

Quand les mathématiques se mettent au service du tennis

Paul-Emile Marcus // Candidat n°31462



Introduction

Classement ATP : l'exemple de Djokovic



	2018	2019	RÉSULTAT
POINTS	180	2000	+ 1820

Figure 1: Résultat Open d'Australie 2019



	2018	2019	RÉSULTAT
POINTS	2000	2000	+ 0

Figure 2: Résultat Wimbledon 2019

Force relative d'un joueur j sur un joueur i

Classement Elo

Calcul du classement Elo

Probabilité de victoire de A sur B :

$$p(D) = \frac{1}{1 + 10^{-\frac{D}{400}}}$$

Le classement Elo est calculé à l'aide de la formule suivante :

$$E_{n+1} = E_n + K \times [W - p(D)]$$

$W \in \{0, 1\}$ représente l'issue du match

Force relative d'un joueur j sur i

Quelques complexités

- Meilleur classement en carrière : $\mathcal{O}(n)$
- Ratio de victoires au cours de la saison/en carrière : $\mathcal{O}(n)$
- Face-à-face entre deux joueurs donnés : $\mathcal{O}(n^2)$

Force relative d'un joueur j sur un joueur i

Prise en compte de plusieurs paramètres

$$\mathbb{P}_{\text{ELO}} = \frac{1}{1 + 10^{-\frac{D_{\text{ELO}}}{400}}}$$

$$\mathbb{P}_{\text{H2H}} = \frac{1}{1 + 10^{-\frac{D_{\text{H2H}}}{10}}}$$

$$\mathbb{P}_{\text{RATIO}} = \frac{1}{1 + 10^{-\frac{D_{\text{RATIO}}}{0.4}}}$$

$$\Rightarrow \mathbb{P}(j/i) = \frac{1}{x_{\text{ELO}} + x_{\text{H2H}} + x_{\text{RATIO}}} \cdot \sum_{\text{paramètres}} x_{\text{param}} \cdot \mathbb{P}_{\text{param}}$$

Force relative d'un joueur j sur un joueur i

Calcul

```
422 ### Cette fonction calcule la proba de victoire d'un joueur j sur un joueur i en se basant sur la formule du  
423 classement élo adaptée à plusieurs statistiques  
424  
425 def calcul_proba(player1, player2):  
426     stats1, stats2 = stats_joueurs_pour_proba_final[player1], stats_joueurs_pour_proba_final[player2]  
427  
428     # Calcul proba classement elo  
429     p1_elo, p2_elo = calcul_proba_elo(stats1[3], stats2[3], 400)  
430  
431     # Calcul proba ratio saison  
432     p1_rs, p2_rs = calcul_proba_elo(stats1[5], stats2[5], 0.4)  
433  
434     # Calcul proba h2h  
435     h2h_tot = stats1[7][stats2[0]] + stats2[7][stats1[0]]  
436     if h2h_tot >= 5:  
437         p1_h2h, p2_h2h = calcul_proba_elo(stats1[7][stats2[0]], stats2[7][stats1[0]], 10)  
438  
439     # Calcul proba finale  
440     p = (50*p1_elo + 25*p1_rs + 25*p1_h2h)/100  
441     else:  
442         p = (70*p1_elo + 30*p1_rs)/100  
443  
444     return (p, 1-p)
```

Figure 3: Calcul de la force relative de j sur i

Matrice de transition

$$\begin{array}{c} \text{Djokovic} \\ \text{Federer} \\ \text{Nadal} \end{array} \begin{pmatrix} & \text{Djokovic} & \text{Federer} & \text{Nadal} & & & \\ \text{Djokovic} & 0 & 0.32 & 0.36 & \dots & * & \dots & * \\ \text{Federer} & 0.68 & 0 & 0.57 & \dots & * & \dots & * \\ \text{Nadal} & 0.64 & 0.43 & 0 & \dots & * & \dots & * \\ & \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ & * & * & * & \dots & * & \dots & * \\ & \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ & * & * & * & \dots & * & \dots & 0 \end{pmatrix}$$

Figure 4: Matrice de transition

Théorème de Perron-Frobenius

Théorème de Perron-Frobenius

Soient $A \in \mathcal{M}_n(\mathbb{R})$ positive et primitive et $r = \rho(A)$ son rayon spectral, alors :

- $r > 0$ est une valeur propre dominante et simple.
- $\exists ! x_0 > 0$ / $\|x_0\|_1 = 1$ et $Ax_0 = rx_0$.

Théorème de Perron-Frobenius

Soit $Y_0 \in \mathcal{M}_{n,1}(\mathbb{R}^+)$, tel que $\|Y_0\|_1 = 1$. Alors, $A^n Y_0 \xrightarrow{n \rightarrow +\infty} Y_\infty$, avec $Y_\infty \in \mathcal{M}_{n,1}(\mathbb{R}_+^*)$, $AY_\infty = rY_\infty$, et $\|Y_\infty\|_1 = 1$.

L'algorithme de Perron-Frobenius

```
473
474 # Initialisation du vecteur
475 Y = np.array([[1/n] for _ in range(n)])
476
477 # Itération jusqu'à convergence
478 for i in range(1000):
479     tmp = np.dot(matrice_transition, Y)
480     Y = tmp / np.linalg.norm(tmp)
481     if i == 9 or i==99 or i==999:
482         print ([Y[0],Y[1],Y[2]])
483
484 Y_normalized = Y / np.sum(Y)
485
486 # On trie les indices des états en fonction des valeurs du vecteur propre
487 classement_indices = np.argsort(Y_normalized, axis=0)[::-1]
488
489 # Création du classement
490 classement = [(indice, Y_normalized[indice][0]) for indice in classement_indices]
491
492 C = []
493 for i in range(len(classement)):
494     C.append(dict_index_nom[classement[i][0][0]])
495
496 return(C)
497
```

Figure 5: Algorithme de Perron Frobenius

Algorithme de Perron Frobenius

Quelques itérations

$$AY_0 = \begin{pmatrix} 0.0235 \\ 0.0252 \\ 0.0258 \\ \vdots \\ * \end{pmatrix}$$

$$A^{10}Y_0 = \begin{pmatrix} 0.0246 \\ 0.0250 \\ 0.0260 \\ \vdots \\ * \end{pmatrix}$$

$$A^{100}Y_0 = \begin{pmatrix} 0.0246 \\ 0.0250 \\ 0.0260 \\ \vdots \\ * \end{pmatrix}$$

Table 1: Mon classement

1	Djokovic N.
2	Alcaraz C.
3	Medvedev D.
4	Sinner J
5	Zverev A.
6	Rublev A.
7	Hurkacz H.
8	Dimitrov G.
9	Fritz T.
10	Rune H.

Table 2: Classement ATP

1	Djokovic N.
2	Alcaraz C.
3	Medvedev D.
4	Sinner J
5	Rublev A.
6	Tsitsipas S.
7	Rune H.
8	Ruud C.
9	Zverev A.
10	Fritz T.

Variations de classement

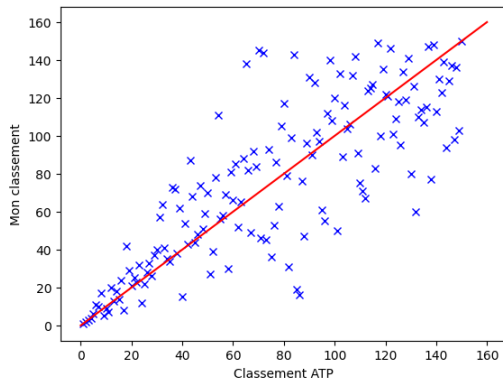


Figure 6: Nouveau classement en fonction du classement ATP

Écart de classement

- Moyenne : 7,79 et Écart-type : 32,4

Prédiction d'un match

Un problème binaire

Classement J1 (ATP)	Classement J2 (ATP)	1er tour	2eme tour	...
14	27	1	0	

**Si le mieux
classé gagne**

Vainqueur
1

**Si le mieux
classé perd**

Vainqueur
0

Distance de Levenshtein

```
315 def distance_levenshtein(mot1, mot2):
316     # On initialise une matrice de taille (longueur_mot1 + 1) x (longueur_mot2 + 1)
317     dist = [[0] * (len(mot2) + 1) for _ in range(len(mot1) + 1)]
318
319     # On initialise la première colonne avec les distances de 0 à la longueur du mot1
320     for i in range(len(mot1) + 1):
321         dist[i][0] = i
322
323     # On initialise la première ligne avec les distances de 0 à la longueur du mot2
324     for j in range(len(mot2) + 1):
325         dist[0][j] = j
326
327     # On calcule la distance de Levenshtein
328     for i in range(1, len(mot1) + 1):
329         for j in range(1, len(mot2) + 1):
330             cout_substitution = 0 if mot1[i - 1] == mot2[j - 1] else 1
331             dist[i][j] = min(dist[i - 1][j] + 1,                # effacement du nouveau caractère de mot1
332                             dist[i][j - 1] + 1,              # insertion dans mot2 du nouveau caractère de mot1
333                             dist[i - 1][j - 1] + cout_substitution) #substitution
334
335     # La distance de Levenshtein est le dernier élément de la matrice
336     return dist[len(mot1)][len(mot2)]
337
338 def trouver_correspondance_nom(nom_inexact, liste_noms):
339     correspondance_score_min = float('inf')
340     nom_correspondant = None
341     for nom_exact in liste_noms:
342         distance = distance_levenshtein(nom_inexact, nom_exact)
343         if distance < correspondance_score_min:
344             correspondance_score_min = distance
345             nom_correspondant = nom_exact
346     return nom_correspondant
```

Figure 7: Implémentation de la distance de Levenshtein

Random Forest

Localité	Superficie	Chambres	Jardin	Parking	Salles de bain
Le Mans	70	2	non	non	1
Sigean	128	4	oui	oui	2
Bayonne	256	5	oui	oui	3
Saumur	28	1	non	oui	1

Figure 8: Division aléatoire du jeu de données

Random Forest

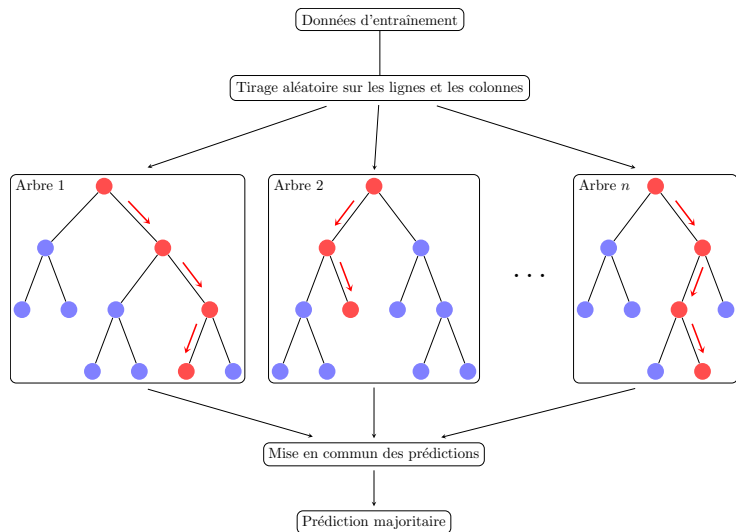


Figure 9: Explication du Random Forest Classifier

Score (précision) du modèle entraîné avec mon classement et ses meilleurs hyperparamètres :

0.68

Score (précision) du modèle entraîné avec le classement ATP et ses meilleurs hyperparamètres :

0.66

Remarque

$$\text{précision} = \frac{\text{nombre de bonnes prédictions}}{\text{nombre total de prédictions}}$$

Evaluation des modèles

Matrices de confusion

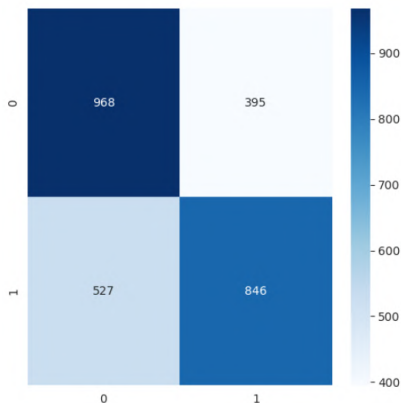


Figure 10: Matrice de confusion classement ATP

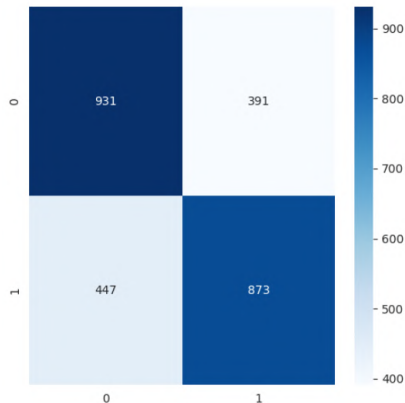


Figure 11: Matrice de confusion non classement

Evaluation des modèles

"Lucidité" des modèles

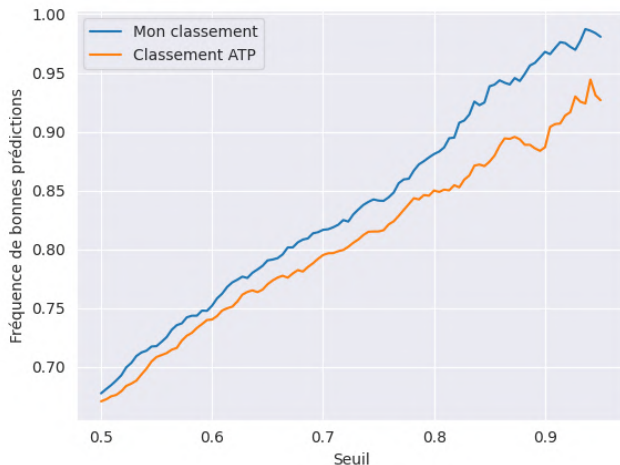


Figure 12: Fréquence de bonnes prédictions en fonction du seuil de probabilité

Le travail de mon binôme

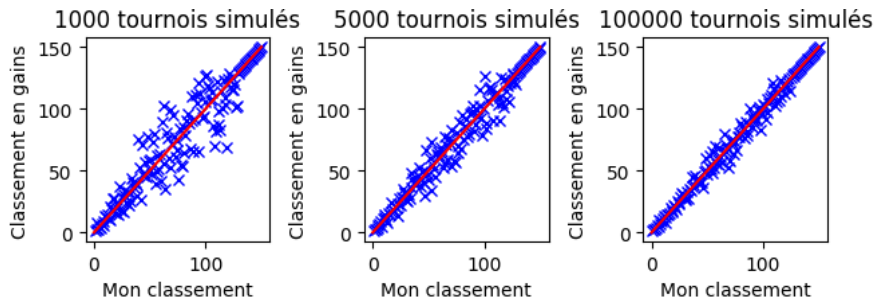


Figure 13: Classement après les tournois en fonction de notre classement

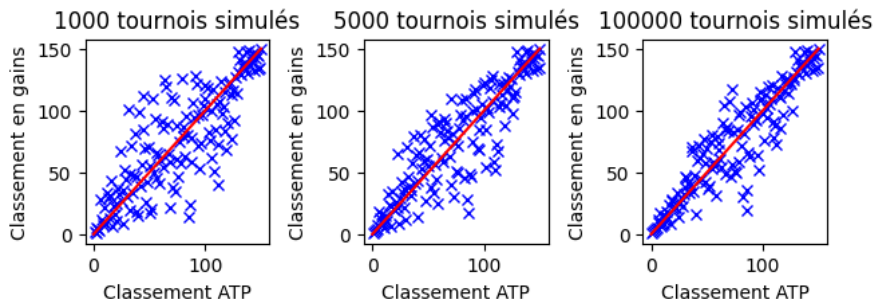


Figure 14: Classement après les tournois en fonction du classement ATP

Conclusion



Soit $P \in \mathcal{GL}_n(\mathbb{C})$ telle que

$$P^{-1}AP = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & \lambda_1 I_{m_1} + N_1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & \lambda_s I_{m_s} + N_s \end{pmatrix}$$

où $\text{Sp}_{\mathbb{C}}(A) = \{1, \lambda_1, \dots, \lambda_s\}$ et m_1, \dots, m_s les multiplicités respectives de $\lambda_1, \dots, \lambda_s$.

$$A^k = P \cdot \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & (\lambda_1 I_{m_1} + N_1)^k & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & (\lambda_s I_{m_s} + N_s)^k \end{pmatrix} \cdot P^{-1}$$

A est stochastique donc $\rho(A) = 1$ et pour tout $i \in \llbracket 1, s \rrbracket$, $|\lambda_i| < 1$.

$$(\lambda_i I_{m_i} + N_i)^k = \sum_{j=0}^{r_i-1} \binom{k}{j} \lambda_i^{k-j} N_i^j \quad \text{où} \quad N_i^{r_i} = 0 \quad \text{et} \quad N_i^{r_i-1} \neq 0$$

$$\text{où} \quad \binom{k}{j} \underset{k \rightarrow +\infty}{\sim} \frac{k^j}{j!}$$

$$\text{Il vient,} \quad \sum_{j=0}^{r_i-1} \binom{k}{j} \lambda_i^{k-j} N_i^j \underset{k \rightarrow +\infty}{\sim} \frac{k^{r_i-1} \lambda_i^{k-(r_i-1)}}{(r_i-1)!} N_i^{r_i-1}$$

Annexe

Preuve de la convergence de l'algorithme

Finalement,

$$A^k \xrightarrow{k \rightarrow +\infty} P \cdot \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix} \cdot P^{-1}$$

D'où,

$$A^k - P \cdot \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix} \cdot P^{-1} \underset{k \rightarrow +\infty}{\sim} C \cdot \lambda_0^k \cdot k^{r_0-1}$$

où $|\lambda_0| = \max_{1 \leq i \leq s} |\lambda_i|$.

```
13 import pandas as pd
14 import numpy as np
15 import matplotlib.pyplot as plt
16 from google.colab import drive
17 from sklearn import preprocessing
18 from sklearn.preprocessing import OneHotEncoder
19 import seaborn as sns
20 from datetime import timedelta
21 from datetime import *
22 from random import *
23 import copy
24 from math import *
25 from time import *
26 from datetime import timedelta, datetime
27 color = sns.color_palette()
28 sns.set_style('darkgrid')
29
30 drive.mount('/content/drive')
31
32 # Chargement du dataset/suppression des colonnes avec des données manquantes
33
34 data = pd.read_csv('/content/drive/MyDrive/TIPE/DATASET_INITIAL.csv')
35 data = data.dropna()
36 data = data.reset_index(drop=True)
37
38 #Conversion des données en type int
39
40 data['Player1Rank'] = pd.to_numeric(data['Player1Rank'])
41 data['Player2Rank'] = pd.to_numeric(data['Player2Rank'])
42 data['Player1Sets'] = pd.to_numeric(data['Player2Sets'])
43 data['Date'] = pd.to_datetime(data.Date)
44
45 winner = data['Winner'].copy()
46
```

Figure 15: Programme 1

```
42 #Player 1 = mieux classé ATP
43
44 def rename(data):
45     for k in range(data.shape[0]):
46         if data.iloc[k,5]>data.iloc[k,12]:
47             data.iloc[k,0],data.iloc[k,1]=data.iloc[k,1],data.iloc[k,0]
48             data.iloc[k,4],data.iloc[k,11]=data.iloc[k,11],data.iloc[k,4]
49             data.iloc[k,5],data.iloc[k,12]=data.iloc[k,12],data.iloc[k,5]
50             data.iloc[k,6],data.iloc[k,13]=data.iloc[k,13],data.iloc[k,6]
51             data.iloc[k,7],data.iloc[k,14]=data.iloc[k,14],data.iloc[k,7]
52     return data
53
54 data = rename(data)
55
56 data['Winner']= winner
57
58 #Colonne Winner : 1 si P1 gagne 0 sino,
59 def win(data):
60     for k in range(data.shape[0]):
61         m = min(data.iloc[k, 5], data.iloc[k, 12])
62         # Si P1 gagne...
63         if data.iloc[k,6]==data.iloc[k, 15]:
64             data.iloc[k,15] = 1
65         else:
66             data.iloc[k, 15] = 0
67     return data
68
69 data = win(data)
70
```

Figure 16: Programme 2

```
71 """## Transformation de certaines variables en variables muettes (dummy variables)"""
72
73 #Transformation des variables round et surface en variables muettes
74 def round_number(x):
75     if x == '1st Round':
76         return 1
77     elif x == '2nd Round':
78         return 2
79     elif x == '3rd Round':
80         return 3
81     elif x == '4th Round':
82         return 4
83     elif x == 'Quarterfinals':
84         return 5
85     elif x == 'Semifinals':
86         return 6
87     elif x == 'The Final':
88         return 7
89
90 data['Round'] = data['Round'].apply(round_number)
91
92 dummy_rounds = pd.get_dummies(data['Round'], prefix='Round')
93
94 data = data.join(dummy_rounds.loc[:, 'Round_1':])
95 data[['Round_1', 'Round_2', 'Round_3',
96       'Round_4', 'Round_5', 'Round_6', 'Round_7']] = data[['Round_1.0', 'Round_2.0', 'Round_3.0', 'Round_4.0',
97       'Round_5.0', 'Round_6.0', 'Round_7.0']].astype('int')
98 data = data.drop(['Round_1.0', 'Round_2.0', 'Round_3.0', 'Round_4.0', 'Round_5.0', 'Round_6.0', 'Round_7.0'], axis =
99 1)
100
101 dummy_surface = pd.get_dummies(data['Surface'], prefix='Surface')
102
103 data = data.join(dummy_surface.loc[:, 'Surface_Clay':])
104 data = data.drop(['Surface'], axis = 1)
105 data[['Surface_Clay', 'Surface_Grass', 'Surface_Hard']] =
106 data[['Surface_Clay', 'Surface_Grass', 'Surface_Hard']].astype('int')
107 data = data.drop(['Surface'], axis = 1)
```

Figure 17: Programme 3

```
106 """***Indexation des joueurs : création d'un dictionnaire contenant en clé le nom de joueur et en valeur :  
107 [index du joueur, son nombre de victoire, total de matchs joués, meilleur classement en carrière]***"  
108 ##peut-être une meilleure manière d'encoder les joueurs  
109 index_joueurs_stats = {}  
110 #Dictionnaire contenant le nom des joueurs comme une clé [player_index,x,y,z,t,w]  
111 #x : nombres de matchs gagnés  
112 #y : nombres de matchs joués  
113 #z : meilleur classement en carrière (on le fixe par défaut à 3000)  
114 # t : date du dernier match  
115 # w : liste des h2h contre les autres joueurs  
116 i=0  
117  
118 for player in data['Player1'].unique():  
119     if player not in index_joueurs_stats.keys():  
120         index_joueurs_stats[player] = [i,0,0,3000,0]  
121         i+=1  
122 for player in data['Player2'].unique():  
123     if player not in index_joueurs_stats.keys():  
124         index_joueurs_stats[player] = [i,0,0,3000,0]  
125         i+=1  
126  
127 print('Nombre exact de joueurs : ',i)  
128
```

Figure 18: Programme 4

```
135 def classement_elo(data,classement,nb_matches):
136
137     # Création d'une colonne elo_ranking
138     elo_p1 = np.zeros((data.shape[0], 1))
139     elo_p2 = np.zeros((data.shape[0], 1))
140     proba_win_best_ranked = np.zeros((data.shape[0], 1))
141
142     # On parcourt tous les matchs, puis on calcule la proba de gain du joueur le
143     # mieux classé contre le moins bien
144     for index, row in data.iterrows():
145         player1, player2 = row[6], row[11]
146         elo1, elo2 = classement[player1], classement[player2]
147         K1, K2 = K[player1], K[player2]
148
149         elo_p1[index, 0] = classement[player1]
150         elo_p2[index, 0] = classement[player2]
151
152
153         D = elo1 - elo2
154         if abs(D) > 400:
155             D = (D/abs(D))*400
156
157         # Proba de gain du joueur le mieux classé sur le pire
158         # Proba de gain de P1 sur P2
159         p = 1/(1 + 10**(-D/400))
160
161         proba_win_best_ranked[index, 0] = round(p,2)
162
163         if nb_matches[player1]<=30:
164             K[player1] = 40
165         elif elo1 >= 2250 or K1==10:
166             K[player1] = 10
167         else:
168             K[player1] = 30
169
170         if nb_matches[player2]<=30:
171             K[player2] = 40
172         elif elo2 >= 2250 or K[player2]==10:
173             K[player2] = 10
174         else:
175             K[player2] = 30
176
```

Figure 19: Programme 5

```
179
180     # Détermination de W, résultat réel
181     W = row[13]
182
183     # Actualisation du classement
184     classement[player1] = round(elo1 + K1*(W-p))
185     classement[player2] = round(elo2 + K2*(p-W))
186
187     nb_matches[player1] += 1
188     nb_matches[player2] += 1
189
190
191
192     return (elo_p1, elo_p2, proba_win_best_ranked)
193
194 elo_p1, elo_p2, proba_win_best_ranked = classement_elo(data,classement,nb_matches)
195
196 # Insertion des colonnes elo_p1 et elo_p2
197 data['elo_P1'] = elo_p1
198 data['elo_P2'] = elo_p2
199 data['proba_elo'] = proba_win_best_ranked
200
```

Figure 20: Programme 6


```
201 """***Colonne meilleur classement en carrière***"""
202
203 data['Best_rank_p1'] = 0
204 data['Best_rank_p2'] = 0
205
206 for c, row in data.iterrows():
207     p1 = row[6]
208     p2 = row[11]
209     if row[5]<row[27]:
210         data.at[c, 'Best_rank_p1'] = row[5]
211         index_joueurs_stats[p1][3] = row[5]
212     else:
213         data.at[c, 'Best_rank_p1'] = index_joueurs_stats[p1][3]
214
215     if row[10]<index_joueurs_stats[p2][3]:
216         data.at[c, 'Best_rank_p2'] = row[10]
217         index_joueurs_stats[p2][3] = row[10]
218     else:
219         data.at[c, 'Best_rank_p2'] = index_joueurs_stats[p2][3]
220
```

Figure 21: Programme 7

```
221 """***Colonnes du pourcentage de victoires (saison en cours)***"""
222
223 index_joueurs_stats_clean = copy.deepcopy(index_joueurs_stats)
224 tableau_ratio = np.zeros((len(data)),2)
225 date_limite = pd.to_datetime('2007-01-01')
226
227 for c,row in data.iterrows():
228     if row[3]>=date_limite:
229         date_limite = pd.to_datetime(date_limite) + pd.DateOffset(years=1)
230         index_joueurs_stats = copy.deepcopy(index_joueurs_stats_clean)
231         score_p1 = index_joueurs_stats[row[6]]
232         score_p2 = index_joueurs_stats[row[11]]
233         if score_p1[2]<5:
234             tableau_ratio[c,0]=0.2
235         else:
236             tableau_ratio[c,0] =round(score_p1[1]/score_p1[2],2)
237         if score_p2[2]<5:
238             tableau_ratio[c,1]=0.2
239         else:
240             tableau_ratio[c,1] = round(score_p2[1]/score_p2[2],2)
241         if row[13]==1:
242             index_joueurs_stats[row[6]][0],index_joueurs_stats[row[6]][1],index_joueurs_stats[row[6]][2] =
score_p1[0],score_p1[1]+1,score_p1[2]+1
243             index_joueurs_stats[row[11]][0],index_joueurs_stats[row[11]][1],index_joueurs_stats[row[11]][2] =
score_p2[0],score_p2[1]+1,score_p2[2]+1
244         else:
245             index_joueurs_stats[row[6]][0],index_joueurs_stats[row[6]][1],index_joueurs_stats[row[6]][2] =
score_p1[0],score_p1[1],score_p1[2]+1
246             index_joueurs_stats[row[11]][0],index_joueurs_stats[row[11]][1],index_joueurs_stats[row[11]][2] =
score_p2[0],score_p2[1],score_p2[2]+1
247
248 data['ratio_p1_saison'] = tableau_ratio[:,0]
249 data['ratio_p2_saison'] = tableau_ratio[:,1]
```

Figure 22: Programme 8

```
251 """## **Colonne du pourcentage de victoire (en carrière)****
252
253 index_joueurs_stats = copy.deepcopy(index_joueurs_stats_clean)
254 tableau_ratio_1 = np.zeros(((len(data)),2))
255
256 for c,row in data.iterrows():
257     score_p1 = index_joueurs_stats[row[6]]
258     score_p2 = index_joueurs_stats[row[11]]
259     if score_p1[2]<5:
260         tableau_ratio_1[c,0]=0.2
261     else:
262         tableau_ratio_1[c,0] =round(score_p1[1]/score_p1[2],2)
263     if score_p2[2]<5:
264         tableau_ratio_1[c,1]=0.2
265     else:
266         tableau_ratio_1[c,1] = round(score_p2[1]/score_p2[2],2)
267     if row[13]==1 :
268         index_joueurs_stats[row[6]][0],index_joueurs_stats[row[6]][1],index_joueurs_stats[row[6]][2] =
269         score_p1[0],score_p1[1]+1,score_p1[2]+1
270         index_joueurs_stats[row[11]][0],index_joueurs_stats[row[11]][1],index_joueurs_stats[row[11]][2] =
271         score_p2[0],score_p2[1],score_p2[2]+1
272     else:
273         index_joueurs_stats[row[6]][0],index_joueurs_stats[row[6]][1],index_joueurs_stats[row[6]][2] =
274         score_p1[0],score_p1[1],score_p1[2]+1
275         index_joueurs_stats[row[11]][0],index_joueurs_stats[row[11]][1],index_joueurs_stats[row[11]][2] =
276         score_p2[0],score_p2[1]+1,score_p2[2]+1
277
278 data['ratio_p1_carrière'] = tableau_ratio_1[:,0]
279 data['ratio_p2_carrière'] = tableau_ratio_1[:,1]
```

Figure 23: Programme 9

```
277 """Fonction qui calcule la proba de gain avec le classement elo"""
278
279 def calcul_proba_elo(pts_elo1, pts_elo2, seuil):
280     D = pts_elo1 - pts_elo2
281     if abs(D) > seuil:
282         D = (D/abs(D))*seuil
283
284     # Proba de gain de P1 sur P2
285     p = 1/(1 + 10**(-D/seuil))
286
287     return (p, 1-p)
288
289 !pip install fuzzywuzzy
290
291 """Fonction basé sur la distance de levenstein pour harmoniser les noms des joueurs entre le classement et le
dataset"""
292
293 from fuzzywuzzy import fuzz
294
295 def abréviation_nom(nom_complet):
296     parties_nom = nom_complet.split()
297     nom_abrégé = parties_nom[-1] + ' ' + ''.join(parties_nom[0][0]) + '.'
298     return nom_abrégé
299
300 def trouver_correspondance_nom(nom_inexact, liste_noms):
301     correspondance_score_max = 0
302     nom_correspondant = None
303     for nom_exact in liste_noms:
304         score = fuzz.partial_ratio(nom_inexact, nom_exact)
305         if score > correspondance_score_max:
306             correspondance_score_max = score
307             nom_correspondant = nom_exact
308     return nom_correspondant
309
```

Figure 24: Programme 10

```

314 ###cette fonction permet de creer le classement pour une date donnée
315
316 def classement_mois(date):
317     # Ici on récupère les données relatives à chacun des joueurs à la date directement disponible dans le dataset
318     # ie : pour un joueur : [classement_elo, ratio_saison, ratio_carrière]
319     stats_joueurs_pour_proba = {}
320     date_limite = pd.to_datetime(date)
321     date_limite_formatee = date_limite.strftime('%d-%m-%Y')
322     data_r = data[(data['Date']>=date_limite-timedelta(days=730)) & (data['Date']<=date_limite_formatee)]
323     data_r = data_r.reset_index(drop=True)
324     data_r['Date'] = data_r['Date'].dt.strftime('%d-%m-%Y')
325     for index, ligne in data_r.iterrows():
326         # Données joueur 1
327         stats_joueurs_pour_proba[ligne[6]] = [ligne[24], ligne[29], ligne[31]]
328
329         # Données joueur 2
330         stats_joueurs_pour_proba[ligne[11]] = [ligne[25], ligne[30], ligne[32]]
331
332
333     # On importe le classement ATP du 30/10/23 et on récupère les données de chaque joueur (100 premiers à l'ATP)
334     # NB : Nécessairement il ne reste donc que les joueurs (ayant joué un match ou plus les deux dernières années)
335     ET (présent dans le top 100 ATP à la date)
336     atp_rkg = pd.read_csv('f:/content/drive/MyDrive/TIPE/classements atp/{date}.csv')
337     atp_rkg = atp_rkg[1:100]
338     stats_joueurs_pour_proba_final = {}
339     liste_joueurs = stats_joueurs_pour_proba.keys()
340     print(liste_joueurs)
341     cpt_index = 0
342     for index, row in atp_rkg.iterrows():
343         name = row[1]
344         l = name.split(' ')
345         if len(l) >= 2:
346             name = trouver_correspondance_nom(abreviation_nom(name), liste_joueurs)
347             if name in liste_joueurs:
348                 f_stats = stats_joueurs_pour_proba[name]
349                 # STATS : [index, points ATP, classement ATP actuel, points elo, meilleur classement carrière ATP, Ratio
350                 saison, Ratio carrière]
351                 stats = [cpt_index, row[4], row[0], f_stats[0], row[5]] + f_stats[1:]
352                 stats_joueurs_pour_proba_final[name] = stats
353                 cpt_index += 1
354                 #print(cpt_index, name, stats, stats_joueurs_pour_proba_final[name])
355                 print(stats_joueurs_pour_proba_final)

```

Figure 25: Programme 11

```
359 ##(Cette fonction ajoute au dictionnaire stats_joueurs_pour_proba_final le nombre de matchs gagnés sur chaque  
360 autre joueur sous forme d'une liste  
361  
362 def get_h2h():  
363     n = len(stats_joueurs_pour_proba_final)  
364     for player1 in stats_joueurs_pour_proba_final:  
365         p1_h2h = [0 for _ in range(n)]  
366         cpt_p2 = 0  
367  
368         for player2 in stats_joueurs_pour_proba_final:  
369             data_pour_h2h = data[((data['Player1'] == player1) & (data['Player2'] == player2)) | ((data['Player1']  
== player2) & (data['Player2'] == player1))]  
370  
371             match_p1 = 0  
372             for index, row in data_pour_h2h.iterrows():  
373                 # S'il est player1, mieux classé et winner = 1 OU player2, moins bien classé et winner = 0  
374                 if (row[6] == player1 and row[13]==1) or (row[11]==player1 and row[13]==0):  
375                     match_p1 += 1  
376  
377             p1_h2h[cpt_p2] = match_p1  
378             cpt_p2 += 1  
379             stats_joueurs_pour_proba_final[player1].append(p1_h2h)  
380  
381 get_h2h()  
382
```

Figure 26: Programme 12

```
384 ### Cette fonction calcule la proba de victoire d'un joueur i sur un joueur j en se basant sur la formule du  
385 classement elo adaptée à plusieurs statistiques  
386  
387 def calcul_proba(player1, player2):  
388     stats1, stats2 = stats_joueurs_pour_proba_final[player1], stats_joueurs_pour_proba_final[player2]  
389  
390     # Calcul proba elo  
391     p1_elo, p2_elo = calcul_proba_elo(stats1[3], stats2[3], 400)  
392  
393     # Calcul proba ATP  
394     p1_atp, p2_atp = calcul_proba_elo(stats1[1], stats2[1], 3000)  
395  
396     # Calcul proba classement ATP  
397     p1_catp, p2_catp = calcul_proba_elo(stats1[2], stats2[2], 40)  
398  
399     # Calcul proba ratio saison  
400     p1_rs, p2_rs = calcul_proba_elo(stats1[5], stats2[5], 0.4)  
401  
402     # Calcul proba h2h  
403     h2h_tot = stats1[7][stats2[0]] + stats2[7][stats1[0]]  
404     if h2h_tot >= 5:  
405         p1_h2h, p2_h2h = calcul_proba_elo(stats1[7][stats2[0]], stats2[7][stats1[0]], 10)  
406  
407     # Calcul proba finale  
408     p = (50*p1_elo + 25*p1_rs + 25*p1_h2h)/100  
409     else:  
410         p = (70*p1_elo + 30*p1_rs)/100  
411  
412     return (p, 1-p)
```

Figure 27: Programme 13

```
413 ### Création de la matrice de transition
414 dict_index_nom = {stats_joueurs_pour_proba_final[player][0]:player for player in stats_joueurs_pour_proba_final}
415 print(dict_index_nom)
416
417 n = len(stats_joueurs_pour_proba_final)
418 matrice_transition = [[0 for _ in range(n)] for _ in range(n)]
419 for j in range(n):
420     somme = 0
421
422     # On remplit la colonne j
423     for i in range(n):
424         if i != j :
425             pji = calcul_proba(dict_index_nom[j], dict_index_nom[i])[0]
426             matrice_transition[i][j] = pji
427             somme += pji
428         else :
429             matrice_transition[i][j] = 0
430
431
432     #On normalise la colonne
433     for i in range(n):
434         matrice_transition[i][j] /= somme
435
436 matrice_transition = np.array(matrice_transition)
437
438
439 n = 100 # dimension de la matrice de transition
440
441
```

Figure 28: Programme 14


```
443 # Initialisation du vecteur
444 Y = np.array([[1/n] for _ in range(n)])
445
446 # Itération jusqu'à convergence
447 for i in range(1000):
448     tmp = np.dot(matrice_transition, Y)
449     Y = tmp / np.linalg.norm(tmp)
450
451 # Calcul de la valeur propre dominante
452 valeur_propre_dominante = np.dot(tmp.T, Y) / np.dot(Y.T, Y)
453 Y_normalized = Y / np.sum(Y)
454
455 # Triez les indices des états en fonction des valeurs du vecteur propre
456 classement_indices = np.argsort(Y_normalized, axis=0)[::-1]
457
458 # Création du classement
459 classement = [(indice, Y_normalized[indice][0]) for indice in classement_indices]
460
461 C = []
462 for i in range(len(classement)):
463     C.append(dict_index_nom[classement[i][0][0]])
464
465 return(C)
466
```

Figure 29: Programme 15

```
469 ### Cette fonction permet finalement creer le dictionnaire souhaité
470 def obtenir_classements_personnels():
471     classements_personnels = {}
472     date_debut = datetime(2015, 1, 1)
473     date_fin = datetime(2023, 10, 1)
474     mois_actuel = date_debut
475
476     while mois_actuel <= date_fin:
477         # Vérifier si le mois actuel est dans la période à sauter
478         if mois_actuel.year == 2020 and mois_actuel.month >= 4 and mois_actuel.month <= 8:
479             # Passer au mois suivant sans rien faire
480             mois_suivant = mois_actuel.replace(month=mois_actuel.month + 1)
481             mois_actuel = mois_suivant
482             continue
483
484         # Calculer le premier jour du mois suivant
485         mois_suivant = mois_actuel.replace(day=1)
486         if mois_suivant.month == 12:
487             mois_suivant = mois_suivant.replace(year=mois_suivant.year + 1, month=1)
488         else:
489             mois_suivant = mois_suivant.replace(month=mois_suivant.month + 1)
490
491         # Obtenir le classement des joueurs pour le premier jour de ce mois
492         date_str = mois_actuel.strftime('%d-%m-%Y')
493         classement = classement_mois(date_str)
494
495         # Ajouter ce classement à la liste des classements mensuels
496         classements_personnels[mois_actuel.strftime('%Y-%m')] = classement
497
498         # Passer au mois suivant
499         mois_actuel = mois_suivant
500
501     return classements_personnels
502
503 # Obtention des classements mensuels
504 classements_personnels = obtenir_classements_personnels()
```

Figure 30: Programme 16

```

506 ### Je ne garde dans le dataset initial que les matchs à partir du 01-01-2015
507
508 data = data[data['Date']>='2015-01-01']
509 data = data.reset_index(drop=True)
510
511 #Suppression des matchs impliquant des joueurs au delà de la 100ème place ATP(ceux qui n'ont pas de classement)
512 for c, row in data.iterrows():
513     if (row[5]>100 or row[10]>100):
514         data = data.drop([c], axis = 0)
515 data = data.reset_index(drop=True)
516
517 *** ##Création d'un nouveau dataset dans lequel les colonnes classement sont remplies avec mon classement****
518
519 import pandas as pd
520 from datetime import datetime
521
522 def mettre_a_jour_classements(data, classements_personnels):
523     # Créer une copie du dataset original pour le nouveau dataset
524     data_my_rank = data.copy()
525
526     # Convertir la colonne de dates en format datetime
527     data_my_rank.iloc[:, 3] = pd.to_datetime(data_my_rank.iloc[:, 3])
528
529     # Fonction pour obtenir le classement d'un joueur à une date donnée
530     def obtenir_classement_joueur(joueur, date_match):
531         date_str = date_match.strftime('%Y-%m')
532         classement_mois = classements_personnels.get(date_str, [])
533         try:
534             classement_joueur = len(classement_mois) - classement_mois.index(joueur)
535         except ValueError:
536             classement_joueur = float('NaN')
537         return classement_joueur

```

Figure 31: Programme 17

```
538
539 # Mettre à jour les colonnes de classement des joueurs dans le nouveau dataset
540 for index, row in data_my_rank.iterrows():
541     date_match = row[3]
542     nom_joueur1 = row[6]
543     nom_joueur2 = row[11]
544     classement_joueur1 = obtenir_classement_joueur(nom_joueur1, date_match)
545     classement_joueur2 = obtenir_classement_joueur(nom_joueur2, date_match)
546     data_my_rank.at[index, 'Player1Rank'] = classement_joueur1
547     data_my_rank.at[index, 'Player2Rank'] = classement_joueur2
548
549     return data_my_rank
550
551 data_my_rank = mettre_a_jour_classements(data, classements_personnels)
552
553 data = data.dropna()
554 data = data.reset_index(drop=True)
555
556 data_my_rank = data_my_rank.dropna()
557 data_my_rank = data_my_rank.reset_index(drop=True)
```

Figure 32: Programme 18

```
27 # Chargement du data_atpset/suppression des colonnes avec des données manquantes
28
29 data_atp = pd.read_csv('/content/drive/MyDrive/TIPE/DATASET_ATP_FINAL.csv')
30 data_atp = data_atp.dropna()
31 data_atp = data_atp.reset_index(drop=True)
32
33 data_myrank = pd.read_csv('/content/drive/MyDrive/TIPE/DATA_MYRANK_FINAL.csv')
34 data_myrank = data_myrank.dropna()
35 data_myrank = data_myrank.reset_index(drop=True)
36
37 #Conversion des données en type int
38
39 data_atp['Player1Rank'] = pd.to_numeric(data_atp['Player1Rank'])
40 data_atp['Player2Rank'] = pd.to_numeric(data_atp['Player2Rank'])
41 data_atp['Player1Sets'] = pd.to_numeric(data_atp['Player2Sets'])
42 data_atp['Date'] = pd.to_datetime(data_atp.Date)
43
44 data_myrank['Player1Rank'] = pd.to_numeric(data_myrank['Player1Rank'])
45 data_myrank['Player2Rank'] = pd.to_numeric(data_myrank['Player2Rank'])
46 data_myrank['Player1Sets'] = pd.to_numeric(data_myrank['Player2Sets'])
47 data_myrank['Date'] = pd.to_datetime(data_myrank.Date)
48
```

Figure 33: Programme 19

```
58 ##peut-être une meilleure manière d'encoder les joueurs
59 index_joueurs_stats = {}
60 #Dictionnaire contenant le nom des joueurs comme une clé [player_index,x,y,z,t,w]
61 #x : nombres_de_matches_gagnés
62 #y : nombres_de_matches_joués
63 #z : meilleur classement en carrière (on le fixe par défaut à 3000)
64 # t : date du dernier match
65 # w : liste des h2h contre les autres joueurs
66 i=0
67
68 for player in data_myrank['Player1'].unique():
69     if player not in index_joueurs_stats.keys():
70         index_joueurs_stats[player] = [1,0,0,3000,0]
71         i+=1
72 for player in data_myrank['Player2'].unique():
73     if player not in index_joueurs_stats.keys():
74         index_joueurs_stats[player] = [1,0,0,3000,0]
75         i+=1
76
77 print('Nombre exact de joueurs : ',i)
78
79 data_myrank['Best_rank_p1'] = 0
80 data_myrank['Best_rank_p2'] = 0
81
82 for c, row in data_myrank.iterrows():
83     p1 = row[6]
84     p2 = row[11]
85     if row[5]<index_joueurs_stats[p1][3]:
86         data_myrank.at[c, 'Best_rank_p1'] = row[5]
87         index_joueurs_stats[p1][3] = row[5]
88     else:
89         data_myrank.at[c, 'Best_rank_p1'] = index_joueurs_stats[p1][3]
90
91     if row[10]<index_joueurs_stats[p2][3]:
92         data_myrank.at[c, 'Best_rank_p2'] = row[10]
93         index_joueurs_stats[p2][3] = row[10]
94     else:
95         data_myrank.at[c, 'Best_rank_p2'] = index_joueurs_stats[p2][3]
96
```

Figure 34: Programme 20

```

97 """On mélange aléatoirement P1 et P2"""
98
99 import pandas as pd
100 import numpy as np
101
102 def randomize_players_and_winner(df):
103     # Sélectionner les colonnes correspondant aux caractéristiques des joueurs
104     player_cols = ['Player1', 'B365Player1', 'B365Player2', 'Player1Rank', 'Player2Rank',
105                   'Player2', 'elo_P1', 'elo_P2', 'Best_rank_p1', 'Best_rank_p2',
106                   'ratio_p1_saison', 'ratio_p2_saison', 'ratio_p1_carrière', 'ratio_p2_carrière']
107
108     # Créer une copie du data_atpFrame pour ne pas modifier l'original
109     df_copy = df.copy()
110
111     # Pour chaque ligne du data_atpFrame
112     for index, row in df_copy.iterrows():
113         # Inverser aléatoirement les colonnes correspondant aux caractéristiques des joueurs
114         if np.random.rand() < 0.5:
115             # Inverser les colonnes
116             df_copy.at[index, 'Winner'] = 1 - df_copy.at[index, 'Winner']
117             for col in player_cols:
118                 col_index = df_copy.columns.get_loc(col)
119                 opponent_col = col.replace('1', '2').replace('2', '1')
120                 opponent_col_index = df_copy.columns.get_loc(opponent_col)
121                 df_copy.iloc[index, col_index], df_copy.iloc[index, opponent_col_index] = df_copy.iloc[index,
122                 opponent_col_index], df_copy.iloc[index, col_index]
123
124     return df_copy
125
126 data_atp = randomize_players_and_winner(data_atp)
127 data_myrank = randomize_players_and_winner(data_myrank)

```

Figure 35: Programme 21

```
128 *****Dernière modifications des datasets...*****
129
130 # Suppression des colonnes qui ne seront plus utiles au modèle
131 data_atp = data_atp.drop(['Player1', 'Player2', 'Player1Sets', 'Player2Sets', 'Date'], axis = 1)
132 data_myrank = data_myrank.drop(['Player1', 'Player2', 'Player1Sets', 'Player2Sets', 'Date'], axis = 1)
133
134 data_atp['Winner'] = pd.to_numeric(data_atp['Winner'])
135 data_myrank['Winner'] = pd.to_numeric(data_myrank['Winner'])
136
137 data_atp = data_atp.drop(['Series'], axis = 1)
138 data_myrank = data_myrank.drop(['Series'], axis = 1)
139
140 data_atp['Player1Pts'] = data_atp['Player1Pts'].astype('int64')
141 data_atp['Player2Pts'] = data_atp['Player2Pts'].astype('int64')
142 data_atp['Player1Rank'] = data_atp['Player1Rank'].astype('int64')
143 data_atp['Player2Rank'] = data_atp['Player2Rank'].astype('int64')
144 data_atp['elo_P1'] = data_atp['elo_P1'].astype('int64')
145 data_atp['elo_P2'] = data_atp['elo_P2'].astype('int64')
146
147
148 data_myrank['Player1Pts'] = data_myrank['Player1Pts'].astype('int64')
149 data_myrank['Player2Pts'] = data_myrank['Player2Pts'].astype('int64')
150 data_myrank['Player1Rank'] = data_myrank['Player1Rank'].astype('int64')
151 data_myrank['Player2Rank'] = data_myrank['Player2Rank'].astype('int64')
152 data_myrank['elo_P1'] = data_myrank['elo_P1'].astype('int64')
153 data_myrank['elo_P2'] = data_myrank['elo_P2'].astype('int64')
154
```

Figure 36: Programme 22


```
155 """***DIVISION DES DATASETS POUR L'ENTRAINEMENT ET LE TEST***  
156  
157 def pre_processing(df):  
158     df = df.dropna()  
159  
160     X = df.drop(['Winner'], axis = 1)  
161     y = df['Winner']  
162  
163     return X, y  
164  
165 from sklearn.model_selection import train_test_split  
166 trainset_atp, testset_atp = train_test_split(data_atp, test_size = 0.2, random_state = 45)  
167 trainset_myrank, testset_myrank = train_test_split(data_myrank, test_size = 0.2, random_state = 45)  
168  
169 X_train_atp, y_train_atp = pre_processing(trainset_atp)  
170 X_train_myrank, y_train_myrank = pre_processing(trainset_myrank)  
171  
172 X_test_atp, y_test_atp = pre_processing(testset_atp)  
173 X_test_myrank, y_test_myrank = pre_processing(testset_myrank)  
174  
175 print(y_test_myrank.shape)  
176  
177 """***MODELISATION***  
178  
179 from sklearn import metrics  
180 from sklearn.model_selection import GridSearchCV  
181 from sklearn.metrics import classification_report  
182 from sklearn.model_selection import GroupKFold  
183 from sklearn.model_selection import RandomizedSearchCV  
184 from sklearn.linear_model import LogisticRegression  
185 from sklearn.ensemble import RandomForestClassifier  
186 from sklearn.ensemble import AdaBoostClassifier  
187 from sklearn.neural_network import MLPRegressor  
188 from sklearn.model_selection import learning_curve  
189 from sklearn.metrics import accuracy_score
```

Figure 37: Programme 23

```
207 from sklearn.tree import DecisionTreeClassifier
208 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
209 from sklearn.naive_bayes import GaussianNB
210 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis, LinearDiscriminantAnalysis
211
212 from sklearn import model_selection
213
214
215 X = X_train_atp
216 Y = y_train_atp
217
218 seed = 42
219
220 models = []
221 models.append(('LR', LogisticRegression()))
222 models.append(('LDA', LinearDiscriminantAnalysis()))
223 models.append(('QDA', QuadraticDiscriminantAnalysis()))
224 models.append(('KNN', KNeighborsClassifier(125, n_jobs=-1)))
225 models.append(('CART', DecisionTreeClassifier(max_depth=10)))
226 models.append(('NB', GaussianNB()))
227
228 models.append(('RandomForest', RandomForestClassifier(n_estimators=100, n_jobs=-1)))
229 models.append(('MLP', MLPClassifier(alpha=0.00001)))
230 models.append(('ADABOOST', AdaBoostClassifier()))
231
232
233 results = {}
234 scoring = {'accuracy': make_scorer(accuracy_score),
235           'precision_score': make_scorer(precision_score),
236           'recall_score': make_scorer(recall_score),
237           'f1_score': make_scorer(f1_score)}
238 names = []
```

Figure 38: Programme 24

```
239 for name, model in models:
240
241     stratifiedKFold = model_selection.StratifiedKFold(n_splits=10, random_state=seed, shuffle = True)
242     cv_results = model_selection.cross_validate(model, X, Y, cv=stratifiedKFold, scoring=scoring)
243     results.append(cv_results)
244     names.append(name)
245     msg =
246     = '-----\n'
247     msg = "Model : %s \n" % (name)
248     msg = msg + '\n'
249     msg = msg + "Accuracy : %f (%f)\n" % (cv_results['test_accuracy'].mean(), cv_results['test_accuracy'].std())
250     msg = msg + "Precision score : %f (%f)\n" % (cv_results['test_precision_score'].mean(),
251     cv_results['test_precision_score'].std())
252     msg = msg + "Recall score : %f (%f)\n" % (cv_results['test_recall_score'].mean(),
253     cv_results['test_recall_score'].std())
254     msg = msg + "F1 score : %f (%f)\n" % (cv_results['test_f1_score'].mean(), cv_results['test_f1_score'].std())
255     msg = msg +
256     '-----\n'
257     print(msg)
```

Figure 39: Programme 25

```
256 *****GRID SEARCH CV**
257
258 **RandomForest**
259 ***
260
261 from sklearn.ensemble import RandomForestClassifier
262 from sklearn.model_selection import GridSearchCV
263 from sklearn.metrics import accuracy_score
264
265 # Définition des hyperparamètres à tester
266 param_grid = {
267     'n_estimators': [50, 100, 200],
268     'max_depth': [None, 10, 20],
269     'min_samples_split': [2, 5, 10],
270     'min_samples_leaf': [1, 2, 4]
271 }
272
273 # Création du premier modèle avec le dataset ATP
274 rf_atp = RandomForestClassifier()
275
276 # Création du grid search pour le premier modèle
277 grid_search_atp = GridSearchCV(estimator=rf_atp, param_grid=param_grid, cv=5, scoring='accuracy', n_jobs=-1)
278
279 # Entraînement du grid search pour le premier modèle
280 grid_search_atp.fit(X_train_atp, y_train_atp)
281
282 # Meilleurs hyperparamètres pour le premier modèle
283 best_params_atp = grid_search_atp.best_params_
284
285 # Meilleur modèle pour le premier modèle
286 best_model_atp = grid_search_atp.best_estimator_
287
288 # Score du meilleur modèle pour le premier modèle
289 best_score_atp = grid_search_atp.best_score_
290
291 # Affichage des résultats pour le premier modèle
292 print("Meilleurs hyperparamètres pour le modèle avec classement ATP:", best_params_atp)
293 print("Score avec les meilleurs hyperparamètres pour le modèle avec classement ATP:", best_score_atp)
294
295 # Création du deuxième modèle avec votre propre classement
296 rf_personnel = RandomForestClassifier()
297
298 # Création du grid search pour le deuxième modèle
299 grid_search_personnel = GridSearchCV(estimator=rf_personnel, param_grid=param_grid, cv=5, scoring='accuracy',
300                                     n_jobs=-1)
```

Figure 40: Programme 26

```
300 -
301 # Entraînement du grid search pour le deuxième modèle
302 grid_search_personnel.fit(X_train_myrank, y_train_myrank)
303
304 # Meilleurs hyperparamètres pour le deuxième modèle
305 best_params_personnel = grid_search_personnel.best_params_
306
307 # Meilleur modèle pour le deuxième modèle
308 best_model_personnel = grid_search_personnel.best_estimator_
309
310 # Score du meilleur modèle pour le deuxième modèle
311 best_score_personnel = grid_search_personnel.best_score_
312
313 # Affichage des résultats pour le deuxième modèle
314 print("\nMeilleurs hyperparamètres pour le modèle avec votre propre classement:", best_params_personnel)
315 print("Score avec les meilleurs hyperparamètres pour le modèle avec votre propre classement:",
316       best_score_personnel)
317
318 random_forest_myrank = RandomForestClassifier(max_depth = 10, min_samples_leaf = 3, min_samples_split = 5,
319                                              n_estimators = 25)
320 random_forest_myrank.fit(X_train_myrank, y_train_myrank)
321 print(random_forest_myrank.score(X_test_myrank, y_test_myrank))
322 print(random_forest_myrank.score(X_train_myrank, y_train_myrank))
323
324 # Prédire les probabilités pour le modèle avec votre propre classement
325 probas_myrank = random_forest_myrank.predict_proba(X_test_myrank)[: , 1]
326 predictions_myrank = random_forest_myrank.predict(X_test_myrank)
327
328 random_forest_atp = RandomForestClassifier(max_depth = 10, min_samples_leaf = 4, min_samples_split = 10,
329                                           n_estimators = 200)
330 random_forest_atp.fit(X_train_atp, y_train_atp)
331 print(random_forest_atp.score(X_train_atp, y_train_atp))
332 print(random_forest_atp.score(X_test_atp, y_test_atp))
333 # Prédire les probabilités pour le modèle avec classement ATP
334 probas_atp = random_forest_atp.predict_proba(X_test_atp)[: , 1]
335 predictions_atp = random_forest_atp.predict(X_test_atp)
```

Figure 41: Programme 27

```
335 """Lucidité des modèles"""
336 """
337
338 import numpy as np
339 import matplotlib.pyplot as plt
340
341 # Définir une liste de seuils
342 seuils = np.linspace(0.5, 0.95, 100)
343
344 # Initialiser une liste pour stocker les fréquences de bonnes prédictions
345 freq_bonnes_predictions_myrank = []
346
347 # Pour chaque seuil, calculer la fréquence de bonnes prédictions
348 for seuil in seuils:
349     # Filtrer les prédictions avec des probabilités éloignées de 0,5
350     indices = ((probas_myrank < 1-seuil) | (probas_myrank > seuil)).flatten()
351
352     # Sélectionner les prédictions et les vraies étiquettes correspondantes
353     score_filtre = accuracy_score(y_test_myrank[indices], predictions_myrank[indices])
354     freq_bonnes_predictions_myrank.append(score_filtre)
355
356
357 # Initialiser une liste pour stocker les fréquences de bonnes prédictions
358 freq_bonnes_predictions_atp = []
359
360 # Pour chaque seuil, calculer la fréquence de bonnes prédictions
361 for seuil in seuils:
362     # Filtrer les prédictions avec des probabilités éloignées de 0,5
363     indices = ((probas_atp < 1-seuil) | (probas_atp > seuil)).flatten()
364
365     # Sélectionner les prédictions et les vraies étiquettes correspondantes
366     score_filtre = accuracy_score(y_test_atp[indices], predictions_atp[indices])
367     freq_bonnes_predictions_atp.append(score_filtre)
368
369
370 freq_bonnes_predictions_atp, freq_bonnes_predictions_myrank = np.array(freq_bonnes_predictions_atp),
371 np.array(freq_bonnes_predictions_myrank)
372 ecart = np.abs(freq_bonnes_predictions_myrank - freq_bonnes_predictions_atp)
373 # Calculer l'écart moyen
374 ecart_moyen = np.mean(ecart)
375
376 print(f"L'écart moyen entre les deux courbes est : {ecart_moyen}")
377
378 # Tracer la courbe de la fréquence de bonnes prédictions en fonction des seuils
379 plt.plot(seuils, freq_bonnes_predictions_myrank, label = 'Mon classement')
380 plt.plot(seuils, freq_bonnes_predictions_atp, label = 'Classement ATP')
381 plt.xlabel('Seuil')
382 plt.ylabel('Fréquence de bonnes prédictions')
383 plt.grid(True)
384 plt.legend()
385 plt.show()
```

Figure 42: Programme 28