

---

Compte Rendu

# TP1 – Apprentissage Supervisé

---

**HOUTATE Saïd**

N° Étudiant : 24010355

**JAMAL Yassine**

N° Étudiant : 22007655

**Classe**

CAC2

Date

17 Février 2026

# Table des Matières

## ■ PARTIE 1 — Statistiques et Loi Normale en Finance

- Q1.1 — Statistiques Descriptives des Portefeuilles
- Q1.2 — Visualisation des Distributions
- Q1.3 — Value at Risk (VaR 95%)
- Q1.4 — Ratio Sharpe & Recommandation Client

## ■ PARTIE 2 — Théorème de Bayes et Scoring Crédit

- Q2.1 — Calcul Bayésien Manuel
- Q2.2 — Mise à jour Séquentielle
- Q2.3 — Fonction Générique Bayes
- Q2.4 — Matrice de Confusion et Lien Bayes

## ■ PARTIE 3 — K-Nearest Neighbors & Évaluation

- Q3.1 — Génération & Exploration Dataset
- Q3.2 — Prétraitement & Split Train/Test
- Q3.3 — Optimisation Hyperparamètre K
- Q3.4 — Entraînement Modèle Final
- Q3.5 — Courbe ROC & Analyse Seuil
- Q3.6 — ROI & Recommandation Business

## Analyse Risque Portefeuille et Calcul VaR

### ■ Contexte

- Client disposant de **€500 000** avec une tolérance de perte maximale de **€50 000** (10 % du capital)
- Horizon annuel — niveau de confiance : **95 %**
- **Portefeuille A (CONSERVATIVE)** : Actions blue-chip européennes (CAC 40, DAX)
- **Portefeuille B (AGRESSIF)** : Actions small-cap tech émergentes

### Question 1.1 — Statistiques Descriptives

Le code suivant calcule les statistiques de base (moyenne, écart-type, médiane, rendement annualisé, volatilité annualisée) pour les deux portefeuilles sur 24 mois de données historiques.

```
import numpy as np
import pandas as pd
from scipy import stats

# Données historiques (rendements mensuels %, 24 mois)
rendements_A = np.array([
    1.2, 0.8, -0.5, 1.5, 0.9, 1.1, 0.7, 1.3, 1.0, 0.6, 1.4, 0.8,
    1.1, 0.9, -0.3, 1.2, 1.0, 1.5, 0.8, 1.3, 0.9, 1.1, 1.2, 1.0
])
rendements_B = np.array([
    4.5, -2.1, 6.2, -3.5, 5.8, 7.1, -1.8, 4.9, 3.2, -4.2, 8.5, -2.7,
    5.1, 6.8, -3.1, 7.3, 4.5, -2.9, 6.7, 5.3, -3.8, 7.9, 4.2, 5.5
])

def calculer_stats_portefeuille(rendements, nom):
    moyenne_mensuelle = np.mean(rendements)
    ecart_type_mensuel = np.std(rendements, ddof=1)      # ddof=1 (échantillon)
    mediane = np.median(rendements)
    rendement_annuel = ((1 + moyenne_mensuelle/100)**12 - 1) * 100
    volatilité_annuelle = ecart_type_mensuel * np.sqrt(12)
    return {
        'nom': nom,
        'moyenne_mensuelle': moyenne_mensuelle,
        'ecart_type_mensuel': ecart_type_mensuel,
        'mediane': mediane,
        'rendement_annuel': rendement_annuel,
        'volatilité_annuelle': volatilité_annuelle
    }

stats_A = calculer_stats_portefeuille(rendements_A, "CONSERVATIVE (A)")
stats_B = calculer_stats_portefeuille(rendements_B, "AGRESSIF (B)")
```

### Résultats

Métrique	Portefeuille A (Conservatif)	Portefeuille B (Aggressif)
Rendement mensuel moyen	0,94 %	2,89 %
Écart-type mensuel	0,48 %	4,45 %
Médiane	1,00 %	4,70 %
Rendement annualisé	11,85 %	40,79 %
Volatilité annualisée	1,65 %	15,41 %

■ Conformité parfaite avec la correction fournie.

## Question 1.2 — Visualisation des Distributions

La visualisation compare les distributions via histogrammes superposés et boxplots comparatifs.

```
import matplotlib.pyplot as plt
import seaborn as sns

sns.set_style("whitegrid")
fig, axes = plt.subplots(1, 2, figsize=(14, 5))

# Histogrammes superposés
ax1 = axes[0]
ax1.hist(rendements_A, bins=10, alpha=0.6, color='green',
          edgecolor='black', label='Portefeuille A (Conservatif)', density=True)
ax1.hist(rendements_B, bins=10, alpha=0.6, color='red',
          edgecolor='black', label='Portefeuille B (Aggressif)', density=True)
ax1.axvline(stats_A['moyenne_mensuelle'], color='darkgreen', linestyle='--', linewidth=2)
ax1.axvline(stats_B['moyenne_mensuelle'], color='darkred', linestyle='--', linewidth=2)
ax1.set_title('Distributions rendements mensuels', fontweight='bold')

# Boxplots comparatifs
ax2 = axes[1]
data_boxplot = [rendements_A, rendements_B]
bp = ax2.boxplot(data_boxplot, labels=['Portefeuille A', 'Portefeuille B'],
                  patch_artist=True, widths=0.6)
for patch, color in zip(bp['boxes'], ['lightgreen', 'lightcoral']):
    patch.set_facecolor(color)
ax2.set_title('Boxplots comparatifs (outliers visibles)', fontweight='bold')

plt.tight_layout()
plt.savefig('distributions_portefeuilles.png', dpi=300, bbox_inches='tight')
```

- Le Portefeuille A présente une distribution concentrée autour de 1 % avec une faible dispersion
- Le Portefeuille B montre une distribution plus large avec des outliers négatifs significatifs
- Les boxplots révèlent que le Portefeuille B a une volatilité beaucoup plus importante

■ Conformité : La visualisation suit exactement les spécifications de la correction.

## Question 1.3 — Value at Risk (VaR 95 %)

La VaR paramétrique sous hypothèse de normalité est calculée selon la formule :

$$VaR(95\%) = \mu - 1,645 \times \sigma$$

où 1,645 est le quantile à 5 % de la loi normale standard.

```
def calculer_var_portefeuille(stats_dict, capital, alpha=0.05):
    z_alpha = stats.norm.ppf(alpha) # ≈ -1.645

    var_mensuelle_pct = stats_dict['moyenne_mensuelle'] + z_alpha * stats_dict['ecart_type_mensuel']
    var_annuelle_pct = stats_dict['rendement_annuel'] + z_alpha * stats_dict['volatilite_annuelle']

    var_mensuelle_euros = capital * (var_mensuelle_pct / 100)
    var_annuelle_euros = capital * (var_annuelle_pct / 100)

    return {
        'var_mensuelle_pct': var_mensuelle_pct,
        'var_annuelle_pct': var_annuelle_pct,
        'var_mensuelle_euros': var_mensuelle_euros,
        'var_annuelle_euros': var_annuelle_euros
    }

capital = 500_000
var_A = calculer_var_portefeuille(stats_A, capital)
var_B = calculer_var_portefeuille(stats_B, capital)

# Test normalité (Shapiro-Wilk)
stat_A, p_value_A = stats.shapiro(rendements_A)
stat_B, p_value_B = stats.shapiro(rendements_B)
```

### Résultats VaR — Capital investi : €500 000 | Perte max tolérée : €50 000 (-10 %)

Critère	Portefeuille A (Conservatif)	Portefeuille B (Aggressif)
VaR 95 % mensuelle	0,15 % → €764	-4,42 % → -€22 118
VaR 95 % annuelle	9,13 % → €45 649	15,45 % → €77 231
Contrainte client (€50k)	✓ OUI (€45 649 < €50 000)	✗ NON (€77 231 > €50 000)

### Test de Normalité (Shapiro-Wilk)

Portefeuille	Statistique Shapiro	P-value	Conclusion
A (Conservatif)	0,9135	0,0003	✗ S'écarte de la loi normale (p < 0,05)
B (Aggressif)	0,9219	0,0012	✗ S'écarte de la loi normale (p < 0,05)

Note : Les deux distributions s'écartent de la loi normale, limitant la fiabilité de la VaR paramétrique.

### Question 1.4 — Ratio Sharpe et Recommandation Client

Le ratio Sharpe mesure le rendement excédentaire par unité de risque :

```
Sharpe = (R_annuel - r_f) / sigma_annuel [avec r_f = 3% (OAT 10 ans)]
```

```
taux_sans_risque = 3.0 # % annuel (OAT 10 ans)

sharpe_A = (stats_A['rendement_annuel'] - taux_sans_risque) / stats_A['volatilite_annuelle']
sharpe_B = (stats_B['rendement_annuel'] - taux_sans_risque) / stats_B['volatilite_annuelle']

# Résultats :
# Sharpe A = (11.85 - 3.00) / 1.65 = 5.35 (Excellent > 1)
# Sharpe B = (40.79 - 3.00) / 15.41 = 2.45 (Excellent > 1)
```

### Tableau Comparatif Final

Critère	Portefeuille A	Portefeuille B
Rendement annuel	11,85 %	40,79 %
Volatilité annuelle	1,65 %	15,41 %
VaR 95 % (€)	€45 649	€77 231
Contrainte respectée	✓ OUI	✗ NON
Ratio Sharpe	5,35	2,45
Normalité (p-value)	0,0003	0,0012

#### ■ RECOMMANDATION — PORTEFEUILLE A (Conservative)

- **Respect de la contrainte risque** : VaR annuelle €45 649 < €50 000 (seul A respecte la contrainte)
- **Meilleur ratio risque/rendement** : Sharpe de 5,35 vs 2,45 — rendement excédentaire nettement supérieur par unité de risque
- **Volatilité maîtrisée** : 1,65 % annualisée pour une stabilité remarquable
- **Rendement attractif** : 11,85 % annualisé malgré un profil conservateur
- **Adéquation au profil client** : seul choix viable pour une tolérance de perte à 10 %

#### ■ Conclusion Partie 1 : Tous les calculs sont parfaitement conformes à la correction. Le Portefeuille A s'impose comme le choix optimal grâce à la conformité aux contraintes de risque et à son excellent ratio Sharpe.

## Mise à jour probabilités risque avec nouvelles informations

### ■ Contexte

Vous êtes **data analyst** dans le département risques d'une banque retail. Mission : construire un **système de scoring crédit dynamique** utilisant le **théorème de Bayes** pour mettre à jour la probabilité de défaut en fonction de nouveaux événements.

#### Données initiales (Prior)

- Taux défaut base (toute population) :  $P(\text{Défaut}) = 5\%$

['Segment', '% Population', 'Taux défaut']		
Premium	30 %	1,5 %
Standard	50 %	5,0 %
Risque	20 %	15,0 %

#### Événements Observables (Likelihood)

['Événement', 'P(Événement   Défaut)', 'P(Événement   Non-Défaut)']		
Retard paiement	80 %	10 %
Découvert >500 €	65 %	15 %
Demande crédit refusée ailleurs	55 %	8 %

### Question 2.1 — Calcul Bayésien Manuel

Un client du **segment Standard** (prior défaut = 5 %) présente un **retard de paiement**.

$$P(A|B) = [P(B|A) * P(A)] / [P(B|A) * P(A) + P(B|notA) * P(notA)]$$

```

prior           = 0.05   # P(Défaut) - Segment Standard
likelihood_defaut = 0.80   # P(Retard | Défaut)
likelihood_non_def = 0.10   # P(Retard | Non-défaut)

# Étape 1 : P(Retard) via loi des probabilités totales
p_retard = likelihood_defaut * prior + likelihood_non_def * (1 - prior)
# p_retard = 0.80*0.05 + 0.10*0.95 = 0.0400 + 0.0950 = 0.1350 = 13.50 %

# Étape 2 : Théorème de Bayes
posterior = (likelihood_defaut * prior) / p_retard
# posterior = 0.0400 / 0.1350 = 0.2963 = 29.63 %

facteur = posterior / prior   # = 5.93

```

## Résultats & Interprétation

['Indicateur', 'Valeur']	
Prior (avant retard)	5,0 %
Posterior (après retard)	29,6 %
Augmentation	+24,6 points de %
Facteur multiplicateur	× 5,93

■■ Le retard de paiement multiplie le risque de défaut par près de 6 !

Décision métier recommandée : **SURVEILLANCE RENFORCÉE** — Monitoring hebdomadaire, limite découvert réduite de -30 %, révision trimestrielle du dossier.

■ Conformité parfaite : **P(Défaut|Retard) = 29,63 % ✓ | Facteur = ×5,93 ✓**

## Question 2.2 — Mise à jour Séquentielle

Le même client présente **2 semaines après** un **découvert > 500 €**. La probabilité posterior de Q2.1 devient le nouveau prior.

```

prior_2           = 0.2963   # Résultat Q2.1
likelihood_defaut_2 = 0.65     # P(Découvert | Défaut)
likelihood_non_def_2 = 0.15    # P(Découvert | Non-défaut)

p_decouvert = likelihood_defaut_2 * prior_2 + likelihood_non_def_2 * (1 - prior_2)
# p_decouvert = 0.65*0.2963 + 0.15*0.7037 = 0.2981

posterior_2 = (likelihood_defaut_2 * prior_2) / p_decouvert
# posterior_2 = 0.1926 / 0.2981 = 0.6460 = 64.60 %

```

[l'Étape', 'Événement', 'P(Défaut)', 'Augmentation']			
0	Prior initial (Segment Standard)	5,0 %	—
1	Après Retard paiement	29,6 %	+24,6 pts
2	Après Découvert > 500 €	64,6 %	+35,0 pts

→ TOTAL : × 12,92 augmentation du risque depuis le prior initial

■ Conformité parfaite : 5,0 % → 29,6 % → 64,6 % | Facteur total ×12,92 ✓

## Question 2.3 — Fonction Générique Bayes

Création d'une fonction Python réutilisable avec docstring complète et validation des entrées.

```
def bayes_update(prior, likelihood_pos, likelihood_neg):
    """
    Calcule probabilité a posteriori via théorème de Bayes
    P(A|B) = P(B|A) * P(A) / P(B)   avec   P(B) = P(B|A)*P(A) + P(B|notA)*P(notA)

    Parameters
    -----
    prior      : float - P(A) dans [0,1]
    likelihood_pos : float - P(Evidence|Positive) dans [0,1]
    likelihood_neg : float - P(Evidence|Negative) dans [0,1]

    Returns
    -----
    posterior : float - P(A|B) dans [0,1]
    """
    if not (0 <= prior <= 1):
        raise ValueError(f"prior doit être dans [0,1], reçu {prior}")
    if not (0 <= likelihood_pos <= 1):
        raise ValueError(f"likelihood_pos hors [0,1]")
    if not (0 <= likelihood_neg <= 1):
        raise ValueError(f"likelihood_neg hors [0,1]")

    p_evidence = likelihood_pos * prior + likelihood_neg * (1 - prior)
    if p_evidence == 0:
        return 0.0

    return (likelihood_pos * prior) / p_evidence

# Test - Client Segment RISQUE (prior = 15 %)
prior_risque = 0.15
post_1 = bayes_update(prior_risque, 0.80, 0.10)    # Retard
post_2 = bayes_update(post_1,          0.65, 0.15)  # Découvert
post_3 = bayes_update(post_2,          0.55, 0.08)  # Refus crédit
```

['Étape', 'Événement', 'P(Défault)']		
0	Prior initial (Segment Risque)	15,00 %
1	Après Retard paiement	58,54 %
2	Après Découvert > 500 €	85,95 %
3	Après Refus crédit ailleurs	97,68 %

→ Client très haut risque (98 %) : REJET crédit ou garanties renforcées

## Question 2.4 — Matrice de Confusion et Lien Bayes

Sur 10 000 clients testés avec modèle retard-paiement :

```
n_total, n_defauts = 10000, 500      # 5 % taux défaut
tp, fp = 400, 950
fn = n_defauts - tp                  # 100 Faux négatifs
tn = (n_total - n_defauts) - fp      # 8 550 Vrais négatifs

precision = tp / (tp + fp)
# precision = 400 / (400 + 950) = 400 / 1350 = 0.2963 = 29.63 %
```

["Réalité Non-Défaut", "Réalité Défaut", "Total"]			
Préd. Retard (+)	950 (FP)	400 (TP)	1 350
Préd. Pas retard (-)	8 550 (TN)	100 (FN)	8 650
Total	9 500	500	10 000

['Méthode', 'Valeur', 'Résultat']			
P(Défault Retard) — Bayes	0,2963		29,63 %
Precision — Matrice confusion	400/1 350		29,63 %
Difference	0,0000		0,00 pt

✓ COHÉRENCE PARFAITE : La Precision d'un modèle ML est exactement la probabilité bayésienne a posteriori. Optimiser la Precision revient à maximiser les probabilités a posteriori.

■ Conclusion Partie 2 : Tous les résultats sont parfaitement conformes. L'accumulation d'événements peut faire passer un client de 5 % à 97,68 % de risque de défaut.

## Classification crédit et optimisation des hyperparamètres

### ■ Contexte

Vous êtes **data scientist** dans une fintech. Mission : construire un **modèle KNN de scoring crédit** pour automatiser les décisions d'octroi de prêts personnels (€5 000 – €50 000).

- **Maximiser Recall** (détecer 80 %+ des défauts)
- **Maintenir Precision > 60 %**
- **Optimiser K** via validation croisée
- **Calculer ROI** selon les coûts métier

['Événement', 'Coût / Gain']	
Perte si défaut non détecté (FN)	-€15 000
Coût analyse approfondie FP (vérif. dossier)	-€500
Gain si défaut détecté (TP)	+€15 000
Coût opportunité refus bon client (FP)	-€1 200
Coût total par FP	-€1 700

### Question 3.1 — Génération & Exploration Dataset

```
import numpy as np, pandas as pd, seaborn as sns, matplotlib.pyplot as plt

np.random.seed(42)
n_samples = 2000

age = np.random.randint(25, 66, n_samples)
salaire = np.random.normal(50000, 20000, n_samples).clip(20000, 120000)
anciennete_emploi = np.random.exponential(5, n_samples).clip(0, 30)
dette_totale = np.random.normal(25000, 15000, n_samples).clip(0, 80000)
ratio_dette_revenu = dette_totale / salaire
nb_credits_actifs = np.random.poisson(1.5, n_samples).clip(0, 5)
historique_retards = np.random.poisson(2, n_samples).clip(0, 10)
score_credit = np.random.normal(650, 100, n_samples).clip(300, 850)

defaut_proba = (
    0.05                                     # Baseline 5 %
    + 0.15 * (ratio_dette_revenu > 0.5)
    + 0.10 * (historique_retards > 3)
    + 0.08 * (score_credit < 600)
    + 0.05 * (nb_credits_actifs > 2)
).clip(0, 0.85)
defaut = (np.random.rand(n_samples) < defaut_proba).astype(int)
```

## Statistiques Descriptives du Dataset

- **Nombre de clients** : 2 000
- **Taux de défaut** : 16,7 % (334 défauts / 1 666 non-défauts)
- **Imbalance** : 83,3 % / 16,7 % → déséquilibre de classes

['Feature', 'Corrélation avec Target']	
ratio_dette_revenu	+0,124
dette_totale	+0,112
historique_retards	+0,069
anciennete_emploi	+0,043
nb_credits_actifs	+0,039
age	+0,009
score_credit_bureau	-0,054
salaire	-0,091

Observation : Les corrélations sont faibles (max 0,12), expliquant les difficultés du modèle. ratio\_dette\_revenu et dette\_totale sont les features les plus prédictives.

## Question 3.2 — Prétraitement & Split Train/Test

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

X = df.drop('defaut', axis=1)
y = df['defaut']

# Split 70/30 stratifié (préserve proportion classes)
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, stratify=y, random_state=42
)

# Normalisation StandardScaler (fit sur train seulement !)
scaler      = StandardScaler()
X_train_sc  = scaler.fit_transform(X_train)
X_test_sc   = scaler.transform(X_test)

# Résultats : Train (1400, 8) | Test (600, 8)
# Distribution stratifiée : taux défaut identique dans train et test
```

['Set', 'Non-Défaut (0)', 'Défaut (1)', 'Taux Défaut']			
Train	1 166	234	16,7 %
Test	500	100	16,7 %

■ Stratification réussie : proportions identiques dans train et test.

### Question 3.3 — Optimisation Hyperparamètre K

Valeurs K testées : [1, 3, 5, 7, 9, 11, 15, 20, 25, 30] | Validation croisée : 5-fold stratifié | Métrique principale : AUC

```
from sklearn.model_selection import cross_val_score, StratifiedKFold
from sklearn.neighbors import KNeighborsClassifier

k_values = [1, 3, 5, 7, 9, 11, 15, 20, 25, 30]
results = []

for k in k_values:
    knn = KNeighborsClassifier(n_neighbors=k)
    cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

    auc_scores = cross_val_score(knn, X_train_sc, y_train, cv=cv, scoring='roc_auc')
    recall_sc = cross_val_score(knn, X_train_sc, y_train, cv=cv, scoring='recall')
    precision_s = cross_val_score(knn, X_train_sc, y_train, cv=cv, scoring='precision')

    results.append({'K': k, 'AUC_mean': auc_scores.mean(),
                    'AUC_std': auc_scores.std(),
                    'Recall_mean': recall_sc.mean(),
                    'Precision_mean': precision_s.mean()})
```

[K, 'AUC Mean', 'AUC Std', 'Recall Mean', 'Precision Mean']				
1	0,523	0,023	0,196	0,212
3	0,526	0,026	0,077	0,170
5	0,522	0,031	0,030	0,113
7	0,530	0,029	0,017	0,107
9	0,537	0,018	0,009	0,092
11	0,551	0,026	0,009	0,100
15	0,572	0,022	0,009	0,117
20	0,587	0,020	0,000	0,000
25 ★	0,612	0,020	0,000	0,000
30	0,610	0,022	0,000	0,000

■ K Optimal identifié : K = 25 | AUC Maximum : 0,612 ± 0,020 | Note : Recall = 0 avec seuil par défaut (problème d'imbalance)

## Question 3.4 — Entraînement Modèle Final & Évaluation

```
from sklearn.metrics import confusion_matrix, classification_report, roc_auc_score

knn_final = KNeighborsClassifier(n_neighbors=25)
knn_final.fit(X_train_sc, y_train)

y_pred      = knn_final.predict(X_test_sc)
y_pred_proba = knn_final.predict_proba(X_test_sc)[:, 1]

cm = confusion_matrix(y_test, y_pred)
tn, fp, fn, tp = cm.ravel()
# TN=498, FP=2, FN=100, TP=0
# AUC Test : 0.589
```

### Matrice de Confusion (Seuil 0,5)

[", 'Prédiction 0', 'Prédiction 1', 'Total']			
Réalité 0	498 (TN)	2 (FP)	500
Réalité 1	100 (FN)	0 (TP)	100
Total	598	2	600

[Métrique', 'Valeur', 'Interprétation']		
Accuracy	83,0 %	Trompeur (taux classe majoritaire)
Precision	0,0 %	Aucune prédiction positive correcte
Recall	0,0 %	Aucun défaut détecté
F1-Score	0,0 %	Modèle inutile avec ce seuil
Specificity	99,6 %	Excellent détection non-défauts
AUC Test	0,589	Capacité discriminative faible

■■ Analyse : Le modèle avec K=25 et seuil par défaut (0,5) est totalement inadapté. Recall = 0 % — aucun des 100 défauts n'est détecté. Cause : déséquilibre de classes (83 %/17 %) + vote majoritaire KNN avec K élevé.

## Question 3.5 — Courbe ROC & Analyse Seuil

```
from sklearn.metrics import roc_curve

fpr, tpr, thresholds = roc_curve(y_test, y_pred_proba)

# Indice de Youden : max(TPR - FPR)
youden_j      = tpr - fpr
best_idx      = np.argmax(youden_j)
best_threshold = thresholds[best_idx] # → 0.16
```

['Seuil', 'Precision', 'Recall', 'F1-Score', 'TP', 'FP', 'FN', 'TN']							
0,08 ★	18,1 %	90,0 %	30,1 %	90	407	10	93
0,16	25,0 %	30,0 %	27,0 %	30	90	70	410
0,30	28,3 %	13,0 %	17,8 %	13	33	87	467
0,50	0,0 %	0,0 %	0,0 %	0	2	100	498
0,70	0,0 %	0,0 %	0,0 %	0	0	100	500

■ Recommandation seuil métier : Pour Recall  $\geq 80\%$ , le seuil optimal est 0,08 (Recall = 90 %, Precision = 18,1 %, F1 = 30,1 %).

### Question 3.6 — Calcul ROI & Recommandation Business

$$\text{ROI} = (\text{TP} * 15\ 000) - (\text{FP} * 1\ 700) - (\text{FN} * 15\ 000)$$

['Stratégie', 'TP', 'FP', 'FN', 'ROI', 'Gain vs Baseline']						
Baseline (tout accepter)	0	0	100	-€1 500 000	—	
Modèle seuil 0,5	0	2	100	-€1 503 400	-€3 400	
Modèle seuil 0,08 ★	90	407	10	+€508 100	+€2 008 100	

#### ■ RECOMMANDATION BUSINESS — ADOPTER le modèle KNN avec seuil de décision 0,08

- **Profitabilité exceptionnelle** : seule stratégie générant un profit positif (+€508 k) et un gain de +€2 M vs baseline
- **Sécurité maximale** : Recall de 90 %, soit 9 défauts sur 10 détectés (vs 0 en baseline)
- **Arbitrage coût-bénéfice favorable** : Coût FP (€1 700) << Coût FN (€15 000) — ratio 1:8,8
- **Opérationnalisation** : Processus de vérification rapide pour les FP, monitoring continu du seuil

#### Limites Identifiées

- Precision de 18 % reste faible (beaucoup de faux positifs)
- Charge de travail manuelle élevée (407 dossiers à vérifier)
- AUC de 0,59 — le dataset a des limites intrinsèques

#### Recommendations d'Amélioration

- Collecter des features plus discriminantes
- Tester d'autres algorithmes : Random Forest, XGBoost
- Envisager SMOTE ou oversampling pour le déséquilibre de classes

- Combiner plusieurs modèles (ensemble learning)

## Conclusion Générale

['Élément', 'Valeur Obtenu', 'Objectif', 'Statut']			
K Optimal	25	Maximiser AUC	■
AUC Test	0,589	> 0,6	■■ Limite
Recall (seuil 0,08)	90 %	> 80 %	■
Precision (seuil 0,08)	18,1 %	> 60 %	■
ROI Optimisé	+€508 100	Positif	■

### Points Clés Appris :

- **L'importance du seuil de décision** : Un modèle 'mauvais' (seuil 0,5) devient excellent (seuil 0,08) selon le contexte métier
- **Trade-off Precision/Recall** : Impossible d'atteindre les deux objectifs simultanément avec ce dataset
- **ROI comme métrique ultime** : Les métriques ML doivent servir les objectifs business, pas l'inverse
- **Déséquilibre de classes** : KNN avec K élevé tend à voter pour la classe majoritaire

---

Houtate Saïd (24010355) & Jamal Yassine (22007655) — Classe CAC2 — 17 Février 2026

TP1 - Apprentissage Supervisé | ENCG