FiveThirtyEight's club football predictions Sports modelling reading group

Clement Lee

2023-04-28

Nate Silver



Nate Silver

- ▶ Predicted the outcomes in all but one states in 2008 U.S. presidential election
 - 2008: All states but one
 - 2012: All states
 - 2016: Least worst of major forecasts
 - ► 2020: All states
- ► 538 electors in the Electoral College

FiveThirtyEight

Does other forecasts as well



Updated April 25, 2023, at 4:56 p.m.

REMIER EAGUE
77%
23%

England

Outline

- 1. Soccer Power Index (SPI)
 - Nate originally devised for ESPN
 - Revised substantially over the years
- 2. Forecasting
 - Poisson model
 - Monte Carlo simulations
- 3. Leagues and tiers
 - For e.g. Champions League

Soccer Power Index (SPI)

What is SPI?

- Overall rating: the percentage of available points expected to take
- Used in turn for the simulations

Pre-season SPI

 \blacktriangleright 2/3 \times end-of-season's SPI + 1/3 \times SPI implied by market value

Reason: market value strongly correlated with end-of-season SPI

SPI during season

- Updated after each match
- > Calculated according to an offensive rating and a defensive rating
 - Opaque about how

TEAM	SPI	OFF. DEF.
💿 Man. City	92.3	3.0 0.3
\delta Real Madrid	86.9	2.7 0.6
🛞 Inter Milan	77.7	2.3 0.7
🛞 AC Milan	74.5	2.1 0.8

Offensive rating

- Goals expected to score against an "average" team
- Using actual goals?
 - Scoreline often disagrees with people's impressions of the quality of each team's play
 - ► Low-scoring nature → prolonged periods of luck (good results despite playing poorly), or vice versa

Offensive rating

Average of

- 1. Goals
- 2. Adjusted goals
- 3. Shot-based expected goals
- 4. Non-shot expected goals

 ${\sf Defensive\ rating\ =\ offensive\ rating\ of\ opposing\ team}$

High defensive rating is bad

Adjusted goals

- Reducing value of goals scored when more players
- Reducing value of goals scored late when leading
- \blacktriangleright Increasing value of other goals to add up to total # actual goals

Shot-based expected goals

- Each shot assigned a probability based on distance & angle
- Part of the body the shot was taken with
- Adjustment for the player

Non-shot expected goals

- Passes, interceptions, take-ons and tackles
- \blacktriangleright Intercepting the ball at opposing team's penalty spot leads to a goal \approx 9% of the time
- \blacktriangleright Pass completed at the center of the six-yard box leads to a goal \approx 14% of the time

SPI pipeline

- 1. End-of-season SPI & market value \rightarrow pre-season SPI
- 2. Use current SPI to forecast match result (next section)
- 3. Use match data to calculate offensive & defensive ratings
- 4. Update current SPI after match using the ratings
 - Unlike ELO, a win doesn't necessarily improve SPI

Forecasting

Simulating number of goals

- Poisson with mean the offensive rating
- Adjusted for
 - league-specific home-field advantage
 - importance of the match to each team
- Two independent Poisson distributions
- Inflated for draws (league-specific, around 9%)

Scoring matrix



Importance, quality, match rating



Monte Carlo simulations

- Using the scoring matrix
- > Teams' SPI can rise and fall depending on simulated matches played
- Widened distribution of possible outcomes

Whole pipeline

- 1. End-of-season SPI & market value \rightarrow pre-season SPI
- 2. Use current SPI to forecast match result (next section)
- 3. Use match data to calculate offensive & defensive ratings
- 4. (Update current SPI after match using the ratings)
- 5. Build scoring matrix using offensive ratings & draw inflation
- 6. Simulate the (remaining) matches of the season

Performance - Ranked Probability Score (RPS)

- Measures how good forecasts are in matching observed outcomes
- \triangleright RPS = 0: wholly accurate; RPS = 1: wholly inaccurate
- Current RPS = 0.1957 after reductions:
 - by 0.0018 using expected-goals metrics
 - by 0.0011 using market values to pre-season SPI
 - by 0.0006 using match importance

Leagues and tiers

Relative strengths between leagues

- ► For predictions of e.g. Champions League
- Originally used a tiered system

European countries' soccer leagues, sorted into six tiers by strength

TIER	LEAGUES
1	England, Germany, Italy, Spain
2	France
3	Portugal
4	Belgium, Czech Republic, Netherlands, Russia, Ukraine
5	Austria, Bulgaria, Croatia, Denmark, Finland, Greece, Hungary, Ireland, Israel, Norway, Poland, Romania, Scotland, Slovakia, Slovenia, Sweden, Switzerland, Turkey
6	Albania, Andorra, Armenia, Azerbaijan, Belarus, Bosnia, Cyprus, Estonia, Faroe Islands, Georgia, Iceland, Kazakhstan, Latvia, Lithuania, Luxembourg, Macedonia, Malta, Moldova, Montenegro, Northern Ireland, Serbia, Wales

Top divisions only. Countries are listed alphabetically within each tier.

Current version

- 1. Calculate SPI using domestic league results only
- 2. Calculate expected scores of inter-league matches using domestic SPI
- 3. Apply Massey's method to actual expected scores to find league strengths
- 4. Regress league strengths toward market-value based ratings
- 5. Run through all matches again, this time with the league strengths

Interpretation: bonus (in goals) given to each team in an inter-league match

Summary

Modularised approach

- 1. Match data to calculate offensive/defensive ratings
- 2. Offensive ratings to calculate scoring matrix, with adjustments
- 3. Scoring matrix to simulate future matches
- 4. Domestic results to compute league strengths

Extra 1: @Experimental361

HOW MUCH THE TABLE CAN CHANGE: PREMIER LEAGUE, 1-3 OCT 2022

OUTLINES SHOW SPREAD OF POSSIBLE MOVEMENT; INNER CHARTS SHOW MODELLED PROBABILITIES



Extra 2: Engsoccerdata GitHub repository

https://github.com/jalapic/engsoccerdata

engsoccerdata

This R package is mainly a repository for complete soccer datasets, along with some built-in functions for analyzing parts of the data. Currently I include three English ones (League data, FA Cup data, Playoff data - described below), several European leagues (Spain, Germany, Italy, Holland, France, Belgium, Portugal, Turkey, Scotland, Greece) as well as South Africa and MLS.

Free to use for non-commerical use. Compiled by James Curley.

Please cite as: James P. Curley (2016). engsoccerdata: English Soccer Data 1871-2016. R package version 0.1.5