

The background of the slide is a vibrant image of a soccer stadium at night. The sky is a deep blue with wispy clouds, and two powerful floodlights on the left and right sides create bright starburst effects. Numerous smaller lights along the stadium's rim cast beams of light across the scene. The soccer field is a lush green, with the center circle clearly visible in the foreground.

# Predicting the Overall Rating of FIFA players with Linear Regression

Jesús Cabezas


Sandra Cunha

# Content



1. Column Selection
2. Data Cleaning
3. Exploratory Data Analysis
4. Correlation Heatmap
5. Features Selection
6. Numerical Features Transformation
7. Categorical Features Transformation
8. Modeling
9. Model Evaluation
10. Conclusion

# Column Selection

- 
- Age
  - OVA
  - Nationality
  - Club
  - BOV
  - Position
  - POT
  - Value
  - Crossing
  - Finishing
  - Heading
  - Short Passing
  - Volley
  - Dribbling
  - Curve
  - FK Accuracy
  - Long Passing
  - Ball Control
  - Acceleration
  - Sprint Speed
  - Agility
  - Reactions
  - Balance
  - Stamina
  - Strength
  - Long Shots
  - Aggression
  - Interceptions
  - Positioning
  - Vision
  - Shot Power
  - Jumping
  - Penalties
  - Marking
  - Standing Tackle
  - Sliding Tackle
  - GK Diving
  - GK Handling
  - GK Kicking
  - GK Positioning
  - GK Reflexes

# Data Cleaning



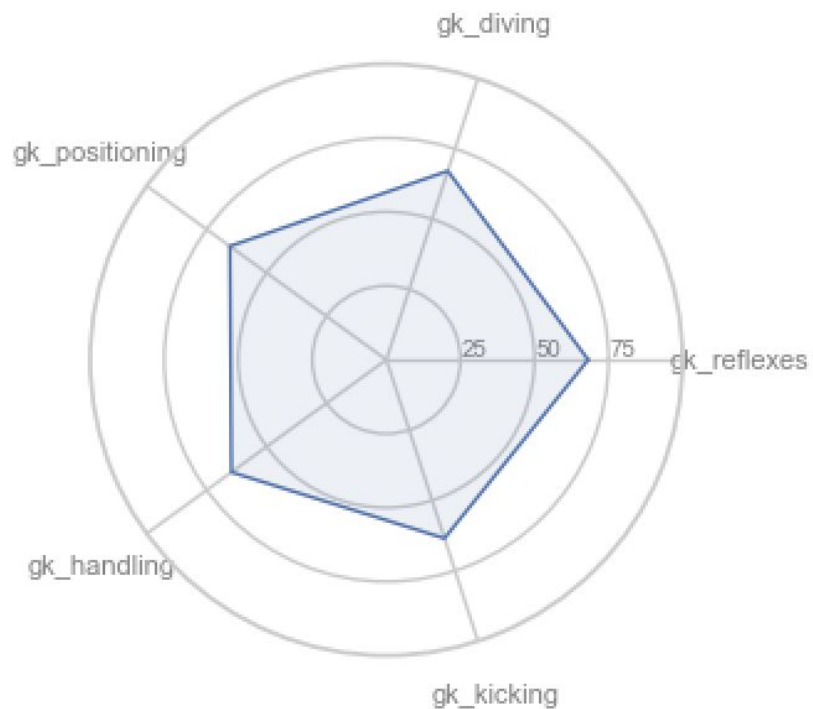
- Standardize header names
- Checking for nulls
- Grouping values in “Position” Column to 'DF', 'MF', 'AT', 'GK'

# Exploratory Data Analysis

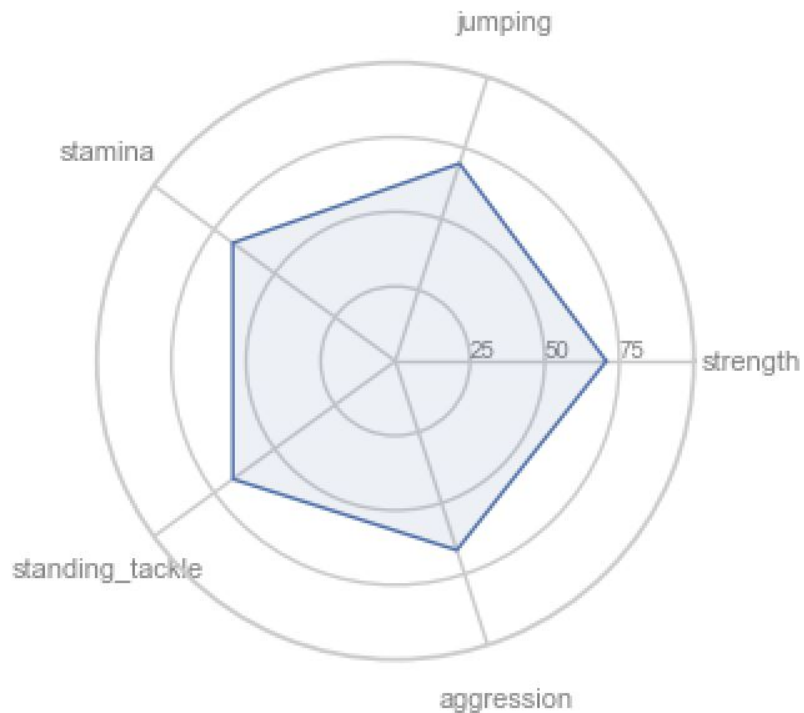


# Top 5 Features by Position

GK

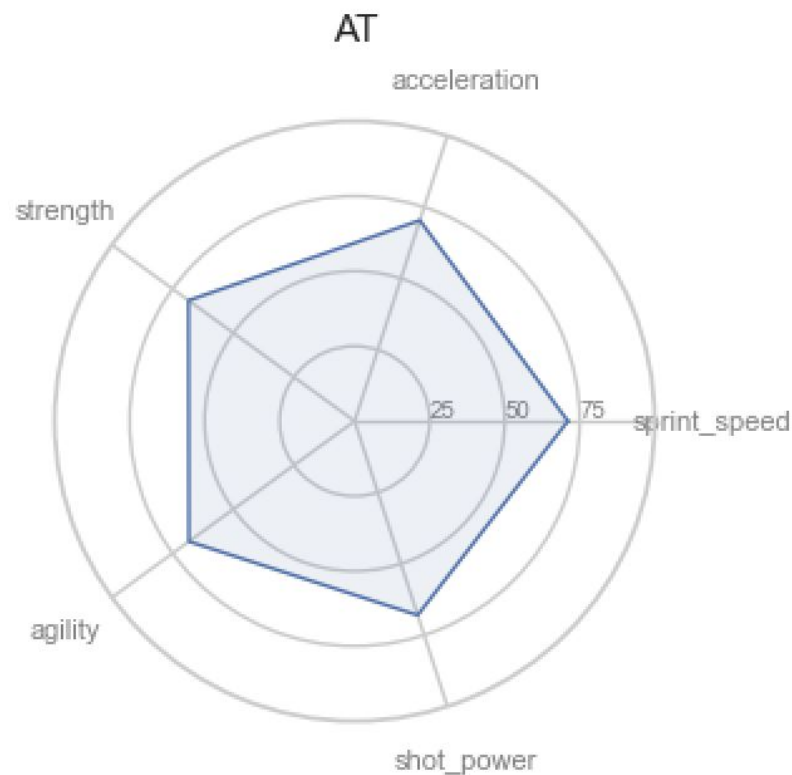
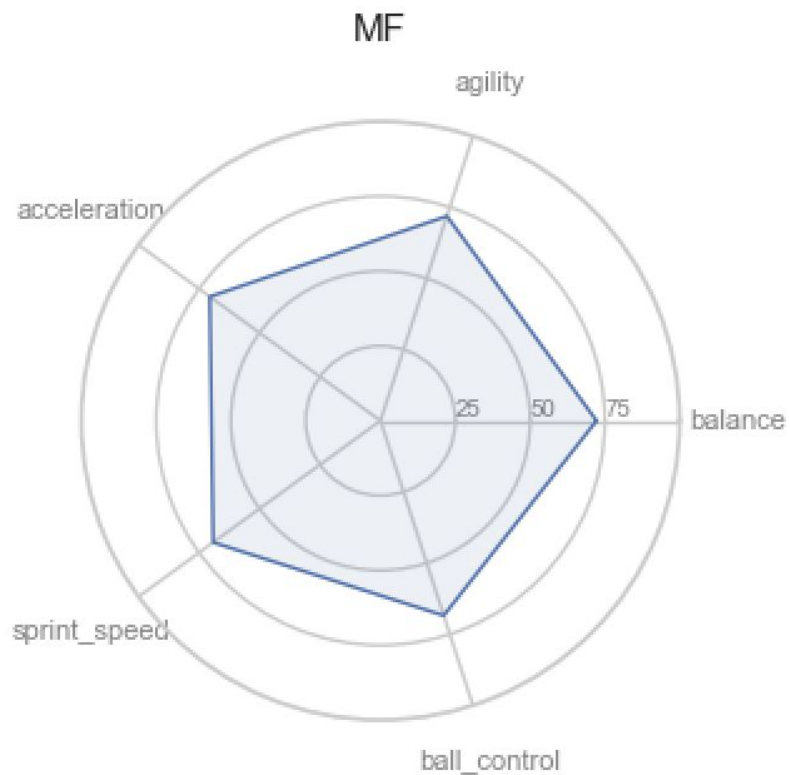


DF

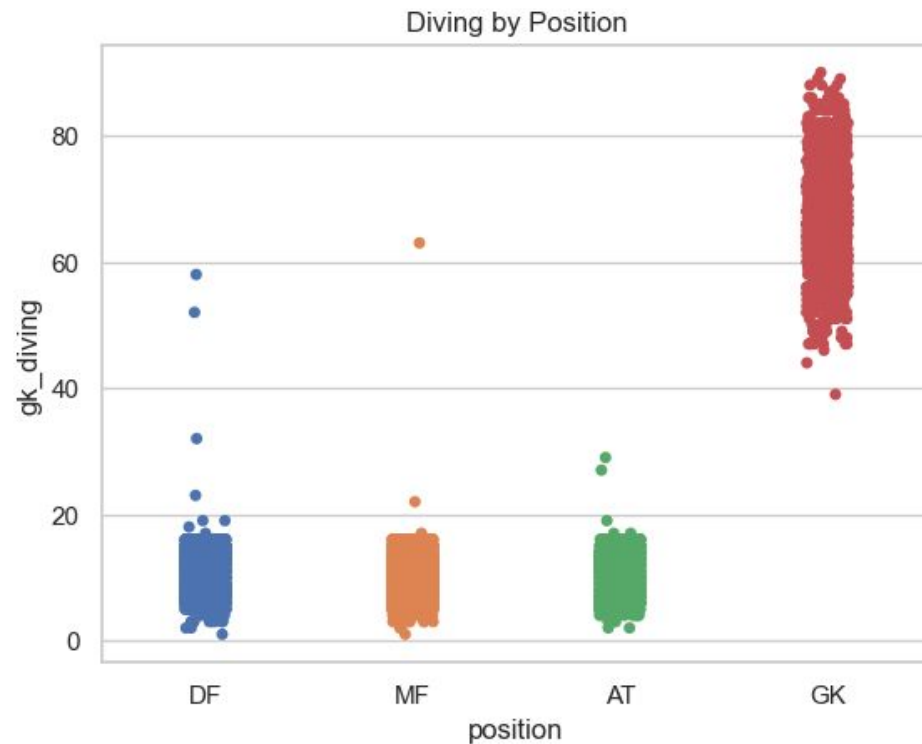
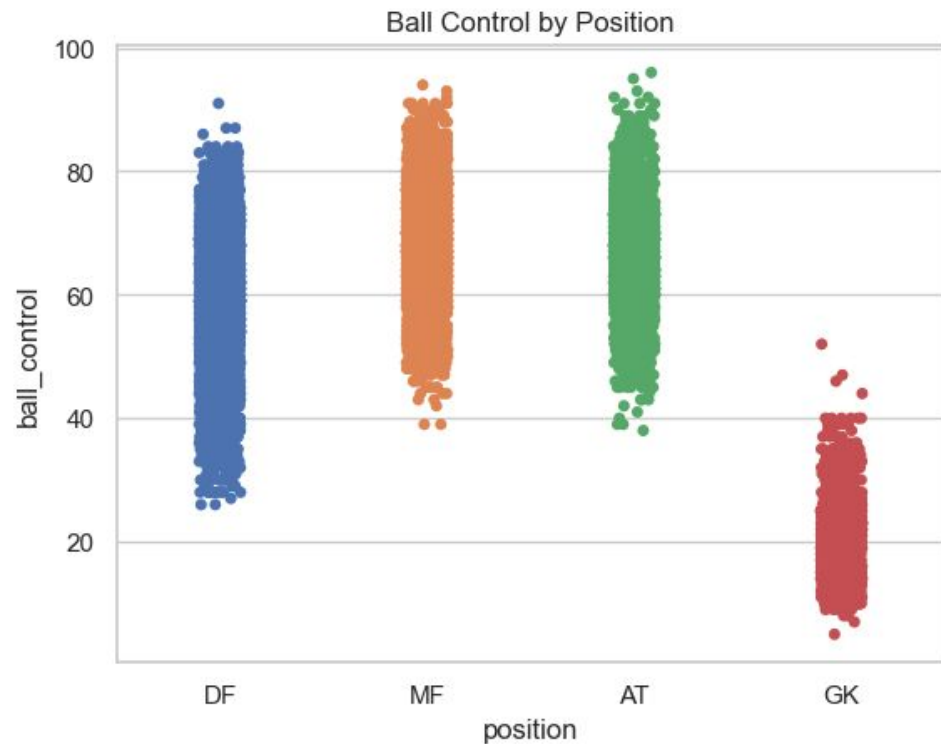




# Top 5 Features by Position

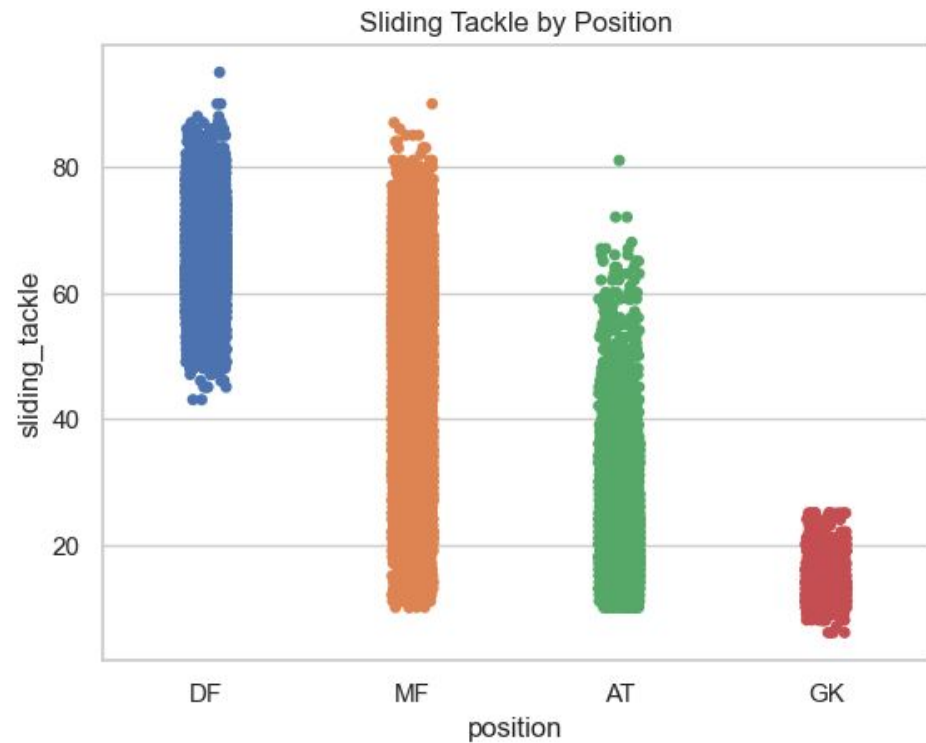
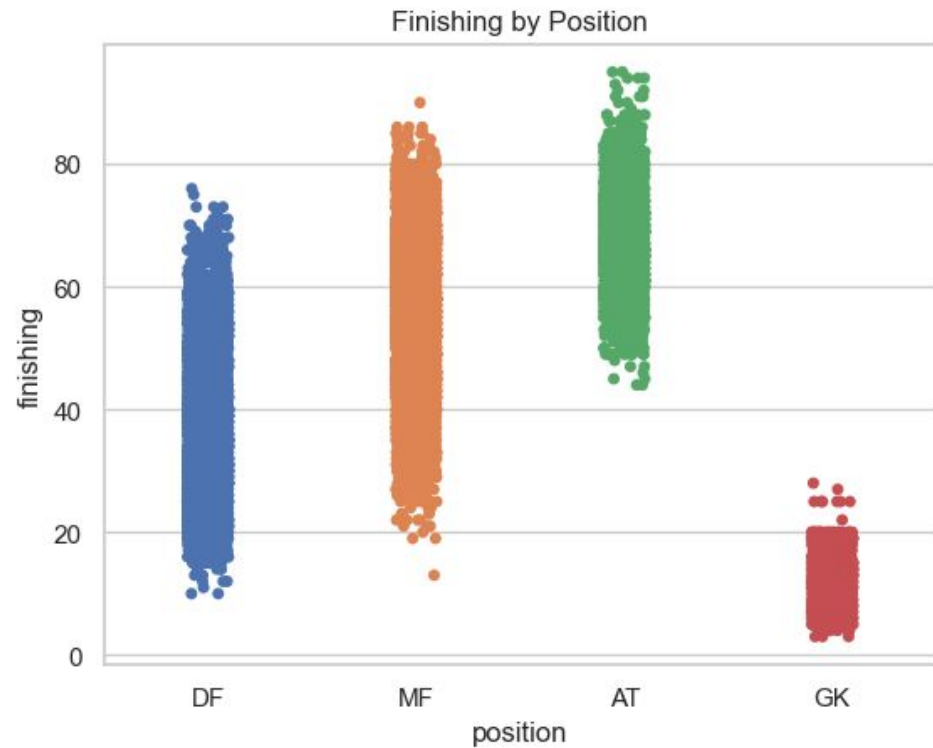


# Some differences in features

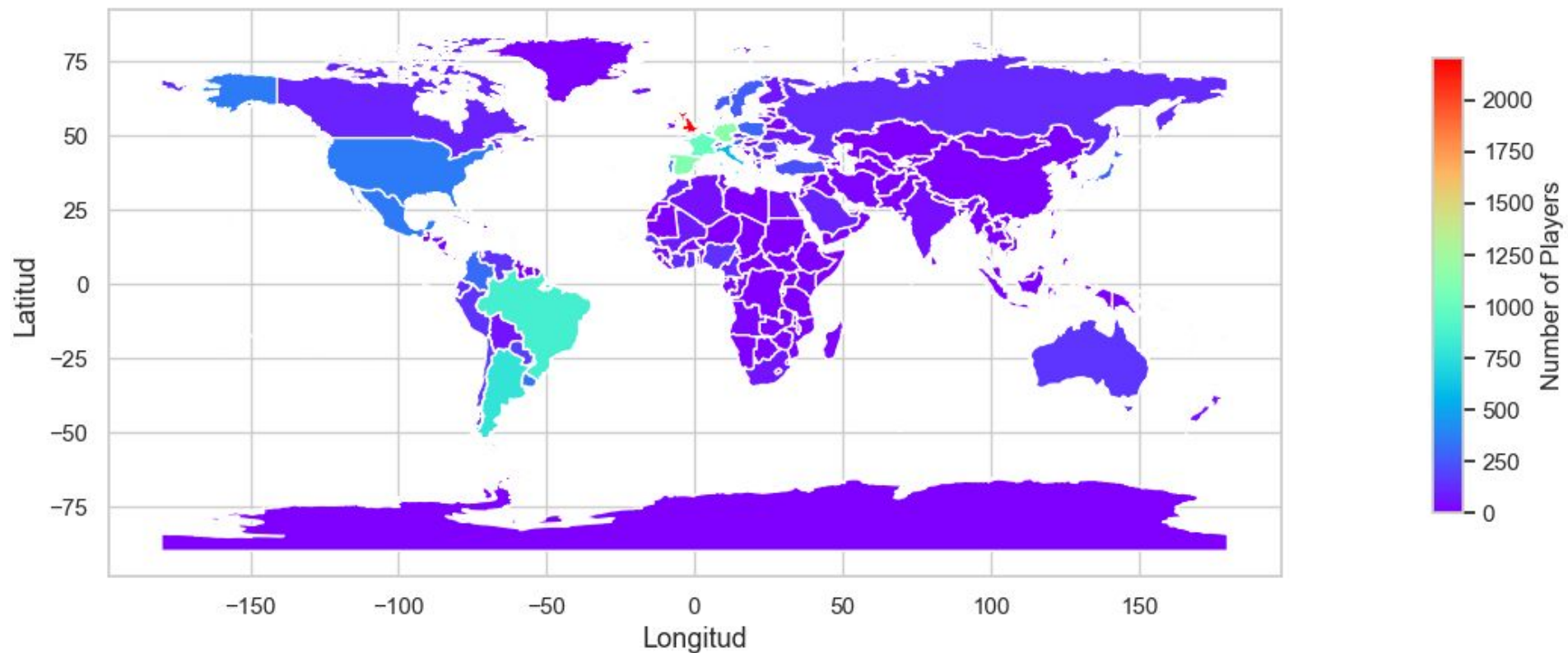




# Some differences in features



# Players per Country

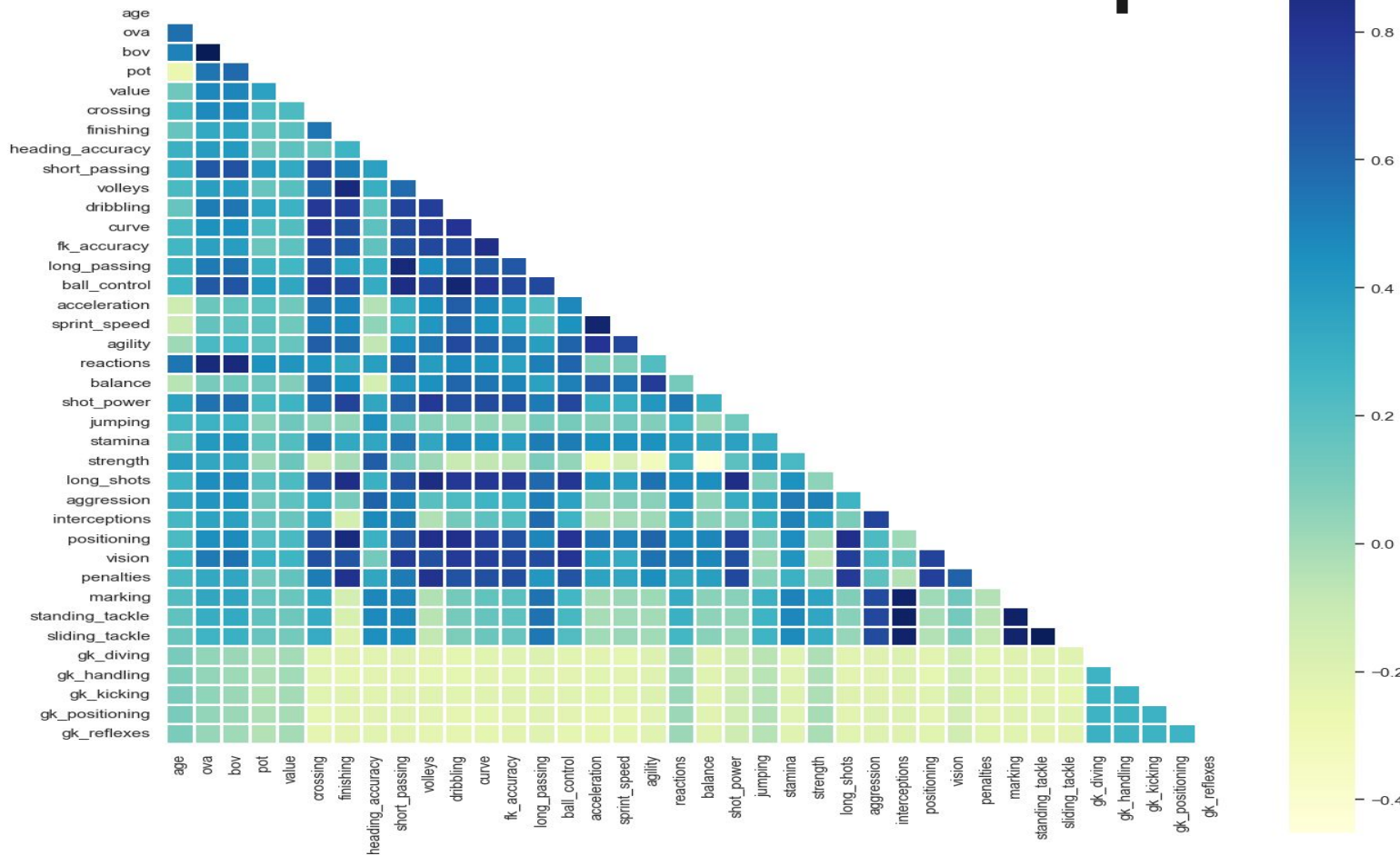




# Correlation Heatmap

- ❖ Understand which variables are related to each other;
- ❖ The relationship with our target: “ova” - Overall Rating;
- ❖ Feature selection: which columns are relevant for our model?

# Correlation Heatmap





## Feature Selection

Moderate to High Correlated columns (above 0.8) with our target “ova”:

- ❖ vision
- ❖ shot\_power
- ❖ reactions
- ❖ ball\_control
- ❖ short\_passing
- ❖ bov

# Preprocessing



## Numerical Features

- Data scaling
- Box-cox

## Categorical Features

- Dummy Encoding

Nationality

Position

# Modeling



## X-y split

`y = finaldata["ova"]` – Our Target “Overall”

`X = finaldata.drop(["ova"], axis=1)`

## Train-test split

30% --> test , 70% --> train



# Model Evaluation



## Results with Box-cox:

- $R^2 = 0.98$
- $RMSE = 0.71$
- $MSE = 0.50$
- $MAE = 0.53$

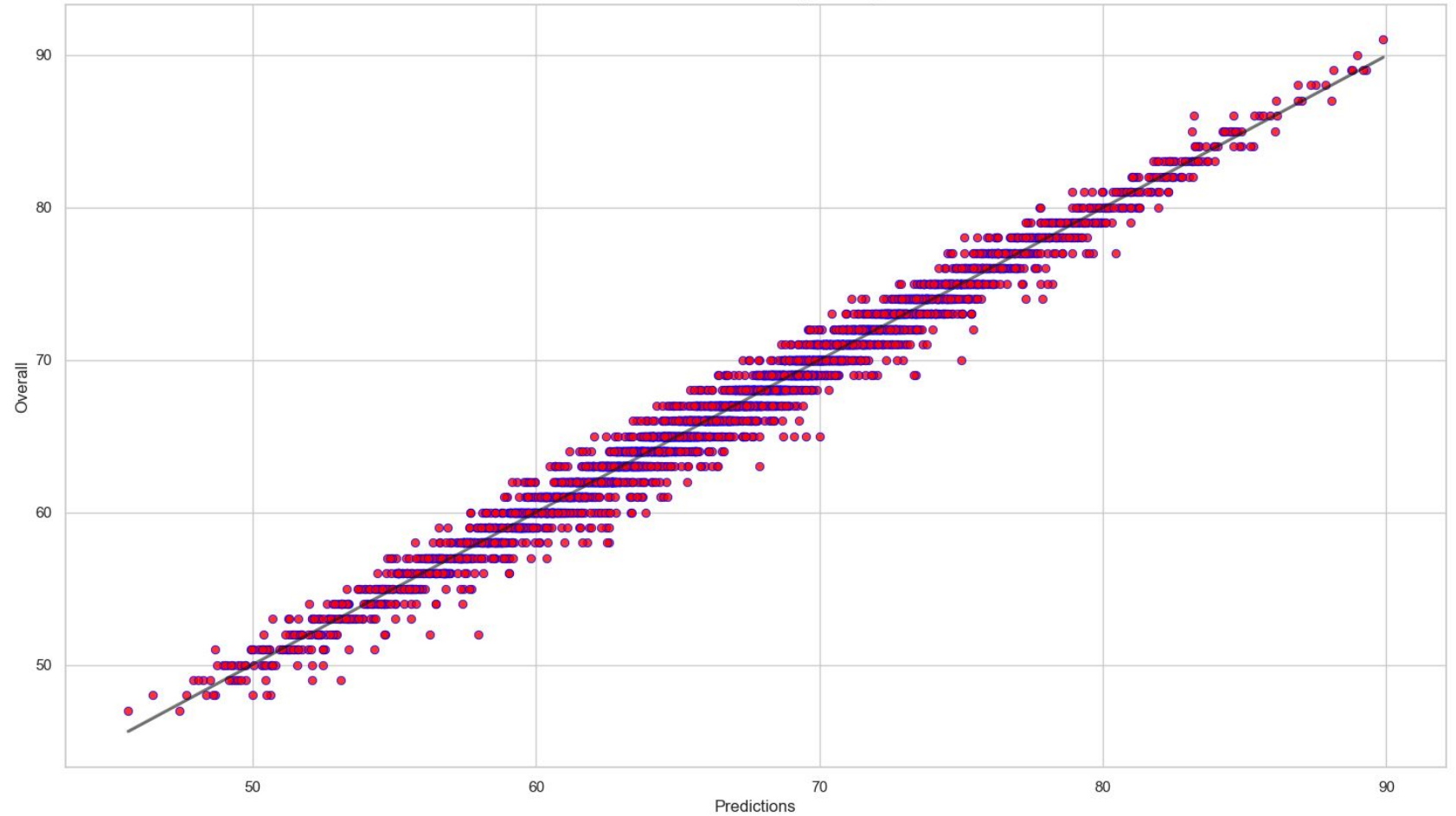
## Results without Box-cox:

- $R^2 = 0.98$
- $RMSE = 0.92$
- $MSE = 0.85$
- $MAE = 0.69$

98% of the variance of the target is explained by the variance of the features

$\pm 0.7 / \pm 0.9$  difference between values predicted by the model and the actual values

Linear Prediction of Player Rating



# Conclusions

- Model performs better without categoricals;
- Slightly better prediction using box-cox transformation;

Thank you!!

