 1 -84.42807 1026 -5   A America/New_York   2001 Fixed wing multi engine
## 2 -71.00518   19 -5   A America/New_York   2002 Fixed wing multi engine
## 3 -81.30899   96 -5   A America/New_York   2002 Fixed wing multi engine
## 4 -87.75242  620 -6   A  America/Chicago   2006 Fixed wing multi engine
## 5 -83.35339  645 -5   A America/New_York   1992 Fixed wing multi engine
## 6 -81.30899   96 -5   A America/New_York   2006 Fixed wing multi engine
##        manufacturer        model engines seats     engine sched_arr_time_minute
## 1            BOEING      737-7H4       2   140 Turbo-fan                     25
## 2            BOEING     A320-232       2   149 Turbo-fan                      0
## 3            AIRBUS     A319-114       2   145 Turbo-fan                     20
## 4            BOEING      737-7H4       2   140 Turbo-fan                      5
## 5 AIRBUS INDUSTRIE     A320-211       2   182 Turbo-jet                     31
## 6            AIRBUS     A320-232       2   200 Turbo-fan                     57
##   sched_arr_time_num_minute sched_dep_time_minute sched_dep_time_num_minute
## 1                       505                     0                       360
## 2                       840                    50                       770
## 3                      1220                    15                      1035
## 4                      1325                    30                      1230
## 5                      1111                    25                       985
## 6                       717                     0                       540
##   sched_air_time dep_delay_flag temp_flag dewp_flag humid_flag wind_dir_flag
## 1            145              0         0         0          0             0
## 2             70              0         0         0          0             0
## 3            185              0         0         0          0             0
## 4             95              0         0         0          0             0
## 5            126              0         0         0          0             0
## 6            177              0         0         0          0             0
##   wind_speed_flag precip_flag pressure_flag visib_flag name_flag lat_flag
## 1               0           0             0          0         0        0
## 2               0           0             0          0         0        0
## 3               0           0             0          0         0        0
## 4               0           0             0          0         0        0
## 5               0           0             0          0         0        0
## 6               0           0             1          0         0        0
##   lon_flag alt_flag tz_flag dst_flag tzone_flag year.y_flag type_flag
## 1        0        0       0        0          0           0         0
## 2        0        0       0        0          0           1         1
## 3        0        0       0        0          0           0         0
## 4        0        0       0        0          0           0         0
## 5        0        0       0        0          0           0         0
## 6        0        0       0        0          0           0         0
##   manufacturer_flag model_flag engines_flag seats_flag engine_flag
```

```
## 1                    0        0        0        0        0
## 2                    1        1        1        1        1
## 3                    0        0        0        0        0
## 4                    0        0        0        0        0
## 5                    0        0        0        0        0
## 6                    0        0        0        0        0
```

```r
library(dplyr)
DF <- DF %>% mutate_if(is.numeric, scale)
head(DF)
```

```
##      month          day carrier flight origin dest   distance       temp
## 1 1.30322 -0.9929373      WN   1716    LGA  ATL -0.3777852  0.3339858
## 2 1.01019  1.6325235      AA    178    JFK  BOS -1.1644742  0.1127815
## 3 1.59625  0.2627179      DL   1585    LGA  MCO -0.1205721 -1.2848272
## 4 1.30322  0.4910188      WN   3494    EWR  MDW -0.4475611 -1.1138966
## 5 1.01019  0.6051693      DL   2231    LGA  DTW -0.7335054  0.3339858
## 6 1.30322 -0.9929373      B6     27    EWR  MCO -0.1383581  0.4144237
##         dewp      humid   wind_dir wind_speed     precip    pressure      visib
## 1  0.7418623  0.9315583  0.07717317  0.4871566 -0.1492223 -0.96830167 0.3664282
## 2  0.2753242  0.2375566  0.36735789 -0.3415806 -0.1492223  1.01596006 0.3664282
## 3 -1.2735821 -0.5125364  0.65754261  1.1087096 -0.1492223  0.28792159 0.3664282
## 4 -1.2269283 -0.7129351 -1.76066338 -0.9631336 -0.1492223  2.54341334 0.3664282
## 5 -0.0325909 -0.7695363 -0.40646802  0.4871566 -0.1492223 -0.12606108 0.3664282
## 6  0.7138700  0.9040226  0.36735789  0.9015253 -0.1492223 -0.04040949 0.3664282
##                              name        lat       lon         alt
## 1     Hartsfield Jackson Atlanta Intl -0.4207546 0.3298674  0.48364704
## 2 General Edward Lawrence Logan Intl  1.1190951 1.2378347 -0.60619417
## 3                    Orlando Intl -1.3395034 0.5408516 -0.52285973
## 4              Chicago Midway Intl  1.0170501 0.1049977  0.04424731
## 5          Detroit Metro Wayne Co  1.0922942 0.4025621  0.07130394
## 6                    Orlando Intl -1.3395034 0.5408516 -0.52285973
##          tz dst           tzone     year.y                        type
## 1  0.6826595   A America/New_York -0.08500492 Fixed wing multi engine
## 2  0.6826595   A America/New_York  0.08617407 Fixed wing multi engine
## 3  0.6826595   A America/New_York  0.08617407 Fixed wing multi engine
## 4 -0.2514221   A  America/Chicago  0.77089000 Fixed wing multi engine
## 5  0.6826595   A America/New_York -1.62561576 Fixed wing multi engine
## 6  0.6826595   A America/New_York  0.77089000 Fixed wing multi engine
##        manufacturer    model     engines      seats    engine
## 1            BOEING  737-7H4 0.05879311 0.02232546 Turbo-fan
## 2            BOEING A320-232 0.05879311 0.15869100 Turbo-fan
## 3            AIRBUS A319-114 0.05879311 0.09808410 Turbo-fan
## 4            BOEING  737-7H4 0.05879311 0.02232546 Turbo-fan
## 5 AIRBUS INDUSTRIE A320-211 0.05879311 0.65869797 Turbo-jet
## 6            AIRBUS A320-232 0.05879311 0.93142905 Turbo-fan
##   sched_arr_time_minute sched_arr_time_num_minute sched_dep_time_minute
## 1            -0.2348938                -1.4325947            -1.36042229
## 2            -1.6716145                -0.3129587             1.23408583
## 3            -0.5222379                 0.9570761            -0.58206985
## 4            -1.3842703                 1.3080068             0.19628258
## 5             0.1099192                 0.5927766            -0.06316823
## 6             1.6041087                -0.7240489            -1.36042229
##   sched_dep_time_num_minute sched_air_time dep_delay_flag temp_flag dewp_flag
## 1              -1.6236293     0.14883297              0         0         0
```

```
## 2                -0.1673447   -0.24317221          0        0        0
## 3                 0.7739125    0.35790240          0        0        0
## 4                 1.4665357   -0.11250381          0        0        0
## 5                 0.5963168    0.04952499          0        0        0
## 6                -0.9842849    0.31608852          0        0        0
##   humid_flag wind_dir_flag wind_speed_flag precip_flag pressure_flag visib_flag
## 1          0             0               0           0             0          0
## 2          0             0               0           0             0          0
## 3          0             0               0           0             0          0
## 4          0             0               0           0             0          0
## 5          0             0               0           0             0          0
## 6          0             0               0           0             1          0
##   name_flag lat_flag lon_flag alt_flag tz_flag dst_flag tzone_flag year.y_flag
## 1         0        0        0        0       0        0          0           0
## 2         0        0        0        0       0        0          0           1
## 3         0        0        0        0       0        0          0           0
## 4         0        0        0        0       0        0          0           0
## 5         0        0        0        0       0        0          0           0
## 6         0        0        0        0       0        0          0           0
##   type_flag manufacturer_flag model_flag engines_flag seats_flag engine_flag
## 1         0                 0          0            0          0           0
## 2         1                 1          1            1          1           1
## 3         0                 0          0            0          0           0
## 4         0                 0          0            0          0           0
## 5         0                 0          0            0          0           0
## 6         0                 0          0            0          0           0
```

```r
DF$dep_delay <- dep_delay_vec
```

#take out extreme departure delays

```r
DF<-DF[DF$dep_delay < 30,]
```

```r
set.seed(42)
DF$flight <- NULL
train_index <- sample(1:nrow(DF),size=2*nrow(DF)/3,replace=FALSE)
train_df <- DF[train_index,]
test_df <- DF[-train_index,]
```

## predicting 0

```r
rmse = mean((test_df$dep_delay-0)^2) %>% sqrt()
rmse
```

```
## [1] 8.30571
```

## predicting the mean

```r
rmse = mean((test_df$dep_delay-mean(train_df$dep_delay))^2)%>% sqrt()
rmse
```

```
## [1] 8.299767
```

## predicting the median

```
rmse = mean((test_df$dep_delay-median(train_df$dep_delay))^2)%>% sqrt()
rmse
```

```
## [1] 8.469257
```

## linear regression

```
model <- lm(dep_delay ~ .-model, data=train_df)

summary <- round(summary(model)$coefficients,6)
sorteddf <- summary[order(summary[,ncol(summary)]),]
head(sorteddf)
```

```
##            Estimate Std. Error   t value Pr(>|t|)
## carrierAA -2.318607   0.269781 -8.594395        0
## carrierAS -3.341812   0.616737 -5.418535        0
## carrierDL -1.392992   0.253561 -5.493710        0
## carrierEV  1.274029   0.187960  6.778204        0
## carrierMQ -2.396562   0.272333 -8.800101        0
## carrierUS -2.323719   0.260042 -8.935921        0
```

```
sorteddf
```

```
##                             Estimate Std. Error    t value
## carrierAA                  -2.318607   0.269781  -8.594395
## carrierAS                  -3.341812   0.616737  -5.418535
## carrierDL                  -1.392992   0.253561  -5.493710
## carrierEV                   1.274029   0.187960   6.778204
## carrierMQ                  -2.396562   0.272333  -8.800101
## carrierUS                  -2.323719   0.260042  -8.935921
## carrierWN                   2.796283   0.325335   8.595102
## originJFK                   0.831817   0.115446   7.205266
## wind_speed                  0.294774   0.027966  10.540539
## precip                      0.208324   0.029304   7.109046
## pressure                   -0.293555   0.027097 -10.833530
## year.y                      0.257028   0.042093   6.106234
## seats                       0.365535   0.068252   5.355693
## sched_arr_time_num_minute   0.241584   0.045987   5.253282
## sched_dep_time_minute       0.132111   0.024902   5.305163
## sched_dep_time_num_minute   1.258505   0.044951  27.996992
## dep_delay_flag1            -3.014132   0.155718 -19.356329
## pressure_flag1              0.606457   0.098873   6.133711
## carrierHA                  -4.771275   1.108689  -4.303528
## destCHO                   -36.123526   9.266129  -3.898449
## destILM                   -31.097284   8.056795  -3.859759
## destPDX                    14.661963   3.837985   3.820224
## destSMF                    16.425104   4.321366   3.800905
## destPHX                     8.055347   2.133058   3.776432
## destCRW                   -31.215951   8.376979  -3.726397
## visib                      -0.127430   0.034660  -3.676623
## destSAN                    13.889006   3.780985   3.673383
```

```
## destJAC         11.591104    3.169410    3.657181
## destSEA         13.325506    3.646194    3.654635
## destLAX         14.348648    3.943394    3.638654
## destCAK        -31.214329    8.580927   -3.637640
## destTUL        -13.660070    3.754857   -3.637973
## destOMA        -15.118697    4.158319   -3.635771
## destSFO         16.632081    4.588360    3.624842
## destBHM        -21.009316    5.837679   -3.598916
## destLGB         14.042295    3.939582    3.564412
## destHNL         66.885768   18.771149    3.563222
## destDSM        -17.187039    4.854017   -3.540786
## destTYS        -25.105682    7.108433   -3.531817
## destBGR        -30.575425    8.693169   -3.517178
## destXNA        -14.613857    4.161947   -3.511303
## destBWI        -34.467570    9.818060   -3.510629
## destLAS          9.278585    2.648563    3.503253
## destAVL        -25.975094    7.418111   -3.501578
## destSNA         13.386324    3.822915    3.501601
## destBDL        -35.271411   10.155622   -3.473092
## destROC        -32.435956    9.347846   -3.469886
## destGSP        -25.376154    7.320086   -3.466647
## destCLE        -29.207281    8.425622   -3.466483
## destRIC        -31.792119    9.183695   -3.461800
## distance       -15.049604    4.352182   -3.457945
## destCVG        -25.692606    7.447115   -3.450008
## destBTV        -31.939098    9.300241   -3.434223
## destDAY        -26.344313    7.671046   -3.434253
## destBOS        -33.302948    9.738780   -3.419622
## destBUF        -31.162718    9.124996   -3.415094
## destSJC         15.585073    4.565151    3.413923
## destOAK         15.741470    4.621128    3.406413
## destPWM        -31.438770    9.240920   -3.402126
## destSLC          4.520401    1.329819    3.399261
## destGRR        -24.648547    7.256882   -3.396575
## destGSO        -27.791331    8.186192   -3.394903
## destIAD        -32.407962    9.550812   -3.393215
## destOKC        -10.626435    3.139739   -3.384496
## destPIT        -30.140890    8.909258   -3.383098
## destDCA        -32.545897    9.625597   -3.381182
## destMCI        -14.798414    4.379165   -3.379277
## destALB        -33.767488    9.997762   -3.377505
## destPHL        -34.872031   10.335187   -3.374108
## destCHS        -24.003318    7.115993   -3.373151
## destIND        -23.582074    6.995318   -3.371123
## destSTL        -19.054904    5.660461   -3.366317
## destPVD        -33.332083    9.908856   -3.363868
## destANC         36.554722   10.888763    3.357105
## destMHT        -32.248146    9.615996   -3.353594
## destSYR        -32.442780    9.675578   -3.353059
## destORF        -30.723323    9.167680   -3.351265
## destCMH        -27.023263    8.067612   -3.349599
## destBNA        -21.381935    6.388039   -3.347183
## destSDF        -23.471800    7.027043   -3.340210
## destRDU        -27.895202    8.360245   -3.336649
```

```
## destTVC                                 -23.763734   7.123077  -3.336161
## destSAV                                 -22.145147   6.646547  -3.331828
## destMDW                                 -22.017052   6.613415  -3.329150
## destACK                                 -32.386998   9.753827  -3.320440
## destMKE                                 -21.556521   6.534602  -3.298827
## destDTW                                 -26.134652   7.928589  -3.296255
## destCLT                                 -25.088232   7.674413  -3.269075
## destJAX                                 -19.502019   5.989172  -3.256213
## destORD                                 -21.244774   6.565174  -3.235981
## destMCO                                 -17.147920   5.306618  -3.231422
## destMSY                                 -12.718760   3.935982  -3.231408
## destMSN                                 -19.745453   6.122372  -3.225131
## destMVY                                 -31.958595   9.911091  -3.224528
## day                                       0.075847   0.023801   3.186663
## destATL                                 -20.349736   6.406683  -3.176329
## destRSW                                 -14.406148   4.545223  -3.169514
## destSRQ                                 -15.035777   4.752239  -3.163935
## destTPA                                 -15.652896   4.951828  -3.161034
## destBUR                                  12.507966   3.973579   3.147783
## destCAE                                 -22.661801   7.428543  -3.050639
## destMSP                                 -14.801747   4.870956  -3.038776
## destMEM                                 -15.840986   5.227838  -3.030122
## destFLL                                 -13.585753   4.566169  -2.975307
## destMIA                                 -13.114128   4.444244  -2.950812
## destPSE                                  -4.940578   1.678873  -2.942794
## destPBI                                 -13.976366   4.805236  -2.908570
## (Intercept)                              16.474044   5.670479   2.905230
## destAUS                                  -5.873576   2.038240  -2.881690
## carrierB6                                -0.654102   0.227007  -2.881420
## destDFW                                  -7.889899   2.756116  -2.862688
## destMYR                                 -21.926351   7.794817  -2.812940
## originLGA                                -0.325794   0.115862  -2.811914
## destIAH                                  -7.258532   2.613026  -2.777827
## destBQN                                  -4.859601   1.759015  -2.762683
## destHOU                                  -6.906619   2.546684  -2.712005
## humid                                     0.417871   0.155521   2.686914
## manufacturerGULFSTREAM AEROSPACE          8.169688   3.182253   2.567265
## manufacturerCESSNA                        4.567610   1.811097   2.522013
## destBZN                                   5.625552   2.281507   2.465718
## precip_flag1                             -9.363342   3.876804  -2.415222
## temp_flag1                                8.517235   3.583399   2.376859
## carrierYV                                -1.410746   0.595048  -2.370812
## manufacturerLEARJET INC                   9.621229   4.098480   2.347511
## destSAT                                  -4.082388   1.761550  -2.317497
## carrierUA                                 0.589784   0.254952   2.313314
## destEYW                                 -11.192190   4.855385  -2.305109
## carrierFL                                 0.948807   0.417663   2.271705
## manufacturerMCDONNELL DOUGLAS CORPORATION 7.130518   3.144929   2.267307
## carrierOO                                -4.676371   2.080031  -2.248222
## month                                     0.055963   0.025119   2.227950
## engineReciprocating                      -4.784441   2.155649  -2.219490
## engineTurbo-jet                          -7.579155   3.523834  -2.150826
## manufacturerBOEING                        6.655788   3.121857   2.131996
## manufacturerBARKER JACK L                 4.156889   1.949770   2.131989
```

14

```
## manufacturerCANADAIR LTD                    10.020777   4.719325   2.123350
## manufacturerDEHAVILLAND                       5.277486   2.489152   2.120194
## destPSP                                       9.000603   4.426308   2.033434
## manufacturerMCDONNELL DOUGLAS                 6.307716   3.121992   2.020414
## manufacturerPIPER                             4.028235   2.033192   1.981236
## manufacturerMCDONNELL DOUGLAS AIRCRAFT CO     6.192869   3.126076   1.981036
## engineTurbo-fan                              -6.970145   3.522649  -1.978666
## manufacturerPAIR MIKE E                       5.954488   3.056590   1.948082
## manufacturerCIRRUS DESIGN CORP                3.649312   1.879395   1.941748
## manufacturerAIRBUS INDUSTRIE                  5.928450   3.122708   1.898496
## destSTT                                      -2.956627   1.572582  -1.880110
## destSJU                                      -2.898323   1.599828  -1.811647
## manufacturerEMBRAER                           5.378862   3.116909   1.725704
## destSBN                                     -15.570339   9.034869  -1.723361
## manufacturerAIRBUS                            5.285531   3.124349   1.691722
## wind_dir_flag1                                0.269529   0.159642   1.688332
## manufacturerLEBLANC GLENN T                   4.512985   2.739959   1.647100
## destMTJ                                       5.184644   3.150881   1.645458
## manufacturerFRIEDEMANN JON                    3.889750   2.392071   1.626102
## manufacturerBOMBARDIER INC                    4.655030   3.120928   1.491553
## destDEN                                      -2.165776   1.518672  -1.426099
## manufacturerCANADAIR                          4.355101   3.145879   1.384383
## manufacturerKILDALL GARY                      2.862667   2.445757   1.170462
## manufacturerAMERICAN AIRCRAFT INC             3.066792   2.683506   1.142830
## manufacturerMARZ BARRY                        2.829639   2.587439   1.093606
## manufacturerAVIONS MARCEL DASSAULT            6.754935   6.699339   1.008299
## manufacturerBEECH                             3.431401   3.551823   0.966096
## manufacturerJOHN G HESS                      -4.810782   5.388356  -0.892811
## manufacturerSIKORSKY                          3.186564   3.678961   0.866159
## manufacturerBELL                              2.610957   3.284993   0.794814
## type_flag1                                   -0.203922   0.271119  -0.752150
## manufacturerDOUGLAS                           4.008806   5.632419   0.711738
## wind_dir                                      0.018733   0.027026   0.693150
## sched_arr_time_minute                        -0.013804   0.024404  -0.565633
## carrierVX                                     0.181052   0.323759   0.559219
## year.y_flag1                                 -0.102593   0.199432  -0.514424
## manufacturerLAMBERT RICHARD                   1.058252   2.476434   0.427329
## manufacturerHURLEY JAMES LARRY               -1.997386   4.928047  -0.405310
## carrierF9                                     0.242613   0.637308   0.380684
## destHDN                                       1.137227   3.190730   0.356416
## engines                                      -0.051607   0.146845  -0.351438
## destEGE                                       0.431238   1.375901   0.313422
## wind_speed_flag1                              0.449983   1.446992   0.310978
## typeFixed wing single engine                  0.606625   2.103014   0.288455
## dewp                                          0.040676   0.302037   0.134670
## manufacturerAVIAT AIRCRAFT INC               -0.501995   3.968705  -0.126488
## temp                                          0.023006   0.260704   0.088244
## manufacturerROBINSON HELICOPTER CO            0.248651   3.245307   0.076619
## typeRotorcraft                               -0.198373   4.095669  -0.048435
##                                            Pr(>|t|)
## carrierAA                                   0.000000
## carrierAS                                   0.000000
## carrierDL                                   0.000000
## carrierEV                                   0.000000
```

```
## carrierMQ                    0.000000
## carrierUS                    0.000000
## carrierWN                    0.000000
## originJFK                    0.000000
## wind_speed                   0.000000
## precip                       0.000000
## pressure                     0.000000
## year.y                       0.000000
## seats                        0.000000
## sched_arr_time_num_minute    0.000000
## sched_dep_time_minute        0.000000
## sched_dep_time_num_minute    0.000000
## dep_delay_flag1              0.000000
## pressure_flag1               0.000000
## carrierHA                    0.000017
## destCHO                      0.000097
## destILM                      0.000114
## destPDX                      0.000133
## destSMF                      0.000144
## destPHX                      0.000159
## destCRW                      0.000194
## visib                        0.000236
## destSAN                      0.000239
## destJAC                      0.000255
## destSEA                      0.000258
## destLAX                      0.000274
## destCAK                      0.000275
## destTUL                      0.000275
## destOMA                      0.000277
## destSFO                      0.000289
## destBHM                      0.000320
## destLGB                      0.000365
## destHNL                      0.000366
## destDSM                      0.000399
## destTYS                      0.000413
## destBGR                      0.000436
## destXNA                      0.000446
## destBWI                      0.000447
## destLAS                      0.000460
## destAVL                      0.000463
## destSNA                      0.000463
## destBDL                      0.000515
## destROC                      0.000521
## destGSP                      0.000527
## destCLE                      0.000528
## destRIC                      0.000537
## distance                     0.000545
## destCVG                      0.000561
## destBTV                      0.000594
## destDAY                      0.000594
## destBOS                      0.000627
## destBUF                      0.000638
## destSJC                      0.000641
## destOAK                      0.000658
```

```
## destPWM                              0.000669
## destSLC                              0.000676
## destGRR                              0.000683
## destGSO                              0.000687
## destIAD                              0.000691
## destOKC                              0.000713
## destPIT                              0.000717
## destDCA                              0.000722
## destMCI                              0.000727
## destALB                              0.000732
## destPHL                              0.000741
## destCHS                              0.000743
## destIND                              0.000749
## destSTL                              0.000762
## destPVD                              0.000769
## destANC                              0.000788
## destMHT                              0.000798
## destSYR                              0.000799
## destORF                              0.000805
## destCMH                              0.000810
## destBNA                              0.000817
## destSDF                              0.000837
## destRDU                              0.000848
## destTVC                              0.000850
## destSAV                              0.000863
## destMDW                              0.000871
## destACK                              0.000899
## destMKE                              0.000971
## destDTW                              0.000980
## destCLT                              0.001079
## destJAX                              0.001129
## destORD                              0.001213
## destMCO                              0.001232
## destMSY                              0.001232
## destMSN                              0.001260
## destMVY                              0.001262
## day                                  0.001440
## destATL                              0.001492
## destRSW                              0.001527
## destSRQ                              0.001557
## destTPA                              0.001573
## destBUR                              0.001646
## destCAE                              0.002284
## destMSP                              0.002376
## destMEM                              0.002445
## destFLL                              0.002928
## destMIA                              0.003170
## destPSE                              0.003253
## destPBI                              0.003632
## (Intercept)                          0.003671
## destAUS                              0.003956
## carrierB6                            0.003960
## destDFW                              0.004201
## destMYR                              0.004910
```

```
## originLGA                                     0.004926
## destIAH                                       0.005473
## destBQN                                       0.005734
## destHOU                                       0.006689
## humid                                         0.007213
## manufacturerGULFSTREAM AEROSPACE              0.010252
## manufacturerCESSNA                            0.011670
## destBZN                                       0.013675
## precip_flag1                                  0.015727
## temp_flag1                                    0.017462
## carrierYV                                     0.017751
## manufacturerLEARJET INC                       0.018901
## destSAT                                       0.020478
## carrierUA                                     0.020707
## destEYW                                       0.021162
## carrierFL                                     0.023106
## manufacturerMCDONNELL DOUGLAS CORPORATION 0.023373
## carrierOO                                     0.024564
## month                                         0.025886
## engineReciprocating                           0.026455
## engineTurbo-jet                               0.031492
## manufacturerBOEING                            0.033009
## manufacturerBARKER JACK L                     0.033010
## manufacturerCANADAIR LTD                      0.033727
## manufacturerDEHAVILLAND                       0.033992
## destPSP                                       0.042011
## manufacturerMCDONNELL DOUGLAS                 0.043343
## manufacturerPIPER                             0.047567
## manufacturerMCDONNELL DOUGLAS AIRCRAFT CO 0.047590
## engineTurbo-fan                               0.047856
## manufacturerPAIR MIKE E                       0.051408
## manufacturerCIRRUS DESIGN CORP                0.052170
## manufacturerAIRBUS INDUSTRIE                  0.057633
## destSTT                                       0.060096
## destSJU                                       0.070043
## manufacturerEMBRAER                           0.084403
## destSBN                                       0.084826
## manufacturerAIRBUS                            0.090702
## wind_dir_flag1                                0.091350
## manufacturerLEBLANC GLENN T                   0.099540
## destMTJ                                       0.099878
## manufacturerFRIEDEMANN JON                    0.103931
## manufacturerBOMBARDIER INC                    0.135819
## destDEN                                       0.153843
## manufacturerCANADAIR                          0.166244
## manufacturerKILDALL GARY                      0.241817
## manufacturerAMERICAN AIRCRAFT INC             0.253112
## manufacturerMARZ BARRY                        0.274130
## manufacturerAVIONS MARCEL DASSAULT            0.313313
## manufacturerBEECH                             0.333998
## manufacturerJOHN G HESS                       0.371960
## manufacturerSIKORSKY                          0.386405
## manufacturerBELL                              0.426724
## type_flag1                                    0.451962
```

```
## manufacturerDOUGLAS                     0.476629
## wind_dir                                0.488217
## sched_arr_time_minute                   0.571644
## carrierVX                               0.576013
## year.y_flag1                            0.606957
## manufacturerLAMBERT RICHARD             0.669141
## manufacturerHURLEY JAMES LARRY          0.685251
## carrierF9                               0.703438
## destHDN                                 0.721530
## engines                                 0.725260
## destEGE                                 0.753960
## wind_speed_flag1                        0.755818
## typeFixed wing single engine           0.772999
## dewp                                    0.892873
## manufacturerAVIAT AIRCRAFT INC          0.899346
## temp                                    0.929683
## manufacturerROBINSON HELICOPTER CO      0.938927
## typeRotorcraft                         0.961370
```

```r
lm_test_df <- test_df

in_test_but_not_train <- setdiff(unique(lm_test_df$model), unique(train_df$model))
lm_test_df <- lm_test_df[ !lm_test_df$model %in% in_test_but_not_train, ]

in_test_but_not_train <- setdiff(unique(lm_test_df$dest), unique(train_df$dest))
lm_test_df <- lm_test_df[ !lm_test_df$dest %in% in_test_but_not_train, ]

preds = predict(model, newdata=lm_test_df)
```

```
## Warning in predict.lm(model, newdata = lm_test_df): prediction from a rank-
## deficient fit may be misleading
```

```r
rmse = sqrt(mean((lm_test_df$dep_delay - preds)^2))
rmse
```

```
## [1] 8.003373
```

# gbm

```r
set.seed(42)
library(gbm)
```
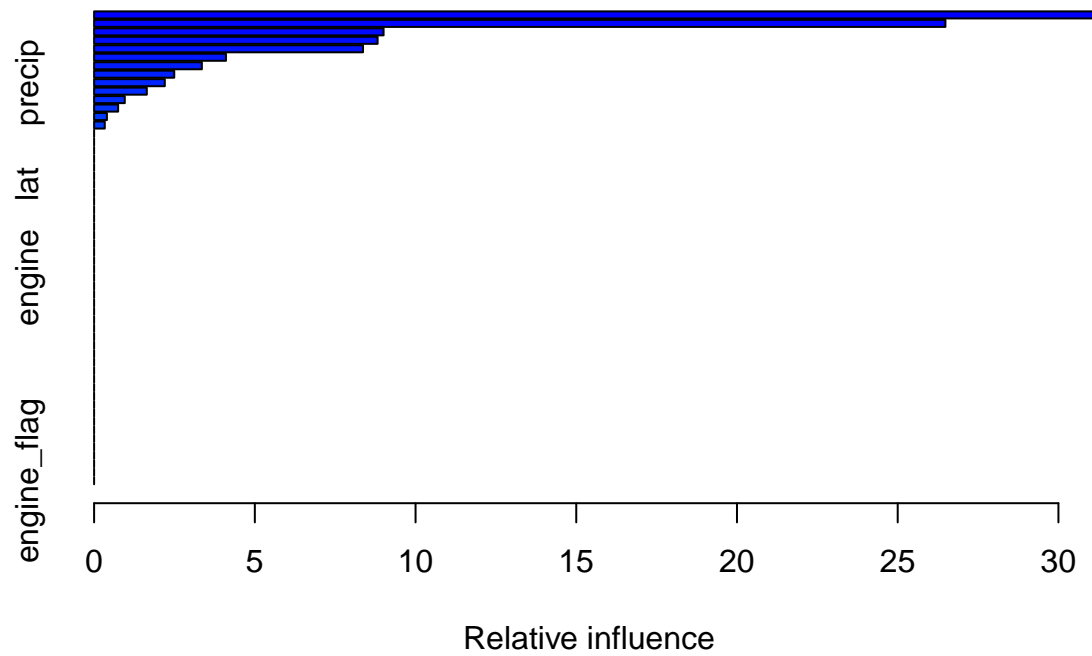
```
## Loaded gbm 2.1.5
```

```r
model <- gbm(dep_delay ~ ., data=train_df,
             n.trees=100, shrinkage=0.1, cv.folds = 3) # default shrinkage = 0.1
```

```
## Distribution not specified, assuming gaussian ...
```

```r
preds = predict(model, newdata=test_df, n.trees=100)
rmse = sqrt(mean((test_df$dep_delay - preds)^2))
rmse
```

```
## [1] 7.978221
```

```
summary(model)
```



Relative influence

```
##                                             var        rel.inf
## sched_dep_time_num_minute sched_dep_time_num_minute 31.1084810
## model                                         model 26.4841554
## dest                                           dest  9.0052224
## carrier                                     carrier  8.8195760
## month                                         month  8.3653954
## dewp                                           dewp  4.1022280
## origin                                       origin  3.3525792
## sched_arr_time_num_minute sched_arr_time_num_minute  2.4930686
## precip                                       precip  2.1992862
## pressure                                   pressure  1.6390248
## dep_delay_flag                       dep_delay_flag  0.9536652
## humid                                         humid  0.7472539
## pressure_flag                         pressure_flag  0.3983933
## temp                                           temp  0.3316705
## day                                             day  0.0000000
## distance                                   distance  0.0000000
## wind_dir                                   wind_dir  0.0000000
## wind_speed                               wind_speed  0.0000000
## visib                                         visib  0.0000000
## name                                           name  0.0000000
## lat                                             lat  0.0000000
## lon                                             lon  0.0000000
## alt                                             alt  0.0000000
## tz                                               tz  0.0000000
## dst                                             dst  0.0000000
## tzone                                         tzone  0.0000000
## year.y                                       year.y  0.0000000
## type                                           type  0.0000000
## manufacturer                           manufacturer  0.0000000
## engines                                     engines  0.0000000
```

```
## seats                                          seats  0.0000000
## engine                                        engine  0.0000000
## sched_arr_time_minute        sched_arr_time_minute  0.0000000
## sched_dep_time_minute        sched_dep_time_minute  0.0000000
## sched_air_time                    sched_air_time  0.0000000
## temp_flag                            temp_flag  0.0000000
## dewp_flag                            dewp_flag  0.0000000
## humid_flag                          humid_flag  0.0000000
## wind_dir_flag                    wind_dir_flag  0.0000000
## wind_speed_flag                wind_speed_flag  0.0000000
## precip_flag                        precip_flag  0.0000000
## visib_flag                          visib_flag  0.0000000
## name_flag                            name_flag  0.0000000
## lat_flag                              lat_flag  0.0000000
## lon_flag                              lon_flag  0.0000000
## alt_flag                              alt_flag  0.0000000
## tz_flag                                tz_flag  0.0000000
## dst_flag                              dst_flag  0.0000000
## tzone_flag                          tzone_flag  0.0000000
## year.y_flag                        year.y_flag  0.0000000
## type_flag                            type_flag  0.0000000
## manufacturer_flag            manufacturer_flag  0.0000000
## model_flag                          model_flag  0.0000000
## engines_flag                      engines_flag  0.0000000
## seats_flag                          seats_flag  0.0000000
## engine_flag                        engine_flag  0.0000000
```

Here, you can see the relative influence for each variable for gbm.

For a gbm, the improvement in the splitting criterion (which is mean squared error for regression) for a given variable is calculated at each step. The relative influence for a given variable is the average of these improvements over all the trees where the aforementioned variable is used.
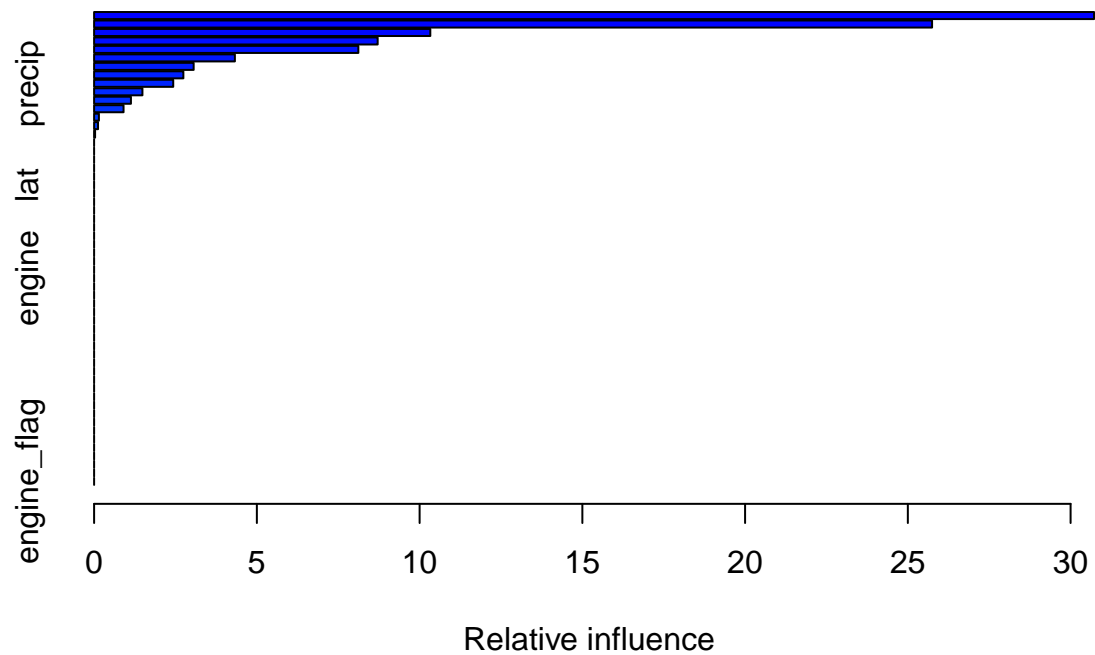
```
model <- gbm(dep_delay ~ ., data=train_df,
             n.trees=1000, shrinkage=0.01) # default shrinkage = 0.1
```

```
## Distribution not specified, assuming gaussian ...
```

```
preds = predict(model, newdata=test_df, n.trees=1000)
rmse = sqrt(mean((test_df$dep_delay - preds)^2))
rmse
```

```
## [1] 7.980782
```

```
summary(model)
```

Relative influence

```
##                                                  var       rel.inf
## sched_dep_time_num_minute sched_dep_time_num_minute 30.7209486
## model                                            model 25.7445886
## dest                                              dest 10.3287677
## carrier                                        carrier  8.7111198
## month                                            month  8.1184266
## dewp                                              dewp  4.3249508
## origin                                          origin  3.0586748
## sched_arr_time_num_minute sched_arr_time_num_minute  2.7422873
## precip                                          precip  2.4312138
## pressure                                      pressure  1.4845190
## humid                                            humid  1.1316047
## dep_delay_flag                          dep_delay_flag  0.9042069
## pressure_flag                            pressure_flag  0.1476206
## temp                                              temp  0.1205245
## day                                                day  0.0305463
## distance                                      distance  0.0000000
## wind_dir                                      wind_dir  0.0000000
## wind_speed                                  wind_speed  0.0000000
## visib                                            visib  0.0000000
## name                                              name  0.0000000
## lat                                                lat  0.0000000
## lon                                                lon  0.0000000
## alt                                                alt  0.0000000
## tz                                                  tz  0.0000000
## dst                                                dst  0.0000000
## tzone                                            tzone  0.0000000
## year.y                                          year.y  0.0000000
## type                                              type  0.0000000
## manufacturer                              manufacturer  0.0000000
## engines                                        engines  0.0000000
## seats                                            seats  0.0000000
## engine                                          engine  0.0000000
```

```
## sched_arr_time_minute          sched_arr_time_minute  0.0000000
## sched_dep_time_minute          sched_dep_time_minute  0.0000000
## sched_air_time                        sched_air_time  0.0000000
## temp_flag                                  temp_flag  0.0000000
## dewp_flag                                  dewp_flag  0.0000000
## humid_flag                                humid_flag  0.0000000
## wind_dir_flag                          wind_dir_flag  0.0000000
## wind_speed_flag                      wind_speed_flag  0.0000000
## precip_flag                              precip_flag  0.0000000
## visib_flag                                visib_flag  0.0000000
## name_flag                                  name_flag  0.0000000
## lat_flag                                    lat_flag  0.0000000
## lon_flag                                    lon_flag  0.0000000
## alt_flag                                    alt_flag  0.0000000
## tz_flag                                      tz_flag  0.0000000
## dst_flag                                    dst_flag  0.0000000
## tzone_flag                                tzone_flag  0.0000000
## year.y_flag                              year.y_flag  0.0000000
## type_flag                                  type_flag  0.0000000
## manufacturer_flag                  manufacturer_flag  0.0000000
## model_flag                                model_flag  0.0000000
## engines_flag                            engines_flag  0.0000000
## seats_flag                                seats_flag  0.0000000
## engine_flag                              engine_flag  0.0000000
```

set.seed(42)

x <- 2^seq(5,14, by=1) rmse_vec <- numeric(length(x)) count <- 1 for (val in x) { hboost <- gbm( dep_delay ~ ., data = train_df, n.trees = val, distribution = 'gaussian', shrinkage = 0.01 ) preds = predict(hboost, n.trees = val, newdata = test_df) mse = mean((test_df$dep_delay - preds) ^ 2) rmse <- sqrt(mse) rmse_vec[count] <- rmse
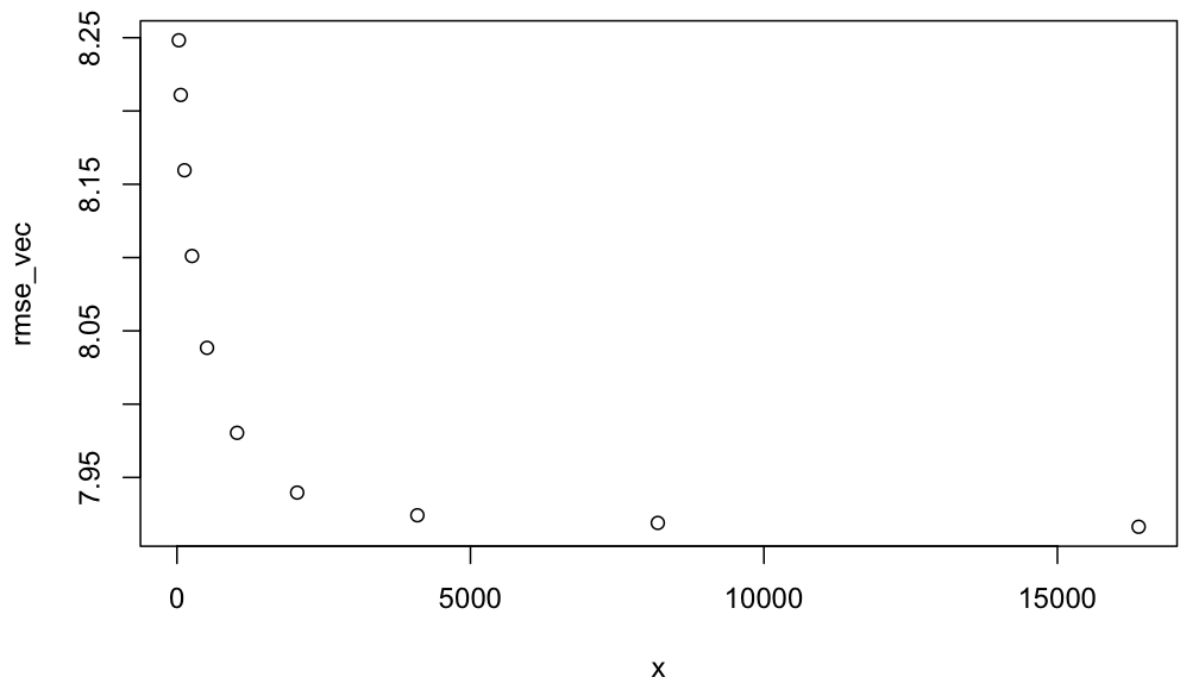
print(val) print(rmse) count = count + 1 }

plot(x, rmse_vec)

summary(hboost) class(summary(hboost)) summary <- summary(hboost) write.csv(summary,'16384trees_gbm.csv')

Analysis:

Tuning gbm

Here I plotted root mean squared error (rmse) vs the number of trees for shrinkage of 0.01 and all other variables as default for gbm. You can see that after around 5000 trees, increasing the number of trees further gives diminishing returns.

#Methods:

## Data Preprocessing

I performed data preprocessing. My data preprocessing steps include the following: - Dropping columns that contain data from after the planes' departure which may leak information about the response variable dep_delay. - Dropping columns with too many NAs. - Impute NAs for the remaining columns. - Scaling the data to work well with methods like lasso regression.

## Modelling

# Basic Models

dep_delay is the number of minutes that the plane either departs early or late. Negative numbers are for early departures and positive numbers are for the number of minutes the plane is late. First, I used a basic model of simply predicting the dep_delay to always be 0. This was done to establish baseline performance. This model had an root mean squared error (RMSE) of

24

# Linear Regression

# GBM

## Basic

Conclusion: In conclusion, out of the methods that we covered in class, I found gradient boosted models to provide the best performance based on having the lowest root mean squared error on the hold out test set.

Based on the relative influence scores provided by the gbm, some of the most important feature variables include dest, model, and sched_dep_time_num_minute.

The dest column contains the airport code for where a given flight is flying to. Based on my run of gbm with a shrinkage of 0.01 and 16834 trees, dest was the most important feature with 49.56 relative influence.("Gradient Boosting Machines · UC Business Analytics R Programming Guide" 2019).

# References

"Gradient Boosting Machines · UC Business Analytics R Programming Guide." 2019. http://uc-r.github.io/gbm_regression#h2o.