

Final Project

Yifei Yang, Yiming Yuan

Read in the data

```
library(skimr)

library(ggplot2)
library(dplyr)

library(readr)

library(MASS)

ufc <- read_csv("data/ufc-fighters-statistics.csv")
```

Introduction and data

Around 300 million individuals worldwide identify themselves as fans of Mixed Martial Arts (MMA), with its popularity peaking in nations such as the United States, the United Kingdom, Brazil, Singapore, and China. The Ultimate Fighting Championship (UFC) is the premier organization in the MMA world. Our motivation is to research on what contributes to the fighters' performance. Thus, the research questions are: How largely do the innate physical attributes affect fighters' performance? What's the most effective stance? How does the way fighters stroke (speed/accuracy/amount/defence success) influence their performance?

Today's data are UFC fighter statistics including:

wins: The number of wins the fighter has in their career.

draws: The number of draws the fighter has in their career.

height_cm: The height of the fighter in centimeters.

weight_in_kg: The weight of the fighter in kilograms.

reach_in_cm: The reach of the fighter in centimeters.

stance: The fighting stance of the fighter (Orthodox/Southpaw/Switch).

significant_strikes_landed_per_minute: The average number of significant strikes landed by the fighter per minute.

significant_striking_accuracy: The percentage of significant strikes that land successfully for the fighter.

significant_strikes_absorbed_per_minute: The average number of significant strikes absorbed by the fighter per minute.

significant_strike_defence: The percentage of opponent's significant strikes that the fighter successfully defends.

average_takedowns_landed_per_15_minutes: The average number of takedowns landed by the fighter per 15 minutes.

takedown_accuracy: The percentage of takedown attempts that are successful for the fighter.

takedown_defense: The percentage of opponent's takedown attempts that the fighter successfully defends.

average_submissions_attempted_per_15_minutes: The average number of submission attempts made by the fighter per 15 minutes.

Our response variable is wins, predictors are

sources:

<https://www.kaggle.com/datasets/aaronfriasr/ufc-fighters-statistics?resource=download>

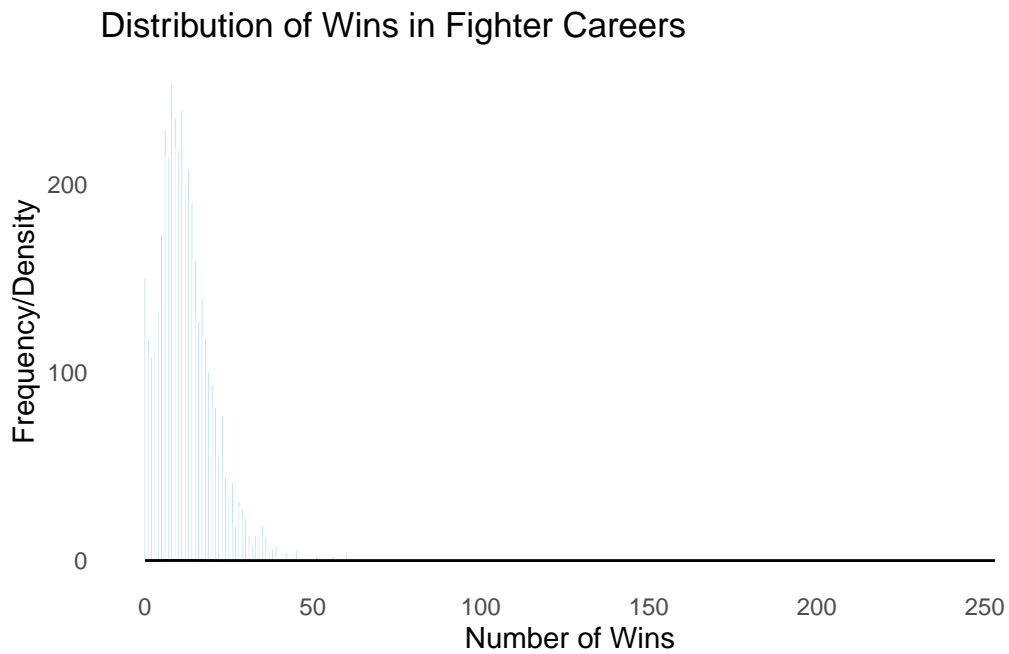
<https://www.euronews.com/business/2023/09/27/the-booming-billion-dollar-business-of-combat-sports>

EDA

```
library(ggplot2)

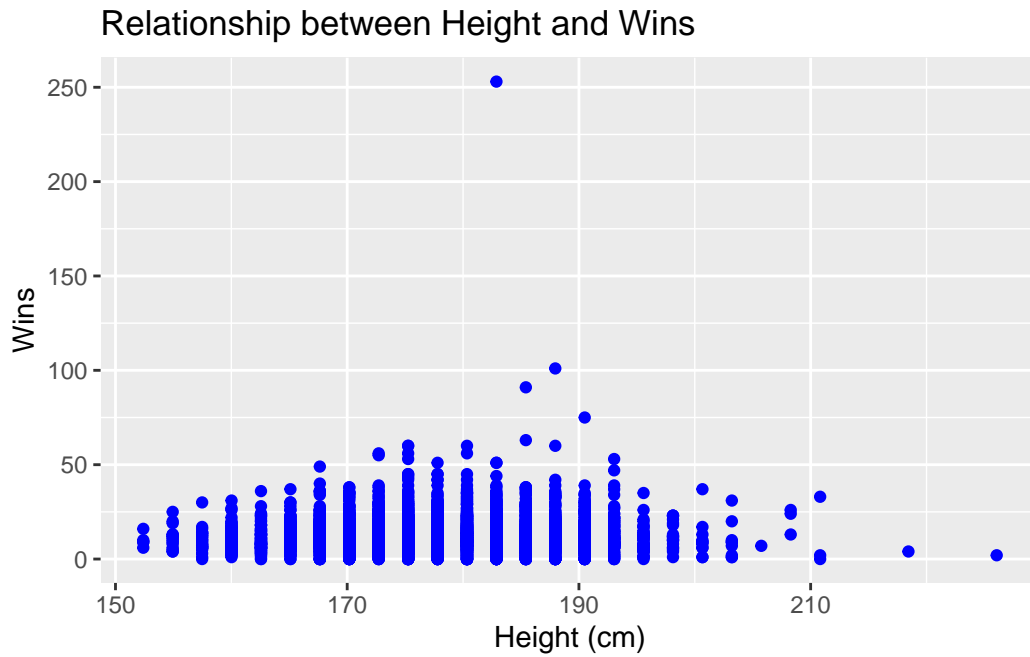
plot <- ggplot(ufc, aes(x = wins)) +
  geom_histogram(binwidth = 1, fill = "skyblue", color = "white") +
  geom_density(alpha = 0.5, fill = "skyblue") +
  labs(x = "Number of Wins", y = "Frequency/Density",
       title = "Distribution of Wins in Fighter Careers") +
  theme_minimal() +
  theme(panel.grid = element_blank())
```

```
print(plot)
```



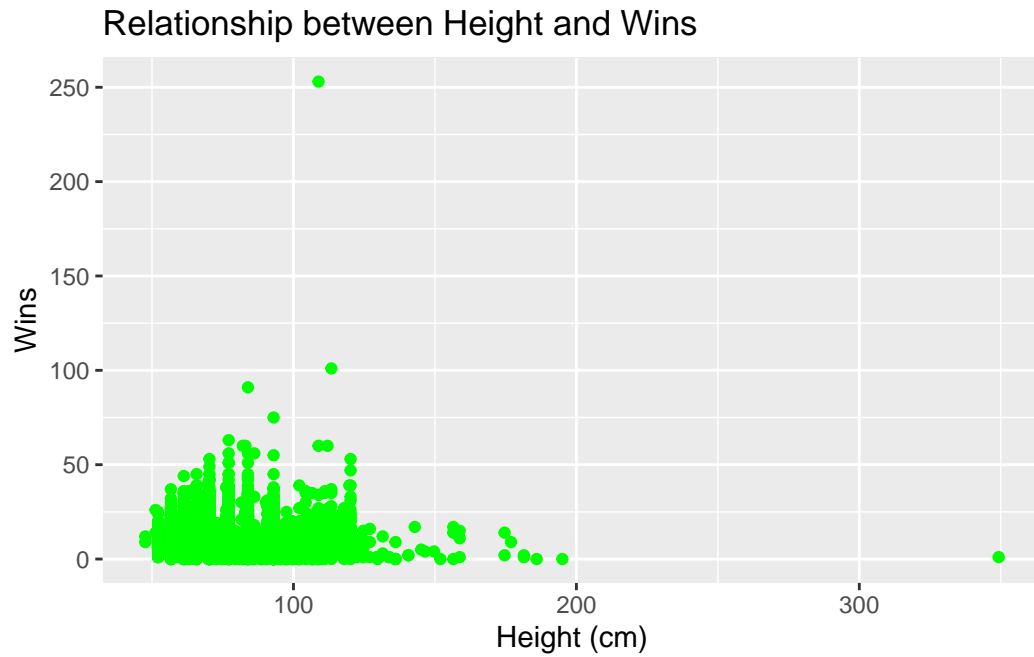
```
plot_height <- ggplot(ufc, aes(x = height_cm, y = wins)) +  
  geom_point(color = "blue") +  
  labs(x = "Height (cm)", y = "Wins",  
        title = "Relationship between Height and Wins")  
  
print(plot_height)
```

Warning: Removed 298 rows containing missing values (`geom_point()`).



```
plot_weight <- ggplot(ufc, aes(x = weight_in_kg, y = wins)) +  
  geom_point(color = "green") +  
  labs(x = "Height (cm)", y = "Wins",  
       title = "Relationship between Height and Wins")  
  
print(plot_weight)
```

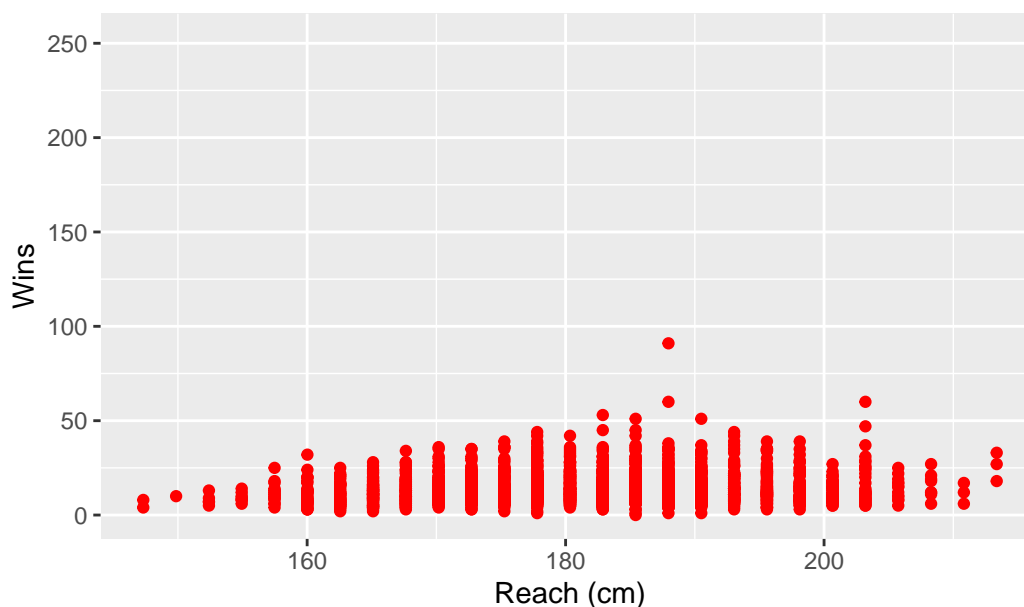
Warning: Removed 87 rows containing missing values (`geom_point()`).



```
plot_reach <- ggplot(ufc, aes(x = reach_in_cm, y = wins)) +  
  geom_point(color = "red") +  
  labs(x = "Reach (cm)", y = "Wins",  
       title = "Relationship between Reach and Wins")  
print(plot_reach)
```

Warning: Removed 1927 rows containing missing values (`geom_point()`).

Relationship between Reach and Wins



Data Preparation

```
ufc$win_ratio <- ufc$wins / (ufc$wins + ufc$losses + ufc$draws)
ufc$age <- as.integer(format(Sys.Date(), "%Y")) - as.integer(format(ufc$date_of_birth, "%Y"))

print(head(ufc))
```

```
# A tibble: 6 x 20
  name      nickname  wins losses draws height_cm weight_in_kg reach_in_cm stance
<chr>    <chr>    <dbl> <dbl> <dbl>   <dbl>      <dbl>    <dbl> <chr>
1 Robert ~ <NA>         7     0     0    190.        93.0         NA Ortho~
2 Daniel ~ The Ani~    15    37     0    185.        83.9         NA <NA>
3 Dan Mol~ <NA>        13     9     0    178.        98.0         NA <NA>
4 Paul Ru~ <NA>         7     4     0    168.        61.2         NA <NA>
5 Collin ~ All In     8     2     0    190.        83.9        193. Ortho~
6 Gerald ~ The Fin~     9     7     0    175.        70.3         NA Ortho~
# i 11 more variables: date_of_birth <date>,
#   significant_strikes_landed_per_minute <dbl>,
#   significant_striking_accuracy <dbl>,
#   significant_strikes_absorbed_per_minute <dbl>,
#   significant_strike_defence <dbl>,
```

```
# average_takedowns_landed_per_15_minutes <dbl>, takedown_accuracy <dbl>,
# takedown_defense <dbl>, ...
```

```
skim(ufc)
```

Table 1: Data summary

Name	ufc
Number of rows	4111
Number of columns	20
Column type frequency:	
character	3
Date	1
numeric	16
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
name	0	1.00	5	27	0	4105	0
nickname	1854	0.55	1	30	0	1784	0
stance	823	0.80	6	11	0	5	0

Variable type: Date

skim_variable	n_missing	complete_rate	min	max	median	n_unique
date_of_birth	1135	0.72	1943-01-25	2004-10-08	1986-11-06	2565

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
wins	0	1.00	12.37	9.37	0.00	7.00	11.00	17.00	253.00	
losses	0	1.00	5.73	5.10	0.00	2.00	5.00	8.00	83.00	
draws	0	1.00	0.26	0.82	0.00	0.00	0.00	0.00	11.00	
height_cm	298	0.93	178.23	8.89	152.40	172.72	177.80	185.42	226.06	
weight_in_kg	87	0.98	77.40	17.98	47.63	65.77	77.11	83.91	349.27	

skim_variable	n_missing	complete_data	mean	sd	p0	p25	p50	p75	p100	hist
reach_in_cm	1927	0.53	181.81	110.68	147.32	175.26	182.88	190.50	213.36	
significant_strikes_landed_per_minute	0	1.00	2.44	1.99	0.00	0.83	2.33	3.60	17.65	
significant_striking_accuracy	0	1.00	35.54	20.40	0.00	27.00	40.00	49.00	100.00	
significant_strikes_absorbed_per_minute	0	1.00	3.15	2.85	0.00	1.55	2.94	4.23	52.50	
significant_strike_defence	0	1.00	42.64	22.32	0.00	36.00	50.00	58.00	100.00	
average_takedowns_landed_per_15_minutes	0	1.00	1.25	1.94	0.00	0.00	0.59	1.94	32.14	
takedown_accuracy	0	1.00	26.30	28.70	0.00	0.00	22.00	45.00	100.00	
takedown_defense	0	1.00	38.96	34.43	0.00	0.00	42.00	66.00	100.00	
average_submissions_attempted_per_15_minutes	0	1.00	0.66	1.51	0.00	0.00	0.00	0.70	21.90	
win_ratio	19	1.00	0.66	0.19	0.00	0.60	0.69	0.78	1.00	
age	1135	0.72	38.59	7.78	20.00	33.00	38.00	44.00	81.00	

Choose Predictors

```

selected_columns <- ufc %>%
  dplyr::select(age, height_cm, weight_in_kg, reach_in_cm, significant_strikes_landed_per_minute,
    significant_striking_accuracy, significant_strikes_absorbed_per_minute,
    significant_strike_defence, average_takedowns_landed_per_15_minutes,
    takedown_accuracy, takedown_defense, average_submissions_attempted_per_15_minutes)

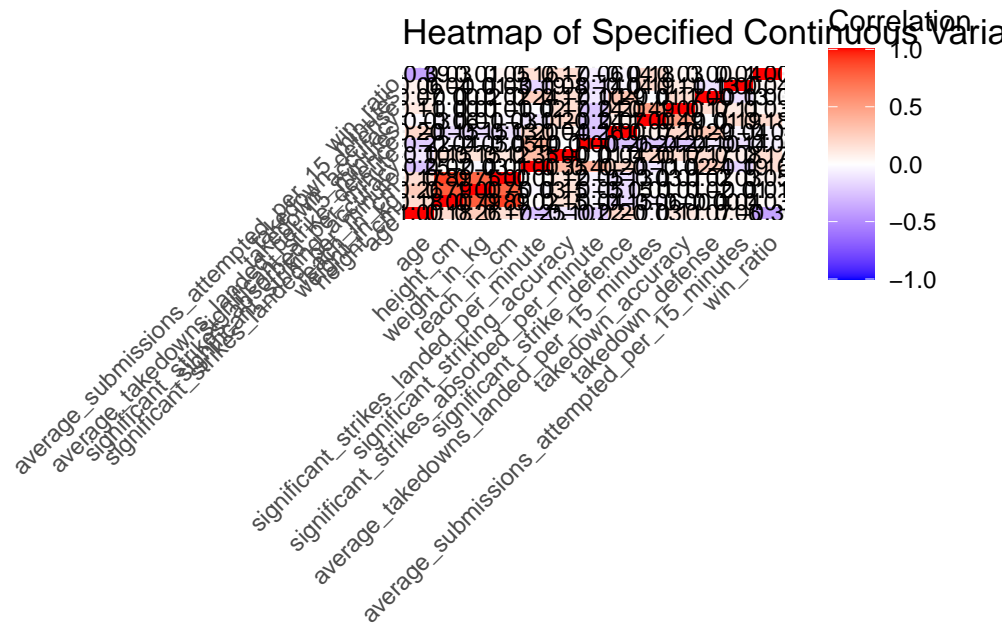
# Compute the correlation matrix
correlation_matrix <- cor(selected_columns, use = "complete.obs") # Handles NA by excluding NA

# Convert the correlation matrix to a long format for ggplot
correlation_data <- as.data.frame(as.table(correlation_matrix))

# Rename columns for clarity
names(correlation_data) <- c("variable1", "variable2", "value")

# Plotting the heatmap with correlation coefficients
ggplot(correlation_data, aes(x = variable1, y = variable2, fill = value)) +
  geom_tile() + # This creates the heatmap tiles
  geom_text(aes(label = sprintf("%.2f", value)), color = "black", size = 3) + # Adds text
  scale_fill_gradient2(low = "blue", high = "red", mid = "white", midpoint = 0, limit = c(-1, 1)) +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
    axis.text.y = element_text(angle = 45, vjust = 1)) + # Adjust text alignment if needed
  labs(title = "Heatmap of Specified Continuous Variables", x = "", y = "")

```

Predictors: age Stance (category) average_takedowns_landed_per_15 minutes Significant_striking_accuracy Significant_strikes_landed_per_minute significant_strikes_absorbed_per_minute significant_strike_defence

Model Fitting

```
#linear
ufc$stance <- as.factor(ufc$stance)

model1 <- lm(win_ratio ~ age + stance + average_takedowns_landed_per_15_minutes + significant_striking_accuracy + significant_strikes_landed_per_minute + significant_strikes_absorbed_per_minute + significant_strike_defence,
             data = ufc)

summary(model1)
```

Call:

```
lm(formula = win_ratio ~ age + stance + average_takedowns_landed_per_15_minutes + significant_striking_accuracy + significant_strikes_landed_per_minute + significant_strikes_absorbed_per_minute + significant_strike_defence,
```

```

data = ufc)

Residuals:
    Min       1Q   Median       3Q      Max
-0.70326 -0.06308  0.00335  0.07501  0.58006

Coefficients:
                Estimate Std. Error t value Pr(>|t|)
(Intercept)      0.9418870   0.0539846   17.447 < 2e-16
age             -0.0082692   0.0003249  -25.453 < 2e-16
stanceOrthodox    0.0179482   0.0502243    0.357  0.721
stanceSideways    0.1840189   0.1322619    1.391  0.164
stanceSouthpaw    0.0236782   0.0504429    0.469  0.639
stanceSwitch      0.0238210   0.0510905    0.466  0.641
average_takedowns_landed_per_15_minutes  0.0095566   0.0013340    7.164 9.86e-13
significant_striking_accuracy  0.0009264   0.0001842    5.029 5.22e-07
significant_strikes_landed_per_minute  0.0063530   0.0016135    3.937 8.42e-05
significant_strikes_absorbed_per_minute -0.0051719   0.0009167   -5.642 1.84e-08
significant_strike_defence  0.0001415   0.0001732    0.817  0.414

(Intercept)          ***
age                  ***
stanceOrthodox
stanceSideways
stanceSouthpaw
stanceSwitch
average_takedowns_landed_per_15_minutes ***
significant_striking_accuracy          ***
significant_strikes_landed_per_minute  ***
significant_strikes_absorbed_per_minute ***
significant_strike_defence
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1224 on 2965 degrees of freedom
(1135 observations deleted due to missingness)
Multiple R-squared:  0.2876,    Adjusted R-squared:  0.2852
F-statistic: 119.7 on 10 and 2965 DF,  p-value: < 2.2e-16

```

Categorical variable:

0.0435211 represents the difference in average win ratio for fighters with Orthodox stance, compared with fighters with Open stance, while holding the other predictors constant.

-0.0893852 represents the difference in average win ratio for fighters with Sideways stance, compared with fighters with Open stance, while holding the other predictors constant.

0.0484787 represents the difference in average win ratio for fighters with Southpaw stance, compared with fighters with Open stance, while holding the other predictors constant.

0.0895293 represents the difference in average win ratio for fighters with Switch stance, compared with fighters with Open stance, while holding the other predictors constant.

We notice that the p-value for Stance variables are all less than 0.05, and thus, there is insufficient evidence to suggest a linear relationship between fighter's stance and wins ratio at 0.05 significance level, while controlling for other predictors

Continuous variable:

While holding the other predictors constant, when the average number of takedowns landed by the fighter per 15 minutes increases by 1 time, the average win ratio of the fighter will increase by 0.0113384.

While holding the other predictors constant, when the significant striking accuracy increases by 1%, the average win ratio of the fighter will increase by 0.0017843.

While holding the other predictors constant, when the average number of significant strikes landed by the fighter per minute increases by 1 time, the average win ratio of the fighter will increase by 0.0199915.

While holding the other predictors constant, when the average number of significant strikes absorbed by the fighter per minute increases by 1 time, the average win ratio of the fighter will decrease by 0.0016777.

While holding the other predictors constant, when the opponent's significant strikes that the fighter successfully defends increases by 1%, the average win ratio of the fighter will increase by 0.0009887.

```
#linear & transformation
ufc <- ufc %>%
  filter(significant_strike_defence > 0)

model2 <- lm(win_ratio ~ age + stance + average_takedowns_landed_per_15_minutes + significant_strikes_landed_per_minute + significant_strikes_absorbed_per_minute + log(significant_strike_defence),
  data = ufc)

summary(model2)
```

Call:

```
lm(formula = win_ratio ~ age + stance + average_takedowns_landed_per_15_minutes +
    significant_striking_accuracy + significant_strikes_landed_per_minute +
    significant_strikes_absorbed_per_minute + log(significant_strike_defence),
    data = ufc)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.62694	-0.06236	0.00217	0.07236	0.43321

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.8201894	0.0682592	12.016	< 2e-16
age	-0.0077436	0.0003230	-23.971	< 2e-16
stanceOrthodox	0.0071090	0.0516415	0.138	0.890519
stanceSideways	0.1736378	0.1260981	1.377	0.168618
stanceSouthpaw	0.0096258	0.0518276	0.186	0.852671
stanceSwitch	0.0050580	0.0524267	0.096	0.923148
average_takedowns_landed_per_15_minutes	0.0095747	0.0012638	7.576	4.76e-14
significant_striking_accuracy	0.0013110	0.0002060	6.365	2.27e-10
significant_strikes_landed_per_minute	0.0060730	0.0016010	3.793	0.000152
significant_strikes_absorbed_per_minute	-0.0036717	0.0009872	-3.719	0.000204
log(significant_strike_defence)	0.0251122	0.0098567	2.548	0.010894

(Intercept)	***
age	***
stanceOrthodox	
stanceSideways	
stanceSouthpaw	
stanceSwitch	
average_takedowns_landed_per_15_minutes	***
significant_striking_accuracy	***
significant_strikes_landed_per_minute	***
significant_strikes_absorbed_per_minute	***
log(significant_strike_defence)	*

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.115 on 2870 degrees of freedom

(521 observations deleted due to missingness)

Multiple R-squared: 0.2789, Adjusted R-squared: 0.2763

F-statistic: 111 on 10 and 2870 DF, p-value: < 2.2e-16

```
#logistic
ufc$win_ratio_binary <- ifelse(ufc$win_ratio > 0.5, 1, 0)

model3 <- glm(win_ratio_binary ~ age + stance + average_takedowns_landed_per_15_minutes +
              significant_striking_accuracy + significant_strikes_landed_per_minute +
              significant_strikes_absorbed_per_minute + significant_strike_defence,
              family = binomial, data = ufc)

summary(model3)
```

Call:

```
glm(formula = win_ratio_binary ~ age + stance + average_takedowns_landed_per_15_minutes +
    significant_striking_accuracy + significant_strikes_landed_per_minute +
    significant_strikes_absorbed_per_minute + significant_strike_defence,
    family = binomial, data = ufc)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	7.430179	1.474741	5.038	4.7e-07
age	-0.149109	0.013137	-11.350	< 2e-16
stanceOrthodox	-0.126691	1.200785	-0.106	0.915974
stanceSideways	13.894061	535.412518	0.026	0.979297
stanceSouthpaw	-0.317526	1.211986	-0.262	0.793331
stanceSwitch	0.037334	1.318228	0.028	0.977406
average_takedowns_landed_per_15_minutes	0.215908	0.072626	2.973	0.002950
significant_striking_accuracy	0.024920	0.007205	3.459	0.000542
significant_strikes_landed_per_minute	0.093100	0.077137	1.207	0.227452
significant_strikes_absorbed_per_minute	-0.037184	0.026051	-1.427	0.153476
significant_strike_defence	0.008041	0.007835	1.026	0.304745

(Intercept)	***
age	***
stanceOrthodox	
stanceSideways	
stanceSouthpaw	
stanceSwitch	
average_takedowns_landed_per_15_minutes	**
significant_striking_accuracy	***
significant_strikes_landed_per_minute	
significant_strikes_absorbed_per_minute	
significant_strike_defence	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1303.0 on 2880 degrees of freedom
Residual deviance: 1030.8 on 2870 degrees of freedom
(521 observations deleted due to missingness)
AIC: 1052.8

Number of Fisher Scoring iterations: 12

We impose a threshold of 0.5 for wins ratio: when it's larger than 0.5, we consider it Satisfactory Performance, and when it's less than 0.5, we consider it Unsatisfactory Performance.

```
#ordinal classified by age

ufc$win_ratio_group <- cut(ufc$win_ratio,
                           breaks = c(0, 0.33, 0.66, 1),
                           labels = c("Low", "Medium", "High"),
                           include.lowest = TRUE)
ufc$win_ratio_group <- ordered(ufc$win_ratio_group)

model4 <- polr(win_ratio_group ~ age + stance + average_takedowns_landed_per_15_minutes +
               significant_striking_accuracy + significant_strikes_landed_per_minute +
               significant_strikes_absorbed_per_minute + significant_strike_defence,
               data = ufc)
summary(model4)
```

Re-fitting to get Hessian

Call:

```
polr(formula = win_ratio_group ~ age + stance + average_takedowns_landed_per_15_minutes +
      significant_striking_accuracy + significant_strikes_landed_per_minute +
      significant_strikes_absorbed_per_minute + significant_strike_defence,
      data = ufc)
```

Coefficients:

	Value	Std. Error	t value
age	-0.114646	0.0070482	-1.627e+01

stanceOrthodox	-0.302536	0.9022047	-3.353e-01
stanceSideways	11.923761	0.0001064	1.121e+05
stanceSouthpaw	-0.027337	0.9069988	-3.014e-02
stanceSwitch	-0.564248	0.9227359	-6.115e-01
average_takedowns_landed_per_15_minutes	0.153758	0.0320932	4.791e+00
significant_striking_accuracy	0.016123	0.0042248	3.816e+00
significant_strikes_landed_per_minute	0.056107	0.0359596	1.560e+00
significant_strikes_absorbed_per_minute	-0.036733	0.0205325	-1.789e+00
significant_strike_defence	0.005212	0.0046364	1.124e+00

Intercepts:

	Value	Std. Error	t value
Low Medium	-8.2170	1.0347	-7.9415
Medium High	-4.3951	1.0181	-4.3167

Residual Deviance: 3383.979

AIC: 3407.979

(521 observations deleted due to missingness)