



MASTER INFORMATIQUE

PARCOURS MACHINE LEARNING

DE LA FOUILLE DE DONNÉES À L'AUTO-ML

Projet FIFA - Part IV: autoML

Auteur:
Selim LAKHDAR

Professeur:
Laetitia JOURDAN



January 10, 2022

Contents

1	Context	2
2	Librairies autoML	2
2.1	TPOT	2
2.2	Hyperopt-Sklearn	3
2.3	Auto-Sklearn	3
3	Dataset	3
3.1	Attributs	3
3.1.1	Non Significatif	3
3.1.2	Numérique	4
3.1.3	Catégoriel	5
3.2	Discrétisation	5
3.2.1	Value	5
3.2.2	Wage	7
4	Regression	7
4.1	Value	7
4.2	Wage	8
4.3	Overall	9
5	Classification	10
5.0.1	DValue: 3 Groups	10
5.0.2	DValue: 4 Groups	11
5.0.3	DValue: 6 Groups	12
5.1	DWage	12
6	GridSearch	13
7	RandomSearch	13
8	Conclusion	13
9	Annexe	14

1 Context

Ce rapport fait suite aux deux premiers rapports: *Analyse descriptive des données* et *Prediction* où nous avons abordé:

- La préparation de données (nettoyage de notre dataset FIFA 2019).
- La création de groupes d'attributs (Discrétisation).
- L'analyse de corrélation.
- La segmentation grâce au clustering.
- De la prediction: Classification et Régression.

Dans cette dernière partie nous allons étudier un mécanisme d'automatisation de la chaine de Machine Learning. Nous allons plus particulièrement étudier la librairie *auto-sklearn* et comparer les méthodes disponibles avec les méthodes classiques de random search et grid search.

2 Bibliothèques autoML

Plusieurs libraires d'autoML sont disponibles dans la littérature. Trois libraires se démarquent. Elles sont toutes utilisées au-dessus de Scikit-Learn que nous avons utilisé lors des derniers rapports.

2.1 TPOT

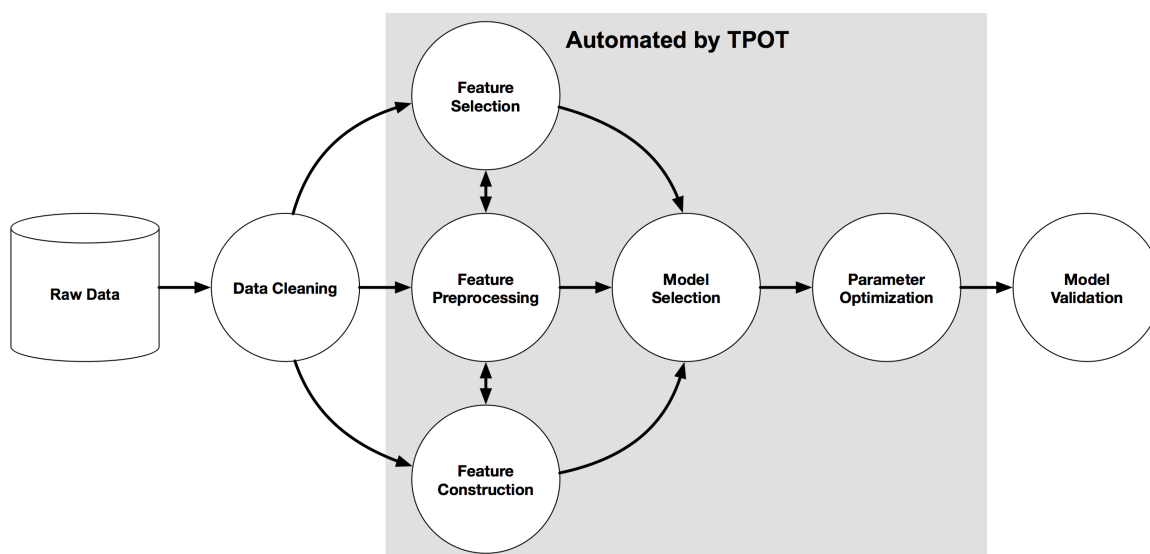


Figure 1: TPOT pipeline

Nous pouvons par exemple citer Tree-based Pipeline Optimization Tool (TPOT) qui utilise une structure d'arbre pour représenter son pipeline comme l'indique la Figure 1. Introduit en 2016 par l'article Evaluation of a Tree-based Pipeline Optimization Tool for Automating Data Science. Proceedings of GECCO 2016, pages 485-492.

2.2 Hyperopt-Sklearn

Hyperopt-Sklearn est une autre librairie open source pour l'optimisation bayésienne. Développé par James Bergstra et introduit en 2013 grâce à l'article Hyperopt-Sklearn: Automatic Hyperparameter Configuration for Scikit-Learn, 2014. Il est utilisé pour l'optimisation à large échelle avec plusieurs centaines de paramètres de façon distribuée.

2.3 Auto-Sklearn

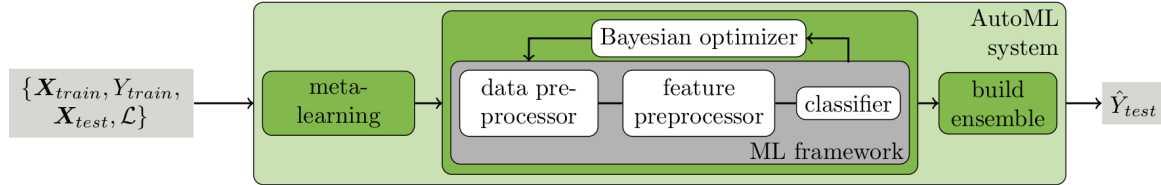


Figure 2: Auto-Sklearn overview

Dans la suite du rapport nous verrons l'utilisation de la librairie Auto-Sklearn. Elle est l'une des plus utilisées. Dans notre cas nous avons utilisé cette librairie car elle est facile à prendre en main et beaucoup d'exemples sont disponible.

Introduit en 2015 par Matthias Feurer et al dans leur article Efficient and Robust Automated Machine Learning. La Figure 1 résume les différentes phases de auto-sklearn. Dans sa première version, Auto-Sklearn utilise 15 classifieurs, 14 méthodes de feature preprocessing et 4 data preprocessing methods, ce qui donne un totale de 110 hyperparamètres.

3 Dataset

Nous reprendrons le dataset que nous avons obtenu après avoir "nettoyer" et discrétiser nos attributs.

3.1 Attributs

Pour résumer notre dataset nous pouvons catégoriser nos attributs ainsi;

3.1.1 Non Significatif

La Figure 3 représente les différents attributs que nous n'utiliserons pas lors de la prédiction, car elles n'apportent pas plus d'informations à notre dataset.

Attribut	Description
Name	Nom du joueur
Photo	Photo du joueur
Club Logo	Logo du Club
Flag	Drapeau du joueur
Real Face	Image du joueur
Name	Nom du joueur
Nationality	Nationalité du joueur
Club	Club du joueur
Joined	Date à la quelle le joueur a rejoint le club
Loaned From	Acheté depuis quel club
Contract Valid Until	durée du contrat

Figure 3: Attributs non significatif

3.1.2 Numérique

La Figure 4 représente les différents attributs numériques de notre dataset.

Attribut	Description
Overall	Score du joueur
Potential	Potentiel du joueur
Value	Valeur du joueur
Wage	Valeur du joueur suivant son age
Special	Statistique Spécial liée au joueur
International Reputation	Réputation du joueur
Weak Foot	Pied faible du joueur
Skill Moves	Statistique de déplacement liée au joueur
Jersey Number	Numéro du joueur
Height	Longueur du joueur
Weight	Poids du joueur
Release Clause	Clause de libération
BMI	Indice Masse Corporel

Figure 4: Attributs Numérique

La Figure 5 regroupe les différents attributs numériques liés aux statistiques des positions.

Attribut	Position	Attribut	Position	Attribut	Position
SW	DEF	RDM	MID	RF	FWD
RWB	DEF	CDM	MID	CF	FWD
RB	DEF	LDM	MID	LF	FWD
RCB	DEF	RM	MID	RW	FWD
CB	DEF	RCM	MID	RS	FWD
LCB	DEF	CM	MID	ST	FWD
LB	DEF	LCM	MID	LS	FWD
LWB	DEF	LM	MID	LW	FWD
-	-	RAM	MID	-	-
-	-	CAM	MID	-	-
-	-	LAM	MID	-	-

Figure 5: Statistique Numérique

La figure 6 représente elle, les différents scores liés aux aptitudes du joueur.

Attribut	Description	Attribut	Description
Crossing	Valeur entre 0 et 100	Finishing	Valeur entre 0 et 100
HeadingAccuracy	Valeur entre 0 et 100	ShortPassing	Valeur entre 0 et 100
Volleys	Valeur entre 0 et 100	Dribbling	Valeur entre 0 et 100
Curve	Valeur entre 0 et 100	FKAccuracy	Valeur entre 0 et 100
LongPassing	Valeur entre 0 et 100	BallControl	Valeur entre 0 et 100
Acceleration	Valeur entre 0 et 100	SprintSpeed	Valeur entre 0 et 100
Agility	Valeur entre 0 et 100	Reactions	Valeur entre 0 et 100
Balance	Valeur entre 0 et 100	ShotPower	Valeur entre 0 et 100
Jumping	Valeur entre 0 et 100	Stamina	Valeur entre 0 et 100
Strength	Valeur entre 0 et 100	LongShots	Valeur entre 0 et 100
Aggression	Valeur entre 0 et 100	Interceptions	Valeur entre 0 et 100
Positioning	Valeur entre 0 et 100	Vision	Valeur entre 0 et 100
Penalties	Valeur entre 0 et 100	Composure	Valeur entre 0 et 100
Marking	Valeur entre 0 et 100	StandingTackle	Valeur entre 0 et 100
SlidingTackle	Valeur entre 0 et 100	GKDivng	Valeur entre 0 et 100
GKHandling	Valeur entre 0 et 100	GKKicking	Valeur entre 0 et 100
GKPositioning	Valeur entre 0 et 100	GKReflexes	Valeur entre 0 et 100
Marking	Valeur entre 0 et 100	-	-

Figure 6: Statistique Numérique

3.1.3 Catégoriel

La Figure 7 représente les différents attributs catégoriels de notre dataset.

Attribut	Description
Age	Groupe d'âge: -20,20-25,25-30,330-35,30+
Preferred Foot	Droit ou Gauche
Work Rate	Niveau: Low/Low, ..., Low/Medium, ..., High/High
Body Type	Type: Normal, Lean, Stocky, Unkown
Position	GK, DEF, MID, FWD

Figure 7: Attributs Catégoriel

3.2 Discrétisation

3.2.1 Value

Nous avons vu dans la première partie qu'il y avait plusieurs façons de discrétiser la valeur Value. En effet, les valeurs ne sont pas répartis uniformément ce qui introduit une variance plus ou moins importante entre les groupes, suivant notre discrétisation. On peut observer cela grâce au traçage de la courbe qui représente Value pour chaque joueur (Figure 8). On remarque bien que la courbe est en dent de scie, et que les écarts entre ses dents de scies sont visibles (pas de continuité).

La Figure 9 représente différent découpage possible de la valeur Value. Nous observerons par la suite le score de classification sur ces groupes. Nous remarquerons qu'au plus il y a de la variance entre les groupes, plus notre classification sera moins précise (on aura des classes non balancées).

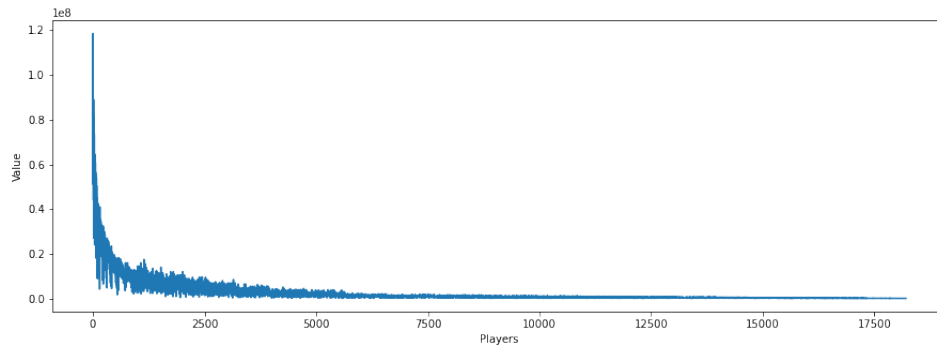


Figure 8: Value plot

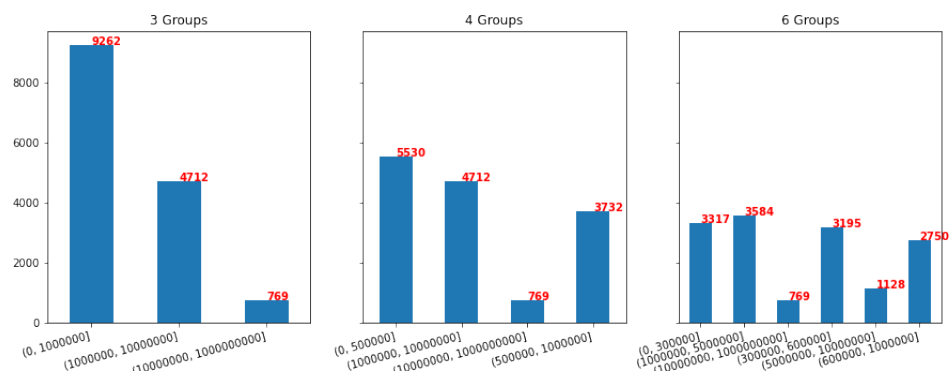


Figure 9: DValue Discretisation

3.2.2 Wage

Contrairement à Value, les valeurs de Wage sont moins dispersées. On peut l'observer grâce à la courbe (Fig 10), mais aussi à la valeur de l'écart type qui est nettement inférieur à celui de Value, **5833752** pour Value contre **22834** pour Wage.

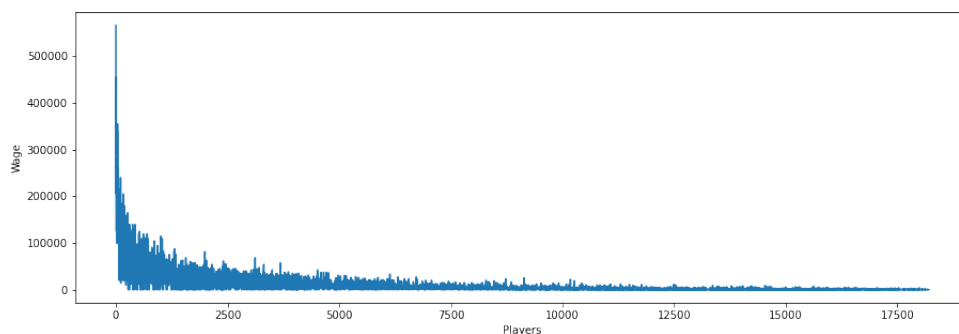


Figure 10: Wage plot

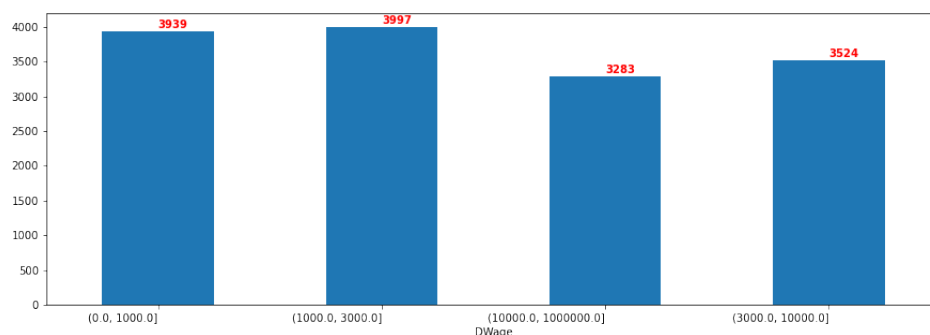


Figure 11: DWage repartition

Grâce à ce découpage (Figure 11) nous observons une bonne répartition des éléments entre les groupes, ce qui facilitera la tâche de notre classifieur.

4 Regression

Nous utiliserons un budget (temps) constant pour le reste des expériences de **5 minutes**. De plus nous utiliserons de la cross validation à 10 pour éviter l'overfitting.

4.1 Value

Dans le rapport précédent nous avons prédit la valeur de Value grâce à des classifieurs linéaires. La Figure 12 résume les différents scores obtenus. Nous avons remarqué des scores négatifs lors de la cross validation ce qui nous laisse penser que ces classifieurs n'étaient pas adaptés à notre dataset.

	R2	R2_mean_CV10	R2_std_CV10
LinearRegression	0.819823	-255.783940	684.525880
LASSO	0.819823	-255.783520	684.525638
Ridge	0.819818	-255.722425	684.361062

Figure 12: Regression Scores for Value

En utilisant AutoSklearnRegressor nous obtenons les résultats qui sont décrits dans la Figure 13.

```
In [34]: print(clf_aml.sprint_statistics())

auto-sklearn results:
Dataset name: d531bc01-7240-11ec-8fb3-b5f4799831ab
Metric: r2
Best validation score: 0.977154
Number of target algorithm runs: 8
Number of successful target algorithm runs: 2
Number of crashed target algorithm runs: 0
Number of target algorithms that exceeded the time limit: 4
Number of target algorithms that exceeded the memory limit: 2
```

Figure 13: auto-sklearn Value statistics

Nous observons que seulement 2 algorithmes ont réussi à se terminer, 4 ont été arrêtés à cause du timeout, et 2 qui ont dépassé la taille mémoire allouée. On observe aussi que le score obtenu est de 0.97 qui correspond au meilleur score R2 avec une cross validation de 10.

Le meilleur modèle obtenu ainsi qu'un résumé des paramètres est décrit par la figure suivante (Table 1). On y observe que le meilleur modèle est Gradient Boosting.

Table 1: Auto-Sklearn Value

model_id	rank	ensemble_weight	type	cost	duration	config_id	train_loss	seed	start_time	end_time	budget	status	data_preprocessors	feature_preprocessors	balancing_strategy	config_origin
5	1	1.0	gradient_boosting	0.022846	32.827184	4	0.000069	0	1.641840e+09	1.641840e+09	0.0	StatusType.SUCCESS	[]	[no_preprocessing]	None	Initial design

Le Tableau 7 (en annexe) résume les différents paramètres de chaque modèle. On peut y observer que le modèle a utilisé du OneHotEncoding sur les attributs catégoriels.

4.2 Wage

Tout come Value, nous allons prédire la valeur de Wage grâce aux mêmes classifieurs linéaires. La Figure 14 représente les différents scores. On remarque de moins bon score que pour Value.

	R2	R2_mean_CV10	R2_std_CV10
LinearRegression	0.794189	-1.285790	2.032606
LASSO	0.794198	-1.260272	1.987027
Ridge	0.794182	-1.283868	2.030639

Figure 14: Regression Scores for Wage

En utilisant AutoSklearnRegressor nous obtenons les résultats qui sont décrits dans la Figure 15.

Nous observons que seulement 2 algorithmes ont réussi à se terminer, 4 ont été arrêtés à cause du timeout, et 1 qui a dépassé la taille mémoire allouée. On observe aussi que le score obtenu est de 0.75 qui correspond au meilleur score R2 avec une cross validation de 10.

```
In [41]: print(clf_aml2.sprint_statistics())

auto-sklearn results:
Dataset name: 60eacd61-7245-11ec-95dc-adff35704611
Metric: r2
Best validation score: 0.750079
Number of target algorithm runs: 7
Number of successful target algorithm runs: 2
Number of crashed target algorithm runs: 0
Number of target algorithms that exceeded the time limit: 4
Number of target algorithms that exceeded the memory limit: 1
```

Figure 15: auto-sklearn wage statistics

Le meilleur modèle obtenu ainsi que un résumé des paramètres est décrit par la figure suivante (Table 2). On y observe que le meilleur modèle est Gradient Boosting.

Table 2: Auto-Sklearn Value

model_id	rank	ensemble_weight	type	cost	duration	config_id	train_loss	seed	start_time	end_time	budget	status	data_preprocessors	feature_preprocessors	balancing_strategy	config_origin
5	1	0.96	gradient_boosting	0.249921	29.304173	4	0.003023	0	1.641840e+09	1.641840e+09	0.0	StatusType.SUCCESS	[]	[no_preprocessing]	None	Initial design
4	2	0.04	ard_regression	1.001049	39.363339	3	1.000000	0	1.641840e+09	1.641840e+09	0.0	StatusType.SUCCESS	[]	[polynomial]	None	Initial design

Le Tableau 8 (en annexe) résume les différents paramètres de chaque modèle. On peut y observer que le modèle a utilisé du OneHotEncoding sur les attributs catégoriels.

La figure 16 représente les scores obtenus sur l'ensemble de train et de test. Ils sont nettement meilleurs que nos classifieurs linéaires.

```
In [58]: train_predictions2 = clf_aml2.predict(X_train_w)
print("Train R2 score:", metrics.r2_score(y_train_w, train_predictions2))
test_predictions2 = clf_aml2.predict(X_test_w)
print("Test R2 score:", metrics.r2_score(y_test_w, test_predictions2))

Train R2 score: 0.9898922712503859
Test R2 score: 0.8023640710253628
```

Figure 16: auto-sklearn wage pred

4.3 Overall

Tout comme Value et Wage nous allons prédire la valeur de Overall grâce aux mêmes classifieurs linéaire. La Figure 17 représente les différents scores.

	R2	R2_mean_CV10	R2_std_CV10
LinearRegression	0.922870	-5.797899	2.934559
LASSO	0.843138	-8.311402	2.550680
Ridge	0.922840	-5.796284	2.934056

Figure 17: Regression Scores Overall

En utilisant AutoSklearnRegressor nous obtenons les résultats qui sont décrits dans la Figure 18.

Nous observons que seulement 1 algorithme a réussi à se terminer, 4 on était arrêté à cause du timeout. On observe aussi que le score obtenu est de 0.99 qui correspond au meilleur score R2 avec une cross validation de 10. Le meilleur modèle obtenu ainsi qu'un résumé des paramètres est décrit par la figure suivante (Table 3). On y observe que le meilleur modèle est Gradient Boosting.

Table 3: Auto-Sklearn Overall

model_id	rank	ensemble_weight	type	cost	duration	config_id	train_loss	seed	start_time	end_time	budget	status	data_preprocessors	feature_preprocessors	balancing_strategy	config_origin
5	1	0.96	gradient_boosting	0.249921	29.304173	4	0.003023	0	1.641840e+09	1.641840e+09	0.0	StatusType.SUCCESS	[]	[no_preprocessing]	None	Initial design
4	2	0.04	ard_regression	1.001049	39.363339	3	1.000000	0	1.641840e+09	1.641840e+09	0.0	StatusType.SUCCESS	[]	[polynomial]	None	Initial design

```
In [60]: print(clf_aml3.sprint_statistics())

auto-sklearn results:
Dataset name: b1640e16-7249-11ec-95dc-adff35704611
Metric: r2
Best validation score: 0.993811
Number of target algorithm runs: 6
Number of successful target algorithm runs: 1
Number of crashed target algorithm runs: 1
Number of target algorithms that exceeded the time limit: 4
Number of target algorithms that exceeded the memory limit: 0
```

Figure 18: auto-sklearn overall statistics

Le Tableau 9 (en annexe) résume les différents paramètres de chaque modèle. On peut y observer que le modèle a utilisé du OneHotEncoding sur les attributs catégoriels.

La figure 19 représente les scores obtenus sur l'ensemble de train et de test. Ils sont nettement meilleurs que nos classifieurs linéaires.

```
In [66]: train_predictions3 = clf_aml3.predict(X_train_o)
print("Train R2 score:", metrics.r2_score(y_train_o, train_predictions3))
test_predictions3 = clf_aml3.predict(X_test_o)
print("Test R2 score:", metrics.r2_score(y_test_o, test_predictions3))

Train R2 score: 0.9993470987493016
Test R2 score: 0.9948201586947482
```

Figure 19: auto-sklearn overall pred

5 Classification

Si on reprend les différentes discrétisations qui ont été faites en divisant Value en 3, 4 et 6 groupes (Voir Figure 9) nous obtenons différents scores.

D'après les différents scores obtenus nous pouvons observer que la discrétisation en 3 groupes (Figure 20) est meilleure que celle en 4 groupes (Figure 23) qui est elle-même meilleure que la discrétisation en 6 groupes (Figure 26).

Concernant Auto-Sklearn, nous observons que nous obtenons toujours de meilleurs scores. Pour la prédiction de DValue 3 et 4, le meilleur classifieur reste Gradient Boosting. Les paramètres sont représentés par la table 11 et 12 en annexe. On remarque toujours de très bon score pour la classification. D'autre part pour la classification de Dvalue_6 Auto-Sklearn n'est pas arrivé à trouver un classifieur exploitable.

5.0.1 DValue: 3 Groups

	accuracy	f1_weighted	recall_weighted	precision_weighted	f1_weighted_mean.CV10	f1_weighted_std.CV10
LogisticRegression	0.963147	0.963052	0.963147	0.963047	0.916721	0.118629
DecisionTreeClassifier	0.968799	0.968787	0.968799	0.968801	0.883725	0.151050
RandomForestClassifier	0.962921	0.962901	0.962921	0.962938	0.860834	0.186419
GaussianNB	0.840606	0.843481	0.840606	0.850207	0.810712	0.176005
KNeighborsClassifier_3	0.867737	0.866033	0.867737	0.867085	0.825663	0.089968
KNeighborsClassifier_4	0.862311	0.856590	0.862311	0.861703	0.816515	0.102246
KNeighborsClassifier_6	0.870224	0.864807	0.870224	0.869804	0.824056	0.107767

Figure 20: Classification DValue Scores

```
In [138]: print(clf_aml5.sprint_statistics())
```

auto-sklearn results:
Dataset name: cb906652-725b-11ec-95dc-adff35704611
Metric: accuracy
Best validation score: 0.980426
Number of target algorithm runs: 6
Number of successful target algorithm runs: 1
Number of crashed target algorithm runs: 0
Number of target algorithms that exceeded the time limit: 5
Number of target algorithms that exceeded the memory limit: 0

Figure 21: auto-sklearn Dvalue 3 pred

Table 4: Auto-Sklearn Dvalue 3

modelId	rank	ensemble_weight	type	cost	duration	config_id	train_loss	seed	start_time	end_time	budget	status	data_preprocessors	feature_preprocessors	balancing_strategy	config_origin
4	1	1.0	gradient_boosting	0.019574	21.680902	3	0.002821	0	1.641850e+09	1.641850e+09	0.0	StatusType.SUCCESS	[]	[no_preprocessing]	weighting	Initial design

	accuracy	f1_macro	f1_micro	f1_weighted	recall_macro	recall_micro	recall_weighted	precision_macro	precision_micro	precision_weighted
train	0.998062	0.994703	0.998062	0.998069	0.998329	0.998062	0.998062	0.991177	0.998062	0.998091
test	0.978747	0.969795	0.978747	0.978845	0.982579	0.978747	0.978747	0.958094	0.978747	0.979167

Figure 22: Auto-Sklearn Classification DValue Scores

5.0.2 DValue: 4 Groups

	accuracy	f1_weighted	recall_weighted	precision_weighted	f1_weighted_mean_CV10	f1_weighted_std_CV10
LogisticRegression	0.905946	0.905821	0.905946	0.905859	0.843669	0.180698
DecisionTreeClassifier	0.931494	0.931511	0.931494	0.931529	0.816784	0.182516
RandomForestClassifier	0.926973	0.926982	0.926973	0.927013	0.773125	0.233648
GaussianNB	0.732308	0.737031	0.732308	0.751656	0.692522	0.189903
KNeighborsClassifier_3	0.730726	0.728015	0.730726	0.729527	0.692400	0.085815
KNeighborsClassifier_4	0.760118	0.750274	0.760118	0.755794	0.705030	0.093633
KNeighborsClassifier_6	0.766448	0.759450	0.766448	0.763807	0.717500	0.105438

Figure 23: Classification DValue Scores

```
In [147]: print(clf_aml6.sprint_statistics())
```

auto-sklearn results:
Dataset name: 7b6c1719-725c-11ec-95dc-adff35704611
Metric: accuracy
Best validation score: 0.956686
Number of target algorithm runs: 7
Number of successful target algorithm runs: 1
Number of crashed target algorithm runs: 0
Number of target algorithms that exceeded the time limit: 4
Number of target algorithms that exceeded the memory limit: 2

Figure 24: auto-sklearn Dvalue 4 pred

Table 5: Auto-Sklearn Dvalue 4

modelId	rank	ensemble_weight	type	cost	duration	config_id	train_loss	seed	start_time	end_time	budget	status	data_preprocessors	feature_preprocessors	balancing_strategy	config_origin
3	1	1.0	gradient_boosting	0.043314	36.786762	2	0.003844	0	1.641850e+09	1.641850e+09	0.0	StatusType.SUCCESS	[]	[no_preprocessing]	weighting	Initial design

	accuracy	f1_macro	f1_micro	f1_weighted	recall_macro	recall_micro	recall_weighted	precision_macro	precision_micro	precision_weighted
train	0.998643	0.996724	0.998643	0.998648	0.998943	0.998643	0.998643	0.994560	0.998643	0.998665
test	0.957495	0.955919	0.957495	0.957625	0.964537	0.957495	0.957495	0.948128	0.957495	0.958055

Figure 25: Auto-Sklearn Classification DValue Scores

5.0.3 DValue: 6 Groups

	accuracy	f1_weighted	recall_weighted	precision_weighted	f1_weighted_mean_CV10	f1_weighted_std_CV10
LogisticRegression	0.840380	0.840079	0.840380	0.840511	0.777335	0.207813
DecisionTreeClassifier	0.892607	0.892292	0.892607	0.892176	0.744778	0.255623
RandomForestClassifier	0.886050	0.886260	0.886050	0.887518	0.713415	0.270342
GaussianNB	0.627854	0.628603	0.627854	0.635709	0.583461	0.186409
KNeighborsClassifier_3	0.600271	0.592759	0.600271	0.596563	0.555849	0.083376
KNeighborsClassifier_4	0.633281	0.621264	0.633281	0.625496	0.577331	0.096768
KNeighborsClassifier_6	0.647072	0.638079	0.647072	0.642696	0.588475	0.105887

Figure 26: Classification DValue Scores

5.1 DWage

Tout comme DValue, nous allons essayer de prédire les valeurs discrètes de DWage que nous avons construite plus haut.

	accuracy	f1	recall	precision	f1_mean_CV10	f1_std_CV10
LogisticRegression	0.609767	0.612708	0.609767	0.618085	0.528905	0.262971
DecisionTreeClassifier	0.531540	0.531562	0.531540	0.531672	0.252740	0.066976
RandomForestClassifier	0.615193	0.615103	0.615193	0.615312	0.351820	0.193055
GaussianNB	0.556636	0.560099	0.556636	0.566438	0.500922	0.178237
KNeighborsClassifier_3	0.508026	0.507893	0.508026	0.517820	0.444934	0.109921
KNeighborsClassifier_4	0.539001	0.534621	0.539001	0.535107	0.468517	0.133113
KNeighborsClassifier_6	0.535835	0.533975	0.535835	0.536262	0.474865	0.141189

Figure 27: Classification DWage Scores

D'après la Figure 27, on observe de mauvais score de classification.

En utilisant AutoSklearnClassifier nous obtenons les résultats qui sont décrits dans la Figure 28.

```
In [112]: print(clf_aml4.sprint_statistics())

auto-sklearn results:
Dataset name: 4fbbf242-7256-11ec-95dc-adff35704611
Metric: accuracy
Best validation score: 0.621609
Number of target algorithm runs: 7
Number of successful target algorithm runs: 2
Number of crashed target algorithm runs: 0
Number of target algorithms that exceeded the time limit: 5
Number of target algorithms that exceeded the memory limit: 0
```

Figure 28: auto-sklearn DWage statistics

Nous observons que seulement 2 algorithmes ont réussi à se terminer, 5 on était arrêté à cause du timeout. On observe aussi que le score obtenu est de 0.62 qui correspond au meilleur score d'accuracy avec une cross validation de 10.

Le meilleur modèle obtenu ainsi qu'un résumé des paramètres est décrit par la figure suivante (Table 6). On y observe que le meilleur modèle est Gradient Boosting.

Table 6: Auto-Sklearn DWage

model_id	rank	ensemble_weight	type	cost	duration	config_id	train_loss	seed	start_time	end_time	budget	status	data_preprocessors	feature_preprocessors	balancing_strategy	config_origin
7	1	0.66	gradient_boosting	0.378391	31.091871	6	0.052466	0	1.641848e+09	1.641848e+09	0.0	StatusType.SUCCESS		[no_preprocessing]	weighting	Initial design
6	2	0.34	passive_aggressive	0.403004	9.427623	5	0.389686	0	1.641848e+09	1.641848e+09	0.0	StatusType.SUCCESS		[select_rates_classification]	weighting	Initial design

Le Tableau 10 (en annexe) résume les différents paramètres de chaque modèle. On peut y observer que le modèle a utilisé du OneHotEncoding sur les attributs catégoriels.

La figure 29 représente les scores obtenus sur l'ensemble de train et de test.

	accuracy	f1_macro	f1_micro	f1_weighted	recall_macro	recall_micro	recall_weighted	precision_macro	precision_micro	precision_weighted
train	0.935756	0.935777	0.935756	0.935793	0.937380	0.935756	0.935756	0.935046	0.935756	0.936680
test	0.617228	0.619844	0.617228	0.615417	0.622011	0.617228	0.617228	0.618231	0.617228	0.614185

Figure 29: Auto-Sklearn Classification DWage Scores

6 GridSearch

Dans cette partie nous allons utiliser GridSearchCV avec le classifieur RandomForestClassifier en variant différents paramètres:

- n_testimator : [200, 500]
- max_features : ['auto', 'sqrt', 'log2']
- max_depth : [4, 8]
- criterion : ['gini', 'entropy']

Pour une question de temps et de ressources nous nous limiterons à la prédiction de DValue_3. À la fin de l'opération nous obtenons comme meilleurs paramètres :

- n_testimator : 500
- max_features : auto
- max_depth : 8
- criterion : entropy

La figure suivante (Figure 30) résume les différents scores sur l'ensemble de train et de test. On observe de très bon score, comparable à Auto-Sklearn en prenant moins de temps (2 minutes à peu près) !

	accuracy	f1_macro	f1_micro	f1_weighted	recall_macro	recall_micro	recall_weighted	precision_macro	precision_micro	precision_weighted
train	0.975291	0.969555	0.975291	0.975305	0.965418	0.975291	0.975291	0.973882	0.975291	0.975359
test	0.952747	0.932226	0.952747	0.952657	0.922958	0.952747	0.952747	0.942356	0.952747	0.952722

Figure 30: GridSearch: RandomForest Classification DValue_3 Scores

7 RandomSearch

Nous répétons la même opération avec RandomSearch ie: avec les mêmes paramètres. Nous obtenons exactement les mêmes meilleurs paramètres que cité plus haut, ainsi que quasiment les mêmes scores que GridSearch (Voire Figure 31).

	accuracy	f1_macro	f1_micro	f1_weighted	recall_macro	recall_micro	recall_weighted	precision_macro	precision_micro	precision_weighted
train	0.975291	0.969555	0.975291	0.975305	0.965418	0.975291	0.975291	0.973882	0.975291	0.975359
test	0.952747	0.932226	0.952747	0.952657	0.922958	0.952747	0.952747	0.942356	0.952747	0.952722

Figure 31: RandomSearch: RandomForest Classification DValue_3 Scores

8 Conclusion

À travers les différentes expérimentations nous avons vu que l'autoML, en particulier Auto-Sklearn, est un outil simple et efficace. Dans la majeure partie des cas (à part pour DValue_6) Auto-Sklearn arrivé à dépasser les scores que nous avons obtenus en utilisant notre propre pipeline. Nous avons aussi vu que dans certains cas un simple GridSearch ou RandomSearch pouvait être aussi performant.

9 Annexe

Table 7: Auto-Sklearn Value parameters

[illegible]

Table 10: Auto-Sklearn DWage parameters

[illegible]

Table 11: Auto-Sklearn DValue 3 parameters

[illegible]

Table 12: Auto-Sklearn DValue 4 parameters

[illegible]