PARTIE 2 : CRÉATION D'UNE BASE DE DONNÉES RELATIONNELLE 1. Import des librairies python et de l'ORM SQLAlchemy Definition ORM: Object-relational mapping (ORM, O/RM, and O/R mapping tool) in computer science is a programming technique for converting data between incompatible type systems using object-oriented programming languages. This creates, in effect, a "virtual object database" that can be used from within the programming language. 1. Import de deux fichiers csv contenant les travaux de la partie 1 2. Préparation des 4 tables demandées avec les colonnes demandées 3. Import de ces tables dans la base de données 4. Je les importe finalement dans ma bdd afin de pouvoir faire les requetes SQL demandées. In [1]: # Import des librairies import pandas as pd import numpy as np from sqlalchemy import create\_engine # permet de se connecter à la base de donnée from IPython.display import display # je cré l'engine avec SQLAlchemy. La bdd nommé P3 2.db sera creer automatiquement dans '/Users/ismail/Projet3' # selon la configuration de mon ordinateur (Macbook Air M1)(ref. voir https://getdoc.wiki/Sqlalchemy-dialects) bdd engine = create engine('sqlite://P3 2.db') In [2]: #visualization du chemin %pwd '/Users/ismail/Projet3' Out[2]: In [3]: # J'importe les données de la 1er partie pour creer les 3 premieres tables. df partiel = pd.read csv("data partie I.csv") # Je renomme les colonnes car SQL ne gere pas les lettres avec des accents ou caracteres sepciaux df partiel.rename(columns={ 'Zone' : 'zone', 'Code zone' : 'code zone', 'Produit' : 'produit', 'Population': 'population', 'Code Produit' : 'code produit', 'Disponibilité alimentaire en quantité (kg/personne/an)' : 'dispo alim kg', 'Disponibilité alimentaire (Kcal/personne/jour)' : 'dispo alim kcal an', 'Disponibilité de protéines en quantité (g/personne/jour)' : 'dispo prot', 'Disponibilité de matière grasse en quantité (g/personne/jour)' : 'dispo mat gr', 'Disponibilité intérieure' : 'dispo int', 'Aliments pour animaux' : 'alim ani', 'Semences' : 'semences', 'Pertes' : 'pertes', 'Traitement' : 'transfo', 'Nourriture' : 'nourriture', 'Autres utilisations (non alimentaire)' : 'autres\_utilisations', inplace=True) # Affichage du dataframe df partiel.head() Out[3]: code produit origine alim\_ani autres\_utilisations dispo\_alim\_kcal\_an dispo\_alim\_kg dispo\_mat\_gr ... **Production** pertes semences transfo produit **0** Afghanistan 2511 Blé vegetale NaN NaN 1369.0 160.23 4.69 ... 775000000.0 5.169000e+09 322000000.0 NaN 1 Afghanistan 2513 Orge vegetale 360000000.0 26.0 2.92 0.24 ... 52000000.0 5.140000e+08 NaN 22000000.0 NaN 2 Afghanistan 2514 Maïs vegetale 20000000.0 21.0 2.50 0.30 ... NaN 31000000.0 3.120000e+08 5000000.0 NaN 1000000.0 1.300000e+07 3 Afghanistan 2517 NaN 3.0 0.40 0.02 ... Millet vegetale NaN 0.0 NaN vegetale **4** Afghanistan NaN 15.0 7.53 0.04 ... 9000000.0 3.030000e+08 23000000.0 NaN NaN de Terre  $5 \text{ rows} \times 25 \text{ columns}$ In [4]: # Je verfie les noms des colonnes df partiel.columns Index(['zone', 'code zone', 'code produit', 'produit', 'origine', 'alim\_ani', Out[4]: 'autres\_utilisations', 'dispo\_alim\_kcal\_an', 'dispo\_alim\_kg', 'dispo\_mat\_gr', 'dispo\_prot', 'dispo\_int', 'Exportations - Quantité', 'Importations - Quantité', 'nourriture', 'pertes', 'Production', 'semences', 'transfo', 'Variation de stock', 'population', 'DispoAlimKcal', 'DispoProtKg', 'Ratio Kcal/KgNourriture', '%proteine'], dtype='object') In [5]: # Création du dataframe df population df population = df partiel[['zone','code zone','population']] # Affichage du dataframe df population.head() Out[5]: zone code zone population O Afghanistan 2 30552000 **1** Afghanistan 2 30552000 **2** Afghanistan 2 30552000 **3** Afghanistan 2 30552000 **4** Afghanistan 2 30552000 #seule ligne par pays df\_population = df\_population.drop\_duplicates(subset=['code zone'],inplace=False) df\_population.head() Out[6]: zone code zone population Afghanistan 30552000 **51** Afrique du Sud 202 52776000 135 Albanie 3173000 4 39208000 200 Algérie 276 82727000 Allemagne # Clef primaire proposée pour la table df\_population: code zone df population['primary key'] = df population['code zone'] In [ ]: # Création du dataframe df dispo alim df\_dispo\_alim = df\_partie1[['zone', 'code zone', 'produit', 'code produit', 'origine', 'dispo\_prot','dispo\_alim\_kcal\_an','dispo\_mat\_gr','DispoAlimKcal', # Affichage du dataframe display(df\_dispo\_alim) produit code produit origine dispo\_prot dispo\_alim\_kcal\_an dispo\_mat\_gr DispoAlimKcal zone code zone DispoProtKg 2 Blé 1.526638e+13 **0** Afghanistan 2511 vegetale 36.91 1369.0 4.116011e+08 2 0.79 26.0 **1** Afghanistan 2513 vegetale 2.899385e+11 8.809669e+06 Orge 2 0.56 **2** Afghanistan 2514 vegetale 21.0 0.30 2.341811e+11 6.244829e+06 Maïs 2 3.0 **3** Afghanistan Millet 2517 vegetale 0.08 3.345444e+10 8.921184e+05 **4** Afghanistan 2 Pommes de Terre 2531 vegetale 0.25 15.0 1.672722e+11 2.787870e+06 11582 Îles Salomon 5.57 7.371540e+09 25 Poissons Marins, Autres 2764 animale 36.0 1.140541e+06 11583 Îles Salomon 0.0 25 2767 animale 0.06 0.000000e+00 1.228590e+04 Mollusques, Autres 11584 Îles Salomon 25 623.0 Riz (Eq Blanchi) 2805 vegetale 10.90 1.275686e+11 2.231938e+06 **11585** Îles Salomon 25 Lait - Excl Beurre animale 1.05 19.0 0.70 3.890535e+09 2.150033e+05 **11586** Îles Salomon 25 0.10 Miscellanees 2899 vegetale 3.0 0.04 6.142950e+08 2.047650e+04 11587 rows × 10 columns #Clef primaire proposée pour la table df\_dispo\_alim: combinaison de code zone + code produit. with pd.option context('mode.chained assignment', None): df\_dispo\_alim['primary\_key'] = df\_dispo\_alim['code zone'].astype(str) + '-' + df\_dispo\_alim['code produit'].astype(str) pd.reset\_option('mode.chained\_assignment') In [10]: # Création du dataframe df equilibre prod df\_equilibre\_prod = df\_partiel[['zone', 'code zone', 'produit', 'code produit', 'dispo\_int', 'alim\_ani', 'semences', 'pertes', 'transfo', 'nourriture', 'a # Affichage du dataframe display(df\_equilibre\_prod) zone code zone produit code produit dispo\_int alim\_ani semences pertes transfo nourriture autres\_utilisations NaN 322000000.0 775000000.0 **0** Afghanistan 2 Blé 2511 5.992000e+09 NaN 4.895000e+09 NaN 2 2513 5.240000e+08 360000000.0 **1** Afghanistan 22000000.0 52000000.0 8.900000e+07 NaN Orge 2 2514 3.130000e+08 200000000.0 **2** Afghanistan Maïs 5000000.0 31000000.0 7.600000e+07 NaN **3** Afghanistan 2 Millet 2517 1.300000e+07 NaN 0.0 1000000.0 1.200000e+07 NaN 4 Afghanistan 2 Pommes de Terre 23000000.0 2531 2.620000e+08 NaN 9000000.0 NaN 2.300000e+08 NaN ••• 11582 Îles Salomon 1.100000e+07 1.100000e+07 25 Poissons Marins, Autres 2764 NaN NaN NaN NaN 0.0 11583 Îles Salomon 25 Mollusques, Autres 1.000000e+06 NaN NaN NaN NaN 1.000000e+06 0.0 11584 Îles Salomon 25 Riz (Eq Blanchi) 2805 4.900000e+07 0.0 0.0 1000000.0 NaN 3.600000e+07 12000000.0 11585 Îles Salomon 25 Lait - Excl Beurre 2848 7.000000e+06 0.0 NaN 0.0 NaN 6.000000e+06 NaN **11586** Îles Salomon 25 Miscellanees 2899 NaN NaN NaN NaN NaN NaN NaN 11587 rows × 11 columns In [11]: #Clef primaire proposée pour la table df equilibre prod: code zone + code produit. with pd.option context('mode.chained assignment', None): df\_equilibre\_prod['primary\_key'] = df\_equilibre\_prod['code zone'].astype(str) + '-' + df\_equilibre\_prod['code produit'].astype(str) pd.reset option('mode.chained assignment') df sousnut = pd.read csv("sousnut partie I.csv") # Création du dataframe df sous nutrition df sousnut1 = df sousnut[['Zone', 'Code zone', 'sousalimentés', 'Population']] df sousnut1.rename(columns={ 'Zone' : 'zone', 'Code zone' : 'code zone', 'Produit' : 'produit', 'sousalimentés' : 'sousalimentes', 'Population': 'population', inplace=True) # Affichage du dataframe display(df sousnut1) /opt/anaconda3/lib/python3.8/site-packages/pandas/core/frame.py:4296: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy return super().rename( zone code zone sousalimentes population **0** Afghanistan 7900000.0 30552000 7900000.0 30552000 **1** Afghanistan **2** Afghanistan 2 7900000.0 30552000 2 7900000.0 **3** Afghanistan 30552000 2 7900000.0 30552000 **4** Afghanistan **15600** Îles Salomon 25 50000.0 561000 **15601** Îles Salomon 25 50000.0 561000 **15602** Îles Salomon 25 50000.0 561000 15603 Îles Salomon 25 561000 50000.0 **15604** Îles Salomon 25 50000.0 561000  $15605 \text{ rows} \times 4 \text{ columns}$ In [13]: #Pour eviter d'avoir des lignes completement identiques pour chaque pays je supprime les doublons pour avoir une #seule ligne par pays df sousnut1 = df sousnut1.drop duplicates(subset=['code zone'],inplace=False) df\_sousnut1.head() zone code zone sousalimentes population Out[13]: 2 Afghanistan 7900000.0 30552000 **60** Afrique du Sud 202 2600000.0 52776000 3 155 Albanie 200000.0 3173000 247 Algérie 4 1700000.0 39208000 79 82727000 340 Allemagne 0.0 In [14]: #Clef primaire proposée pour la table df sousnut1: code zone df sousnut1['primary key'] = df sousnut1['code zone'] In [15]: #Transfer des tables dans la bdd df population.to sql( name='population', # nom de la table apres transfer dans la bdd con = bdd engine, # connexion utilisé (ici, celle donnée par SQLAlchemy) index=False, # on ne veut pas écrire les index dans la table if\_exists='replace' # si la table existe déjà, on l'écrase df\_dispo\_alim.to\_sql( name='dispo alim', con = bdd engine, index=False, if\_exists='replace' df equilibre prod.to sql( name='equilibre prod', con = bdd\_engine, index=False, if exists='replace' df sousnut1.to sql( name='sous\_nutrition', con = bdd engine, index=False, if\_exists='replace' #Question 19 : Écrivez les requêtes SQL permettant de connaître... #Les 10 pays ayant le plus haut ratio disponibilité alimentaire/habitant en termes de protéines (en kg) #par habitant, puis en termes de kcal par habitant. #I. en kg top10 dispo prot = pd.read sql query(con=bdd engine, sql = f""" SELECT sum(dispo prot) as dispoprot FROM dispo alim GROUP BY zone ORDER BY dispoprot DESC LIMIT 10""") top10 dispo prot.index += 1 print('Top 10 des pays avec le plus de disponibilité alimentaire par habitant et par an (en kg de protéines)') display(top10 dispo prot) Top 10 des pays avec le plus de disponibilité alimentaire par habitant et par an (en kg de protéines) zone dispoprot 130.91 1 Islande 2 Chine - RAS de Hong-Kong 129.01 3 127.98 Israël Lituanie 124.17 5 Maldives 119.16 Finlande 117.41 Luxembourg Pays-Bas 111.45 9 Albanie 111.14 10 Portugal 110.86 In [17]: #II. en Kcal top10\_dispo\_kcal = pd.read\_sql\_query(con=bdd\_engine, sql = f""" SELECT sum(dispo\_alim\_kcal\_an) dispoalimkcal FROM dispo alim GROUP BY zone ORDER BY dispoalimkcal DESC LIMIT 10""") top10\_dispo\_kcal.index += 1 print('Top 10 des pays avec le plus de disponibilité alimentaire par habitant et par an (en kcal)') display(top10\_dispo\_kcal) Top 10 des pays avec le plus de disponibilité alimentaire par habitant et par an (en kcal) zone dispoalimkcal 3769.0 1 Autriche 2 Belgique 3737.0 3 3708.0 Turquie **4** États-Unis d'Amérique 3682.0 5 3610.0 Israël 6 Irlande 3596.0 7 Italie 3578.0 8 Égypte 3518.0 9 3503.0 Allemagne 10 Canada 3499.0 In [18]: # les 10 pays ayant le plus faible ratio disponibilité alimentaire/habitant #en termes de protéines (en kg) par habitant. top10\_faible\_dispo\_prot = pd.read\_sql\_query(con=bdd\_engine, sql = f""" SELECT zone, sum(dispo\_prot) as dispoprot FROM dispo alim GROUP BY zone ORDER BY dispoprot ASC LIMIT 10""") top10\_faible\_dispo\_prot.index += 1 print('Top 10 des pays avec le moins de disponibilité alimentaire par habitant et par an (en kg de protéines)') display(top10\_faible\_dispo\_prot) Top 10 des pays avec le moins de disponibilité alimentaire par habitant et par an (en kg de protéines) zone dispoprot 37.62 1 Libéria 2 Guinée-Bissau 43.91 3 Mozambique 45.67 4 République centrafricaine 46.03 5 46.68 Madagascar 47.69 6 Haïti 7 Zimbabwe 48.32 8 Congo 51.30 9 Sao Tomé-et-Principe 51.65 10 Ouganda 52.64 In [19]: #La quantité totale (en kg) de produits perdus par pays pertes par pays = pd.read sql query(con=bdd engine, sql = f""" SELECT zone, sum(pertes) FROM equilibre\_prod GROUP BY zone ORDER BY zone""") print('Quantité des pertes par pays (en kg)') display(pertes\_par\_pays) Quantité des pertes par pays (en kg) sum(pertes) zone 0 Afghanistan 1.135000e+09Afrique du Sud 2.193000e+09 1 2 Albanie 2.750000e+08 3 Algérie 3.747000e+09 Allemagne 3.781000e+09 4 ••• 169 Émirats arabes unis 7.050000e+08 170 Équateur 7.040000e+08 **171** États-Unis d'Amérique 7.162000e+09 172 Éthiopie 2.256000e+09 173 Îles Salomon 6.000000e+06 174 rows × 2 columns In [20]: #Les 10 pays pour lesquels la proportion de personnes sous-alimentées est la plus forte top10\_pays\_sous\_nutrition = pd.read\_sql\_query(con=bdd\_engine, sql=f""" SELECT zone, ROUND(AVG((sousalimentes / population) \* 100),2) AS prop\_p\_ss\_alim FROM sous nutrition NATURAL JOIN sous nutrition GROUP BY zone ORDER BY prop\_p\_ss\_alim DESC LIMIT 10""") top10\_pays\_sous\_nutrition.index += 1 print('Top 10 des pays avec la proportion de personnes sous-alimentées la plus importante (en %)') display(top10\_pays\_sous\_nutrition) Top 10 des pays avec la proportion de personnes sous-alimentées la plus importante (en %) zone prop\_p\_ss\_alim 1 Dominique 69.44 2 Haïti 50.40 3 Kiribati 49.02 48.15 4 Zambie 5 Zimbabwe 46.64 Saint-Vincent-et-les Grenadines 6 45.87 7 République centrafricaine 43.33 8 République populaire démocratique de Corée 42.58 9 Congo 40.47 10 Tchad 38.21 In [21]: #Les 10 produits pour lesquels le ratio Autres utilisations/Disponibilité intérieure est le plus élevé. top10\_produit\_autres\_util = pd.read\_sql\_query(con=bdd\_engine, sql=f""" SELECT produit AS Produit, ROUND(AVG(autres utilisations / dispo int)\*100,2) AS Part autres utilisations FROM equilibre prod GROUP BY Produit ORDER BY Part autres utilisations DESC LIMIT 10""") top10 produit autres util.index += 1 print('Top 10 des produits pour lesquels le ratio Autres utilisations/Disponibilité intérieure est le plus élevé') display(top10\_produit\_autres\_util) Top 10 des produits pour lesquels le ratio Autres utilisations/Disponibilité intérieure est le plus élevé **Produit Part\_autres\_utilisations** Alcool, non Comestible 98.25 2 73.91 **Piments** 3 Huile de Palmistes 72.65 4 Huile de Palme 62.17 5 **Palmistes** 57.56 6 Huile de Colza&Moutarde 57.36 7 Huile de Son de Riz 50.30 **Huil Plantes Oleif Autr** 8 49.66 9 Huile de Coco 48.99 10 Sucre non centrifugé 46.51 #Question 20: pour quelques uns des produits identifiés dans cette dernière requête SQL, supposez quelles sont ces "autres utilisations" possibles (recherchez sur internet!) 1. Piments: Utilisé comme parfum ou insecticide: https://fr.wikipedia.org/wiki/Piment#Parfums\_et\_cosmétiques https://fr.wikipedia.org/wiki/Piment#Comme\_insecticide 2. Alcool, non comestible: https://ec.europa.eu/taxation\_customs/denatured-alcohol-not-human-consumption\_fr Alcool dénaturé non destiné à la consommation humaine L'alcool peut être rendu impropre à la consommation humaine par l'ajout de produits ayant très mauvais goût et/ou très mauvaise odeur et (de préférence) d'un marqueur chimique. Le procédé permettant de rendre l'alcool impropre à la consommation humaine est appelé dénaturation (totale ou partielle) de l'alcool. Lorsqu'il a été dénaturé conformément à la réglementation de l'UE, l'alcool peut être exonéré de droits d'accise. Cette règle s'applique

notamment à l'alcool utilisé à des fins industrielles, pour produire des denrées alimentaires et des médicaments ou pour fabriquer des cosmétiques, des biocarburants, des peintures, ainsi que des

produits antigel et de nettoyage. 3. Huile de Palme: pour la production des cosmetiques et des agrocarburants https://fr.wikipedia.org/wiki/Huile\_de\_palme#Oléochimie

https://fr.wikipedia.org/wiki/Huile\_de\_palme#Agro-carburants