

## Instructions

You're almost done with the semester! Take a second to congratulate yourself on getting here. As a reminder, this final project is simply an (imperfect) way of measuring what you have learned throughout the semester. So take a deep breath and do your best, but also remember that it doesn't determine your value as a human being.

The exam is split into 4 sections: Module 1, 2 and 3 (6 questions), Modules 4 and 5 (3 questions), Module 6 (2 questions) and the final project. Most of the questions on this exam are short answers. You don't need to write out an overly long response (a sentence or so for each part of the question should be fine), but you should be specific in explaining your response. For example, if there is a question about whether the assumptions are reasonable. You shouldn't just say "from the plot we can see that the linearity assumption is (or is not) reasonable," but instead you should explain specifically why the plot leads you to believe the linearity assumption is (or is not) reasonable.

The exam is open notes so you **can** use any of the material or any of the notes you have taken throughout the class. You **cannot** discuss the exam (while it is in progress) with anyone else. You also **cannot** use any generative AI tools. Submissions will be sent by e-mail to **nbb45@cornell.edu** before **May 14th 11:59pm**.

## Module 1, 2, and 3

In the questions for Modules 1, 2, and 3, we will look at data from SNCF, France's national railway. The data has been cleaned and made easily available by TidyTuesday. In particular, we have data on train delays from each month between 2015-2018 for each train route (i.e., from city A to city B). So each observation (i.e., row in the data) corresponds to a specific route in a specific year and month. In the dataset, we will be particularly interested in the following variables

For each row in the data, we have the following variables

- year : year of observation (2015, 2016, 2017 or 2018)
- month : month of observation (1, 2, ..., 12)
- departure\_station : station where the route begins (e.g., “PARIS NORD” or “MONTPELLIER”)
- arrival\_station : station where the route ends (e.g., “PARIS NORD” or “MONTPELLIER”)
- journey\_time\_avg : average journey time in minutes for the route for that year and month
- avg\_delay\_all\_departing : average delay in minutes for all departures for the route for that year and month (i.e., how many minutes the train was late to leave departure station)
- avg\_delay\_all\_arriving : average delay in minutes for all arrivals for the route for that year and month (i.e., how many minutes the train was late to arrive at the arrival\_station)

In the following questions, the model you fit or consider may change from question to question.

```
## Load in data and remove some outliers
train_data <- read.csv("https://raw.githubusercontent.com/rfordatascience/tidytuesday/master/data/2019/
# removing some outliers
train_data <- train_data[-which(train_data$avg_delay_all_arriving < -30),]
train_data <- train_data[-which(train_data$avg_delay_all_departing > 30),]
# make month and year factors
train_data$month <- as.factor(train_data$month)
train_data$year <- as.factor(train_data$year)
```

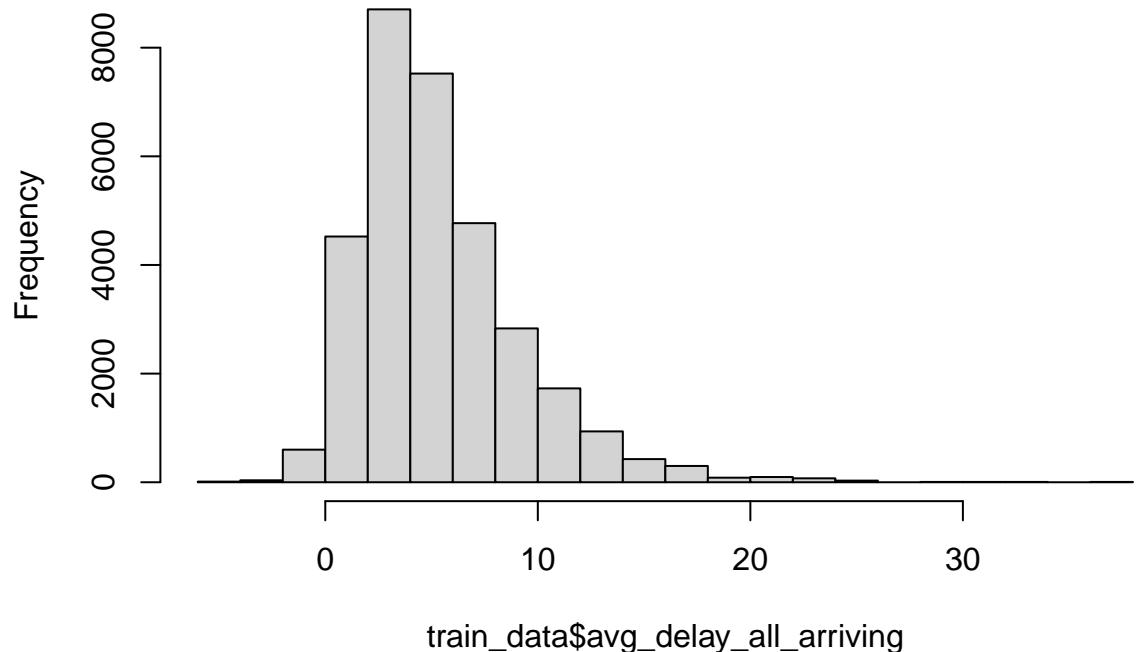
### Question 1 (2 pts)

Suppose we are interested in modeling the average delayed arrival; i.e., avg\_delay\_all\_arriving is the outcome variable. Specifically, we would like to investigate the association between average delayed arrival and journey time (journey\_time\_avg) when controlling for the average departure delay (avg\_delay\_all\_departing).

Fit the relevant linear model below and write 1 sentence interpreting the estimated coefficient for journey\_time\_avg.

```
hist(train_data$avg_delay_all_arriving)
```

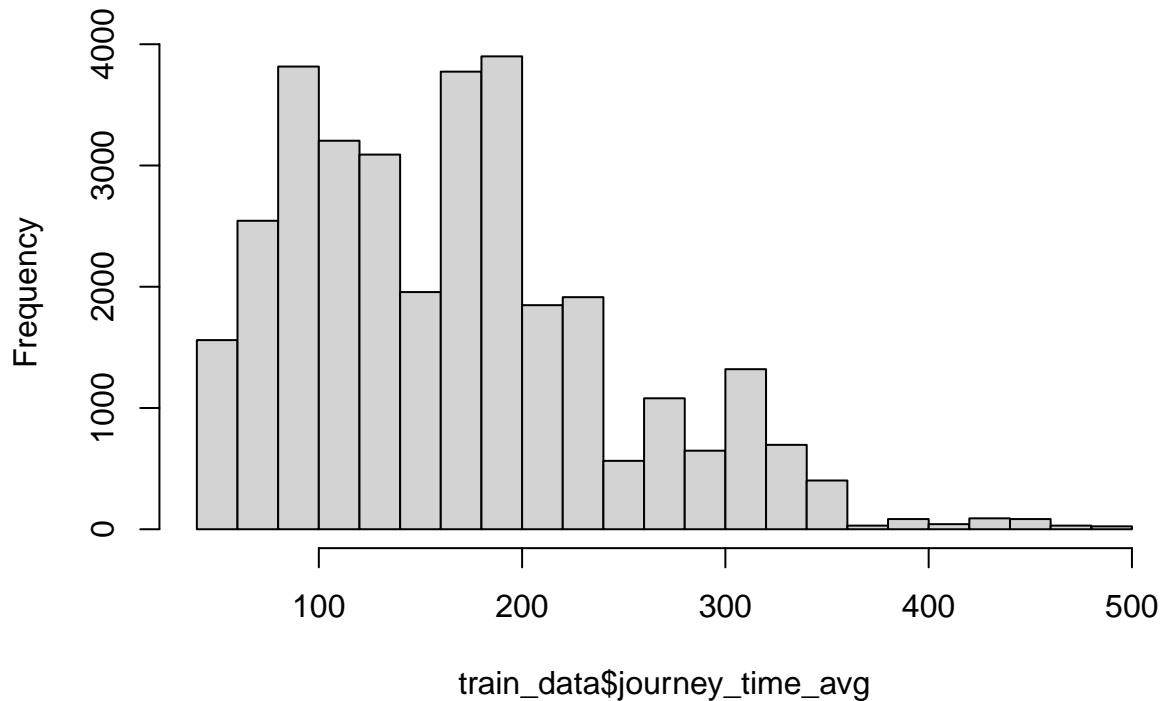
**Histogram of train\_data\$avg\_delay\_all\_arriving**



Question 1 Answer

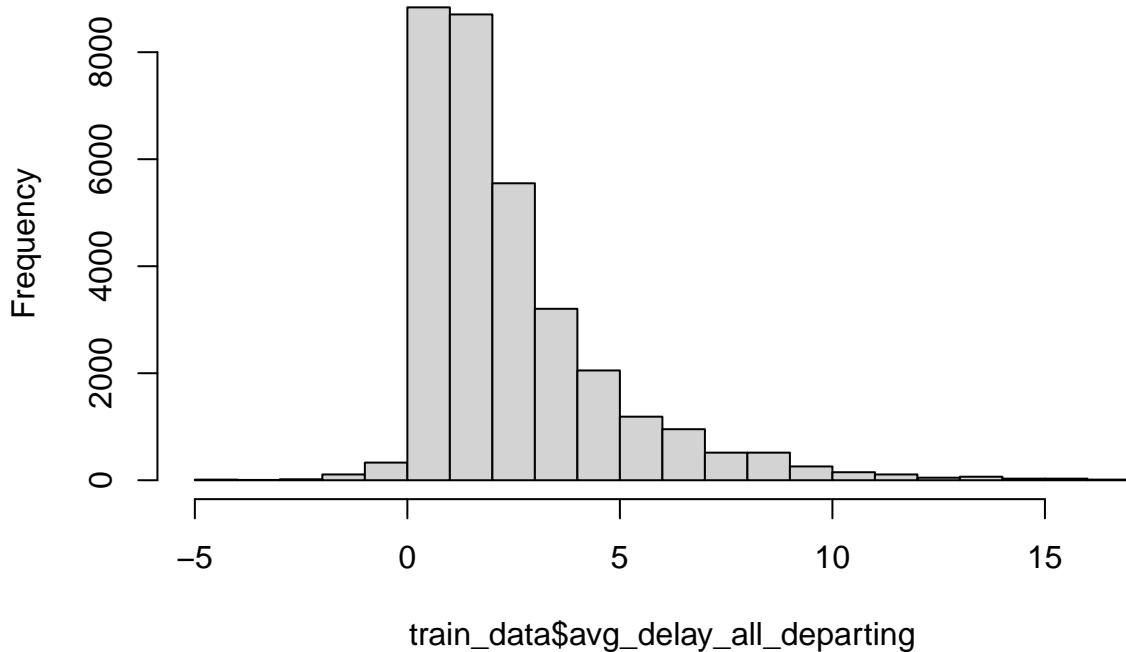
```
hist(train_data$journey_time_avg)
```

### Histogram of train\_data\$journey\_time\_avg



```
hist(train_data$avg_delay_all_departing)
```

## Histogram of train\_data\$avg\_delay\_all\_departing



```
mod <- lm(avg_delay_all_arriving ~ journey_time_avg + avg_delay_all_departing, data=train_data)
summary(mod)
```

```
##
## Call:
## lm(formula = avg_delay_all_arriving ~ journey_time_avg + avg_delay_all_departing,
##      data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.9335 -1.6684 -0.1168  1.3580 22.8919
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)             -0.3513436  0.0404302 -8.69   <2e-16 ***
## journey_time_avg         0.0220767  0.0001989 111.00   <2e-16 ***
## avg_delay_all_departing  0.8528395  0.0068434 124.62   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.839 on 32697 degrees of freedom
## Multiple R-squared:  0.4544, Adjusted R-squared:  0.4543
## F-statistic: 1.361e+04 on 2 and 32697 DF,  p-value: < 2.2e-16
```

The estimated coefficient for journey\_time\_avg is 0.022, which means that for every one minute increase

in the length of the journey, there is a 0.022 minute increase in the delay in arrival at the destination when there is no change in the delay in departing.

## Question 2 (2 pts)

Some output for a **different model** is shown below. Using the output, predict the average arrival delay for a train route which has an average journey time of 200 minutes, has an average departure delay of 3 minutes, and took place in January (i.e., month == 1).

```
##                               Estimate Std. Error   t value   Pr(>|t|) 
## (Intercept)           -0.89153617  0.0625821018 -14.245865 6.529044e-46
## journey_time_avg      0.02215535  0.0001925167 115.082758 0.000000e+00
## avg_delay_all_departing 0.79854766  0.0067729269 117.902891 0.000000e+00
## month2                 0.45637989  0.0735194856  6.207605 5.444405e-10
## month3                 -0.47362582  0.0734957887 -6.444258 1.177825e-10
## month4                 -0.08049415  0.0735030454 -1.095113 2.734751e-01
## month5                  0.32511688  0.0735458330  4.420602 9.874291e-06
## month6                  1.86150670  0.0739347360 25.177701 1.475638e-138
## month7                  1.94714571  0.0742472465 26.225157 4.916823e-150
## month8                  0.75296105  0.0737335963 10.211913 1.908528e-24
## month9                  0.61577106  0.0735351836  8.373829 5.797172e-17
## month10                 0.93242762  0.0734938111 12.687158 8.508994e-37
## month11                 1.29220290  0.0736160461 17.553278 1.160711e-68
## month12                 0.13960420  0.0810379232  1.722702 8.495186e-02
```

**Question 2 Answer** The MLR finds the solution to the model:  $y_{\text{hat}} = b_0(\text{month}) + b_1x_1 + b_2x_2 + e$ . I will solve for the predicted average arrival, ignoring the error term. The month 1 is the reference term so we will just use the intercept

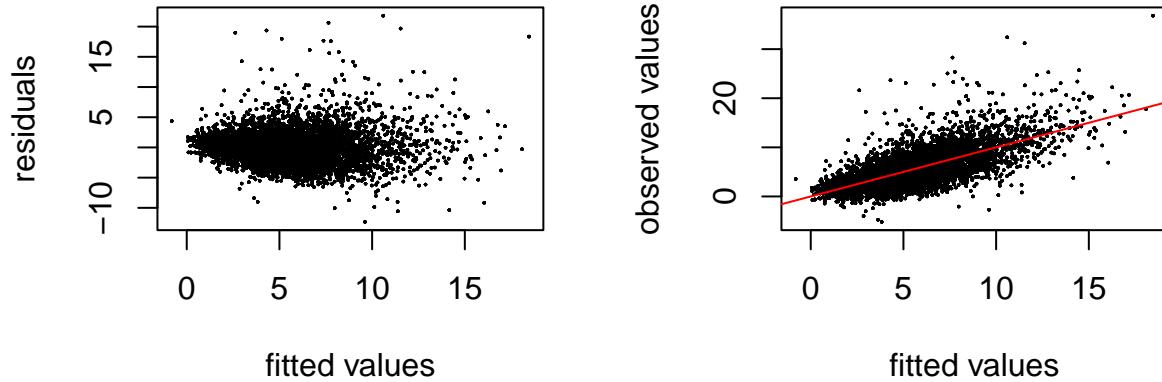
```
y_hat = -0.89153617 + 0.02215535 * 200 + 0.79854766 * 3
y_hat
```

```
## [1] 5.935177
```

The predicted average delay in arriving is 5.93 minutes.

## Question 3 (6 pts)

Do the assumptions for linear regression seem reasonable for the model fit in Question 2? Explain why or why not? You should use the plots below to justify your answer.



**Question 3 Answer** The assumptions to test are really the independence of errors. Here I list the key assumptions and describe my thoughts on them 1. independence of observations - the sampling of the data is independent train trips, so I think this is okay. 2. independence of errors- here, we see the errors generally fan out in a blob around the data. It looks like there is a slight increase in spread towards the greater fitted values, but I think that the data points affected by this are few enough that it does not violate the assumption. The majority of the weight of the data experiences equal variance of error 3. normality (which I looked at above) - while the distribution is not perfectly symmetric, it looks pretty plausible (smooth, unimodal) and does not violate the assumption

#### Question 4 (2 pts)

Suppose you think the association between arrival delay and journey time (i.e., the slope of journey time) may change from year to year. Fit a linear model below which would allow for that. For this problem, you **do not** need to consider adjusting for other variables in the model.

Let's introduce year as an interaction term.

```
mod <- lm(avg_delay_all_arriving ~ journey_time_avg * year, data=train_data)
summary(mod)
```

```
##
## Call:
## lm(formula = avg_delay_all_arriving ~ journey_time_avg * year,
##      data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.2122  -1.9821  -0.5073   1.4446  28.2437
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)               0.6554726  0.0872504  7.513 5.95e-14 ***
## journey_time_avg          0.0225342  0.0004863 46.341 < 2e-16 ***
```

```

## year2016          0.4856883  0.1231418   3.944 8.03e-05 ***
## year2017          1.5495532  0.1207531   12.832 < 2e-16 ***
## year2018          2.8680510  0.1179537   24.315 < 2e-16 ***
## journey_time_avg:year2016 -0.0021768  0.0006853  -3.176  0.00149 **
## journey_time_avg:year2017 -0.0027683  0.0006659  -4.157 3.23e-05 ***
## journey_time_avg:year2018 -0.0010376  0.0006418  -1.617  0.10597
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.27 on 32692 degrees of freedom
## Multiple R-squared:  0.2764, Adjusted R-squared:  0.2762
## F-statistic:  1784 on 7 and 32692 DF,  p-value: < 2.2e-16

```

#### Question 4 Answer

Here we see that year had a significant but small impact on the interaction between `journey_time_avg` and the delay in arrival in 2016 and 2017, but not significant in 2018. The effect of the interaction is much smaller than the effect of year or `journey_time_avg` independently, so we can conclude that the interaction effect is minimal.

#### Question 5 (3 pts)

Below, we fit a model which includes the covariates `journey time`, average departing delay and month. Suppose we want to test if the average arrival delay is associated with month after adjusting for `journey time` and average departure delay. For this problem, you don't need to consider interaction terms and you don't need to include other covariates. Describe how you would test this hypothesis. You don't need to actually perform any calculations or write any code, but specify which function in R you would use and be specific about what the inputs would be.

```

mod_year <- lm(avg_delay_all_arriving ~ journey_time_avg + avg_delay_all_departing + month,
                 data = train_data)
summary(mod_year)

```

```

##
## Call:
## lm(formula = avg_delay_all_arriving ~ journey_time_avg + avg_delay_all_departing +
##     month, data = train_data)
##
## Residuals:
##      Min       1Q       Median      3Q      Max 
## -12.3284  -1.6465  -0.1308   1.3595  21.8094 
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)    
## (Intercept)           -0.8915362  0.0625821 -14.246 < 2e-16 ***
## journey_time_avg       0.0221554  0.0001925 115.083 < 2e-16 ***
## avg_delay_all_departing 0.7985477  0.0067729 117.903 < 2e-16 ***
## month2                  0.4563799  0.0735195   6.208 5.44e-10 ***
## month3                 -0.4736258  0.0734958  -6.444 1.18e-10 ***
## month4                 -0.0804941  0.0735030  -1.095   0.273  
## month5                  0.3251169  0.0735458   4.421 9.87e-06 ***
## month6                  1.8615067  0.0739347  25.178 < 2e-16 ***

```

```

## month7          1.9471457  0.0742472  26.225 < 2e-16 ***
## month8          0.7529611  0.0737336  10.212 < 2e-16 ***
## month9          0.6157711  0.0735352   8.374 < 2e-16 ***
## month10         0.9324276  0.0734938  12.687 < 2e-16 ***
## month11         1.2922029  0.0736160  17.553 < 2e-16 ***
## month12         0.1396042  0.0810379   1.723   0.085 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.748 on 32686 degrees of freedom
## Multiple R-squared:  0.4892, Adjusted R-squared:  0.489
## F-statistic:  2408 on 13 and 32686 DF,  p-value: < 2.2e-16

```

**Question 5 answer** I would use this model exactly, which controls for journey time and delay in departing and tests the significance of the impact of month on arrival delay when there is no change in journey time or avg delay in departure. I would use the model provided, and read the p-values to make my decision as to the association of month. Since all months except 4 have a significant association, I conclude there is a significant effect.

## Question 6 (2 pt)

Suppose we fit the model below where we have used the log of journey\_time\_avg. Write 1 sentence interpreting the coefficient for journey time.

```

mod_log <- lm(avg_delay_all_arriving ~ log(journey_time_avg) ,
                data = train_data)
summary(mod_log)

```

```

##
## Call:
## lm(formula = avg_delay_all_arriving ~ log(journey_time_avg) ,
##      data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11.8915  -2.2847  -0.4961   1.5921  29.8838
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)           -11.06500   0.19401 -57.03 <2e-16 ***
## log(journey_time_avg)  3.29684   0.03868  85.23 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.477 on 32698 degrees of freedom
## Multiple R-squared:  0.1818, Adjusted R-squared:  0.1817
## F-statistic:  7264 on 1 and 32698 DF,  p-value: < 2.2e-16

```

```
3.29684 * log(1.01)
```

## Question 6 answer

```
## [1] 0.03280465
```

Since we log transform journey time, but not the response, I interpret this to mean two observations which differ by 1% journey time have expected avg delay of arriving that differs by 3.28%.

## Module 4 and 5

### Question 7 (3 pts)

In the model you fit in Question 1, each observation in the dataset corresponds to a specific route observed in a specific month and year. Thus each route appears in the data multiple times. Explain why this might violate an assumption for linear regression. How could you fix this? If your suggestion involves additional covariates or a different modeling assumption, be specific about what you mean (i.e., say what covariates would you include, or what model you would fit). There is more than 1 reasonable answer for this question, but just pick one.

**Question 7 answer** This might violate the assumption of independent observations, since there could be factors specific to each route that affect all observations made on that route, introducing covariance where none is assumed.

I could fix this by including a random effect for each departure or arrival station - or, perhaps, by creating a new variable composed of departure and arrival station, and including it as a random effect in the model.

```
train_data$route <- as.factor(paste(train_data$departure_station, train_data$arrival_station))
length(unique(train_data$route))

## [1] 130

library(lme4)

## Loading required package: Matrix

model <- lmer(avg_delay_all_arriving ~ journey_time_avg + avg_delay_all_departing + (1|route),
               data=train_data)
summary(model)

## Linear mixed model fit by REML ['lmerMod']
## Formula: avg_delay_all_arriving ~ journey_time_avg + avg_delay_all_departing +
##           (1 | route)
## Data: train_data
##
## REML criterion at convergence: 149862.3
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max 
## -6.3108 -0.5283 -0.0676  0.4201  9.9764 
## 
## Random effects:
```

```

## Groups     Name      Variance Std.Dev.
## route     (Intercept) 5.799    2.408
## Residual          5.600    2.366
## Number of obs: 32700, groups: route, 130
##
## Fixed effects:
##                               Estimate Std. Error t value
## (Intercept)           1.111310  0.241575   4.60
## journey_time_avg     0.007810  0.000655  11.92
## avg_delay_all_departing 1.282857  0.009630 133.22
##
## Correlation of Fixed Effects:
##             (Intr) jrnny_--
## jrnny_tm_vg -0.471
## avg_dly_ll_ -0.121  0.043

coef(model)$route

##                                         (Intercept) journey_time_avg
## AIX EN PROVENCE TGV PARIS LYON      -1.488149269  0.007810064
## ANGERS SAINT LAUD PARIS MONTPARNASSE   0.179147658  0.007810064
## ANGOULEME PARIS MONTPARNASSE        -0.761130301  0.007810064
## ANNECY PARIS LYON                   4.085535110  0.007810064
## ARRAS PARIS NORD                  -0.175240775  0.007810064
## AVIGNON TGV PARIS LYON            -1.834609379  0.007810064
## BARCELONA PARIS LYON              6.778021501  0.007810064
## BELLEGARDE (AIN) PARIS LYON       -1.247976395  0.007810064
## BESANCON FRANCHE COMTE TGV PARIS LYON -0.336430657  0.007810064
## BORDEAUX ST JEAN PARIS MONTPARNASSE   0.456095474  0.007810064
## BORDEAUX ST JEAN PARIS VAUGIRARD     0.261774948  0.007810064
## BORDEAUX ST JEAN TOURCOING         7.481812156  0.007810064
## BREST PARIS MONTPARNASSE            2.642236911  0.007810064
## CHAMBERY CHALLES LES EAUX PARIS LYON -2.593360664  0.007810064
## DIJON VILLE PARIS LYON             -3.431400378  0.007810064
## DOUAI PARIS NORD                  0.729029165  0.007810064
## DUNKERQUE PARIS NORD              1.605447341  0.007810064
## FRANCFORPT PARIS EST             -0.171297350  0.007810064
## GENEVE PARIS LYON                -0.290887161  0.007810064
## GRENOBLE PARIS LYON              0.576661743  0.007810064
## ITALIE PARIS LYON                2.060498890  0.007810064
## LA ROCHELLE VILLE PARIS MONTPARNASSE 2.159624441  0.007810064
## LAUSANNE PARIS LYON              -1.111183346  0.007810064
## LAVAL PARIS MONTPARNASSE          -0.301981789  0.007810064
## LE CREUSOT MONTCEAU MONTCHANIN PARIS LYON -4.096875133  0.007810064
## LE MANS PARIS MONTPARNASSE         -0.659816217  0.007810064
## LILLE LYON PART DIEU              0.898191094  0.007810064
## LILLE MARSEILLE ST CHARLES        5.414014320  0.007810064
## LILLE PARIS NORD                 -0.285447096  0.007810064
## LYON PART DIEU LILLE              -2.884296515  0.007810064
## LYON PART DIEU MARNE LA VALLEE    -2.853879689  0.007810064
## LYON PART DIEU MARSEILLE ST CHARLES -0.184104564  0.007810064
## LYON PART DIEU MONTPELLIER        -0.205448878  0.007810064
## LYON PART DIEU PARIS LYON          -1.868213779  0.007810064
## LYON PART DIEU RENNES              -1.048553258  0.007810064

```

## MACON LOCHE PARIS LYON	-5.310687509	0.007810064
## MADRID MARSEILLE ST CHARLES	8.603719226	0.007810064
## MARNE LA VALLEE LYON PART DIEU	-3.723563705	0.007810064
## MARNE LA VALLEE MARSEILLE ST CHARLES	0.243508411	0.007810064
## MARSEILLE ST CHARLES LILLE	6.075460948	0.007810064
## MARSEILLE ST CHARLES LYON PART DIEU	3.660179266	0.007810064
## MARSEILLE ST CHARLES MADRID	-0.191621964	0.007810064
## MARSEILLE ST CHARLES MARNE LA VALLEE	3.689143142	0.007810064
## MARSEILLE ST CHARLES PARIS LYON	0.359866966	0.007810064
## MARSEILLE ST CHARLES TOURCOING	3.516004943	0.007810064
## METZ PARIS EST	-0.663149314	0.007810064
## MONTPELLIER LYON PART DIEU	3.538498214	0.007810064
## MONTPELLIER PARIS LYON	-0.219732397	0.007810064
## MULHOUSE VILLE PARIS LYON	-1.861456723	0.007810064
## NANCY PARIS EST	-1.426431930	0.007810064
## NANTES PARIS MONTPARNASSE	1.511739470	0.007810064
## NANTES PARIS VAUGIRARD	0.938436601	0.007810064
## NANTES STRASBOURG	1.195196144	0.007810064
## NICE VILLE PARIS LYON	3.082660101	0.007810064
## NIMES PARIS LYON	-1.703155151	0.007810064
## PARIS EST FRANCFORT	2.104798785	0.007810064
## PARIS EST METZ	-0.269939388	0.007810064
## PARIS EST NANCY	-0.725092846	0.007810064
## PARIS EST REIMS	0.378502350	0.007810064
## PARIS EST STRASBOURG	0.003507406	0.007810064
## PARIS EST STUTTGART	1.716298163	0.007810064
## PARIS LYON AIX EN PROVENCE TGV	2.047981090	0.007810064
## PARIS LYON ANNECY	2.616783632	0.007810064
## PARIS LYON AVIGNON TGV	1.829245615	0.007810064
## PARIS LYON BARCELONA	6.988785874	0.007810064
## PARIS LYON BELLEGARDE (AIN)	3.132265040	0.007810064
## PARIS LYON BESANCON FRANCHE COMTE TGV	1.690872385	0.007810064
## PARIS LYON CHAMBERY CHALLES LES EAUX	2.220165128	0.007810064
## PARIS LYON DIJON VILLE	-0.436962144	0.007810064
## PARIS LYON GENEVE	1.962819271	0.007810064
## PARIS LYON GRENOBLE	1.825336453	0.007810064
## PARIS LYON ITALIE	3.028263884	0.007810064
## PARIS LYON LAUSANNE	2.974626392	0.007810064
## PARIS LYON LE CREUSOT MONTCEAU MONTCHANIN	1.057780242	0.007810064
## PARIS LYON LYON PART DIEU	0.166068235	0.007810064
## PARIS LYON MACON LOCHE	0.367223274	0.007810064
## PARIS LYON MARSEILLE ST CHARLES	2.371761586	0.007810064
## PARIS LYON MONTPELLIER	2.856321407	0.007810064
## PARIS LYON MULHOUSE VILLE	0.846402132	0.007810064
## PARIS LYON NICE VILLE	5.452120562	0.007810064
## PARIS LYON NIMES	2.967678664	0.007810064
## PARIS LYON PERPIGNAN	4.296214396	0.007810064
## PARIS LYON SAINT ETIENNE CHATEAUCREUX	2.890566545	0.007810064
## PARIS LYON TOULON	4.221430115	0.007810064
## PARIS LYON VALENCE ALIXAN TGV	1.163575878	0.007810064
## PARIS LYON ZURICH	1.607921974	0.007810064
## PARIS MONTPARNASSE ANGERS SAINT LAUD	2.161682774	0.007810064
## PARIS MONTPARNASSE ANGOULEME	1.805261518	0.007810064
## PARIS MONTPARNASSE BORDEAUX ST JEAN	2.498672313	0.007810064

## PARIS MONTPARNASSE BREST	1.737981149	0.007810064
## PARIS MONTPARNASSE LA ROCHELLE VILLE	1.434159087	0.007810064
## PARIS MONTPARNASSE LAVAL	1.109929645	0.007810064
## PARIS MONTPARNASSE LE MANS	0.773693378	0.007810064
## PARIS MONTPARNASSE NANTES	2.304359004	0.007810064
## PARIS MONTPARNASSE POITIERS	0.584874871	0.007810064
## PARIS MONTPARNASSE QUIMPER	1.635317932	0.007810064
## PARIS MONTPARNASSE RENNES	1.692933746	0.007810064
## PARIS MONTPARNASSE ST MALO	1.353109437	0.007810064
## PARIS MONTPARNASSE ST PIERRE DES CORPS	0.535431390	0.007810064
## PARIS MONTPARNASSE TOULOUSE MATABIAU	4.382524357	0.007810064
## PARIS MONTPARNASSE TOURS	1.530746817	0.007810064
## PARIS MONTPARNASSE VANNES	2.551685188	0.007810064
## PARIS NORD ARRAS	0.152288339	0.007810064
## PARIS NORD DOUAI	0.535664394	0.007810064
## PARIS NORD DUNKERQUE	-0.099736922	0.007810064
## PARIS NORD LILLE	-0.137818489	0.007810064
## PARIS VAUGIRARD BORDEAUX ST JEAN	-0.624507998	0.007810064
## PARIS VAUGIRARD NANTES	5.399235582	0.007810064
## PARIS VAUGIRARD RENNES	0.871104531	0.007810064
## PERPIGNAN PARIS LYON	2.763425909	0.007810064
## POITIERS PARIS MONTPARNASSE	-1.695897246	0.007810064
## QUIMPER PARIS MONTPARNASSE	2.538321307	0.007810064
## REIMS PARIS EST	-0.601128276	0.007810064
## RENNES LYON PART DIEU	3.150265019	0.007810064
## RENNES PARIS MONTPARNASSE	0.633931284	0.007810064
## RENNES PARIS VAUGIRARD	0.984902210	0.007810064
## SAINT ETIENNE CHATEAUCREUX PARIS LYON	2.262999970	0.007810064
## ST MALO PARIS MONTPARNASSE	2.255778745	0.007810064
## ST PIERRE DES CORPS PARIS MONTPARNASSE	-0.476725537	0.007810064
## STRASBOURG NANTES	3.528164241	0.007810064
## STRASBOURG PARIS EST	-1.449974136	0.007810064
## STUTTGART PARIS EST	-2.422669924	0.007810064
## TOULON PARIS LYON	-1.182688102	0.007810064
## TOULOUSE MATABIAU PARIS MONTPARNASSE	4.784046219	0.007810064
## TOURCOING BORDEAUX ST JEAN	4.265119762	0.007810064
## TOURCOING MARSEILLE ST CHARLES	1.228460819	0.007810064
## TOURS PARIS MONTPARNASSE	2.308276346	0.007810064
## VALENCE ALIXAN TGV PARIS LYON	-4.186364637	0.007810064
## VANNES PARIS MONTPARNASSE	1.301072081	0.007810064
## ZURICH PARIS LYON	-1.873405101	0.007810064
## avg_delay_all_departing		
## AIX EN PROVENCE TGV PARIS LYON	1.282857	
## ANGERS SAINT LAUD PARIS MONTPARNASSE	1.282857	
## ANGOULEME PARIS MONTPARNASSE	1.282857	
## ANNECY PARIS LYON	1.282857	
## ARRAS PARIS NORD	1.282857	
## AVIGNON TGV PARIS LYON	1.282857	
## BARCELONA PARIS LYON	1.282857	
## BELLEGARDE (AIN) PARIS LYON	1.282857	
## BESANCON FRANCHE COMTE TGV PARIS LYON	1.282857	
## BORDEAUX ST JEAN PARIS MONTPARNASSE	1.282857	
## BORDEAUX ST JEAN PARIS VAUGIRARD	1.282857	
## BORDEAUX ST JEAN TOURCOING	1.282857	

## BREST PARIS MONTPARNASSE	1.282857
## CHAMBERY CHALLES LES EAUX PARIS LYON	1.282857
## DIJON VILLE PARIS LYON	1.282857
## DOUAI PARIS NORD	1.282857
## DUNKERQUE PARIS NORD	1.282857
## FRANCFORPT PARIS EST	1.282857
## GENEVE PARIS LYON	1.282857
## GRENOBLE PARIS LYON	1.282857
## ITALIE PARIS LYON	1.282857
## LA ROCHELLE VILLE PARIS MONTPARNASSE	1.282857
## LAUSANNE PARIS LYON	1.282857
## LAVAL PARIS MONTPARNASSE	1.282857
## LE CREUSOT MONTCEAU MONTCHANIN PARIS LYON	1.282857
## LE MANS PARIS MONTPARNASSE	1.282857
## LILLE LYON PART DIEU	1.282857
## LILLE MARSEILLE ST CHARLES	1.282857
## LILLE PARIS NORD	1.282857
## LYON PART DIEU LILLE	1.282857
## LYON PART DIEU MARNE LA VALLEE	1.282857
## LYON PART DIEU MARSEILLE ST CHARLES	1.282857
## LYON PART DIEU MONTPELLIER	1.282857
## LYON PART DIEU PARIS LYON	1.282857
## LYON PART DIEU RENNES	1.282857
## MACON LOCHE PARIS LYON	1.282857
## MADRID MARSEILLE ST CHARLES	1.282857
## MARNE LA VALLEE LYON PART DIEU	1.282857
## MARNE LA VALLEE MARSEILLE ST CHARLES	1.282857
## MARSEILLE ST CHARLES LILLE	1.282857
## MARSEILLE ST CHARLES LYON PART DIEU	1.282857
## MARSEILLE ST CHARLES MADRID	1.282857
## MARSEILLE ST CHARLES MARNE LA VALLEE	1.282857
## MARSEILLE ST CHARLES PARIS LYON	1.282857
## MARSEILLE ST CHARLES TOURCOING	1.282857
## METZ PARIS EST	1.282857
## MONTPELLIER LYON PART DIEU	1.282857
## MONTPELLIER PARIS LYON	1.282857
## MULHOUSE VILLE PARIS LYON	1.282857
## NANCY PARIS EST	1.282857
## NANTES PARIS MONTPARNASSE	1.282857
## NANTES PARIS VAUGIRARD	1.282857
## NANTES STRASBOURG	1.282857
## NICE VILLE PARIS LYON	1.282857
## NIMES PARIS LYON	1.282857
## PARIS EST FRANCFORPT	1.282857
## PARIS EST METZ	1.282857
## PARIS EST NANCY	1.282857
## PARIS EST REIMS	1.282857
## PARIS EST STRASBOURG	1.282857
## PARIS EST STUTTGART	1.282857
## PARIS LYON AIX EN PROVENCE TGV	1.282857
## PARIS LYON ANNECY	1.282857
## PARIS LYON AVIGNON TGV	1.282857
## PARIS LYON BARCELONA	1.282857
## PARIS LYON BELLEGARDE (AIN)	1.282857

## PARIS LYON BESANCON FRANCHE COMTE TGV	1.282857
## PARIS LYON CHAMBERY CHALLES LES EAUX	1.282857
## PARIS LYON DIJON VILLE	1.282857
## PARIS LYON GENEVE	1.282857
## PARIS LYON GRENOBLE	1.282857
## PARIS LYON ITALIE	1.282857
## PARIS LYON LAUSANNE	1.282857
## PARIS LYON LE CREUSOT MONTCEAU MONTCHANIN	1.282857
## PARIS LYON LYON PART DIEU	1.282857
## PARIS LYON MACON LOCHE	1.282857
## PARIS LYON MARSEILLE ST CHARLES	1.282857
## PARIS LYON MONTPELLIER	1.282857
## PARIS LYON MULHOUSE VILLE	1.282857
## PARIS LYON NICE VILLE	1.282857
## PARIS LYON NIMES	1.282857
## PARIS LYON PERPIGNAN	1.282857
## PARIS LYON SAINT ETIENNE CHATEAUCREUX	1.282857
## PARIS LYON TOULON	1.282857
## PARIS VALENCE ALIXAN TGV	1.282857
## PARIS LYON ZURICH	1.282857
## PARIS MONTPARNASSE ANGERS SAINT LAUD	1.282857
## PARIS MONTPARNASSE ANGOULEME	1.282857
## PARIS MONTPARNASSE BORDEAUX ST JEAN	1.282857
## PARIS MONTPARNASSE BREST	1.282857
## PARIS MONTPARNASSE LA ROCHELLE VILLE	1.282857
## PARIS MONTPARNASSE Laval	1.282857
## PARIS MONTPARNASSE LE MANS	1.282857
## PARIS MONTPARNASSE NANTES	1.282857
## PARIS MONTPARNASSE POITIERS	1.282857
## PARIS MONTPARNASSE QUIMPER	1.282857
## PARIS MONTPARNASSE RENNES	1.282857
## PARIS MONTPARNASSE ST MALO	1.282857
## PARIS MONTPARNASSE ST PIERRE DES CORPS	1.282857
## PARIS MONTPARNASSE TOULOUSE Matabiau	1.282857
## PARIS MONTPARNASSE TOURS	1.282857
## PARIS MONTPARNASSE VANNES	1.282857
## PARIS NORD ARRAS	1.282857
## PARIS NORD DOUAI	1.282857
## PARIS NORD DUNKERQUE	1.282857
## PARIS NORD LILLE	1.282857
## PARIS VAUGIRARD BORDEAUX ST JEAN	1.282857
## PARIS VAUGIRARD NANTES	1.282857
## PARIS VAUGIRARD RENNES	1.282857
## PERPIGNAN PARIS LYON	1.282857
## POITIERS PARIS MONTPARNASSE	1.282857
## QUIMPER PARIS MONTPARNASSE	1.282857
## REIMS PARIS EST	1.282857
## RENNES LYON PART DIEU	1.282857
## RENNES PARIS MONTPARNASSE	1.282857
## RENNES PARIS VAUGIRARD	1.282857
## SAINT ETIENNE CHATEAUCREUX PARIS LYON	1.282857
## ST MALO PARIS MONTPARNASSE	1.282857
## ST PIERRE DES CORPS PARIS MONTPARNASSE	1.282857
## STRASBOURG NANTES	1.282857

```

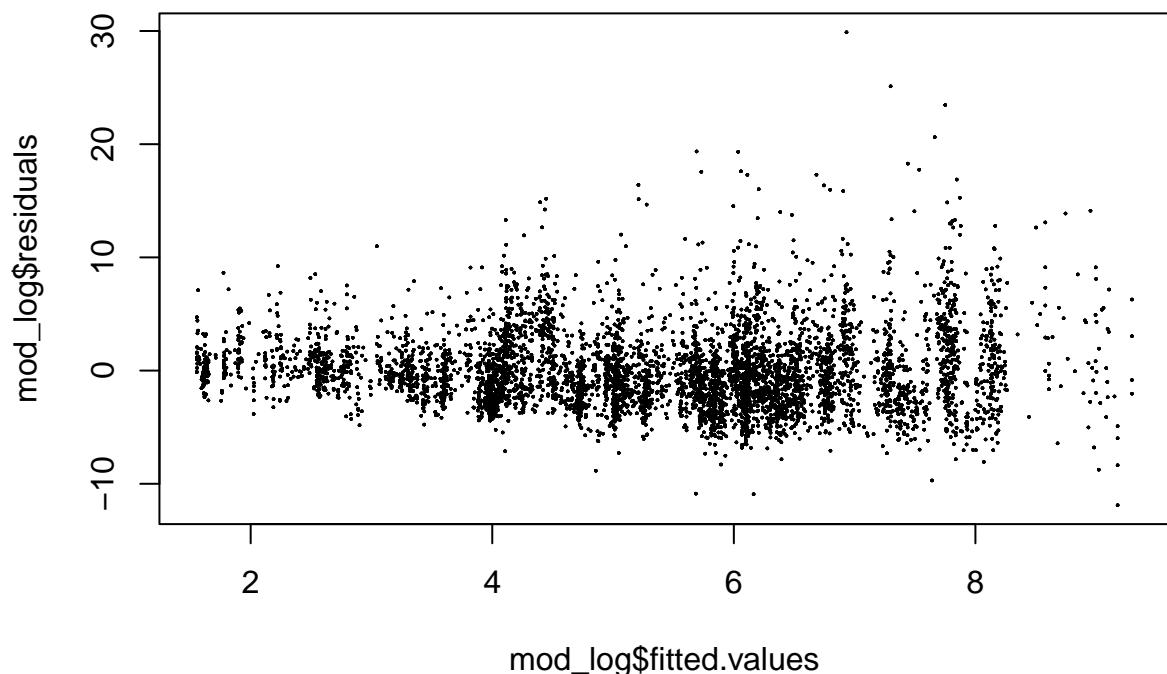
## STRASBOURG PARIS EST           1.282857
## STUTTGART PARIS EST           1.282857
## TOULON PARIS LYON             1.282857
## TOULOUSE MATABIAU PARIS MONTPARNASSE 1.282857
## TOURCOING BORDEAUX ST JEAN      1.282857
## TOURCOING MARSEILLE ST CHARLES 1.282857
## TOURS PARIS MONTPARNASSE       1.282857
## VALENCE ALIXAN TGV PARIS LYON   1.282857
## VANNES PARIS MONTPARNASSE      1.282857
## ZURICH PARIS LYON              1.282857

```

And it looks like there is quite a large degree of variation within these clusters.

### Question 8 (3 pts)

Using the model from Question 5, we plot the fitted values vs the residuals below. Explain why you might want to use robust standard errors. What might be the advantages and disadvantages of using the robust standard errors as opposed to the model based errors (the ones that come out of `summary`)?



```

##
## Call:
## lm(formula = avg_delay_all_arriving ~ log(journey_time_avg),
##     data = train_data)
##
## Residuals:

```

```

##      Min     1Q   Median     3Q    Max
## -11.8915 -2.2847 -0.4961  1.5921 29.8838
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)           -11.06500  0.19401 -57.03 <2e-16 ***
## log(journey_time_avg) 3.29684  0.03868  85.23 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.477 on 32698 degrees of freedom
## Multiple R-squared:  0.1818, Adjusted R-squared:  0.1817
## F-statistic:  7264 on 1 and 32698 DF, p-value: < 2.2e-16

```

**Question 8 answer** Here, the error structure shows a banding pattern, as well as a mild fanning pattern (the spread of the residuals increases as the fitted values increase). The violation of the assumption of homoskedasticity can lead to increased type I error rates and incorrect inference procedures (since our standard errors become biased) if we do not account for it.

The disadvantages are that with large sample sizes, hypothesis tests and CIs are valid even with the heteroskedasticity, so this may add compute and complexity unnecessarily. We might also lose power unnecessarily if our model is correct, and if the error in model set up results in bias in parameter estimates, then the robust estimators will not help, since they affect the standard error, not the coefficient estimates.

### Question 9 (3 pts)

Suppose you are taking a train tomorrow from Lille to Paris Nord and want to predict the delay in arrival. You want to be very sure about the prediction, so you gather data for 1000 different variables you think might be relevant (temperature, whether it is raining, GDP of France per month/year, the win/loss record of the soccer team in Lille, etc). You then regress average arrival delay onto all of those variables, and use it to predict the arrival delay for tomorrow's train. Explain why this might not give a good prediction. What might you do instead? 2-3 sentences for this answer is fine.

**Question 9 answer** Here, I introduce an over-fitting problem. My model might essentially learn the regression on the present variables really well, but not generalize well. What I would do instead is 1. used a penalized regression method like RIDGE or LASSO - with 1k parameters it is likely not all have equal predictive weight, so I would want to limit the impact of some. This would also help with model interpretability. I might also try a variable selection method, like forward or backward search. And 2. I would use cross-validation, repeatedly holding out subsets of the data, and measure MSE to determine which model is likely to generalize best.

## Module 6

For the following questions, suppose we are analyzing data for Big Red Airlines, Cornell's latest idea for getting people to and from Ithaca. The dependent variable is whether or not a flight took off on time. In the `OnTime` variable: 1 indicates that the flight took off on time, 0 indicates that it was delayed. The covariates we have recorded include Temperature (in degrees), `TimeOfDay` (Evening, Midday, Morning), and `Rain` (FALSE, TRUE).

```

airlineData <- read.csv("https://raw.githubusercontent.com/ysamwang/btry6020_sp22/main/lab11/airline.csv")
names(airlineData)

## [1] "OnTime"      "Temperature"   "TimeOfDay"     "Rain"

```

## Question 10 (2 pts)

What is the appropriate type of regression for modeling the binary data? What is being predicted by the linear model we are fitting? i.e., if the model we set up is

$$\hat{y} = b_0 + b_1 X_{1,i} + b_2 X_{2,i} \dots$$

what is on the left side of the equation (you can write it out in words instead of typing out the math)?.

**Question 10 answer** We are setting up a logistic regression, which aims to predict a 1 or 0 value of whether or not something will occur. It does this by taking the log of the exponential function:  $\log(\theta/(1-\theta))$ , where  $\theta$  is the result of the linear model. We are modeling the log odds of the event happening.

## Question 11 (2 pts)

We fit the model below. How would you interpret the coefficient associated with `Temperature`?

```

mod <- glm(OnTime ~ Temperature + TimeOfDay + Rain,
           data = airlineData, family = "binomial")
summary(mod)

##
## Call:
## glm(formula = OnTime ~ Temperature + TimeOfDay + Rain, family = "binomial",
##      data = airlineData)
##
## Deviance Residuals:
##    Min      1Q   Median      3Q      Max
## -1.3319 -0.6108 -0.4108 -0.2530  2.6456
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) 1.46094   0.37880  3.857 0.000115 ***
## Temperature -0.05248   0.00652 -8.050 8.29e-16 ***
## TimeOfDayMidday 0.18066   0.22978  0.786 0.431717
## TimeOfDayMorning -0.59611   0.25310 -2.355 0.018511 *
## RainTRUE      -0.52725   0.19977 -2.639 0.008308 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 765.46 on 892 degrees of freedom
## Residual deviance: 671.25 on 888 degrees of freedom
## AIC: 681.25
##
## Number of Fisher Scoring iterations: 5

```

$\exp(-0.05248)$

### Question 11 answer

```
## [1] 0.9488733
```

For this regression, the coefficient of temperature indicates that for an observation that only differs by an increase of 1 degree in temperature from another observation, the odds ratio of the second event occurring vs the first event is 0.9488. So there is a slight decrease in plane departure on time probability for an increase in temperature.

# **Final Project (30 pts)**

## **Introduction**

This final project is designed to demonstrate your mastery of linear regression techniques on real-world data. You will apply the theoretical concepts we've covered in class to a dataset of your choice, perform a comprehensive analysis, and present your findings in a professional format suitable for showcasing to potential employers.

## **Objectives**

By completing this project, you will:

- Apply linear regression techniques to solve real-world problems
- Demonstrate your ability to verify and address regression assumptions
- Perform meaningful feature selection and hypothesis testing
- Communicate the practical significance of your statistical findings
- Create a professional portfolio piece for future employment opportunities

## **Project Requirements**

### **Dataset Selection**

1. Choose a dataset from Kaggle
2. Your dataset must have a continuous target variable suitable for linear regression
3. The dataset should contain multiple potential predictor variables
4. Choose a dataset that interests you and has meaningful real-world applications

### **Analysis Requirements**

Your analysis must include the following components:

### **Exploratory Data Analysis**

- Summary statistics of variables
- Visualization of distributions and relationships
- Identification of missing values and outliers
- Data cleaning and preprocessing steps

### **Regression Assumptions Verification**

- Linearity assessment
- Normality of residuals
- Homoscedasticity (constant variance of residuals)
- Independence of observations
- Multicollinearity assessment

## **Assumption Violation Handling**

- Apply appropriate transformations when assumptions are violated
- Document your approach to each violation
- Compare models before and after corrections

## **Variable Selection & Hypothesis Testing**

- Implement at least two different variable selection techniques
- Perform hypothesis tests on coefficients
- Assess model performance with metrics ( $R^2$ , adjusted  $R^2$ , RMSE, etc.)
- Validate your model using appropriate cross-validation techniques

## **Feature Impact Analysis**

- Quantify and interpret the impact of each feature on the target
- Provide confidence intervals for significant coefficients
- Explain the practical significance of your findings in the context of the dataset

**Deliverables** GitHub Repository containing:

- All code (well-documented Rmd files)
- README.md with clear instructions on how to run your analysis
- Data folder (or instructions for accessing the data)
- Requirements.txt or environment.yml file

**Final Report (PDF) containing:**

- Introduction: dataset description and problem statement
- Methodology: techniques used and justification
- Results: findings from your analysis
- Discussion: interpretation of results and limitations
- Conclusion: summary and potential future work
- References: cite all sources used

## **Evaluation Criteria**

Your project will be evaluated based on:

- Correctness of statistical analysis and procedures
- Proper handling of regression assumptions
- Quality of variable selection and hypothesis testing
- Clarity of interpretation and insights
- Organization and documentation of code
- Professional presentation of findings

## **Timeline and Submission**

- Release Date: May 5th, 2025
- Due Date: Wednesday, May 14th, 2025 (11:59 PM EST)
- Submission: Email your GitHub repository link and PDF report to nbb45@cornell.edu with the subject line "Final Project - [Your Name]"

## **Resources**

- Course materials and lecture notes
- Kaggle Datasets
- GitHub tutorial and GitHub documentation for repository setup.

## **Academic Integrity**

This is an individual project. While you may discuss general concepts with classmates, all submitted work must be your own. Proper citation is required for any external resources used.

Good luck with your project! This is an opportunity to demonstrate your skills and create a valuable addition to your professional portfolio.

## **Finished**

You're done, congratulations!