

Caracterización de Redes y Topologías Biológicas

Resumen

En esta asignatura se estudian los principales tipos de conectividad que se pueden dar en una red biológica. Se describen además cuales puede ser la mejor estrategia de conexión entre los elementos de una red sujetos a una determinada dinámica. También se proporcionan métodos para calcular los principales parámetros topológicos y de rendimiento de una red dada. Además se estudian redes resistentes a una determinada estrategia de ataque o frente a errores en la red.

Índice general

I	Introducción y descripción de algunas redes reales	2
I.1	Qué es una red	2
I.2	Algunos ejemplos de redes y algunas de sus propiedades	2
I.2.1	World Wide Web	2
I.2.2	Internet	3
I.2.3	Red de actores	4
I.2.4	Red de colaboración científica	4
I.2.5	Red de contactos sexuales	4
I.2.6	Red de llamadas telefónicas	4
I.2.7	Redes lingüísticas	5
I.2.8	Redes eléctricas	5
I.3	Algunos ejemplos de redes biológicas y algunas de sus propiedades	5
I.3.1	Redes de ecología	5
I.3.2	Redes celulares	5
I.3.3	Redes neuronales	6
I.3.4	Redes de interacción de proteínas	6
I.3.5	Redes genéticas	6
II	Teoría de grafos y métricas	7
II.1	Introducción a la teoría de grafos	7
II.2	Bucles y ramas paralelas	9
II.3	Grafos dirigidos y ponderados	9
II.4	Grado de un nodo	10
II.5	Subgrafos	11
II.6	Paseos, caminos, circuitos y ciclos	11
II.7	Medidas de centralidad, betweeness y closeness	12
II.8	Conexidad	13
II.9	Bosques y árboles	14

Capítulo I

Introducción y descripción de algunas redes reales

Aunque lo vayamos a utilizar como sinónimos, un grafo y una red no es lo mismo; el grafo es la representación matemática de la red. En una red aleatoria, no hay que medir nada; si una red biológica sale aleatoria, se ha medido mal. Las redes biológicas son todas de mundo pequeño. Además, casi todas son libres de escala.

I.1. Qué es una red

Una red es un conjunto de elementos (personas, ciudades, proteínas, especies animales, productos químicos, etc) de las cuales algunas están conectadas con otras y otras no. Se puede representar en bolas que se unen con líneas con otras líneas. Las bolitas se denominan como nodos.

Las redes se estudian con NetworkX y Cytoscape.

I.2. Algunos ejemplos de redes y algunas de sus propiedades

I.2.1. World Wide Web

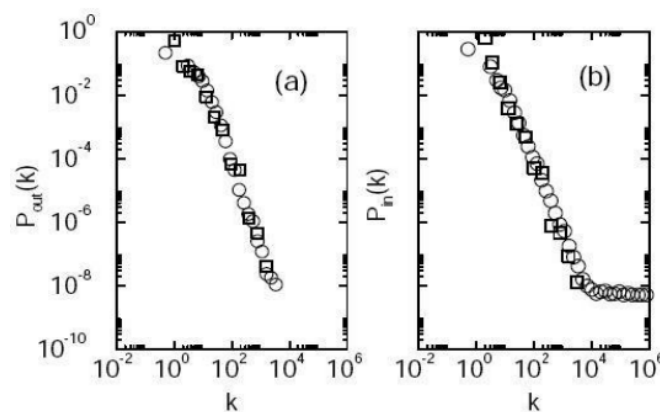
La World Wide Web es la mayor red para la cual existe información topológica. Los nodos de la red son los documentos, y las ramas son los enlaces (hyperlinks) entre documentos. El tamaño actual de esta red es de más de 1000 millones de nodos. Esta red es dirigida: la página A apunta a la página B, pero sin que la página B apunte a la página A. El grupo CAIDA se dedica a analizar la red. Esta red es enorme, pudiendo dibujar solo a un nivel muy alto.

La distribución del grado de las páginas web tiene una distribución libre de escala tanto en los enlaces de salida como en los enlaces de entrada. Esto es una distribución de probabilidad. Por ejemplo, en el queso de Gruyere, los agujeros son de distinto tamaño, los cuales tienen una distribución de tamaño. Se llama libre de escala porque

se pueden encontrar diez veces más los agujeros de un tamaño mayor y diez veces menos los agujeros de tamaño pequeño. Así, no hay una escala fija de la distribución (no se puede representar con ninguna escala, ni logarítmica ni nada). Esto con el queso manchego no pasa. Si en la WWW vemos cuántas páginas web tienen 100 enlaces de salida, 10, 1000, etc, y se dibuja en escala logarítmica logarítmica, sale una recta. Esto pasa también con los enlaces de entrada.

La distancia entre dos páginas de la WWW es pequeña (entre 11 y 16). Los nodos de la WWW están muy clusterizados.

La cola de la derecha parece que rompe la recta. Las redes libres de escala se pueden producir por muchas razones, pero al utilizar un proceso evolutivo en el que cada tiempo se generan nuevos nodos y tienen mayor preferencia para conectarse a otros nodos, e incluso pueden desaparecer algunos nodos antiguos. Esto produce las colas residuales.



Las redes libres de escala son muy resistentes a ataques aleatorios (fallos en la red) en cuanto a la conectividad, por lo que hay una razón evolutiva por la que las redes biológicas son libres de escala. La red regulatoria de P53 está muy estudiada y caracterizada. Uno de los elementos más importantes es MDM5.

1.2.2. Internet

Internet es la red de enlaces físicos entre ordenadores u otros servicios de comunicación. La topología de internet se suele estudiar a dos niveles: Enrutadores y Sistemas autónomos. Los enrutadores son las máquinas que mandan los "paquetes" a otros enrutadores. Hay algoritmos de enrutación que deciden hacia dónde enviar las cosas. Los sistemas autónomos son conjuntos de máquinas que organizan y gestionan otras máquinas.

Para ambos tipos de red (enrutadores y sistemas autónomos) el grado de cada nodo seguía una distribución libre de escala. De nuevo la red está altamente clusterizada (coeficiente de clustering entre 0,18 y 0,3) y los caminos entre nodos son cortos (aproximadamente 9).

El **índice de clusterización** es una medida de la probabilidad de que los dos vecinos de un nodo sean vecinos entre sí, favoreciendo la creación de triángulos. Es decir, en redes sociales, que mis amigos también sean amigos entre sí. En biología, si dos

proteínas son expresadas por una tercera proteína, las dos mantienen una relación entre sí (aunque puede no pasar). Los vecinos de un mismo nodo tienen una probabilidad alta de ser vecinos entre sí. En una red aleatoria, los vecinos de un nodo dependen de la probabilidad de rama de que esos nodos también sean vecinos entre sí (como cualquier otro).

La métrica de caminos cortos o largos se hace en comparación con el grafo aleatorio con el mismo número de nodos y ramas. En biología, los caminos también suelen ser cortos, y si son largos se puede deber a una enfermedad o patología.

1.2.3. Red de actores

Los nodos son actores, y dos de ellos están conectados si han participado juntos en alguna película. Actualmente, la red consta de unos 450.000 actores. La distancia media entre actores es 3,65. La red está altamente clusterizada (100 veces más que un grafo aleatorio). La distribución de grados sigue una ley de potencias (libre de escala).

1.2.4. Red de colaboración científica

Los nodos están constituidos por científicos. Dos nodos están conectados si alguna vez publicaron un trabajo en común. La red de nuevo presenta una distribución libre de escala, caminos cortos entre los nodos y una alta clusterización.

El **centro de la red** es el nodo que está a una menor distancia promedio del resto de nodos de la red. Este centro lo tiene un científico húngaro llamado Paul Erdős que trabajaba en teoría de grafos.

Para una red de citaciones científicas, los nodos de la red son artículos científicos. Las ramas son citaciones entre artículos. Se tiene una base de datos de unos 750.000 artículos. Tanto los grados de entrada como los de salida siguen una distribución libre de escala.

1.2.5. Red de contactos sexuales

Los nodos y las ramas tienen una definición obvia. Tiene interés por la difusión de enfermedades (especialmente aquellas de transmisión sexual como el SIDA). Presenta una distribución libre de escala. Se sospecha que los datos de esta red no son totalmente fiables (es defectuosa al tener muchos datos falsos). Entre un 10-15 % es falsa.

Se define como k-core un grafo no dirigido creado a partir de un grafo más grande en el que se crean jerarquías o grupos en el que los nodos están separados por k vértices.

1.2.6. Red de llamadas telefónicas

Los nodos son números de teléfono. Las ramas son llamadas de larga distancia entre nodos. De nuevo la red presenta una distribución libre de escala.

I.2.7. Redes lingüísticas

Los nodos son palabras. Dos nodos están conectados si están juntas en alguna frase y hay solamente una palabra entre ambas. Un estudio realizado en inglés sobre 440.902 palabras presentó una distancia media de 2,62 y un índice de clusterización de 0,43.

Otra red lingüística considera de nuevo los nodos como palabras. Dos nodos están conectados si se considera que ambas palabras son sinónimas (de acuerdo con el Merrian Webster Dictionary). El camino medio es de 4,7, el índice de clusterización es de 0,7 y los nodos presentan una distribución libre de escala.

En la red semántica, cada nodo es un objeto o un concepto. Dos nodos se relacionan entre sí, si existe una relación de la forma "es un" o "tiene un" entre ambos nodos. Se ha estudiado poco, pero parece presentar un camino medio corto, alta clusterización y una distribución de nodos libre de escala.

I.2.8. Redes eléctricas

La red eléctrica del Oeste de los Estados Unidos está compuesta por nodos (generadores, transformadores y subestaciones) y ramas (cables físicos entre nodos). La red tiene 4.941 nodos y un grado medio por nodo de 2,41. Esta red se aparta del patrón habitual teniendo una estructura muy jerárquica y en forma de estrella. Esto hace que sea muy frágil y condicionada a cuestiones económicas y políticas. Ocurre de forma similar con las redes de internet. No se utiliza el camino más rápido o corto, si no el camino más barato (como a la hora de buscar vuelos).

I.3. Algunos ejemplos de redes biológicas y algunas de sus propiedades

I.3.1. Redes de ecología

En las redes alimentarias, los nodos de la red son especies, y las ramas relaciones predador-presa entre especies. Las distancias son cortas entre los elementos de la red. En general, son redes con pocos nodos.

Al ser redes pequeñas es difícil dibujar la distribución del grado de los nodos. Parecen presentar una distribución libre de escala, con un exponente inusualmente pequeño. Esta red es dirigida (aunque pueda haber dobles ramas).

I.3.2. Redes celulares

Se presentan al estudiar el metabolismo de organismos. Los nodos son sustratos químicos (ATP, ADP, etc), y las ramas presentan reacciones químicas entre los sustratos. Esta red va de arriba a abajo, empezando con unos productos de entrada de la célula y terminando con productos de salida que la célula no puede descomponer más.

I.3.3. Redes neuronales

Cada nodo es una neurona (biológica o artificial), y las ramas son conexiones sinápticas entre neuronas. La primera red estudiada de este tipo es la del gusano *Caenorhabditis elegans*, del cual se tiene el mapa neuronal completo.

Las redes neuronales artificiales están ahora en auge para las inteligencias artificiales al utilizarse para el aprendizaje profundo.

I.3.4. Redes de interacción de proteínas

Cada nodo es una proteína. Las ramas representan relaciones de expresión entre las proteínas. Una de las redes más importantes es la red p53 de control de crecimiento del cáncer. Un paper muy bueno es [Surfing the p53 network \(DOI 10.1038/35042675\)](https://doi.org/10.1038/35042675).

Esta es la red en la que más se trabaja en biología. Se buscan los efectos entre los nodos (aumenta la expresión, inhibe), los componentes clave, los parámetros, etc.

I.3.5. Redes genéticas

Cada nodo es expresión genética (nucleótidos). Las ramas conectan los nucleótidos que presentan un alto índice de similitud entre ambas. Una vez representada la red, se buscan familias o grupos de genes similares. Hay que diferenciar identidad con similitud (sobre todo con desajuste de fase). Se utiliza programación dinámica para calcular la mayor longitud de subsecuencia idéntica, como por ejemplo con el algoritmo Soldier's Walk.

Si clusterizamos y obtenemos 2 cluster, cada cluster indica un gen con errores, o dos individuos distintos. Luego hay que interpretar por qué hay ese número de cluster. Normalmente hay muchos clusters que se quieren clasificar, y en cada cluster suele aparecer el mismo gen que se ha mutado.

Las máquinas de microarrays ahora dan un conjunto de nucleótidos muy grandes, pero antes se obtenían fragmentos que había que unir. Para ello, se debían utilizar algoritmos sobre grafos para calcular cadenas largas a partir de las cadenas cortas, pero ahora ya no se usa por las mejoras tecnológicas.

Capítulo II

Teoría de grafos y métricas

II.1. Introducción a la teoría de grafos

La teoría de grafos ha sido utilizada recientemente para:

- Clasificación automática de secuencias de proteínas.
- Detección de jerarquías de proteínas.
- Análisis de redes genéticas.
- Reconstrucción de redes genéticas grandes obtenidas mediante modificación de genes.

Un grafo G es un par de conjuntos (V, E) donde $V = \{v_1, v_2, \dots, v_n\}$ es el conjunto de vértices o nodos y $E = \{(v_i, v_j), (v_{i'}, v_{j'}), \dots\}$ es un conjunto de pares no ordenados de elementos de V y se denomina conjunto de ramas del grafo. El número de nodos se denomina **orden** del grafo, y el número de ramas es el **tamaño** del grafo.

Pregunta de test: define orden y tamaño, dado un grafo dar el orden y tamaño, etc.

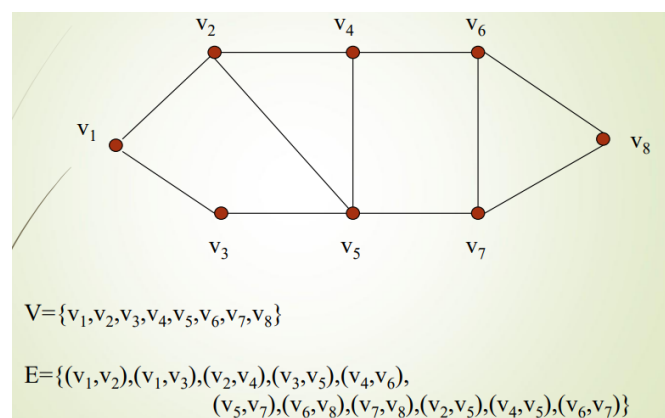
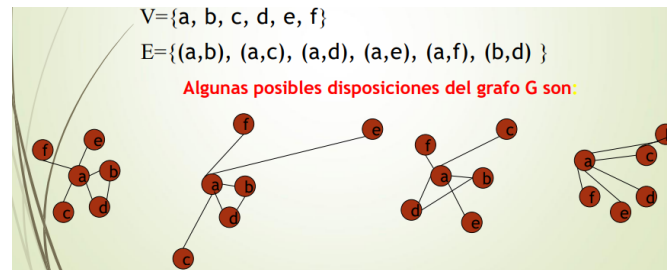


Figura II.1: Ejemplo de grafo de orden 8 y tamaño 11.

Para una red de proteínas, cada proteína sería un nodo del grafo, y una rama indicaría interacción entre ambas proteínas.

Una disposición (layout) es una posible colocación de los nodos y las ramas en un espacio 2D o 3D. Un mismo grafo puede tener múltiples colocaciones. Ejemplo, consideremos el grafo $G=(V,E)$.



Existen programas de ordenador que nos permiten obtener colocaciones predefinidas (Gephy, Pajek). Cuando no se especifica ninguna colocación, se entiende que los nodos se sitúan aleatoriamente sobre el plano o espacio. Algunos de los tipos más habituales de colocaciones son:

- Colocaciones regulares
- Basadas en la física (atracción-repulsión)
- Basadas en propiedades topológicas (jerarquías, número de vecinos, etc)

Un hipergrafo H es un también par de conjuntos (V,E) donde $V = \{v_1, v_2, \dots, v_n\}$ es el conjunto de vértices o nodos y $E = \{(v_{i1}, v_{i2}, \dots), (v_{i'1}, v_{i'2}, \dots), \dots\}$ es una familia de subconjuntos no ordenados de elementos de V . E se denomina conjunto de hiperramas o hiperaristas del hipergrafo. El número de hiperramas $|E|$ se denomina cardinalidad del hipergrafo. El valor $|E| * |V|$ se denomina tamaño o volumen del grafo. Si tenemos un grafo de n nodos, ¿cuántas parejas podemos tener como máximo? $(n \cdot n - 1)/2$ Por tanto, en un grafo con n nodos, ¿cuántas ramas puede tener? Igual, $(n \cdot n - 1)/2$

Pregunta
examen

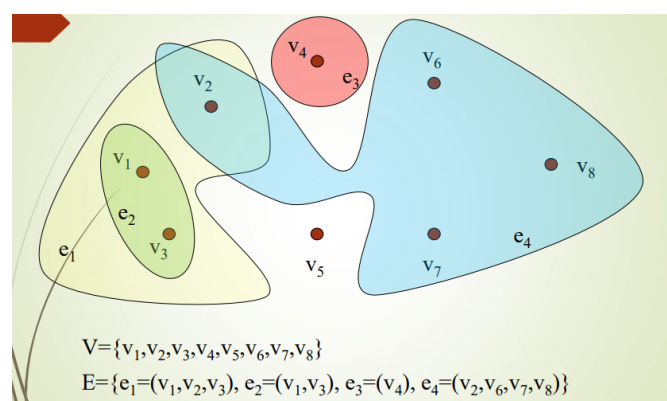


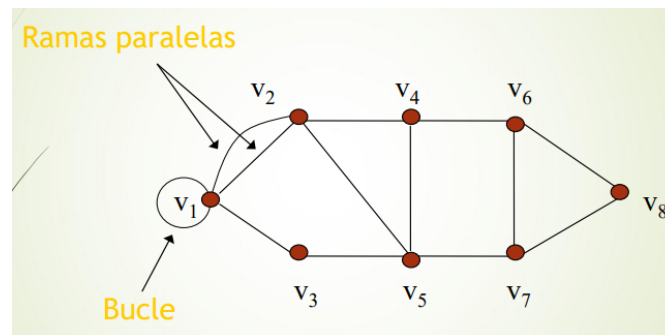
Figura II.2: Ejemplo de hipergrafo de cardinalidad 4 y tamaño 32.

Un hipergrafo H se dice que es **propio** si no es vacío ($V \neq \emptyset$) y no contiene ninguna arista vacía. Un hipergrafo H se dice que tiene **dominio completo** si todos los nodos están en al menos una arista, en caso contrario se dice que tiene **dominio parcial**. Si en un hipergrafo todas las hiperramas tienen el mismo número de nodos, entonces se denomina **hipergrafo k-uniforme**.

Ejercicio: Indicar si el hipergrafo del ejemplo anterior es propio, tiene dominio completo y si es k uniforme. Es propio (el conjunto de vértices tiene 8 elementos y todas las ramas e tienen vértices dentro), es de dominio parcial (v_5 no está en ninguna rama) y no es k -uniforme (e_1 tiene 3 elementos, e_2 tiene 2, e_3 tiene 1 y e_4 tiene 4).

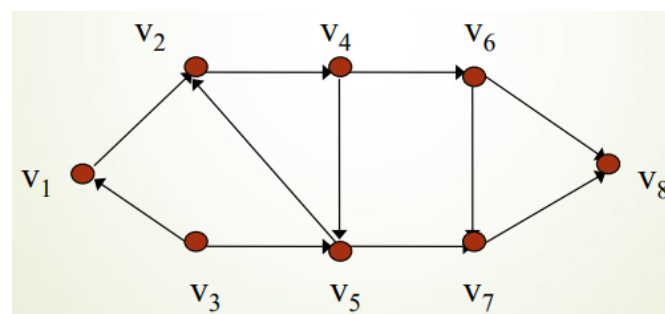
II.2. Bucles y ramas paralelas

Un bucle es una rama que empieza y termina en el mismo nodo (v_i, v_i). Cuando dos ramas conectan el mismo par de vértices se denominan paralelas. Un grafo con bucles se denomina pseudografo. Un grafo con ramas paralelas pero sin bucles se denomina multigrafos. Un grafo sin bucles ni ramas paralelas se denomina grafo simple.

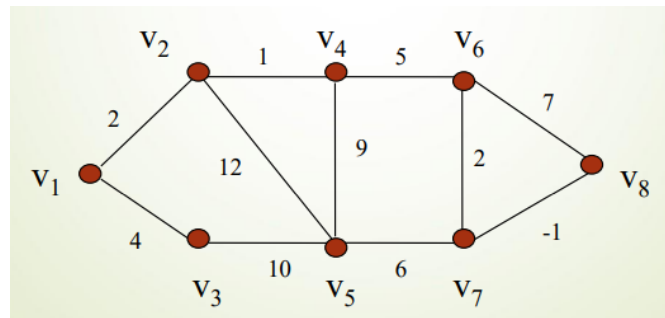


II.3. Grafos dirigidos y ponderados

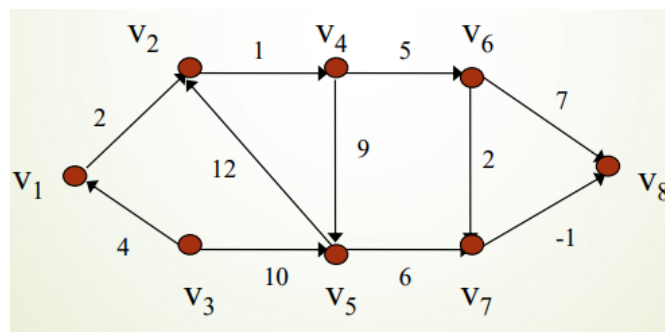
Se puede considerar que los enlaces entre nodos son dirigidos (v_i, v_j) = (v_j, v_i). Los grafos dirigidos se denominan también **digrafos**.



En los grafos ponderados, a cada rama del grafo se le puede asociar un número. El número asociado a cada rama puede indicar entre otras cosas una distancia, una capacidad, un valor temporal, etc.

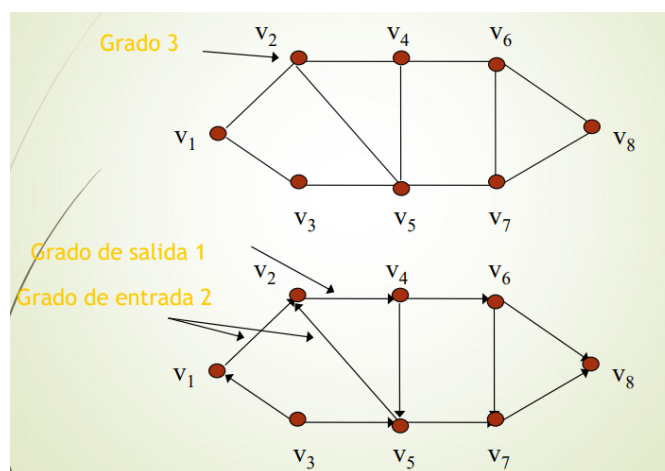


Los grafos dirigidos y ponderados poseen ramas dirigidas a las que se asocia un número.



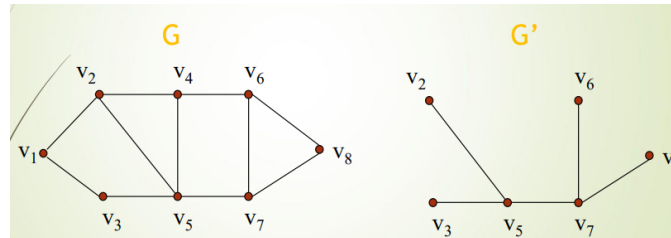
II.4. Grado de un nodo

Dos nodos de un grafo son **vecinos o adyacentes** si existe una rama que los conecta. El **grado** de un nodo es el número de vecinos que tiene dicho nodo. En los grafos dirigidos se calcula el **grado de entrada** y el **grado de salida**. En los grafos ponderados, el grado se puede promediar por el número asociado a las ramas. Un grafo se dice que es **regular** si todos los nodos tienen el mismo grado.



II.5. Subgrafos

Un grafo $G'=(V',E')$ es un subgrafo de un grafo $G=(V,E)$ si V' es un subconjunto de V y E' es un subconjunto de E . En otras palabras, un subgrafo es un trozo de un grafo más grande.



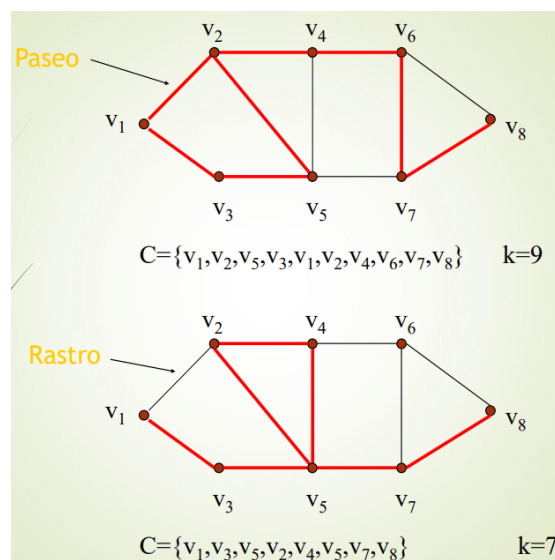
Un subgrafo $G'=(V',E')$ de un grafo $G=(V,E)$ se dice que es **abarcador** si $V=V'$, es decir, si están todos los nodos, pero faltan algunas ramas.

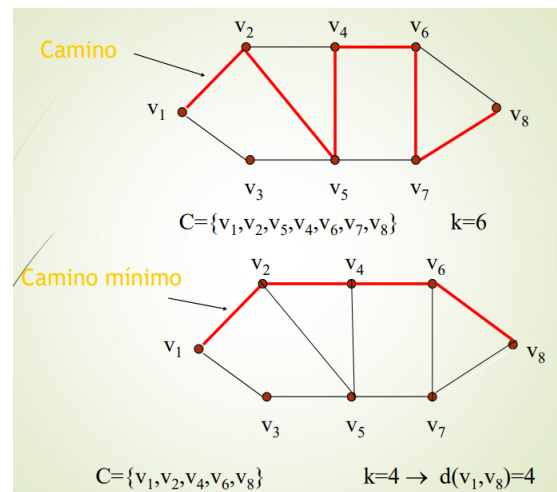
Un grafo es un subgrafo de sí mismo. Además, un grafo vacío es un subgrafo de cualquier grafo.

II.6. Paseos, caminos, circuitos y ciclos

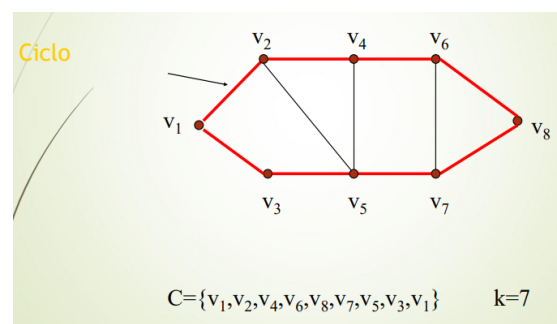
Un **paseo** de un nodo u a un nodo v es una secuencia de vértices $\{v_0, v_1, \dots, v_k\}$ con $v_1 = uv_k = v$ y (v_{i-1}, v_i) rama del grafo. El número de ramas del paseo es su **longitud**. Un paseo en el cual no se repiten ramas se denomina **rastro**. Un paseo en el cual todos los vértices $\{v_0, v_1, \dots, v_k\}$ son distintos se denomina **camino**. Un camino siempre debe ser un rastro y un paseo. Si algo no es rastro, no puede ser camino, y si no es paseo, no puede ser ni rastro ni camino. Cada uno es cada vez más restrictivo.

Entre dos nodos, puede haber varios caminos posibles. Un **camino mínimo** entre dos nodos es aquel de menor longitud de entre todos los posibles caminos entre ambos nodos. La **distancia** entre dos nodos del grafo se define como la longitud de cualquier camino mínimo que los una.





Un **paseo cerrado** es un paseo $\{v_0, v_1, \dots, v_k\}$ tal que $v_0 = v_k$. Un paseo cerrado en el que no se repiten ramas es un **circuito**. Un **ciclo** es un circuito en el que no se repiten vértices. Los ciclos son importantes, porque las redes biológicas tienen ciclos (que suelen ser largos), pero en las redes aleatorias no aparecen ciclos, o éstos son muy pequeños.



El nodo con menor distancia entre los demás es muy importante, denominándose como **centro del grafo**.

Para un grafo con excesivos nodos, los caminos mínimos y las distancias se calculan con un algoritmo. Si el grafo es no ponderado, se utiliza el algoritmo búsqueda en anchura, mientras que si es ponderado, utiliza Dijkstra.

II.7. Medidas de centralidad, betweenness y closeness

Pregunta examen:
Betweenness/Closeness/Farness se define como...

Dado un nodo v_i se define su **betweenness** $C_B(v_i)$ como la fracción de caminos mínimos que hay entre el resto de nodos del grafo y que pasan por el nodo v_i . Es decir, se hacen parejas de todos los nodos del grafo excluyendo el nodo de interés, y se calculan los caminos mínimos. Algunos pasarán por el nodo de interés, que son los que nos quedamos. Con eso se evalúa el cociente (los que pasan por ese nodo entre todos), que será el betweenness (un valor entre 0 y 1). La centralidad de un nodo es muy costosa de calcular, usualmente se emplean algoritmos aproximados.

Dado un nodo v_i se define su **lejanía o farness** $C_F(v_i)$ como la suma de las distancias de v_i al resto de nodos del grafo.

Dado un nodo v_i se define su **cercanía o closeness** $C_C(v_i)$ como la inversa de su lejanía $C_C(v_i) = 1/C_F(v_i)$.

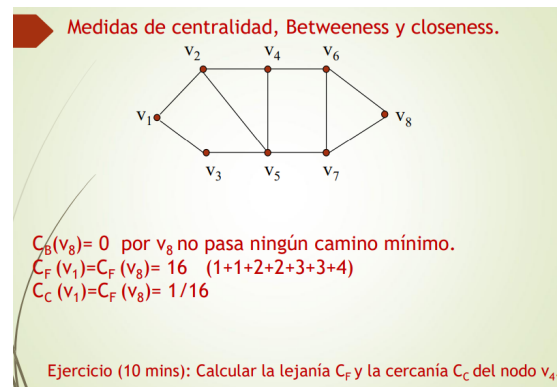


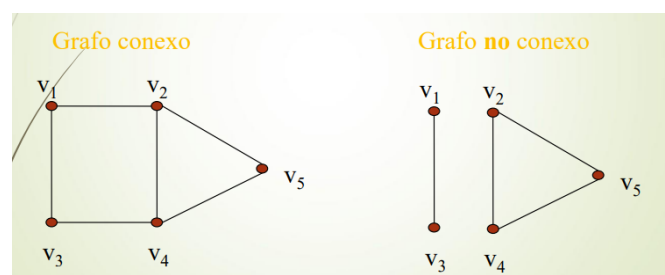
Figura II.3: Respuesta al ejercicio: Cogiendo v_4 , la lejanía será $2+1+2+1+1+2+2 = 11$, y la cercanía $1/11$.

La cercanía y lejanía tiene un problema: su valor numérico depende del orden del grafo. Por tanto, sirve para comparar dentro del mismo grafo, pero no entre grafos. Para eso, habría que normalizar dividiendo por el número total de nodos. A esto se le conoce como **camino característico**.

Pregunta
examen:
Calcular
camino
caracterís-
tico

II.8. Conexidad

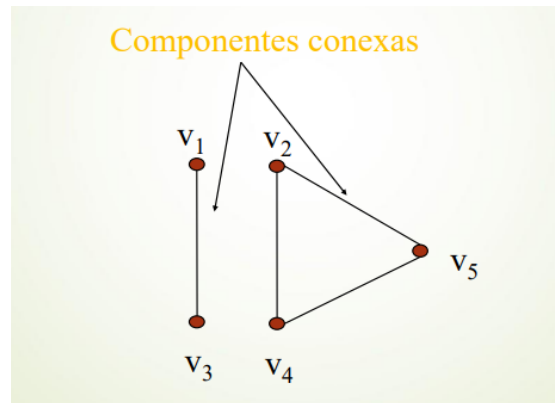
Un grafo es **conexo** si para cada par de nodos del grafo existe al menos un camino que los une. En otras palabras, que no esté separado en distintos trozos.



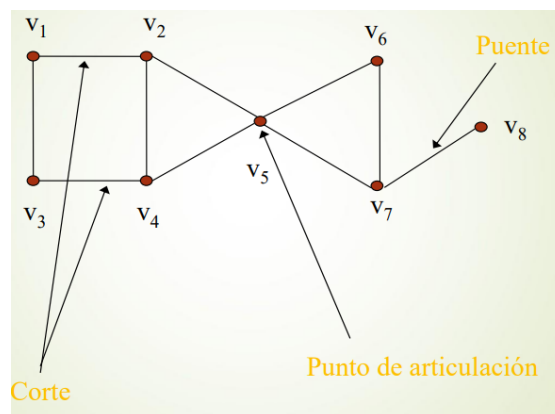
Hay un algoritmo muy rápido y eficiente que calcula si un grafo es conexo o no.

Una **componente conexa** de un grafo es cada uno de los subgrafos maximales conexos. Esto quiere decir que el subgrafo no puede ser más grande, que no se le puede añadir más nodos.

Un **punto de articulación** es un nodo que desconecta un grafo conexo. Un **corte** es un conjunto de ramas que desconecta un grafo conexo. Si un corte está compuesto por una única rama, se denomina **punto de articulación**. Un **corte mínimo** de un grafo es el mínimo número de ramas que al ser eliminadas desconectan el grafo.



El algoritmo CLICK (CLuster Identification via Connectivity Kernels) calcula una aproximación al corte mínimo. Esto lo hacían cogiendo los dos nodos más lejanos. Los puentes suelen ser muy malos para la conexidad de los grafos.

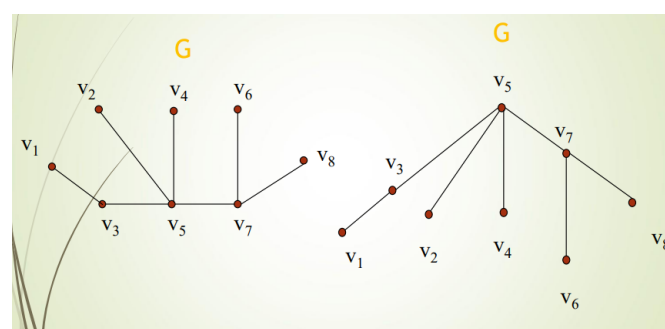


La máxima distancia entre cualquier par de nodos se denomina como diámetro.

El corte mínimo entre dos nodos es siempre mayor que el corte mínimo de todo el grafo.

II.9. Bosques y árboles

Un grafo sin ciclos (acíclico) se denomina bosque. Un árbol es un grafo acíclico conexo. Cada componente conexas de un bosque es un árbol.



Un subgrafo abarcador acíclico de un grafo G se denomina un **bosque abarcador**. Un subgrafo abarcador conexo acíclico de un grafo G se denomina un **árbol abarcador**.

