

○○○○

PRESENTATION

ACADEMY

○○○○

SOMMAIRE

- Problématique liée à l'entreprise
- Validation de la qualité des données
- Description du Dataset principal
- Sélection des indicateurs via les informations pertinentes
- Description des indicateurs retenus
- Statistiques des indicateurs retenus
- Top 10 des pays sélectionnés
(Score synthétique)



PROFIL ACADAMY

Qui somme nous ?

L'entreprise est une organisation dont l'objectif est d'étendre ses activités à l'international (niveau lycée et études supérieurs).

Notre mission

Notre mission est d'utiliser l'analyse des données sur l'éducation pour informer et faciliter l'expansion internationale de l'entreprise, en identifiant des opportunités de marché prometteuses.

Notre vision

Notre vision est de devenir un leader mondial dans les services éducatifs en ligne en exploitant intelligemment les données pour améliorer l'accès à l'éducation à l'échelle mondiale.

Our Goals

Notre objectif est de trouver des pays attractifs afin de contribuer à l'amélioration de l'accès à l'éducation à travers le monde.

Quels sont les pays à fort
potentiel pour nos services
d'éducation en ligne ?
Comment les départager ?

PROBLÉMA TIQUE



THE WORLD BANK
IBRD • IDA

**PRÉ-ANALYSE DES DONNÉES DU PORTAIL
EDUCATION STATISTICS (EDSTATS)**

- 🔍 Indicateurs pertinents.
- 📊 Indicateurs statistiques
- 🎯 Top 10 pays

VALIDATION DE LA QUALITÉ DES DONNÉES



EdStatsCountry-Series.csv

Les jeux de données
EdstatsCountry-Series
et EdstatsFootNote sont
exclus de notre analyse
en raison de leur
moindre pertinence et
de leur apport limité à
nos objectifs d'analyse.

EdStatsFootNote.csv



EdStatsData.csv

se compose de

- 886 930 lignes
- 70 colonnes.

Le taux moyen de complétion des données est de 13,9%.
Pas de doublons, evolution des indicateurs par pays (1970 - 2100)
c'est notre jeu de données principal

EdStatsCountry.csv

se compose de

- 241 lignes
- 32 colonnes

Pas de doublons.
Informations géographiques,
données économique.

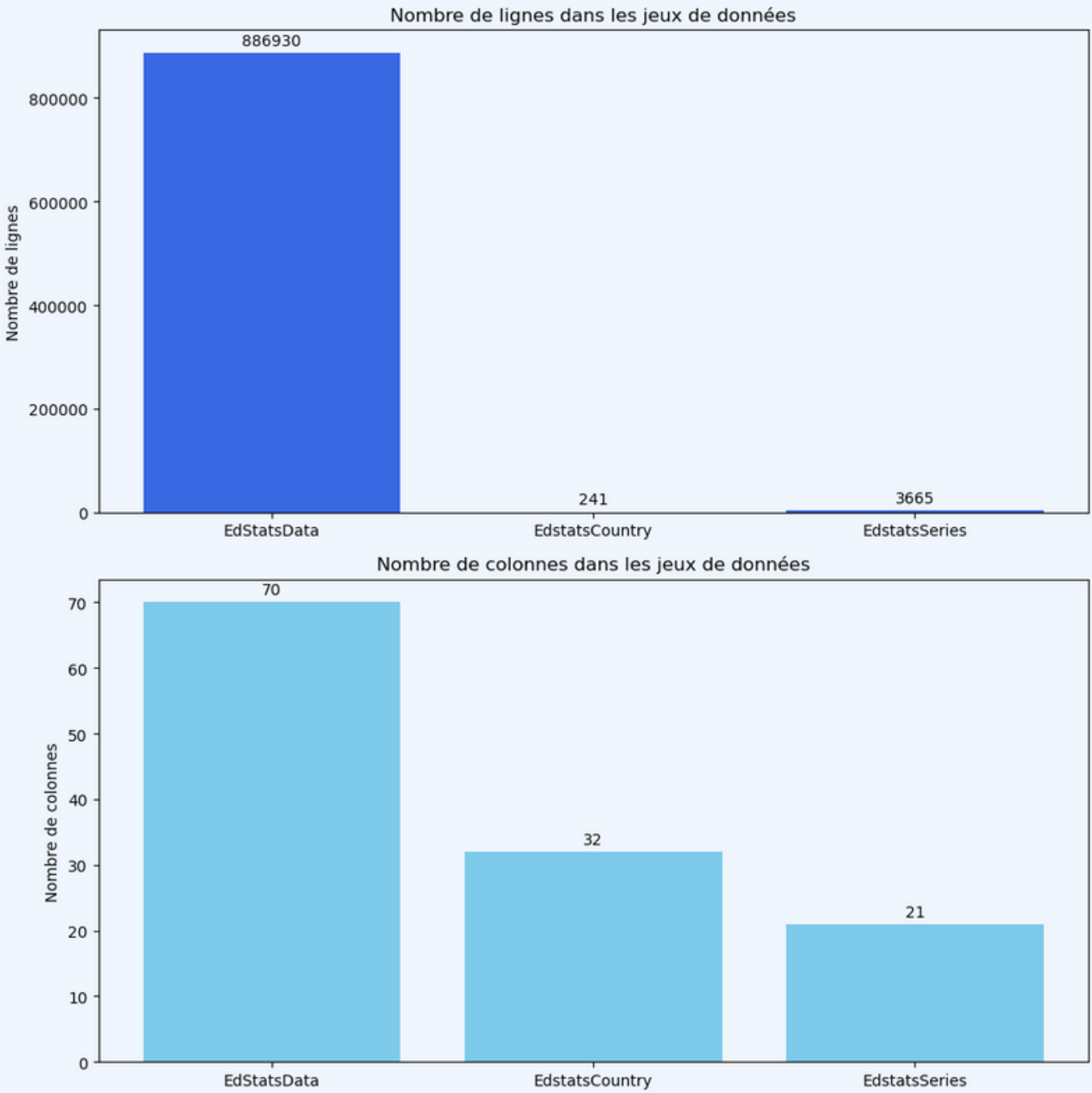
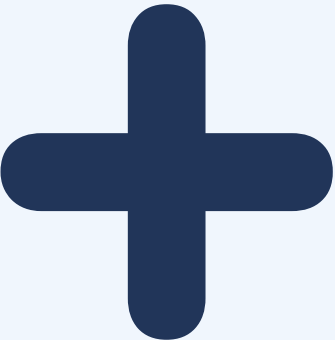
EdStatsSeries.csv

se compose de

- 3665 lignes
- 21 colonnes

Il présente une valeur intéressante à savoir les Topic (thème de l'indicateur)

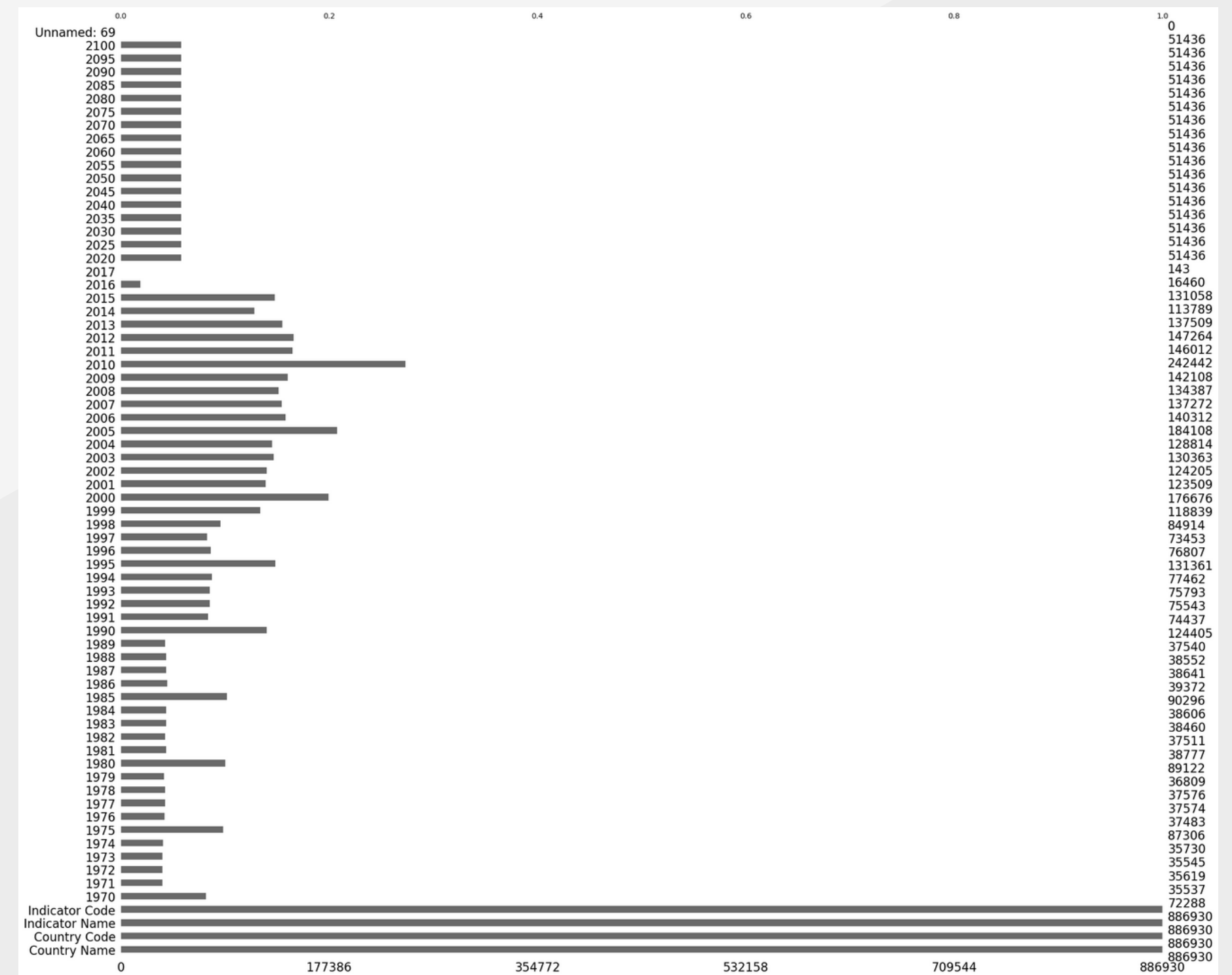
VALIDATION DE LA QUALITÉ DES DONNÉES





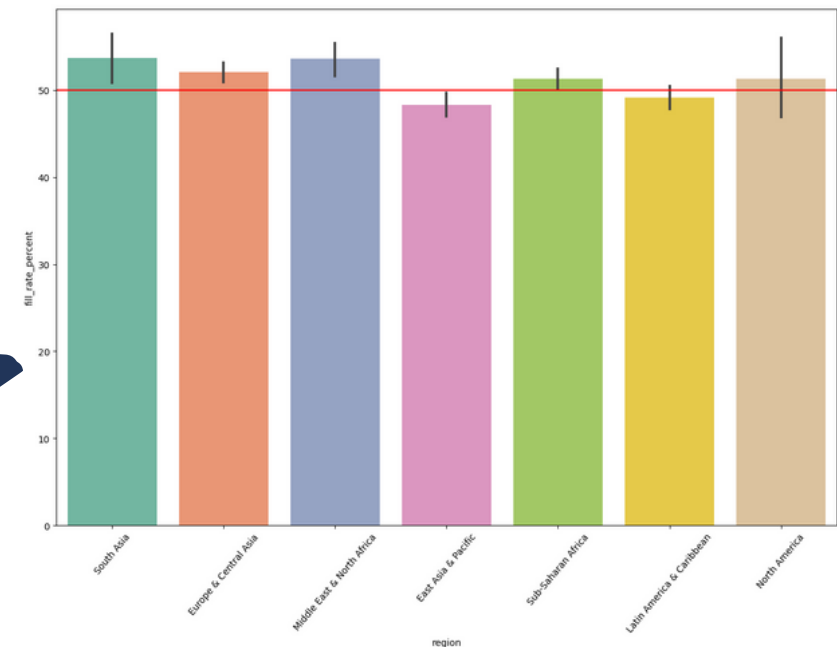
NOTRE JEU DE DONNÉE PRINCIPAL

- Il est essentiel de souligner que les valeurs manquantes sont entre 0 et 95 % dans ce dataset
- Le data set est composé comme tel :
Une ligne représente un indicateur pour son pays par année (1970 - 2100)
- Informations utiles : une plage temporelle pour connaître l'évolution dans le temps, sélectionner des indicateurs pertinents.



SELECTIONS DES INFORMATIONS PERTINENTES

```
1 # Liste des indicateurs restants
2 indicateurs_restants = indicator_list_deleted
3
4 # Liste des indicateurs à filtrer
5 # Indicateurs potentiels que nous souhaitons sélectionner
6 indicateurs_potential = ['GDP', 'age 15-19', 'age 20-24', 'Population, ages 15-24, total',
7 'Population of the official age for secondary education, both sexes (number)',
8 'Population of the official age for tertiary education, both sexes (number)',
9 'computers', 'Gross enrolment ratio, secondary, both sexes (%)',
10 'Gross enrolment ratio, secondary, gender parity index (GPI)',
11 'Gross enrolment ratio, tertiary, both sexes (%)',
12 'Gross enrolment ratio, tertiary, gender parity index (GPI)',
13 'DHS: Secondary completion rate',
14 'Enrolment in tertiary education per 100,000 inhabitants, both sexes',
15 'Internet users']
16
17 # Filtrer Les indicateurs qui sont dans la liste indicateurs_potential
18 indicateurs_pertinents = [indicator for indicator in indicateurs_restants if any(ind in indicator for ind in indicateurs_potential)]
19
20 # Créer un DataFrame avec Les indicateurs pertinents
21 data_potential_indicators = data_EdStats[data_EdStats['Indicator Name'].isin(indicateurs_pertinents)]
22
23 print(len(data_potential_indicators['Indicator Name'].unique().tolist()))
24 data_potential_indicators['Indicator Name'].unique().tolist()
25
```



Step 2

Recherche d'indicateurs pertinents et compréhension des termes pour les indicateurs : <https://datatopics.worldbank.org/education/indicator>

Nb : 68

Step 4

Analyser le taux de remplissage avec pour seuil le quartile 3 = 50% .

Nb : 35

1

Nb : 3665

Step 1

Brainstorming

- Population visée
- Aspects économique
- Education : taux de scolarisation, % d'élèves diplômés..

2

3

Nb : 35

Step 3

Analyse par plage temporelle et par thème (Topic)

4

5

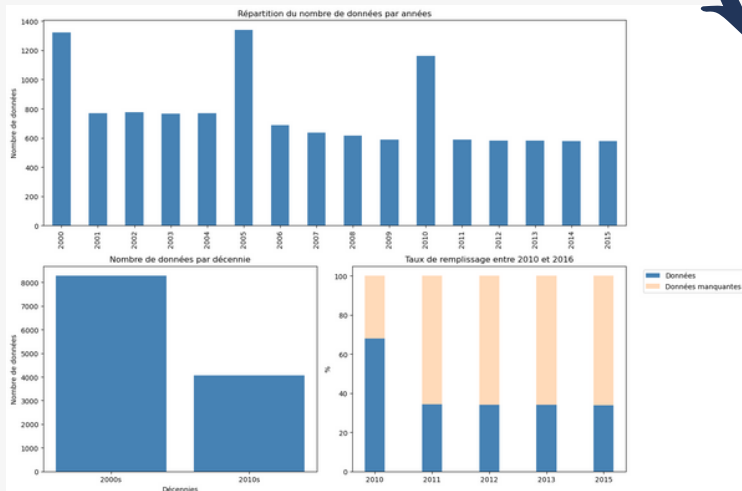
Nb : 8

Step 5

Affiner la sélection avec la recherche par mot clefs.

```
1 # Liste de tout les indicateurs
2 indicator=data_EdStats['Indicator Name'].unique()
3 nb_indicator=len(indicator)
4 country_=data_EdStats['Country Name'].unique()
5 nb_country=len(country_)
6
7
8 print('Nombre d indicateurs : ')
9 print(nb_indicator)
10 print('Nombre de pays : ')
11 print(nb_country)
12
13 nb_country*nb_indicator==data_EdStats.shape[0]
14 indicator.tolist()
```

Nombre d indicateurs :
3665
Nombre de pays :
242



```
1 # Liste des topics souhaités
2 topics_souhaites = ["Attainment",
3 "Education Equality",
4 "Tertiary",
5 "Economic Policy & Debt: National accounts: US$ at current prices",
6 "Secondary", "Population", "Infrastructure: Communications"]
7
8 # Filtrer le DataFrame pour ne conserver que Les lignes avec Les topics souhaités
9 data_academy = data_academy[data_academy['topic'].isin(topics_souhaites)]
10
11 # Obtenir La Liste des noms des indicateurs correspondants
12 noms_des_indicateurs = data_academy['indicator_name'].unique()
13
14 # Afficher La Liste des noms des indicateurs correspondants
15 print(len(noms_des_indicateurs))
16 print(noms_des_indicateurs)
```

```
1 # Liste des mots-clés pour la filtration (en minuscules)
2 mots_cles = ['gdp',
3 'internet',
4 'computer',
5 'population, ages 15-24, total',
6 'completed secondary',
7 'completed tertiary']
8
9 # Filtrer Les indicateurs en utilisant une compréhension de liste
10 indicators_filtered = [indicator for indicator in list_indicators_refined if any(mot in indicator.
11
12 # Afficher La Liste des indicateurs filtrés
13 print(indicators_filtered)
14 print(len(indicators_filtered))
15
```

['Barro-Lee: Percentage of population age 15-19 with secondary schooling. Completed Secondary', 'Barro-Lee: Percentage of population age 15-19 with tertiary schooling. Completed Tertiary', 'Barro-Lee: Percentage of population age 20-24 with secondary schooling. Completed Secondary', 'Barro-Lee: Percentage of population age 20-24 with tertiary schooling. Completed Tertiary', 'GDP per capita (current US \$)', 'Internet users (per 100 people)', 'Personal computers (per 100 people)', 'Population, ages 15-24, total']

LES INDICATEURS RETENUS

- Education
- Economie
- Numérique
- Démographique

BAR.SEC.CMPT.1519.ZS

Pourcentage de la population âgée de 15 à 19 ans ayant suivi un enseignement secondaire.

BAR.TER.CMPT.1519.ZS

Pourcentage de la population âgée de 15 à 19 ans ayant suivi un enseignement supérieur.

BAR.SEC.CMPT.2024.ZS

Pourcentage de la population âgée de 20 à 24 ans ayant suivi un enseignement secondaire.

BAR.TER.CMPT.2024.ZS

Pourcentage de la population âgée de 20 à 24 ans ayant fait des études supérieures.

NY.GDP.PCAP.CD

PIB par habitant (USD courants)

IT.NET.USER.P2

Utilisateurs d'Internet (pour 100 personnes)

IT.CMP.PCMP.P2

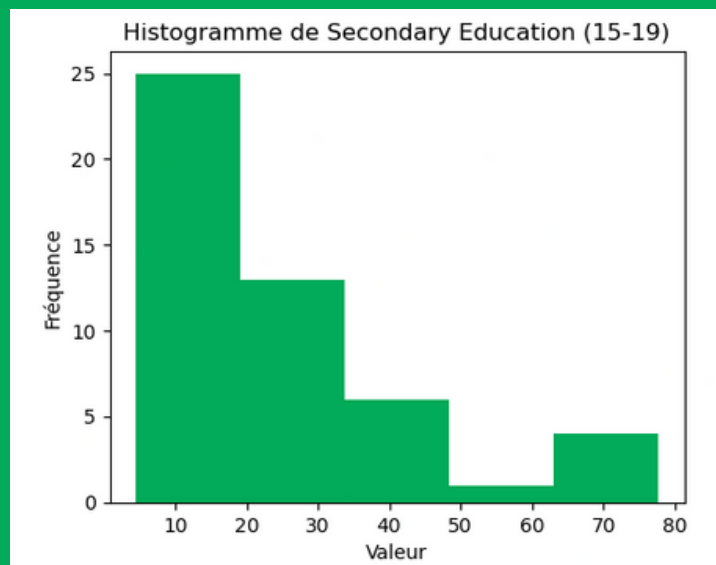
Ordinateurs personnels (pour 100 personnes)

SP.POP.1524.TO.UN

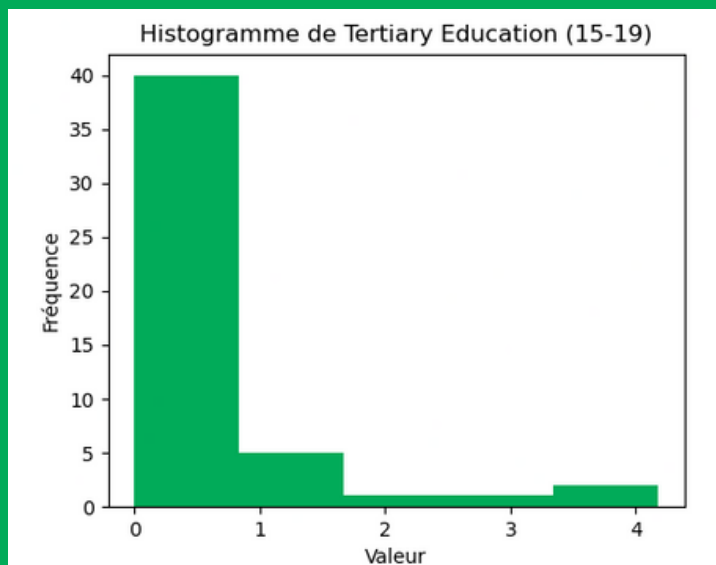
Population âgée de 15 à 24 ans, total

STATISTIQUE DES INDICATEURS RETENUS

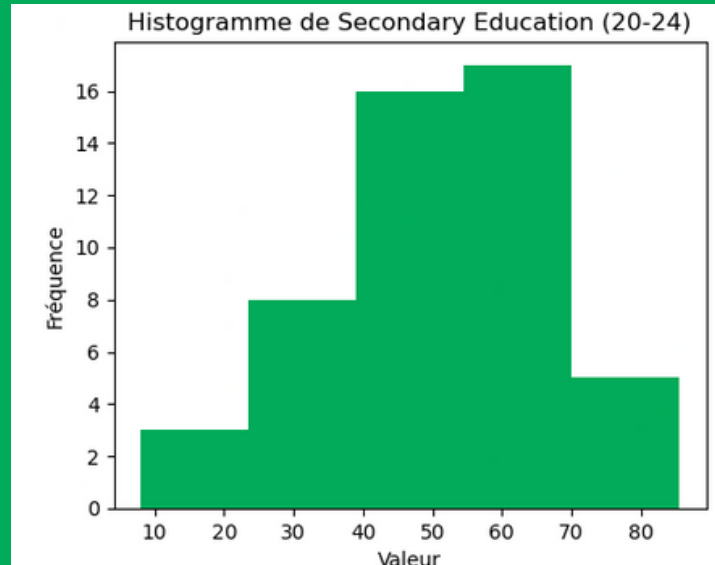
BAR.SEC.CMPT.1519.ZS



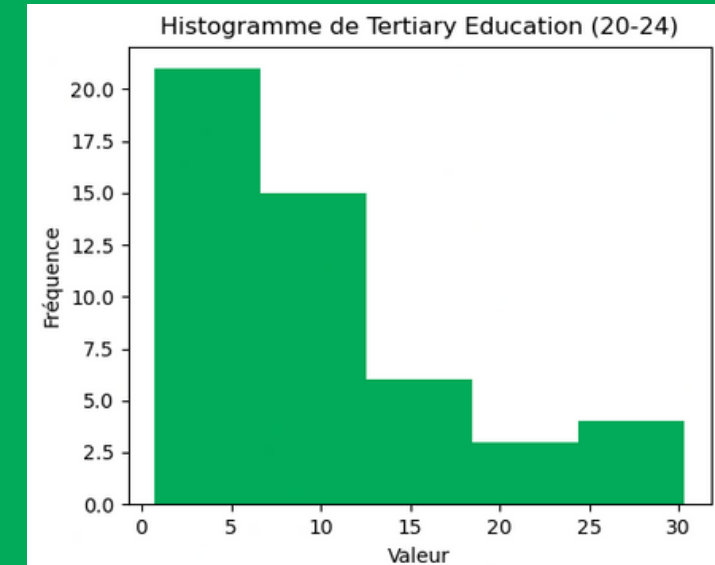
BAR.TER.CMPT.1519.ZS



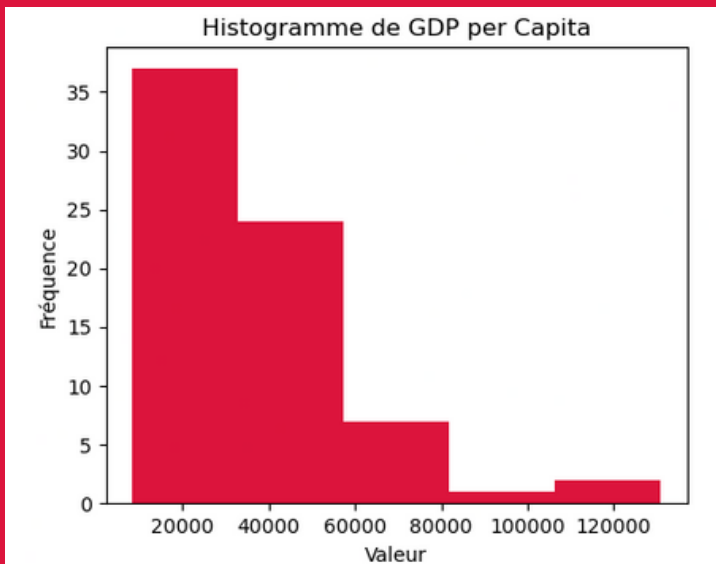
BAR.SEC.CMPT.2024.ZS



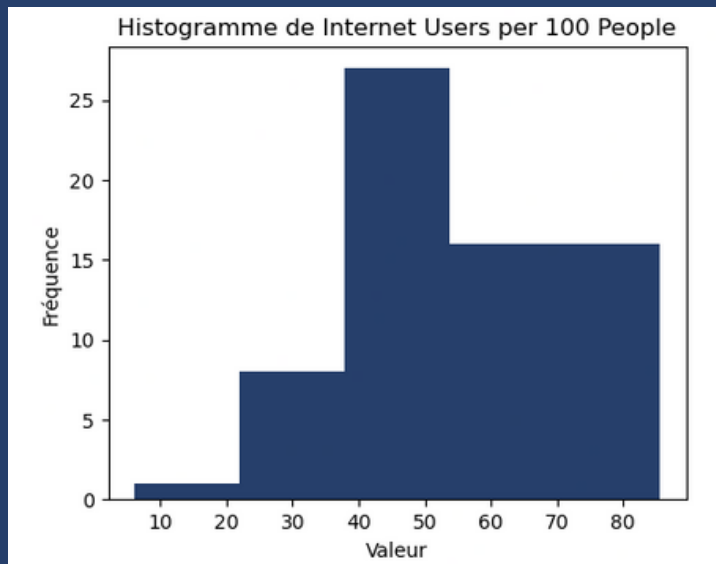
BAR.TER.CMPT.2024.ZS



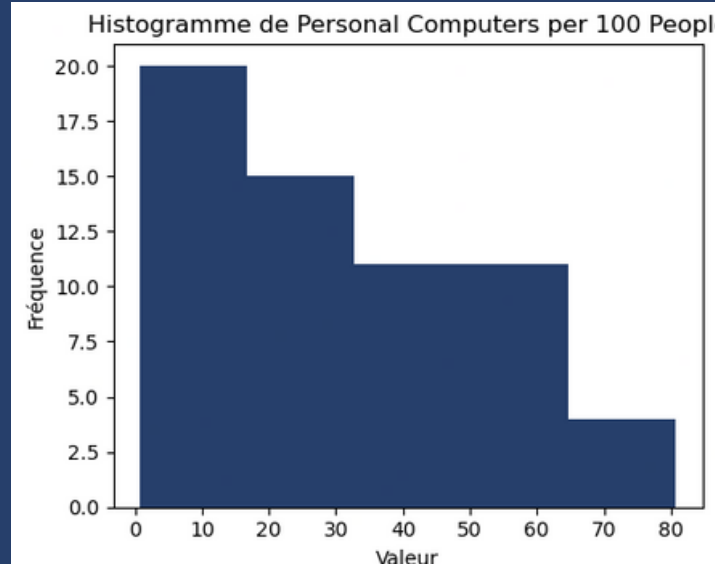
NY.GDP.PCAP.CD



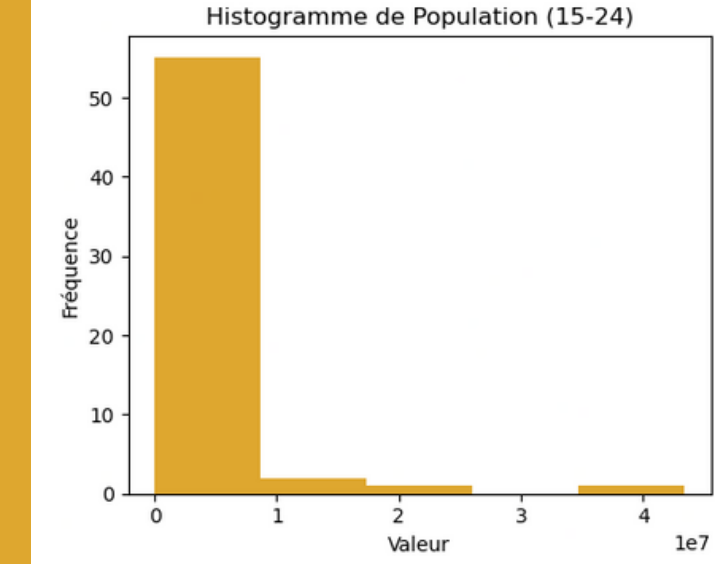
IT.NET.USER.P2



IT.CMP.PCMP.P2



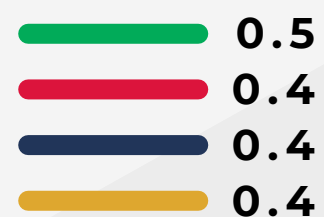
SP.POP.1524.TO.UN





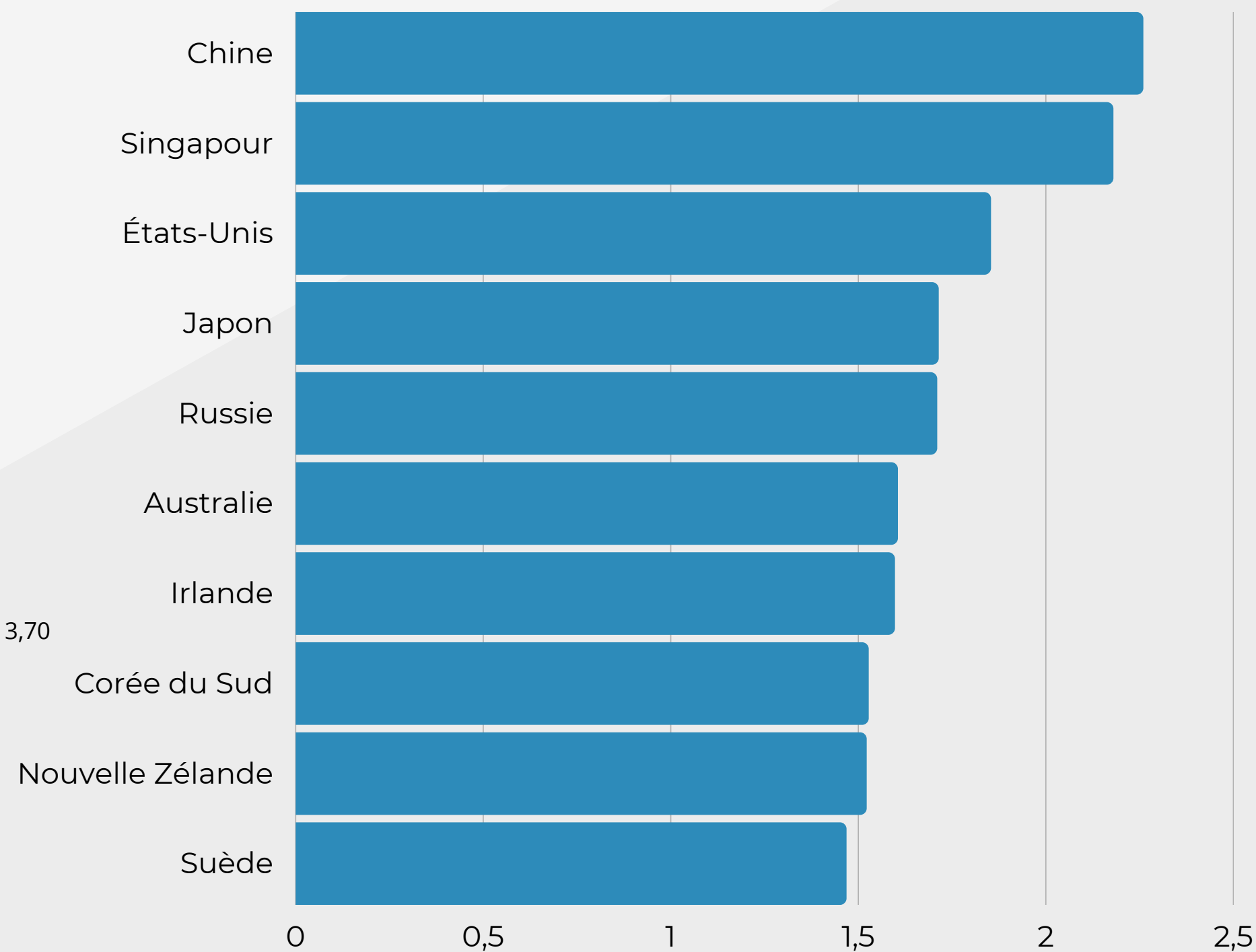
TOP 10 DES PAYS RETENUS

(SCORE SYNTHÉTIQUE AVEC PONDÉRATION)



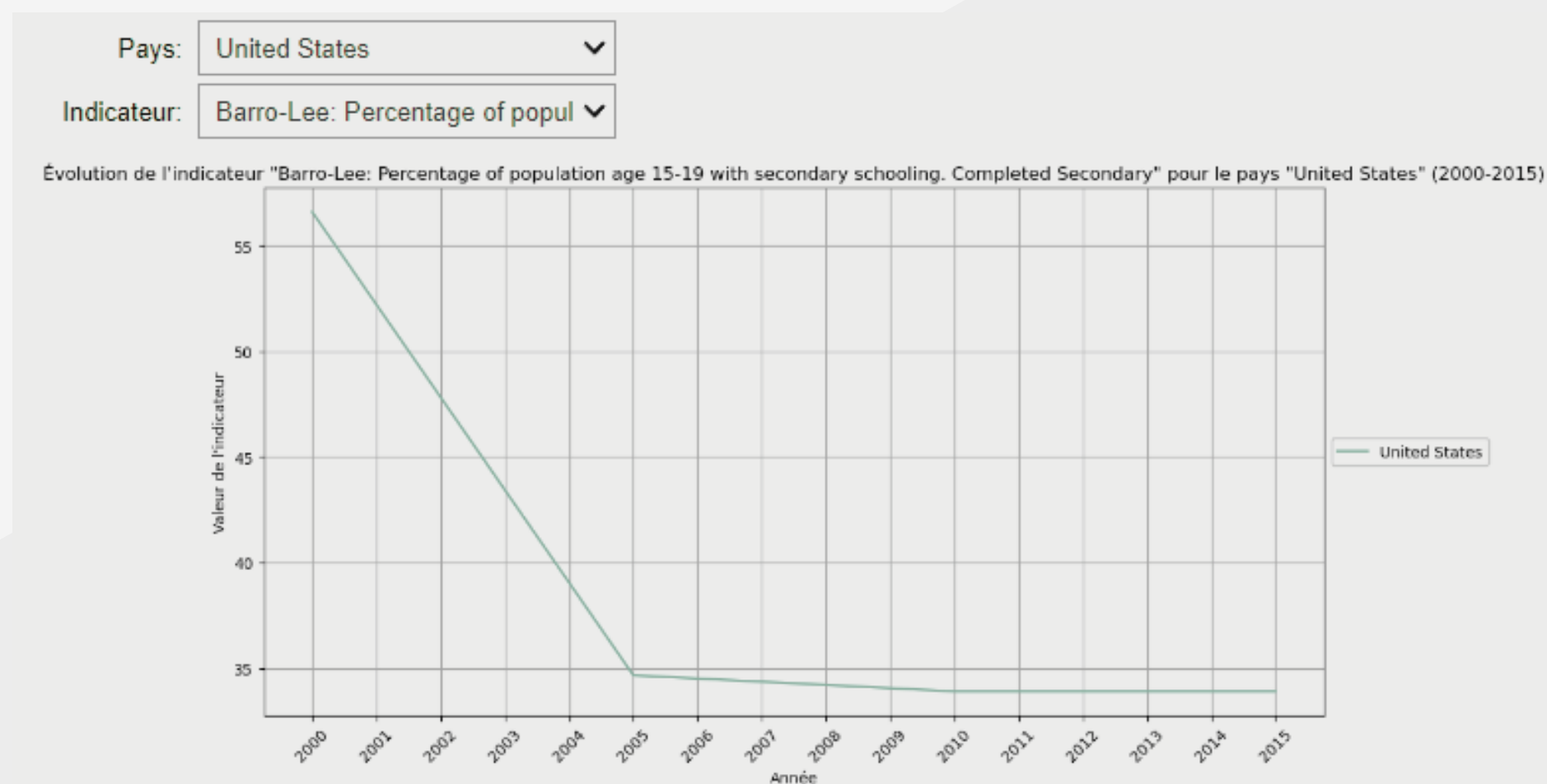
la valeur maximale peut atteindre 3,70

Top 10 des Pays avec les Scores Synthétiques les Plus Élevés





EVOLUTION DANS LE TEMPS (INDICATEURS PAR PAYS)





MERCI

Des question ?

