

Que révèlent vos requêtes, votre réseau social et vos mouvements à votre propos?

Sébastien Gambs

Chaire de recherche en sécurité des systèmes d'information
(Université de Rennes 1 - INRIA)
IRISA

28 septembre 2010

Introduction

Protection de la vie privée

Le **droit au respect de la vie privée** est un des droits fondamentaux des individus.

- ▶ Déclaration universelle des droits de l'homme des Nations Unies (article 12), 1948.
- ▶ Directive européenne 95/46/EC sur la protection des données personnelles.

Un des principaux défis de la "Société de l'information".

Exemples :

- ▶ Adresse IP \Rightarrow localisation, identifiant, contenu.
- ▶ Historique de requêtes \Rightarrow centre d'intérêts.
- ▶ Connaissance du réseau social \Rightarrow inférences sur opinions politiques, religion, hobbies, ...

Danger : utilisation de cette information à des fins frauduleuses.

Exemples : spam ciblé, usurpation d'identité, profilage.

Problématique de l'usurpation d'identité en France

Près de 213.000 cas d'usurpation d'identité en France chaque année, selon le CREDOC

AP | 06.10.2009 | 19:12

Près de 213.000 cas d'usurpation d'identité avérés sont enregistrés chaque année en France, d'après une étude publiée mardi par le CREDOC. Comptes débités, emprunts indûment contractés, prestations sociales détournées: les escroqueries liées au vol de données personnelles coûteraient en moyenne près de 2.229 euros à chaque victime, soit un total de 474 millions d'euros.

L'enquête a été réalisée par le Centre de recherche pour l'étude et l'observation des conditions de vie (CREDOC) pour Fellowes, une entreprise qui produit notamment des machines de destruction de documents.

Selon les calculs du CREDOC, 212.762 personnes sont victimes d'usurpation d'identité en France chaque année, phénomène qui serait ainsi plus fréquent que les cambriolages au domicile principal (153.000 environ). Sur 2.007 personnes interrogées, 4,2% ont déclaré avoir été victimes d'une usurpation d'identité au cours des dix dernières années.

Dans plus de la moitié des cas, les usurpateurs débitent le compte bancaire (51,9% sur 300 personnes victimes d'une usurpation). Autre escroquerie fréquente, ils contractent un ou plusieurs emprunts au nom de leur victime (21,7%) ou bénéficient de prestations sociales à leur place (11,1%). Mais ils peuvent aussi causer des infractions au code de la route (16,9%), prendre un emploi (1,7%) ou même réaliser un mariage (1,7%) en se faisant passer pour leur victime.

La grande majorité des personnes visées découvrent qu'on leur a volé leur identité en vérifiant leur relevé bancaire (41,4%), en recevant des factures ou procès-verbaux indus (18,7%).

Le coût d'une usurpation atteint en moyenne 2.228,7 euros, entre le montant des détournements (1.520,4 euros), les démarches administratives et judiciaires (584,5 euros) et les coûts supplémentaires comme les frais médicaux ou postaux (142,8 euros). Sur ce total, le montant des remboursements atteint 661,8 euros en moyenne, 1.566,9 euros restant à la charge de la victime.

(Extrait d'un article du nouvel observateur paru le 6 octobre 2009)

Problématique de l'usurpation d'identité au Canada

ID theft scams target Canada's health-care system

Criminals are exploiting lax security in government databases to assume false identities and take advantage of Canada's health-care system, warns a leading expert in identity fraud.

BY THE OTTAWA CITIZEN

NOVEMBER 3, 2008

BE THE FIRST TO POST A COMMENT

...

For a time, "the competition for dead children was so fierce that we had situations where two and three different crime groups were accessing the same dead child's identity," said Mr. Pendleton.

Bris de vie privée

Remplacer le nom d'une personne par un **pseudonyme** ⇒
protection de la vie privée d'un individu

A Face Is Exposed for AOL Searcher No. 4417749

The New York Times

August 8, 2006

What Revealing Search Data Reveals

AOL posted, but later removed, a list of the Web search inquiries of 658,000 unnamed users on a new Web site for academic researchers. An interview with one of those unnamed users, Thelma Arnold, combined with her data reveal what she was searching for, why and on which Web sites.

A sample of Thelma Arnold's search data released by AOL

4417749	swing sets	2006-04-24	15:39:30	4	http://www.byswingset.com
4417749	swing sets	2006-04-24	15:39:30	9	http://www.buychoice.com
4417749	swing sets	2006-04-24	15:39:30	10	http://www.creativeplaythings.com
4417749	swing sets	2006-04-24	15:39:30	5	http://www.childlife.com
4417749	swing sets	2006-04-24	15:39:30	6	http://www.planitplay.com
4417749	that do not shed	2006-04-28	9:05:54	2	http://www.gopetsamerica.com
4417749	dog who urinate on everything	2006-04-28	13:24:07	6	http://www.dogdaysusa.com
4417749	walmart	2006-04-28	14:07:32	1	http://www.walmart.com
4417749	womens underwear	2006-04-28	14:12:28	10	http://www.bizrate.com
4417749	jcpenny	2006-04-28	14:16:05		
4417749	jcpenny	2006-04-28	14:16:49	1	http://www.jcpenny.com
4417749	tortus and turtles	2006-04-29	13:12:47		
4417749	manchester terrier	2006-05-02	9:05:31	1	http://www.manchestertierrier.com
4417749	delta	2006-05-02	13:48:26		
4417749	fingers going numb	2006-05-02	17:35:47		
4417749	dances by laura	2006-05-02	17:59:32		
4417749	dances by lori	2006-05-02	17:59:57		
4417749	single dances	2006-05-02	18:00:18	1	http://solosingles.com
4417749	single dances in atlanta	2006-05-02	18:01:13		
4417749	single dances in atlanta	2006-05-02	18:01:50		
4417749	dry mouth	2006-05-06	16:49:14	2	http://www.mayoclinic.com
4417749	dry mouth	2006-05-06	16:49:14	8	http://www.wrongdiagnosis.com
4417749	tyroid	2006-05-06	16:59:34		
4417749	thyroid	2006-05-06	16:55:44		
4417749	competitive market analysis of homes in illburn	2006-05-14	12:14:52		
4417749	competitive market analysis of homes in illburn	2006-05-14	12:16:17		
4417749	competitive market analysis of homes in illburn	2006-05-14	12:16:43		

Why the search

"I was thinking about my grandchildren"

"I was looking for some."

"A woman was in the [public] bathroom crying. She was going through a divorce. I thought there was a place called 'Dances by Lori,' for singles."

"I wanted to find out what my house was worth."



Erik S. Lesser for The New York Times

Thelma Arnold's identity was betrayed by AOL records of her Web searches, like ones for her dog, Dudley, who clearly has a problem.

AOL posted, but later removed, a list of the Web search inquiries of 658,000 unnamed users on a new Web site for academic researchers. An interview with one of those unnamed users, Thelma Arnold, combined with her data reveal what she was searching for, why and on which Web sites.

Technologies respectueuses de la vie privée

Technologies respectueuses de la vie privée (*Privacy Enhancing Technologies* ou *PET* en anglais) : ensemble de techniques et d'applications qui permettent à un individu de protéger ses informations personnelles pendant qu'il est en ligne.

Exemples : anonymisateur, communication anonyme, pseudonyme.

Deux principes fondamentaux derrière les PETs :

- ▶ **Minimisation des données** : seule l'information nécessaire pour compléter une application particulière devrait être collectée/révélee (et pas plus).
- ▶ **Souveraineté** : permettre à l'individu de garder le contrôle sur ses données personnelles et sur comment elles sont collectées et diffusées.

Classification proposée des méthodes de protection

Dimensions proposées pour classifier les méthodes :

1. Moment de la protection :


- ▶ Protection *en-ligne* (lorsque l'utilisateur est physiquement connecté) ou
- ▶ protection *hors-ligne* (lors d'un accès futur aux données enregistrées).

2. But de la protection :

- ▶ Protéger un individu à un niveau local ou
- ▶ un groupe de personnes à un niveau global.

3. Technique de protection utilisée :

- ▶ Perturbation de l'information.
- ▶ Primitives cryptographiques et calcul multipartit sécuritaire.
- ▶ Accès restreint ou limité à l'information (par exemple au niveau des requêtes).

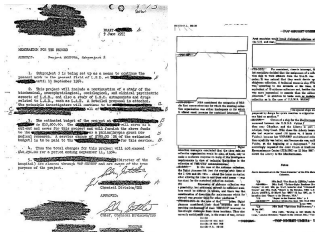
Autres dimensions à considérer : localité, interactivité, adaptivité. 

Assainissement

Assainissement (*sanitisation* en anglais) : *processus qui accroît l'incertitude dans les données afin de protéger la vie privée.*

⇒ Compromis inhérent entre le niveau de protection de la vie privée désirée et l'utilité des données "assainie".

Exemple typique d'utilisation : rendre public des données.

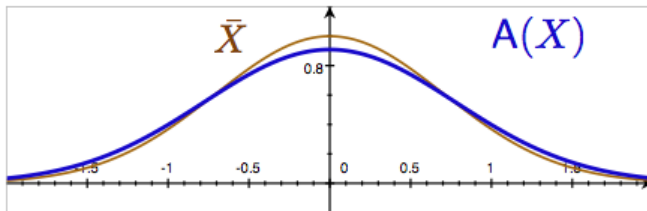


Exemples tirés de l'entrée "sanitization" sous Wikipedia

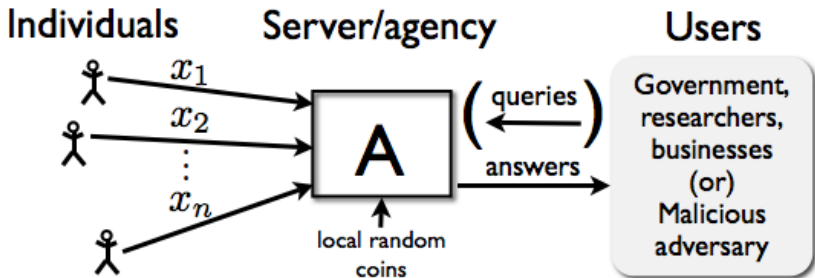
Méthodes de randomisation

Randomisation: ajout de bruit indépendant (comme gaussien ou uniforme) aux valeurs des informations transmises (par exemple la localisation ou les données enregistrées).

But: cacher les valeurs spécifiques des attributs tout en préservant la distribution jointe des données.



Modèle possible pour les méthodes de randomisation



Extrait d'un tutoriel de Adam Smith sur la protection de la vie privée dans les bases de données (mars 2008)

Attaque par inférence

- ▶ **Attaque par inférence** : l'adversaire prend en entrée un ensemble de données assaini (et possiblement des connaissances auxiliaires) et essaye d'inférer de nouvelles informations personnelles.
- ▶ **Attaque par chaînage** : l'adversaire essaye de relier ensemble les enregistrements de deux ensembles de données différents contenant une fraction d'individus en commun.
- ▶ Le *risque de divulgation par chaînage* mesure la probabilité de succès de cette attaque.
- ▶ **Défi principal** : être capable de donner des garanties de protection de vie privée même contre un adversaire ayant des connaissances auxiliaires.
- ▶ Cependant, il est possible que cette connaissance a priori ne puisse pas être modélisée.

Attaque par chaînage originelle de Sweeney

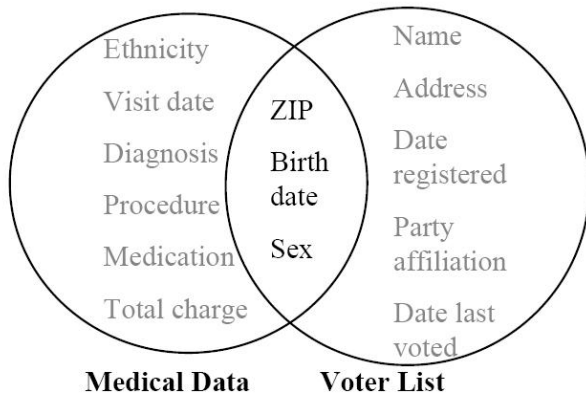


Figure 1 Linking to re-identify data

k-anonymité (Sweeney 02)

- ▶ **Garantie**: dans chaque groupe de l'ensemble de données assaini, chaque individu est identique à au moins $k - 1$ autres.
- ▶ Atteint par une combinaison de généralisation et suppression.
- ▶ **Exemple d'utilisation**: assainissement de données médicales.

	Non-Sensitive			Sensitive
	Zip Code	Age	Nationality	Condition
1	13053	28	Russian	Heart Disease
2	13068	29	American	Heart Disease
3	13068	21	Japanese	Viral Infection
4	13053	23	American	Viral Infection
5	14853	50	Indian	Cancer
6	14853	55	Russian	Heart Disease
7	14850	47	American	Viral Infection
8	14850	49	American	Viral Infection
9	13053	31	American	Cancer
10	13053	37	Indian	Cancer
11	13068	36	Japanese	Cancer
12	13068	35	American	Cancer

Figure 1. Inpatient Microdata

	Non-Sensitive			Sensitive
	Zip Code	Age	Nationality	Condition
1	130**	< 30	*	Heart Disease
2	130**	< 30	*	Heart Disease
3	130**	< 30	*	Viral Infection
4	130**	< 30	*	Viral Infection
5	1485*	≥ 40	*	Cancer
6	1485*	≥ 40	*	Heart Disease
7	1485*	≥ 40	*	Viral Infection
8	1485*	≥ 40	*	Viral Infection
9	130**	3+	*	Cancer
10	130**	3+	*	Cancer
11	130**	3+	*	Cancer
12	130**	3+	*	Cancer

Figure 2. 4-anonymous Inpatient Microdata

- ▶ **Défi principal**: extraire des connaissances utiles tout en préservant la confidentialité des données sensibles.

Quelques autres métriques de protection de la vie privée

- ▶ *l*-diversité (MKG^V 07) : maintient de la diversité dans chaque groupe par rapport aux valeurs possibles des attributs sensibles.
- ▶ Peut être instanciée par une métrique se basant sur l'entropie.
- ▶ Protège contre des attaques basées sur l'homogénéité et certaines autres attaques.
- ▶ *t*-proximité (LLV² 07) : la distribution des attributs dans chaque groupe doit être proche de celle de la population globale.
- ▶ *t* est un seuil à ne pas dépasser et qui représente la proximité entre distributions.

¹Machanavajjhala, Gehrke, Kifer et Venkatasubramanian.

²Li, Li et Venkatasubramanian.

Attaque par composition

- **Question** : supposons que le patron d'Alice sache qu'elle a 28 ans, qu'elle habite dans le quartier ayant le code postal 13012 et qu'elle est visitée les deux hôpitaux. Qu'est ce qu'il peut apprendre?

	Non-Sensitive			Sensitive
	Zip code	Age	Nationality	Condition
1	130**	<30	*	AIDS
2	130**	<30	*	Heart Disease
3	130**	<30	*	Viral Infection
4	130**	<30	*	Viral Infection
5	130**	>40	*	Cancer
6	130**	>40	*	Heart Disease
7	130**	>40	*	Viral Infection
8	130**	>40	*	Viral Infection
9	130**	3*	*	Cancer
10	130**	3*	*	Cancer
11	130**	3*	*	Cancer
12	130**	3*	*	Cancer

(a)

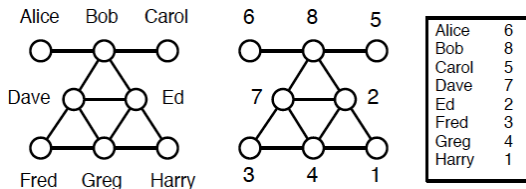
	Non-Sensitive			Sensitive
	Zip code	Age	Nationality	Condition
1	130**	<35	*	AIDS
2	130**	<35	*	Tuberculosis
3	130**	<35	*	Flu
4	130**	<35	*	Tuberculosis
5	130**	<35	*	Cancer
6	130**	<35	*	Cancer
7	130**	>35	*	Cancer
8	130**	>35	*	Cancer
9	130**	>35	*	Cancer
10	130**	>35	*	Tuberculosis
11	130**	>35	*	Viral Infection
12	130**	>35	*	Viral Infection

(b)

Attaques par inférence dans les réseaux sociaux

Anonymisation d'un graphe social

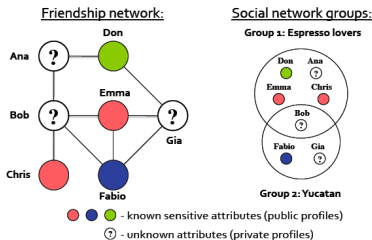
- ▶ Anonymiser un graphe (social) peut être une tâche très difficile car certains motifs dans le graphe peuvent être uniques.
- ▶ **Connaissance possible:** vous êtes le seul dans le réseau à avoir 47 amis et qui a 3 amis ayant chacun 52 amis.
- ▶ **Connaissance plus structurée:**



- ▶ **Conséquence:** anonymiser le graphe en enlevant les étiquettes des noeuds et des arêtes n'est pas suffisant.

Ce que vos amis révèlent à votre propos même si votre profil est privé

- ▶ **Idée principale** (Zheleva et Ghetoor 09): même si votre profil est privé la connaissance de votre réseau social + certains des attributs de vos amis peut être utilisée pour inférer certains de vos attributs personnels.
- ▶ Modélisé comme une tâche d'apprentissage semi-supervisé et ensuite un algorithme de propagation d'information est utilisé.



Project Gaydar

- ▶ Expérience menée par deux étudiants du MIT.
- ▶ **Hypothèse principale**: les préférences sexuelles de vos amis peuvent révéler de l'information à propos de vos propres préférences sexuelles.



- ▶ **Exemple** : si un individu a un nombre important d'amis gays alors il a une probabilité non-négligeable d'être gay lui aussi.
- ▶ **Attaque par inférence**:
 - ▶ consulter les pages Facebook des amis qui ont un profil public et qui ont déclaré explicitement leur orientation sexuelle.
 - ▶ construire un classificateur prédisant si oui ou non un individu est gay.

Compétition de Netflix

- ▶ Compétition mise en place par Netflix et proposé comme un défi à la communauté d'apprentissage machine pour améliorer la précision du système de recommandation de films.
- ▶ Ensemble de données de 35000 individus (lignes) où chaque dimension (colonne) est un film.
- ▶ La cellule $C_{i,j}$ contient le score donné par l'individu i au film j .
- ▶ **Caractéristiques principales:** les données sont de haute dimensionnalité (milliers de films) et éparées (la plupart des individus n'ont pas voté pour plus de 100-200 films).
- ▶ **Remarque:** la k -anonymité et autres techniques similaires ne fonctionnent pas dans ce contexte.

Ensemble de données publié et dé-anonymisation

► Forme de l'ensemble de données:

	Item 1	Item 2	Item M		
User 1	👍	👎	👍		
User 2		👍			
	👍		👎	👍	👍
	👍		👎		
		👍	👎	👎	
User N		👎	👍		

- Aucun identifiant n'a été utilisé.
- Semble fournir une certaine forme d'anonymat mais...
- la dé-anonymisation a été possible pour un nombre important d'enregistrements à l'aide d'une attaque par inférence utilisant Internet Movie DataBase (IMDB) comme information auxiliaire (Narayanan et Shmatikov 08).

Exemple d'un enregistrement IMDB

seaview1

Send an [IMDb private message](#) to this author or view their [message board profile](#).

Jump To: [A](#) [B](#) [C](#) [D](#) [E](#) [F](#) [G](#) [H](#) [I](#) [J](#) [K](#) [L](#) [M](#) [N](#) [O](#) [P](#) [Q](#) [R](#) [S](#) [T](#) [U](#) [V](#) [W](#) [X](#) [Y](#) [Z](#)

Page 1 of 3: [\[1\]](#) [\[2\]](#) [\[3\]](#) ▶

104 reviews in total

[Expanded](#) | [Show all](#) | [Show summaries](#) | [Alphabetical](#) | [Chronological](#) | [Useful](#)

1. [3:10 to Yuma](#) (2007) 13 September 2007
2. [A History of Violence](#) (2005) 12 October 2005
3. [An Education](#) (2009) 5 March 2010
4. [Angels & Demons](#) (2009) 25 May 2009
5. [A Serious Man](#) (2009) 5 March 2010
6. [Atonement](#) (2007) 8 February 2008
7. [Avatar](#) (2009) 5 March 2010
8. [AVP: Alien vs. Predator](#) (2004) 12 August 2004
9. [Babel](#) (2006) 3 November 2006
10. [Batman Begins](#) (2005) 27 July 2005
11. [Be Cool](#) (2005) 3 March 2005
12. [Beyond the Sea](#) (2004) 22 December 2004
13. [Blade: Trinity](#) (2004) 2 December 2004
14. [Bridget Jones: The Edge of Reason](#) (2004) 4 November 2004
15. [Broken Flowers](#) (2005) 12 August 2005
16. [Capitalism: A Love Story](#) (2009) 14 November 2009
17. [Charlie's Angels](#) (2000) 4 November 2000
18. [Christmas with the Kranks](#) (2004) 23 November 2004
19. [Cinderella Man](#) (2005) 4 July 2005
20. [Collateral](#) (2004) 6 August 2004

Quand la recherche sur la vie privée a un impact tangible

...

FRIDAY, MARCH 12, 2010

Netflix Prize Update

This is Neil Hunt, Chief Product Officer for Netflix.

About five months ago we announced that Netflix would sponsor a sequel to the Netflix Prize. We've given a lot thought to how to sponsor a contest that discovers more about the predictability of Netflix members' movie watching behavior while always ensuring we protect Netflix members' privacy.

In the past few months, the Federal Trade Commission (FTC) asked us how a Netflix Prize sequel might affect Netflix members' privacy, and a lawsuit was filed by KamberLaw LLC pertaining to the sequel. With both the FTC and the plaintiffs' lawyers, we've had very productive discussions centered on our commitment to protecting our members' privacy.

We have reached an understanding with the FTC and have settled the lawsuit with plaintiffs. The resolution to both matters involves certain parameters for how we use Netflix data in any future research programs.

In light of all this, we have decided to not pursue the Netflix Prize sequel that we announced on August 6, 2009.

Protection des requêtes

Ce que vos requêtes disent à votre propos

- ▶ **Leçon apprise de l'incident AOL**: vos requêtes révèlent beaucoup d'information à votre propos.
- ▶ Expérimentation de Yahoo conduite sur les logs de requêtes de 65000 utilisateurs (Jones, Kumar, Pang et Tomkins 07).
- ▶ **But**: prédire des attributs sensibles tel que l'âge, le sexe ou le code postal.
- ▶ Le profil de recherche d'un utilisateur est représenté sous la forme d'un vecteur (*bag-of-words* en anglais).
- ▶ Algorithmes d'apprentissage : classifieur bayésien et machines à vecteurs de support.
- ▶ Un dixième de l'ensemble de données est utilisée pour l'entraînement alors que le reste est utilisé pour le test.

Résultats de l'étude

- ▶ **Prédiction du sexe:** 83.3% de précision sur l'ensemble de test contre 57% si on prédisait la classe majoritaire (mâle).
- ▶ **Stéréotypes de requêtes pour la classe mâle:** *fanfiction, bridal, makeup, womens, knitting, hair, ecards, glitter, yoga, and diet.*
- ▶ **Stéréotypes de requêtes pour la classe femelle:** *nfl, poker, espn, ufc, railroad, prostate, football, golf, male, wrestling.*
- ▶ **Prédiction de l'âge:** prédiction avec une précision d'un maximum de 10 ans de différence avec l'âge réel.
- ▶ **Stéréotypes de requêtes pour la classe "jeunes":** *myspace, pregnancy, wikipedia, lyrics, quotes, apartments, torrent, baby, wedding, mall, soundtrack.*
- ▶ **Stéréotypes de requêtes pour la classe "vieux":** *aarp, telephone, lottery, amazon.com, retirement, funeral, senior, mapquest, medicare, newspapers, repair.*

Attaque sur une personne

- ▶ **But de l'adversaire**: identifier l'enregistrement d'un individu particulier dans un ensemble de données publié.
- ▶ L'adversaire peut utiliser ses connaissances auxiliaires pour essayer d'inférer des attributs sensibles.
- ▶ **Exemple** : votre voisin a une connaissance partielle à propos de vos loisirs et vos intérêts et essaye de vous identifier parmi un ensemble de logs de requêtes.

	Common	Rare
Cars	volkswagen beetle (478) honda odyssey (1504) toyota prius (1070)	triumph tr3 (23) e-type jaguar (5)
Sports	skiing (9618) football (123802)	bassmaster (388) skulling (17)
Food	pizza (104,888) italian restaurant (4998) brie (39,325)	assam (747)
Books	harry potter (27,838) danielle steele (238) freakonomics (574)	holly lisle (20) elizabeth moon (27)

Attaque sur une personne

- **Cadre de l'expérimentation:** organisation en paquets correspondant à une combinaison de plusieurs requêtes singletons.
- **Exemples de paquets :**

Query set	Bin size
harry potter, pizza	4855
football, skiing	2430
italian restaurant, pizza	1441
harry potter, volkswagen beetle	27
honda odyssey, italian restaurant	20
football, skiing, toyota prius	9
football, triumph tr3	4
football, harry potter, volkswagen beetle	3
pizza, triumph tr3	2
danielle steele, volkswagen beetle	1
brie, holly lisle, pizza	1

# users in bin	100+	51-99	26-50	6-25	3-5	2	1	< 100
# bins	51	13	17	65	44	31	99	320

Résultats de recherche personnalisés

privac

[Advanced Search](#)
[Language Tools](#)

privacy [Remove](#)

privacy enhancing technologies symposium 2010 [Remove](#)

privacy guard

privacy act

privacy assist

privacy policy

privacy assist bank of america

privacy center

privacy yacht

privacy center virus removal

History suggestions

Generic suggestions

Google Search I'm Feeling Lucky

Ingénierie inverse sur l'historique des requêtes

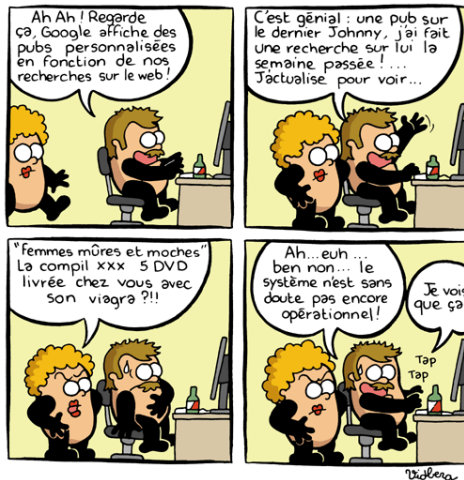
- ▶ **Google Web History**: résultats de recherche personnalisés par rapport à l'historique des requêtes d'un utilisateur et son comportement de navigation.
- ▶ **Historiographieur** (CCP³ 10): à partir des résultats de recherche personnalisés retournés par Google, il est possible de reconstruire une proportion importante de son historique de requêtes.
- ▶ Les résultats de recherche personnalisés nécessaires pour mener cette attaque peuvent être obtenus en :
 - ▶ détournement de session d'un utilisateur et faire des recherches à sa place ou
 - ▶ écoute passive du réseau de communication pour les services non-protégées par HTTPS.

³Castelluccia, de Cristofaro et Perito.

État actuel de certains services Google

Service Name	Default Connect.	HTTPS Support	Domain specific cookie	Purpose
Search	HTTP	no	no	Web search
Maps	HTTP	no	no	Maps search
Reader	HTTP	yes	no	RSS/Atom feed reader
Contacts	HTTP	yes	no	Address book manager
History	HTTP	yes	no	Search history manager
Gmail	HTTPS	mand.	no	Web mail application
Accounts	HTTPS	mand.	no	Google account manager
News	HTTP	no	no	News aggregator
Bookmarks	HTTP	yes	no	Bookmark manager
Docs	HTTP	yes	yes	Office application
Calendar	HTTP	yes	yes	Calendar application
Groups	HTTP	yes	yes	Discussion groups application
Books	HTTP	no	no	Personalized digital library

Peut on faire le même type d'attaque sur de la publicité ciblée?



Une alternative possible à Google



recherche anonyme - confidentielle - sécurisée - rien à télécharger



Protection de la vie privée dans Gossple

- ▶ **Gossple** : moteur de recherche décentralisé (réseau pair à pair) qui offre des résultats de recherche personnalisés par rapport au voisinage sémantique d'un individu.



- ▶ **But principal en terme de protection de la vie privée** : identifier les voisins sémantiques d'un individu en utilisant un protocole de type *gossip* et sans avoir à échanger en clair le profil des individus.
- ▶ Travail conjoint avec Anne-Marie Kermarrec (INRIA - Rennes) et Mohammad Nabil Alaggan.
- ▶ **Solution possible** : combinaison de techniques cryptographiques + *differential privacy*.

Géolocalisation et protection de la vie privée

Géo-localisation et protection de la vie privée

- ▶ La **géo-localisation** associe une localisation géographique à un objet tel qu'un téléphone portable, un ordinateur ou un véhicule équipé d'un GPS.
- ▶ Ces objets sont souvent personnels et associés à un individu spécifique.
- ▶ Si divulgué à des entités non-autorisées, ces informations peuvent amener à un **bris de vie privée** de la même manière que l'historique des achats d'un individu ou ses requêtes personnelles.



Défi principal : concilier géo-localisation et respect de la vie privée dans le cas d'applications liées aux données spatio-temporelles d'un individu.

Geo-privacy

La **geo-privacy**⁴ cherche à *empêcher une entité non-désirée d'apprendre la localisation géographique passée, présente et future d'un individu* (Beresford et Stajano 03).

Remarque : les **données personnelles spatio-temporelles** d'un individu peuvent jouer le rôle de *quasi-identificateurs*.

Ainsi, les données d'un individu peut permettre d'inférer :

- ▶ son lieu d'habitation et de travail,
- ▶ son identité,
- ▶ ses centres d'intérêts,
- ▶ ses habitudes ou
- ▶ une déviation par rapport à son comportement habituel.

⇒ **Bris de vie privée**

⁴Aussi appelé parfois "*locational privacy*" en anglais.

INRIA Alumni

Coordonnées personnelles

Caché aux autres membres

Email :

N° de téléphone :

N° de mobile :

Adresse (rue) :

Ville :

Code Postal :

Pays :


Géolocalisation : Je veux être géo-localisable.

Vos coordonnées géographiques : Longitude : Latitude :

1/ Se localiser à partir de l'adresse saisie :

Se localiser directement sur la carte en déplaçant le marqueur

2/



Abus lié au contrôle de la localisation géographique

Exemple :

GPS Vehicle Tracking Device

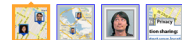
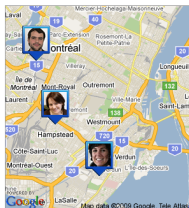
TravelEyes2® is a powerful vehicle tracking device that utilizes the GPS (Global Positioning System) technology developed by the Department of Defense. With amazing accuracy, it gathers miles driven on a daily basis, maintains time for billing records and produces comprehensive reports on all vehicle travel activities. It's ideal for the busy executive, a concerned parent with a watchful eye to track their teenager's late night activity or to track suspicious spouse, TravelEyes2® can work for you. It's a real time and money-saver -- and in some cases, a lifesaver!



Price \$ 199.00

Pas si
différent
de ...

Google latitude



Fred wants to hang out with his friends, and checks to see where they are.

[Learn more about Google Latitude](#)

See where your friends are right now

Enjoy Google Latitude on your phone, computer, or both.

Start using it on your phone

See your friends' locations and share yours with them.

Enter your number or visit google.com/latitude on your mobile browser.

> [Will it work with my phone?](#)

View it on your computer

See your friends' locations on a full screen even without a compatible phone or data plan.

[Add Latitude to iGoogle](#)



This service is free from Google; carrier charges may apply.

Please rob me



PLEASE ROB ME

Listing all those empty homes out there

Check out the same results on [Twitter search](#).

Next step

 We at Forthehack have been thinking about how we want to continue pleaserobme.com. It has received a lot of attention and it's time for a next step. We want to offer this website to a professional foundation, agency or company that focuses on raising awareness, helping people understand and provide answers to online privacy related issues.

If you're such a foundation, agency or company, [contact us](#).

More Info

[Home](#)

[Why](#)

[About](#)

Made Possible By

Forthehack

Connaissances *a priori*

L'adversaire peut avoir des connaissances *a priori* lui permettant d'essayer d'obtenir un bris de vie privée.

Exemples de connaissances potentielles :

- ▶ présence d'un individu parmi les données anonymisées,
- ▶ connaissance partielle de ses attributs (comme adresse de la maison ou du lieu de travail),
- ▶ modèle de ses habitudes,
- ▶ connaissance de son réseau social,
- ▶ connaissance de la distribution des attributs parmi la population,
- ▶ connaissance géographique des routes et du relief,
- ▶ ...

Exemple de connaissance géographique : Google Earth

The screenshot displays the Google Earth interface. On the left, the search panel shows results for 'IRISA'. The search bar contains 'IRISA'. Below it, a list of results is shown:

- IRISA (1 - 10)
- Fiches d'entreprises fournies par [PageJaune.it](#), Fiches d'entreprises fournies par [eniro.se](#), Fiches d'entreprises fournies par
- A** [Irisa hotel](#)
24, Banu Manta Avenue, District 1, Bucharest / Bukarest, Romania
- B** [Inria Rennes - Bretagne Atlantique](#)
Campus de Beaulieu, 263 avenue du Général Leclerc, 35042 Rennes
- C** [Irisa](#)

The main map area shows a 3D satellite view of the Inria Rennes campus. A red location pin labeled 'B' is placed on the main building, with a label 'Inria Rennes - Bretagne Atlantique (Irisa)'. The map includes navigation controls on the right side, such as a compass and a vertical slider. The bottom right corner of the map area shows the copyright notice '© 2009 Tele Atlas' and the Google logo.

Connaissances pouvant être inférées à partir de traces

À partir simplement des **traces de mobilité et de contact** d'une personne (mesurées par son téléphone), on peut parfois inférer :

- ▶ la maison et le lieu de travail, les horaires de travail ainsi que l'itinéraire utilisé,
- ▶ quelles machines/personnes elle a croisé,
- ▶ si elle a une Freebox ou une Livebox,
- ▶ à quelle heure elle a couché son fils et le film qu'elle a regardé,
- ▶ son heure de réveil,
- ▶ quels amis elle a visité,
- ▶ quand elle voulait être tranquille (pas d'enregistrements),
- ▶ quand la secrétaire du labo est allée prendre sa pause café,
- ▶ ...

Identification de lieu

Supposons qu'on dispose des traces GPS de la voiture d'un individu où le nom de la personne a été remplacé par un pseudonyme généré aléatoirement.

Heuristique pour identifier la maison de cet individu :

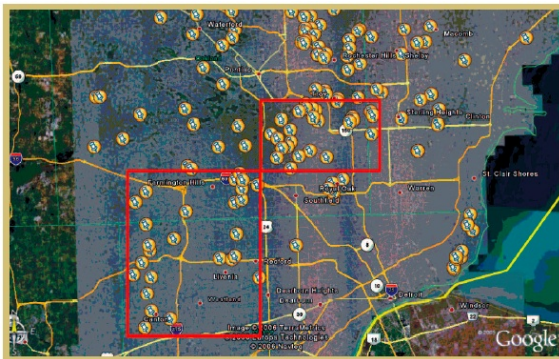
- ▶ Choisir le dernier arrêt avant minuit.

Heuristique pour identifier le lieu de travail :

- ▶ Choisir l'endroit où il y a le moins de déplacement dans la journée.

Géocodage inverse : transforme les coordonnées d'une localisation en une adresse physique.

Illustration de l'identification de maison



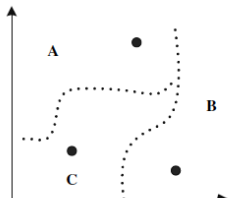
Les auteurs ont ciblés deux zones et manuellement localisés 65 foyers potentiels⁵. Dans 85% des cas, l'algorithme d'identification a retourné les mêmes localisations.

⁵L'identité exacte des conducteurs était gardée secrète.

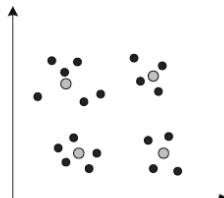
Assainissement de données géolocalisées

- ▶ **Masques géographiques** : *modifier la localisation géographique afin de préserver la vie privée d'un individu.*
- ▶ **Exemples de modifications** : aggrégation, perturbation aléatoire, randomisation en fonction de la densité.
- ▶ **Autres transformations possibles** : **sous-échantillonner**, **échanger** certaines traces de deux pseudonymes différents, **rajouter** des enregistrements artificiels.
- ▶ **Couverture spatiale** (Gruteser et Grunwald 03) : faire en sorte qu'à chaque unité de temps, chaque individu soit dans une zone qui est partagée par au moins $k - 1$ autres individus.
- ▶ **Mix-zone** (Beresford et Stajano 03) : zone de l'espace où
 - ▶ aucune observation n'est produite et
 - ▶ telle qu'un nouveau pseudonyme sera généré à la sortie qui est différent de celui de l'entrée.

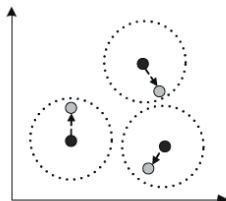
Exemples de masques géographiques



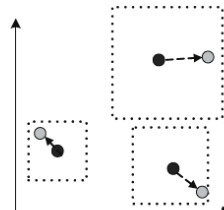
(a.) Aerial aggregation



(b.) Point aggregation




(c.) Random perturbation



(d.) Density sensitive

(Extrait de Armstrong, Rushton et Zimmerman 99)

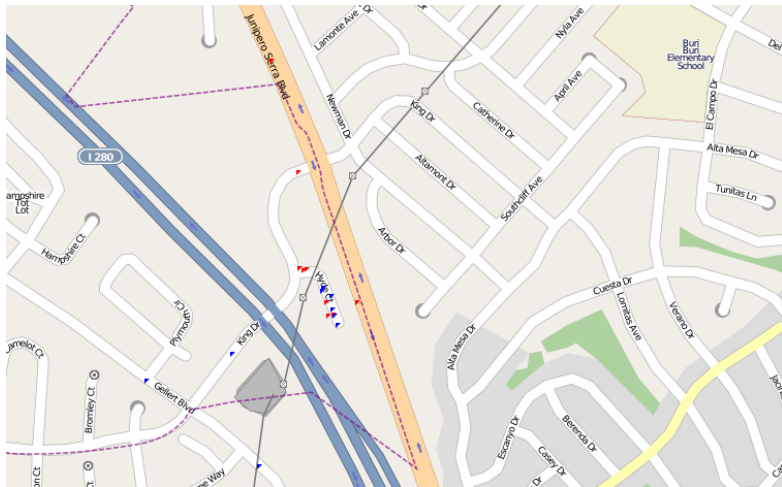
Limites des masques géographiques et de la couverture spatiale

- ▶ Si l'adversaire dispose de connaissances géographiques sur la zone sanitisée, il peut écarter certaines hypothèses qui semblent peu probables.
- ▶ **Exemple** : si suite à une perturbation aléatoire la localisation retournée est au milieu d'une zone difficilement accessible comme une falaise ou une rivière \Rightarrow l'adversaire peut facilement abandonner cette hypothèse pour la zone accessible la plus proche.
- ▶ **Risque de chaînabilité** : même s'il est impossible d'identifier exactement un individu, il est parfois possible de chaîner les actions d'un groupe d'individus.
- ▶ **Exemple d'inférence** : à chaque étape de temps, je suis capable de suivre le déplacement d'un groupe d'une zone à l'autre. 

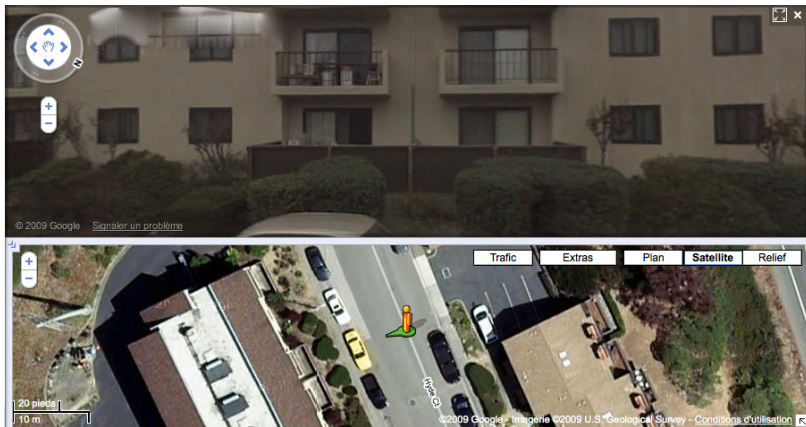
GEPETO

- ▶ Travail en cours conjoint avec Marc-Olivier Killijian (LAAS-CNRS) et Miguel Núñez del Prado (INSA, Toulouse).
- ▶ **GEPETO** (*GEoPrivacy-Enhancing TOolkit*): logiciel permettant de manipuler des données géolocalisées et de les visualiser, sanitiser, faire des attaques d'inférence et mesurer l'utilité.
- ▶ **Algorithmes de sanitisation** actuellement implémentés : pseudonymisation, sous-échantillonnage, perturbation de la position, aggrégation spatiale.
- ▶ **Algorithme d'inférence** : localisation de la maison et du lieu de travail.
- ▶ Approche testé sur un ensemble de données public de déplacement de taxis à San Francisco.
- ▶ **Connaissance géographique** : Google Maps et StreetView.

Maison de taxi identifié (vue de GEPETO)



Maison de taxi identifié (vue de GoogleMaps et StreetView)

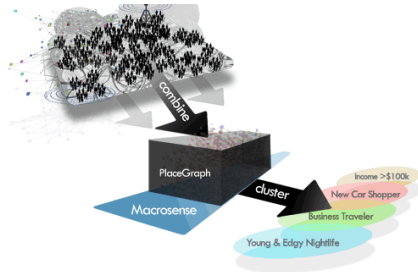


Un dernier exemple : Sense Networks



The image shows the header and main content area of the Sense Networks website. The header is a dark blue/black bar with the 'Sense Networks' logo on the left and navigation links 'Macrosense', 'Citysense', 'Technology', and 'Principles' on the right. Below the header is a large, high-angle photograph of a dense city skyline under a blue sky with white clouds. A semi-transparent white box is overlaid on the city image, containing the text: 'Indexing the real world using location data for predictive analytics.' At the bottom of the image, there is a row of small navigation icons including arrows, a search icon, and a refresh icon.

Utilisation des attaques par inférence pour faire du profilage



MacroSense enables companies to:

- Better understand customers using existing data, without requiring any change in behavior
- Segment and cluster customers into marketing groups based on actual unbiased behavior with unprecedented accuracy and relevance
- Personalize recommendations and advertisements based on popularity with "people like me"
- Automatically find and present the most relevant suggestions to a particular audience
- Identify group influencers

Conclusion

Au fur et à mesure que la frontière entre les domaines devient de plus en plus floue. . .

2 réponse(s) pour "Sarkozy Nicolas; France entière"

Voir aussi 2 réponse(s) trouvée(s) sur : Copains d'avant, Facebook, LinkedIn, Trombi, Twitter, Viadeo

Nouveau !

1 

✓ 2 réponses exactes | Réponses 1 à 2

✓ HAUTE GARONNE (31)

Sarkozy Nicolas

› fax

fax : 09 55 11 13 33

Mail : services@web-diffusion-france.com

Envoyer vers : [mobile](#) | [mail](#)

Ajouter au : [Carnet d'adresses](#)

✓ PARIS (75)

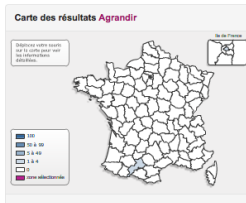
Sarkozy N

›

mobile  : 06 78 29 57 06

Envoyer vers : [mobile](#) | [mail](#)

Ajouter au : [Carnet d'adresses](#)



Nouveau ! pour : Nicolas Sarkozy, France entière

2 réponse(s) sur :

 (0)

Copains d'avant

 (1)

Facebook

 (0)

LinkedIn

 (0)

Trombi

 (1)

Twitter

 (0)

Viadeo

le travail de l'adversaire devient de plus en plus facile

pagesjaunes

les professionnels PagesBlanches à qui est ce numéro qui porte ce nom petites annonces

Revenir à la page de réponse(s) PagesBlanches

Recherche pour : **Sarkozy Nicolas PARIS ILE-DE-FRANCE** Nouveau !

Bêta 2 réponse(s) sur :

- Tous les profils
- Copains d'avant (0)
- Facebook (1)**
- LinkedIn (0)
- Trombi (0)
- Twitter (1)
- Viadeo (0)

Affiner par mot(s)-clé(s)

- Profils avec photo (1)

 **Nicolas Sarkozy**
Mot(s)-clé(s) : *Aucun*

facebook
[voir son profil](#)

Incitations à développer la protection de la vie privée

- ▶ **Prise de conscience**: effort important à faire au niveau de l'éducation du grand public afin d'atteindre une prise de conscience des risques.
- ▶ Risque malheureusement d'arriver suite à des scandales liées à de graves brèches de vie privée.
- ▶ **Cadre législatif**: avoir des lois ou des directives qui encadrent la collecte et la protection des données personnelles.
- ▶ **Exemple**: loi sur la protection des données personnelles, avis de la CNIL ou directive européenne.
- ▶ **Incitation commerciale**: le respect de la vie privée pourrait être un argument commercial.
- ▶ **Exemple**: produit labellisé vs produit non-labellisé.
- ▶ **Autre incitatif**: amende à payer pour les entreprises si fuite avérée d'informations personnelles.

C'est la fin !

Merci pour votre attention.
Questions?