

# Lecture 10: Multiagent Evaluation

Jun Wang  
UCL

# Table of Content

- ▶ Elo rating system
- ▶ Glicko
- ▶ TrueSkill
- ▶ Multidimensional Elo (mElo2k)
- ▶  $\alpha$ -rank

# Importance of evaluation

- ▶ The third wave of AI was driven by *ImageNet* benchmark in Computer Vision
- ▶ In Natural Language Processing we have *SuperGLUE Benchmark*
- ▶ In Reinforcement learning, we have *atari* and *OpenAI Gym*
- ▶ How about multiagent learning?
  - ▶ What would be the tasks?
  - ▶ Difficult as require evaluation against another agent

# The Elo Rating System<sup>1</sup>

- ▶ Elo rating assigns each player a numerical rating score representing their ability from historical plays
  - ▶ developed by Arpad Elo for chess competition
  - ▶ widely used for many games such as chess, football, tennis, and video games
- ▶ It is computed based on the aggregated wins and losses weighted by how likely the win or loss was based on the opponent's ability
- ▶ Elo evaluation system has various purposes:
  - ▶ match players against each other when players have similar rating score
  - ▶ use skill rating to have cut-offs for different skill tiered tournaments

---

<sup>1</sup>Arpad E Elo. *The rating of chessplayers, past and present*. Arco Pub., 1978.

# The Elo Rating System

- ▶ Elo assigns a rating  $r$  to each player
- ▶ According to rating  $r$ , the predicted winning probability is:

$$\hat{p}_{ij} = \sigma(\alpha r_i - \alpha r_j) = \frac{1}{1 + e^{-\alpha(r_i - r_j)}}$$

- ▶ Minimizing cross entropy of the predicted  $\hat{p}$  and the true  $p$  gives

$$L_{Elo}(p_{ij}, \hat{p}_{ij}) = -p_{ij} \log(\hat{p}_{ij}) - (1 - p_{ij}) \log(1 - \hat{p}_{ij})$$

- ▶ *stochastic gradient descent* up finishing a game at time  $t$ :

$$r_i^{t+1} \leftarrow r_i^t + \eta(S_{ij}^t - \hat{p}_{ij}^t) = r_i^t + \eta(S_{ij}^t - \hat{p}_{ij}^t)$$

where  $S_{ij}^t$  is a binary value, i winning j gives 1 otherwise 0

- ▶ Elo ratings are stationary at:

$$\sum_j \hat{p}_{ij} = \sum_j \bar{p}_{ij} = \sum_j \frac{S_{ij}}{N_{ij}}$$

# Explanation of super-parameters in ELO

- ▶  $\alpha$  is an additional parameter that represent the scaling of the skill rating being used.

$$\alpha = \frac{\ln(10)}{400}, \quad \hat{p}_{ij} = \frac{1}{1 + e^{-\alpha(r_i - r_j)}} = \frac{1}{1 + 10^{-(r_i - r_j)/400}}$$

For example, if we want to  $\hat{p}_{ij} = 0.8$  then  $r_i - r_j \approx 240$ .

- ▶  $\eta$  is a constant, frequently 32 in chess, used to control the change magnitude of ratings at every step.

$$r_i^{t+1} \leftarrow r_i^t + \eta(S_{ij}^t - \hat{p}_{ij}^t) = r_i^t + \eta(S_{ij}^t - \hat{p}_{ij}^t)$$

- ▶ In addition, according to the above formula, the two players participating in a competition will increase or decrease the same score.

# Glicko<sup>2</sup>

- ▶ Glicko extends Elo, considering the rating reliability
  - ▶ Suppose two players both rated 1700
  - ▶ The first defeating the second. In ELO the first player gains 16 rating points; the second loses 16 rating points
  - ▶ The first player last played at several years ago, but the second plays every week. So the first player's 1700 rating score is less reliable than the second. Intuitively,
  - ▶ The first player should add more than 16 points, because his rating is not believable and he defeated a player with a fairly precise rating of 1700.
  - ▶ The second player should decrease less than 16 points because his rating is already precisely measured to be near 1700, and he loses to a player with unreachable rating.
- ▶ Glicko adds ratings deviation(RD) to measure the uncertainty of a player's rating.

---

<sup>2</sup>Mark E Glickman. "The glicko system". In: (1995).

# Glicko Formulation

- ▶ Rating
  - ▶ Rating  $r$  measures the mean skill of each player and is updated from game outcomes.
- ▶ RD(Rating Deviation)
  - ▶ High RD correspond to unreliable rating. Higher RD indicates a player has only participated small number of games, or not play frequently
  - ▶ Changes from game outcomes and from the passage of time when not playing. Playing games always decreases player's RD, and time passing without competing always increases player's RD.
- ▶ Rating changes in Glicko system are not balanced as they usually are in the Elo system.
- ▶ Player's strength with 95 % confidence in interval  $[r - 2 * RD, r + 2 * RD]$ .



# Glicko Rating Algorithm

A rating period consists of many pairwise competitions, such as tournaments.

Step 1: at beginning of a rating period:

- ▶ New player: initialize rating = 1500, RD = 350.
- ▶ Old player:  $RD = \min(\sqrt{RD_{old}^2 + c^2 t}, 350)$ , where  $c$  is a constant and  $t$  represents the number of rating periods without playing the game.

Step 2: at end of a rating period, updating ratings  $r$  and RD's according to game results

# Updating ratings and RDs

- ▶ Suppose a player has competed with  $m$  opponents in this rating period, with ratings  $r_1, r_2, \dots, r_m$ , and rating deviations  $RD_1, RD_2, \dots, RD_m$ .
- ▶ Let  $s_1, s_2, \dots, s_m$  are the results of  $m$  competitions, with value 0, 0.5, or 1 for a loss, draw and win.
- ▶ Updating rating and RD:

$$r' = r + \frac{q}{1/RD^2 + 1/d^2} \sum_{j=1}^m g(RD_j)(s_j - E(s|r, r_j, RD_j))$$

$$RD' = \sqrt{\left(\frac{1}{RD^2} + \frac{1}{d^2}\right)^{-1}}, \text{ where}$$

$$q = \frac{\ln 10}{400}, \quad g(RD) = \frac{1}{1 + 3q^2(RD^2)/\pi^2}$$

$$E(s|r, r_j, RD_j) = \frac{1}{1 + 10^{-g(RD_j)(r-r_j)/400}}$$

$$d^2 = \left( q^2 \sum_{j=1}^m g(RD_j)^2 E(s|r, r_j, RD_j)(1 - E(s|r, r_j, RD_j)) \right)^{-1}$$

# TrueSkill<sup>3</sup>

- ▶ Use Gaussian to model player's skill and performance in game which has mean and variance, not just a skill as Elo.

Suppose each player  $i$  has skill  $s_i \sim \mathcal{N}(s_i; \mu_i, \sigma_i)$ , performance  $p_i \sim \mathcal{N}(p_i; s_i, \beta)$  with a fixed variance  $\beta$ . The probability of  $i$  beating can be denoted as

$$P(p_1 > p_2 | s_1, s_2) > \Phi\left(\frac{s_1 - s_2}{\sqrt{2\beta}}\right),$$

where  $\Phi$  denotes the cumulative density of standard normal distribution.

- ▶ It can deal with any number of competing entities and can infer individual skills from team results

---

<sup>3</sup>Ralf Herbrich, Tom Minka, and Thore Graepel. "TrueSkill™: a Bayesian skill rating system". In: *Proceedings of the 19th international conference on neural information processing systems*. 2006, pp. 569–576.

# Transitive and non-transitive games

- ▶ Elo rating assumes that relative skill is transitive
- ▶ In Rock-paper-scissors game, the maxent Nash equilibrium<sup>4</sup> is  $\mathbf{p}_A^* = (1/3, 1/3, 1/3)^T$  and the **Nash average** is  $\mathbf{n}_A := \mathbf{A} \cdot \mathbf{p}_A^* = (0, 0, 0)^T$ .

$\mathbf{A}$	R	P	S
R	0	-1	1
P	1	0	-1
S	-1	1	0

**Table:** Intransitive case, thus anti-symmetric matrix

- ▶ It's easy to check that the unique maxent Nash for  $\mathbf{A}'$  is  $\mathbf{p}_A^* = (1, 0, 0)^T$ , Nash average is  $\mathbf{n}_A = (0, -1, -2)^T$

$\mathbf{A}'$	a	b	c
a	0	1	2
b	-1	0	1
c	-2	-1	0

**Table:** Transitive case

---

<sup>4</sup>with greater entropy than any other Nash equilibrium

# Why Elo can only solve transitive problems?

- ▶ Elo rating cannot solve intransitive problems.
- ▶ For example, rock-paper-scissors (e.g. paper beats rock with  $p = 1$ , but the predicted  $\hat{p} = 0.5$ )

**Hodge decomposition**<sup>5</sup>: any antisymmetric matrix decomposes as

$$\mathbf{A} = \text{transitive component} + \text{cyclic component} = \text{grad}(\mathbf{r}) + \text{rot}(\mathbf{A})$$

- ▶ The divergence of  $\mathbf{A}$  is  
 $\mathbf{r} = \text{div}(\mathbf{A}) = \frac{1}{n} \mathbf{A} \cdot \mathbf{1}$ .  $\text{grad}(\mathbf{r}) = \mathbf{r} \mathbf{1}^T - \mathbf{1} \mathbf{r}^T$
- ▶ The rotation of  $\mathbf{A}$  is  
 $\text{rot}(\mathbf{A})_{ij} = \frac{1}{n} \sum_{k=1}^n \text{curl}(\mathbf{A})_{ijk}$   $\text{curl}(\mathbf{A})_{ijk} = \mathbf{A}_{ij} + \mathbf{A}_{jk} - \mathbf{A}_{ik}$
- ▶ There is an Elo rating that generates probabilities  $\mathbf{P}$  iff  
 $\text{curl}(\text{logit } \mathbf{P}) = 0$

---

<sup>5</sup>David Balduzzi Karl Tuyls Julien Perolat Thore Graepel.

“Re-evaluating Evaluation”. In: (NIPS 2018); X. Jiang. “Statistical ranking and combinatorial Hodge theory”. In: (2011).

# Multidimensional Elo( $m\text{Elo}_{2k}$ )<sup>6</sup>

$m\text{Elo}_{2k}$  can solve intransitive problems because it considers the low-rank( $2k$ ) approximation of the cyclic component of the matrix  $\text{logit}(\mathbf{P})$ .

- ▶  $\mathbf{C}^T \mathbf{\Omega} \mathbf{C}$  is an approximation to the **Schur decomposition** of the cyclic component  $\tilde{\mathbf{A}}$ .

$$\mathbf{A}_{n \times n} = \text{grad}(\mathbf{r}) + \tilde{\mathbf{A}} \approx \text{grad}(\mathbf{r}) + \mathbf{C}^T \begin{pmatrix} 0 & 1 & & \\ -1 & 0 & & \\ & & \ddots & \end{pmatrix} \mathbf{C} =: \text{grad}(\mathbf{r}) + \mathbf{C}_{n \times 2k}^T \mathbf{\Omega}_{2k \times 2k} \mathbf{C}_{2k \times n}$$

- ▶ Let  $m\text{Elo}_{2k}$  assign each player Elo rating  $r_i$  and  $2k$ -dimensional vector  $\mathbf{c}_i$ .

$$\mathbf{mElo}_{2k}: \hat{p}_{ij} = \sigma\left(r_i - r_j + \mathbf{c}_i^T \cdot \mathbf{\Omega}_{2k \times 2k} \cdot \mathbf{c}_j\right) \text{ where } \mathbf{\Omega}_{2k \times 2k} = \sum_{i=1}^k (\mathbf{e}_{2i-1} \mathbf{e}_{2i}^T - \mathbf{e}_{2i} \mathbf{e}_{2i-1}^T).$$

---

<sup>6</sup>David Balduzzi Karl Tuyls Julien Perolat Thore Graepel.  
“Re-evaluating Evaluation”. In: (NIPS 2018).

## $\alpha$ -rank<sup>7</sup>

$\alpha$ -rank provides an alternative solution concept for multi-player normal form game

		II		
		L	C	R
I	U	2, 1	1, 2	0, 0
	M	1, 2	2, 1	1, 0
	D	0, 0	0, 1	2, 2

- ▶ In this case, there are two players, player 1's strategy set is U, M, D, player 2's strategy set is L, C, R, and there are  $3 \times 3 = 9$  strategy profiles (U, L), (U, C), (U, R), ....

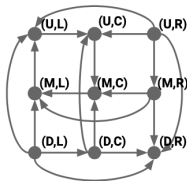
---

<sup>7</sup>Shayegan Omidshafiei et al. “ $\alpha$ -rank: Multi-agent evaluation by evolution”. In: *Scientific reports* 9.1 (2019), pp. 1–29.

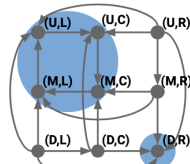
$\alpha$ -rank

		II		
		L	C	R
I	U	2, 1	1, 2	0, 0
	M	1, 2	2, 1	1, 0
	D	0, 0	0, 1	2, 2

(a)



(b)



(c)

- ▶ According to the payoff table M in Fig(a), construct the Response graph
- ▶ In response graph nodes correspond to strategy profiles, and directed edges if the deviating player's new strategy is a better-response.<sup>8</sup> For example in Fig(b), from (U,L) to (U,C) player 2's payoff increases from 1 to 2.
- ▶ The sink strongly connected components (SSCC) of the response graph are illustrated in Fig(c) in blue. Desirable strategy profiles are all in SSCC, because the edges always point to the increasing direction.

<sup>8</sup>Karl Tuyls Mark Rowland Shayegan Omidshafiei. “Multiagent Evaluation under Incomplete Information”. In: (NIPS 2019).



# $\alpha$ -rank

- In order to be able to find SSCC and its stationary distribution-Rank creates a so-called, Markov-Conley chain, where the edges are “soft”.

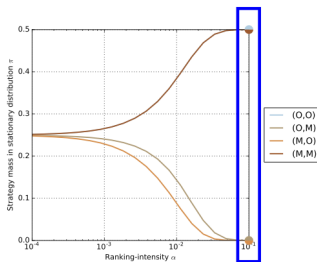
$$\mathbf{C}_{s,\sigma} = \begin{cases} \eta \frac{1 - \exp(-\alpha(\mathbf{M}^k(\sigma) - \mathbf{M}^k(s)))}{1 - \exp(-\alpha m(\mathbf{M}^k(\sigma) - \mathbf{M}^k(s)))} & \text{if } \mathbf{M}^k(\sigma) \neq \mathbf{M}^k(s) \\ \frac{\eta}{m} & \text{otherwise,} \end{cases} \quad \text{and} \quad \mathbf{C}_{s,s} = 1 - \sum_{\substack{k \in [K] \\ \sigma | \sigma^k \in S^k \setminus \{s^k\}}} \mathbf{C}_{s,\sigma},$$

- The stationary distribution computed by  $\alpha$ -Rank corresponds to a ranking of strategy profiles:

$$\pi = \mathbf{C}^T \pi$$

# $\alpha$ -rank results - Battle of the Sexes

- As  $\alpha$  increases, two sink chain components emerge corresponding to strategy profiles  $(O, O)$  and  $(M, M)$ .



(b) Ranking-intensity sweep.

Agent	Rank	Score
$(O, O)$	1	0.5
$(M, M)$	1	0.5
$(O, M)$	2	0.0
$(M, O)$	2	0.0

(c)  $\alpha$ -Rank results.

	$O$	$M$
$O$	$(3, 2)$	$(0, 0)$
$M$	$(0, 0)$	$(2, 3)$

# References I



Arpad E Elo. *The rating of chessplayers, past and present*. Arco Pub., 1978.



Mark E Glickman. “The glicko system”. In: (1995).



David Balduzzi Karl Tuyls Julien Perolat Thore Graepel. “Re-evaluating Evaluation”. In: (NIPS 2018).



Ralf Herbrich, Tom Minka, and Thore Graepel. “TrueSkill™: a Bayesian skill rating system”. In: *Proceedings of the 19th international conference on neural information processing systems*. 2006, pp. 569–576.

# References II



X. Jiang. “Statistical ranking and combinatorial Hodge theory”. In: (2011).



Karl Tuyls Mark Rowland Shayegan Omidshafiei. “Multiagent Evaluation under Incomplete Information”. In: (NIPS 2019).



Shayegan Omidshafiei et al. “ $\alpha$ -rank: Multi-agent evaluation by evolution”. In: *Scientific reports* 9.1 (2019), pp. 1–29.