

Flow of knowledge in information networks

Jean-François BAFFIER¹

¹Japan Society for the Promotion of Science
French National Center for Scientific Research
Hosted at the Tokyo Institute of Technology

Monday, December 18, 2018

Network Optimization course

Table of contents

- 1 Introduction
- 2 Information Networks
- 3 Flow of knowledge
- 4 Multiplex Networks
- 5 Cycles, connected component and ethics
- 6 Conclusion

Contextualization

A not very funny joke

Three scientists, one in computer science, one in humanities, and one in biology enter a café...

Contextualization

A not very funny joke

Three scientists, one in computer science, one in humanities, and one in biology enter a café...

⇒ They have the same number of citations !

J.-F. Baffier (2018)

Contextualization

A not very funny joke

Three scientists, one in computer science, one in humanities, and one in biology enter a café...

⇒ They have the same number of citations!

J.-F. Baffier (2018)

The analysis of academic citation networks depends strongly on

- Research field, era, host institution, fundings, ...
- Network size, its shape, ... so its structure
- Analytical tools used!

Comment en suis-je arrivé là ?



Table of contents

- 1 Introduction
- 2 Information Networks**
- 3 Flow of knowledge
- 4 Multiplex Networks
- 5 Cycles, connected component and ethics
- 6 Conclusion

The value of the transmission of knowledge

Education is the transmission of civilization.

(Areil and Will DURANT — 1968)

*Structure is more important than content in the transmission of information.
(The medium is the message)*

(Abbie HOFFMAN — 1968)

An academic citations network

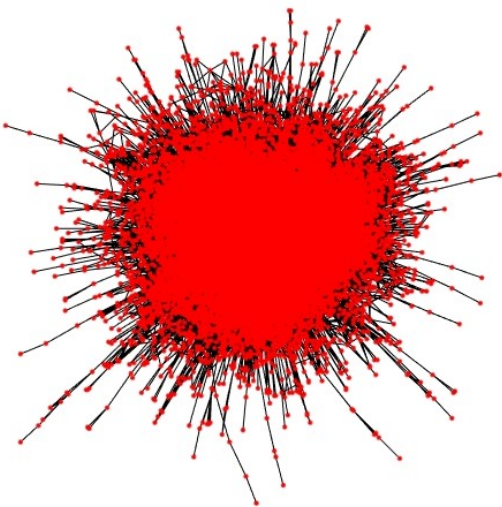


FIGURE — Arxiv HEP-TH (high energy physics theory) : 27770 nodes (articles) — 352807 arcs (citations)

An academic citations network

Expected properties :

- **Directed**
- (Pseudo-) **Acyclic**
- **Dynamic**
- **Stable** nodes (articles) : the papers are not modified after publication
- **Large et complex**
 - According to the University of Ottawa, in 2009 the amount of 50 millions of academic publications was reached (starting from 1665). Currently there is about 2.5 millions of new publications per year.
 - With the constant growth in the number of researcher and the global growing population, it is believed that the total number of publications will double every 9 years.
- **Unweighted** (on the links)
- **Information** is generated on each node (or at least on many)

Other information networks

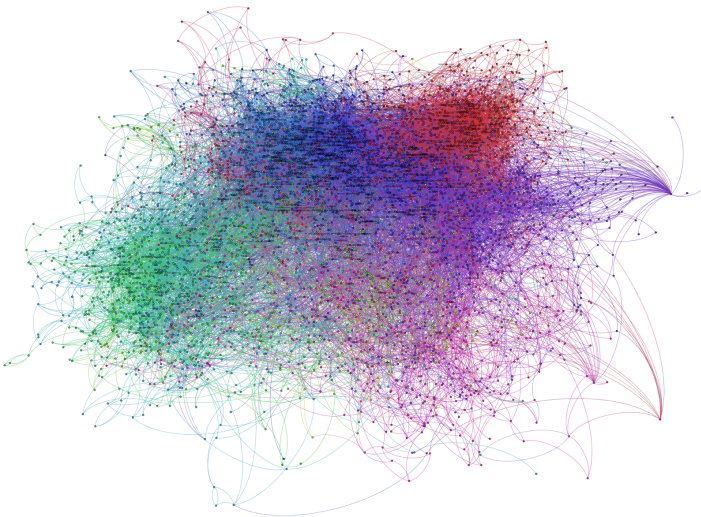


FIGURE – Réseaux d'interconnexions des domaines et disciplines dans Wikipédia

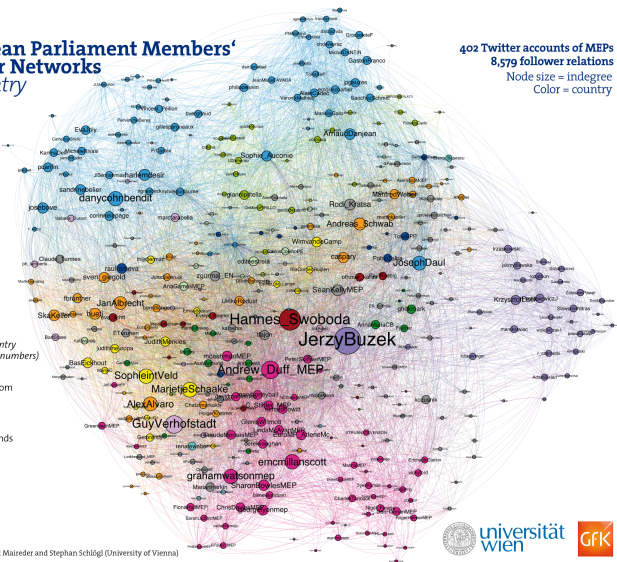
Other information networks

European Parliament Members' Twitter Networks by country

402 Twitter accounts of MEPs
8,579 follower relations
Node size = indegree
Color = country

Accounts by country
(in order of user numbers)

- France
- United Kingdom
- Germany
- Poland
- Italy
- The Netherlands
- Sweden
- Spain
- Belgium
- Portugal
- Romania
- Austria
- Other



CC BY-SA 4.0 — Axel Mairreder and Stephan Schögl (University of Vienna)



universität
wien



Other information networks

Game of Thrones Family Ties

- Node color:
- House Baratheon
 - House Frey
 - House Greyjoy
 - House Lannister
 - House Martell
 - House Stark
 - House Targaryen
 - House Tully
 - House Tyrell
- Edge color:
- father
 - brother/sister
 - mother
 - spouse

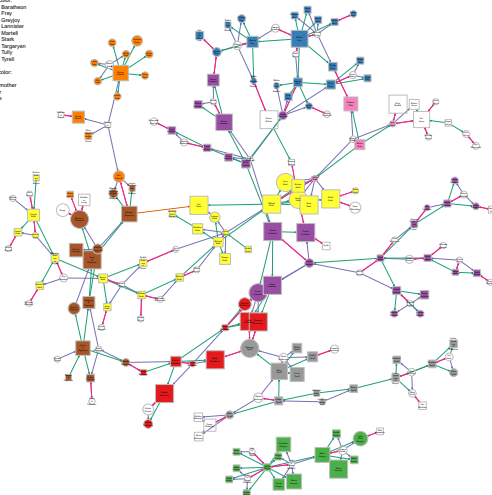


Table of contents

- 1 Introduction
- 2 Information Networks
- 3 Flow of knowledge**
- 4 Multiplex Networks
- 5 Cycles, connected component and ethics
- 6 Conclusion

The h -index

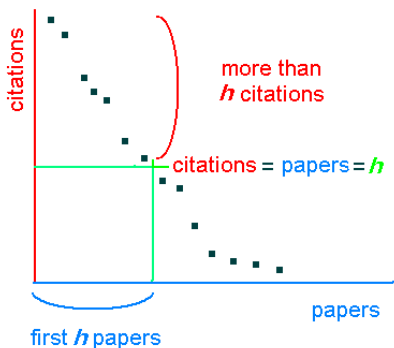


FIGURE – Short exercise : what is the value of h here ?

The Hirsch index aims to quantify the impact of a publication (or author).

- Quick computation
- Improvement over the simple metric of the number of direct citations
- Reference metric (good metric... on average)

The h -index

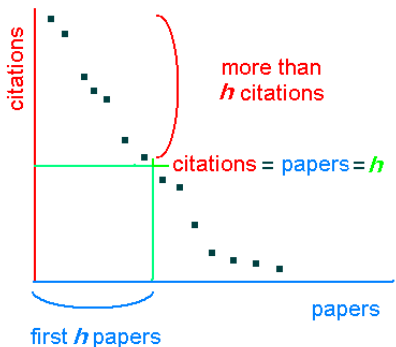


FIGURE – Short exercise : what is the value of h here ?

The Hirsch index aims to quantify the impact of a publication (or author).

- Quick computation
- Improvement over the simple metric of the number of direct citations
- Reference metric (good metric... on average)
- Only consider a small part of the citation network
- Is it actually useful to evaluate (to grade) the agents creating knowledge ?

Strahler's flow

The flow of Strahler is a first simple way to consider the influence (here the river flow) of an agent on the others.

How to extend it to the context of the transmission of knowledge ?

- Information should possibly be evaluated on an ascendant way
- Classic metrics (i.e. direct citation number, *h*-index, etc.) should be compatible
- Algorithmic cost should be reasonable

Size/Speed ratio			
n	n^2	n^3	2^n
10	100	1000	≈ 1000
100	10^4	10^6	$\approx 10^{30}$
1000	10^6	10^9	$\approx 10^{300}$

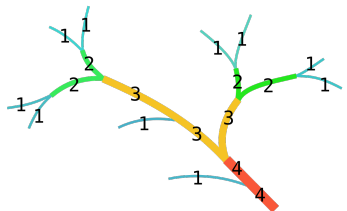


FIGURE – A Strahler flow example to evaluate the influence of each tributary river in the formation of the main river. Strahler number (4 here) measure the complexity of the branching factor of the river.

A general framework

Definition

For a publication p , its **in-neighborhood** $\mathcal{N}^-(p)$ is the set of all the publications referring to p . The size of $\mathcal{N}^-(p)$ is simply its in-degree $d^-(p)$. The corollary implies that $\mathcal{N}^+(p)$, the **out-neighborhood** of p , corresponds to all publications to which p is referring to (with a size $d^+(p)$, its out-degree).

$$K(p) = \begin{cases} \lambda, & \text{if } \mathcal{N}^+(p) = \emptyset \\ F(K(s_1), \dots, K(s_{d^+(p)})), & \text{otherwise,} \end{cases} \quad (1)$$

where λ designates a constant for terminal cases (leafs, often $\lambda = 1$), $s_i \in \mathcal{N}^+(p)$ represents the successors of node p , and F is an application depending on the values $K(s_1), \dots, K(s_{d^+(p)})$. To simplify the notations, we denote $F(\mathcal{N}^+(p)) = F(K(s_1), \dots, K(s_{d^+(p)}))$.

A general framework

Strahler Flow

$$F(\mathcal{N}^-(p)) = \begin{cases} 1, & \text{if } d^-(p) = 0 \\ \max_{q \in \mathcal{N}^-(p)} (K(q)) + \begin{cases} d^-(p) - 1 & \text{if all values } K(q) \text{ are equal} \\ d^-(p) - 2 & \text{otherwise} \end{cases} \end{cases} \quad (1)$$

h-index

$$F(\mathcal{N}^+(p)) = \begin{cases} 0, & \text{if } d^+(p) = 0 \\ \max_{X \subset \mathcal{N}^+(p)} \min_{q \in X} (d^+(q), |X|) \end{cases} \quad (2)$$

Definitions

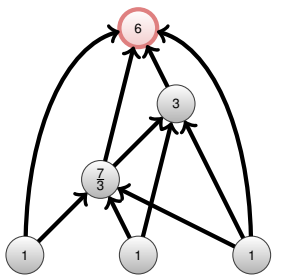
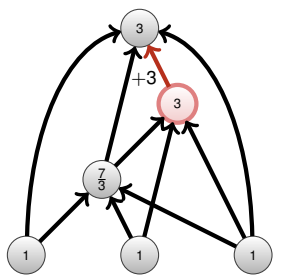
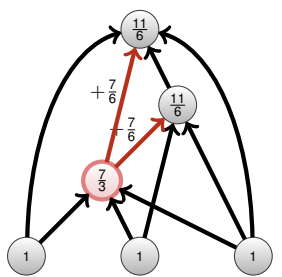
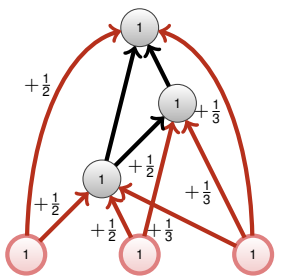
Definition (Related)

Two articles p and q are said to be related if and only if there exist a path from p to q or from q to p . They are k -related if they are related and if the shortest path between them is at most of length k .

Definition (k -diffuse)

A measure of a node p is k -diffuse when it limits its computation to a subgraph composed of the k -related nodes of p

A first straightforward ascending flow



A first straightforward ascending flow

ALGORITHM 1: ascending flow

input : A citation network with nodes (articles) and arcs (citations)

An empty dequeue Q (FIFO)

output: The ascending flow on each node (article) and each arc (citation)

```

1 Initialize each article  $v$  with flow value  $\alpha_v = 1$ 
2 Color each arc in white
3 Add all leaves in  $Q$ 
4 while  $Q$  is not empty do
5    $v \leftarrow pop\_first(Q)$ 
6   for each  $w$  son of  $v$  do
7     Color each  $(v, w)$  in blue
8      $\alpha_w \leftarrow \alpha_w + \alpha_v / d^-(v)$ 
9     if all incoming arcs of  $w$  are blue then
10       $Q \leftarrow push\_last(w)$ 
11    end
12  end
13 end

```

A first straightforward ascending flow

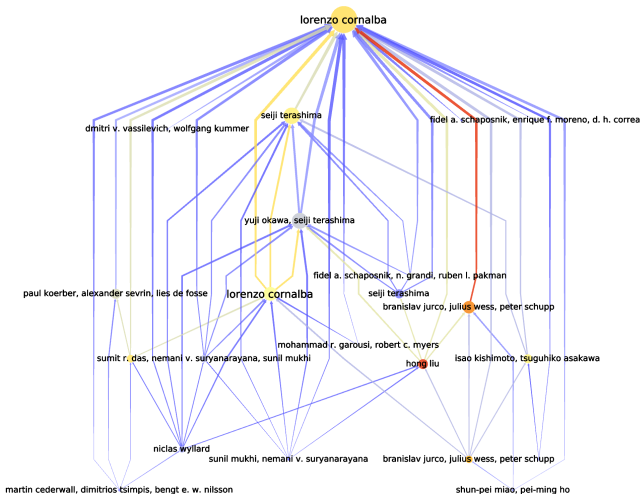
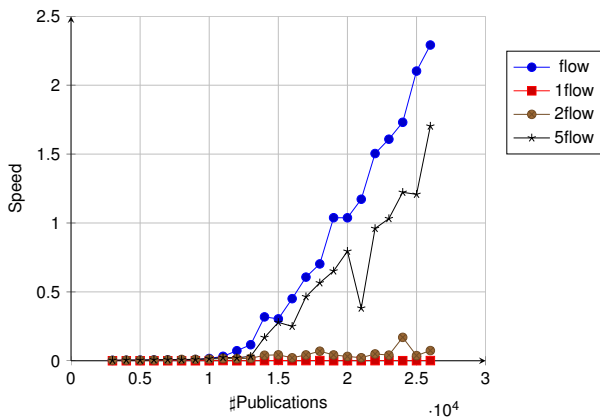


FIGURE – An ascending flow example (publications with a high h -index)

A first straightforward ascending flow



In terms of computation, from $k = 2$, the ranks obtained by the k -flow are $r_s = 0.99$ similar of those of the regular flow so when a gain of computation is needed, one can use k -diffuse version of the algorithm .

Table of contents

- 1 Introduction
- 2 Information Networks
- 3 Flow of knowledge
- 4 Multiplex Networks**
- 5 Cycles, connected component and ethics
- 6 Conclusion

The intrication of data

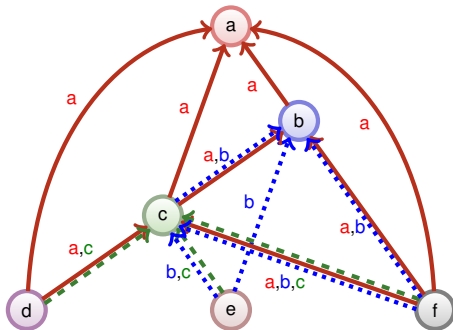


FIGURE – Illustration of the previous citation network transformed to a multiplex network, each citation creates its own layer of interaction.

The intrication of data

Definition

A citation subgraph $G_p = (V_p, E_p)$ is induced by a node p such as $V_p = p \cup \mathcal{N}^-(p) \subset V$, and $E_p \subset E$ is such that, $\forall (q, r) \in V_p^2, \exists a(q, r) \in E$.

Let consider, for each publication p , its induced citation subgraph $G_p = (V_p, E_p)$ of G , the multiplex citation network \mathcal{G} results in combining all individual subgraphs G_p together.

Definition

A multiplex network $\mathcal{G} = (V, \mathcal{E}, \mathcal{L})$ connects nodes (p, q) on different layers l such as arcs $a(p, q, l) \in \bigcup_{l \in \mathcal{L}} \mathcal{E}_l$. A multiplex citation network $\mathcal{G} = (V, \mathcal{E}, \mathcal{L})$ is defined such as $\mathcal{G} = \bigcup_p^{p \in V} G_p$, hence $\mathcal{E} = \bigcup_p^{p \in V} E_p$ and $\mathcal{L} = V$.

Note that an arc $a(p, q, l)$ exists if and only if both p and q cite l or if $l = q$. As a consequence, the multiplex network once “flattened”, has the exact same topology as the original citation network. The difference lies in the multiple edges.

The intrication of data

In our multiplex citation network, the notion of neighborhood remains the same as in the monoplex case. However we can refer to a different notion of multiplex degrees $\delta^+(p)$ and $\delta^-(p)$ that takes into account the number of arcs connecting a node to its neighborhood.

Definition

Denote the multiplex out-degree $\delta^+(p)$ (respectively the multiplex in-degree $\delta^-(p)$) a node p in the multiplex network \mathcal{G} .

$\delta^+(p) = |a(p, q, l)|, \forall q \in V, \forall l \in \mathcal{L}, \text{ s.a. } \exists a(p, q, l) \in \mathcal{E}$ (respectively

$\delta^-(p) = |a(q, p, l)|)$

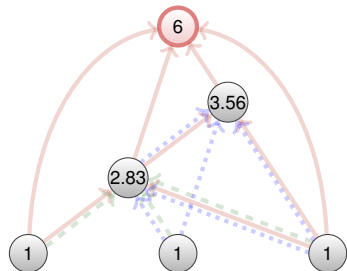
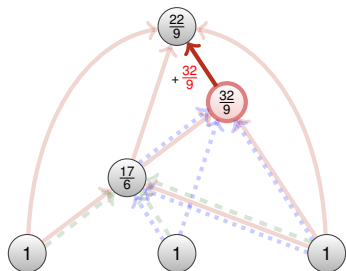
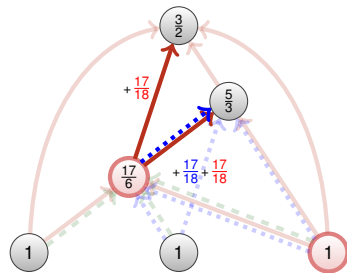
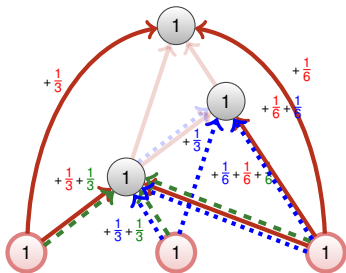
The degrees $d^+(p)$ and $d^-(p)$ still refer to the degree in the monoplex network G , **i.e.** the size of the neighborhood. We then introduce the degrees $d_l^+(p)$ and $d_l^-(p)$ corresponding to the degree in the subgraph G_l , **i.e.** the number of arc adjacent to p on the layer l .

Definition

Denote the layer out-degree $\delta_l^+(p)$ (respectively the layer in-degree $\delta_l^-(p)$) a node p in the subgraph \mathcal{G}_\uparrow . $\delta_l^+(p) = |a(p, q, l)|, \forall q \in V, \text{ s.a. } \exists a(p, q, l) \in \mathcal{E}$ (respectively

$\delta_l^-(p) = |a(q, p, l)|)$

Aggregated Flow



Aggregated Flow

ALGORITHM 2: Aggregated flow (multiplex)

input : A citation network with nodes (articles) and arcs (citations)
 An empty dequeue Q (FIFO)
 A function α_{init} over the nodes

output: The ascending flow on each node (article) and each arc (citation)

```

1 Initialize each article  $p$  with a flow value  $\lambda_p = \lambda_{init}(p)$  (= 1 by default)
2 Color each arc in white
3 Add all leaves in  $Q$ 
4 while  $Q$  is not empty do
5    $p \leftarrow pop\_first(Q)$ 
6   for each  $q$  son of  $p$  do
7     for each layer  $l \in \mathcal{L}_p$  do
8       Color  $a(p, q, l)$  in blue
9        $\lambda_q \leftarrow \lambda_q + \lambda_p / \delta^+(p)$ 
10    end
11    if all incoming arcs of  $q$  are blue then
12       $Q \leftarrow push\_last(q)$ 
13    end
14  end
15 end
```

Aggregated Flow

General framework formulation

$$F(\mathcal{N}^-(p)) = \sum_{q \in \mathcal{N}^-(p)} K(k_q) / \delta^+(k_q) + \lambda_p \quad (3)$$

The complexity of the monoplex ascending flow is $\Theta(m)$ where m is the number of citations. Since the input is a multiplex network, in the worst case the number of links in the networks is equal to the number of nodes (layers) times the number of citations. Thus, the time-complexity of this aggregated flow is $\Theta(mn)$.

Sum Flow

This second extension of the ascending flow to a multiplex network consists in combining multiple monoplex versions of the ascending flow.

Recall the definition of the ascending flow, adapted to a subgraph G_l :

$$F_{G_l}(\mathcal{N}^-(p)) = \sum_{q \in \mathcal{N}^-(p)} K_{G_l}(k_q) / d_{G_l}^+(k_q) + \lambda_{(p,l)} \quad (4)$$

The sum flow will simply sum ascending flows over all subgraphs composed of one layer, determined for a node p by :

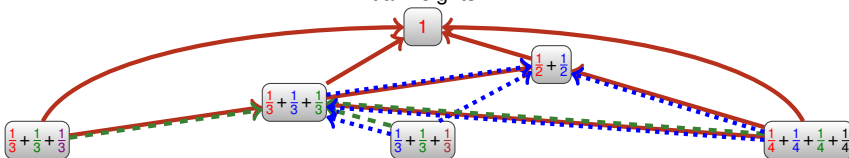
$$F_G(p) = \sum_l^{l \in \mathcal{L}(p)} F_{G_l}(\mathcal{N}^-(p)) \quad (5)$$

$$\lambda_{(p,l)} = \frac{\lambda_p}{|\mathcal{L}(p)|}$$

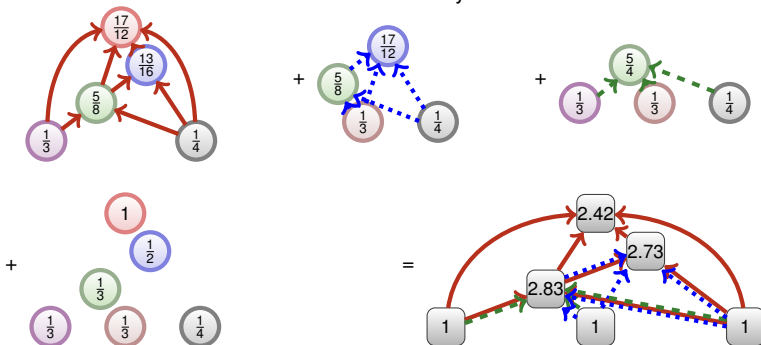
If we set the parameter $\lambda_{(p,l)} = 1$, the contribution of one publication to the whole system will be exactly its number of citations, and a publication that cites lots of work will produce a lot of flow. In order to maintain constant the unit of contribution of a publication of a work, we set $\lambda_p = 1$ hence $\sum_l^{l \in \mathcal{L}(p)} \frac{1}{|\mathcal{L}(p)|} = 1$ such as a publication brings exactly 1 unit of contribution to the system.

Sum Flow

Initial weights



Sum over each layer



Sum Flow

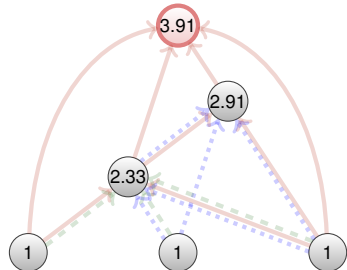
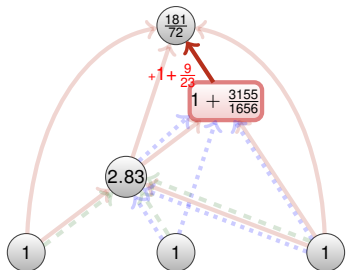
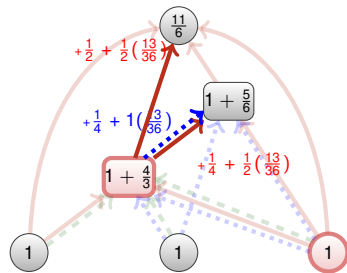
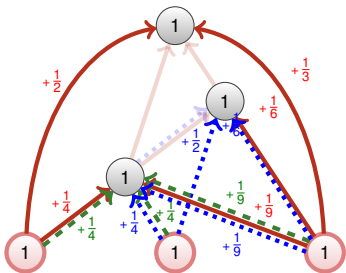
ALGORITHM 3: Sum flow (multiplex)

input : A citation network with nodes (articles) and arcs (citations)

output: The sum flow on each node (article) and each arc (citation)

- 1 Initialize each article p with a flow value $\lambda_p = 0$
 - 2 **for each article** p **do**
 - 3 Construct the citation subgraph G_p
 - 4 $\mathcal{F} \leftarrow \text{ascending-flow}(G_p)$ with $\forall q \in V_p, \lambda_{init}(q) = \frac{1}{d^+(q)+1}$
 - 5 **for each article** $q \in V_p$ **do**
 - 6 $\lambda_q \leftarrow \lambda_q + \mathcal{F}(q)$
 - 7 **end**
 - 8 **end**
-

Selective Flow



Selective Flow

input : A citation network with nodes (articles) and arcs (citations)

An empty dequeue Q (FIFO)

output: The selective flow on each node (article) and each arc (citation)

1 Initialize each article p with the color **white** and the value $\lambda_p = 1$

2 Color all leaves in **red** and push them into Q

3 **while** Q is not empty **do**

4 $q \leftarrow \text{pop_first}(Q)$

5 **for each** r son of q **do**

6 **if** r is white **then**

7 Construct G_r and initialize all $\lambda(a(q, r, l))$ to 0

8 $\text{push_last}(Q, r)$

9 Color r in **red**

10 **end**

11 **end**

12 **for each layer** $l \in L_q$ **do**

13 Initialize $\lambda_{l,q} = 0$

14 **for each parent** p of q **do**

15 $\lambda_{l,q} \leftarrow \lambda_{l,q} + \lambda(a(p, q, l))$

16 **end**

17 **for each son** r of q **do**

18 $\lambda(a(p, q, l)) \leftarrow \frac{1}{d^+(q) \times (1 + |L(q, r)|)} + \frac{\lambda_{l,q}}{d_l^+(q)}$

19 **end**

20 $\lambda_q \leftarrow \lambda_q + \lambda_{l,q}$

21 **end**

22 **end**

Selective Flow

$$F_{G_l}(\mathcal{N}^-(p)) = \lambda_{(p,l)} + \sum_{q \in \mathcal{N}^-(p)} \frac{K_{G_l}(k_q)}{d_{G_l}^+(k_q)}$$
$$F_G(p) = \sum_l F_{G_l}(\mathcal{N}^-(p)) \quad (4)$$
$$\alpha_{(p,l)} = \sum_j \frac{\lambda_p}{d^+(p) \times |\mathcal{L}(p,j)|}$$

Table of contents

- 1 Introduction
- 2 Information Networks
- 3 Flow of knowledge
- 4 Multiplex Networks
- 5 Cycles, connected component and ethics**
- 6 Conclusion

What ethics in influence measures

We can define some criteria to guarantee an ethic that we hope flows naturally

- **Equity** between articles
- Metrics with bounded values in practice
- A network (or sub-network) with **consistent** content : for instance articles from a same field or institute.

These criteria are put to the test by the mere presence of **cycles** within the network. That is, the existence of a couple of nodes (items here) that are descended from each other.

We say that a sub-network is a (strongly) connected component if every vertex is reachable from every other vertex.

A costly solution, decycling (with drawings)

A costly solution, decycling (with drawings)

A costly solution, decycling (with drawings)

Solving the time/accuracy ratio in Information Networks

Conclusion

Solutions to the Big Data Science problems are inherently cross-domain.

(Jean-françois BAFFIER — 2018)