

# A Survey of Advanced Modeling Techniques for Forecasting College Football Game Outcomes

Charles South<sup>1</sup>, PhD and Edward Egros<sup>2</sup>, MS

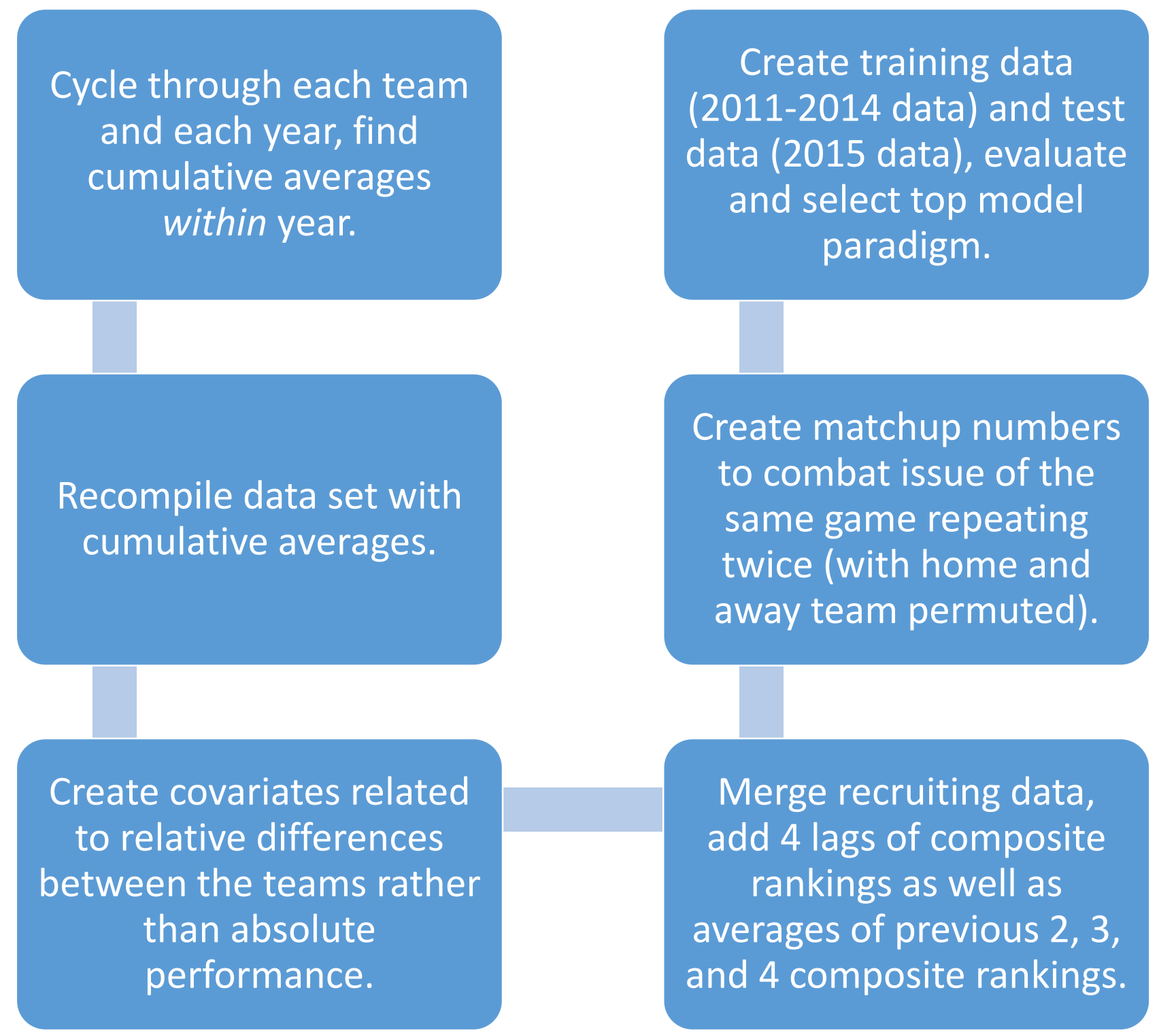
## Purpose:

Can we use box score, point spread, and recruiting data to accurately forecast outcomes of college football games using modern machine learning and Bayesian modelling approaches?

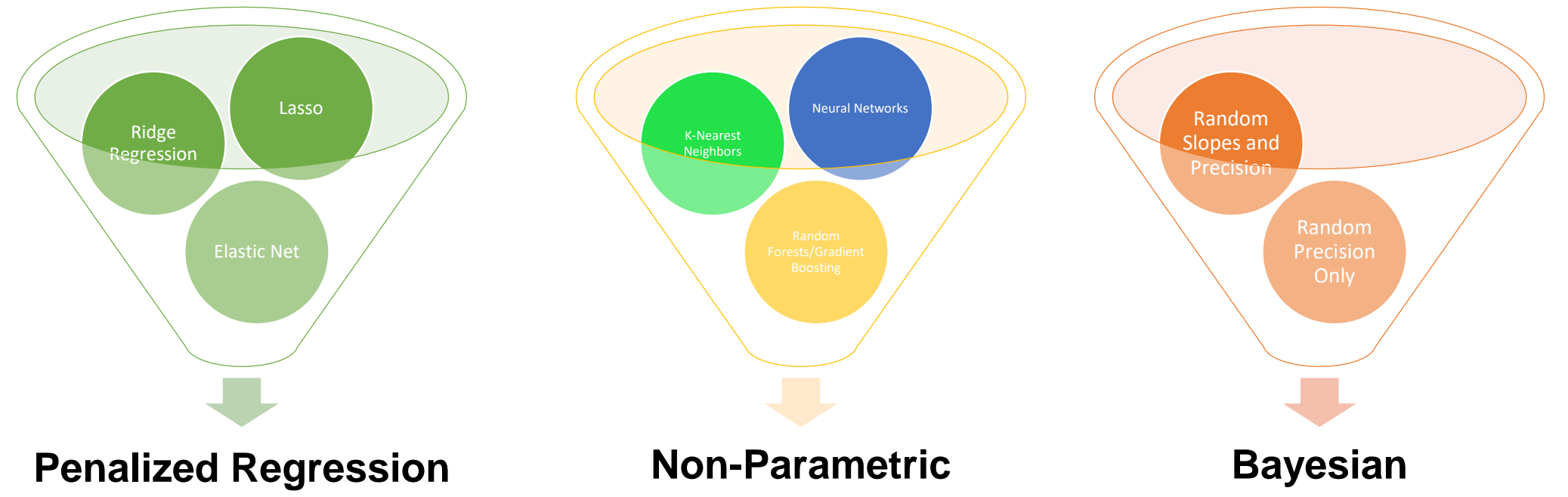
## Data:

The data used for this college football forecast consisted of 4,339 games between Football Bowl Subdivision (FBS) teams between the 2011 and 2016 seasons<sup>3</sup>. The data set was augmented with composite team ratings<sup>4</sup> for each recruiting class dating back to 2008 to account for seniors in the 2011 season.

## Data Prep/Manipulation Process:



## Modelling Paradigms:



## Modelling Details:

- Dependent Variable: point differential
- R packages utilized: *glmnet*, *randomForest*, *FNN*, *nnet*, *xgboost*.
- Where applicable, the *caret* package was used to tune the models.
- For the Bayesian models, a the model was compressed to only include variables selected by the lasso.
- Bayesian framework:

$$\begin{aligned}
 & \text{predicted point difference for team } i \text{ on game } j \\
 & y_{ij} \sim \text{Normal}(\mu_{ij}, \tau_i), \\
 & \mu_{ij} = \beta_0 + \beta_1 X_{1ij} + \dots + \beta_p X_{pij}, \\
 & \beta_p \sim \text{Normal}(0.001, 0.001), \\
 & \tau_i \sim \text{Gamma}(1, 1)
 \end{aligned}$$

random precision for team *i*

Lasso Retained Variables		
Yards Per Pass Attempt (YPPA)	Yards Per Rush Attempt (YPPRA)	Rush Attempts
Total Yards	Yards Per Play (YPP)	Turnovers (TO)
Opponent Points Scored	Opponent YPPRA	Opponent Total Yards
Opponent Turnovers	Opponent Penalty Yards	Average Point Differential
Opponent Offense Passing Yards	Opponent Offense YPPRA	Opponent Offense Total Yards
Opponent Offense YPP	Opponent Defense Total Rush Yds	Opponent Defense YPPRA
Opponent Defense Total Yards	Opponent Def YPP	Opponent Defense TO
Opponent Average Points Differential	Difference in Win Percentage	Home Field Advantage
Composite Ranking, Lag 2 (CR)	Average CR (Last 2 Years)	Average CR (Last 3 Years)

## Final Predictions:

All generated models made predictions for the same game twice. Our decision rule was designed to be as intuitive as possible – if both predictions resulted in the same team being favored, the favored team was deemed the predicted winner. If the predictions diverged, the team with the larger predicted point differential was deemed the predicted winner.

## Summary of Results:

Model	Overall Prediction	PPV	NPV
Lasso	75.0%	77.3%	72.4%
Random Forest	74.1%	78.9%	69.0%
Gradient Boosting	72.9%	76.9%	68.5%
K-Nearest Neighbors	69.0%	71.3%	66.3%
Neural Network	71.6%	75.8%	66.4%
Bayesian	74.1%	76.5%	71.5%

## Conclusions/Future Work:

- The lasso was the top performer, slightly edging out the Bayesian model and random forests.
- In terms of variable importance, composite recruit rankings are important, though it appears it takes at least 2 years for players to develop.
- Volume and efficiency are important – as is home field advantage.
- Of the 27 predictors identified by the lasso, 15 were directly related to opponent strength. Schedule matters.
- An ensemble method may improve predictive ability.

<sup>3</sup>[http://www.seldomusedreserve.com/?page\\_id=8805](http://www.seldomusedreserve.com/?page_id=8805)  
<sup>4</sup><http://247sports.com/Article/247Rating-Explanation-81574>

<sup>1</sup>South Statistical Consulting and Analytics, LLC  
<sup>2</sup>Fox 4 News (Dallas), InsideSportsAnalytics.com