

Projet 3 : Concevez une application au service de la santé publique

Mentor : Morgan MOISON
Etudiant : Marin DUCHEMIN



Plan de la Présentation

I ~ Présentation de l'idée d'application

II ~ Préparation des données

III ~ Analyse exploratoire effectuée

IV ~ Faisabilité de l'application

I ~ Présentation de l'idée d'application

Générer un panier de provisions alimentaires pour une semaine.

Basée sur l'idée du «Panier Confinement » de certaines grandes enseignes commerciales

I ~ Présentation de l'idée d'application

Proposer un choix de types de panier différents selon :

- Ses préférences
- Sa localisation

**Basée sur la base de données
OpenFoodFact**



I ~ Présentation de l'idée d'application

Exemple de panier de confinement (Carrefour)

2 briques de Jus d'orange avec pulpe pur jus, Carrefour, 1L

2 paquets de Biscottes froment, Cracotte, 250g

1 paquet de Céréales trésor chocolat noisette, Kellogg's, 375g

1 pot de Confiture de fraise, Carrefour, 370g

1 pot de Pâte à tartiner noisettes et cacao, Nutella, 400g

1 sachet de Café soluble, Nescafe, 100g

1 boîte de Thé vert menthe, Carrefour, 25 sachets

...

Définition des catégories suite à cet exemple

Jus de Fruits, Tartines, Céréales, Laits, Confitures et Pâtes à Tartiner, Cafés, Thés, Biscuits Sucrés, Desserts, Soupes, Pates, Riz, Fruits, Légumes, Sources de Protéines, Sauces, Biscuits Apéritifs, Matières Grasses

I ~ Présentation de l'idée d'application

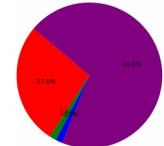
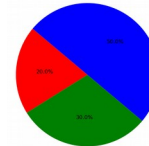
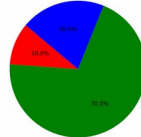
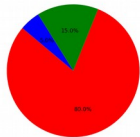
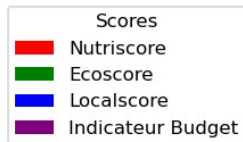
Construction des Profils :

- Choix des indicateurs d'intérêt :
 - Récupérés dans la base de données :
 - le Nutriscore (allant de -15 à +40, classé en lettre de A à E),
 - L'Ecoscore (allant de 0 à +100, classé en lettre de A à E),
 - Créés par la suite :
 - un Localscore (qui sera créé en fonction de l'emplacement du consommateur),
 - un Indicateur Budget (un filtre selon les marques de distribution pour sélectionner les produits moins onéreux)

I ~ Présentation de l'idée d'application

Construction des Profils :

- Choix de différents profils à partir de nos indicateurs correspondant à différents besoins (nutritionnels, écologiques, régionaux ou économiques)



II ~ Préparation des données

Jeu de données très conséquent impliquant :

- Importation au sein de la base de données NoSQL MongoDB
- Filtres successifs pour réduire aux produits d'intérêt :
 - Au départ : 1 601 557 produits.
 - Appel de Santé Publique France donc réduction aux produits vendus en France : 425 015 (*filtre de requête MongoDB*)
 - Pour le localscore, besoin de connaître le lieu de manufacture : 193 747 (*filtre de requête MongoDB*)
 - Produits avec données utilisables : 44 214 (*filtre RegEx sur les noms, lieux de manufacture, codes barres*)

II ~ Préparation des données

Plusieurs étapes d'uniformisation des données

- Noms des produits
- Noms des lieux
- etc

```
1 # Uniformisation de la colonne ecoscore_data
2
3 for i in range(len(name_df)):
4     if type(name_df.loc[i,'ecoscore_data']) == dict:
5         if name_df.loc[i,'ecoscore_data']['agribalyse']:
6             name_df.loc[i,'ecoscore_data'] = name_df.loc[i,'ecoscore_data']['agribalyse']['name_fr']
7         else:
8             name_df.loc[i,'ecoscore_data'] = np.nan
9     else:
10        continue
```

```
1 # Uniformisation des noms
2
3 for i in range(len(name_df)) :
4     if type(name_df.loc[i,'generic_name_fr']) == str and len(name_df.loc[i,'generic_name_fr']) == 0:
5         name_df.loc[i,'generic_name_fr'] = np.nan
6     elif type(name_df.loc[i,'generic_name']) == str and len(name_df.loc[i,'generic_name']) == 0:
7         name_df.loc[i,'generic_name'] = np.nan
8     elif type(name_df.loc[i,'product_name_fr']) == str and len(name_df.loc[i,'product_name_fr']) == 0:
9         name_df.loc[i,'product_name_fr'] = np.nan
10    elif type(name_df.loc[i,'product_name']) == str and len(name_df.loc[i,'product_name']) == 0:
11        name_df.loc[i,'product_name'] = np.nan
12    elif type(name_df.loc[i,'ecoscore_data']) == str and len(name_df.loc[i,'ecoscore_data']) == 0:
13        name_df.loc[i,'ecoscore_data'] = np.nan
14    else:
15        continue
```

II ~ Préparation des données

Enrichissement avec des coordonnées GPS :

La fréquence des 20 premiers lieux est :

| | |
|-------------------|----------|
| France | 0.386620 |
| Italie | 0.041865 |
| Belgique | 0.020871 |
| Espagne | 0.018331 |
| Allemagne | 0.017599 |
| France, Bretagne | 0.013576 |
| Bretagne, France | 0.010934 |
| France | 0.010934 |
| Pays-Bas | 0.007926 |
| Union Européenne | 0.006320 |
| Suisse | 0.005426 |
| Thaïlande | 0.004349 |
| Chine | 0.003780 |
| Normandie, France | 0.003374 |
| Royaume-Uni | 0.003353 |
| Bretagne | 0.003048 |
| Guérou, France | 0.002906 |
| Pologne | 0.002886 |
| Maroc | 0.002744 |
| Bretagne | 0.002703 |

Première étape de normalisation des lieux de manufacture

(n'avoir que des Pays et, pour la France, des Régions, Villes ou Code Postaux)

Recherche des coordonnées via des bases de données disponibles sur le site du gouvernement et Kaggle

- <https://www.data.gouv.fr/fr/datasets/regions-departements-villes-et-villages-de-france-et-doutre-mer/>
- <https://www.kaggle.com/paultimothymooney/latitude-and-longitude-for-every-country-and-state>

II ~ Préparation des données

Etape de standardisation des lieux de manufacture :

```
1 # Actualisation du tableau
2
3 for i in range(len(coordonnees_df)):
4     # Analyse de la cellule
5     if coordonnees_df.loc[i, 'Pays'] in ['france'] and len(coordonnees_df.loc[i, 'Pays']) == 1:
6         coordonnees_df.loc[i, 'Pays'] = 'FRANCE'
7     else:
8         test = [j for j in coordonnees_df.loc[i, 'Pays'] if(j not in ['france'])]
9
10    # Tri des Pays
11    pays = [k for k in paysList if (k in test)]
12    pays = list(set(s for s in pays if len(s) == max(len(s) for s in pays)))
13    if len(pays) != 0:
14        coordonnees_df.loc[i, 'Pays'] = pays[0].upper()
15
16    # Tri des Regions
17    else:
18        region = [k for k in regionList if (k in test)]
19        region = list(set(s for s in region if len(s) == max(len(s) for s in region)))
20        if len(region) != 0:
21            coordonnees_df.loc[i, 'Pays'] = 'FRANCE'
22            coordonnees_df.loc[i, 'Region'] = region[0]
23
24    # Tri des Codes Postaux
25    else:
26        codePostal = [k for k in codePostalList if (k in test)]
27        codePostal = list(set(s for s in codePostal if len(s) == max(len(s) for s in codePostal)))
28        if len(codePostal) != 0:
29            coordonnees_df.loc[i, 'Pays'] = 'FRANCE'
30            coordonnees_df.loc[i, 'Region'] = codePostal[0]
31
32    # Tri des Villes
33    else:
34        ville = [k for k in villeList if (k in test)]
35        ville = list(set(s for s in ville if len(s) == max(len(s) for s in ville)))
36        if len(ville) != 0:
37            coordonnees_df.loc[i, 'Pays'] = 'FRANCE'
38            coordonnees_df.loc[i, 'Region'] = ville[0]
39        else:
40            if np.vectorize(lambda x: bool(re.compile(r".*france.*").match(x)))(coordonnees_df.loc[i, 'Pays']).any():
41                coordonnees_df.loc[i, 'Pays'] = 'FRANCE'
```

Suivie d'une étape d'uniformisation pour n'avoir que des pays ou des regions

II ~ Préparation des données

Comment associer les catégories aux données:

- Sélection de mots-clés correspondant aux catégories définies (*diapo 5*)

```
1 # On liste le type d'ingrédients pouvant aller dans le panier de confinement
2
3 jusList = ['jus d'orange', 'jus de pomme', 'jus multifruit', 'jus d'ananas', 'jus de goyave']
4 tartineList = ['pain de mie', 'biscotte', 'pain tranche']
5 cerealeList = ['cereales', 'petales de ble', 'petales', 'muesli', 'petales']
6 laitList = ['lait entier', 'lait demi-ecreme', 'lait-ecreme', 'lait de coco', 'lait d'amande', 'lait de riz']
7 tartinableList = ['confiture', 'pate a tartiner']
8 cafeList = ['cafe']
9 theList = ['the']
10 biscuitList = ['cookie', 'biscuit sec chocolate', 'biscuit sec petit beurre', 'biscuit sec feuilleté', 'sable']
11 dessertList = ['compote', 'yaourt', 'dessert']
12 soupeList = ['soupe']
13 pateList = ['pates', 'pates fraiches', 'pates seches']
14 rizList = ['riz basmati', 'riz thai', 'riz gluant', 'riz melange', 'riz complet', 'riz blanc']
15 fruitList = ['orange', 'pomme, crue', 'poire', 'banane, crue', 'ananas, cru', 'peche', 'abricot', 'cerise', 'framboise', 'mure', 'prune']
16 legumeList = ['carotte', 'patate', 'oignon', 'haricot', 'maïs', 'champignon', 'tomate', 'concombre', 'epinard']
17 proteineList = ['poulet', 'boeuf', 'jambon', 'sardine', 'thon', 'colin', 'oeuf']
18 sauceList = ['pesto', 'sauce tomate', 'ketchup', 'mayonnaise']
19 apertifList = ['cacahuete', 'pistache', 'chips', 'noix de cajou']
20 matiereList = ['beurre doux', 'beurre demi-sel', 'beurre demi sel', 'beurre sale', 'creme fraiche', 'huile d'olive', 'huile de colza', 'huile de tournesol']
21
```

- Recherche de ces mots-clés parmi les noms des produits

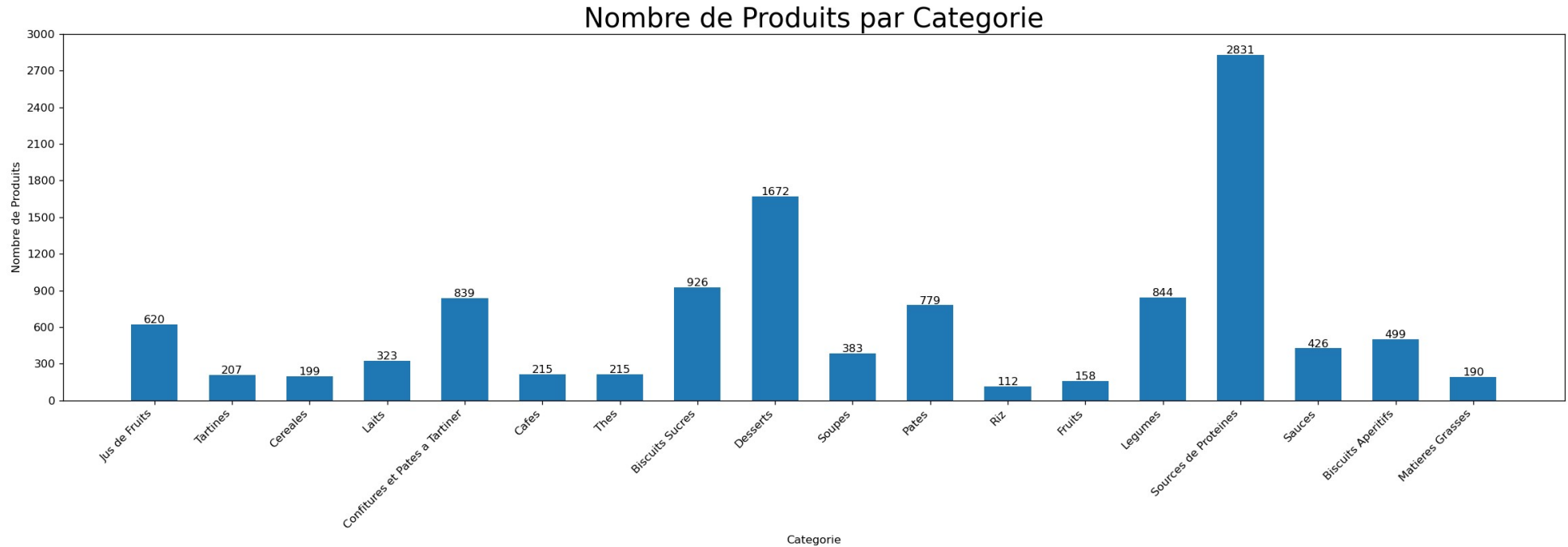
II ~ Préparation des données

- Construction des catégories en jeux de données :

| | Code Barre | Nom du Produit | Pays de Manufacture | Region de Manufacture | Latitude GPS | Longitude GPS | Nutriscore | Ecoscore | Marque de Distribution |
|---|--------------|----------------------------|---------------------|-----------------------|--------------|---------------|------------|-----------|------------------------|
| 0 | 25110001525 | Huile d'olive vierge extra | Italie | NaN | 41.871940 | 12.567380 | 6.0 | 49.531594 | Non |
| 1 | 319110413922 | Huile d'olive vierge extra | France | NaN | NaN | NaN | NaN | 65.766726 | Non |
| 2 | 637913429156 | Huile de colza | France | Hauts-de-France | 50.614407 | 3.118079 | NaN | 47.450168 | Non |
| 3 | 764549058679 | Huile d'olive vierge extra | Espagne | NaN | 40.463667 | -3.749220 | 6.0 | 64.472459 | Non |
| 4 | 11972464 | Huile d'olive vierge extra | Italie | NaN | 41.871940 | 12.567380 | NaN | 41.481594 | Non |

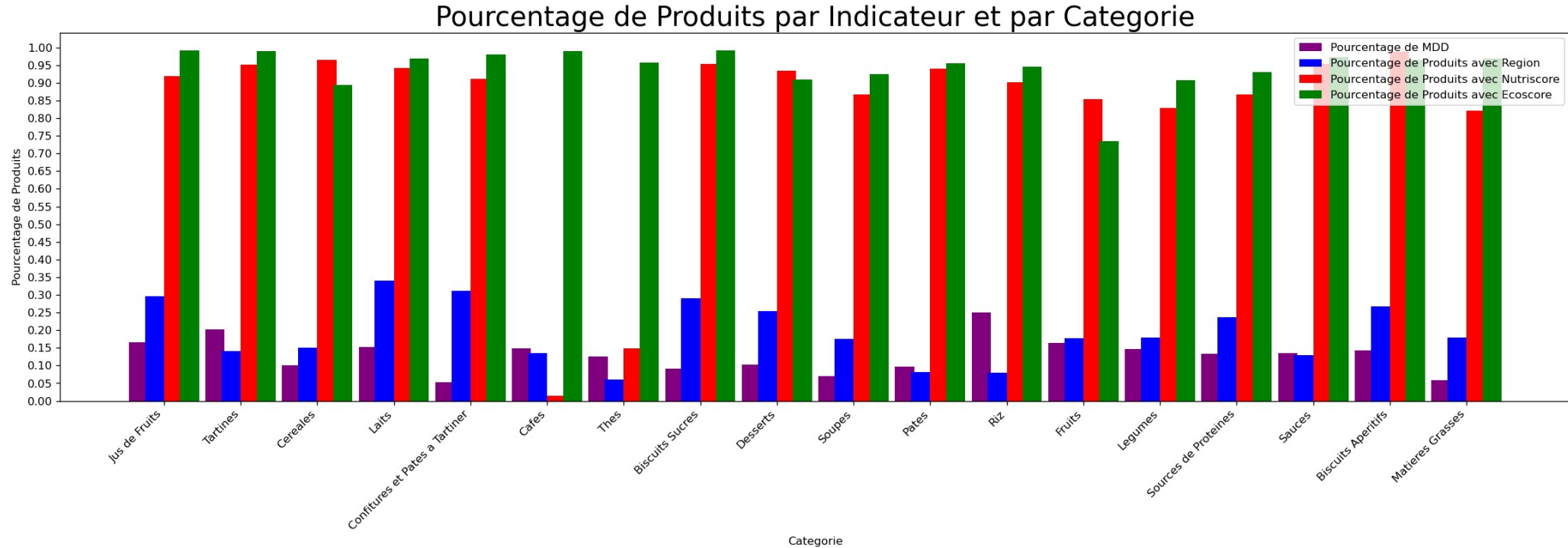
- Enregistrement sous format csv pour les réutiliser dans les autres NoteBooks.

III ~ Analyse exploratoire effectuée

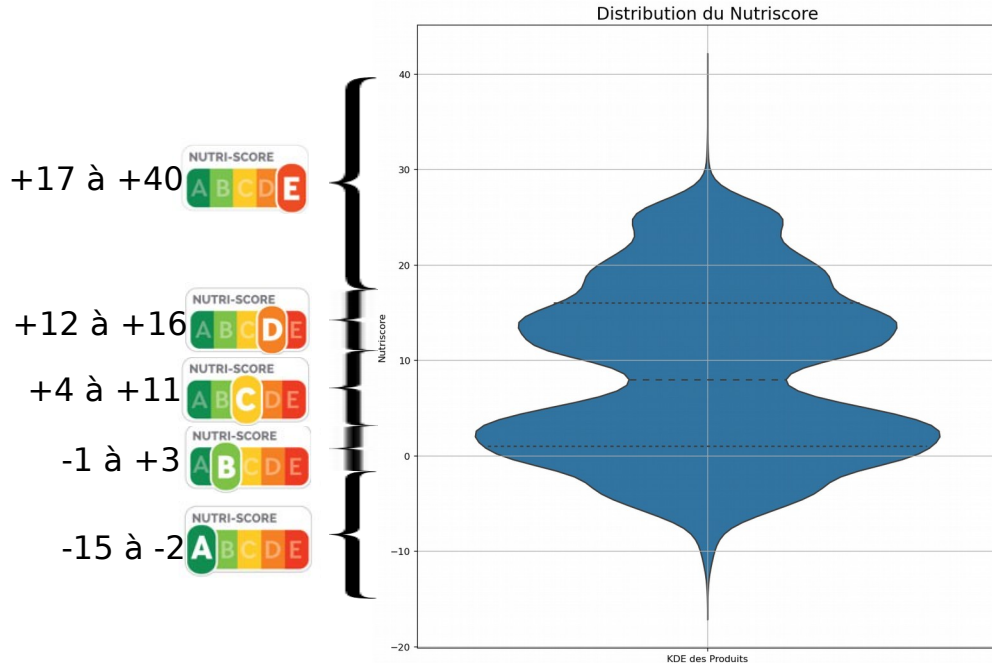


Nombre total de produits conservés : 11 438

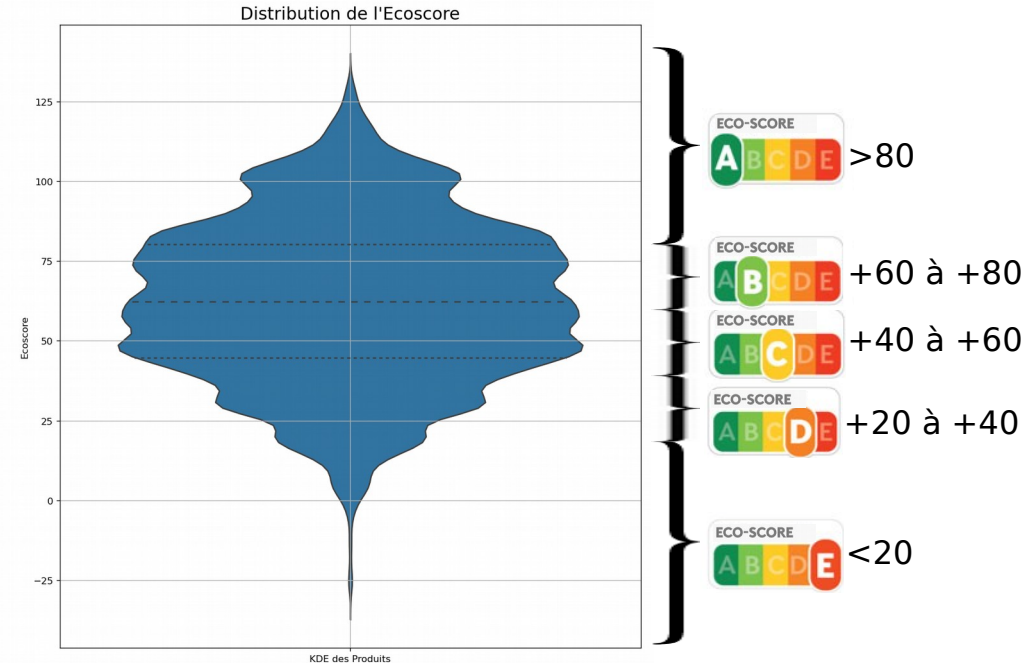
III ~ Analyse exploratoire effectuée



III ~ Analyse exploratoire effectuée



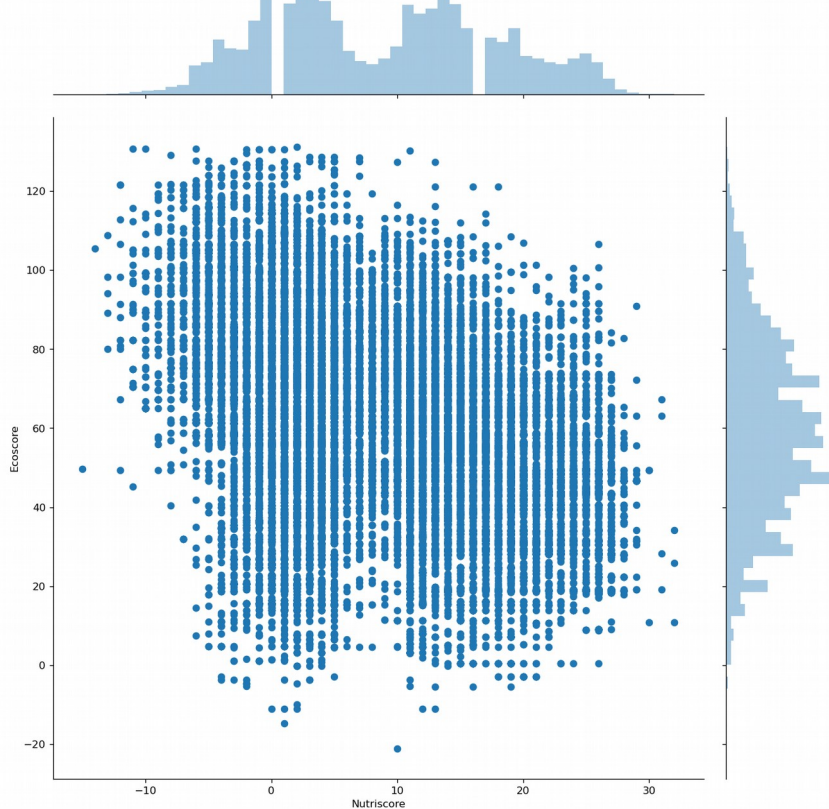
Distribution bimodale
du Nutriscore



Distribution monomodale
de l'Ecoscore

III ~ Analyse exploratoire effectuée

Recherche de corrélation entre le Nutriscore et l'Ecoscore



**Aucune corrélation visible
entre ces deux scores.**

III ~ Analyse exploratoire effectuée

Synthèse après la partie exploratoire:

- Données suffisantes et uniformes pour exploiter les indicateurs
- Catégories établies, créées et fonctionnelles
- Indicateurs récupérés non-corrélés et exploitables et
Indicateurs non récupérés prêt à être dynamiquement créés

Le projet d'application paraît donc faisable.

IV ~ Faisabilité de l'application

Mise en pratique pour tester un prototype d'application.

IV ~ Faisabilité de l'application

Fonction pour trouver le score de proximité

```
1 # Fonction pour trouver la region la plus lointaine (DOM TOM exclus)
2
3 def find_further_region(region):
4     distanceList = []
5
6     for item in regionList:
7         distLatCarre = np.square(latitudeDict[item] - latitudeDict[region])
8         distLongCarre = np.square(longitudeDict[item] - longitudeDict[region])
9         distanceList.append(np.sqrt(distLatCarre + distLongCarre))
10
11     proximiteDict = dict(zip(regionList, distanceList))
12
13     global furtherRegion
14     furtherRegion = (list(proximiteDict.keys())[list(proximiteDict.values()).index(max(proximiteDict.values()))], max(proximiteDict.values()))
15     return furtherRegion

```

```
1 # Fonction pour construire le localscore en fonction de la region ou on se trouve
2
3 def find_localscore(region, categorie):
4     categorie['Localscore'] = ''
5     for i in range(len(categorie)):
6         # Si on a l'information de la region ou l'information du pays (autre que la France)
7         if type(categorie.loc[i, 'Region de Manufacture']) == str or categorie.loc[i, 'Pays de Manufacture'] not in ['France']:
8             distLatCarre = np.square(categorie.loc[i, 'Latitude GPS'] - latitudeDict[region])
9             distLongCarre = np.square(categorie.loc[i, 'Longitude GPS'] - longitudeDict[region])
10             categorie.loc[i, 'Localscore'] = (np.sqrt(distLatCarre + distLongCarre))
11
12         # Si le produit est manufacture en France
13         else:
14             categorie.loc[i, 'Localscore'] = find_further_region(region)[1]
15

```

IV ~ Faisabilité de l'application

Classement des produits selon le profil choisi

```
1 # Fonction pour classer les produits selon le profil
2
3 def find_rank(profil, categorie):
4     profilRatio = [profilDict[profil]['Nutriscore'], profilDict[profil]['Ecoscore'], profilDict[profil]['Localscore']]
5
6     # Pour ne pas avoir de nutriscore negatif ce qui nous embeterait pour le calcul du ratio
7     categorie['Nutriscore'] = categorie['Nutriscore'].add(abs(categorie['Nutriscore'].min()))
8
9     # Creation du nouveau score global
10    for i in range(len(categorie)):
11        if categorie.loc[i, 'Nutriscore'] == 0:
12            scoreNut = 1
13        else:
14            scoreNut = categorie['Nutriscore'].min() / categorie.loc[i, 'Nutriscore']
15
16        scoreEco = categorie.loc[i, 'Ecoscore'] / categorie['Ecoscore'].max()
17
18        if categorie.loc[i, 'Localscore'] == 0:
19            scoreLoc = 1
20        else:
21            scoreLoc = categorie['Localscore'].min() / categorie.loc[i, 'Localscore']
22
23        categorie.loc[i, 'Score pour le Profil'] = np.average([scoreNut, scoreEco, scoreLoc], weights=profilRatio)
24
25
26 # Reorganisation du jeu de donnees
27 if profil in ['Petit Prix']:
28     categorie.sort_values(['Marque de Distribution', 'Score pour le Profil'], ascending = False, ignore_index = True, inplace = True)
29 else:
30     categorie.sort_values(['Score pour le Profil'], ascending = False, ignore_index = True, inplace = True)
```

IV ~ Faisabilité de l'application

Création d'un panier selon le profil et la région

```
1  # Fonction pour faire un panier
2
3  def create_basket(region, profil):
4
5      for cat in categ_df_list:
6          find_localscore(region, cat)
7          find_rank(profil, cat)
8
9      global basket_df
10     basket_df = pd.DataFrame(columns = list(categ_df_list[0].columns))
11
12     for i in range(len(categ_df_list)):
13         basket_df.loc[i] = dict(categ_df_list[i].loc[0])
14
15
```

IV ~ Faisabilité de l'application

Test en direct parce que c'est plus sympa.

**Merci de votre
Attention !**