

Projet d'expansion à l'international d'Academy

Analyse des données de systèmes éducatifs
(données de la banque mondiale)



Contexte de ce travail

Academy souhaite étendre son activité à l'international. Pour cela, il est nécessaire de repérer les pays les plus propices à cette expansion.

Les données de la banque mondiale semblent être un bon point de départ pour mettre en place une stratégie.

Le but de ce travail est donc de mettre en forme ces données pour en extraire un classement de pays. Ce dernier permettra d'élaborer un début de réponse quant à l'implantation d'un service d'éducation en ligne pour des niveaux lycée et université.

Description de la base de données de la banque mondiale

Ensemble de données sur le champ de l'éducation regroupé en 5 dataframes



<u>Nom du jeu de données</u>	Description du jeu de données
EdStatsCountrySeries	Représente une collection de variables disponibles pour chaque pays, à partir desquelles les données sont extraites. Chaque ligne correspond à une variable décrite et contient le code pays, le nom de la variable et la description du mode de recueil.
EdStatsCountry	Contient un ensemble d'informations concernant les différents pays. Chaque ligne représente un pays avec plusieurs variables décrivant les données du pays, leur ancienneté, le système de récupération des données, etc.
EdStatsData	Chaque ligne représente la valeur d'un indicateur pour un pays et pour une année donnée.
EdStatsFootNote	Chaque ligne représente un code de série avec, en indication, le pays auquel il est rattaché, une description et une année.
EdStatsSeries	Chaque ligne représente un nom de série avec les informations qui lui sont rattachées pour la décrire.

Code utilisé pour le traitement de données : Python
IDE: JupyterLab

Description de la base de données de la banque mondiale

Ensemble de données sur le champ de l'éducation regroupé en 5 dataframes

<u>Nom du jeu de données</u>	Description du jeu de données
EdStatsCountrySeries	613 lignes et 4 colonnes.
EdStatsCountry	241 lignes et 32 colonnes.
EdStatsData	886930 lignes et 70 colonnes.
EdStatsFootNote	643638 lignes et 5 colonnes.
EdStatsSeries	3665 lignes et 21 colonnes.

Code utilisé pour le traitement de données : Python
IDE: JupyterLab



academy

Vérifications réalisées

- Vérification de la présence de doublons : Aucun doublon détecté.
- Recherche des colonnes vides et suppression de ces colonnes :
 - EdStatsCountrySeries : “Unnamed: 3”
 - EdStatsCountry: “Unnamed: 31”
 - EdStatsData : “Unnamed: 69”
 - EdStatsFootNote : “Unnamed: 4”.
 - EdStatsSeries : “Other web links”, “Unnamed: 20”, “License Type”, “Notes from original source”, “Unit of measure”, “Related indicators”.
- Description de données qualitatives : l’intégralité des colonnes est décrite dans le code mis à disposition sur GitHub
- Suppression de données “pays” non valides dans le fichier EdStatsCountry

Liste des données “pays” non valides

Dans la serie pays du fichier EdStatsCountry un certain nombre d'items ne semblent pas valides et ont été supprimés. Par exemple :

- Des données correspondant à des régions (“Arab World”, “East Asia & Pacific”, “Europe & Central Asia”, ...)
- Des regroupements selon une typologie (“High income”, “Low income”, ...)

L'ensemble de la liste des éléments supprimés est disponible dans le code fourni sur le GitHub en fin de document.

Obtention d'une liste de variables utilisables

- À partir de la liste des pays, un nettoyage du dataframe `EdStatsCountrySeries`, grâce à un merge avec le dataframe `EdStatsCountry` nettoyé, a permis un nettoyage de l'ensemble des dataframes.
- Une sélection des topics intéressants pour une approche métier a été effectuée sur le dataframe `EdStatsSeries`
- La sélection a été faite pour retenir les sujets suivants : le niveau scolaire, les volumes de populations, les volumes financiers, les capacités de communication pour l'aspect distanciel et la place de l'enseignement privé.
- Grâce au lien fait via les codes de series du dataframe `EdStatsCountrySeries`, le dataframe `EdStatsData` a été réduit aux topics pertinents, limitant le nombre de variables disponibles.

Sélection des années pour l'analyse et mise en forme du dataframe EdStatsData

- Le dataframe contient de nombreuses années, mais toutes ne disposent pas d'un volume de données suffisant ; un certain nombre ont donc été supprimées.
- Les années 1999 à 2015, avec le plus grand nombre de données, ont été conservées.
- Le dataframe EdStatsData a été mis en forme selon la structure suivante : une ligne par pays et une colonne par indicateur. Les valeurs correspondent à la moyenne des données sur l'ensemble des années sélectionnées.

Sélection des variables pour l'analyse

2 phases ont été réalisées pour la sélection des variables d'intérêts :

- 1) Une approche métier qui contient des variables de comme les niveaux de population, les niveaux scolaires visés, le niveau d'utilisation d'internet et les données économiques.
- 2) Une approche statistique avec des matrices de corrélation pour ne sélectionner qu'un nombre réduit de variable.

Matrice de corrélation entre les indicateurs (moyenne 1999-2015)

Indicator Name

Enrolment in secondary education, private institutions, both sexes (number)	1.00	0.75	0.62	-0.08	-0.15	-0.12	0.20	0.70	0.72	0.73	0.72	0.74	0.70	0.70	0.60	0.05	-0.23
Enrolment in secondary general, both sexes (number)	0.75	1.00	0.78	-0.13	-0.13	-0.08	0.02	0.91	0.92	0.93	0.91	0.91	0.84	0.85	0.77	0.14	-0.13
Enrolment in tertiary education, all programmes, both sexes (number)	0.62	0.78	1.00	0.15	0.01	0.19	-0.21	0.67	0.68	0.71	0.69	0.68	0.70	0.72	0.69	0.49	-0.01
GDP per capita (constant 2005 US\$)	-0.08	-0.13	0.15	1.00	0.17	0.88	-0.44	-0.28	-0.26	-0.22	-0.22	-0.27	-0.24	-0.21	-0.03	0.71	0.07
Government expenditure on education as % of GDP (%)	-0.15	-0.13	0.01	0.17	1.00	0.20	-0.21	-0.22	-0.21	-0.20	-0.21	-0.22	-0.20	-0.19	-0.18	0.22	0.08
Internet users (per 100 people)	-0.12	-0.08	0.19	0.88	0.20	1.00	0.55	-0.24	-0.22	-0.19	-0.20	-0.25	-0.20	-0.17	-0.00	0.77	0.10
Population growth (annual %)	-0.20	0.02	-0.21	-0.44	-0.21	-0.55	1.00	0.16	0.15	0.14	0.14	0.20	0.11	0.08	-0.04	-0.64	-0.24
Population of the official age for lower secondary education, both sexes (number)	0.70	0.91	0.67	-0.28	-0.22	-0.24	0.16	1.00	0.99	0.98	0.98	0.98	0.86	0.85	0.78	-0.04	-0.18
Population of the official age for secondary education, both sexes (number)	0.72	0.92	0.68	-0.26	-0.21	-0.22	0.15	0.99	1.00	0.99	0.99	0.99	0.86	0.86	0.79	-0.03	-0.19
Population of the official age for tertiary education, both sexes (number)	0.73	0.93	0.71	-0.22	-0.20	-0.19	0.14	0.98	0.99	1.00	0.98	0.99	0.87	0.87	0.80	0.01	-0.20
Population of the official age for upper secondary education, both sexes (number)	0.72	0.91	0.69	-0.22	-0.21	-0.20	0.14	0.98	0.99	0.98	1.00	0.98	0.86	0.86	0.79	-0.01	-0.18
Population of the official entrance age to secondary general education, both sexes (number)	0.74	0.91	0.68	-0.27	-0.22	-0.25	0.20	0.98	0.99	0.99	0.98	1.00	0.87	0.87	0.79	-0.05	-0.20
Population, ages 12-18, total	0.70	0.84	0.70	-0.24	-0.20	-0.20	0.11	0.86	0.86	0.87	0.86	0.87	1.00	1.00	0.80	-0.00	-0.14
Population, ages 15-24, total	0.70	0.85	0.72	-0.21	-0.19	-0.17	0.08	0.85	0.86	0.87	0.86	0.87	1.00	1.00	0.80	0.03	-0.14
Population, ages 15-64, total	0.60	0.77	0.69	-0.03	-0.18	-0.00	-0.04	0.78	0.79	0.80	0.79	0.79	0.80	0.80	1.00	0.16	-0.15
School life expectancy, tertiary, both sexes (years)	-0.05	0.14	0.49	0.71	0.22	0.77	-0.64	-0.04	-0.03	0.01	-0.01	-0.05	-0.00	0.03	0.16	1.00	0.14
Unemployment, total (% of total labor force)	-0.23	-0.13	-0.01	0.07	0.08	0.10	-0.24	-0.18	-0.19	-0.20	-0.18	-0.20	-0.14	-0.14	-0.15	0.14	1.00



Indicator Name

Sélection des variables pour l'analyse

4 variables sont gardées et utilisées pour la construction d'un classement de pays d'intérêt :

- 1) "Enrolment in secondary general, both sexes (number)" : utilisé dans l'analyse principale.
- 2) "School life expectancy, tertiary, both sexes (years)" : utilisé dans l'analyse principale.
- 3) "Unemployment, total (% of total labor force)" : utilisé dans une analyse secondaire.
- 4) "Population growth (annual %)" : utilisé dans une analyse secondaire.

Les deux premières variables sont utilisées pour établir un classement des pays qui sont le mieux dotés pour une implantation allant du lycée jusqu'à l'université.

Les deux variables suivantes permettent un éclairage des conditions d'évolution à prendre en compte et de la stratégie possible concernant le chômage et l'apport de la formation à la problématique. Cette analyse a été appliquée pour les 20 premiers pays du classement

Classement des pays selon les variables d'intérêts

Country Name	Enrolment in secondary general, both sexes (number)	School life expectancy, tertiary, both sexes (years)_y	Classement moyen	Rang final	Note croissance	Note chômage
United States	3.0	3.0	3.0	1	Faible	Tension modérée
Russian Federation	8.0	14.0	11.0	2	Negatif	Tension modérée
Korea, Rep.	26.0	1.0	13.5	3	Faible	Tension faible
Ukraine	25.0	12.0	18.5	4	Negatif	Tension modérée
Argentina	23.0	16.0	19.5	5	Stable	Accès difficile à l'emploi

Classement des pays selon les variables d'intérêts

Selon l'analyse réalisée, les pays à privilégier sont les suivants :

United States, Russian Federation, Korea Rep., Ukraine, Argentina, Germany, Japan, United Kingdom, France, Spain, Australia, Italy, Canada, Poland, Turkey, Thailand, Belarus, Chile, Greece, Brazil

Bien entendu, les données étant datées de 2015, certains pays peuvent avoir connu une variation selon plusieurs variables non répertoriées ici (par exemple le contexte géopolitique).

Ressources



Lien Github vers le code d'analyse

- <https://github.com/Thomas-Auvin/Projet2>

Lien vers les jeux de données

- <https://datacatalog.worldbank.org/search/dataset/0038480>

