



DAVID WISMER

# UFC SCORE PREDICTIONS

LINEAR REGRESSION

# INTRODUCTION

---

- ▶ Project Goal
  - ▶ Create a multivariate linear regression model to determine the winner of a UFC fight using fight metrics (striking statistics, grappling statistics, control time, etc).
- ▶ Ultimate Fighting Championship UFC Judging
  - ▶ There are 3 to 5 rounds per fight.
  - ▶ For each round, each judge gives 10 points to the winner. The loser receives 9 points, or in some cases, 8 points. Both fighters get 10 points in the event of a tie.
  - ▶ Judges score fights based on a hierarchy of criteria: 1) effective striking and grappling, 2) effective aggression, and 3) control of the fight area.

# MODEL DATA

---

## Dataset

2,258 fights and 4,516 records. Each record represents a fighter\_id/fight\_id combination

## Features

Individual statistics were converted to disparity between fighter\_1 and fighter\_2

## Train Test Split

60% training, 20% validation, and 20% testing data using a time-based split.

Significant Strikes Landed, Significant Strikes Attempted, Takedowns Landed, Takedowns Attempted, Control Time, Leg Strikes, Body Strikes, Head Strikes, Submission Attempts, Total Strikes Landed, Total Strikes Attempted

# DATA SOURCES AND TOOLS

Data Sources:



Web Scraping:



Data Storage:



Database Interactions:



Data Cleaning:



Regression Model:



Visualization:



# WEB SCRAPING - UFCSTATS.COM

1

### Events & Fights

Enter Event Name...

Completed      Upcoming

NAME/DATE	LOCATION
NEXT → UFC Fight Night: Rozenstruik vs. Sakai June 05, 2021	Las Vegas, Nevada, USA
<b>UFC Fight Night: Font vs. Garbrandt</b> May 22, 2021	Las Vegas, Nevada, USA
UFC 262: Oliveira vs. Chandler May 15, 2021	Houston, Texas, USA
UFC Fight Night: Rodriguez vs. Waterson May 08, 2021	Las Vegas, Nevada, USA

3

### UFC Fight Night: Font vs. Garbrandt

L Felicia Spencer "FEENOM"  
W Norma Dumont "THE IMMORTAL"

#### WOMEN'S FEATHERWEIGHT BOUT

METHOD: Decision - Split ROUND: 3 TIME: 5:00 TIME FORMAT: 3 Rnd (5-5-5) REFEREE: Chris Tognoni  
DETAILS: Sal D'Amato 27 - 30. Brian Miner 29 - 28. Junichiro Kamijo 28 - 29.

2

### UFC Fight Night: Font vs. Garbrandt

DATE: May 22, 2021 LOCATION: Las Vegas, Nevada, USA

Click on a row below to see in-depth event stats. Fight, Perf, Sub, and KO of the Night Bonuses: **\*FIGHT \*PERF \*SUB \*KO**

W/L	FIGHTER	KD	STR	TD	SUB	WEIGHT CLASS	METHOD	ROUND	TIME
WIN	Rob Font	0	176	2	0	Bantamweight	U-DEC	5	5:00
WIN	Cody Garbrandt	0	63	3	0				
WIN	Carla Esparza	0	27	3	0	Women's Strawweight	KO/TKO		
WIN	Yan Xiaonan	0	5	0	0		Punches	2	2:58
WIN	Jared Vaudreuil	0	121	0	0	Heavyweight	U-DEC	3	5:00
WIN	Justin Tafa	0	74	0	0				
WIN	Norma Dumont	0	68	1	0	Women's Featherweight	S-DEC	3	5:00
WIN	Felicia Spencer	0	47	0	0				

### SIGNIFICANT STRIKES

FIGHTER	SIG. STR	SIG. STR. %	HEAD	BODY	LEG	DISTANCE	CLINCH	GROUND
Felicia Spencer	47 of 114	41%	27 of 76	14 of 27	6 of 11	38 of 102	9 of 12	0 of 0
Norma Dumont	68 of 135	50%	40 of 100	14 of 20	14 of 15	59 of 123	9 of 12	0 of 0
PER ROUND ▾								
ROUND 1								
Felicia Spencer	11 of 37	29%	7 of 25	2 of 8	2 of 4	8 of 34	3 of 3	0 of 0
Norma Dumont	21 of 46	45%	11 of 35	3 of 3	7 of 8	18 of 43	3 of 3	0 of 0
ROUND 2								
Felicia Spencer	20 of 51	39%	12 of 35	5 of 10	3 of 6	19 of 49	1 of 2	0 of 0
Norma Dumont	35 of 60	58%	23 of 45	6 of 9	6 of 6	33 of 57	2 of 3	0 of 0
ROUND 3								
Felicia Spencer	16 of 26	61%	8 of 16	7 of 9	1 of 1	11 of 19	5 of 7	0 of 0
Norma Dumont	12 of 29	41%	6 of 20	5 of 8	1 of 1	8 of 23	4 of 6	0 of 0

# WEB SCRAPING - MMADECISIONS.COM

1

EVENTS													
2021	2020	2019	2018	2017	2016	2015	2014	2013	2012	2011	2010	2009	2008
2007	2006	2005	2004	2003	2002	2001	2000	1999	1998	1997	1996	1995	1994
<b>Date</b>												#	
May 22, 2021	UFC on ESPN+ 46: Font vs. Garbrandt											8	
May 21, 2021	Bellator 259: Cyborg vs. Smith 2											6	
May 21, 2021	Invicta FC: Rodriguez vs. Torquato											4	
May 15, 2021	UFC 262: Oliveira vs. Chandler											5	
May 08, 2021	UFC on ESPN 24: Rodriguez vs. Waterson											6	
May 07, 2021	Bellator 258: Archuleta vs. Pettis											5	
May 01, 2021	UFC on ESPN 23: Reyes vs. Prochazka											8	
Apr 24, 2021	KSW 60: De Fries vs. Narkun 2											3	
Apr 24, 2021	UFC 261: Usman vs. Masvidal 2											4	
Apr 17, 2021	UFC on ESPN 22: Whittaker vs. Gastelum											8	

2

UFC on ESPN+ 46: Font vs. Garbrandt			
UFC Apex			
Las Vegas, Nevada, USA			
May 22, 2021			
<b>Fight</b>		<b>Scores</b>	
Font def. Garbrandt Unanimous	CLEARY 48 - 47	D'AMATO 50 - 45	KAMIJO 50 - 45
Vanderaa def. Tafa Unanimous	BELL 30 - 27	COLÓN 29 - 28	WEEKS 30 - 27
Dumont def. Spencer Split	D'AMATO 30 - 27	KAMIJO 29 - 28	MINER 28 - 29
Ramos def. Algeo Unanimous	BELL 30 - 27	CLEARY 29 - 28	COLÓN 30 - 27
Hermansson def. Shahbazyan Unanimous	BYRD 29 - 27	D'AMATO 29 - 27	KAMIJO 29 - 27
McGee def. Silva Unanimous	BELL 29 - 27	BYRD 30 - 26	CLEARY 30 - 26
Culibao def. Nuerdanbieke Unanimous	CLEARY 29 - 28	HAGEN 29 - 28	MINER 29 - 28
Ismagulov def. Alves Unanimous	BELL 29 - 28	HAGEN 29 - 28	KAMIJO 29 - 28

3

**Norma Dumont**  
defeats  
**Felicia Spencer**

**SPLIT DECISION**

**UFC on ESPN+ 46: Font vs. Garbrandt**  
May 22, 2021  
Las Vegas, Nevada, USA  
REFEREE: Unknown

**Sal D'Amato**

ROUND	Dumont	Spencer
1	10	9
2	10	9
3	10	9
<b>TOTAL</b>	<b>30</b>	<b>27</b>

**Junichiro Kamijo**

ROUND	Dumont	Spencer
1	-	-
2	-	-
3	-	-
<b>TOTAL</b>	<b>29</b>	<b>28</b>

**Bryan Miner**

ROUND	Dumont	Spencer
1	-	-
2	-	-
3	-	-
<b>TOTAL</b>	<b>28</b>	<b>29</b>

**MEDIA SCORES**

Dayne Fox <i>BloodyElbow.com</i>	<b>29-28</b>	Dumont
Matthew Wells <i>TheBodyLockMMA.com</i>	<b>29-28</b>	Dumont
Ryan Frederick <i>WrestlingObserver.com</i>	<b>29-28</b>	Dumont
Drake Riggs <i>TheBodyLockMMA.com</i>	<b>29-28</b>	Dumont
Rob Tatum <i>CombatPress.com</i>	<b>29-28</b>	Dumont
Marcel Dorff <i>mmdna.nl</i>	<b>29-28</b>	Dumont
Shawn Bitter <i>Cageside Press</i>	<b>29-28</b>	Dumont
MMAJunkie.com	<b>29-28</b>	Dumont
MMANmania.com	<b>29-28</b>	Dumont
Jay Petry <i>Sherdog.com</i>	<b>28-29</b>	Spencer
Ben Duffy <i>Sherdog.com</i>	<b>28-29</b>	Spencer
Tyler Treese <i>Sherdog.com</i>	<b>28-29</b>	Spencer

**YOUR SCORECARD**

ROUND	Dumont	Spencer
1	- ▼	- ▼
2	- ▼	- ▼
3	- ▼	- ▼
<b>TOTAL</b>	-	-

[Show Results](#)

# CROSS VALIDATION - R-SQUARED

	Linear Regression	LASSO Regression	Ridge Regression	SGD Regression	Random Forest
Unaltered Data	0.6832	0.6832	0.6832	N/A	0.6955
MinMaxScaler (Normalized)	0.6832	0.6831	0.6832	0.6661	0.6954
StandardScaler (Standardized)	0.6832	0.6832	0.6832	0.6832	<b>0.6958</b>

- ▶ The best R-Squared in the models tested was approximately **0.6958**.
- ▶ Random Forest had the highest R-Squared, but also had the largest generalization error from training to validation. Hyperparameter tuning only partially fixed overfitting.
- ▶ Normalizing with MinMaxScaler and standardizing with StandardizedScaler had minimal impact on performance.
- ▶ SGDRegressor performed differently each iteration. Unaltered data produced nonsensical results.

# CROSS VALIDATION – RMSE

	Linear Regression	LASSO Regression	Ridge Regression	SGD Regression	Random Forest
Unaltered Data	0.2025	0.2025	0.2025	N/A	<b>0.1984</b>
MinMaxScaler (Normalized)	0.2025	0.2025	0.2025	0.2075	0.1987
StandardScaler (Standardized)	0.2025	0.2025	0.2025	0.2026	0.1986

- ▶ The minimum RSME in the models tested was approximately **0.1984**.
- ▶ The reasonable actual per round score range is 8.0 - 10.0. The RMSE results range from 9.4% - 10.4% of this range.
- ▶ Normalizing with MinMaxScaler and standardizing with StandardizedScaler had minimal impact on performance.
- ▶ SGDRegressor performed differently each iteration. Unaltered data produced nonsensical results.

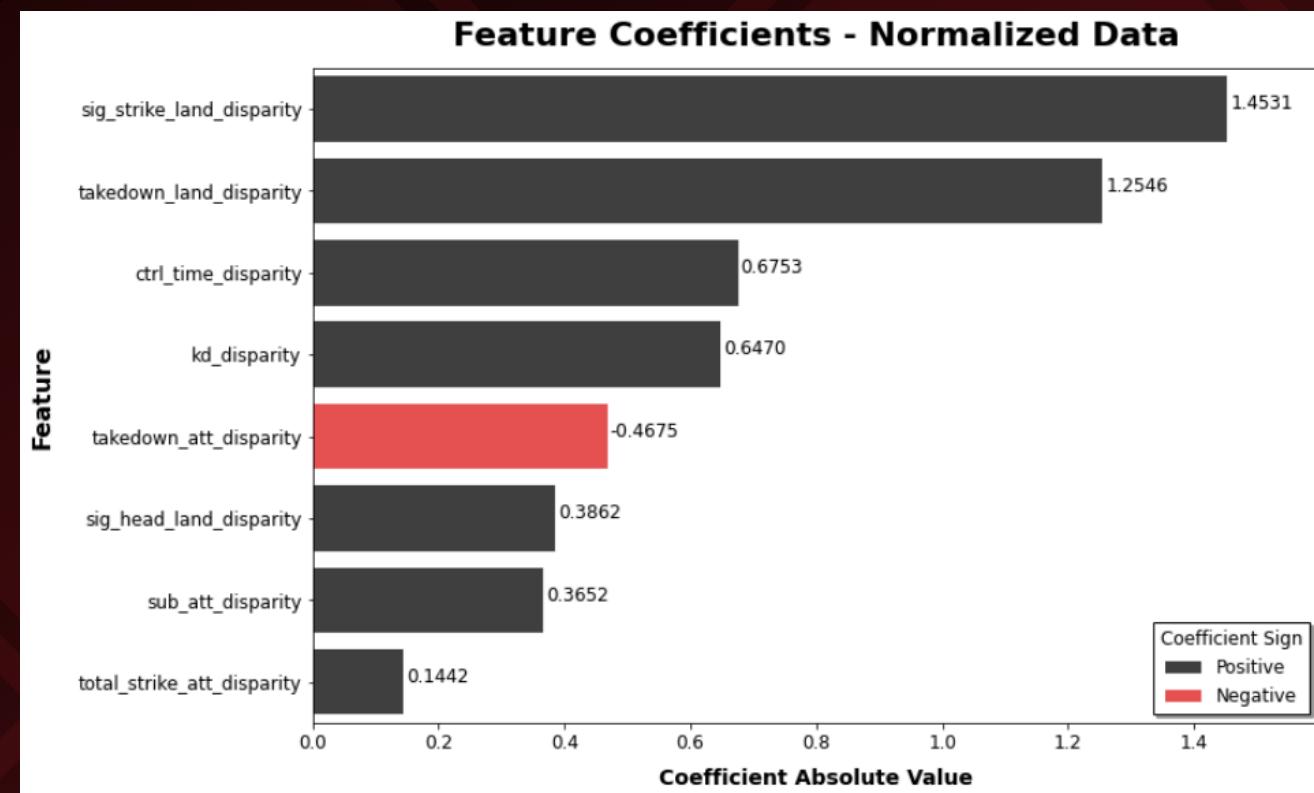
# TEST RESULTS - TOP 3 MODELS

	Linear Regression Normalized	SGD Regression Normalized	Random Forest Regression Standardized
R-Squared	0.7070	0.7085	<b>0.7152</b>
RMSE	0.2078	0.2072	<b>0.2053</b>
Prediction Accuracy	85.5%	<b>86.4%</b>	84.4%

- ▶ My preferred regression model is **Linear Regression with Normalized Data**
- ▶ Linear Regression is easily interpretable and normalized features allow for a more natural comparison of feature coefficients
- ▶ All models performed very similarly in cross validation, simple validation, and in testing
- ▶ Random Forest Regression had strong cross validation results, but even with hyperparameter tuning, the model was overfit to the training data.

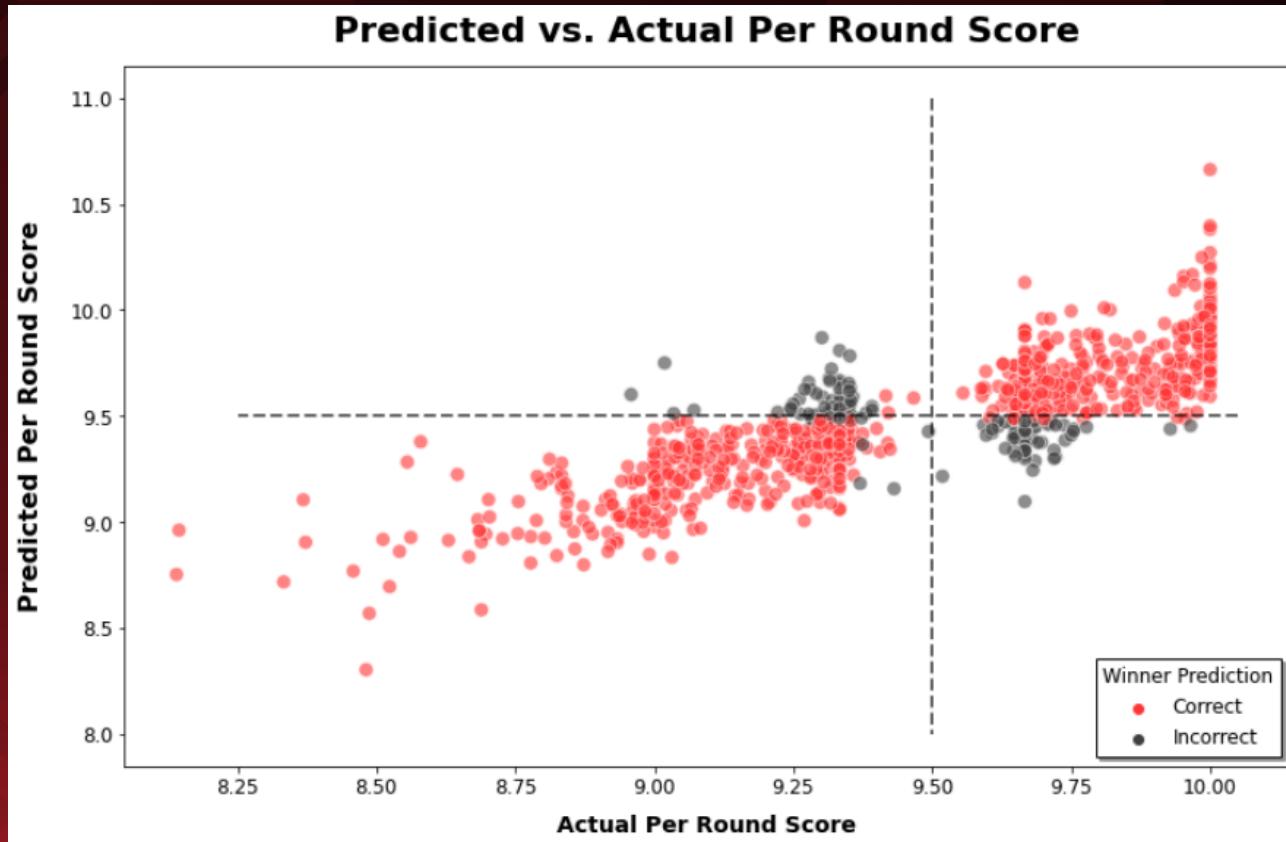
# FEATURE IMPORTANCE - LINEAR REGRESSION, NORMALIZED DATA

Per Round Score =  $1.4531X_1 + 1.2546X_2 + 0.6753X_3 + 0.6470X_4$   
 $- 0.4675X_5 + 0.3865X_6 + 0.3652X_7 + 0.1442X_8 + 7.2543$



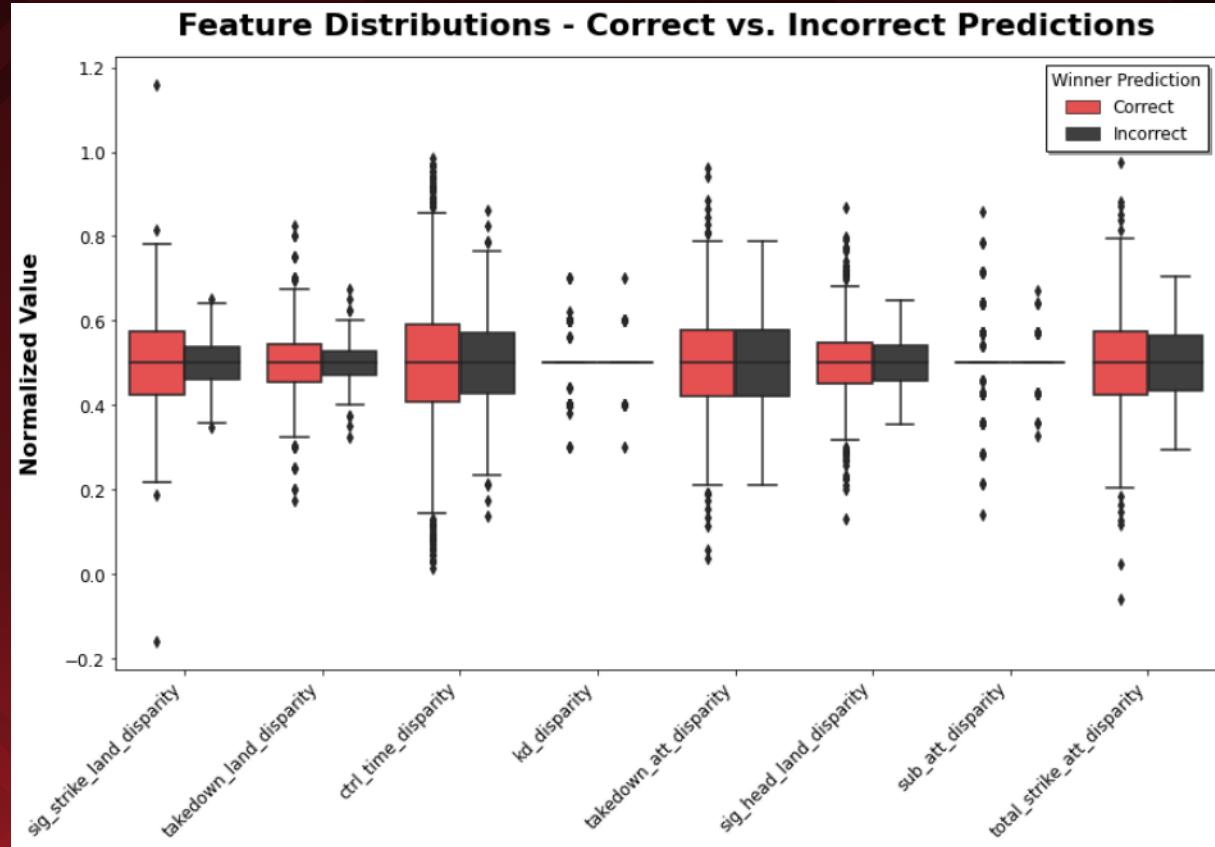
- ▶ The metrics with the greatest impact on predictions are in line with intended UFC judging criteria:
  - ▶ 1) effective striking and grappling
  - ▶ 2) aggression
  - ▶ 3) fight area control
- ▶ The error term for this model is **7.2543** points per round.
- ▶ Takedown attempts have a negative impact on projected score. For every takedown attempt, the fighter is penalized. For each successful takedown, the reward exceeds the penalty.

# WINNER PREDICTION - LINEAR REGRESSION, NORMALIZED DATA



- ▶ Winner Prediction Accuracy: **85.5%**
- ▶ Typically, the model is correct when both predicted and actual per round score fall on the same side of 9.5 points per round.
- ▶ The model predicts the winner more accurately the further the actual score gets from 9.5 points per round.

# WHERE STATISTICS FALL SHORT



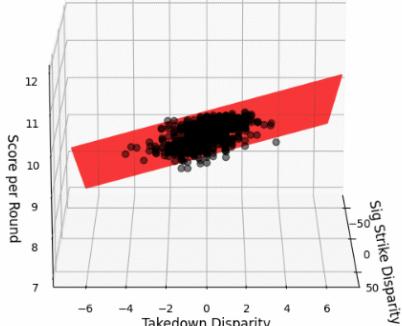
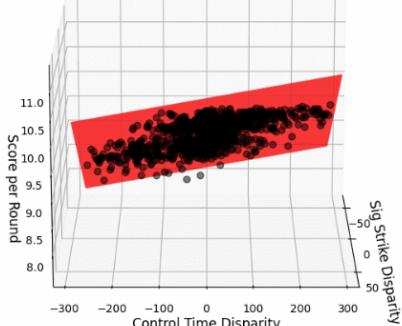
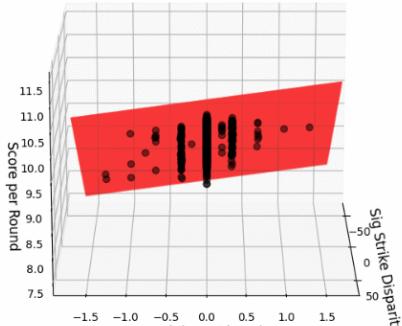
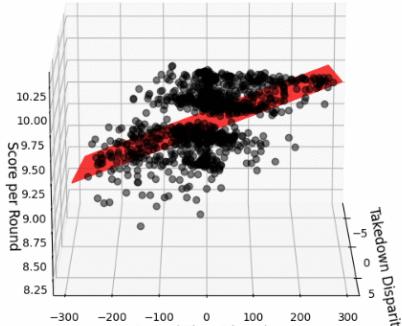
- ▶ Fight metric distributions were narrower in fights where the model made the wrong prediction. Closer fights are harder to judge.
- ▶ Statistics struggle to capture aggression and damage inflicted by fighters, both key factors in judging close fights.
- ▶ Knockdowns are much more rare than, for example, significant strikes. Any knockdown disparity is rewarded heavily, as even the 3rd quartile for knockdown disparity per round is zero. Knockdown disparity has the 4th largest coefficient in the regression model.

## FUTURE WORK

---

- ▶ Model Enhancements
  - ▶ Logistic Regression on a per round basis
- ▶ Analysis of Fighter Styles
  - ▶ Which styles match up most effectively with other styles?
- ▶ Judging the Judges
  - ▶ Which judges historically deviate furthest from the regression model?

# THANK YOU!

<b>Feature 1</b>	Sig Strike Disparity		<b>Feature 1</b>	Sig Strike Disparity
<b>Feature 2</b>	Takedown Disparity		<b>Feature 2</b>	Control Time Disparity
<b>R-Squared</b>	0.6190		<b>R-Squared</b>	0.6624
<b>Feature 1</b>	Sig Strike Disparity		<b>Feature 1</b>	Takedown Disparity
<b>Feature 2</b>	Knockdown Disparity		<b>Feature 2</b>	Control Time Disparity
<b>R-Squared</b>	0.5071		<b>R-Squared</b>	0.3283

# **APPENDICES**

# FEATURE IMPORTANCE COMPARISON

Linear Regression			Random Forest		
Rank	Feature	Coefficient	Rank	Feature	Importance
1	sig_strike_land_disparity	1.4531	1	sig_strike_land_disparity	0.6165
2	takedown_land_disparity	1.2546	2	ctrl_time_disparity	0.2280
3	ctrl_time_disparity	0.6753	3	takedown_land_disparity	0.0515
4	kd_disparity	0.6470	4	sig_head_land_disparity	0.0410
5	takedown_att_disparity	-0.4675	5	total_strike_att_disparity	0.0264
6	sig_head_land_disparity	0.3862	6	takedown_att_disparity	0.0180
7	sub_att_disparity	3652	7	kd_disparity	0.0094
8	total_strike_att_disparity	0.1442	8	sub_att_disparity	0.0094

The diagram illustrates the relationship between the features identified by the two models. Red arrows point from the top four features of the LR model to the top four features of the RF model. A blue arrow points from the fifth feature of the LR model to the fifth feature of the RF model. A green arrow points from the sixth feature of the LR model to the sixth feature of the RF model. A purple arrow points from the seventh feature of the LR model to the seventh feature of the RF model.

# MOST CONTROVERSIAL FIGHTS OF 2020

Fighter 1	Fighter 2	Judges	Media	Judges + Media	Model
Jon Jones	Dominick Reyes	Jones 9.67	Jones 9.60	Jones 9.61	Reyes 9.53
Paul Felder	Dan Hooker	Hooker 9.53	Hooker 9.60	Hooker 9.59	Hooker 9.54
Stipe Miocic	Daniel Cormier	Miocic 9.73	Miocic 9.64	Miocic 9.65	Miocic 9.58
Pedro Munhoz	Frankie Edgar	Munhoz 9.53	Edgar 9.67	Edgar 9.65	Munhoz 9.57
Israel Adesanya	Yoel Romero	Adesanya 9.67	Adesanya 9.60	Adesanya 9.61	Adesanya 9.52
Alexander Volkanovski	Max Holloway	Volkanovski 9.53	Volkanovski 9.60	Volkanovski 9.59	Volkanovski 9.64
Robert Whittaker	Darren Till	Whittaker 9.60	Whittaker 9.63	Whittaker 9.63	Whittaker 9.62
Dan Ige	Edson Barboza	Ige 9.56	Ige 9.67	Ige 9.65	Barboza 9.53
Song Yadong	Marlon Vera	Yadong 9.67	Yadong 9.68	Yadong 9.68	Yadong 9.49
Max Holloway	Alexander Volkanovski	Volkanovski 9.73	Volkanovski 9.74	Volkanovski 9.74	Volkanovski 9.53

Source: [MMAOddsBreaker.com](https://MMAOddsBreaker.com)

# MUNHOZ V. EDGAR – 2020

Fighter 2	Munhoz	Edgar
Significant Strikes Landed Disparity	4.6	-4.6
Takedowns Landed Disparity	0.0	0.0
Control Time Disparity (Seconds)	8.6	-8.6
Knockdown Disparity	0.0	0.0
Takedowns Attempted Disparity	0.8	-0.8
Significant Head Strikes Landed Disparity	-3.6	3.6
Submissions Attempted Disparity	0.0	0.0
Total Strikes Attempted Disparity	-0.2	0.2

## Judges

- ▶ Edgar wins by split decision
- ▶ Total points per round in favor of Munhoz with 9.53 points per round

## Media

- ▶ Edgar wins with 9.67 points per round

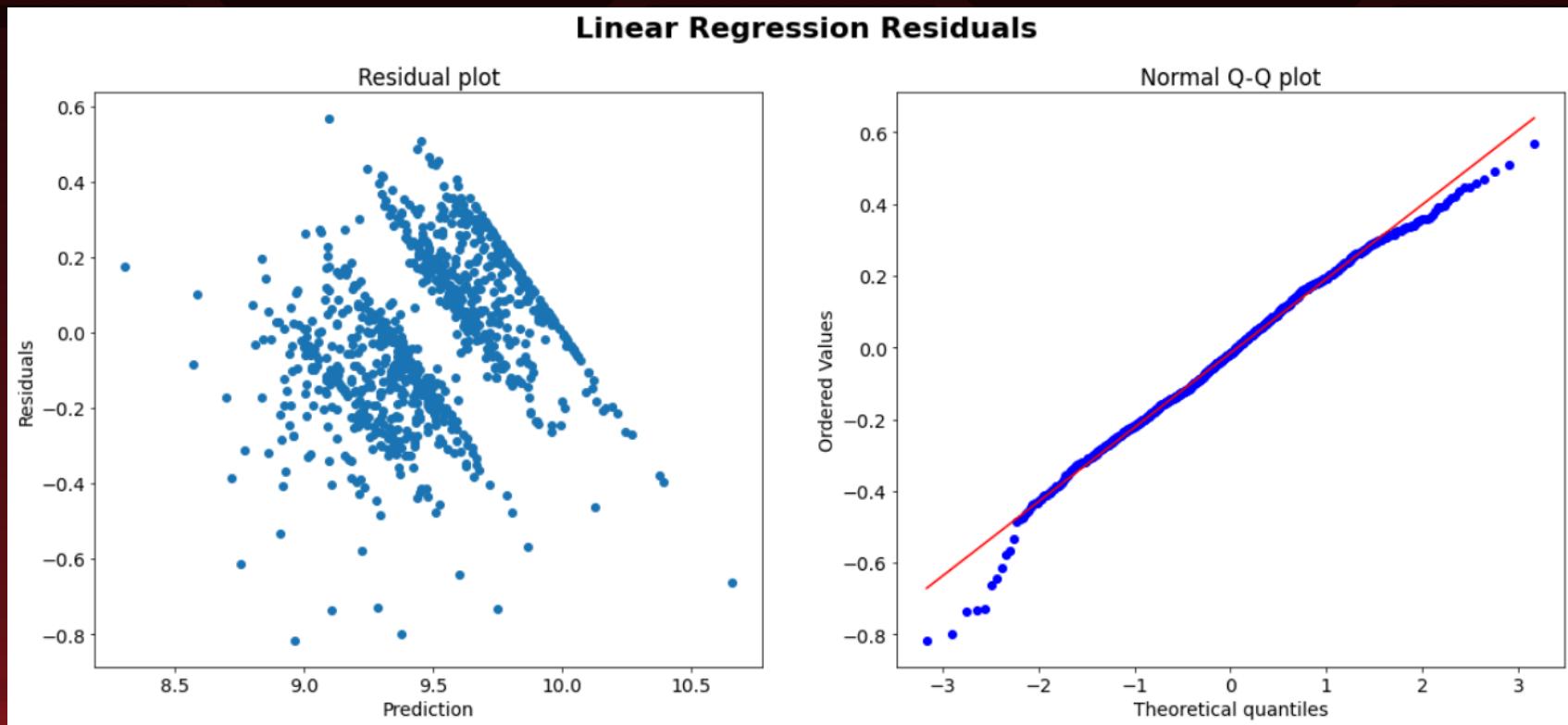
## Judges & Media Combined

- ▶ Edgar wins with 9.65 points per round

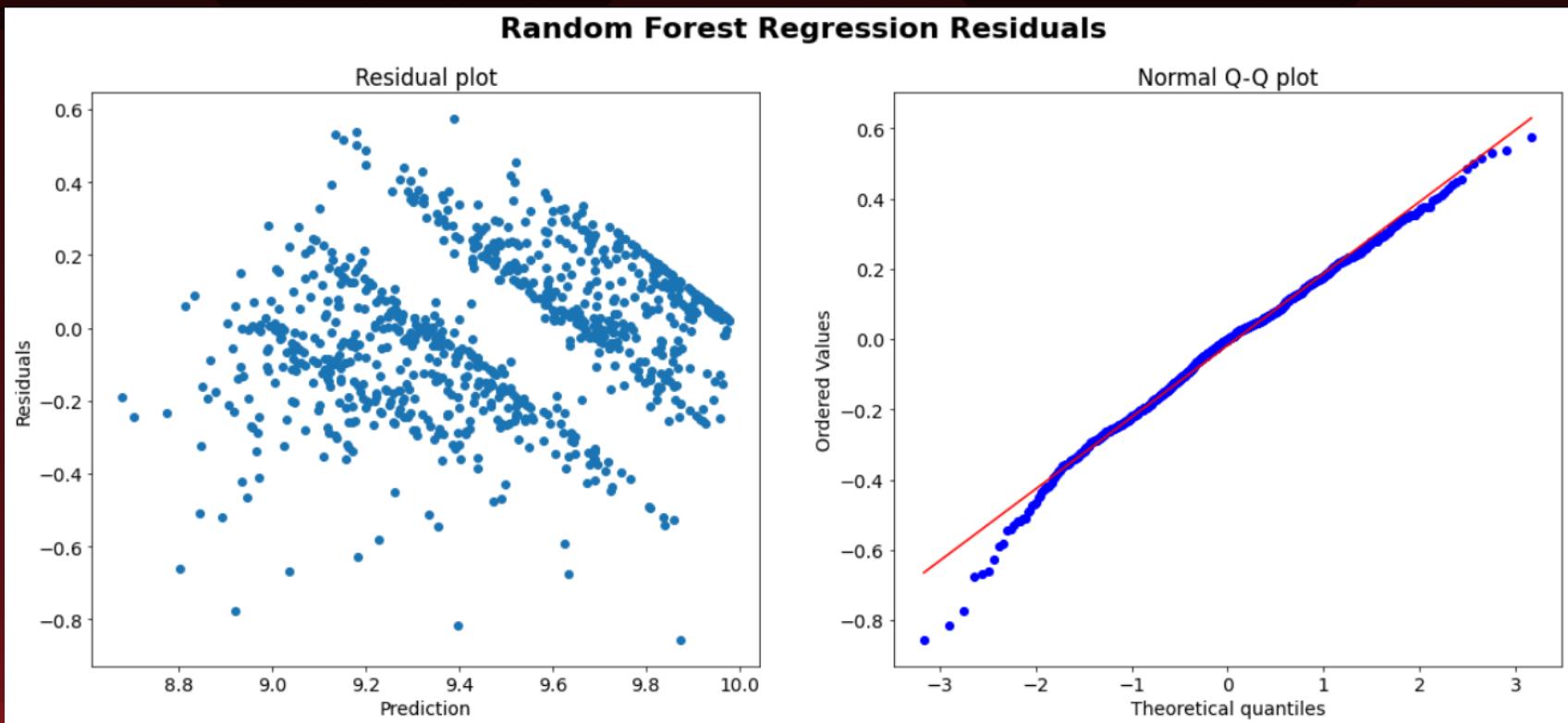
## Regression Model

- ▶ Munhoz wins with 9.57 points per round

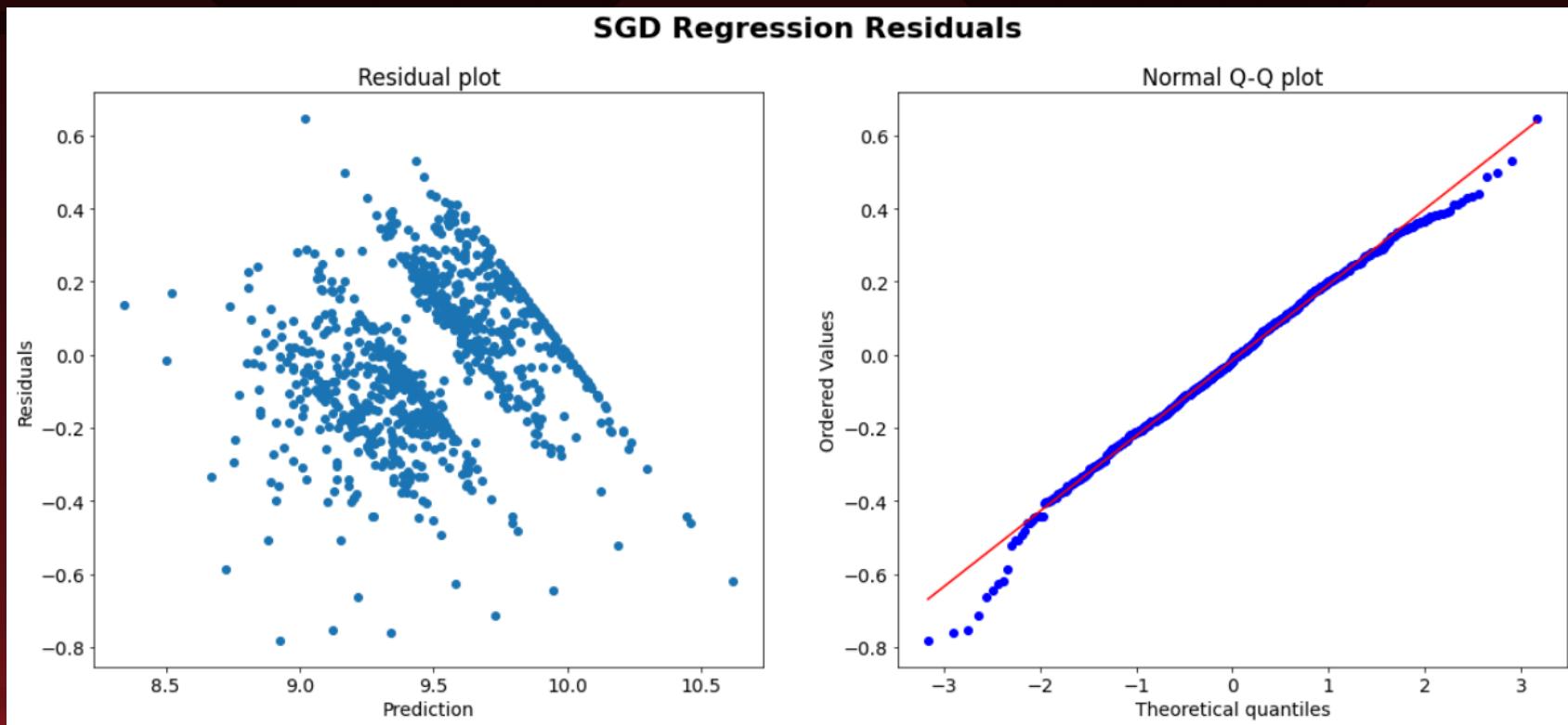
# RESIDUAL PLOTS - LINEAR REGRESSION



# RESIDUAL PLOTS - RANDOM FOREST



# RESIDUAL PLOTS – STOCHASTIC GRADIENT DESCENT



# SIMPLE VALIDATION - WINNER PREDICTION

	Linear Regression	LASSO Regression	Ridge Regression	SGD Regression	Random Forest
Unaltered Data	85.2%	85.2%	85.2%	47.2%	84.4%
MinMaxScaler (Normalized)	85.2%	85.8%	85.2%	<b>86.6%</b>	84.0%
StandardScaler (Standardized)	85.2%	85.2%	85.2%	85.1%	84.4%

- ▶ The best winner prediction accuracy in the models tested was approximately **86.6%**.
- ▶ SGDRegressor was worse than 50/50 using unaltered data. When data was normalized, it produced the best results of all models tested in simple validation.
- ▶ Normalizing with MinMaxScaler and standardizing with StandardizedScaler had minimal impact on prediction accuracy.
- ▶ Random Forest Regression had the lowest predictive ability, which makes intuitive sense given the established hierarchy of judging criteria.

# VALIDATION MODEL COMPARISON - R-SQUARED

---

Model	Data	Training R2	Simple Val R2	Simple Val Diff	Cross Val R2	Cross Val Diff
Linear Regression	Unaltered	0.699158	0.716776	0.017618	0.683173	-0.015985
Linear Regression	Normalized	0.699158	0.716776	0.017618	0.683173	-0.015985
Linear Regression	Standardized	0.699158	0.716776	0.017618	0.683173	-0.015985
LASSO Regression	Unaltered	0.699130	0.716628	0.017498	0.683170	-0.015960
LASSO Regression	Normalized	0.698653	0.715392	0.016738	0.683096	-0.015558
LASSO Regression	Standardized	0.699119	0.716667	0.017549	0.683171	-0.015948
Ridge Regression	Unaltered	0.699158	0.716775	0.017618	0.683173	-0.015985
Ridge Regression	Normalized	0.699158	0.716775	0.017618	0.683173	-0.015985
Ridge Regression	Standardized	0.699158	0.716776	0.017618	0.683173	-0.015985
SGD Regression	Normalized	0.687999	0.720987	0.032988	0.668185	-0.019814
SGD Regression	Standardized	0.698567	0.716564	0.017997	0.682459	-0.016109
Random Forest Regression	Unaltered	0.801955	0.715137	-0.086818	0.690881	-0.111074
Random Forest Regression	Normalized	0.802300	0.715964	-0.086336	0.691025	-0.111275
Random Forest Regression	Standardized	0.801614	0.715161	-0.086452	0.691330	-0.110284

# VALIDATION MODEL COMPARISON – RMSE

---

Model	Data	Training RMSE	Simple Val RMSE	Simple Val Diff	Cross Val RMSE	Cross Val Diff
Linear Regression	Unaltered	0.197553	0.204277	0.006723	0.202475	0.004922
Linear Regression	Normalized	0.197553	0.204277	0.006723	0.202475	0.004922
Linear Regression	Standardized	0.197553	0.204277	0.006723	0.202475	0.004922
LASSO Regression	Unaltered	0.197562	0.204330	0.006768	0.202476	0.004914
LASSO Regression	Normalized	0.197719	0.204775	0.007056	0.202504	0.004785
LASSO Regression	Standardized	0.197566	0.204316	0.006750	0.202476	0.004910
Ridge Regression	Unaltered	0.197553	0.204277	0.006723	0.202475	0.004922
Ridge Regression	Normalized	0.197553	0.204277	0.006723	0.202475	0.004922
Ridge Regression	Standardized	0.197553	0.204277	0.006723	0.202475	0.004922
SGD Regression	Normalized	0.200283	0.202763	0.002480	0.207538	0.007255
SGD Regression	Standardized	0.198540	0.206520	0.007979	0.202568	0.004027
Random Forest Regression	Unaltered	0.160121	0.204813	0.044692	0.198378	0.038257
Random Forest Regression	Normalized	0.160120	0.204869	0.044749	0.198724	0.038604
Random Forest Regression	Standardized	0.160181	0.205574	0.045393	0.198638	0.038456