A person is seen from the side, sitting at a wooden desk and using a laptop. The laptop screen displays a video of a woman and a child, with the text 'ONLINE CLASS' overlaid. The person's hand is on the laptop keyboard. On the desk, there are various items: a pair of glasses, a notebook, a pen, and a cup of pencils. The background is slightly blurred, showing a desk lamp and some papers.

Analysez des données de systèmes éducatifs

Mitra Dadgar _ Data Scientist

1

PRÉSENTATION DU CONTEXTE

2

PRÉSENTATION DES DONNÉES

3

QUALITÉ DES DONNÉES

4

SÉLECTION DES INDICATEURS

5

ANALYSE UNIVARIÉE

6

SCORING

7

CLASSEMENT FINAL

8

ÉVOLUTION DU POTENTIEL

9

CONCLUSION

● ● ● Présentation du contexte

- ❖ **Academy (start-up de la EdTech) :**
propose des contenus de la formation en ligne niveau lycée et université.
- ❖ ***l'objectif*** : Analyser la possibilité d'un projet d'expansion à l'international de l'entreprise
- ❖ Ce jeu de données peut informer les décisions d'ouverture vers de nouveaux pays:
 - Quels sont les pays avec un fort potentiel de clients pour nos services ?
 - Pour chacun de ces pays, quelle sera l'évolution de ce potentiel de clients ?
 - Dans quels pays l'entreprise doit-elle opérer en priorité ?



● ● ● Présentation des données

- ❖ **La source du dataset:** Banque mondiale
- ❖ **Nombre de lignes :** 886930 lignes
- ❖ **Nombre de colonnes :** 101 colonnes
- ❖ **Nombre de pays :** 1302 pays
- ❖ **Nombre d'indicateurs:** 3665 indicateurs

Comment sélectionner les indicateurs intéressants ?
sélectionner les indicateurs les mieux remplis



● Avant nettoyage

```
# Taille de notre dataframe  
df.shape
```

```
(886930, 101)
```

THE WORLD BANK

Help us improve this section of the site. Can we get your feedback? [Click here](#)

Sign In

DataBank Education Statistics - All Indicators

Table Chart Map Metadata

Clear Selection | Add Country (242) | Add Series (1468) | Add Time (7)

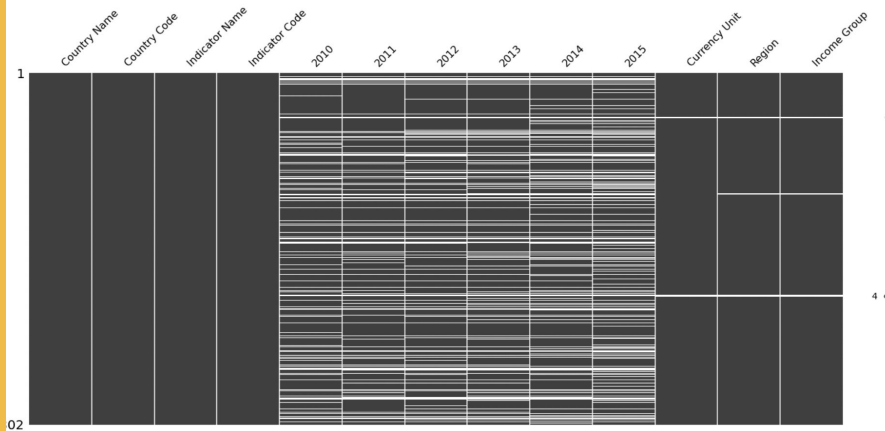
Internet users (per 100 people)

	2010	2014	2015	2016	2017	2018	2019
Afghanistan	4.0	7.0	8.3	11.2	13.5
Albania	45.0	60.1	63.3	66.4	71.8
Algeria	12.5	29.5	38.2	42.9	47.7	59.6	..
American Samoa
Andorra	81.0	95.9	96.9	97.9	91.6
Angola	2.8	21.4	12.4	13.0	14.3
Antigua and Barbuda	47.0	67.8	70.0	73.0	76.0
Arab World	26.7
Argentina	45.0	64.7	68.0	71.0	74.3
Armenia	25.0	54.6	59.1	64.3	64.7
Aruba	62.0	83.8	88.7	93.5	97.2
Australia	76.0	84.0	84.6	86.5	86.5
Austria	75.2	81.0	83.9	84.3	87.9	87.7	..
Azerbaijan	46.0	75.0	77.0	78.2	79.0	79.8	..
Bahamas, The	43.0	76.9	78.0	80.0	85.0
Bahrain	55.0	90.5	93.5	98.0	95.9	98.6	..
Bangladesh	3.7	13.9	14.4	18.0	15.0
Barbados	65.1	75.2	76.1	79.5	81.8

Source: Education Statistics - All Indicators. Click on a metadata icon for original source information to be used for citation.

● ● ● Qualité des données

- Comporte-t-il beaucoup de données manquantes?
- Existe-t-il des valeurs duplicates dans ce jeu de données ?
Non
- Identifier les indicateurs en fonction de la qualité



● Après nettoyage

```
df2.shape
```

```
(875, 14)
```

Taux de valeurs manquantes par année

2017	99.983877
2016	98.144160
1971	95.993258
1973	95.992356
1972	95.984012
1974	95.971497
1979	95.849842
1976	95.773849
1982	95.770692
1989	95.767422
1977	95.763589
1978	95.763364
1983	95.663694
1988	95.653321
1984	95.647233
1987	95.643286
1981	95.627953
1986	95.560867
2020	94.200670
2025	94.200670
2030	94.200670
2035	94.200670
2040	94.200670
2045	94.200670
2050	94.200670
2055	94.200670
2060	94.200670
2065	94.200670
2070	94.200670
2075	94.200670
2080	94.200670
2085	94.200670
2090	94.200670
2095	94.200670
2100	94.200670
1970	91.849639
1997	91.718287
1991	91.607342
1992	91.482642
1993	91.454455
1996	91.340128
1994	91.266278
1998	90.426076
1975	90.156382
1980	89.951631
1985	89.819264
2014	87.170464

● ● ● Sélection des Indicateurs

les indicateurs importants pour l'étude de marché :

→ Quel est le nombre de clients potentiels dans le pays
(ceux qui ont envie de se former)
niveau lycée et université.



Enrolment in tertiary education, all programmes,
both sexes (number)



Population, ages 15-24, total



Population of the official entrance age to
secondary general education, both sexes
(number)



Population growth (annual %)

→ Les prospects peuvent-ils accéder à notre service
(ordinateur/smartphone/internet)



Internet users (per 100 people)

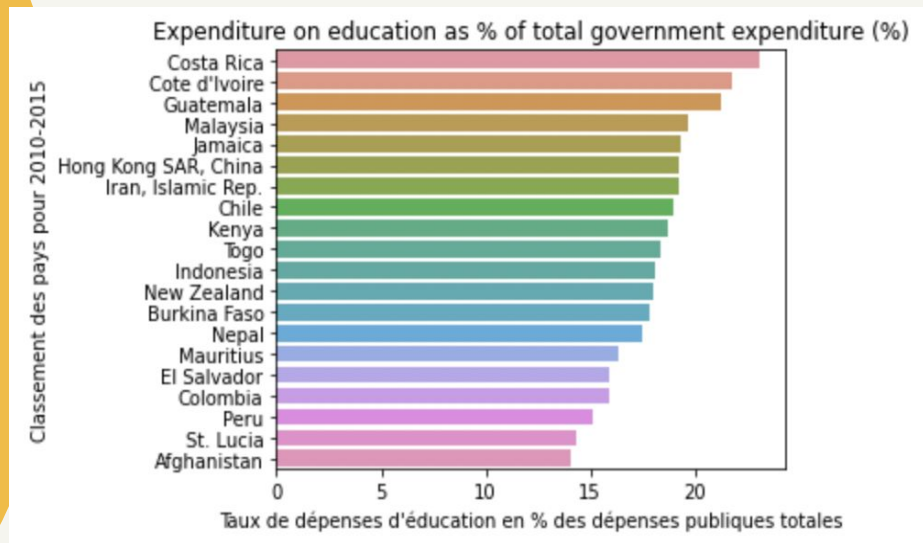
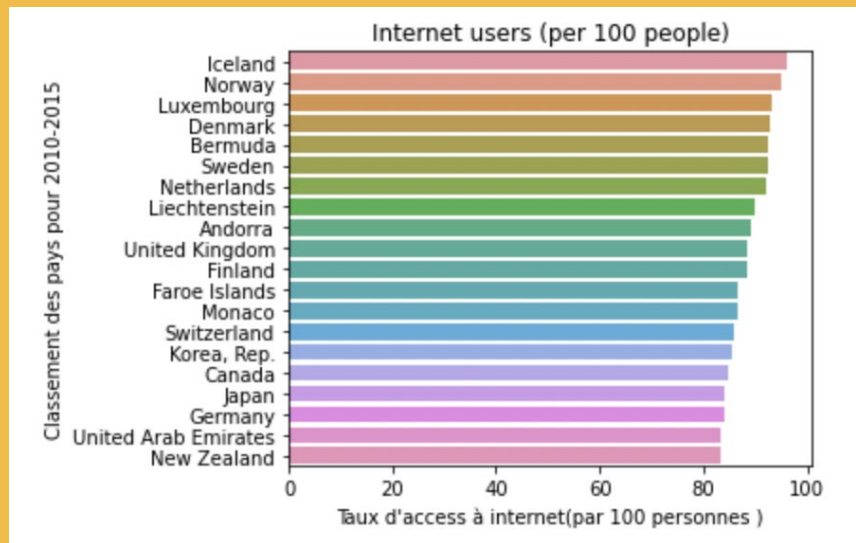
→ Les prospects peuvent-ils payer notre service



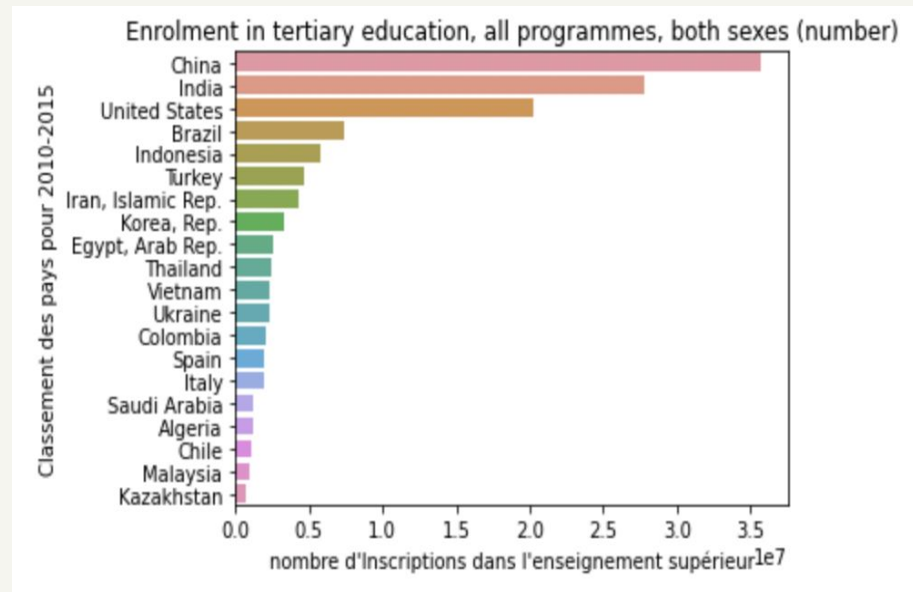
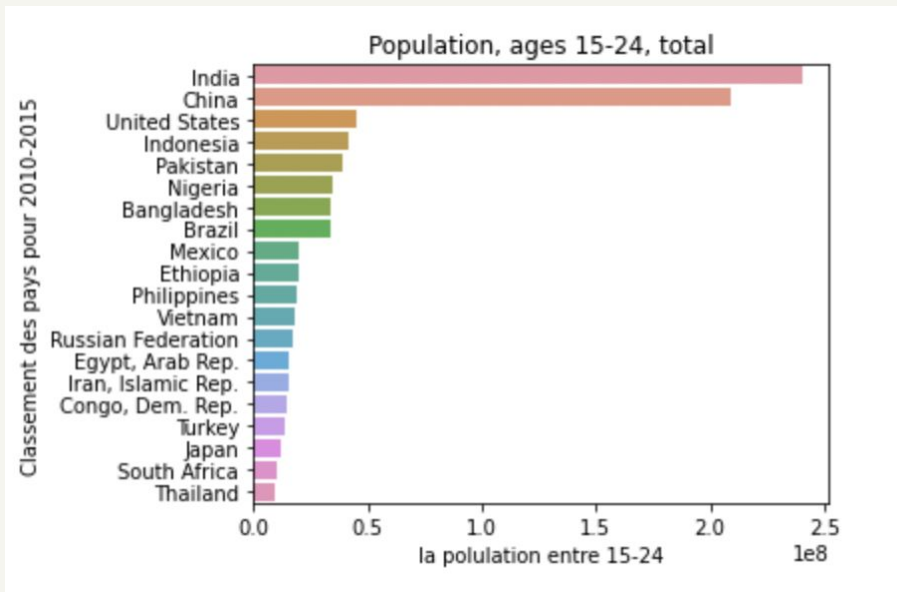
Expenditure on education as % of total
government

Accès à Internet et dépense dans l'éducation

❖ Classement des pays par indicateurs pour 2010-2015



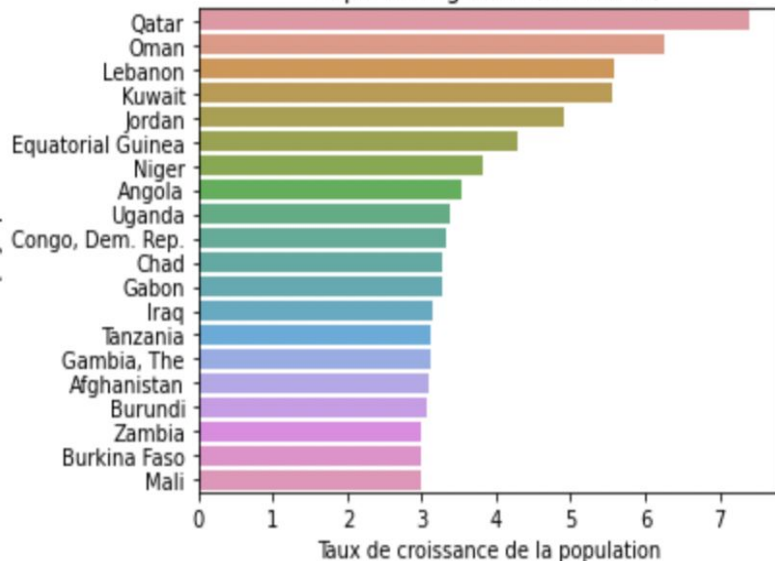
Démographie et éducation



Démographie et âge d'entrée dans le secondaire

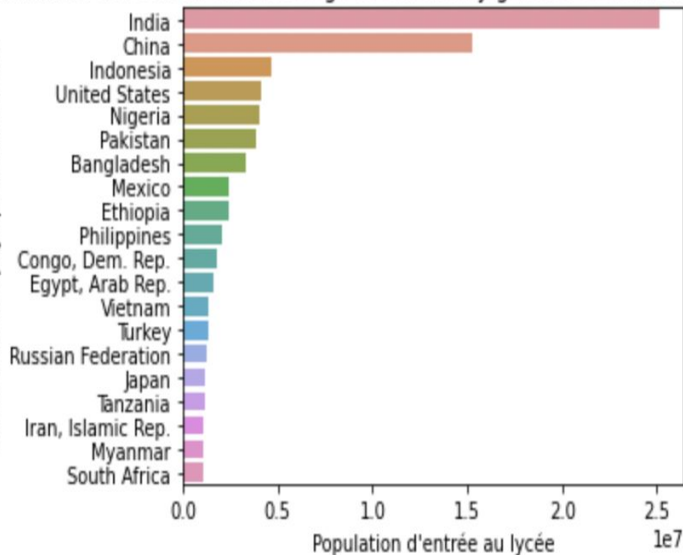
Population growth (annual %)

Classement des pays pour 2010-2015



Population of the official entrance age to secondary general education, both sexes (number)

Classement des pays pour 2010-2015



● ● ● Scoring

On utilise un classement par quartiles un score entre 0 et 3 pour chaque indicateur afin de noter le critère indicateur de chaque pays (3 pays pour la plus potentiel, 0 pour la moins potentiel).

	Average Years	Score_Internet_User	Score_Population_Growth	Score_Population_Secondary	Score_Population_15_24	Score_Enrolment_Tertiary	Score
Country Name							
Afghanistan	7.483954e+06	0.0	3.0	3.0	3.0	0.0	
Albania	7.862124e+05	2.0	0.0	1.0	1.0	1.0	
Algeria	8.665805e+06	1.0	2.0	3.0	3.0	3.0	
American Samoa	-2.057934e-01	0.0	0.0	0.0	0.0	0.0	
Andorra	8.788385e+01	3.0	0.0	0.0	0.0	0.0	
...
Virgin Islands (U.S.)	4.243283e+01	2.0	0.0	0.0	0.0	0.0	
West Bank and Gaza	1.288042e+06	2.0	3.0	1.0	1.0	1.0	
Yemen, Rep.	6.307480e+06	1.0	3.0	3.0	3.0	0.0	
Zambia	3.252853e+06	1.0	3.0	2.0	2.0	0.0	
Zimbabwe	3.715917e+06	0.0	2.0	2.0	2.0	0.0	

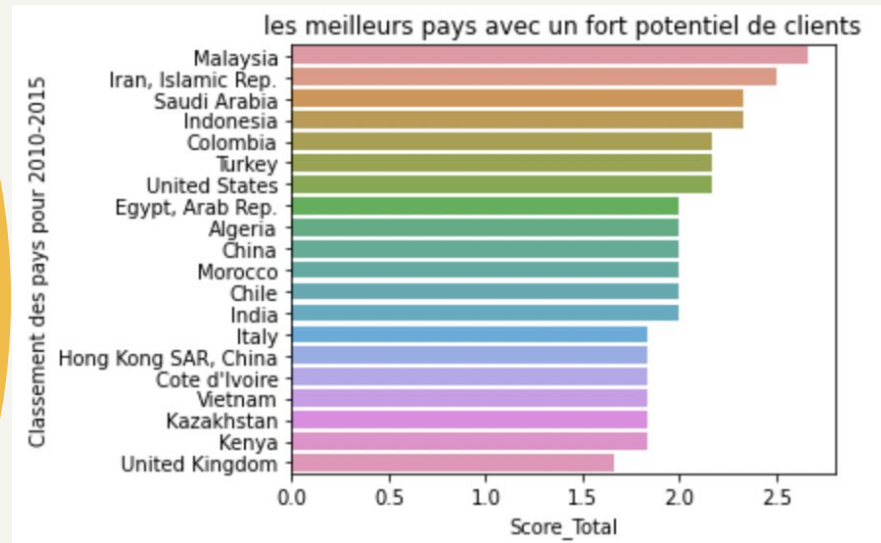
```
df3['Score_Internet_User'] = pd.qcut(df_Indict_Internet['Average Years'], 4, labels = False)
df3['Score_Population_Growth'] = pd.qcut(df_indic_Population_Growth['Average Years'], 4, labels = False)
df3['Score_Population_Secondary'] = pd.qcut(df_Population_Secondary['Average Years'], 4, labels = False)
df3['Score_Population_15_24'] = pd.qcut(df_Population_15_24['Average Years'], 4, labels = False)
df3['Score_Enrolment_Tertiary'] = pd.qcut(df_Enrolment_Tertiary['Average Years'], 4, labels = False)
df3['Score_Indic_Expenditure'] = pd.qcut(df_Indic_Expenditure['Average Years'], 4, labels = False)
```

● ● ● Classement final

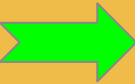
Sélectionner les pays pour lesquels tous les scores sont supérieures à une valeur (le 3ème quartile),

Par exemple un pays qui est dans le 3.0 sur plusieurs indicateurs aura un meilleur score qu'un pays qui n'y apparait qu'une fois (les meilleurs pays sont ceux pour lesquels la valeur de tous leurs indicateurs est supérieure au 3e quartile).

Donc pour réaliser ça nous ajoutons une colonne qui a fait une moyennes des score des indicateurs de chaque pays et après nous faisons un tri par cette valeur pour afficher les 30 pays avec un fort potentiel



● ● ● Evolutions du potentiel



Méthode 1: Trouver les indicateurs qui ont été renseignés dans les années 2020 à 2100

Méthode 2: Utiliser la librairie Scipy pour extrapoler des NAN

Résultat:

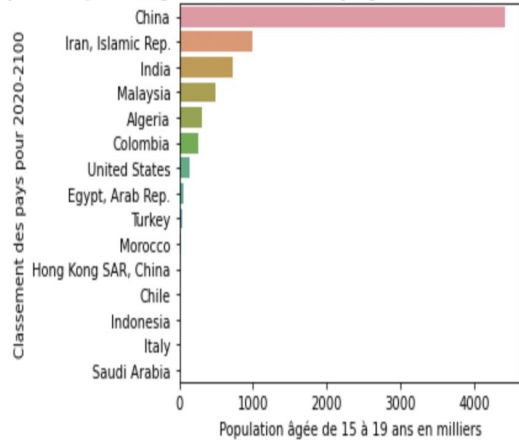
China', 'India', 'Malaysia', 'United States', 'Colombia', 'Egypt, Arab Rep.', 'Iran, Islamic Rep.', 'Algeria', 'Indonesia', 'Turkey', 'Hong Kong SAR, China', 'Morocco', 'Saudi Arabia', 'Chile', 'Italy'

Méthode 1:

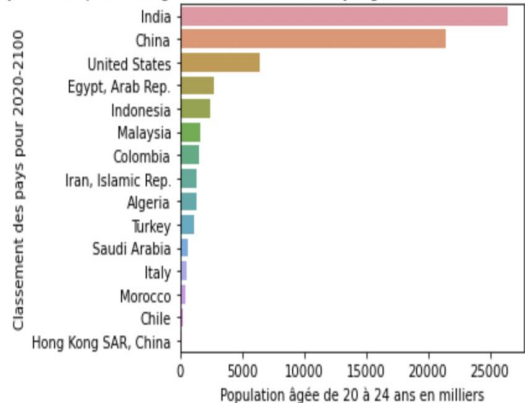
1. Trouver les indicateurs les mieux renseignés
2. Choisir 2 indicateurs :
 - Wittgenstein Projection: Population age 20-24 in thousands by highest level of educational attainment. Post Secondary. Total
 - Wittgenstein Projection: Population age 15-19 in thousands by highest level of educational attainment. Post Secondary. Total
3. Trouver les meilleurs pays pour chaque un des deux indicateur
4. Classement par quartiles un score entre 0 et 3 pour chaque indicateur
5. 10 top pays qui ont les meilleurs scores pour nos deux indicateurs

Scoring sur deux indicateurs prédictifs pour les années 2020-2100

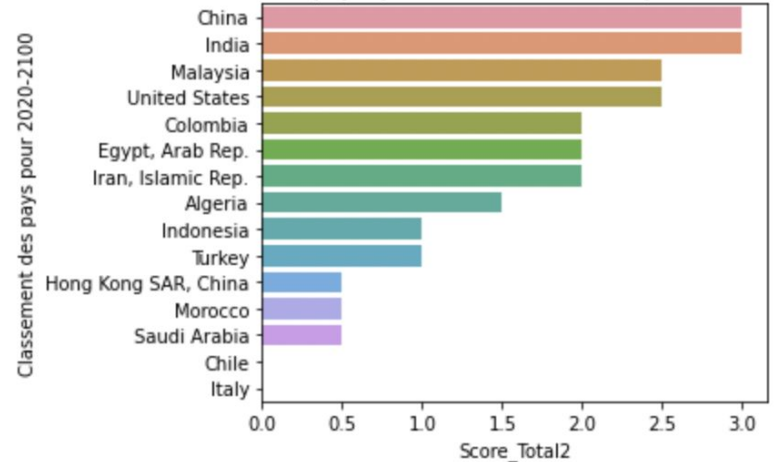
Wittgenstein Projection: Population age 15-19 in thousands by highest level of educational attainment. Post Secondary. Total



Wittgenstein Projection: Population age 20-24 in thousands by highest level of educational attainment. Post Secondary. Total



les meilleurs pays qui ont l'évolution de ce potentiel de clients



● ● ● Conclusion :

les pays commun dans les deux scoring précédents
pour les années 2010_2015 et 2020 _2100

les pays qui ont eu un fort potentiel pour devenir des
clients dans le passé et qui ont un fort potentiel pour
devenir des clients à l'avenir

**l'entreprise doit opérer en priorité dans 10 pays
suivants:**

- ★ United States
- ★ China
- ★ Iran, Islamic Rep.
- ★ Colombia
- ★ Algeria
- ★ Egypt, Arab Rep.
- ★ Indonesia
- ★ Turkey
- ★ Malaysia

● ● ● Pistes d'amélioration

- Création d'un dashboard interactif
- Visualisation du classement sur une carte interactive
- etc



Thank you

@ Mitra.dadgar1367@gmail.com

