



# Predicting Tennis Player Performance

From historical match data.



# Interesting questions to ask:

---

Which players are on the rise?

Who will unexpectedly perform super well next season?

What are the other, unexpected indicators of ranking, wins, or lifetime prize money?

How does performance in one year serve as a leading indicator of performance in the next year?

Which players had bad or good luck this season, given other stats?

Hypothetical presentation audience:

---

*ESPN is looking to predict winners in advance of matches, to highlight anomalies, or identify underdogs with promise.*

# Measures of overall player performance:

---

Win-loss ratio:

- All time
- Per year

Prize money won

- All time
- Per year

Rank (integer values)



# Demonstrated useful stats to use in our model:

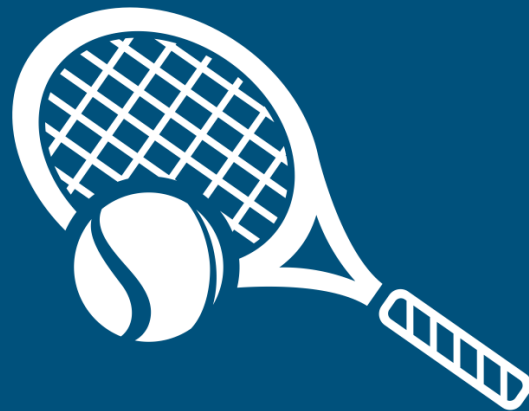
---

## More useful:

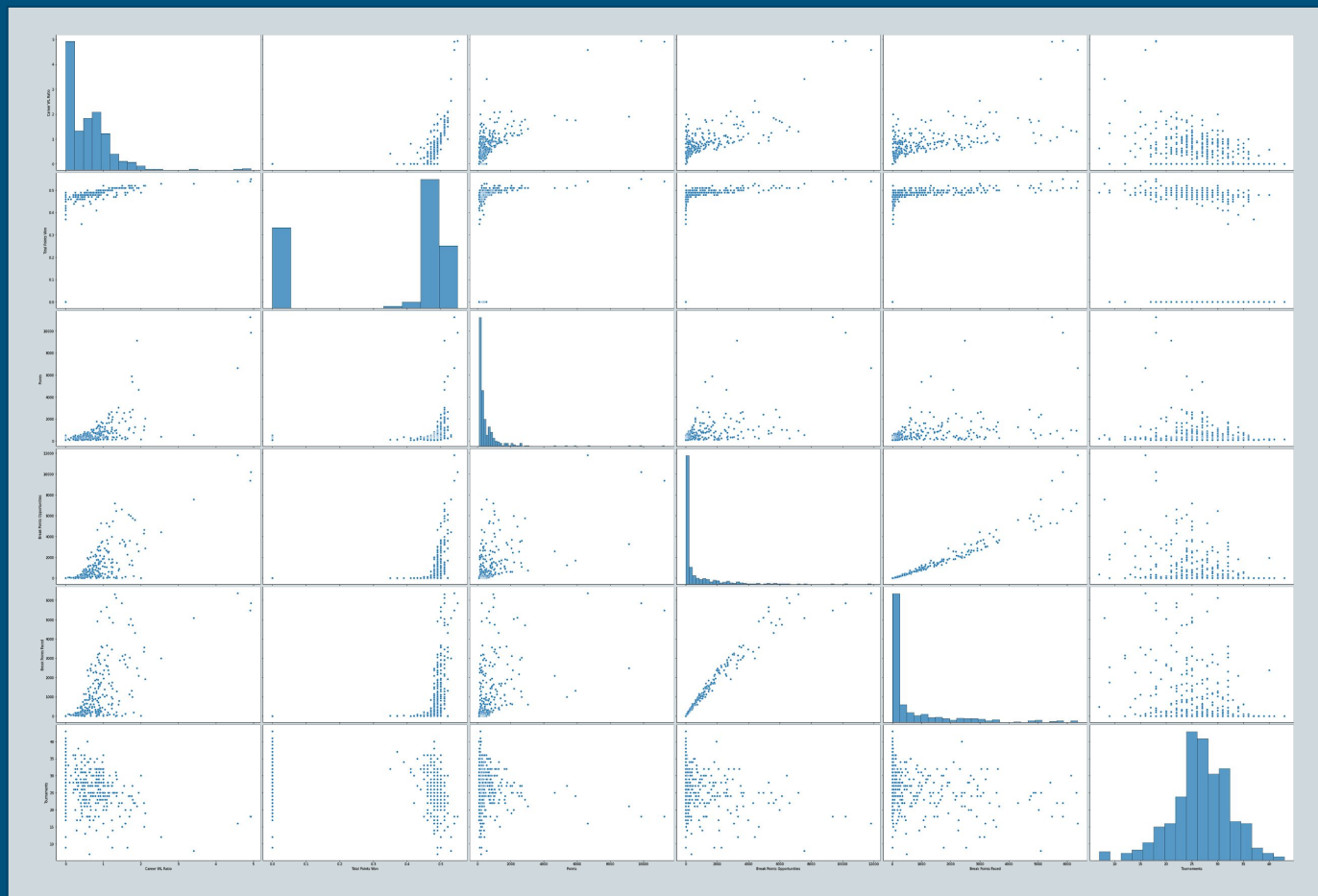
- Total points won
- Break points won

## Less useful:

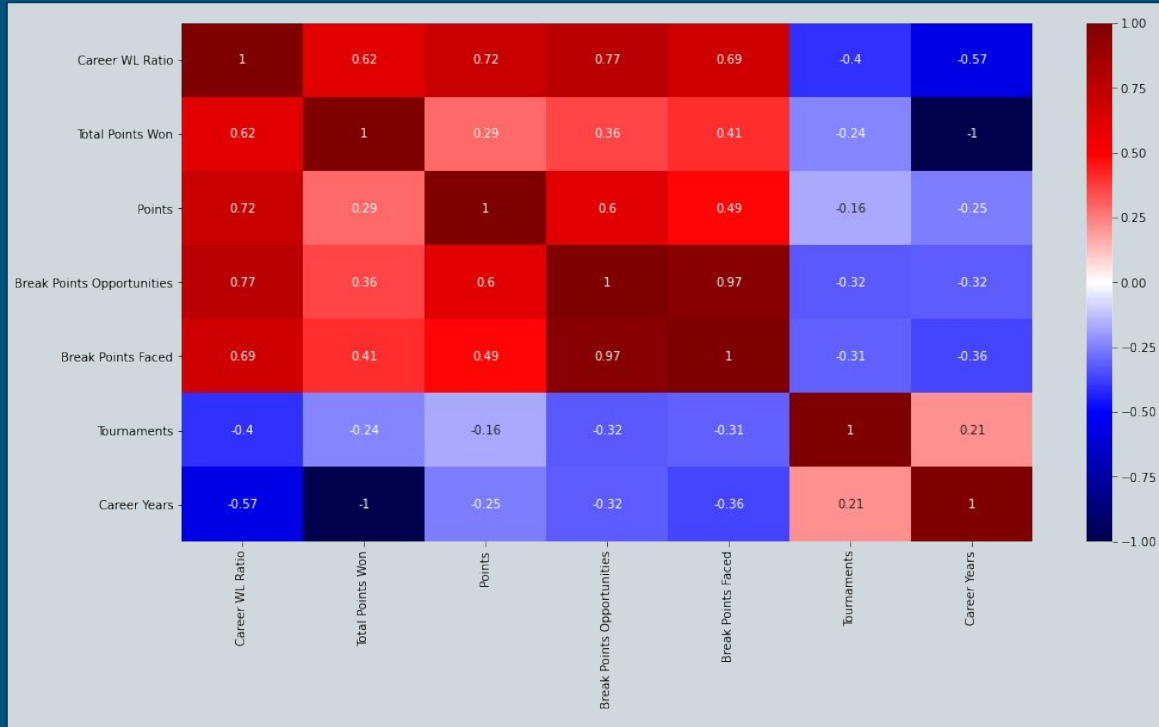
- Tournaments played (negative correlation!)
- Double faults



# Pairplots



# Correlation heatmap (post feature reduction)



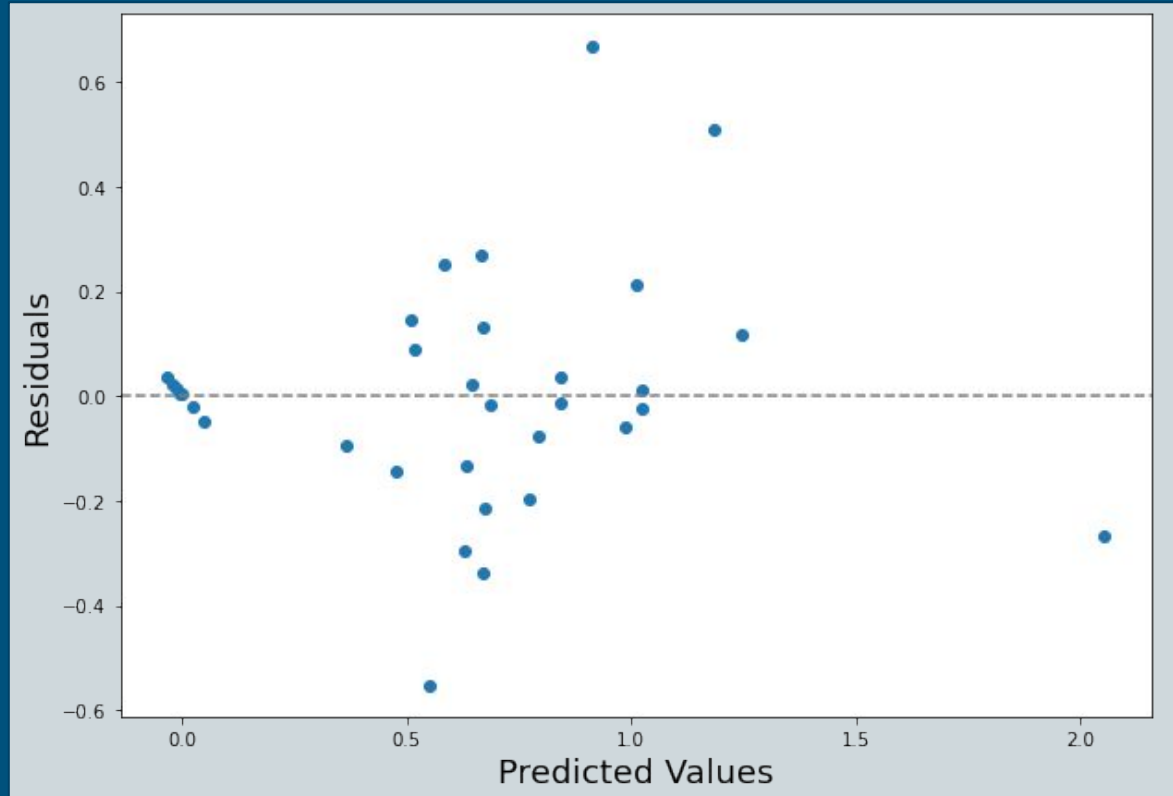
# How well do our models work?

100 players			262 players (4-way fold)		
Model	R <sup>2</sup>	RMSE	Model	R <sup>2</sup>	RMSE
<i>Std reg</i>	0.94	0.17	<i>Standard reg</i>	$0.847 \pm 0.081$	$0.214 \pm 0.013$
<i>Lasso</i>	0.94	0.26	<i>LassoCV</i>	$0.810 \pm 0.099$	$0.233 \pm 0.012$
<i>Ridge</i>	0.98	0.17	<i>RidgeCV</i>	$0.842 \pm 0.083$	$0.215 \pm 0.014$
			<i>Feat-reduced standard reg</i>	$0.855 \pm 0.048$	$0.212 \pm 0.015$



# Residual plot for final model

---



# Methodology:

---

- Web scraping with Selenium and BeautifulSoup
- Linear regression
- Regularization
- Test-train split, with holdout for validation
- Lasso and Ridge regressions
  - Now with cross-validation
- Dimensionality reduction
- Model simplification

# Appendix

---

# Challenges:

---

- CAPTCHA inhibited scraping at random intervals
  - HTML DOM for ATP Tour website was inconsistent
- Higher  $R^2$  values on smaller sample size of 100 players
- Although 400 players scraped, only 262 had detailed stats
- For next time:
  - Make scraping process more robust
  - Break out error calculation and charting into separate helper functions
  - Predict rank, even though it is integer
  - Look at intra-year data to predict break-out players
  - Logarithmic transforms to account for high-performers skewing the data