

SCORE_Project

Avi Benki and Ryan Henry

2024-04-21

1. Learning Goals To conduct a multivariable analysis of a binary outcome, we must be able to: a.

Read data from a CSV file

b. Select key variables and format them to be computed

c. Create a logistic regression model

d. Analyze p-values and summarize results

2. Introduction Can someone predetermine which fighter will win a fight in UFC? Is it stance, height, weight, age, number of strikes landed, or some combination? Every fighter and their style is unique so what plays the most significant role in helping a fighter win a fight? Analyzing each UFC fight and its winner, we will use a logistic regression model to explain what factor(s) play the most significant role in determining who is more likely to win a given fight.

3. Data The data set that will be used is the “fights” data set. This includes over 6000 fights dating back to 1994 and summarizes all stats for each fighter during the fight such as height, weight, reach, stance, style as well as each fighters previous average number of take downs, significant strikes which is further broken down to shots to the head, body, shots in the clinch, etc. It also gives how each fight ended whether that be by knockout, split decision, unanimous decision, how many rounds the fight was scheduled for and what round it ended in. a few rows have been included below.

Downloading data

Use the code below to download the the UFC data to your computer. We will save the fight data in a data frame called fights, and information about the fighters in a data set called fighters. For simplicity the following tables make up all of the stats for the “Blue” fighter, the data also includes all of the same stats for the “Red” fighter.

```
library(tidyverse)
library(broom)
library(knitr)
library(rsample)
fights <- read.csv("C:/Users/avaneesh.benki/Documents/data.csv")
```

| R_fighter | B_fighter | Referee | date | location | Winner | title_bout | weight_class |
|--------------|---------------|---------------|------------|------------------------|--------|------------|--------------|
| Adrian Yanez | Gustavo Lopez | Chris Tognoni | 2021-03-20 | Las Vegas, Nevada, USA | Red | False | Bantamweight |
| Trevin Giles | Roman Dolidze | Herb Dean | 2021-03-20 | Las Vegas, Nevada, USA | Red | False | Middleweight |

| R_fighter | B_fighter | Referee | date | location | Winner | title_bout | weight_class |
|-------------|-----------------|-----------|------------|------------------------|--------|------------|--------------|
| Tai Tuivasa | Harry Hunsucker | Herb Dean | 2021-03-20 | Las Vegas, Nevada, USA | Red | False | Heavyweight |

| B_avg_KD | B_avg_opp_KD | B_avg_SIG_STR_pct | B_avg_opp_SIG_STR_pct | B_avg_TD_pct |
|----------|--------------|-------------------|-----------------------|--------------|
| 0.0 | 0 | 0.42 | 0.50 | 0.33 |
| 0.5 | 0 | 0.66 | 0.31 | 0.30 |
| NA | NA | NA | NA | NA |

| B_avg_opp_TD_pct | B_avg_SUB_ATT | B_avg_opp_SUB_ATT | B_avg_REV | B_avg_opp_REV |
|------------------|---------------|-------------------|-----------|---------------|
| 0.36 | 0.5 | 1 | 0 | 0 |
| 0.50 | 1.5 | 0 | 0 | 0 |
| NA | NA | NA | NA | NA |

| B_avg_SIG_STR_att | B_avg_SIG_STR_landed | B_avg_opp_SIG_STR_att | B_avg_opp_SIG_STR_landed |
|-------------------|----------------------|-----------------------|--------------------------|
| 50.0 | 20 | 84 | 45.0 |
| 65.5 | 35 | 50 | 16.5 |
| NA | NA | NA | NA |

| B_avg_TOTAL_STR_att | B_avg_TOTAL_STR_landed | B_avg_opp_TOTAL_STR_att |
|---------------------|------------------------|-------------------------|
| 76.5 | 41.0 | 114.0 |
| 113.5 | 68.5 | 68.5 |
| NA | NA | NA |

| B_avg_opp_TOTAL_STR_landed | B_avg_TD_att | B_avg_opp_TD_landed | B_avg_opp_TD_att |
|----------------------------|--------------|---------------------|------------------|
| 64 | 1.5 | 6.5 | 9.0 |
| 29 | 2.5 | 0.5 | 0.5 |
| NA | NA | NA | NA |

| B_avg_HEAD_att | B_avg_HEAD_landed | B_avg_opp_HEAD_att | B_avg_opp_HEAD_landed |
|----------------|-------------------|--------------------|-----------------------|
| 39.5 | 11 | 63 | 27.5 |
| 46.0 | 20 | 36 | 7.5 |
| NA | NA | NA | NA |

| B_avg_BODY_att | B_avg_BODY_landed | B_avg_opp_BODY_att | B_avg_opp_BODY_landed |
|----------------|-------------------|--------------------|-----------------------|
| 7.5 | 7 | 12 | 9 |
| 12.0 | 8 | 8 | 3 |
| NA | NA | NA | NA |

| B_avg_LEG_att | B_avg_LEG_landed | B_avg_opp_LEG_att | B_avg_opp_LEG_landed |
|---------------|------------------|-------------------|----------------------|
| 3.0 | 2 | 9 | 8.5 |
| 7.5 | 7 | 6 | 6.0 |
| NA | NA | NA | NA |

| B_avg_GROUND_att | B_avg_GROUND_landed | B_avg_opp_GROUND_att | B_avg_opp_GROUND_landed |
|------------------|---------------------|----------------------|-------------------------|
| 4.5 | 3.0 | 36.5 | 24.5 |
| 7.0 | 4.5 | 1.5 | 0.5 |
| NA | NA | NA | NA |

| B_avg_CTRL_time.seconds. | B_avg_opp_CTRL_time.seconds. | B_total_time_fought.seconds. |
|--------------------------|------------------------------|------------------------------|
| 34.0 | 277.5 | 531.5 |
| 219.5 | 24.5 | 577.5 |
| NA | NA | NA |

| B_total_rounds_fought | B_total_title_bouts | B_current_win_streak | B_current_lose_streak |
|-----------------------|---------------------|----------------------|-----------------------|
| 4 | 0 | 0 | 1 |
| 4 | 0 | 2 | 0 |
| 0 | 0 | 0 | 0 |

| B_longest_win_streak | B_wins | B_losses | B_draw | B_win_by_Decision_Majority |
|----------------------|--------|----------|--------|----------------------------|
| 1 | 1 | 1 | 0 | 0 |
| 2 | 2 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 |

| B_win_by_Decision_Split | B_win_by_Decision_Unanimous | B_win_by_KO.TKO | B_win_by_Submission |
|-------------------------|-----------------------------|-----------------|---------------------|
| 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 |

| B_win_by_TKO_Doctor_Stoppage | B_Stance | B_Height_cms | B_Reach_cms | B_Weight_lbs |
|------------------------------|----------|--------------|-------------|--------------|
| 0 | Orthodox | 165.10 | 170.18 | 135 |
| 0 | Orthodox | 187.96 | 193.04 | 205 |
| 0 | Orthodox | 187.96 | 190.50 | 241 |

4. Methods and Instructional Content It will be useful while constructing our SCORE Module to reference other work that has been done on teaching statistics, especially in mixed martial arts. The first source that we found is a Kaggle tutorial at this link : <https://www.kaggle.com/code/olawaleoladipo/ufc-data-analysis-visualization-beginner> . Although the tutorial works through python and not R, it is useful as an example of how to introduce statistical concepts to individuals who might not have much experience with statistics, or even MMA. A very interesting part of the tutorial worth looking at is the last section in which many

different MMA statistics are visualized in a variety of models. This will be useful to emulate when making our SCORE module so that the student can help portray their findings in a digestible way.

Another useful reference to use, can be found at this link : <https://www.datacamp.com/tutorial/generalized-linear-models>. This is a very accessible and interesting tutorial on Generative Linear Models in R that starts with the expectation that one knows almost nothing about them. It is written and implemented in a very understandable and fun way, which is how we want our own module to end up as well.

The fights data set has a fight in each of row, with fight statistics given in terms of Blue corner fighter (B_fighter) and red corner fighter (R_fighter). The statistics are given as an average of how they have done in their fights historically as well as how well their opponents have done against them, so if the Blue corner fighter landed 8,4, and 6 head kicks each in their only three fights, which option below would we expect to see in the data frame?

- a. B_avg_opp_HEAD_landed = 4
- b. B_avg_HEAD_att = 6
- c. B_avg_HEAD_landed = 6
- d. B_avg_HEAD_landed = 8,4,6

The answer is c.

Data cleaning

Now that we have a basic understanding of how the data is laid out, lets combine the datasets to create a tailored data frame to make our analysis easier. The first thing we can do is remove every time the fight ended in a draw or a no contest. This can be done by using the filter function.

First we need to figure out what the possible options are in the “Winner category”

```
fights %>%  
  select(Winner) %>%  
  unique()
```

```
##      Winner  
## 1      Red  
## 4      Blue  
## 13     Draw
```

Now we can take out the Draws.

```
fights_new <- fights %>%  
  filter(Winner != "Draw")
```

In order to continue to hone our dataset, lets restrict the data set to only include the columns that we will be using in our analysis. This can be done using the select function.

```
fights_new <- fights_new %>%  
  select(B_fighter, R_fighter, Winner, date, B_avg_SIG_STR_pct, B_avg_TD_pct, B_avg_SIG_STR_landed, B_a
```

Lets take a look at the reduced data set

```

tfnew1 <- fights_new %>%
  select(B_fighter, R_fighter, Winner, date, B_avg_SIG_STR_pct, B_avg_TD_pct) %>%
  head(5) %>%
  kable()
tfnew2 <- fights_new %>%
  select(B_avg_SIG_STR_landed, B_avg_opp_SIG_STR_landed, B_avg_TD_landed, B_wins) %>%
  head(5) %>%
  kable()
tfnew3 <- fights_new %>%
  select(B_losses, B_Stance, B_Height_cms, B_Reach_cms, B_Weight_lbs, R_avg_SIG_STR_pct) %>%
  head(5) %>%
  kable()
tfnew4 <- fights_new %>%
  select(R_avg_TD_pct, R_avg_TD_landed, R_wins, R_losses, R_Stance, R_Height_cms) %>%
  head(5) %>%
  kable()
tfnew5 <- fights_new %>%
  select(R_Reach_cms, R_Weight_lbs, B_avg_SIG_STR_landed, B_avg_opp_SIG_STR_landed) %>%
  head(5) %>%
  kable()
tfnew1

```

| B_fighter | R_fighter | Winner | date | B_avg_SIG_STR_pct | B_avg_TD_pct |
|-------------------|---------------|--------|------------|-------------------|--------------|
| Gustavo Lopez | Adrian Yanez | Red | 2021-03-20 | 0.420000 | 0.330 |
| Roman Dolidze | Trevin Giles | Red | 2021-03-20 | 0.660000 | 0.300 |
| Harry Hunsucker | Tai Tuivasa | Red | 2021-03-20 | NA | NA |
| Montserrat Conejo | Cheyenne Buys | Blue | 2021-03-20 | NA | NA |
| Macy Chiasson | Marion Reneau | Blue | 2021-03-20 | 0.535625 | 0.185 |

tfnew2

| B_avg_SIG_STR_landed | B_avg_opp_SIG_STR_landed | B_avg_TD_landed | B_wins |
|----------------------|--------------------------|-----------------|--------|
| 20.0000 | 45.0000 | 1.0 | 1 |
| 35.0000 | 16.5000 | 1.5 | 2 |
| NA | NA | NA | 0 |
| NA | NA | NA | 0 |
| 57.9375 | 28.4375 | 1.5 | 4 |

tfnew3

| B_losses | B_Stance | B_Height_cms | B_Reach_cms | B_Weight_lbs | R_avg_SIG_STR_pct |
|----------|----------|--------------|-------------|--------------|-------------------|
| 1 | Orthodox | 165.10 | 170.18 | 135 | 0.5000000 |
| 0 | Orthodox | 187.96 | 193.04 | 205 | 0.5768750 |
| 0 | Orthodox | 187.96 | 190.50 | 241 | 0.5389063 |
| 0 | Southpaw | 152.40 | 154.94 | 115 | NA |
| 1 | Orthodox | 180.34 | 182.88 | 135 | 0.4030762 |

tfnew4

| R_avg_TD_pct | R_avg_TD_landed | R_wins | R_losses | R_Stance | R_Height_cms |
|--------------|-----------------|--------|----------|----------|--------------|
| 0.0000000 | 0.000000 | 1 | 0 | Orthodox | 170.18 |
| 0.4062500 | 0.781250 | 4 | 2 | Orthodox | 182.88 |
| 0.0000000 | 0.000000 | 4 | 3 | Southpaw | 187.96 |
| NA | NA | 0 | 0 | Switch | 160.02 |
| 0.5117188 | 1.261719 | 5 | 6 | Orthodox | 167.64 |

tfnew5

| R_Reach_cms | R_Weight_lbs | B_avg_SIG_STR_landed | B_avg_opp_SIG_STR_landed |
|-------------|--------------|----------------------|--------------------------|
| 177.80 | 135 | 20.0000 | 45.0000 |
| 187.96 | 185 | 35.0000 | 16.5000 |
| 190.50 | 264 | NA | NA |
| 160.02 | 115 | NA | NA |
| 172.72 | 135 | 57.9375 | 28.4375 |

Which of the following statistics is not considered by our new dataset?

- a. average takedowns by the opponents of the Blue Corner fighter
- b. Red corner fighter's average reversals
- c. Blue corner fighter's stance
- d. Average significant strike percentage by red corner fighter

answer b.

In order to simplify our analysis, we look at the data as differences between the red and blue corner fighters. Using significant strikes as an example, if the red corner averages 30 per fight and the blue corner averages 45, the new data point would be -15. This will be useful when trying to find a regression line that considers the winner of the fight with the differences in different aspects of fighting. By analyzing the differences in statistics between fighters, we make it easier to determine if there is a positive or negative linear correlation to winning. The assumption is that typically the taller fighter with more reach who lands the most significant strikes will likely be the winner. We'll look at the differences of Red - Blue and see if as the height gap, gap between the fighters number of take downs, and other variables grows, so does the probability of the red fighter winning. One important thing to mention is that for all the fights recorded in this data set is that typically the current champion or fighter with more wins is represented by the red corner. To do this, we will utilize the mutate function.

```
fightes_new <- fights_new %>%  
  mutate(diff_avg_SIG_STR_pct = R_avg_SIG_STR_pct - B_avg_SIG_STR_pct,  
         diff_avg_TD_pct = R_avg_TD_pct - B_avg_TD_pct,  
         diff_avg_SIG_STR_landed = R_avg_SIG_STR_landed - B_avg_SIG_STR_landed,  
         diff_avg_opp_SIG_STR_landed = R_avg_opp_SIG_STR_landed - B_avg_opp_SIG_STR_landed,  
         diff_avg_TD_landed = R_avg_TD_landed - B_avg_TD_landed,  
         diff_age = R_age - B_age,  
         diff_Reach_cms = R_Reach_cms - B_Reach_cms,  
         diff_height = R_Height_cms - B_Height_cms,  
         diff_wins = R_wins - B_wins,
```

```

    diff_losses = R_losses - B_losses,
    diff_fights = (R_wins + R_losses) - (B_wins + B_losses))

fights_diff <- fights_new %>%
  select(R_fighter, B_fighter, Winner, date, diff_avg_SIG_STR_pct, diff_avg_TD_pct, diff_avg_TD_landed,

```

We now have all the same data, but now in the form of red fighter minus blue fighter. By running a logistic regression on all the times that the red fighter won, we will be able to discover how much the differences between the fighters impact a win.

Lets make one final change and remove all the rows that contain N/A for any value.

```
Final_df <- na.omit(fights_diff)
```

Lets take one last look at the data set.

```

tfinal1 <- Final_df %>%
  select(R_fighter, B_fighter, Winner, date, diff_avg_SIG_STR_pct, diff_avg_TD_pct) %>%
  head(5) %>%
  kable()
tfinal2 <- Final_df %>%
  select(diff_avg_SIG_STR_landed, diff_avg_opp_SIG_STR_landed, diff_avg_TD_landed) %>%
  head(5) %>%
  kable()
tfinal3 <- Final_df %>%
  select(diff_age, diff_Reach_cms, diff_height, diff_wins, diff_losses, diff_fights) %>%
  head(5) %>%
  kable()
tfinal1

```

| | R_fighter | B_fighter | Winner | date | diff_avg_SIG_STR_pct | diff_avg_TD_pct |
|---|-----------------|---------------|--------|------------|----------------------|-----------------|
| 1 | Adrian Yanez | Gustavo Lopez | Red | 2021-03-20 | 0.0800000 | -0.3300000 |
| 2 | Trevin Giles | Roman Dolidze | Red | 2021-03-20 | -0.0831250 | 0.1062500 |
| 5 | Marion Reneau | Macy Chiasson | Blue | 2021-03-20 | -0.1325488 | 0.3267188 |
| 6 | Leonardo Santos | Grant Dawson | Blue | 2021-03-20 | 0.0501563 | -0.0979688 |
| 7 | Song Kenan | Max Griffin | Blue | 2021-03-20 | 0.0338477 | -0.3221875 |

```
tfinal2
```

| | diff_avg_SIG_STR_landed | diff_avg_opp_SIG_STR_landed | diff_avg_TD_landed |
|---|-------------------------|-----------------------------|--------------------|
| 1 | -3.00000 | -39.000000 | -1.0000000 |
| 2 | 8.15625 | 11.093750 | -0.7187500 |
| 5 | -13.57520 | 56.117188 | -0.2382812 |
| 6 | 1.43750 | 6.921875 | -1.9921875 |
| 7 | -26.61328 | -24.048828 | -1.4531250 |

tfinal3

| | diff_age | diff_Reach_cms | diff_height | diff_wins | diff_losses | diff_fights |
|---|----------|----------------|-------------|-----------|-------------|-------------|
| 1 | -4 | 7.62 | 5.08 | 0 | -1 | -1 |
| 2 | -4 | -5.08 | -5.08 | 2 | 2 | 4 |
| 5 | 14 | -10.16 | -12.70 | 1 | 5 | 6 |
| 6 | 14 | 7.62 | 5.08 | 3 | 1 | 4 |
| 7 | -4 | -12.70 | 2.54 | 0 | -5 | -5 |

Nice! Now that we have our data all cleaned up, lets begin analyzing what qualities makes a fighter more likely to win. Look at all the columns and make a prediction as to what factor(s) are the most significant to winning a UFC fight.

ANS: (ex. height and significant strikes landed)

One thing statisticians do when analyzing data is splitting it into training and test groups. This means that a model is trained with only the data in the training group, and then its accuracy is checked against the test group. For this dataset, we will be using a 75/25 split, where the data will be trained on a randomly selected 75% of the rows.

```
dt = sort(sample(nrow(Final_df), nrow(Final_df)*.75))
traindiff<-Final_df[dt,]
testdiff<-Final_df[-dt,]
```

To analyze the data with the training set, we will use a linear regression model. This means that we are assuming all dependent variables have a linear relationship with winning. For example, as the difference in height increases so does the likelihood of a fighter winning. Here is an example using the linear regression model for one variable.

```
height_model<- traindiff %>%
  glm(Winner == "Red" ~ diff_height,
      data = .,
      family = "binomial")
height_model %>% tidy() %>% kable()
```

| term | estimate | std.error | statistic | p.value |
|-------------|-----------|-----------|-----------|-----------|
| (Intercept) | 0.5269950 | 0.0385790 | 13.660146 | 0.0000000 |
| diff_height | 0.0131381 | 0.0061424 | 2.138927 | 0.0324415 |

What does the p-value say about the difference in height? Provide a brief description below on whether a difference in height leads to a better chance of winning a UFC fight.

ANS

It appears that there is very strong evidence that a positive difference in height contributes to a better chance of winning because of the low p-value and the positive slope.

Now lets look at takedowns and strikes (qualities not controlled by genetics.)


```
skill_model <- traindiff %>%
  glm(Winner == "Red" ~ diff_avg_SIG_STR_pct + diff_avg_TD_pct + diff_avg_SIG_STR_landed + diff_avg_opp
      data = .,
      family = "binomial")
summary(skill_model)
```

```
##
## Call:
## glm(formula = Winner == "Red" ~ diff_avg_SIG_STR_pct + diff_avg_TD_pct +
##      diff_avg_SIG_STR_landed + diff_avg_opp_SIG_STR_landed + diff_avg_TD_landed,
##      family = "binomial", data = .)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      0.528312   0.039036  13.534 < 2e-16 ***
## diff_avg_SIG_STR_pct      0.173526   0.252514   0.687  0.49196
## diff_avg_TD_pct          0.018061   0.136082   0.133  0.89441
## diff_avg_SIG_STR_landed    0.009157   0.001847   4.958 7.11e-07 ***
## diff_avg_opp_SIG_STR_landed -0.009726   0.001904  -5.108 3.26e-07 ***
## diff_avg_TD_landed        0.087446   0.027464   3.184  0.00145 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 3809.0  on 2885  degrees of freedom
## Residual deviance: 3749.2  on 2880  degrees of freedom
## AIC: 3761.2
##
## Number of Fisher Scoring iterations: 4
```

Does anything stand out to you? Take a closer look at the p-value for average significant strikes landed by an opponent. Does the significance of this variable surprise you?

Lets create a model for the genetic variables and the fighter's current record at the time of the fight.

```
gen_model <- traindiff %>%
  glm(Winner == "Red" ~ diff_age + diff_Reach_cms + diff_height + diff_wins + diff_losses,
      data = .,
      family = "binomial")
summary(gen_model)
```

```
##
## Call:
## glm(formula = Winner == "Red" ~ diff_age + diff_Reach_cms + diff_height +
##      diff_wins + diff_losses, family = "binomial", data = .)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      0.562567   0.041740  13.478 < 2e-16 ***
## diff_age         -0.063153   0.008482  -7.446 9.65e-14 ***
## diff_Reach_cms    0.020353   0.006258   3.252  0.00114 **
## diff_height       -0.016278   0.008178  -1.991  0.04653 *
```

```
## diff_wins      0.034895    0.013795    2.529  0.01142 *
## diff_losses    -0.091790    0.020378   -4.504 6.66e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 3809.0  on 2885  degrees of freedom
## Residual deviance: 3677.9  on 2880  degrees of freedom
## AIC: 3689.9
##
## Number of Fisher Scoring iterations: 4
```

These variable seem to carry much more significance than the previous variables. Why do you think that is? Finally lets complete out model with all 11 variables.

```
final_model <- traindiff %>%
  glm(Winner == "Red" ~ diff_avg_SIG_STR_pct + diff_avg_TD_pct + diff_avg_SIG_STR_landed + diff_avg_opp_SIG_STR_landed +
    data = .,
    family = "binomial")
summary(final_model)
```

```
##
## Call:
## glm(formula = Winner == "Red" ~ diff_avg_SIG_STR_pct + diff_avg_TD_pct +
##      diff_avg_SIG_STR_landed + diff_avg_opp_SIG_STR_landed + diff_avg_TD_landed +
##      diff_age + diff_Reach_cms + diff_height + diff_wins + diff_losses,
##      family = "binomial", data = .)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      0.564435   0.042018  13.433 < 2e-16 ***
## diff_avg_SIG_STR_pct      0.165892   0.257595   0.644 0.519573
## diff_avg_TD_pct      -0.033059   0.138449  -0.239 0.811273
## diff_avg_SIG_STR_landed    0.006263   0.001914   3.272 0.001069 **
## diff_avg_opp_SIG_STR_landed -0.006244   0.001975  -3.161 0.001571 **
## diff_avg_TD_landed      0.081354   0.028250   2.880 0.003979 **
## diff_age      -0.059018   0.008542  -6.909 4.87e-12 ***
## diff_Reach_cms      0.019840   0.006318   3.140 0.001689 **
## diff_height      -0.014354   0.008267  -1.736 0.082487 .
## diff_wins      0.024447   0.014036   1.742 0.081543 .
## diff_losses      -0.070067   0.020844  -3.361 0.000775 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 3809.0  on 2885  degrees of freedom
## Residual deviance: 3648.6  on 2875  degrees of freedom
## AIC: 3670.6
##
## Number of Fisher Scoring iterations: 4
```

Provide a brief summary of which variables are significant and why you think these variables are important. Look at both the p-value and the estimated effect on the model.

Now that we have generated a model from our data, let's try to utilize it to predict winners of fights that have not happened. Before we do this, it might be helpful to visualize what our models means. Let's take a look at the first model that we created which tried to determine the winner based only on height. By creating a list of several possible height differences, and utilizing our model to predict the odds of victory for these heights, we will get a look at what the prediction curve looks like. <https://www.theanalysisfactor.com/r-glm-plotting/>
Let's make a list of height differences from -15 cm to 15 cm, with a step size of 0.1 cm.

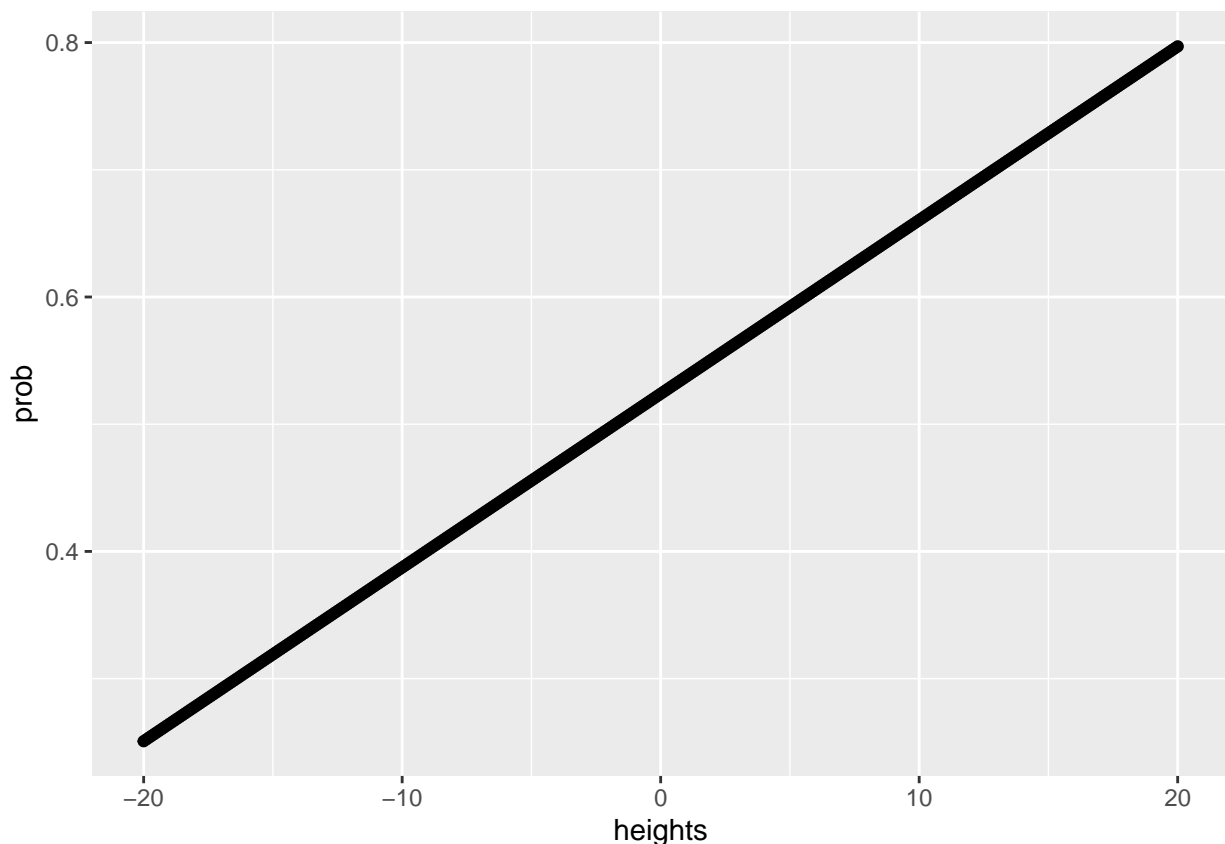
```
heights <- seq(-20,20,0.1) %>% data.frame(.) %>% rename(.,heights = .)
```

The predict function will use the intercept and slope that we found above to determine a win probability for each of these heights.

```
predheights <- heights %>% mutate(prob = 0.5239208 + 0.0136545*heights)
```

now we can plot the results

```
ggplot(predheights, aes(heights, prob))+ geom_point()
```



As expected, as the height difference goes up, so does the expected probability of a win. Furthermore, a height difference of 0 gives a probability of around 0.5 which makes sense as neither fighter has the advantage.

Okay, now let's try to complete the same process for the model that we created that looks at all the factors. We will create a data frame to simulate a mismatch, in which one fighter will have a big advantage in every category that we are looking at. Let's see how it goes.

```
mismatchdata <- data.frame(diff_avg_SIG_STR_pct = 0.2, diff_avg_TD_pct = 0.2, diff_avg_SIG_STR_landed = 0.2,
                           diff_avg_TD_landed = 2, diff_age = -7, diff_Reach_cms = 6, diff_height = 6, diff_weight_lbs = 10)
```

Now lets predict!

```
mismatchprediction <- predict(final_model, newdata = mismatchdata, type = "response")
mismatchprediction
```

```
##          1
## 0.894925
```

Wow! It looks like our model predicts an 88% chance for the fighter with all the advantages to win! Now lets go back and use the test data set to find how many times the model's predicted winner matched up with the actual winner. The code below will create a new dataframe containing the red fighters win probability based on the chosen variables.

```
testing <- final_model %>%
  augment(type.predict = "response",
          newdata = testdiff)
```

Now, we will assign a 1 for every time the actual winner of the fight matched the fighter which our model gave a better chance of winning.

```
testing2 <- testing %>%
  mutate(PredWin = ifelse((.fitted >= .5), "Red", "Blue"))
testing2 <- testing2 %>%
  mutate(correct = ifelse(Winner == PredWin, 1, 0))
```

Now we can find a proportion of times that our model gave the advantage to the real winner.

```
score <- testing2 %>% pull(correct)
correct <- sum(score)
total <- length(score)
correct/total
```

```
## [1] 0.6476091
```

It seems like our model was right about 66% of the time. Not Bad!

To wrap up our module we will do some quick analysis on which of the models that we created were the most accurate. If you need a refresher, go back above and take a look at the four models that we trained based on different combinations of variables. What were these four models?

ANS:

height, genetic/record, skill, final/everything

The two metrics that we will use for this analysis are AICs and BICs. AICS (Akaike Information Criterion) and BICs (Bayesian Information Criterion) both measure how well the model fits the data and put this analysis into a single score. The difference between the two is the BICs more heavily penalizes complexity in the model since complexity can introduce more opportunities for error. A lower AICS and BICS score indicates a better fitting/ more accurate model.

```
AICs = c("AICheight" = AIC(height_model), "AICgenetics" = AIC(gen_model), "AICskills" = AIC(skill_model), "AICfinal" = AIC(final_model))
BICs = c("BICheight" = BIC(height_model), "BICgenetics" = BIC(gen_model), "BICskills" = BIC(skill_model), "BICfinal" = BIC(final_model))
```

```
##      AICheight AICgenetics AICskills AICfinal
##      3808.376   3689.881   3761.171   3670.555
```

```
BICs
```

```
##      BICheight BICgenetics BICskills BICfinal
##      3820.312   3725.687   3796.976   3736.199
```

Based on the results, which of our models were the most fitting?

ANS: the final models had the lowest scores for both.

6. Conclusion In examining fights throughout the history of the UFC, we sought to determine what factors makes a fighter more likely to win. We observed 11 different factors and found the most significant were the fighters current number of wins and losses, height, age, and the average number of significant strikes an opponent lands on them per fight. We found that you fighters who have a height advantage against their opponent are typically more successful. A fighters ability to avoid strikes also leads to a fighters likelihood of winning. Most importantly a fighters current record is the most relevant variable in determining which fighter will win a UFC fight. The more wins the better the chance and the less losses the better the chance. We found that the final model that included the most factors was the most accurate, and gave the advantage to the actual winner about 66% of the time.