

Statistical modelling in soccer

Hans Van Eetvelde

Napels, April 8, 2018

Overview

Modelling soccer games

- Directly Modelling soccer outcomes

- Modelling scores

Estimating parameters

- Maximum Likelihood

- Weighted Maximum Likelihood

Applications

- Prediction

- Ranking

Bradley Terry

First assume that there are only two outcomes: a home win (H) or an away win (A)

We model the chances on these outcomes as follows:

$$P_H = \frac{H}{H + A};$$

$$P_A = \frac{A}{H + A};$$

where H is the strength of the home team and A is the strength of the away team.

Possibility of a draw

Of course, in soccer we have also the possibility of a draw, which we will model as follows:

$$P_H = \frac{p_H}{H + d' \frac{p_H}{H A + A}};$$

$$P_D = \frac{d' \frac{p_H}{H A + A}}{H + d' \frac{p_H}{H A + A}};$$

$$P_A = \frac{p_A}{H + d' \frac{p_A}{H A + A}};$$

Exercise: Assume p_H is constant. Prove that $P_D(A)$ is maximal for $A = H$.

Adding the Home advantage

A first simple addition to this model is adding a home effect. The strength of the home team is multiplied with h :

$$P_H = \frac{h_H}{h_H + d} \frac{p_H}{h_H A + A};$$

$$P_D = \frac{d}{h_H + d} \frac{p_H A}{h_H A + iA};$$

$$P_A = \frac{p_A}{h_H + d} \frac{A}{h_H A + A};$$

Disadvantage of the Bradley-Terry Model

It does only consider the outcomes of the matches (Home win, Draw, away win), but not the scores.

Winning a game by 5-0 or by 1-0 is not the same!

) We will try to model the scores

Modelling scores

We can indirectly model the outcome of the game by modelling the home and away score. G_H is the number of goals made by the home team and G_A is the number of goals made by the away team, then we have

$$P_H = P(G_H > G_A)$$

$$P_D = P(G_H = G_A)$$

$$P_A = P(G_H < G_A)$$

Poisson distribution

- | How to model a number of goals?
- | Most used 'count distribution': Poisson distribution

$$P(G = g) = \frac{\mu^g}{g!} e^{-\mu}$$

where the parameter μ is the mean.

- | In this case: μ is the average number of goals scored

Soccer scores

Let G_H be the number of goals made by the home team and G_A the number of goals scored by the away team, then we have for the Poisson model:

$$P(G_H = x) = \frac{\lambda_H^x}{x!} e^{-\lambda_H}$$
$$P(G_A = x) = \frac{\lambda_A^x}{x!} e^{-\lambda_A}$$

where λ_H is the expected number of goals scored by the home team and λ_A is the expected number of goals scored by the away team.

From scores to outcome

If the variables G_H and G_A are Poisson distributed with parameters μ_H and μ_A , then the distribution of the difference $G_H - G_A$ is a known distribution, named the 'Skellam' distribution. Using this distribution we can easily compute the chances on every outcome:

$$P_H = P(G_H - G_A > 0)$$

$$P_D = P(G_H - G_A = 0)$$

$$P_A = P(G_H - G_A < 0)$$

Poisson model

Idea from Maher (1982)

$$H = C \frac{H}{A}$$
$$A = C \frac{A}{H}$$

where H is again the strength of the home team, A is the strength of the away team and C is the expected number of goals for teams with equal strengths.

Adding the home effect (bis)

Idea from Maher (1982)

$$H = c + h + \frac{H}{A}$$
$$A = c + \frac{A}{H}$$

where H is again the strength of the home team, A is the strength of the away team and c is the expected number of goals for teams with equal strengths.

Attack and defence strengths

Instead of using one strength parameter per team, we could use separate parameters for a team's attacking strength and a team's defensive strength:

$$H = c \cdot h \cdot \frac{H;\text{att}}{A;\text{def}}$$

$$A = c \cdot \frac{A;\text{att}}{H;\text{def}}$$

Correlation between home and away score

Consider again the model with one strength parameter per team

$$H = c \cdot h \cdot \frac{H}{A}$$

$$A = c \cdot \frac{A}{H}$$

There is clearly a (negative) correlation between the home and the away score.

Should we allow for other types of correlation)? Bivariate Poisson

Bivariate Poisson Model

In the simple Poisson model we have

$$\begin{aligned}
 P(G_H = x; G_A = y) &= P(G_H = x) P(G_A = y) \\
 &= \frac{\lambda_H^x}{x!} e^{-\lambda_H} \frac{\lambda_A^y}{y!} e^{-\lambda_A} \\
 &= \frac{\lambda_H^x \lambda_A^y}{x! y!} e^{-(\lambda_H + \lambda_A)}
 \end{aligned}$$

In the bivariate model, we allow for a positive covariance between G_H and G_A

$$P(G_H = x; G_A = y) = \frac{\lambda_H^x \lambda_A^y}{x! y!} e^{-(\lambda_H + \lambda_A + c)} \sum_{k=0}^{\min(x,y)} \frac{c^k}{k!} \frac{c}{\lambda_H \lambda_A}$$

Bivariate Poisson Model

- | The bivariate Poisson distribution is derived from the simple Poisson distribution as follows: X_H , X_A en X_C are Poisson distributed with parameters μ_H , μ_A en μ_C , then $G_H = X_H + X_C$ en $G_A = X_A + X_C$ are bivariate poisson distributed.
- | So the number of goals for both teams now consist out of a team-specific component and a common component.

How to determine c_C ?

We will determine b_H and b_A as before, for c_C we have several options:

1. $c_C = b_C = \text{constant}$
2. $c_C = b_C \quad b_H$ depending on the home team
3. $c_C = b_C \quad b_A$ depending on the away team
4. $c_C = b_C \quad b_H \quad b_A$ depending on both teams

Other ways to model soccer goals

- | One can use other count distributions instead of the Poisson distribution. (For example the Negative binomial or the Weibull Count Distribution are used)
- | Other ways to model the dependency between home and away scores?
- | We look for a function c for which:

$$P(G_H = x; G_A = y) = c(P(G_H = x); P(G_A = y))$$

This kind of functions are called copula (density) functions.

Overview

Modelling soccer games

- Directly Modelling soccer outcomes

- Modelling scores

Estimating parameters

- Maximum Likelihood

- Weighted Maximum Likelihood

Applications

- Prediction

- Ranking

Maximum Likelihood

- | When we know the values of the parameters $\theta_H(A, h, c, \dots)$ we can calculate the chances for the outcome of a game
- | but how do we get these parameters?
- | We can estimate them by using maximum likelihood!
- | We look for the optimal parameters using numerical optimization

Maximum likelihood Bradley Terry

Imagine we have collected the outcome of H; D; A of N past games, including T teams

$$L(\theta_1; \dots; \theta_T; d; h) = \prod_{i=1}^N P(O_i = o_i):$$

$$P(O_i = H) = \frac{h_{iH}}{h_{iH} + d_i p \frac{h_{iH}}{h_{iH} + i_A}};$$

$$P(O_i = D) = \frac{d_i p \frac{h_{iH}}{h_{iH} + i_A}}{h_{iH} + d_i p \frac{h_{iH}}{h_{iH} + i_A}};$$

$$P(O_i = A) = \frac{i_A}{h_{iH} + d_i p \frac{h_{iH}}{h_{iH} + i_A}};$$

Maximum likelihood Poisson

Let now x_{iH} en x_{iA} be the number of goals scored by the home and away team.

For example, if the result of the game is 1-2, $x_{iH} = 1$ and $x_{iA} = 2$.

$$L(x_1; \dots; x_T; c; h) = \prod_{i=1}^T P(G_{iH} = x_{iH}) P(G_{iA} = x_{iA}):$$

$$P(G_{iH} = g) = \frac{c_{iH}^g}{g!} e^{-c_{iH}} \quad P(G_{iA} = g) = \frac{c_{iA}^g}{g!} e^{-c_{iA}}$$

$$c_{iH} = c + h \frac{iH}{iA}$$

$$c_{iA} = c - \frac{iA}{iH}$$

Time decay parameter

- | How many matches do we consider for our MLE-estimates?
- | Is every game equally important?

We introduce a weight such that more recent games have a higher weight in the maximum likelihood

$$w_{\text{time}}(t_i) = \frac{1}{2} \frac{t_i}{\text{Half Period}};$$

where t_i is the time (in days) that has passed since the game is played and 'Half Period' is a parameter whose value is to be chosen. (see later)

When we use the time weights, our likelihood formulas for the Bradley-Terry and the simple Poisson model change to

$$L = \prod_{i=1}^N P(O_i = o_i)^{w_{\text{time}}(t_i)}:$$

$$L = \prod_{i=1}^N (P(G_{iH} = x_{iH})P(G_{iA} = x_{iA}))^{w_{\text{time}}(t_i)}$$

Importance weights

We could also use weights to give important games more weights than less important games. For example, if we would consider the national teams, we could use the weights used by the FIFA:

- | World Cup: $w_{\text{importance}} = 4$
- | Confederations Championships (e.g. ECFA): $w_{\text{importance}} = 3$
- | Qualification matches for tournaments above: $w_{\text{importance}} = 2$
- | Friendlies and small tournaments: $w_{\text{importance}} = 1$

E.g. for Bradley-Terry the likelihood function changes to:

$$L = \prod_{i=1}^N P(O_i = o_i)^{w_{\text{time}}(t_i)w_{\text{importance}}}$$

Overview

Modelling soccer games

- Directly Modelling soccer outcomes

- Modelling scores

Estimating parameters

- Maximum Likelihood

- Weighted Maximum Likelihood

Applications

- Prediction

- Ranking

Prediction

- | Since our models can produce probabilities for (future) soccer games, it is straightforward that they can be used for prediction.
- | We repeatedly estimate our parameters on a number of games in a certain time period and we then predict the outcome of the next game.
- | We measure if a model is good in predicting by using a loss function

loss functions

- | Rank Probability Score:

$$\frac{1}{N} \sum_{i=1}^N \frac{(P(O_i = H) - \mathbb{I}(o_i = H))^2 + (P(O_i = A) - \mathbb{I}(o_i = A))^2}{2}$$

The best model is the model with the lowest RPS.

Results Premier League

Table: Comparison table for the best performing models of each of the considered classes with respect to the RPS criterion. All of the second season half English Premier League matches in the period between the seasons 2000-2001 and 2016-2017 are considered.

Model Class	Parameters per team	Half Period best model	Lowest RPS
Independent Poisson,	1	HP = 200	0.1979
Bivariate Poisson, $c = \text{const.}$	1	HP = 180	0.1981
Bivariate Poisson, $c = \text{const.}$	2	HP = 240	0.1986
Independent Poisson	2	HP = 240	0.1988
Bradley-Terry	1	HP = 540	0.2020

Ranking

Another application of this modelling, is the making of rankings. Since we have estimated the strengths s_i for all teams, we can order them by (estimated) strength.

Example ranking national teams

	teams	points		team	points
1	Brazil	4.32	11	England	3.15
2	Spain	4.18	12	Uruguay	2.98
3	Germany	4.13	13	Denmark	2.97
4	Argentina	4.13	14	Croatia	2.95
5	France	3.53	15	Sweden	2.93
6	Belgium	3.50	16	Poland	2.92
7	Colombia	3.42	17	Italy	2.91
8	Chile	3.37	18	Peru	2.86
9	Netherlands	3.37	19	Ecuador	2.61
10	Portugal	3.36	20	Switzerland	2.60
131	Luxembourg	1.12	212	Vatican	0.29

