



**UNIVERSIDAD NACIONAL AUTÓNOMA DE
MÉXICO**

**POSGRADO EN CIENCIA E INGENIERÍA DE LA
COMPUTACIÓN**

**EXTRACCIÓN DE INTERACCIONES
SOCIALES REALES A PARTIR DE
PLATAFORMAS VIRTUALES**

TESIS

QUE PARA OPTAR POR EL GRADO DE

MAESTRO EN CIENCIA E INGENIERÍA DE LA COMPUTACIÓN

PRESENTA:

ALBERTO GARCÍA RODRÍGUEZ

MIEMBROS COMITÉ TUTOR

Dr. CARLOS GERSHENSON GARCÍA
IIMAS-UNAM

Dr. RAFAEL ÁNGEL BARRIO PAREDES
IF-UNAM

Ciudad Universitaria, Ciudad de México Junio 2022



Universidad Nacional
Autónoma de México

Dirección General de Bibliotecas de la UNAM

Biblioteca Central



UNAM – Dirección General de Bibliotecas
Tesis Digitales
Restricciones de uso

DERECHOS RESERVADOS ©
PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL

Todo el material contenido en esta tesis esta protegido por la Ley Federal del Derecho de Autor (LFDA) de los Estados Unidos Mexicanos (México).

El uso de imágenes, fragmentos de videos, y demás material que sea objeto de protección de los derechos de autor, será exclusivamente para fines educativos e informativos y deberá citar la fuente donde la obtuvo mencionando el autor o autores. Cualquier uso distinto como el lucro, reproducción, edición o modificación, será perseguido y sancionado por el respectivo titular de los Derechos de Autor.

Agradecimientos

Quisiera agradecer primeramente a mí familia, por su amor, su apoyo incondicional y por haber fomentado en mí la curiosidad que a día de hoy me acompaña.

A mis amistades, porque a pesar de la distancia o del tiempo que pase sin podernos ver, han estado para mí en los momentos cruciales.

Por supuesto a mis asesores Carlos Gershenson y Rafael Barrio, porque a pesar de su insaciable sed de aprender y sus capacidades técnicas, son unas magníficas personas.

Al IIMAS, al Instituto de Física; a los profesores, de quienes aprendí tanto estos años; a mis compañeros, gente de muy diversas áreas, regiones y hasta países.

Agradezco también a CONACYT, por el apoyo económico que me fue otorgado, me fue de gran ayuda.

Por último, quisiera agradecer a la Asociación de Estudiantes Sudcalifornianos en México, por haberme brindado la oportunidad de alojarme en su casa de estudiantes y haberme dado la oportunidad de conocer estudiantes y personas excepcionales.

Índice general

Resumen	4
Introducción	4
1. Teoría de redes	10
1.1. Inicios	13
1.1.1. 7 puentes de Königsberg	13
1.1.2. Redes aleatorias	14
1.1.3. Mundo pequeño	15
1.1.4. Modelo de Watts-Strogatz	15
1.1.5. Modularidad y comunidades	16
1.1.6. Modelo de enlace preferencial	17
1.2. Bases	18
1.2.1. Matriz de adyacencia	19
1.2.2. Redes pesadas	19
1.2.3. Medidas	20
1.2.4. Caminos y más	23
1.2.5. Hipergráficas	23
1.2.6. Gráficas potencia	24
1.2.7. Redes bipartitas	24
1.2.8. Redes multirelacionales	25
1.2.9. Redes multicapa	26
1.3. Antecedentes	27
1.3.1. Detección de comunidades	27
1.3.2. Redes coordinadas	28
1.3.3. Homofilia	29
1.3.4. Cámaras de eco	29
2. Metodología	31
2.1. Recopilación de datos	31

2.2. Primer acercamiento	33
2.3. Red de co–eventos	36
2.4. Análisis de susceptibilidad	38
3. Resultados	40
3.1. Red de co–eventos	40
3.2. Susceptibilidad	44
3.3. Agrupamientos a diferentes niveles de la red de co–eventos .	46
4. Conclusiones	51
A. Twitter y características	53
B. Método de Louvain	55
C. Identificación de usuarios	56
C.1. Bots y usuarios relacionados por afinidad política	56
C.2. Usuarios con relación laboral	59

Resumen

El análisis de redes sociales de las plataformas digitales Facebook ¹, Twitter ², Reddit ³, sólo por mencionar algunas, ha tenido un gran desarrollo en los últimos años debido al interés en identificar bots y usuarios coordinados han aplicado métodos de aprendizaje automático, procesamiento del lenguaje natural, minería de textos, aprendizaje profundo y demás técnicas sofisticadas, pues en recientes estudios se ha encontrado que estos interfieren sustancialmente en las interacciones humanas, pueden modificar nuestra percepción de los eventos ocurridos en estas plataformas y hasta fuera de ellas. La mayoría de estas técnicas dependen en gran medida del conjunto de datos previamente etiquetados (para los que hacen uso de aprendizaje supervisado), del idioma y el modelo del lenguaje en el cual fueron basados (para aquellos que hacen uso del texto) el problema de ser tan sensibles a estas situaciones, es que entonces estas técnicas serían limitadas a un uso geográfico o a la cantidad y calidad de datos etiquetados. Por ello, en el presente trabajo se proponen una serie de técnicas para atacar este problema, pero desde otra perspectiva detectando redes sociales reales, así como la detección de bots o usuarios cuyo propósito es difundir masivamente información, en plataformas virtuales tipo Twitter con la ventaja de que las técnicas empleadas son independientes al idioma y sin la necesidad del uso de aprendizaje supervisado. Al final del trabajo mostramos como resultados algunos bots detectados, redes sociales donde la agrupación y conexión fue dada por una afinidad política y sub-redes donde los usuarios de estas mantienen un vínculo por relacion laboral así como una red que puede ser segmentada de cierta forma en la que se aprecian distintos niveles de vinculación.

¹<https://www.facebook.com/>

²<https://twitter.com/>

³<https://www.reddit.com/>

Introducción

Este último siglo la producción de artículos científicos con la temática de *redes* ha tenido un crecimiento exponencial parte de ese crecimiento, se lo debemos al desarrollo y masificación de las computadoras personales pero también a la versatilidad que tienen las redes para describir una vasta variedad de sistemas complejos, como: la interconexión de las páginas web, las redes telefónicas, la red eléctrica, la red de transporte, las redes de distribución, la estructura de las citas en artículos científicos, las redes de amistad, los procesos metabólicos, la regulación genética, las estructuras cerebrales, las interacciones entre especies, las epidemias, entre muchos otros.

Aunque es claro que cada uno de estos fenómenos son particularmente distintos entre sí, al transformarlos a *nodos* unidos por *líneas*, donde los nodos representan los elementos del sistema que se desean modelar y las líneas las interacciones entre ellos, dicho de otra forma, al transformar estos fenómenos a redes vemos que comienzan a surgir varias propiedades en común, muchas de ellas debidas a la topología de la red y otras debidas a su evolución en el tiempo.

A pesar de que muchos de los sistemas que son modelados con redes comparten propiedades en común, existen algunos enfoques tales como el estudio de redes biológicas, sociales o de transporte, donde existirán muchas más similitudes cuando las redes pertenecen al mismo tipo. Para este trabajo haremos uso del marco teórico para el análisis de redes sociales (SNA por sus siglas en inglés). En lo general éste marco teórico se enfoca en el estudio del comportamiento individual: los patrones de las relaciones, dados por la estructura de la red y las interacciones entre ambos. Debido a que las redes sociales son causa y resultado de las acciones individuales pero a su vez, la estructura de la red puede condicionar o limitar las acciones de los individuos. De tal forma que aquí se puede observar causalidad hacia arriba y causalidad hacia abajo.

Seleccionamos la plataforma Twitter debido a que esta plataforma cuenta con su propia API (Interfaz de Programación de Aplicaciones) para acceder a sus datos, lo cual facilita la recopilación de datos, búsquedas de datos almacenados y extracción de mensajes liberados en tiempo real. Además esta red se caracteriza por contener usuarios activos en discusiones altamente polarizadas, a diferencia de otras redes donde el tema principal es generar contactos y compartir con ellos cierta información, aquí se centran más en entrar en discusiones y peleas muy agresivas. De tal modo que supusimos sería más sencillo recopilar datos de discusiones.

Nuestro principal objetivo con este trabajo será proponer la construcción de una o más redes a partir de los datos recabados de Twitter con la intención de obtener de ellas interacciones sociales reales, es decir: extraer un conjunto de usuarios que estén vinculados por razones políticas, laborales, familiares, intelectuales y no sólo vinculados de forma circunstancial ⁴.

Twitter es una plataforma virtual de tipo ‘microbloggin’, la cual a grandes rasgos consiste en una interfaz en la que uno puede escribir textos cortos (280 caracteres máximo), publicar imágenes o vídeos. El contenido que cada usuario genere, podrá ser visto por cualquier otro usuario (a no ser que el usuario propietario coloque alguna restricción) sin embargo, sólo se mostrará de forma automática en la interfaz de los usuarios que sean seguidores del usuario en cuestión. La relación entre cada usuario es dirigida, la comunicación por lo general no es bidireccional. Cada mensaje puede ser relanzado sin alteración (retweeteado) por otros usuarios, o puede ser referenciado agregándole un extra (cita o ‘quote’), se contabiliza y expone cada vez que el mensaje es relanzado o citado, lo cual agrega un valor extra al mensaje. Los mensajes más citados o relanzados serán expuestos en una zona especial de la interfaz de Twitter, otorgando una mayor visibilidad.

Además, el texto seguido de un ‘#’ (hashtag) lo contará internamente la plataforma y los más populares también serán expuestos en la zona donde se colocan las tendencias.

Todo esto supone una serie de problemas, como por ejemplo: las relaciones circunstanciales (relaciones de muy corto plazo) son muy propicias ante eventos ‘espontáneos’ ya que la plataforma impulsa los mensajes y hashtags

⁴Decimos que un vínculo es circunstancial cuando, la relación entre dos o más usuarios se forma únicamente por circunstancias exclusivas de ese momento particular y no necesariamente de ideas o circunstancias que proliferen en un mediano o largo plazo.

más replicados. Como una bola de nieve, los mensajes más populares al ser expuestos a todos los usuarios, tendrán más oportunidad de crecer en popularidad y enconces perpetuar el ciclo. Otro problema es que las relaciones al ser direccionales, no implican que ambos usuarios tengan conocimiento del motivo de la relación. Por ejemplo, el presidente de México emite un tweet, en cuestión de minutos ya miles de usuarios habrán retweeteado el mensaje, pues tiene miles de usuarios que lo siguen sin embargo, el presidente de México sólo sabrá de aquellos tweets que hayan emitido los usuarios que él sigue.

Al existir las ‘quotes’, se le agrega incertidumbre al sentido del mensaje pues sin procesamiento del lenguaje natural no se puede saber si el mensaje lleva un sentido de apoyo o ataque.

Dadas las problemáticas anteriores, encontrar relaciones fuertes, relaciones de mediano–largo plazo y no sólo relaciones circunstanciales en Twitter, no es un problema trivial.

En principio: las “social media” plataformas digitales que facilitan a los usuarios compartir y crear contenido con otros usuarios fueron diseñadas para permitir a sus usuarios expresar sus opiniones sobre diversos temas, hacer nuevos amigos y/o contactos. Sin embargo, no pasó mucho tiempo después de su creación para que compañías y/o gobiernos se dieran cuenta del enorme potencial que tienen estas plataformas para poder influir⁵ en los usuarios que comenzaron a invertir dinero y recursos humanos para lograr sus fines políticos, económicos y demás. La gama de plataformas virtuales diseñadas para permitir estas interacciones es muy amplia y su popularidad puede variar en función de la edad de sus usuarios y hasta la zona geográfica a la que pertenezca, sólo por mencionar algunos ejemplos Facebook, Instagram, WhatsApp, todas las plataformas derivadas de Meta⁶, Twitter⁷, Discord⁸, Wechat⁹, V Kontakte¹⁰, Telegram¹¹, Sina Weibo¹², Viber¹³, QQ¹⁴, cuya función principal es comunicar mediante mensajes de texto, existen otras en las cuales su principal función es compartir fotografías o vídeos como, Snap-

⁵<https://www.theguardian.com/news/series/cambridge-analytica-files>

⁶<https://developers.facebook.com/>

⁷<https://developer.twitter.com/en>

⁸<https://discord.com/>

⁹<https://www.wechat.com/>

¹⁰<https://vk.com/>

¹¹<https://web.telegram.org/>

¹²<https://weibo.com/>

¹³<https://www.viber.com/>

¹⁴<https://international.qq.com/>

chat¹⁵, Youtube¹⁶, Tiktok¹⁷, Kuaishou¹⁸, además existen las plataformas orientadas al arte como, DeviantArt¹⁹, Behance²⁰, Artstation²¹, VSCO²², para citas, Hinge²³, Tinder²⁴, Feeld²⁵, Bumble²⁶, adicionalmente podríamos agregar a LinkedIn²⁷ que se usa más para cuestiones laborales y Reddit²⁸ que es de tipo blog. La mayoría de ellas cuentan con la posibilidad de generar interacciones no humanas, además de las características que ofrece cada una de las plataformas para poder realizar ciertas acciones de forma automatizada, es sabido que también se pueden transgredir dichas políticas de uso, mediante el uso de algunas técnicas computacionales, esto supone un gran problema, pues abre paso a una amplia gama de opciones para interferir en las relaciones humanas.

Hoy en día, es sabido que las redes sociales tienen impacto en las decisiones para realizar políticas (Keller et al., 2020; Levy and Razin, 2019), cuestiones de imagen pública, estrategias de mercadotecnia, etc., es importante tomar en cuenta esto para el entendimiento de las dinámicas y mecanismos subyacentes a las plataformas de “social media”.

Se sabe que es muy complicado determinar si un usuario de Twitter es humano (Mirko Lai and Rosso, 2017; Dutta et al., 2020). Algunas estrategias de inteligencia artificial se han implementado para tratar de realizar esta tarea (Wang, 2010; Yang et al., 2022), el problema es que la eficiencia obtenida ha sido baja y limitada. Usualmente estas tareas requieren la intervención humana para el etiquetado de los datos y el etiquetado es fuertemente dependiente del contexto, requiere manejar sutilezas como sarcasmo, humor negro y demás (Mirko Lai and Rosso, 2017). Se alcanza a observar que los usuarios de Twitter tienden a encapsularse en burbujas ideológicas²⁹ (Pa-

¹⁵<https://www.snapchat.com/>

¹⁶<https://www.youtube.com/>

¹⁷<https://www.tiktok.com/>

¹⁸<https://www.kuaishou.com/>

¹⁹<https://www.deviantart.com/>

²⁰<https://www.behance.net/>

²¹<https://www.artstation.com/>

²²<https://vsco.co/>

²³<https://www.hinge.co/>

²⁴<https://tinder.com/>

²⁵<https://feeld.co/>

²⁶<https://bumble.com/>

²⁷<https://www.linkedin.com/>

²⁸<https://www.reddit.com/>

²⁹Con el desarrollo de los algoritmos de aprendizaje de máquina, se a vuelto muy sencillo que las empresas vinculadas a las tecnologías de la información, reconozcan nuestros gustos y con ello crean categorías de usuarios con gustos en común. Esto aplicado en las plata-

riser, 2011) muy bien definidas de tal forma que la mayoría de usuarios en estas burbujas están de acuerdo con el tema en cuestión y prácticamente no hay debate interno. En este trabajo implementaremos algunas técnicas de análisis de redes para identificar distintas agrupaciones políticas, características interesantes de estas redes y relaciones sociales de mediano-largo plazo.

Por lo tanto, una pregunta interesante sería saber si podemos encontrar interacciones de humanos reales a partir de la topología de la red, sin el uso de métodos que se basen en el lenguaje, humor, etc., u otros contextos difíciles de evaluar. La cual intentaremos responder a lo largo de este trabajo.

En el primer capítulo, se presenta una breve introducción a la teoría de redes, empezando por sus antecedentes históricos, seguido de un compendio de bases matemáticas requerido para el seguimiento del trabajo y concluyendo con los trabajos que antecedentes relacionados al estudio de redes sociales, redes politizadas y coordinación en redes.

El segundo capítulo, aborda la metodología que se siguió, desde la recopilación de datos hasta la construcción de las redes de co-eventos. En el capítulo tercero, se muestran los resultados obtenidos de cada uno de los análisis realizados, en él se aborda tanto el resultado estadístico como la información que esto nos proporciona sobre las relaciones sociales por medios virtuales. Por último, el capítulo cuarto, aborda las conclusiones de este trabajo, se realiza una recopilación de las metas alcanzadas, las limitaciones de este trabajo y posibles trabajos a futuro.

formas virtuales de redes sociales, tiende a generar filtros de información a modo, de tal forma que terminamos comunicandonos sólo con usuarios afines a nuestros pensamientos. A esto se le conoce también como burbujas ideológicas.

Capítulo 1

Teoría de redes

En este capítulo se presentan de forma introductoria los orígenes de los sistemas complejos y la teoría de redes, así como el fundamento matemático de teoría de redes con la finalidad de que al lector se le pueda facilitar la comprensión de este trabajo.

Al paso del tiempo, la forma en que se ha hecho ciencia ha ido evolucionando, cada vez ha tomado mayor formalidad y criterios más objetivos, en busca de disminuir los sesgos con la finalidad de obtener mejores interpretaciones de los descubrimientos. La mayor parte de los grandes avances en la ciencia moderna los debemos al “enfoque reduccionista”, el cual sostiene básicamente que “el todo no es más que la suma de las partes”, lo cual quiere decir que podríamos explicar un fenómeno complejo, al entender las características que poseen las partes fundamentales que lo componen. La mecánica clásica particularmente fue desarrollada haciendo uso exhaustivo del reduccionismo, empleando como su marco teórico y sólo por mencionar algunos de sus frutos, están las tres leyes del movimiento de Newton (1726), además de esto, también se encontraba el **demonio de Laplace** el cual es un “determinismo científico” en el que se argumenta que si se conociera la ubicación precisa y el momento de cada átomo en el universo, se podría determinar cualquier valor pasado o futuro, a partir de las leyes de la mecánica clásica. Entonces, se generó una visión de cómo debían analizarse los fenómenos y el poder de las matemáticas para describirlos, parecía ser que no había nada que no pudiera ser descrito matemáticamente y con esas descripciones se podría describir el pasado y futuro de cualquier cosa. Sin embargo, con el surgimiento de las computadoras y nuevos descubrimientos, fue claro que estos métodos tendrían sus limitaciones y que éstas serían bastante relevantes.

El reduccionismo científico no permite la causalidad hacia abajo mientras que en la biología y los sistemas sociales sí es permitido, los fenómenos emergentes contradicen la esencia del reduccionismo. La emergencia se plantea como propiedades de los sistemas complejos que no pueden ser predichas conociendo únicamente las características individuales de los elementos que los forman. El descubrimiento del caos en sistemas dinámicos (Lorenz, 1963) acabó con el determinismo de Laplace.

Entonces ¿este fue el fin del reduccionismo?, ¿ya no se hace ciencia con un enfoque reduccionista?, ¿surgió un nuevo enfoque para realizar ciencia?, ¿qué sigue?. El método reduccionista para realizar investigación científica se sigue y seguirá empleando, el reduccionismo ha probado ser un enfoque útil para atacar ciertos problemas, por ello es que continuará en uso, sin embargo, al identificar sus limitaciones es que se plantearon nuevas formas de estudiar los sistemas. El enfoque con mayor uso en la actualidad es el de las llamadas “ciencias de la complejidad”, este enfoque sirve particularmente para el análisis de los llamados “sistemas complejos”, los cuales son sistemas donde sus elementos mantienen interacciones que hacen muy difícil separarlos para analizarlos individualmente, en los sistemas complejos se plantea que “el todo es más grande que la suma de sus partes”, contradiciendo al reduccionismo.

Aunque a día de hoy no se tiene una definición muy clara de ¿qué es un sistema complejo, o qué son las ciencias de la complejidad?, podemos decir que las ciencias de la complejidad ¹ estudian cómo los sistemas compuestos por elementos que a bajas escalas interaccionan entre sí, pueden de forma espontánea auto-organizarse exhibiendo estructuras y comportamientos no triviales a escalas más grandes, cuyos comportamientos emergentes no pueden ser predichos ignorando estas dinámicas y centrándose exclusivamente en el estudio exhaustivo sus elementos de forma aislada. Para abordar estos sistemas, se requiere de nuevos marcos matemáticos y métodos científicos.

Por ello se contemplan por lo menos 7 marcos teóricos para el estudio de los sistemas complejos, mismos que son mostrados en la Fig. 1.1 y como vemos las redes constituyen uno de los 7 bloques.

Y la razón es clara, la teoría de redes ha mostrado ser una herramienta muy útil a la hora de estudiar las interacciones entre los elementos que forman un sistema complejo, los patrones que emergen, las topologías que se forman debido a las interacciones y cómo estas topologías propician ciertas

¹<https://complexityexplained.github.io/>

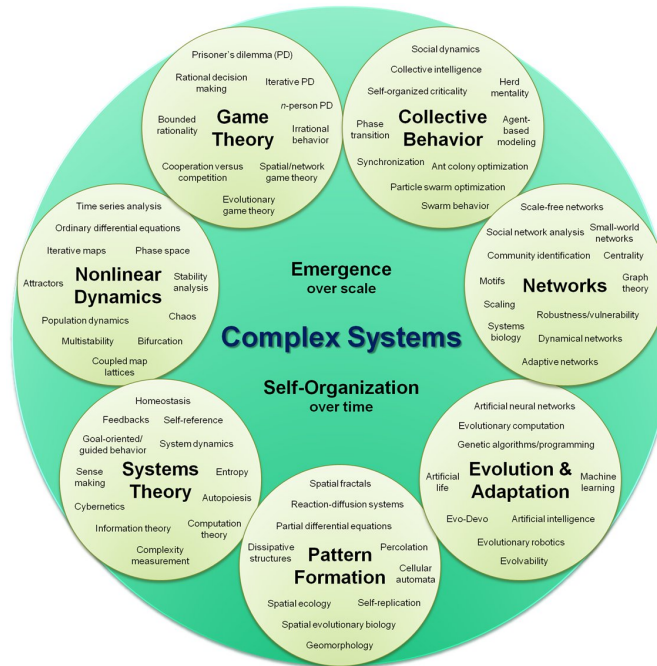


Figura 1.1: Mapa organizacional de los sistemas complejos dividido en 7 subgrupos, elaborado por Hiroki Sayama. Obtenido de Wikipedia.

dinámicas. Pero entonces, ¿qué son las redes? la red, en su forma más elemental, es una colección de elementos los cuales están conectados en pares por una línea, a los elementos se les conoce como *nodos o vértices* y a las líneas como *ligas o enlaces* (Newman, 2010). Podría parecer demasiado ambigua esta definición, pero esta flexibilidad es la que nos permite trabajar las redes en prácticamente cualquier área de estudios, como en sociología, biología, física, economía, entre otras, pues en cada una de estas áreas las relaciones entre los elementos, pueden ser diferentes, pero el hecho de que exista una relación es suficiente para generar el enlace.

La teoría de redes ha sido un área de estudio que ha crecido significativamente los últimos años, debido en gran medida al desarrollo de las computadoras, internet, las bases de datos y la tecnología en general, su impacto en la ciencia y tecnología ha sido abrumador, la teoría de redes ha resultado tener un sin fin de aplicaciones a problemas prácticos y eso fue muy bien visto por los gobiernos y por las empresas.

Por ello a continuación esbozaremos brevemente un poco de su historia y su avance hasta nuestros tiempos.

1.1. Inicios

1.1.1. 7 puentes de Königsberg

La teoría de redes se construye a partir de la teoría de gráficas y esta nace por los años de 1730 en Königsberg, capital de Prusia oriental y una ciudad comercial muy activa. En esos tiempos el comercio se realizaba mayormente por medio de los barcos, controlando rutas marítimas y demás. Dado que también era el caso de esta ciudad y se encontraban en una buena época comercial, decidieron construir 7 puentes sobre el río Pregel para interconectar la ciudad. 5 de ellos intercomunicaban con la isla de Kneiphof y los otros 2 cruzaban las dos ramificaciones del río. Esto lo podemos ver representado en la Fig. 1.2.

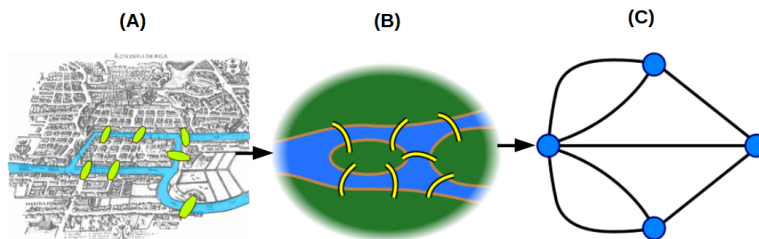


Figura 1.2: 7 Puentes de Königsberg, donde (A) representa el mapa de los 7 puentes con mayor detalle, (B) el mapa simplificado de los 7 puentes de Königsberg y (C) es una representación en forma de gráfica. Obtenido de Wikipedia.

De estas conexiones (puentes) surgió la siguiente cuestión ¿Uno puede recorrer los 7 puentes sin tener que cruzar por uno de ellos más de una vez?, esta pregunta permaneció sin respuesta hasta que el matemático Leonard Euler ofreció una demostración rigurosa para responder el problema. Euler generó una gráfica como la mostrada en la Fig. 1.2 (panel C) y de ella comenzó a realizar sus análisis y deducciones. Euler demostró que no es posible y de hecho concluyó que no es posible en ninguna gráfica que tenga más de dos nodos con un número impar de enlaces.

1.1.2. Redes aleatorias

El estudio de las redes aleatorias logró alcanzar gran visibilidad gracias al trabajo de Pál Erdős y Alfréd Rényi, ellos juntaron teoría de probabilidad y combinatoria con teoría de gráficas en una serie de 8 artículos publicados entre 1959 y 1968 estableciendo la teoría de redes aleatorias como una nueva rama de las matemáticas.

Una red aleatoria consiste en N nodos donde cada par de nodos es conectado con la misma probabilidad p .

En sus dos principales trabajos (Erdős and Rényi, 1959, 1960) plantearon un modelo de redes aleatorias y su evolución, curiosamente Gilbert (1959) también desarrolló y publicó un modelo de redes aleatorias el mismo año, sin embargo su impacto no fue tan significativo. Existen sutiles diferencias entre ambos modelos, a continuación mostraremos una breve definición de cada modelo:

- Sea $\mathbf{G}(\mathbf{N}, \mathbf{L})$ una gráfica con N nodos conectados con L enlaces colocados de forma aleatoria, modelo de Erdős and Rényi (1959).
- Sea $\mathbf{G}(\mathbf{N}, \mathbf{p})$ una gráfica donde cada par de N nodos es enlazado con una probabilidad p , modelo de Gilbert (1959).

Como se puede apreciar, el modelo $G(N, p)$ fija la probabilidad p de que dos nodos sean conectados, en cambio $G(N, L)$ fija el número total de enlaces L . Cuando queremos calcular algunas características de la red, como la distribución del grado, su dispersión, entre otras, es más sencillo deducirlas por medio del modelo de Gilbert, la distribución del grado en una red aleatoria, está definida por la distribución binomial

$$p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k}, \quad (1.1)$$

donde k representa el grado en cuestión, esta distribución puede ser aproximada por una distribución de Poisson siempre y cuando $k \ll N$, definida por

$$p_k = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}. \quad (1.2)$$

El estudio de las redes aleatorias es de gran utilidad a la hora de realizar análisis sobre redes *reales*, pues haciendo uso de las propiedades de las redes aleatorias, tales como la distribución de grado, grado promedio, entre

otras, son muy útiles para contrastarlas con las de nuestra red real, de esta forma podemos saber si las propiedades identificadas en nuestra red, son particulares de la red o siguen propiedades de una red aleatoria.

1.1.3. Mundo pequeño

Probablemente, uno de los primeros trabajos que llamó la atención de muchos científicos que estaban inursionando en la teoría de redes fue el experimento de Travers and Milgram (1969) sobre el “problema de mundo pequeño”, el cual a grandes rasgos consistió en seleccionar arbitrariamente un conjunto de individuos a los cuales se les pidió contactaran a una persona apoyándose de su red de contactos de tal forma que aún si no lo conocían directamente, contactaran a la persona que creyeran tuviera algún conocido que sí la conociera directamente. Como resultado de este experimento, se obtuvo que un subconjunto del conjunto inicial sí logró hacer contacto con la persona en cuestión y en promedio tomó “5.5 pasos”, lo cual fue bastante sorprendente, pues eso traería consigo muchas preguntas, sobre la forma en la que estamos relacionados, nuestras redes de contacto y demás. De los resultados de este experimento emergió la frase “seis grados de separación”, la cual hace alusión a que no importa a quién deseamos contactar, entre cualquier persona y nosotros sólo estamos a 6 personas de separación. De hecho, en años recientes se han realizado pruebas sobre algunas plataformas de redes sociales tales como Facebook en Backstrom et al. (2012) se observó que la distancia promedio era de 4,74 la cual corresponde a 3.74 intermediarios, sin embargo en 2016 Facebook realizó nuevamente experimentos ² y encontró que el promedio de intermediarios se redujo a 3,57, lo cual indica que hubo una reducción considerable respecto al anterior (Backstrom et al., 2012).

1.1.4. Modelo de Watts-Strogatz

Pasaron más de 20 años para que se formularan modelos que pudieran explicar las características que producían el fenómeno de *mundo pequeño* con el modelo de Watts and Strogatz (1998) se entendió que las redes de “mundo pequeño” se encontraban en un punto medio entre las redes regulares y las redes aleatorias. A grandes rasgos el modelo inicia con un anillo de nodos, donde cada nodo es conectado a sus vecinos inmediatos, con esa

²<https://research.facebook.com/blog/2016/2/three-and-a-half-degrees-of-separation/>

configuración el *coeficiente de agrupamiento*³ es $C = 3/4$. A partir de aquí se inicia con una reconexión aleatoria para cada enlace, dada por una probabilidad p . Para una p pequeña la red mantiene un alto coeficiente de clustering promedio, pero las distancias promedio a nodos decrecen drásticamente, lo cual induce el efecto de *mundo pequeño*. Sin embargo para valores cercanos a 1, básicamente obtenemos una red aleatoria. Lo podemos ver de forma gráfica en las Figs. 1.3 y 1.4.

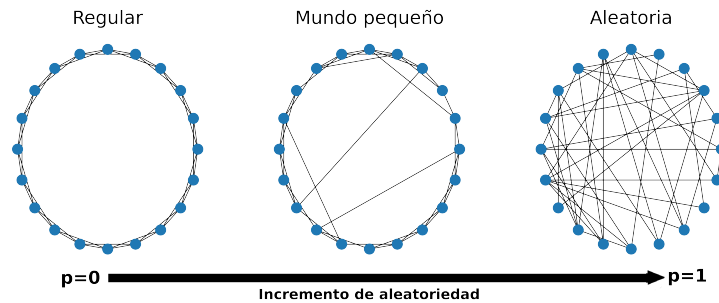


Figura 1.3: Fases de la red al realizar una reconexión aleatoria sobre los enlaces de una red regular, variando la aleatoriedad de 0 a 1 de forma gradual.

1.1.5. Modularidad y comunidades

Años más tarde del experimento de Milgram surgió otro artículo bastante interesante en el que se analizaba una red de un club de karate (Zachary, 1977) en el que por medio del uso del algoritmo de “Ford–Fulkerson” logró detectar comunidades dentro de la red, es decir, para este problema logró detectar que existían subgrupos mayormente cohesionados que con el resto de la red. Esto abrió una brecha para comenzar a crear algoritmos que fueran capaces de identificar grupos internos en las redes. Más tarde Girvan and Newman (2002) lograron diseñar un algoritmo capaz de detectar comunidades dentro de una red. Pocos años más tarde Newman (2006) planteó una nueva medida, denominada *modularidad* la cual, serviría como base para identificar qué comunidades están mejor cohesionadas en función de sus enlaces.

³El coeficiente de agrupamiento mide la fracción de caminos de tamaño dos que están cerrados, es decir, caminos de tamaño dos que al incorporarle un enlace más ya existente, da lugar a un ciclo de tamaño tres. Para mayor detalle, consultar el apartado de medidas 1.2.3

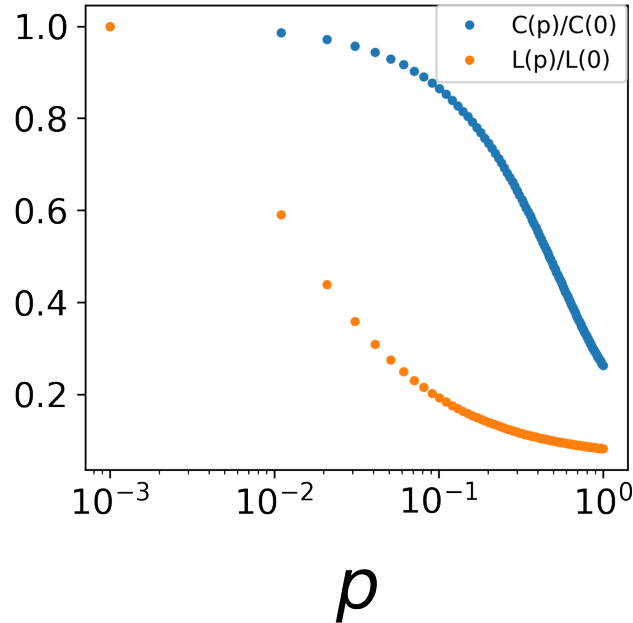


Figura 1.4: Variaciones del coeficiente de agrupamiento ($C(p)$) y la longitud de camino mas corto promedio ($L(p)$) en función de la aleatoriedad p , ambos normalizados por $C(0)$ y $L(0)$ respectivamente.

1.1.6. Modelo de enlace preferencial

Sin duda el modelo de Barabasi and Albert (1999) ha jugado un papel fundamental en estos tiempos, porque nos deja ver que muchas de las actividades sociales están llenas de comportamientos que siguen distribuciones de leyes de potencias, la popularidad de los usuarios en las plataformas virtuales, la concentración de riqueza, etc. y no sólo eso, a medida que pasa el tiempo concentrarán más recursos, generando mayor desigualdad. Su modelo plantea una red que va evolucionando con el tiempo, inicia con m_0 nodos con enlaces entre ellos escogidos arbitrariamente, pero con la condición de que cada nodo tenga al menos un enlace. La red se desarrolla siguiendo los siguientes pasos:

1. En cada paso de tiempo agregamos un nuevo nodo con m ($\leq m_0$) enlaces que conectan al nuevo nodo a m nodos existentes en la red.

2. La probabilidad $p(k)$ de que uno de los enlaces del nuevo nodo se conecte al nodo i depende del grado k_i como se define a continuación:

$$p(k_i) = \frac{k_i}{\sum_j k_j} \quad (1.3)$$

Este modelo reúne las condiciones mínimas para generar redes cuya distribución de grado sigue una ley de potencias, conocidas también como redes libres de escala. El término *libre de escala* hace referencia a la divergencia que se presenta al momento de calcular cualquier n momento de la distribución de grado en una red que sigue una distribución de ley de potencias. Como recordamos, la desviación estándar nos da información sobre que tan dispersos están los datos respecto a la media, entonces el grado de un nodo elegido al azar lo podemos definir de la siguiente forma

$$k = \langle k \rangle \pm \sigma_k, \quad (1.4)$$

donde σ_k es la desviación estandar del grado k , para el caso de de una red que sigue una distribución de Poisson (una red aleatoria) sería $k = \langle k \rangle \pm \langle k \rangle^{1/2}$ aquí podemos tomar a $\langle k \rangle$ como escala, sin embargo para una red que sigue una distribución de ley de potencias sería $k = \langle k \rangle \pm \infty$, sería un sin sentido tomar $\langle k \rangle$ como escala.

1.2. Bases

Una red es una abstracción de un sistema, un modelo en el cual se enlistan los elementos de un sistema, normalmente a esos elementos se les llama *nodos* o *vértices* y estos están vinculados por alguna característica, a dicho vínculo se le conoce como *enlace* o *arista* (Newman, 2010; Barabási and Pósfai, 2016). Más formalmente, una red consiste en un conjunto de vértices \mathcal{N} y un conjunto de enlaces, $\mathcal{L} \subset \mathcal{N} \times \mathcal{N}$. Los enlaces pueden ser dirigidos o no dirigidos. Un enlace no dirigido es un conjunto de dos nodos distintos, $\{i, j\}$. Para el caso en que el enlace sea dirigido, entonces el enlace sería descrito como un par ordenado (i, j) cuya dirección va del nodo i al nodo j . En este último caso al nodo i se le conoce como fuente y al j objetivo.

Al número de nodos lo definiremos como $N = |\mathcal{N}|$, el cual representa el número de elementos del sistema y también denota el tamaño de la red. Al número de enlaces lo denotaremos por $L = |\mathcal{L}|$, este representma el número de interacciones entre los nodos.

1.2.1. Matriz de adyacencia

Como vimos anteriormente, un enlace dirigido puede ser representado por un par ordenado de nodos, si quisieramos representar un enlace no dirigido como una lista de pares ordenados, sólo bastaría con duplicar los pares de nodos pero con orden inverso. Una forma muy sencilla de mostrar el registro de enlaces de una red, sería mostrar la lista de enlaces que la constituyen, ejemplo: $\{(i_1, j_1), (i_1, j_2), \dots, (i_n, j_n)\}$. Sin embargo existe otra forma de representar las redes, que nos proporciona facilidad para realizar operaciones matemáticas sobre ella. A esta representación se le conoce como *matriz de adyacencia* (Newman, 2010; Barabási and Pósfai, 2016) y consiste en representar de N nodos en una matriz con N filas y N columnas, donde usaremos a i como un subíndice para denotar la fila y a j para denotar la columna. Los elementos de la matriz serán:

$A_{ij} = 1$ en caso de existir un enlace que apunte de i a j .

$A_{ij} = 0$ en caso contrario.

Es claro que para el caso en que la red sea *no dirigida* entonces $A_{ij} = A_{ji}$, generando una matriz simétrica. Hasta aquí las redes anteriormente descritas también se les conoce como *redes simples*, existen otras denominadas *multigráficas* (Newman, 2010; Barabási and Pósfai, 2016) en las cuales se permiten enlaces en paralelo y autoenlaces. Para estas redes, su matriz de adyacencia se expresa como en el caso de las redes pesadas descritas a continuación.

1.2.2. Redes pesadas

Hasta ahora hemos definido cómo expresar una red en forma matricial, colocamos $A_{ij} = 1$ si existe un enlace y $A_{ij} = 0$ en caso de que no exista dicha conexión, pero esto asume que todas las conexiones valen lo mismo, ¿en el mundo real todas las interacciones valen lo mismo, todos los vínculos valen lo mismo?. La respuesta a esto es *no*, y como hay contextos donde vale la pena los distintos grados de relación, es que existen las *redes pesadas* (Newman, 2010; Barabási and Pósfai, 2016). Estas siguen la misma idea que las redes sin pesos, lo único que modificaremos es el valor del peso introducido en la matriz de adyacencia $A_{ij} = w_{ij}$ donde w_{ij} es el peso que tendrá dicho enlace.

En la Fig. 1.5 podemos observar una *red no dirigida* (A) construida a partir de la matriz 1.5, también de la misma figura podemos ver la *red dirigida* (B) formada por la matriz 1.6 y por último la *red pesada dirigida*

(C) la cual fue construida a partir de la matriz 1.7.

$$A_{ij} = \begin{bmatrix} 0. & 1. & 1. & 1. \\ 1. & 0. & 1. & 0. \\ 1. & 1. & 0. & 0. \\ 1. & 0. & 0. & 0. \end{bmatrix} \quad (1.5)$$

$$A_{ij} = \begin{bmatrix} 0. & 1. & 0. & 0. \\ 0. & 0. & 1. & 0. \\ 1. & 0. & 0. & 0. \\ 1. & 0. & 0. & 0. \end{bmatrix} \quad (1.6)$$

$$A_{ij} = \begin{bmatrix} 0. & 1. & 0. & 0. \\ 0. & 0. & 0,5 & 0. \\ 2. & 0. & 0. & 0. \\ 1. & 0. & 0. & 0. \end{bmatrix} \quad (1.7)$$

1.2.3. Medidas

En una red no dirigida podemos expresar la cantidad de enlaces que tiene con otros nodos como el *grado* (Newman, 2010; Barabási and Pósfai, 2016) y para denotarlo en función de los i nodos que existen en la red, podemos expresarlo como k_i . El número de enlaces (L) para este tipo de red puede ser expresado como:

$$L = \frac{1}{2} \sum_{i=1}^N k_i \quad (1.8)$$

El grado promedio puede ser definido por

$$\langle k \rangle \equiv \frac{1}{N} \sum_{i=1}^N k_i = \frac{2L}{N} \quad (1.9)$$

Para redes dirigidas el grado depende de la dirección y decimos que a los enlaces que inciden el enlace en el nodo (cuando la flecha apunta al nodo) cuentan como grado de entrada o k_i^{in} y en caso contrario sería k_i^{out} o grado de salida. De tal forma que el grado en el nodo i está dado por

$$k_i = k_i^{in} + k_i^{out} \quad (1.10)$$

Por lo tanto aquí el número de enlaces estaría representado por

$$L = \sum_{i=1}^N k_i^{in} = \sum_{i=1}^N k_i^{out} \quad (1.11)$$

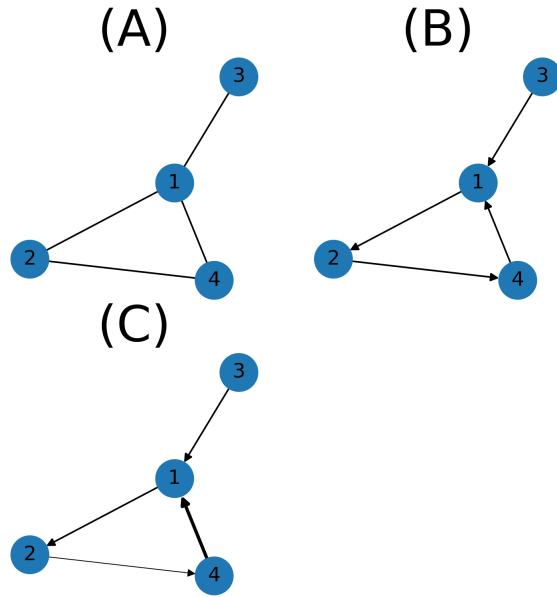


Figura 1.5: (A) Red no dirigida, donde cada enlace es trazado por una línea que une a un par de nodos.
 (B) Red dirigida, donde se puede apreciar que cada enlace direcciona en sentido de la flecha.
 (C) Red pesada y dirigida, aquí podemos apreciar que los enlaces varían su grosor en función del peso.

Para este caso el grado promedio sería

$$\langle k^{in} \rangle = \frac{1}{N} \sum_{i=1}^N k_i^{in} = \langle k^{out} \rangle = \frac{1}{N} \sum_{i=1}^N k_i^{out} = \frac{L}{N} \quad (1.12)$$

Lo anterior lo podemos expresar en forma matricial, por lo tanto el grado en una red no dirigida sería

$$k_i = \sum_{j=1}^N A_{ij} = \sum_{i=1}^N A_{ij} \quad (1.13)$$

Para una red dirigida sería

$$k_i^{in} = \sum_{j=1}^N A_{ij}, k_i^{out} = \sum_{i=1}^N A_{ij} \quad (1.14)$$

El número de enlaces

$$2L = \sum_{i=1}^N k_i^{in} = \sum_{j=1}^N k_i^{out} = \sum_{ij} A_{ij} \quad (1.15)$$

El *coeficiente de agrupamiento* mide la fracción de triadas cerradas con respecto al total de triadas en la red. Donde una triada o triada abierta es un camino de tamaño 2 y se dice que es una triada cerrada, camino cerrado o triángulo si existe un enlace que cierre este camino, dando lugar a un ciclo de tamaño 3.

$$C = \frac{\text{número de triángulos}}{\text{número máximo posible de ciclos de tamaño tres}} \quad (1.16)$$

El coeficiente de agrupamiento también se puede calcular de forma local,

$$C_i = \frac{2L_i}{k_i(k_i - 1)} \quad (1.17)$$

Donde L_i son los enlaces entre los vecinos de i .

Imaginemos que quisieramos identificar grupos de nodos que tengan una relación más estrecha que otros, porque a pesar de que podamos tener una red conexas, tal vez no todos los nodos están igualmente relacionados. Para esto existe una medida denominada modularidad Q la cual mide qué tan similares son los nodos de la misma clase, en función de sus enlaces. La modularidad (Newman, 2006) la podemos definir por

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{g_i g_j}, \quad (1.18)$$

donde m es el número de enlaces y $\delta_{g_i g_j}$ vale 1 siempre y cuando el nodo i y el nodo j se ubiquen en el mismo grupo (clase). De tal modo que obtendremos valores positivos y más cercanos a 1 cuanto más similares sean los nodos de ese grupo (en función de sus conexiones) y negativos cuanto menos similares sean. Existen diversos métodos para la selección de esos grupos, a los cuales a partir de ahora denominaremos *comunidades*.

1.2.4. Caminos y más

En redes, nos podría interesar medir distancias, en este caso la distancia estaría dada por la longitud del camino, donde un **camino** (Newman, 2010; Barabási and Pósfai, 2016) es el recorrido que se realiza sobre una serie de enlaces interconectados. Por ejemplo de la Fig. 1.5, un camino sobre la red no dirigida (panel A) podría ser: $P = [(3, 1), (1, 2), (2, 4)]$, su longitud la podemos definir como $l = |P|$, aplicado al camino definido anteriormente su longitud sería $l = 3$.

Sin embargo, es claro que habrá mejores caminos que otros, para expresar el **camino mas corto** de un nodo a otro, lo podemos hacer con d_{ij} . En este sentido el camino más corto para la misma red de la Fig. 1.5 sería $d_{3,4} = 2$ dado por el camino $P = [(3, 1), (1, 4)]$.

El **diametro** d_{max} es el más grande de los caminos más cortos existentes en una red.

Un **ciclo** es un camino que inicia y termina en el mismo nodo.

Los **caminos eulerianos** son caminos que pasan por cada enlace una única vez.

Los **caminos hamiltonianos** son caminos que pasan por cada nodo una única vez.

Se dice que una red es **conexa**, si para cualquier par de nodos existe un camino que los conecte. Para el caso contrario sería **disconexa**. En una red conexa existe una única componente, la cual contiene todos los enlaces y nodos de la misma, pero para una red disconexa pueden existir tantas componentes como nodos. Más formalmente una **componente** es un subconjunto de nodos en la cual existe un camino al resto de nodos existentes en dicha componente.

Para el caso de que la red sea dirigida, decimos que la componente es **fuertemente conectada** si para cualquier par i, j de nodos existe un camino tanto de i a j como de j a i . En caso de que esto no se cumpla, la componente sería una componente **débilmente conectada**. Los elementos i, j de la matriz de adyacencia elevada a una potencia n , es el número de caminos que conectan i y j .

1.2.5. Hipergráficas

Existen casos donde necesitamos unir a más de dos nodos con un mismo enlace, imaginemos que tenemos una familia de nodos que comparten una

característica en común, en este caso no sería muy útil poder enlazar a todos estos nodos con un mismo enlace. A estos enlaces se les conocen como *hiperenlaces* y a las redes con hiperenlaces se les denomina *hipergrafos* (Newman, 2010; Thurner Hanel, Rudolf A., Klimek, Peter, Oxford University Press., 2019). Donde los elementos de \mathcal{L} , son elementos del conjunto potencia de \mathcal{N} , $P(\mathcal{N})$. Recordar que el conjunto potencia se define como el conjunto de todos los subconjuntos el conjunto original, en este caso sería \mathcal{N} .

1.2.6. Gráficas potencia

Aquí los enlaces están definidos entre conjuntos de nodos. Supongamos el conjunto potencia $P(\mathcal{N})$ del conjunto de nodos \mathcal{N} . En las gráficas de potencia los nodos son elementos de $P(\mathcal{N})$, llamados nodos potencia. Esos nodos potencia, son conectados por enlaces potencia (Thurner Hanel, Rudolf A., Klimek, Peter, Oxford University Press., 2019). Dado que los nodos provienen del conjunto potencia, las dimensiones de la matriz de adyacencia sería $2^N \times 2^N$. Las gráficas de potencia pueden verse como una generalización de los hipergrafos y son útiles cuando necesitamos mapear un conjunto de elementos a un conjunto diferente de elementos.

1.2.7. Redes bipartitas

En busca evitar trabajar con redes tan gigantescas como las gráficas de potencia o hipergrafos, es que salen al rescate las *redes bipartitas*. Las *redes bipartitas* (Newman, 2010; Thurner Hanel, Rudolf A., Klimek, Peter, Oxford University Press., 2019; Barabási and Pósfai, 2016) consisten en dos conjuntos disjuntos de nodos, supongamos los conjuntos A y B , donde cada enlace conecta un nodo de A a uno de B y no existe enlace que conecte a dos nodos del mismo conjunto.

Para el caso de las redes bipartitas, la representación matricial se le conoce como matriz de incidencia. Esta se puede definir así: para n elementos del primer conjunto y g elementos del segundo conjunto, entonces la matriz de incidencia \mathbf{B} es una matriz de $g \times n$. Donde sus elementos B_{ij} valdrían $B_{ij} = 1$ si el elemento j del primer conjunto enlaza al i elemento del segundo conjunto y $B_{ij} = 0$ en caso contrario.

Dado que una red bipartita está formada por dos conjuntos disjuntos de nodos, uno podría tener interés en saber cómo están relacionados los nodos de un conjunto en particular, aunque esa relación fuera indirecta, esto lo

podemos hacer al realizar una proyección sobre uno de los conjuntos. Para realizar la proyección debemos escoger el conjunto sobre el cual la realizaremos, supongamos que A es el primer conjunto y B el segundo, para realizar la proyección sobre el conjunto A , debemos conectar todas las parejas de nodos que pertenecen al conjunto A en caso de que estos enlacen al mismo nodo perteneciente al conjunto B . Para el caso de la proyección sobre B enlazamos aquellas parejas de B que enlacen a un mismo nodo del conjunto A .

Además podemos generar una red bipartita pesada si variamos los pesos en función de las veces en que las parejas de nodos de A (o el conjunto sobre el cual se realiza la proyección) coinciden enlazando a un nodo distinto de B (o el conjunto en cuestión).

Las proyecciones pueden ser expresadas matemáticamente como $\mathbf{P} = \mathbf{B}^T \mathbf{B}$, de la misma forma para la proyección del otro conjunto $\mathbf{P}' = \mathbf{B} \mathbf{B}^T$.

A continuación, se muestra un ejemplo de una matriz de incidencia para la red bipartita mostrada en la Fig. 1.6 (en el panel A) y sus proyecciones respectivas (panel B y C).

$$\mathbf{B} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix} \quad (1.19)$$

1.2.8. Redes multirelacionales

Hay veces en que necesitamos expresar diferentes tipos de relación en una misma red, imaginemos que quisieramos mapear la red de contactos de un grupo de personas, pero resulta que de ese grupo algunos se contactan por teléfono, otros por Facebook, otros por WhatsApp y así. Una opción que tendríamos, sería generar una red por cada medio de comunicación, pero probablemente eso nos deje una serie de redes con algunos nodos aislados. Una forma de solucionarlo, podría ser la implementación de una red multirelacional, la cual además del conjunto de nodos \mathcal{N} y el conjunto de enlaces \mathcal{L} agrega un tercer conjunto \mathcal{A} referente al tipo de enlaces. Por lo tanto, en las redes multirelacionales (Thurner Hanel, Rudolf A., Klimek, Peter, Oxford University Press., 2019) cada enlace es definido por la tripleta $(n \in \mathcal{N}, l \in \mathcal{L}, \alpha \in \mathcal{A})$. Las redes multirelacionales pueden ser descritas por tensores de adyacencia. Supongamos i y j están enlazados en la capa α , entonces habría una entrada distinta de cero para el tensor M_{ij}^α .

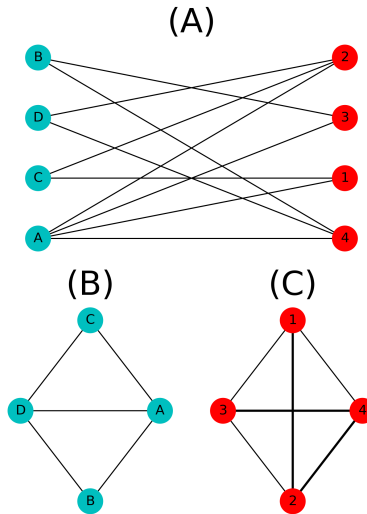


Figura 1.6: (A) Ejemplo de red bipartita con un conjunto coloreado en *azul* que contiene letras y un conjunto en *rojo* el cual contiene números.
 (B) Proyección de la red bipartita sobre el conjunto *azul*.
 (C) Proyección pesada de la red bipartita sobre el conjunto *rojo*.

1.2.9. Redes multicapa

Como su nombre lo dice, las redes multicapa (Newman, 2010; Thurner Hanel, Rudolf A., Klimek, Peter, Oxford University Press., 2019) tienen la capacidad de albergar múltiples capas, las cuales las denotaremos por \mathcal{M} . Cada nodo puede estar presente en cualquier subconjunto de capas. Cada enlace conecta a un nodo en la capa α con un nodo en la capa β , donde $\alpha, \beta \in \mathcal{M}$. La gran ventaja de las redes multicapa es que nos permite dinámicas sobre el tiempo, si secuenciamos los conjuntos de capas, donde cada conjunto corresponde a un instante de tiempo específico, entonces podríamos almacenar la dinámica. En términos de tensores lo podríamos describir así: $M_{ij}^\alpha(t)$, donde $t = 1, \dots, T$.

1.3. Antecedentes

1.3.1. Detección de comunidades

Una de las primeras cosas que se realizan al trabajar con redes sociales es identificar grupos con características en común y aunque esto lo podemos hacer de forma natural cuando trabajamos con redes pequeñas, sin embargo a medida que trabajamos con redes más grandes ya no parece útil realizar esta clasificación de forma manual, es cuando hacemos uso de la computadora y alguna metodología para realizar esta tarea.

Claro que la metodología cambia según la información que poseemos de la red y de cómo hayamos generamos dicha red, por ejemplo si trabajáramos con datos de Twitter: podemos hacer uso de técnicas de análisis de sentimientos para encontrar clasificar a los usuarios según su afinidad política (Caetano et al., 2018); haciendo uso de las interacciones entre usuarios y sus listas de seguidores para encontrar comunidades y hasta su función en el debate como en (Amor et al., 2016; Chamberlain et al., 2016) o hasta encontrar las comunidades en tiempo real respectivamente; haciendo uso de las interacciones entre los usuarios y la teoría de balance de Heider⁴ se puede identificar comunidades políticamente afines (Harmon, 1959; Amor et al., 2016); también podemos hacer primero la detección de comunidades dadas la red de interacciones y posteriormente extraer las temáticas en en las cuales son afines (Kim et al., 2013); haciendo uso del aprendizaje supervisado podemos entrenar a la computadora para reconocer el grupo político al que pertenece cada usuario, según los hashtags que emplea, palabras que usa y demás (Lai et al., 2017; Yuan and Crooks, 2018); empleando los hashtags que usan, las menciones y demás interacciones que tienen con otros usuarios (Darmon et al., 2015).

⁴Modelo propuesto por Heider F., el cual se define por una tríada cuyos actores son una persona P , otra persona O y un objeto X en el cual tanto la persona P como la persona O sienten afinidad (+) o rechazo (-). El objeto X podría ser cualquier cosa, material o inmaterial, incluso una tercera persona. Y para que esta triada esté balanceada todos los enlaces deberán ser positivos o un positivo y dos negativos, de cualquier otra forma la triada estaría desbalanceada. En el contexto social se puede ver este balance en el proverbio árabe “*el enemigo de mi enemigo es mi amigo*”.

1.3.2. Redes coordinadas

El ‘astroturfing’ es un termino referido al intento de crear la impresión de que existe un genuino apoyo popular generalizado hacia una política, individuo o producto, donde en realidad esto no ocurre, o al menos no en esa magnitud. Para lograr esto en las plataformas virtuales actuales, se utilizan múltiples identidades en línea y grupos de presión falsos para engañar al público haciéndole creer que la posición del ‘astroturfer’ es la opinión de las masas.

Dada esta situación es que se han hecho algunos estudios para tratar de identificar la coordinación de cuentas cuyo fin sea colocar algún mensaje dentro de las tendencias de las plataformas en cuestión. Por ejemplo el artículo de (Keller et al., 2020) propone generar una red de retweets para identificar las cuentas con mayor concentración de retweets y los grupos de usuarios que los soportan. También proponen la red de co-tweets en la cual uno dos usuarios si ambos usuarios publicaron el mismo mensaje en una ventana de tiempo de 1 minuto. Trabajan también con la red de co-retweet uniendo a dos usuarios, siempre y cuando ambos hayan retweeteado el mismo mensaje en el mismo minuto. Todo esto al final lo comparan con tendencias generales de frecuencias de retweets y con frecuencias de palabras (en los tweets publicados) para tratar de identificar la existencia de coordinación en estas redes.

En otro trabajo para la detección de redes coordinadas (Pacheco et al., 2021) proponen la intervención de las redes bipartitas para la detección de varios tipos de comportamientos coordinados, por ejemplo: usar el mismo `screen_name` (el nombre que funge como identificador único), imágenes publicadas, hashtags, co-retweets, acciones sincronizadas. De tal forma que aquí no sólo se limitaron a vincular a los usuarios por los textos publicados, también por las imágenes y por las ventanas de tiempo en que publicaban. Unos de las mejoras que se añaden en este artículo es que además de utilizar redes bipartitas, lo hacen pensando en la proyección posterior y el modo de re-pesarlas, como la distancia de jaccard, similaridad de coseno y otras de tal forma que pueden hacer uso de técnicas de clusterización para generar los grupos según su grado de similaridad.

También (Torres-Lugo et al., 2020) propone analizar una de las formas en que se construyen de las cámaras de eco, básicamente lo que hacen revisar los mensajes de menciones, únicamente estos y entonces ir vinculando el emisor con los mencionados, de tal forma que termina con una red de usuarios orquestada por unos cuantos. En este caso particular se ve lo fácil que

es generar una cámara de eco, los coordinadores comienzan a citar a su red de fieles seguidores para que estos participen apoyando en la discusión y a su vez convoquen a más usuarios para posicionar el mensaje en la plataforma.

1.3.3. Homofilia

Homofilia es el principio que propicia que la gente tienda a juntarse con sus similares (en ciertas cualidades o rasgos) y con mucha menor frecuencia a juntarse con personas que son más distintas (Miller et al., 2001).

1.3.4. Cámaras de eco

Por lo general en plataformas virtuales como Twitter los usuarios tienden a seguir sólo a las cuentas con las cuales concuerdan en algunos puntos de vista y no tienen conflictos fuertes en ningún tema. Esto tiende a ocasionar que los usuarios sólo se rodean de gente con la que concuerdan y no se entre en ningún especie de debate interno. De tal forma que cuando existen ciertas discusiones por temas políticos o ideas que tienden a generar conflicto de opiniones se tienden a fragmentar en 2 o más grupos (según la cantidad de vertientes que existan sobre el tema en cuestión) y en cada uno pareciera existir una opinión en común que resuena en todo el grupo. A esto se le conoce como cámaras de eco, las cuales pueden hacer pensar a los usuarios que todos piensan igual que ellos, reforzando la idea de que su punto de vista es el correcto, sin embargo esto sólo es un sesgo ocasionado por la forma en la que generó su red social, a este sesgo cognitivo se le conoce como el ‘sesgo de confirmación’ el cual refuerza las ideas que confirman sus puntos de vista y disminuye desproporcionadamente la consideración a posibles alternativas. Sobre esto se han realizado varios artículos para entender cuáles son las razones que originan y magnifican estos comportamientos. Por ejemplo, tenemos el artículo de (Sasahara et al., 2020) en el que muestran cómo el repetir acciones por influencia social y el dejar de seguir a aquellos usuarios que tienen opiniones contrarias a las propias propician una acelerada emergencia de las cámaras de eco y a su vez polarizan la red pues estas acciones tienden a aislar a los grupos según sus ideologías políticas, creencias, entre otros. En otro artículo de (Torres-Lugo et al., 2020) analizan un método que se usa para generar estas cámaras de eco, básicamente consiste en que un pequeño grupo de usuarios cita a otros (por lo general de su misma ideología política) en un tweet, esto con la intención de convocarlos a que participen, estos a su vez repiten la acción y citan a otros. De tal forma que en poco tiempo

ya han formado grupos masivos que opinarán en la misma sintonía sobre un tema dado, creando de forma no espontanea las cámaras de eco y a su vez sirve como ‘astroturfing’ pues ya con esta masa pueden colocar un mensaje en la opinión pública, sin importar que esta acción no haya sido espontanea.

Dados estos antecedentes y la breve introducción matemática a teoría de redes, podemos continuar al siguiente capítulo, donde abordaremos la forma en la cual recolectamos los datos, así como la metodología que seguimos para realizar nuestros análisis. Recordar que este trabajo tiene como objetivo la extracción de un conjunto de usuarios que estén vinculados por razones políticas, laborales, familiares, intelectuales y no sólo vinculados de forma circunstancial a partir de un conjunto de datos recabados de Twitter.

Capítulo 2

Metodología

En este capítulo se abordan los procedimientos seguidos para la recopilación de datos, la caracterización de la red de retweets en Twitter hasta la construcción de las redes de co-eventos

2.1. Recopilación de datos

Para la recolección de los datos se utilizó la librería `tweepy`¹ de python debido a la facilidad con la que se puede acceder a los “tweets”, perfiles de usuario, búsqueda de “hashtags”. Se recolectaron tweets por un periodo de aproximadamente 3 semanas durante el mes de mayo de 2020 (específicamente del 6-21 de mayo), dicho periodo de tiempo fue seleccionado fue seleccionado así debido a restricciones computacionales y las fechas en las que se ejecutó la recolección de datos fue propuesta, debido a que al rededor de esas fechas el gobierno federal de México estimaba que la transmisión del SARS-CoV-2 disminuiría a tal grado que algunos municipios podrían comenzar con la reactivación económica, debido a esto supusimos que en esas fechas habría una fuerte actividad en Twitter. Para la captura de tweets en tiempo real, `tweepy` te permite capturar tweets por zona geográfica o aquellos que coincidan con una lista de palabras (que uno genere previamente), el problema de utilizar el filtro de la zona geográfica es que sólo al rededor del 1% de los usuarios tiene activada su geolocalización (Schlosser et al., 2021; Zohar, 2021; Sloan and Morgan, 2015). Por ello decidimos usar el filtro por palabras clave, de tal forma que sólo nos quedaríamos con aquellos usuarios que emitieran mensajes con un conjunto de palabras específicas.

¹<https://www.tweepy.org/>

Dado que en 2020 México se encontraba en plena pandemia ocasionada por el virus SARS–CoV–2 y además en Twitter se vivían discusiones muy álgidas¹ en torno al sistema de salud mexicano, se optó por recopilar tweets, retweets, quotes cuya temática estuviera de alguna forma relacionada con la emergencia sanitaria y la respuesta del gobierno de México, de tal forma que los datos obtenidos tenderían a seguir una discusión mayormente política.

Entonces primero generamos el conjunto de palabras relacionadas con figuras políticas o instancias relacionadas al gobierno de México.

1. **amlo, lopezobrador** López Obrador actual presidente de México (2018–2024) y miembro del partido MORENA.
2. **morena** Partido que actualmente cuenta con la mayoría parlamentaria en México.
3. **4t** La autodenominada “cuarta transformación”, se refiere a los cambios promovidos por el actual gobierno.
4. **gatell** Hugo López–Gatell Ramírez, Actual Subsecretario de Prevención y Promoción de la Salud, a cargo de la estrategia y respuesta al COVID–19 de México.
5. **FelipeCalderon** Es la cuenta de Twitter de Felipe Calderón, expresidente de México del 2006 al 2012 y muy crítico del actual gobierno en las redes sociales.
6. **SSalud_mx** Cuenta oficial de Twitter del sistema de salud mexicano.
7. **insabi** Instituto de Salud para el Bienestar.

Con la intención de que se hiciera referencia a la pandemia, el mensaje también deberá involucrar por lo menos una de las siguientes palabras clave:

2.1 encasa, sanitaria, sarampion, cuarentena, salud, fase3 Palabras relacionadas a mensajes emitidos por la Secretaría de Salud. A estas palabras las denominaremos “*etiquetas de salud*”.

2.2 covid, coronavirus, corona, pulmonia, neumonia, sarscov2, respir A estas las denominaremos “*etiquetas de COVID–19*”.

¹<https://www.eluniversal.com.mx/tag/insabi>

Al final de la recolección de los datos realizamos un primer filtro para quedarnos únicamente con aquellos usuarios que hayan emitido por lo menos una vez un mensaje que involucrara al menos una de las palabras perteneciente a alguna de las figuras políticas 1–7 y al menos otra palabra perteneciente a cualesquiera de las dos etiquetas antes mencionadas. De modo tal que esto nos aseguraría que cada usuario en cuestión recopilado de esta forma habría emitido alguna opinión y/o mensaje que tuviera que ver con la salud y México. Al final de esto nos quedamos con 2,950,080 tweets.

2.2. Primer acercamiento

Como primer acercamiento, primero graficamos la red de discusión que se genera tomando en cuenta los retweets, quotes y comments Fig.2.1, de forma tal que generamos un enlace de un usuario “A” a un usuario “B” si el usuario “A” retweeteo, comentó o realizó un quote a usuario “B”, el problema es que no se alcanza a percibir de forma clara comunidades bien definidas. Lo que sí se alcanza a apreciar con colores más claros, son los usuarios que tienen más conexiones, las cuales se pueden deber a que han sido más retweeteados, citados con las “quotes”, o bien, etiquetados en un comentario, en azul marcamos los enlaces bidireccionales de peso superior a 8 lo cual implica que en esos usuarios hay una comunicación mucho más intensa que en el resto.

Debido a que en este primer acercamiento no logramos identificar de forma clara comunidades con ideologías políticas afines o algún interés común bien definido optamos por conservar únicamente los retweets (Rt), consideramos que el acto de retweetear un mensaje es un voto de apoyo, de alguna forma se concuerda con el mensaje y si esto lo ampliamos a un espacio de tiempo extenso, probablemente obtendríamos una red que se acopla según intereses comunes de tal manera que podríamos apreciar grupos más conectados que otros debido a que tienen un grado de afinidad mayor o existe una coordinación detrás. Para la red de retweets el grado de entrada (k_{in}), puede ser interpretado como el grado de popularidad de cada usuario y el grado de salida (k_{out}), como la actividad generada por cada usuario (a cuantos usuarios retweeteo).

En la Fig. 2.2, se muestra la red resultante. Con diferentes colores, se pintaron las distintas comunidades detectadas, estas fueron detectadas ha-

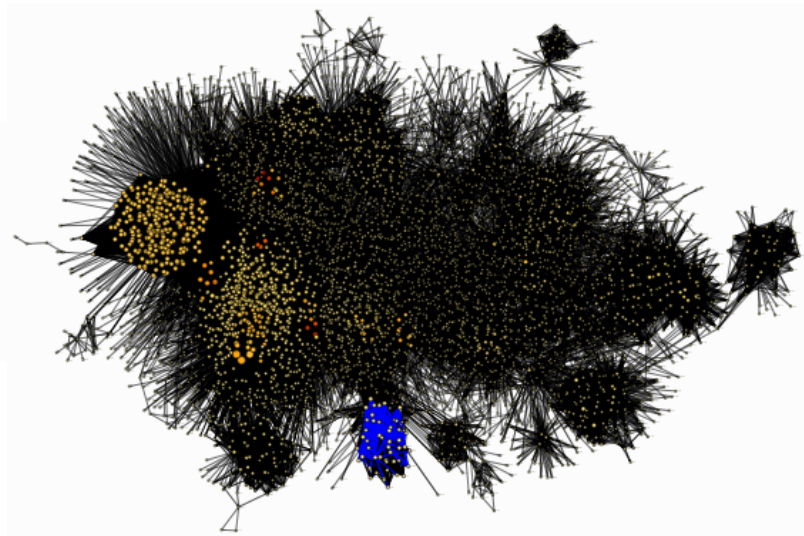


Figura 2.1: Figura de una discusión en Twitter. Los enlaces entre nodos representan retweets, quotes o menciones de un usuario a otro. Se colocó en color azul los enlaces bidireccionales cuyo peso era superior a 8 debido a que esto nos indicaría cuales son los usuarios que tienen una comunicación más intensa que el resto de la población. También colocamos en distintas tonalidades de naranja los usuarios con mayor grado de entrada para apreciar a los usuarios más populares.

ciendo uso del algoritmo de Louvain² (Blondel et al., 2008), el cual es un algoritmo de modularidad fuertemente optimizado para redes grandes. Aquí, se pueden identificar dos grupos polarizados, cada uno con ideologías políticas opuestas, hemos distinguido el grupo de los “pro-gobierno” en color guinda y el grupo de los “opositores” con color azul. La asociación a estos grupos fue realizada extrayendo una muestra de usuarios para ambas comunidades e identificando temas en común entre estos usuarios, tales como las imágenes compartidas, los tweets compartidos, la autodescripción de los usuarios y demás características personales de cada usuario. Sin embargo, en la red mostrada no existe una separación clara entre ambos grupos, se aprecia una conexión no despreciable. Esto se lo atribuimos a que existe una gran cantidad de usuarios que retweetean cosas que no tienen un rasgo político tan definido, por ejemplo noticias, estadísticas sobre la pandemia y

²Para mayor detalle sobre el algoritmo de Louvain revisar el Ap. B.

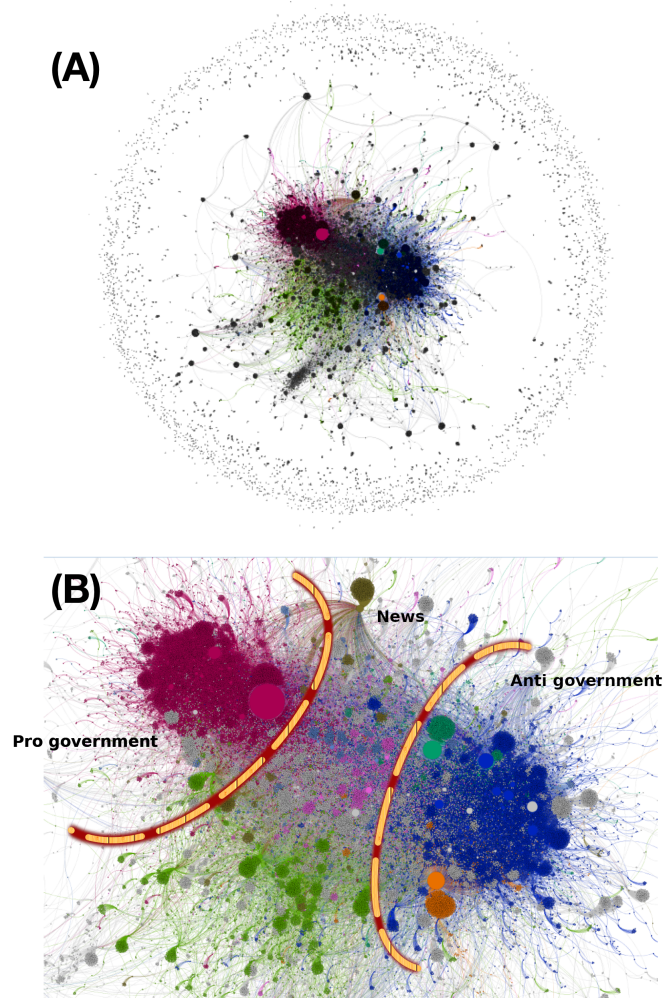


Figura 2.2: (A) Comunidades identificadas en discusiones políticas mostradas con diferentes colores. Las comunidades fueron detectadas haciendo uso del algoritmo de Louvain (Blondel et al., 2008). El tamaño de cada nodo está asociado al grado de entrada. La red está polarizada principalmente en dos grupos con puntos de vista opuestos relacionados con los temas de COVID, de ún lado están los “pro–gobierno” con color guinda y del otro los “opositores” con color azul. Los links van en sentido a las manecillas del reloj son links que salen y vice versa. (B) Es un zoom de la misma red.

otros temas en los que ambas comunidades participan activamente (puntos

en común). Aquellos nodos que no entran dentro de la componente gigante es porque no retweetearon un tweet o tema perteneciente al centro de la discusión “trending topic” o en un caso muy extremo, son usuarios cuya red de seguidores tienen una conexión muy débil o nula con el resto de la sociedad.

Como primer análisis sobre nuestra red de Rt calculamos las distribuciones $P(k_{in} + 1)$ and $P(k_{out} + 1)$. Es importante notar que se aplicó el desplazamiento en k por 1 para poder graficar aquellos nodos con $k_{in} = k_{out} = 0$; debido a que ellos conforman la mayor parte de la red, pues la mayoría de los usuarios a pesar de que generan tweets, lo normal es que no reciban retweets de vuelta pues su grado de popularidad no es muy alto. 98% de los nodos pertenecen a este tipo. En la Fig. 2.3 se muestran esas distribuciones y como era de esperarse siguen un comportamiento de ley de potencias: $P(k_{in}) \sim (k_{in} + 1)^{-2,0}$ y $P(k_{out}) \sim (k_{out} + 1)^{-4,4}$, respectivamente.

Como se puede apreciar, en la Fig. 2.3 el 99,9% de los usuarios no retweetean a más de 20 usuarios distintos, lo cual implica que los usuarios son muy selectivos a la hora de seleccionar al usuario al cual se le va a apoyar con un retweet. Sin embargo esos 20 apoyos son suficientes para generar a los dichos “influencer” (usuarios que concentran una gran cantidad de seguidores en comparación con el resto de los usuarios), de hecho se puede apreciar que hay algunos usuarios que recolectan más de 1000 retweets a pesar de la limitante previamente mencionada. Cuando analizamos a los usuarios en función del grado de entrada *vs.* el grado de salida para cada usuario (esto no se muestra en la gráfica), podemos notar que los usuarios más populares reciben más de 500 retweets y no generan más de 4 retweets, mientras que los usuarios con un grado de entrada $k_{in} \leq 10$ envían el 99,5% de todos los retweets.

2.3. Red de co–eventos

Como lo mencionamos al inicio, lo primero que quisieramos es obtener comunidades bien definidas en las que los usuarios se agrupen según sus ideologías políticas o algún interés común y entonces podamos identificar grupos opositores y demás. Pero ¿cómo podríamos obtener esas redes con los datos previamente recabados?, pues en algunos artículos ya se propone la aplicaciones de redes bipartitas para la identificación de grupos coordinados o con comportamientos similares (Keller et al., 2020; Pacheco et al., 2021). En nuestro caso hemos diseñado una red de “co–eventos”, la cual se forma de la siguiente forma: generamos un enlace entre los usuarios el retweet

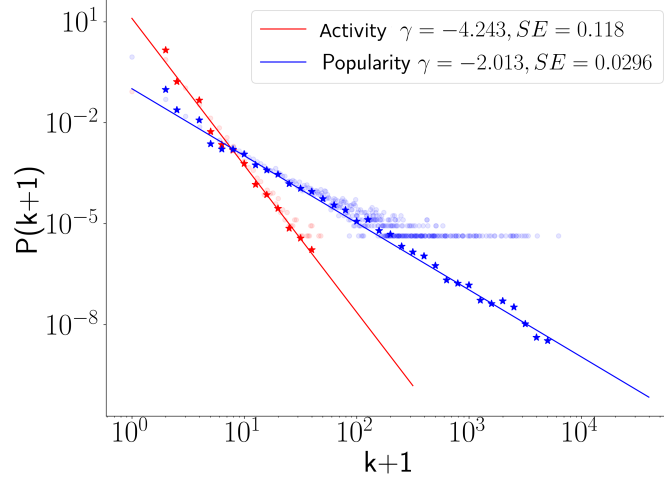


Figura 2.3: Distribuciones en log–log de la popularidad (“popularity” azul) y el número de usuarios a los que se le dio apoyo con un retweet (“activity” rojo). En colores menos intensos presentamos los datos crudos y en colores más intensos los datos al aplicarles un agrupamiento logaritmico sobre estos trazando líneas rectas correspondientes a los mejores ajustes para estas leyes de potencias. Favor de observar cómo la popularidad decrece con menor pendiente que la actividad.

en cuestión si dos usuarios retweetearon el mismo mensaje en la misma ventana de tiempo (la ventana que empleamos es de una hora), en caso de que coincidan en más de una ocasión incrementaremos el peso del enlace en función de las veces que coincidieron, de la misma forma esto se puede aplicar para el caso de los “hashtags”. Una vez obtenida esta red bipartita, formada por los conjuntos disjuntos de los mensajes y los usuarios podemos generar una proyección pesada (Barabási, 2013) hacia los usuarios, de tal modo que obtendremos una red de los usuarios en función de sus coincidencias con otros usuarios. Definimos $w_{ij}(l)$ entre dos usuarios i y j , los cuales tienen un retweet o hashtag en común en la ventana de tiempo que inicia en $l\Delta t$ donde $\Delta t = 1$ es una hora y l es el número de la ventana. Finalmente, el peso total entre los usuarios i y j es,

$$w_{ij} = \sum_{l=1}^L w_{ij}(l) \quad (2.1)$$

donde $L = T/\Delta t$ y T es el periodo de observación, que para este caso es de 3 semanas.

También definimos $k_w(i)$ donde,

$$k_w(i) = \sum_{j \neq i} w_{ij} \quad (2.2)$$

Al generar la proyección de la bipartita sobre los usuarios es que nosotros obtenemos nuestra red de co-eventos, co-retweet para el caso de los retweets y co-hashtag para el caso de los hashtags. Ya con ella, se decidió eliminar todos los enlaces con peso 1 pues el hecho de que se tenga una única coincidencia con otro usuario, no debería implicar que de alguna forma pertenecen al mismo grupo ni podría implicar un grado de coordinación o algo similar, así que eliminando estos enlaces acabamos con la mayoría de las interacciones dadas simplemente por casualidad. Antes de remover los enlaces con $w_{ij} = 1$, el total de nodos era de 45,454 y al haberlos removido se redujo a 8,105.

2.4. Análisis de susceptibilidad

Lo siguiente que proponemos es un análisis de percolación sobre las redes de co-retweet y co-hashtag. Para realizar esto, nosotros calculamos el tamaño promedio de las componentes (quitando previamente la componente gigante), también asociado como el segundo momento del número n_s de s tamaño de la componente (Onnela et al., 2007),

$$\langle s \rangle = \frac{\sum'_s n_s s^2}{\sum'_s n_s s} \quad (2.3)$$

donde las sumas primarias van sobre las componentes de todos los tamaños excluyendo la componente gigante, la cual es la componente más grande (Iñiguez et al., 2009).

Para ver cómo se comporta la susceptibilidad en función de las componentes de la red, renormalizamos los pesos w_{ij} de las redes de co-retweet y co-hashtag como una función de sus coincidencias promedio y el grado de participaciones de la siguiente forma:

$$w_{ij}^* = \frac{2w_{ij}}{k_w(i) + k_w(j)}, \quad (2.4)$$

Usando esta definición, $w_{ij}^* = 1$ es porque ambos usuarios i y j han coincidido en todos sus mensajes. De tal forma que podríamos ver el peso

del enlace como el grado de homofilia que existe entre ambos usuarios, siendo 1 el grado máximo y 0 ningún parecido entre ambos.

Ahora podemos estudiar la susceptibilidad de la red generada por los pesos w_{ij}^* . Primero, removeremos de forma progresiva los enlaces de forma descendente, partiendo del enlace con peso máximo w_{ij}^* y terminando con el enlace de peso mínimo w_{ij}^* . Posteriormente lo haremos al revés, con el objeto de ver si existe una transición de fase de percolación en el proceso.

Capítulo 3

Resultados

Una vez descrita la metodología, en este capítulo procederemos a mostrar los resultados obtenidos de las simulaciones para las redes de co- eventos.

3.1. Red de co-eventos

La Fig. 3.1 muestra la red resultante de co-retweets, En el Panel (A) es fácil apreciar que la red está separada en tres grupos principales¹, de izquierda a derecha encontramos al grupo de los pro-gobierno, los opositores y el grupo covid (usuarios que básicamente sólo retweetearon noticias sobre el covid). Cada comunidad está coloreada con un color distinto, en el Panel (A) podemos apreciar los tres principales bloques en los que se agrupan los usuarios y en el Panel (B) podemos observar como son las interacciones entre las comunidades que forman cada uno de ellos. Es importante resaltar que las interacciones dentro de las comunidades, comprenden el 96 % del total de los enlaces. Del 4 % restante podemos decir que la mayoría de los enlaces están dentro de las comunidades pertenecientes al bloque pro-gobierno. Las comunidades encontradas en la red de co-retweet podrían indicar algún tipo de organización, la existencia de intereses en común o metas en común. Por lo tanto, esto confirma que el COVID-19 como tópico fue politizado y trajo consigo un debate altamente polarizado en una situación de emergencia tal como ha sido la pandemia.

Observamos en la Fig. 3.1 perteneciente al cúmulo rosa del Panel (A), tres nodos pertenecientes a la comunidad de color rosa parecieran concentrar

¹Para más detalles acerca de cómo identificamos la relación en común de estos grupos revisar el Ap.C

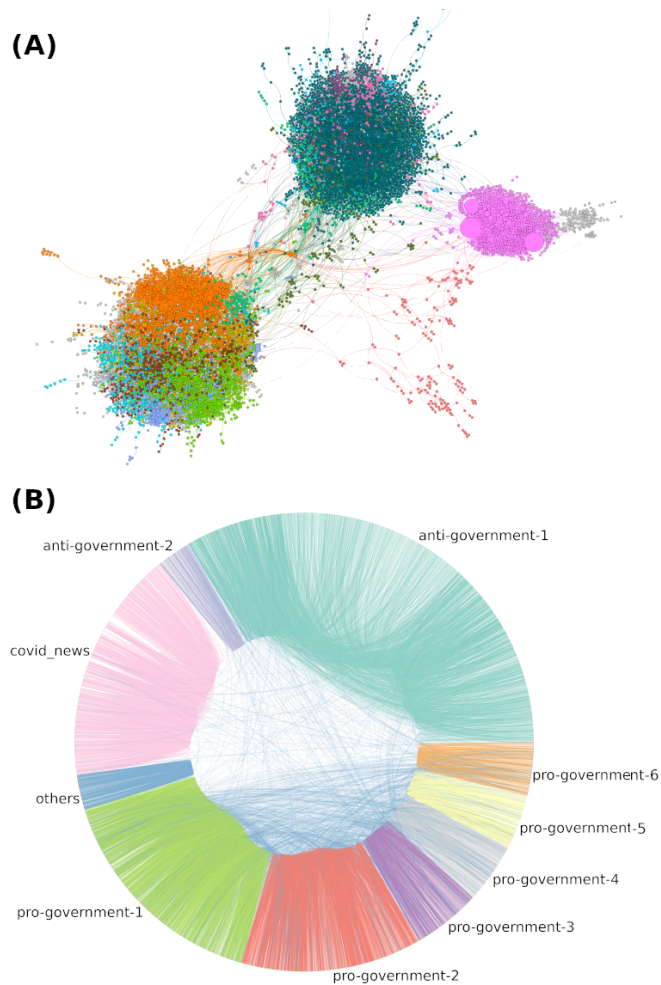


Figura 3.1: Panel (A), Red de co-retweet con una ventana de tiempo de una hora. Colores diferentes fueron asignados para distinguir a cada comunidad, las cuales fueron encontradas haciendo uso del método de Louvain. Se pueden apreciar 3 grandes cúmulos, de izquierda a derecha encontramos al grupo de los pro-gobierno, los opositores y el grupo covid (usuarios que básicamente sólo retweetearon noticias sobre el covid). El tamaño del nodo es proporcional al grado del mismo. Panel (B), es la misma red de co-retweets, pero ploteada en un diseño circular. Los enlaces dentro de la misma comunidad tienen el color de la comunidad, caso contrario son coloreados con un azul oscuro.

una gran cantidad de los enlaces, para comprobar esto se optó por generar una gráfica con la distribución de grado y obtuvimos que 3 usuarios tuvieron un comportamiento atípico al resto de usuarios en la red, pues concentraron una red de conexiones muy por encima del resto.

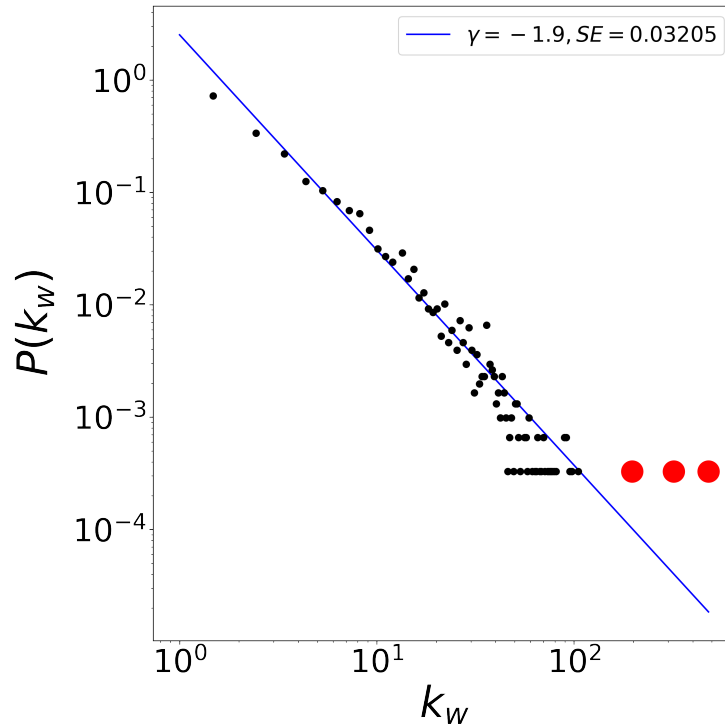


Figura 3.2: Distribución de grado para la red de co–retweet para ventanas de tiempo de una hora. Hay tres hubs, los cuales se distinguen con círculos de color rojo, con un claro comportamiento de valor atípico. Estos tres hubs se ubican en la comunidad de color rosa en la Fig. 3.1.

Estos usuarios fueron identificados en la Fig. 3.2 con estrellas. Realizando una búsqueda más en profundidad sobre sus perfiles de usuario, encontramos que CoronaUpdateBot está configurado para retweetear todo lo que encuentre sobre coronavirus, el segundo world_new_seng retweeteaba noticias que tienen que ver con la pandemia y el tercero BillEsteem es una cuenta que ya no existe, sin embargo estaba muy activa durante el periodo de recolección de datos. Vista esta red como una red de contagios tenemos que los tres usuarios cumplen con la definición de super propagadores, pues gracias a

todas esas conexiones y actividad pueden propagar información muy rápidamente a una multitud de usuarios.

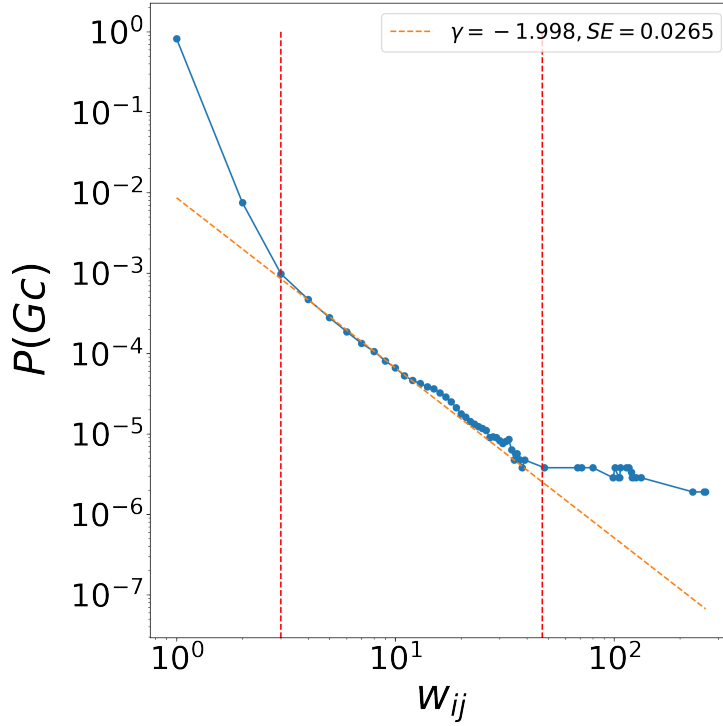


Figura 3.3: Distribución del tamaño de la componente gigante al ir removiendo los enlaces pesados (no normalizados) de forma progresiva de forma ascendente w_{ij} . Es importante notar que hemos encontrado una región en los intervalos que van de $w_{ij} = 3$ al $w_{ij} = 47$ que se puede aproximar siguiendo una ley de potencias. Una vez que los datos se apartan de la ley de potencias, $P(G_c)$ se alínea casi de forma horizontala partir de los filtros $w_{ij} \geq 48$ indicando algún tipo de organización.

Para descubrir más información sobre la red social contenida en la red de co-retweet, removimos los enlaces de forma progresiva de forma ascendente w_{ij} para ver cómo se comporta el tamaño de la componente gigante. Encontramos que cuando los enlaces pesados 1 al 3 son removidos, el tamaño de la componente gigante decrece abruptamente (ver Fig.3.3).

Cuando eliminamos los enlaces del $k_w = 3$ al $k_w = 48$, el tamaño de

la componente más grande pareciera seguir una disminución siguiendo una ley de potencias. En este intervalo la componente gigante es mucho más grande que el resto de componentes. Como se aprecia en la última parte de la Fig.3.3, desde el enlace $k_w = 48$ en adelante la disminución del tamaño de la componente más grande es mucho más lento. En este intervalo, los componentes obtenidos para cada uno de los pesos removidos son mucho más homogéneos, pareciera no existir gran diferencia entre la componente de tamaño mayor que el resto de las componentes. Por lo tanto, después del enlace con peso 48 las coincidencias con otro usuario deben reflejar algún tipo de organización o un tema de interés común que genera un alto grado de cohesión.

3.2. Susceptibilidad

En la Fig. 3.4, se muestra $\langle s \rangle$ en función de la remoción de enlaces progresivamente en forma ascendente respecto a w_{ij}^* . Para ambas redes co-retweet y co-hashtag se presentan picos en la susceptibilidad. La Fig. 3.5 muestra $\langle s \rangle$ pero removiendo los enlaces pesados en sentido opuesto (de forma descendente).

El objetivo de este análisis es observar si la susceptibilidad varía de forma similar al experimento realizado por Onnela et al. (2007), es decir, se esperaba que al eliminar progresivamente los enlaces de forma descendente, la susceptibilidad decayera suavemente y al realizarlo de forma ascendente, esperábamos ver una transición de fase como la figura 17 del experimento de las redes de telefonía celular (Onnela et al., 2007). Esto sugiere la importancia de los enlaces débiles en esta red, la muestran como una red social con un fuerte comportamiento real (no virtual). Al identificar los enlaces débiles eliminados justo antes de que la red se pulverice en muchos grupos pequeños, se puede detectar el esqueleto de la red que permite que el sistema se mantenga unido, lo cual es información extremadamente importante ya que estas conexiones son vitales para mantener cohesionadas las redes grandes (redes sociales). Estas conexiones son sumamente importantes en las redes sociales, en experimentos realizados por Granovetter (1973), se han identificado algunas de las cualidades de los enlaces débiles la más relevante para nuestro experimento, es que se esos enlaces nos conectan a grupos sociales a los que no tenemos acceso vía nuestros enlaces fuertes, es decir, aquellos contactos ocasionales son los que nos conectan a otros círculos sociales, propiciando así el fenómeno de mundo pequeño.

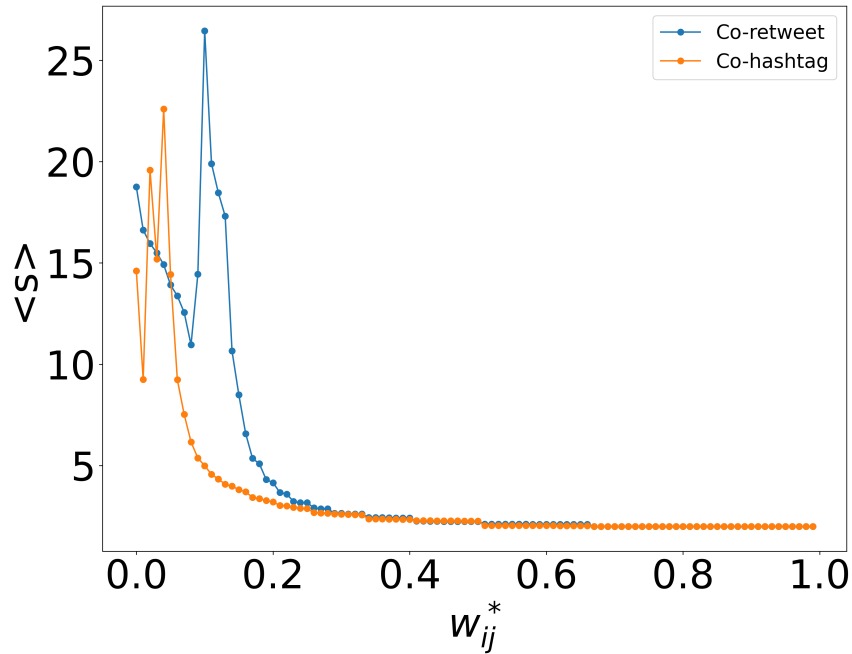


Figura 3.4: Susceptibilidad en función de la remoción ascendente de los enlaces pesados w_{ij}^* para la redes de co-retweet y co-hashtag. Notar que ambas redes presentan un pico, indicando una alta susceptibilidad por remover los enlaces debiles, una característica que es usualmente vista en las redes sociales físicas. Las líneas entre los puntos fueron colocadas para ayuda visual.

Observando las Figs. 3.4,3.5 podemos observar que en general, la susceptibilidad se comporta como se esperaba. Eso nos dice que en general, su comportamiento pertenece a una red social, lo cual era de esperarse, pues Twitter es una plataforma virtual donde interactuan seres humanos, además de bots. Sin embargo, observando con más detalle la Fig. 3.5 y comparándola con la figura 17 de Onnela et al. (2007) es claro que la susceptibilidad de nuestras redes de co-eventos es mucho más alta que la obtenida por ellos debido a que el grado de conexiones es bastante más alto que el que se da en telefonía celular. Si comparamos la red de co-retweet con la co-hashtag vemos que la de co-hashtag cae mucho más rápido que la de co-retweet, esto se debe a que a diferencia de los retweets, los usuarios se definen mejor por una variedad de hashtags, generando grupos bien definidos por afinidades muy peculiares. En cambio para los retweets las personas los grupos se dan por una afinidad más general, en nuestro caso particular esa afinidad fue política.

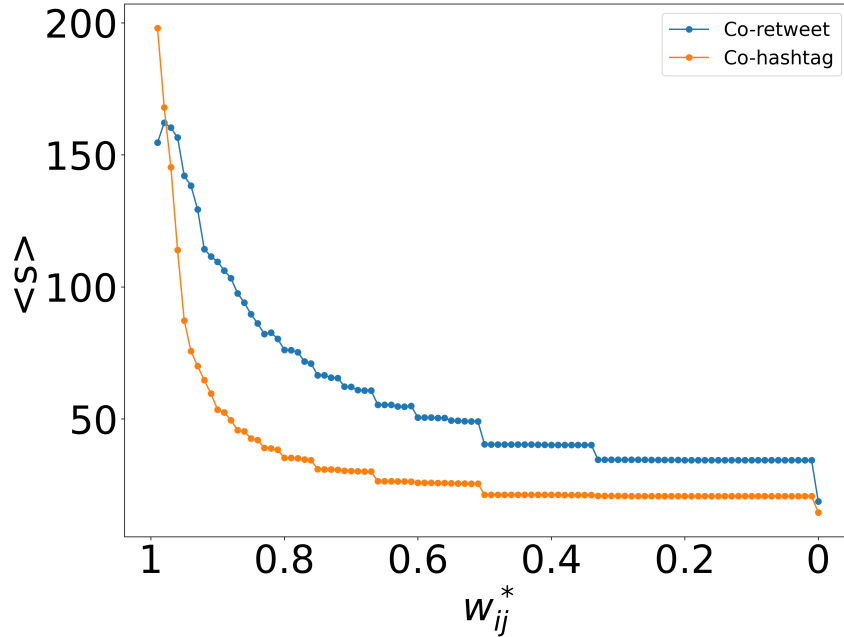


Figura 3.5: susceptibilidad en función de la remoción descendente de los enlaces pesados de forma progresiva w_{ij}^* para las redes de co-retweet y co-hashtag. Las líneas entre los puntos fueron colocadas para ayuda visual.

Analizando la Fig. 3.4 observamos que la transición de fase se da muy cerca del $w_{ij}^* \approx 0,1$ para el caso de la red de co-retweet y $w_{ij}^* \approx 0,02$ para la red de co-hashtag, sumado a esto, vemos que la susceptibilidad de ambas redes decae mucho más rápido que en la Fig. 3.5 y pasando los enlaces de $w_{ij}^* \approx 0,3$ en ambas redes, pareciera que llega a una región bastante estable, con subredes de tamaño 3 en promedio y posteriormente decrece un poco. Esto nos dice que en general los grupos sociales mejor cohesionados en Twitter, están formados por triadas y parejas.

3.3. Agrupamientos a diferentes niveles de la red de co-eventos

A partir del análisis realizado sobre la Fig. 3.4 identificamos la transición de fase y con su ubicación colocamos el límite de nuestro primer segmento de red $w_{ij}^* \leq 0,1$ ($w_{ij}^* \leq 0,02$ para el caso de la red de co-hashtag), como

límite superior de nuestro segundo segmento proponemos $w_{ij}^* \approx 0,3$ pues al rededor de este punto pareciera que la susceptibilidad no varía mucho. Este mismo límite sera utilizado como límite inferior de nuestro tercer segmento de red. Dados los límite propuestos, procedimos a realizar una inspección visual a los perfiles de usuarios de una muestra aleatoria de 50 usuarios para cada segmento de red. Como resultado de dichas observaciones hemos identificado que existen 3 tipos de usuarios en la red w_{ij}^* , según el estrato de la red:

Nivel bajo $w_{ij}^* \leq 0,1$: Aquí encontramos a los bots o propagadores de noticias, son los usuarios cuya única finalidad es difundir una serie de mensajes, su tasa de publicación es considerablemente más alta que el

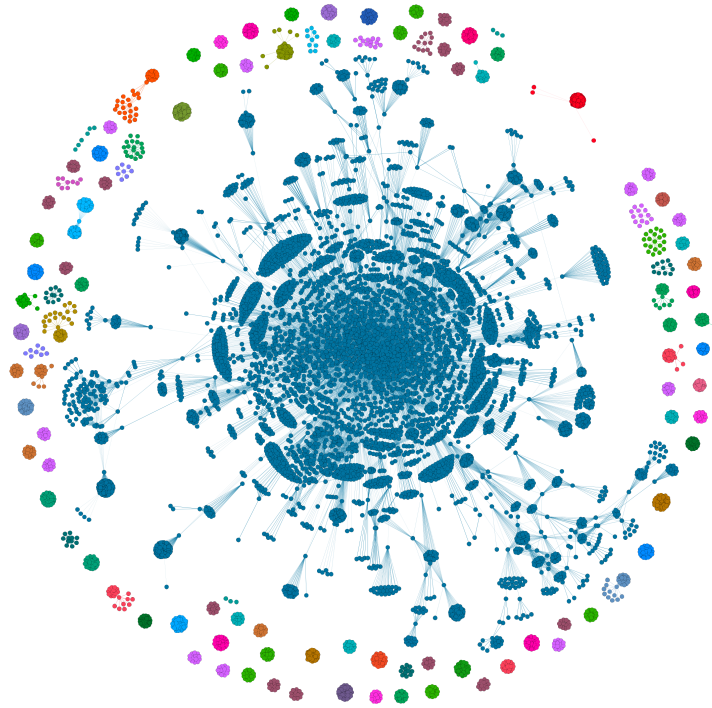


Figura 3.6: Relaciones al grado más bajo 0,0-0,1, donde cada color refiere a un tamaño de componente específico, la mayoría de los las componentes ronda por los 25 usuarios. En la componente gigante se concentran la mayoría de los bots o propagadores de noticias y concentra un total de 5946 usuarios.

de la mayoría de los usuarios, tienden a encontrarse en el extremo derecho de la gráfica de la Fig.3.1, es decir son los grandes concentradores de contactos y forman el grupo mayoritario de la red, para el caso particular de esta red de co-retweet fueron identificados como la componente gigante de la Fig.3.6. Una cosa que es importante resaltar, es que para el caso de la red de co-retweet sí se encontraron 3 bots que distaban por mucho del resto, como se puede observar en la Fig.3.2, sin embargo en la red de co-hashtag no se detectaron hubs que se separaran tanto del resto.

Nivel más alto $w_{ij}^* \geq 0,3$: Aquí detectamos usuarios altamente coordinados, estos son usuarios que tienden a retweetearse mutuamente o pu-

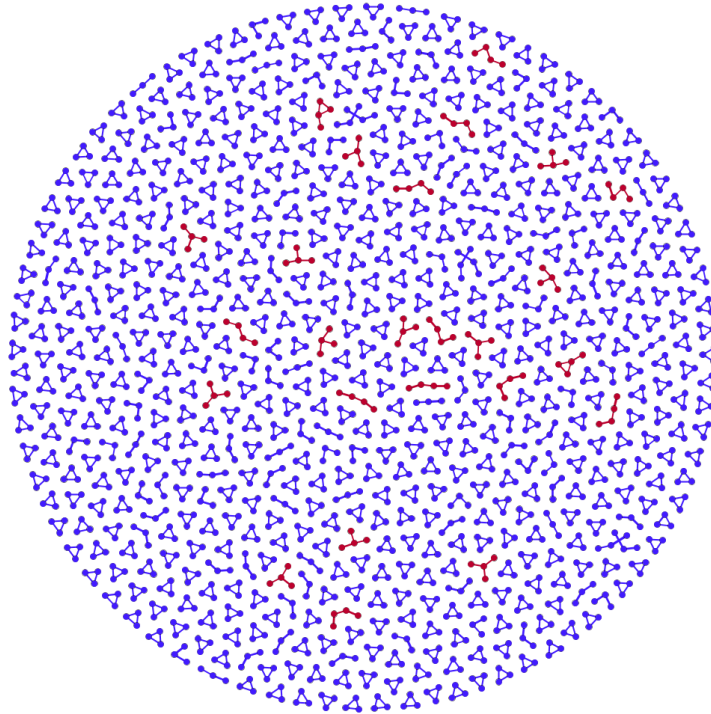


Figura 3.7: Relaciones sólidas 0,3-1, cada color refiere a un tamaño de componente específico. Aquí se observan los usuarios altamente coordinados, fueron eliminados los enlaces de parejas simplemente por estética, sin embargo las parejas concentran el 80 % de las componentes, mientras que las componentes de tamaño tres concentran casi el 20 % restante.

blicar los mismos hashtags en los mismos periodos de tiempo, normalmente son grupos de no más de 3 usuarios que tratan de darse auto-difusión. Este bloque posee 6858 usuarios. Las diferencias entre la red de co-retweet y co-hashtag fueron considerables, aunque no se dieron a nivel topológico si no a nivel social, curiosamente para el caso de la red de co-retweet no se encontró un lazo social definido entre los usuarios enlazados, a diferencia de las conexiones que se encontraron en la red de co-hashtag, aquí prácticamente todas las conexiones implicaban una relación laboral. Los resultados de esta red los podemos observar en la Fig.3.7.

Nivel intermedio $0,1 < w_{ij}^* < 0,3$: Encontramos usuarios normales, estos usuarios no presentan un patrón bien definido como el de los anteriores, parecieran estar relacionados simplemente por algunos puntos en común con los usuarios con quienes comparten conexiones, como en un grupo social real. Cuenta con 6373 usuarios y los podemos observar en la Fig. 3.8.

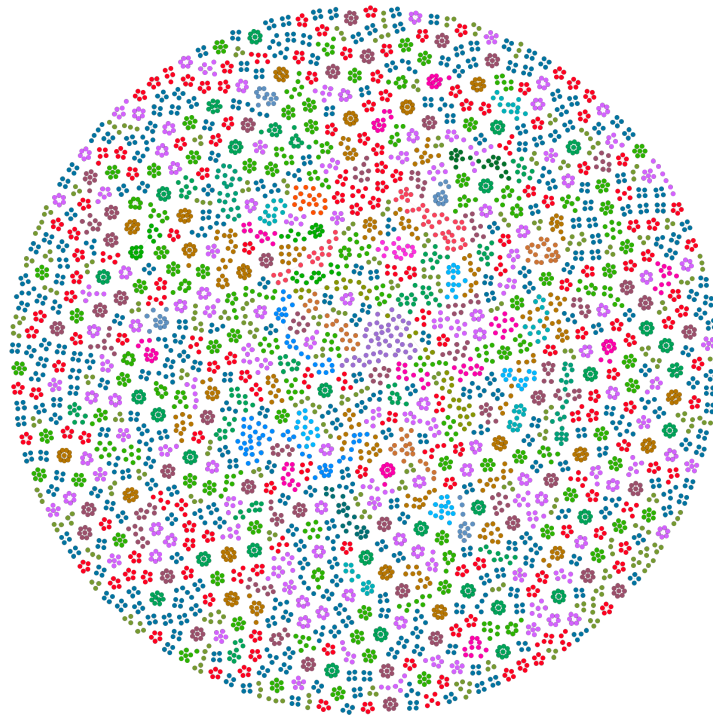


Figura 3.8: Relaciones intermedias 0,1-0,3, cada color refiere a un tamaño de componente específico. Aquí no se aprecia la existencia de una componente gigante y el tamaño de predominante de las componentes ronda por los 6 usuarios. En esta sección se observan a los usuarios ‘normales’ que se relacionan con otros (parte de la misma componente) por algunos puntos en común.

Capítulo 4

Conclusiones

En este trabajo se presentaron las redes de co–eventos con los pesos w_{ij} y w_{ij}^* , además de una serie de análisis hechos sobre las mismas con el objetivo de extraer de ellas un conjunto de usuarios que estén vinculados por razones políticas, laborales, familiares, intelectuales y no sólo vinculados de forma circunstancial.

Twitter da la posibilidad de automatizar cuentas, esto puede darse para una variedad muy amplia de propósitos, siempre y cuando no se violen sus lineamientos. Sumado a esto, sabemos que las plataformas virtuales de redes sociales juegan un papel muy importante a la hora de influir en el comportamiento colectivo, muestras de ello se han visto en campañas políticas, comerciales, cuestiones relacionadas a la salud, entre otras, por ello es importante tener un mejor entendimiento del rol que juegan ciertas dinámicas, así como los mecanismos que se emplean para influir en la conversación pública. Como se puede ver en este trabajo, el empleo de algunos métodos estadísticos han sido de utilidad para la identificación de comportamientos anómalos.

Como pudimos observar, la red de co–eventos con los pesos w_{ij} es de utilidad para identificar los bloques de comunidades que conforman la red general, los cuales pueden ser identificados al realizar mediciones sobre la proporción de enlaces internos y externos o aplicando algoritmos de visualización de tipo repulsivos como el *ForceAtlas2* (Jacomy et al., 2014). Aquí identificar las agrupaciones por ideologías a fines es muy sencillo. Para el caso específico de la red de co–retweet con los pesos w_{ij} podemos realizar la identificación de bots o usuarios super–propagadores de información.

Por otro lado, la red de co–eventos con los pesos w_{ij}^* permite tener in-

formación sobre el grado de semejanza entre los usuarios, esta ponderación contempla todas las acciones que tuvieron cada uno y en cuantas coincidieron ambos. Esto nos permitió lograr identificar 3 variedades de comportamientos en los usuarios de Twitter, los usuarios ‘normales’, los usuarios altamente coordinados y los bots o propagadores de noticias. En la sección de alta coordinación tenemos que, con la red de co-retweet no logramos identificar relaciones de mediano o largo plazo sin embargo, para la red de co-hashtag sí encontramos dichas relaciones, la mayoría de ellas fueron relaciones laborales. Nuestra interpretación a esto, es que al menos para la coordinación laboral el uso de *hashtags* es de mayor utilidad para posicionar la marca o noticia que darle popularidad a un mensaje en específico.

Los bots o usuarios super-propagadores de información que se encontraron en la red de co-eventos con pesos w_{ij}^* , fueron identificados únicamente en el nivel más bajo de pesos $w_{ij}^* < 0,1$. Probablemente esto se deba a que este tipo de usuarios, tienen como propósito difundir información, no interactuar con el resto de la población, o a mostrar rasgos de personalidad.

Hemos identificado algunas limitaciones en este trabajo, entre ellas tenemos que las redes de co-eventos con los pesos w_{ij} no son de utilidad para identificar relaciones de mediano, largo plazo en subredes. Esto se lo adjudicamos a que las conexiones están pesadas por el grado de coincidencias mutuas, pero esto no nos proporciona información sobre la variedad de actividad que ejerció cada usuario.

Otra limitante que tenemos es que las redes de co-eventos construidas, son independientes entre sí y en este trabajo no realizamos análisis para ver cómo varían en el tiempo, esto no lo implementamos por limitaciones de tiempo sin embargo, proponemos como trabajo a futuro un acercamiento, haciendo uso de redes multicapa para construir redes de co-eventos a varias capas y analizar cómo varían en el tiempo. Con esto esperamos obtener una mejor segregación de redes sociales de mediano y largo plazo.

En suma se ha mostrado que, la construcción de las redes de co-eventos con sus variaciones respectivas sirvieron para identificar relaciones sociales de mediano y largo plazo, específicamente relaciones laborales entre grupos de usuarios y relaciones por ideales políticos. Es importante resaltar que los métodos aquí propuestos no hacen uso del procesamiento de lenguaje natural, ni en técnicas de aprendizaje automático o aprendizaje profundo, lo cual permite poder trasladarlo a cualquier otro lenguaje o cualquier otro tema, sin la necesidad de tener que realizar etiquetado de datos ni entrenamiento de alguna red neuronal.

Apéndice A

Twitter y características

Twitter es una red social que nació con la idea de comunicarse por medio de mensajes muy cortos, en sus inicios tenía el límite de 140 caracteres, actualmente ya se extendió hasta 280. La red social de Twitter funciona más o menos de la siguiente forma: cada usuario que ingresa a esta red comienza a seguir a los usuarios que desea, algo muy importante en esta red es que es una red dirigida, por lo tanto a diferencia de Facebook y otras redes donde cuando sigues a un usuario este también te sigue, aquí la relación puede ser unidireccional. Cada usuario puede subir cuantos “tweets” (mensajes con las restricciones mencionadas) desee, estos tweets pueden ser visto por cualquier otro usuario de Twitter, pues a menos que el usuario coloque restricciones de privacidad, por defecto cada tweet emitido es público. Sin embargo Twitter manda notificaciones y coloca con un orden de prioridad los tweets emitidos por los usuarios a los que sigues y a pesar de que los tweets de los usuarios que uno no sigue son públicos, para llegar a ellos uno tendría que realizar una búsqueda un poco más específica. Por lo tanto, si la idea es que un mensaje sea visto por mucha gente, lo mejor sería que quien lo emitiera fuera una persona a la cual lo siguen muchos usuarios. De tal forma que esta red social tendrá características como las redes de conexión preferencial.

Cada usuario tiene acceso al “Time line”, la cual es una sección donde el usuario puede ver en orden cronológico la serie de eventos que han acontecido dentro de su propia red (los usuarios a los que sigue), en esta red aparecerán todos los eventos donde los usuarios a los que sigue hayan interactuado de alguna forma. Las formas de interactuar en Twitter son las siguientes:

like Todo usuario puede emitir un “like” sobre cualquier tweet emitido y esto es básicamente como un voto de aprobación hacia el mensaje

retweet Todo usuario puede tomar un tweet y emitirlo de nuevo sin ninguna

alteración, el tweet saldrá sin ningún cambio, por lo tanto seguirá diciendo quién fue su creador, la hora que se creó y los demás datos propios del tweet. Un retweet básicamente sirve para promocionar un tweet de forma altruista, pues no queda ninguna intervención en el mensaje y uno simplemente ayuda a propagarlo a su red de contactos.

quote Para este último uno toma un tweet existente y le añade un comentario, al final queda el tweet original envuelto por el comentario en cuestión, de tal forma que ahora se propagará tanto el mensaje original como el añadido en el comentario.

reply Los usuarios pueden emitir opiniones sobre los tweets generados sin la necesidad de usar el “quote”, simplemente en una subsección de replicas que tiene por defecto cada tweet. La diferencia del “reply” con el “quote” es que en el “quote” uno vuelve a difundir el mensaje original enmascarado con nuestra intervención y en un “reply” uno sólo emite el comentario y no se enmascara ni redifunde el mensaje original, es un mensaje un poco más directo sin intención de que tus seguidores promocionen tu intervención.

mention En cada tweet se puede etiquetar a uno o más usuarios, de forma tal que esto emitirá una notificación a los usuarios etiquetados para que ellos puedan interactuar con el tweet en cuestión.

Apéndice B

Método de Louvain

El método de Louvain es un algoritmo de tipo voraz. Este algoritmo se divide en dos fases que se repiten de forma iterativa. Tiene como supuesto que iniciamos con una red pesada de N nodos. Primero, nosotros asignamos una comunidad diferente a cada nodo de la red. Así que, en esta partición inicial existen tantas comunidades como nodos en la red. Ahora, para cada node i consideramos los vecinos j de i y evaluamos la ganancia de modularidad que obtendríamos de remover i de su comunidad y colocarlo en la comunidad de j . El nodo i es colocado en la comunidad para la cual su ganancia sea máxima, siempre y cuando la ganancia sea positiva. En caso de no ser positiva la ganancia, i permanecerá en su comunidad original. Este proceso es aplicado repetidas veces sobre cada uno de los nodos hasta que ya no se pueda obtener alguna mejora, es en este momento en que la fase uno es completada. La segunda fase del algoritmo consiste en la construcción de una nueva red donde sus nodos ahora se encuentran en las comunidades encontradas durante la fase previa. Para realizar esto, los pesos de los enlaces entre los nuevos nodos están dados por la suma de los pesos entre los enlaces en las dos comunidades correspondientes. Los enlaces entre nodos de la misma comunidad conducen a auto-bucles hacia esta comunidad en la nueva red. Una vez completada la segunda fase, es posible repetir la primera fase sobre la nueva red pesada generada por la fase dos. Denotaremos como "paso" a la combinación de la fase uno con la fase dos. EL número de comunidades irá decreciendo en cada paso, los pasos son iterados hasta que no existan más cambios y un máximo de modularidad sea alcanzado. para más detalles es recomendable leer el artículo original de Blondel et al. (2008).

Apéndice C

Identificación de usuarios

Con la finalidad de que el lector pueda observar algunos ejemplos de los usuarios identificados como bots, usuarios con una relación de afinidad política y aquellos identificados con una relación de tipo laboral, mostramos en este apéndice algunos perfiles de usuario, así como una breve explicación de cómo se realizó el proceso de identificación.

Es importante aclarar que en este trabajo, primero le dimos al algoritmo de Louvain las redes y una vez obtenidas las comunidades, procedimos a encontrar la relación que existía entre los usuarios pertenecientes a una misma comunidad. Esto lo hicimos de forma muestral, eligiendo de forma aleatoria 50 usuarios de cada comunidad, para posteriormente ingresar a sus perfiles de usuario e identificar su relación en común.

Para el caso particular de las relaciones laborales (para estas no se necesitó realizar una identificación de comunidades), realizamos una muestra única de 50 componentes sobre la red resultante al aplicar el filtro de $w_{ij}^* \geq 0,5$ sobre la red de co-hashtags y con ellos realizar la inspección sobre sus perfiles de usuario.

C.1. Bots y usuarios relacionados por afinidad política

Como lo vimos en la Fig. 3.1, esta red se subdivide básicamente en tres cúmulos principales de usuarios, los cuales están relacionados por tener una afinidad política (para el caso de los identificados como pro-gobierno o los identificados como opositores) o los que se relacionan por compartir noticias acerca del covid. En esta sección mostramos algunos ejemplos de usuarios

pertenecientes al grupo que comparten noticias sobre el covid, los pertenecientes al grupo de los *pro-gobierno* y al grupo de los *opositores*. En las Figs. C.1 y C.2 se muestran al usuario CoronaUpdateBot ¹ y world_news_eng ² respectivamente, ambos pertenecientes al grupo de usuarios relacionados por compartir noticias referentes al COVID.



Figura C.1: CoronaUpdateBot es un usuario de tipo bot, cuyo propósito es retweetear noticias relacionadas al coronavirus.



Figura C.2: world_news_eng es un usuario de tipo bot. Actualmente es una cuenta suspendida por violar las reglas de Twitter

En el caso del usuario BillEsteem no mostramos una imagen de su cuenta pues este usuario realizó un cambio a su nombre de usuario. Dentro de

¹La cuenta CoronaUpdateBot fue diseñada por el usuario Plutonicindia y este ha sido suspendido por violar las reglas de Twitter.

²La cuenta del usuario world_news_eng actualmente ha sido suspendida por violar las reglas de Twitter.

las políticas de Twitter, el nombre de usuario ³ (también es conocido como handle) debe ser único, cada usuario tiene la posibilidad de cambiar su nombre de usuario a cualquier otro que no esté siendo usado en ese momento. Estos tres usuarios anteriormente mencionados, aparecen como los nodos más grandes en el cúmulo rosa (en el panel (A)) de la Fig.3.1, también están identificados con estrellas de color rosa en la Fig.3.2 y en la red de co-eventos con pesos w_{ij}^* se encuentran en el segmento de $w_{ij}^* \leq 0,1$.

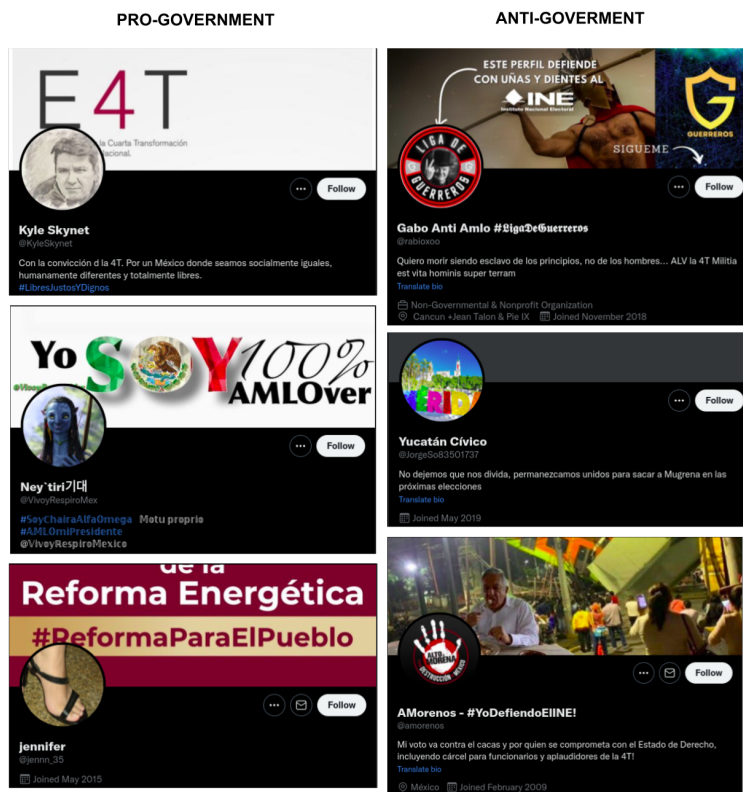


Figura C.3: Muestra de cuentas relacionadas por cuestiones políticas, del lado izquierdo se muestran aquellas con una afinidad pro-gobierno y del lado derecho las identificadas como opositores.

En la Fig.C.3 se muestran tres ejemplos de usuarios pertenecientes al grupo con afinidad al gobierno, conocidos como los pro-gobierno y tres ejem-

³<https://help.twitter.com/en/managing-your-account/change-twitter-handle>

plos de los usuarios pertenecientes al grupo de los opositores. Estos usuarios fueron obtenidos de la red de co-retweet para w_{ij} , también mostrada en la Fig.3.1.

C.2. Usuarios con relación laboral

En la Fig. C.4, se muestran algunos ejemplos de relaciones humanas de tipo laboral, encontradas al filtrar los enlaces de la red de co-hashtag $w_{ij}^* \geq 0,5$.

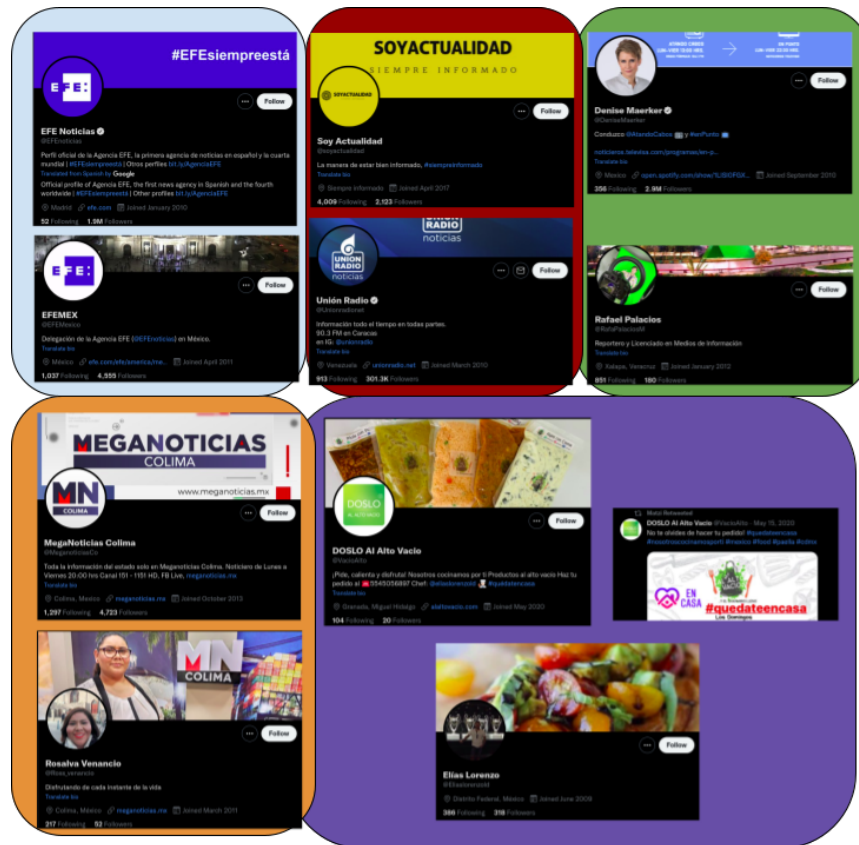


Figura C.4: Muestra de cuentas relacionadas por cuestiones laborales, donde las cuentas con un mismo color de fondo indican una relación.

Bibliografía

- Amor, B. R., Vuik, S. I., Callahan, R., Darzi, A., Yaliraki, S. N., and Barahona, M. (2016). Community detection and role identification in directed networks: Understanding the twitter network of the care.data debate. *Dynamic Networks and Cyber-Security*, 1:111–136.
- Backstrom, L., Boldi, P., Rosa, M., Ugander, J., and Vigna, S. (2012). Four degrees of separation. *Proceedings of the 4th Annual ACM Web Science Conference, WebSci'12*, volume:33–42.
- Barabási, A.-L. (2013). Network science. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 371(1987):20120375.
- Barabasi, A.-L. and Albert, R. (1999). Albert, r.: Emergence of scaling in random networks. *science* 286, 509-512. *Science (New York, N.Y.)*, 286:509–12.
- Barabási, A.-L. and Pósfai, M. (2016). *Network science*. Cambridge University Press, Cambridge.
- Blondel, V. D., Guillaume, J.-L., Lambiotte, R., and Lefebvre, E. (2008). Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10):P10008.
- Caetano, J. A., Lima, H. S., Santos, M. F., and Marques-Neto, H. T. (2018). Using sentiment analysis to define twitter political users' classes and their homophily during the 2016 American presidential election. *Journal of Internet Services and Applications*, 9(1).
- Chamberlain, B., Humby, C., and Deisenroth, M. P. (2016). *Real-Time Community Detection in Large Social Networks on a Laptop*. Number 1.

- Darmon, D., Omodei, E., and Garland, J. (2015). Followers are not enough: A multifaceted approach to community detection in online social networks. *PLoS ONE*, 10(8):1–20.
- Dutta, H. S., Dutta, V. R., Adhikary, A., and Chakraborty, T. (2020). Hawkeseye: Detecting fake retweeters using hawkes process and topic modeling. *IEEE Transactions on Information Forensics and Security*, 15:2667–2678.
- Erdős, P. and Rényi, A. (1959). On random graphs i. *Publicationes Mathematicae Debrecen*, 6:290.
- Erdős, P. and Rényi, A. (1960). On the evolution of random graphs. *Publ. Math. Inst. Hungary. Acad. Sci.*, 5:17–61.
- Gilbert, E. N. (1959). Random graphs. *The Annals of Mathematical Statistics*, 30(4):1141–1144.
- Girvan, M. and Newman, M. E. (2002). Community structure in social and biological networks. *Proceedings of the National Academy of Sciences of the United States of America*, 99(12):7821–7826.
- Granovetter, M. S. (1973). The strength of weak ties. *American Journal of Sociology*, 78(6):1360–1380.
- Harmon, J. (1959). *The Psychology of Interpersonal Relations*. By Fritz Heider. New York: John Wiley and Sons, Inc., 1958. 322 pp. *Social Forces*, 37(3):272–273.
- Iñiguez, G., Kertész, J., Kaski, K. K., and Barrio, R. A. (2009). Opinion and community formation in coevolving networks. *Phys. Rev. E*, 80:066119.
- Jacomy, M., Venturini, T., Heymann, S., and Bastian, M. (2014). ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. *PLoS ONE*, 9(6):1–12.
- Keller, F. B., Schoch, D., Stier, S., and Yang, J. (2020). Political astroturfing on twitter: How to coordinate a disinformation campaign. *Political Communication*, 37(2):256–280.
- Kim, Y. H., Seo, S., Ha, Y. H., Lim, S., and Yoon, Y. (2013). Two applications of clustering techniques to twitter: Community detection and issue extraction. *Discrete Dynamics in Nature and Society*, 2013.

- Lai, M., Hernández Farías, D. I., Patti, V., and Rosso, P. (2017). Friends and enemies of clinton and trump: Using context for detecting stance in political tweets. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 10061 LNAI:155–168.
- Levy, G. and Razin, R. (2019). Echo Chambers and Their Effects on Economic and Political Outcomes. *Annual Review of Economics*, 11:303–328.
- Lorenz, E. N. (1963). Deterministic Nonperiodic Flow. *Journal of Atmospheric Sciences*, 20(2):130–141.
- Miller, M., Lynn, S.-L., and James, M. C. (2001). Birds of a Feather: Homophily in Social Networks. *Annual Review of Sociology*, 27:415–444.
- Mirko Lai, Delia Irazú Hernández Farías, V. P. and Rosso, P. (2017). *Friends and Enemies of Clinton and Trump: Using Context for Detecting Stance in Political Tweets*. Book Series: Lecture Notes in Computer Science. Springer International Publishing.
- Newman, M. (2010). *Networks: An Introduction*. Oxford University Press, Inc., USA.
- Newman, M. E. J. (2006). Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23):8577–8582.
- Newton, I. (1726). *Philosophiæ naturalis principia mathematica*. Editio tertia, aucta emendata. Londini : Apud G. J. Innys, 1726.
- Onnela, J. P., Saramäki, J., Hyvönen, J., Szabó, G., De Menezes, M. A., Kaski, K., Barabási, A. L., and Kertész, J. (2007). Analysis of a large-scale weighted network of one-to-one human communication. *New Journal of Physics*, 9.
- Pacheco, D., Hui, P.-M., Torres-Lugo, C., Truong, B. T., Flammini, A., and Menczer, F. (2021). Uncovering coordinated networks on social media. In *Proc. AAAI International Conference on Web and Social Media (ICWSM)*. Forthcoming. Preprint arXiv:2001.05658.
- Pariser, E. (2011). *The Filter Bubble: What the Internet Is Hiding from You*. Penguin Group, The.
- Sasahara, K., Chen, W., Peng, H., Ciampaglia, G. L., Flammini, A., and Menczer, F. (2020). Social influence and unfollowing accelerate the emergence of echo chambers. *Journal of Computational Social Science*.

- Schlosser, S., Toninelli, D., and Cameletti, M. (2021). Comparing methods to collect and geolocate tweets in Great Britain. *Journal of Open Innovation: Technology, Market, and Complexity*, 7(1):1–20.
- Sloan, L. and Morgan, J. (2015). Who tweets with their location? Understanding the relationship between demographic characteristics and the use of geoservices and geotagging on twitter. *PLoS ONE*, 10(11):1–15.
- Thurner Hanel, Rudolf A., Klimek, Peter, Oxford University Press., S. (2019). *Introduction to the theory of complex systems*. Oxford University Press, Oxford.
- Torres-Lugo, C., Yang, K.-C., and Menczer, F. (2020). The Manufacture of Political Echo Chambers by Follow Train Abuse on Twitter.
- Travers, J. and Milgram, S. (1969). An Experimental Study of the Small World Problem. *Sociometry*, 32(4):425–443.
- Wang, A. H. (2010). Detecting spam bots in online social networking sites: A machine learning approach. In Foresti, S. and Jajodia, S., editors, *Data and Applications Security and Privacy XXIV*, pages 335–342, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Watts, D. J. and Strogatz, S. H. (1998). Collective dynamics of ‘small-world’ networks. *Nature*, 393(6684):440–442.
- Yang, K.-C., Ferrara, E., and Menczer, F. (2022). Botometer 101: Social bot practicum for computational social scientists. pages 1–14.
- Yuan, X. and Crooks, A. T. (2018). Examining online vaccination discussion and communities in Twitter. *ACM International Conference Proceeding Series*, pages 197–206.
- Zachary, W. W. (1977). An Information Flow Model for Conflict and Fission in Small Groups. *Journal of Anthropological Research*, 33(4):452–473.
- Zohar, M. (2021). Geolocating tweets via spatial inspection of information inferred from tweet meta-fields. *International Journal of Applied Earth Observation and Geoinformation*, 105:102593.