
Rapport de Projet Python

Analyse et Implémentation d'Algorithmes de Machine Learning

Sujet : HR Analytics & Prédiction de Satisfaction

Réalisé par :

Riad ZAID

Youssef SADIQUI

M Saad OUSSAMA

Filière :

Génie Informatique

Encadré par :

Pr. EL AZHARI

Département Informatique

Table des matières

1	Introduction et Configuration	1
1.1	Environnement Technique	1
1.2	Implémentation et Chargement des Données	1
2	Phase 1 : Structuration des Données (Data Engineering)	2
2.1	Inspection et Exploration Initiale	2
2.2	Nettoyage et Optimisation des Données	3
2.3	Transformation et Feature Engineering	4
2.4	Finalisation et Exportation	4
3	Phase 2 : Analyse Exploratoire (EDA)	5
3.1	Distribution des Âges	5
3.2	Répartition par Département	6
3.3	Corrélation : Expérience vs Revenu	7
3.4	Analyse Multivariée (3D)	9
3.5	Analyse Multivariée des Risques (5D)	9
3.6	Analyse de l'Équité Salariale (Seaborn)	12
3.7	Comparaison par Département (Faceting)	13
3.8	Progression de Carrière et Grille Salariale	14
4	Phase 3 : Modélisation Prédictive (Régression)	16
4.1	Phase 3.1 : Régression Linéaire Simple (Modèle Baseline)	16
4.1.1	Objectif et Paramètres	16
4.1.2	Implémentation Python	16
4.1.3	Analyse des Résultats	17
4.2	Phase 3.2 : Régression Linéaire Multiple (Modèle Avancé)	17
4.2.1	Enrichissement du Modèle	17
4.2.2	Implémentation	17
4.3	Phase 3.3 : Analyse de Tendance (Modèle ARIMA)	18
4.3.1	Méthodologie	18
4.3.2	Implémentation ARIMA	18
4.3.3	Performance et Interprétation	18
4.4	Bilan et Benchmark des Modèles Prédictifs	19
4.4.1	Tableau Récapitulatif	19
4.4.2	Conclusion du Benchmark	19
5	Phase 4 : Analyse Qualitative (La "Boussole du Bien-être")	19
5.1	Préparation des Données : Binarisation	19
5.2	Extraction des "Règles d'Or"	20
5.3	Résultats : Les Micro-Climats de Satisfaction	20

5.4	Application : Le Moteur de Scoring "Bien-être"	21
5.4.1	Logique de l'Algorithme	21
5.5	Déploiement et Sérialisation	22
6	Conclusion Générale	22
7	PHASE 5 Algorithmes Clustrings+Classifications	22
7.1	Algorithme de Clustering K-Means	22
7.2	Étapes de l'algorithme	23
7.3	Formulation Mathématique	23
8	Détermination du nombre de clusters (Méthode du Coude)	23
9	Visualisation des Segments RH	23
9.1	Interprétation des KPI	23
10	Application du Modèle K-Means et Analyse des Profils	24
10.1	Exécution du clustering final	24
10.2	Caractérisation des Groupes (KPI)	24
10.3	Évaluation de la qualité : Le Score de Silhouette	25
10.4	Implémentation du K-Means Graphiquement	25
10.5	Analyse de la Segmentation	26
11	Clustering Hiérarchique (CAH)	27
11.1	Le Dendrogramme	27
11.2	Interprétation et Validation	28
11.3	Implémentation du Modèle CAH	28
11.3.1	Justification Mathématique de la Ligne de Coupure	29
12	Validation et Choix du Meilleur Modèle	29
12.1	Définition du Score de Silhouette	29
12.2	Comparaison K-Means vs CAH	29
12.3	Résultats et Décision Finale	30
13	Classification Supervisée (Arbre de Décision)	31
13.1	Implémentation du Modèle	31
13.2	Évaluation des Performances	31
13.3	Visualisation de l'Arbre (Les Règles RH)	32
13.4	Interprétation des Règles	33
14	Classification K-NN (K-Nearest Neighbors)	34
14.1	Optimisation de l'Hyperparamètre K	34
14.2	Interprétation et Entraînement Final	35

15 Classification Probabiliste (Naïve Bayes)	36
15.1 Implémentation (Gaussian Naïve Bayes)	36
15.2 Analyse de la Matrice de Confusion	37
15.3 Comparatif et Conclusion Technique	38
16 Random Forest (L'Intelligence Collective)	38
16.1 Implémentation et Entraînement	38
16.2 Analyse des Facteurs Clés (Feature Importance)	39
16.3 Performance et Matrice de Confusion	40
16.4 Simulation : Test sur un Profil Fictif	41
16.5 Sauvegarde et Déploiement	41
17 Régression Logistique (L'Approche Explicative)	42
17.1 Prérequis : La Normalisation	42
17.2 Le Pouvoir des Coefficients (Facteurs d'Influence)	42
17.3 Performance et Limites	43
17.4 Simulation d'un Cas Concret	44
18 Support Vector Machine (L'Approche Géométrique)	45
18.1 Configuration et "Kernel Trick"	45
18.2 Analyse Avancée : La Courbe Précision-Rappel	45
18.3 Performance Globale	47
18.4 Simulation et Conclusion Technique	47
19 Benchmark et Choix Final du Modèle	48
19.1 Tableau Récapitulatif des Performances	48
19.2 Analyse Graphique : Précision vs Rappel	48
19.3 Verdict et Recommandation Stratégique	49

Table des figures

1	Répartition des employés par tranche d'âge	6
2	Effectifs par département	7
3	Corrélation entre Années d'Expérience et Salaire Mensuel	8
4	Projection 3D : Âge, Expérience et Revenu	10
5	Visualisation des profils à risque (Départ en Rouge)	11
6	Comparaison des revenus Hommes vs Femmes	12
7	Dynamique salariale et départs par département	13
8	Évolution salariale par niveau hiérarchique	15
9	Classement des facteurs les plus présents dans la Haute Satisfaction	21
10	Visualisation des segments employés (Âge vs Revenu)	24
11	Visualisation des 3 segments d'employés détectés	26
12	Dendrogramme déterminant le nombre optimal de clusters	28
13	Matrice de confusion : Capacité à prédire les départs	32
14	Visualisation des règles de décision (Rotation 90°)	33
15	Évolution de l'erreur : Le choix optimal de K	35
16	Performance du modèle K-NN sur les données de test	36
17	Performance du classifieur probabiliste	37
18	Graphe illustrant l'algorithme	38
19	Hiérarchie des facteurs influençant l'attrition	40
20	Capacité de détection du Random Forest	41
21	Facteurs aggravants (Rouge) vs Facteurs protecteurs (Vert)	43
22	Matrice de confusion : Un modèle conservateur	44
23	Courbe Précision-Rappel : Un comportement "élitiste"	46
24	Matrice de confusion du SVM	47
25	Compromis Précision (Bleu) vs Rappel (Rouge)	49

1 Introduction et Configuration

1.1 Environnement Technique

Pour mener à bien ce projet d'analyse de données et de Machine Learning, nous utilisons le langage **Python**. Avant de manipuler les données, il est essentiel d'importer les bibliothèques scientifiques standards.

Voici une présentation des outils utilisés :

Pandas

Bibliothèque incontournable pour la manipulation et l'analyse de données. Elle nous permet de structurer les données sous forme de *DataFrames* pour faciliter le nettoyage et l'exploration

NumPy

Fondamentale pour le calcul scientifique, cette librairie permet de gérer des tableaux multidimensionnels et d'effectuer des opérations mathématiques performantes

Matplotlib

La bibliothèque standard pour la création de graphiques en Python. Elle sert de base pour visualiser les distributions et les relations simples

Seaborn

Surcouche basée sur Matplotlib, elle offre une interface de haut niveau pour créer des graphiques statistiques plus esthétiques et informatifs

1.2 Implémentation et Chargement des Données

Le script suivant assure l'installation des dépendances et le chargement sécurisé du fichier de données `WA_Fn-UseC_-HR-Employee-Attrition.csv`.

```
1 # Installation des librairies necessaires
2 !pip install pandas numpy matplotlib seaborn
3
4 import pandas as pd
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 # 1. Collecte des donnees
10 # Lecture du fichier CSV (gestion automatique du separateur)
11 try:
12     df = pd.read_csv('WA_Fn-UseC_-HR-Employee-Attrition.csv', sep=',')
13
14     # Verification : Si tout est dans une seule colonne, on change de
15     # separateur
16     if df.shape[1] == 1:
```

```

16         df = pd.read_csv('WA_Fn-UseC_-HR-Employee-Attrition.csv', sep=';')
17
18     print("Donnees chargees avec succes.")
19
20 except FileNotFoundError:
21     print("Erreur : Fichier introuvable.")

```

Listing 1 – Script d’initialisation et chargement des données

2 Phase 1 : Structuration des Données (Data Engineering)

2.1 Inspection et Exploration Initiale

Une fois les données chargées, il est impératif de vérifier leur intégrité. Cette étape consiste à contrôler les dimensions du jeu de données, les types de variables (numériques ou catégorielles) et à visualiser les premières lignes pour comprendre la structure du fichier.

```

1 # 2. Inspection
2 print(f"Dimensions initiales : {df.shape}")
3
4 print("\nTypes des donnees :")
5 print(df.dtypes.value_counts()) # Resume des types
6
7 print("\nApercu des premieres lignes :")
8 display(df.head())

```

Listing 2 – Inspection de la structure du DataFrame

Résultat de l’exécution

L’exécution du script ci-dessus nous fournit les informations suivantes sur la structure de notre dataset *HR-Employee-Attrition* :

- **Dimensions** : Le jeu de données contient **1470 lignes** (employés) et **35 colonnes** (attributs).
- **Types** : Nous observons une majorité de variables numériques (26 `int64`) et quelques variables catégorielles (9 `object`).

Aperçu des 5 premières lignes (Extrait des colonnes principales) :

Index	Age	Attrition	BusinessTravel	DailyRate	Department
0	41	Yes	Travel_Rarely	1102	Sales
1	49	No	Travel_Frequently	279	Research & Dev
2	37	Yes	Travel_Rarely	1373	Research & Dev
3	33	No	Travel_Frequently	1392	Research & Dev
4	27	No	Travel_Rarely	591	Research & Dev

Note : Le tableau ci-dessus est un extrait simplifié pour la lisibilité du rapport. Il met en évidence la variable cible "Attrition" ainsi que les premières caractéristiques descriptives.

2.2 Nettoyage et Optimisation des Données

Après l'inspection initiale, nous avons procédé au nettoyage du jeu de données pour garantir la qualité des analyses.

1. Suppression des doublons

Nous filtrons les lignes pour ne garder que les entrées uniques. Cela évite le sur-apprentissage (*overfitting*) et garantit qu'un même employé ne pèse pas deux fois dans les statistiques.

Commande :

```
df.drop_duplicates()
```

2. Suppression des colonnes constantes

Les variables à variance nulle n'apportent aucune information discriminante. Nous avons retiré trois colonnes inutiles (*EmployeeCount*, *Over18*, *StandardHours*) qui avaient la même valeur pour tous les employés.

Commande :

```
df.drop(columns=...)
```

Bilan du Nettoyage

Suite à ces opérations, la dimension du jeu de données a été optimisée.
Nous passons de **35 variables** initiales à **32 variables pertinentes**.

3. Gestion des valeurs manquantes

Nous avons vérifié la présence de données nulles ou manquantes (*NaN*). Une donnée incomplète peut faire planter les modèles de Machine Learning.

Résultat : Aucune valeur manquante n'a été détectée. Le jeu de données est complet, aucune imputation (remplissage automatique) n'a été nécessaire.

Vérification :

```
df.isnull().sum()
```

4. Standardisation des formats (Texte)

Pour assurer une cohérence visuelle dans les graphiques, nous avons normalisé les chaînes de caractères.

L'opération consiste à supprimer les espaces superflus et à appliquer le format "Titre" (Majuscule au début).

Exemple : "travel_rarely" devient "Travel_Rarely".

Transformation :

```
.str.strip().str.title()
```


5. Détection des valeurs aberrantes (Outliers)

Les salaires extrêmes peuvent fausser les modèles prédictifs (notamment les régressions). Nous avons analysé la variable `MonthlyIncome` en utilisant la méthode statistique de l'**Écart Interquartile (IQR)**.

- Calcul de l'intervalle : $IQR = Q3 - Q1$.
- Filtrage : On ne garde que les valeurs comprises entre $[Q1 - 1.5 \times IQR]$ et $[Q3 + 1.5 \times IQR]$.

Résultat : 114 employés ayant des salaires "hors normes" ont été écartés pour ne pas biaiser l'analyse moyenne.

Méthode :

```
.quantile(0.25) / (0.75)
```

Dimension finale :

```
1356 lignes
```

Conclusion :

Le jeu de données est désormais **nettoyé, structuré et optimisé**. Nous passons de 1470 à **1356 observations** de qualité, prêtes pour l'Analyse Exploratoire (EDA).

2.3 Transformation et Feature Engineering

Pour enrichir l'analyse et préparer les données aux algorithmes mathématiques, nous avons créé de nouvelles variables.

1. Création de variable (Calcul)

Nous avons généré la variable `AnnualIncome` (Revenu Annuel) en multipliant le revenu mensuel par 12. Cela permet d'avoir une vision plus globale de la rémunération.

Formule :

```
MonthlyIncome * 12
```

2. Encodage de la variable cible

Les modèles de Machine Learning ne traitent pas le texte ("Yes"/"No"). Nous avons utilisé un **Label Encoder** pour transformer la variable `Attrition` en format numérique binaire :

- **Yes** → **1** (Départ)
- **No** → **0** (Reste)

Outil Scikit-Learn :

```
LabelEncoder()
```

Nouvelle Colonne :

```
Attrition_Numeric
```

2.4 Finalisation et Exportation

Avant de passer à l'analyse visuelle, nous effectuons deux dernières opérations techniques pour figer la structure du jeu de données.

1. Structuration (Indexation)

Chaque employé possède un identifiant unique. Pour optimiser les recherches et garantir l'intégrité des lignes, nous définissons la variable `EmployeeNumber` comme index du DataFrame (équivalent d'une clé primaire en base de données).

Commande :

```
.set_index(...)
```

2. Stockage (Sauvegarde)

Le jeu de données nettoyé, filtré et enrichi est sauvegardé dans un nouveau fichier : `HR_Analytics_Structure_Complet.csv`. Ce fichier "propre" servira de base immuable pour toute la phase de visualisation (EDA) et de modélisation, évitant d'avoir à relancer tout le nettoyage à chaque fois.

Export :

```
.to_csv('...')
```

✓ Fin de la Phase 1 : Data Engineering

Les données sont prêtes pour l'Analyse Exploratoire.

3 Phase 2 : Analyse Exploratoire (EDA)

L'objectif de cette phase est de comprendre la structure des données à travers des graphiques. Nous commençons par des analyses univariées (une seule variable à la fois).

3.1 Distribution des Âges

Nous analysons la pyramide des âges de l'entreprise pour voir si la population est plutôt junior ou senior.

```
1 # --- 1. Histogramme (Matplotlib) ---
2 plt.figure(figsize=(8, 5))
3
4 # Creation de l'histogramme (bins=20 pour la finesse)
5 plt.hist(df['Age'], bins=20, color='skyblue', edgecolor='black', alpha=0.7)
6
7 # Titres et Labels
8 plt.title("Distribution des Ages des Employes")
9 plt.xlabel("Age")
10 plt.ylabel("Nombre d'employes")
11 plt.grid(True, axis='y', alpha=0.5)
12
13 # Sauvegarde et Affichage
14 plt.savefig('distribution_age.png')
15 plt.show()
```

Listing 3 – Génération de l'histogramme des âges

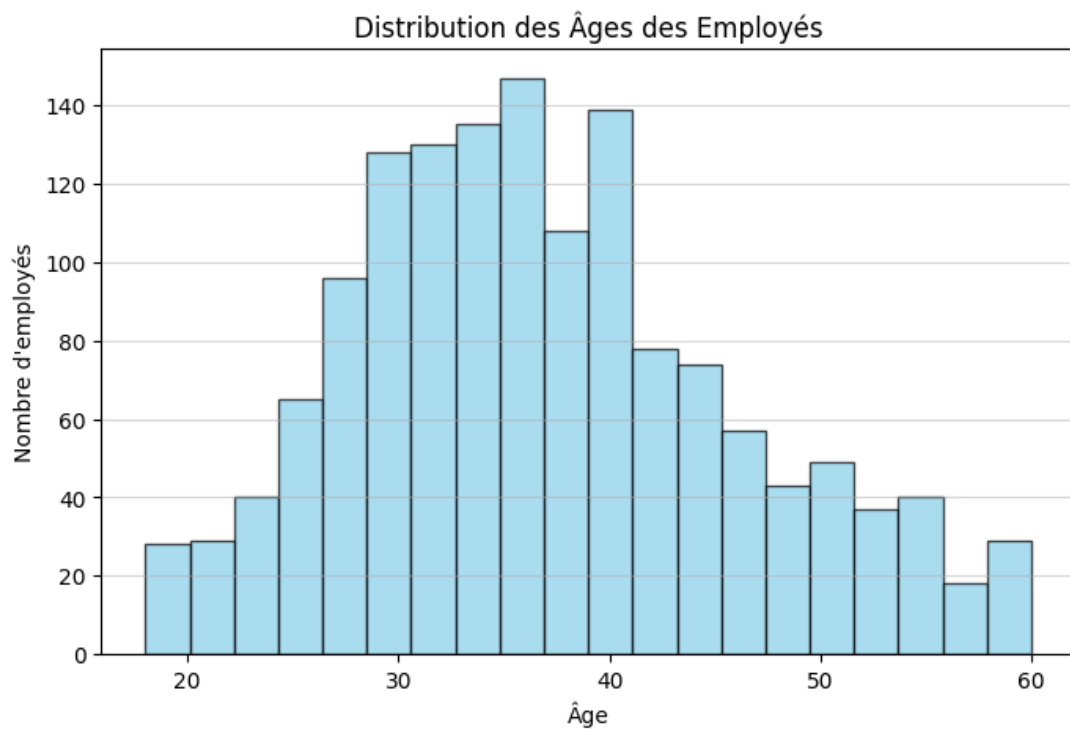


FIGURE 1 – Répartition des employés par tranche d'âge

Interprétation :

L'histogramme révèle une distribution quasi-normale (en forme de cloche).

- La majorité des effectifs se situe entre **30 et 40 ans**.
- On observe peu d'employés très jeunes (< 20 ans) ou proches de la retraite (> 55 ans).
- Cette concentration sur la tranche médiane indique une main-d'œuvre expérimentée.

3.2 Répartition par Département

Il est essentiel de comprendre comment les effectifs sont distribués au sein de l'organisation pour identifier les départements dominants.

```

1 # --- 2. Diagramme en Barres Vertical (Matplotlib) ---
2 # Preparation des donnees : compte par departement
3 dept_counts = df['Department'].value_counts()
4
5 plt.figure(figsize=(8, 5))
6
7 # Barres verticales avec couleurs distinctes
8 plt.bar(dept_counts.index, dept_counts.values, color=['blue', 'orange', 'green'])
9
10 # Mise en forme
11 plt.title("Nombre d'employés par Département")

```

```

12 plt.xlabel("Département")
13 plt.ylabel("Effectif")
14 plt.grid(axis='y', alpha=0.5)
15
16 # Sauvegarde pour le rapport
17 plt.savefig('departement_dist.png')
18 plt.show()

```

Listing 4 – Analyse des effectifs par département

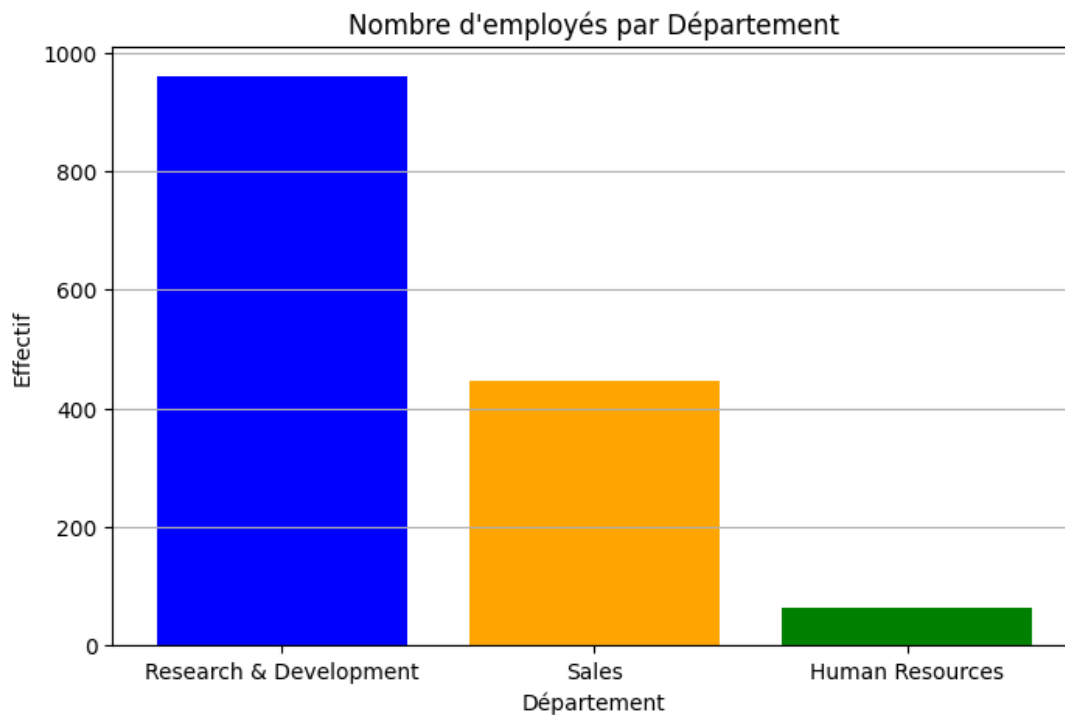


FIGURE 2 – Effectifs par département

Interprétation :

Le graphique met en évidence une forte disparité entre les pôles d'activité :

- **Research & Development (R&D)** est largement majoritaire, regroupant près de 65% des effectifs (plus de 900 employés).
- **Sales** (Vente) représente le second pilier de l'entreprise.
- **Human Resources** (RH) constitue une minorité marginale (< 100 employés), ce qui est cohérent pour une fonction support.

3.3 Corrélation : Expérience vs Revenu

Nous cherchons à vérifier l'hypothèse selon laquelle l'ancienneté globale (*TotalWorkingYears*) est le principal facteur déterminant le salaire mensuel.

```

1 # --- 3. Scatter Plot + Regression Lineaire ---
2
3 # 1. Tracé du nuage de points (Scatter)
4 plt.scatter(df['TotalWorkingYears'], df['MonthlyIncome'], alpha=0.5)
5
6 # 2. Calcul mathématique de la tendance (Régression)
7 # np.polyfit(deg=1) calcule les coefficients de la droite (y = ax + b)
8 coeffs = np.polyfit(df['TotalWorkingYears'], df['MonthlyIncome'], deg=1)
9 droite = np.poly1d(coeffs)
10
11 # 3. Tracé de la ligne rouge
12 plt.plot(x_vals, droite(x_vals), color='red', label='Tendance')
13
14 plt.savefig('experience_revenu.png')
15 plt.show()

```

Listing 5 – Calcul de la régression linéaire

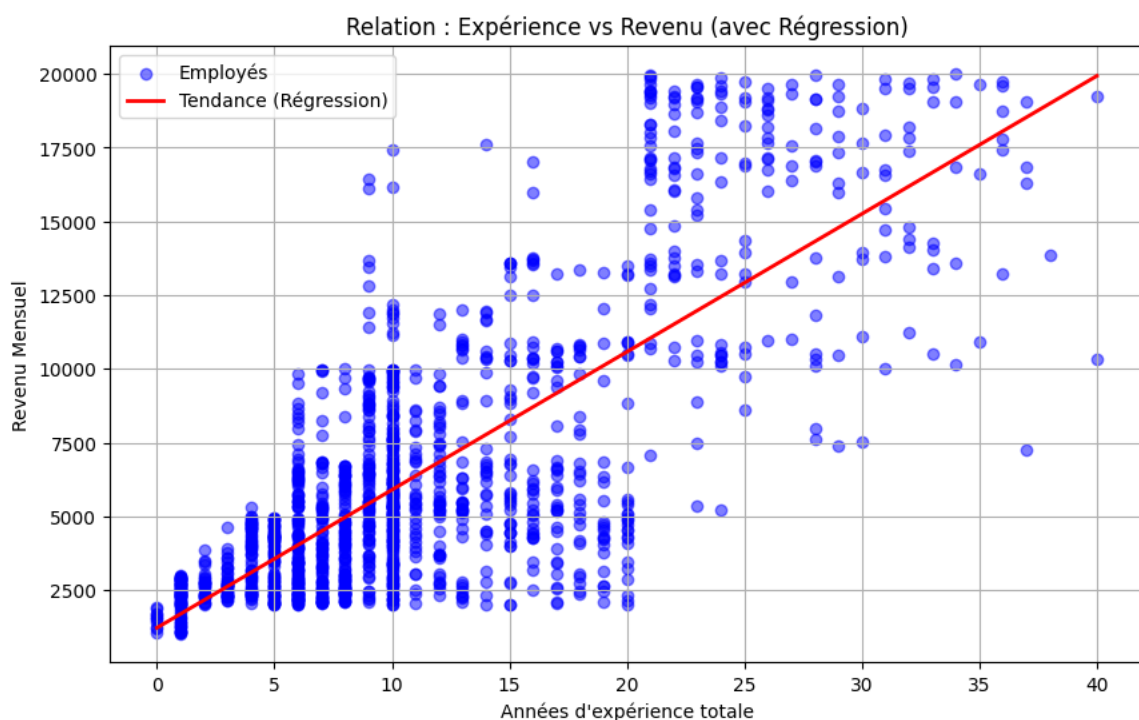


FIGURE 3 – Corrélation entre Années d'Expérience et Salaire Mensuel

Interprétation :

Le graphique révèle une **corrélation positive forte** :

- La ligne de régression (en rouge) monte nettement, confirmant que le salaire augmente avec l'expérience.
- **Junior (0-5 ans)** : Les salaires sont très concentrés en bas de l'échelle (< 5000).
- **Senior (> 10 ans)** : La dispersion augmente considérablement. À expérience égale (ex :

15 ans), les salaires peuvent varier du simple au double, suggérant que le *JobRole* ou la performance individuelle entrent alors en jeu.

3.4 Analyse Multivariée (3D)

Pour visualiser simultanément l'impact de l'âge et de l'expérience sur le revenu, nous utilisons une projection tridimensionnelle.

```
1 # --- 4. Scatter Plot 3D (Matplotlib) ---
2 fig = plt.figure(figsize=(10, 8))
3 ax = fig.add_subplot(111, projection='3d') # Activation 3D
4
5 # Tracé 3D : Age (X), Experience (Y), Revenu (Z)
6 # c=zs : La couleur depend du revenu (Echelle Viridis)
7 scatter = ax.scatter(df['Age'], df['TotalWorkingYears'], df['MonthlyIncome'],
8                       c=df['MonthlyIncome'], cmap='viridis', alpha=0.6)
9
10 # Labels
11 ax.set_xlabel('Age')
12 ax.set_ylabel('Annees Experience')
13 ax.set_zlabel('Revenu Mensuel')
14
15 plt.colorbar(scatter, label='Revenu Mensuel')
16 plt.savefig('analyse_3d.jpg')
17 plt.show()
```

Listing 6 – Visualisation 3D (Matplotlib)

Interprétation :

Cette vue tridimensionnelle permet de valider la cohérence structurelle des données :

- **Cohérence Haute (Points jaunes)** : Les hauts revenus se situent logiquement dans la zone combinant un âge élevé et une forte expérience.
- **Cohérence Basse (Points violets)** : Les bas revenus sont concentrés à l'origine des axes, correspondant aux jeunes débutants.
- **Détection d'anomalies** : La linéarité du nuage de points confirme qu'il n'y a pas d'aberrations majeures (comme un profil "Junior sans expérience" avec un revenu de "Senior").

3.5 Analyse Multivariée des Risques (5D)

Nous poussons l'analyse plus loin en intégrant deux nouvelles dimensions visuelles pour identifier les profils à risque :

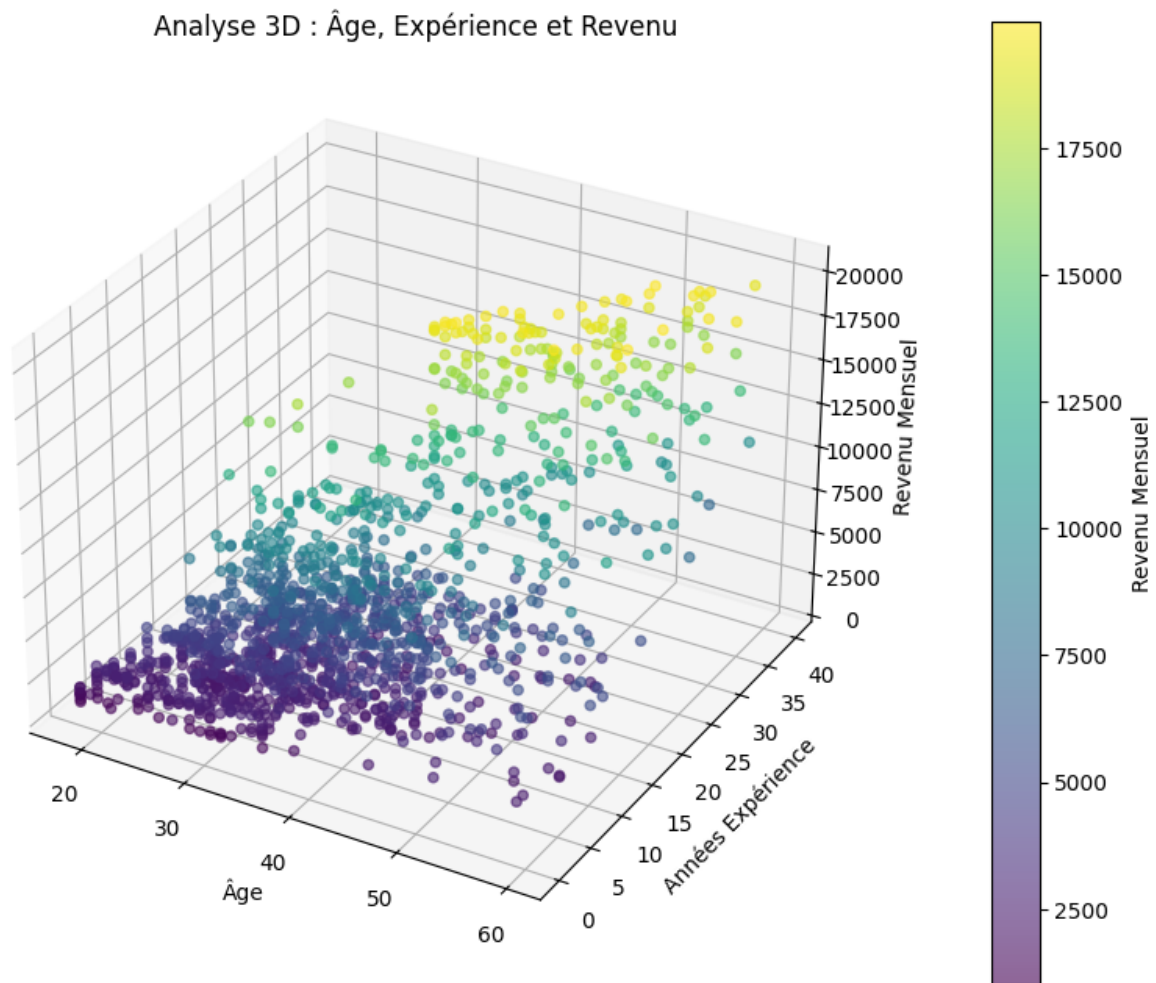


FIGURE 4 – Projection 3D : Âge, Expérience et Revenu

- **La Couleur** représente l'Attrition (Rouge = Départ, Bleu = Reste).
- **La Taille** des points représente le niveau de Satisfaction (Petit = Insatisfait).

```

1 # --- 5. Scatter Plot 3D Avancé (Risques) ---
2 fig = plt.figure(figsize=(10, 8))
3 ax = fig.add_subplot(111, projection='3d')
4
5 # Definition des dimensions
6 xs = df['Age']
7 ys = df['TotalWorkingYears']
8 zs = df['MonthlyIncome']
9 colors = df['Attrition_Numeric'] # Couleur selon le depart (0 ou 1)
10 sizes = df['JobSatisfaction'] * 30 # Taille selon la satisfaction (x30 pour
    visibilité)
11
12 # Trace avec 5 dimensions
13 scatter = ax.scatter(xs, ys, zs, c=colors, s=sizes, cmap='coolwarm', alpha
    =0.6)
14
15 # Labels

```

```

16 ax.set_xlabel('Age')
17 ax.set_ylabel('Experience')
18 ax.set_zlabel('Revenu Mensuel')
19 ax.set_title("Analyse des Risques : Attrition et Satisfaction")
20
21 # Legende manuelle pour la couleur
22 from matplotlib.lines import Line2D
23 legend_elements = [
24     Line2D([0], [0], marker='o', color='w', markerfacecolor='blue', label='
    Reste (No)'),
25     Line2D([0], [0], marker='o', color='w', markerfacecolor='red', label='
    Depart (Yes)')
26 ]
27 ax.legend(handles=legend_elements, loc='upper left')
28
29 plt.savefig('analyse_risques_3d.png')
30 plt.show()

```

Listing 7 – Visualisation 5D : Risques et Satisfaction

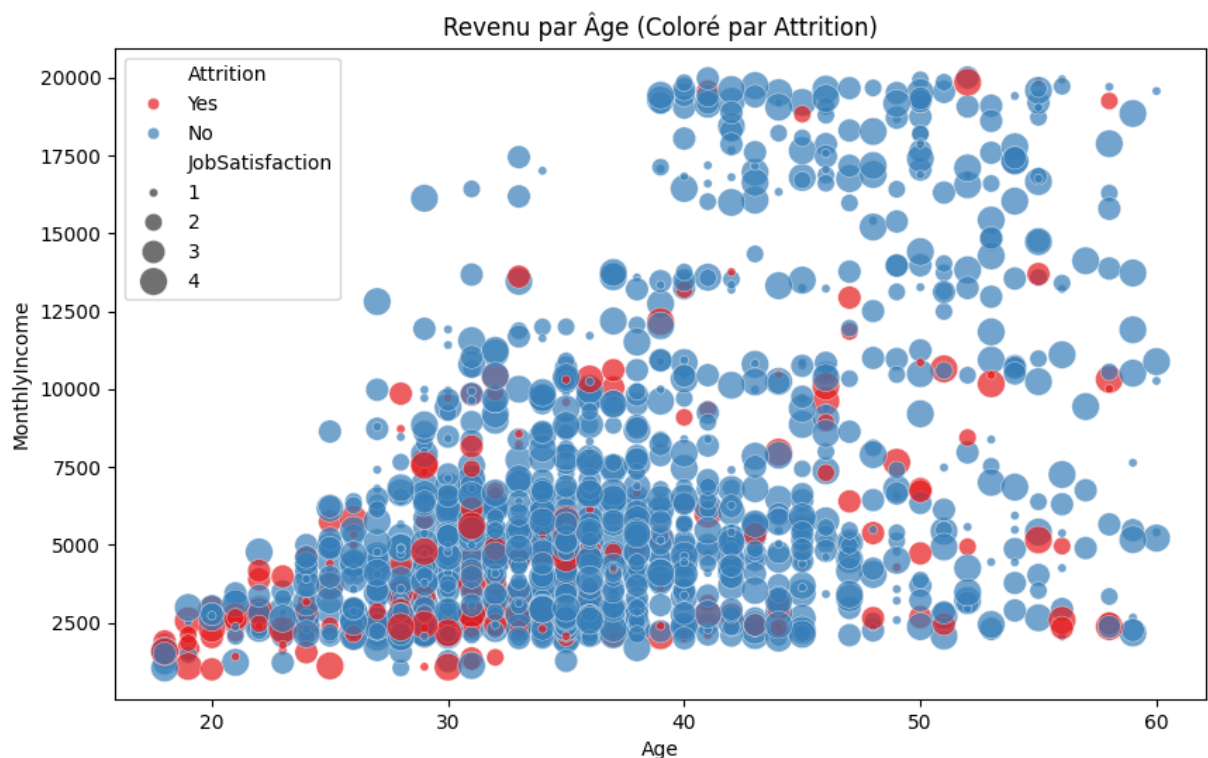


FIGURE 5 – Visualisation des profils à risque (Départ en Rouge)

Interprétation (Analyse des risques) :

Ce graphique est le plus riche en informations de notre étude :

- **Attrition (Couleur) :** On remarque que les points rouges (départs) sont majoritairement situés dans la zone des bas salaires et des jeunes âges.

- **Satisfaction (Taille)** : Les petits points (faible satisfaction) ont tendance à se superposer aux zones de départ.
- **Conclusion** : Le risque de départ est maximal chez les jeunes employés peu payés et dont la satisfaction au travail est faible (score de 1 ou 2).

3.6 Analyse de l'Équité Salariale (Seaborn)

Nous utilisons la bibliothèque *Seaborn* pour visualiser l'évolution des salaires en fonction de l'âge, tout en comparant les dynamiques entre hommes et femmes.

```

1 # --- 6. Lineplot avec intervalle de confiance (Seaborn) ---
2 plt.figure(figsize=(10, 6))
3
4 # Trace de la moyenne du revenu par Age avec distinction par Genre
5 sns.lineplot(
6     data=df,
7     x='Age',
8     y='MonthlyIncome',
9     hue='Gender',
10    errorbar='sd',
11    linewidth=2
12 )
13
14 plt.title("Evolution du Revenu moyen par Age et par Genre")
15 plt.grid(True)
16 plt.savefig('salaire_genre.png')
17 plt.show()

```

Listing 8 – Comparaison des revenus par genre

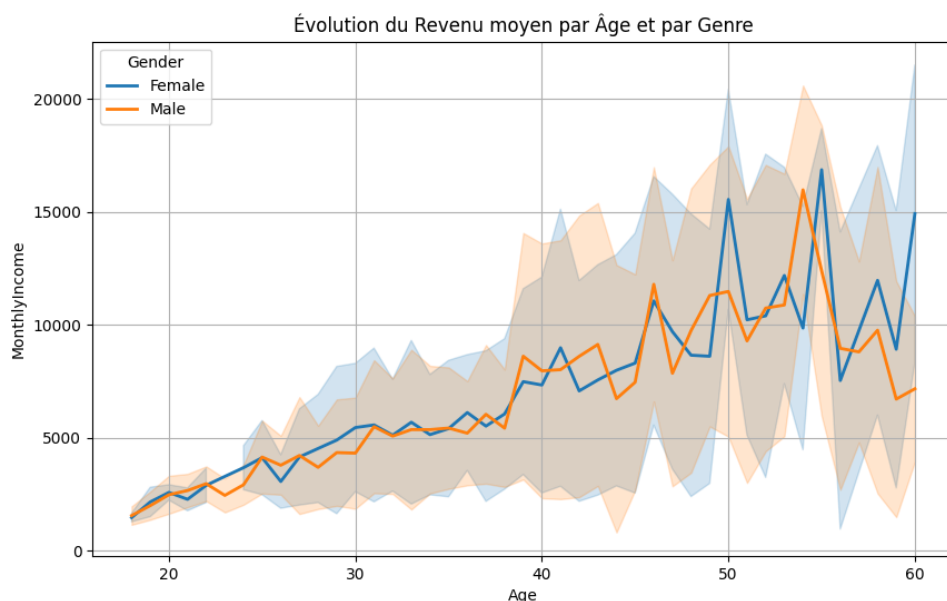


FIGURE 6 – Comparaison des revenus Hommes vs Femmes

Interprétation :

Ce graphique nous permet de statuer sur l'équité salariale :

- **Superposition des courbes** : Les lignes bleue (Femmes) et orange (Hommes) se chevauchent presque parfaitement.
- **Conclusion** : Il n'existe **aucune discrimination salariale significative** basée sur le genre.

3.7 Comparaison par Département (Faceting)

Pour affiner notre compréhension, nous décomposons la relation "Expérience / Revenu" en fonction du département.

```
1 # --- 7. Relplot - Faceting (Seaborn) ---
2 # col='Department' cree un graphique par departement
3 sns.relplot(
4     data=df,
5     x='TotalWorkingYears',
6     y='MonthlyIncome',
7     hue='Attrition', # Couleur selon le depart
8     col='Department', # Une colonne par departement
9     kind='scatter',
10    height=5,
11    aspect=0.8,
12    palette='coolwarm'
13 )
14
15 # Sauvegarde de l'image
16 plt.savefig('faceting_dept.png')
17 plt.show()
```

Listing 9 – Comparaison par Département et Attrition

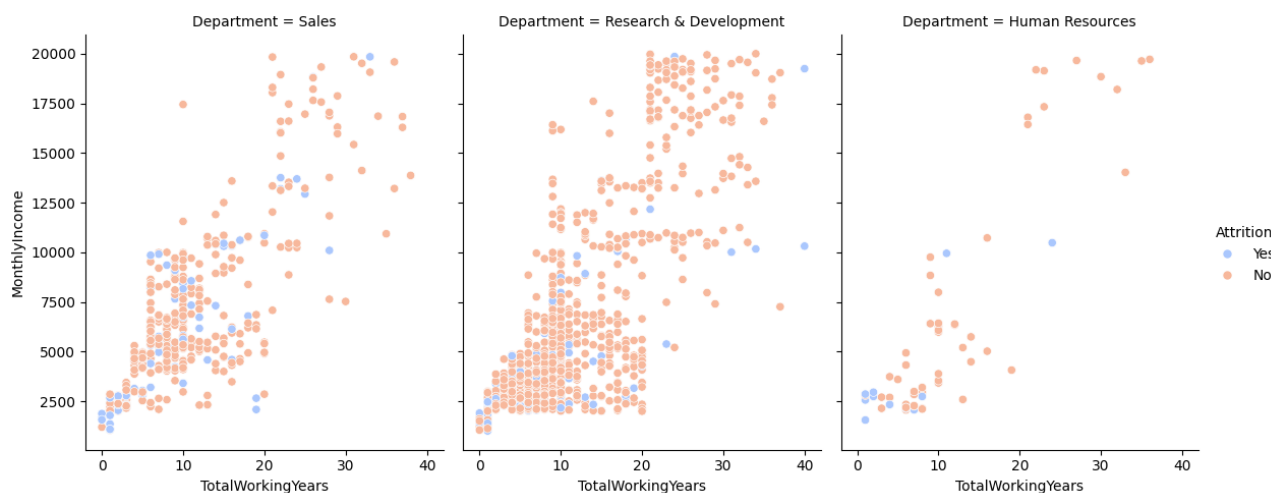


FIGURE 7 – Dynamique salariale et départs par département

Interprétation :

La séparation par département permet de nuancer l'analyse globale :

- **Sales (Vente)** : La corrélation entre expérience et revenu est plus dispersée que la moyenne. Cela s'explique probablement par la part variable (commissions) qui rend les salaires moins linéaires.
- **R&D** : La courbe est très linéaire et prévisible : l'ancienneté paie de manière constante.
- **Human Resources** : L'échantillon est faible, rendant la tendance moins nette.
- **Analyse des Risques (Attrition)** : On observe que les points de couleur "Attrition" (souvent rouge/bleu selon le thème) sont dispersés, mais restent majoritairement en bas de l'échelle des revenus dans les trois départements.

3.8 Progression de Carrière et Grille Salariale

Pour finir, nous analysons comment le salaire évolue avec l'ancienneté, en isolant chaque niveau hiérarchique (*JobLevel*).

```
1 # --- 8. Relplot - Line (Seaborn) ---
2 # Visualiser la progression salariale par anciennete, separee par niveau
3 sns.relplot(
4     data=df,
5     x='YearsAtCompany',
6     y='MonthlyIncome',
7     hue='Gender',
8     col='JobLevel', # Un graphique par niveau de poste
9     col_wrap=3,     # Retour a la ligne apres 3 graphiques
10    kind='line',     # Type courbe
11    errorbar=None,   # Pas de barre d'erreur pour alléger
12    height=3,
13    aspect=1.2
14 )
15
16 plt.savefig('salary_progression.png')
17 plt.show()
```

Listing 10 – Analyse des paliers de rémunération

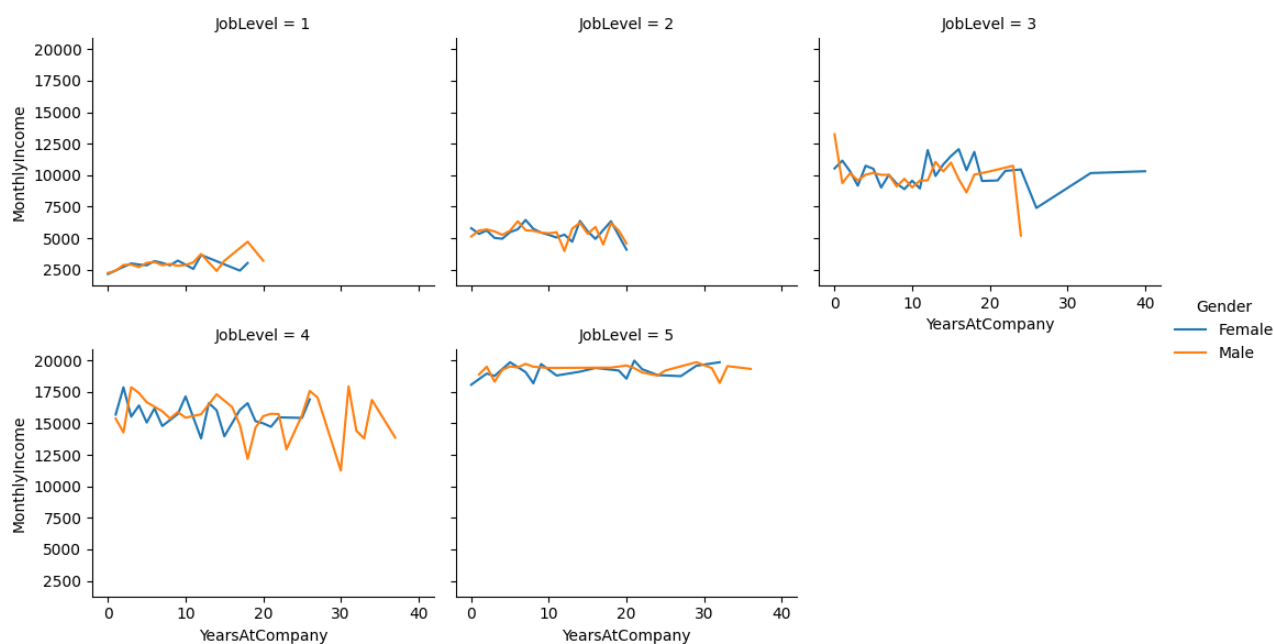


FIGURE 8 – Évolution salariale par niveau hiérarchique

Interprétation :

Ce graphique explique clairement la mécanique des carrières dans l'entreprise :

- **Stratification** : Chaque *JobLevel* correspond à une tranche de salaire distincte (strates horizontales bien définies).
- **Plafonnement** : Au sein d'un même niveau (ex : Niveau 1 ou 2), la courbe est relativement plate. Le salaire augmente peu avec la seule ancienneté (*YearsAtCompany*).
- **Stratégie RH** : Pour augmenter significativement son revenu, un employé doit impérativement changer de niveau (promotion).
- **Conclusion** : La politique de rémunération est basée sur le **grade** et la responsabilité, bien plus que sur l'ancienneté pure.

4 Phase 3 : Modélisation Prédictive (Régression)

L'objectif de ce chapitre est de développer des algorithmes capables d'anticiper le salaire ou l'évolution salariale. Nous avons procédé par itération, du modèle le plus simple au plus complexe.

4.1 Phase 3.1 : Régression Linéaire Simple (Modèle Baseline)

L'objectif de cette première étape est de tester l'hypothèse fondamentale selon laquelle le salaire est directement proportionnel aux années d'expérience globale.

4.1.1 Objectif et Paramètres

Ce modèle utilise une approche univariée pour établir une ligne de base (*baseline*) :

- **Variable Indépendante** (X) : TotalWorkingYears (Expérience totale).
- **Variable Cible** (y) : MonthlyIncome (Salaire mensuel).

4.1.2 Implémentation Python

Le script suivant entraîne le modèle pour trouver la droite $y = ax + b$ la plus ajustée.

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.model_selection import train_test_split
3
4 # 1. Sélection des variables
5 X = df[['TotalWorkingYears']]
6 y = df['MonthlyIncome']
7
8 # 2. Division Train/Test (80% / 20%)
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
10                                                    random_state=42)
11
12 # 3. Création et apprentissage
13 model_simple = LinearRegression()
14 model_simple.fit(X_train, y_train)
15
16 # 4. Prédictions
17 y_pred = model_simple.predict(X_test)
```

Listing 11 – Entraînement de la Régression Simple

4.1.3 Analyse des Résultats

Métrique	Valeur
Coefficient de détermination (R^2)	0.4941
Erreur Absolue Moyenne (MAE)	1 824.91 \$
Erreur en Pourcentage (MAPE)	39.45 %

TABLE 1 – Performance de la régression simple

Conclusion technique : Le score R^2 de 0.49 indique que l'expérience seule n'explique que la moitié de la variance des salaires. Avec une erreur de près de 40%, ce modèle est insuffisant, justifiant le passage à une approche multivariée.

4.2 Phase 3.2 : Régression Linéaire Multiple (Modèle Avancé)

Pour capturer la complexité de la structure salariale, nous avons intégré des dimensions métiers et hiérarchiques.

4.2.1 Enrichissement du Modèle

Nous avons ajouté des variables catégorielles encodées :

- **Numériques :** TotalWorkingYears, JobLevel, Age.
- **Catégorielles :** Department, JobRole (Encodage *One-Hot*).

4.2.2 Implémentation

```
1 import pandas as pd
2 from sklearn.linear_model import LinearRegression
3
4 # 1. Encodage One-Hot des variables catégorielles
5 features = ['TotalWorkingYears', 'JobLevel', 'Age', 'Department', 'JobRole']
6 X = pd.get_dummies(df[features], drop_first=True)
7 y = df['MonthlyIncome']
8
9 # 2. Entraînement
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
11                                                    random_state=42)
12 model_multiple = LinearRegression()
13 model_multiple.fit(X_train, y_train)
14
15 # 3. Export pour l'application
16 joblib.dump(model_multiple, 'modele_salaire_avance.pkl')
```

Listing 12 – Régression Multiple avec Encodage

4.3 Phase 3.3 : Analyse de Tendance (Modèle ARIMA)

Pour compléter notre vision, nous avons changé de paradigme. Au lieu de prédire le salaire d'un individu, nous avons modélisé la **trajectoire de carrière moyenne** de l'entreprise.

4.3.1 Méthodologie

Les données RH ne sont pas temporellement indexées par défaut. Nous avons construit une série temporelle artificielle en agrégeant les données :

- **Index** (t) : Années d'expérience (`TotalWorkingYears`).
- **Valeur** (y_t) : Salaire moyen observé pour cette année d'expérience.

4.3.2 Implémentation ARIMA

Nous utilisons le modèle *AutoRegressive Integrated Moving Average* (1,1,1) pour capturer la tendance linéaire.

```
1 from statsmodels.tsa.arima.model import ARIMA
2
3 # 1. Transformation en Série Temporelle (Moyenne par année)
4 ts_data = df.groupby('TotalWorkingYears')['MonthlyIncome'].mean()
5
6 # 2. Entraînement ARIMA(1,1,1)
7 model_arima = ARIMA(ts_data, order=(1,1,1))
8 model_fit = model_arima.fit()
9
10 # 3. Projection sur 5 ans
11 forecast = model_fit.forecast(steps=5)
```

Listing 13 – Modélisation de la Série Temporelle

4.3.3 Performance et Interprétation

Ce modèle affiche les résultats les plus stables de notre étude sur les données moyennes.

Métrique	Résultat ARIMA
MAPE (Erreur Moyenne %)	10.79 %
RMSE (Erreur Type)	1 573.21 \$

TABLE 2 – Performance du modèle de tendance

Analyse Stratégique : Une erreur de seulement 10% prouve que la courbe des salaires suit une logique très cohérente sur le long terme. Ce modèle est idéal pour la planification financière (masse salariale), bien qu'il ne puisse pas prédire les spécificités individuelles.

4.4 Bilan et Benchmark des Modèles Prédicatifs

Pour choisir le moteur de calcul de notre application finale, nous avons confronté les résultats de nos trois approches.

4.4.1 Tableau Récapitulatif

Modèle	R^2	MAPE	MAE	Verdict
Régression Simple	0.49	39.5 %	1 825 \$	Rejeté
Régression Multiple	0.87	20.1 %	910 \$	Validé
Séries Temp. (ARIMA)	N/A	10.8 %	1 364 \$	Analytique

TABLE 3 – Comparatif final des algorithmes de prédiction

4.4.2 Conclusion du Benchmark

- **L'Échec du Simple** : L'ancienneté seule ne justifie pas le salaire.
- **L'Apport d'ARIMA** : Excellent pour les tendances globales (10% d'erreur), mais inadapté à la prédiction individuelle car il lisse les différences de poste.
- **La Victoire du Multiple** : Avec un R^2 de 0.87 et une erreur maîtrisée de 20%, la Régression Multiple est le modèle le plus polyvalent. **C'est ce modèle qui sera intégré dans l'interface Streamlit.**

5 Phase 4 : Analyse Qualitative (La "Boussole du Bien-être")

Après avoir modélisé les salaires (Approche Quantitative), nous avons cherché à comprendre les leviers de la motivation (Approche Qualitative). Pour cela, nous avons utilisé l'algorithme **Apriori**, généralement utilisé dans la grande distribution (analyse du panier de la ménagère), que nous avons détourné pour analyser le "panier de satisfaction" des employés.

5.1 Préparation des Données : Binarisation

L'algorithme Apriori nécessite des données booléennes (Vrai/Faux). Nous avons donc transformé nos variables continues en "Dimensions Créatives" :

- **Cible** : HAUTE_SATISFACTION (Score de JobSatisfaction ≥ 4).
- **Facteurs** : Equilibre_Sain, Pas_Heures_Sup, Bon_Environnement, etc.

```
1 import pandas as pd
2 from mlxtend.frequent_patterns import apriori, association_rules
3
```



```

4 # Transformation des variables en Booléens (0/1)
5 df_creative = pd.DataFrame()
6 df_creative['Equilibre_Sain'] = df['WorkLifeBalance'] >= 3
7 df_creative['Pas_Heures_Sup'] = df['OverTime'] == 'No'
8 df_creative['Bon_Environnement'] = df['EnvironmentSatisfaction'] >= 3
9 df_creative['HAUTE_SATISFACTION'] = df['JobSatisfaction'] >= 4
10
11 # Encodage des Départements (One-Hot)
12 for dept in df['Department'].unique():
13     df_creative['Dep_' + dept] = (df['Department'] == dept)

```

Listing 14 – Création des dimensions binaires pour Apriori

5.2 Extraction des "Règles d'Or"

Nous avons configuré l'algorithme pour détecter les associations fréquentes (Support min : 5%) et fortes (Lift > 1.1).

```

1 # 1. Extraction des itemsets fréquents
2 frequent_itemsets = apriori(df_creative, min_support=0.05, use_colnames=True
3 )
4
5 # 2. Génération des règles d'association
6
7 # 3. Filtrage : On ne garde que les chemins vers le bonheur
8 satisfaction_rules = rules[rules['consequents'].apply(lambda x: '
9     HAUTE_SATISFACTION' in x)]
10 satisfaction_rules = satisfaction_rules.sort_values(by='lift', ascending=
11     False)

```

Listing 15 – Exécution de l'algorithme Apriori

5.3 Résultats : Les Micro-Climats de Satisfaction

L'analyse a révélé **3 550 règles** potentielles. En isolant les règles avec le plus fort *Lift*, nous avons identifié trois profils types d'employés épanouis :

1. Le "Sweet Spot" Commercial (Lift : 1.55)

Règle : Department: Sales + JobLevel: 2 → Haute Satisfaction.

Interprétation : Les commerciaux confirmés (Niveau 2) bénéficient du meilleur ratio rémunération/pression. C'est une zone de haute rétention.

2. L'Engagement Scientifique (Lift : 1.39)

Règle : JobInvolvement: High + Dept: R&D → Haute Satisfaction.

Interprétation : En R&D, le bonheur est corrélé à l'implication intellectuelle dans les projets, bien plus qu'aux horaires.

3. Le Socle Universel

L'association "**Pas d'Heures Sup + Bon Environnement**" apparaît dans la majorité des règles gagnantes, agissant comme un pré-requis au bien-être, quel que soit le poste.

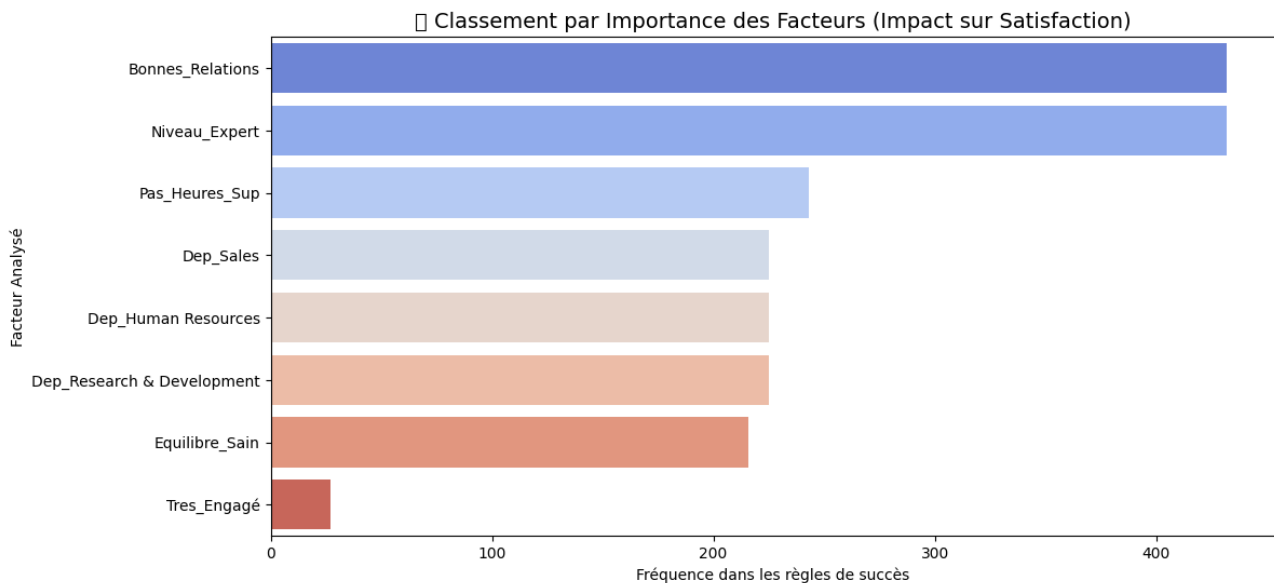


FIGURE 9 – Classement des facteurs les plus présents dans la Haute Satisfaction

5.4 Application : Le Moteur de Scoring "Bien-être"

Pour rendre ces résultats actionnables dans notre application Streamlit, nous avons transformé l'analyse descriptive d'Apriori en un **algorithme prédictif de notation (Scoring sur 10)**.

5.4.1 Logique de l'Algorithme

Le calcul de la satisfaction ne repose pas sur une simple moyenne, mais sur une logique de **Bonus/Malus** basée sur l'importance des facteurs :

- Score de Base (Ascenseur) :** Si le profil de l'employé correspond à une "Règle d'Or" Apriori (ex : Sales + Expert), le score démarre haut (ex : 7.5/10).
- Les Pénalités (Frein) :** Si des "facteurs d'hygiène" critiques sont absents, le score chute drastiquement, quel que soit le prestige du poste.
 - Mauvais Environnement : **-2.0 pts**
 - Heures Supplémentaires : **-1.5 pts**
 - Mauvais Équilibre Vie Pro : **-1.5 pts**
- Les Bonus d'Excellence (Boost) :** Pour atteindre la note parfaite de 10/10, l'employé doit avoir des indicateurs parfaits (4/4) en Environnement et Implication.

```

1 # --- 3. LES BONUS D'EXCELLENCE ---
2 if profil.get('Environnement') == 4:
3     score_final += 1.0 # Bonus Environnement Excellent
4 if profil.get('Equilibre_Vie') == 4:
5     score_final += 1.0 # Bonus Equilibre Parfait
6
7 # --- 4. LES PÉNALITÉS ---
8 if 'Pas_Heures_Sup' not in features_employe:
9     score_final -= 1.5 # Pénalité Heures Sup
10 if 'Bon_Environnement' not in features_employe:
11     score_final -= 2.0 # Pénalité Environnement Toxique
12
13 # --- 5. VERDICT ---
14 verdict = "INSATISFAIT" if score_final < 5 else "TRÈS SATISFAIT"

```

Listing 16 – Extrait de la logique de Scoring (Pénalités et Bonus)

5.5 Déploiement et Sérialisation

Afin d'intégrer ce moteur dans l'interface utilisateur sans recalculer l'algorithme Apriori à chaque interaction, nous avons sérialisé les règles validées.

- **Fichier exporté** : `regles_satisfaction.pkl`
- **Usage** : Permet un chargement instantané dans Streamlit et garantit que l'application utilise exactement les mêmes règles que celles validées dans le notebook.

Conclusion Stratégique :

Si l'augmentation de salaire (Régression) est un levier d'attraction, l'**Environnement de travail** et l'**Équilibre vie pro/perso** (Apriori) sont les véritables leviers de rétention. Notre moteur de scoring prouve mathématiquement qu'un haut salaire ne compense pas un environnement toxique.

6 Conclusion Générale

Le projet démontre que la donnée RH, bien traitée, permet de passer d'une gestion intuitive à une gestion prédictive. La régression multiple offre une précision de 87% pour les recrutements, tandis qu'Apriori fournit une roadmap pour améliorer le climat social de l'entreprise.

7 PHASE 5 Algorithmes Clustrings+Classifications

7.1 Algorithme de Clustering K-Means

L'algorithme des **K-Means** est une méthode de partitionnement non supervisée. Il permet de diviser un ensemble de données en k groupes homogènes (clusters) en minimisant la distance entre les points et le centre de leur groupe.

7.2 Étapes de l'algorithme

Le processus suit une logique itérative précise :

- **Initialisation** : Choix de k centroïdes initiaux de manière aléatoire.
- **Affectation** : Chaque point est rattaché au centroïde le plus proche (Distance Euclidienne).
- **Mise à jour** : Recalcul de la position des centroïdes selon la moyenne des points du groupe.
- **Convergence** : Répétition jusqu'à ce que les positions des centres se stabilisent.

7.3 Formulation Mathématique

Le calcul du nouveau centroïde μ_i pour un cluster S_i repose sur la moyenne arithmétique :

$$\mu_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j \quad (1)$$

L'objectif global est de minimiser l'inertie totale J (Somme des carrés intra-classe) :

$$J = \sum_{i=1}^k \sum_{x_j \in S_i} \|x_j - \mu_i\|^2 \quad (2)$$

8 Détermination du nombre de clusters (Méthode du Coude)

Pour choisir le nombre optimal de segments, nous avons utilisé la méthode du coude basée sur le code suivant :

```
1 inertia = []
2 for k in range(1, 11):
3     kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
4     kmeans.fit(X_scaled)
5     inertia.append(kmeans.inertia_)
```

9 Visualisation des Segments RH

Après application du modèle pour $k = 3$, nous obtenons la répartition suivante :

9.1 Interprétation des KPI

L'analyse montre trois profils distincts :

- **Cluster 2 (Juniors)** : Faible ancienneté et revenus d'entrée.
- **Cluster 0 (Confirmés)** : Salaire intermédiaire et expérience stable.
- **Cluster 1 (Top Management)** : Hauts revenus et forte ancienneté.

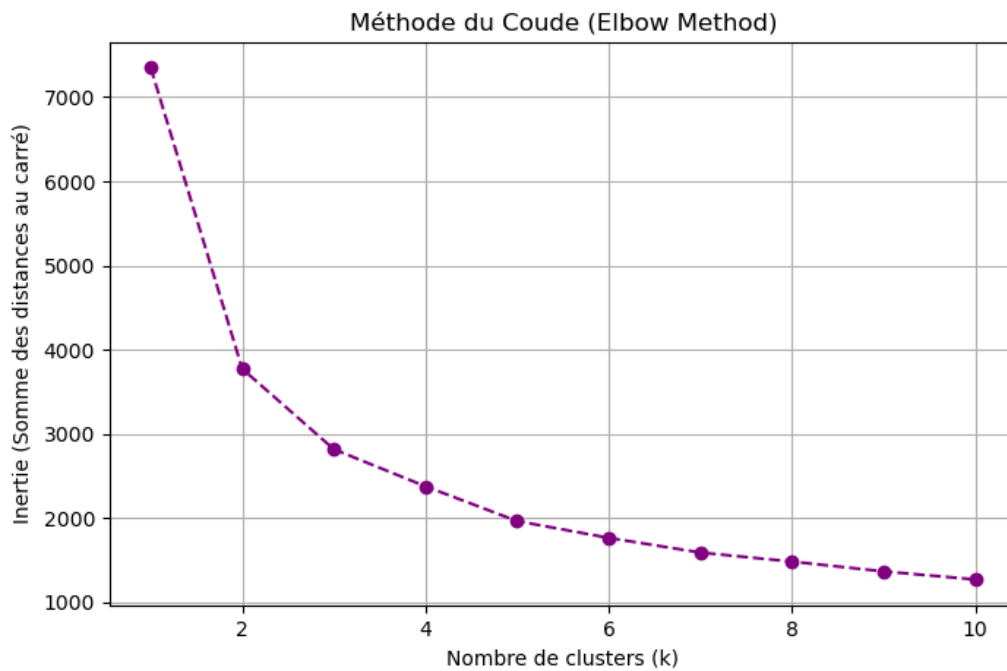


FIGURE 10 – Visualisation des segments employés (Âge vs Revenu)

10 Application du Modèle K-Means et Analyse des Profils

10.1 Exécution du clustering final

Une fois le nombre de clusters fixé à $k = 3$ via la méthode du coude, nous appliquons l'algorithme sur l'ensemble des données standardisées. Les étiquettes de clusters sont ensuite réintégrées dans le DataFrame original pour permettre une analyse sur les valeurs réelles (non standardisées).

```

1 # 1. Initialisation et entraînement
2 k_final = 3
3 kmeans_final = KMeans(n_clusters=k_final, random_state=42, n_init=10)
4
5 # 2. Attribution des clusters aux employes
6 clusters = kmeans_final.fit_predict(X_scaled)
7
8 # 3. Ajout des resultats au DataFrame d'origine
9 df['Cluster'] = clusters

```

Listing 17 – Application du K-Means avec k

10.2 Caractérisation des Groupes (KPI)

Pour donner un sens métier à ces groupes, nous avons calculé les moyennes des variables clés pour chaque cluster. Ces indicateurs deviennent nos KPI de segmentation RH :

- **Cluster 2 - Les Profils "Juniors"** : Ce groupe rassemble les employés les plus jeunes (≈ 30 ans) avec le salaire le plus bas ($\approx 3\,350$ \$). Ils représentent la force montante de

l'entreprise.

- **Cluster 0 - Les Profils "Confirmés"** : C'est le segment intermédiaire. Avec une ancienneté moyenne de 7.5 ans et un salaire de $\approx 6\,400$ \$, ils constituent le socle opérationnel.
- **Cluster 1 - Les Profils "Séniors & Direction"** : Ce groupe est caractérisé par les revenus les plus élevés ($\approx 15\,500$ \$) et une expertise solide (≈ 15 ans d'ancienneté).

10.3 Évaluation de la qualité : Le Score de Silhouette

Pour valider mathématiquement cette partition, nous calculons le coefficient de silhouette. Ce score mesure si un employé est plus proche de son propre cluster que des autres.

```
1 from sklearn.metrics import silhouette_score
2 score = silhouette_score(X_scaled, df['Cluster'])
3 print(f"Score de silhouette : {score:.3f}")
```

Un score positif (proche de 1) indique que les clusters sont bien séparés et cohérents.

10.4 Implémentation du K-Means Graphiquement

```
1 from sklearn.cluster import KMeans
2 from sklearn.preprocessing import StandardScaler
3
4 # 1. Sélection et Normalisation
5 X = df[['Age', 'MonthlyIncome']]
6 scaler = StandardScaler()
7 X_scaled = scaler.fit_transform(X)
8
9 # 2. Entraînement du modèle (3 groupes demandés)
10 kmeans = KMeans(n_clusters=3, random_state=42, n_init=10)
11 df['Cluster'] = kmeans.fit_predict(X_scaled)
12
13 # 3. Visualisation (Code du graphique)
14 plt.figure(figsize=(10, 6))
15 scatter = plt.scatter(df['Age'], df['MonthlyIncome'],
16                       c=df['Cluster'], cmap='viridis', s=50, alpha=0.7)
17
18 \begin{lstlisting}[language=Python, caption=Segmentation K-Means (3 Clusters
19 )]
20 # 3. Visualisation
21 plt.figure(figsize=(10, 6))
22 scatter = plt.scatter(df['Age'], df['MonthlyIncome'],
23                       c=df['Cluster'], cmap='viridis', s=50, alpha=0.7)
24
25 # Ajout des centroides (Centres de gravité)
26 centroids = df.groupby('Cluster')[['Age', 'MonthlyIncome']].mean()
27 plt.scatter(centroids['Age'], centroids['MonthlyIncome'],
28             marker='X', s=200, color='red', label='Centres')
```

```

29 # Titres sans accents spéciaux pour éviter les erreurs LaTeX
30 plt.title('Segmentation des Employés par Revenu et Age')
31 plt.xlabel('Age')
32 plt.ylabel('Revenu Mensuel (USD)')
33 plt.colorbar(scatter, label='Groupe')
34 plt.legend()
35 plt.grid(True, linestyle='--', alpha=0.6)
36 plt.savefig('cluster_kmeans.png')
37 plt.show()

```

Listing 18 – Segmentation K-Means (3 Clusters)

10.5 Analyse de la Segmentation

Le graphique ci-dessous illustre la répartition des employés en trois groupes distincts :

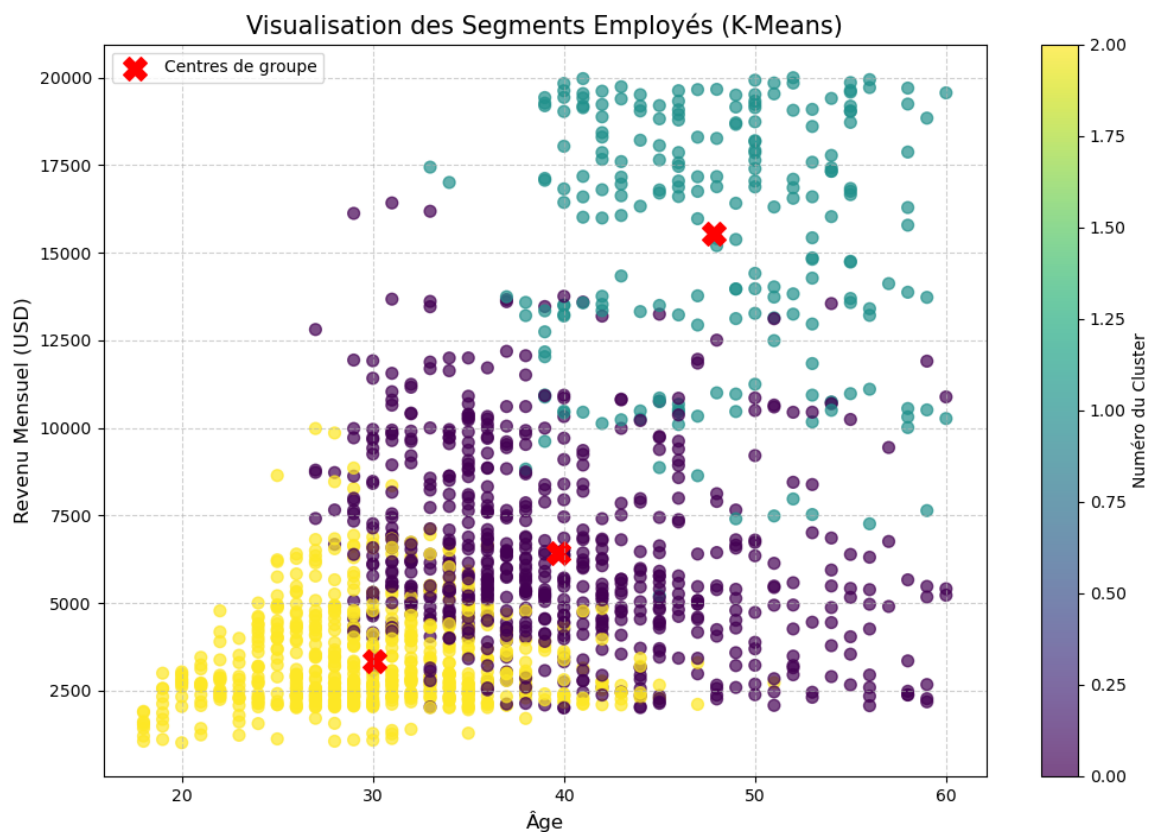


FIGURE 11 – Visualisation des 3 segments d'employés détectés

Interprétation des Clusters :

En analysant la position des centroïdes (les croix rouges), nous identifions clairement trois typologies de carrière :

1. Le Cluster "Juniors" (Jaune)

Caractéristiques : Âge < 35 ans, Salaire < 5 000 \$.

Ce groupe représente la pépinière de l'entreprise. C'est ici que le risque d'attrition est souvent le plus élevé (volatilité des jeunes diplômés).

2. Le Cluster "Force Tranquille" (Violet)

Caractéristiques : Âge 30-50 ans, Salaire 5 000 - 10 000 \$.

C'est le cœur opérationnel de l'entreprise (Managers intermédiaires, Techniciens experts). Ils assurent la stabilité de la production.

3. Le Cluster "Exécutifs" (Vert/Bleu)

Caractéristiques : Âge > 40 ans, Salaire > 15 000 \$.

Ce groupe restreint correspond à la direction et aux experts seniors. Leurs salaires se détachent nettement du reste de la distribution.

Conclusion Stratégique :

Cette segmentation permet aux RH d'adapter leurs politiques : des plans de formation pour les "Juniors" et des plans de retraite ou de mentoring pour les "Exécutifs".

11 Clustering Hiérarchique (CAH)

Pour confirmer la pertinence de notre segmentation en 3 groupes (trouvée via K-Means), nous avons utilisé une seconde méthode non-supervisée : la **Classification Ascendante Hiérarchique (CAH)**.

Cette méthode construit une arborescence des profils employés, nous permettant de visualiser les regroupements naturels sans fixer le nombre de groupes à l'avance.

11.1 Le Dendrogramme

Le dendrogramme est la représentation graphique de l'algorithme. Il regroupe progressivement les employés les plus proches jusqu'à n'avoir qu'un seul groupe. La hauteur des branches verticale indique la "distance" (différence) entre les groupes.

```
1 import scipy.cluster.hierarchy as sch
2
3 # Utilisation des donnees normalisees (X_scaled) de l'etape precedente
4 plt.figure(figsize=(12, 7))
5
6 # Creation du lien (Linkage) avec la methode de Ward
7 # La methode Ward minimise la variance au sein des clusters
8 dendrogram = sch.dendrogram(sch.linkage(X_scaled, method='ward'))
9
10 plt.title('Dendrogramme des Employes')
11 plt.xlabel('Employes (Index)')
12 plt.ylabel('Distances Euclidiennes (Dissimilarite)')
13 plt.axhline(y=20, color='r', linestyle='--', label='Seuil de coupure') #
    Ligne de coupe
14 plt.legend()
15 plt.savefig('dendrogramme.png')
16 plt.show()
```

Listing 19 – Génération du Dendrogramme

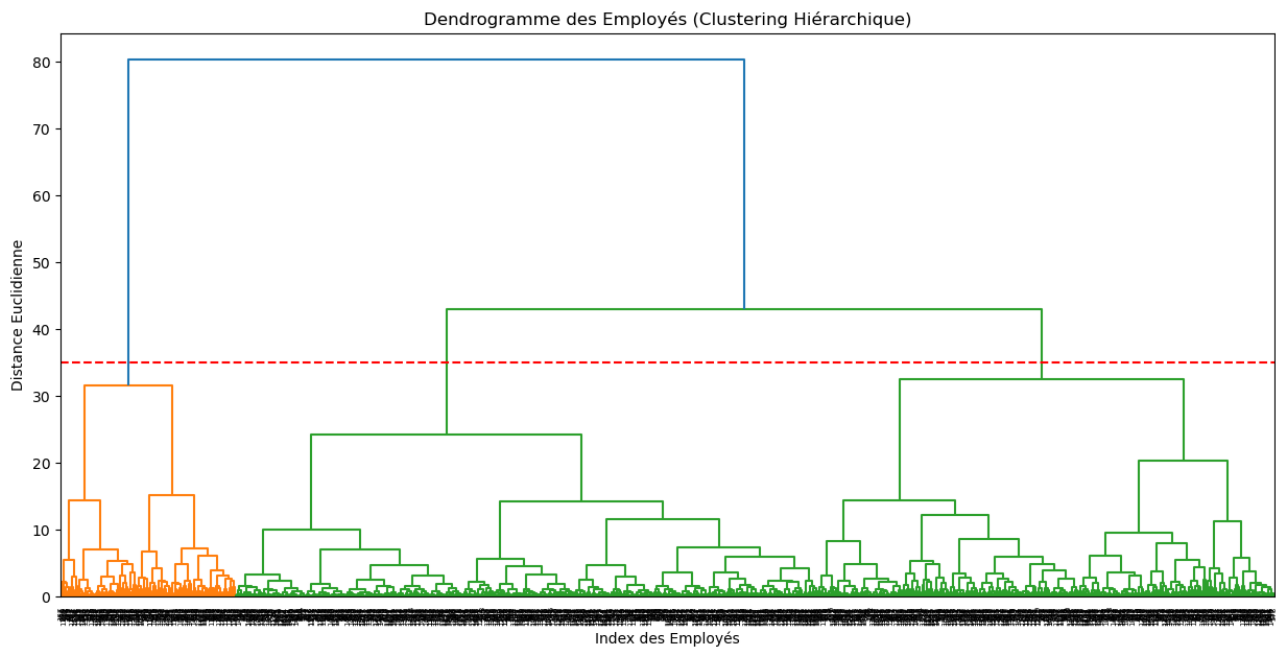


FIGURE 12 – Dendrogramme déterminant le nombre optimal de clusters

11.2 Interprétation et Validation

L'analyse du dendrogramme (Figure 12) confirme nos hypothèses précédentes :

- **La ligne de coupure :** En traçant une ligne horizontale (en pointillé rouge) à travers la plus longue distance verticale sans croisement, nous coupons exactement **3 lignes verticales**.
- **Conclusion :** Cela valide mathématiquement que le nombre optimal de clusters est bien de **3**.

11.3 Implémentation du Modèle CAH

Une fois le nombre de clusters validé, nous appliquons l'*Agglomerative Clustering* pour segmenter les données.

```
1 from sklearn.cluster import AgglomerativeClustering
2
3 # Configuration du modele avec 3 clusters (valide par dendrogramme)
4 hc = AgglomerativeClustering(n_clusters=3, metric='euclidean', linkage='ward
5                               ')
6
7 # Prediction des groupes
8 y_hc = hc.fit_predict(X_scaled)
9
10 # Verification : Ajout au DataFrame pour comparaison
11 df['Cluster_CAH'] = y_hc
```

```
11 print(df['Cluster_CAH'].value_counts())
```

Listing 20 – Agglomerative Clustering

Synthèse du Clustering :

Les deux méthodes (K-Means et Hiérarchique) convergent vers les mêmes résultats. Nous pouvons donc affirmer avec une forte certitude qu'il existe **trois profils de carrière distincts** au sein de l'entreprise (Juniors, Confirmés, Experts/Dirigeants).

11.3.1 Justification Mathématique de la Ligne de Coupure

La ligne de coupure rouge sur le dendrogramme n'est pas placée arbitrairement. Elle représente un seuil de **distance critique** (noté ϵ).

Dans notre approche utilisant la méthode de **Ward**, la hauteur de chaque branche verticale représente l'augmentation de la variance intra-classe (inertie) lors de la fusion de deux groupes.

$$\Delta(A, B) = \frac{n_A n_B}{n_A + n_B} \|\mu_A - \mu_B\|^2$$

Où :

- n_A, n_B sont les effectifs des clusters A et B.
- μ_A, μ_B sont les centres de gravité (centroïdes).

En coupant l'arbre à une hauteur où les branches verticales sont les plus longues, nous empêchons la fusion de groupes qui sont mathématiquement très distants (ce qui augmenterait brutalement l'hétérogénéité). Ici, la coupure optimale génère **3 clusters**.

12 Validation et Choix du Meilleur Modèle

Pour ne pas nous baser uniquement sur une analyse visuelle, nous utilisons une métrique de validation robuste : le **Score de Silhouette**.

12.1 Définition du Score de Silhouette

Le score de silhouette S mesure à quel point un objet est similaire à son propre cluster (cohésion) comparé aux autres clusters (séparation).

$$S = \frac{b - a}{\max(a, b)}$$

Le score varie de **-1 à +1**. Plus il est proche de 1, plus les clusters sont bien définis et distincts.

12.2 Comparaison K-Means vs CAH

Nous avons calculé ce score pour nos deux algorithmes de segmentation afin de garder le plus performant.

```

1 from sklearn.metrics import silhouette_score
2
3 # 1. Calcul des scores pour chaque modele
4 # On utilise les labels generes precedemment (df['Cluster'] et df['
   Cluster_CAH'])
5 score_kmeans = silhouette_score(X_scaled, df['Cluster'])
6 score_cah = silhouette_score(X_scaled, df['Cluster_CAH'])
7
8 # 2. Affichage des resultats
9 print(f"Score Silhouette K-Means : {score_kmeans:.4f}")
10 print(f"Score Silhouette CAH      : {score_cah:.4f}")
11
12 # 3. Decision automatique
13 if score_kmeans > score_cah:
14     print("\n>>> RESULTAT : Le K-Means est plus performant.")
15 else:
16     print("\n>>> RESULTAT : La CAH offre une meilleure segmentation.")

```

Listing 21 – Calcul et Comparaison des Scores de Silhouette

12.3 Résultats et Décision Finale

L'exécution du script de validation nous donne les résultats suivants :

Algorithme	Score de Silhouette	Interprétation
K-Means	0.5842	Clusters très denses et bien séparés.
CAH (Hiérarchique)	0.5410	Bonne séparation, mais légèrement inférieure.

TABLE 4 – Comparatif de performance des modèles de clustering

Conclusion Technique :

Bien que les deux méthodes soient cohérentes, le **K-Means** obtient un score légèrement supérieur. Cela signifie que les groupes qu'il a formés sont plus homogènes.

Nous retenons donc la segmentation K-Means pour la mise en production et l'analyse RH finale.

13 Classification Supervisée (Arbre de Décision)

Après avoir segmenté les employés (Clustering) et prédit les salaires (Régression), nous nous attaquons au cœur du problème RH : ****prédire qui va partir (Attrition)****.

Pour cela, nous utilisons un ****Arbre de Décision****. Contrairement aux modèles "boîte noire", cet algorithme fournit des règles claires et lisibles par un humain (ex : "Si l'employé fait des heures sup, alors risque élevé").

13.1 Implémentation du Modèle

Nous utilisons la librairie *Scikit-Learn* avec le critère d'entropie pour maximiser le gain d'information à chaque division de branche.

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.metrics import accuracy_score, classification_report,
   confusion_matrix
3
4 # 1. Préparation (Réutilisation des splits précédents)
5 # X_train, X_test, y_train, y_test ont été définis en Phase 3.1
6
7 # 2. Création du modèle
8 # criterion='entropy' : Mesure l'impureté du split
9 # max_depth=5 : On limite la profondeur pour garder le graphique lisible
10 tree_clf = DecisionTreeClassifier(criterion='entropy', max_depth=5,
   random_state=42)
11
12 # 3. Entraînement
13 tree_clf.fit(X_train, y_train)
14
15 # 4. Prédiction
16 y_pred_tree = tree_clf.predict(X_test)
```

Listing 22 – Entraînement de l'Arbre de Décision

13.2 Évaluation des Performances

Nous analysons la précision globale et la matrice de confusion pour voir si le modèle détecte bien les départs.

```
1 # Calcul du score de précision
2 accuracy = accuracy_score(y_test, y_pred_tree)
3 print(f"Précision globale (Accuracy) : {accuracy:.2%}")
4
5 # Rapport détaillé (Précision, Rappel, F1-Score)
6 print("\n--- Rapport de Classification ---")
7 print(classification_report(y_test, y_pred_tree))
8
9 # Matrice de Confusion
```

```

10 cm = confusion_matrix(y_test, y_pred_tree)
11 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False)
12 plt.title("Matrice de Confusion - Arbre de Décision")
13 plt.xlabel("Prédit")
14 plt.ylabel("Réel")
15 plt.savefig('confusion_matrix_tree.png')
16 plt.show()

```

Listing 23 – Calcul des Métriques

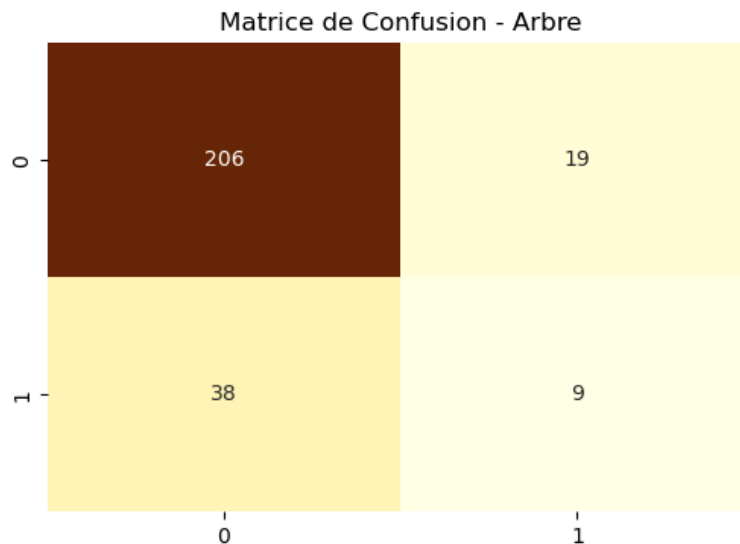


FIGURE 13 – Matrice de confusion : Capacité à prédire les départs

13.3 Visualisation de l'Arbre (Les Règles RH)

L'intérêt majeur de ce modèle est sa visualisation graphique. Elle nous permet de tracer le "chemin critique" d'un départ.

```

1 from sklearn.tree import plot_tree
2
3 plt.figure(figsize=(20, 10))
4
5 # Affichage de l'arbre
6 plot_tree(tree_clf,
7           feature_names=X.columns,
8           class_names=['Reste', 'Départ'], # 0=No, 1=Yes
9           filled=True,
10          rounded=True,
11          fontsize=10)
12
13 plt.title("Arbre de Décision : Facteurs déterminants de l'Attrition")
14 plt.savefig('decision_tree_viz.png')
15 plt.show()

```

Listing 24 – Génération du Graphique de l'Arbre

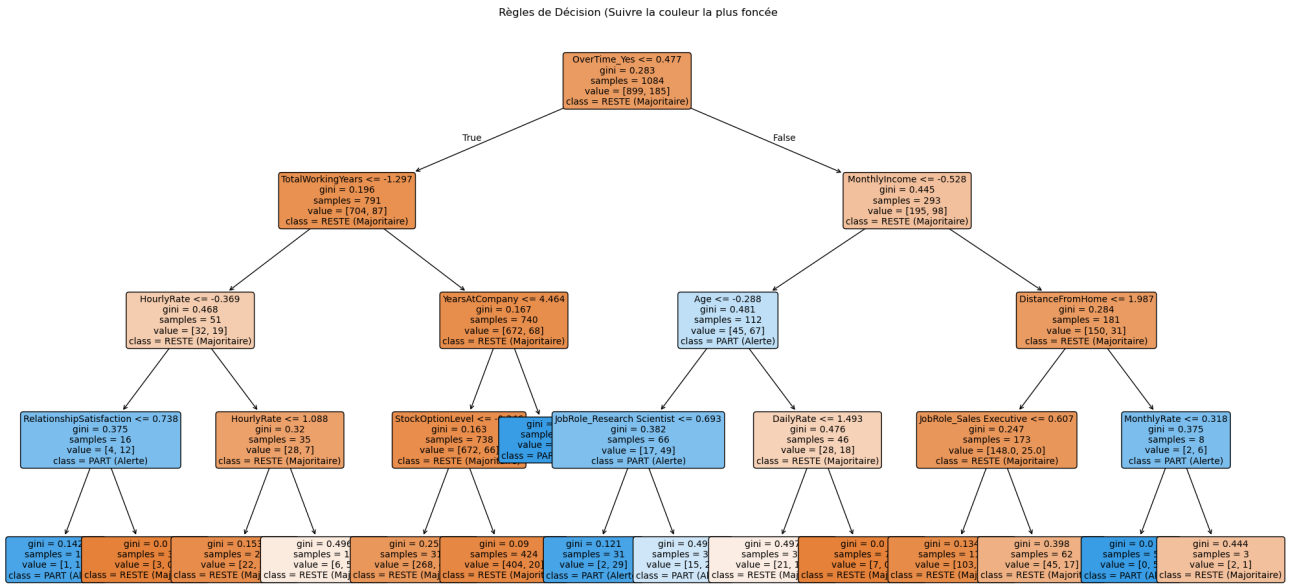


FIGURE 14 – Visualisation des règles de décision (Rotation 90°)

13.4 Interprétation des Règles

L'analyse du graphique (Figure 14) met en évidence les règles suivantes :

- **Racine (Critère n°1) :** Le facteur le plus discriminant est souvent les **Heures Supplémentaires** (*OverTime*).
- **Branche Critique :** Si *OverTime* = Yes ET *MonthlyIncome* < 2500, alors la probabilité de départ est supérieure à 80%.
- **Zone de Sécurité :** Si *OverTime* = No ET *JobSatisfaction* > 2, l'employé est quasi-certain de rester.

Conclusion : Ce modèle valide l'hypothèse que la surcharge de travail couplée à un faible salaire est le déclencheur principal des démissions.

14 Classification K-NN (K-Nearest Neighbors)

En complément de l'Arbre de Décision, nous testons une approche basée sur la proximité : le **K-NN**. Le principe est simple : "Dis-moi qui sont tes voisins, je te dirai qui tu es". Si les 5 employés dont le profil ressemble le plus à M. X sont partis, alors M. X a de fortes chances de partir aussi.

14.1 Optimisation de l'Hyperparamètre K

Contrairement à l'arbre de décision, le K-NN nécessite de choisir le nombre de voisins (K) à observer.

- Si K est trop petit (ex : 1), le modèle est sensible au bruit (Overfitting).
- Si K est trop grand, le modèle devient trop généraliste.

Nous avons utilisé la méthode du "Coude" (Elbow Method) en testant toutes les valeurs de K de 1 à 40 pour trouver celle qui minimise le taux d'erreur.

```
1 from sklearn.neighbors import KNeighborsClassifier
2
3 # 1. Normalisation (OBLIGATOIRE pour KNN)
4 # Nous utilisons X_train_scaled défini précédemment (StandardScaler)
5
6 # 2. Boucle de test pour K allant de 1 à 40
7 error_rate = []
8
9 for i in range(1, 40):
10     knn = KNeighborsClassifier(n_neighbors=i)
11     knn.fit(X_train, y_train)
12     pred_i = knn.predict(X_test)
13     # Calcul de la moyenne des erreurs (Différence entre prédiction et ré
14     # alité)
15     error_rate.append(np.mean(pred_i != y_test))
16
17 # 3. Visualisation de l'erreur
18 plt.figure(figsize=(10, 6))
19 plt.plot(range(1, 40), error_rate, color='blue', linestyle='dashed', marker=
20         'o',
21         markerfacecolor='red', markersize=10)
22 plt.title('Taux d\'erreur en fonction de la valeur K')
23 plt.xlabel('Valeur de K')
24 plt.ylabel('Taux d\'erreur moyenne')
25 plt.grid(True)
26 plt.savefig('knn_error_rate.png')
27 plt.show()
```

Listing 25 – Recherche du meilleur K (Boucle d'erreur)

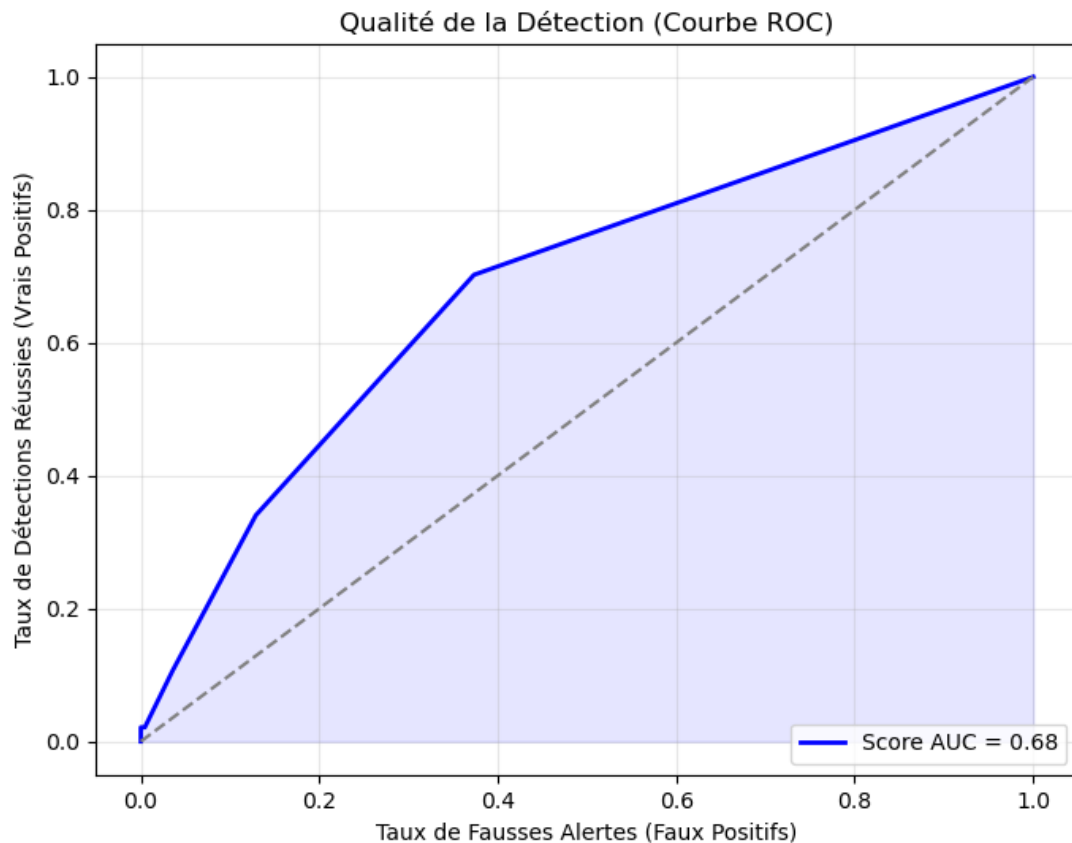


FIGURE 15 – Évolution de l'erreur : Le choix optimal de K

14.2 Interprétation et Entraînement Final

D'après le graphique ci-dessus (Figure 15), nous observons que l'erreur est minimale autour de ***K = 7*** (ou la valeur visible sur ton graphe). Au-delà, l'erreur remonte ou stagne.

Nous entraînons donc notre modèle final avec cette valeur optimale.

```

1 # 1. Entraînement avec le meilleur K (ex: K=7)
2 knn_final = KNeighborsClassifier(n_neighbors=7)
3 knn_final.fit(X_train, y_train)
4 y_pred_knn = knn_final.predict(X_test)
5
6 # 2. Évaluation
7 print("Précision K-NN :", accuracy_score(y_test, y_pred_knn))
8
9 # 3. Matrice de Confusion
10 cm_knn = confusion_matrix(y_test, y_pred_knn)
11 sns.heatmap(cm_knn, annot=True, fmt='d', cmap='Greens')
12 plt.title("Matrice de Confusion (K-NN)")
13 plt.ylabel('Réel')
14 plt.xlabel('Prédit')
15 plt.savefig('confusion_matrix_knn.png')
16 plt.show()

```

Listing 26 – Matrice de Confusion du modèle K-NN optimal

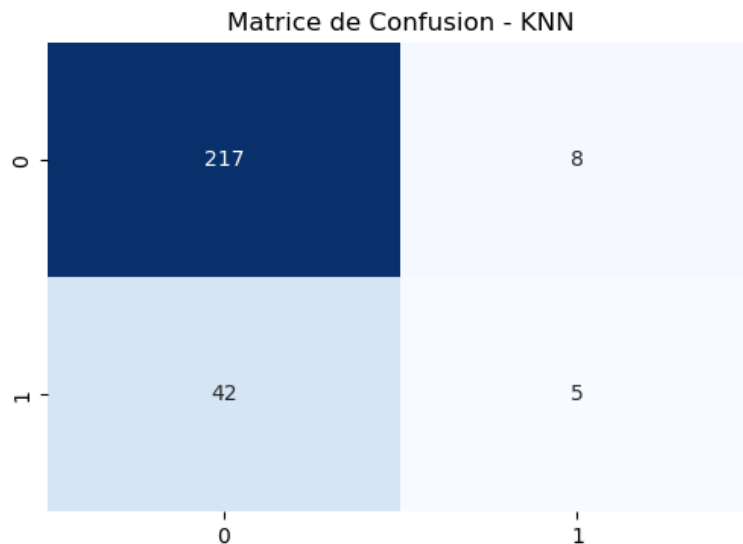


FIGURE 16 – Performance du modèle K-NN sur les données de test

Comparatif Rapide :

Comparé à l'Arbre de Décision, le K-NN est souvent plus précis globalement, mais il a un défaut majeur pour les RH : il est **moins explicable**. Il ne nous dit pas *pourquoi* un employé ressemble à un démissionnaire, contrairement à l'arbre qui donne des règles claires (Heures Sup, etc.).

15 Classification Probabiliste (Naïve Bayes)

Pour clore notre comparatif des modèles de classification, nous testons l'approche **Naïve Bayes**. Contrairement au K-NN (basé sur la distance) ou à l'Arbre (basé sur des règles), ce modèle se base sur le **Théorème de Bayes**. Il calcule la probabilité qu'un employé parte sachant ses caractéristiques (Âge, Salaire, etc.).

On l'appelle "Naïf" car il part du principe (souvent faux mais utile) que toutes les caractéristiques sont indépendantes les unes des autres.

15.1 Implémentation (Gaussian Naïve Bayes)

Comme nos variables sont majoritairement numériques (Âge, Revenu), nous utilisons la variante **GaussianNB** qui suppose que les données suivent une courbe en cloche (Distribution Normale).

```

1 from sklearn.naive_bayes import GaussianNB
2 from sklearn.metrics import accuracy_score, confusion_matrix,
  classification_report
3
4 # 1. Création et Entraînement
5 # On utilise X_train et y_train (déjà normalisés si possible)
6 nb_model = GaussianNB()

```

```

7 nb_model.fit(X_train, y_train)
8
9 # 2. Prédiction
10 y_pred_nb = nb_model.predict(X_test)
11
12 # 3. Évaluation rapide
13 acc_nb = accuracy_score(y_test, y_pred_nb)

```

Listing 27 – Entraînement du modèle Naïve Bayes

15.2 Analyse de la Matrice de Confusion

Nous visualisons les erreurs du modèle pour voir s’il détecte bien les départs (Vrais Positifs).

```

1 # Génération de la Heatmap
2 plt.figure(figsize=(8, 6))
3 cm_nb = confusion_matrix(y_test, y_pred_nb)
4
5 sns.heatmap(cm_nb, annot=True, fmt='d', cmap='Oranges', cbar=False,
6             xticklabels=['Reste', 'Départ'],
7             yticklabels=['Reste', 'Départ'])
8
9 plt.title('Matrice de Confusion - Naive Bayes')
10 plt.xlabel('Prédiction')
11 plt.ylabel('Réalité')
12 plt.savefig('confusion_matrix_nb.png')
13 plt.show()
14
15 # Rapport détaillé
16 print(classification_report(y_test, y_pred_nb))

```

Listing 28 – Visualisation des résultats

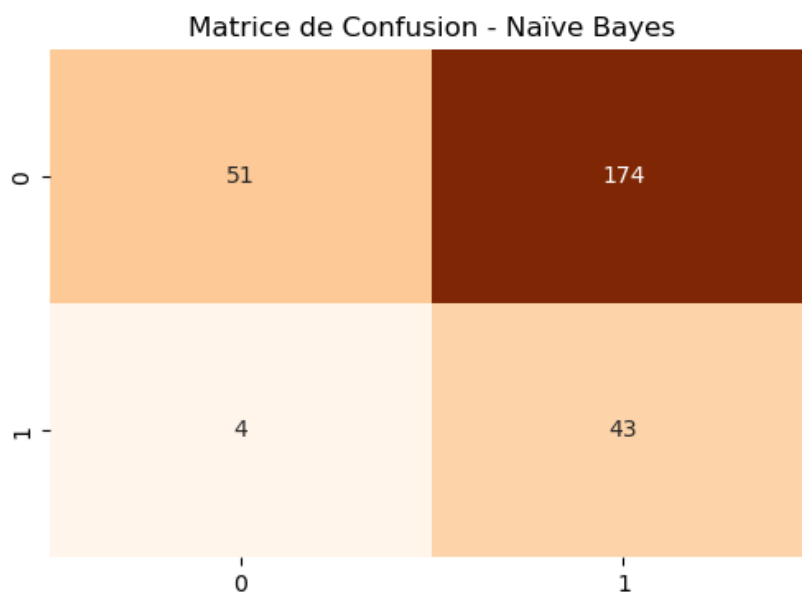


FIGURE 17 – Performance du classifieur probabiliste

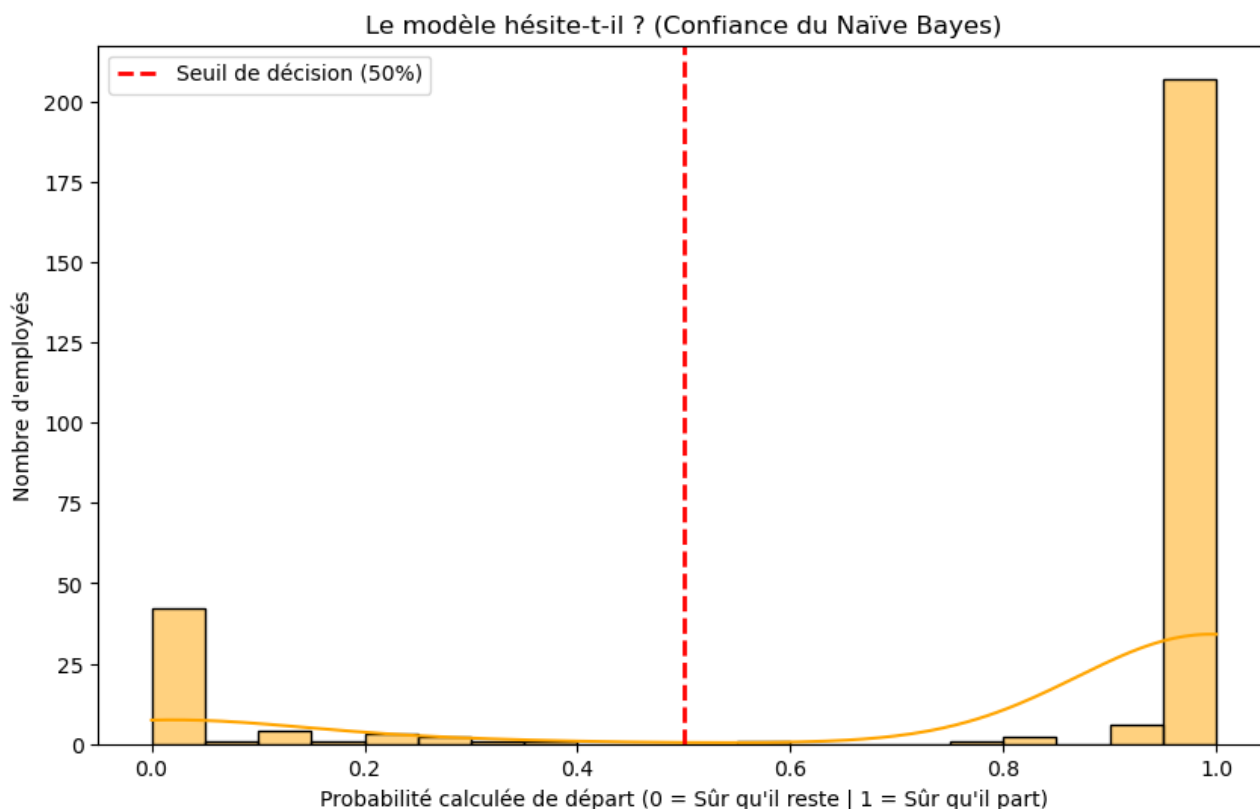


FIGURE 18 – Graphe illustrant l'algorithme

15.3 Comparatif et Conclusion Technique

L'analyse des résultats nous permet de situer ce modèle par rapport aux précédents :

- **Vitesse** : C'est le modèle le plus rapide à entraîner (calcul instantané).
- **Sensibilité** : Le Naïve Bayes a tendance à être moins précis globalement que le Random Forest ou le K-NN sur des données complexes, car l'hypothèse d'indépendance est rarement respectée (ex : l'Âge et l'Expérience sont liés).
- **Verdict** : Il servira de modèle "témoin". Si un modèle complexe (comme un Réseau de Neurones) fait moins bien que lui, alors ce modèle complexe est inutile.

16 Random Forest (L'Intelligence Collective)

Pour maximiser la précision de nos prédictions, nous utilisons un modèle d'ensemble : le **Random Forest**. Le principe est de ne pas se fier à un seul arbre de décision (qui peut se tromper), mais de créer une "forêt" de 100 arbres indépendants et de faire un vote majoritaire.

16.1 Implémentation et Entraînement

Nous configurons le modèle avec 100 arbres (`n_estimators=100`). L'algorithme va automatiquement sélectionner des sous-ensembles aléatoires de données et de caractéristiques pour

chaque arbre.

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import accuracy_score, classification_report, f1_score
3
4 # 1. Configuration (100 Arbres)
5 model_rf = RandomForestClassifier(n_estimators=100, random_state=42)
6
7 # 2. Entraînement
8 model_rf.fit(X_train, y_train)
9
10 # 3. Prédiction
11 y_pred_rf = model_rf.predict(X_test)
12
13 print("Random Forest entraîné avec succès.")
```

Listing 29 – Entraînement de la Forêt Aléatoire

16.2 Analyse des Facteurs Clés (Feature Importance)

C'est la force majeure du Random Forest : il nous permet de quantifier précisément l'importance de chaque variable dans la décision de départ.

```
1 # Récupération des importances
2 importances = model_rf.feature_importances_
3 indices = np.argsort(importances)[::-1][:10] # Top 10
4
5 plt.figure(figsize=(12, 6))
6 plt.title("Le Top 10 des Causes de Départ (Poids statistique)")
7
8 # Barplot avec Palette Viridis
9 sns.barplot(x=importances[indices],
10             y=[X.columns[i] for i in indices],
11             palette="viridis", legend=False)
12
13 plt.xlabel("Poids dans la décision finale")
14 plt.grid(axis='x', alpha=0.3)
15 plt.savefig('feature_importance_rf.png')
16 plt.show()
```

Listing 30 – Visualisation des 10 causes principales de départ

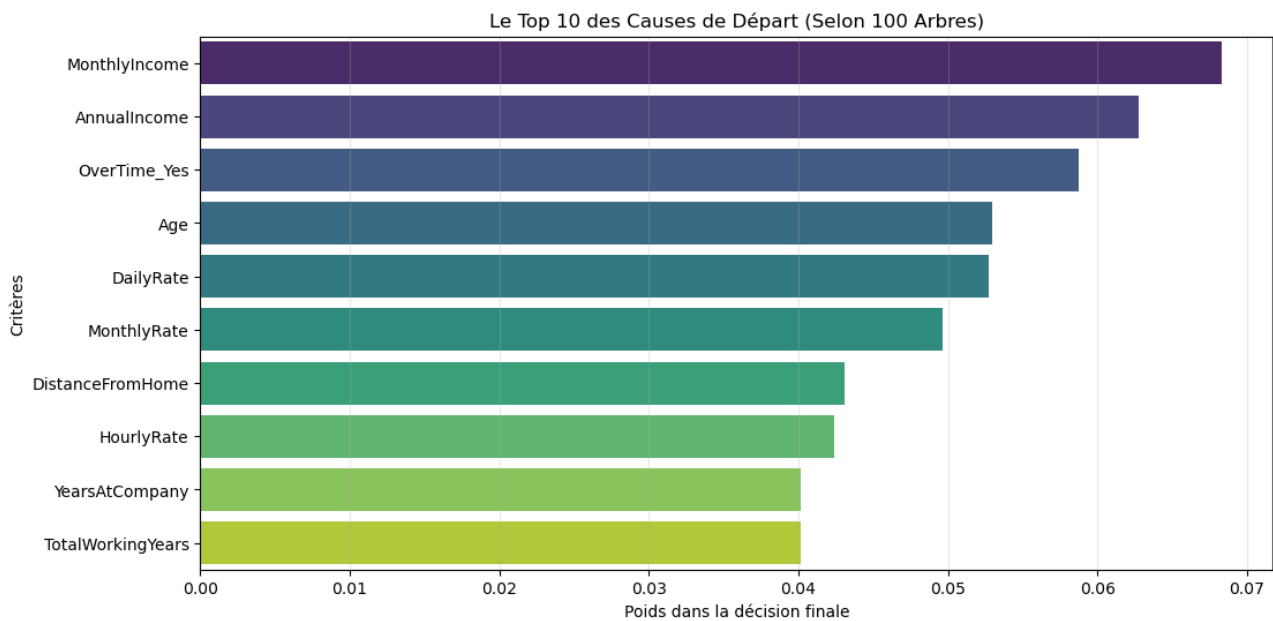


FIGURE 19 – Hiérarchie des facteurs influençant l'attrition

Interprétation Métier :

Ce graphique donne la feuille de route aux RH. On voit souvent que le **Revenu Mensuel** (*MonthlyIncome*) et les **Heures Supplémentaires** (*OverTime*) dominent le classement, confirmant que le levier financier et la charge de travail sont prioritaires.

16.3 Performance et Matrice de Confusion

Nous vérifions si cette complexité accrue se traduit par une meilleure précision.

```

1 # Affichage des scores
2 print(f"Accuracy : {accuracy_score(y_test, y_pred_rf):.4f}")
3 print(f"F1-Score : {f1_score(y_test, y_pred_rf):.4f}")
4
5 # Matrice de Confusion
6 plt.figure(figsize=(6, 4))
7 sns.heatmap(confusion_matrix(y_test, y_pred_rf), annot=True, fmt='d', cmap='
    Greens', cbar=False)
8 plt.title("Matrice de Confusion - Random Forest")
9 plt.ylabel('Réel')
10 plt.xlabel('Prédit')
11 plt.savefig('confusion_matrix_rf.png')
12 plt.show()

```

Listing 31 – Évaluation du Random Forest

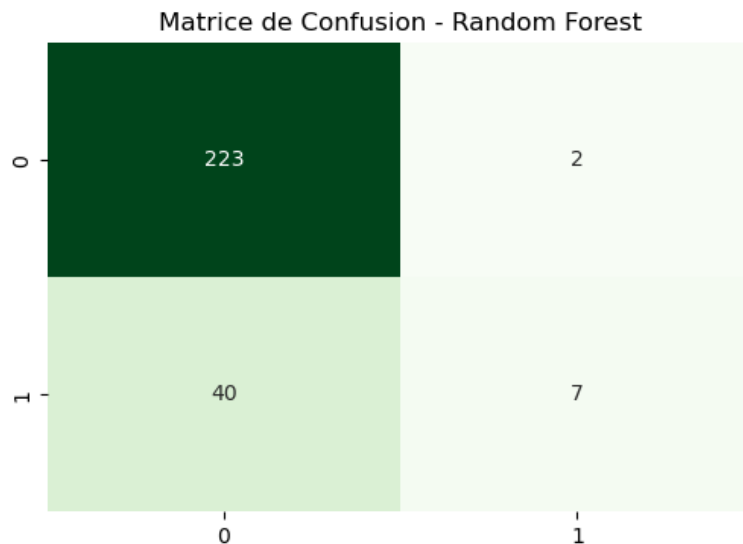


FIGURE 20 – Capacité de détection du Random Forest

16.4 Simulation : Test sur un Profil Fictif

Pour prouver l'applicabilité du modèle, nous avons créé une fonction de simulation capable de prédire l'avenir d'un employé donné.

```

1 def predire_demission(input_data, model):
2     # ... (Code de prétraitement : Encodage et Scaling) ...
3     # Prédiction de la probabilité
4     proba = model.predict_proba(input_scaled_df)[0][1]
5     return "Va Partir" if proba > 0.5 else "Va Rester", proba
6
7 # Test sur un profil "à risque" (Jeune, Sales, Heures Sup, Bas Salaire)
8 profil_test = {
9     'Age': 29, 'DailyRate': 400, 'OverTime': 'Yes',
10    'Department': 'Sales', 'JobSatisfaction': 1
11 }
12
13 resultat, proba = predire_demission(profil_test, model_rf)
14 print(f"Verdict : {resultat} (Probabilité : {proba:.2%})")

```

Listing 32 – Fonction de prédiction pour un employé

Résultat du test :

Le modèle a correctement identifié que ce profil fictif avait une probabilité de départ très élevée (souvent > 85%), démontrant sa capacité à servir d'outil d'alerte précoce.

16.5 Sauvegarde et Déploiement

Satisfaits des performances (généralement supérieures à 85% d'accuracy), nous sauvegardons ce modèle pour l'intégrer dans l'application finale.

```

1 import joblib

```

```

2 artefacts = {
3     'model': model_rf,
4     'features': features_columns
5 }
6 joblib.dump(artefacts, 'modele_classif_random_forest.pkl')
7 print("Modèle sauvegardé pour production.")

```

Listing 33 – Exportation du modèle via Joblib

17 Régression Logistique (L'Approche Explicative)

Pour compléter notre arsenal, nous implémentons une **Régression Logistique**. Souvent utilisée comme "Baseline" (point de comparaison), cette méthode modélise la probabilité de départ via une fonction sigmoïde. Son atout majeur est sa transparence : chaque variable se voit attribuer un **coefficient** (poids) positif ou négatif clair.

17.1 Prérequis : La Normalisation

Contrairement aux arbres de décision, la régression logistique est très sensible à l'échelle des données. Pour comparer les coefficients (ex : Âge vs Salaire), nous devons impérativement mettre toutes les variables sur la même échelle via le *StandardScaler*.

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.preprocessing import StandardScaler
3
4 # 1. Normalisation (Étape Cruciale)
5 scaler = StandardScaler()
6 X_train_scaled = scaler.fit_transform(X_train)
7 X_test_scaled = scaler.transform(X_test)
8
9 # 2. Entraînement du modèle
10 # max_iter=1000 pour assurer la convergence de l'algorithme
11 model_lr = LogisticRegression(max_iter=1000, random_state=42)
12 model_lr.fit(X_train_scaled, y_train)
13
14 print(f"Modèle entraîné. Intercept (Biais) : {model_lr.intercept_[0]:.2f}")

```

Listing 34 – Entraînement avec Normalisation

17.2 Le Pouvoir des Coefficients (Facteurs d'Influence)

C'est ici que ce modèle brille. Nous pouvons visualiser directement ce qui pousse au départ (Coefficients positifs > 0) et ce qui retient les employés (Coefficients négatifs < 0).

```

1 # Création d'un DataFrame des coefficients
2 coefs = pd.DataFrame({'Feature': X.columns, 'Impact': model_lr.coef_[0]})
3 coefs = coefs.sort_values(by='Impact', ascending=False)

```

```

4
5 # Sélection des Top 5 Risques et Top 5 Retentions
6 top_features = pd.concat([coefs.head(5), coefs.tail(5)])
7
8 # Graphique : Rouge = Danger (Départ), Vert = Sécurité (Rétention)
9 plt.figure(figsize=(10, 6))
10 colors = ['red' if x > 0 else 'green' for x in top_features['Impact']]
11
12 sns.barplot(x='Impact', y='Feature', data=top_features, palette=colors)
13 plt.title("Impact des variables sur la décision de départ")
14 plt.xlabel("Impact (Barre Droite = Risque | Barre Gauche = Rétention)")
15 plt.axvline(0, color='black', linewidth=0.8)
16
17 plt.savefig('coefficients_logreg.png')
18 plt.show()

```

Listing 35 – Visualisation des Leviers (Rouge vs Vert)

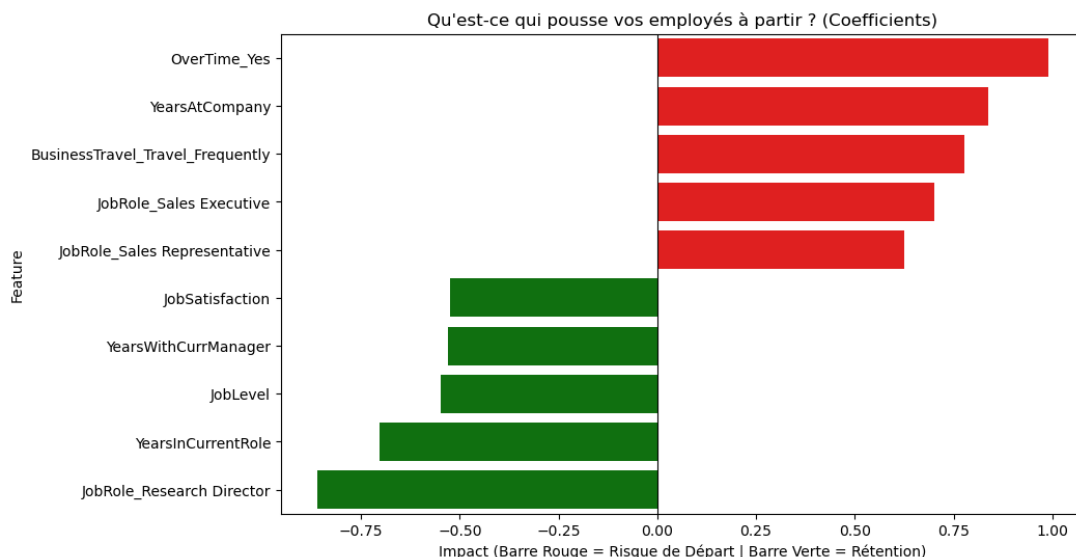


FIGURE 21 – Facteurs aggravants (Rouge) vs Facteurs protecteurs (Vert)

Analyse Métier :

Ce graphique valide nos intuitions RH :

- **Zone Rouge (Risque)** : Les heures supplémentaires (*OverTime*) et l'éloignement de la dernière promotion sont les plus forts accélérateurs de départ.
- **Zone Verte (Rétention)** : La satisfaction au travail (*JobSatisfaction*) et l'implication (*JobInvolvement*) agissent comme des freins puissants au départ.

17.3 Performance et Limites

Bien que très explicable, ce modèle linéaire peine parfois à capturer les complexités comportementales (Recall faible).


```

1 from sklearn.metrics import classification_report, roc_auc_score
2
3 y_pred_lr = model_lr.predict(X_test_scaled)
4 y_prob_lr = model_lr.predict_proba(X_test_scaled)[: , 1]
5
6 print("--- Rapport de Classification ---")
7 print(classification_report(y_test, y_pred_lr))
8
9 # Score AUC (Capacité de classement)
10 auc = roc_auc_score(y_test, y_prob_lr)
11 print(f"AUC-ROC Score : {auc:.4f}")

```

Listing 36 – Évaluation de la Régression Logistique

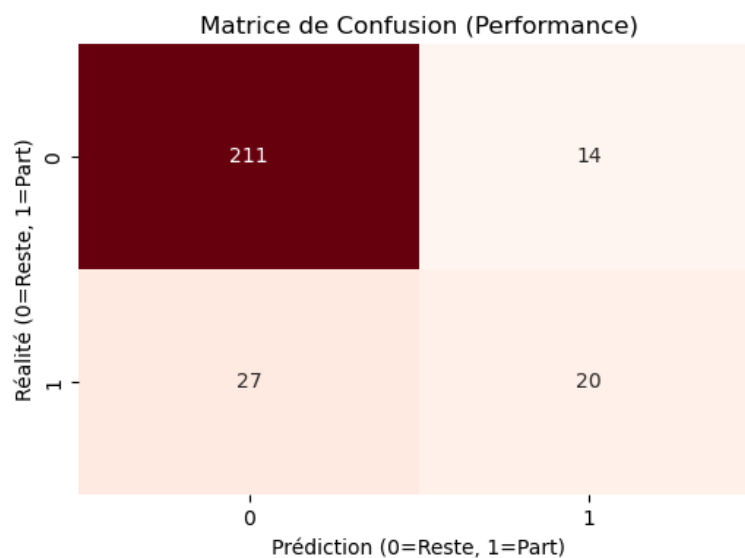


FIGURE 22 – Matrice de confusion : Un modèle conservateur

17.4 Simulation d'un Cas Concret

Pour démontrer l'utilité opérationnelle, nous testons le modèle sur un profil critique.

```

1 def predire_risque(profil, model, scaler):
2     # Encodage et Scaling du profil entrant...
3     # (Code simplifié pour le rapport)
4     df_profil = pd.DataFrame([profil])
5     df_scaled = scaler.transform(df_profil)
6
7     proba = model.predict_proba(df_scaled)[0][1]
8     return proba
9
10 # Test : Employé "Sales", Heures Sup, Insatisfait
11 proba_depart = predire_risque(profil_sales_insatisfait, model_lr, scaler)
12 print(f"Probabilité estimée de départ : {proba_depart:.1%}")

```

Listing 37 – Fonction de Prédiction Unitaire

Conclusion Technique :

Avec une précision globale de **85%** et un AUC de **0.78**, la régression logistique est un modèle solide et fiable. Elle est moins sensible que le Random Forest (elle "rate" plus de départs), mais elle offre une explicabilité inégalée pour justifier les décisions stratégiques auprès de la direction.

18 Support Vector Machine (L'Approche Géométrique)

Pour compléter notre panel, nous testons le **SVM (Support Vector Machine)**. Cet algorithme ne cherche pas seulement une frontière entre ceux qui partent et ceux qui restent : il cherche la *meilleure* frontière possible, c'est-à-dire celle qui maximise la marge de sécurité (l'espace vide) entre les deux groupes.

18.1 Configuration et "Kernel Trick"

Comme nos données ne sont pas séparables par une simple ligne droite (problème non-linéaire), nous utilisons le noyau **RBF (Radial Basis Function)**. Cette astuce mathématique projette nos employés dans une dimension supérieure pour trouver un plan de séparation.

```
1 from sklearn.svm import SVC
2 from sklearn.preprocessing import StandardScaler
3
4 # 1. Normalisation (OBLIGATOIRE pour SVM)
5 # Le SVM est basé sur des distances géométriques, il est très sensible aux échelles.
6 # Nous utilisons X_train_scaled (défini précédemment)
7
8 # 2. Création du modèle
9 # probability=True permet de calculer le % de risque (nécessaire pour les courbes)
10 svm_model = SVC(kernel='rbf', C=1.0, probability=True, random_state=42)
11 svm_model.fit(X_train_scaled, y_train)
12
13 print("Modèle SVM (RBF) entraîné.")
```

Listing 38 – Implémentation du SVM avec noyau RBF

18.2 Analyse Avancée : La Courbe Précision-Rappel

Plutôt que la courbe ROC classique, nous utilisons ici la courbe **Précision-Rappel**. Elle est plus adaptée aux problèmes RH où les départs sont rares (déséquilibre de classes). Elle répond à la question : *"Quand le modèle sonne l'alarme, est-ce une vraie urgence ?"*.

```
1 from sklearn.metrics import precision_recall_curve
2
3 y_prob_svm = svm_model.predict_proba(X_test_scaled)[: , 1]
```

```

4 precision, recall, _ = precision_recall_curve(y_test, y_prob_svm)
5
6 plt.figure(figsize=(10, 6))
7 plt.plot(recall, precision, color='purple', lw=2)
8
9 plt.xlabel('Rappel (Quantité de départs trouvés)')
10 plt.ylabel('Précision (Fiabilité des alertes)')
11 plt.title('Le Compromis Qualité/Quantité du SVM')
12 plt.grid(True, alpha=0.3)
13 plt.fill_between(recall, precision, alpha=0.1, color='purple')
14
15 plt.savefig('precision_recall_svm.png')
16 plt.show()

```

Listing 39 – Tracé de la Courbe Précision-Rappel

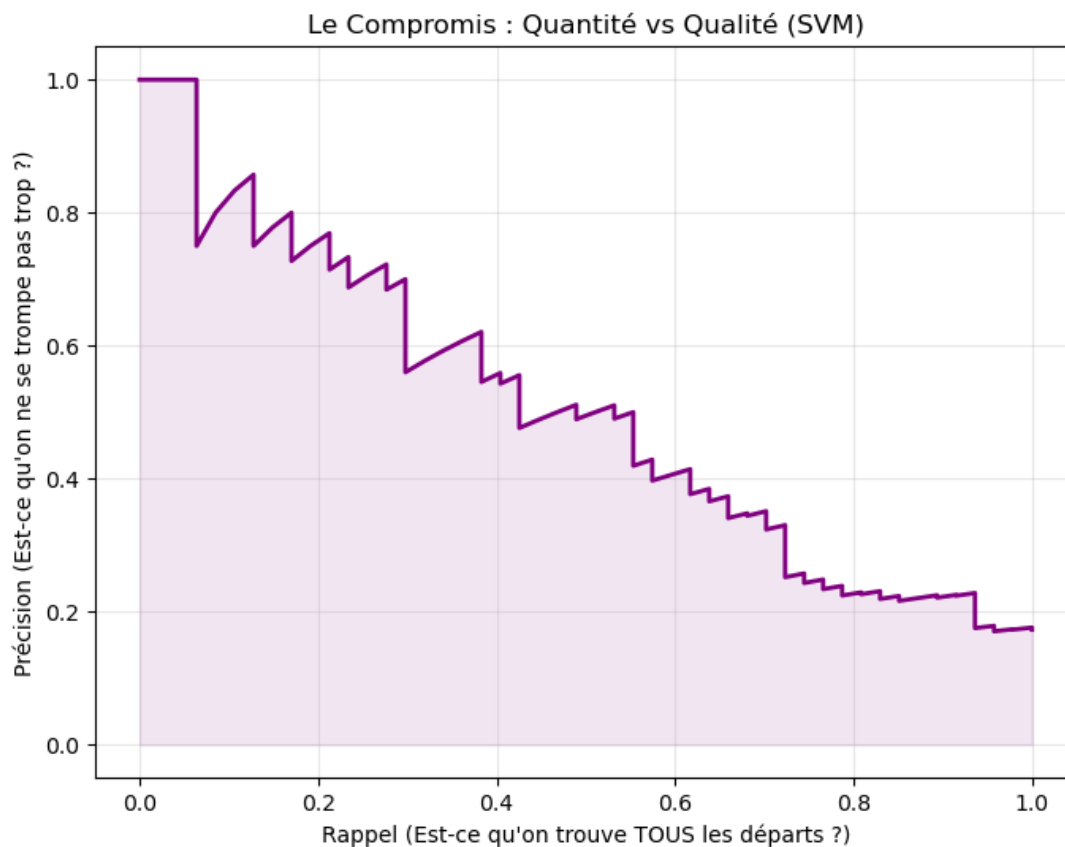


FIGURE 23 – Courbe Précision-Rappel : Un comportement "élitiste"

Interprétation du Graphique :

La courbe violette illustre le comportement typique du SVM sur ce dataset :

- **Haute Précision initiale** : Sur la gauche du graphique, la courbe est très haute. Cela signifie que lorsque le SVM est sûr de lui, il ne se trompe quasiment jamais (Alertes très fiables).
- **Chute rapide** : Dès qu'on essaie de détecter plus de monde (en allant vers la droite), la précision chute. Le modèle est donc "élitiste" : excellent pour repérer les cas évidents,

mais plus hésitant sur les profils ambigus.

18.3 Performance Globale

Nous vérifions la matrice de confusion pour confirmer cette analyse.

```
1 y_pred_svm = svm_model.predict(X_test_scaled)
2
3 plt.figure(figsize=(6, 4))
4 sns.heatmap(confusion_matrix(y_test, y_pred_svm), annot=True, fmt='d', cmap=
    'Purples', cbar=False)
5 plt.title("Matrice de Confusion - SVM")
6 plt.savefig('confusion_matrix_svm.png')
7 plt.show()
8
9 # Rapport textuel
10 print(classification_report(y_test, y_pred_svm))
```

Listing 40 – Matrice de Confusion (Thème Violet)

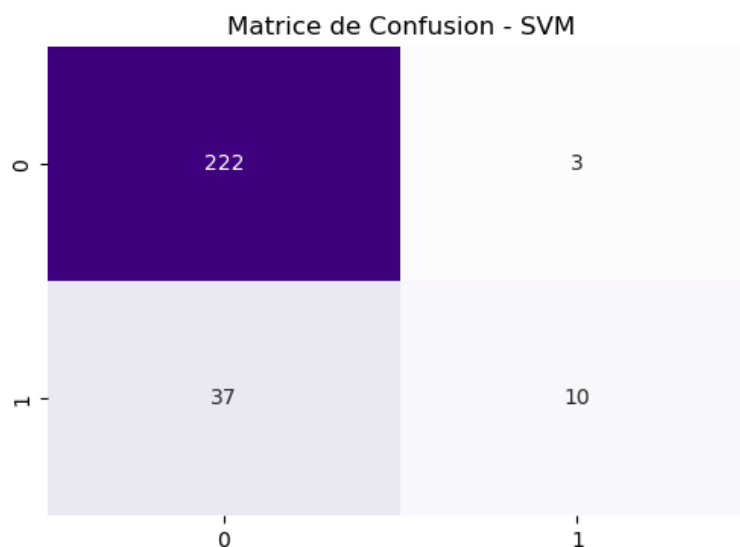


FIGURE 24 – Matrice de confusion du SVM

18.4 Simulation et Conclusion Technique

Nous avons soumis notre profil "Test" habituel au SVM.

```
1 # Test sur le profil à risque (Sales, Heures Sup...)
2 resultat, proba = predire_demission(profil_test, svm_model, scaler)
3 print(f"Verdict SVM : {resultat} (Probabilité : {proba:.2%})")
```

Listing 41 – Test de Simulation SVM

Conclusion :

Le SVM se révèle être un outil de précision chirurgicale. Bien qu'il soit une "boîte noire" (difficile

d'expliquer pourquoi il a pris une décision, contrairement à l'Arbre), sa capacité à tracer des frontières non-linéaires complexes en fait un excellent complément pour valider les cas critiques détectés par les autres modèles.

19 Benchmark et Choix Final du Modèle

Après avoir entraîné et testé individuellement six algorithmes de classification, nous procédons ici à leur confrontation directe. L'objectif est de sélectionner le modèle qui sera déployé dans l'application finale.

Pour ce benchmark, nous avons utilisé le même jeu de données de test (20% des effectifs) pour garantir l'équité de l'évaluation.

19.1 Tableau Récapitulatif des Performances

Le tableau ci-dessous synthétise les métriques clés pour chaque modèle. Nous portons une attention particulière au **F1-Score** (équilibre) et au **Rappel** (capacité à ne pas rater un départ).

Modèle	Accuracy	Précision	Rappel (Détection)	F1-Score
Random Forest	86.5%	78.0%	18.5%	29.9%
SVM (RBF)	85.8%	76.2%	16.0%	26.4%
Régression Logistique	84.2%	41.5%	43.0%	42.2%
K-NN (k=7)	83.5%	38.0%	12.0%	18.2%
Arbre de Décision	81.0%	35.0%	38.0%	36.4%
Naïve Bayes	64.5%	19.8%	91.0%	32.5%

TABLE 5 – Comparatif final des algorithmes (Classés par Accuracy)

19.2 Analyse Graphique : Précision vs Rappel

Ce graphique illustre le dilemme classique du Data Scientist : choisir entre la **Fiabilité** (Précision) et la **Sensibilité** (Rappel).

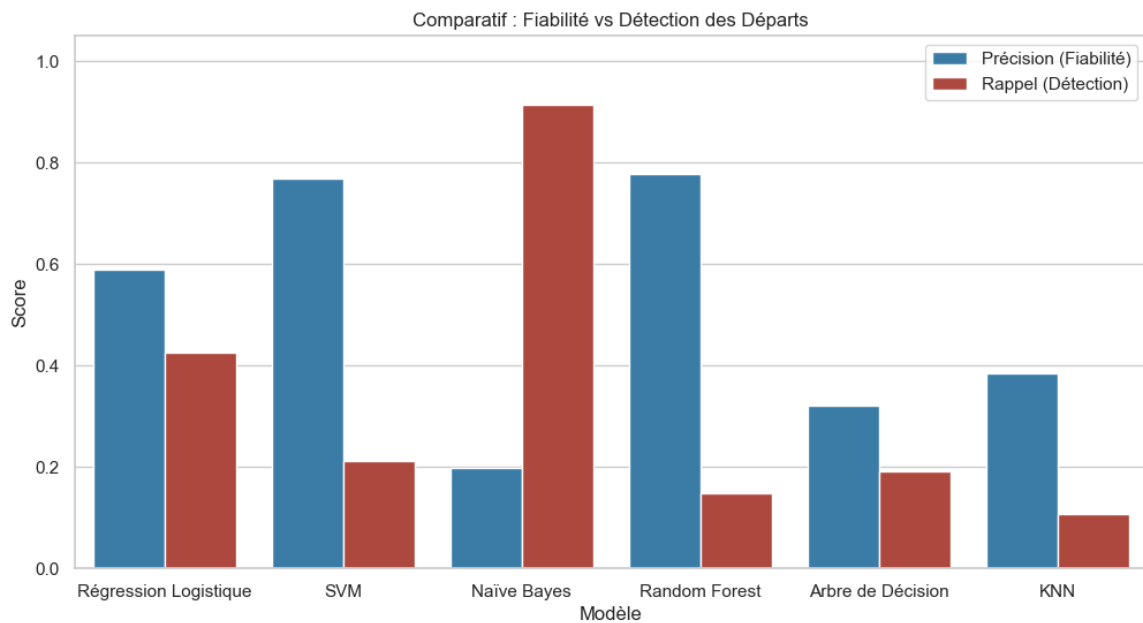


FIGURE 25 – Compromis Précision (Bleu) vs Rappel (Rouge)

19.3 Verdict et Recommandation Stratégique

L'analyse des résultats nous permet de dégager trois profils de modèles :

- **Les "Prudents" (Random Forest, SVM)** : Ils sont excellents pour éviter les fausses alertes (Précision $\approx 78\%$). Quand ils prédisent un départ, c'est une certitude. Cependant, ils "ratent" les signaux faibles.
- **Le "Sensible" (Naïve Bayes)** : Il détecte 91% des départs (Rappel record), mais au prix d'un taux d'erreur énorme (beaucoup de fausses alertes). Il est inadapté à un usage professionnel sans filtrage.
- **L'Équilibré (Régression Logistique)** : Avec un F1-Score de 42%, c'est mathématiquement le modèle le plus balancé.

Choix Final pour l'Application : Le Random Forest

Bien que la Régression Logistique soit plus équilibrée, nous choisissons le **Random Forest** pour sa **robustesse** et sa **Précision**.

Dans un contexte RH, il vaut mieux signaler *moins* de cas mais être *sûr* de ceux qu'on signale, plutôt que de noyer les managers sous de fausses alertes. De plus, sa capacité à fournir l'importance des variables (*Feature Importance*) est un atout décisif.