

Navegando las redes sociales



Desde hace algunos años los usuarios de la Web, nos hemos visto expuestos a una inmensidad de información generada por otros humanos. En particular, por lo que se denomina *Social Media Data*, o Datos de Medios Sociales, que están compuestos por mensajes (texto), fotos, vídeos, etc. El éxito de cientos de aplicaciones sociales, entre las cuales destacan Facebook, Twitter, YouTube, Foursquare, Flickr e Instagram, ha llevado a que millones de personas se vuelvan usuarios activos en la generación y publicación de contenido en la Web. Toda esta información, sin duda alguna, tiene un valor incalculable para investigadores de muchas áreas, como la Sociología, la Computación, la Medicina y la Biología, por sólo nombrar algunas.

Esta revolución del contenido digital ha cambiado el mundo, pasando por cómo nos informamos hasta cómo nos relacionamos con otras personas. Si antes íbamos a Google News o a la CNN para tener las noticias más recientes (o

breaking news), ahora todos los usuarios de Twitter saben que la forma más rápida de informarse es a través de este medio. Twitter, incluso, ha sido incorporado en los medios tradicionales de noticias y en buscadores Web, como fuente primaria de información en tiempo real.

A su vez, el uso de redes sociales está transformando el mundo hacia sociedades más cosmopolitas, permitiendo establecer vínculos sociales activos con personas localizadas en lugares geográficos distantes. Ya la distancia no es un impedimento para la interacción, sino que está dada principalmente por la similitud de intereses entre los usuarios de las diferentes plataformas sociales.

Desde el punto de vista computacional, el análisis de Datos de Medios Sociales, presenta diversos desafíos de gran complejidad:

- Gran volumen de datos, que ascienden a Gb o Tb diarios incluso, dependiendo del tipo de medio.



Bárbara Poblete

Profesora Asistente DCC Universidad de Chile. PhD en Computación, Universitat Pompeu Fabra (2009); Magíster en Ciencias mención Computación, Universidad de Chile (2004); Ingeniero Civil en Computación, Universidad de Chile (2004). Líneas de Especialización: Minería de grandes volúmenes de datos; Minería de Logs de Buscadores; Privacidad de Datos en la Web; Análisis de Redes Sociales en línea.
bpoblete@dcc.uchile.cl

- Análisis en tiempo real de un *stream* o flujo constante de información.
- Crear algoritmos que sean capaces de descubrir información valiosa e interesante a partir de datos de baja calidad o ruidosos.
- Preservar la privacidad de las personas, dependiendo de qué información compartan públicamente.
- El trabajo interdisciplinario que involucra ser un “Data Scientist” al servicio de variadas áreas de investigación.

En mi investigación, me ha tocado pasearme por estos diferentes desafíos y todos me resultan extremadamente interesantes. Lejos de estar resueltas, estas problemáticas aumentan día a día su complejidad en la medida que aumenta la información de los medios sociales.

En particular, en este artículo describiré parte de mi trabajo realizado sobre la red social Twitter, además de otros trabajos que me parecen destacables. Todos con un tema en común: la explotación de los datos generados masivamente por usuarios como fuente de riqueza nueva.

ANÁLISIS DE REDES SOCIALES O “ME ACABO DE COMER UN SÁNDWICH”

Twitter es una red social que permite la publicación de mensajes cortos, denominados *tweets*, con un largo máximo de 140 caracteres. Es reconocida como una plataforma de *microblogging* que además permite a los usuarios suscribirse a lo que dicen otros.

Me ha tocado en más de una oportunidad encontrarme con científicos seniors y muy serios, que me dicen que no entienden cómo se puede hacer Ciencia usando información de redes sociales como Twitter. Argumentan que cómo es posible sacar algo útil de un comentario como: “Me acabo de comer un sándwich”. Es verdad, Twitter a pesar de ser caracterizado como un medio informativo, también está lleno de comentarios personales que carecen de interés general. ¿Pero en realidad son inútiles? Muchos pensamos que no.

Existen muchas cosas interesantes que se pueden sacar a partir de un comentario personal y aparentemente superficial, como por ejemplo: si el comentario trae alguna referencia a un lugar o una etiqueta de geolocalización, podremos saber dónde ocurrió el evento descrito. Luego si muchos usuarios de la red deciden comer sándwiches en ese lugar, quizás sea una buena recomendación para entregar a otros usuarios que viven/trabajan en la zona. También podemos mirar si el comentario está acompañado de algún adjetivo positivo o negativo, como “me acabo de comer un sándwich delicioso” o “me acabo de comer un sándwich malísimo”. Este sentimiento, junto con otros más, se puede analizar automáticamente por medio de heurísticas, para así establecer un sentimiento general con respecto a un lugar, evento o situación. Quizás esto a pequeña escala pueda parecer anecdótico, pero las redes sociales y la automatización de los procesos de análisis, nos permiten llegar a observar a millones de usuarios que publican miles de millones de mensajes. Es así como podríamos generar automáticamente una lista de los mejores restaurantes de América Latina, de Santiago o del mundo. Todo esto sin pedirle ayuda a un crítico gastronómico ni tener que lidiar con los costos de realizar encuestas a personas.

Ahora, hay proyectos mucho más ambiciosos en los que podemos pensar al tener a la mano tanta información. Muchos científicos computacionales, se han ido especializando en el área de la Sociología Computacional y a su vez, sociólogos se han asociado con expertos computacionales. El resultado de

esta incursión es una serie de publicaciones en el área de Minería de Datos, que analizan el comportamiento humano, a partir de datos agregados de redes sociales. Dentro de esto, un estudio que me parece interesante de compartir es el trabajo publicado en Science por científicos de Cornell [Golder & Macy, Science 2011] en el que se hace un análisis a nivel mundial, de usuarios de Twitter en 84 países. Para esto recolectaron información de 2,4 millones de usuarios, a través de unos 500 millones de mensajes en Twitter. El objetivo era monitorear el estado anímico de las personas y cómo éste variaba durante el transcurso del día. Este estudio se expandió por un período de dos años en que los investigadores establecieron que el trabajo, las horas de sueño y la cantidad de luz solar, influyen en emociones cíclicas como el entusiasmo, el miedo, el enojo, el estado de alerta y la satisfacción. Todo esto puede parecer obvio, ya que hace años que se habla de los ritmos circadianos existentes en los estados anímicos, pero hasta ahora esas teorías sólo habían sido evaluadas en laboratorios usando pequeños grupos homogéneos de individuos y por cortos períodos de tiempo. Sin embargo, por medio de investigaciones como éstas, se ha podido validar ésta y otras teorías existentes. En este caso, se ha observado que las mismas emociones cíclicas se repiten en todo el mundo, en forma independiente a la cultura y país de las personas. También se ha concluido que los distintos ciclos se desplazan en función de las horas de luz solar de cada región geográfica y de los días que son considerados como laborales en cada país.

Recopilamos datos de la red social por un lapso de un año, durante 2010, obteniendo información de aproximadamente cinco millones de usuarios activos y cinco mil millones de mensajes, transformando este estudio en el más grande de estas características que se ha publicado.

En una línea similar, junto con colegas de Yahoo! Research Barcelona y la Universidad Técnica Federico Santa María, llevamos a cabo un caracterización de gran escala a nivel de los diez países más activos en Twitter del mundo [Poblete, B. et al. CIKM 2011]. Para esto recopilamos datos de la red social por un lapso de un año, durante 2010, obteniendo información de aproximadamente cinco millones de usuarios activos y cinco mil millones de mensajes, transformando este estudio en el más grande de estas características que se ha publicado. El objetivo de esta publicación fue mostrar un análisis preliminar de diferencias y similitudes entre países, en términos de actividad en la red social, sentimiento, uso del lenguaje y características de sus estructuras de las redes. A continuación resumo a grandes rasgos algunos de nuestros hallazgos:

Actividad. En la Figura 1 se puede apreciar la cantidad de actividad, o mensajes por usuario, observada para cada uno de los países más activos. Se debe notar, que los países más activos no son necesariamente los que tienen usuarios más activos. La Figura 2 muestra los países con más mensajes por usuario, indicando que hay países que tienen muchos menos usuarios en Twitter, pero que estos participan mucho más activamente que en otros países donde la comunidad es más extensa.

Idiomas. Utilizando un clasificador de idiomas, clasificamos cada uno de los mensajes de nuestro set de datos, con un resultado del 99% de los mensajes clasificados en 69 idiomas. En la Figura 3, se pueden ver los idiomas más comunes, siendo el inglés el idioma más popular presente en el 53% de los mensajes. La Figura 4 nos permite apreciar los tres idiomas más utilizados en los diez países más activos de la red social, en donde el inglés continúa apareciendo siempre entre los idiomas más utilizados. A pesar de tener un buen desempeño, la tarea de la clasificación de idiomas no es trivial en los mensajes o tweets, ya que están plagados de abreviaciones, modismos y faltas de ortografía.

Sentimiento. También estudiamos el componente de sentimiento de los mensajes, utilizando la métrica de *happiness* (o felicidad) acuñada por Dodds et al. [DODDS],

Figura 1

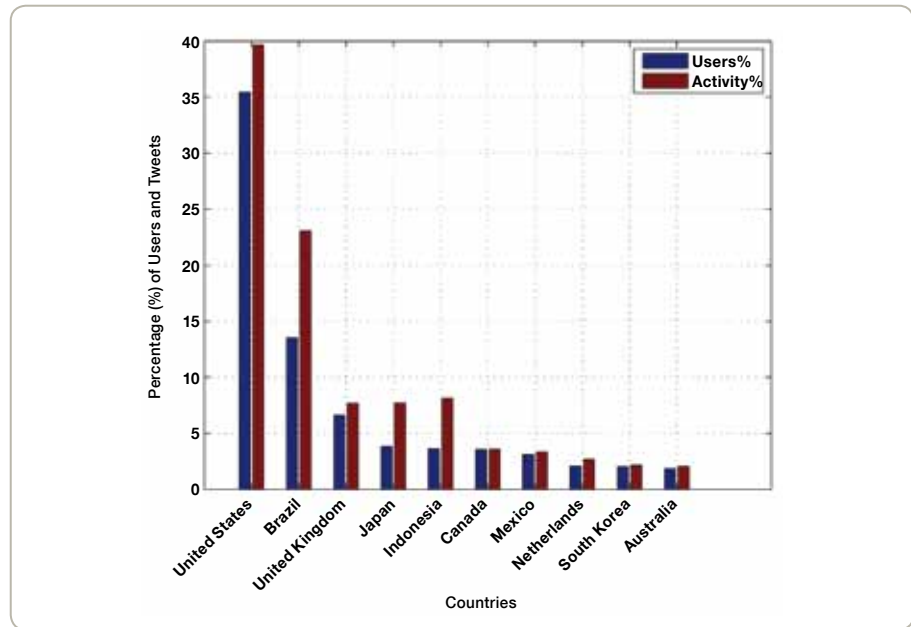
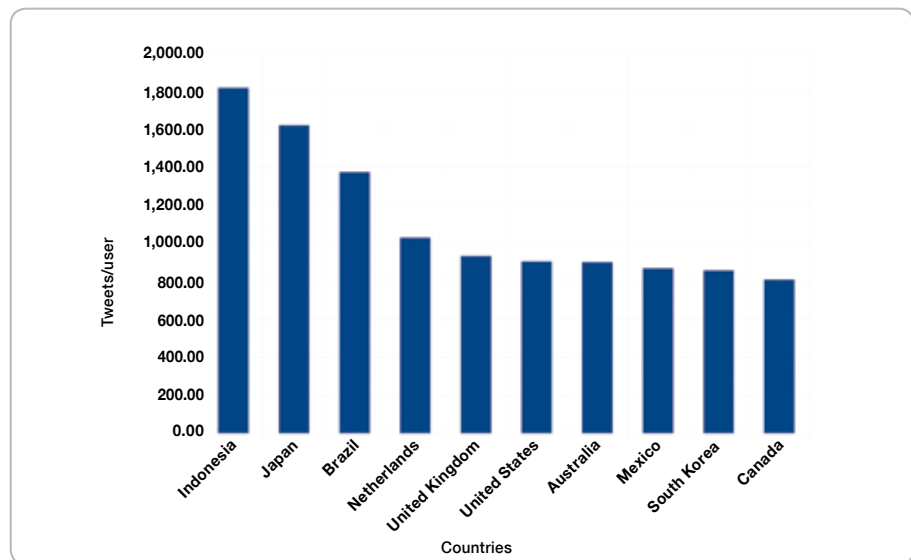


Figura 2



que es comúnmente conocida como *valencia*. Este valor representa la reacción psicológica que tienen los humanos a palabras específicas, de acuerdo a una escala que varía desde “feliz” a “infeliz”. En particular nos remitimos a analizar la valencia de los mensajes clasificados como escritos en inglés y español, para los cuales existen listas estándares de palabras y sus valencias [Bradley and Lang, Redondo et al.]. Los resultados obtenidos se pueden observar en la Figura 5, en donde se ve que los niveles de felicidad aumentan hacia fin

de año y, sin mucha sorpresa, se aprecia que Brasil presenta los valores más altos de felicidad para cada mes. Sin embargo, al igual que en el análisis de idiomas, existen factores que pueden afectar la certeza de esta métrica, como el uso de ironías en las conversaciones.

Contenido de los mensajes. Analizamos brevemente algunas propiedades de los mensajes o tweets de cada uno de los diez países más activos. Entre estas características se encuentran:

Figura 3

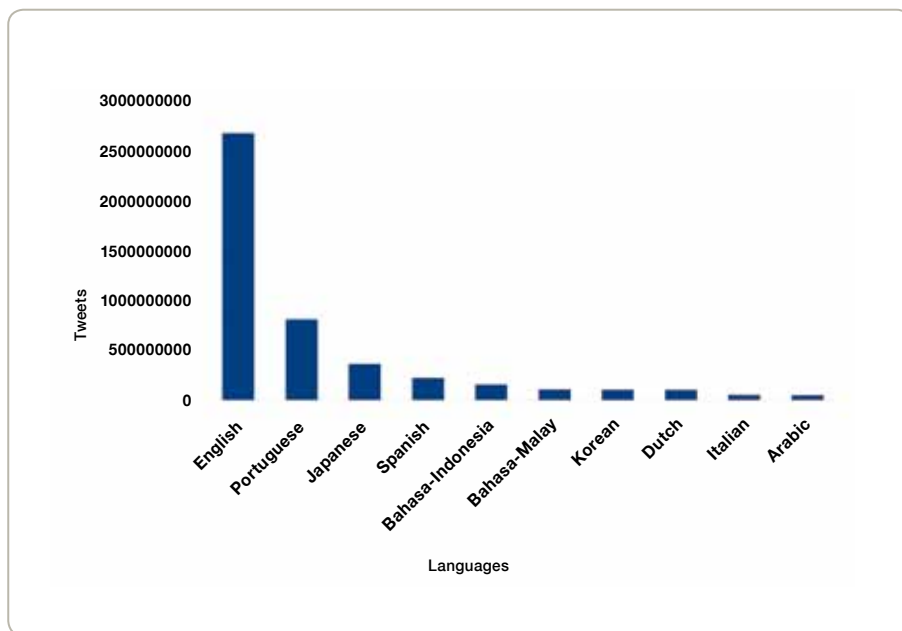
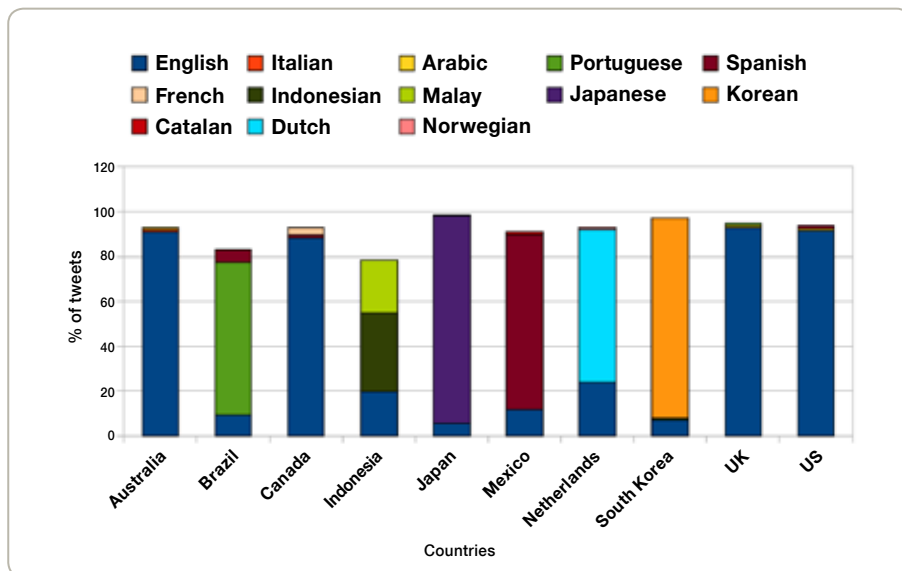


Figura 4



- **#:** este símbolo indica si el mensaje ha sido denominado con un tema en particular o *hashtag*.
- **RT:** nos indica si un mensaje contiene el texto "RT", que indica que corresponde a una republicación de un mensaje o *retweet*.
- **@:** indica si el mensaje contiene un símbolo '@', que se usa junto a un nombre de usuario para indicar una mención a éste.

- **URL:** Indica si el mensaje contiene una URL o no.

Utilizando estas características se calculó su promedio para cada país, que se muestra en la Tabla 1. En esta Tabla se observa que Indonesia y Corea del Sur tienen el porcentaje más alto de menciones de usuarios, contrastando con Japón que tiene el más bajo. Además Japón es el país en que se hace el menor porcentaje de propagación de mensajes (RT). Países Bajos es donde

más se utilizan los hashtags por usuario, y en Estados Unidos se utilizan más URLs en los mensajes, lo cual indicaría un mayor uso de Twitter como medio para difundir información formal.

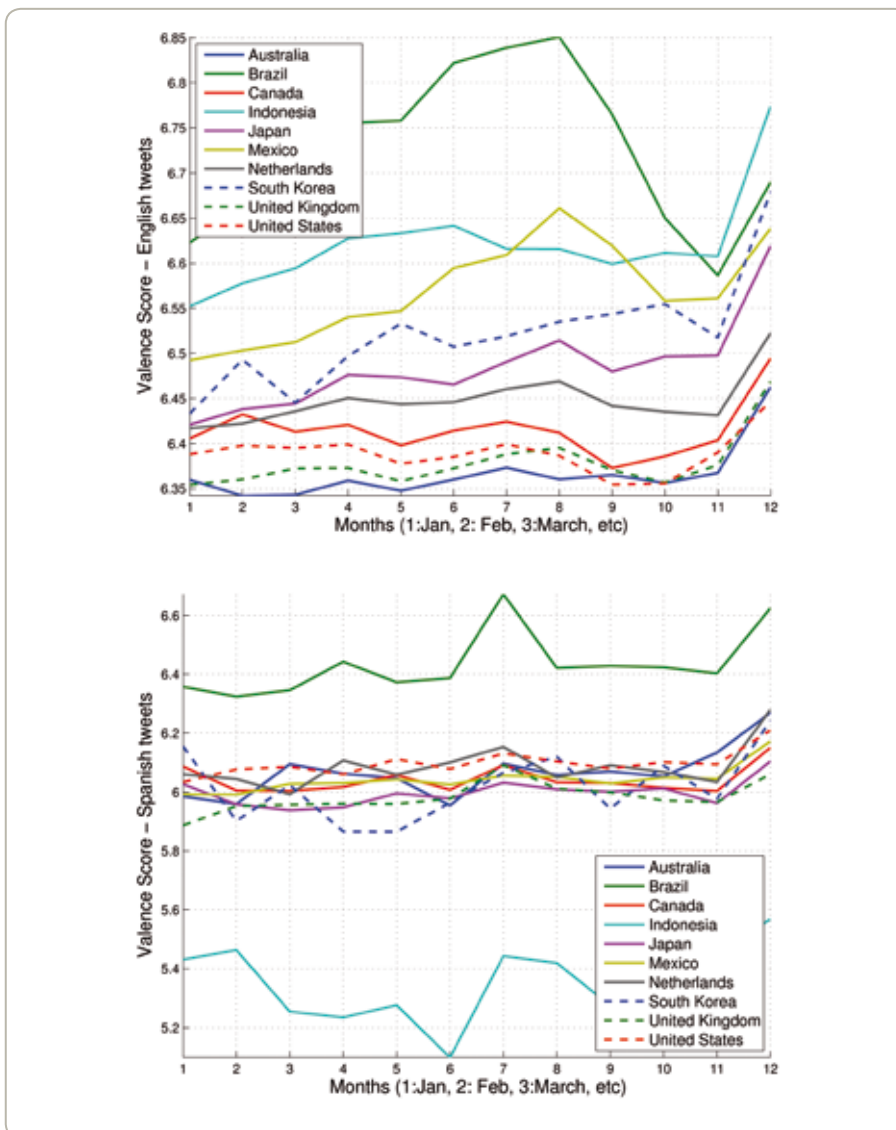
Estructura de red. Twitter provee una red social para sus usuarios, es decir, les permite conectarse por enlaces dirigidos que posibilitan relaciones no recíprocas entre personas (al contrario de Facebook en donde los enlaces son recíprocos). En Twitter, los usuarios tienen la opción de seguir o no seguir a otros usuarios, suscribiendo a los mensajes que estos publican. Estas conexiones entre usuarios pueden visualizarse como un gran grafo dirigido.

Las estructuras de las redes sociales en Twitter son altamente dinámicas, es decir, están cambiando en todo momento. Debido a esto, elegimos realizar un análisis estático de sus características usando una captura o *snapshot* de su estructura para una ventana de tiempo acotada. Para analizar la red de cada país se consideraron los usuarios pertenecientes a éste y los enlaces autocontenidos en este conjunto, sin considerar enlaces externos. En este sentido, un trabajo anterior realizado por Mislove et al. [MISLOVE, 2007], indica que muestras parciales de estos grafos pueden llevar a subestimar medidas como distribución del grado, pero continúan preservando otras métricas, como la densidad, reciprocidad y conectividad. Por lo tanto, al preservar en nuestro estudio el componente activo del grafo, estamos analizando la parte más relevante de su estructura social. Las estructuras resultantes mostraron diferencias notorias entre los países y sus organizaciones sociales. Por ejemplo, observamos que para algunos países la reciprocidad entre usuarios es mucho más importante que en otros, como es el caso de Japón, Corea del Sur, Indonesia y Canadá. La naturaleza simétrica de estas conexiones afecta la estructura de sus redes, aumentando la conectividad y reduciendo su diámetro. Con respecto al grado de los usuarios (o grado de los nodos del grafo) se observa que Estados Unidos y Corea del Sur son los países con mayor grado promedio. Esto significa que los usuarios de estos países tienden a tener más seguidores y a seguir más gente que

Tabla 1:
Average usage of features per user for each country

Country	Tweets Users	(URL)%	(#)%	(@)%	(RT)%
Indonesia	1813.53	14.95	7.63	58.24	9.71
Japan	1617.35	16.30	6.81	39.14	5.65
Brazil	1370.27	19.23	13.41	45.57	12.80
Netherlands	1026.44	24.40	18.24	42.33	9.12
UK	930.58	27.11	13.03	45.61	11.65
US	900.79	32.64	14.32	40.03	11.78
Australia	897.41	31.37	14.89	43.27	11.73
Mexico	865.7	17.49	12.38	49.79	12.61
S. Korea	853.92	19.67	5.83	58.02	9.02
Canada	806	31.09	14.68	42.50	12.50

Figura 5



en otros países. Por otra parte Indonesia presenta el grado promedio más bajo por nodo, a pesar de ser una comunidad altamente activa.

Entre muchas otras cosas, en nuestro análisis evaluamos métricas como, la densidad de las redes de los países, que nos indica qué tan bien conectados están los usuarios entre sí. De aquí se desprende que Estados Unidos es el país con menor densidad y Corea del Sur es el país con mayor densidad, junto con Australia y Países Bajos. Esto indica que las comunidades más pequeñas están mejor conectadas que las más grandes. También estudiamos el coeficiente de *clustering* promedio, observando que los países que poseen alto nivel de clustering pero baja reciprocidad, muestran indicios de ser comunidades más jerarquizadas (i.e., dos usuarios que comparten un enlace recíproco entre sí, siguen a otro usuario que no es recíproco).

La reciprocidad nos habla sobre el nivel de cohesión, confianza y capital social en sociología [REF7]. En este contexto la tendencia de algunas sociedades humanas es hacia las conexiones recíprocas. Sin embargo, las redes de Twitter tienden a un equilibrio que no es recíproco y a una estructura jerarquizada. Por lo tanto sigue un modelo en que existen autoridades, las cuales reciben muchos seguidores, pero no corresponden al resto de la misma forma. Por otra parte, observamos que países que tienen alta reciprocidad también tienen usuarios más activos, que participan de comunidades más pequeñas y locales. Otra observación interesante es que las comunidades que son más recíprocas muestran más altos valores en sus niveles de felicidad, así como las comunidades en las cuáles hay más conversación entre usuarios (@). Esto es razonable, ya que niveles más altos de conversación traen más comunicación informal entre los usuarios. Esto es lo opuesto a lo que se observa en países altamente jerarquizados como Estados Unidos, donde se privilegia la comunicación formal y el uso de Twitter como un medio de información más que de conversación. Esto junto con un bajo nivel de clustering y de densidad, nos indica que Estados Unidos posee una comunidad Twitter más globalizada y con menor interacción en

comunidades pequeñas. En la Figura 6 se puede observar la dirección y porcentaje de las conexiones externas de cada país, la mayoría de éstas conexiones son dirigidas hacia Estados Unidos.

ANÁLISIS DE REDES SOCIALES O CÓMO EXTRAER "LAS NOTICIAS DE VERDAD"

Por otra parte las redes sociales no sólo proveen una forma de interactuar informalmente entre usuarios. Sino que algunas, como Twitter, se han convertido en un medio importante de información en tiempo real. Este uso, que se aprovecha día a día, es aún más relevante cuando se

presentan situaciones de crisis como desastres naturales o emergencias de cualquier tipo.

Por esto, con mi equipo de trabajo, hemos profundizado en los temas de:

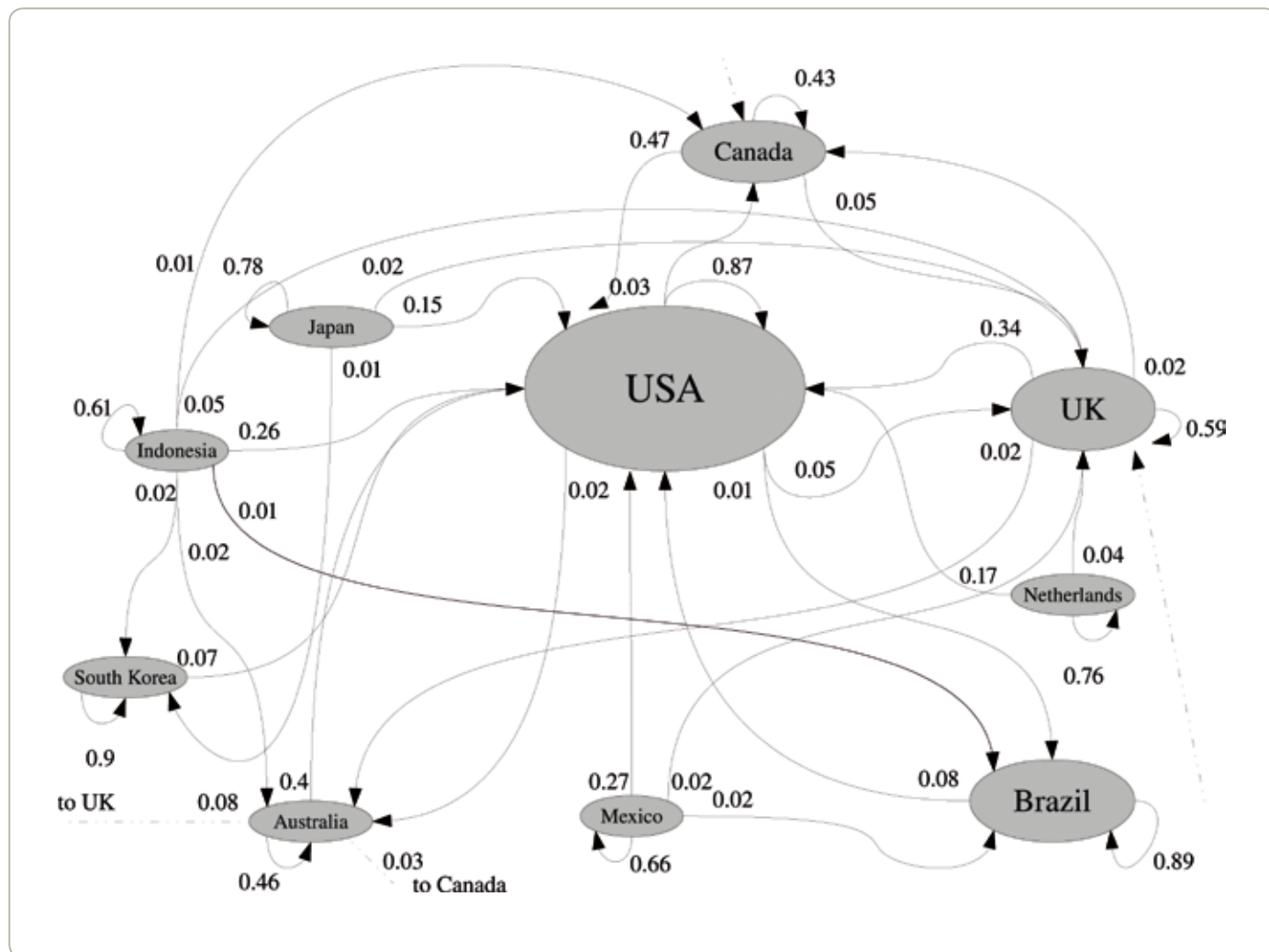
- 1) ¿Cómo separar los tópicos de discusión en Twitter, en *informativos* y en *conversación*? Siendo estos primeros, los de mayor interés para la diseminación de información.
- 2) ¿Cómo saber si un tópico es *creíble* o si se percibe como rumor inverosímil?

Con el afán de explorar estos temas hemos realizado dos trabajos [SOMA] y [WWW2011], junto con un artículo que se encuentra bajo revisión en una revista. En estos trabajos llevamos a cabo un estudio sobre el caso del terremoto

ocurrido en Chile el 27 de febrero de 2010, en el cual pudimos observar una serie de rumores que se propagaron por las redes sociales. Algunos de los cuales resultaron ser falsos y que contribuyeron a generar una sensación de caos en la comunidad. También en este trabajo recopilamos una serie de características comunes que nos hicieron sospechar que la posibilidad de automatizar el proceso de asignar un nivel de credibilidad a una información en Twitter es posible.

Motivados por este descubrimiento, nos dedicamos a trabajar en los dos temas planteados anteriormente. Para esto construimos un clasificador automático de temas *noticiosos*, que permitiese descartar temas de conversaciones sin importancia de

Figura 6



La Web Social ha llegado para quedarse y existen infinitas posibilidades para la información que continuamente está siendo publicada en este medio. Como disciplina, a la Computación le corresponde hacerse cargo de este problema y preparar profesionales capacitados para lidiar con altos volúmenes de información y con conocimientos en Minería de Datos.

los temas que son informativos en Twitter. Además construimos un clasificador que nos permitiese saber si una noticia es percibida como *creíble* o no creíble por los usuarios de la Red. Para poder entrenar y evaluar estos clasificadores creamos un set de datos de gran escala monitoreando temas emergentes de conversación en Twitter por dos meses a nivel mundial. Luego utilizando herramientas de *crowdsourcing*, que permiten acceder a

miles de evaluadores humanos a bajo costo, se etiquetaron manualmente cada uno de los temas, marcando si éste es noticioso o no, y si es creíble o no. A partir de esta investigación, pudimos demostrar que existen características que permiten obtener buenos resultados de clasificación en ambos casos, lo que de ser aplicado puede ser de mucha utilidad en situaciones como la vivida en Chile durante el terremoto.

REFLEXIONES Y LÍNEAS DE INVESTIGACIÓN FUTURAS

La Web Social ha llegado para quedarse y existen infinitas posibilidades para la información que continuamente está siendo publicada en este medio. Como disciplina, a la Computación le corresponde hacerse cargo de este problema y preparar profesionales capacitados para lidiar con altos volúmenes de información y con conocimientos en Minería de Datos. A su vez, nuestros profesionales deben estar dispuestos a trabajar sobre problemas interdisciplinarios y dar soporte científico al análisis de los datos.

El trabajo que actualmente realizo en conjunto con mis alumnos y colaboradores, está enfocado en la recuperación de información multimedia utilizando datos de redes sociales. Todo esto enfocado en la organización y comprensión de los grandes volúmenes de información que se generan día a día en la Web. BITS

REFERENCIAS

[Golder & Macy, Science 2011] Diurnal and seasonal mood vary with work, sleep, and daylength across diverse cultures. Golder, S.A. and Macy, M.W. *Science*, Vol. 333, No. 6051, pp 1878--1881, 2011, American Association for the Advancement of Science.

[Poblete, B. et al. CIKM 2011] Bárbara Poblete, Ruth García, Marcelo Mendoza, and Alejandro Jaimes. 2011. Do all birds tweet the same?: characterizing twitter around the world. In *Proceedings of the 20th ACM international conference on Information and knowledge management (CIKM '11)*, Bettina Berendt, Arjen de Vries, Wenfei Fan, Craig Macdonald, Iadh Ounis, and Ian Ruthven (Eds.). ACM, New York, NY, USA, 1025-1030. DOI=10.1145/2063576.2063724 <http://doi.acm.org/10.1145/2063576.2063724>

[DODDS] P. S. Dodds, K. D. Harris, I. M. Kloumann, C. A. Bliss, and C.

M. Danforth. Temporal patterns of happiness and information in a global social network: Hedonometrics and Twitter. *Computing Research Repository abs/1101.5120v3[physics.soc-ph]*, Feb. 2011.

[Bradley and Lang] M. M. Bradley and P. J. Lang. Affective norms for english words (ANEW): Stimuli, instruction manual, and affective ratings. In *Technical Report C-1, The Center for Research in Psychophysiology*, Gainesville, Florida, 1999.

[Redondo et al.] J. Redondo, I. Fraga, I. Padrn, and M. Comesaa. The spanish adaptation of a new (Affective Norms for English Words). In *Volumne 39*, number 3, pages 600–605. Psychonomic Society Publications, 2007.

[MISLOVE, 2007] A. Mislove, M. Marcon, P. K. Gummadi, P. Druschel, and B. Bhattacharjee. *Measurement and analysis*

of online social networks. In Internet Measurement Conference, pages 29–42, 2007.

[REF7] R. Hanneman and M. Riddle. *Introduction to social network methods*. University of California Riverside, CA, 2005.

[SOMA] Marcelo Mendoza, Bárbara Poblete, and Carlos Castillo. 2010. Twitter under crisis: can we trust what we RT? In *Proceedings of the First Workshop on Social Media Analytics (SOMA '10)*. ACM, New York, NY, USA, 71-79. DOI=10.1145/1964858.1964869 <http://doi.acm.org/10.1145/1964858.1964869>

[WWW2011] C. Castillo, M. Mendoza, and B. Poblete. Information credibility on twitter. In *Proceedings of the 20th international conference on World wide web, WWW '11*, pages 675–684, New York, NY, USA, 2011. ACM.