
Collecte de Données et Data Visualization

Compréhension des objectifs, présentation de la démarche

Plan

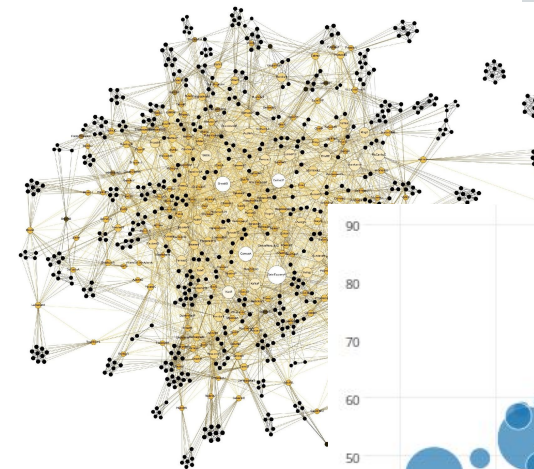
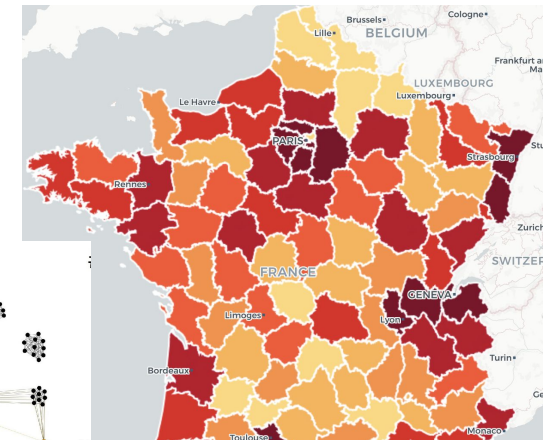
- Contexte & Définitions
- Objectifs
- Démarche & livrables
- Planning
- Organisation

Contexte & Définitions

Data Visualization : représentation visuelle des données pour une **meilleure compréhension**

- Graphiques
- Cartographies
- Listes / Chronologie
- Organigramme

Revenu fiscal de référence par tranche (en euros)	Nombre de foyers fiscaux	Revenu fiscal de référence des foyers fiscaux	Impôt net (total)	Nombre de foyers fiscaux imposés	Revenu fiscal de référence des foyers fiscaux imposés
0 à 10 000	8 779 578	37 017 353	-120 471	64 840	288 206
10 001 à 12 000	1 141 466	23 577 329	-52 668	5 831	63 710
12 001 à 15 000	3 415 487	48 459 161	-97 920	236 299	3 491 006
15 001 à 20 000	907 523	102 764 312	1 757 683	130 770	63 996 113
20 001 à 30 000	6 830 393	167 947 233	5 647 709	507 008	98 705 054
30 001 à 50 000	6 553 966	250 560 654	13 309 362	513 297	109 437 842
50 001 à 100 000	3 305 940	217 391 802	20 630 441	131 015	205 724 469
Plus de 100 000 dont :	740 153	140 578 578	28 027 762	72 817	38 965 445
100 001 à 200 000	597 946	78 515 622	12 812 674	579 100	76 088 011
200 001 à 300 000	88 163	21 094 844	4 697 851	86 452	20 285 145
300 001 à 400 000	28 254	9 670 624	2 422 510	27 654	9 537 660
400 001 à 650 000	2 254	2 254 254	2 254 254	2 254	2 254 254
650 001 à 1 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
1 000 001 à 1 500 000	2 254	2 254 254	2 254 254	2 254	2 254 254
1 500 001 à 2 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
2 000 001 à 3 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
3 000 001 à 4 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
4 000 001 à 5 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
5 000 001 à 6 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
6 000 001 à 7 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
7 000 001 à 8 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
8 000 001 à 9 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
9 000 001 à 10 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
10 000 001 à 11 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
11 000 001 à 12 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
12 000 001 à 13 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
13 000 001 à 14 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
14 000 001 à 15 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
15 000 001 à 16 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
16 000 001 à 17 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
17 000 001 à 18 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
18 000 001 à 19 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
19 000 001 à 20 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
20 000 001 à 21 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
21 000 001 à 22 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
22 000 001 à 23 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
23 000 001 à 24 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
24 000 001 à 25 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
25 000 001 à 26 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
26 000 001 à 27 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
27 000 001 à 28 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
28 000 001 à 29 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254
29 000 001 à 30 000 000	2 254	2 254 254	2 254 254	2 254	2 254 254

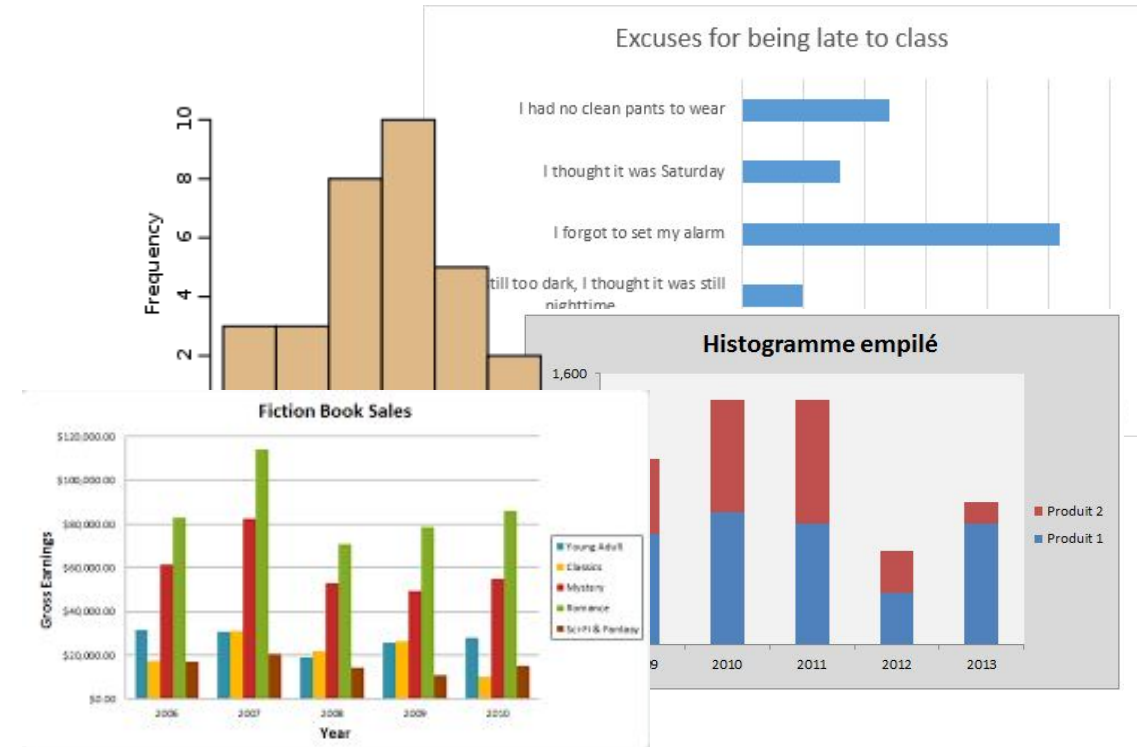


	015170, "01023", "Amia res-sus-nae", "B.	015180, "01024", "Atignat", "TATIGNAT	015190, "01025", "Bacq-d-le-Ville", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015210, "01027", "Balan", "BALAN, BAL	015220, "01028", "Baniac", "BANIAC, BANES"	015230, "01029", "Beaupont", "BEAUPO
	015170, "01023", "Amia res-sus-nae", "B.	015180, "01024", "Atignat", "TATIGNAT	015190, "01025", "Bacq-d-le-Ville", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015210, "01027", "Balan", "BALAN, BAL	015220, "01028", "Baniac", "BANIAC, BANES"	015230, "01029", "Beaupont", "BEAUPO
	015170, "01023", "Amia res-sus-nae", "B.	015180, "01024", "Atignat", "TATIGNAT	015190, "01025", "Bacq-d-le-Ville", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015210, "01027", "Balan", "BALAN, BAL	015220, "01028", "Baniac", "BANIAC, BANES"	015230, "01029", "Beaupont", "BEAUPO
015170, "01023", "Amia res-sus-nae", "B.	015170, "01023", "Amia res-sus-nae", "B.	015180, "01024", "Atignat", "TATIGNAT	015190, "01025", "Bacq-d-le-Ville", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015210, "01027", "Balan", "BALAN, BAL	015220, "01028", "Baniac", "BANIAC, BANES"	015230, "01029", "Beaupont", "BEAUPO
015180, "01024", "Atignat", "TATIGNAT	015180, "01024", "Atignat", "TATIGNAT	015180, "01024", "Atignat", "TATIGNAT	015190, "01025", "Bacq-d-le-Ville", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015210, "01027", "Balan", "BALAN, BAL	015220, "01028", "Baniac", "BANIAC, BANES"	015230, "01029", "Beaupont", "BEAUPO
015190, "01025", "Bacq-d-le-Ville", "B.	015190, "01025", "Bacq-d-le-Ville", "B.	015190, "01025", "Bacq-d-le-Ville", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015210, "01027", "Balan", "BALAN, BAL	015220, "01028", "Baniac", "BANIAC, BANES"	015230, "01029", "Beaupont", "BEAUPO
015200, "01026", "Bacq-d-le-Châtel", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015200, "01026", "Bacq-d-le-Châtel", "B.	015210, "01027", "Balan", "BALAN, BAL	015210, "01027", "Balan", "BALAN, BAL	015220, "01028", "Baniac", "BANIAC, BANES"	015230, "01029", "Beaupont", "BEAUPO	
015210, "01027", "Balan", "BALAN, BAL	015210, "01027", "Balan", "BALAN, BAL	015210, "01027", "Balan", "BALAN, BAL	015220, "01028", "Baniac", "BANIAC, BANES"	015220, "01028", "Baniac", "BANIAC, BANES"	015230, "01029", "Beaupont", "BEAUPO		
015220, "01028", "Baniac", "BANIAC, BANES"	015220, "01028", "Baniac", "BANIAC, BANES"	015220, "01028", "Baniac", "BANIAC, BANES"					
015230, "01029", "Beaupont", "BEAUPO	015230, "01029", "Beaupont", "BEAUPO	015230, "01029", "Beaupont", "BEAUPO					

Contexte & Définitions

Limites :

- Multitudes de graphiques disponibles et imaginables en fonction du jeu de données
- Nommage de ces graphiques non arrêté :
 - un graphe renvoie à plusieurs mot-clés (histogram / bar chart)
 - un mot-clé renvoie à plusieurs graphiques (bar chart)



Comment peut-on labéliser / normer les visuels de data-visualization ?

Objectifs

Collecter un jeu de données de data vizualisation et proposer une **méthode de classification** automatique

Créer un jeu de données

Classifier les visuels

Automatiser la classification
(ML)

Etendre l'étude à de plus
grands jeux de données

Démarche & Livrable

1

Créer un jeu de données

Collecter un jeu de données

Script de scraping
Google Image

Connection API Bing

Obtenir un jeu de données propre

Sélection d'image

Redimensionnement,
cropper



2

Classer les visuels

Sélectionner des visuels

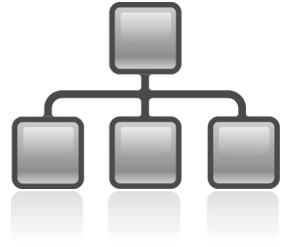
Recherche des visuels
les plus récurrents

Sélection de visuels
analysables (pour ML)

Annoter les visuels

Définition de labels
précis

Nommer les classes



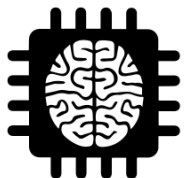
3

Méthode de classification automatique (Machine Learning)

Rechercher les algorithmes existants (Google, etc.)

Choix de l'algorithme et des méthodes les plus
pertinents (KMeans,...)
Définition des seuils de validation

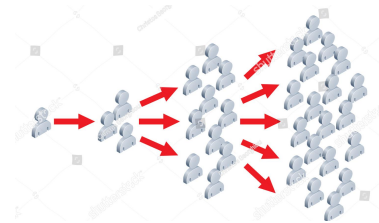
Classification binaire & classification
multi-catégories



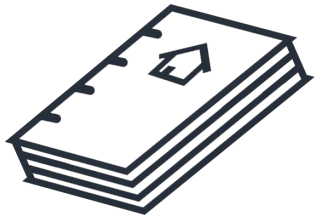
Démarche & Livrable

4

Étendre à de plus grands jeux de données, rendre public les résultats



Livrables :



- Une base de données de data visualizations
- Un programme de classification des data visualizations
- Une documentation explicative du projet détaillant:
 - La démarche du projet
 - Le fonctionnement du programme

Important et limites possibles:

identifier l'aspect légal du traitement des images, notamment par rapport aux droits des images

Planning

Décembre 2017

Constitution de l'équipe et prise de contact avec le tuteur

Jusqu'au 18 janvier

Phase de cadrage : définition des rôles, des livrables, des méthodes de travail
Identification des problématiques non-techniques pouvant impacter le projet(légales par exemple)
Premier reporting le 18 janvier

Jusqu'au 30 janvier

Constitution de la base de données de data vizualisation et premiers traitements en vue de l'obtention du jeu de données propres
Objectif : premier livrable pour le 30 janvier

Jusqu'au 15 février

Classification des visuels sur la base d'une méthodologie de sélection et de normage
Second reporting le 15 février

Jusqu'au 15 mars

Application des méthodes de machine learning afin d'obtenir notre programme de classification de data vizualisation adapté au grand jeu de données
Objectif : second livrable pour le 15 mars

Jusqu'au 29 mars

Préparation de la documentation, préparation de la présentation finale, rattrapage d'un éventuel retard et correction des possibles défauts
Troisième reporting le 29 mars

Organisation

Chef de projet



**Alexis MARTIN
- DELAHAYE**

Fonctionnel



Marc ARNAL

Dev Team



Luca GUÉRY



Louis KRAEMER



Arnaud BRUGIERE

Des questions ?
