

# 30Music

## Finding similarities between songs

Clément BAFFOS

Holly HEALEY

Victor LALEUW

Ping'an YANG

01

Data cleaning

02

Partitioning

03

Density-based

04

Graph-based

05

Organisation et conclusion

# Présentation de 30Music

---

## Entities

Albums
Playlist
Tags
Tracks
Users

## Relations

Events
Love
Sessions

# Extraction des données

Format **.idomaar** : utilisé pour représenter des entités et leurs relations

- “users”

1	user	1	1116715959	{"lastfm_username":"000123","gender":"f","age":24,"country":"US","playcount":221012,"playlists":2,"subscribertype":"base"}
2	user	2	1163123792	{"lastfm_username":"000333","gender":"m","age":39,"country":"CZ","playcount":217535,"playlists":9,"subscribertype":"base"}
3	user	3	1184426573	{"lastfm_username":"00elen","gender":"f","age":,"country":"","playcount":49733,"playlists":2,"subscribertype":"base"}
4	user	4	1123157597	{"lastfm_username":"00Eraser00","gender":"m","age":32,"country":"DE","playcount":168054,"playlists":2,"subscribertype":"base"}
5	user	5	1171302116	{"lastfm_username":"00fieldsy","gender":"m","age":23,"country":"UK","playcount":45700,"playlists":2,"subscribertype":"base"}

- “love”

1	preference	1	-1	{"value":"love"} {"subjects":[{"type":"user","id":44542}], "objects":[{"type":"track","id":2785601}]}
2	preference	2	-1	{"value":"love"} {"subjects":[{"type":"user","id":44542}], "objects":[{"type":"track","id":2785590}]}
3	preference	3	-1	{"value":"love"} {"subjects":[{"type":"user","id":44542}], "objects":[{"type":"track","id":143076}]}
4	preference	4	-1	{"value":"love"} {"subjects":[{"type":"user","id":44542}], "objects":[{"type":"track","id":143037}]}
5	preference	5	-1	{"value":"love"} {"subjects":[{"type":"user","id":44542}], "objects":[{"type":"track","id":143052}]}

→ Nécessité de parser chaque fichier

Nettoyage préliminaire des données :

- suppression des colonnes inutiles issues du parsing des données
- suppression des lignes avec des données manquantes ou aberrantes (0, NaN ou null)

	age	id	country	gender	lastfm_username	playcount	playlists	subscribertype
0	24	1116715959	US	f	000123	221012	2	base
1	39	1163123792	CZ	m	000333	217535	9	base
2	32	1123157597	DE	m	00Eraser00	168054	2	base
3	23	1171302116	UK	m	00fieldsy	45700	2	base
4	48	1273363051	UK	m	01dela	3869	7	base
5	24	1229981595	UK	m	01srainey	77775	4	base

*Table "user" nettoyée*

Nettoyage en vue du partitioning :

- Transformation de certaines données qualitatives en données quantitatives pour pouvoir les exploiter

Variable qualitative	Variable quantitative
pays d'origine	proportion d'utilisateurs pour chacun des 6 pays les plus représentés
genre	proportion de femmes
âge	proportion d'utilisateurs par classe d'âge

01

Data cleaning

02

**Partitioning**

03

Density-based

04

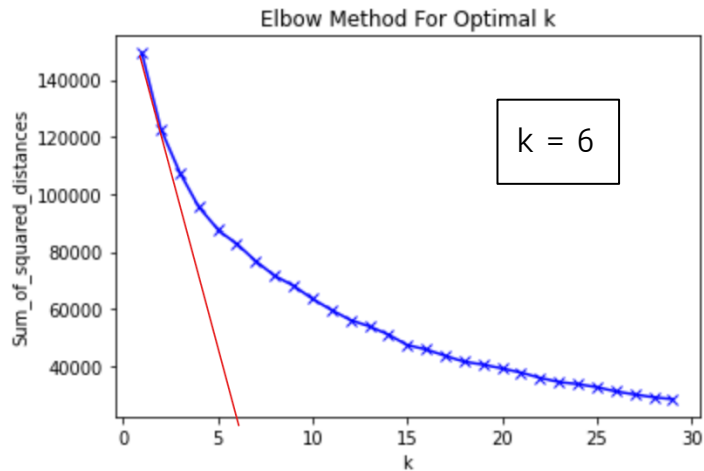
Graph-based

05

Organisation et conclusion

# Partitioning : K-means

- Préparation de la donnée :
  - variables quantitatives uniquement
  - création d'un dictionnaire track-id/track\_title pour interpréter les clusters
- Nombre de clusters déterminé par la méthode du coude



cluster\_0

track_id	track_name
2748152	Tara+Putra/_/Dubland+Coastline
3732410	Kalabi/_/Last+Words
1316489	Gil+Scott-Heron/_/Evolution+(And+Flashback)
1419055	Jacco+Gardner/_/Clear+the+Air
843899	Deacon+Blue/_/Fergus+Sings+The+Blues
1155214	Faithless/_/Postcards
1674649	Lyambiko/_/Feeling+Good
476434	Boyd+Rice+and+Fiends/_/Watery+Leviathan
3015283	The+Zombies/_/If+It+Don%27t+Work+Out
2087924	Orchestral+Manoeuvres+in+the+Dark/_/The+Misund...



## Analyse des clusters obtenus

---

### cluster\_0

- Musiques principalement écoutées en **Allemagne** (proportion moyenne d'Allemands sur ces chansons de 90.8%)
- Pas de classe d'âge spécifique

### cluster\_1

- Musiques principalement écoutées aux **Etats-Unis** par la tranche d'âge **31-40** (proportions supérieures à 70% pour ces 2 critères dans ce cluster)

### cluster\_2

- Musiques principalement écoutées aux **Etats-Unis** par la tranche d'âge **21-30** (proportions supérieures à 88% pour ces 2 critères dans ce cluster)

## Analyse des clusters obtenus

---

### cluster\_3

- Musiques principalement écoutées par la tranche d'âge **41-60** (proportions supérieures à 90%)
- Pas de pays spécifique mais légère prédominance aux **Etats-Unis** (59%)

### cluster\_4

- Musiques principalement écoutées par la tranche d'âge **21-30**
- Musiques écoutées dans tous les pays de façon relativement uniforme

### cluster\_5

- Musiques principalement écoutées au **Royaume-Uni** par les tranches d'âge **21-30** et **31-40**

01

Data cleaning

02

Partitioning

03

**Density-based**

04

Graph-based

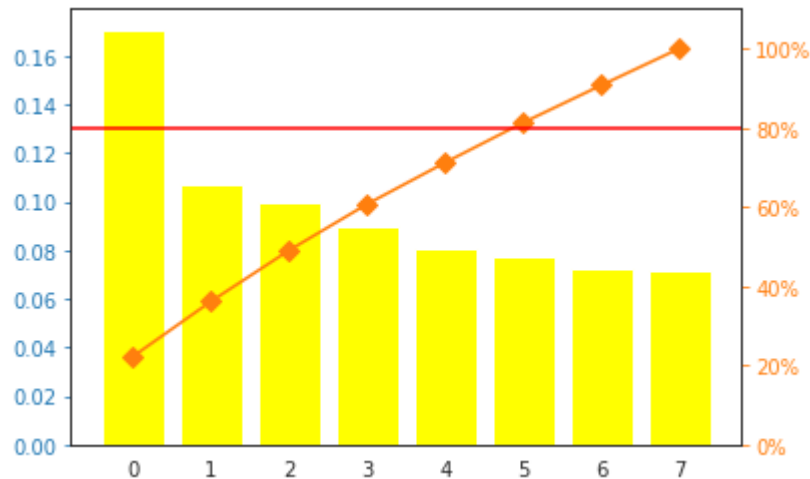
05

Organisation et conclusion

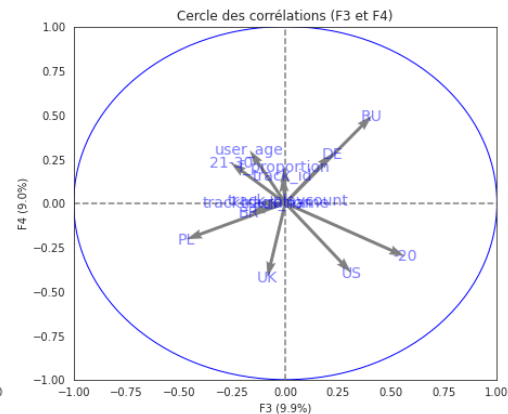
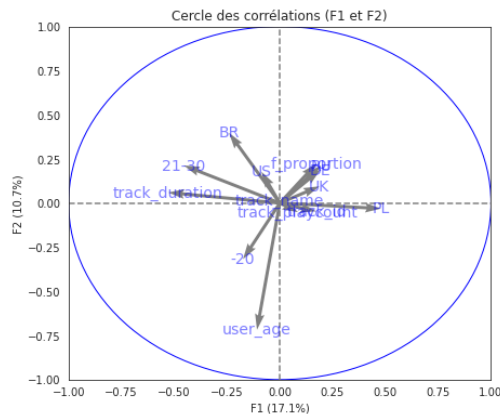
# Analyse en composante principale

- Les méthodes density-based ne permettent pas de traiter un grand nombre de dimension

Part de variabilité expliquée par les composantes



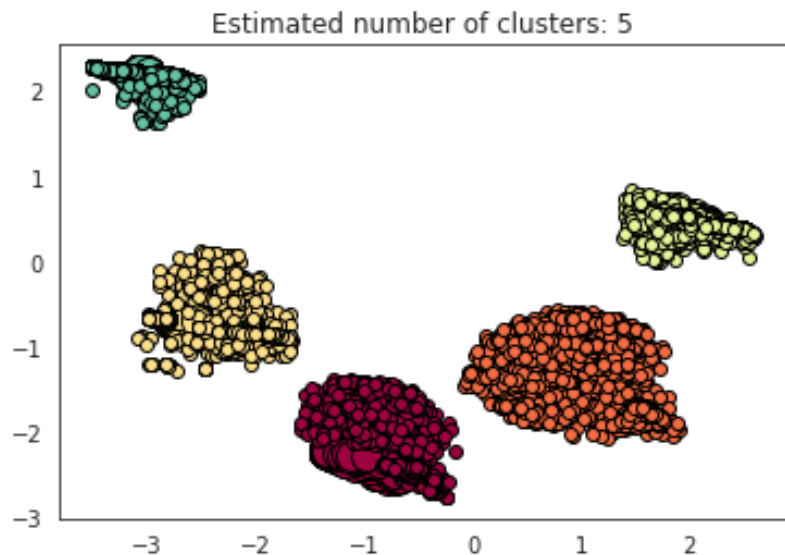
On conserve 4 composantes qui expliquent 60% de la variance



# Density-based : DBSCAN

Méthode DBSCAN :

- $\text{eps} = 0,5$
- $\text{min\_samples} = 5000$



- Nombre de clusters : 5
- Proportion de noise points : 52%

cluster\_0

	track_id	track_name
0	107	009+Sound+System/_/Born+To+Be+Wasted
1	122	009+Sound+System/_/Number+Two
16	1378	10,000+Maniacs/_/City+of+Angels
28	1497	1000+Funerals/_/Igneous+Lips
29	1499	1000+Funerals/_/Night%27s+Dew+(Shape+of+Despai...
30	1503	1000+Funerals/_/Sutured+Lips
31	1505	1000+Funerals/_/Your+Fancy
43	2253	108/_/Deathbed
44	2255	108/_/Declarations+On+A+Grave
66	2667	10+Ft.+Ganja+Plant/_/Midnight+Landing

## Analyse des clusters obtenus

---

### cluster\_0

- Musiques principalement écoutées aux **Etats-Unis** par la tranche d'âge **31-40** (proportions supérieures à 97% pour ces 2 critères dans ce cluster)
- **21%** de femmes

### cluster\_1

- Musiques principalement écoutées aux **Etats-Unis** par la tranche d'âge **21-30** (proportions supérieures à 95% pour ces 2 critères dans ce cluster)
- **31%** de femmes

### cluster\_2

- Musiques principalement écoutées aux **Etats-Unis** par la tranche d'âge **41-60** (proportions supérieures à 98% pour ces 2 critères dans ce cluster)
- **18%** de femmes

## Analyse des clusters obtenus

---

### cluster\_3

- Musiques principalement écoutées aux **Royaume-Uni** par la tranche d'âge **21-30** (proportions supérieures à 97% pour ces 2 critères dans ce cluster)
- **29%** de femmes

### cluster\_4

- Musiques principalement écoutées en **Allemagne** par la tranche d'âge **41-60** (proportions supérieures à 99% pour ces 2 critères dans ce cluster)
- **96%** d'hommes

01

Data cleaning

02

Partitioning

03

Density-based

04

Graph-based

05

Organisation et conclusion



- Création de la matrice avec les musiques et les utilisateurs
  - Ligne : Musique
  - Colonne : Musique
  - Valeur : Le nombre d'utilisateurs qui aiment la Ligne et la Colonne en même temps

[illegible]

# Spectral Clustering

	node_1	node_2	distance
0	60007	14820	100.0
1	71340	14820	100.0
2	71340	60007	100.0
3	74094	14820	100.0
4	74094	60007	100.0
...	...	...	...
124745	4770627	3755697	100.0
124746	4770627	3756296	100.0
124747	4770627	3770868	100.0
124748	4770627	3770874	100.0
124749	4770627	4768900	100.0

Préparation des données:

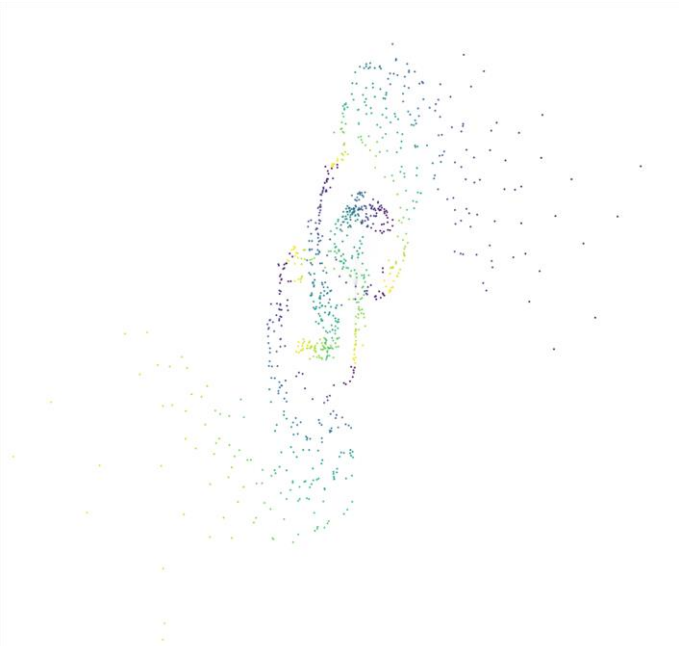
- Musique comme les nodes
- L'inverse des valeurs de matrice comme la distance entre deux nodes
- Pour les valeurs zéros on met 100 comme la distance
- \*S'il y a des utilisateurs aiment les deux musiques, la distance est  $\leq 1$ , sinon, c'est 100
- \*Il n'y pas de duplications des combinaisons entre "node\_1" et "node\_2"

# Spectral Clustering

---

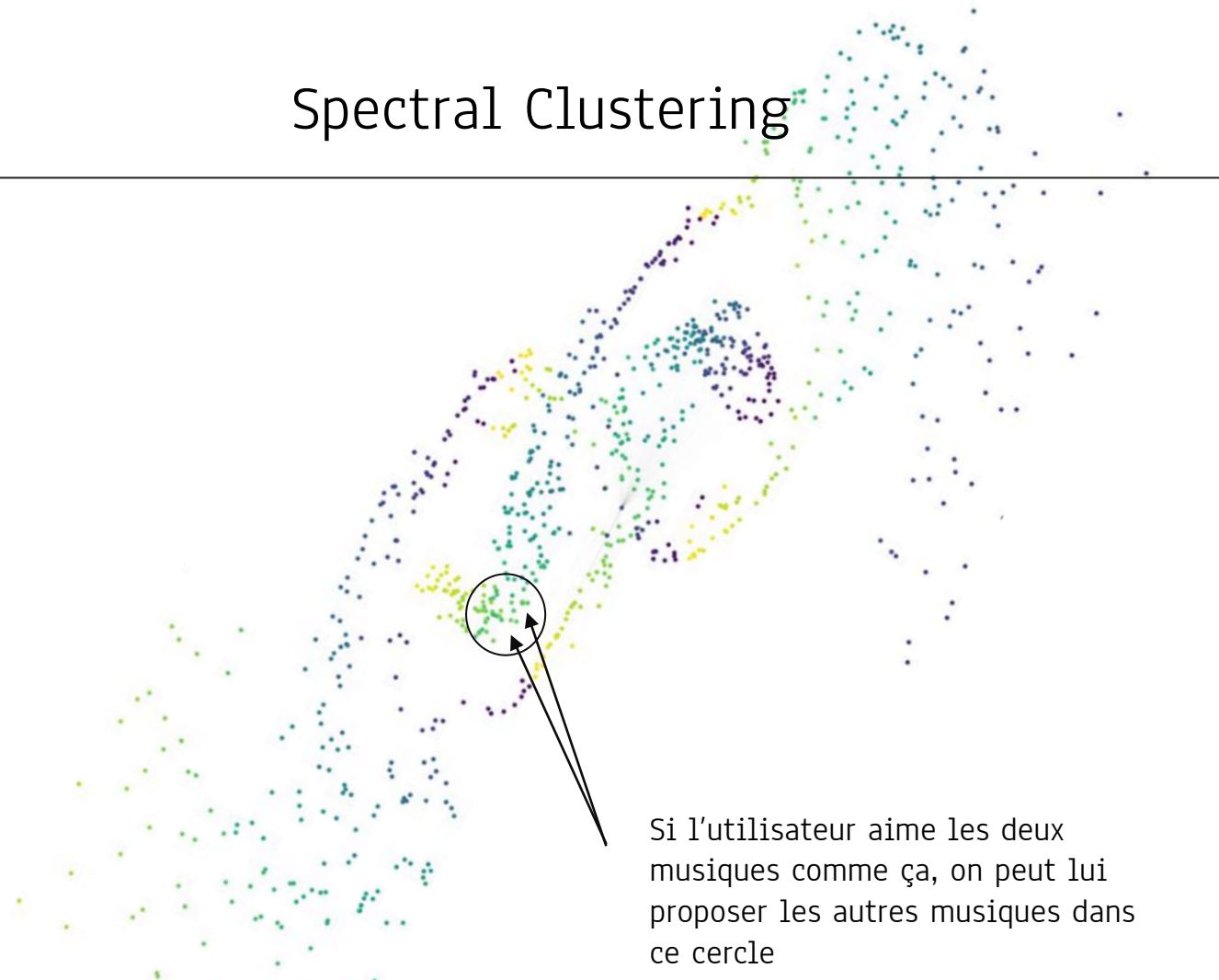
Spectral Clustering:

- Basé sur les librairies networkx et community
- Un point correspond à une musique
- La distance entre deux points correspond au nombre d'utilisateurs qui aiment les deux (juste pour la distance moins de 1)



# Spectral Clustering

---



01

Data cleaning

02

Partitioning

03

Density-based

04

Graph-based

05

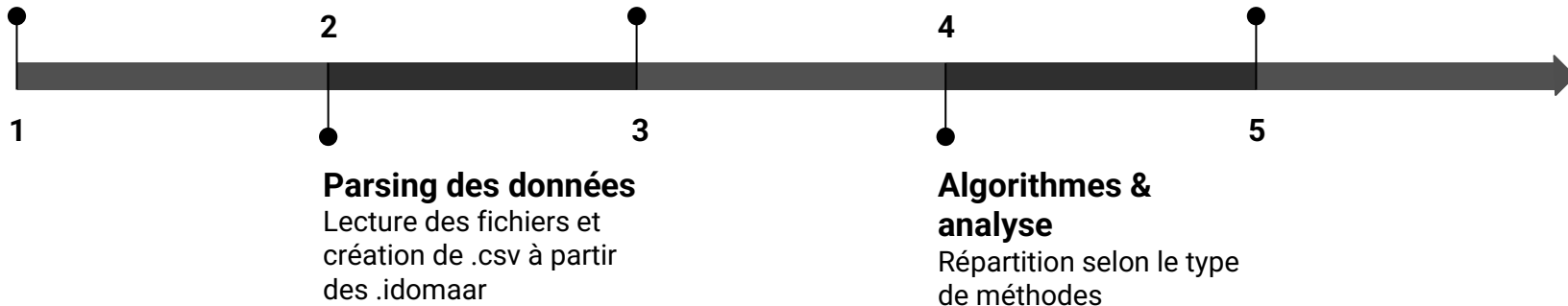
Organisation et conclusion

# Organisation

---

## Récupération des données

Base de données accessible en ligne et mise en place d'un drive collaboratif



## Business use (type spotify)

---

But : satisfaire les envies musicales des utilisateurs pour les fidéliser.

### Algorithmes :

- Proposant des musiques écoutées par les personnes du même groupe que l'utilisateur.
- Proposant des musiques d'un autre cluster (mais proche) que celui formé par les méthodes graph-based.
- Création de playlists

### Utilité :

- Proposant des musiques à écouter pour un client basé sur ses précédentes écoutes et les musiques qu'il a aimées en fonction des préférences des autres personnes de son cluster.
- Découverte des musiques d'un autre groupe proche du sien.
- Complémentaire avec les utilisations précédentes.

## Business use (plus original)

---

### Algorithmes :

- Met en relation différentes personnes d'un même cluster
- Déterminer le genre de musique qui plaît le plus

### Utilité :

- Permet aux utilisateurs de rencontrer des personnes écoutant le même type de musique pour divers buts (concerts...)
- Créer des musiques qui vont plaire à un large groupe (but commercial)



## Pour aller plus loin ...

---

### Avec le même Dataset

Méthode graph based :

- Prise en compte de plus de musiques dans le graph (seulement 500 ici)
- Sélection plus pertinente des musiques (ici celle qui sont le plus aimées)

Autres idées de graph possible :

- Nombre de playlist dans lesquelles les 2 musiques sont présentes
- Nombre de session d'écoute dans lesquelles les 2 musiques sont présentes
- Nombre de tags en commun

### Avec un autre Dataset

- Données labellisées pour pouvoir faire de l'apprentissage supervisé
- Données audio pour pouvoir analyser les similitudes

Merci pour votre  
attention

---