# Le point de vue d'un algorithmicien sur Science des Données

Michel Habib habib@liafa.univ-paris-diderot.fr http://www.liafa.univ-paris-diderot.fr/~habib

Réunion des 4 Pôles USPC, 30 novembre 2015

Le point de vue d'un algorithmicien sur Science des Données

# Plan de l'exposé

Remarques liminaires

Remarques liminaires

Une évolution naturelle de l'informatique

Remarques liminaires

Une évolution naturelle de l'informatique

La rupture

Remarques liminaires

Une évolution naturelle de l'informatique

La rupture

Statistiques sur les graphes

Remarques liminaires

Une évolution naturelle de l'informatique

La rupture

Statistiques sur les graphes

Conclusions

Le point de vue d'un algorithmicien sur Science des Données Lemarques liminaires

#### Remarques liminaires

Une évolution naturelle de l'informatique

La rupture

Statistiques sur les graphes

Conclusions

Le point de vue d'un algorithmicien sur Science des Données — Remarques liminaires

➤ Science des données : éviter les recettes de cuisine et les effets de mode, trouver les concepts durables.

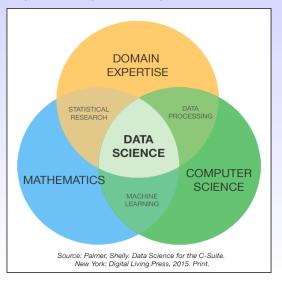
- ► Science des données : éviter les recettes de cuisine et les effets de mode, trouver les concepts durables.
- i.e., l'apprentissage de technologies trop rapidement périssables
   (cf. en micro-électronique selon pôle emploi, un ingénieur

nouvellement formé n'est compétent que pendant 5 ans)

- Science des données : éviter les recettes de cuisine et les effets de mode, trouver les concepts durables.
- i.e., l'apprentissage de technologies trop rapidement périssables
   (cf. en micro-électronique selon pôle emploi, un ingénieur nouvellement formé n'est compétent que pendant 5 ans)
- L'outil est éphémère le concept est durable ...

Remarques liminaires

### Une définition pas très pluridisciplinaire



Le point de vue d'un algorithmicien sur Science des Données Lune évolution naturelle de l'informatique

Remarques liminaires

Une évolution naturelle de l'informatique

La rupture

Statistiques sur les graphes

Conclusions

Le point de vue d'un algorithmicien sur Science des Données Lune évolution naturelle de l'informatique

> S'adapter aux données dont les tailles grandissent vertigineusement

Le point de vue d'un algorithmicien sur Science des Données

Lune évolution naturelle de l'informatique

- S'adapter aux données dont les tailles grandissent vertigineusement
- Cela fait 70 ans que les informaticiens le font.

### Un disque de 5 méga-octets en 1956



Le point de vue d'un algorithmicien  $\mbox{ sur Science des Données } \bigsqcup_{\mbox{ } \mbox{ } \$ 

Remarques liminaires

Une évolution naturelle de l'informatique

La rupture

Statistiques sur les graphes

Conclusions

Le point de vue d'un algorithmicien  $\,$  sur Science des Données  $\,$   $\,$   $\,$   $\,$  La rupture

► Toutefois la croissance exponentielle de la collecte des données a changé la donne.

Le point de vue d'un algorithmicien sur Science des Données La rupture

- ► Toutefois la croissance exponentielle de la collecte des données a changé la donne.
- Il faut une nouvelle informatique pour traiter ces données gigantesques.

Le point de vue d'un algorithmicien sur Science des Données La rupture

### Le blues d'un agent de la NSA

Avalanche de données sur l'internet en 1s :

▶ 220 000 messages sur Whatsapp

- 220 000 messages sur Whatsapp
- 23 000 minutes de vidéos échangées sur Skype

- 220 000 messages sur Whatsapp
- 23 000 minutes de vidéos échangées sur Skype
- ▶ 52 000 likes sur Facebook

- 220 000 messages sur Whatsapp
- 23 000 minutes de vidéos échangées sur Skype
- ▶ 52 000 likes sur Facebook
- On peut parler de masse de données.

- 220 000 messages sur Whatsapp
- 23 000 minutes de vidéos échangées sur Skype
- ▶ 52 000 likes sur Facebook
- On peut parler de masse de données.
- Sans compter les messages électroniques, toutes les images des caméras et autres satellites

Le point de vue d'un algorithmicien sur Science des Données La rupture

# Ce qu'il faudrait développer tant en recherche qu'en enseignement

 De très grandes machines parallèles (peu de recherche en informatique sur ce point au sein d'USPC (C. Cerin P13), mais sujet présent à INRIA)

# Ce qu'il faudrait développer tant en recherche qu'en enseignement

- De très grandes machines parallèles (peu de recherche en informatique sur ce point au sein d'USPC (C. Cerin P13), mais sujet présent à INRIA)
- Des moyens de communication pour échanger ces données (les débits de communication ne suivent pas la croissance des données)

# Ce qu'il faudrait développer tant en recherche qu'en enseignement

- De très grandes machines parallèles (peu de recherche en informatique sur ce point au sein d'USPC (C. Cerin P13), mais sujet présent à INRIA)
- Des moyens de communication pour échanger ces données (les débits de communication ne suivent pas la croissance des données)
- ▶ De nouvelles méthodes de stockage, de compression et d'accès à ces données (par exemple données distribuées accessibles à l'aide de Hadoop, Map Reduce) Peu d'expérience sur ce sujet dans le monde académique.

# Ce qu'il faudrait développer tant en recherche qu'en enseignement

- De très grandes machines parallèles (peu de recherche en informatique sur ce point au sein d'USPC (C. Cerin P13), mais sujet présent à INRIA)
- Des moyens de communication pour échanger ces données (les débits de communication ne suivent pas la croissance des données)
- De nouvelles méthodes de stockage, de compression et d'accès à ces données (par exemple données distribuées accessibles à l'aide de Hadoop, Map Reduce) Peu d'expérience sur ce sujet dans le monde académique.
- Les méthodes d'analyse des données : fouille de données (Themis Palpanas P5) et l'apprentissage (Yann Chevaleyre P13, Mohamed Nadif P5)

Le point de vue d'un algorithmicien sur Science des Données La rupture

### Mais aussi et la liste n'est pas ordonnée

Les méthodes de visualisation de données

Le point de vue d'un algorithmicien sur Science des Données La rupture

- Les méthodes de visualisation de données
- L'algorithmique, telle par exemple l'algorithmique à la volée (streaming, Frédéric Magniez au LIAFA)

- Les méthodes de visualisation de données
- L'algorithmique, telle par exemple l'algorithmique à la volée (streaming, Frédéric Magniez au LIAFA)
- La protection des données : les méthodes de cryptographie et d'anonymisation (est-ce vraiment possible?)

- Les méthodes de visualisation de données
- L'algorithmique, telle par exemple l'algorithmique à la volée (streaming, Frédéric Magniez au LIAFA)
- ► La protection des données : les méthodes de cryptographie et d'anonymisation (est-ce vraiment possible ?)
- Les libertés individuelles (Pourquoi Google archive tout sur nous?) Moi j'ai plein de choses à cacher.

- Les méthodes de visualisation de données
- ► L'algorithmique, telle par exemple l'algorithmique à la volée (streaming, Frédéric Magniez au LIAFA)
- ► La protection des données : les méthodes de cryptographie et d'anonymisation (est-ce vraiment possible?)
- Les libertés individuelles (Pourquoi Google archive tout sur nous?) Moi j'ai plein de choses à cacher.
- Quels sont les modèles économiques compatibles avec la démocratie?

- Les méthodes de visualisation de données
- L'algorithmique, telle par exemple l'algorithmique à la volée (streaming, Frédéric Magniez au LIAFA)
- ► La protection des données : les méthodes de cryptographie et d'anonymisation (est-ce vraiment possible?)
- Les libertés individuelles (Pourquoi Google archive tout sur nous?) Moi j'ai plein de choses à cacher.
- Quels sont les modèles économiques compatibles avec la démocratie?

### Big Brother est déjà installé



Le point de vue d'un algorithmicien sur Science des Données La rupture

### Pourquoi je me sens concerné par le sujet

 Spécialiste d'algorithmique de graphes je m'intéresse aux algorithmes sur les très grands graphes (par ex. graphe des amis de Facebook)

### Pourquoi je me sens concerné par le sujet

- Spécialiste d'algorithmique de graphes je m'intéresse aux algorithmes sur les très grands graphes (par ex. graphe des amis de Facebook)
- Calcul de diamètre et centres de grands graphes https://who.rocq.inria.fr/Laurent.Viennot/road/

## Pourquoi je me sens concerné par le sujet

- Spécialiste d'algorithmique de graphes je m'intéresse aux algorithmes sur les très grands graphes (par ex. graphe des amis de Facebook)
- Calcul de diamètre et centres de grands graphes https://who.rocq.inria.fr/Laurent.Viennot/road/
- Participation au projet ANR Algopol avec le sociologue Dominique Cardon (Orange)

## Pourquoi je me sens concerné par le sujet

- Spécialiste d'algorithmique de graphes je m'intéresse aux algorithmes sur les très grands graphes (par ex. graphe des amis de Facebook)
- Calcul de diamètre et centres de grands graphes https://who.rocq.inria.fr/Laurent.Viennot/road/
- Participation au projet ANR Algopol avec le sociologue Dominique Cardon (Orange)
- 2 cours en Master Informatique P7 sur des sujets Data Science Grands Réseaux d'interaction, Méthodes et Algorithmes pour l'accès à l'Information Numérique (Moteur de recherche, algorithmes de recommandation . . . )



Remarques liminaires

Une évolution naturelle de l'informatique

La rupture

Statistiques sur les graphes

Conclusions

Echantillonnage de graphes, "Property testing" = statistiques sur des graphes plutôt que sur des nombres (points forts actuels : US et Israel).

- ► Echantillonnage de graphes, "Property testing" = statistiques sur des graphes plutôt que sur des nombres (points forts actuels : US et Israel).
- ► Limite de graphes : un sujet difficile quelques résultats fondamentaux : Livre de L. Lovász

- ► Echantillonnage de graphes, "Property testing" = statistiques sur des graphes plutôt que sur des nombres (points forts actuels : US et Israel).
- ▶ Limite de graphes : un sujet difficile quelques résultats fondamentaux : Livre de L. Lovász
- Lemme de régularité de E. Szemerédi, Prix Abel 2012.

- ► Echantillonnage de graphes, "Property testing" = statistiques sur des graphes plutôt que sur des nombres (points forts actuels : US et Israel).
- ▶ Limite de graphes : un sujet difficile quelques résultats fondamentaux : Livre de L. Lovász
- Lemme de régularité de E. Szemerédi, Prix Abel 2012.
- Une théorie des graphes clairsemés ou peu denses, P. Ossona de Mendes et J. Nesestril (Collaboration avec Prague, LEA CNRS)

- ► Echantillonnage de graphes, "Property testing" = statistiques sur des graphes plutôt que sur des nombres (points forts actuels : US et Israel).
- ► Limite de graphes : un sujet difficile quelques résultats fondamentaux : Livre de L. Lovász
- Lemme de régularité de E. Szemerédi, Prix Abel 2012.
- Une théorie des graphes clairsemés ou peu denses, P. Ossona de Mendes et J. Nesestril (Collaboration avec Prague, LEA CNRS)
- Modèles de graphes aléatoires (apports récents très important des physiciens)

Le point de vue d'un algorithmicien  $\,$  sur Science des Données  $\,$   $\,$   $\,$  Conclusions

Remarques liminaires

Une évolution naturelle de l'informatique

La rupture

Statistiques sur les graphes

Conclusions

# Science des données : un sujet pluridisciplinaire, cf. Sergueï Brin et Larry Page (Google)

 Sociologie: Ils se sont inspiré d'une idée de sociologues sur les indices de citations (G. Pinsky et F. Narin 1976), qu'ils ont appliqué aux pages web.

# Science des données : un sujet pluridisciplinaire, cf. Sergueï Brin et Larry Page (Google)

- Sociologie: Ils se sont inspiré d'une idée de sociologues sur les indices de citations (G. Pinsky et F. Narin 1976), qu'ils ont appliqué aux pages web.
- Mathématiques: Ils ont formalisé cette idée et l'ont ramené à un calcul de valeurs propres, ont croisé le théorème de Perron-Frobenius (1900), les chaînes de Markov.

# Science des données : un sujet pluridisciplinaire, cf. Sergueï Brin et Larry Page (Google)

- Sociologie: Ils se sont inspiré d'une idée de sociologues sur les indices de citations (G. Pinsky et F. Narin 1976), qu'ils ont appliqué aux pages web.
- Mathématiques: Ils ont formalisé cette idée et l'ont ramené à un calcul de valeurs propres, ont croisé le théorème de Perron-Frobenius (1900), les chaînes de Markov.
- Informatique: Puis ils se sont construit un système informatique distribué, efficace et adapté.

Le point de vue d'un algorithmicien sur Science des Données La Conclusions

#### Deux notes optimistes

 Les mathématiciens et informaticiens de P7 sont prêts à avoir un rôle moteur dans un projet de master Science des données USPC.

#### Deux notes optimistes

- Les mathématiciens et informaticiens de P7 sont prêts à avoir un rôle moteur dans un projet de master Science des données USPC.
- 2. Dans le cadre des projets USPC, nous avons obtenu une allocation pluridisciplinaire (Math-Info) de doctorat pour étudier avec E. Birmele (P5): Les algorithmes de classification de grandes masses de graphes avec comme premier jeu de données des graphes venant de la biologie.