

Evaluación de Protección de Privacidad de una Herramienta de Navegador Web

Estrada J.*; Rodríguez A.*

*Escuela Politécnica Nacional, Facultad de Ingeniería Eléctrica y Electrónica
Quito, Ecuador (e-mail: jose.estrada@epn.edu.ec; : ana.rodriguez@epn.edu.ec)

Resumen: Internet ha construido un escenario en el que la información del usuario es utilizada para identificarlo y clasificarlo. Muy pocas herramientas han sido propuestas para proteger al usuario ante este inminente riesgo de privacidad. Una de ellas es TrackMeNot, que implementa un mecanismo de perturbación de las consultas de búsqueda del usuario para ofuscar su perfil. Lamentablemente, no existen muchos estudios que determinen si este mecanismo es efectivo y en qué grado. En este trabajo se evalúa TrackMeNot, midiendo su efectividad en base a métricas justificadas de privacidad. Encontramos que, frente a ataques sencillos de identificación, TrackMeNot mejora la privacidad del usuario, pero que, frente a ataques más sofisticados, este mecanismo implementado de ofuscación no tiene ningún éxito.

Palabras clave: profiling, extensión de navegador, métricas de privacidad, ofuscación de consultas, perfil de usuario, evaluación de privacidad.

Abstract: Internet has constructed a scenario where the information of a user is used to identify and classify his interests. There are only a few tools that have been proposed to protect the user's privacy. TrackMeNot is one of these tools. It implements a perturbation mechanism of the user's search queries in order to obfuscate his profile. There is no certainty, however, about the effectiveness of such mechanism, neither about the extent to which privacy is being protected. An evaluation of TrackMeNot is done in this work, by measuring its effectiveness based on justified privacy metrics. We found that it successfully enhances user's privacy against identification attacks but completely fails on protecting privacy in front of classification attacks.

Keywords: profiling, browser extension, privacy metrics, query obfuscation, user profile, privacy evaluation.

1. INTRODUCCION

La concepción de los servicios en Internet ha cambiado significativamente desde que los usuarios son los principales generadores de contenido. Este fenómeno se debe, entre otras cosas, al auge de aplicaciones orientadas a la colaboración y que facilitan la reproducción de las experiencias de los usuarios en Internet.

Gracias al *Big Data* y a las avanzadas técnicas de análisis de datos, actividades como el *profiling* (obtención de perfiles) y la clasificación de usuarios se han vuelto prácticas comunes, llevadas a cabo por sistemas de personalización de contenido que se alimentan de toda la información que entrega el usuario, sin que este último sea consciente del riesgo que esto implica.

La personalización de contenido incentiva a los usuarios a entregar cada vez más información propia para que mejore su experiencia con el servicio. Dicha personalización se alimenta de perfiles de usuario creados a partir de la recopilación de patrones de navegación. El precio de obtener esta personalización es muy alto: la privacidad del usuario, en especial cuando existen múltiples fuentes cuyos datos (rastros) combinados con otros de distintos orígenes podrían revelar

información sensible relacionada con preferencias personales [2], [21].

La información personal sujeta al análisis de terceros es muy variada e involucra todos los datos originados por las interacciones del usuario: desde el contenido de las páginas visitadas, el tiempo consumido en un sitio web, el número de clics, las consultas a un motor de búsqueda, los datos entregados en formularios, y las cookies, hasta la configuración particular del navegador [25].

En ese entorno existe, por lo tanto, una amplia gama de posibles atacantes: motores de búsqueda, sistemas de recomendación, redes sociales, sistemas de etiquetado, etc. Sin embargo, los proveedores de servicios de Internet son entidades que tienen acceso a toda esa información relacionada con la actividad del usuario y en muchos casos ésta es también comercializada con compañías de publicidad o directamente utilizada para alimentar una plataforma de anuncios, sin considerar la privacidad de los dueños de esos datos [9].

Se debe dar particular atención a los motores de búsqueda y las redes sociales que se han convertido en los *gateways* por excelencia para llegar a servicios básicos como blogs o sitios de noticias, ya que los usuarios acceden a estos servicios consultando a Google o simplemente mediante los enlaces

publicados en sus perfiles de Facebook. En general, cualquier servicio web ahora recolecta consultas de búsqueda, etiquetas, clics y más metadatos que pueden mapearse a identidades de manera relativamente sencilla [21]. Aunque estos datos se hayan sometido a procesos de “anonimizado”, varios estudios ([2] y [32]) muestran que la privacidad aún puede comprometerse, es decir que no está garantizada. Esto es especialmente cierto en estos tiempos en los que se descubre que los gobiernos tienen acceso directo a información privada de los usuarios almacenada en los motores de búsqueda y las redes sociales más grandes. El material que se recolecta comprende: historiales de búsqueda, mensajes de correo electrónico, transferencias de archivos e incluso conversaciones en línea.

Además, existe mucha presión sobre las empresas que manejan información personal ([20] y [29]), para que apliquen fuertes políticas de privacidad con el fin de proteger los datos sensibles. Parece, sin embargo, que la presión externa (desde gobiernos por ejemplo) por revelar este rastro digital puede resultar mucho mayor. Las políticas de privacidad son instrumentos muy difundidos y populares para regular y comunicar la forma como se manejan los datos privados. Lamentablemente no son efectivas ya que los usuarios no las toman en cuenta, pero, eso sí, las aceptan sin ninguna reflexión. Esto demuestra una clara falta de consciencia con respecto a los riesgos a los que están expuestos los usuarios en Internet.

La privacidad se enmarca en un contexto sumamente complejo en el que puede depender incluso de los intereses individuales del usuario, por lo que el camino más corto para protegerla es incrementar el nivel de consciencia del usuario, evidenciando las debilidades y fortalezas de su conducta en Internet. El problema es que no existen herramientas que entreguen esta información (nivel de privacidad). Ciertamente, hay algunas que implementan medidas de ofuscación o bloqueo, pero no está claro cuál es su efectividad real. Esto se debe a que la privacidad y, en consecuencia, las herramientas que la protegen están relacionadas con una noción multidimensional y, por tanto, difícil de medir o evaluar. No está claro, entonces, si las múltiples herramientas existentes reducen realmente estos riesgos, razón por la que medir la efectividad de protección de estos mecanismos es imprescindible para poder compararlos y decidir cuál usar en determinados entornos de usuario.

Medir el nivel de privacidad es el primer paso para implementar medidas efectivas de protección pero también para evaluar mecanismos ya existentes.

1.1 Contribución

Considerando la falta de información sobre la efectividad de las herramientas disponibles de protección de privacidad, se propone un mecanismo para evaluar la ganancia de privacidad que se obtiene al instalar una conocida herramienta de protección, a nivel de navegador, y mediante el uso de métricas de privacidad planteadas en un trabajo anterior. Este mecanismo de evaluación puede integrarse con la herramienta de medición de privacidad propuesta

en [15] para que el usuario pueda determinar continuamente, en su navegador, el incremento de privacidad que obtiene gracias a esta herramienta.

Esta evaluación del nivel de privacidad obtenido al implementar un mecanismo de protección podría ayudar a los usuarios a decidir si dicho mecanismo es conveniente para sus intereses.

Nuestra evaluación aprovecha el módulo de *profiling* de otro *add-on* de Mozilla Firefox llamado Adnostic [33]. En concreto, dicho módulo nos permite obtener un perfil de usuario en base al cual determinamos varios niveles de riesgo de privacidad, mediante la utilización de métricas justificadas en conceptos de teoría de la información.

1.2 Organización

Este artículo se ha organizado de la siguiente manera. La Sec. 2 destaca algunas de las tecnologías y herramientas disponibles para la protección de la privacidad. La Sec. 3 resume los modelos de atacante y las métricas de privacidad utilizadas para la evaluación de la herramienta de protección de privacidad. La Sec. 4 describe el proceso de evaluación de un mecanismo de protección de privacidad en el navegador llamado TrackMeNot y, finalmente, en la Sec. 5 se mencionan las conclusiones de este trabajo.

2. ESTADO DEL ARTE

2.1 Tecnologías para el mejoramiento de la privacidad

Las tecnologías para el mejoramiento de la privacidad (PETs, por sus siglas en inglés) son medios técnicos para proteger la privacidad del usuario [34]. La privacidad es un concepto amplio que involucra varios enfoques, desde las características del tráfico de comunicaciones hasta el contenido de los mensajes transmitidos. Así, se puede clasificar a las PETs en tecnologías básicas *anti-tracking*, métodos criptográficos, enfoques basados en terceros (TTP), mecanismos colaborativos y técnicas de perturbación de datos.

Tecnologías Básicas Anti-tracking. El *tracking* es un mecanismo mediante el que una entidad identifica a otra en un proceso de comunicación. Esto es vital para los servicios personalizados ya que les permite mapear identidades con sus correspondientes preferencias. Para ello existen diferentes parámetros que permiten identificar a una entidad, por ejemplo una dirección IP o una *cookie*.

Bloquear o “esconder” estos parámetros de identificación es parte de las medidas básicas *anti-tracking*. El problema es que al hacerlo, varios servicios en Internet no se pueden ofrecer.

Recuperación de información de manera privada (PIR, por sus siglas en inglés).

La PIR permite a un usuario obtener información de una base de datos sin que el proveedor conozca el contenido recuperado [28]. Una solución sencilla, pero no muy práctica, podría ser que el usuario descargue la base de datos completa y luego, localmente, acceda al contenido de interés.

Otra opción, propuesta para sistemas de recomendación es la de no revelar perfiles individualmente sino uno agregado a partir de los perfiles de un grupo de usuarios.

Mecanismos basados en un tercero de confianza (TTP).

Un TTP es un intermediario que recibe las comunicaciones del usuario y las envía en su nombre al destino para proveerle de privacidad. Los mensajes en el destino aparecen originados en el TTP y, por tanto, no se podrían vincular con el usuario. La desventaja clara de estos mecanismos está en el retardo agregado por el TTP.

Los *Mixes* [5] son implementaciones de TTP que reciben un mensaje y lo reenvían a su destino de manera que el evento de llegada no pueda asociarse al de salida. Esto evita el rastreo que puede hacer un atacante al “escuchar” los eventos de reenvío para seguir un mensaje desde el origen a su destino.

Onion Routings también un mecanismo basado en TTP que consiste en enviar un mensaje que va cifrado varias veces desde el enrutador en el origen y que se va descifrando por capas hasta llegar al destino.

La *Colaboraciones* una tecnología que plantea la participación cooperativa entre usuarios del sistema para obtener privacidad. *Crowds*[23] y *LBS* [3] son protocolos que aprovechan la participación de varias entidades para encaminar la información de manera impredecible y consiguientemente anónima para el adversario.

En [7] se propone un mecanismo para privacidad en búsquedas web que consiste en que los usuarios intercambien porciones de sus búsquedas antes de enviar las suyas, con el fin de ofuscar sus perfiles de interés ante atacantes externos.

Perturbación de datos.

Consiste en obstruir al atacante en su afán de construir un perfil preciso de usuario; por ejemplo, enviando datos falsos junto con los reales.

El fraguado de consultas es una aplicación de esta tecnología, en la que se generan consultas forjadas desde el cliente, de modo que el motor de búsqueda no pueda obtener un perfil preciso ya que las consultas que recibe están ofuscadas.

TrackMeNot[4] es una implementación muy conocida de forjado de consultas. Es una extensión de navegador que genera consultas falsas y las envía a distintos motores de búsqueda desde el navegador del usuario. La generación de estas consultas se alimenta de contenido RSS alojado en distintas fuentes de información.

Otra propuesta a nivel de aplicación y en el navegador es *GooPIR*[14]. Ésta herramienta ofusca directamente cada consulta que hace el usuario a Google y utiliza para ello palabras obtenidas de una fuente local. Resulta bastante complicado, sin embargo, ofuscar consultas sensibles relacionadas, por ejemplo, con condiciones de salud o afinidad política.

En el campo de los sistemas de recomendación (basados en la compartición de ratings, por ejemplo) existen también propuestas de perturbación, por ejemplo en [27] donde se propone un algoritmo para enviar ratings perturbados al sistema de recomendación.

2.2 Herramientas orientadas a la protección de la privacidad

Actualmente existen algunas herramientas que tratan de proteger la privacidad del usuario en Internet, esencialmente mediante el bloqueo de funciones en el navegador y que facilitan la entrega de información personal. Estos mecanismos, generalmente basados en la heurística, no miden el riesgo de privacidad del usuario ni evalúan el nivel de protección que ofrecen; simplemente aplican una metodología intuitiva.

Adnostic[33] es un *add-on* desarrollado para el navegador Mozilla Firefox que implementa una arquitectura para desplegar publicidad personalizada, sin comprometer la privacidad del usuario, ya que se decide en el navegador qué anuncios mostrar, en función de un perfil calculado localmente. Este perfil se obtiene a partir del procesamiento de las consultas que realiza el usuario y del contenido de las páginas que visita. Luego, esta información es clasificada utilizando procesamiento natural de lenguaje dentro del navegador. Los anuncios, que forman parte de un conjunto previamente descargado, se despliegan dependiendo de los intereses del usuario.

REPRIV [31] es otro sistema propuesto para trabajar en el navegador que ofrece una personalización mejorada de contenido y un mecanismo de control del usuario sobre la información que entrega a terceros. Usa la información de navegación del usuario para descubrir cuáles son sus intereses, y comunicarlos a terceros para que estos últimos puedan ajustar el contenido en base a esas preferencias. Propone interfaces para sitios web de terceros para los protocolos de comunicación de información personal que funcionan sobre HTTP. Promete una mejora importante en la provisión de contenido a medida, gracias al gran detalle de la información del navegador, pero el control de privacidad podría verse afectado por la falta de usabilidad de las políticas de protección que se implementen y que un usuario promedio tendría que gestionar. Además, no muestran ninguna métrica que indique al usuario el nivel de privacidad que posee. Propuesto por Microsoft, *REPRIV* es un planteamiento interesante, aunque la protección de privacidad se aborda como una derivación del servicio de personalización que ofrece y que ya supone la confianza en un tercero.

En relación específica con la medición de privacidad, existen un par de estudios en [12] y [22] sobre herramientas para redes sociales (Facebook en los dos casos) que determinan el riesgo de privacidad del usuario en función de la cantidad de información que de éste se puede inferir a partir de sus relaciones con otros usuarios. También implementan acciones de protección de privacidad bloqueando estos usuarios, analizando la configuración de privacidad de la cuenta o detectando y eliminando aplicaciones que poseen demasiados permisos.

TrackMeNot[4] es otra herramienta para protección de privacidad a nivel de navegador que propone ofuscar el flujo de consultas que envía un usuario a motores de búsqueda, mediante la generación de consultas falsas. Ha recibido muchas críticas respecto de su eficacia, aunque no se han

propuesto muchos mecanismos para evaluar sus bondades. En [30] se muestra que estas consultas falsas podrían ser identificadas con relativa facilidad utilizando clasificadores basados en inteligencia artificial. Sin duda, la falta de una herramienta de medida de privacidad le impide al usuario valorar su condición de riesgo antes y después de aplicar una estrategia de protección como ésta.

Google Sharing[11] es otra herramienta que implementa un mecanismo de protección de privacidad al prevenir el rastreo de usuario que realiza Google mediante las consultas al motor de búsqueda. El mecanismo consiste en que el usuario envía sus peticiones a un proxy externo que gestiona un grupo de identidades asociadas a *cookies*. Estas *cookies* reemplazan las *cookies* de las peticiones, enmascarando la identidad del usuario, y luego se reenvían con la petición original a Google. Aun cuando permite enviar peticiones cifradas desde el usuario, la privacidad del usuario puede comprometerse si hay colusión entre Google y el servidor proxy.

Ghostery[10] es otro *add-on* de Firefox orientado a proteger la privacidad del usuario mediante la detección y bloqueo de “*trackers*” y otros objetos dedicados al rastreo de la actividad del usuario. Es una herramienta bastante completa y muy popular, con varios módulos que implementan los mecanismos de protección en distintos navegadores.

El modo de “navegación privada” es también una opción de protección de privacidad en los navegadores más conocidos. Ésta deshabilita el almacenamiento de información local (historial, imágenes, videos, *cookies*, etc.) durante la navegación web. Esto complica significativamente el acceso a muchos sitios en Internet, por lo que quienes usan este modo lo hacen durante intervalos de tiempo muy cortos. La protección se limita al ámbito local pues externamente existen otros mecanismos para identificar y clasificar al perfil del usuario.

El bloqueo o desactivación de ciertas características del navegador web es una medida común implementada por varias soluciones en forma de *plug-ins* de navegador (NoScript [24], AdBlock Plus [26], DoNotTrackMe [8]), y evitan que se libere información que pueda usarse para identificar al usuario.

Sin embargo, ninguna de estas herramientas evalúa el nivel de privacidad del usuario. *TrackMeNot* da un paso importante al implementar una medida proactiva para proteger la privacidad, pero no se plantea una valoración de su efectividad, lo que desincentiva su uso. En general, se considera como posibles adversarios únicamente a los anunciantes o a los servicios de redes sociales pero no a los proveedores de servicios de Internet (ISP, *Internet Service Provider*) que son las entidades que más información poseen sobre los usuarios, aunque se podría restringir enormemente este riesgo si todas las conexiones cliente-servidor estuvieran cifradas con HTTPS. En realidad los ISPs tienen acceso a toda la información que el usuario envía a Internet, y el gran detalle de la misma representa un enorme incentivo para su comercialización.

En [27], [17], [19] y [16] se abordan en algunos mecanismos que podrían emplearse para la protección de la privacidad del usuario en entornos donde éste hace consultas o etiqueta contenido; considerando también el costo de estas estrategias que se refleja en la pérdida de utilidad de los datos, la pérdida de funcionalidad de un servicio o el consumo adicional de recursos. Se incluye entre estos mecanismos la falsificación de consultas o la supresión de etiquetas con el fin de mostrar una versión distorsionada del perfil del usuario que el atacante no pueda explotar. La optimización de estos mecanismos así como su impacto son también sujetos de estudio.

3. MODELOS DE ATACANTE Y MÉTRICAS DE PRIVACIDAD

En esta sección se presenta los dos modelos de atacante considerados en este artículo, así como las métricas que permiten evaluar el nivel de privacidad de un usuario. Los modelos de usuario y las métricas de privacidad son ampliamente justificados en [18].

3.1 Modelos de Atacante

Los criterios de privacidad se plantean inicialmente, asumiendo que el perfil del usuario es modelado como una función de masa de probabilidad o un histograma de frecuencias relativas de datos de usuario a lo largo de un conjunto preestablecido de categorías de interés. Este modelo supone una representación muy habitual en los servicios de información personalizada.

El modelo de adversario permite definir las propiedades del atacante, considerando como tal a cualquier entidad capaz de tener acceso a información de usuario, con el objetivo de obtener su perfil, con el riesgo a la privacidad que esto implica.

Conocer al adversario es importante ya que la privacidad del usuario se mide respecto a éste. En función de las propiedades del adversario, el usuario podría implementar medidas de protección de su privacidad que, por ejemplo, modifiquen su perfil de intereses.

Esencialmente, se contemplan dos objetivos del atacante, en función de sus capacidades, y que definen el modelo de adversario; *identificación* y *clasificación*.

- *Identificación*, cuando el atacante intenta distinguir al usuario del resto de la población, detectando desviaciones de sus intereses respecto del perfil promedio de la población.
- *Clasificación*, cuando el atacante intenta clasificar al usuario en un grupo de población, comparando el perfil del usuario con el perfil representativo del grupo.

3.2 Métricas de Privacidad

En [18] se justifica la entropía de Shannon y la divergencia de Kullback-Leibler (KL) como medidas de privacidad. Las interpretaciones de estas medidas dependerán de las hipótesis

que se hagan, fundamentalmente respecto del modelo de adversario.

Otra métrica más general, no limitada a la privacidad de

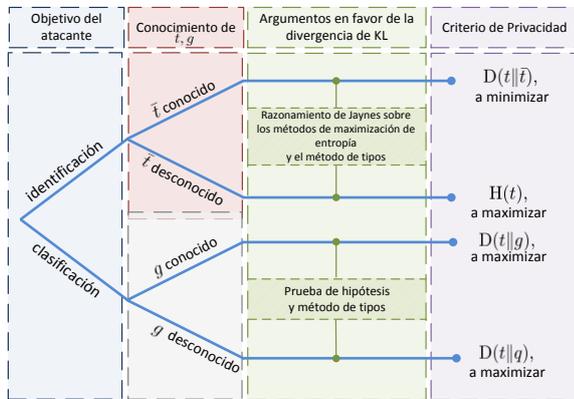


Figura 1. Resumen de las interpretaciones de la entropía de Shannon y de la divergencia de KL como métricas de privacidad, de acuerdo con la justificación presentada en [18].

perfiles, es la propuesta en [6]. En este trabajo, los autores proponen medir la privacidad como el error de estimación de un adversario, e interpretan, mediante argumentos de teoría de la información y teoría de decisión Bayesiana, otras métricas del estado del arte como casos particulares de la suya.

Para facilitar la comprensión, además, se resume a continuación las principales definiciones propuestas para la justificación de estas métricas de privacidad. Se revisa además la interpretación de estas dos cantidades de teoría de la información como métricas de privacidad de perfiles de usuario.

El símbolo H denotará la entropía de Shannon y D denotará la divergencia de KL. La entropía $H(p)$ de una variable aleatoria discreta X con distribución de probabilidad p es una medida de su incertidumbre, definida como

$$H(X) = -E \log p(X) = -\sum_x p(x) \log p(x).$$

La divergencia de KL o entropía relativa $D(p \parallel q)$ entre dos distribuciones de probabilidad $p(x)$ y $q(x)$ sobre el mismo alfabeto se define como

$$D(p \parallel q) = E_p \log \frac{p(X)}{q(X)} = \sum_x p(x) \log \frac{p(x)}{q(x)}.$$

La divergencia de KL es una medida de discrepancia entre distribuciones de probabilidad, garantizando que $D(p \parallel q) \geq 0$, con igualdad si, y sólo si, $p=q$. Consecuentemente se deduce que la entropía $H(p)$ alcanza su valor máximo en $H(u) = \log n$, siendo n la cardinalidad del alfabeto finito sobre el que se calcula $D(p \parallel u)$, para una distribución uniforme u :

$$D(p \parallel u) = \log n - H(p).$$

En concreto, acorde con el análisis en [18] tenemos que la maximización de la entropía resulta ser un caso especial de la minimización de la divergencia, alcanzada idealmente cuando la distribución a optimizar es idéntica a la de referencia.

Sea q el perfil de interés de un usuario, t una versión perturbada o modificada del mismo y \bar{t} la distribución del perfil de la población. En la Fig. 1 se muestran las

interpretaciones de la entropía de Shannon y de la divergencia de KL como medidas de privacidad. Éstas se explican a continuación, de acuerdo al objetivo del atacante.

3.2.1 Métricas contra identificación

En caso de que el objetivo del atacante sea identificar al usuario, el razonamiento de Jaynes acerca de los métodos de maximización de la entropía permite justificar la divergencia y la entropía como medidas de privacidad.

La entropía del perfil aparente del usuario, que es el perfil observado por el atacante, es justificada en [18] como una medida de la probabilidad de este perfil perturbado, en el sentido de frecuencia de aparición de dicho perfil en la población de usuarios. Considerando esta probabilidad del perfil de usuario como una medida razonable de su anonimato (o privacidad), en [18] se justifica también la entropía como una métrica de privacidad. En concreto, mientras mayor sea la entropía de este perfil, mayor es su probabilidad, y por tanto mayor es el número de usuarios que se comportan de acuerdo con este perfil, haciéndolo más privado.

Además, como se puede observar en la primera rama de la Fig. 1, si la distribución del perfil de la población \bar{t} es conocida, se utiliza la divergencia entre el perfil del usuario t y el perfil de la población como métrica de privacidad, de manera que, cuanto más pequeña sea esa divergencia, más privado se puede considerar el perfil.

En definitiva, la elección de perfiles aparentes que conduzcan a la minimización de la divergencia de KL mejora el anonimato. En términos más simples, una menor divergencia corresponde a una mayor frecuencia de ocurrencia de dicho perfil, permitiendo al usuario pasar más desapercibido. En el caso de un perfil de referencia de la población uniforme, esto equivale a la maximización de la entropía de Shannon.

3.2.2 Métricas contra clasificación

Si el objetivo del atacante es clasificar al usuario como miembro de un grupo en particular, se utiliza la divergencia como métrica de privacidad, de acuerdo al análisis realizado en [18], a partir del *test* de hipótesis y el método de tipos. Como se indica en la Fig. 1, en la segunda rama, si el perfil del grupo g es desconocido en el lado del usuario, la opción es maximizar la divergencia entre el perfil real q y el perfil observado (aparente) t , con el fin de evitar ser clasificado de acuerdo a su perfil original.

Nótese que en el problema de clasificación, al contrario de lo que ocurre en el problema de identificación, se busca maximizar la divergencia de KL, en lugar de minimizarla. La intuición subyacente al análisis citado es que se desea agrandar la distancia entre el perfil aparente del usuario y el perfil real, o el representativo del grupo en el que se desea evitar la categorización.

4. EVALUACIÓN DE UN MECANISMO DE OFUSCACIÓN DE PERFIL DE USUARIO

Existen innumerables propuestas para proteger la privacidad pero no muchas herramientas que las implementen a nivel de

usuario. Existen sí un conjunto de herramientas que, basadas en métodos heurísticos, tratan de mitigar los riesgos de privacidad mediante el bloqueo de ciertas interacciones del usuario con la Web.

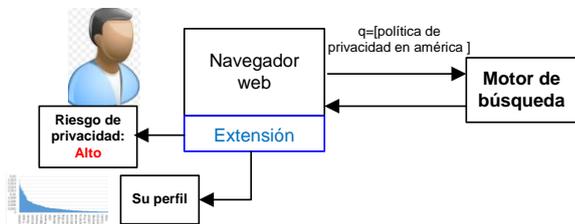


Figura 2. Arquitectura de la aplicación propuesta en [15] para medir la privacidad del usuario en el navegador.

Un paso importante para el desarrollo de aplicaciones más eficientes de protección de privacidad es, por tanto, evaluar las existentes. Esto permitiría al menos aprovechar las estrategias que mayores beneficios representen y descartar los mecanismos que no obtengan resultados satisfactorios.

Esta evaluación de las aplicaciones de protección de privacidad consiste en comparar los niveles de privacidad obtenidos antes y después de implementar una estrategia determinada. Así se podría calcular una suerte de ganancia de privacidad, interpretada como el beneficio potencial de aplicar una herramienta dada.

La medición de privacidad es, como se ha dicho ya, una tarea compleja ya que debe considerarse múltiples variables en los distintos enfoques propuestos tanto teóricamente como en la práctica.

En esta sección se explica la metodología utilizada para para evaluar una de las herramientas de protección de privacidad disponibles, mediante la utilización de las métricas de privacidad propuestas en [18] y algunos de los módulos de medición ya implementados en nuestra extensión de Firefox.

4.1 Consultas de búsqueda: información que identifica al usuario

La búsqueda de información es una actividad muy común durante la navegación en Internet. Es, además, evidente que las consultas de búsqueda de un usuario reflejan de manera muy fiel sus intereses, preocupaciones o problemas. Estas consultas, combinadas, por ejemplo, con interacciones en redes sociales y actividades de etiquetado, podrían revelar de manera muy precisa la identidad de este usuario. Esto ya ha sido demostrado, por ejemplo, cuando el New York Times expuso la identidad de una mujer, usando registros aparentemente “anonimizados” de búsqueda en AOL [21].

Los servicios de información personalizada son extremadamente lucrativos debido a su gran efectividad. Es así que gran parte de las ganancias que se obtienen en Internet provienen de la publicidad personalizada. Existe, sin embargo, un riesgo latente y, en muchos casos, comprobado, de que esta información es compartida con gobiernos u otras entidades externas.

La perturbación de datos

es

una técnica de preservación de privacidad que se puede aplicar del lado del usuario (como se explicó en la Sec. 2). En el caso de las consultas de búsqueda, la perturbación u ofuscado de consultas es un mecanismo que ha sido implementado mediante la extensión de Firefox *TrackMeNot*, con el fin de proteger a los usuarios frente actividades de *profiling*.

4.2 Medición de la privacidad de usuario en el navegador

En nuestro trabajo previo en [15] se propuso una extensión para el navegador Firefox que permita la medición de la privacidad del usuario en función, esencialmente, de las consultas de búsqueda que el usuario realiza a los distintos motores en Internet.

Este *add-on* detecta las palabras asociadas a las consultas de búsqueda para procesarlas en el navegador web y crear en base a ellas un perfil de usuario. Este perfil es luego utilizado como insumo de los módulos de medición de privacidad, de acuerdo a los mismos criterios expuestos en la Sec. 3. Este proceso se ilustra en la Fig. 2.

En base a esta medición de privacidad de perfiles, se plantea una evaluación de esta herramienta para determinar su efectividad en la protección de la privacidad del usuario.

4.3 TrackMeNot

TrackMeNot (TMN) es una extensión del navegador Firefox cuyo objetivo es ofuscar el perfil del usuario, introduciendo material forjado en el flujo de consultas de búsqueda. Este mecanismo se basa en la generación artificial de conjuntos de palabras (que aparentan ser consultas de usuario) que luego son enviadas a un motor de búsqueda mediante peticiones HTTP. TMN implementa además algunas estrategias para evitar que los motores de búsqueda detecten las solicitudes automatizadas y su arquitectura se resume en la Fig. 3. Algunos de sus elementos se describen a continuación.

4.3.1 Listas Dinámicas de consultas

TMN utiliza fuentes públicas de información (*feeds RSS*) para obtener una semilla a partir de la cual se construyen las consultas falsas. Esta lista de palabras es dinámica ya que cambia con el tiempo debido a que la información con la que se construye proviene de periódicos en línea cuyo contenido varía diariamente.

4.3.2 Conciencia de búsqueda en tiempo real

Puesto que varios motores de búsqueda son capaces de detectar un comportamiento automatizado, TMN soluciona esto enviando consultas falsas solamente cuando el usuario envía una consulta real.

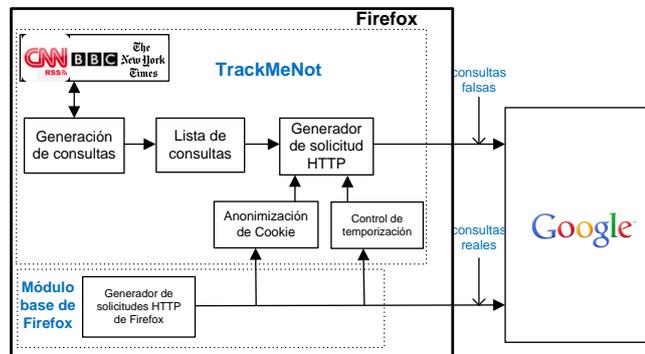


Figura 3. Arquitectura de TrackMeNot

4.3.3 Mapas de Cabeceras

Con el fin de simular mejor un comportamiento de usuario, TMN adapta las cabeceras de los mensajes HTTP que envía, de acuerdo a las cabeceras utilizadas por el usuario en la última consulta realizada.

4.3.4 Consultas en modo ráfaga

Este modo permite a TMN enviar varias consultas en el momento en el que el usuario genera las suyas. Así, el comportamiento del usuario se emula de mejor manera pues normalmente un usuario envía varias consultas en un corto período de tiempo y luego se detiene.

4.3.5 Anonimización de cookies

TMN toma las *cookies* de las consultas de usuario y las agrega a las consultas automáticamente generadas, también con el fin de simular mejor el comportamiento del usuario. TMN es una herramienta muy conocida en el ámbito de la privacidad y utiliza un mecanismo innovador para ofuscar el perfil de usuario visto por los motores de búsqueda. No existe, sin embargo, otras herramientas que implementen métodos similares aunque sí hay varias contribuciones teóricas en ese sentido.

4.4 Metodología de evaluación de TMN

La evaluación de TMN consiste en la comparación de los niveles de privacidad obtenidos antes y después de que TMN es utilizado, de modo que un valor de ganancia de privacidad se pueda calcular. Si la ganancia es positiva, significará que el mecanismo es efectivo. Si existe esta ganancia, se determinará si es suficiente para proteger a los usuarios frente a adversarios externos (es decir, si es eficiente).

Los principales pasos para estimar la eficiencia de TMN son los siguientes:

- Obtener una muestra significativa de registros (*logs*) de consultas de usuario,
- Generar consultas falsas, utilizando el mecanismo ofrecido por TMN,
- Ofuscar las consultas reales del usuario, mezclándolas con consultas falsas,
- Obtener los perfiles real y ofuscado de los usuarios, a partir de los *logs* de consulta,
- Medir la privacidad de los perfiles real y aparente, mediante las métricas antes descritas y,
- Determinar la ganancia de privacidad de los perfiles de usuario luego del proceso de ofuscado.

En la Fig. 3 se observa que el proceso de evaluación se inicia con la obtención de un número significativo de consultas de usuarios. Se usó para esto un *dataset* de AOL [21] liberado en 2006 que contiene cerca de 20 millones de consultas de 650 mil usuarios. La muestra tomada de este *dataset* pertenece a 6674 usuarios con 501 a 1000 consultas cada uno. De estos usuarios, se escogió a 50 con la mayor cantidad de consultas.

Las consultas falsas se obtuvieron mediante el módulo de generación de TMN. Se reutilizó el código pertinente de esta extensión para generar 700 consultas de búsqueda a partir de

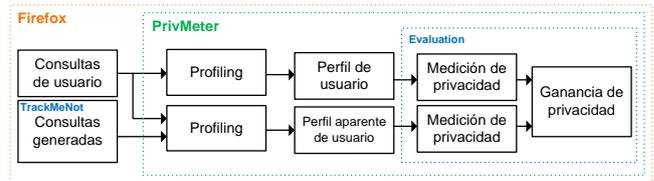


Figura. 4. Arquitectura de evaluación de TrackMeNot

los *feeds*RSS configurados por defecto en TMN (sitios de la CNN, BBC, *TheRegister* y *The New York Times*).

La ofuscación de las consultas reales de cada usuario se realiza mediante la mezcla de éstas con las consultas falsas previamente generadas con TMN.

Tanto las consultas reales como las ofuscadas fueron sujetas del proceso de *profiling* para obtener los perfiles reales y ofuscados (aparentes) de cada usuario en el experimento. Los primeros (perfiles reales) constituyen aquella información a proteger y los segundos son los perfiles obtenidos luego de utilizar TMN, es decir aquellos derivados de las consultas falsas mezcladas con las reales.

El siguiente paso es medir la privacidad de los perfiles de usuario de cada grupo, de acuerdo con las métricas explicadas en la Sec. 3. Finalmente se comparan los valores de privacidad obtenidos para determinar un valor de ganancia que permita juzgar si el mecanismo de protección de TMN es efectivo.

Algunos de estos pasos se explican con más detalle en los siguientes párrafos.

4.5 Obtención de perfiles a partir de información de usuario

Para obtener los perfiles de usuario a partir de los datos de consultas de búsqueda utilizados del *dataset* de AOL, se reutilizó el módulo de *profiling* de la extensión Adnestic. Se modificó este componente, de manera que construya los perfiles de usuario a partir de la categorización continua de las búsquedas del *dataset*. Cabe resaltar que entre las categorías en las que el texto de las consultas puede ser clasificado no se incluye categorías sensibles relacionadas con salud, racismo o pornografía. El esquema de categorización (602 categorías) está basado en la jerarquía que usa Google para clasificar a sus usuarios.

Modelar la estrategia de obtención del perfil de usuario revela el tipo de atacante que se considera en el análisis de privacidad. En este caso el adversario frente al que se mide la privacidad sería, evidentemente, Google.

4.6 Medición de privacidad

Como se muestra en el esquema de la Fig. 4, una vez que los perfiles real y aparente de cada usuario están disponibles, el siguiente paso es medir su privacidad.

Las métricas de privacidad usadas para este experimento se describen en la Sec. 3y, en resumen, son las siguientes:

- La entropía del perfil de usuario y,
- La divergencia de Kullback-Leibler (entropía relativa) del perfil de cada usuario con respecto al perfil de un grupo de población predefinido. En este caso el perfil

promedio de la población, obtenido a partir de la herramienta Google Ad Planer.

4.7 Ganancia de privacidad

Tal como se analizó en la Sec. 3, un valor más alto de entropía significa un valor más alto de privacidad de usuario. Además, si disminuye el valor de divergencia del perfil del usuario con respecto al perfil medio de la población, se obtiene una ganancia de privacidad. En nuestro trabajo, el nivel de ganancia de privacidad se calcula comparando el nivel de privacidad antes y después de que el mecanismo de ofuscación de TMN se ha ejecutado. Para ello se utilizan las siguientes expresiones.

Sea q el perfil de usuario, t el perfil ofuscado de usuario, p el perfil medio de la población, y M el nivel de ganancia de privacidad, tenemos:

- Para la privacidad medida como como la entropía del perfil de usuario, $H(q)$, la ganancia de privacidad relativa al nivel inicial de privacidad se define como

$$M_H = \frac{H(t) - H(q)}{H(q)}$$

- Y para la privacidad medida como la divergencia del perfil de usuario con respecto al perfil promedio de la población $D(q/p)$

$$M_D = \frac{D(q|p) - D(t|p)}{D(q|p)}$$

Estos valores darán una idea inicial de cuánto se está mejorando la privacidad cuando se usa la ofuscación de consultas. Asimismo, el cálculo de percentiles a los que estos valores pertenecen dentro de la población, brindará una medida más realista del nivel de privacidad.

4.8 Entorno de evaluación

Es importante considerar los elementos del entorno de evaluación para determinar cuál de ellos influye decisivamente en la ganancia de privacidad que obtiene el usuario a partir de la configuración de TMN.

Como se explicó previamente, las consultas falsas que se utilizaron fueron generadas artificialmente a partir de TMN. Modificamos TMN para que genere automáticamente diferentes cantidades de consultas falsas, durante 5 días, para luego verificar cómo influye en la privacidad la proporción de consultas falsas agregadas con respecto a las reales. Así, siendo ρ la relación entre el número de consultas falsas con respecto al número de consultas reales, determinamos cómo se incrementa la privacidad en función de ρ .

$$\rho = \frac{\text{Número de consultas falsas}}{\text{Número de consultas reales}}$$

Se realizaron algunos experimentos también modificando los feeds RSS que TMN usa como fuente para generar las consultas falsas.

4.9 Resultados de la medición de privacidad

4.9.1 Análisis de consultas falsas

Las 700 consultas falsas que se obtuvieron (durante 5 días) a partir de TMN se analizan inicialmente en esta sección.

Se obtuvieron los perfiles a partir de las consultas falsas para analizar cómo éstos podrían impactar en los perfiles de usuario. Cada uno de estos perfiles estuvo formado de aproximadamente unas 140 categorías de las 602 disponibles en el “categorizador” de Adnestic.

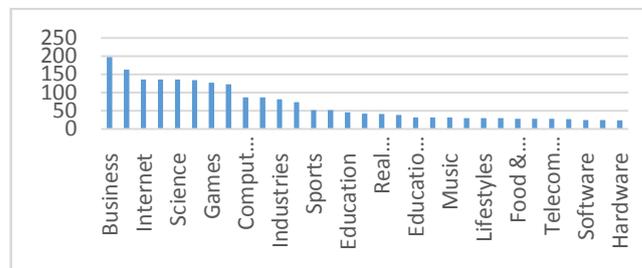


Figura. 5. Histograma de perfiles a partir de las consultas falsas generadas por TrackMeNot- Día 1

La Fig. 5 muestra uno de los histogramas de las consultas falsas obtenidos durante 5 días (sólo se muestran las categorías más populares).

Pudimos observar que, durante el experimento, las categorías a las que las consultas falsas de TMN pertenecían eran básicamente las mismas y su impacto (popularidad) similar durante los 5 días. Además, la categoría “References” aparecía repetidamente, teniendo una importante influencia en estos histogramas. A partir de nuestra experiencia con el módulo de *profiling* de Adnestic, observamos que “Reference” es una especie de categoría genérica utilizada (aunque no siempre) para clasificar todas las consultas relacionadas con información sensible (e.g. condición de salud). Esto minimiza la influencia de tales categorías en el perfil de usuario, por lo que no provee mucha información.

Los perfiles “falsos” obtenidos son muy similares entre sí, lo que significa que, al menos durante estos 5 días, los tópicos de las consultas falsas no cambiaron mucho. Esto facilitaría el trabajo de un adversario al intentar separar la influencia de estas consultas en el perfil ofuscado, en su afán de obtener el perfil real del usuario.

4.9.2 Ganancia de privacidad usando los feeds RSS por defecto en TMN

En el experimento determinamos que existía una ganancia de privacidad gracias a la ofuscación de las consultas de usuario. Esta ganancia se reflejó tanto en términos de la entropía del perfil de usuario como en términos de la divergencia de este perfil con respecto al perfil medio de la población.

En términos de entropía, la ganancia media obtenida fue de 14.58% y de 20.17% en términos de divergencia. Estos resultados se obtuvieron durante los 5 días del experimento y utilizando 100% de consultas falsas (e.g. el mismo número de consultas reales que falsas).

Los valores de entropía de los perfiles de usuario luego de ser ofuscados, también se compararon con los valores originales mediante el uso de medidas de posición no central (percentiles en este caso). Se encontró que, luego del proceso de ofuscación, el valor de privacidad de los usuarios

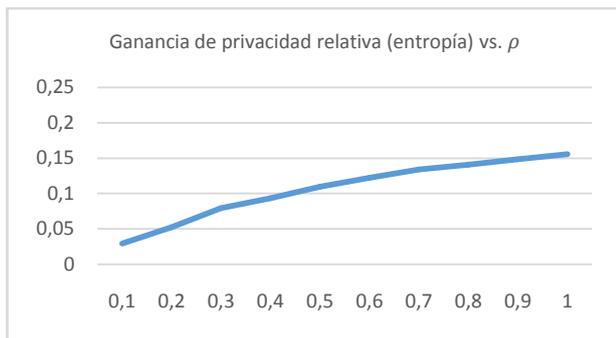


Figura. 6. Incremento de la ganancia de privacidad relativa (en términos de entropía) en función de ρ .

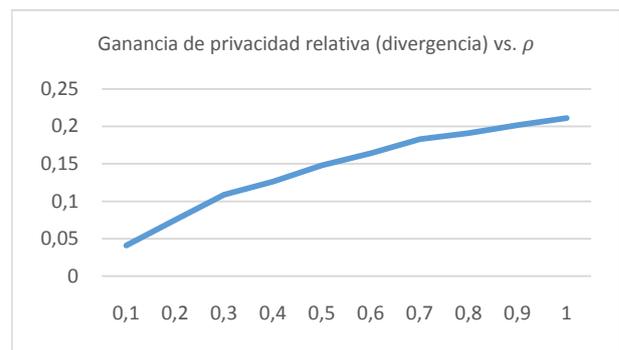


Figura.7. Incremento de la ganancia de privacidad relativa (en términos de divergencia) en función de ρ .

(entropías de sus perfiles) se incrementó en cerca de 50 percentiles de acuerdo a la distribución de entropías obtenida de la población de donde se obtuvo la muestra de consultas. Esto representa una ganancia considerable ya que un usuario que originalmente estaba en el cuarto percentil pasaba, luego de ofuscar su perfil, al quincuagésimo percentil.

4.9.3 Ganancia de privacidad con respecto a ρ

Se midió también la ganancia de privacidad con respecto a ρ (descrita en la Sec. 4.7). Se observó cómo la ganancia de privacidad se incrementa con el porcentaje de consultas falsas, tanto en términos de entropía como de divergencia. Estas mediciones se ilustran en Fig. 6 y Fig. 7.

4.9.4 Ganancia de privacidad contra ataques de clasificación

Tal como se menciona en la Sec. 3 la divergencia de KL podría ser usada como una métrica de privacidad frente a ataques de clasificación. Si disponemos del perfil de un grupo de población, podemos clasificar al perfil de un usuario como parte de éste grupo si la divergencia entre los dos perfiles es suficientemente pequeña.

Medimos, entonces, la divergencia de KL de los perfiles de los usuarios con respecto a los perfiles de algunos grupos de población (grupos de edad y género) para clasificar a los usuarios en algunos de esos grupos.

El resultado fue que casi ningún usuario se clasificó de manera distinta a la categoría asignada antes de la ofuscación del perfil, lo que significa que este mecanismo no fue exitoso frente a ataques de clasificación.

El impacto de las consultas falsas de TMN sí modifica los perfiles de usuario, incrementando, por tanto, la discrepancia con las categorías en las que los usuarios fueron clasificados originalmente. Pero este impacto no es suficiente para clasificar a los usuarios en una categoría diferente. Esto sugiere que las consultas falsas no deberían generarse aleatoriamente pues si es así, aparentemente, la influencia que tienen puede dispersarse tanto, a lo largo de múltiples categorías, que el efecto final es casi nulo. En lugar de eso, una estrategia más dirigida, con el objetivo específico de provocar una clasificación del usuario en una categoría predefinida, nos daría mejores resultados.

Es así que el reto consiste en implementar un proceso de generación más inteligente de consultas falsas, capaz de

ofuscar eficientemente los perfiles de los usuarios frente a ataques de clasificación. Esto significa que el proceso de ofuscación debería adaptarse a las necesidades de privacidad específicas de los usuarios. Sin duda, el proceso de obtención de palabras (luego combinadas para generar consultas falsas) relacionadas con un tópico específico, sin necesidad de apoyarse en terceros para ello, es un campo de investigación muy interesante.

4.9.5 Ganancia de privacidad, usando feeds RSS más específicos

La ganancia de privacidad se midió también usando feeds RSS con contenido más específico para alimentar la generación de consultas falsas. Los feeds RSS que están configurados por defecto en TMN pertenecen a sitios web de periódicos muy conocidos y apuntan a secciones donde se publica un resumen de las noticias más importantes.

En lugar de eso, configuramos en TMN feeds RSS que apunten a secciones de estos periódicos relacionadas con "Deportes", por ejemplo. La categoría "Deportes" tiende a ser una categoría de mayor interés para el público masculino. La intención, en ese sentido, fue ofuscar los perfiles de usuario con consultas falsas 'masculinas' de modo que los perfiles, inicialmente clasificados como de inclinación femenina, pudiesen ser categorizados como perfiles 'masculinos' luego de la ofuscación.

Esta estrategia modificó efectivamente los perfiles femeninos, reduciendo su divergencia con respecto al perfil medio masculino. Una vez más, esto no fue suficiente para cambiar la categoría en la que el usuario fue originalmente clasificado.

4.10 Discusión sobre los parámetros involucrados en la evaluación de TMN

Varios elementos de configuración de TMN se tomaron en cuenta para evaluar su funcionamiento. Desde el punto de vista de la entropía y la divergencia de los perfiles de usuario respecto al perfil medio de la población, la privacidad se mejora, aunque el precio es la generación de tráfico adicional. Mientras más privacidad se necesite, más consultas falsas debe generar el sistema. Esto podría tener también un impacto importante en los servicios personalizados ofrecidos por los motores de búsqueda.

La forma en la que se modelan los usuarios es también un factor importante cuando se mide la privacidad (ver Sec. 3), pero depende de las capacidades del adversario para obtener el perfil del usuario. Así, la evaluación de un mecanismo de mejoramiento de privacidad mediante el uso de distintos métodos de *profiling* podría darnos una mejor idea del rendimiento del mecanismo frente a diversos adversarios. La tarea de obtener un perfil de usuario; sin embargo, no es trivial, y no lo es tampoco saber, de manera precisa, qué técnicas usan los motores de búsqueda para este efecto.

El procesamiento de las consultas es parte de las capacidades de *profiling* del adversario que debemos simular, tanto para medir, como para evaluar un mecanismo de protección. La idea es, nuevamente, imitar las intenciones del adversario de descubrir los intereses del usuario, aunque las consultas estén escritas con errores de tipografía. El hecho de que cerca del 20% de las consultas ofuscadas no fueron exitosamente clasificadas, indica que se puede mejorar significativamente su procesamiento con el fin de obtener perfiles de usuario de mejor calidad a partir de sus correspondientes consultas.

La disponibilidad de información real de usuario, en forma de consultas o etiquetas ayuda de manera importante a interpretar los resultados de la evaluación de este mecanismo de mejoramiento de privacidad. Claramente, se puede hacer una evaluación más realista si la privacidad se mide como un valor relativo a los valores en una población.

Finalmente, esta evaluación muestra que la ofuscación de consultas no es suficientemente eficiente contra ataques más sofisticados (e.g. ataques de clasificación, ver Sec. 3) si las consultas falsas son aleatoriamente generadas. La solución radicaría en la implementación de una estrategia de generación de consultas basada en el perfil mismo de usuario y en sus necesidades particulares.

5. CONCLUSIONES

La información es actualmente una mercancía muy codiciada, y este fenómeno pone en grave riesgo a la privacidad de los usuarios. Los usuarios comunes; sin embargo, no tienen consciencia de estos riesgos y las aplicaciones que se han propuesto para proteger su privacidad se concentran, en su mayoría, en bloquear ciertas interacciones de usuario para evitar la fuga de información.

En este trabajo se ha propuesto evaluar la ganancia de privacidad que ofrece TrackMeNot a su usuario, en base a la metodología de medición de privacidad planteada en nuestro trabajo anterior en [15].

El mecanismo de perturbación de consultas de TMN ofusca el perfil del usuario efectivamente frente a ataques de identificación pero no logra resultados satisfactorios frente a ataques más sofisticados como los de clasificación.

En cuanto a los parámetros de configuración de TMN, mientras mayor es el número de consultas falsas con que se ofusca el perfil de usuario, mayor es la privacidad obtenida para éste.

El trabajo futuro involucra, por ejemplo, la recopilación de información relacionada con preferencias globales de grupos de población, que nos ayuden a medir la privacidad considerando ataques de clasificación.

Además, es necesario un método más inteligente de generación de consultas en TMN, un mecanismo capaz de cambiar su comportamiento en función de las particularidades del perfil de usuario.

Es conveniente, también, utilizar distintos modelos de *profiling* para simular varios atacantes y no sólo Google.

Otro trabajo que queda pendiente es el análisis del *tradeoff* existente entre incremento de privacidad mediante ofuscación y el consecuente incremento de tráfico generado.

REFERENCIAS

- [1] A. Erola, J. Castellà-Roca, A. Viejo, & J. Mateo-Sanz, "Exploiting social networks to provide privacy in personalized web search". *Journal of Systems and Software*, 84(10), 1734-1745, 2011.
- [2] A. Narayanan y V. Shmatikov, "Robust De-anonymization of Large Sparse Datasets", en *Security and Privacy*, 2008. SP 2008. IEEE Symposium on, C1, 2008.
- [3] C. Chow, M. Mokbel, & X. Liu, "A peer-to-peer spatial cloaking algorithm for anonymous location-based service", en *Proceedings of the 14th annual ACM international symposium on Advances in geographic information systems* (pp. 171-178). ACM, 2006.
- [4] D. Howe, & H. Nissenbaum, "TrackMeNot: Resisting surveillance in web search". Lessons from the Identity Trail: Anonymity, Privacy, and Identity in a Networked Society, 417-436, 2009.
- [5] D. L. Chaum, "Untraceable electronic mail, return addresses, and digital pseudonyms", en *Communications of the ACM*, 24(2), 84-90, 1981.
- [6] D. Rebollo-Monedero, J. Parra-Arnau, Claudia Diaz and J. Forné, "On the Measurement of Privacy as an Attacker's Estimation Error", en *Springer, International Journal of Information Security*, vol. 12, n. 2, pp. 129-149, 2013.
- [7] D. Rebollo-Monedero, J. Forné, y J. Domingo-Ferrer, "Query Profile Obfuscation by Means of Optimal Query Exchange between Users", en *IEEE Trans. Depend., Secure Comput.*, 2012.
- [8] DoNotTrackMe, [Online] Disponible: <https://addons.mozilla.org/en-US/firefox/addon/donottrackplus/>
- [9] E. Pfanner, "Internet Providers in Deal for Tailored Ads". *En The New York Times*, Technology. [Online] Disponible: http://www.nytimes.com/2008/02/18/technology/18target.html?_r=2&oref=slogin&, Feb. 2008.
- [10] Ghostery. [Online] Disponible: <http://www.ghostery.com/>
- [11] Google Sharing. [Online] Disponible: <https://addons.mozilla.org/en-us/firefox/addon/googlesharing/>
- [12] J. Becker y H. Chen, "Measuring Privacy Risk in Online Social Networks". En *Proceedings of W2SP 2009: Web 2.0 Security and Privacy*, 2009.
- [13] J. Canny, "Collaborative filtering with privacy", en *Security and Privacy*, 2002. Proc. 2002 IEEE Symposium on (pp. 45-57). IEEE, 2002.
- [14] J. Domingo-Ferrer, A. Solanas, & J. Castellà-Roca, "*k*-private information retrieval from privacy-uncooperative queryable databases", en *Online Information Review*, 33(4), 720-744, 2009.
- [15] J. Estrada-Jiménez, "Implementation of a Firefox Extension that Measures User Privacy Risk in Web Search", Master Thesis, Universitat Politècnica de Catalunya, 2013.
- [16] J. Parra-Arnau, A. Perego, E. Ferrarí, J. Forné y D. Rebollo-Monedero, "Privacy-Preserving Enhanced Collaborative Tagging", en *IEEE Trans. Knowl. Data Eng.*, 2012.
- [17] J. Parra-Arnau, D. Rebollo-Monedero y J. Forné, "A Privacy-Preserving Architecture for the Semantic Web based on Tag Suppression", en *Proc. Int. Conf. Trust, Priv., Secur., Digit. Bus.*, Bilbao, España, pp. 58-68, 2010.
- [18] J. Parra-Arnau, D. Rebollo-Monedero, J. Forné, "Measuring the Privacy of User Profiles in Personalized Information Systems", en *Future Generation Computer Systems*, 2013.
- [19] J. Parra-Arnau, D. Rebollo-Monedero, J. Forné, J. L. Muñoz y O. Esparza, "Optimal tag suppression for privacy protection in the

- semantic Web”, en *Data, Knowl. Eng.*, vol. 81-82, pp. 46-66, 2012.
- [20] K. Hafner, “Google Resists U.S. Subpoena of Search Data”, en *The New York Times*, Technology. [Online] Disponible: http://www.nytimes.com/2006/01/20/technology/20google.html?_r=1, Enero 2006.
- [21] M. Barbaro y T. Zeller Jr., “A Face Is Exposed for AOL Searcher No. 4417749”, en *The New York Times*, Technology. [Online] Disponible: <http://www.nytimes.com/2006/08/09/technology/09aol.html?pagewanted=all>, Agosto 2006.
- [22] M. Fire, D. Kagan, A. Elishar, y Y. Elovici, “Social Privacy Protector - Protecting User Privacy in Social Networks”.
- [23] M. Reiter, & A. Rubin, “Crowds: Anonymity for web transactions”, en *ACM Transactions on Information and System Security (TISSEC)*, 1(1), 66-92, 1998.
- [24] Maone, Giorgio. NoScript. [Online] Disponible: <http://noscript.net>, 2009.
- [25] P. Eckersley, “How Unique Is Your Web Browser?”. [Online] Disponible: <https://panopticklick.eff.org/>
- [26] Palant, Wladimir: Adblock Plus: Save your time and traffic. [Online] Disponible: <http://adblockplus.org/>.
- [27] Polat, H., & Du, W., “Privacy-preserving collaborative filtering using randomized perturbation techniques”, en *Data Mining, 2003. ICDM 2003. Third IEEE International Conference on* (pp. 625-628). IEEE, 2003.
- [28] R. Ostrovsky, W. Skeith III, “A survey of single-database private information retrieval: Techniques and applications”, en *Public Key Cryptography-PKC 2007* (pp. 393-411). Springer Berlin Heidelberg, 2007.
- [29] RussiaToday, “Google se enfrenta al FBI para no revelar datos privados de los usuarios”. [Online] Disponible: <http://actualidad.rt.com/actualidad/view/90908-google-fbi-revelar-datos-usuarios>, Abril 2013.
- [30] S. TejaPeddinti y N. Saxena, “On the Privacy of Web Search Based on Query Obfuscation: A Case Study of TrackMeNot”. En *10th International Symposium, PETS*, 2010.
- [31] S. Vi-a, G. News, S. Vi-b, I. Browsing, and I. Explorer, “REPRIV: Re-Envisioning In-Browser Privacy.”
- [32] TechCrunch, “AOL Proudly Releases Massive Amounts of Private Data”. [Online] Disponible: <http://techcrunch.com/2006/08/06/aol-proudly-releases-massive-amounts-of-user-search-data/>.
- [33] V. Toubiana, D. Boneh, H. Nissenbaum, y S. Barocas, “Adnostic: Privacy Preserving Targeted Advertising *”. Proc. of the 17th Annual Network and Distributed System Security Symposium (NDSS), 2009.
- [34] Y. Wang, & A. Kobsa, “Privacy-enhancing technologies. Social and Organizational Liabilities”, en *Information Security*, 203-227, 2006.