

Système Avancé de Prédition de la Performance des Candidats

Analyse Comparative d'Algorithmes de Machine Learning pour
l'Optimisation du Recrutement

Taha Bchir - Hadi chouaieb

ISET Sfax

DSI 34

Rapport de Projet de Machine Learning

Table des matières

Résumé Étendu	5
1 Introduction	6
1.1 Contexte et Motivation	6
1.2 Problématique	6
1.3 Objectifs Spécifiques	6
1.4 Cadre Théorique	6
2 Description et Analyse des Données	6
2.1 Source et Généalogie des Données	7
2.2 Analyse des Valeurs Manquantes	7
2.3 Analyse du Déséquilibre des Classes	7
3 Méthodologie Détaillée	7
3.0.1 Techniques d'Imputation	7
3.0.2 Encodage des Variables Catégorielles	8
3.0.3 Scaling et Normalisation	8
3.1 Algorithmes Implémentés : Analyse Théorique	8
3.1.1 1. Régression Logistique	8
3.1.2 2. Forêts Aléatoires (Random Forest)	9
3.1.3 3. XGBoost (Extreme Gradient Boosting)	9
3.1.4 4. LightGBM	9
3.1.5 5. Gradient Boosting Classifier	10
3.1.6 6. AdaBoost	10
3.1.7 7. SVM (Support Vector Machine)	10
3.1.8 8. Extra Trees (Extremely Randomized Trees)	11
3.1.9 9. Réseau de Neurones (MLP)	11

3.2	Optimisation des Hyperparamètres	11
3.2.1	Stratégie de Recherche	11
3.2.2	Hyperparamètres Optimisés	11
3.3	Validation et Évaluation	12
3.3.1	Stratégie de Validation Croisée	12
3.3.2	Métriques d'Évaluation	13
4	Résultats Détaillés	13
4.1	Performances Globales	13
4.2	Importance des Caractéristiques	14
4.3	Analyse des Erreurs	14
5	Implémentation Technique Avancée	14
5.1	Pipeline de Production	15
5.1.1	Code Principal d'Entraînement	15
5.2	Application Streamlit	17
5.3	Intégration Continue et Monitoring	19
6	Discussion et Analyse Critique	19
6.1	Interprétation des Résultats	19
6.1.1	Supériorité du Gradient Boosting	19
6.1.2	Importance des Caractéristiques	19
6.2	Limitations et Biais Potentiels	19
6.2.1	Limitations Techniques	20
6.2.2	Biais Éthiques Potentiels	20
6.3	Validation Business	20
6.3.1	Impact sur les Métriques RH	20
7	Conclusion et Perspectives	20
7.1	Synthèse des Contributions	20

7.1.1	Contributions Techniques	21
7.1.2	Contributions RH	21
7.2	Recommandations Stratégiques	21
7.2.1	À Court Terme (0-6 mois)	21
7.2.2	À Moyen Terme (6-18 mois)	21
7.2.3	À Long Terme (18+ mois)	21
7.3	Directions Futures de Recherche	22
7.4	Déclaration Éthique	22
Remerciements		22
Déclaration d'Intérêts		22
8 Annexes Techniques		23
8.1	Code Complet de Prétraitement	23
8.2	Configuration d'Environnement	25
8.3	Documentation de l'API	26
8.4	Plan de Déploiement	27

Liste des tableaux

1	Analyse des valeurs manquantes par variable	7
2	Répartition de la variable cible	7
3	Comparaison des techniques d'encodage testées	8
4	Stratégies d'optimisation des hyperparamètres	11
5	Détail des hyperparamètres optimisés pour chaque algorithme	12
6	Métriques d'évaluation utilisées	13
7	Performances comparées des algorithmes (Validation Croisée 5-Fold)	13
8	Importance des caractéristiques - Top 15	14
9	Système de monitoring et validation	19
10	Analyse des limitations techniques	20
11	Impact business potentiel du système	20
12	Plan de déploiement en production	27

Résumé Étendu

Ce rapport présente une étude exhaustive sur la prédiction de la performance des candidats dans le domaine des ressources humaines, utilisant des techniques avancées de machine learning. Le système développé vise à assister les professionnels RH dans leur processus de sélection en fournissant des prédictions basées sur des données historiques.

Méthodologie

Nous avons implémenté et comparé neuf algorithmes de classification différents, chacun avec ses caractéristiques et avantages spécifiques. Le pipeline de traitement comprend :

- Un nettoyage rigoureux des données pour traiter les valeurs manquantes et aberrantes
- Une transformation appropriée des variables catégorielles et numériques
- Une optimisation systématique des hyperparamètres pour chaque algorithme
- Une évaluation robuste via validation croisée stratifiée

Résultats Principaux

- Meilleur modèle : Gradient Boosting Classifier (F1-score : 0.665)
- Variables les plus importantes : Score test technique (33.9%), Années d'expérience (32.6%)
- Performance sur test : F1-score de 0.640 avec une AUC de 0.620

Applications Pratiques

Le système a été intégré dans une application Streamlit interactive permettant :

- Des prédictions en temps réel pour des candidats individuels
- Le traitement par lots de fichiers CSV
- La visualisation des facteurs d'influence sur les décisions

1 Introduction

1.1 Contexte et Motivation

Dans le paysage concurrentiel actuel du recrutement, les organisations cherchent constamment à optimiser leurs processus de sélection. L'identification des candidats qui réussiront à long terme représente un défi complexe, souvent basé sur des critères subjectifs et des jugements intuitifs. Ce projet s'inscrit dans le mouvement croissant de la "HR Analytics", qui vise à apporter une approche data-driven aux décisions de ressources humaines.

1.2 Problématique

Le recrutement traditionnel présente plusieurs limitations :

- **Subjectivité** : Les biais cognitifs influencent les décisions
- **Inefficacité** : Processus longs et coûteux
- **Risque d'erreur** : Mauvais recrutements entraînant des coûts importants
- **Manque de prédictibilité** : Difficulté à anticiper la performance future

1.3 Objectifs Spécifiques

1. Développer un modèle prédictif robuste pour identifier les candidats performants
2. Analyser l'importance relative des différents facteurs de recrutement
3. Créer une interface utilisateur intuitive pour les praticiens RH
4. Établir un cadre méthodologique reproductible pour l'analyse RH
5. Fournir des recommandations actionnables pour l'amélioration des processus

1.4 Cadre Théorique

Le projet s'appuie sur plusieurs domaines théoriques :

- **Psychologie organisationnelle** : Théories de la performance au travail
- **Statistiques multivariées** : Techniques d'analyse des données
- **Machine Learning** : Algorithmes de classification supervisée
- **Data Mining** : Extraction de connaissances à partir de données

2 Description et Analyse des Données

2.1 Source et Généalogie des Données

Le jeu de données utilisé a été généré synthétiquement pour reproduire les caractéristiques d'un ensemble de données RH réaliste. La génération a respecté les contraintes suivantes :

- Distribution d'âge réaliste (loi normale tronquée entre 18 et 70 ans)
- Corrélations entre variables cohérentes avec la littérature
- Introduction de bruits réalistes (valeurs manquantes, aberrantes)
- Déséquilibre modéré des classes (53% vs 47%)

2.2 Analyse des Valeurs Manquantes

TABLE 1 – Analyse des valeurs manquantes par variable

tableheaderVariable	tableheaderNbre Manquants	tableheaderPourcentage	tableheaderType	tab
score_test_technique	45	4.5%	Numérique	
score_softskills	38	3.8%	Numérique	
secteur_précédent	67	6.7%	Catégorique	
langues_parlées	22	2.2%	Numérique	
mobilité	15	1.5%	Binaire	

2.3 Analyse du Déséquilibre des Classes

TABLE 2 – Répartition de la variable cible

tableheaderClasse	tableheaderEffectif	tableheaderPourcentage	tableheaderDescription
Non Performant	470	47.0%	Candidats n'ayant pas satisfait aux
Performant	530	53.0%	Candidats ayant réussi leur intégration
Total	1000	100%	Ensemble complet

3 Méthodologie Détaillée

3.0.1 Techniques d’Imputation

Nous avons utilisé plusieurs stratégies d'imputation :

- **Imputation par médiane** : Pour les variables numériques avec valeurs aberrantes

- **Imputation par KNN** : Testé pour les variables corrélées
- **Imputation par mode** : Pour les variables catégorielles
- **Indicateurs de valeurs manquantes** : Création de variables binaires supplémentaires

3.0.2 Encodage des Variables Catégorielles

TABLE 3 – Comparaison des techniques d'encodage testées

tableheaderMéthode	tableheaderDescription	tableheaderPerformance
One-Hot Encoding	Création de variables binaires pour chaque catégorie	
Label Encoding	Transformation en valeurs numériques séquentielles	
Target Encoding	Encodage basé sur la moyenne de la variable cible	
Binary Encoding	Conversion en code binaire puis séparation	
Choix final	One-Hot Encoding	avec
		drop_first=True

3.0.3 Scaling et Normalisation

Trois méthodes ont été comparées :

- **StandardScaler** : Standardisation Z-score (moyenne=0, écart-type=1)
- **MinMaxScaler** : Normalisation entre 0 et 1
- **RobustScaler** : Utilisation de médiane et IQR, robuste aux outliers

3.1 Algorithmes Implémentés : Analyse Théorique

3.1.1 1. Régression Logistique

Principe : Modèle linéaire qui modélise la probabilité d'appartenance à une classe via la fonction logistique.

Avantages :

- Interprétabilité des coefficients
- Rapidité d'entraînement et de prédiction
- Faible risque d'overfitting avec régularisation

Limitations :

- Hypothèse de linéarité peu réaliste
- Sensible aux variables corrélées
- Performance limitée sur relations complexes

3.1.2 2. Forêts Aléatoires (Random Forest)

Principe : Ensemble d'arbres de décision indépendants, chacun entraîné sur un sous-échantillon bootstrap des données avec sélection aléatoire de features.

Avantages :

- Robustesse au bruit et aux outliers
- Estimation naturelle de l'importance des variables
- Pas besoin de scaling des features
- Bonne performance sur données non linéaires

Limitations :

- Modèle "black-box" moins interprétable
- Computationnally intensive pour grands datasets
- Peut overfitter sur données bruitées

3.1.3 3. XGBoost (Extreme Gradient Boosting)

Principe : Algorithme de boosting gradient optimisé avec régularisation L1/L2 et techniques avancées de gestion des missing values.

Avantages :

- Excellentes performances compétitives
- Gestion native des valeurs manquantes
- Régularisation intégrée contre l'overfitting
- Optimisation parallèle efficace

Limitations :

- Sensible aux hyperparamètres
- Temps d'entraînement potentiellement long
- Requiert un tuning minutieux

3.1.4 4. LightGBM

Principe : Framework de gradient boosting utilisant l'algorithme Gradient-based One-Side Sampling (GOSS) et Exclusive Feature Bundling (EFB).

Avantages :

- Entraînement extrêmement rapide
- Faible consommation mémoire
- Bonne performance sur grands datasets
- Supporte les données catégorielles natives

Limitations :

- Peut overfitter sur petits datasets
- Sensible aux hyperparamètres
- Moins robuste aux outliers que Random Forest

3.1.5 5. Gradient Boosting Classifier

Principe : Algorithme séquentiel construisant un ensemble d'arbres faibles, où chaque nouvel arbre corrige les erreurs des précédents.

Avantages :

- Excellente performance prédictive
- Flexibilité avec différentes fonctions de perte
- Peut modéliser des relations complexes
- Interprétabilité via feature importance

Limitations :

- Sensible au bruit dans les données
- Temps d'entraînement significatif
- Requiert un tuning soigné

3.1.6 6. AdaBoost

Principe : Adaptive Boosting, ajuste les poids des instances mal classées à chaque itération.

Avantages :

- Simple et efficace
- Réduit à la fois bias et variance
- Peu d'hyperparamètres à tuner

3.1.7 7. SVM (Support Vector Machine)

Principe : Recherche de l'hyperplan optimal séparant les classes avec marge maximale.

Avantages :

- Efficace en haute dimension
- Bonne généralisation
- Contrôle explicite de la complexité

3.1.8 8. Extra Trees (Extremely Randomized Trees)

Principe : Variante de Random Forest avec sélection aléatoire des seuils de coupure.

Avantages :

- Plus rapide que Random Forest
- Réduction de la variance
- Moins sensible aux hyperparamètres

3.1.9 9. Réseau de Neurones (MLP)

Principe : Perceptron multicouches avec fonctions d'activation non linéaires.

3.2 Optimisation des Hyperparamètres

3.2.1 Stratégie de Recherche

TABLE 4 – Stratégies d'optimisation des hyperparamètres

tableheaderMéthode	tableheaderDescription	tableheaderUtilisation
Grid Search	Exploration exhaustive de toutes les combinaisons	Modèles simples
Random Search	Échantillonnage aléatoire dans l'espace des paramètres	Modèles complexes
Bayesian Optimization	Optimisation séquentielle basée sur les performances passées	XGBoost, LightGBM
Choix	Grid Search pour comparaison équitable	Tous modèles

3.2.2 Hyperparamètres Optimisés

TABLE 5 – Détail des hyperparamètres optimisés pour chaque algorithme

tableheaderAlgorithme	tableheaderHyperparamètres	tableheaderValeurs Testées
Régression Logistique	C, penalty, solver	C : [0.01, 0.1, 1, 10, 100]
tableheader Random Forest	n_estimators, max_depth, min_samples_split	n_estimators : [100, 200, 300]
XGBoost	learning_rate, max_depth, n_estimators	learning_rate : [0.01, 0.05, 0.1]
tableheader LightGBM	num_leaves, learning_rate, feature_fraction	num_leaves : [31, 63, 127]
Gradient Boosting	n_estimators, learning_rate, max_depth	n_estimators : [100, 200, 300]
tableheader AdaBoost	n_estimators, learning_rate	n_estimators : [50, 100, 200]
SVM	C, kernel, gamma	C : [0.1, 1, 10], kernel : ['linear', 'rbf']
tableheader Extra Trees	n_estimators, max_features	n_estimators : [100, 200, 300]
MLP	hidden_layer_sizes, activation, alpha	hidden_layer_sizes : [(50,), (100,), (50, 100)], activation : ['relu', 'tanh'], alpha : [0.001, 0.01, 0.1, 1, 10]

3.3 Validation et Évaluation

3.3.1 Stratégie de Validation Croisée

- **K-Fold Stratifié** : 5 folds avec préservation de la distribution des classes
- **Train-Test Split** : 70% entraînement, 30% test
- **Validation Hold-out** : Séparation temporelle simulée

3.3.2 Métriques d'Évaluation

TABLE 6 – Métriques d'évaluation utilisées

tableheaderMétrique	tableheaderFormule	tableheaderInterprétation
Accuracy	$\frac{TP+TN}{TP+TN+FP+FN}$	Performance globale
Precision	$\frac{TP}{TP+FP}$	Qualité des prédictions positives
Recall	$\frac{TP}{TP+FN}$	Capacité à détecter les positifs
F1-Score	$2 \times \frac{Precision \times Recall}{Precision + Recall}$	Moyenne harmonique
AUC-ROC	$\int_0^1 TPR(FPR)dFPR$	Performance globale à tous les seuils
Métrique Principale	F1-Score	Équilibre précision/rappel

4 Résultats Détaillés

4.1 Performances Globales

TABLE 7 – Performances comparées des algorithmes (Validation Croisée 5-Fold)

tableheaderAlgorithme	tableheaderF1-Score	tableheaderAccuracy	tableheaderPrecision
gradient boosting	0.665 ± 0.015	0.652 ± 0.012	0.647 ± 0.018
AdaBoost	0.662 ± 0.014	0.648 ± 0.013	0.643 ± 0.017
XGBoost	0.660 ± 0.016	0.645 ± 0.014	0.640 ± 0.019
Random Forest	0.657 ± 0.013	0.642 ± 0.011	0.638 ± 0.016
LightGBM	0.655 ± 0.017	0.640 ± 0.015	0.635 ± 0.020
Extra Trees	0.648 ± 0.014	0.633 ± 0.012	0.629 ± 0.017
SVM	0.642 ± 0.012	0.627 ± 0.010	0.623 ± 0.015
MLP	0.635 ± 0.018	0.620 ± 0.016	0.616 ± 0.021
Régression Logistique	0.630 ± 0.011	0.615 ± 0.009	0.611 ± 0.014

4.2 Importance des Caractéristiques

TABLE 8 – Importance des caractéristiques - Top 15

tableheaderRang	tableheaderCaractéristique	tableheaderImportance (%)	tableheaderInterprétation
1	Score test technique	33.9	Compétence technique dominant
2	Années d'expérience	32.6	Expérience professionnelle cond. facteur
3	Âge	7.8	Maturité, étape de carrière
4	Score softskills	7.0	Compétences relationnelles émotionnelles
5	Langues parlées	5.2	Compétences linguistiques nationales
6	Mobilité	4.7	Flexibilité géographique
7	Disponibilité immédiate	4.5	Réactivité organisationnelle
8	Niveau études (Bac+5)	1.2	Éducation supérieur
9	Spécialité (Informatique)	1.1	Domaine technique spécialisé
10	Secteur précédent (Tech)	1.0	Expérience sectorielle
11	Niveau études (Doctorat)	0.8	Éducation avancée
12	Spécialité (Data Science)	0.7	Domaine émergent
13	Secteur précédent (Consulting)	0.5	Expérience conseil
14	Spécialité (Finance)	0.4	Domaine financier
15	Mobilité (Non)	0.2	Limitation géographique

4.3 Analyse des Erreurs

tableheaderType Erreur	tableheaderEffectif	tableheaderPourcentage	tableheaderCoût Relatif
Faux Positifs	85	28.3%	Élevé
Faux Négatifs	72	24.0%	Modéré
Vrais Positifs	183	61.0%	-
Vrais Négatifs	160	53.3%	-
Total Erreurs	157	52.3%	-

5 Implémentation Technique Avancée

5.1 Pipeline de Production

5.1.1 Code Principal d'Entraînement

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split, GridSearchCV
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.ensemble import GradientBoostingClassifier
6 import pickle
7 import json
8
9 class HRPredictionPipeline:
10     def __init__(self, model_path='models/'):
11         self.model_path = model_path
12         self.scaler = StandardScaler()
13         self.model = None
14         self.feature_names = None
15
16     def load_and_preprocess(self, filepath):
17         """Chargement et prétraitement des données"""
18         df = pd.read_csv(filepath)
19
20         # Nettoyage des données
21         df_clean = self._clean_data(df)
22
23         # Encodage des variables catégorielles
24         df_encoded = self._encode_categorical(df_clean)
25
26         # Séparation features/target
27         X, y = self._prepare_features_target(df_encoded)
28
29         return X, y, df_encoded
30
31     def _clean_data(self, df):
32         """Tâches de nettoyage"""
33         # Filtrage
34         df = df[(df['ge'] >= 18) & (df['ge'] <= 70)]
35
36         # Imputation valeurs manquantes
37         numeric_cols = df.select_dtypes(include=np.number).columns
38         for col in numeric_cols:
39             if df[col].isnull().any():
40                 df[col].fillna(df[col].median(), inplace=True)
41
42         return df
43
44     def train_model(self, X_train, y_train):
45         """Entraînement avec optimisation hyperparamètres"""
46         # Définition de la grille de recherche
47         param_grid = {

```

```

48         'n_estimators': [100, 200, 300],
49         'max_depth': [3, 5, 7],
50         'learning_rate': [0.01, 0.1, 0.2],
51         'subsample': [0.8, 0.9, 1.0]
52     }
53
54     # Initialisation mod le
55     gb_model = GradientBoostingClassifier(random_state=42)
56
57     # Recherche grille
58     grid_search = GridSearchCV(
59         gb_model,
60         param_grid,
61         cv=5,
62         scoring='f1',
63         n_jobs=-1,
64         verbose=1
65     )
66
67     # Entrainement
68     grid_search.fit(X_train, y_train)
69
70     # Recuperation meilleur mod le
71     self.model = grid_search.best_estimator_
72     self.best_params = grid_search.best_params_
73
74     return self.model
75
76 def save_model(self, filename='hr_gradient_boosting.pkl'):
77     """Sauvegarde du mod le et artefacts"""
78     model_data = {
79         'model': self.model,
80         'scaler': self.scaler,
81         'feature_names': self.feature_names,
82         'best_params': self.best_params,
83         'version': '2.0',
84         'date': pd.Timestamp.now().isoformat()
85     }
86
87     with open(f'{self.model_path}{filename}', 'wb') as f:
88         pickle.dump(model_data, f)
89
90     # Sauvegarde de donnees
91     metadata = {
92         'model_type': 'GradientBoostingClassifier',
93         'performance': self.evaluate_model(),
94         'features_importance': self.get_feature_importance(),
95         'training_date': pd.Timestamp.now().isoformat()
96     }
97
98     with open(f'{self.model_path}metadata.json', 'w') as f:
99         json.dump(metadata, f, indent=4)

```

Listing 1 – Pipeline complet d'entraînement

5.2 Application Streamlit

```
1 import streamlit as st
2 import pandas as pd
3 import numpy as np
4 import pickle
5 import plotly.graph_objects as go
6 from plotly.subplots import make_subplots
7
8 class HRPredictionApp:
9     def __init__(self):
10         self.load_model()
11         self.setup_page()
12
13     def setup_page(self):
14         st.set_page_config(
15             page_title="Système Prédiction Performance RH",
16             page_icon="",
17             layout="wide"
18         )
19
20         st.title("Système Intelligent de Prédiction RH")
21         st.markdown("""
22             Cette application utilise un modèle de Machine Learning pour
23             prédire
24             la performance des candidats après 6 mois d'intégration.
25             """
26         )
27
28     def load_model(self):
29         """Chargement du modèle entraîné"""
30         try:
31             with open('models/hr_gradient_boosting.pkl', 'rb') as f:
32                 model_data = pickle.load(f)
33
34                 self.model = model_data['model']
35                 self.scaler = model_data['scaler']
36                 self.feature_names = model_data['feature_names']
37
38         except FileNotFoundError:
39             st.error("Modèle non trouvé. Veuillez entraîner le
40                     modèle d'abord.")
41             st.stop()
42
43     def display_prediction_dashboard(self):
44         """Interface principale"""
```

```

42     tab1, tab2, tab3 = st.tabs([
43         "Prédiction Individuelle",
44         "Prédiction par Lot",
45         "Analyse du Modèle"
46     ])
47
48     with tab1:
49         self.single_prediction_interface()
50
51     with tab2:
52         self.batch_prediction_interface()
53
54     with tab3:
55         self.model_analysis_interface()
56
57 def single_prediction_interface(self):
58     """Interface pour prédiction unique"""
59     col1, col2 = st.columns([2, 1])
60
61     with col1:
62         st.header("Saisie des Informations du Candidat")
63
64         # Formulaire interactif
65         with st.form("candidate_form"):
66             age = st.slider("âge", 18, 70, 30)
67             experience = st.slider("Années d'expérience", 0, 50,
68                                    5)
69             tech_score = st.slider("Score test technique", 0, 100,
70                                    70)
71             soft_score = st.slider("Score soft skills", 0, 100, 65)
72
73             submitted = st.form_submit_button("Prédire la
74                                         Performance")
75
76             if submitted:
77                 # Préparation des données
78                 candidate_data = self.prepare_candidate_data(
79                     age, experience, tech_score, soft_score
80                 )
81
82                 # Prédiction
83                 prediction, probability = self.predict(
84                     candidate_data)
85
86                 # Affichage résultats
87                 with col2:
88                     self.display_prediction_results(
89                         prediction, probability, candidate_data
90                     )

```

Listing 2 – Application Streamlit interactive

5.3 Intégration Continue et Monitoring

TABLE 9 – Système de monitoring et validation

Composant	Description	Fréquence
Validation des données	Vérification format et plages des inputs	En temps réel
Dérive des données	Monitoring distribution features vs entraînement	Quotidien
Performance modèle	Tracking métriques sur nouvelles données	Hebdomadaire
Business metrics	Impact sur taux de rétention et performance	Mensuel
Ré-entraînement	Mise à jour modèle avec nouvelles données	Trimestriel

6 Discussion et Analyse Critique

6.1 Interprétation des Résultats

6.1.1 Supériorité du Gradient Boosting

La supériorité du Gradient Boosting Classifier peut s'expliquer par plusieurs facteurs :

- **Capture de relations non linéaires** : Capacité à modéliser des interactions complexes entre variables
- **Robustesse au bruit** : Approche séquentielle réduisant l'impact des outliers
- **Équilibre biais-variance** : Boosting réduit le biais tout en contrôlant la variance
- **Flexibilité** : Support de différentes fonctions de perte et paramètres

6.1.2 Importance des Caractéristiques

L'analyse d'importance révèle plusieurs insights RH importants :

1. **Dominance des compétences techniques** : Le score technique représente 34% de l'importance
2. **Valeur de l'expérience** : Les années d'expérience sont presque aussi importantes
3. **Facteurs soft sous-estimés** : Les softskills ont une importance modérée (7%)
4. **Caractéristiques démographiques** : L'âge a une influence significative

6.2 Limitations et Biais Potentiels

6.2.1 Limitations Techniques

TABLE 10 – Analyse des limitations techniques

tableheaderLimitation	tableheaderImpact	tableheaderGravité
Taille échantillon limitée	Risque d'overfitting, généralisation limitée	Élevée
Variables incomplètes	Facteurs importants potentiellement omis	Moyenne
Bruit dans les données	Performance sous-optimale	Faible
Variables corrélées	Instabilité des coefficients	Moyenne
Déséquilibre classes	Biais en faveur de la classe majoritaire	Faible

6.2.2 Biais Éthiques Potentiels

- **Biais d'âge** : Le modèle pourrait discriminer les candidats plus âgés
- **Biais éducatif** : Privilège des diplômes supérieurs
- **Biais sectoriel** : Avantage pour certains secteurs d'activité
- **Biais de mesure** : Les scores tests pourraient être biaisés

6.3 Validation Business

6.3.1 Impact sur les Métriques RH

TABLE 11 – Impact business potentiel du système

tableheaderMétrique	tableheaderAvant	tableheaderAprès (Est.)	tableheaderImpact
Taux succès recrutement	65%	75%	+10 points
Cout par recrutement	15k€	12k€	-20%
Délai intégration	6 mois	5 mois	-17%
Rétention 1 an	80%	85%	+5 points
Satisfaction manager	70%	80%	+10 points

7 Conclusion et Perspectives

7.1 Synthèse des Contributions

Ce projet apporte plusieurs contributions significatives :

7.1.1 Contributions Techniques

- Implémentation et comparaison systématique de 9 algorithmes ML
- Développement d'un pipeline de production robuste
- Création d'une application interactive pour utilisateurs finaux
- Méthodologie de validation rigoureuse

7.1.2 Contributions RH

- Identification des facteurs clés de performance
- Outil d'aide à la décision objectif
- Réduction des biais de recrutement
- Optimisation des processus de sélection

7.2 Recommandations Stratégiques

7.2.1 À Court Terme (0-6 mois)

1. Déployer le système en phase pilote sur un département
2. Former les recruteurs à l'interprétation des résultats
3. Établir un processus de feedback pour amélioration continue
4. Monitorer les impacts sur la diversité et l'inclusion

7.2.2 À Moyen Terme (6-18 mois)

1. Intégrer des sources de données supplémentaires
2. Développer des modèles spécifiques par métier
3. Implémenter un système de recommandation de formations
4. Étendre à la prédiction de la rétention

7.2.3 À Long Terme (18+ mois)

1. Crée un écosystème complet de talent analytics
2. Intégrer l'IA explicative pour la transparence
3. Développer des modèles prédictifs pour le développement de carrière
4. Contribuer à la recherche académique en HR analytics

7.3 Directions Futures de Recherche

- **Apprentissage par transfert** : Adaptation de modèles entre organisations
- **Modèles temporels** : Prédiction de l'évolution de la performance
- **Analyse de réseaux** : Influence des relations interpersonnelles
- **NLP avancé** : Analyse sémantique des CV et entretiens
- **Apprentissage par renforcement** : Optimisation dynamique des processus

7.4 Déclaration Éthique

Ce système doit être utilisé selon les principes suivants :

- **Transparence** : Les candidats doivent être informés de l'utilisation d'IA
- **Contrôle humain** : Les prédictions doivent rester des recommandations
- **Équité** : Surveillance continue des biais potentiels
- **Vie privée** : Protection des données personnelles
- **Responsabilité** : Désignation claire de la responsabilité décisionnelle

Remerciements

Je tiens à exprimer ma gratitude aux experts RH qui ont contribué à la validation des hypothèses de ce projet, ainsi qu'à la communauté open-source pour les outils et bibliothèques qui ont rendu ce travail possible.

Déclaration d'Intérêts

L'auteur déclare n'avoir aucun conflit d'intérêts financier ou personnel qui pourrait influencer les résultats présentés dans ce rapport.

8 Annexes Techniques

8.1 Code Complet de Prétraitement

```

1 """
2 Script de prétraitement complet pour les données RH
3 Auteur: Taha Bchir
4 Date: \today
5 """
6
7 import pandas as pd
8 import numpy as np
9 from sklearn.preprocessing import StandardScaler, OneHotEncoder
10 from sklearn.impute import KNNImputer
11 import warnings
12 warnings.filterwarnings('ignore')
13
14 class HRDataPreprocessor:
15     """Classe pour le prétraitement des données RH"""
16
17     def __init__(self, config=None):
18         self.config = config or self.default_config()
19         self.scaler = StandardScaler()
20         self.imputer = KNNImputer(n_neighbors=5)
21         self.encoder = OneHotEncoder(drop='first', sparse=False)
22         self.feature_names = None
23
24     @staticmethod
25     def default_config():
26         return {
27             'age_range': (18, 70),
28             'score_range': (0, 100),
29             'imputation_strategy': 'median',
30             'encoding_method': 'onehot',
31             'scaling_method': 'standard'
32         }
33
34     def full_pipeline(self, df):
35         """Exécute le pipeline complet"""
36
37         # 1. Nettoyage initial
38         df_clean = self.clean_data(df)
39
40         # 2. Traitement valeurs manquantes
41         df_imputed = self.handle_missing_values(df_clean)
42
43         # 3. Encodage variables catégorielles
44         df_encoded = self.encode_categorical(df_imputed)
45
46         # 4. Scaling des variables numériques

```

```

47     df_scaled = self.scale_features(df_encoded)
48
49     # 5. Feature engineering
50     df_engineered = self.feature_engineering(df_scaled)
51
52     # 6. Selection features
53     df_final = self.select_features(df_engineered)
54
55     return df_final
56
57 def clean_data(self, df):
58     """ tapes de nettoyage """
59
60     # Suppression des doublons
61     df = df.drop_duplicates()
62
63     # Filtrage des ges invalides
64     min_age, max_age = self.config['age_range']
65     df = df[(df['ge'] >= min_age) & (df['ge'] <= max_age)]
66
67     # Correction des scores hors limites
68     min_score, max_score = self.config['score_range']
69     df['score_test_technique'] = df['score_test_technique'].clip(
70         min_score, max_score)
71     df['score_softskills'] = df['score_softskills'].clip(min_score,
72         max_score)
73
74     # Standardisation du texte
75     text_cols = ['niveau_tudes', 'sp_cialit', 'secteur_pr_c_dent']
76     for col in text_cols:
77         if col in df.columns:
78             df[col] = df[col].str.lower().str.strip()
79
80     return df
81
82 def handle_missing_values(self, df):
83     """Gestion des valeurs manquantes"""
84
85     # Strat gies par type de variable
86     strategies = {
87         'numeric': self.config['imputation_strategy'],
88         'categorical': 'mode',
89         'binary': 'mode'
90     }
91
92     for col in df.columns:
93         if df[col].isnull().any():
94             col_type = self.get_column_type(df[col])
95             strategy = strategies[col_type]
96
97             if strategy == 'median':

```

```

96         df[col].fillna(df[col].median(), inplace=True)
97     elif strategy == 'mean':
98         df[col].fillna(df[col].mean(), inplace=True)
99     elif strategy == 'mode':
100        df[col].fillna(df[col].mode()[0], inplace=True)
101    elif strategy == 'knn':
102        # Imputation KNN pour variables corr l es
103        pass
104
105    return df
106
107 def get_column_type(self, series):
108     """D termine le type de colonne"""
109     if series.dtype in ['int64', 'float64']:
110         if series.nunique() == 2:
111             return 'binary'
112         else:
113             return 'numeric'
114     else:
115         return 'categorical'

```

Listing 3 – Script complet de prétraitement des données

8.2 Configuration d’Environnement

```

1  # Core Data Science
2  numpy==1.21.0
3  pandas==1.3.0
4  scikit-learn==0.24.2
5  scipy==1.7.0
6
7  # Machine Learning Algorithms
8  xgboost==1.4.2
9  lightgbm==3.2.1
10 catboost==0.26.1
11
12 # Visualization
13 matplotlib==3.4.2
14 seaborn==0.11.1
15 plotly==5.3.1
16
17 # Web Application
18 streamlit==1.8.1
19 flask==2.0.1
20 dash==2.0.0
21
22 # Utilities
23 joblib==1.0.1
24 imbalanced-learn==0.8.0

```

```

25 optuna==2.10.0
26 shap==0.39.0
27
28 # Documentation
29 jupyter==1.0.0
30 notebook==6.4.0
31 pytest==6.2.5

```

Listing 4 – Fichier requirements.txt

8.3 Documentation de l'API

```

1 """
2 API Documentation - HR Prediction System
3 Base URL: https://api.hr-predict.com/v1
4 """
5
6 # Endpoints disponibles
7 ENDPOINTS = {
8     'single_prediction': '/predict',
9     'batch_prediction': '/predict/batch',
10    'model_info': '/model/info',
11    'model_retrain': '/model/retrain',
12    'monitoring': '/monitoring/metrics'
13}
14
15 # Exemple de requête
16 import requests
17 import json
18
19 def predict_candidate(candidate_data):
20     """Prédiction pour un candidat unique"""
21
22     url = "https://api.hr-predict.com/v1/predict"
23     headers = {
24         'Content-Type': 'application/json',
25         'Authorization': 'Bearer YOUR_API_KEY'
26     }
27
28     response = requests.post(
29         url,
30         headers=headers,
31         data=json.dumps(candidate_data)
32     )
33
34     if response.status_code == 200:
35         return response.json()
36     else:
37         raise Exception(f"API Error: {response.status_code}")

```

```

38
39 # Rponse type
40 SAMPLE_RESPONSE = {
41     "success": True,
42     "prediction": 1,
43     "probability": 0.78,
44     "confidence_interval": [0.72, 0.84],
45     "key_factors": [
46         {"feature": "score_test_technique", "importance": 0.34},
47         {"feature": "annees_experience", "importance": 0.27}
48     ],
49     "model_version": "2.0",
50     "timestamp": "2024-01-15T10:30:00Z"
51 }
```

Listing 5 – Documentation API REST

8.4 Plan de Déploiement

TABLE 12 – Plan de déploiement en production

Phase	Durée	Responsable	Activités
Pré-déploiement	2 semaines	Équipe Data	Tests de charge, validation données
Déploiement pilote	1 mois	Équipe RH	Formation utilisateurs, collecte feedback
Déploiement complet	3 mois	Équipe IT	Intégration systèmes existants
Monitoring	Continue	Équipe Ops	Surveillance performance, maintenance
Optimisation	Continue	Équipe Data	Amélioration continue modèle