

Análisis y detección de correlaciones en relevamientos transcripcionales de
gran escala

Tesis de Licenciatura en Ciencias Físicas

Facultad de Ciencias Exactas y Naturales

Universidad de Buenos Aires

Andrés Rabinovich

Marzo 2016

Hoja a completar por los jurados

Resumen

Abstract

Dedicatoria

Índice general

1. Introducción biológica	14
1.1. Información hereditaria: ADN	14
1.2. Transcripción y traducción: dogma central de la biología molecular	15
2. Materiales y Metodos	18
2.1. Micromatrices de ADN	18
2.2. Cambios transcripcionales en respuesta a stress abiotico en <i>A.thaliana</i>	20
2.3. PIN - Redes de interacción de proteínas	21
2.3.1. PIN AI1 y LCI binaria	21
2.4. KEGG - Vías metabólicas	22
2.5. GO - Ontología genética	24
3. Métodos de agrupamiento de datos	29
3.1. Similaridad, distancia y disimilaridad	29
3.1.1. Medidas de distancia	31
3.1.2. Similaridad semántica	34
3.2. Estrategias de agrupamiento	36
3.3. Agrupamientos no jerárquicos	36
3.3.1. K-means	36
3.3.2. PAM	37
3.4. Agrupamientos jerárquicos	38
3.4.1. Representación de un agrupamiento jerárquico - dendrogramas	38
3.5. Detectando grupos en el agrupamiento jerárquico	39
3.5.1. Corte de árbol estático	40
3.5.2. Corte de árbol dinámico híbrido	41
3.6. Agrupamiento y similaridad en redes	42
3.6.1. Solapamiento topológico	42
3.6.2. Detección de comunidades en redes	43
4. Análisis de conjunto de datos transcripcionales de respuesta a stress abiótico en plantas	44
4.1. Proceso de filtrado	44
4.2. Agrupamiento con k-means	46
4.3. Agrupamiento con corte de árbol dinámico	49
4.4. Comparación de escalas de resolución de los métodos	52

4.5. Discusión	52
5. Congruencia biológica de agrupamientos transcripcionales	55
5.1. Densidades de interacción	55
5.2. Índice de homogeneidad biológica	56
5.2.1. Modificaciones al Índice de homogeneidad biológica	58
5.3. Congruencia biológica de las particiones transcripcionales	59
6. Coherencia entre métrica transcripcional y espacio de conocimiento GO	63
6.1. Alineamiento de núcleo-objetivo	63
6.2. Espacio de expresión y GO	64
6.3. Alineamiento de núcleo-objetivo local	66
6.3.1. Caracterización de las redes	66
6.3.2. KTA local red 30kmnn	69
7. Métrica mixta	71
7.1. Hacia una métrica mixta	71
7.1.1. Modificación de la similaridad de correlación	71
7.2. Método heurístico	73
7.3. Resultados	76
8. Conclusiones y perspectivas	82

Motivaciones y objetivos

La genómica funcional es un campo de la biología molecular que hace extenso uso de datos genómicos y transcriptómicos para estudiar, describir y responder preguntas acerca de la expresión, función e interacción de genes y proteínas en una escala global (a lo largo de todo el genoma), en contraposición con los métodos más tradicionales de estudio que se realizan gen por gen.

Desde principios del año 2000, a partir de la aparición de tecnologías experimentales modernas, tales como la tecnología de Micromatrices de ADN (DNA microarray), secuenciadores de nueva generación (NGS) o secuenciadores de ARN (RNAseq), es posible relevar el estado transcripcional de una célula de forma global, es decir, cuantificar los niveles de todo el RNA mensajero que está siendo exportado en un dado momento desde el núcleo celular hacia el citoplasma con el fin de producir determinadas proteínas.

La realización de este tipo de estudios posee un potencial enorme, con aplicaciones tanto en áreas de investigación básica como aplicada, investigaciones biomédicas, farmacológicas y de la salud.

En particular, en relevamientos transcripcionales de gran escala es posible obtener información sobre el nivel de activación de miles de genes, para decenas o cientos de condiciones ambientales/experimentales diferentes. Para ganar conocimiento biológico a partir de la cantidad enorme de datos que estos relevamientos generan, es necesario implementar estrategias de búsqueda de correlaciones en espacios de alta dimensionalidad.

Para ello, es de fundamental importancia el estudio e implementación de procedimientos de búsqueda de estructuras aplicables a este tipo de relevamientos, cobrando predominancia técnicas estadísticas y técnicas de aprendizaje automático no supervisado, tales como las técnicas de agrupamiento o “clustering”, que permitan reconocer subconjuntos de genes que evidencien patrones de coexpresión similares a lo largo de conjuntos específicos de condiciones experimentales. [1]

Objetivos y organización de la tesis

El presente trabajo tiene como objetivo analizar la coherencia entre la métrica transcripcional y la inferida a partir de otros espacios de conocimiento, como ser redes de interacción de proteínas (PIN por sus siglas en inglés), redes inferidas de literatura curada (LCI), vías metabólicas (KEGG) y ontología genética (GO).

Vamos a hacerlo cuantitativamente y tratar de incorporar lo encontrado en la elaboración de métricas mixtas que permitan agrupar perfiles de expresión y obtener estructuras compactas (coherentes) en varios espacios. Para esto utilizaremos el conjunto de datos Wiegel, un

exhaustivo estudio de expresión del transcriptoma de *Arabidopsis thaliana*, dos PINs, AI1 y LCI de [2], una red KEGG de [3] y una base de datos de anotaciones GO de [4].

Esta tesis está organizada de la siguiente forma. En el capítulo 1 se introducirán los conceptos biológicos necesarios para comprender y motivar los datos presentados y analizados. En el capítulo 2 introduciremos los materiales y métodos utilizados a lo largo del trabajo. Describiremos la composición y funcionamiento de los métodos de obtención de los cuatro tipos de datos que analizaremos (micromatrices de ADN, redes de interacción de proteínas, redes de vías metabólicas y ontología GO). En el capítulo 3 presentaremos en detalle los métodos de agrupamiento de datos utilizados en este trabajo y analizaremos las problemáticas asociadas a cada uno. En el capítulo 4 analizaremos los datos mediante los métodos presentados en el capítulo 3 y los caracterizaremos buscando información biológica en los mismos. En el capítulo 5 utilizaremos la información obtenida en el capítulo 4 para proponer una métrica mixta que permita aumentar la cantidad de información biológica conseguida previamente. Finalmente, en el último capítulo analizaremos los resultados obtenidos y plantearemos futuras líneas de estudio.

Capítulo 1

Introducción biológica

Este capítulo tiene por objetivo el introducir al lector en los conceptos biológicos básicos necesarios para comprender y motivar los datos presentados y analizados en este trabajo. El lector que desee profundizar sobre los mismos puede remitirse a [5, 6].

1.1. Información hereditaria: ADN

Las células y los organismos pueden ser divididos en dos ramas, procariotas (como las bacterias) y eucariotas (como las plantas, hongos y animales). En las procariotas, el material genético no ocupa una región definida dentro de la célula, sino que se encuentra dispersa en el citoplasma, mientras que en las eucariotas, el material genético se encuentra separado del citoplasma en una región denominada núcleo (figura ??).

Todas las células vivas de La Tierra transmiten su información genética hereditaria por medio del ADN (ácido desoxirribonucleico). El ADN es una molécula unidimensional formada por dos hebras enrolladas una alrededor de la otra en una estructura de doble hélice (figura: ??). Las hebras son cadenas largas de polímeros formadas por monómeros (los nucleótidos), que consisten en dos partes: una columna conformada por un azúcar (desoxirribosa) con un grupo fosfato adherido, y una base nitrogenada, que puede ser adenina (A), guanina (G), citosina (C) o tiamina (T). Cada azúcar se conecta al siguiente mediante un grupo fosfato, con una protuberancia formada por la base, creando de esta manera una cadena polimérica. En principio, es posible extender la cadena de ADN agregando cualquier monómero al final de la misma. Sin embargo, el ADN no se sintetiza como una única hebra, sino a partir de una hebra preexistente, por lo que cada nucléotido debe conectarse mediante puentes de hidrógeno con un nucléotido de la hebra preexistente siguiendo unas reglas estrictas definidas por la estructura complementaria de las bases: A se conecta con T (mediante dos puentes de hidrógeno) y C con G (mediante tres). De esta manera, se forma la estructura de doble hélice con hebras complementarias del ADN. Las uniones entre las bases son mucho más débiles que entre los azúcares y los grupos fosfato, lo que permite a las hebras separarse sin que se rompan.

Un gen es un segmento de ADN que contienen la información necesaria para la síntesis de una proteína en particular. Las proteínas son las moléculas que llevan a cabo casi todos los procesos dentro de una célula, y están compuestas por hasta 20 aminoácidos diferentes. Un gen es por lo tanto una receta que indica el orden en que se colocarán estos aminoácidos,

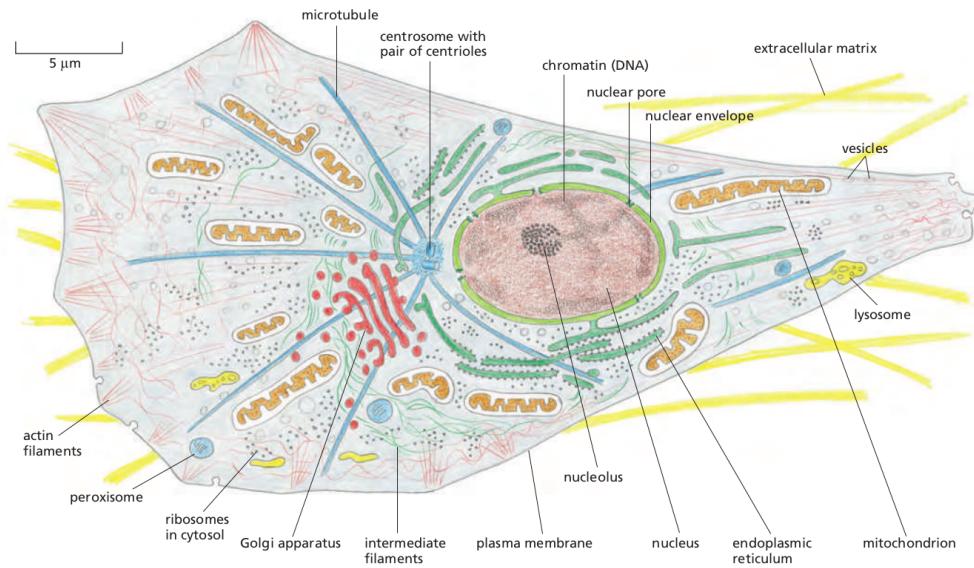


Figura 1.1: La célula eucariota y sus principales características (fuente: [6])

codificada en la secuencia lineal de bases en la molécula de ADN. El genoma es la colección de todos los genes que codifican todas las proteínas que un organismo requiere para vivir. El genoma de un organismo sencillo como el de la levadura contiene alrededor de 6000 genes, mientras que el del humano contiene entre 30000 y 40000. La mayor parte del ADN humano (un 98 %) contiene regiones no codificantes, es decir, hebras que no codifican ninguna proteína en particular. Aunque en un principio se pensaba que esta enorme proporción del genoma no cumplía funcionalidad alguna (se hablaba del genoma basura o “garbage genome” en inglés), estudios recientes sugieren que, al menos en algunos casos, podría jugar un rol regulatorio en la síntesis de ciertas proteínas.

1.2. Transcripción y traducción: dogma central de la biología molecular

Para poder llevar a cabo la función de transmitir información, el ADN debe poder hacer algo más que replicarse. Debe poder expresar esa información que permite la síntesis de otras polímeros: el ARN y las proteínas.

El proceso para la síntesis de una proteína se conoce como transcripción y comienza con la síntesis de una molécula más corta de un polímero llamado ARN (ácido ribonucléico). En el ARN, la columna está conformada por el azúcar ribosa y cuatro bases, uracilo (U) en lugar de timina, y las otras tres bases A, C y G son las mismas que en el ADN, apareándose cada una con su respectiva base complementaria. Durante la transcripción, las hebras de ADN se separan en la región a ser copiada y los monómeros que conforman el ARN son conectados con sus bases complementarias en el ADN (figura ??). La molécula de ARN final es una secuencia que reproduce fielmente la información del gen copiado, donde cada triplete de bases consecutivas (llamados codones) codifica cada aminoácido de la proteína a sintetizar,

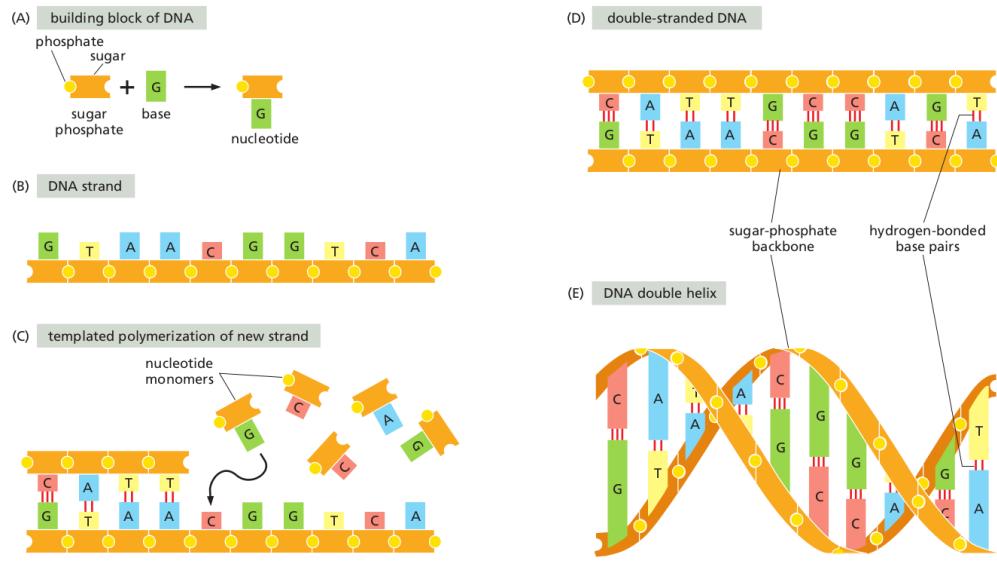


Figura 1.2: Información hereditaria en el ADN. (A) Bloques constitutivos del ADN, esqueleto azúcar-fosfato y base nitrógenada. (B) Una hebra de ADN compuesta por el esqueleto de azúcares-fosfatos y las bases. (C) Polimerización de una hebra a partir de otra que funciona como plantilla. (D) ADN completo con las dos hebras. (E) Forma final del ADN con las hebras en configuración de doble hélice. (fuente: [6])

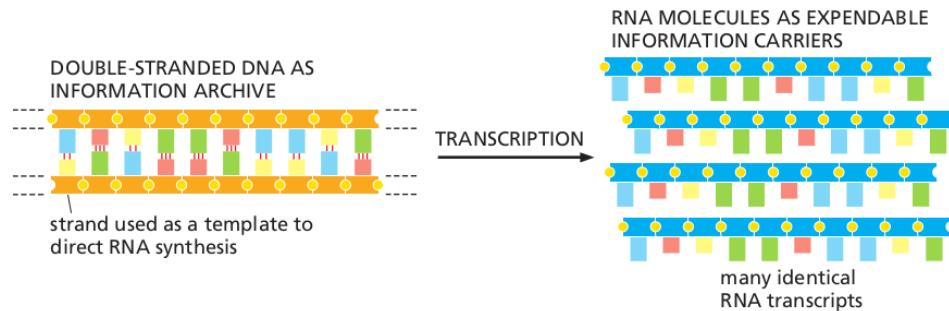


Figura 1.3: Proceso de transcripción de información genética (fuente: [6])

y es esta molécula la que es exportada desde el núcleo al citoplasma en forma de ARNm, dejando la información original intacta dentro del núcleo celular.

Esta molécula de ARNm será luego utilizada por el ribosoma, una maquinaria catalítica compleja consistente en más de 50 proteínas ribosomales diferentes y varias moléculas de ARN ribosomal, para sintetizar la proteína codificada por el gen, en un proceso llamado traducción. Todo el proceso completo de transcripción y traducción se conoce como dogma central de la biología molecular. Si bien cada célula que compone un organismo complejo posee el mismo ADN, células tomadas de distintos órganos realizan diferentes funciones, al igual que las proteínas en los mismos. Por ejemplo, las células de la retina requieren moléculas fotosensibles, mientras que las células que componen el hígado no las requieren. Existe entonces un proceso conocido como diferenciación dentro de cada célula. En lugar de sintetizar todas las posibles proteínas, la célula regula los niveles de transcripción y traducción de los genes que codifican las proteínas necesarias para la misma y únicamente esas proteínas

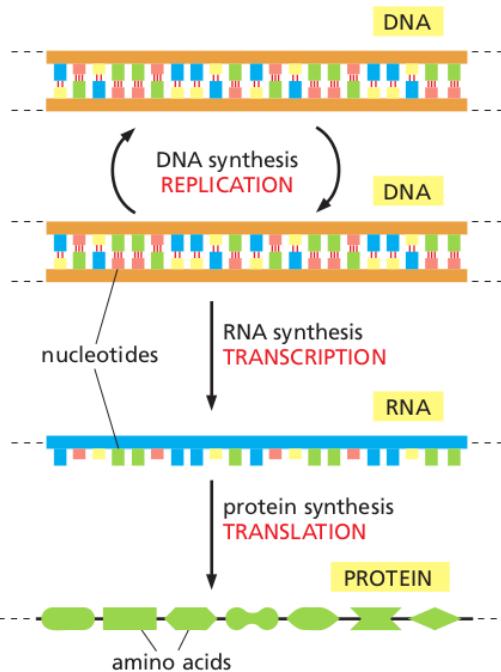


Figura 1.4: Proceso de transcripción y traducción de información genética: dogma central de la biología molecular (fuente: [6])

son las que serán sintetizadas.

En un dado momento, la célula puede requerir muchas proteínas de un tipo y pocas de otro, es decir, en un dado momento cada gen individual puede expresarse a niveles diferentes. La transcripción de un gen (la orden de comenzar a copiarlo o de finalizar la copia) es regulada por proteínas especiales llamadas factores de transcripción, que se ligan a regiones específicas del ADN fuera de la región codificante, que inicia o suprime la transcripción. Esto lleva a la asunción en que se basa el análisis de expresión genética: el estado biológico de una célula queda determinado por su perfil de expresión, es decir, los niveles de expresión de cada gen individual en el genoma, que pueden ser inferidos a partir de las concentraciones de ARNm. Conocer cuáles son los genes que se expresan frente a determinados estímulos, puede brindar información sobre la función que realizan las proteínas codificadas por los mismos en el organismo, información clave para comprender las bases de enfermedades complejas como el cáncer.

Capítulo 2

Materiales y Métodos

Las técnicas de relevamientos transcripcionales de gran escala, tales como las micromatrizes de ADN y secuenciamientos de ARN por medio de secuenciadores de nueva generación, permiten el monitoreo en paralelo del nivel de actividad transcripcional de la totalidad del genoma. En este capítulo daremos una introducción al funcionamiento del primer tipo de tecnologías, que será la que usaremos extensivamente en este trabajo, y a los diferentes conjuntos de datos que utilizaremos. [7]

2.1. Micromatrizes de ADN

La tecnología de micromatrizes de ADN se constituyó como una herramienta indispensable para el monitoreo de niveles de expresión a lo largo de todo el genoma de un organismo, estimando la concentración de ARNm que está siendo exportado desde el núcleo celular hacia el citoplasma para la síntesis de determinadas proteínas.

Una micromatriz es típicamente un portaobjetos de vidrio u otra superficie sólida a la cual se le adosan de forma ordenada y en lugares específicos (llamados sondas) moléculas de ADN. Un mismo sitio puede contener varios millones de copias de moléculas idénticas de ADN de composición conocida (tanto genómico como hebras cortas de oligo-nucleótidos) que se corresponden de forma única con un gen. Una micromatriz de ADN puede medir en simultáneo los niveles de expresión de hasta 40000 genes distintos.

En la actualidad, la aparición de tecnologías más rápidas (y cada vez más económicas) de secuenciamiento, conocidas colectivamente como Secuenciamiento de próxima generación (Next-generation sequencing) y RNA-seq, están comenzando a dejar obsoleta la tecnología de micromatrizes. Sin embargo, las mismas siguen siendo una herramienta útil en el estudio de los perfiles de expresión genética.

Dependiendo de la tecnología utilizada, las micromatrizes pueden ser de canal único o de doble canal.

En las micromatrizes de un solo canal, las moléculas de ARNm son extraídas de las células de interés del organismo y mediante diversas técnicas son transcritas inversamente a ADN. Luego, el ADN es transcrita nuevamente a ARNm utilizando ARN marcado con un compuesto fluorescente (biotina). Estas copias marcadas y aumentadas son luego colocadas en la micromatriz, permitiendo que el ARNm se difunda por toda la misma.

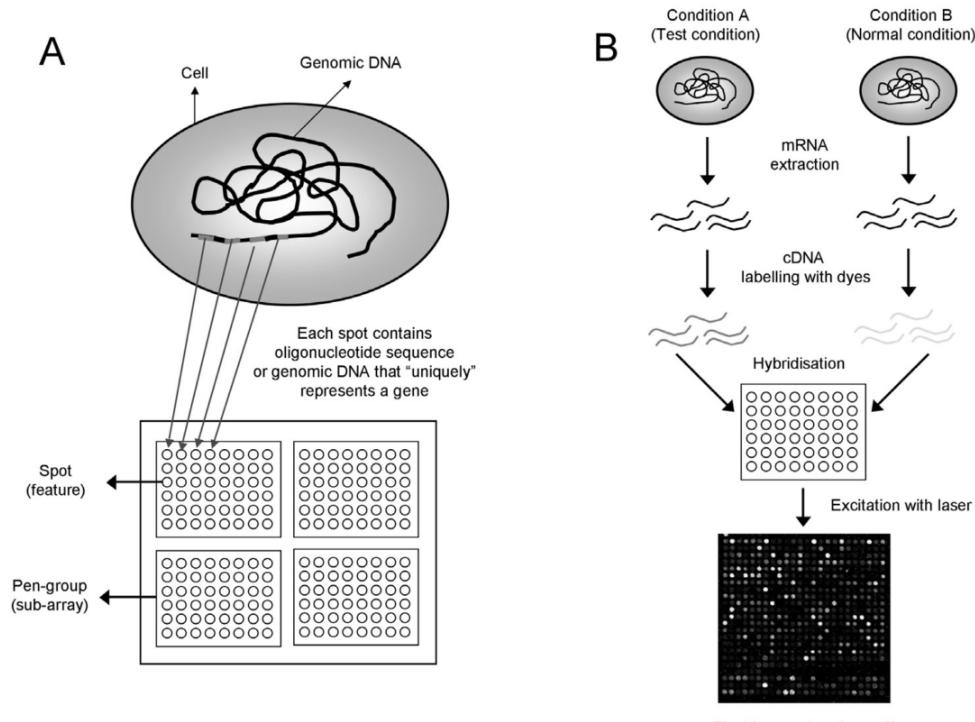


Figura 2.1: Funcionamiento básico de una micromatriz de ADN [Hacer esta figura nuevamente](#)

Cuando el ARNm encuentra una sonda que contiene su copia complementaria, se hibridiza con la misma, es decir se pega con una afinidad mucho mayor con la que se puede pegar a cualquier otra. Al lavarse la solución de ARNm, solo aquellos que se hibridizaron con la copia complementaria se mantienen unidos. Finalmente, se ilumina la micromatriz con luz laser de longitud adecuada y se mide la cantidad de fluorescencia emitida por cada sonda. Esta cantidad está asociada a la cantidad de ácido nucléico que se ligó a una dada sonda y eso a su vez será proporcional a la concentración de ese ARNm particular en el tejido de interés. El resultado de un experimento con micromatrices es una tabla o matriz de expresión de $N_g \times N_m$ donde cada fila corresponde a los niveles de expresión de cada gen particular (N_g del orden de 10000 genes), y cada columna a cada muestra (N_m del orden de 10 muestras) de tejido tomada.

Estos tecnologías plantean entonces el problema de como analizar bastas cantidades de datos para obtener información de interés, como ser:

1. La identificación de los genes que forman parte de algún proceso biológico
2. Agrupar tumores para su clasificación clínica
3. Proveer evidencia de la función de proteínas cuyo rol en el organismo se desconoce

En este trabajo, analizaremos el conjunto de datos Wiegel, datos obtenidos mediante esta tecnología, que detallaremos a continuación. [5,8,9]

2.2. Cambios transcripcionales en respuesta a stress abiotico en A.thaliana

Debido a su naturaleza sésil, las plantas están constantemente sometidas a una variedad de estímulos y stress de origen biótico (virus, interacción con animales, etc) y abiótico (temperaturas extremas, salinidad del suelo, etc). Para sobrevivir en este entorno, estos organismos han desarrollado una compleja red de señalización y respuesta transcripcional que les permiten elaborar las respuestas adecuadas. Por ejemplo, en respuesta a sequía o bajas temperaturas la respuesta transcripcional de estos organismos incluye un incremento en la expresión de genes que codifican para proteínas hidrofílicas, con el objeto de prevenir daños en la membrana. Así mismo, la expresión de genes de proteínas “chaperonas”, que ayudan a estabilizar la forma de otras proteínas en situaciones de congelamiento o deshidratación también es incrementada.

En el marco del proyecto AtGenExpress, (un esfuerzo multinacional desarrollado para descubrir el transcriptoma del organismo modelo multicelular *Arabidopsis thaliana*), el grupo *Weigel & Lohmann*, de Alemania, realizó en el año 2004 un exhaustivo estudio de expresión del transcriptoma de *Arabidopsis thaliana* utilizando las micromatrizes de ADN Affymetrix ATH1, con el objetivo de comprender las complejas redes de genes que, se conjectura, controlan la tolerancia de la planta al estrés. Para ello, se sometió a plantas de *Arabidopsis*, de idéntico genotipo y de idénticas condiciones de crecimiento, a diversos tratamientos de estrés. Los tratamientos de estrés se realizaron de tal forma de excluir efectos circadianos (oscilaciones de las variables biológicas en intervalos regulares de tiempo asociadas con un cambio ambiental rítmico), tomando muestras o de la raíz (root) o del tallo (shoot) en dos réplicas biológicas cada 0 minutos, 30 minutos, 1 hora, 3 horas, 6 horas, 12 horas y 24 horas luego del comienzo del tratamiento. En algunos de los tratamientos, por ejemplo, de luz ultravioleta, las alteraciones transcripcionales producto del estrés fueron tan rápidas que se tomaron además muestras a los 15 minutos del comienzo del mismo. Las muestras de control se tomaron de plantas no sometidas a ningún tratamiento de estrés de la misma forma que con las plantas tratadas.

Además de un tratamiento de control, se realizaron los siguientes tratamientos de estrés:

Tratamiento de frío

Las cajas conteniendo las plantas fueron transferidas a hielo para un rápido enfriamiento y mantenidas a 4°C en un cuarto frío hasta la cosecha.

Tratamiento de calor

Las cajas fueron transferidas a una incubadora y sometidas a una temperatura de 38°C durante 3 horas antes de la cosecha.

Tratamiento osmótico y de sal

Se removieron las balsas de polipropileno que sostienen a las plantas y se agregaron a la solución acuosa, mannitol y NaCl en una concentración final de 300 mM y 150 mM

respectivamente. Luego, se devolvieron las balsas a su lugar hasta la cosecha.

Tratamiento de heridas

Se hirió a las plantas utilizando un elemento punzante consistente en 16 agujas, tres veces por hoja, dejando en promedio entre 3 y 4 agujeros.

Tratamiento de sequía

Las plantas fueron expuestas a una corriente de aire durante 15 minuto, lapso durante el cual perdieron un 10% de su peso. Luego, se las devolvió a la cámara de cultivo hasta la cosecha.

Tratamiento con luz ultravioleta B

Se irradió a las plantas durante 15 minutos con luz ultravioleta B. Bajo estas condiciones se induce una respuesta de la planta tanto para daño por radiación de onda corta como para radiación ultravioleta. [10, 11]

2.3. PIN - Redes de interacción de proteínas

Las redes son construcciones útiles para esquematizar la organización de las interacciones en distinto tipo de sistemas. Las redes permiten tener una vision global de como estan organizadas dichas interacciones.

La mayor parte de las funciones biológicas en una célula es llevada a cabo por proteínas a través de procesos de interacción física entre ellas, por ejemplo formando complejos proteicos. Por lo tanto, es de fundamental importancia conocer no solo los niveles de expresión de una dada proteína, sino también, en simultaneo, las interacciones físicas que la misma podria llevar a cabo con otras proteínas. El registro en forma global de estas interacciones conforma lo que se denomina red de interacción de proteínas o PIN, y si la misma contempla la totalidad de las proteínas de una dada especie, la PIN correspondiente se conoce como interactoma completo.

2.3.1. PIN AI1 y LCI binaria

A lo largo de esta tesis se analizaron dos redes de interacción de proteínas con el objetivo de utilizarlas como referencia.

La primera, una red experimental de interacciones binarias de alta confianza establecida entre 2700 proteinas [2] que reporta 5700 interacciones entre las mismas. Para generar este interactoma, el Consorcio de Mapéo del Interactoma de Arabidopsis utilizó una colección de aproximadamente 8000 marcos abiertos de lectura (secuencias de ARN comprendidas entre un codón de inicio de traducción y un codón de terminación) representando alrededor del 30% de los genes codificantes. Probaron todas las interacciones de a pares con un método conocido como *Sistema de doble híbrido* (Y2H por sus siglas en inglés), consistente en la activación de un gen reporter mediante la acción de un factor de transcripción sobre la secuencia

regulatoria. En esta técnica, el factor de transcripción es separado en dos fragmentos, uno que reconoce la secuencia regulatoria y otro que promueve la activación de la transcripción. Estos dos fragmentos son luego conectados cada uno a cada una de las dos proteínas (llamadas carnada y presa) que se desean analizar. Si las dos proteínas son capaces de interactuar físicamente, el factor de transcripción se reconstituirá y se activará el gen reportero, visualizándose como crecimiento en un medio específico o una reacción con cambio de color. [12] Utilizando los pares obtenidos confeccionaron un conjunto de datos consistente en 5664 interacciones binarias entre 2661 proteínas, llamado Arabidopsis Interactome versión 1 “main screen”, que llamaremos $AI1_{main}$.

La segunda red utilizada fue una red binaria de interacción de proteínas, que llamaremos, $LCI_{binaria}$, obtenida de [2], material suplementario, tabla 4, consistente en aproximadamente 4300 interacciones entre alrededor de 2200 proteínas de *Arabidopsis*. La misma fue obtenida mediante curado manual de literatura, es decir, en lugar de realizar ensayos de alto rendimiento en busca de pares de proteínas interactuantes, se realiza una revisión exhaustiva de la literatura existente en busca de interacciones que aparezcan en ensayos de pequeña escala previamente realizados sobre pocas proteínas y motivados por hipótesis previas (hypothesis-driven en inglés), ensayos altamente fiables. [13]

El solapamiento observado entre ambas se encuentra en el rango esperado dado la cobertura del proteoma que hacen estas redes, como muestra el diagrama de la figura ??.

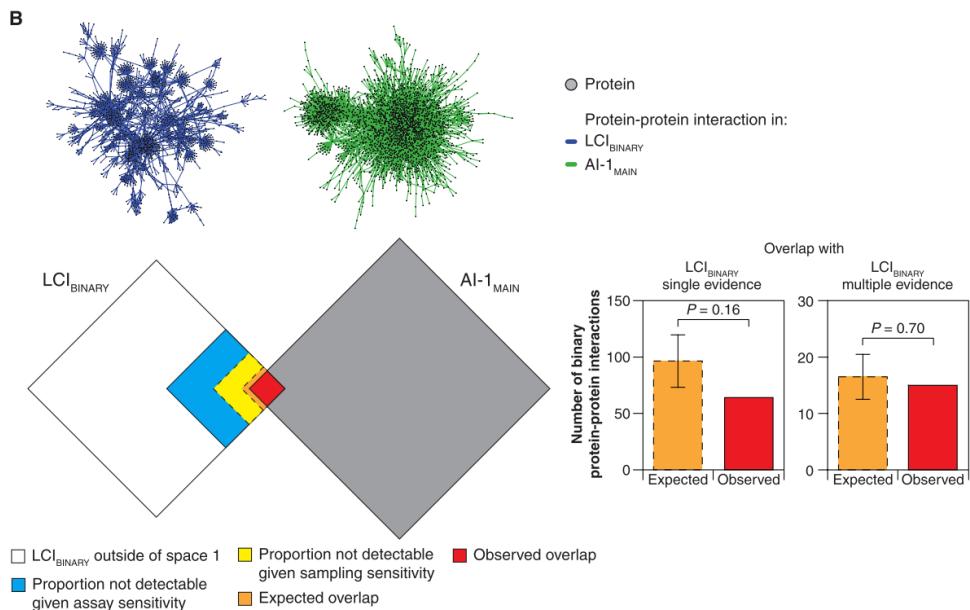


Figura 2.2: (Arriba) Representación de las redes $LCI_{binaria}$ (azul) y $AI1_{main}$ (verde). (Abajo a la izquierda) Los conjuntos de datos son representados por diagramas de Venn cuadrados. El tamaño es proporcional al número de interacciones. (Abajo a la derecha) Los solapamientos observados y esperados de $AI1_{main}$ con $LCI_{binaria}$ (fuente: [2])

2.4. KEGG - Vías metabólicas

Una vía metabólica es un conjunto de reacciones químicas que suceden dentro de una célula. En una vía, la sustancia química inicial, llamada metabolito, es modificada por una

serie de reacciones químicas catalizadas por enzimas. En estas reacciones, el producto de una enzima es utilizado como substrato por la siguiente enzima y así hasta alcanzar un producto final, que puede usarse inmediatamente, almacenarse o iniciar una nueva vía metabólica. El metabolismo de una célula consiste en una red o vías interconectadas que permiten la síntesis o ruptura de las moléculas, acciones conocidas como anabolismo y catabolismo, respectivamente. Descubrir este tipo de redes es fundamental para obtener una imagen global de la actividad celular (figura ??).

La Enciclopedia de Genes y Genomas de Kyoto (KEGG, por sus siglas en inglés), es una base de datos sobre diversos genomas, vías biológicas, enfermedades, drogas y sustancias químicas.

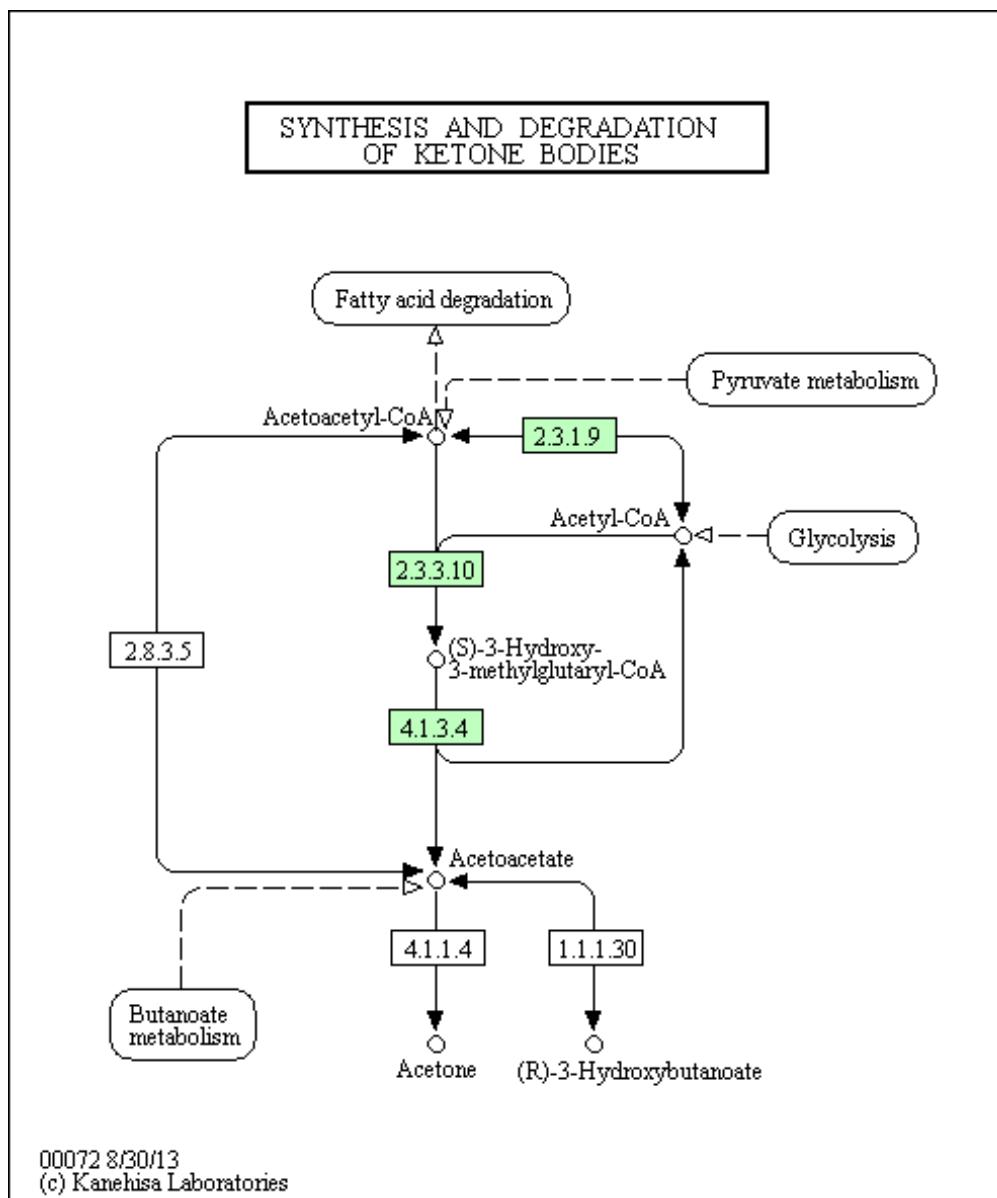


Figura 2.3: Mapa KEGG de la vía metabólica de *Arabidopsis Thaliana* Synthesis and degradation of ketone bodies

La misma provee de una base de datos de vías metabólicas que contiene recursos para la representación de procesos celulares tales como el metabolismo, transducción de señales y ciclo

celular. La figura ?? muestra un mapa KEGG de la vía metabólica de *Arabidopsis Thaliana* “Synthesis and degradation of ketone bodies”. En la misma, se puede observar que a priori, la estructura de una vis metabólica excede el lenguaje de redes. Las interacciones metabólicas suelen involucrar sustratos, productos y enzimas en relaciones que son difícil de capturar utilizando únicamente interacciones binarias. En nuestro trabajo utilizamos el abordaje de Gabriele Sales y colaboradores [14] para mapear vias metabólicas en redes.

Se conformó una red uniendo todas las vías metabólicas presentes en la base de datos y teniendo en cuenta solamente aquellos genes presentes en el conjunto de datos Weigel, obteniéndose una red de 1992 nodos y 23009 arcos. [3, 15]

2.5. GO - Ontología genética

Poder comparar y clasificar entidades es un mecanismo fundamental de las ciencias biológicas. El advenimiento de tecnologías de alto rendimiento hace que sea necesario adoptar sistemas de representación del conocimiento que sean objetivos y estandarizados. Esto llevó al desarrollo de diversas ontologías para anotación de genes y de sus productos, y en particular, al desarrollo de la Ontología Génica (Gene ontology, GO por sus siglas en inglés). El proyecto de Ontología Génica (GO) es un consorcio que intenta mantener un vocabulario y una descripción consistente de conceptos biológicos a lo largo de distintas bases de datos. Esta ontología provee un vocabulario controlado de términos definidos para caracterizar las propiedades de productos génicos (proteínas y secuencias de ARN, por ejemplo).

El proyecto GO consta de tres ontologías estructuradas que describen los productos génicos en términos de sus procesos biológicos asociados (ontología *Biological Process*, BP), de sus componentes celulares (ontología *Cellular Component*, CC) y de sus funciones moleculares (ontología *Molecular Function*, MF).

Un término de un proceso biológico (BP) describe una serie de eventos realizados por uno o varios grupos de eventos moleculares con un comienzo y un fin definidos, por ejemplo, “proceso celular fisiológico” o “transducción de señal”. Un proceso biológico no es equivalente a una vía metabólica ya que no intenta representar la dinámica o dependencias de la misma. Un término de componente celular (CC) describe un componente de una célula que es parte de un objeto mayor, como ser una estructura anatómica (por ejemplo, retículo endoplasmático rígido, núcleo, etc.) o un grupo de productos génicos (por ejemplo, ribosoma, proteasoma, etc.).

Finalmente, los términos de función molecular (MF) describen las actividades que ocurren a nivel molecular, por ejemplo, “actividad catalítica” o “actividad de transporte”.

Cada una de estas tres ontologías está estructurada como un grafo acíclico dirigido (DAG por sus siglas en inglés).

Cada nodo representa un término que describe alguna función. Los términos se unen entre si mediante relaciones direccionales del tipo “es un” o “es parte de”, donde el primero expresa una relación de clase-subclase y el segundo una relación de parte-todo (figura ??). Cuando un producto genético es descrito por un termino GO, se dice que el mismo está anotado en ese término, ya sea de forma directa o a través de herencia, ya que estar anotado en un término implica estar anotado en todos los términos ancestrales, regla conocida como *regla del camino verdadero*.

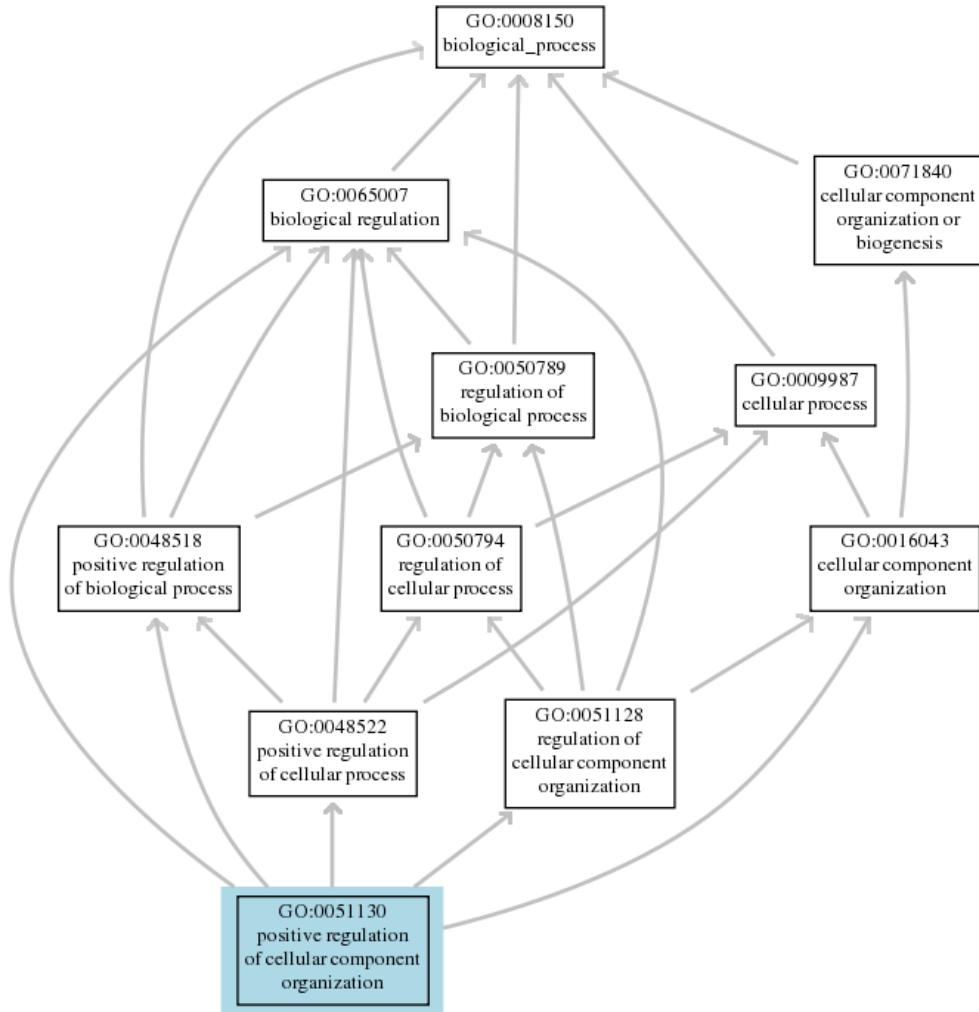


Figura 2.4: Subgrafo de la ontología BP de GO mostrando el proceso biológico “Positive regulation of cellular component organization”.

Formalmente, podemos describir estas relaciones de la ontología GO de la siguiente manera: Sea $C = \{c_i / 1 \leq i \leq N\}$ un conjunto ordenado finito de conceptos que representan términos GO. Los mismos se relacionan entre sí a través de las relaciones antes consignadas, de tal forma que $c_i \rightarrow c_j$ denota que c_i es un/es parte de c_j . Basado en esto, es posible definir una relación binaria sobre C , denotada por \preceq , tal que $c_i \preceq c_j$, es decir c_j es un ancestro de c_i en la jerarquía GO. Notar entonces que si $c_k \preceq c_i$ y $c_i \preceq c_j \Rightarrow c_k \preceq c_j$ (regla del camino verdadero). En cada grafo existe un término raíz de la jerarquía r , tal que $c_i \preceq r \forall c_i \in C$.

Los conceptos más generales se hallarán más próximos al término raíz, mientras que los más específicos e informativos se alejarán del mismo. La anotación de un gen o producto génico se realiza siempre al nodo más específico, pudiendo ser anotado además en varios conceptos biológicos a la vez.

Una anotación en GO para un dado producto génico consiste en un término GO junto con una referencia que describe el tipo de trabajo o análisis que se realizó para asociar un gen con un término específico. Cada anotación debe además incluir un código de evidencia que

indica la forma en que se justifica la anotación a un término particular, lo que le confiere un grado de fiabilidad.

En particular, existen dos grupos de anotaciones, aquellas que fueron curadas manualmente y aquellas que fueron inferidas de anotaciones electrónicas (IEA). Este último tipo de anotaciones funcionales se realiza de forma automatizada sin que intermedie un curador e involucran comparaciones por similitud de secuencia o anotaciones transferidas de bases de datos y por lo tanto poseen una baja calidad y una gran cobertura, formando alrededor del 40% de las anotaciones totales. Además, dentro del grupo de las anotaciones que fueron curadas manualmente, se tienen aquellas que fueron inferidas por medio de experimentos (IDA, IEP, IGI, IMP, IPI), aquellas que fueron inferidas por medio de análisis computacional (IBA, ISS, RCA, ISM). La tabla ?? muestra una descripción de cada código de anotación.

La figura ?? muestra la fracción que representa cada tipo de anotación para cada ontología. Las cantidad total de anotaciones para BP, MF y CC, sin tener en cuenta aquellas pertenecientes a la categoría IEA, totalizan 2540816, 207087 y 1043851 anotaciones respectivamente.

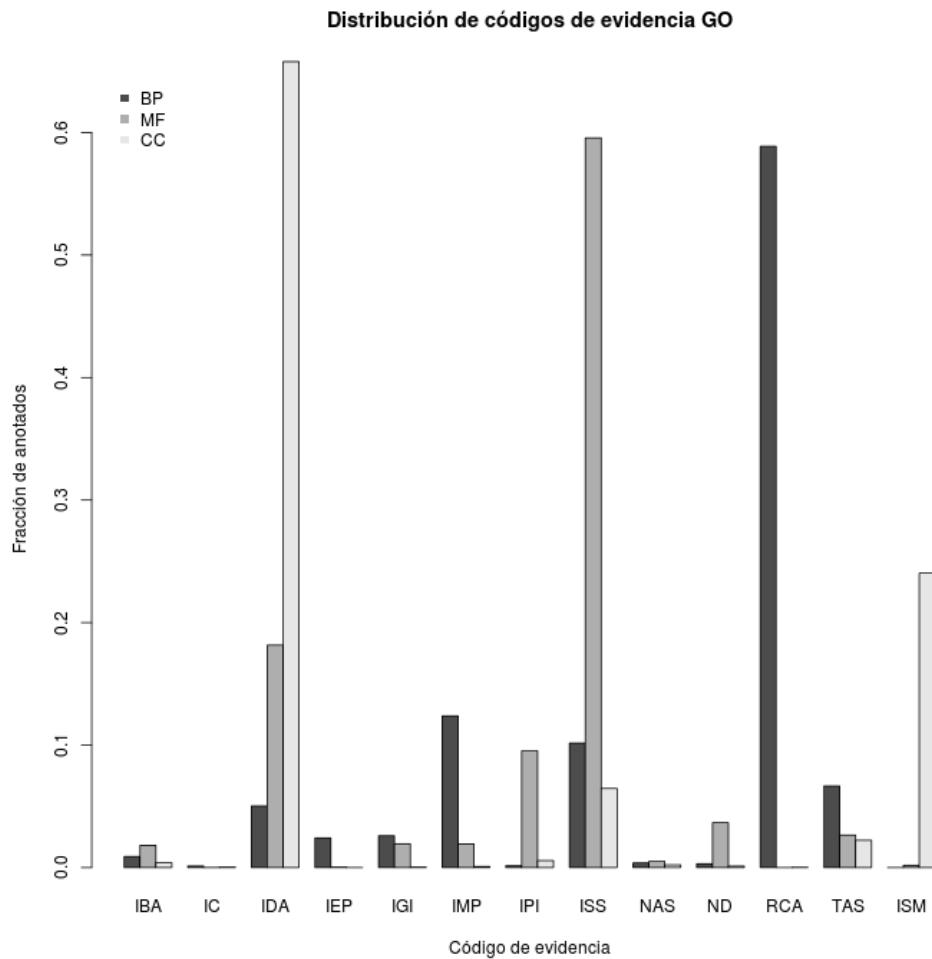


Figura 2.5: Códigos de evidencia en cada una de las ontologías y la fracción del total que representan.

En particular, en este trabajo se tuvieron en cuenta únicamente las evidencias obtenidas experimentalmente. Para ello, se tomaron dos subconjuntos de anotaciones de la ontología

Anotación (inferida de)	Siglas	Tipo	Descripción
Ensayo directo	IDA	Experimental	Indica que se llevó a cabo un ensayo directo para determinar la función, proceso o componente indicado por el término GO.
Patrón de expresión	IEP	Experimental	Cubre casos donde la anotación fue inferida por el tiempo o lugar de expresión de un gen.
Interacción genética	IGI	Experimental	Incluye cualquier combinación de alteraciones en la secuencia (mutaciones) o la expresión de más de un gen/producto génico.
Fenotipo mutante	IMP	Experimental	Cubre los casos donde la función, proceso o localización celular de un producto génico es inferido basado en diferencias en la función, proceso o localización celular entre dos alelos diferentes del gen correspondiente.
Interacción física	IPI	Experimental	Cubre interacciones físicas entre la entidad de interés y otra molécula (como ser proteínas, iones o complejos).
Aspecto biológico de un ancestro	IBA	Análisis comp.	Tipo de evidencia filogenética en donde un aspecto de un descendiente es inferido a través de la caracterización de un aspecto de un gen ancestral.
Secuencia o similaridad estructural	ISS	Análisis comp.	Se utiliza cuando un análisis basado en secuencia forma la base de la anotación y la revisión de la misma se realiza de forma manual.
Ánalysis computacional revisado	RCA	Análisis comp.	Se utiliza para anotaciones realizadas en base a predicciones de análisis computacional de conjuntos de datos de experimentos de gran escala.
Modelo de secuencia	ISM	Análisis comp.	Se utiliza cuando la evidencia de algún tipo de modelo estadístico es usada para realizar predicciones sobre un producto génico.

Cuadro 2.1: Códigos de evidencia GO

BP, que llamaremos BPA, consistente en las anotaciones IDA, IPI, IGI, IMP, con un total de 512235 anotaciones y BPB, consistente en las anotaciones IDA, IPI, IGI, IMP y IEP, con un total de 573688. La diferencia entre estos dos conjuntos es la consideración o no de anotaciones IEP, obtenidas a partir de relevamientos de expresión. Es decir, de ensayos de la misma naturaleza que los que deseamos contextualizar en este trabajo. Además, se utilizó un subconjunto de la ontología CC, consistente en las anotaciones IDA, IPI, IGI, IMP, con un total de 693991 anotaciones. [7, 16–20]

Capítulo 3

Métodos de agrupamiento de datos

Un método de agrupamiento de datos o método de “clustering”, es un método de clasificación no supervisado que permite la partición de un conjunto de N objetos en K grupos o clases, de tal forma que los objetos miembro de un grupo sean más similares entre si (en algún sentido a definir) que entre los miembros de otros grupos.

Son métodos no supervisados ya que en un proceso de agrupación no existen clases definidas previamente ni ejemplos de que tipo de relaciones se desea encontrar entre los objetos, por lo que el mismo proceso debe generar las clases iniciales a las cuales asignar los objetos en el proceso de clasificación.

Estas técnicas permiten el descubrimiento o identificación de distribuciones y patrones subyacentes en los datos, posibilitando obtener conclusiones sobre los mismos, lo que las hace una de las herramientas más útiles en procesos de minería de datos y aprendizaje automatizado en campos tan diversos como las ciencias sociales, las ciencias médicas y la ingeniería.

Dependiendo de los criterios utilizados para realizar la partición, un proceso de agrupamiento puede resultar en diferentes particiones. Como ejemplo de esto podemos tomar el conjunto de números $\{-5, -3, -2, 2, 3\}$. Si decidimos agruparlos por su módulo, obtendremos los conjuntos $\{-5\}$, $\{-3, 3\}$, $\{-2, 2\}$, mientras que si decidimos agruparlos por positividad o negatividad, obtendremos los conjuntos $\{-5, -3, -2\}$ y $\{2, 3\}$. También podríamos haber optado por agrupar por paridad, si son o no primos, etc. Como se observa de un ejemplo tan sencillo, es de fundamental importancia la elección de las propiedades de los objetos a partir de las cuales realizar el agrupamiento.

En el presente trabajo nos interesará agrupar y caracterizar conjuntos de genes de un organismo modelo, la planta *Arabidopsis thaliana*, en base a sus perfiles de expresión génica a lo largo de los diversos tratamientos presentados en el capítulo anterior. [21–23]

Discutiremos a continuación diferentes metodologías y criterios de similaridad que pueden ser considerados para ello.

3.1. Similaridad, distancia y disimilaridad

Las distancias y similaridades tienen un rol preponderante en el análisis de agrupamiento de datos y por regla general son conceptos recíprocos.

Una medida de similaridad o coeficiente de similaridad se utiliza para indicar de forma

cuantitativa la fuerza de la relación entre dos objetos del conjunto. Los $i = 1, 2, \dots, N$ objetos de un conjunto E pueden ser definidos en términos de las coordenadas \vec{X}^i de sus puntos representativos en un espacio $d - dimensional$. Sean $\vec{x} = \{x_0, x_1, \dots, x_d\}$ e $\vec{y} = \{y_0, y_1, \dots, y_d\}$ dos puntos $d - dimensionales$. Entonces, el coeficiente de similaridad entre ambos será una función de sus atributos:

$$s(\vec{x}, \vec{y}) = s((x_0, x_1, \dots, x_d), (y_0, y_1, \dots, y_d)) \quad (3.1)$$

con s una función simétrica, es decir, $s(\vec{x}, \vec{y}) = s(\vec{y}, \vec{x})$. Cuanto mayor es el coeficiente de similaridad, mayor es la similaridad entre ambos.

Por otro lado, las medidas de disimilaridad o de distancia se comportan de forma inversa, a mayor distancia o disimilaridad, más diferentes son dos puntos. Una métrica de distancia es una función $d \in R$ definida sobre un conjunto E que cumple las siguientes propiedades:

1. No-negatividad: $d(\vec{x}, \vec{y}) \geq 0$
2. Reflexividad: $d(\vec{x}, \vec{y}) = 0 \iff \vec{x} = \vec{y}$
3. Conmutatividad: $d(\vec{x}, \vec{y}) = d(\vec{y}, \vec{x})$
4. Desigualdad triangular: $d(\vec{x}, \vec{y}) \leq d(\vec{x}, \vec{z}) + d(\vec{z}, \vec{y})$

con $\vec{x}, \vec{y}, \vec{z}$ objetos arbitrarios del conjunto.

Una medida de disimilaridad es una métrica si cumple con las propiedades antes enunciadas.

Aunque no pareciera existir una definición formal de métrica de similaridad, Chen y colaboradores definen una métrica de similaridad como una función s que cumple:

1. $s(\vec{x}, \vec{y}) = s(\vec{y}, \vec{x})$
2. $s(\vec{x}, \vec{x}) \geq 0$
3. $s(\vec{x}, \vec{x}) \geq s(\vec{x}, \vec{y})$
4. $s(\vec{x}, \vec{x}) = s(\vec{y}, \vec{y}) = s(\vec{x}, \vec{y}) \iff x = y$
5. $s(\vec{x}, \vec{y}) + s(\vec{y}, \vec{z}) \leq s(\vec{x}, \vec{z}) + s(\vec{y}, \vec{y})$

La condición 5 indica que la similaridad entre \vec{x} y \vec{z} a través de \vec{y} no es mayor que la similaridad directa entre \vec{x} y \vec{z} sumada a la autosimilaridad de \vec{y} . Esta propiedad es el equivalente de la desigualdad triangular para una distancia métrica.

Si bien es deseable que una similaridad o disimilaridad sea una métrica, existen muchas medidas de similaridad o disimilaridad que dan excelentes resultados en técnicas de agrupamiento de datos sin ser métricas, es decir, sin que necesariamente cumplan la desigualdad triangular o el ítem 5 de métrica de similaridad. [24]

Finalmente, los objetos del conjunto pueden ser especificados por medio de una “matriz de distancia” de $N \times N$ cuyos elementos d_{ij} indican la disimilaridad entre los puntos i y j . [21–23, 25]

3.1.1. Medidas de distancia

El análisis de datos de expresión genética se basa principalmente en la comparación de perfiles de expresión génica. Para poder comprarlos, se requiere una medida que cuantifique cuán similares o disimilares son los objetos considerados. La elección de una medida de distancia será entonces de fundamental importancia para lograr agrupamientos que tengan sentido en el contexto de los datos analizados. En las siguientes secciones se listarán las medidas de distancia más comúnmente utilizadas en el agrupamiento de datos (no necesariamente de datos de perfiles de expresión).

Distancia euclíadiana

La distancia euclíadiana es probablemente la distancia más utilizada en el contexto de datos numéricos. Para dos puntos \vec{x} e \vec{y} en un espacio $d - dimensional$, la distancia euclíadiana se define como:

$$d_{euc}(\vec{x}, \vec{y}) = \left[\sum_{i=1}^d (x_i - y_i)^2 \right]^{\frac{1}{2}} = [(\vec{x} - \vec{y})(\vec{x} - \vec{y})^T]^{\frac{1}{2}} \quad (3.2)$$

con x_i e y_i los valores de la i esima componente de \vec{x} e \vec{y} respectivamente.

Distancia Manhattan o Taxicab

La distancia Manhattan o taxicab es llamada así por ser la distancia que debería recorrer un taxi en una ciudad para ir de un punto a otro, suponiendo la ciudad como una cuadrícula perfecta. Para dos puntos \vec{x} e \vec{y} en un espacio $d - dimensional$, la distancia Manhattan se define como:

$$d_{man}(\vec{x}, \vec{y}) = \sum_{i=1}^d |(x_i - y_i)| \quad (3.3)$$

Distancia máxima

Para dos puntos \vec{x} e \vec{y} en un espacio $d - dimensional$, la distancia máxima se define como:

$$d_{max}(\vec{x}, \vec{y}) = \max_{1 \leq i \leq n} |x_i - y_i| \quad (3.4)$$

Distancia de Minkowsky

Para dos puntos \vec{x} e \vec{y} en un espacio $d - dimensional$, la distancia de Minkowsky se define como:

$$d_{mink}(\vec{x}, \vec{y}) = \left[\sum_{i=1}^d (x_i - y_i)^r \right]^{\frac{1}{r}}, r \geq 1 \quad (3.5)$$

r es el orden de la distancia de Minkowsky. Notar que si tomamos $r = 2, 1, \infty$ obtenemos la distancia euclíadiana, la Manhattan y la máxima, respectivamente.

Coeficiente de correlación de Pearson

Una de las métricas más utilizadas para medir similaridad entre perfiles de expresión, como los presentados en la figura ??, es el coeficiente de correlación de Pearson [8].

El coeficiente de correlación fue desarrollado por Karl Pearson basado en ideas introducidas por Francis Galton alrededor del año 1880.

Para dos puntos \vec{x} e \vec{y} en un espacio $d - dimensional$, representando el perfil de expresión de dos genes a lo largo de un dado tratamiento, el CCP se define como:

$$r(\vec{x}, \vec{y}) = \frac{\sum_{i=1}^d (x_i - \bar{x})(y_i - \bar{y})}{[\sum_{i=1}^d (x_i - \bar{x})^2]^{\frac{1}{2}} [\sum_{i=1}^d (y_i - \bar{y})^2]^{\frac{1}{2}}} \quad (3.6)$$

o de forma equivalente:

$$r(\vec{x}, \vec{y}) = \frac{\frac{1}{d-1} \sum_{i=1}^d (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y} \quad (3.7)$$

con s la desviación estandar de la muestra. El centrar alrededor de la media permite comparar la forma de ambos perfiles, en lugar de su magnitud.

Valores altos de r implican que las fluctuaciones respecto de la media de las respectivas componentes se encuentran “en sincronía”. En caso de no estarlo el valor esperado tiende a cero. Finalmente si las fluctuaciones tienden a ocurrir “sincronizadamente” pero en sentidos opuestos $r \Rightarrow -1$.

En nuestro caso, $r = +1$ corresponde a genes que están siendo coexpresados por la maquinaria celular, mientras que $r = -1$ a genes anti-coexpresados. En el primer caso los perfiles de ambos genes (apropiadamente reescalados) coinciden perfectamente, mientras que en el segundo son perfectamente opuestos.

La correspondiente medida de distancia puede ser calculada como [26]:

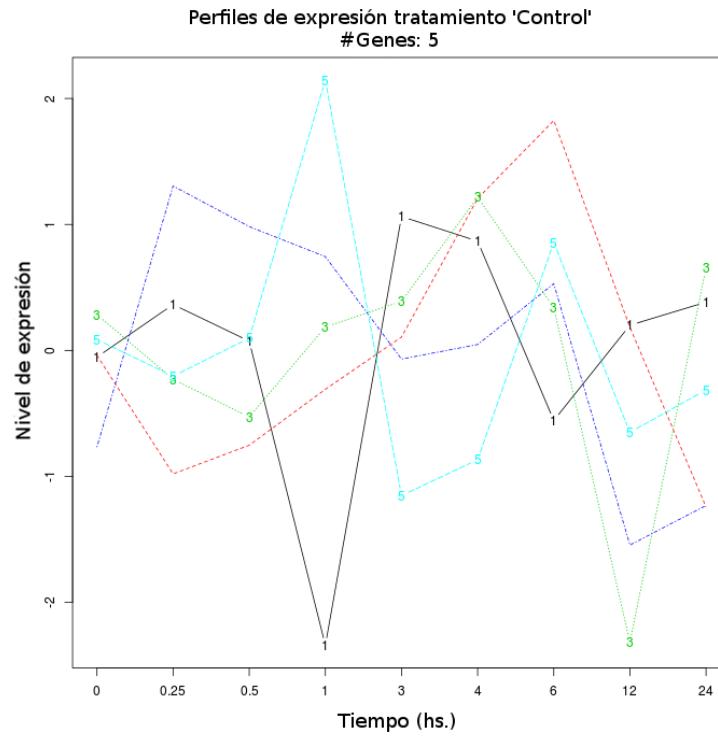
$$d_{ccp}(\vec{x}, \vec{y}) = 1 - r(\vec{x}, \vec{y}) \quad (3.8)$$

o alternativamente:

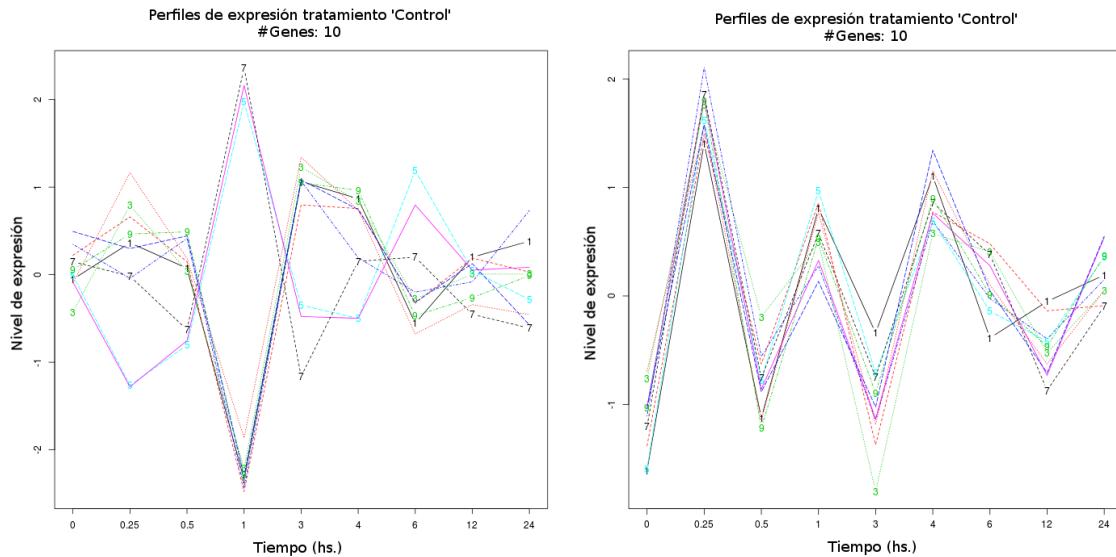
$$d_{ccp}(\vec{x}, \vec{y}) = 1 - |r(\vec{x}, \vec{y})| \quad (3.9)$$

En el caso de la distancia definida en ??, al tomar el valor absoluto del CCP, genes cuyos perfiles son iguales pero opuestos (están anti-coexpresados) pueden encontrarse más cerca en el sentido de d_{ccp} que aquellos que son expresados hacia arriba o abajo pero en distintas magnitudes. Por lo tanto, esta distancia permite encontrar grupos de genes que son coexpresados, sin importar en qué sentido (Figura ??) sean coexpresados. En el caso de la distancia definida en ??, solamente se consideran cercanos aquellos genes cuyos perfiles sean coexpresados o bien hacia arriba o bien hacia abajo (Figura ??). [8, 21, 25, 27]

En el presente trabajo se utilizará como distancia la definida en ?? para encontrar grupos de genes que únicamente se hayan coexpresado o bien hacia arriba o bien hacia abajo. [28]



(a) Perfiles de expresión para el tratamiento Control de cinco genes tomados al azar.



(b) Perfiles de expresión para el tratamiento Control de 10 genes que están co-regulados y anti co-regulados.

Figura 3.1: Distintos grupos de perfiles de expresión

3.1.2. Similaridad semántica

La adopción de ontologías provee los medios para comparar aspectos de entidades que de otra forma no podrían ser comparados. Por ejemplo, si dos productos génicos son anotados dentro del mismo esquema, es posible compararlos mediante el análisis de los términos en los cuales están anotados de forma explícita utilizando medidas de similaridad semántica. Se define una medida de similaridad semántica como una función tal que dados dos términos de la ontología o un conjunto de términos en los que dos genes están anotados, la función devuelve un escalar que refleja la cercanía de sentido entre ellos.

Es posible cuantificar la similaridad semántica en una ontología representada por un grafo como GO, mediante diversas estrategias.

Comparación de términos en GO

Sea C el conjunto de todos los términos de una ontología GO, con un número total $\#C$ de anotaciones. Un término c_i tendrá $\#c_i$ anotaciones, ya sea directamente o por intermedio de cualquiera de sus hijos. La probabilidad de que un gen tomado al azar, sin otro tipo de información, se encuentre anotado al concepto c_i será entonces $P(c_i) = \frac{\#c_i}{\#C}$, con $P : C \Rightarrow [0 : 1]$.

Se define el contenido de información de c_i como $IC = -\log_2(P(c_i))$, cantidad en el intervalo $(0, -\log_2[\frac{1}{\#C}])$, que indica cuán específico e informativo es un término de la ontología. Para un c_i y c_j tales que $c_i \preceq c_j$, se tiene que $IC(c_i) \geq IC(c_j)$. Cuanto más específico sea un término, es menos probable que un gen dado esté anotado en el mismo, y por lo tanto, su contenido de información es mayor. El nodo raíz de la ontología tiene un contenido de información nulo, ya que es el ancestro de todos los términos de la misma y por lo tanto, saber que un concepto está anotado a la raíz no aporta información.

Es posible definir una medida entre pares de términos utilizando el IC. Una de las medidas de similaridad semántica más comúnmente utilizadas es la medida de similaridad semántica introducida por Resnik en [17], que consiste en asignar como la medida de similaridad entre dos términos, el contenido de información del ancestro en común más informativo (el MICA):

$$Sim_{res}(c_i, c_j) = \max_{c \in S(c_i, c_j)} (-\log_2[P(c)]) = IC(MICA[c_i, c_j]) \quad (3.10)$$

Con $S(c_i, c_j)$ el conjunto de ancestros comunes de c_i y c_j . De esta manera, para cuantificar la información compartida (y estimar entonces su similaridad semántica) se considera el contenido de información de los ancestros en común que dos términos poseen.

A modo de ejemplo, tomemos el DAG de la figura ??, con 9 términos o conceptos: $C = \{R, c_0, \dots, c_7\}$ y con 5 entidades mapeadas (genes anotados): $g_1 = \{5, 6, 2, 0, r\}$, $g_2 = \{5, 4, 2, 3, 0, r\}$, $g_3 = \{7, 1, r\}$, $g_4 = \{4, 3, 0, r\}$ y $g_5 = \{2, 0, r\}$. Podemos calcular la similaridad semántica de Resnik entre los términos c_4 y c_5 , por ejemplo, sabiendo que $\#C = 5$ y que el ancestro común más informativo de ambos es c_0 , con $\#c_0 = 4$. Se tiene entonces que $Sim_{res}(c_4, c_5) = IC(MICA) = IC(c_0) = -\log_2(\frac{\#c_0}{\#C}) = -\log_2(\frac{4}{5}) = 0,32$. Si quisieramos calcular ahora la similaridad semántica Resnik entre c_5 y c_6 , obtendríamos $Sim_{res}(c_5, c_6) = IC(c_2) = -\log_2(\frac{3}{5}) = 0,73$. Por lo tanto, para Resnik, los conceptos c_5 y c_6 son entre sí, más similares que los conceptos c_4 y c_5 .

Existen más de dos docenas de medidas de similaridad entre términos GO, y no siempre es

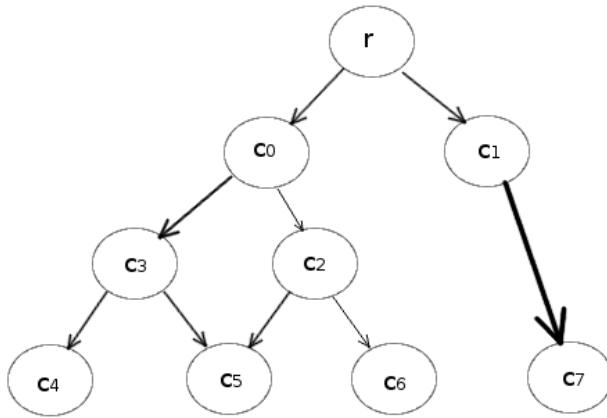


Figura 3.2: DAG con 9 términos o conceptos: $C = \{R, c_0, \dots, c_7\}$ y con 5 entidades mapeadas (genes anotados): $g_1 = \{5, 6, 2, 0, r\}$, $g_2 = \{5, 4, 2, 3, 0, r\}$, $g_3 = \{7, 1, r\}$, $g_4 = \{4, 3, 0, r\}$ y $g_5 = \{2, 0, r\}$

claro cuál es el mejor para un dado propósito. Sin embargo, generalmente la elección de una medida por defecto es suficiente [7]. En este trabajo utilizaremos la Sim_{res} , por tratarse de una medida simple y efectiva.

Una vez establecida una medida de similaridad semántica entre término GO, existen distintas formas para extender esta idea y definir una similaridad semántica entre genes. Básicamente existen 2 estrategias. La primera se basa en medidas globales (groupwise en inglés), que comparan globalmente los conjuntos de términos en los que dos genes están anotados, $GO(g_1)$ y $GO(g_2)$, por ejemplo, contando cuantos términos comparten: $|GO(g_1) \cap GO(g_2)|$. [29]

La segunda estrategia se basa en medidas de a pares (pairwise en inglés), calculando la similaridad semántica término a término de cada uno de los conjuntos $GO(g_1)$ y $GO(g_2)$ y luego aplicando sobre esta similaridad alguna operación para obtener una medida de similaridad entre estos genes.

El primer paso para esto es calcular una matriz de similaridad S de $N \times M$ que contenga la similaridad de a pares, entre todos los pares de términos de estos conjuntos, con $N = |GO(g_1)|$ y $M = |GO(g_2)|$, utilizando alguna de las medidas de similaridad semántica entre términos presentadas anteriormente (Sim_{res} , Sim_{lin} , etc.):

$$S_{ij} = Sim(GO(g_1^i), GO(g_2^j)), \forall i \in \{1, \dots, N\} \text{ y } \forall j \in \{1, \dots, M\} \quad (3.11)$$

Notar que esta matriz puede no ser simétrica.

Cada una de las N filas corresponde a la similaridad entre la anotación i – *esima* del gen 1 y todas las M anotaciones del gen 2 y cada una de las M columnas corresponde a la similaridad entre la anotación j – *esima* del gen 2 y todas las N anotaciones del gen 1.

A partir de S_{ij} es posible definir tres métodos para obtener una medida de similaridad entre genes. El primer método, propuesto en [30], consiste en tomar como similaridad, la máxima similaridad entre todos los términos:

$$Sim_{max}(GO(g_i), GO(g_j)) = \max\{S_{ij}\} \quad (3.12)$$

El segundo método, propuesto en [31], consiste en tomar el valor medio de todos los valores de la matriz S_{ij} :

$$Sim_{med}(GO(g_i), GO(g_j)) = \frac{1}{N \cdot M} \sum_{i,j} S_{ij} \quad (3.13)$$

Finalmente, el tercer método, propuesto en [18], implica tomar el valor medio de los máximos de cada fila, el valor medio de los máximos de cada columna, y quedarse con el máximo de esos dos valores. Este criterio de similaridad se conoce como *remax*:

$$Sim_{remax}(GO(g_1), GO(g_2)) = \max\left\{\frac{1}{N} \sum_i \max_{1 \leq j \leq M} S_{ij}, \frac{1}{M} \sum_j \max_{1 \leq i \leq N} S_{ij}\right\} \quad (3.14)$$

Como muchos genes están anotados en conceptos muy diversos por participar en procesos biológicos muy distintos, e incluso puede haber genes que no están anotados en ningún concepto, la medida de similaridad Sim_{med} tiende a dar valores más bajos que otros métodos. Por el contrario, la medida Sim_{max} tiende a dar valores más altos, por ser una medida más optimista. En este trabajo utilizaremos el tercer método, Sim_{remax} , por ser un compromiso entre ambos casos extremos. [7, 17, 18, 30–33]

3.2. Estrategias de agrupamiento

En la sección anterior abordamos distintas metodologías para cuantificar la noción de similaridad en diversos espacios.

En lo que sigue introduciremos las diferentes estrategias de agrupamiento de datos utilizadas en este trabajo, tanto para agrupamiento de perfiles transcripcionales como de armado de comunidades en las redes presentadas anteriormente.

Es posible distinguir dos tipos de agrupamientos, conocidos como agrupamiento duro (hard clustering en inglés), y agrupamiento difuso (fuzzy clustering en inglés). En el primer caso, el de agrupamiento duro, cada objeto del conjunto de datos es asignado a un y solo un grupo, mientras que en el segundo caso, el de agrupamiento difuso, un elemento del conjunto puede pertenecer a varios grupos, con distinta probabilidad. En este trabajo utilizaremos únicamente métodos de agrupamiento duro.

3.3. Agrupamientos no jerárquicos

Además de la distinción mencionada más arriba, los métodos de agrupamiento pueden dividirse (entre otros) fundamentalmente entre agrupamientos jerárquicos y agrupamientos no jerárquicos. Las dos estrategias de agrupamientos no jerárquicos que se presentan a continuación fueron utilizadas en el desarrollo de este trabajo.

3.3.1. K-means

K-means es un método usual de agrupamiento no jerárquico en donde cada observación pertenece al grupo con la media más cercana a la observación, y donde K , el número de

grupos, debe ser fijado a priori y puede ser determinado arbitrariamente o estimado a partir de algún factor de mérito apropiado. [34]

Más formalmente, sea un conjunto de observaciones $\{\vec{x}_1, \dots, \vec{x}_n\}$, k-means construye una partición de las observaciones en k grupos con $k \leq n$ a fin de minimizar una función de costo, como ser la suma de los cuadrados dentro de cada grupo $G = \{g_1, \dots, g_k\}$:

$$C = \operatorname{argmin} \sum_{i=1}^k \sum_{x_j \in g_i} \|x_j - \mu_i\|^2 \quad (3.15)$$

Con μ_i el valor medio de los elementos del grupo g_i . La figura ?? muestra un conjunto de observaciones y los grupos que se obtienen fijando $k = 2$ y $k = 5$, junto con sus respectivos centroides. Se observa que dependiendo del k utilizado, el algoritmo encuentra particiones con

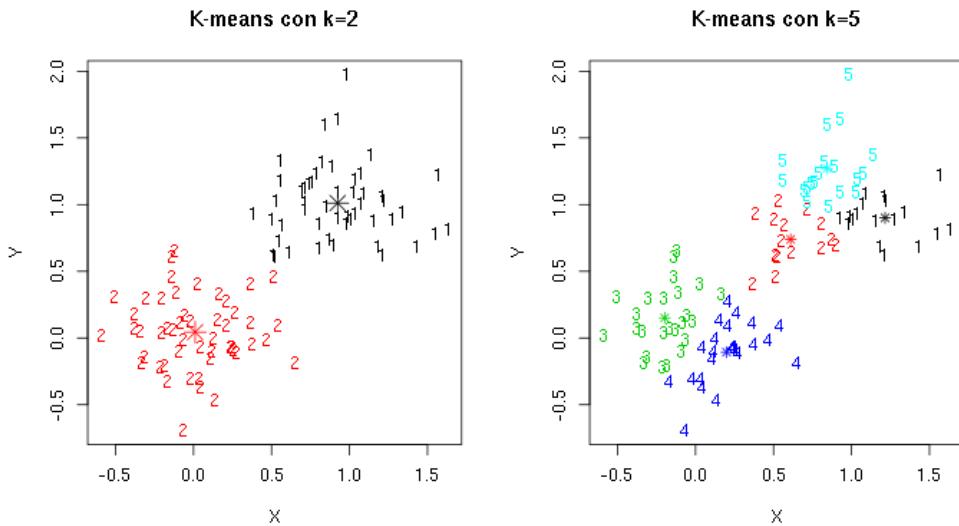


Figura 3.3: Agrupamiento utilizando k-means con $k = 2$ y $k = 5$. mejorar epígrafe

mayor o menor nivel de *resolución*. Volveremos sobre el tema de la resolución más adelante. [35]

3.3.2. PAM

Si bien k-means es uno de los métodos de partición más utilizados ya que es muy eficiente en términos de tiempo computacional, el mismo es muy sensible a observaciones aisladas. Por esta razón, en algunos métodos se reemplazan los centroides, que son puntos no necesariamente pertenecientes al conjunto de observaciones, por medoides, que son los objetos más centrales dentro del grupo (se reemplaza k-means por k-medoids). Esto hace que el método sea insensible a observaciones aisladas.

Particionar alrededor de medoides (Partitioning around medoids en inglés) es uno de los métodos más conocidos que hace uso de este concepto, buscando minimizar la función de

costo:

$$C = \operatorname{argmin}_{i=1}^k \sum_{x_j \in g_i} d(x_j, m_i) \quad (3.16)$$

Con m_i el medoide del grupo i y $d(x_j, m_i)$ la distancia entre el objeto x_j del grupo i y el medoide del mismo grupo. [36, 37]

3.4. Agrupamientos jerárquicos

Existen dos acercamientos distintos para realizar un agrupamiento jerárquico: se puede ir “desde abajo hacia arriba”, agrupando grupos más chicos en grupos más grandes, lo que se conoce como agrupamiento aglomerativo, o se puede ir “desde arriba hacia abajo”, dividiendo grupos más grandes en grupos más chicos, lo que se conoce como agrupamiento divisivo. En este trabajo nos interesarán únicamente trabajar con agrupamientos aglomerativos.

Un agrupamiento jerárquico aglomerativo comienza con cada objeto en un grupo separado. Luego, se unen los dos grupos más cercanos de acuerdo a algún criterio definido generando un nuevo grupo a partir de ambos. Al nuevo grupo se le asignará una distancia al resto de los grupos de acuerdo a cierto criterio. Esto se repite hasta que solo quede un único grupo. Es un tipo de procedimiento determinista y voraz (greedy en inglés), ya que realiza las decisiones tomando en cuenta los óptimos locales en cada etapa, esperando obtener con esto un óptimo global.

Se dice que una partición es más fina (o un refinamiento) de otra partición, si cada grupo de una partición más fina está contenido dentro de un grupo de la partición más gruesa, es decir, cada grupo de la partición más fina es un sub-grupo de un grupo de la partición más gruesa. El agrupamiento jerárquico es un método cuyo resultado es un conjunto de particiones anidadas P_n, P_{n-1}, \dots, P_1 cada vez más gruesas, donde cada nivel más alto une dos grupos de una partición de un nivel más bajo.

Para poder realizar este procedimiento, es necesario definir cuán cercanos son dos grupos. Existen diferentes criterios para ello, entre los que cabe mencionar el método de Ward, el de enlace único y el de enlace promedio. En este trabajo, utilizaremos el método de enlace completo (o complete-linkage en inglés), que consiste en tomar la distancia entre dos grupos como el máximo de la distancia entre sus puntos:

$$D_{cl}(C_i, C_j) = \max_{\vec{x} \in C_i, \vec{y} \in C_j} d(\vec{x}, \vec{y}) \quad (3.17)$$

con $d(\vec{x}, \vec{y})$ la función de distancia utilizada para calcular la matriz de disimilaridad entre los elementos.

Este método permite el manejo de grupos con formas complejas y es invariante ante transformaciones monótonas (como una transformación logarítmica) [38] y considera solamente la separación entre elementos, dejando de lado la compacidad o el balance en los grupos. [21, 39]

3.4.1. Representación de un agrupamiento jerárquico - dendrogramas

Un agrupamiento jerárquico puede representarse como un árbol, llamado dendrograma, que permite una rápida interpretación. En un dendrograma, cada nodo está asociado con una altura h , tal que si A y B son dos nodos del dendrograma, h cumple:

$$h(A) \leq h(B) \Leftrightarrow A \subseteq B \quad (3.18)$$

A modo ilustrativo, la figura ?? muestra el agrupamiento jerárquico realizado sobre 10 puntos colocados de forma aleatoria en el plano, agrupados utilizando la distancia euclídea y mediante los tres métodos vistos anteriormente (Ward, enlace único y enlace completo). De

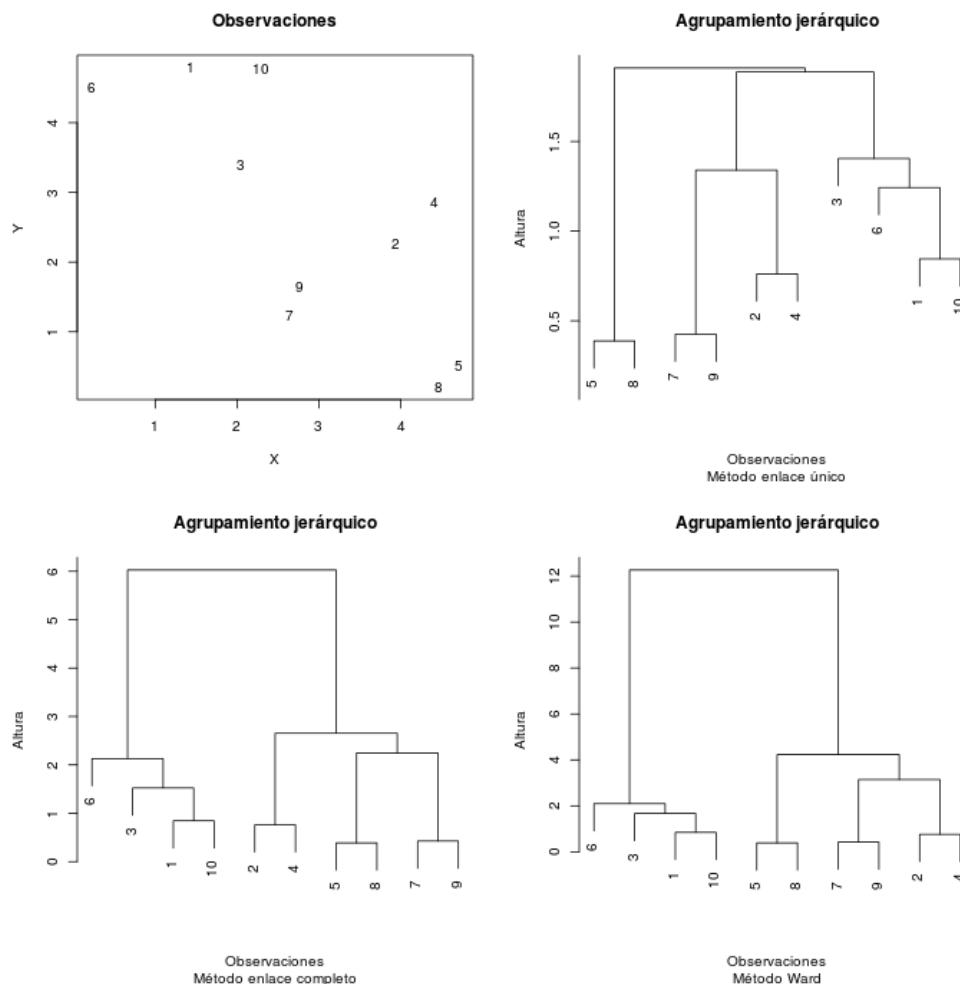


Figura 3.4: Ejemplos de agrupamientos jerárquicos utilizando el mismo conjunto de datos pero distintos métodos de distancia entre grupos. **mejorar epigrafe**

estos gráficos es claro que cada método produce una secuencia diferente de particiones, y dependerá de la aplicación que se requiera, cual de los métodos utilizar.

3.5. Detectando grupos en el agrupamiento jerárquico

El agrupamiento jerárquico organiza los objetos en árboles (dendrogramas) cuyas ramas son los grupos deseados. El proceso de detección de grupos a partir de esta estructura se conoce como corte de árbol, corte de ramas o podado de ramas.

3.5.1. Corte de árbol estático

El método más sencillo de podado es conocido como corte de árbol estático, y funciona definiendo cada rama contigua debajo de una altura fija de corte, como un grupo separado. La cantidad de grupos obtenidos por éste método depende fuertemente de la altura de corte elegida. La figura ?? muestra dos alturas de corte posibles y los grupos que se obtienen a partir de cada una de ellas. Al cortar el árbol en $h = 3$, se obtienen dos grupos, el grupo g_1 ,

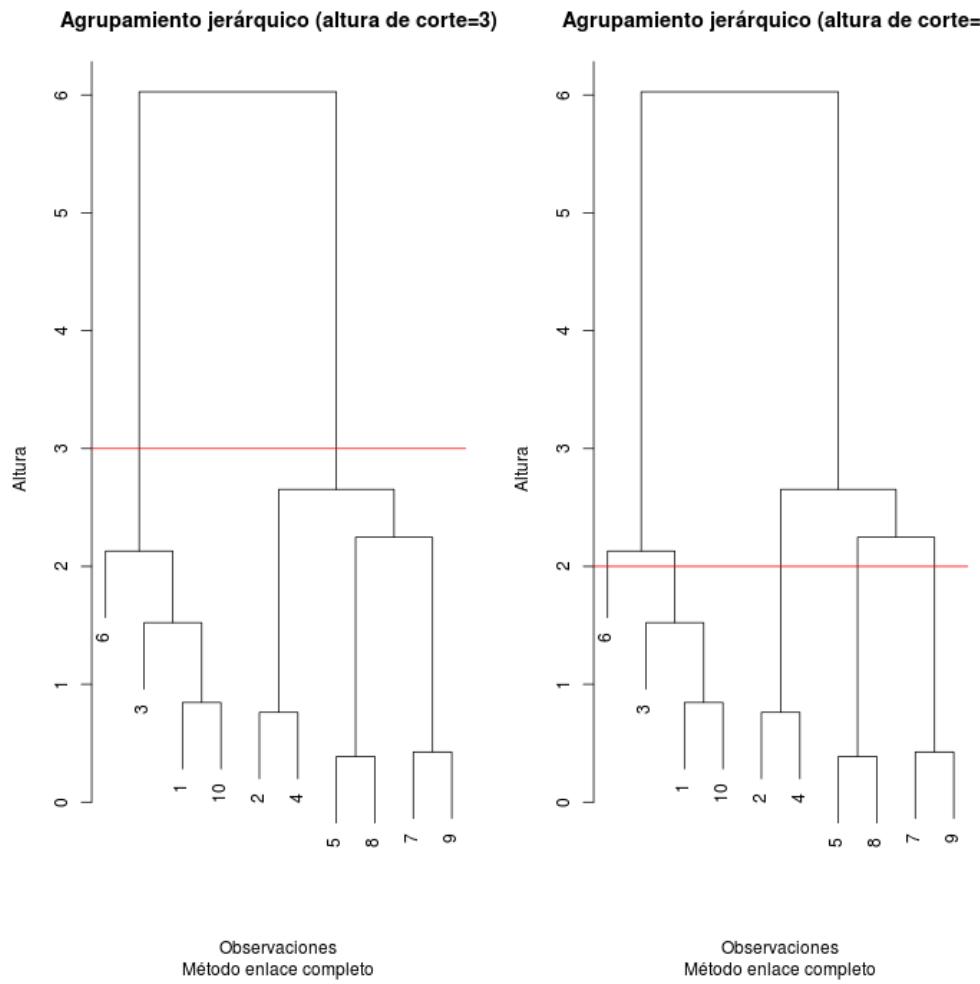


Figura 3.5: Corte de árbol a dos alturas diferentes. mejorar epígrafe

que contiene a las observaciones $\{6, 3, 1, 10\}$ y el grupo g_2 que contiene a las observaciones $\{2, 4, 8, 5, 7, 9\}$, mientras que al cortarlo en $h = 2$, se obtienen cuatro grupos, g'_1 con la observación $\{6\}$, g'_2 con las observaciones $\{3, 1, 10\}$, g'_3 con las observaciones $\{2, 4\}$ y g'_4 con

las observaciones $\{5, 8, 7, 9\}$.

Este ejemplo tan sencillo permite poner de manifiesto que el problema del agrupamiento es un problema “mal planteado”, es decir, cualquier conjunto de puntos puede ser agrupado de maneras drásticamente distintas, sin que exista a priori un único criterio para preferir uno u otro agrupamiento. La fuente de ambigüedades a este respecto más importante, es que la forma en que los datos deberían ser agrupados, depende fuertemente de la *resolución* deseada. Lo que parece una única nube de puntos puede resultar ser, al analizar los datos con mayor resolución, una partición compuesta de muchos grupos. Cada tarea deberá encontrar el nivel adecuado de resolución para obtener la cantidad “correcta” de grupos. [23] [40]

3.5.2. Corte de árbol dinámico híbrido

Si bien es posible detectar grupos distintos en el dendrograma a partir de una inspección visual, utilizar una técnica de corte de árbol estático de forma programática no siempre logra identificar adecuadamente los grupos, ya que al poseer grupos anidados, un solo corte a una altura prefijada no será capaz de detectarlos todos. El método de corte de árbol dinámico híbrido ataca este problema analizando la forma de las ramas del dendrograma en lugar de una altura absoluta [40]. El mismo construye los grupos de abajo hacia arriba en dos pasos. En el primer paso, se detectan las ramas que satisfacen un criterio específico para ser grupos. Este paso de poda está basado en la información de unión del dendrograma. En el segundo paso, se miden cuán cerca de los grupos detectados en el primer paso están todos los objetos no asignados previamente. Si un objeto está suficientemente cerca de un grupo, es asignado a ese grupo. En este paso, se ignora el dendrograma y se utiliza únicamente la información de disimilitud. Este paso puede considerarse un método modificado de particionado alrededor de medoides (modified Partitioning Around Medoids o mPAM, en inglés). Por eso el nombre de *híbrido*, al tratarse de una mezcla entre agrupamiento jerárquico y no jerárquico. Los criterios específicos para la detección de grupos se basan en los siguientes cuatro criterios de la forma de las ramas:

1. Un grupo debe tener una cantidad mínima de objetos.
2. Los objetos que están muy lejos del grupo son excluidos del grupo aunque pertenezcan a la misma rama del dendrograma.
3. Cada grupo debe estar separado de su entorno por una brecha o espacio vacío.
4. El núcleo de cada grupo (el conjunto de objetos con menor altura de unión en el grupo) debe estar fuertemente conectado.

O más formalmente, dado un núcleo de un grupo, llamamos d al promedio de las disimilitudes de a pares entre objetos del núcleo, es decir, a su dispersión y definimos la brecha g de un grupo como la diferencia entre d y la altura donde el grupo se une al resto del dendrograma y entonces, una rama se considera un grupo si:

1. Tiene al menos N_0 objetos.
2. Todas las alturas de unión son a lo sumo de h_{max} .

3. La brecha g del grupo es mayor que un g_{min} .
4. La dispersión d del núcleo es a lo sumo d_{max} .

Los parámetros N_0 , h_{max} , g_{min} y d_{max} son parámetros ajustables del método. La figura ??

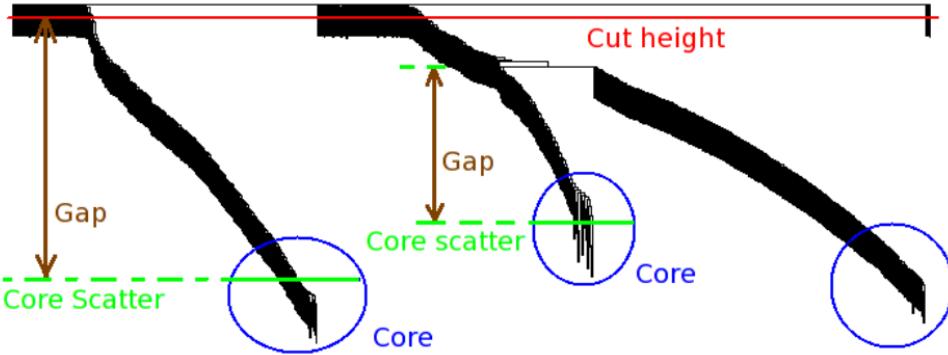


Figura 3.6: Dendrograma simulado con tres ramas con alturas de unión diferentes. La altura de corte corresponde a h_{max} (fuente: [40]). poner en castellano

muestra un ejemplo de los parámetros utilizados para definir los grupos en el paso 1. Para el paso 2, de tipo PAM, los objetos no asignados (o aquellos grupos que no cumplen tener al menos N_0 objetos) son asignados al grupo más cercano si la disimilitud correspondiente es más pequeña que una disimilitud máxima definida previamente, o si es más pequeña que el “radio” del grupo. El “radio” se define como la máxima de las disimilitudes del medoide del grupo al resto de los objetos del mismo.

Es posible controlar la sensibilidad de las divisiones de los grupos mediante el parámetro *deepSplit*, que puede tomar los valores de 1 a 4. Para un *deepSplit* = 1, el método producirá relativamente pocos grupos, de muchos elementos y bien definidos, mientras que para *deepSplit* = 4, el método producirá más grupos pero con una dispersión mayor en el núcleo y separado por brechas más pequeñas.

Para una descripción más detallada del algoritmo, el lector interesado puede referirse a [40], [41].

3.6. Agrupamiento y similaridad en redes

Como se desarrolló en la sección ??, las redes son construcciones útiles para esquematizar la organización de las interacciones en distintos tipos de sistemas. En particular, a partir de estas construcciones es posible estimar nociones de distancia y similaridad entre nodos.

3.6.1. Solapamiento topológico

Una red puede representarse con una matriz de adyacencia $A = [a_{ij}]$ que codifica que pares de nodos están conectados. Para el caso de enlaces no dirigidos, A es una matriz simétrica donde cada a_{ij} puede tomar un valor entre $[0, 1]$. Para una red no pesada, 0 indica que dos

nodos no están conectados, mientras que 1 indica que si lo están. Para una red pesada, el elemento de matriz indica la fuerza de la conexión, tomando un valor entre [0, 1].

A partir de la matriz de adyacencia, es posible construir una matriz de solapamiento topológico $T = [t_{ij}]$ (TOM por sus siglas en inglés) que puede utilizarse como una medida de similaridad para redes biológicas y está definida como:

$$t_{ij} = \begin{cases} \frac{l_{ij} + a_{ij}}{\min\{k_i, k_j\} + 1 - a_{ij}} & i \neq j \\ 1 & i = j \end{cases} \quad (3.19)$$

donde $l_{ij} = \sum_u a_{iu}a_{uj}$, $k_i = \sum_u a_{iu}$ y u es un índice que recorre todos los nodos de la red.

El solapamiento topológico de dos nodos refleja su similaridad en términos de los nodos en común que conectan. Básicamente, t_{ij} es un indicador del acuerdo entre el conjunto de nodos vecinos a i con el conjunto de nodos vecinos a j . Utilizando esta similaridad, se obtiene una matriz de disimilaridad $D = 1 - TOM$ y con esto es posible realizar agrupamientos utilizando, entre otras, alguna de las técnicas antes mencionadas. Además de TOM, es posible definir una matriz de solapamiento topológico generalizada de orden m , $T[m] = [t[m]_{ij}]$ (GTOMm), tal que mida el acuerdo entre los nodos que son accesibles por i y por j en m pasos. [42]

3.6.2. Detección de comunidades en redes

En muchos casos, el análisis de patrones de interacción entre partes de un sistema, relevadas en redes complejas permite reconocer la existencia de una organización modular entre sus componentes. Esto también ocurre en muchas redes de origen biológico [43]. De hecho, en los últimos años, muchos trabajos han explorado el vínculo existente entre comunidades de nodos en una red que representa entidades biológicas y la noción de “módulos de funcionalidad biológica” de Hartwell [44] que postula la existencia de agrupamientos de componentes moleculares, y sus interacciones, capaces de llevar adelante una función biológica específica. En este trabajo utilizaremos dos métodos de reconocimiento de comunas en redes, Infomap [45] y CNM [46].

El método o algoritmo Infomap hace uso de criterios de optimización basados en teorías de información, donde los módulos se definen de tal forma que la longitud media de la descripción de un proceso de paseo al azar en el grafo sea mínima, mientras que el desarrollado por Clauset, Newman y Moore que denominaremos CNM, a partir de ciertas heurísticas, busca particiones de la red optimizando directamente una función de calidad Q .

Ambos métodos serán utilizados en este trabajo con el fin de comparar los resultados obtenidos para las comunidades Infomap y CNM con los obtenidos para los métodos de agrupamiento usados. [19, 45]

Capítulo 4

Análisis de conjunto de datos transcripcionales de respuesta a stress abiótico en plantas

En este capítulo analizaremos el conjunto de datos transcripcionales Wiegel & Lohmann para la planta *Arabidopsis thaliana* presentados en la sección ??, utilizando para ello los métodos de agrupamiento k-means (sección ??) y corte de árbol dinámico híbrido (sección ??) introducidos en el capítulo ?? para obtener grupos en el espacio de expresión.

Una vez obtenidos los grupos en el espacio de expresión, utilizaremos los índices BHI e Interacting Densities para cuantificar el grado de coherencia entre estas estructuras y el conocimiento almacenado en el espacio GO.

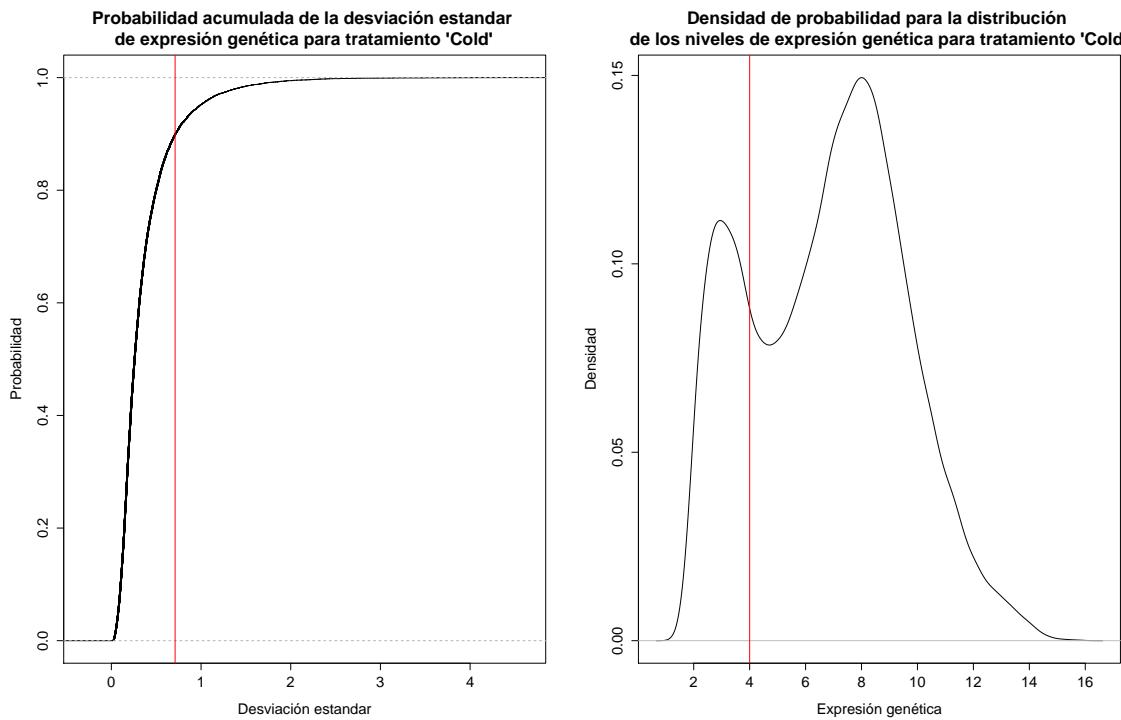
Luego, analizaremos la coherencia de los resultados obtenidos en el espacio de expresión con la de resultados obtenidos en otros espacios de conocimiento, como GO (sección ??), PIN (sección ??) y KEGG (sección ??), esperando que estos conocimientos sean diferentes pero no ortogonales.

4.1. Proceso de filtrado

El conjunto de datos transcripcionales utilizado consta de los niveles de expresión de 22810 sondas que se mapean a 20149 genes a lo largo de 11 tratamientos diferentes y con entre 4 y 9 muestreos en dos réplicas. Para comenzar nuestro análisis realizamos un proceso de filtrado no-específico de manera de focalizarnos en aquellos genes que estuvieran en principio siendo regulados a lo largo del experimento.

Para ello, se aplicaron dos tipos de filtros por tratamiento, por desviación estándar y por tipo “*KsobreA*”. Para el primero, se calculó la desviación estándar por gen a lo largo de todo el tratamiento y se decidió tomar los genes cuya desviación estándar se encontrara en el cuantil 0.9, es decir, utilizar el 10 % de los genes con mayor desviación estándar, considerando estos como los que formaron parte de la respuesta biológica al tratamiento. La figura ?? muestra la distribución de probabilidad acumulada (empírica) de la desviación estándar para los genes del tratamiento “Frío”.

Una vez aplicado este filtro por desviación estándar, se aplicó un filtro de tipo “*KsobreA*”,



(a) Distribución de probabilidad acumulada de la desviación estándar para los genes del tratamiento Frío. (b) distribución de probabilidad para los niveles de expresión para el tratamiento Frío. La recta vertical roja muestra el valor a partir del cual se descartados.

Figura 4.1: Funciones de distribución de probabilidad para perfiles de expresión

que toma únicamente con aquellos genes que tengan al menos K datos por encima del valor A . En nuestro caso, decidimos utilizar como valor de K , la mitad de las mediciones que tuviera el tratamiento. Si el tratamiento tenía mediciones a los 0 minutos, 30 minutos, 1 hora, 3 horas, 6 horas, 12 horas y 24, es decir, 6 mediciones en total, se tomó $K = 3$. Para A , se decidió utilizar una medida usual de $A = 4$, ya que valores de señal menores a 4 no se distinguen del ruido. La figura ?? muestra la distribución de probabilidad para los niveles de expresión para el tratamiento “Frío”. La tabla ?? muestra los filtros aplicados y la cantidad de genes finales por tratamiento. Una vez aplicados los filtros y obtenido los genes de mayor variabilidad en su expresión, se estandarizaron los datos obtenidos para poner a todos los genes en igualdad de condiciones y pesarlos de la misma forma en el agrupamiento. Un procedimiento normal de estandarización de genes para que cada gen tenga media cero y varianza unitaria implica realizar la transformación:

$$\tilde{x}_i = \frac{x_i - \bar{x}}{s_x} \quad (4.1)$$

Con x_i cada observación del gen x a lo largo del tiempo para un determinado tratamiento. Una vez realizado el filtrado y estandarizado procedimos a agrupar los datos mediante los diferentes métodos mencionados en el capítulo 3.

Tratamiento	σ	A	Cantidad de genes
Control	0.37	4	1885
Frío	0.71	3	1955
Osmótico	0.71	3	1923
Sal	0.88	3	1927
Sequía	0.54	4	1870
Genotóxico	0.46	3	1899
Oxidativo	0.41	3	1880
UV-B	0.51	4	1872
Heridas	0.41	4	1877
Calor	0.75	2	1960
Calor y recuperación	0.65	2	1944

Cuadro 4.1: Cantidad de genes y filtros utilizados por tratamiento.

4.2. Agrupamiento con k-means

El método de agrupamiento k-means hace uso de la distancia euclíadiana para minimizar la suma de los cuadrados. Si los datos están estandarizados y centrados, es posible relacionar la distancia euclíadiana d con el coeficiente de correlación mediante la fórmula:

$$d(\vec{x}, \vec{y}) = \sqrt{2(n - 1)(1 - r(\vec{x}, \vec{y}))} \quad (4.2)$$

con n la dimensión del espacio y por lo tanto, para datos estandarizados, la distancia euclíadiana se comportará de forma similar a la distancia de correlación y podremos utilizar el método k-means.

Para decidir el k a utilizar en el método, se realizó un barrido variando k entre $k = 2$ y $k = 30$ con pasos de 1. Al tratarse de un método heurístico, no existe garantía de convergencia al óptimo global y el resultado del mismo puede entonces depender de los grupos iniciales. Por lo tanto, para cada k , se realizaron cien agrupamientos y se midieron los índices de validación internos Calinski-Harabasz y Dunn en cada uno, definidos respectivamente como:

$$CH_k = \frac{SS_B}{SS_W} \frac{n - k}{n - 1} \quad (4.3)$$

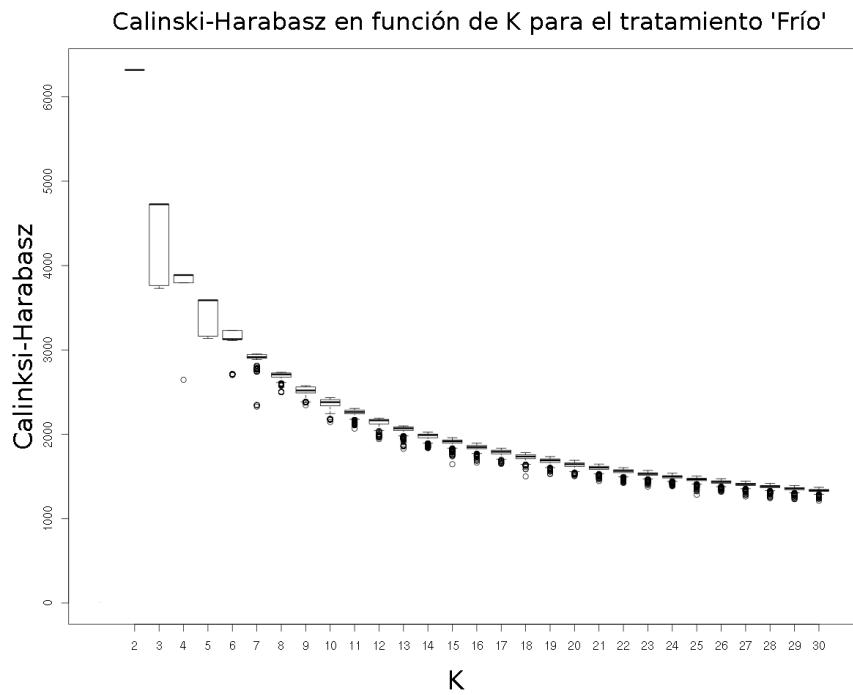
con SS_B el promedio de la varianza entre grupos, SS_W el promedio de la varianza intra grupos, k es el número de grupos y n el número de observaciones y:

$$DI = \frac{\min \delta}{\max \Delta} \quad (4.4)$$

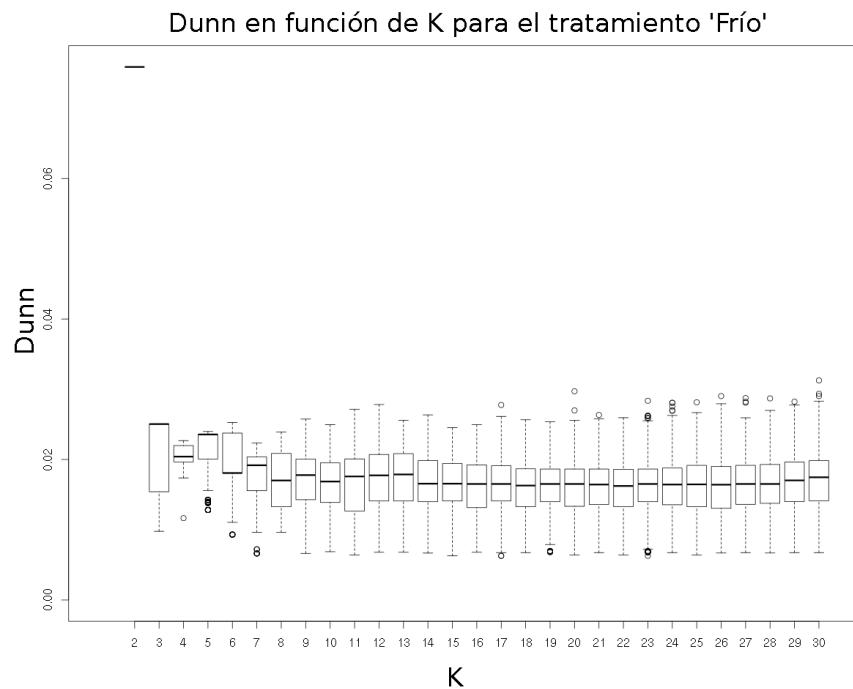
con δ la menor de las de distancias entre grupos y Δ la mayor de las distancias intra grupos. Grupos bien definidos tendrán distancias grandes entre ellos comparados con las distancias intra grupos, por lo que a mayor CH o DI , mejor definidos estarán los grupos.

Las figuras ?? y ?? muestran un gráfico de caja (o boxplot en inglés), para el índice CH y Dunn respectivamente para cada uno de los k en el barrido. Un boxplot consiste en una caja con una linea horizontal que indica el segundo cuartil, es decir, la mediana del conjunto

de datos, y dos líneas verticales llamadas bigotes (o whiskers en inglés) que se extiende una desde el primer cuartil hasta el valor más pequeño del conjunto (con excepción de puntos aislados) y la otra desde el tercer cuartil hasta el valor más grande. Los puntos aislados se grafican de forma separada en el gráfico.



(a) Índice CH de particiones realizadas con k-means para k entre 2 y 30.



(b) Índice Dunn de particiones realizadas con k-means para k entre 2 y 30.

Figura 4.2: Índices de validación interna para particiones realizadas con k-means

Se observa que la cantidad de grupos que maximiza estos índices es 2. Se realizó entonces un agrupamiento con $k = 2$, obteniéndose los perfiles que muestra la figura ??, con una correlación media de $\rho = 0,74$ para el primero y de $\rho = 0,79$ para el segundo, con aproximadamente

el 50 % de los genes en cada grupo. Estas estructuras tan grandes son de difícil interpretación

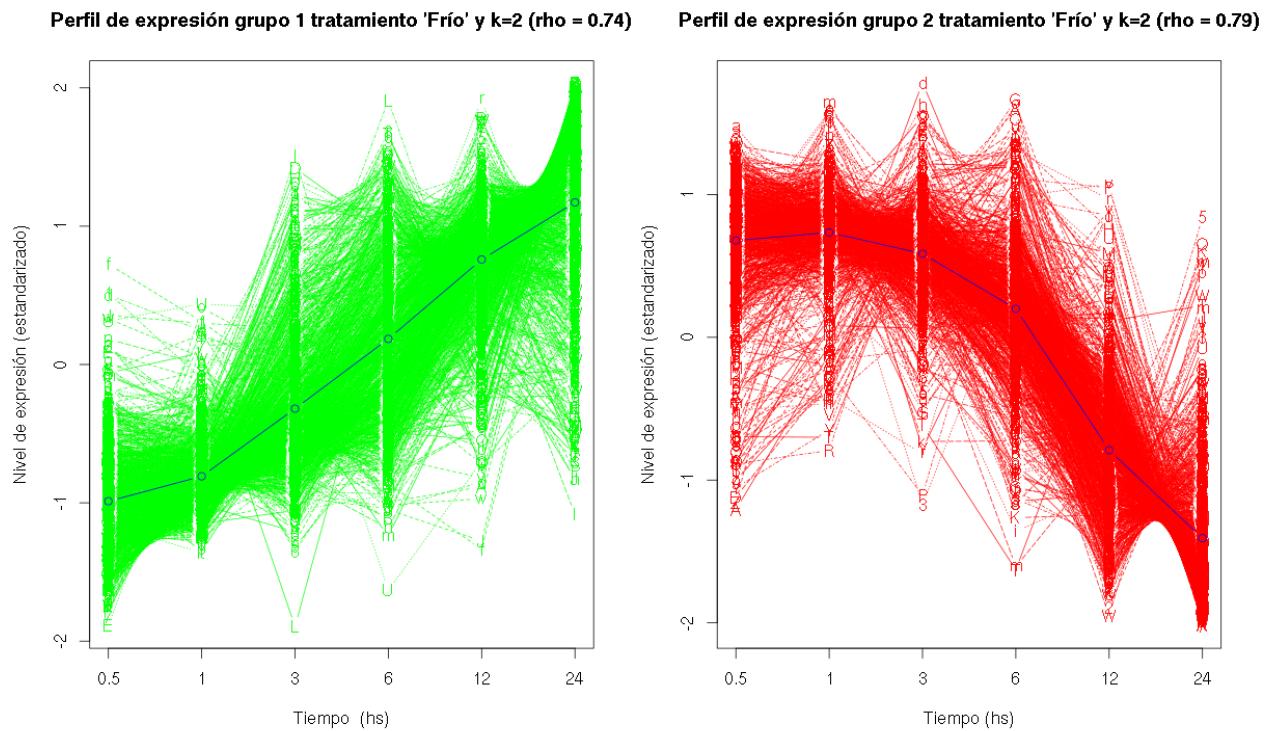


Figura 4.3: Perfiles de expresión génica obtenidos con el método k -means ($k=2$) para el tratamiento 'Frío'. En azul, el valor medio de cada grupo.

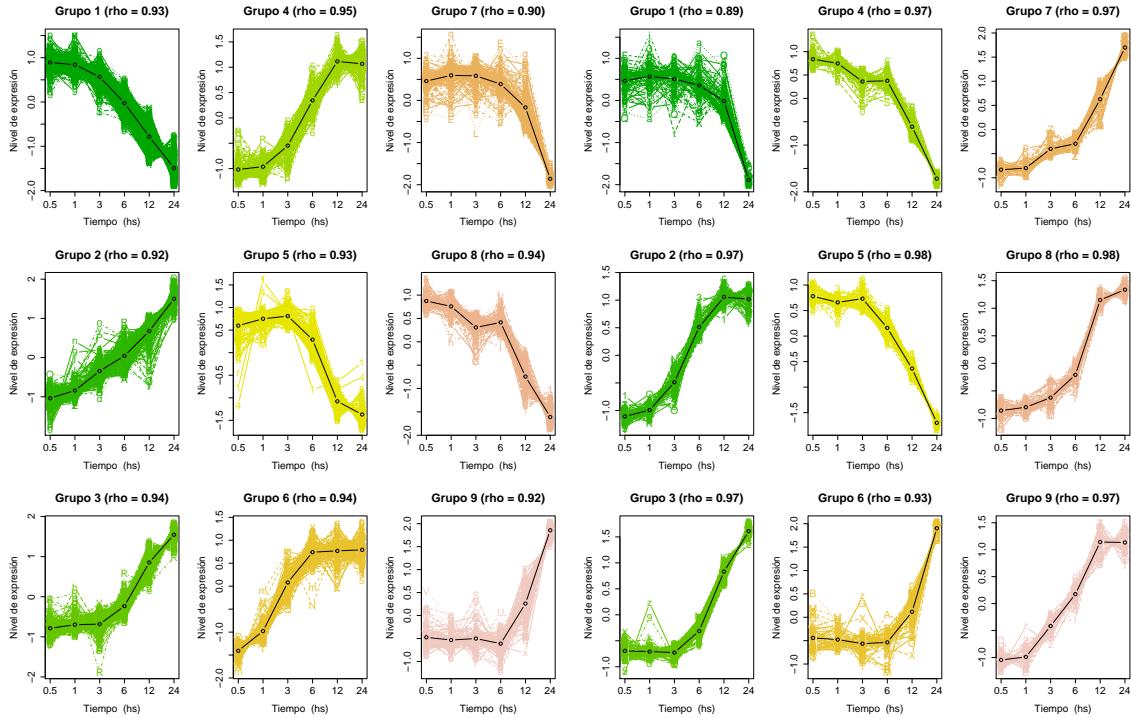
biológica, ya que si bien las respuestas de expresión dentro de cada grupo son similares, existe mucha heterogeneidad en las funciones biológicas de los genes que los componen. El método k -means está entonces trabajando a una escala que no permite extraer información biológica de los grupos. Será necesario entonces aumentar la granularidad mediante la consideración de otros métodos de agrupamiento.

4.3. Agrupamiento con corte de árbol dinámico

Utilizando el método de corte de árbol dinámico se realizó un agrupamiento para cada tratamiento, utilizando alternativamente los parámetros $deepSplit = 1$ (que llamaremos $ds1$, de menor granularidad) y $deepSplit = 4$ (que llamaremos $ds4$, de mayor granularidad). Las figuras ?? y ?? muestran algunos de los perfiles obtenidos con cada parámetro respectivamente para el tratamiento "Frío".

En general, para todos los tratamientos, los grupos obtenidos por este método tienen mayor correlación media (ρ) que los obtenidos por el método k -means, obteniéndose una mayor cantidad de grupos con el $ds4$ que con $ds1$.

Para cada parámetro, cada tratamiento y cada grupo, se realizó un control nulo consistente en tomar la misma cantidad de genes presentes en el grupo, pero de forma aleatoria,



(a) Perfiles obtenidos con ds1.

(b) Perfiles obtenidos con ds4.

Figura 4.4: Perfiles de expresión génica obtenidos con el método corte de árbol dinámico para ds1 y ds4 para el tratamiento 'Frío'. En negro, el valor medio de cada grupo.

del conjunto de genes que formaban el tratamiento, y medir su correlación media. Esto se realizó 1000 veces para cada grupo. Las figuras ?? y ?? muestran la correlación media por tamaño de grupo y el control nulo para ds1 y ds4 respectivamente. Los grupos fueron agrupados por tamaño de a 10 genes, donde los colores más claros indican mayor cantidad de grupos que los oscuros. El gráfico tiene además la media, en negro, y el segundo y tercero cuartil, en gris, para la distribución del control nulo.

Se observa que la correlación media de los grupos es en todos los casos notoriamente superior a la del control nulo. Esto muestra que existe estructura no trivial en los grupos hallados para ambos parámetros.

Por otro lado, en las figuras ?? y ?? se observa que ds1 llega a tener grupos de mayor tamaño que ds4. Esto es esperable ya que cada parámetro aumenta o disminuye la granularidad del método. Estos grupos de ds1 comparativamente grandes tienen alta correlación. Sin embargo hay una menor correlación en los grupos pequeños para ds1 que para ds4. Una posible explicación para esto es que para que exista un grupo grande, es necesario que el mismo tenga alta correlación. De lo contrario, el método buscará partirlo en grupos más chicos hasta maximizar la correlación de cada grupo.

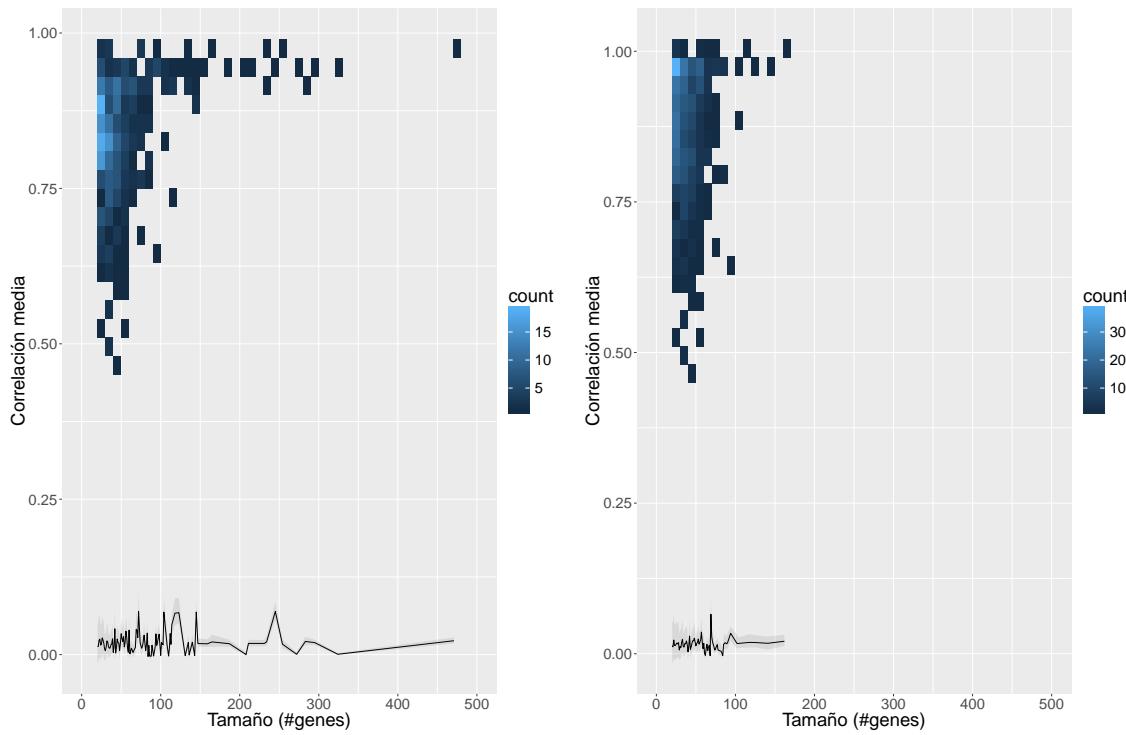
4.4. Comparación de escalas de resolución de los métodos

Otra forma de visualizar la diferencia en los tamaños de los grupos que obtiene cada método es mediante la función de distribución acumulada empírica que se observa en la figura ???. En la misma se observa que corte de árbol dinámico con *ds4* produce la mayor cantidad de grupos con los menores tamaños, seguida por la misma técnica pero con *ds1* y finalmente por k-means con solamente dos grupos muy masivos. Por otro lado, para poder caracterizar la granularidad de las particiones halladas con cada método, podemos calcular la fracción de los grupos más grandes de una partición en otra partición. La figura ?? muestra como están distribuidos los dos grupos de la partición k-means en los grupos de la partición corte de árbol dinámico con *ds1* y los 4 grupos más grandes de la partición corte de árbol dinámico con *ds1* en la de *ds4*.

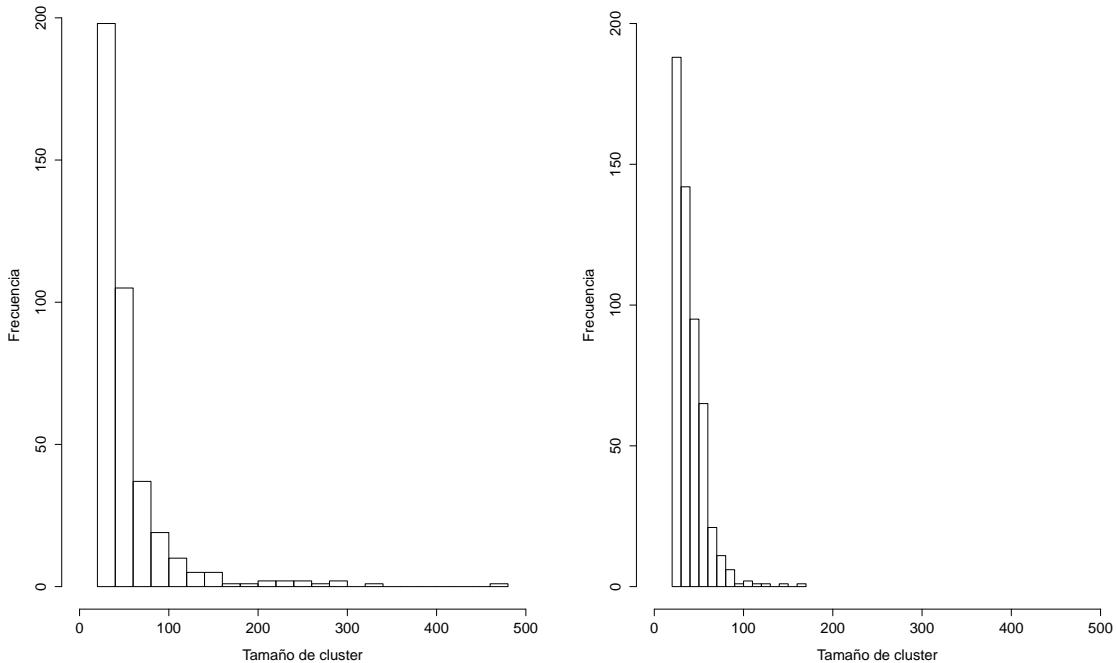
Aquellos subgrupos (los grupos de una partición contenidos en el grupo de otra partición) que están contenidos en más (menos) de un 50 % en el grupo aparecen en rojo (negro). La linea punteada indica el porcentaje del grupo que representaría el subgrupo, si el subgrupo estuviera contenido completamente en el grupo. En cada caso, se observa que los grupos más grandes de una partición con menor granularidad se parten en grupos más pequeños en otra partición con mayor granularidad, es decir, la partición *ds4* está contenida en la partición *ds1* (es un refinamiento de la misma) y esta a su vez está contenida en la partición k-means.

4.5. Discusión

En el presente análisis de estructura de los grupos obtenidos por medio de los métodos k-means, *ds1* y *ds4*, encontramos que todos los métodos producen particiones altamente coherentes, con el método k-means generando las particiones más gruesas y los métodos subsiguientes, refinamientos de las mismas. La alta coherencia detectada es indicativo de que cada método logra hallar estructuras en el espacio de expresión génica, aunque no siempre es factible realizar una interpretación biológica de las estructuras encontradas. Sobre todo en los grupos encontrados con k-means, que solamente toman en cuenta la expresión o inhibición de los genes. En los siguientes capítulos introduciremos algunas herramientas que nos permitirán cuantificar la homogeneidad biológica de las particiones para encontrar la escala óptima en el análisis de expresión.



(a) Correlación media por tamaño de grupo para los grupos obtenidos con ds1 para todos los tratamientos. (b) Correlación media por tamaño de grupo para los grupos obtenidos con ds4 para todos los tratamientos.



(c) Histograma por tamaño de grupo para los grupos obtenidos con ds1 para todos los tratamientos. (d) Histograma por tamaño de grupo para los grupos obtenidos con ds4 para todos los tratamientos.

Figura 4.5: Correlación media por tamaño de grupo para los grupos obtenidos por corte de árbol dinámico con ds1, ds4 y control nulo para todos los tratamientos y sus respectivos histogramas

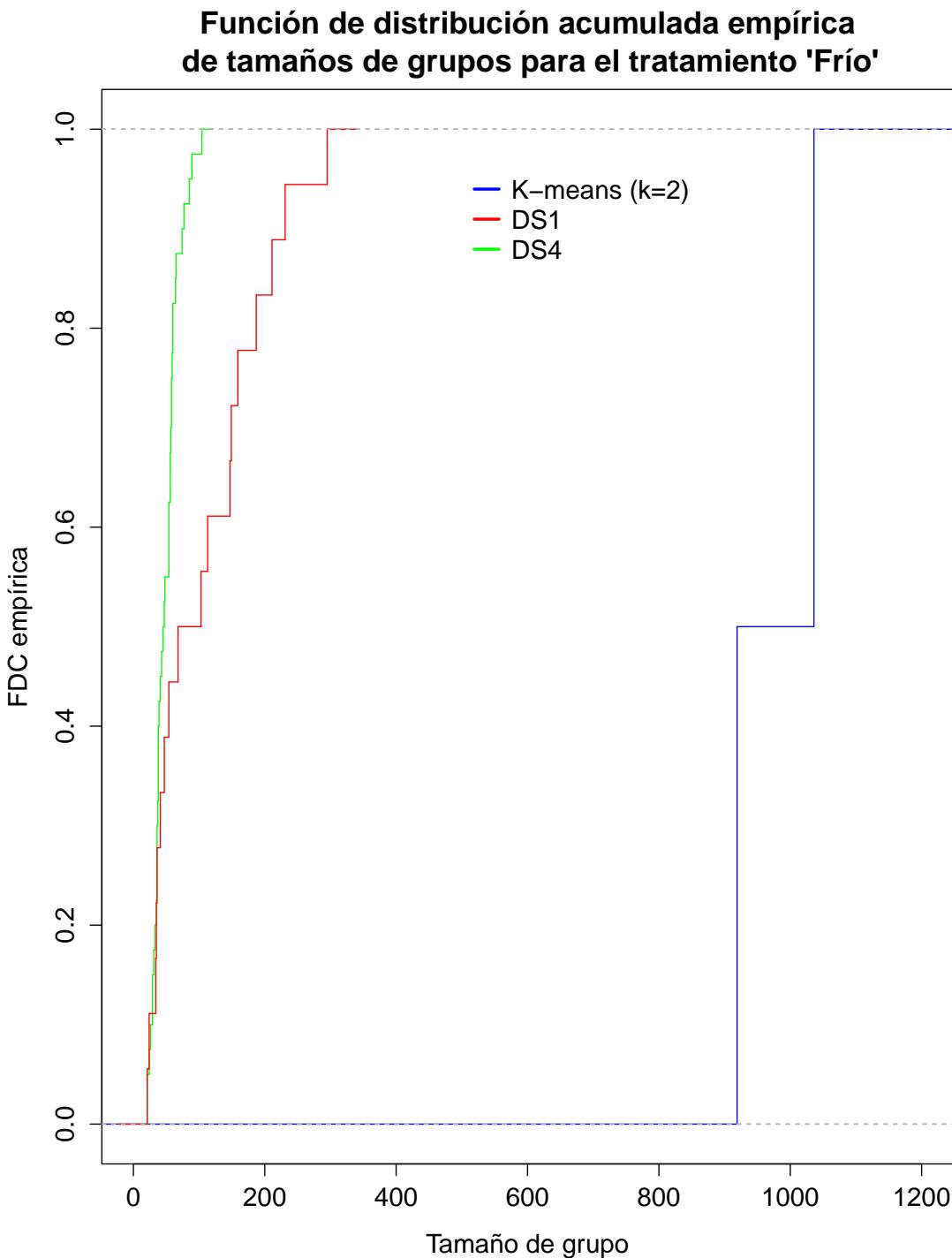


Figura 4.6: Función de distribución acumulada empírica para los métodos k -means ($k=2$) y corte de árbol dinámico ($ds=1$ y $ds=4$) para el tratamiento 'Frío'.

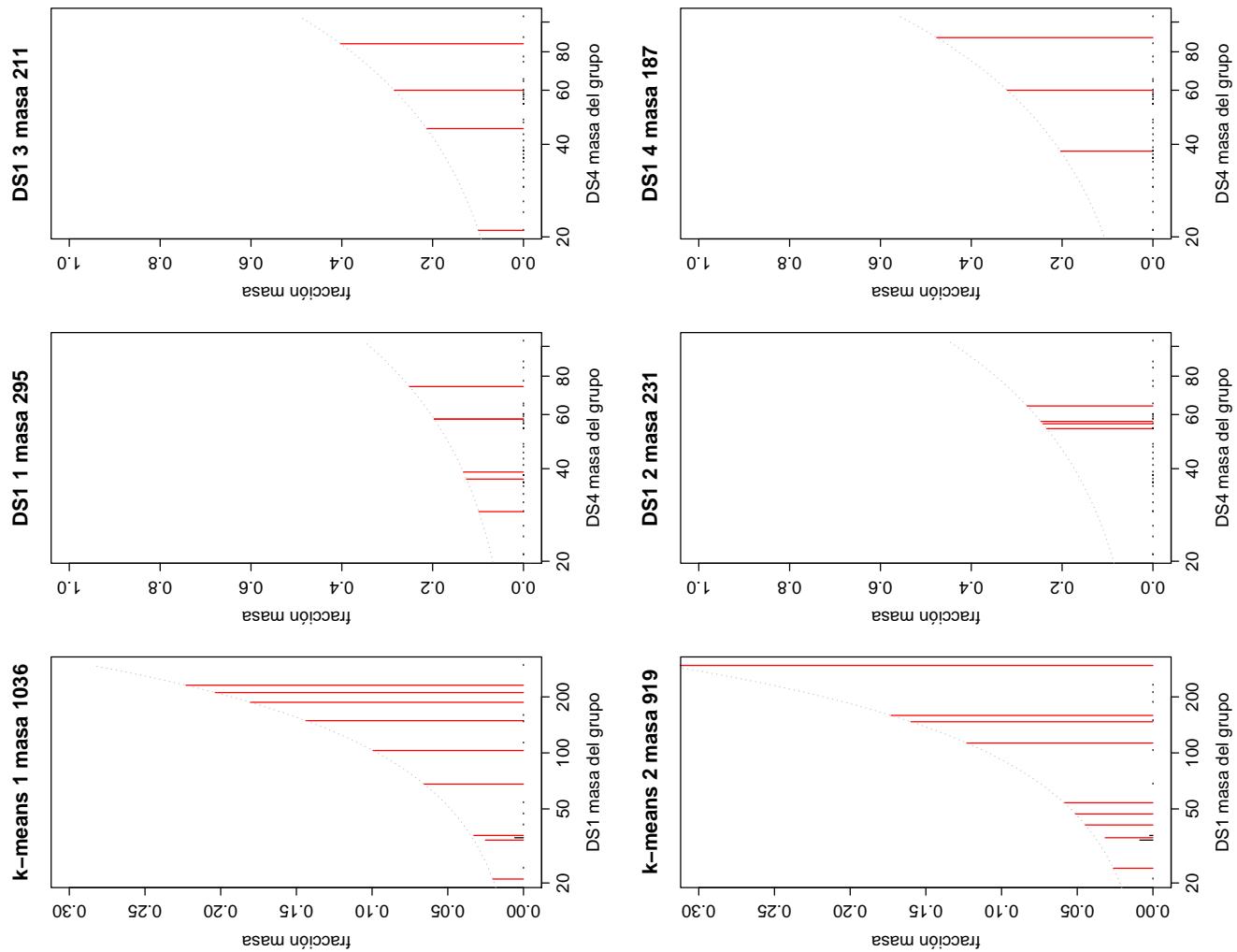


Figura 4.7: Fracción de grupos de una partición más fina dentro de grupos en una partición más gruesa para el tratamiento 'Frío', con $ds1$, $ds4$ y k -means. En rojo, aquellos subgrupos que están contenidos en más de un 50% en el grupo. La linea punteada marca el porcentaje del grupo que representa el total del subgrupo.

Capítulo 5

Congruencia biológica de agrupamientos transcripcionales

Las respuestas transcripcionales relevadas en los datos que estamos analizando ocurren en respuesta a determinadas situaciones ambientales. Las mismas tienen como objeto orquestar la maquinaria celular para llevar adelante determinada funcionalidad biológica que asegure la sobrevida del organismo. Esperamos que los conocimientos (entendidos como nociones de similitud) de los distintos espacios (el de expresión y el biológico) sean diferentes pero no ortogonales. Por lo tanto, una vez detectadas las estructuras en distintas resoluciones en el espacio de expresión, nos interesará cuantificar la congruencia biológica de las mismas. Para ello haremos uso de varios índices, BHI, BHI_{IC} , BHI_{Resnik} , zBHI e ID que servirán como criterios biológicos de validación externos.

5.1. Densidades de interacción

El índice de Densidades de interacción, o ID, introducido por Dutkowski en [47] es un observable que cuantifica el grado en que los genes de una partición comparten anotaciones en GO y además forman parte del mismo grupo en el espacio transcripcional. El mismo se define para un término j en una ontología GO (utilizaremos las definidas en la sección , GO BPA, GO BPB y GO CC) y una partición como:

$$ID(GO_j) = \frac{NE(GO_j)}{N(GO_j)} \quad (5.1)$$

Con $NE(GO_j)$ la cantidad de pares de genes anotados en GO_j que se encuentran juntos en un mismo grupo transcripcional C_x y $N(GO_j)$ la cantidad de pares de genes anotados en GO_j .

Por ejemplo, para un término GO_j con 20 genes anotados y una partición de 3 grupos transcripcionales, donde en el primer grupo hay 5 genes que están anotados en GO_j , en el segundo hay dos y en el tercero hay tres, obtenemos que $ID(GO_j) = \frac{\binom{5}{2} + \binom{2}{2} + \binom{3}{2}}{\binom{20}{2}} = \frac{14}{190} \approx 0,07$.

Como es de nuestro interés establecer un escenario de funcionalidad biológica para los agrupamientos detectados en el espacio transcripcional, nos interesará cotejar la densidad de

interacciones para términos de las ontologías BPA y BPB. Al mismo tiempo analizaremos términos de la ontología CC, ya que asumiremos que para cumplir una determinada función, las proteínas deben colocalizar en alguna estructura subcelular, y por lo tanto estar anotadas en las mismas categorías CC (o similares).

Para tener una idea del orden de magnitud de los indices obtenidos para el caso de estructuras detectadas en el espacio transcripcional, calculamos ademas este indice para las comunidades encontradas en: las redes de interacción de proteínas, y de reacciones metabólicas presentadas en la sección ???. Para ello consideramos las respectivas particiones Infomap (ver sección ??). Finalmente, realizamos un control nulo de tipo 2 para cada tratamiento, reordenando las etiquetas de la partición de forma aleatoria 1000 veces y calculando el *ID* en cada caso.

Para cada ontología utilizada, se calculó la media de *ID* de todos los tratamientos agrupados por cantidad de anotaciones por término, al igual que para cada una de las redes y para el control nulo. Estos valores pueden observarse en las figuras ???. Las escalas utilizadas son escalas logarítmicas. Encontramos que:

1. En todos los casos, términos mas específicos (es decir con menos anotaciones) presentan mayor ID, de nivel notablemente superior que el control nulo, lo que sugiere que los diferentes agrupamientos hasta un tamaño relativamente grande, del orden de 100 genes, en general correlacionan con la información biológica embebida en la ontología respectiva.
2. Las estructuras detectadas en la red de interacciones metabólicas son las que presentan mayor congruencia con las ontologías. Esto es en cierta manera razonable ya que las interacciones de dicha red reflejan interacciones bioquímicas asociadas a una funcionalidad metabólica concreta.
3. Las comunidades encontradas en la red *AI1* y luego las estructuras detectadas en el espacio transcripcional son las que siguen en cuanto a grado de congruencia general con el espacio de conocimiento de las ontologías. Esto es razonable ya que las primeras surgen de información curada de interacción de proteínas, mientras que en el caso transcripcional, la interacción entre proteínas es inferida a partir de su expresión, sin tener en cuenta posibles capas regulatorias postranscripcionales.
4. Las estructuras detectadas en el espacio transcripcional con resolución *ds1* presentan mayor congruencia biológica, según este indicador, que aquellas obtenidas con resolución *ds4*. Esto sería un indicio acerca de cual sería la escala de granularidad apropiada para interpretar las estructuras encontradas.
5. Finalmente queremos mencionar como resultado interesante la baja coherencia encontrada para el caso de la red de interacción de proteínas *LCI*, aunque encontrar las razones de este resultado exceden los alcances de este trabajo.

5.2. Índice de homogeneidad biológica

El índice de homogeneidad biológica (o BHI por sus siglas en inglés) de una partición, introducido por Datta [48] es un observable que cuantifica el grado en que una partición

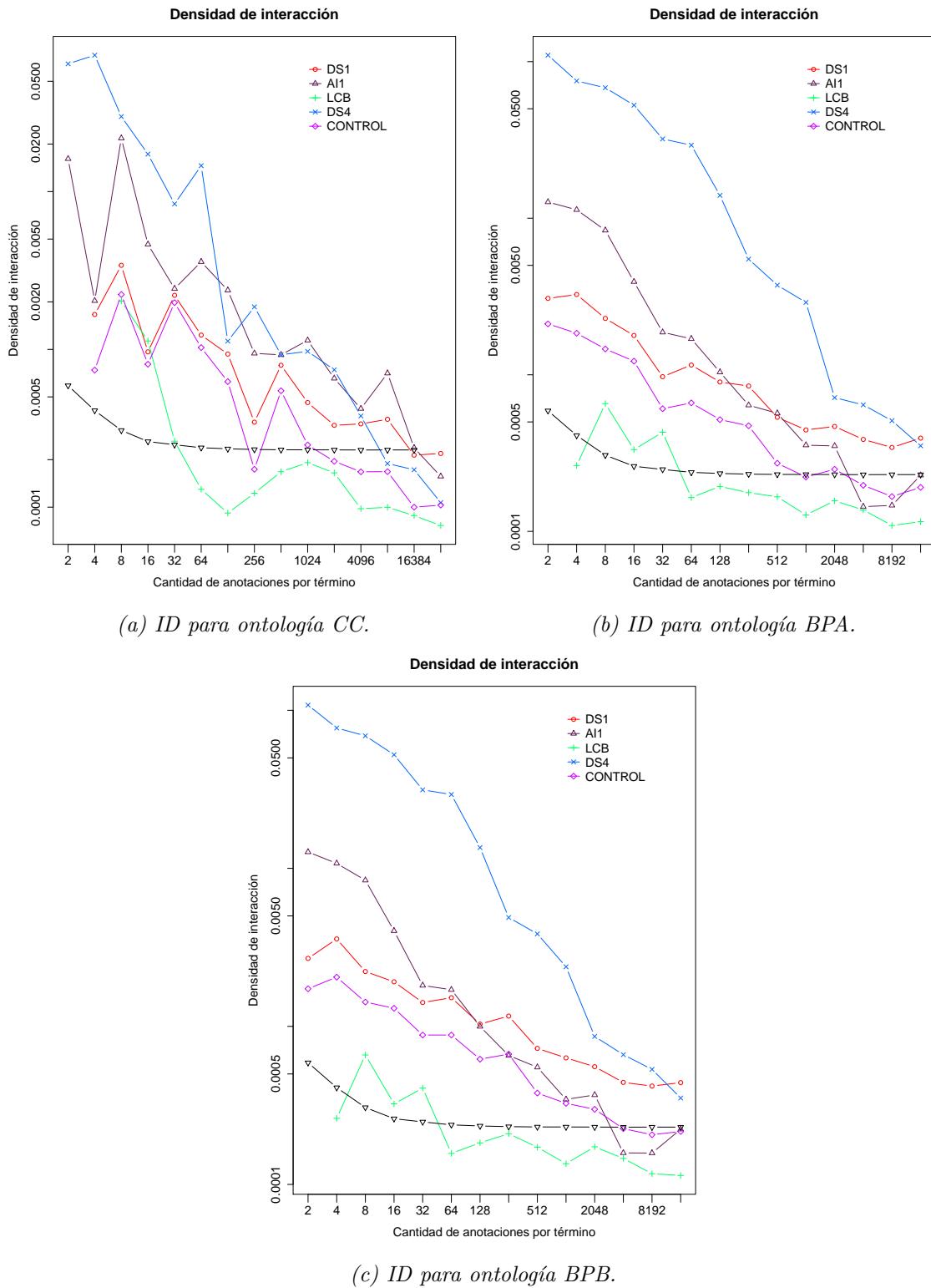


Figura 5.1: Índice de densidad de interacción, ID, para distintas redes de proteínas, vías metabólicas y particiones de expresión.

presenta grupos biológicamente homogéneos, reportando, para cada grupo, la máxima proporción de pares de genes agrupados que comparten una misma clase funcional de Ontología Génica. Consideremos dos genes x e y que pertenecen a un mismo grupo D de una partición dada, con un total de k grupos, y sean $C(x)$ y $C(y)$ los conjuntos de todas las clases funcionales que tienen anotados a los genes x e y respectivamente. Sea además la función indicadora $I(C(x) = C(y))$ que toma el valor 1 si hay al menos una clase en donde ambos genes estén anotados, y 0 en caso contrario. Entonces, el índice de homogeneidad biológica queda definido como:

$$BHI = \frac{1}{k} \sum_{j=1}^k \frac{1}{n_j(n_j - 1)} \sum_{x \neq y \in D_j} I(C(x) = C(y)) \quad (5.2)$$

con n_j la cantidad de genes anotados en el grupo D_j .

5.2.1. Modificaciones al Índice de homogeneidad biológica

Presentaremos a continuación dos variantes del BHI que modificarán la función indicadora para hacer uso de la similaridad semántica y del contenido de información génico.

El índice de homogeneidad biológica con contenido de información (BHI_{IC}) para un grupo se define como:

$$BHI_{IC} = \frac{1}{k} \sum_{j=1}^k \frac{1}{n_j(n_j - 1)} \sum_{x \neq y \in D_j} I(C(x) = C(y)) IC(C(x) \cap C(y)) \quad (5.3)$$

donde el $IC(C(x) \cap C(y))$ es el contenido de información del término más informativo en el que ambos genes se encuentran anotados.

Este índice permite pesar la homogeneidad biológica de un grupo utilizando la especificidad de los conceptos en coincidencia.

Por otro lado, el índice de homogeneidad biológica Resnik para un grupo, BHI_{Resnik} queda definido como:

$$BHI_{Resnik} = \frac{1}{k} \sum_{j=1}^k \frac{1}{n_j(n_j - 1)} \sum_{x \neq y \in D_j} Sim_{rcmax}(C(x), C(y)) \quad (5.4)$$

donde $Sim_{rcmax}(C(x), C(y))$ es la similaridad biológica entre los genes x e y .

Este índice pesa la homogeneidad biológica de un grupo por la similaridad semántica de los genes que lo componen.

Finalmente, el índice de homogeneidad biológica estandarizado para un grupo, $zBHI$, utiliza como referencia los valores obtenidos en situaciones triviales, sin estructura y se define como:

$$zBHI = \frac{BHI - \langle BHI_r \rangle}{s(BHI_r)} \quad (5.5)$$

donde $\langle BHI_r \rangle$ es el valor medio del conjunto de valores del BHI del grupo para un control nulo de 1000 reasignaciones de las etiquetas de la partición y $s(BHI_r)$ es la desviación

estandar de la muestra para el mismo conjunto.

Se realizaron además dos tipos de controles nulos.

El primero, un control nulo que llamaremos “control nulo 1”, se realizó considerando agrupamientos aleatorios (de tamaños variables) a partir de un conjunto de los aproximadamente 6000 genes que pasaron el filtrado en al menos un tratamiento analizado. Para cada tamaño de agrupamiento analizado (entre 2 y 500 genes), se realizaron 1000 grupos aleatorios y se calculó su BHI. Se encontró que el valor medio de los ensambles se mantenía aproximadamente constante, mientras que existía una dependencia de la desviación estandar con el tamaño de los grupos. Se realizaron dos ajustes por funciones de ley de potencias, para tamaños entre 1 y 50 y de 51 en adelante. Las funciones halladas permiten rápidamente obtener el BHI aleatorio medio para una partición de cualquier tamaño y su desviación estandar.

El segundo, que llamaremos “control nulo 2”, consistió en realizar 1000 reasignaciones aleatorias de las etiquetas de cada partición y calcular el BHI de cada grupo de la misma. Encontramos que la media de BHI calculada de esta manera coincidía con la del control nulo anterior, pero no así su desviación estandar. Concluimos que la diferencia fundamental se basa en que en el segundo caso, en la reasignación de etiquetas, se mantiene siempre la estructura de tamaños de la partición, mientras que en el primer caso, cada grupo fue tomado por separado.

Para caracterizar el comportamiento de cada uno de estos índices se midieron los mismos para cada tratamiento y se calculó su correlación de a pares de índices. La figura ?? muestra las distribuciones y correlaciones de a pares para estos índices en el tratamiento “Frío”. Se encuentra que los índices modificados tienen una alta correlación entre ellos y con BHI. Esto parecería indicar que no aportan más información que la que se obtiene a través del índice original. Por ser el más sencillo de calcular, es el que utilizaremos como criterio de validación externa de la calidad de una partición.

5.3. Congruencia biológica de las particiones transcripcionales

Los valores de BHI calculados para cada uno de los grupos del tratamiento “frío” en las particiones k-means (puntos rojos), *ds1* (triángulos verdes) y *ds4* (cuadrados azules) se presentan en las figuras ??, con control nulo 1 y ?? con control nulo 2. Los grupos fueron ordenados según su masa de forma creciente.

Se observa que de los dos grupos de kmeans, solo uno presenta un BHI superior a una desviación estandar para el control nulo 1 y del tercer cuartil para el control nulo 2, mientras que para *ds1*, el 45 % de los grupos superan una desviación estandar para control nulo 1 y el 30 % de los grupos superan el tercer cuartil para el control nulo 2.

Finalmente, para la partición *ds4*, aproximadamente el 40 % de los grupos presenta un BHI por sobre una desviación estandar para el control nulo 1 y un 50 % presenta un BHI por sobre el tercer cuartil del control nulo 2.

Esta baja calidad en el índice BHI de las particiones se encontró de forma similar a lo largo de todos los tratamientos. Esto sugiere que si bien el aumentar la granularidad de la partición con el método corte de árbol dinámico resulta en un aumento de la consistencia biológica global de las estructuras observadas, esto no implica que las resoluciones utilizadas sean las

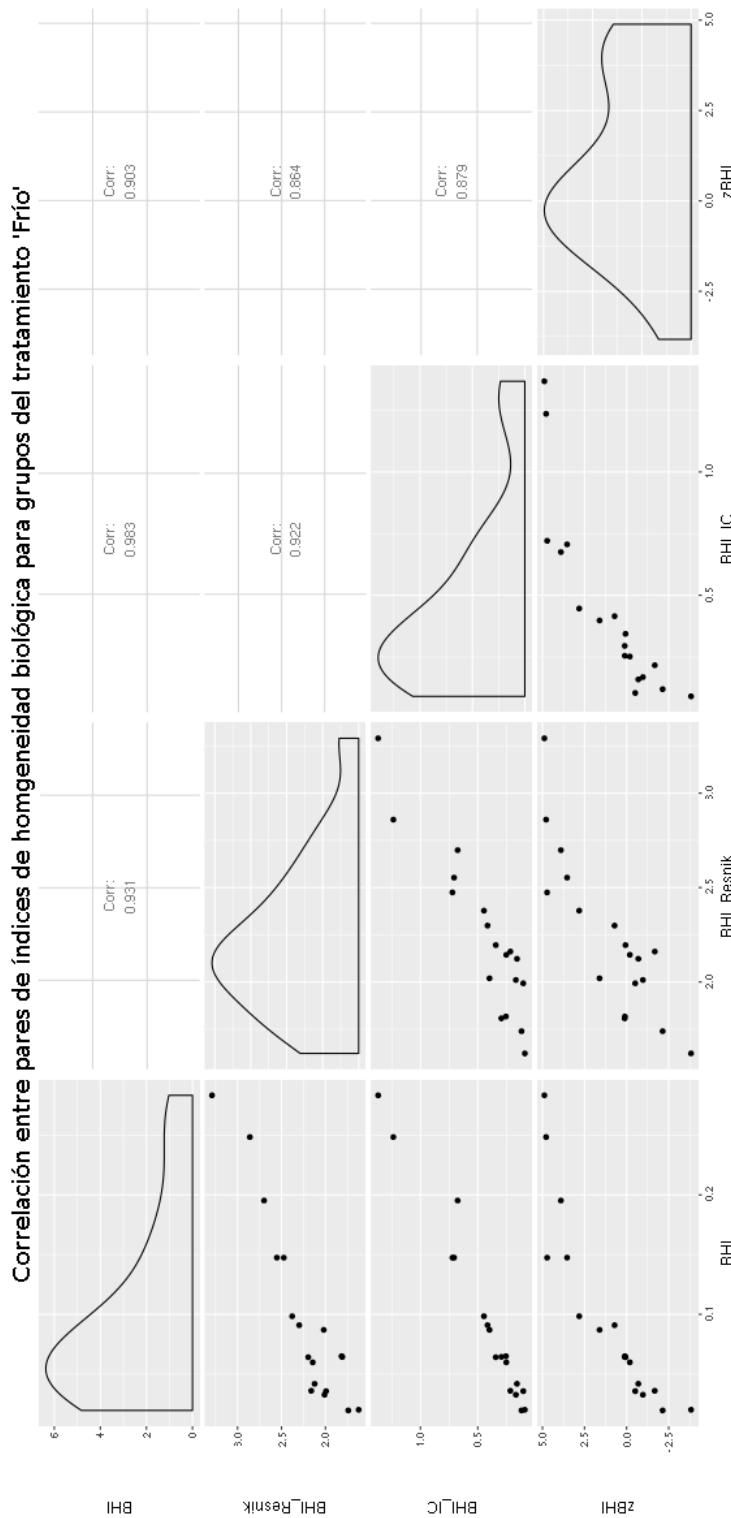
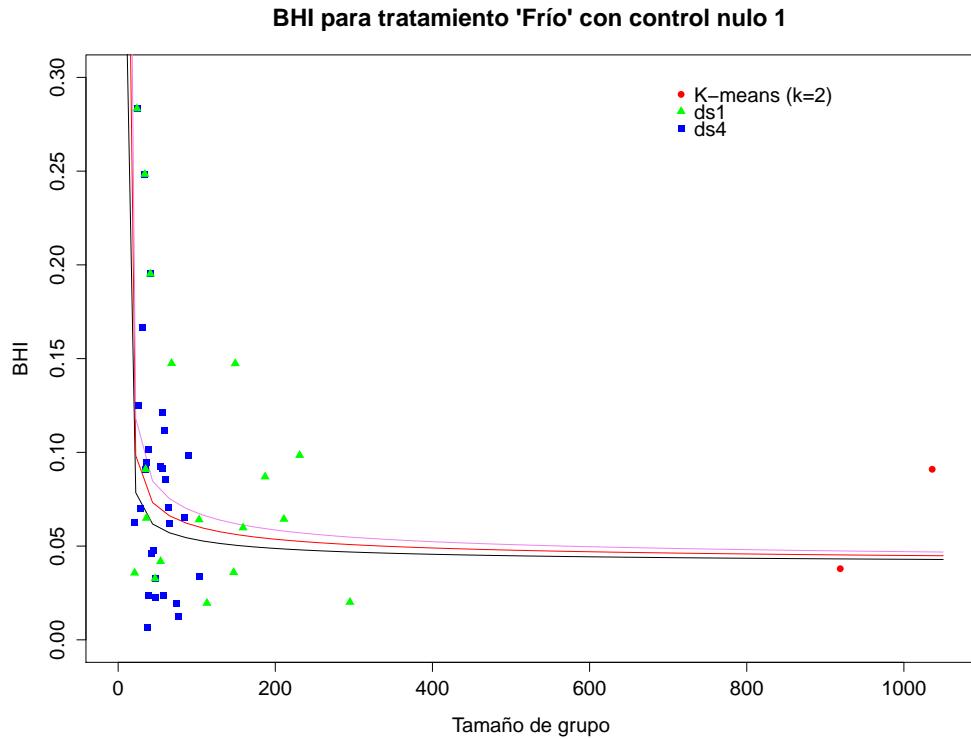


Figura 5.2: Correlación de a pares para los distintos índices de homogeneidad biológica presentados para cada uno de los grupos del tratamiento 'Frío' obtenidos con $ds = 1$. Se observa que todos los índices tienen una alta correlación entre sí.

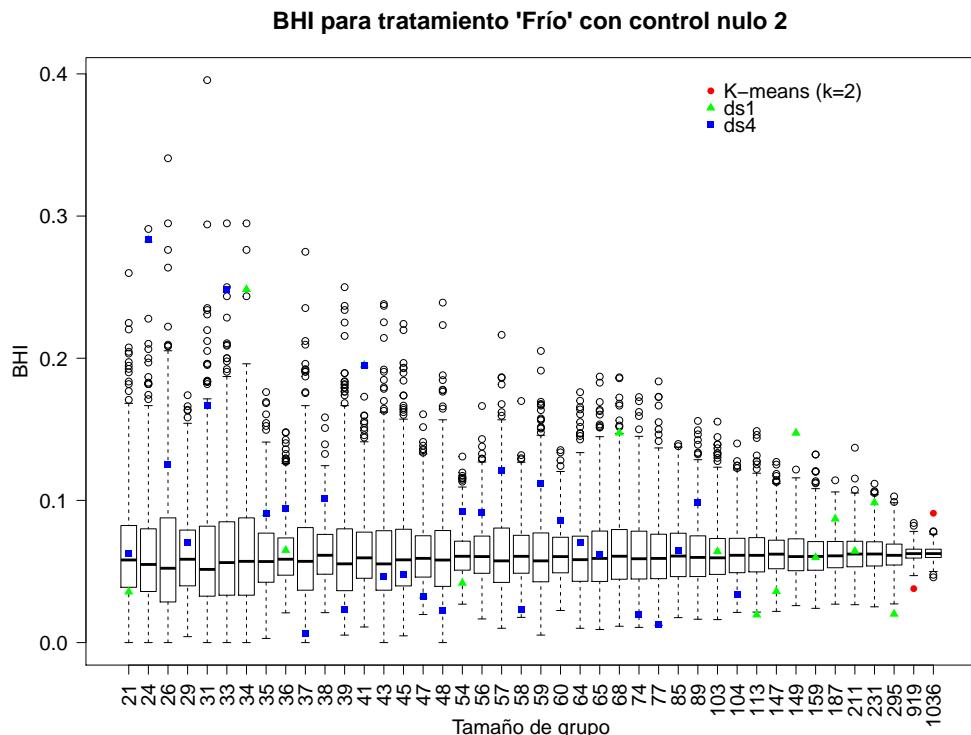
óptimas, ya que por lo general el BHI no es superior al del control nulo.

A pesar de que en el análisis de estructura de los grupos obtenidos por medio de los métodos k-means, *ds1* y *ds4* encontramos que todos los métodos producen particiones altamente coherentes, el análisis de BHI indica que las particiones halladas no pueden ser fácilmente interpretadas a la luz del conocimiento biológico almacenado en GO.

En el próximo capítulo buscaremos cuantificar la coherencia entre los espacios de expresión genética y de conocimiento biológico desde una perspectiva diferente: desde la métrica en lugar de desde las agrupaciones.



(a) *BHI para cada uno de los grupos del tratamiento 'Frío' obtenidos con kmeans, ds1 y ds4 y control nulo 1. En rojo, el valor medio del BHI, en violeta una desviación estándar por sobre el valor medio y en negro una desviación estandar por debajo del mismo.*



(b) *BHI para cada uno de los grupos del tratamiento 'Frío' obtenidos con kmeans, ds1 y ds4, y control nulo 2.*

Figura 5.3: Índice de Homogeneidad Biológica, BHI, para cada uno de los grupos del tratamiento 'Frío' obtenidos con kmeans, ds1 y ds4 y controles nulos.

Capítulo 6

Coherencia entre métrica transcripcional y espacio de conocimiento GO

En los capítulos precedentes cuantificamos por medio de diversos índices la congruencia biológica de los grupos encontrados en el espacio de expresión génica. En este capítulo buscaremos cuantificar la coherencia entre los espacios de expresión génica y de conocimiento biológico desde una óptica diferente: desde la métrica en lugar de desde las agrupaciones.

6.1. Alineamiento de núcleo-objetivo

Una matriz de núcleo o matriz de Gram o matriz de kernel K puede ser pensada informalmente como una matriz semidefinida positiva de similaridad de a pares entre puntos de un conjunto de datos. Para un conjunto de datos $\{x_1, \dots, x_m\}$ esta noción de similaridad esta dada en términos de una función k llamada kernel tal que:

$$K = K_{ij} = k(x_i, x_j) \quad (6.1)$$

con $i, j = \{1, \dots, m\}$ y $k : \mathbb{R}^m \times \mathbb{R}^m \Rightarrow \mathbb{R}$. Una función $k(x, y)$ es un kernel si y solo si para cualquier conjunto finito de datos $C = \{x_1, \dots, x_m\}$ y para cualquier conjunto $\{a_1, \dots, a_m\} \in \mathbb{R}^m$ se tiene que:

$$\sum_{i,j=1}^m a_i a_j k(x_i, x_j) \geq 0 \quad (6.2)$$

Se puede demostrar que esto implica que K debe ser semidefinida positiva (SDP), es decir, $K = \sum_i \lambda_i v_i v_i'$, con $\lambda_i \geq 0$ los autovalores de la matriz K y v_i sus autovectores.

Intuitivamente, un kernel es una transformación que mapea pares de puntos en un espacio de alta dimensionalidad a un indice de similaridad entre los mismos mediante el uso de un producto interno.

Existen multiplicidad de kernels disponibles y para cada aplicación será necesario encontrar el adecuado.

Es de esperar que si es posible extraer información biológica del espacio de expresión genética,

entonces dos puntos que son similares (en algún sentido a definir por el kernel elegido) en el espacio de expresión, también lo sean en el espacio GO (nuevamente, en algún sentido a definir por el kernel elegido). Para cada espacio habrá que definir un kernel adecuado.

Una forma de cuantificar la similaridad entre estos dos espacios es mediante una cantidad conocida como alineamiento núcleo-objetivo o KTA. El KTA de un kernel k_1 con respecto a un kernel k_2 del conjunto C está definido como:

$$\hat{A}(S, k_1, k_2) = \frac{\langle K_1, K_2 \rangle_F}{\sqrt{\langle K_1, K_1 \rangle_F \langle K_2, K_2 \rangle_F}} \quad (6.3)$$

Donde $\langle K_1, K_1 \rangle_F = \sum_{i,j=1}^m K_1(x_i, x_j)K_2(x_i, x_j)$ es el producto de interno de Frobenius entre matrices y K_i son las matrices de kernels simétricas y semidefinida positivas de los espacios a comparar. Este índice tiene un rango entre $[0, 1]$. [49]

Es posible extender este concepto a matrices simétricas indefinidas (no SDP) S mediante diversas técnicas que consisten en transformar S para obtener una S' SDP. Esto es relevante para nosotros, porque matrices de similaridad basadas en similaridad semántica no suelen ser semidefinidas positiva. La que utilizaremos en este trabajo se conoce como *corrimiento del espectro*. Si S es simétrica entonces admite una descomposición en autovalores y autovectores tal que $S = U\Lambda U^T$ con U una matriz ortogonal y Λ una matriz diagonal de autovalores reales, es decir, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$. Entonces, el corrimiento del espectro consiste en correr todo el espectro de S por el mínimo necesario:

$$S_{corrida} = U(\Lambda + |\min \lambda_{min}(S), 0|I)U^T \quad (6.4)$$

Decidimos utilizar este método porque el mismo solo aumenta las autosimilaridades, sin modificar la similaridad entre dos puntos distintos, preservando la estructura de grupo al agrupar datos no necesariamente métricos. [50]

Notar que esta medida es una medida global, ya que toma en cuenta todas las similaridades para calcular KTA.

6.2. Espacio de expresión y GO

Para cuantificar la coherencia métrica entre el espacio de expresión de cada tratamiento y las ontologías GOBPA, GOBPB y GOCC, utilizamos como kernel de espacio de expresión, K_x , la similaridad derivada de la correlación:

$$k_x(g_i, g_j) = \left(\frac{\text{correlacion}(g_i, g_j) + 1}{2} \right) \quad (6.5)$$

con g_i y g_j genes pertenecientes al tratamiento en cuestión.

Para el kernel del espacio de ontologías, utilizamos la similaridad definida en la ecuación ?? y transformamos la matriz en SDP por medio de ???. La matriz se construyó tomando en cuenta todos los genes del tratamiento. Si un gen del tratamiento no se encontraba anotado en la ontología, se lo anotaba al nodo raíz y por lo tanto su similaridad con el resto de los genes era cero. Se calculó entonces para cada tratamiento y cada ontología, el KTA y se

construyó además un control nulo de tipo 2, realizando 1000 reordenamientos aleatorios de las etiquetas de la matriz K_x .

Las figuras ??, ?? y ?? presentan en un boxplot, para cada tratamiento y cada ontología, el KTA del control nulo, junto con un punto rojo para el KTA de expresión.

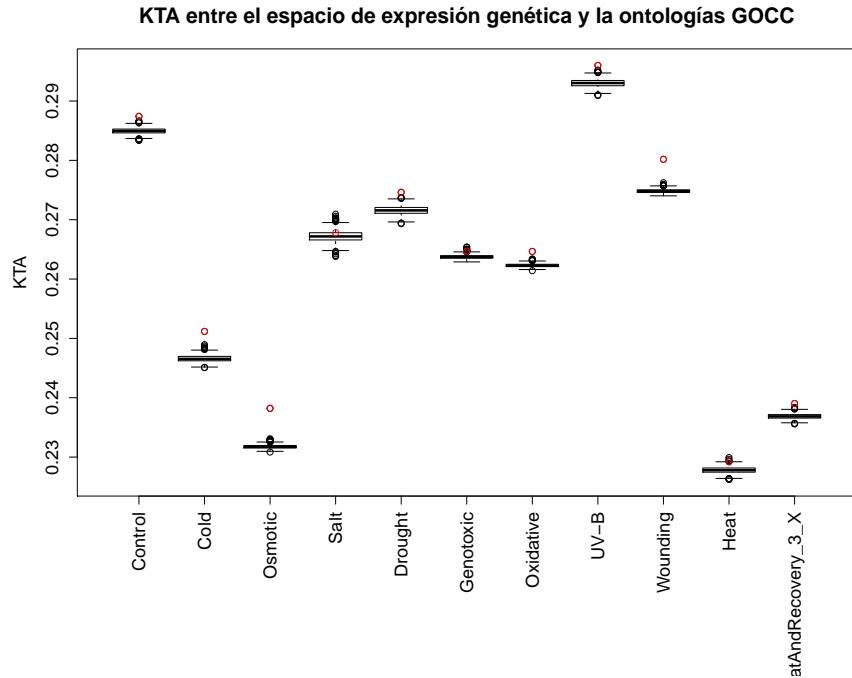


Figura 6.1: KTA para distintos tratamientos entre espacio de expresión y ontología CC.

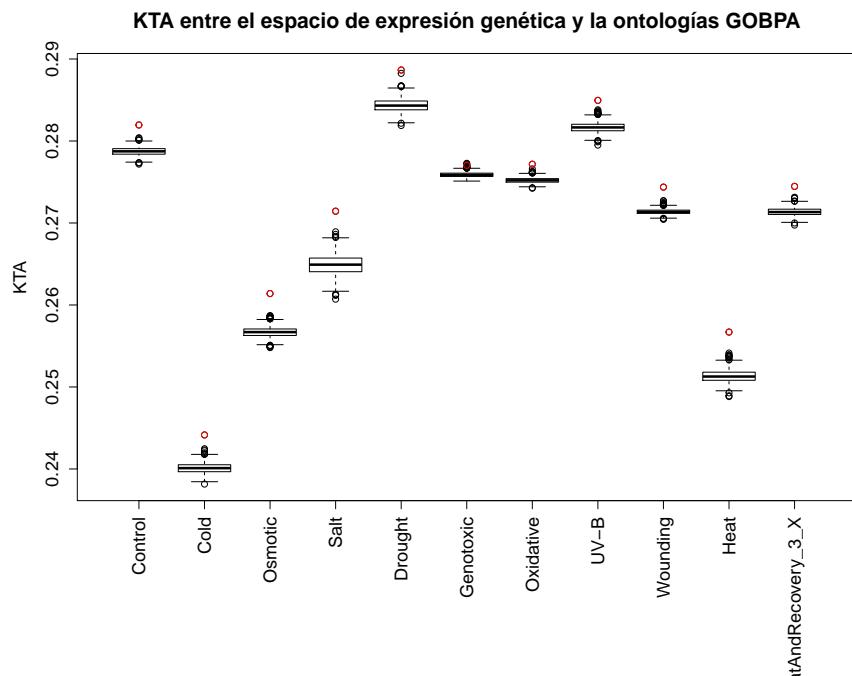


Figura 6.2: KTA para distintos tratamientos entre espacio de expresión y ontología BPA.

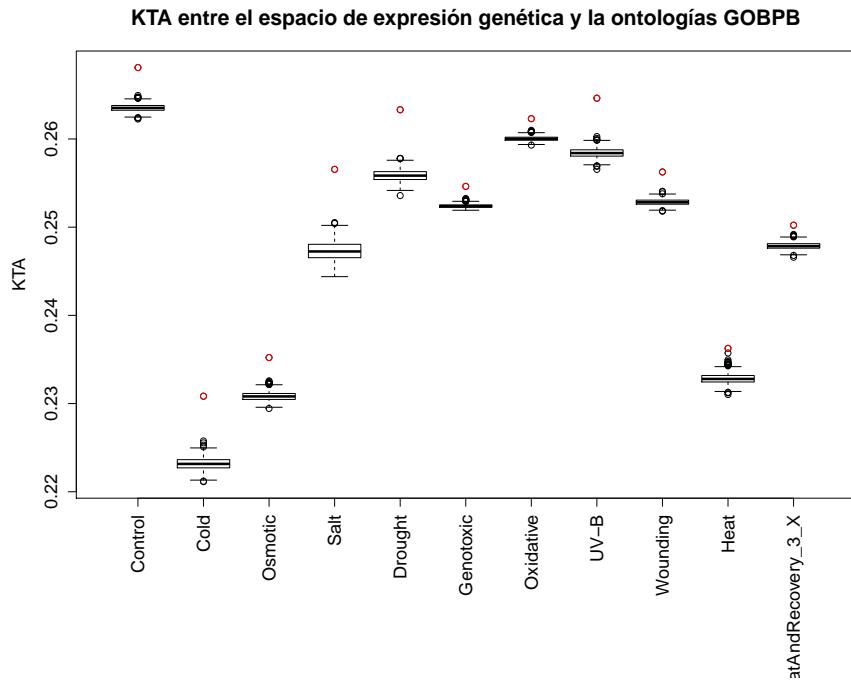


Figura 6.3: KTA para distintos tratamientos entre espacio de expresión y ontología BPB.

En todos los casos encontramos que el KTA de expresión supera todos los valores del KTA de control nulo, lo que soporta la idea de que existe un grado de coherencia no trivial, de orden global, entre la métrica de expresión transcripcional y la métrica del espacio de conocimiento biológico, inferido a partir de la estructura de GO.

6.3. Alineamiento de núcleo-objetivo local

En la sección anterior presentamos una forma de cuantificar globalmente la coherencia entre la métrica transcripcional y el espacio GO mediante el índice KTA. Es posible redefinir este índice para obtener una medida de alineamiento de estos dos espacios pero de forma local, en vecindades transcripcionales. La idea de esta nueva medida apunta a cuantificar e identificar regiones acotadas del espacio transcripcional (i.e. grupos de genes) donde exista un alineamiento espacio-transcripcional vs. espacio-GO diferente a la media. Este índice nos será de gran utilidad en un análisis posterior para caracterizar la congruencia biológica en grupos de genes de similar comportamiento transcripcional y para definir una métrica mixta, que permita definir agrupamientos de genes que presente alto grado de coherencia no solo a nivel transcripcional, sino también en el espacio de conocimiento biológico.

Para definir vecindades transcripcionales consideramos redes de k-primeros-vecinos-mutuos, donde cada nodo es un gen y cada arista tiene un peso w_{ij} entre [0, 1] dado por la similaridad de correlación entre los dos genes g_i y g_j que son unidos por esa arista. Construimos una red $k = 30$, que llamamos $30kmnn$, y estudiamos su topología ¹.

¹Exploramos diferentes alternativas en este punto que incluye el uso de redes de k primeros vecinos en lugar de k primeros vecinos mutuos, y de un número variable de primeros vecinos. Encontramos que la red $30kmnn$ presentada en el trabajo era un buen compromiso entre cobertura y tiempo de computo.

6.3.1. Vecindades transcripcionales

Como primer paso para entender alguna de las particularidades y limitaciones de nuestra caracterización de las vecindades transcripcionales estudiaremos a la red $30kmnn$ mediante dos observables topológicos, la distribución de grado y la intermediación central o *betweenness centrality*.

Además, compararemos la red con dos modelos de redes aleatorias, el modelo de Erdős Renyi, que consiste en construir una red aleatoria con N nodos, conectando pares elegidos al azar y omitiendo múltiples conexiones entre dos mismos nodos hasta alcanzar una cantidad total de K aristas impuestas a priori [51], y el modelo configuracional, que genera una red aleatoria a partir de un recableado de conexiones que mantiene fijo el grado de cada nodo, es decir, no altera la distribución de grado $P(k)$ de la red original [52].

Distribución de grado

El grado k_i de un nodo de la red es la cantidad de primeros vecinos que tiene el nodo. La distribución de grado $P(k)$ es entonces la probabilidad de un que nodo i tomado al azar tenga grado k .

La figura ?? muestra la distribución de grado en escala logarítmica de la red de 30 primeros vecinos mutuos y un modelo nulo de Erdős-Rényi con la misma cantidad de nodos y aristas para el tratamiento frió (1951 nodos y 18436 aristas). Se observa que el grado máximo que alcanzan estas distribuciones está relacionado directamente con el k utilizado, ya que a lo sumo un nodo tendrá k primeros vecinos (k aristas). El hecho de tener un gran número de nodos con k máximo sugiere que ya alcanzamos la “saturación” en una representación basada en interacciones de vecinos mutuos. Vemos de la figura que la distribución obtenida se aleja de la Poissoniana correspondiente al modelo nulo Erdős-Rényi.

Intermediación central o betweenness centrality

La longitud de un camino entre dos nodos se define como la cantidad de aristas que se recorren para llegar de un nodo al otro. El camino (o caminos) más corto es aquél camino cuya longitud es la menor entre todos los caminos. La longitud de un camino más corto se conoce como distancia geodésica.

El betweenness centrality de un nodo i es igual a la cantidad de caminos más cortos desde todos los nodos a todos los otros nodos que pasan por el nodo i . Es una medida de la influencia del nodo i en la red, ya que un nodo con alto betweenness centrality recibirá una gran parte de la carga de la red, suponiendo que la carga se distribuye a través de los caminos más cortos.

Muchas redes reales presentan algunos nodos de alta conectividad, llamados conectores o *hubs*, por donde pasa la mayor parte de la carga de la red.

La figura ?? presenta en escala logarítmica la relación entre el betweenness y el grado k de los nodos de la red $30kmnn$ obtenida para el tratamiento “Frío” (1951 nodos, 18436 aristas). En círculos grises, la de Erdős-Rényi, en triángulos amarillos y la configuracional, en rombos

rojos.

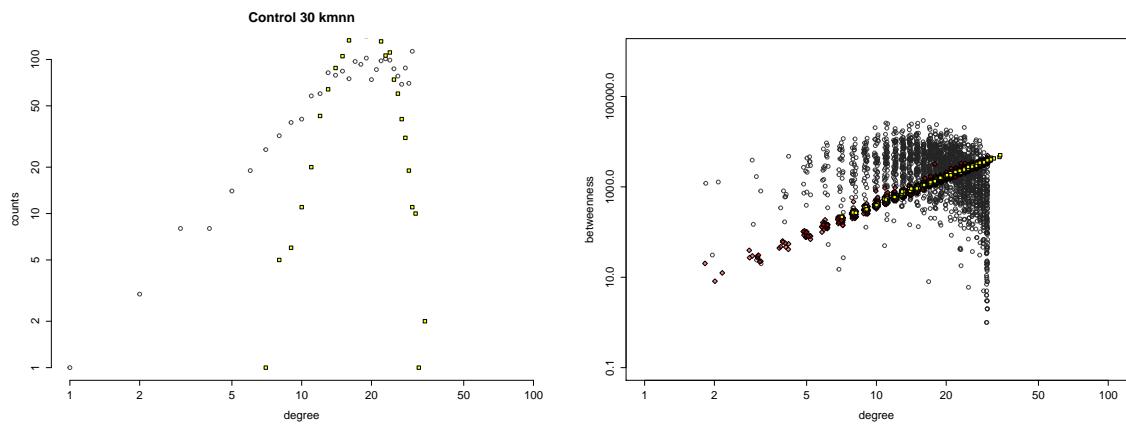
En todos los casos se observa una correlación positiva entre el betweenness y el grado de la red, pero cabe destacar que para un grado fijo k , la red $30kmnn$ presenta una mayor dispersión en el betweenness que las otras. Existen en la red $30kmnn$ nodos de grados intermedios con altos valores de betweenness, a la vez que hubs con bajísimos valores de betweenness. Esto es compatible, por ejemplo, con la existencia dentro de la red de módulos, con hubs locales, interconectados entre si mediante nodos de grado medio-bajo. En todo caso, estos resultados sugieren la existencia de estructura y patrones de conectividad no triviales en la red $30kmnn$.

6.3.2. KTA local red $30kmnn$

La red $30kmnn$ tiene embebida en su estructura de interconectividades la información sobre los entornos locales de cada gen. Para obtener una caracterización del grado de coherencia local entre espacios transcripcionales y GO, para cada arista de la red $30kmnn$ de cada tratamiento, encontramos su vecindario local a primeros vecinos (los n_x primeros vecinos de los nodos unidos por la arista) y construimos la matriz de similaridad de correlación reducida, consistente en la similaridad de correlación entre esos nodos y sus primeros vecinos. Generamos además la matriz de similaridad semántica reducida usando únicamente los genes anteriores, para cada una de las ontologías GOBPA, GOBPB y GOCC (aquellos genes que no estaban anotados en las ontologías, fueron anotados en la raíz de cada una respectivamente). Nos interesa estimar la coherencia media en el espacio GO de la vecindad de cada arista de la red transcripcional. La vecindad de una arista o simplemente vecindad se define como el conjunto de nodos formado por los dos nodos unidos por la arista junto con los primeros vecinos de cada nodo. Llamaremos n_x a la cantidad de nodos vecinos en la red transcripcional que tienen los dos nodos unidos por una arista, y llamaremos n_y a la cantidad de nodos vecinos en la red transcripcional que tienen los dos nodos unidos por la arista y que además esos nodos están anotados en la ontología. De la figura ?? se desprende que n_y depende linealmente de n_x .

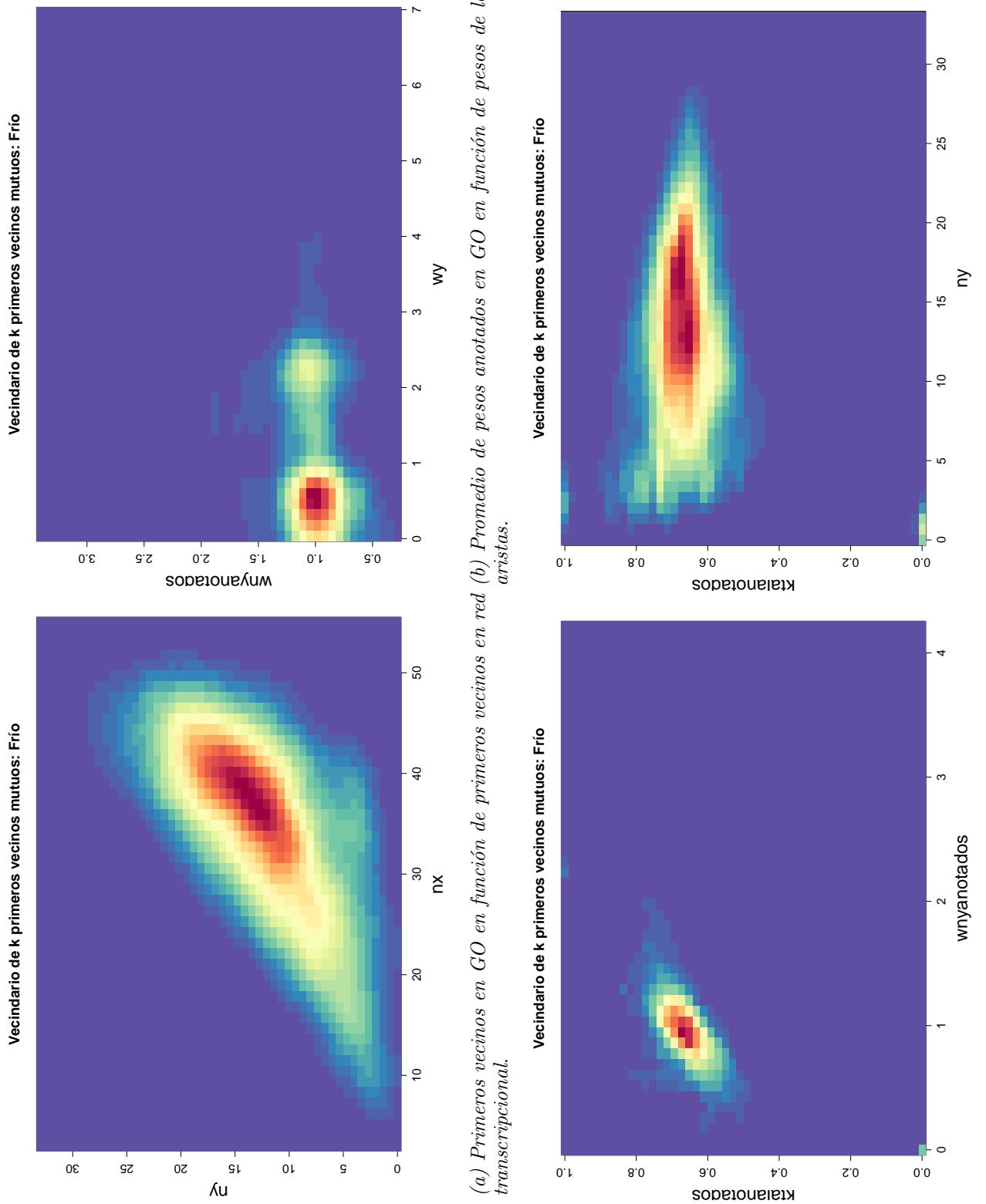
Por otro lado, para estimar el comportamiento de una vecindad, podemos tomar el promedio de los pesos en GO de las aristas en una vecindad, cantidad que llamaremos wyn , junto con el promedio de los pesos de las aristas en una vecindad cuyos nodos se encuentran anotados en la ontología, cantidad que llamaremos $wyn_{anotados}$. De la figura ?? observamos que $wyn_{anotados}$ presenta una gran dispersión en su relación con los pesos de las aristas. Esto implica que el comportamiento en una vecindad no puede ser predicho a partir del conocimiento de una arista, sino que es necesario tomar toda la vecindad en su conjunto. A partir de estas relaciones, introduciremos el índice KTA local como el índice KTA aplicado a las matrices reducidas de ambos espacios, previamente transformadas en SDP. Llamamos matrices reducidas a las matrices de similaridad en el espacio transcripcional y en el espacio GO que solo tienen información de una vecindad. Para caracterizar su comportamiento nos interesa observar primero que el comportamiento de este índice depende linealmente de $wyn_{anotados}$, como muestra la figura ??, siendo entonces coherente con el promedio de los pesos de las aristas de su vecindad, lo que da cuenta que la información contenida en la ontología GO está relacionada con el índice KTA local. Por otro lado, de la figura ?? podemos verificar que no existe una dependencia fuerte del KTA local con la cantidad de vecinos anotados, que llamamos $ny_{anotados}$. Por lo tanto, el KTA local presenta valores más altos para vecinos con

alta coherencia biológica, independientemente de la cantidad de vecinos que tenga.



(a) Distribución de grado para los nodos de la red 30kmnn para el tratamiento 'Frío' (círculos vacíos) y para una red de Erdős-Rényi (cuadrados amarillos) con idéntica cantidad de nodos y configuración. (b) Distribución de betweenness centrality en función del grado de la red para los nodos de las redes vacías y para la red 30kmnn para el tratamiento 'Frío', Erdős-Rényi y ambos con idéntica cantidad de nodos y configuración.

Figura 6.4: Distribución de grado y de betweenness centrality para los nodos de la red 30kmnn y modelos nulos



(a) *Primeros vecinos en GO en función de primeros vecinos en red* (b) *Promedio de pesos anotados en GO en función de pesos de las aristas.*

(c) *KTA local de los anotados en función del promedio de pesos* (d) *Vecindario de k primeros vecinos mutuos: Frio*

Figura 6.5: Caracterización de KTA local para tratamiento 'Frio'.

Capítulo 7

Métrica mixta

Como vimos en el capítulo 4 existen diferentes algoritmos de reconocimiento de agrupamientos que pueden ser utilizados para detectar, con diferente grado de granularidad, grupos de genes coexpresados.

Dichos algoritmos reportan conjuntos de genes que muestran un nivel de correlación transcripcional alto, pero pueden diferir notablemente en el grado de coherencia biológica que es posible asociar a dichas estructuras.

En este capítulo utilizaremos los resultados obtenidos en el capítulo 6 para desarrollar un procedimiento que integre la información contenida en el espacio de expresión con la información contenida en el espacio de conocimientos GO, para favorecer la búsqueda de estructuras biológicamente coherentes. La idea principal es la de montar sobre la heterogeneidad topológica embebida en la estructura de aristas de la red inferida a partir de información transcripcional, un desorden adicional modificando el valor de los pesos de dichas aristas en función de la coherencia biológica local hallada en cada caso. De esta manera, pretendemos obtener una métrica mixta que contenga información de ambos espacios.

7.1. Hacia una métrica mixta

Existen muchas formas de mezclar las métricas del espacio de expresión y del espacio GO. Una de las métricas mixtas investigadas previamente por el grupo [53] fue la de mezcla convexa de distancias, donde se define una nueva distancia d_{mix} a partir de las distancias de expresión y GO y un parámetro α que controla la contribución de cada métrica al algoritmo:

$$d_{mix} = \sqrt{\alpha d_x^2 + (1 - \alpha) \tilde{d}_{GO}^2} \quad (7.1)$$

con $\tilde{d}_{GO} = \frac{\langle d_x \rangle}{\langle d_{GO} \rangle} d_{GO}$, de tal manera que coincidan los valores medios de ambas distribuciones de distancia.

Esta métrica busca el consenso a partir del parámetro global α , que parametriza de manera continua una métrica en contraposición con la otra.

Para este trabajo, implementamos una métrica mixta que en lugar de buscar el consenso entre las métricas, penaliza las correlaciones no soportadas por las distintas ontologías e incentiva aquellas que si lo están.

7.1.1. Modulación de heterogeneidades transcripcionales con información biológica

Partiendo de la red de 30 primeros vecinos mutuos, $30kmnn$, con información transcripcional contenida en los patrones de conectividad de la red, buscamos incluir información biológica en el esquema de conexiónado de la misma. Para ello, modificaremos los pesos de las aristas utilizando la información de coherencia biológica presente en KTA, tomando cada tratamiento y calculando el KTA promedio de toda la red y el KTA local de cada arista. Estas dos cantidades, el KTA promedio, que llamaremos KTA_{fondo} , y el KTA local de la arista entre dos nodos i y j , que llamaremos $KTAl_{ij}$, se relacionan en una cantidad que llamaremos *stress* y que utilizaremos como criterio para identificar si debemos penalizar o incentivar una correlación en el espacio de expresión. El *stress* se define como:

$$\text{stress} = \frac{KTAl_{ij}}{KTA_{fondo}} \quad (7.2)$$

Si $KTAl_{ij} > KTA_{fondo}$, se obtiene un *stress* < 1 , mientras que si $KTAl_{ij} < KTA_{fondo}$, obtenemos que *strees* > 1 . Por lo tanto, si realizamos la transformación no lineal de la similaridad de correlación:

$$w_{ij} = \text{simcor}_{ij}^{\text{stress}} \quad (7.3)$$

obtenemos unos nuevos pesos para la similaridad en donde aquellas relaciones que son similares en GO ($KTAl_{ij} > KTA_{fondo}$) serán incentivadas en el espacio de expresión y a la inversa, las que no son similares en GO serán penalizadas en el espacio de expresión. Los valores típicos del *stress* oscilan entre (0,8 y 1,2). Si graficamos la distribución de fuerza o *strength* para los nodos de la red con los pesos modificados, definida como:

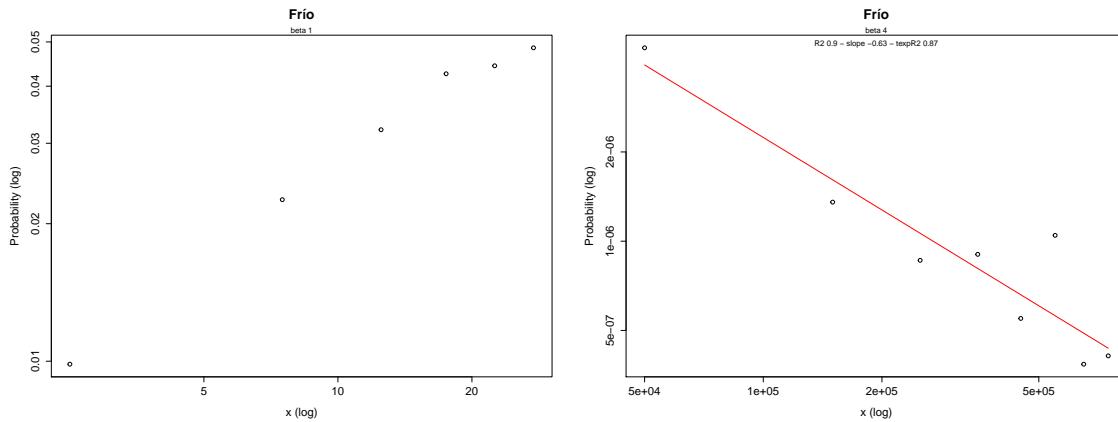
$$S_i = \sum_j w_{ij} \quad (7.4)$$

obtenemos la que se observa en la figura ???. Esta distribución es muy similar a la distribución de grado, $P(k)$, y esto es razonable, ya que los pesos w_{ij} son muy cercanos a 1 debido a la alta correlación entre vecinos de un nodo y los valores típicos de *stress*. Para poder maximizar efectivamente la heterogeneidad a nivel de pesos con información biológica, buscaremos un parámetro β al cual elevar los pesos w_{ij} tal que la distribución de *strength* de la red siga una distribución de tipo ley de potencias, como propone Horvath, en [54].

Buscamos el parametro β adecuado realizando un barrido para β desde 1 hasta 4 con pasos de 1 y desde 15 hasta 65 con pasos de 5, transformando los pesos de la red w_{ij} mediante la ecuación ?? y realizando un ajuste lineal para el logaritmo de la distribución de *strength*. Tomamos el parámetro β que mejor ajustaba cada tratamiento en el sentido de R cuadrado y lo utilizamos para modificar los pesos w_{ij} . Por ejemplo, para el tratamiento 'Frío', obtuvimos $\beta = 4$ como muestra la figura ??.

7.1.2. Heterogeneidades en escala de grupo

La idea de nuestro enfoque es partir de los grupos detectados en el espacio transcripcional y estudiar como se comportan en términos de su coherencia biológica. Para ello pretendemos



(a) Logaritmo de la distribución de strength para la red 30knmm con pesos w_{ij} y $\beta = 1$ del tratamiento 'Frío'.
(b) Logaritmo de la distribución de strength para la red 30knmm con pesos w_{ij} y $\beta = 4$ del tratamiento 'Frío'. En rojo ajuste lineal (R^2 cuadrado = 0.9).

Figura 7.1: Logaritmo de la distribución de strength para la red 30knmm con pesos w_{ij} y distintos β del tratamiento 'Frío'.

analizar heterogeneidades transcripcionales, en la escala interna al grupo, modificando la similaridad de correlación con la información obtenida mediante el índice KTA local en las redes de 30 primeros vecinos mutuos.

Para ello, tomamos cada grupo de cada tratamiento y calculamos el KTA de todo el grupo, además del KTA local de cada arista del grafo 30kmnn correspondiente al grupo. Con estas dos cantidades redefiniremos el *stress* para que sea un *stress* de grupo. La figura ?? presenta la distribución de *stress* para el tratamiento 'Frío' en función del tamaño de grupo. Se observa que el 92 % de los valores de *stress* se ubica entre 0,8 y 1,2, como habíamos mencionado previamente.

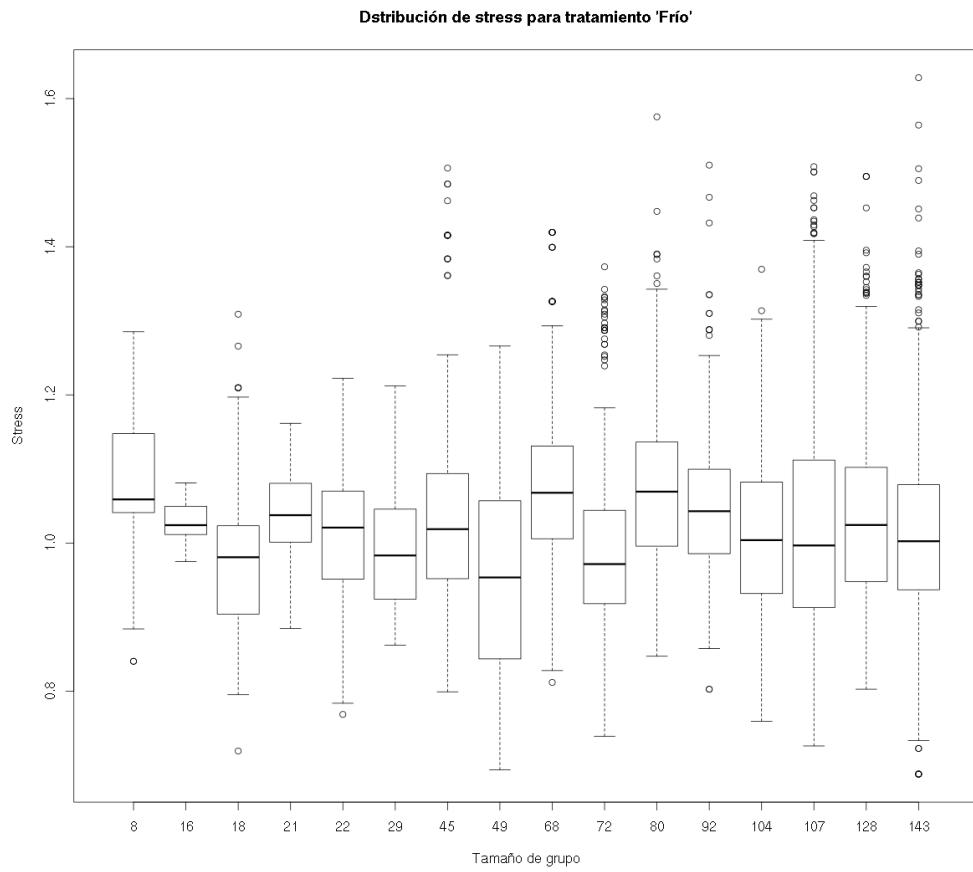


Figura 7.2: Distribución de stress para el tratamiento 'Frio' y ontología CC en función del tamaño de grupo.

Una vez obtenidos los parámetros para la métrica mixta, desarrollamos un método heurístico para poder aplicar esta nueva métrica mixta en el proceso de detección de agrupaciones de genes.

7.2. Método heurístico

El método heurístico desarrollado consiste en tomar cada grupo de una partición realizada previamente con corte de árbol dinámico e intentar encontrar subestructura dentro del mismo a partir de la métrica mixta arriba presentada de tres formas distintas.

La primera forma consiste en aplicar sobre cada grupo nuevamente un corte de árbol dinámico, utilizando la ahora métrica mixta en la confección del dendrograma, es decir: penalizando (incentivando) las conexiones más importantes, que son las que aparecen en la red de k primeros vecinos mutuas con alto (bajo) stress. A este método que llamaremos *lkta.dtc*. Por otro lado, la segunda y la tercera formas, que llamaremos *lkta.infomap* y *lkta.cnm*, consisten en obtener comunidades mediante infomap y cnm, respectivamente, en la red de los genes del grupo original, teniendo en cuenta los pesos modificados por la métrica mixta.

En todos los casos, una vez obtenida una partición del grupo, calculamos el índice BHI de los subgrupos obtenidos que no alcanzaron un nivel de significancia biológica suficiente y volvemos a unir aquellos subgrupos que se encuentren por debajo de una desviación estandar

del control nulo 1 presentado en la sección ???. Esto lo hacemos mediante un algoritmo voraz o greedy, que en cada paso busca los dos subgrupos tales que el BHI de los dos juntos sea superior al BHI de cada uno por separado. Cuando el algoritmo no consigue unir dos subgrupos para mejorar el BHI, se detiene.

Finalmente, todos los subgrupos que todavía quedan por debajo de una desviación estandar del control nulo 1 son unidos entre sí en un único subgrupo.

Además de las 3 metodologías descriptas, utilizaremos un cuarto método, llamado *insideX*, a modo de control. El mismo consiste en volver a particionar cada grupo usando corte de árbol dinámico con *deepsplit* = 4, pero sin cambiar la métrica por la métrica mixta. Esto nos permitirá controlar la efectividad de la métrica mixta para encontrar mayor resolución en las particiones.

Para cuantificar el cambio en la información biológica que brinda la nueva partición, elegimos tomar el valor medio del BHI de los subgrupos, $\langle BHI \rangle$, y el valor medio del BHI de los subgrupos que superan el control nulo, $\langle BHI \rangle_+$. Esta ultima cantidad refleja la situación de aquellas subestructuras del grupo original que se 'beneficiaron' desde el punto de vista del BHI, con la mayor granularidad de la nueva partición.

A modo de ejemplo, presentamos en las figuras ?? y ?? los resultados de esta heurística aplicada a los grupos 2 y 9 del tratamiento 'Frío'. En las mismas, la curva roja indica el valor medio para la distribución del BHI del control nulo 1, la verde una desviación estandar por sobre la media y la negra, una desviación estandar por debajo. El punto lleno representa el grupo original y su BHI. En verde si su BHI supera el control nulo y en rojo si no. Los puntos vacíos en gris representan los subgrupos nuevos. En la leyenda, el primer número es $\langle BHI \rangle$ y el segundo, entre paréntesis, es $\langle BHI \rangle_+$.

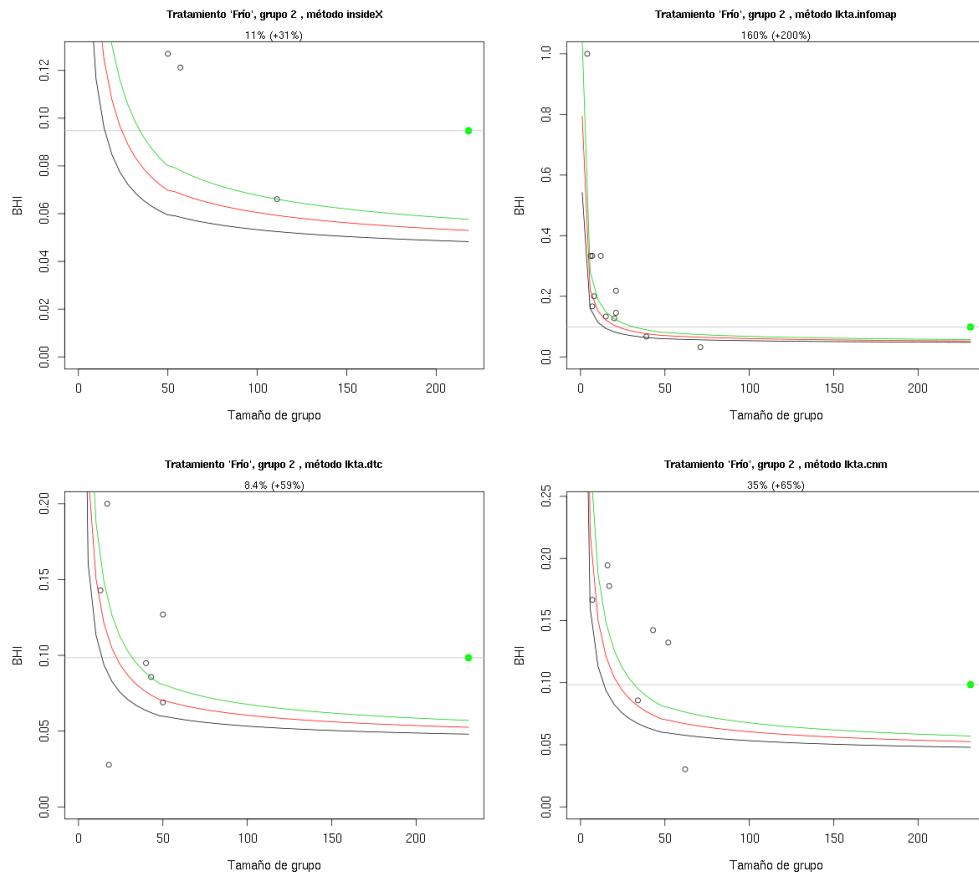


Figura 7.3: Métodos con métrica mixta y control para el grupo 2 del tratamiento 'Frío'.

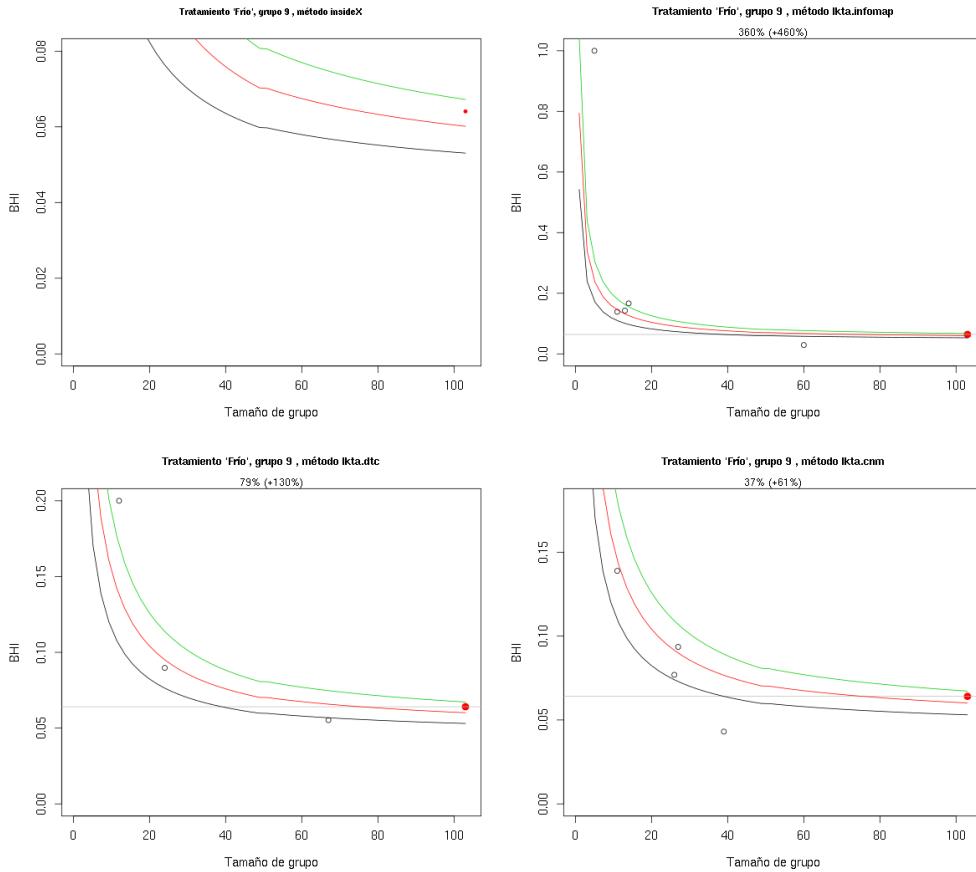


Figura 7.4: Métodos con métrica mixta y control para el grupo 9 del tratamiento 'Frío'.

Se observa que para el grupo 2, un grupo con un BHI por sobre el control nulo y el segundo grupo más grande de la partición, el método *insideX* logró encontrar heterogeneidades y partir el grupo, obteniendo tres subgrupos de más de 50 genes cada uno y una mejora del 11 % en $\langle BHI \rangle$, y una del 31 % para $\langle BHI \rangle_+$. Por otro lado, los tres métodos *lkta* lograron superar la mejora obtenida por *insideX* en $\langle BHI \rangle_+$, y tanto *lkta.infomap* como *lkta.cnm* lograron mejorar además el obtenido en $\langle BHI \rangle$, siendo *lkta.infomap* el que mejores puntajes logró, llegando a obtener incluso un grupo con un BHI de 0,97.

Para el grupo 9, un grupo con un BHI por debajo del control nulo y relativamente pequeño, el método *insideX* no logró encontrar heterogeneidades dentro del grupo y por lo tanto no pudo particionarlo, sin lograr ningún tipo de mejora en el mismo. Sin embargo, todos los métodos *lkta* lograron particionar el grupo y mejorar ambos índices, siendo nuevamente *lkta.infomap* el que logró la mejora más significativa, con un 360 % de incremento en $\langle BHI \rangle$ y un 460 % de incremento en $\langle BHI \rangle_+$.

Finalmente, realizamos este análisis para todos los tratamientos y todos los grupos en cada tratamiento. Para cada método, graficamos el $\langle BHI \rangle$ de los subgrupos en función del BHI del grupo original en la figura ???. En rojo, graficamos una recta de pendiente unitaria. Se observa que todos los métodos que utilizan la métrica mixta consiguieron grupos que no solo presentan altos valores de coherencia transcripcional sino que tienen BHI superiores a los que tenían originalmente. En particular, para *lkta.infomap* el 92 % de los grupos fue dividido en subgrupos que poseen un BHI promedio superior al de su grupo original, mientras que lo

mismo sucede en un 91 % de los grupos para lkta.cnm y un 87 % para lkta.dtc. Finalmente, para el control insideX, solo el 10 % de los grupos mejoró su BHI promedio al ser particionado nuevamente.

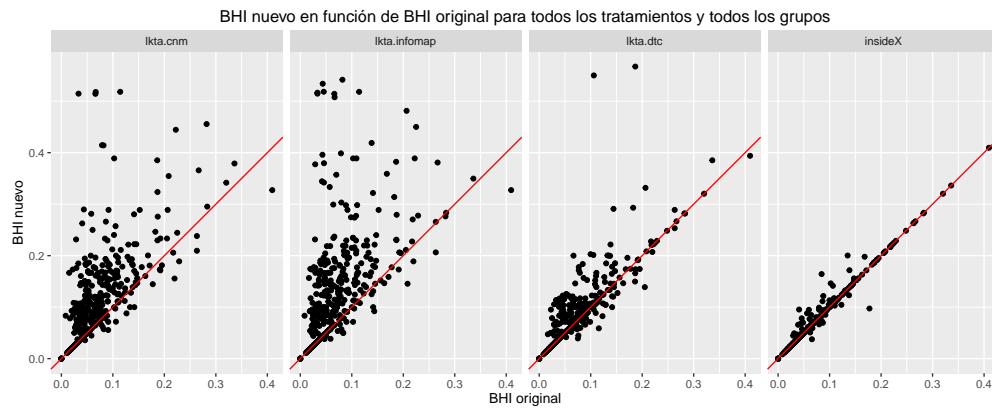


Figura 7.5: BHI nuevo en función de BHI original para todos los tratamientos y todos los grupos.

7.3. Interpretación biológica

Una vez obtenidos los nuevos grupos mediante los distintos métodos, buscaremos ganar conocimiento sobre el significado biológico de cada grupo y en particular estudiaremos si los genes dentro de un mismo grupo participan de algún proceso biológico común. Para ello, recurriremos al análisis de sobreexpresión de conceptos GO dentro de un dado agrupamiento de genes.

El mismo consiste en considerar un grupo de genes y analizar si existen categorías de una ontología (BP u CC), sobrerepresentadas en el mismo, utilizando una prueba de Fisher exacta con un p-valor $< 10e - 4$.

Ilustraremos los resultados de esta tipo de análisis considerando los grupos 2, 6 y 13 del tratamiento 'Frío'. En la figura ?? se observan los grupos originales y los subgrupos obtenidos por medio de cada método junto con el control nulo 1 y los correspondientes valores de BHI, mientras que en las figuras ??, ?? y ?? se visualizan los resultados de la prueba exacta de Fisher para cada uno de estos grupos.

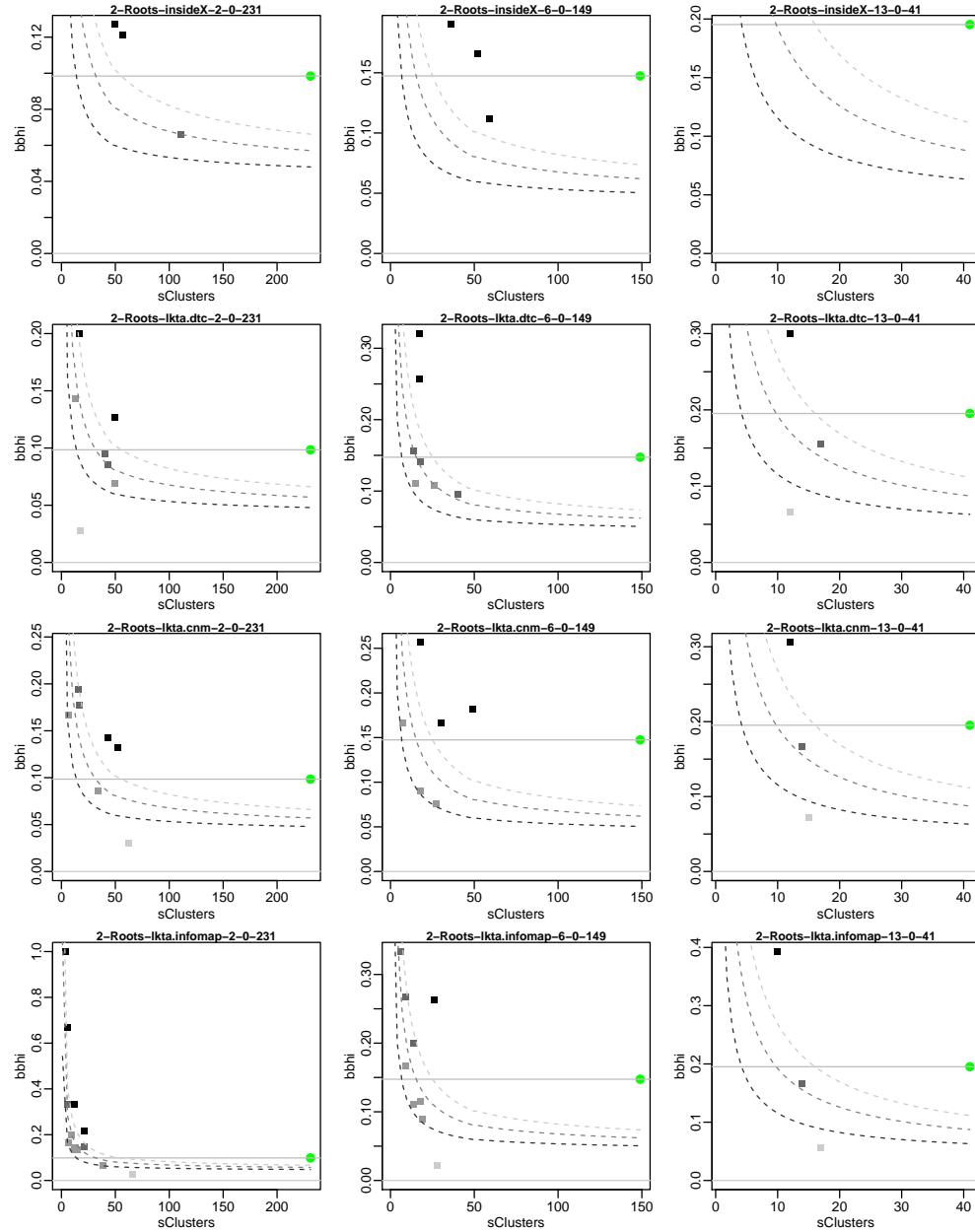


Figura 7.6: En la primera segunda y tercera columna mostramos los resultados del análisis para clusters transcripcionales 2, 6 y 13 del tratamiento Frio respectivamente. La primera, 2da y 3era fila... dtc, cnm, infomp Circulo cluster original. Cuadrados son subgrupos. Nivel de grisado depende de alejamiento de curvas de control. Las curvas punteadas son 1sigma, 2sigma y 3sigma...etc el grisado de los c.

En el Panel izquierdo: mapa de colores o *heatmap* de p-valores donde las filas representan categorías de la ontología GOBP, y las columnas subgrupos del grupo transcripcional 2 detectados por la metodología insideX y las 3 variantes lkta.

El color corresponde al nivel de significancia obtenido en el test de sobrerepresentación

Para una mejor interpretabilidad se utilizaron los 50 conceptos GO relevantes (sobrerepresentados) más específicos y se ordenaron según un dendrograma generado a partir de la similaridad semántica de los mismos.

Se incluye en el panel central detalle sobre el contenido de información, el nivel de profundidad

dad en la jerarquía de la ontología y la descripción de los conceptos GO respectivos. En el panel derecho: conceptos GO de menor especificidad para mejorar la interpretabilidad.

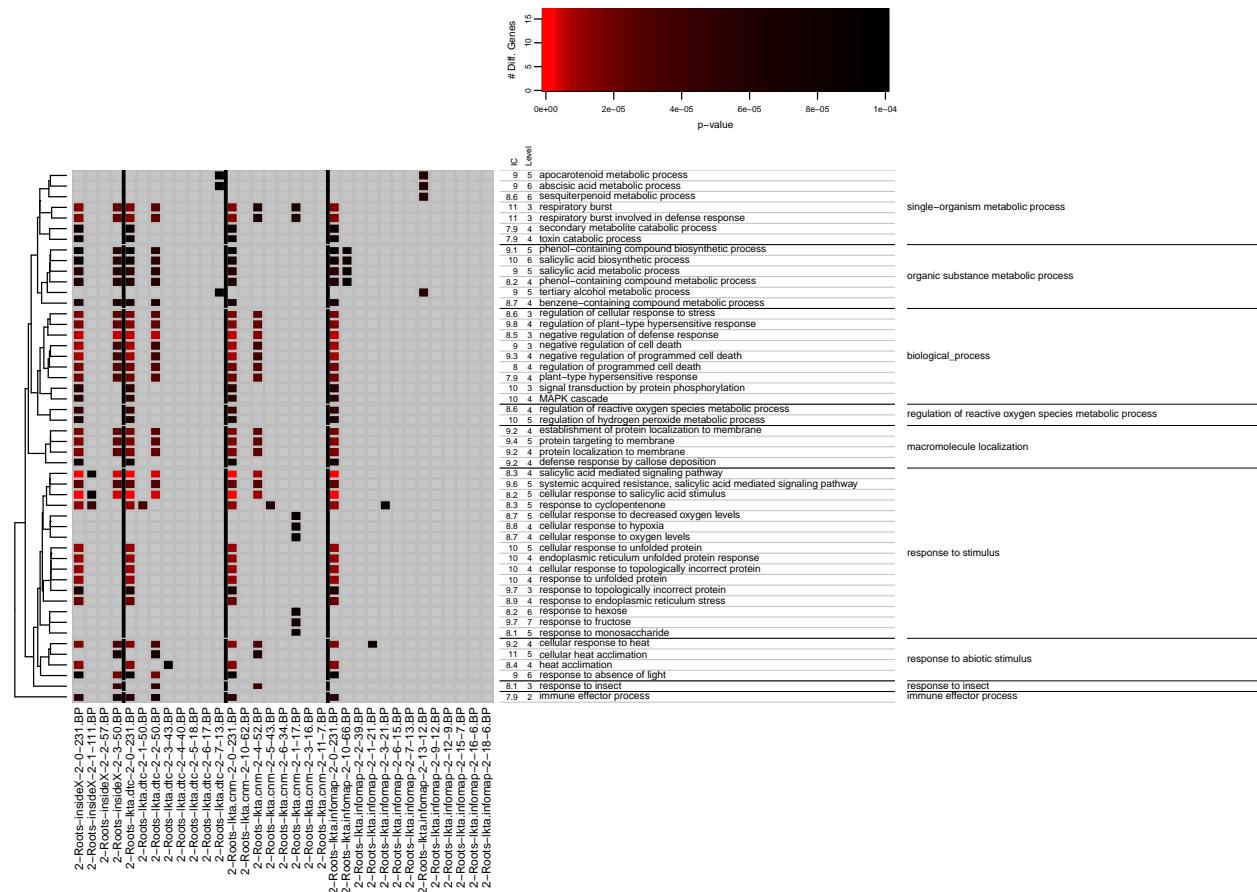


Figura 7.7: Resultado del test de sobreexpresión de nodos GO para el grupo 2 del tratamiento 'Frío' para los métodos con métrica mixta y control.

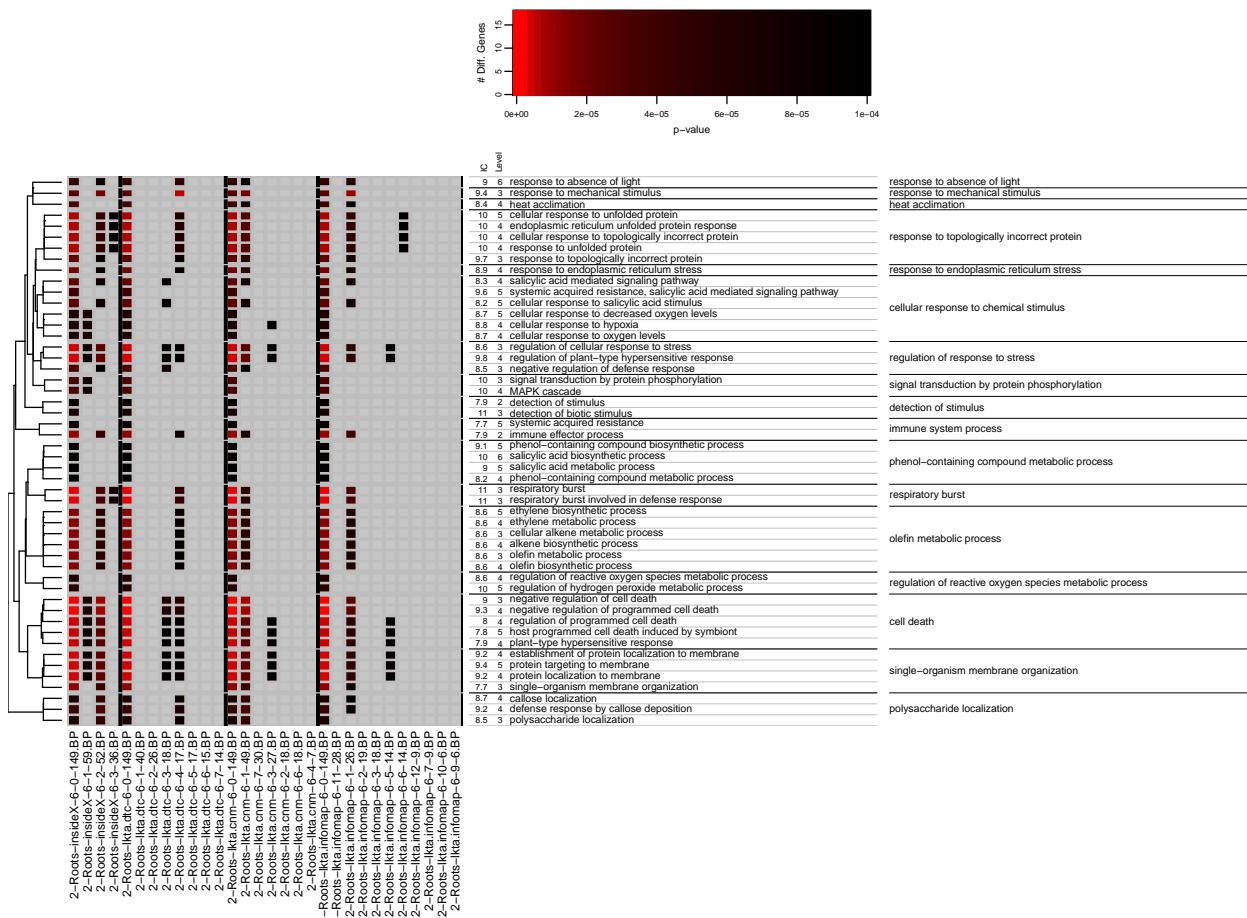


Figura 7.8: Resultado del test de sobreexpresión de nodos GO para el grupo 6 del tratamiento 'Frío' para los métodos con métrica mixta y control.

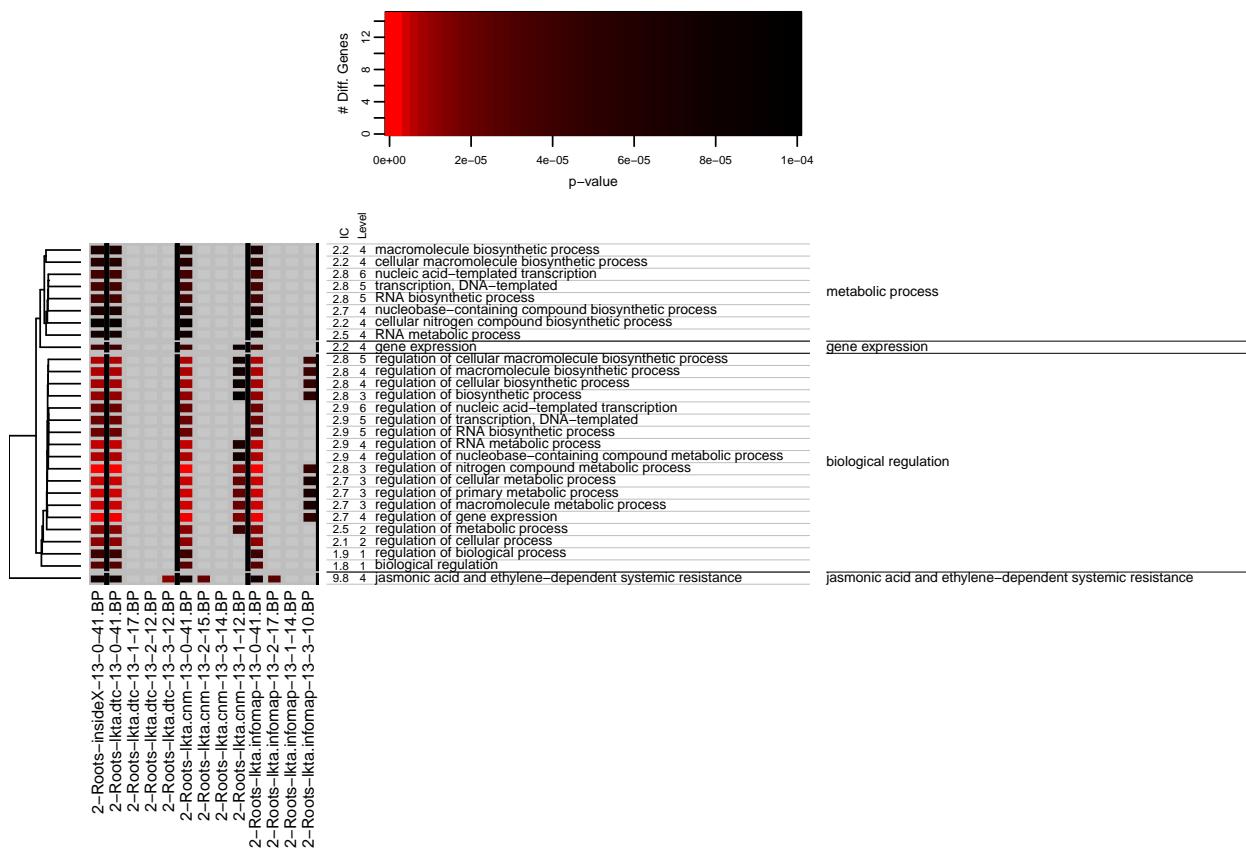


Figura 7.9: Resultado del test de sobreexpresión de nodos GO para el grupo 13 del tratamiento 'Frío' para los métodos con métrica mixta y control.

Para el grupo 2, el segundo grupo más grande de la partición, de 231 genes, y el método insideX, se observa que uno de los dos subgrupos con ganancia en BHI, el de 50 genes, explica casi en su totalidad los procesos biológicos del grupo original, mientras que no se encontró en el otro subgrupo con ganancia nodos sobreexpresados. Por otro lado, para lkta.dtc y lkta.cnmse se obtuvo un subgrupo de 50 genes con ganancia en BHI que explica los procesos del grupo original mientras que el resto de los subgrupos con ganancia no realiza aportes significativos al conocimiento biológico. Finalmente, para lkta.infomap, el de mayor ganancia en el índice BHI, ninguno de los subgrupos logra explicar los nodos sobreexpresados.

Parecería que dentro del grupo transcripcional original (231 genes) un subgrupo de 50, detectado por 2 de las variantes lkta y la puramente transcripcional (insideX) es el responsable del enriquecimiento. El hecho de que insideX lo haya detectado sugiere que aun heterogeneidades puramente transcripcionales, puestas de manifiesto en la escala del grupo 2, correlacionan con la congruencia biológica del subgrupo.

Es interesante notar que en este caso la alta granularidad de la partición infomap parece excesiva en términos de interpretación biológica.

Para el grupo 6, un grupo de tamaño intermedio, con 149 genes, el método insideX consigue nuevamente dos subgrupos con ganancia en BHI, de 59 y 52 genes respectivamente, que explican en conjunto la mayor parte de los nodos sobreexpresados del grupo original. Para lkta.dtc se obtienen dos subgrupos informativos de tamaño pequeño, de 18 y 17 genes respectivamente. Para lkta.cnm y lkta.infomap se obtuvieron un subgrupo informativo de 49 y otro de 26 genes respectivamente.

En este caso los 3 métodos lkta detectan estructuras que se pueden asociar a los mismos conceptos biológicos.

La resolución de la partición insideX parece detectar estructuras más grandes, pero aun permite identificar conceptos biológicos similares en relación a los subgrupos encontrados.

Finalmente, para el grupo 13, un grupo pequeño de 41 genes, el método insideX no logró producir subgrupos, mientras que solo para lkta.cnm y lkta.infomap se consiguieron grupos explicativos para un grupo de 12 y otro de 10 genes respectivamente.

En este caso, solo el uso de la metodología lkta pudo detectar subestructura relevante.

Cabe notar que en este caso, se trata de conjuntos muy chicos de genes, por lo que es de esperar una mayor divergencia entre la estimación de congruencia biológica dada a partir de significancia de test de sobrerepresentación y la asociada al índice BHI (que solo exige coincidencia en un grupo, sin importar su especificación).

En líneas generales, observamos que en varios casos una gran parte de los nodos sobreexpresados del grupo original podían ser explicados por uno o dos subgrupos de tamaño reducido, lo que implica que existen efectivamente heterogeneidades en los grupos originales que no podía detectar el método de agrupamiento inicial *ds1*. Por otro lado, las estructuras detectadas fueron de tamaños intermedios, de alrededor de 50 o 60 genes. Esto da un indicio de que esta debería ser la resolución o escala a la cual deberíamos trabajar con los métodos de agrupamiento. Finalmente, este análisis nos permite ver que si bien lkta.infomap genera las particiones de grupos con mayor ganancia en el índice BHI, los subgrupos obtenidos no necesariamente brindan información biológica. Esto podría deberse al tamaño tan reducido de los subgrupos, que no facilita su interpretación en términos de GO.

Capítulo 8

Conclusiones y perspectivas

La tecnología de micromatrices de ADN se constituyó como una herramienta indispensable para el monitoreo de niveles de expresión a lo largo de todo el genoma de un organismo. Las misma permiten estudiar los perfiles de expresión transcripcionales de miles de genes simultáneamente, estimando la concentración de ARNm que está siendo exportado desde el núcleo celular hacia el citoplasma para la síntesis de determinadas proteínas.

Mediante técnicas de agrupamiento no supervisado de datos es posible encontrar grupos de genes con perfiles de expresión correlacionados, y en particular, utilizando la distancia de correlación lineal considerada en ??, es posible buscar estructuras linealmente correlacionadas. La importancia desde el punto de vista biológico de este tipo de análisis radica en que grupos de genes correlacionados son un primer indicador de que dichos genes forman parte de un mismo proceso biológico.

Distintos métodos de agrupamiento darán como resultado particiones diferentes en función de la resolución o granularidad que puede alcanzar cada método. Encontrar la resolución óptima es entonces de vital importancia para lograr grupos biológicamente coherentes. En esa linea, caracterizamos dos métodos de agrupamiento, el método k-means y el método corte de árbol dinámico.

El método k-means aplicado a los datos de cambios transcripcionales en respuesta a estrés abiótico en *A. thaliana* presenta una muy baja resolución obteniendo únicamente grupos del orden de 1000 genes, con una coherencia media del orden de 0,7. Estos grupos tan grandes son de difícil interpretación biológica.

El método corte de árbol dinámico, por otro lado, es de especial interés por la granularidad que logra alcanzar en el análisis de datos de expresión génica. Es posible regular la resolución del método a partir de variar el parámetro *deepsplit*, con un rango entre *deepsplit* = 1 (*ds1*) para particiones más gruesas, y *deepsplit* = 4 (*ds4*) para particiones más finas. Encontramos grupos de entre 20 y 450 genes para *ds1* y grupos entre 20 y 150 genes para *ds4*, con coherencias medias elevadas en todos los casos de entre 0,85 y 0,97. Esta coherencia tan elevada da cuenta de estructura en los grupos, estructura a la que el método k-means no logra acceder. Para poder dar una interpretación biológica a las particiones obtenidas en el espacio de expresión transcripcional nos valimos de la información contenida en la ontología génica GO. Utilizamos el índice densidad de información (ID) para cuantificar el grado en que los genes de una partición comparten anotaciones en GO y además forman parte del mismo grupo en el espacio transcripcional. Medimos además este índice para tres redes distintas, una red de

vías metabólicas (KEGG) y dos redes de interacción de proteínas (AI1 y LCI), junto con el índice para un control nulo. Este índice nos permitió verificar que es posible inferir información biológica tanto de las particiones *ds1* como las *ds4*, siendo las de *ds1* las que obtuvieron mejor ID entre las dos. Fue entonces de particular interés para nosotros caracterizar las particiones obtenidas con *ds1*.

Para ello, cuantificamos la homogeneidad biológica de las particiones *ds1* por medio del índice de homogeneidad biológica (BHI), que reporta, para cada grupo, la máxima proporción de pares de genes agrupados que comparten una misma clase funcional de Ontología Génica. El análisis de BHI indicó que las particiones halladas no pueden ser fácilmente interpretadas a la luz del conocimiento biológico almacenado en GO.

Buscamos entonces cuantificar la coherencia entre los espacios de expresión genética y de conocimiento biológico desde una perspectiva diferente: desde la métrica en lugar de desde las agrupaciones. Para ello, hicimos uso del alineamiento de núcleo-objetivo (KTA), que da una medida en que dos métricas están alineadas globalmente. Utilizamos la similaridad de correlación definida en ?? como medida de similaridad para el espacio de expresión y la similaridad semántica entre dos genes *rcmax* definida en ?? para el espacio GO. Encontramos que para todos los tratamientos el índice KTA era superior al del control nulo, indicando una coherencia global entre la métrica de expresión transcripcional y la métrica del espacio GO. Conociendo esto, buscamos entonces la forma de definir agrupamientos de alta coherencia no solo en el espacio transcripcional sino también en el espacio de conocimiento biológico. Para ello, definimos una métrica mixta a partir de la congruencia biológica de grupos de genes en una vecindad transcripcional (por medio del índice KTA local), y definimos vecindades transcripcionales a partir de una red de 30 primeros vecinos mutuos. Encontramos que el comportamiento del índice KTA local depende linealmente de *wyn_{anotados}*, lo que da cuenta que la información contenida en la ontología GO está relacionada con este índice.

Para obtener una métrica que permita detectar heterogeneidades dentro de los grupos, modificamos la similaridad de correlación a partir del *stress*, definido en ?? y un parámetro β , lo que permite penalizar o incentivar una correlación en el espacio de expresión dependiendo de su similaridad en el espacio GO. A partir de esta métrica desarrollamos tres métodos heurísticos, lkta.dtc, lkta.infomap y lkta.cnm, para buscar estructura dentro de los grupos obtenidos por medio de *ds1*. Además, realizamos un control, llamado *insideX*, utilizando el método *ds4* sobre los grupos de *ds1*. Obtuvimos para todos los casos de lkta una mejora sustancial en el índice BHI de cada partición, sobre todo con lkta.infomap. Con los nuevos grupos, realizamos una prueba de sobrerepresentación por medio de un test exacto de Fisher, y encontramos que de todos los métodos, el método *insideX* obtenía subgrupos de tamaño intermedio, del orden de 50 genes, que podían explicar el contenido biológico del grupo original por si solos, lo que implica la existencia de heterogeneidades no detectadas por *ds1* en el agrupamiento inicial. Por otro lado, el resto solo obtenía subgrupos de este tipo en contadas ocasiones, y más notoriamente para el caso de lkta.infomap, los subgrupos encontrados eran tan pequeños que no se encontraban sobrerepresentados en ninguna categoría significativa de GO. Esto da un indicio de cual debería ser la resolución óptima para trabajar con este tipo de datos.

Proponemos también, como proyección a futuro...que proponemos

Bibliografía

- [1] NATURE.COM. *Functional genomics*. Accedido: 2016-01-13.
URL <http://www.nature.com/subjects/functional-genomics>
- [2] ARABIDOPSIS-INTERACTOME-MAPPING-CONSORTIUM. *Evidence for Network Evolution in an Arabidopsis Interactome Map*. Annual review of plant biology **10** (2013) 161.
- [3] M. KANEHISA. *Yeast Biochemical Pathways. KEGG: Kyoto encyclopedia of genes and genomes*. Nucleic Acids Res **28** (2000) 27.
URL <http://pathway.yeastgenome.org/biocyc/>
- [4] ARABIDOPSIS.ORG. *org.At.tair.db*. <https://www.arabidopsis.org/biocyc/>.
- [5] E. DOMANY. *Cluster Analysis of Gene Expression Data 1* **110** (2003) 1117.
- [6] B. ALBERTS. *Molecular Biology of The Cell*, volume 6 (2015).
- [7] B. BOSE. *In Vitro Differentiation of Pluripotent Stem Cells into Functional B Islets Under 2D and 3D Culture Conditions and In Vivo Preclinical Validation of 3D Islets*. Methods in Molecular Biology (2016) 257.
- [8] M. BABU. *An Introduction to Microarray Data Analysis*. Computational Genomics: Theory and Application (2004) 225.
URL <http://www.mrc-lmb.cam.ac.uk/genomes/madanm/microarray/chapter-final.pdf>
- [9] A. SCHULZE. *Navigating gene expression using microarrays: a technology review*. Nature cell biology **3** (2001) E190.
- [10] ARABIDOPSIS.ORG. *Microarray data from AtGenExpress*. <https://www.arabidopsis.org/portals/expression/microarray/ATGenExpress.jsp>.
- [11] J. KILIAN *et al.* *The AtGenExpress global stress expression data set: Protocols, evaluation and model data analysis of UV-B light, drought and cold stress responses*. Plant Journal **50** (2007) 347.
- [12] A. BRÜCKNER *et al.* *Yeast two-hybrid, a powerful tool for systems biology*. International Journal of Molecular Sciences **10** (2009) 2763.
- [13] M. E. CUSICK *et al.* *NIH Public Access*. Nature Methods **6** (2009) 39.

- [14] G. SALES *et al.* *graphite: GRAPH Interaction from pathway Topological Environment* (2015). R package version 1.16.0.
- [15] E. SEGAL *et al.* *Discovering molecular pathways from protein interaction and gene expression data.* Bioinformatics **19** (2003).
- [16] J. PANDEY *et al.* *Functional coherence in domain interaction networks.* Bioinformatics **24** (2008) 28.
- [17] P. RESNIK. *Using Information Content to Evaluate Semantic Similarity in a Taxonomy.* roceedings of the 14th international joint conference on Artificial intelligence - Volume 1 - IJCAI'95 **1** (1995) 6.
URL <http://arxiv.org/abs/cmp-lg/9511007>
- [18] C. PESQUITA *et al.* *Semantic similarity in biomedical ontologies.* PLoS Computational Biology **5** (2009).
- [19] A. BERENSTEIN. *Análisis de redes complejas en sistemas biomoleculares* (2014).
- [20] ASHBURNER. *Gene ontology: tool for the unification of biology.* Nat Genet **25** (2000).
- [21] G. GAN *et al.* *Data Clustering: Theory, Algorithms, and Applications*, volume 20 (2007).
- [22] M. HALKIDI *et al.* *On clustering validation techniques.* Journal of Intelligent Information Systems **17** (2001) 107.
- [23] E. DOMANY. *Superparamagnetic clustering of data—the definitive solution of an ill-posed problem.* Physica A: Statistical Mechanics and its Applications **263** (1999) 158.
URL <http://www.sciencedirect.com/science/article/pii/S0378437198004944>
- [24] S. CHEN *et al.* *On the similarity metric and the distance metric.* Theoretical Computer Science **410** (2009) 2365.
URL <http://dx.doi.org/10.1016/j.tcs.2009.02.023>
- [25] L. W. KHENG. *Image Registration* (2010).
- [26] P. D'HAESELEER. *How does gene expression clustering work?* Nat Biotech **24** (2005).
- [27] C. HENNIG. *How to find an appropriate clustering for mixed-type variables with application to socio-economic stratification.* Journal of the Royal Statistical Society. Series C: Applied Statistics **62** (2013) 309.
- [28] M. EISEN. *Cluster analysis and display of genome-wide expression patterns.* Proceedings of the National Academy of Sciences of the United States of America **95** (1998) 14863.
- [29] H. K. LEE *et al.* *Coexpression Analysis of Human Genes Across Many Microarray Data Sets* (2004) 1085.
- [30] J. SEVILLA. *Correlation between gene expression and go semantic similarity.* In: IEEE/ACMTransactions on Computational Biology and Bioinformatics (2005).

- [31] P. LORD. *Investigating semantic similarity measures across the Gene Ontology: the relationship between sequence and annotation.* Bioinformatics (2003).
- [32] D. LIN. *An Information-Theoretic Definition of Similarity.* In: Proc. of the 15th International Conference on Machine Learning (1998) 296.
- [33] J. JIANG. *Semantic Similarity Based on Corpus Statistics and Lexical Taxonomy.* Proceedings of International Conference Research on Computational Linguistics (1997) 19.
- [34] J. HARTIGAN. *A K-Means Clustering Algorithm.* Journal of the Royal Statistical Society **28** (1979) 100.
- [35] J. KOGAN. *Introduction to Clustering Large and High-Dimensional Data* (2006).
- [36] H. S. PARK. *A simple and fast algorithm for K-medoids clustering.* Expert Systems with Applications **36** (2009) 3336.
- [37] L. IBRAHIM. *Using Modified Partitioning Around Medoids Clustering Technique in Mobile Network Planning* **9** (2012) 299.
- [38] J. STEPHEN. *Hierarchical clustering schemes.* Psychometrika (1967).
- [39] C. SHALIZI. *Distances between Clustering , Hierarchical Clustering.* Data Mining (2009) 36.
- [40] P. LANGFELDER *et al.* *Defining clusters from a hierarchical cluster tree: The Dynamic Tree Cut package for R.* Bioinformatics **24** (2008) 719.
- [41] P. LANGFELDER *et al.* *Dynamic Tree Cut : in-depth description , tests and applications* (2007) 1.
- [42] S. HORVATH. *The Generalized Topological Overlap Matrix For Detecting Modules in Gene Networks.* bioinformatics (2007).
- [43] A. BARABÁSI. *Network biology: understanding the cell's functional organization.* Nat. Rev.Genet. (2004) 101.
- [44] L. HARTWELL *et al.* *From molecular to modular cell biology.* Nature (1999) 47.
- [45] M. ROSVALL. *Maps of random walks on complex networks reveal community structure.* Proceedings of the National Academy of Sciences of the United States of America **105** (2008) 1118.
- [46] A. CLAUSET *et al.* *Finding community structure in very large networks.* Phys. Rev. E **70** (2004) 66111.
URL <http://prola.aps.org/abstract/PRE/v70/i6/e066111>
- [47] J. DUTKOWSKI *et al.* *A gene ontology inferred from molecular networks.* **31** (2013) 38.
URL <http://www.nature.com/nbt/journal/v31/n1/abs/nbt.2463.html>\n<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3603333/>

- [48] S. DATTA. *Methods for evaluating clustering algorithms for gene expression data using a reference set of functional classes.* BMC bioinformatics **7** (2006) 397.
URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1590054&tool=pmcentrez>
- [49] N. CRISTIANINI *et al.* *On kernel target alignment.* Studies in Fuzziness and Soft Computing **194** (2006) 205.
- [50] Y. CHEN *et al.* *Learning kernels from indefinite similarities.* Proceedings of the 26th Annual International Conference on Machine Learning - ICML 2009 (2009) 1.
- [51] P. ERDOS *et al.* *On random graphs.* Publicationes Mathematicae (1959) 290.
- [52] M. MOLLOY *et al.* *A critical-point form random graphs with a given degree sequence.* Random structures and Algorithms (1995) 161.
- [53] A. J. BERENSTEIN. *Técnicas de Mecánica estadística para la detección de correlaciones en perfiles de expresión génica.* Tesis de Licenciatura en Ciencias Físicas (2010).
- [54] S. HORVAT. *A general framework for weighted gene co-expression network analysis.* Stat. Appl. Genet. Mol. Biol. (2005).