

RCP216

Fouille de graphes et réseaux sociaux

Auteurs: Raphaël Fournier-S'niehotta, Michel Crucianu, Marin Ferecatu
(fournier@cnam.fr, michel.crucianu@cnam.fr, marin.ferecatu@cnam.fr)

Département d'informatique
Conservatoire National des Arts & Métiers, Paris, France

Plan

- 1 Introduction
- 2 Analyse
- 3 Modélisation
- 4 Mesure
- 5 Algorithmique
- 6 GraphX

Plan du cours

1 Introduction

- Expérience de Milgram
- Exemples de réseaux/graphes
- Éléments de théorie des graphes
- Présentation du cours

Plan du cours

1 Introduction

- Expérience de Milgram
- Exemples de réseaux/graphes
- Éléments de théorie des graphes
- Présentation du cours

Expérience de Milgram (1967)

Stanley Milgram (1933-1984), psychologue social américain. Connu notamment pour les expériences de soumission à l'autorité.



- Objectif de l'expérience : faire transiter une lettre de Omaha, NE à Boston, MA

Règle :

- une personne initie la chaîne
- transition de la main à la main à des personnes que l'on connaît, chacune étant supposée se rapprocher de la destination



Expérience de Milgram (1967)

■ Résultats

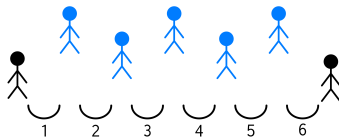
- 44 lettres sur 160 arrivent
- Chemins avec 5 intermédiaires en moyenne.

■ Remarques :

- Chemin interrompu \neq Il n'existe pas de chemin.
- Chemin de longueur $x \neq$ Il n'existe pas de chemin de longueur $< x$

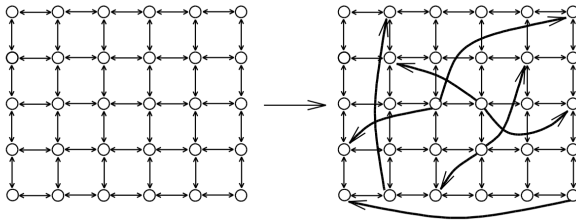
■ Conclusions :

- Il existe des chemins courts.
- Les intermédiaires arrivent à les trouver sans connaissance globale du réseau.



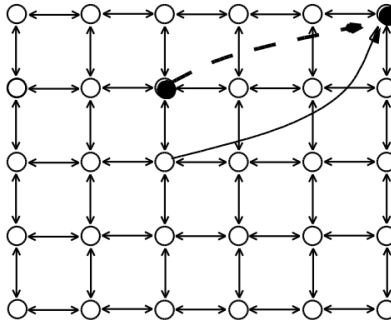
Expérience de Milgram : modélisation

- Objectif : formaliser l'expérience de Milgram
- Travail de D. Watts/S. Strogatz, puis de J. Kleinberg
- Initialement une grille (amis proches).
- On ajoute q voisins quelconques à chaque sommet (amis lointains).



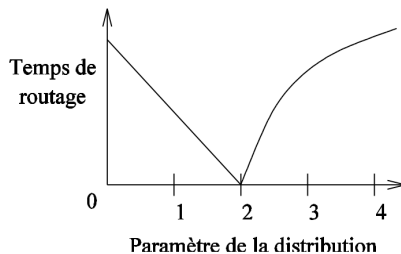
Expérience de Milgram : modélisation

- Un sommet connaît :
 - Sa position, celle de ses voisins, celle de la destination.
 - Il envoie le message à son voisin le plus proche de la destination.



Expérience de Milgram : modélisation

- Un seul lien supplémentaire pour chaque sommet u .
- La destination choisie avec une probabilité dépendant de sa distance à u .
- Dans la majorité des cas, pas de chemins courts



Plan du cours

1 Introduction

- Expérience de Milgram
- Exemples de réseaux/graphes
- Éléments de théorie des graphes
- Présentation du cours

Individus : nombre d'Erdős

Paul Erdős (1913-1996), mathématicien hongrois très prolifique et qui eut plus de 500 collaborateurs directs.



- Graphe de collaboration :
 - Deux scientifiques sont connectés s'ils ont co-écrit un article
 - Chaque scientifique à un nombre d'Erdős :
 - $0 = \text{Erdős}$
 - $1 = \text{collaborateurs d'Erdős}$
 - $2 = \text{collaborateurs de collaborateurs d'Erdős}$
 - Erdős Number project : <http://www.oakland.edu/enp/>
- Récupération de la liste des co-auteurs de tous les articles scientifiques
- Ensuite il ne reste qu'à faire des calculs de plus courts chemins d'Erdős vers les autres chercheurs.

Individus : Kevin Bacon Game

Kevin Bacon (1958–), acteur américain, qui a joué dans plus de 75 films.



- Graphe d'acteurs

- Deux acteurs sont reliés s'ils ont joué dans un même film.

- Distance entre acteurs ?

- <http://oracleofbacon.org/>

- Distance entre Tom Cruise et Clint Eastwood ? 2 (acteur commun entre Space Cowboys et Eyes Wide Shut)

- Distance entre Mickey Mouse et Omar Sy ? 4

- graphes constructible à partir de <http://www.imdb.com/interfaces>

- calculs de plus courts chemins

Individus : possesseurs de fichiers P2P

- Propagation d'un fichier d'utilisateurs en utilisateurs

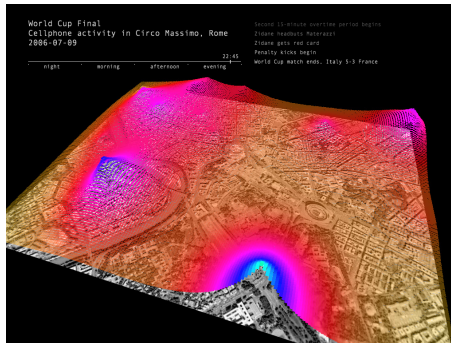
♠ video

- Problèmes et biais de mesure

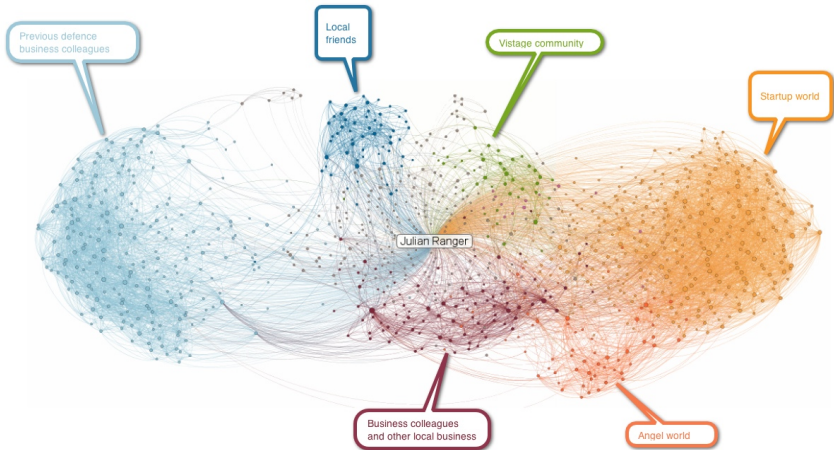
- dynamicité du réseau
- parcours non exhaustif et depuis une source

Individus : communications téléphoniques

- Suivi de communications :
 - Date, heure, durée, type, correspondant
 - Type d'appelant, mobilité, ...
 - <http://senseable.mit.edu>

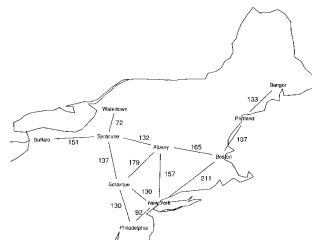


Réseau personnel : LinkedIn Maps



Réseau routier national

- Noeuds : les villes / communes
- Arêtes : (auto)routes
- Valuation possible : distance, ou temps de parcours
- Orientation possible



Des questions :

- quel est le plus court chemin passant par des villes données ?
- quel chemin traverse le moins de villes pour aller d'un point à un autre ?
- peut-on passer par toutes les villes sans passer deux fois par la même route ?
(*voyageur de commerce*)

Autres types de réseaux étudiés

- informatique : pages Web, routeurs, P2P, etc.
- biologie : protéines, neurones cérébraux, etc.
- sciences sociales : amitiés, collaboration, contacts sexuels, etc.
- économie : échanges financiers
- histoire : mariages
- linguistique : synonymie, co-occurrence
- transports : réseau aérien, électrique

Propriétés et problématiques communes

Plan du cours

1 Introduction

- Expérience de Milgram
- Exemples de réseaux/graphes
- Éléments de théorie des graphes
- Présentation du cours

Définitions

Un graphe est défini par un couple $G = (V, E)$ tel que :

- V (pour l'anglais *vertices*) est un ensemble fini de sommets
- E (pour l'anglais *edges*) est un ensemble fini de arêtes

Un graphe peut être orienté, ou non :

- si oui, les couples $(v_i, v_j) \in E$ sont ordonnés, v_i est le sommet initial, v_j est le sommet terminal.
- on appelle alors le couple (v_i, v_j) un *arc*, représenté graphiquement par $v_i \rightarrow v_j$.
- si non, les couples ne sont pas orientés et (v_i, v_j) est équivalent à (v_j, v_i) , et on l'appelle *arête*, représenté par $v_i - v_j$

Terminologie

- l'**ordre** d'un graphe, c'est son nombre de sommets (souvent désigné par n).
- une **boucle** est un arc/une arête reliant un sommet à lui-même
- un graphe dépourvu de boucle est dit **élémentaire**
- un graphe **simple** ne comporte pas de boucle et au plus une arête entre deux sommets
- un graphe **partiel** est le graphe obtenu en supprimant certains arcs ou arêtes
- un **sous-graphe** est le graphe obtenu en supprimant certains sommets et tous les arcs/arêtes incidents aux sommets supprimés.
- un graphe est dit **complet** s'il comporte une arête (v_i, v_j) pour toute paire de sommets $(v_i, v_j) \in E^2$.
- un sommet v_i est dit **adjacent** (familièrement on parle de **voisins**) à un autre s'il existe une arête entre eux.
- le **degré** d'un sommet est le nombre de d'arêtes incidentes à ce sommet.

Plan du cours

1 Introduction

- Expérience de Milgram
- Exemples de réseaux/graphes
- Éléments de théorie des graphes
- **Présentation du cours**

Objectifs

Comprendre le comportement des entités
qui interagissent dans le système étudié, et les lois qui les gouvernent

- On cherche donc :
 - quelle est la structure des graphes
 - quelle est l'évolution de cette structure
 - quels sont les phénomènes reposant sur l'existence de ce réseau

Applications

- Informatique
 - Réseaux : routage, protocoles, sécurité
 - P2P : conception de systèmes, déviations
 - Web : indexation, moteurs de recherche
 - Dessin de graphes
- Sociologie :
 - Diffusion d'innovations, rumeurs
 - Identification de communautés
- Épidémiologie
 - Diffusion de virus, vaccination

Méthodologie

- Outils formels
 - Théorie des graphes
 - Analyse statistique
 - Modélisation probabiliste
- Études expérimentales
 - Simulation
 - Utilisation de données réelles
- Étudier des applications
 - Comprendre en profondeur certains réseaux
 - Extraction de concepts généraux

Ce cours

- Problématiques classées dans 4 grandes catégories :
 - Mesure
 - Comment mesurer les réseaux réels ?
 - Modélisation
 - A quoi ressemblent-ils ?
 - Analyse
 - Peut-on créer des réseaux artificiels similaires ?
 - Algorithmique
 - Comment calculer des choses sur ces grands graphes ?
- Détection de communautés (clustering)
- Réputation, prédiction, innovations et leaders

Plan du cours

2 Analyse

- Propriétés classiques
- Étude de cas

Plan du cours

2 Analyse

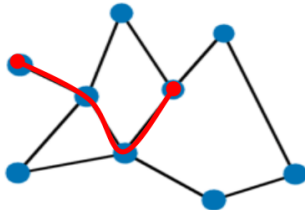
- Propriétés classiques
- Étude de cas

Analyse ?

- Objectifs de l'analyse (statistique) :
 - Description (statistique)
 - Obtenir de l'information pertinente
 - Interprétation des résultats obtenus
- Comment ?
 - Propriétés connues
 - Définition de propriétés (statistiques) pertinentes
 - Corrélations entre ces propriétés
 - Comparaison avec des graphes aléatoires
 - Observation de la croissance des graphes, etc.

Propriétés classiques

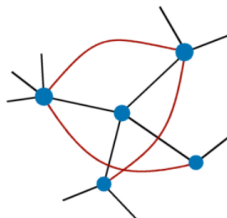
- Longueur des chemins : distance moyenne



Propriétés classiques

■ Clustering

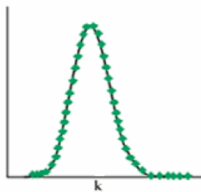
- densité de liens autour d'un nœud
- comparé à la densité globale



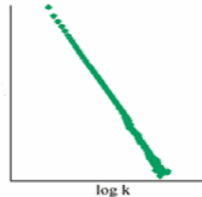
$$c(i) = \frac{2 * |(x,y) \in E, x,y \in N(i)|}{k_i(k_i-1)} \quad (\text{ou } 0 \text{ si } d(i) < 2)$$

Propriétés classiques

- Distribution de degrés
 - Taille ou salaire des individus



$$P_d \sim e^{-\lambda} \cdot \frac{\lambda^d}{d!}$$



$$P_d \sim d^{-\alpha}$$

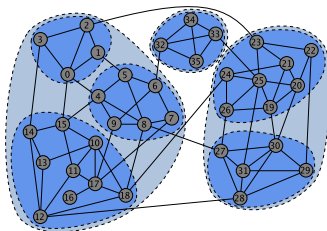
Propriétés classiques

- Composantes connexes
 - Ensemble maximal de sommets tel qu'il existe un chemin entre toute paire de sommets de l'ensemble
 - Graphe connexe = une seule composante connexe

Propriétés classiques

■ Communautés

- ensemble de nœuds très densément liés
- peu de connexion en dehors de l'ensemble



Propriétés classiques

- Autres propriétés
 - Centralité
 - Nombre de plus courts chemins passant par un sommet, etc.
 - Taille des cliques

Propriétés des réseaux réels

- faible densité
- fort clustering
- faible distance moyenne
- distribution de degré fortement hétérogène
- composante géante
- présence de communautés

propriétés différentes de celles des graphes aléatoires

Plan du cours

2 Analyse

- Propriétés classiques
- Étude de cas

Exemple d'analyse : réseau de contacts

- Nombreux équipements avec capacités sans-fil :
 - Ordinateurs, téléphones, PDA, GPS, cartes Navigo
 - Réseaux sans-fils de plus en plus omniprésents
- Contacts physiques ou virtuels permanents :
 - Rencontres physiques, appels téléphoniques, envoi de mails
- Objectifs :
 - Tirer parti des contacts naturels des individus
 - Transmission de l'information de proche en proche
 - Réseau dynamique, non connexe : problèmes de routage ...

Proximité physique ou radio

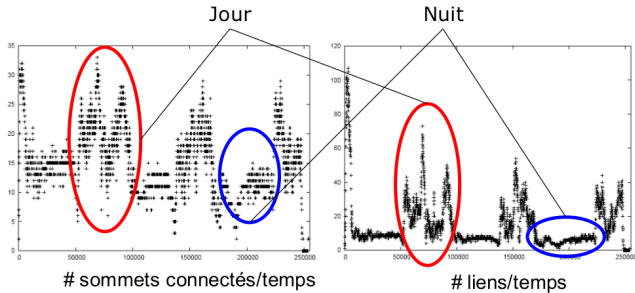
- Quels contacts entre individus ?
 - Physique
 - proximité géographique
 - déplacements
- mesure de la mobilité
 - suivi de déplacements
 - géolocalisation : coûteux, dur à mettre en uvre
 - équipement de chaque individu
- application informatique/télécom : déploiement de réseau dans des environnements "hostiles" (zones militaire, forêts)
- Étude de cas
 - 41 capteurs pendant 3 jours
 - propriétés dynamique du réseau

Étude de cas

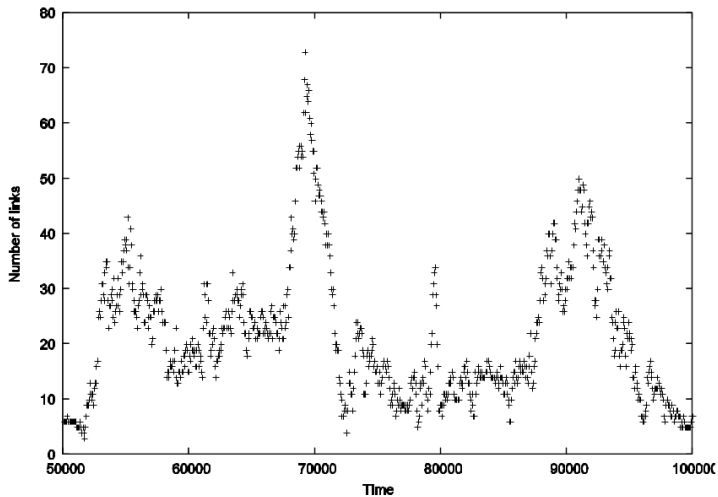
- Conférence INFOCOM 2005, dans un hôtel à Miami (USA)
- 54 capteurs Bluetooth initialement (perte, pannes)
- Fonctions :
 - recherche de contact (5s)
 - attente (110s env)
 - pas de géolocalisation
- données
 - ensemble de liens à chaque instant
 - liens non symétriques
 - <http://plausible.lip6.fr>

Étude de cas

- Effets sociologiques :
 - jour/nuit, repas, pauses, etc.
 - beaucoup de petites variations
 - 50% de sommets isolés
 - max 34 sommets connectés

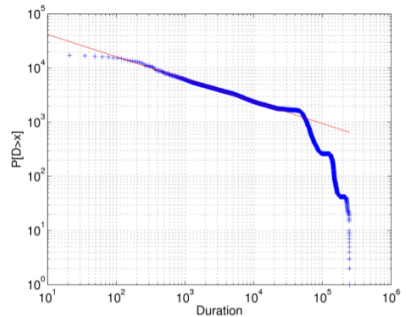
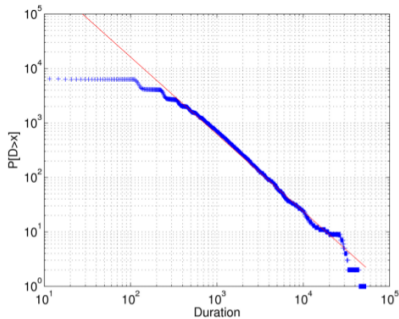


Étude de cas

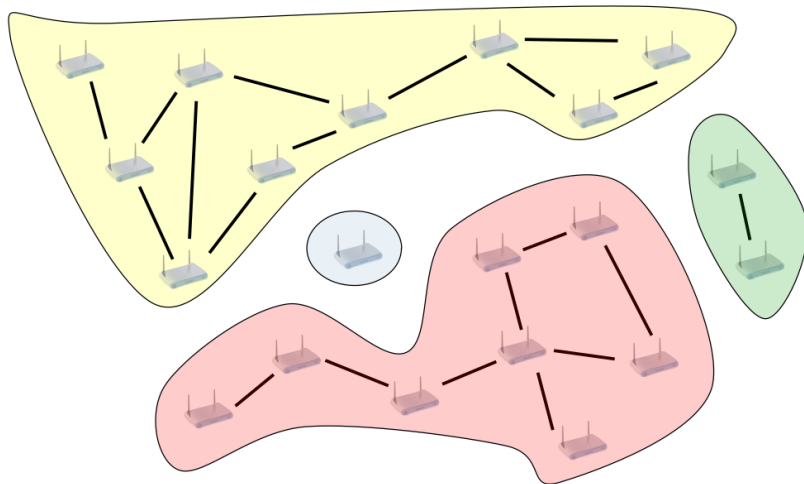


Durée de contacts

- distribution en loi de puissance
- certains liens sont fréquents, d'autres pas
- liens non fréquents pour atteindre des zones spécifiques

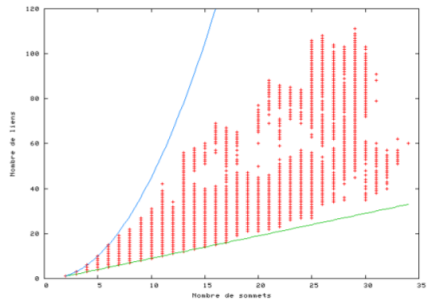
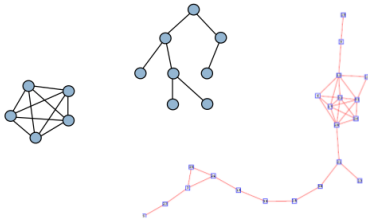


Composantes connexes



Composantes connexes

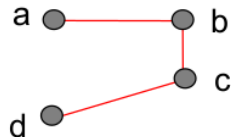
- Petites composantes : densité variable.
- Grosse composantes : faible densité ($\max(nb_liens) \sim 4.5 \times nb_sommets$)



Approche fouille de données

- Graphe dynamique (liens x temps)
 - Rectangles maximaux de 1
 - Calcul exhaustif ?
 - Graphes fréquents : seuils sur la durée.
 - Graphes significatifs : seuils sur le nombre de liens

	t1	t2	t3	t4	t5	t6
a-b	1	1	1	1	0	0
a-c	0	0	0	0	0	0
a-d	0	0	0	1	1	1
b-c	1	1	1	1	0	0
b-d	0	0	0	1	1	1
c-d	1	1	1	1	1	1



Plan du cours

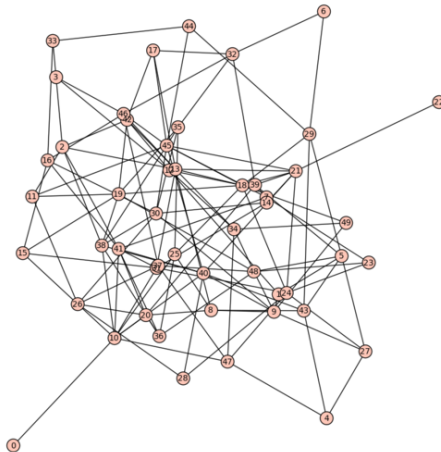
- 3 Modélisation
 - Application : robustesse

Générer des graphes réalistes

- Est-ce que les propriétés observées sur les graphes réels sont “normales”
 - On peut comparer avec un graphe aléatoire ayant certaines propriétés
- Simulation de phénomènes (attaques, diffusion, etc.)
- Évaluation de protocoles
- Compréhension
- Prévion

Tout aléatoire

- Créer n sommets/nœuds
- Ajouter au hasard m liens ($m \leq n^2$)



Propriété attendue

- Graphe aléatoire, $n = m = 4950$
- Graphe réel : clique de 100 sommets, autres noeuds de degré 0
- Probable ?
 - proba degré 0 : $p = (1 - \frac{2}{n})^n \sim 0.14$
 - on attend donc : $n \times p \sim 683$ sommets de degré 0
 - graphe réel peu probable

Propriétés observées

- densité fixée
- Connexité : composante géante de taille $O(n)$
- Distance moyenne, diamètre $\sim \log(n)$
- Distribution des degrés homogène
- Clustering proche de 0
- Pas de structure communautaire

Basé sur la distribution de degrés

- Attachement préférentiel
 - ajout de sommets un à un
 - ajout de lien à des sommets déjà connectés
- Modèle configurationnel (*configuration model*)
 - on prend n sommets
 - on fixe le degré de chaque sommet
 - on ajoute des liens au hasard en respectant les degrés
- ne génèrent pas de clustering

Basé sur le Coefficient de clustering

- Mélanger un graphe très rigide :
 - Donne du clustering et une distance moyenne courte
 - Ne donne pas de degrés hétérogènes !



régulier

$$p = 0$$



$$p = 0.25$$



$$p = 0.5$$



$$p = 0.75$$



aléatoire

$$p = 1$$

Plan du cours

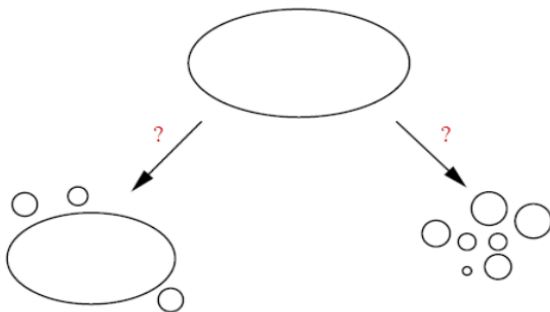
- 3 Modélisation
 - Application : robustesse

Application : robustesse

- Étude des phénomènes visant des sommets :
 - Internet : pannes ou attaques sur routeurs.
 - Réseaux sociaux : maladies, rumeurs,
 - Échanges de-mails : virus informatiques.
- Deux types d'atteintes
 - Pannes : aléatoires.
 - Attaques : ciblées.
- But : Comprendre ces phénomènes pour pouvoir :
 - Prédire.
 - Construire des stratégies d'attaque/défense.

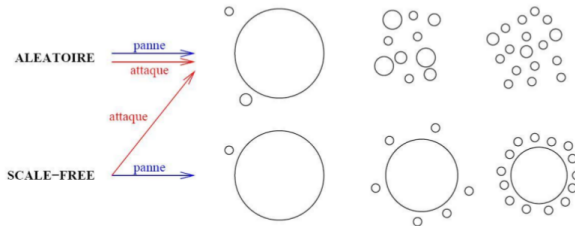
Impact d'une panne/attaque

- Critères :
 - Basés sur la distance.
 - Tailles des composantes connexes.
 - etc.



Résultats

- Suppression :
 - Panne = aléatoire
 - Attaque = ciblée (plus fort degré d'abord)
- Question : qui vacciner pour limiter une épidémie ?



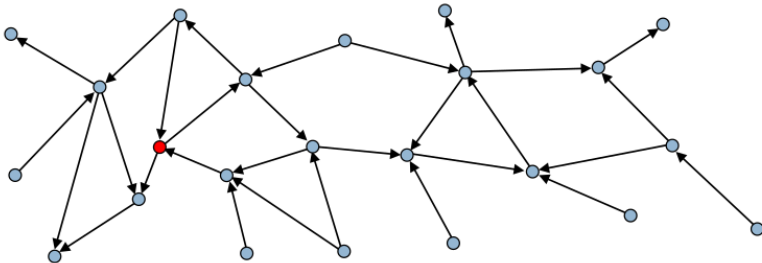
Plan du cours

4 Mesure

- Métrologie : exemple de l'Internet

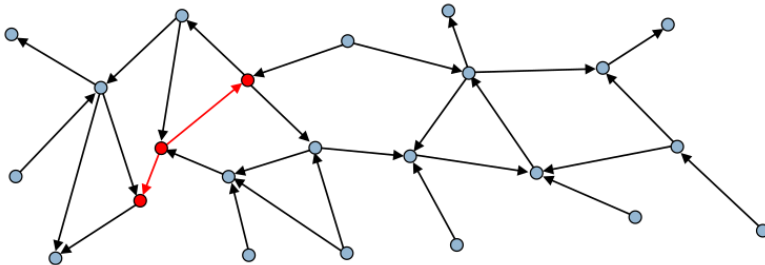
Mesure de l'Internet

- Processus de mesure par parcours en largeur depuis plusieurs sources
- Réseau : orienté, non connexe, dynamique



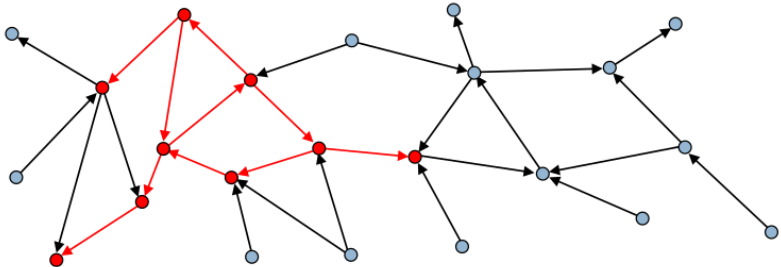
Mesure de l'Internet

- Processus de mesure par parcours en largeur depuis plusieurs sources
- Réseau : orienté, non connexe, dynamique



Mesure de l'Internet

- Processus de mesure par parcours en largeur depuis plusieurs sources
- Réseau : orienté, non connexe, dynamique

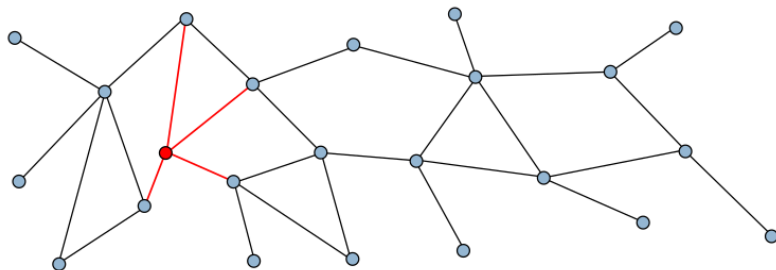


Mesure de réseaux sociaux

Processus de mesure :

- Réseau égo-centrés
- Listes de diffusion, communautés

Réseau : orienté, non connexe, dynamique

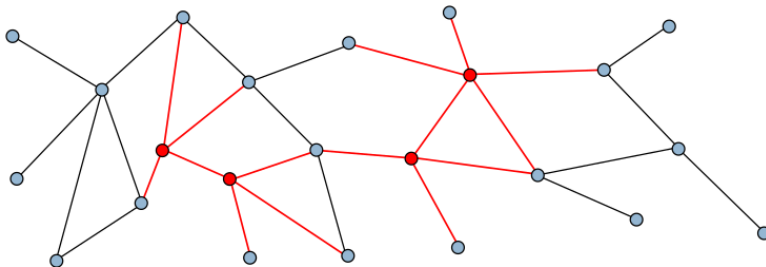


Mesure de réseaux sociaux

Processus de mesure :

- Réseau égo-centrés
- Listes de diffusion, communautés

Réseau : orienté, non connexe, dynamique

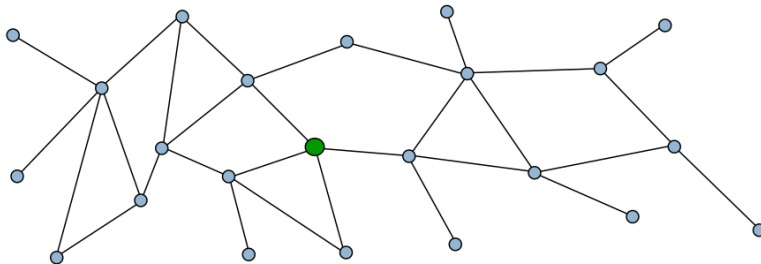


Métrologie des réseaux d'échanges

Processus de mesure

- trafic passant par un sommet

Réseau orienté, pondéré

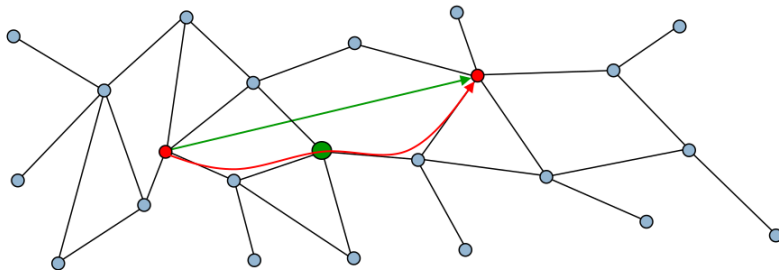


Métrologie des réseaux d'échanges

Processus de mesure

- trafic passant par un sommet

Réseau orienté, pondéré

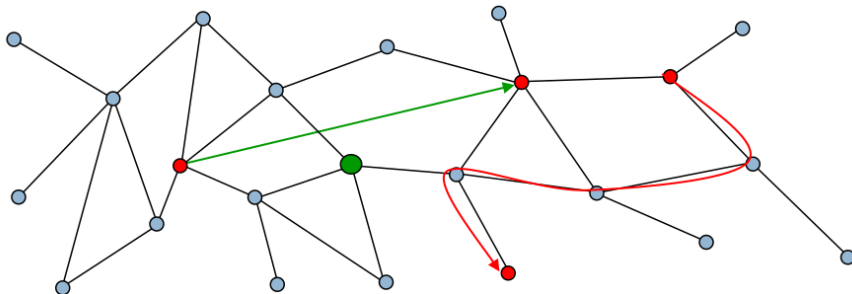


Métrologie des réseaux d'échanges

Processus de mesure

- trafic passant par un sommet

Réseau orienté, pondéré

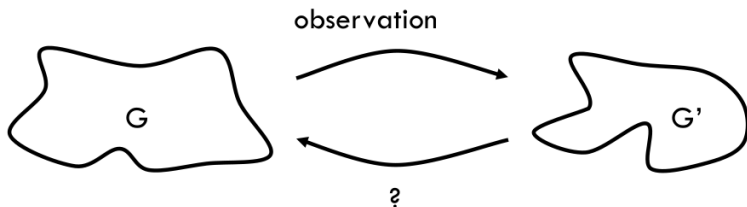


Métrologie

- En général : impossibilité d'étudier l'objet réel, seulement une mesure
- Questions :
 - qui a fait la mesure ?
 - quelle proportion a été mesurée ?
 - combien de temps la mesure a-t-elle duré ?
 - quelles étaient les contraintes / biais ?
 - la mesure peut-elle être reproduite ?

Métrologie

- Étude du biais introduit par l'observation
- Que dire de l'objet réel à partir de l'observation ?
- Nouveaux protocoles de mesures, etc.



- Évaluer la représentativité des "cartes"

Une approche

- On simule la mesure sur un graphe aléatoire
- Modélisation du processus de mesure :
 - Internet : traceroute = chemins courts
 - Web : crawl = parcours en largeur
- Modélisation du réseau :
 - Graphes aléatoires
 - Respect des degrés, du clustering, etc.

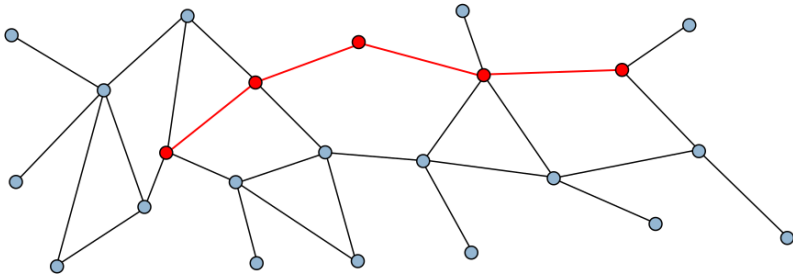
Plan du cours

4 Mesure

- Métrologie : exemple de l'Internet

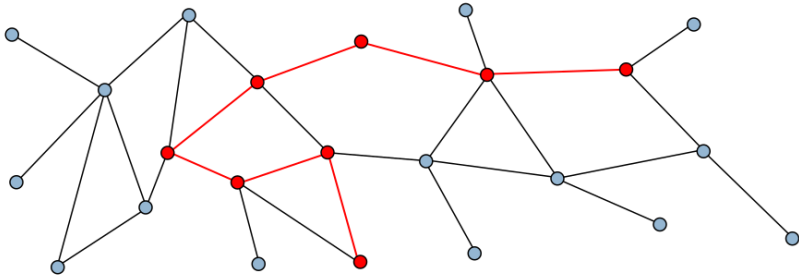
Métrologie de l'Internet

- Processus de mesure :
 - Traceroute, plus courts chemins de plusieurs sources vers plusieurs destinations
- Réseau : (non) orienté, pondéré (RTT,...)



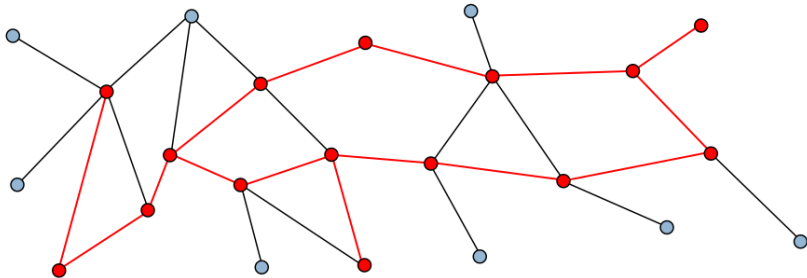
Métrologie de l'Internet

- Processus de mesure :
 - Traceroute, plus courts chemins de plusieurs sources vers plusieurs destinations
- Réseau : (non) orienté, pondéré (RTT,...)



Métrologie de l'Internet

- Processus de mesure :
 - Traceroute, plus courts chemins de plusieurs sources vers plusieurs destinations
- Réseau : (non) orienté, pondéré (RTT,...)

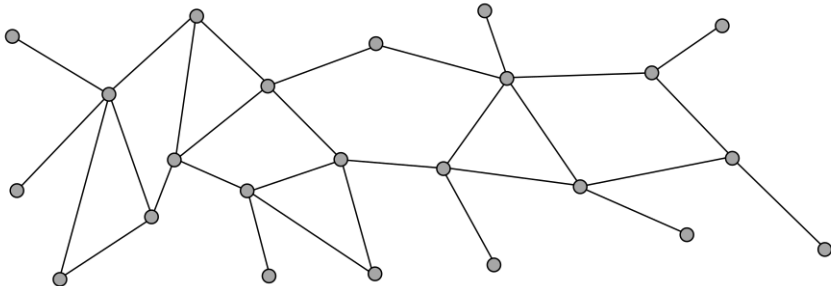


Questions

- Influence sur le résultat de :
 - Nombre de sources et destinations
 - Propriétés du réseau
 - Localisation des sources et destinations
- Modélisation :
 - Traceroute = plus courts chemins (un ou tous)
 - Graphe = graphe aléatoire (modèle à choisir)

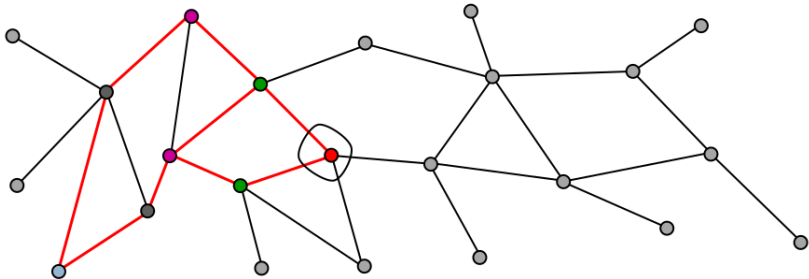
Que voit-on ?

- D'une source vers tout le monde
 - liens rouges découverts (sur plus courts chemins)
 - on répète pour les autres destinations
 - liens noirs invisibles



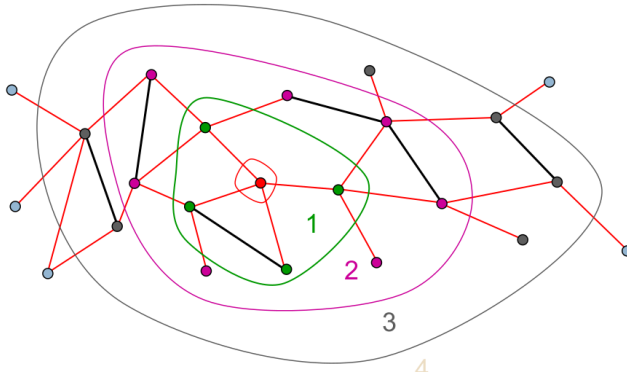
Que voit-on ?

- D'une source vers tout le monde
 - liens rouges découverts (sur plus courts chemins)
 - on répète pour les autres destinations
 - liens noirs invisibles



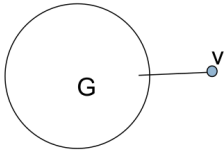
Que voit-on ?

- D'une source vers tout le monde
 - liens rouges découverts (sur plus courts chemins)
 - on répète pour les autres destinations
 - liens noirs invisibles

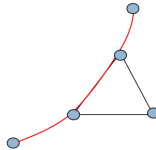


Zones dures à mesurer

- Sommet de degré 1 : uniquement visible si source ou destination

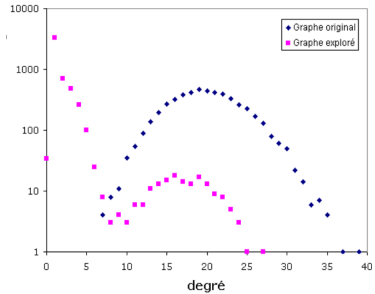
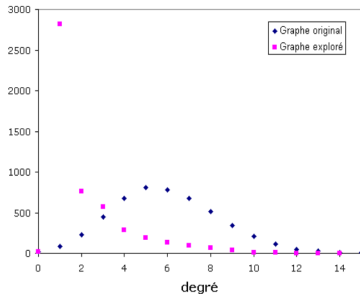


- graphe complet : visiter tous les liens



Distribution de degrés

- différences entre original et mesuré
 - beaucoup de sommets de faible degré
 - peu de sommets de fort degré
- mauvaise estimation de la propriété réelle



Plan du cours

5 Algorithmique

Besoin d'algorithmes spécifiques

- Gros problème = taille :
 - Internet = Millions de sommets (routeurs)
 - Facebook = plus de 800 millions d'utilisateurs actifs
 - Web = Google connaît plus de 1000 milliards d'URL distinctes
- **il est non trivial** de
 - stocker le graphe en mémoire
 - faire des calculs sur le graphe

Exemples

- Compter les triangles d'un graphe (clustering) :
 - naïvement $O(n^3)$
 - $O(m * n^{1/a})$ si distribution des degrés en loi de puissance d'exposant a .
- Diamètre :
 - complexité théorique : $O(nm)$
 - approximation en $O(m)$
- Problèmes NP-complets
- Beaucoup de problèmes spécifiques aux graphes réels (détection de communautés).
Approximation (non prouvée) linéaire.

Plan du cours

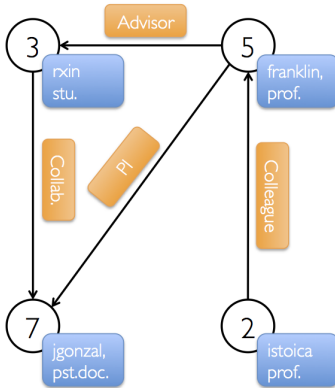
6 GraphX

GraphX

- librairie de Spark pour gérer les calculs sur les graphes
- en particulier, le parallélisme
- introduit une abstraction Graph (au-dessus de RDD) :
 - un multigraphe orienté, avec des propriétés attachées à chaque sommet et chaque arête
 - facilite les cas où il y a plusieurs arêtes entre des noeuds
- <https://spark.apache.org/docs/latest/graphx-programming-guide.html>

GraphX

Property Graph



Vertex Table

Id	Property (V)
3	(rxin, student)
7	(jgonzal, postdoc)
5	(franklin, professor)
2	(istoica, professor)

Edge Table

SrcId	DstId	Property (E)
3	7	Collaborator
5	3	Advisor
2	5	Colleague
5	7	PI

GraphX

```
val sc: SparkContext
// Create an RDD for the vertices
val users: RDD[(VertexId, (String, String))] =
  sc.parallelize(Array((3L, ("rxin", "student")), (7L, ("jgonzal", "postdoc")),
    (5L, ("franklin", "prof")), (2L, ("istoica", "prof"))))
// Create an RDD for edges
val relationships: RDD[Edge[String]] =
  sc.parallelize(Array(Edge(3L, 7L, "collab"), Edge(5L, 3L, "advisor"),
    Edge(2L, 5L, "colleague"), Edge(5L, 7L, "pi")))
// Define a default user in case there are relationship with missing user
val defaultUser = ("John Doe", "Missing")
// Build the initial Graph
val graph = Graph(users, relationships, defaultUser)
```

GraphX : opérateurs

```
val graph: Graph[(String, String), String]  
// Use the implicit GraphOps.inDegrees operator  
val inDegrees: VertexRDD[Int] = graph.inDegrees
```

D'autres opérateurs :

- numEdges/numVertices
- collectNeighbors
- subgraph
- connectedComponents
- triangleCount

Références

- Ce cours repose sur les travaux de :
 - l'équipe ComplexNetworks du LIP6 (UPMC), <http://www.complexnetworks.fr> (membres passés et présents)
 - en particulier les cours de Jean-Loup Guillaume (PR, U. de La Rochelle) et de Clémence Magnien
 - le livre *Mining Massive datasets* (<http://www.mmms.org>), de Jure Leskovec, Anand Rajaraman, Jeff Ullman