

# Exploitation d'une Base de Données

## Cours 3 - *Open data* & visualisation

Anaïs Durand

28 Février 2023

# Open data

# Open data, données ouvertes/libres

- ▶ Données *accessibles* et *utilisables* librement
- ▶ *Origine* :
  - ▷ privée (entreprises, organisations à but non lucratif, ...), ou
  - ▷ publique (collectivités, établissements publics, ...)

# Open data, données ouvertes/libres

- ▶ Données *accessibles* et *utilisables* librement
- ▶ *Origine* :
  - ▷ privée (entreprises, organisations à but non lucratif, ...), ou
  - ▷ publique (collectivités, établissements publics, ...)

## Exemples :

- ▶ <https://ressources.data.sncf.com/> :  
données de la SNCF (trafic, retards, ...)
- ▶ <https://opendata.clermontmetropole.eu/> :  
données de la métropole de Clermont-Ferrand (monuments historiques, équipements sportifs, indice de qualité de l'air, ...)

# Intérêts de l'*open data*

- ▶ Pour la *recherche* : libre accès aux résultats scientifiques, reproductibilité des expériences, vérification des résultats ...
- ▶ Pour la *transparence/le contrôle des institutions* : vérifications des dépenses publiques, contrôle de l'action du gouvernement, ...
- ▶ Pour l'*intelligence artificielle (IA)* : plus de données = meilleur apprentissage
- ▶ ...

# Incitations à l'*open data*

## ► *Dans le monde :*

- ▷ 2004 : déclaration signées par les pays membres de l'OCDE pour rendre libres les publications scientifiques financées par des fonds publiques.

# Incitations à l'open data

## ► Dans le monde :

- ▷ 2004 : déclaration signées par les pays membres de l'OCDE pour rendre libres les publications scientifiques financées par des fonds publiques.

## ► En Europe :

- ▷ 2003, 2008, et 2019 : directives du conseil de l'UE pour inciter à l'open data surtout pour les données publiques
- ▷ 2015 : ouverture du portail européen  
<https://data.europa.eu/>



# Incitations à l'*open data*

## ► *En France :*

- ▷ 2011 : ouverture du portail <https://www.data.gouv.fr/>
- ▷ 2014 : mission etalab chargée de promouvoir et d'encadrer l'*open data* pour les données publiques
- ▷ 2016 : loi pour une République numérique  
⇒ par défaut toute donnée publique doit être ouverte

# Incitations à l'open data

## ► En France :

- ▷ 2011 : ouverture du portail <https://www.data.gouv.fr/>
- ▷ 2014 : mission etalab chargée de promouvoir et d'encadrer l'open data pour les données publiques
- ▷ 2016 : loi pour une République numérique  
⇒ par défaut toute donnée publique doit être ouverte

## ► Exemple récent :

CovidTracker (<https://covidtracker.fr/>)



# Utilisation des données ouvertes

# Récupérer des données ouvertes

► *Où ?* sur une banque de données en lignes :

- ▷ <https://www.data.gouv.fr/>
- ▷ <https://data.europa.eu/>
- ▷ <https://www.data.gov/>
- ▷ <https://www.kaggle.com/datasets>
- ▷ ...

# Récupérer des données ouvertes

► *Où ?* sur une banque de données en lignes :

- ▷ <https://www.data.gouv.fr/>
- ▷ <https://data.europa.eu/>
- ▷ <https://www.data.gov/>
- ▷ <https://www.kaggle.com/datasets>
- ▷ ...

► *Sous quel format ?*

un ou plusieurs fichiers au format .csv, .json, .xls, ...

# Préparation des données

## *Problèmes possibles :*

données non structurées, incomplètes, incohérentes ...

⇒ Besoin de préparer les données avant de les exploiter

- ▶ Structuration en tables
- ▶ Suppression des doublons
- ▶ Suppression ou complétion des données manquantes
- ▶ Rectification des incohérences
- ▶ ...

# Outils utilisés dans ce cours

- ▶ SGBD *PostgreSQL*
- ▶ Langage de programmation *Python*

# Outils utilisés dans ce cours

- ▶ SGBD *PostgreSQL*
- ▶ Langage de programmation *Python*
  - ▷ *psycopg2* : connexion à la base PostgreSQL
  - ▷ *pandas* : manipulation des données et interrogation de la base
  - ▷ *numpy* : manipulation des données



# Outils utilisés dans ce cours

- ▶ SGBD *PostgreSQL*
- ▶ Langage de programmation *Python*
  - ▷ psycopg2 : connexion à la base PostgreSQL
  - ▷ pandas : manipulation des données et interrogation de la base
  - ▷ numpy : manipulation des données



- ▶ Exemples dans le polycopié de cours + documentation des bibliothèques

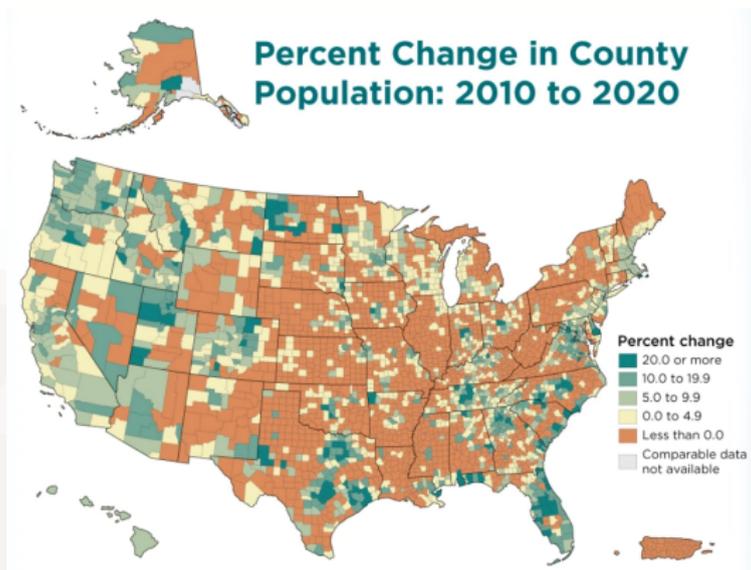
# Visualisation

# Visualisation de données

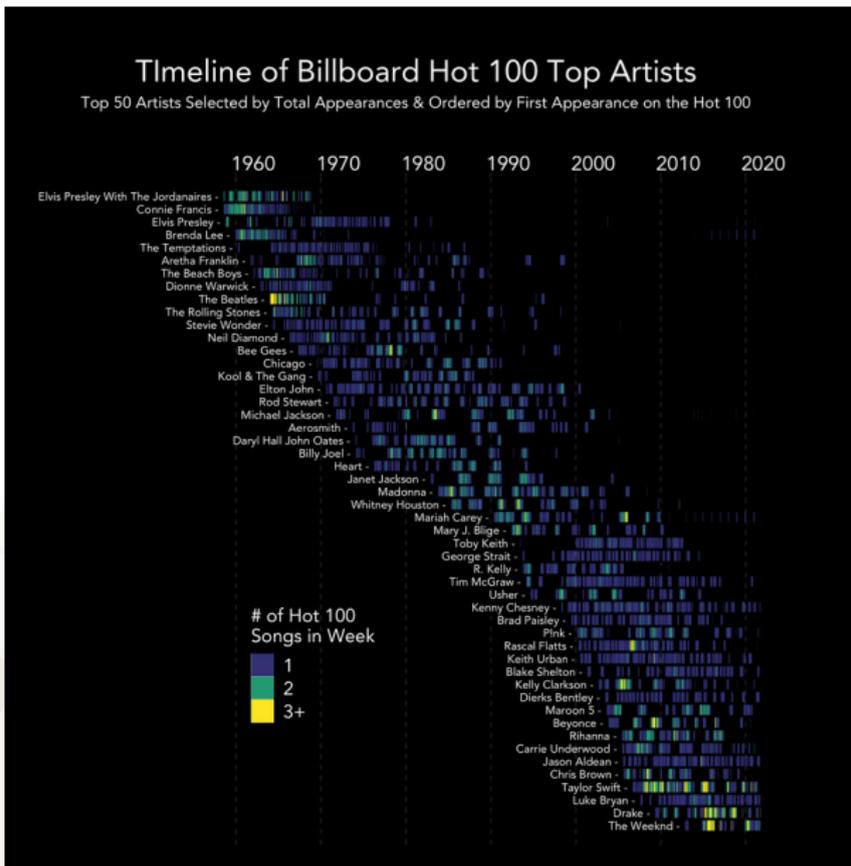
Réprésentation graphique des données

⇒ plus *lisible* que les données brutes

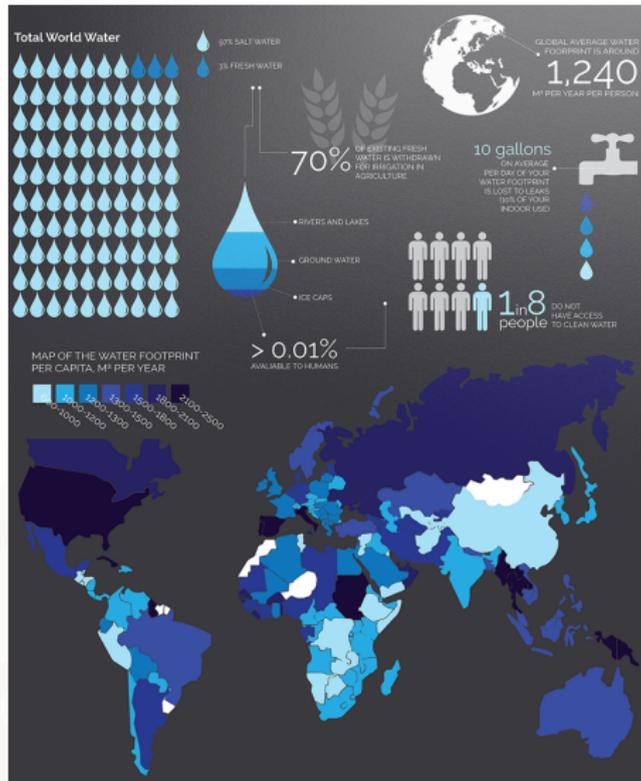
Exemple :



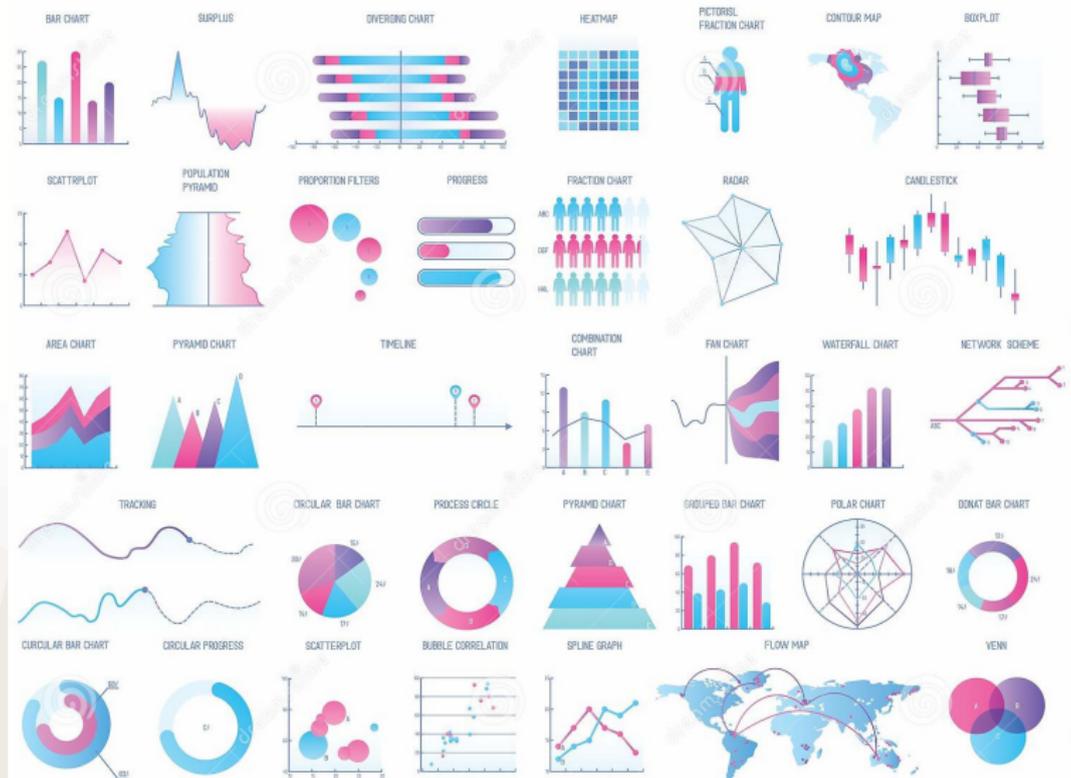
# Exemple :



# Exemple :



# Choix de la représentation



# Choix de la représentation

- Montrer une évolution dans le temps

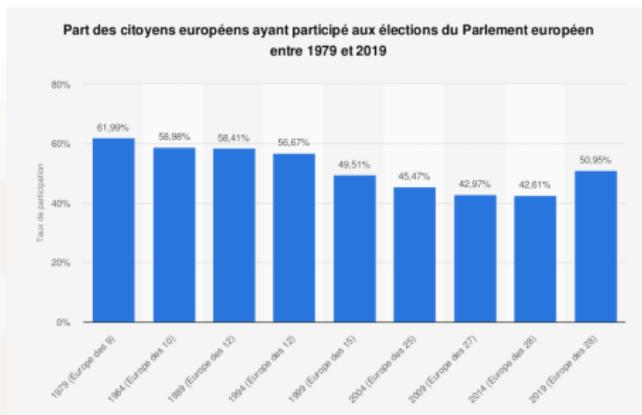
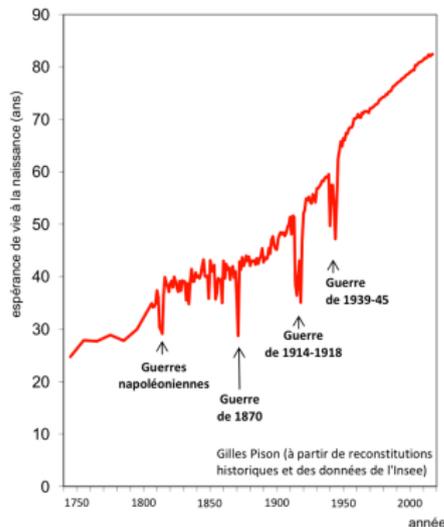
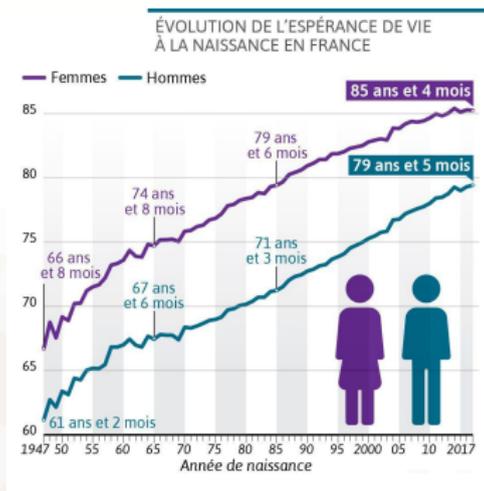


Figure 1. Évolution de l'espérance de vie à la naissance en France de 1740 à 2017

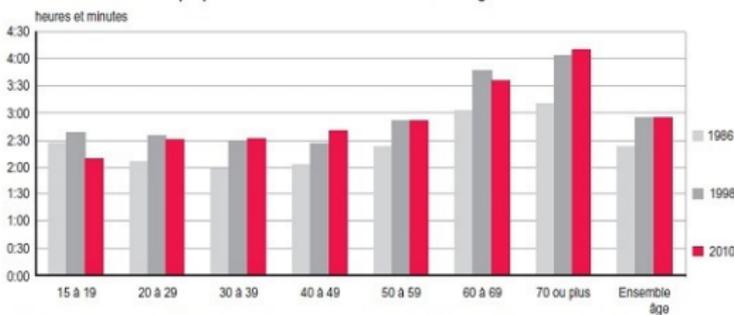


# Choix de la représentation

## ► Comparer des valeurs



## ① Évolution du temps passé devant la télévision selon l'âge



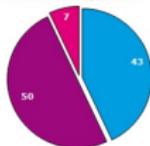
# Choix de la représentation

## ► Montrer une répartition/des proportions

Répartition du temps passé chaque semaine sur écran entre télévision, internet et ordinateur

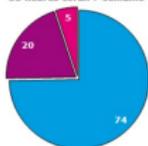
Le cas des 12-17 ans :

31 heures écran / semaine



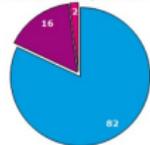
Le cas des 60-69 ans :

33 heures écran / semaine



Le cas des non-diplômés :

33 heures écran / semaine



Le cas des diplômés du supérieur :

40 heures écran / semaine

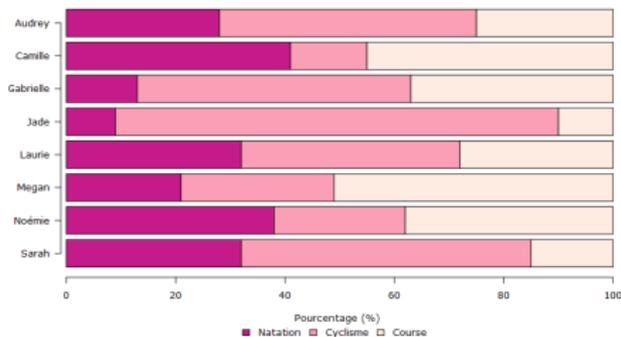


Télévision

Internet

Ordinateur

Triathlon à l'école Rousseau, pourcentage du temps consacré à chaque sport par les compétitrices





# Préparer les données à visualiser

- ▶ Filtrer les données  
⇒ choix des données à afficher, du format
- ▶ Trier les données  
⇒ afficher dans l'ordre le plus lisible (par valeur croissante, ...)
- ▶ Calculer des statistiques  
⇒ moyenne, médiane, écart-type ...

# Outils utilisés dans ce cours

- ▶ SGBD *PostgreSQL*
- ▶ Langage de programmation *Python*
  - ▷ pandas : manipulation des données, interrogation de la base, affichage
  - ▷ matplotlib : création de graphiques



