# Augmentez vos données avec les algorithmes de graphes

•••

Estelle Scifo
PyConFR 2019, Bordeaux

# About me

- Physicist & Data scientist (Luxembourg/World/Remote)
- Graph enthousiast
    - Last project: neomap, visualization tool for Neo4j (written in React)
- Slides and code samples available on my github
                    github.com/stellasia/pyconfr19
- Get in touch via twitter/linkedin!

*twitter: @st3llasia*
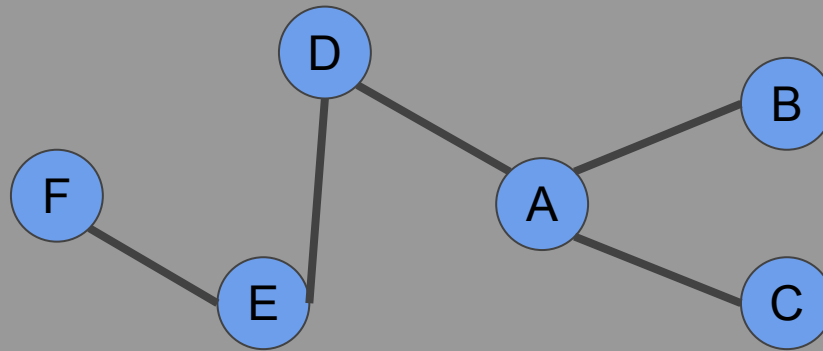*github: stellasia*
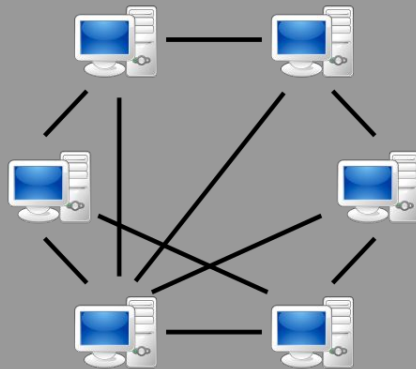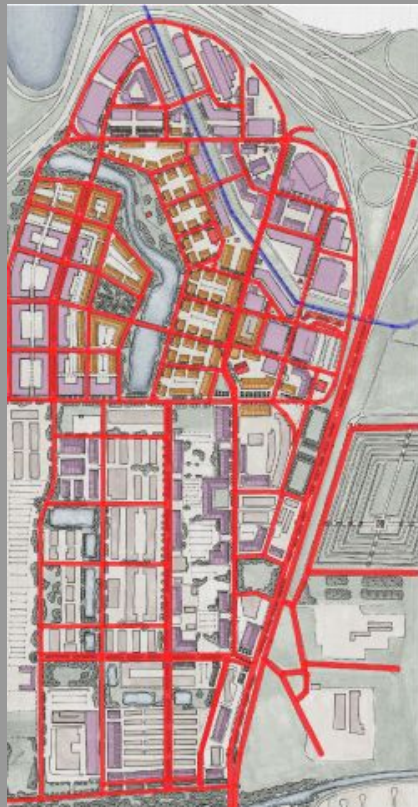*linkedin: estellescifo*

# Graphes

https://www.vox.com/2018/12/13/18106455/best-of-2018-data-charts-tech-end-year-list-amazon-facebook-juul-moviepass-elon-musk

# Graphes

https://www.vox.com/2018/12/13/18106455/best-of-2018-data-charts-tech-end-year-list-amazon-facebook-juul-moviepass-elon-musk
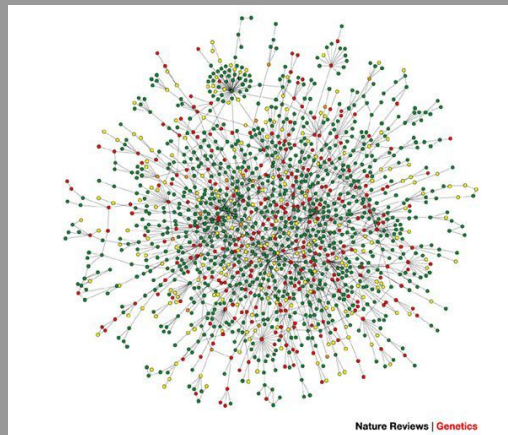
# Graphes

# Des réseaux...

# Mais aussi…



Proteins interaction network
https://www.nature.com/articles/nrg1272



Fig. 1. Minimum spanning tree (MST) obtained by considering the daily logarithmic return time series of all stocks traded in the Greek stock market (Athens Stock Exchange, ASE) in year 1997. The different colours represent the different sectors of economic activity. The existence of stocks that behave as hubs, having more connections than the average, is evident. Inset: MST obtained using all stocks traded in ASE for the year 2000, which is in the middle of the 1999–2001 crisis period.

Stock market
http://kelifos.physics.auth.gr/publications/pdf/p133.pdf



Nodes2019 conference
@WilliamLyon

# Pourquoi maintenant ?



**Complete trend, starting with January 2013**

Legend:
- Graph DBMS
- Document stores
- Search engines
- Key-value stores
- Time Series DBMS
- Wide column stores
- RDF stores
- Native XML DBMS
- Object oriented DBMS
- Multivalue DBMS
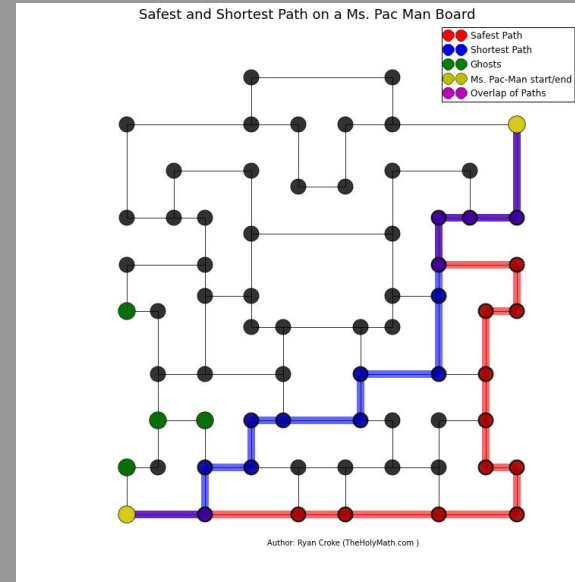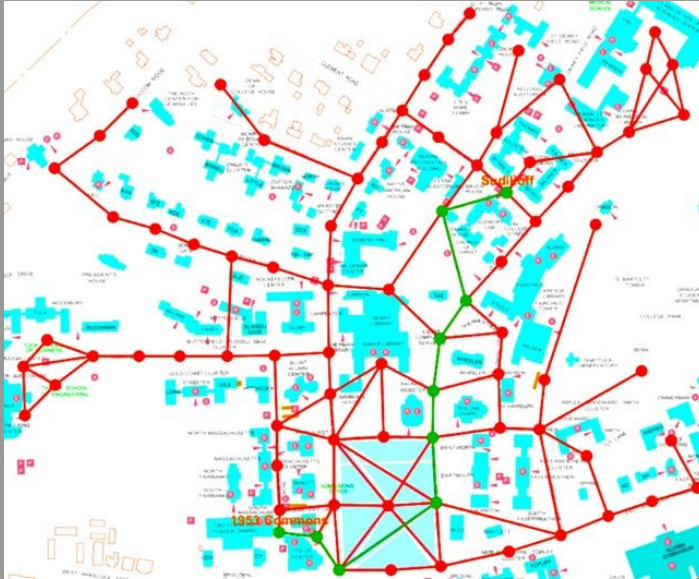- Relational DBMS

© 2019, DB-Engines.com

# Utilisations

- Extraction d'information (description)

- Extraction de "**graphy**" features pour ML

- Node/Graph embedding

- Prédiction de liens

# 1. Algorithmes
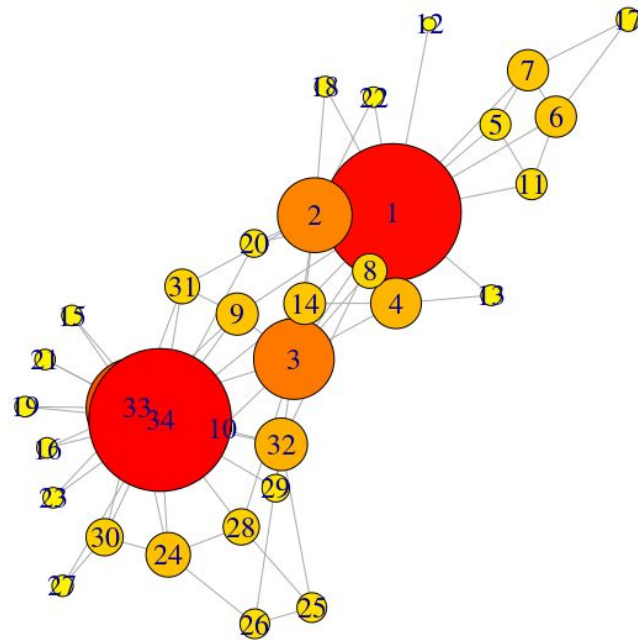
# Plus court chemin
# (Shortest Path)

# Shortest Path



Safest and Shortest Path on a Ms. Pac Man Board

Author: Ryan Croke (TheHolyMath.com)

# Importance (Centrality)

# PageRank

"Nombre et importance des noeuds liés à A"

```
PR(A) =
(1-d) + d Σ (PR(i) / Ni)
```
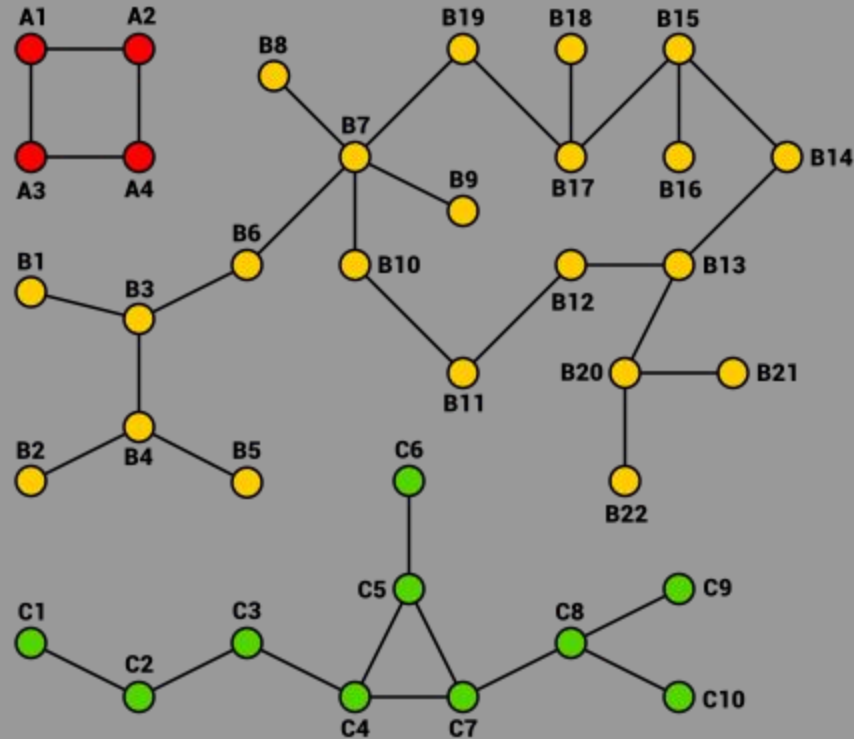
# Betweenness

"Nombre de plus courts chemins entre deux noeuds i et j qui passent par A"
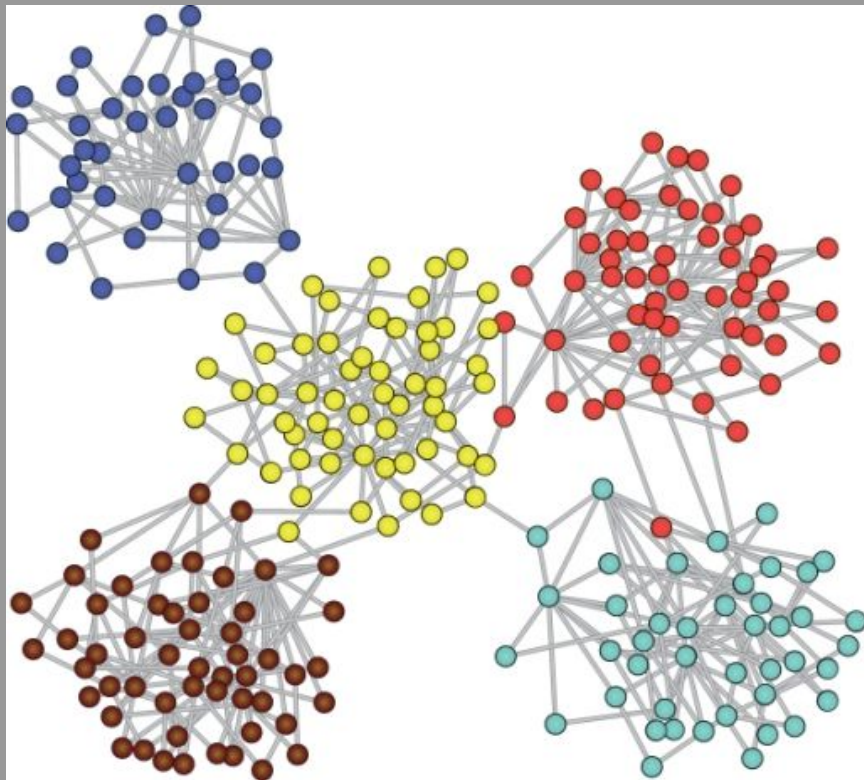
$\Rightarrow$ Noeud "critique"
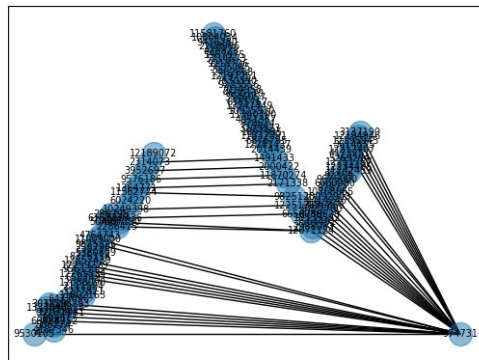
# Communauté (Clustering)

# Communauté - Connected components
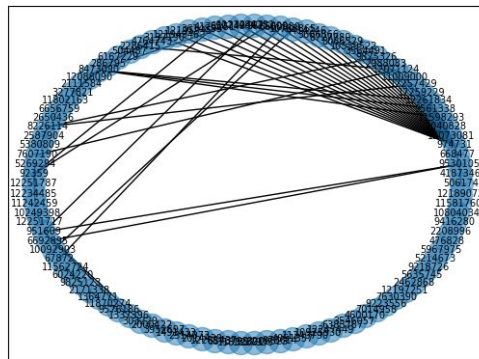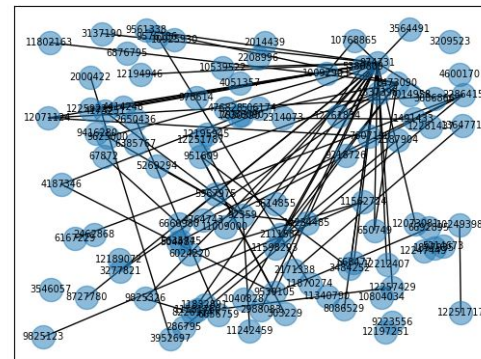
# Communauté

# Visualisation



Layout: planar

Layout: random

Layout: circular
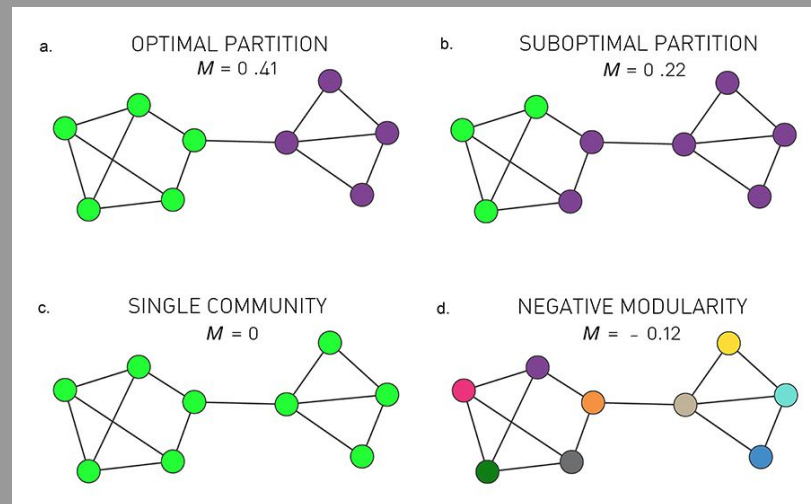
# Communauté - Louvain

- Modularité :

  *Mesure la densité de lien à l'intérieur d'une communauté*

  *VS liens entre communautés*

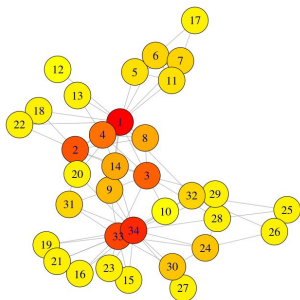$$Q = \frac{1}{2m} \sum_{ij} \left[ A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j),$$



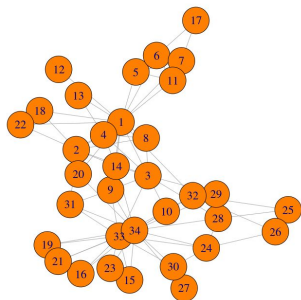- ~ KMeans

# Et d'autres…

## Triangles



Triangles

Number of triangles using this node

Probability that two neighbours are also connected

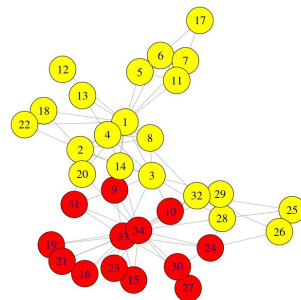## Components



Strongly connected components

Paths going through the node

(weak: one direction
strong: both directions)

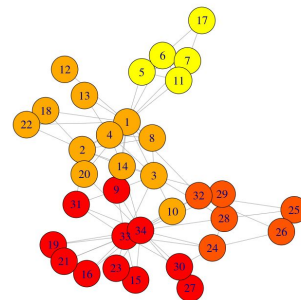## Label propagation



Label propagation

Label propagation to neighbours using the "majority vote" rule

## Louvain



Louvain

More relationships among nodes within a cluster than with nodes outside the cluster

# Prédiction de lien
# (Link Prediction)

# Prédiction de lien

- Combler un manque d'information :
    - Étude de réseaux criminels
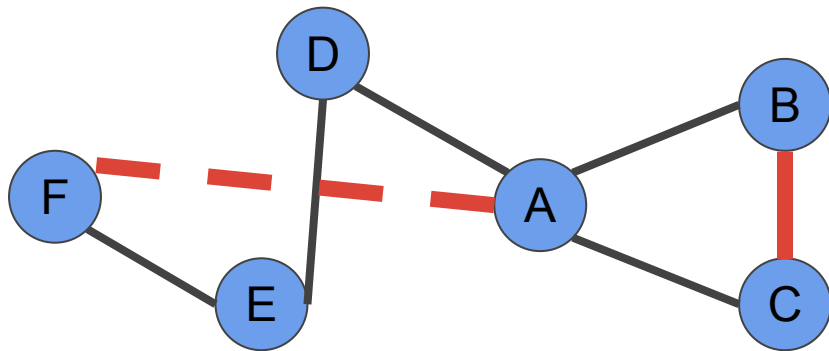
- Prédire des liens futurs :
    - Recommendation

# Metrics for link prediction

- Common neighbours

⇒ *Two people who have a friend in common are more likely to be introduced than those who don't have any friends in common.*

⇒ Size of the set of common friends between u and v:

$$C(u, v) = |\mathcal{N}(u) \cap \mathcal{N}(v)|$$
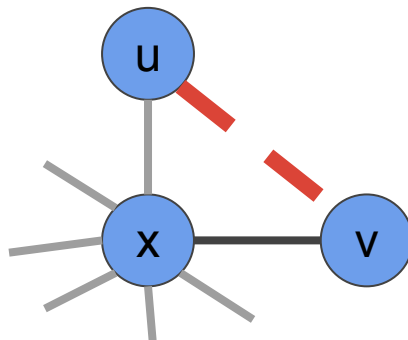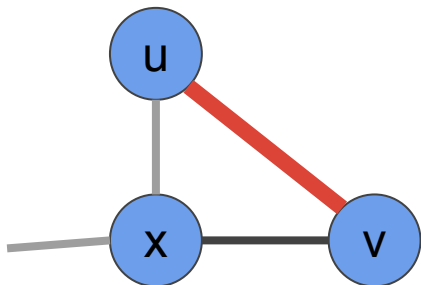
# Metrics for link prediction

- Common neighbours
- Adamic-Adar

  $\Rightarrow$ *Rare connections are giving more informations than common ones*

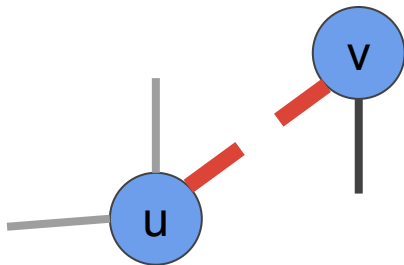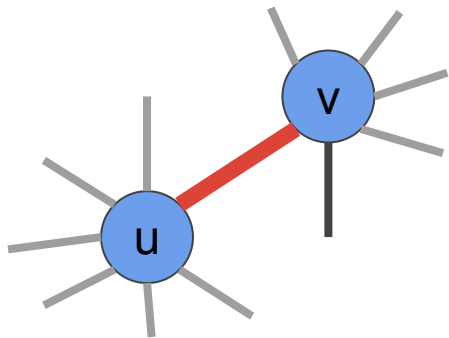$$A(u,v) = \sum_{x \in \mathcal{N}(u) \cap \mathcal{N}(v)} \frac{1}{log|\mathcal{N}(x)|}$$

# Metrics for link prediction

- Common neighbours
- Adamic-Adar
- Total neighbours

  $\Rightarrow$ *the more connected a node is, the more social it is, the more likely it is to receive new links*

$$T(u,v) = |\mathcal{N}(u) \cup \mathcal{N}(v)|$$

# 2. Exemple d'utilisation en python

# Comment ?

- Base de données graphes :
  - Incluant des implémentations des principaux algorithmes (Neo4j)

- Package python : `networkx`

Features
- Data structures for graphs, digraphs, and multigraphs
- Many standard graph algorithms
- Network structure and analysis measures
- Generators for classic graphs, random graphs, and synthetic networks
- Nodes can be "anything" (e.g., text, images, XML records)
- Edges can hold arbitrary data (e.g., weights, time-series)
- Open source 3-clause BSD license
- Well tested with over 90% code coverage
- Additional benefits from Python include fast prototyping, easy to teach, and multi-platform

# Export des données

- Créer un graphe :
    - A partir de fichiers csv

```
nodes.csv:
node
2014439
12281437
6876795
506174
11832891
12212407
668477
4764743
```
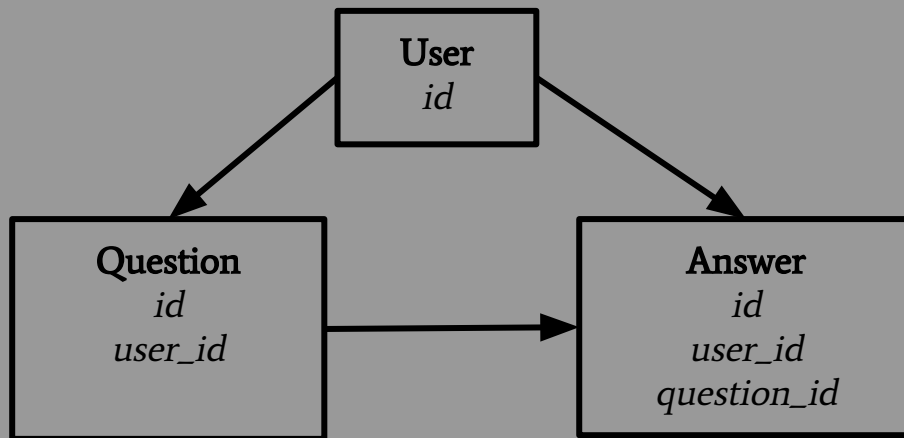
```
edges.csv:
node1, node2
506174,4187346
668477,974731
4764743,2286415
12073081,974731
1040828,974731
11598293,8473090
11598293,974731
12088090,8473090
```
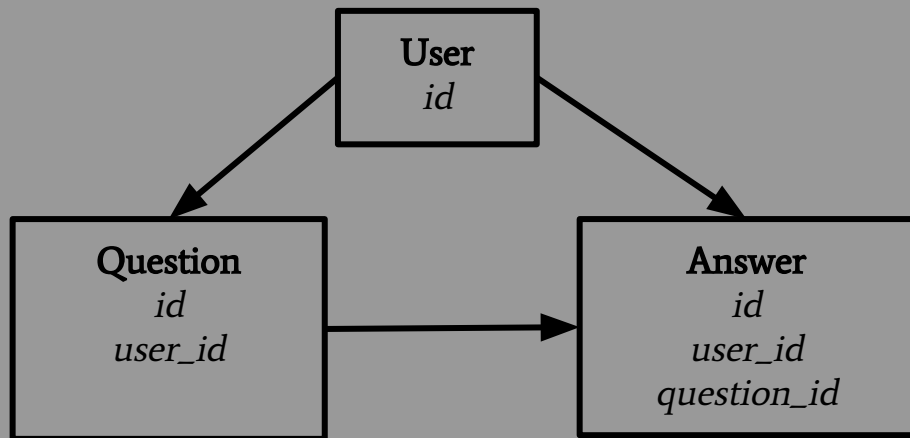
# Export des données

- Créer un graphe :
    - Ex: base de données type StackOverflow:
        - **Liens entre les utilisateurs ?**

```
          ┌─────────────┐
          │    User     │
          │     id      │
          └─────────────┘
          ↙             ↘
┌──────────────┐      ┌──────────────┐
│   Question   │      │    Answer    │
│     id       │ ───→ │     id       │
│   user_id    │      │   user_id    │
│              │      │  question_id │
└──────────────┘      └──────────────┘
```

# Export des données

- Créer un graphe :
  - Ex: base de données type StackOverflow:
    - User; Question: Answer
    - **Liens entre les utilisateurs ?**



```
nodes.csv:


SELECT u.id
FROM users
```

```
edges.csv:


SELECT DISTINCT u1.id, u2.id
FROM user u1
JOIN question q ON q.user_id = u1.id
JOIN answer a ON a.question_id = q.id
JOIN user u2 ON a.user_id = u2.id
```

# Import dans networkx

```
import networkx as nx

G = nx.Graph()

edges = nx.read_edgelist('edges.csv', delimiter=",")
nodes = nx.read_adjlist("nodes.csv")

G.add_edges_from(edges.edges())
G.add_nodes_from(nodes)
```
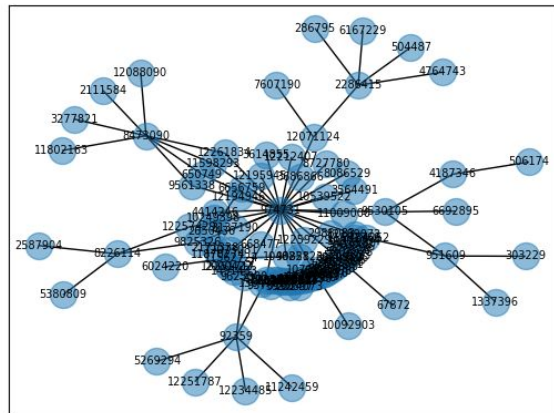
# Visualisation

```python
import networkx as nx
import matplotlib.pyplot as plt

pos = nx.kamada_kawai_layout(G)
f = plt.figure()

plt_nodes = nx.draw_networkx_nodes(
    G,
    pos,
    nodelist=G.nodes,
    alpha=0.5
)
nx.draw_networkx_edges(G, pos)
nx.draw_networkx_labels(G, pos, font_size=7)
```
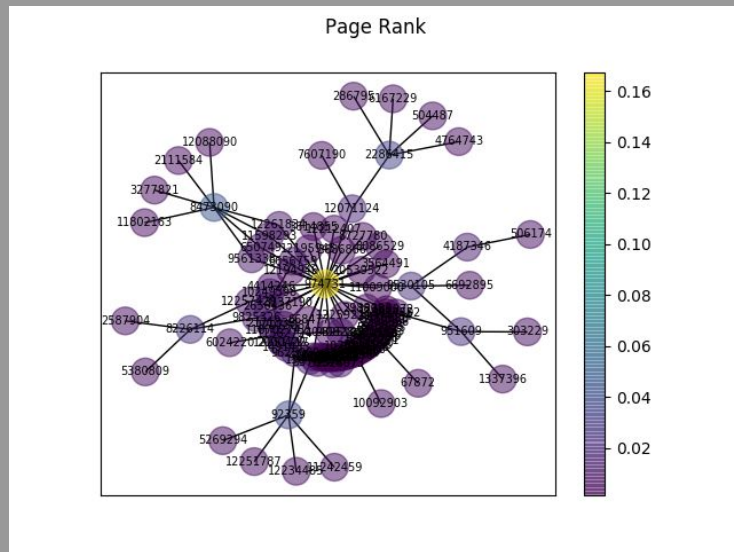
# 2.1 Description des données

# Algorithme : importance (PageRank)

```
res = nx.pagerank(G, alpha=0.85)
res_nodes = list(res.keys())
res_values = list(res.values())

f = plt.figure()
plt_nodes = nx.draw_networkx_nodes(
    G,
    pos,
    nodelist=res_nodes,
    node_color=res_values,
)
nx.draw_networkx_edges(G, pos)
nx.draw_networkx_labels(G, pos, font_size=7)
```
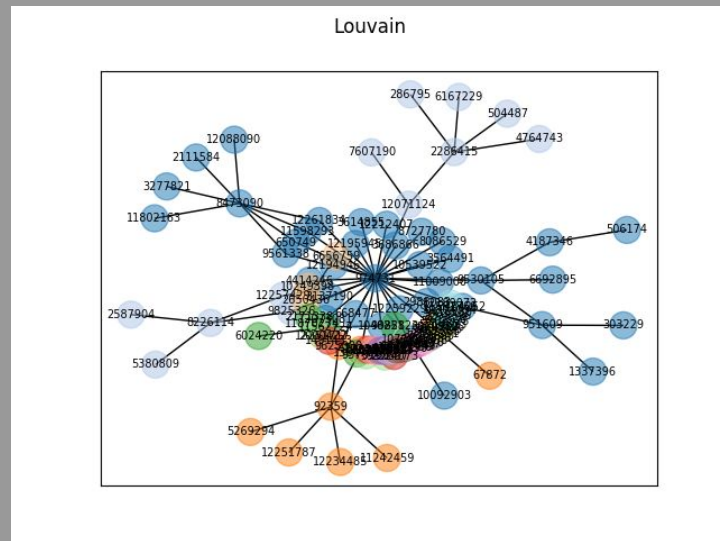


Page Rank

# Algorithme : communauté (Louvain)

```
import community

res = community.best_partition(G)
res_nodes = list(res.keys())
res_values = list(res.values())

f = plt.figure()
plt_nodes = nx.draw_networkx_nodes(
    G,
    pos,
    nodelist=res_nodes,
    node_color=res_values,
)
nx.draw_networkx_edges(G, pos)
nx.draw_networkx_labels(G, pos, font_size=7)
```
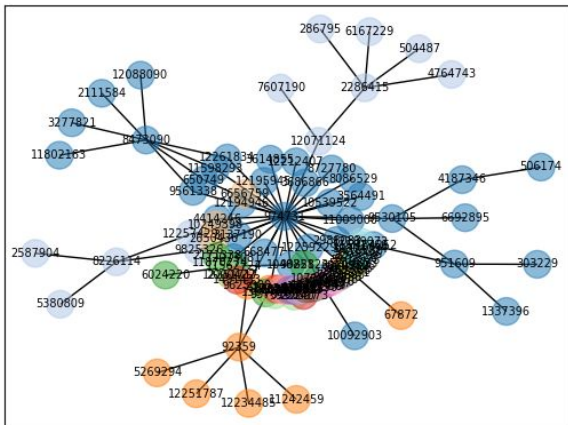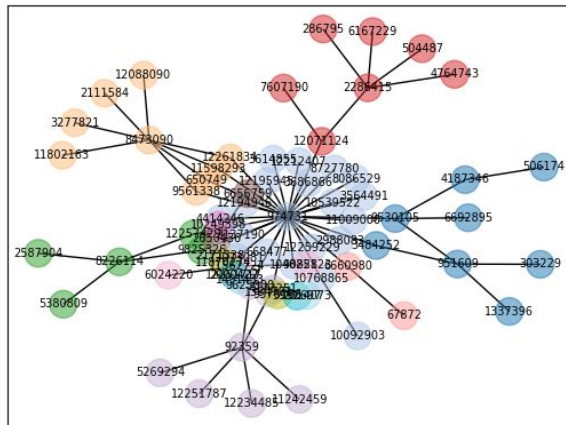


Louvain

# Noeuds isolés

```
import networkx as nx


G.remove_nodes_from(list(nx.isolates(G)))
```
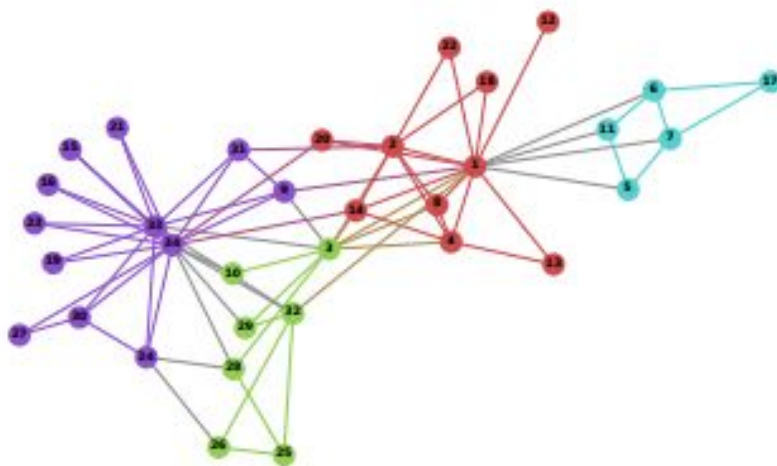
# 2.2 Prédiction
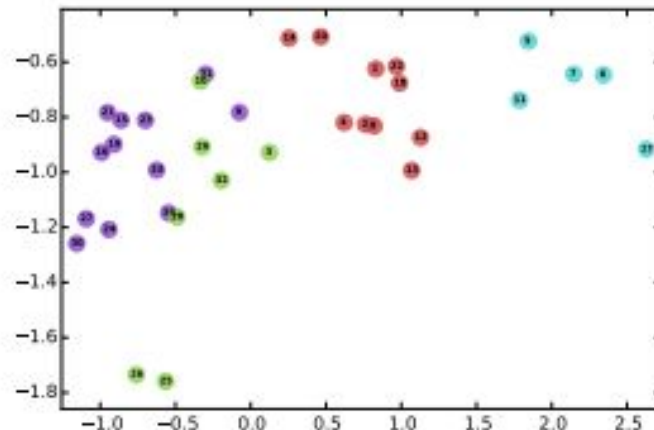
# "Graphy" features

- Problème :
    - Est-ce qu'une personne donnée va venir à la prochaine PyConFR en 2020 ?

| name (str) | age (int) | is_afpy_member (bool) | was_present_in_2019 (bool) | ... | Importance (float) | Community (int) |
|---|---|---|---|---|---|---|
| Estelle | 30 | false | true | ... | 0.01 | 4 |
| Pierre | 22 | true | true | ... | 1.4 | 1 |
| …. | …. | …. | …. | ... | ... | ... |

# Node embeddings



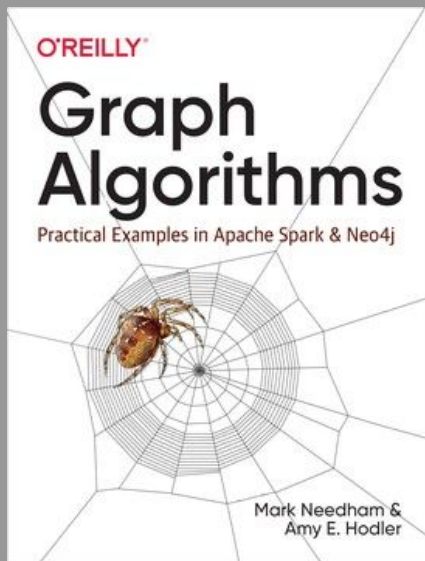(a) Input: Karate Graph

(b) Output: Representation

# Conclusion







*twitter: @st3llasia*

*github: stellasia*

*linkedin: estellescifo*