

Différentes méthodes de classement

Les systèmes Elo et TrueSkill

v3.1 (30/05/2011)

Introduction De nombreux jeux et sports utilisent des classements pour déterminer le niveau des participants. Nous allons étudier 2 systèmes de classements, le système Elo et le système TrueSkill, en essayant de déterminer dans quelles conditions ces systèmes sont le plus efficace (c'est-à-dire dans quels cas le classement obtenu permet-il de bien pronostiquer le résultat de la rencontre entre 2 joueurs à partir de leur classement respectif).

1 Classement Elo

Présentation Le classement Elo est connu pour être utilisé pour les échecs et sert à classer les millions de joueurs pratiquants, des débutants aux grands maîtres internationaux. Chaque joueur possède un score (le score Elo). Plus ce score est élevé, plus le joueur a un haut niveau (par exemple, aux échecs, ce score varie de 1000 pour un joueur débutant à environ 2800 pour les meilleurs joueurs mondiaux). Ce score augmente après une victoire et diminue après une défaite, cette variation dépendant de la différence de niveau entre le joueur et son adversaire.

Calcul du score Elo Pour obtenir le nouveau score Elo d'un joueur après un match, on utilise la formule $S' = S + K \cdot (W - p(D))$, où S désigne l'ancien score du joueur et S' le nouveau. W désigne le résultat de la partie : 1 pour une victoire, 0,5 pour un match nul et 0 pour une défaite. $p(D)$ est le résultat attendu du match (voir ci-dessous). Enfin, K est un facteur permettant de dilater les valeurs obtenues.

Calcul de $p(D)$

- ▶ **Force relative** : Soit $p_{A,B}$ la probabilité pour un joueur A de battre un joueur B. On note alors $f(p_{A,B}) = \frac{p_{A,B}}{1-p_{A,B}}$ la force relative de A contre B. Ainsi, si $f(p_{A,B}) > 1$, A est censé être plus fort que B, et inversement. On peut donc ainsi calculer $f(p_{A,C})$ à partir de $f(p_{A,B})$ et $f(p_{B,C})$: $f(p_{A,C}) = f(p_{A,B}) \cdot f(p_{B,C})$, et donc $p_{A,C} = \frac{f(p_{A,C})}{1+f(p_{A,C})}$.
- ▶ **Passage à un système additif** : On utilise un logarithme pour passer à un système additif. On pose alors $\Delta(p_{A,B}) = 400 \cdot \log[f(p_{A,B})]$ (le facteur 400 utilisé dans le système Elo permet d'étendre la plage de valeurs). $\Delta(p_{A,B})$ représente alors l'écart entre les deux joueurs A et B dans le classement.

On définit alors la fonction $p(D)$ comme étant la fonction réciproque de $\Delta(p)$ et permettant ainsi de définir la probabilité de gain en fonction de la différence D entre les deux joueurs au classement. On a donc $\log \frac{p(D)}{1-p(D)} = \frac{D}{400}$ d'où $p(D) = \frac{1}{1+10^{\frac{D}{400}}}$ (le logarithme utilisé est le logarithme décimal).

Exemple Un joueur A ayant un classement de 1800 points Elo bat un joueur B ayant 1400 points Elo.

On a $D = 1800 - 1400 = 400$ donc $p(D) = \frac{1}{1+10^{\frac{400}{400}}} = 0,91$. On obtient alors, en prenant $K = 15$,

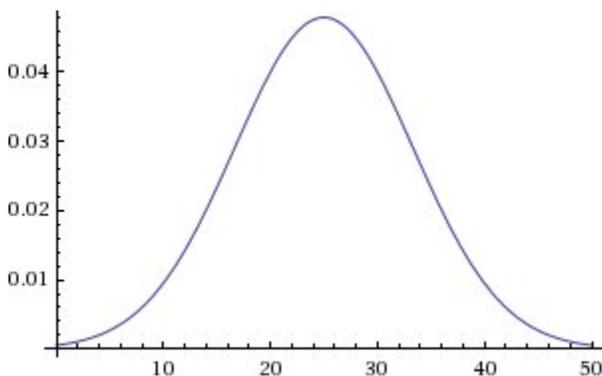
$$S' = S + K.(W - p(D)) = 1800 + 15.(1 - 0,91) = 1801,35.$$

Si A avait battu un joueur B' ayant 1900 points Elo, on aurait alors obtenu $D = 1800 - 1900 = -100$, $p(D) = \frac{1}{1+10^{\frac{100}{400}}} = 0,36$ et donc $S' = 1800 + 15.(1 - 0,36) = 1809,6$.

2 TrueSkill

Présentation Le système TrueSkill est utilisé par Microsoft afin d'évaluer le niveau des joueurs sur les différents jeux disponibles sur sa console de jeu, la Xbox 360. Il définit pour chaque joueur 2 paramètres : μ , représentant le niveau supposé du joueur, et σ , représentant l'incertitude de ce score. Alors que μ augmente après une victoire et diminue après une défaite, σ diminue après chaque partie. Le rang du joueur est alors représenté par $R = \mu - 3\sigma$. Les valeurs de départ de μ et σ pour un nouveau joueur sont $\mu = 25$ et $\sigma = \frac{25}{3}$. Pour un joueur débutant, on a donc $R=0$. C'est ce rang R que l'on utilise pour établir le classement.

Ce système permet de représenter le niveau du joueur par une distribution normale : $\mathcal{N}(\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$



$\mathcal{N}(\mu, \sigma^2)$ pour un joueur débutant ($\mu = 25$ et $\sigma = \frac{25}{3}$)

Calcul du nouveau score TrueSkill après un match Deux joueurs A et B s'affrontent. Leurs μ et σ valent respectivement (μ_A, σ_A) et (μ_B, σ_B) . On supposera que A bat B.

Quelques définitions

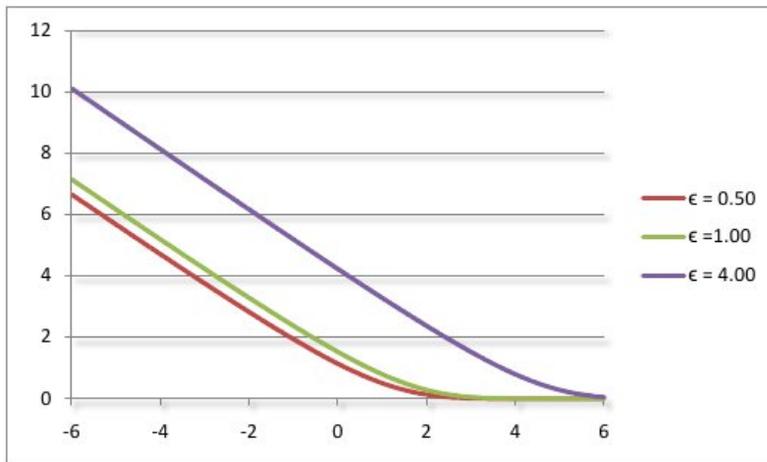
- ▶ β est un paramètre dépendant de chaque jeu. Plus β est faible, plus le jeu nécessite de dextérité. À l'opposé, un jeu reposant sur le hasard aura un β élevé. La valeur standard de β est 0,5.
- ▶ ε dépend de la probabilité que le match se termine par un match nul. Plus ε est grand, plus cette probabilité est grande.
- ▶ On détermine pour chaque rencontre une valeur c telle que $c^2 = 2\beta^2 + \sigma_A^2 + \sigma_B^2$.
- ▶ On pose $t = \frac{\mu_A - \mu_B}{c}$, en considérant toujours que A bat B. Ainsi, si A était le favori de ce match, on aura $t > 0$. À l'inverse, si B était favori, t sera négatif. La présence de c au dénominateur permet d'augmenter ou de diminuer la valeur de t , grâce à β (il est moins surprenant que le joueur le moins bien côté gagne si le jeu donne une plus grande part au hasard) et aux valeurs de σ (si les σ sont grands, cela signifie que les valeurs de μ sont peu représentatives du niveau réel des joueurs).

Modification de μ Les formules de mise à jour des valeurs de μ_A et μ_B sont

$$\mu_A \leftarrow \mu_A + \frac{\sigma_A^2}{c} \cdot v\left(t, \frac{\varepsilon}{c}\right)$$

$$\mu_B \leftarrow \mu_B - \frac{\sigma_B^2}{c} \cdot v\left(t, \frac{\varepsilon}{c}\right)$$

v est une fonction dépendant de t et de ε . Voici la représentation de $v(t, \varepsilon)$ en fonction de t pour quelques valeurs de ε .



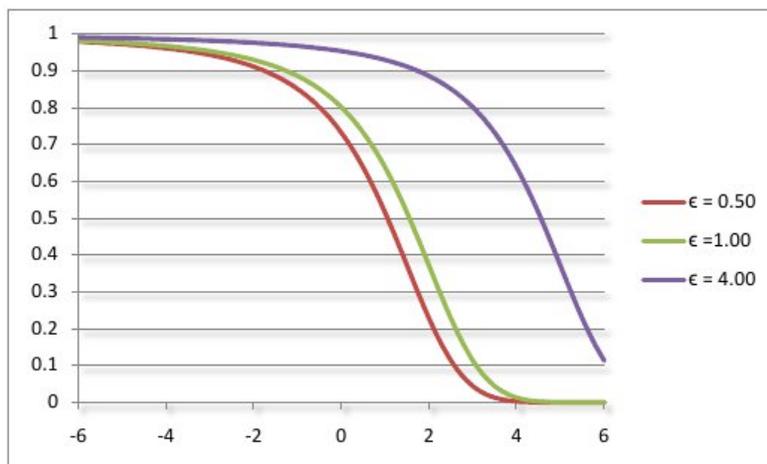
On peut donc remarquer que, vu que v est positive, μ baisse forcément après une défaite et augmente toujours après une victoire. De plus, plus le résultat est inattendu (c'est à dire t petit), plus la valeur de v en valeur absolue sera importante et les variations de μ seront donc plus importantes.

Modification de σ Les formules de mise à jour des valeurs de σ_A et σ_B sont

$$\sigma_A^2 \leftarrow \sigma_A^2 \cdot \left[1 - \frac{\sigma_A^2}{c^2} \cdot w\left(t, \frac{\varepsilon}{c}\right)\right]$$

$$\sigma_B^2 \leftarrow \sigma_B^2 \cdot \left[1 - \frac{\sigma_B^2}{c^2} \cdot w\left(t, \frac{\varepsilon}{c}\right)\right]$$

w est elle aussi une fonction dépendant de t et de ε . Voici la représentation de $w(t, \varepsilon)$ en fonction de t pour quelques valeurs de ε .



On a donc w comprise entre 0 et 1. Vu que $\frac{\sigma_A^2}{c^2} = \frac{\sigma_A^2}{2\beta^2 + \sigma_A^2 + \sigma_B^2} \in [0, 1]$, on aura donc $1 - \frac{\sigma_A^2}{c^2} \cdot w(t, \frac{\epsilon}{c}) \in [0, 1]$. La nouvelle valeur de σ sera donc comprise entre 0 et l'ancienne valeur de σ : on obtient une baisse de σ . Cette baisse sera plus importante si $\frac{\sigma_A^2}{c^2}$ est grand (c'est à dire quand le σ du joueur est encore grand) et si w est proche de 1 (cela correspond à t petit, ce qui signifie donc que le résultat obtenu est inattendu).

3 Comparaison des deux systèmes

Pour comparer les systèmes Elo et TrueSkill, nous allons les tester dans des situations où ils sont tous les deux utilisables : les matchs en un contre un. Pour chaque rencontre, nous déterminerons, à partir des classements obtenus, un favori et nous comparerons le résultat de notre pronostic au résultat du match.

Exemple 1 : championnat de France de volley-ball 2010/2011 Pour tester les systèmes Elo et TrueSkill, nous allons nous intéresser au championnat de France de volley-ball 2010/2011, qui regroupe 14 équipes s'affrontant en matchs aller/retour, ce qui représente un total de 182 matchs. Pour chaque match, on détermine le favori, qui sera notre pronostic et que l'on comparera au résultat du match. On calcule également les nouveaux scores Elo ou TrueSkill. Les tableaux ci-dessous donnent les classements obtenus avec ces systèmes, ainsi que le classement officiel du championnat, et le pourcentage de bons pronostics obtenus.

Elo		
1	Tours	1216,3
2	Cannes	1121,36
3	Poitiers	1098,87
4	Sète	1050,94
5	Rennes	1049,8
6	Nantes Rezé	1036,35
7	Paris	1020,72
8	Montpellier	1001,08
9	Tourcoing	999,45
10	Beauvais	993,66
11	Ajaccio	923,04
12	Toulouse	916,96
13	Nice	834,36
14	Saint-Quentin	737,04

Pourcentage de bons pronostics
63,75%

TrueSkill		
1	Rennes	24,74
2	Cannes	24,72
3	Tours	24,62
4	Nantes Rezé	24,53
5	Sète	24,35
6	Poitiers	24,33
7	Beauvais	24,2
8	Montpellier	23,85
9	Paris	23,45
10	Tourcoing	23,28
11	Toulouse	22,8
12	Ajaccio	22,56
13	Nice	22,3
14	Saint-Quentin	12,43

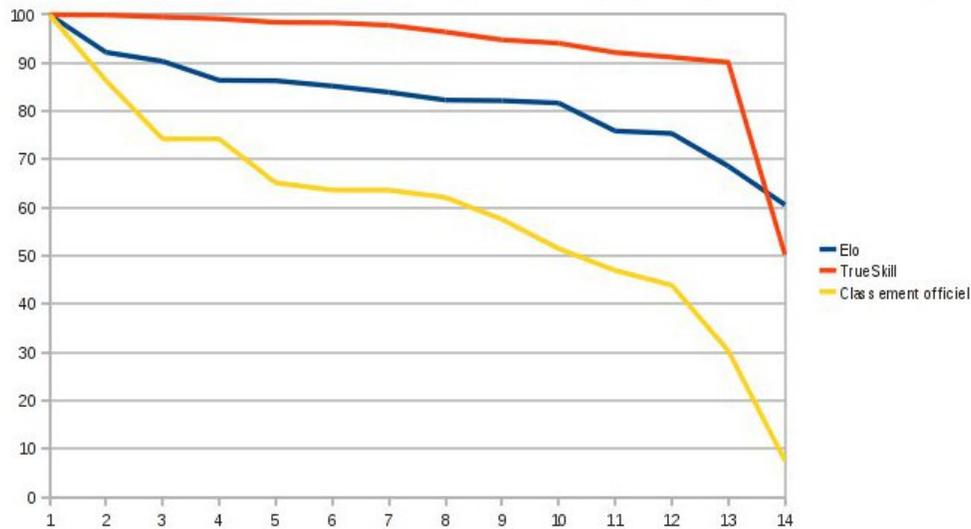
Pourcentage de bons pronostics
61,25%

Classement officiel		
1	Tours	66
2	Poitiers	57
3	Cannes	49
4	Sète	49
5	Rennes	43
6	Nantes Rezé	42
7	Paris	42
8	Montpellier	41
9	Tourcoing	38
10	Beauvais	34
11	Ajaccio	31
12	Toulouse	29
13	Nice	20
14	Saint-Quentin	5

Pourcentage de bons pronostics
61,67%

On peut ainsi remarquer que les deux systèmes obtiennent un taux de bons pronostics équivalents. On obtient également le même pourcentage de bons pronostics si l'on utilise le classement officiel. On peut donc se demander si ces deux systèmes sont utiles dans le cadre d'un championnat.

Le graphique ci-dessous représente les différents scores obtenus avec les 3 systèmes de classements, en fixant dans chaque cas le score de l'équipe en première position à 100 et en déduisant ensuite le score correspondant pour les autres équipes par proportionnalité.



On peut remarquer que les scores obtenus grâce au système TrueSkill sont peu dispersés (le score de l'avant-dernier correspond à 90% du score du premier, alors que pour le système Elo le score du dernier équivaut à un peu moins de 70% du score du premier). Il faut cependant étudier le cas de Saint-Quentin, classé dernier pour chaque classement. On remarque un très grand écart entre Saint-Quentin et l'équipe qui la précède, contrairement à l'Elo où l'écart est moindre.

Exemple 2 : championnat à 16 équipes avec plusieurs groupes de niveaux Pour ce deuxième exemple, nous allons simuler un championnat à 16 équipes, réparties en 4 groupes de niveaux comportant chacun 4 équipes : A, B, C et D, avec les équipes du groupe A étant sensées être les plus faibles et celles du groupe D les plus fortes. Les équipes sont représentées par le nom du groupe de niveau auquel elles appartiennent suivi d'un chiffre entre 1 et 4 (par exemple, les équipes du groupe A sont notées équipes A1, A2, A3 et A4).

Simulation des matchs À chaque groupe de niveau correspond un "coefficient", donné dans le tableau ci-dessous :

Groupe de niveaux	Coefficient
A	1
B	10
C	100
D	1000

Pour chaque match et pour chacune des 2 équipes participantes, on tire au hasard un réel entre 0 et la valeur du coefficient correspondant au groupe de niveau de l'équipe. L'équipe ayant obtenu le nombre le plus élevé est désignée gagnante du match.

Résultats Voici les classements obtenus avec les systèmes Elo et TrueSkill pour une simulation :

Elo		
1	D4	1147,04
2	D3	1129,89
3	D1	1127,42
4	D2	1123,61
5	C3	1065,65
6	C4	1046,61
7	C2	1045,66
8	C1	1024,19
9	B2	968,91
10	B4	952,5
11	B3	951,24
12	B1	942,78
13	A3	877,9
14	A1	873,82
15	A4	864,67
16	A2	858,01

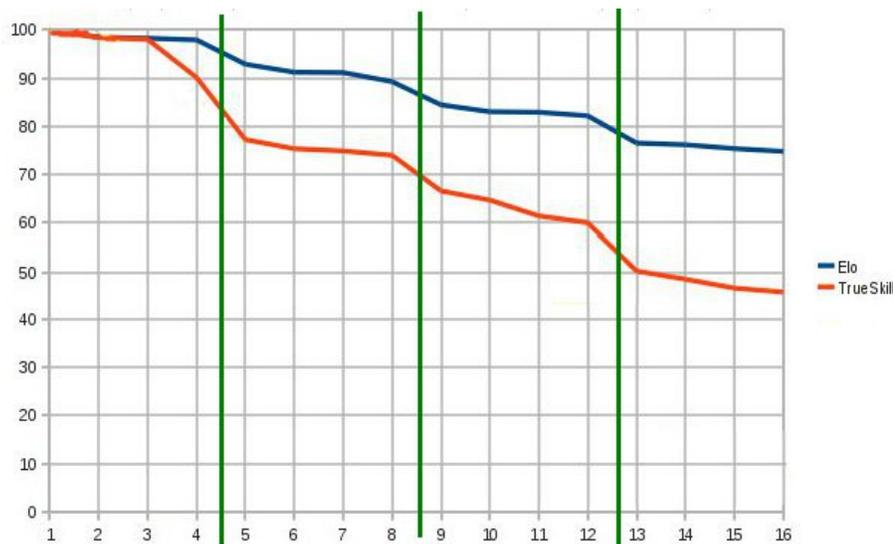
Pourcentage de bons pronostics
81,67%

TrueSkill		
1	D3	31,5
2	D2	31,01
3	D4	30,89
4	D1	28,4
5	C4	24,33
6	C1	23,74
7	C3	23,59
8	C2	23,31
9	B2	20,98
10	B1	20,38
11	B3	19,35
12	B4	18,91
13	A3	15,83
14	A4	15,3
15	A1	14,72
16	A2	14,46

Pourcentage de bons pronostics
80,42%

On constate ainsi pour les deux systèmes un taux de bons pronostics proche de 80%, soit environ 20% de plus que dans l'exemple précédent.

Le graphique ci-dessous représente les différents scores obtenus avec les 2 systèmes de classements, en fixant dans chaque cas le score de l'équipe en première position à 100 et en déduisant ensuite le score correspondant pour les autres équipes par proportionnalité.



(Les lignes verticales vertes marquent les séparations entre groupes de niveaux)

On constate que, contrairement à l'exemple précédent, c'est le TrueSkill qui "dilata" le plus les valeurs. Cela s'explique par le fait qu'il introduise un plus gros écart entre les groupes de niveaux (environ 10% en moyenne, contre 5% pour le système Elo).

Exemple 3 : comparaison Elo/TrueSkill à une plus grande échelle - Beta-test de Halo 2

Halo 2 est un FPS (jeu de tir à la première personne) édité par Microsoft et sorti en 2004. Ce jeu dispose d'un mode multijoueur où peuvent s'affronter deux joueurs ou deux équipes de 4 ou 8 joueurs. Les créateurs du système TrueSkill ont profité des phases de test de ce jeu pour comparer le TrueSkill à l'Elo. Voici les résultats qu'ils ont obtenus :

Mode de jeu	Elo	TrueSkill
1 contre 1	59,43%	69,17%
4 contre 4	57,45%	62,83%
8 contre 8	55,88%	70,06%

Pourcentage de bons pronostics pour les systèmes Elo et TrueSkill dans chacun des 3 modes de jeu multijoueurs de Halo 2.

On constate donc une nette domination du système TrueSkill face à l'Elo dans chaque cas de figure.

4 Conclusion

- Les systèmes Elo et TrueSkill semblent peu utiles pour les championnats sportifs : il apparaît en effet que ces systèmes n'apportent pas plus de précision qu'un système de classement plus basique, qui sera quand à lui plus facilement compréhensible par le grand public et modulable selon les souhaits de l'organisateur de la compétition (par exemple en donnant des points supplémentaires à une équipe gagnant avec un certain écart afin d'inciter une équipe ayant match gagné à continuer d'attaquer et ainsi améliorer l'aspect spectaculaire du jeu). De plus, ces championnats sportifs sont organisés en divisions successives, ce qui fait que le niveau des équipes dans une division est assez homogène, ce qui permet d'obtenir plus de résultats imprévus.
- Dans le cas de disciplines regroupant un très grand nombre de joueurs, il devient impossible d'organiser un championnat à cause du grand nombre de matchs nécessaires (pour que n équipes s'affrontent toutes entre elles une fois, il faut organiser $\frac{n(n-1)}{2}$ matchs) : il est donc nécessaire d'utiliser un système de type Elo ou TrueSkill. Dans ces conditions, il apparaît que le système TrueSkill semble être le plus performant, ce qui justifie son utilisation dans un grand nombre de jeux en ligne multijoueurs.

Remarque Il existe d'autres systèmes de classements, dont certains sont dérivés des systèmes Elo et TrueSkill. On peut ainsi citer en exemple le système Glicko, reprenant le système Elo en y ajoutant un paramètre permettant de moduler la variation du score d'un joueur en fonction de la fiabilité du score des joueurs participants, à la manière du paramètre σ utilisé dans le système TrueSkill. On peut également citer le MatchMaking Rating, utilisé notamment dans le jeu de stratégie Starcraft 2, basé sur le système TrueSkill mais ayant comme principale différence le fait que la valeur de σ d'un joueur peut augmenter si le résultat d'un match ne correspond pas au résultat attendu, le système remettant alors en doute le niveau du joueur, qu'il juge alors plus incertain.