

Rapport de Création du Data Pool

Prédiction de la Performance des Étudiants Marocains
Lycées Qualifiants - Tous les Niveaux (Tronc Commun, 1BAC, 2BAC)

Projet : Prédiction de la Performance des Étudiants
Date : 5 Février 2026

Résumé

Ce rapport documente le processus de création d'un ensemble de données synthétiques complet pour la prédiction de la performance des étudiants marocains du cycle secondaire qualifiant. Le dataset généré contient **10 000 enregistrements** d'étudiants avec **286 variables (features)** couvrant tous les aspects pertinents pour l'analyse et la prédiction de la réussite scolaire.

Table des matières

1	Introduction	3
1.1	Contexte du Projet	3
1.2	Objectif du Data Pool	3
2	Description du Dataset	3
2.1	Caractéristiques Générales	3
2.2	Catégories de Variables	3
2.2.1	1. Informations Personnelles et Démographiques (20 variables)	3
2.2.2	2. Informations Scolaires (15 variables)	4
2.2.3	3. Informations Familiales (25 variables)	4
2.2.4	4. Conditions Socio-économiques (20 variables)	4
2.2.5	5. Informations de Santé (15 variables)	5
2.2.6	6. Performance Académique (50+ variables)	5
2.2.7	7. Habitudes d'Étude (25 variables)	5
2.2.8	8. Activités Parascolaires (20 variables)	6
2.2.9	9. Orientation et Aspirations (15 variables)	6
2.2.10	10. Facteurs Psychologiques (20 variables)	6
2.2.11	11. Mode de Vie (15 variables)	6
2.2.12	12. Compétences et Aptitudes (20 variables)	7
2.2.13	13. Compétences Linguistiques (10 variables)	7
3	Couverture Géographique	7
4	Filières Couvertes	8
5	Méthodologie de Génération	8
5.1	Approche Utilisée	8
6	Analyse et Nettoyage des Données	8
6.1	Analyse Initiale	8
6.2	Actions de Nettoyage	8
7	Utilisation du Dataset	9
7.1	Applications Possibles	9
7.2	Variable Cible	9
7.3	Modèles Recommandés	9
8	Conclusion	9

1 Introduction

1.1 Contexte du Projet

Le projet de pr ediction de la performance des  tudiants vise  a d velopper des mod les de machine learning capables d'identifier les facteurs influen ant la r ussite scolaire des lyc ens marocains et de pr dire leur probabilit  de r ussite au baccalaur at.

1.2 Objectif du Data Pool

La cr eation de ce data pool r pond  a plusieurs objectifs :

- Fournir une base de donn es r aliste et repr sentative de la population estudiantine marocaine
- Inclure une diversit  de profils socio- conomiques et acad miques
- Couvrir toutes les r gions du Maroc
- Int grer des variables pr dictives pertinentes bas es sur la litt rature scientifique

2 Description du Dataset

2.1 Caract ristiques G n erales

Caract�ristique	Valeur
Nombre d'enregistrements	10 000 �tudiants
Nombre de variables	286 features
Format du fichier	CSV (UTF-8)
Taille du fichier	~17 MB
Langue des donn�es	Fran�ais
Couverture g�ographique	12 r�gions du Maroc

TABLE 1 – Caract ristiques g n erales du dataset

2.2 Cat gories de Variables

Le dataset est organis  en plusieurs cat gories th matiques :

2.2.1 1. Informations Personnelles et D mographiques (20 variables)

- Identifiant  tudiant (id_ tudiant)
- Pr nom, Nom, Nom complet
- Sexe (M/F)
- Date de naissance, Âge
- Code Massar (identifiant national)
- R gion, Province, Commune

- Zone (Urbain/Semi-Urbain/Rural)
- Adresse, Code postal
- T l phone, Email

2.2.2 2. Informations Scolaires (15 variables)

- Nom de l tablissement
- Type d tablissement (Lyc  Qualifiant)
- Secteur (Public/Priv )
- Acad mie r gionale
- Direction provinciale
- Niveau (Tronc Commun / 1BAC / 2BAC)
- Fili re et Sp cialit 
- Classe
- Ann e dscription et ann e scolaire

2.2.3 3. Informations Familiales (25 variables)

- Statut des parents (vivant/d c d )
- Informations sur le p re : pr nom, nom, niveau dducation, profession, secteur dactivit , revenu
- Informations sur la m re : pr nom, nom, niveau dducation, profession, secteur dactivit , revenu
- Information sur le tuteur
- Revenu familial total et source de revenu
- Nombre de fr res et surs, rang dans la fratrie
- Statut parental (Mari s/Divorc s/Veuf)

2.2.4 4. Conditions Socio- conomiques (20 variables)

- Type de logement (Villa/Appartement/Maison/Maison Traditionnelle)
- Statut de propri t  (Propri taire/Locataire)
- Nombre de pi ces
- Acc s aux services (lectricit , eau, internet)
-  quipements disponibles (chambre personnelle, bureau, ordinateur, laptop, tablette, smartphone)
- Livres ´ la maison
- Distance ´ lcole et moyen de transport
- Bourses et aides (Tayssir, bourses sociales, bourses dexcellence)

2.2.5 5. Informations de Sant  (15 variables)

- Assurance maladie et type
- Maladies chroniques
- Handicap
- Port de lunettes, probl mes auditifs, allergies
-  tat de sant  g n ral
- IMC, taille, poids

2.2.6 6. Performance Acad mique (50+ variables)

Notes pour chaque mati re (Semestre 1, Semestre 2, Moyenne annuelle) :

- Arabe
- Fran ais
- Anglais
- Math matiques
- Physique-Chimie
- Sciences de la Vie et de la Terre (SVT)
- Histoire-G ographie
-  ducation Islamique
- Philosophie
-  ducation Physique
- Informatique
-  conomie, Comptabilit , Gestion (pour les fili res  conomiques)

Indicateurs globaux :

- Moyenne g n rale (S1, S2, Annuelle)
- Rang dans la classe
- Effectif de la classe
- Absences (totales, justifi es, non justifi es)
- Retards
- Avertissements et sanctions
- Comportement

2.2.7 7. Habitudes d' tude (25 variables)

- Heures d' tude par jour et weekend
- Cours particuliers et mati res
- Co t du soutien scolaire
- Lieu et moment pr f r  d' tude
-  tude en groupe
- Utilisation des ressources en ligne
- Taux de remise des devoirs
- Participation en classe
- Prise de notes
- Utilisation de la biblioth que

2.2.8 8. Activites Parascolaires (20 variables)

- Activites sportives et type de sport
- Activites artistiques
- Clubs et associations
- Benevolat
- Participation aux evenements scolaires
- Membre du conseil des el`es

2.2.9 9. Orientation et Aspirations (15 variables)

- Aspiration de carriere
- Universite souhaitee
- Domaine d’etude souhaite
- Projet d’etudes a l’etranger
- Pays cible
- Connaissance de l’orientation
- Seances d’orientation suivies

2.2.10 10. Facteurs Psychologiques (20 variables)

- Niveau de motivation
- Confiance en soi
- Niveau de stress et d’anxiete
- Anxiete aux examens
- Satisfaction scolaire
- Relations avec les pairs
- Soutien familial
- Implication parentale
- Attentes et pression parentales

2.2.11 11. Mode de Vie (15 variables)

- Petit dejeuner et repas par jour
- Heures de sommeil (semaine/weekend)
- Heure de coucher et de reveil
- Activite physique
- Temps d’cran
- Reaux sociaux
- Jeux video
- Lecture
- Consommation de cafeine
- Travail a temps partiel

2.2.12 12. Comp tences et Aptitudes (20 variables)

- Efficacit  d'auto-apprentissage
- Capacit    fixer des objectifs
- Gestion du temps
- Organisation
- R solution de probl mes
- Pens e critique
- Communication
- Style d'apprentissage (Visuel/Auditif/Kinesth sique)
- Intelligence dominante
- Cr ativit , Adaptabilit , R silience

2.2.13 13. Comp tences Linguistiques (10 variables)

- Niveau en arabe (C2-A1)
- Niveau en fran ais (C1-A2)
- Niveau en anglais (B2-A1)
- Niveau en espagnol
- Niveau en allemand
- Autres langues
- Locuteur amazigh

3 Couverture G ographique

Le dataset couvre les 12 r gions administratives du Maroc :

1. Rabat-Sal -K nitra
2. Casablanca-Settat
3. F es-Mekn s
4. Marrakech-Safi
5. Tanger-T touan-Al Hoceima
6. Souss-Massa
7. Oriental
8. Beni Mellal-Khenifra
9. Draa-Tafilalet
10. Laayoune-Sakia El Hamra
11. Guelmim-Oued Noun
12. Dakhla-Oued Ed-Dahab

Fili�re	Sp�cialit�s
Sciences Exp�rimentales	Sciences Physiques (PC), SVT
Sciences Math�matiques	Maths A (SMA), Maths B (SMB)
Sciences �conomiques	Gestion Comptable (SGC), Sciences �conomiques (SE)
Lettres et Sciences Humaines	Lettres, Sciences Humaines
Sciences et Technologies	Sciences M�caniques (STM), Sciences �lectriques (STE)

TABLE 2 – Fili res et sp cialit s du baccalaur at couvertes

4 Fili res Couvertes

5 M thodologie de G n ration

5.1 Approche Utilis e

Les donn es ont  t  g n r es de mani re synth tique en utilisant Python avec les principes suivants :

- **R alisme** : Les valeurs sont bas es sur des distributions r alistes correspondant au contexte marocain
- **Coh rence** : Les corr lations logiques entre variables sont respect es (ex : revenu familial et  quipements disponibles)
- **Diversit ** : Large  ventail de profils socio- conomiques et acad miques
- **Repr sentativit ** : Distribution proportionnelle entre zones urbaines, semi-urbaines et rurales

6 Analyse et Nettoyage des Donn es

Avant de proc der   l’exploration des donn es, une phase initiale d’analyse et de nettoyage a  t  r alis e pour garantir la qualit  des donn es.

6.1 Analyse Initiale

L’examen initial du dataset a r v l  les caract ristiques suivantes :

- **Dimensions** : 10 000 enregistrements et 286 variables.
- **Types de donn es** : 65 variables flottantes (float64), 48 entiers (int64) et 173 objets (object).
- **Valeurs manquantes** : Identification de colonnes enti rement vides (100% manquantes) telles que `type_handicap`, `economie_s1`, `remarques`, et d’autres partiellement renseign es comme `etablissement_precedent` (5925 manquantes) et `annees_redoublees` (4537 manquantes).

6.2 Actions de Nettoyage

Les  tapes suivantes ont  t  appliqu es pour produire le fichier final :

`Morocco_Student_Data_Cleaned.csv`

1. **Suppression des colonnes vides** : Retrait des 17 colonnes ne contenant aucune donn ee.
2. **Traitement des valeurs manquantes** : Imputation des ann es de redoublement par 0 et des activit s par "Aucun".
3. **Formatage des donn ees** : Standardisation du format des dates de naissance et de collecte.

7 Utilisation du Dataset

7.1 Applications Possibles

- D閎veloppement de mod les de pr ediction de r eussite scolaire
- Identification des facteurs de risque d' echec
- Analyse des in galit es  dues
- Conception de syst mes d'alerte pr oce
-  tudes sur l'impact des facteurs socio- conomiques

7.2 Variable Cible

La variable cible   pr dire est **moyenne_annuelle** : la moyenne g n rale de l' tudiant pour l'ann e scolaire en cours. Il s'agit d'un probl me de **r egression**.

7.3 Mod les Recommand s

- **R egression Lin aire** : Mod le de base pour  tablir une r f rence
- **Random Forest Regressor** : Capture les relations non-lin aires
- **XGBoost / Gradient Boosting** : Performance  lev e pour les donn es tabulaires
- **Deep Learning (MLP)** : R seaux de neurones pour les relations complexes

8 Conclusion

Ce data pool constitue une ressource compl te et r aliste pour le d veloppement de mod les de pr dition de la performance des  tudiants marocains. Avec ses 10 000 enregistrements et 286 variables (269 apr s nettoyage), il offre une base solide pour l'analyse exploratoire, le feature engineering, et l'entra nement de mod les de machine learning.

Fichiers g n r s :

`Morocco_Student_Data_Pool.csv` (Donn es brutes)
`Morocco_Student_Data_Cleaned.csv` (Donn es nettoy es)
`analyze_data.ipynb` (Notebook d'analyse)
`clean_data.ipynb` (Notebook de nettoyage)

Date de cr ation : 5 F vrier 2026

Encodage : UTF-8