

Haciendo breve el asunto, el análisis de redes sociales comenzó creciendo lentamente, a veces a grandes saltos (ver [9, 12]), y demostrando cosas bien peculiares en sociología de las organizaciones, etnias y contactos sexuales [6]. Y ahora ha invadido la investigación sobre la Web y de cómo se hace ciencia [3].

¿Por qué estudiar o investigar redes sociales? Hay muchas razones, entre las cuales podemos considerar el interés por:

1. El estudio de La Web, que concierne directamente a la gente de ciencia de la computación. El análisis de redes sociales ha permitido descubrir propiedades de la Web, como el *Efecto Mateo* (ver más adelante) en la distribución de los vínculos. En una línea similar, veremos la relación entre PageRank [5] y el análisis de redes sociales.
2. El estudio de cómo se hace y mide la ciencia (*scientometrics, epistemometria*), que es de interés para gran parte de la comunidad científica.
3. Las ciencias sociales, ya que el análisis de redes sociales se utiliza activamente en sociología, antropología, ciencia política, gestión organizacional, medios de comunicación (*social media analysis*), etc.
4. La teoría de grafos y la matemática discreta, el fundamento técnico del análisis de redes sociales.

5. El modelamiento, la simulación y el diseño de algoritmos, habilidades y conocimientos claves que han hecho que el análisis de redes sociales escale en tamaño.

En esta exposición se presentan varias oportunidades específicas de investigación en el análisis de redes sociales.

¿CÓMO ESTUDIAMOS UNA RED SOCIAL?

Primero necesitamos recopilar información fiable y expresarla como un grafo o sociograma. Luego, analizamos el grafo para determinar propiedades de la red social original. Aquí veremos cuatro maneras de analizar esta información.

Estudiando las características generales

Es posible que existan muchos tipos de redes por ahí que faltan ser clasificadas, pero ya se observan fenómenos bien frecuentes en ellas.

Redes de mundo pequeño (small world networks). ¿Quién no ha escuchado hablar de los *seis grados de separación* o que estamos a seis personas de distancia de todo el mundo? Este fenómeno se conoce como

mundo pequeño (small world) y ocurre en las redes que tienen una conectividad especial que hace que la distancia promedio, entre dos actores cualquiera, sea muy pequeña en comparación con el tamaño (número de actores) de la red.

Stanley Milgram, un importante psicólogo norteamericano, realizó un experimento que medía la distancia promedio entre personas, en redes de contacto [9, 12]. Eligió personas de Ohama (Nebraska) y Wichita (Kansas) para que se contactaran con personas de Boston (Massachusetts). La gente de Ohama y Wichita indicaba sí conocían, o no, a las personas de Boston. En caso contrario, remitían a un contacto que pudiera conocerlas, con las que se repetía el proceso. Milgram, informado de todo esto, pudo medir cuál era el largo de los caminos recorridos. El resultado: 6 personas de distancia en promedio.

Milgram realizó muchos otros experimentos de conectividad. Muchos criticaron sus procedimientos, calificándolos como mitos urbanos, incluso recientemente [2]. Sin embargo, esta propiedad se sigue observando una y otra vez [3]. Su aparición es tan recurrente que se considera conocimiento general, e incluso es parte de teorías de innovación [13] y propuestas de abandono a los medios de comunicación masivos [10].

Redes libres de escala (scale free networks).

Muchas redes, como las de citación de artículos científicos, la Web, y muchas otras, tienen una distribución de grados que sigue una ley de potencias o similar [3]. A esas redes las llamamos *redes libres de escala* (scale free networks), porque si tomamos un subgrafo de esta red lo más probable es que los grados se sigan distribuyendo como ley de potencias.

Las redes libres de escala son interesantes porque son regidas por leyes de potencia, que se repiten en otros casos como en la distribución del ingreso; *el rico se vuelve más rico mientras el pobre se hace más pobre*, efecto tan típico que hasta aparece en la Biblia. En efecto, se le llama Efecto Mateo, y es reconocido en sociología, economía y comunicaciones, tal como lo indica el famoso sociólogo estructuralista Robert Merton [8].

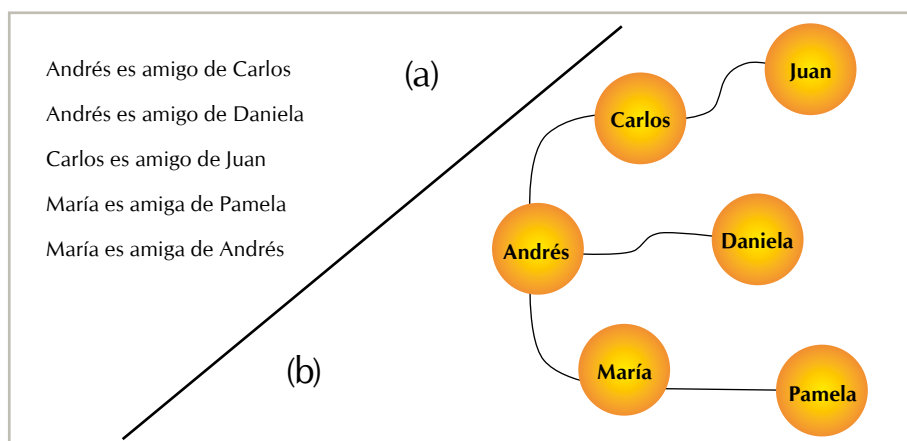


Fig. 1 Especificando una red social.

Red social que modela una situación de amistades. Cada (a) relación de amistad se traduce en un (b) vínculo en el grafo. Los vínculos no tienen dirección, porque la amistad es una relación simétrica o refleja (A es amigo/a de $B \Leftrightarrow B$ es amigo/a de A).

Estudiando la posición de los actores

El concepto tras la posición o localidad de un actor en una red corresponde al acceso que tiene al resto de la red. En principio, sabemos que dos actores ocupan el mismo lugar en la red si comparten los mismos vecinos (*equivalencia estructural*, una versión local de *isomorfismo* de vértices). Pero en general deseamos ir más lejos. En esta necesidad definimos las medidas de *centralidad*, que *miden* la posición de un actor en una red de acuerdo a ciertos criterios.

Centralidad de grado (degree centrality). Un hombre o mujer popular es aquel que tiene muchos amigos o conocidos, ¿no? Con esta simple intuición, definimos nuestra primera medida de centralidad: la centralidad de grado (degree centrality).

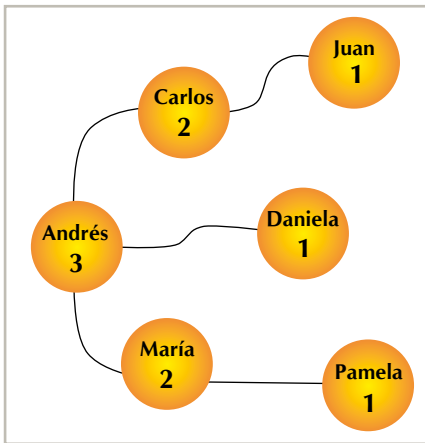


Fig. 2 Centralidad de grado. La centralidad de cada actor se calcula como su número de vecinos.

En términos de grafos, la centralidad de grado de un actor se calcula como su número de vecinos. Si estamos modelando una red social de amigos, la centralidad de cada actor consiste en su número de amigos. Un buen ejemplo se puede apreciar en la figura 2.

Centralidad $c(\beta)$ de Bonacich (Bonacich's $c(\beta)$). Hay gente cuya favorable posición en una red social les permite iniciar procesos influenciales como la transmisión de

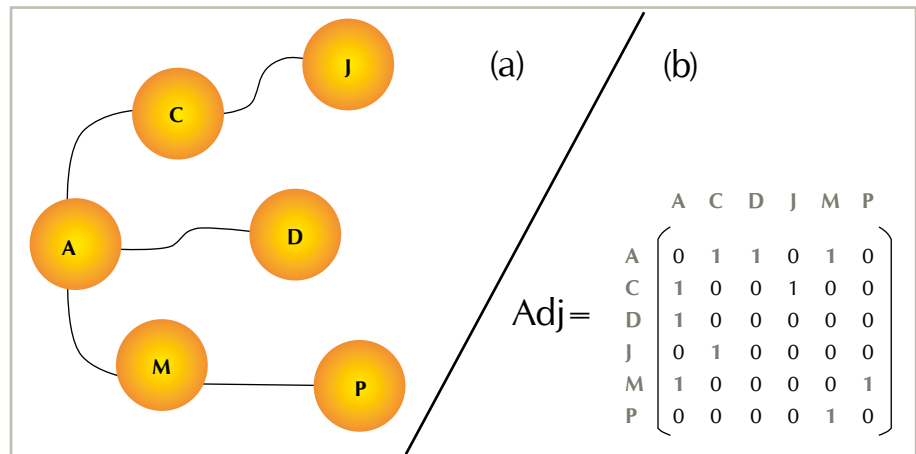


Fig. 3 Grafo y su matriz de adjacencia. Como vemos, los actores están asociados a las filas y a las columnas de la matriz de adjacencia. Si hay un vínculo entre el i -ésimo y el j -ésimo actor, entonces la componente i, j de la matriz de adjacencia será 1. De lo contrario, será 0.

creencias, chismes (*gossip*), publicidad viral, etc. En estos casos, el proceso empieza en un actor y se distribuye a su vecinos, los cuales redistribuyen a sus propios vecinos, sucesivamente. Entonces, podemos proponer una medida de centralidad para tales situaciones, que consista en contar los caminos.

Sin embargo, las redes con ciclos nos dan problemas pues tienen infinitos caminos. Por eso, no podemos contar los caminos así nada más. La solución práctica a este dilema es *atenuar* los caminos usando una *tasa de descuento*, tal como se usa en la evaluación económica, las series de potencias, etc. Así, los caminos más largos se suman como números más pequeños, y los infinitos son cero.

Usar una tasa de descuento que hace más pequeños los caminos más grandes tiene ventajas conceptuales. Los procesos influenciales como los chismes, las creencias, etc. pueden ser muy efectivos en distancias cortas, pero su difusión se hace menos efectiva (o lenta) a grandes distancias. Ajustando la tasa de descuento se puede simular cuán rápido se atenúa un proceso de difusión.

Pasando a lo matemático, definimos la centralidad $c(\beta)$ de Bonacich como $c(\beta) = (\sum_{k=1} \beta^k A^k) \cdot \vec{1}$, donde β es la tasa de descuento y A^k cuenta los caminos de

largo k entre cualquier par de actores. Esta es una propiedad de la matriz de adjacencia, la cual especifica al grafo como una matriz (vea la figura 3). Esta centralidad ha sido llamada centralidad $c(\beta)$ de Bonacich por su creador. Sin embargo, la idea es bastante antigua, casi tanto como el análisis de redes sociales.

Un ejemplo práctico de la centralidad $c(\beta)$ se muestra en la figura 4, en la cual se ilustra la función $c(\beta)$ evaluada en 0,5.

Notemos que la serie $\sum_{k=1} \beta^k A^k$ converge (si lo hace) a $\beta A(1-\beta A)^{-1}$, donde I es la matriz identidad, así que podemos calcular $c(\beta)$ como $c(\beta) = \beta A(1-\beta A)^{-1} \vec{1}$.

¿Qué valor de β usar? Obviamente, un β menor que uno, pero eso no garantiza convergencia. La solución que asegura la convergencia es usar un β menor que el valor propio de A que tiene mayor norma (recordemos álgebra lineal). Un método para calcular el vector y el valor propio más grandes se presenta en el siguiente punto.

Centralidad de vector propio (eigenvector centrality). ¿Qué tal si definimos la centralidad de manera recursiva? Digamos que un actor central es aquel que tiene un vecindario con buena centralidad. Por ejemplo, un feriante puede conocer a mucha gente común mientras que un político puede

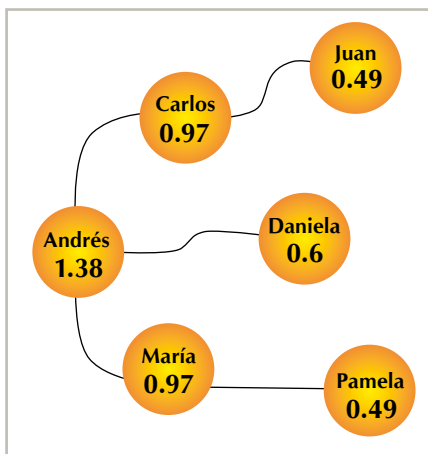


Fig. 4 Centralidad $c(0;5)$.

En el grafo, se ve la centralidad $c(\beta)$ para cada actor, evaluado con $\beta = 0;5$. Se puede ver cómo la centralidad de un actor es influenciada por los vecinos.

conocer menos gente, pero que son personas influyentes. Al final, el político está mejor posicionado en influencia que el feriante, aunque conozca menos gente (¡pero no conoce poca!). Ahora, no todo el vecindario de un actor contará con la misma centralidad, por lo que hay que considerar que los actores con mayor centralidad son más influyentes que el resto.

Siguiendo el concepto de centralidad recursiva, diremos que la centralidad de un actor es *proporcional* a la suma de las centralidades de sus vecinos en el grafo. Matemáticamente, esto se expresa como $c_i = \lambda \sum_j a_{ij} c_j$, donde λ es la constante de proporcionalidad y a_{ij} es el elemento de la fila i y la columna j de la matriz de adjacencia A de nuestra red social. En álgebra lineal, $\vec{c} = \lambda A \vec{c}$, o sea, \vec{c} es un vector propio de A (y λ^{-1} es su correspondiente valor propio). ¡Ahora queda claro el nombre de la centralidad de vector propio!

Pero el lector experto notará que, desafortunadamente, este problema no tiene solución única; en general, una matriz tiene varios vectores propios diferentes. ¿Cuál usar? Cualquier valor propio positivo tiene sentido en la ecuación, ¿no? Bueno, los investigadores decidieron dejar todo en el valor propio más alto, por una propiedad muy sencilla: su vector propio

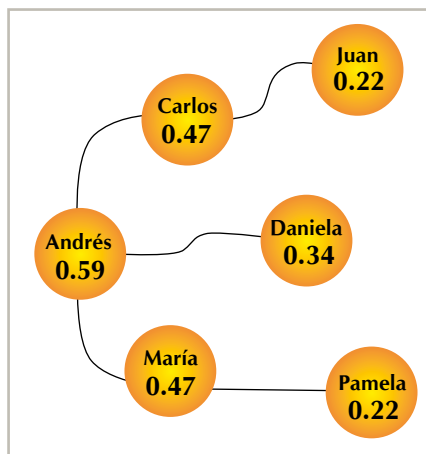


Fig. 5 Centralidad de vector propio.

En el grafo, se ve la centralidad de vector propio, calculada con $\vec{c} = (A^{50} \vec{1}) / \|A^{50} \vec{1}\|$. Se puede ver cómo la centralidad de un actor es influenciada por sus vecinos.

asociado es fácil de calcular. La sucesión $\vec{c}_{k+1} = A \vec{c}_k / \|A \vec{c}_k\|$, que inicia $\vec{c}_0 = \vec{1}$, nos permite obtener rápidamente este valor.

De manera más sencilla, podemos aproximar la centralidad de vector propio como $\vec{c} \approx A^k \vec{1} / \|A^k \vec{1}\|$, para un k adecuado tal que la aproximación varíe poco (para k y $k+1$ hay diferencias minúsculas). Esto se puede asegurar con un valor de k grande (por ejemplo, 50). Un ejemplo práctico que usa esta aproximación se ve en la figura 5.

Centralidad de cercanía (closeness centrality). Otra medida de posición sale de considerar la distancia promedio al resto de la red. El actor que está más cerca de todo otro elemento de la red es el más central, versa la idea tras la centralidad de cercanía (*closeness centrality*). Pero nosotros no medimos “cercanías”, sino “distancias”; o sea, lo contrario.

Matemáticamente podemos expresar la centralidad de cercanía como el inverso a la suma de las distancias, o sea, $c_i = \frac{1}{\sum_j d_{ij}}$, donde d_{ij} es la distancia entre el actor i y el actor j . Obviamente hablamos de distancias euclídeas o rutas mínimas.

Centralidad de intermediación (betweenness centrality). Digamos que un grupo terrorista

se toma Linares (¡en Linares no pasa nada!) e impide el movimiento de camiones entre las zonas central y sur del país. ¿Qué hacen? Desconectan a Chile. Linares, por inocente y tranquilo que parezca, es una ciudad clave en la red de suministro de Chile, pues es camino obligado. O sea, todos los caminos pasan por Linares. (La realidad es que es fácil hacerle el quite, pero éste es un ejemplo.) Esto nos inspira a repensar las medidas de centralidad.

En estrategias militares y terrorismo, es importante distinguir los actores claves. Si son atacados, desconectan una red, o interrumpen sustancialmente los flujos que se pudieran producir en ésta. Estos actores, que son objeto de ataque y defensa, se pueden descubrir contando cuántas rutas mínimas pasan por ellos; o sea, por su calidad de intermediarios o puntos intermedios. Por eso definimos la centralidad de intermediación como el número de rutas mínimas en las que el actor participa.

Estudiando los grupos que tiene

La detección de comunidades, grupos, cliques (grupos exclusivos), etc. es tema de alto interés en redes sociales. El asunto es complicado pues no es fácil definir un grupo. La definición es fácil cuando hablamos de una estructura formal, cuando existe un grupo definido y un grupo de adherentes que dice ser parte del grupo. Por ejemplo, Chile y los chilenos. Pero todo se vuelve complicado, oscuro, hasta esotérico cuando hablamos de la estructura *informal*. Un grupo de amigos es un montón de gente que son todos o casi todos amigos entre sí, pero ellos a su vez tienen varios amigos comunes... ¿Cuáles pertenecen al grupo y cuales no?

Técnicas de detección. Técnicas para detectar grupos hay muchas; hay muchos algoritmos, con muchas velocidades diferentes, que obedecen a diferentes ideas de cómo se detecta un grupo, situación muy opuesta a la de las medidas de centralidad.

Una manera tradicional consiste en reducir la detección de grupos a una clasificación o *clustering*. Dentro de estas técnicas están el

tradicional k-Means, los algoritmos genéticos, el análisis de *modularidad* (el número de vínculos entre grupos es pequeño, dentro de grupos es alto), etc. Adicionalmente, estas técnicas son parametrizables (i.e. número de clases en k-Means, modularidad mínima, etc.), lo que permite analizar la calidad de la clasificación. Aquí se hace posible usar *árboles jerárquicos* para decidir cuándo la clasificación es buena.

Otra manera tradicional consiste en ver el problema como uno de teoría de grafos. Por ejemplo, la *coloración* es una forma tradicional de hacer clasificación en grafos. En este caso, también es posible ver el problema como uno de *equivalencia estructural* transformado a uno de *equivalencia regular*: “en un grupo de amigos, los amigos compartimos los mismos amigos” (esto define un algoritmo iterativo). Adicionalmente, se pueden buscar *cliques* y *k-Cliques* para encontrar los grupos.

Maneras más novedosas incluyen el uso de las medidas de centralidad: “en un grupo, todos los actores son cercanos” (centralidad de cercanía), “un grupo es una red más o menos aislada del resto” (centralidad de intermediación). También se incluyen medidas *democráticas*, que consisten en consensuar dos o más criterios diferentes de detección de grupos.

Una guía a la historia de los algoritmos de clustering la hace Freeman [4]. Otra revisión más sintética la entrega Boccaletti [3].

Los problemas. Aún queda mucha investigación por hacer en el tema de detección de comunidades. Aquí listo algunos de esos desafíos.

Entre los problemas conceptuales, nos encontramos con: ¿Qué es un grupo o una comunidad? El caso formal es sencillo, pues los actores declaran su pertenencia a un grupo, pero el caso informal es bastante complejo. Asimismo, tenemos el siguiente problema: ¿Cuándo es conveniente comparar con un caso formal? También debemos considerar que los actores pueden pertenecer a varios grupos, cosa que muchos algoritmos no admiten.

Entre los problemas de eficiencia (rapidez), notamos que hay muchos algoritmos que son NP-HARD por tratar de cumplir exactamente una definición, lo que motiva a usar aproximaciones. Pero debemos ser aún más exigentes: si hablamos de cientos, miles o millones de vértices, un algoritmo de orden $Q(n^3)$ puede ser desastrosamente lento. ¡No podemos conformarnos sólo con estar en la clase P (tiempo polinomial)! Las redes son cada vez más grandes, y la necesidad de algoritmos más rápidos es cada vez mayor.

Como indicamos, se hace necesario aproximar en muchos casos. Sin embargo, esto supone nuevos desafíos: ¿Cuándo es bueno usar algoritmos aleatorios? ¿Cuándo es bueno usar heurísticas? Más aun, debemos tener en cuenta que pueden haber actores y vínculos que no consideramos en la construcción del grafo, por lo que nuestros algoritmos deben ser precisos incluso cuando la información escasea o falla.

Personalmente tengo las siguientes interrogantes, las que, de ser contestadas, podrían dar origen a varias publicaciones:

1. ¿Puede un grupo contener otros grupos? Esto ocurre en las estructuras sociales formales.
2. ¿Cómo generar algoritmos de clustering para nuevas representaciones gráficas?
3. ¿Cómo podemos definir clusters cuando hay grafos dirigidos? ¿Qué ocurre en el caso de grafos con pesos?
4. Si consideramos un algoritmo iterativo, ¿podemos usar un resultado aproximado para generar otro más preciso? (Combinar algoritmos.)

Visualización

La visualización de las redes sociales también sirve como método para descubrir propiedades de ésta, aunque tiene menos peso teórico en el análisis. Pero cuenta con la ventaja de alimentar rápidamente la intuición del investigador.

Visualizar redes complejas es un gran desafío; por lo general, se busca presentar gran cantidad de información de forma



Fig. 6 Visualización de una red social.

estética. Se busca la claridad y la simpleza, pese a la gran complejidad de los datos, como se ilustra en el ejemplo de la figura 6. Y hay que considerar que hay muchas potenciales vistas de los datos, que pueden ilustrar propiedades diferentes: centralidad, comunidades, jugadores clave (que, si desaparecen, desconectan la red), etc.

Tal como en la detección de comunidades, existe una gran variedad de algoritmos para visualizar redes sociales. Cada uno obedece a una idea u objetivo diferente, aunque muchas veces se busca la presentación instantánea.

¿CÓMO OBTENEMOS REDES SOCIALES?

Uno de los desafíos con las hipótesis sobre redes sociales es reproducir los fenómenos que se dice que ocurren.

La Web. Una de las redes más estudiadas en computación es la Web, cuya relevancia al área es clarísima. La Web es una red gigante, masiva, en donde participan millones de páginas con vínculos entre sí. Notemos que hay muchos tipos diferentes de páginas web; estáticas y dinámicas (que se generan en el vuelo), que se actualizan, algunas que se borran, otras se crean; hay buscadores con vínculos a grandes porciones de la red, catálogos, sitios de noticias, blogs, foros, sitios de fotografías, de vídeos, bibliotecas, sitios corporativos, etc. los cuales están llenos de páginas y vínculos.

Los desafíos que pone la Web para su estudio caen en los problemas de recuperación y consulta de información. Este es el problema tradicional de los buscadores, que deben buscar y buscar páginas web, siguiendo vínculos, y deben recuperar y clasificar su contenido. Luego, deben explotar la bases de datos construidas para responder las consultas.

Una de las grandes aplicaciones de las medidas de centralidad se da justamente en la red; Google, en vez de buscar páginas por la calidad de su contenido, las busca por su fama. Cada página tiene una nota, un ranking, que sale de la fama de las páginas que la referencian. Este algoritmo, llamado Page Rank, es justamente una aproximación de la *centralidad de vector propio*. Veamos la propia explicación que da Google:

(...) En lugar de contar los vínculos directos, Page-Rank interpreta un vínculo de la página A a la B como un voto para la página B por parte de A. (...) Esta tecnología también tiene en cuenta la importancia de cada página que efectúa un voto, dado que los votos de algunos se consideran de mayor valor, con lo que incrementan el valor de la página a la que enlazan. [5]

Estudiar la Web es un asunto colosal; la Web es demasiado grande, por eso se hacen estudios locales. Por ejemplo, en Chile se puede explotar el registro de NIC, con lo que se ha realizado el Estudio de la Web Chilena [1].

Sitios de redes sociales, Web 2.0. Muchos de los datos de redes sociales son recopilados de los sitios sociales, sitios extremadamente populares cuya estructura no está completamente predefinida sino que se construye dinámicamente de acuerdo a las acciones de sus usuarios. En estos sitios se suele cumplir con los 5 ó 6 grados de separación de Milgram.

Entre los sitios sociales se cuentan: Facebook, Wikipedia, Fotolog, Habbo, Last.fm, Orkut, Youtube, MySpace, Xing, Flickr, Piccasa, Advogato, SourceForge, MyHeritage, aSmallWorld, Broadcaster.com, Classmates.com, DeviantART, Twitter, Sonico.com, etc. (Más de alguno de estos sitios debiera sonar conocido).

Otros registros digitales. Los seres humanos solemos dejar huellas de nuestras interacciones en los medios digitales, más allá de los sitios sociales. Por ejemplo, foros, news, irc, email, CVS-SVN, comercio electrónico, telefonía IP, etc. son evidencias digitales de interacciones humanas. Todas éstas están sujetas a estudio. Sin embargo, aparecen los dilemas éticos de la información confidencial que los sitios sociales expresamente hacen pública.

Encuestas. Este es el método de recuperación de información social más usado en el estudio de las redes sociales tradicionales que están fuera de la Web. Sin embargo, suele ser muy caro realizar estudios de este tipo, sobre todo en papel. Mas aun, muchas veces se requieren investigadores en terreno que supervisen el correcto proceder de las encuestas. Versiones más baratas son las encuestas por Teléfono e Internet, aunque su validez es limitada.

Simulación. Esta es una técnica muy usada por bastantes científicos sociales que trabajan con comunidades artificiales. Sin embargo, su uso aparece más útil en la comprobación de hipótesis que en el análisis desde cero. Usando simulación se puede responder a una pregunta como "¿este proceso social genera redes con estas características?". Luego, las redes artificiales y las reales se comparan con las técnicas de análisis presentadas previamente, y se puede concluir el alcance de una hipótesis.

CONCLUSIONES

El análisis de redes sociales es un área que presenta muchas oportunidades de investigación para la gente de ciencia de la computación. Vimos los exigentes desafíos algorítmicos que propone el área, las diversas métricas que se obtienen de los grafos, su incidencia en el estudio de La Web y el diseño de buscadores, incluso ligeramente el uso de la simulación en el área (que se vio muy poco para lo extendido que es su uso). El área va mucho más allá de lo que son la informática y computación sociales; da a la computación y la matemática discreta un lugar privilegiado en la teoría social. BITS

REFERENCIAS

- [1] R. Baeza-Yates, C. Castillo, E. Graells. "Características de la Web Chilena 2006". Centro de Investigación de la Web. 2006. http://www.ciw.cl/material/web_chilena_2006/index.html
- [2] Blastland interviewing Kleinfield. "Connecting with people in six steps". More or Less. BBC News. http://news.bbc.co.uk/1/hi/programmes/more_or_less/5176698.stm
- [3] S. Boccaletti et al. "Complex networks: Structure and dynamics". Physics Reports 424, 2006, 175-328.
- [4] L. Freeman. "Finding Social Groups: A Meta-Analysis of the Southern Women Data". <http://moreno.ss.uci.edu/85.pdf>
- [5] Información corporativa de Google: tecnología. Visto en Octubre de 2008. <http://www.google.cl/corporate/tech.html>
- [6] E. Lawmann. "A 45-year retrospective on doing networks". Connections 27(1), 65-90. 2006.
- [7] "mc-50 map of FlickrLand: flickr's social network". <http://www.flickr.com/photos/gustavog/4499404/in/set-113313/>
- [8] R. Merton. "The Matthew Effect in Science". Science, 159 (3810): 56-63, Enero 5, 1968. <http://www.garfield.library.upenn.edu/merton/matthew1.pdf>
- [9] S. Milgram. "The Small World Problem". Psychology Today, 1967, Vol. 2, 60-67.
- [10] L. De Rossi. "The Power Of Open Participatory Media And Why Mass Media Must Be Abandoned". Robin Good, Master New Media. Visto en Octubre de 2008. http://www.masternewmedia.org/news/2006/03/20/the_power_of_open_participatory.htm
- [11] J. Scott. "Social Network Analysis: A Handbook". Sage Publications. 2000.
- [12] Travers, Jeffrey, and S. Milgram. "An Experimental Study of the Small World Problem". Sociometry 32, 1969, 425-443.
- [13] D. Ward. "Knock, Knock, Knocking on Newton's Door: Building Collaborative Networks for Innovative Problem Solving". Defense AT&L Journal. Marzo-Abril de 2005. http://www.dau.mil/pubs/dam/03_04_2005/war-ma05.pdf