## A New Bayesian Rating System for Team Competitions

Sergey Nikolenko<sup>1,2</sup> Alexander Sirotkin<sup>1,3</sup>

<sup>1</sup>St. Petersburg Academic University
<sup>2</sup>Steklov Mathematical Institute, St. Petersburg
<sup>3</sup>St. Petersburg Institute for Informatics and Automation of the RAS

#### ICML 2011, June 30, 2011

Sergey Nikolenko, Alexander Sirotkin A New Bayesian Rating System for Team Competitions

### Outline

#### TrueSkill and its problems

- TrueSkill
- Motivation and TrueSkill problems

#### Problems and solutions

- Undersized teams
- Multiway ties and the new factor graph

- In probabilistic rating models, Bayesian inference aims to find a linear ordering on a certain set given noisy comparisons of relatively small subsets of this set.
- Useful whenever there is no way to compare a large number of entities directly, but only partial (noisy) comparisons are available.
- We will stick to the metaphor of matches and players.

### Introduction

- Elo rating system: first probabilistic rating model (chess: two players).
- Bradley–Terry models: assume that each player has a "true" rating γ<sub>i</sub>, and the win probability is proportional to this rating: γ<sub>1</sub> wins over γ<sub>2</sub> with probability <sup>γ<sub>1</sub></sup>/<sub>γ<sub>1</sub>+γ<sub>2</sub></sub>.
- Inference: fit this model to the data from matches played.
- Several extensions, but large matches are hard for Bradley–Terry models.

・ロット 全部 マート・ キャー

### Introduction

- TrueSkill was initially developed in MS Research for Xbox Live gaming servers [Graepel, Minka, Herbrich, 2007].
- Given results of team competitions, learn the ratings of players of these teams.
- Direct application matchmaking: find interesting opponents for a player or team.
- [Graepel et al., 2010]: AdPredictor. Predicts CTRs of advertisements based on a set of features: the features are a team, and the team wins whenever a user clicks the ad.

TrueSkill Motivation and TrueSkill problems

#### TrueSkill factor graph



Sergey Nikolenko, Alexander Sirotkin A New Bayesian Rating System for Team Competitions

### TrueSkill

- Layers of TrueSkill factor graph:
  - $s_{i,j}$  skill of player *i* from team *j*; normally distributed around  $\mu_{i,j}$  with variance  $\sigma_{i,j}$ , where  $(\mu_{i,j}, \sigma_{i,j})$  are prior ratings;
  - *p<sub>i,j</sub>* performance of player *i* from team *j* in this match; conditionally normally distributed around the skill *s<sub>i,j</sub>* with variance β (a global model parameter);
  - *t<sub>j</sub>* − performance of team *j*; in TrueSkill<sup>TM</sup>, team performance is the sum of player performances;
  - $d_j$  difference in performance between teams who took neighboring places in the tournament; a tie between two teams corresponds to  $|d_j| \le \varepsilon$  for some model parameter  $\varepsilon$ , and a win corresponds to  $d_j > \varepsilon$ .

《曰》 《部》 《글》 《글》 - 글

### TrueSkill

- There is no evidence per se, it is incorporated in the structure of the graph, we just have to marginalize by message passing.
- The marginalization problem is complicated by the step functions at the bottom; solved with Expectation Propagation [Minka, 2001]:
  - approximate messages from  $\mathbb{I}(d_i > \epsilon)$  and  $\mathbb{I}(|d_i| \le \epsilon)$  to  $d_i$  with normal distributions;
  - repeat message passing on the bottom layer of the graph until convergence.

(日) (同) (日) (日)

TrueSkill Motivation and TrueSkill problems

#### Example: a match of four players



4 ロ ト 4 合 ト 4 き ト 4 き ト き うへの A New Bayesian Rating System for Team Competitions

### Motivation

- We started with a practical problem: we tried to apply TrueSkill<sup>TM</sup> to a Russian game "What? Where? When?".
- $\bullet\,$  Teams of  $\leq 6$  players answer questions, whoever gets the most correct answers wins.
- It turned out that TrueSkill<sup>TM</sup> works poorly on this dataset because of its properties:
  - large multiway ties are common; it is common to have 30–40 different places (because there were 35-50 questions in total) in a tournament with a thousand teams;
  - teams vary in size (max 6 players, but often incomplete).
- Why is it bad for TrueSkill<sup>™</sup> and what do we do about it?

イロト 不得 とうき イロト

### Outline

#### TrueSkill and its problems

- TrueSkill
- Motivation and TrueSkill problems

#### Problems and solutions

- Undersized teams
- Multiway ties and the new factor graph

### Variable team size

- An undesired feature of TrueSkill<sup>TM</sup> is the assumption that a team's performance is the sum of player performances.
- In many competitions (and comparison problems), an undersized team stands a very good chance against a full one, and it would be an unfair boost for the smaller team.
- To alleviate the team performance formula problem, we simply select a different function.
- We can very easily use any affine function, e.g., average (but it would be unfair for smaller teams now).

#### Variable team size

- Moreover, there is a simple way to approximate nonlinear functions: replace player performances with their estimates provided by the prior ratings µ<sub>i</sub>.
- For instance, to approximate a team performance function

$$t=p_1^2+p_2^2+\ldots+p_n^2$$

we replace it with

$$t = \mu_1 p_1 + \mu_2 p_2 + \ldots + \mu_n p_n$$

(here  $p_i$  are model variables, and  $\mu_i$  are constants fixed before inference and equal to prior ratings).

#### Variable team size

- I don't know a universally good team performance function, I can only encourage you to try different ones.
- In the end, for our dataset the function that worked best was (assuming m<sub>i,j</sub>'s are sorted)

$$t_{i} = \begin{cases} \frac{\sum\limits_{j=1}^{n_{i}} p_{i,j}}{n_{i}} \cdot (0.88 + 0.02n_{i}), & \text{if } n_{i} \leq 6, \\ \sum\limits_{j=1}^{n_{i}} p_{i,j} \cdot \frac{\sum\limits_{j=1}^{6} \mu_{i,j}}{6\sum\limits_{j=1}^{n_{i}} \mu_{i,j}}, & \text{if } n_{i} > 6. \end{cases}$$

where  $n_i$  is number of players in team *i*.

• Obviously, it wouldn't work for other applications.

### Multiway ties

- Large multiway ties are deadly for TrueSkill<sup>TM</sup>. Consider four teams in a tournament with performances  $p_1, \ldots, p_4$ .
- Team 1 has won, while teams 2–4, listed in this order, drew behind the first.
- Then the factor graph tells us that

$$t_2 < t_1 - \varepsilon, \quad |t_2 - t_3| \leq \varepsilon, \quad |t_3 - t_4| \leq \varepsilon.$$

• Team 3's performance t<sub>3</sub> may actually nearly equal t<sub>1</sub>, and t<sub>4</sub> may exceed t<sub>1</sub>!

• Moreover, these boundary cases are realized in practice when unexpected results occur.

$$t_2 < t_1 - \varepsilon, \quad |t_2 - t_3| \leq \varepsilon, \quad |t_3 - t_4| \leq \varepsilon.$$

- Suppose the winning team  $t_1$  was an underdog, and its prior distribution fell behind the priors of  $t_2$ ,  $t_3$ , and  $t_4$ ,  $t_4$  being the prior leader of all four.
- Then the maximum likelihood value of  $t_4$  is likely to exceed  $t_1$ .

・ロット 全部 マート・ キャー

### Changes in the factor graph

- For the multiway tie problem, we add another layer in the factor graph, namely the layer of *place performances l<sub>i</sub>*.
- Each team performs in the  $\epsilon$ -neighborhood of its place performance, and place performances relate to each other with strict inequalities like  $l_2 < l_1 2\epsilon$ .
- Then it's inference as usual. We have not experienced any slowdown in convergence.

Undersized teams Multiway ties and the new factor graph

#### New factor graph



Sergey Nikolenko, Alexander Sirotkin A New Bayesian Rating System for Team Competitions

Undersized teams Multiway ties and the new factor graph

#### Experimental results



### Thank you!

# Thank you for your attention!

Sergey Nikolenko, Alexander Sirotkin A New Bayesian Rating System for Team Competitions