

Partie 1 : Questions à réponse courte

1.1 Définition du problème

Problème hypothétique

Prédire le risque d'abandon scolaire chez les élèves du secondaire dans les zones rurales.

Objectifs

- Identifier les élèves à risque d'abandon avant la fin de l'année scolaire.
- Proposer des interventions ciblées (soutien psychologique, tutorat, aide financière).
- Réduire le taux global d'abandon scolaire dans les établissements concernés.

Parties prenantes

- Direction des établissements scolaires
- Ministère de l'Éducation

KPI

Taux de précision du modèle dans la détection des élèves à risque.

1.2 Collecte et prétraitement des données

Sources de données

- Dossiers scolaires historiques
- Enquêtes socio-économiques locales

Biais potentiel

Sous-représentation des élèves issus de milieux défavorisés ou non scolarisés régulièrement.

Étapes de prétraitement

- Imputation des données manquantes
- Encodage des variables catégorielles
- Normalisation des variables numériques

1.3 Développement du modèle

Modèle choisi

Random Forest — robuste, interprétable, adapté aux données tabulaires.

Division des données

- 70% entraînement
- 15% validation
- 15% test

Hyperparamètres à régler

- Nombre d'arbres : influence la stabilité du modèle
- Profondeur maximale : contrôle le surapprentissage

1.4 Évaluation et déploiement

Mesures d'évaluation

- Score F1 : équilibre entre précision et rappel
- AUC-ROC : performance globale du modèle

Dérive conceptuelle

Changement dans la relation entre les variables et la cible.

Surveillance via des audits réguliers et des alertes de performance.

Défi technique

Scalabilité — adapter le modèle à un grand volume de données en temps réel.

Partie 2 : Étude de cas

Scénario

Un hôpital souhaite disposer d'un système d'IA pour prédire le risque de réadmission des patients dans les 30 jours suivant leur sortie.

2.1 Portée du problème

Problème

Anticiper les réadmissions pour améliorer les soins et réduire les coûts.

Objectifs

- Identifier les patients à risque
- Optimiser les ressources hospitalières

Parties prenantes

- Médecins
- Personnel hospitalier
- Patients

2.2 Stratégie de données

Sources

- Dossiers de sortie (DSE)
- Données démographiques

Préoccupations éthiques

- Confidentialité des données (RGPD, HIPAA)
- Risque de discrimination algorithmique

Pipeline de prétraitement

- Nettoyage des doublons
- Création de variables dérivées (durée d'hospitalisation, nombre de visites)
- Encodage des diagnostics

2.3 Développement de modèles

Modèle choisi

Gradient Boosting — performant sur données tabulaires, gère les interactions complexes.

Matrice de confusion (hypothétique)

$TP = 80, FP = 20, FN = 30, TN = 70$

$Précision = 80 / (80 + 20) = 0.80$

$Rappel = 80 / (80 + 30) = 0.73$

2.4 Déploiement

Étapes

- Intégration via API dans le système hospitalier
- Formation du personnel
- Tests en environnement réel

Conformité

- Chiffrement des données
- Journalisation des accès
- Audit régulier

2.5 Optimisation

Méthode

Validation croisée pour détecter et limiter le surapprentissage.

Partie 3 : Pensée critique

3.1 Éthique et biais

Impact

Un modèle biaisé pourrait sous-estimer le risque chez certains groupes (ex. : minorités), entraînant des soins inadaptés.

Stratégie d'atténuation

- Échantillonnage équilibré
- Analyse de l'équité des prédictions
- Révision régulière des données d'entraînement

3.2 Compromis

Interprétabilité vs Précision

- Modèles simples (arbres de décision) : faciles à expliquer mais moins précis
- Modèles complexes (réseaux de neurones) : plus performants mais opaques

Ressources limitées

- Choisir un modèle léger (ex. : régression logistique)
- Utiliser des services cloud pour le calcul

Partie 4 : Réflexion et diagramme

4.1 Réflexion

La partie la plus difficile a été l'analyse éthique, car elle nécessite une compréhension multidisciplinaire. Avec plus de temps, nous aurions approfondi les tests de robustesse et impliqué des experts médicaux.

4.2 Diagramme du flux de travail IA



