



UNIVERSIDAD NACIONAL AUTÓNOMA DE MÉXICO
PROGRAMA DE MAESTRÍA Y DOCTORADO EN CIENCIAS MATEMÁTICAS Y
DE LA ESPECIALIZACIÓN EN ESTADÍSTICA APLICADA

Inferencia de redes de interacciones causales en sistemas no lineales

TESIS
QUE PARA OPTAR POR EL GRADO DE:
MAESTRO EN CIENCIAS

PRESENTA:
José María Ibarra Rodríguez

DIRECTOR
Dr. Marco Tulio Angulo Ballesteros
CONACyT - Instituto de Matemáticas, UNAM.

QUERÉTARO, QUERÉTARO (26/05/2020)



Universidad Nacional
Autónoma de México

Dirección General de Bibliotecas de la UNAM

Biblioteca Central



UNAM – Dirección General de Bibliotecas
Tesis Digitales
Restricciones de uso

DERECHOS RESERVADOS ©
PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL

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

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

JOSÉ MARÍA IBARRA RODRÍGUEZ

INFERENCIA DE REDES DE INTERACCIONES
CAUSALES EN SISTEMAS NO LINEALES

José María Ibarra Rodríguez: *Inferencia de redes de interacciones causales en sistemas no lineales*. Enero 2020.

The miracle of the appropriateness of the language of mathematics for the formulation of the laws of physics is a wonderful gift which we neither understand nor deserve. We should be grateful for it and hope that it will remain valid in future research and that it will extend, for better or for worse, to our pleasure, even though perhaps also to our bafflement, to wide branches of learning.

— **Eugene Wigner** [33]

AGRADECIMIENTOS

Gracias a la comunidad del Instituto de Matemáticas UNAM Campus Juriquilla. Gracias a cada uno de los maestros, en especial a mi tutora, la Dra. Adriana Hansberg. Gracias también a mis compañeros y amigos, por las continuas tardes de discutir y contrastar ideas. Gracias Zamantha, Ruth, César, Sebastián por los momentos de trabajo compartido.

Gracias Esteban, Carlos, Fernando, Abraham y a toda la familia C3. Este proyecto sólo fue posible gracias al apoyo económico, la confianza y la opción constante que hacen en creer en cada uno de sus integrantes.

Gracias a mi familia. Gracias a mis papás y a mis hermanos por animarme siempre.

Gracias a mis maestros. A Francisco Juan, a Joel, a Víctor y a Marco.

Gracias en especial a Marco. Gracias por la oportunidad de trabajar juntos. Gracias por tu tacto pedagógico, tu compromiso, por la confianza y la libertad creativa.

Gracias Ale. Gracias por ser mi esposa, compañera y amiga. Gracias por siempre creer en mí. Gracias por tu entrega constante, por tu alegría y entusiasmo. Gracias porque vamos juntos construyendo sueños.

ÍNDICE GENERAL

INTRODUCCIÓN	1
Organización de la tesis	2
1 PRELIMINARES	5
1.1 Sistemas dinámicos y redes complejas	6
1.2 Inferencia de redes	7
1.3 Objetivos de la tesis	10
2 INFERENCIA DE REDES CAUSALES	13
2.1 Caracterización de todas las redes admisibles	14
2.2 Algoritmo para la inferencia de redes admisibles	19
2.3 Indicadores de la inferencia	29
2.4 Limitaciones fundamentales	31
2.5 Parámetro de tolerancia	35
3 VALIDACIONES NUMÉRICAS	37
3.1 Caso de Estudio: Ecosistema Aleatorio	38
3.2 Sensibilidad del metodo a muestras con ruido	46
3.3 Caso de Estudio: Microbiota intestinal humano	51
4 DISCUSIÓN	57
4.1 Comentarios finales	57
BIBLIOGRAFÍA	63

ÍNDICE DE FIGURAS

Figura 2.1	Conteo de vectores de signos contenidos en $S(y, z)$ en función de las entradas no cero de $w(y, z)$ para un sistema con $N = 12$ agentes.	26
Figura 2.2	Conteo de vectores de signos contenidos en $S(y, z)$ en función de las entradas no cero de $w(y, z)$ para un sistema con $N = 12$ agentes comparado con el número máximo de patrones(verde oscuro) y el tamaño promedio de los patrones (verde claro) que se alcanzaron al analizar muestras provenientes de un sistema de 12 agentes cuyos detalles se exponen en el siguiente capítulo.	28
Figura 2.3	Ejemplo en dimensión $N = 2$ de que al considerar sólo muestras en equilibrio no es posible alcanzar unicidad 1. (a) Muestra de un sistema de dos dimensiones que consiste en los tres puntos en equilibrio no triviales. (b) Ortantes considerados para cada una de las especies. (c) Conteo de los ortantes compatibles para cada una de las especies.	33
Figura 2.4	Ejemplo en dimensión $n = 2$ de que al considerar sólo muestras en movimiento no es posible alcanzar unicidad 1. (a) Muestra de un sistema de dos dimensiones que consiste en los cuatro puntos en movimiento. (b) Ortantes compatibles para cada par de muestras y cada una de las especies. (c) Conteo de los ortantes compatibles para cada una de las especies. Se observa que las muestras en movimiento hacen que cada especie tenga tres ortantes compatibles.	35

Figura 2.5	<p>Ejemplo en dimensión $N = 2$ de que al considerar un conjunto de datos que contiene tanto muestras en equilibrio como muestras en movimiento. (a) Conjunto de datos de un sistema de dos dimensiones que consiste en los cuatro puntos en movimiento y los tres equilibrios no triviales. (b) Ortantes compatibles para cada par de muestras y cada una de las especies. (c) Conteo de los ortantes compatibles para cada una de las especies. Al observar la frecuencias se observa que si la muestra contiene datos en equilibrio y en movimiento, entonces si es posible alcanzar unicidad 1 para ambas especies. 36</p>
Figura 3.1	<p><i>Simulación ecosistema aleatorio con cinco especies.</i> Series de tiempo correspondientes a las abundancias de cada especie. El color de las líneas indica la especie. Las muestras en movimiento se tomaron en la primera parte del proceso (línea gruesa), y las <i>muestras en equilibrio</i> se tomaron cuando las series cambiaban minimamente (Puntos a la derecha). Se observa como las diferentes configuraciones de especies determinan los valores de equilibrio a los que puede llegar cada especie. Resaltadas en negro se observan las abundancias correspondientes a la configuración con las cinco especies presentes. 40</p>
Figura 3.2	<p><i>Indicadores del proceso de inferencia.</i> Gráfica que muestra el cambio en los indicadores al agregar al análisis las muestras en movimiento. Suponiendo como punto de partida que ya se tomaron en cuenta todas las muestras en equilibrio disponibles. 41</p>
Figura 3.3	<p><i>Comparación de la matriz interacciones y la matriz de interacciones inferidas</i> Comparación de la matriz de parámetros originales (Izquierda) contra la matriz de interacciones inferidas (Derecha). La matriz de interacciones inferidas codifica en el color de cada entrada la cantidad de interacciones admisibles para cada entrada a_{ij}. En las entradas unívocamente determinadas se muestra el correspondiente valor inferido. . . . 41</p>
Figura 3.4	<p><i>Medias de los indicadores para distintos valores de ϵ_T.</i> Valores promedio de los indicadores del proceso de inferencia para diferentes valores del parámetro de tolerancia ϵ_T. 43</p>

Figura 3.5	<p><i>Heatmap Error vs. Unicidad para distintos valores de ϵ_T.</i> Se muestra el resultado para distintos valores del parámetro de tolerancia ϵ_T. Como se corrobora en la Tabla 3.3, el único valor que obtuvo error cero es $\epsilon_T = 10^{-6}$. Es de destacar que, independientemente del valor ϵ_T, la gran mayoría de las inferencias que tienen unicidad no cero, presentan un error pequeño.</p>	44
Figura 3.6	<p><i>Indicadores del proceso de inferencia.</i> Se muestra el cambio en los indicadores al agregar muestras en movimiento al análisis. En cada caso se tiene como punto de partida el resultado de analizar las muestras en equilibrio disponibles.</p>	45
Figura 3.7	<p><i>Unicidad vs. Número de configuraciones.</i> Cada punto corresponde al resultado final de la inferencia para un conjunto de datos aleatorio. Como se observa, la composición del conjunto de datos, indicada en el eje x y en el color, influye fuertemente en la unicidad que se alcanza. Notemos que, cuando no hay muestras en movimiento, la unicidad queda acotada por el número de ceros en la matriz de interacciones, tal como se espera.</p>	46
Figura 3.8	<p><i>Unicidad en distintas configuraciones de muestras.</i> Se observa como el valor de unicidad alcanzado está relacionado con el número de configuraciones y el número de muestras en movimiento consideradas. Para pocas configuraciones, la unicidad parece estancarse lejos del 1 independientemente del número de muestras en movimiento que se consideren. En cambio, cuando se consideran la mayoría de las configuraciones, incluso con muy pocas muestras en movimiento se alcanzan unicidades altas.</p>	47
Figura 3.9	<p><i>Unicidad en distintas configuraciones de muestras.</i> Se observa como el valor de unicidad alcanzado está relacionado con el número de configuraciones y el número de muestras en movimiento consideradas. Para pocas configuraciones, la unicidad parece estancarse lejos del 1 independientemente del número de muestras en movimiento que se consideren. En cambio, cuando se consideran la mayoría de las configuraciones, incluso con muy pocas muestras en movimiento se alcanzan unicidades altas.</p>	48

Figura 3.10	<i>Error vs Unicidad para distintos parámetros de tolerancia para muestras con diferentes niveles de ruido.</i> Mapas de calor de los conteos de los resultados de la inferencia. Los renglones corresponden a los distintos parámetros de ruido σ_R , mientras que los renglones corresponden a los parámetros de tolerancia ϵ_T considerados. En cada gráfica, el eje x corresponde a la tasa de error y el eje y a la unicidad obtenidas. Se busca que la inferencia maximice la unicidad minimizando la tasa de error. Por lo tanto, lo ideal es obtener la mayor concentración de resultados lo más cercano a la esquina superior izquierda. Por otro lado, si un conjunto se vuelve incompatible, o no tiene suficiente información, obtendrá unicidades bajas o cero.	50
Figura 3.11	<i>Medias de los indicadores para distintos valores de ϵ_T.</i> Valores promedio de los indicadores del proceso de inferencia para diferentes valores del parámetro de tolerancia ϵ_T . Observemos que $\epsilon_T = 10^{-4}$ es el único valor que no presenta ningún error, asimismo, presenta una informatividad más baja que los ϵ_T menores. Esto sugiere que el alto valor de unicidad que alcanzan con los valores de ϵ_T menores a él se debe a la presencia de errores.	53
Figura 3.12	<i>Heatmap Error vs. Unicidad para distintos valores de ϵ_T.</i>	54
Figura 3.13	<i>Matriz de interacciones inferidas.</i> Resultado de la inferencia de un conjunto de datos con 79 muestras en movimiento y $\epsilon_T = 10^{-4}$. Este conjunto de datos muestra un comportamiento típico, alcanzó una unicidad de 0.1944 y cero errores.	54
Figura 3.14	<i>Indicadores del proceso de inferencia.</i> Se muestra el cambio en los indicadores al agregar muestras en movimiento al análisis. En cada caso se tiene como punto de partida el resultado de analizar las muestras en equilibrio disponibles.	55
Figura 3.15	<i>Matriz de interacciones inferidas.</i> Resultado de la inferencia de un conjunto de datos con 158 muestras, 79 en equilibrio y 79 en movimiento. Este conjunto de datos muestra un comportamiento típico, alcanzó una unicidad de 0.94444 y cero errores.	55

ÍNDICE DE CUADROS

Cuadro 1.1	Ejemplos de sistemas dinámicos. En todos ellos, la red de interacciones subyacente queda codificada en la matriz de interacciones $A = (a_{ij})$. El modelo Lotka-Volterra generalizado [36], con respuesta funcional Holling II, es ampliamente utilizado para modelar ecosistemas. En el contexto de regulación genética, es usual considerar el modelo Michaelis-Menten [26]. Las poblaciones de neuronas suelen modelarse utilizando el modelo Wilson-Cowan [35].	12
Cuadro 3.1	Parámetros aleatorios para las interacciones del ecosistema simulado.	39
Cuadro 3.2	<i>Tabla con los valores utilizados para los parámetros de muestreo en la exploración inicial.</i> Se realizaron 20 repeticiones para cada combinación de los parámetros, dando un total de 2400 realizaciones.	42
Cuadro 3.3	<i>Análisis Exploratorio para seleccionar un parámetro de tolerancia</i> Valores promedio de los indicadores del proceso de inferencia para diferentes valores del parámetro de tolerancia ϵ_T	43
Cuadro 3.4	<i>Tabla con los valores de los parámetros de muestreo utilizados para ϵ_T fijo.</i> Se realizaron 30 repeticiones para cada combinación de los parámetros, dando un total de 900 realizaciones.	44
Cuadro 3.5	<i>Tabla con los valores de los parámetros de muestreo utilizados para ϵ_T fijo.</i> Se realizaron 10 repeticiones para cada combinación de los parámetros, dando un total de 7200 realizaciones.	49
Cuadro 3.6	<i>Parámetros para el modelo Lotka-Volterra generalizado para doce especies del microbiota intestinal humano.</i> Las doce especies consideradas son: <i>Prevotella copri</i> (PC), <i>Bacteroides vulgatus</i> (BV), <i>Bacteroides uniformis</i> (BU), <i>Bacteroides ovatus</i> (BO), <i>Bacteroides thetaiotaomicron</i> (BT), <i>Faecalibacterium prausnitzii</i> (FP), <i>Blautia hydrogenotrophica</i> (BH), <i>Eubacterium rectale</i> (ER), <i>Collinsella aerofaciens</i> (CA), <i>Eggerthella lenta</i> (EL), <i>Desulfovibrio piger</i> (DP) y <i>Clostridium hiranonis</i> (CH).	51

Cuadro 3.7	<i>Tabla con los valores utilizados para los parámetros de muestreo. Se realizaron 20 repeticiones para cada combinación de los parámetros, dando un total de 600 realizaciones.</i>	52
Cuadro 3.8	<i>Análisis Exploratorio para seleccionar un parámetro de tolerancia. Valores promedio de los indicadores del proceso de inferencia para diferentes valores del parámetro de tolerancia ϵ_T.</i>	53
Cuadro 3.9	<i>Tabla con los valores de los parámetros de muestreo utilizados para ϵ_T fijo. Se realizaron 20 repeticiones para cada combinación de los parámetros, dando un total de 100 realizaciones.</i>	53

INTRODUCCIÓN

Esta tesis presenta los resultados teóricos que sustentan un método para inferir de manera completa y libre de modelos las interacciones causales subyacentes a un sistema dinámico.

El concepto abstracto de red se ha vuelto una idea clave para entender sistemas complejos. Las redes permiten integrar el conocimiento de los distintos agentes de un sistema al considerar como se relacionan entre ellos.

Conocer la red de interacciones causales de un sistema dinámico nos ayuda a entender, predecir y controlar su comportamiento. Sin embargo, en una gran variedad de disciplinas científicas desde ecología hasta genética y neurociencias no es posible observar estas interacciones directamente. Para estos casos, es necesario desarrollar técnicas que permitan realizar la inferencia de las interacciones causales a partir de observaciones de la actividad de los agentes del sistema. Sin embargo, sigue siendo un problema abierto el realizar esta inferencia usando una metodología que tenga el soporte teórico necesario para garantizar la confiabilidad de la misma.

Uno de los aspectos centrales que dificultan la inferencia es el hecho de que un mismo conjunto de datos observado puede ser explicado igual de bien por distintas redes. A pesar de ello, todos los métodos de inferencia *libres de modelo* hasta el momento obtienen como resultado de la inferencia una única red de interacciones. Esta situación puede originar errores fundamentales de modelación, ya que al elegir una única matriz *arbitraria* se ignoran por completo las demás redes admisibles, y por lo tanto, se descartan posibles modelos alternativos.

Más aún, la gran mayoría de los métodos de inferencia de redes existentes requieren conocer las funciones específicas que describen la dinámica de las interacciones de un sistema. Incluso, una vez que se conoce la dinámica del modelo, existen métodos de inferencia bayesiana de parámetros que permiten obtener un conjunto (distribución) de redes posibles. Sin embargo, el gran limitante sigue siendo conocer la dinámica del modelo. Este requerimiento implica que estos métodos sólo pueden aplicarse en contextos muy específicos y haciendo fuertes hipótesis sobre los datos.

En esta tesis, desarrollamos los resultados teóricos que sustentan un método de inferencia libre de modelos y completo. Decimos que es un método libre de modelos ya que puede ser utilizado en una gran variedad de sistemas sin necesidad de conocer su dinámica, siempre que pueda ser descrita por un sistema de ecuaciones diferenciales or-

dinarias. Asimismo, decimos que el método es completo en el sentido de que infiere EXACTAMENTE todas las redes que son admisibles para el conjunto de datos dado, dentro de una amplia familia de modelos (que corresponde a aquellos cuya dinámica es continua con respecto a la actividad de los agentes).

Al ser libre de modelos, el método puede ser utilizado en una amplia variedad de sistemas naturales y tecnológicos. La principal novedad de nuestro método es que al ser completo, permite identificar las interacciones comunes a todas las redes admisibles para un conjunto de datos. De esta manera, al distinguir las interacciones que son unívocamente determinadas por el conjunto de datos de las que no lo son, se proporciona un nivel de confianza en la inferencia que no es alcanzable por ningún otro método existente.

Por otro lado, nuestros resultados permiten establecer condiciones necesarias que un conjunto de datos debe satisfacer para que, como resultado de la inferencia, se les pueda asociar una única red. Esto es un primer paso crucial para el diseño de experimentos que generen datos que determinen de manera única todas las entradas de la matriz de interacciones causales.

ORGANIZACIÓN DE LA TESIS

En el primer capítulo presentamos el problema de reconstrucción de redes causales y el objetivo del nuevo enfoque aquí presentado. Comenzamos planteando el problema y mencionando distintas metodologías existentes para afrontar el mismo. Concluimos esta parte definiendo de manera formal las hipótesis necesarias para nuestro método, así como los objetivos de esta tesis.

En el segundo capítulo presentamos los resultados matemáticos centrales que permiten definir las limitaciones y el alcance del método. Presentamos primero las definiciones y los resultados preliminares necesarios. Sorpresivamente nuestros resultados no requieren un andamiaje matemático sofisticado, sino que están basados en el Teorema del valor medio para espacios vectoriales. La aportación principal es un teorema que caracteriza todas las redes admisibles para un conjunto de datos. Se incluyen también resultados que nos permiten construir un algoritmo eficiente para obtener estas redes. A continuación, se definen nuevos indicadores para medir el desempeño de la inferencia y procedemos a realizar observaciones sobre las limitaciones fundamentales de la informatividad de un conjunto de datos. Concluimos con un resultado que prueba que, para poder obtener una inferencia completa, es necesario tener disponibles tanto muestras en movimiento como muestras en estado estable.

En el tercer capítulo realizamos validaciones numéricas del método de inferencia propuesto usando modelos de ecosistemas microbianos. Mostramos primero el desempeño del método con un ecosistema

aleatorio, lo que permite describir su comportamiento al ajustar el parámetro de tolerancia ante distintos niveles de ruido. Se realizó también una prueba del método con un modelo de microbiota intestinal humano que ilustra su utilidad de manera práctica.

Por último, el cuarto capítulo contiene las conclusiones del trabajo. Asimismo, se presenta una lista de posibles líneas de investigación para extender y profundizar en los resultados presentados en esta tesis.

PRELIMINARES

En la última década, la noción matemática de red se ha consolidado como una herramienta fundamental para modelar una gran cantidad de sistemas naturales y tecnológicos en diversas disciplinas [28]. Un modelo de red describe no sólo los *agentes* o componentes de un sistema (*e.g.*, genes), sino también las *relaciones* que existen entre ellas (*e.g.*, activación o inhibición). Este enfoque le da una flexibilidad única a los modelos basados en redes ya que la interpretación y alcance de los mismos depende de las relaciones específicas codificadas por las aristas de la red. Es decir, al considerar diferentes tipos de relaciones entre los componentes de un sistema, se pueden obtener redes distintas que codifican información diferente del mismo.

Una de las preguntas elementales que se hace al estudiar un sistema y sus componentes, y que resulta fundamental para comprender muchos fenómenos, es tratar de entender cuales de las unidades del sistema actúan *causal* y *directamente* sobre otras [1]. En este contexto, se puede asociar una red dirigida que modele las *interacciones causales* entre los agentes de un sistema. Es decir, se puede asociar una red al sistema donde los nodos correspondan a los *agentes* y las aristas dirigidas indiquen interacciones causales directas entre ellos.

Conocer las interacciones causales de un sistema proporciona información muy útil para entender, predecir y controlar su dinámica [32]. Estudios recientes muestran un avance significativo para entender las implicaciones que la topología de la red de interacciones tiene sobre la dinámica y la función de un sistema [17].

Así, resulta de especial importancia poder obtener la red de interacciones causales de un sistema de interés. Sin embargo, no siempre es posible observar directamente estas interacciones. Esto sucede, por ejemplo, en sistemas de regulación genética o ecosistemas. En estos casos sólo es posible medir la *actividad* de cada uno de los agentes del sistema (*e.g.*, nivel de expresión de genes o abundancia de especies ecológicas), y con esa información se intenta *inferir* la red subyacente de interacciones. Esta inferencia resulta ser un proceso difícil, y se han desarrollado diversas técnicas de *reconstrucción de redes* para resolverlo. Las aplicaciones de esta inferencia van desde las ciencias naturales y sociales hasta la ingeniería [30].

En general, cuando se tienen interacciones causales en un sistema, la actividad de cada agente del mismo se modela mediante un sistema dinámico. Para este sistema dinámico se tiene una red subyacente que codifica la presencia (o ausencia) de cada interacción y la "naturaleza",

positiva o negativa, de cada una de ellas. Con lo cual, la presencia de una arista dirigida en la red corresponde a inferir que un agente dado influye, promoviendo o inhibiendo, la actividad de otro.

En este capítulo se presenta el problema de *inferencia de interacciones causales*, a partir de observaciones de la actividad de los agentes. Como conclusión al capítulo, se discuten las limitaciones fundamentales que motivan los resultados de esta tesis.

1.1 SISTEMAS DINÁMICOS Y REDES COMPLEJAS

En general, los modelos que se utilizan para describir la *actividad* de los agentes de un sistema asumen que es posible describir su dinámica mediante un sistema de ecuaciones diferenciales. Así, si el sistema tiene n agentes, se considera que la actividad $x_i(t) \in \mathbb{R}$ asociada al i -ésimo agente al tiempo t es influida por los otros agentes a través de la ecuación diferencial

$$\frac{dx_i}{dt} = F_i(x_1(t), \dots, x_n(t)), \quad i = 1, \dots, n. \quad (1.1)$$

Aquí, $F_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, \dots, n$, es una función que representa la dinámica del i -ésimo agente.

En este caso, inferir la red de interacciones significa determinar cuales de las variables x_j aparece explícitamente en F_i . Este tipo de modelos se utilizan en una gran variedad de disciplinas, donde la actividad de los agentes puede ser, por ejemplo, la abundancia de especies de un ecosistema, el nivel de expresión de genes de una red de regulación genética, ó el nivel de concentración de reactivos en un sistema bioquímico. Este concepto de causalidad se presenta a detalle en [29].

Más aún, estamos interesados en conocer la *naturaleza* de las interacciones. Esto es, queremos conocer si el j -ésimo agente promueve o inhibe al i -ésimo agente. Esto se puede definir de manera natural, ver [1], como sigue.

DEFINICIÓN 1.1 (INTERACCIONES CAUSALES) Dado un sistema con n agentes como en la Ecuación 1.1 se tiene que el agente j tiene una interacción causal positiva directa, o que *promueve*, al agente i si

$$\frac{\partial F_i(x)}{\partial x_j} > 0, \quad \forall x \in \mathbb{R}^n.$$

De la misma manera se dice que el agente j tiene una interacción causal negativa directa, o que *inhibe*, al agente i si

$$\frac{\partial F_i(x)}{\partial x_j} < 0, \quad \forall x \in \mathbb{R}^n.$$

Dicho esto, dado un sistema, puede inferirse la red subyacente del sistema en dos niveles. Un primer nivel corresponde a determinar cuales interacciones causales existen, es decir determinar si $\frac{\partial F_i}{\partial x_j} \neq 0$, $i, j = 1, \dots, n$. Un segundo nivel de la inferencia corresponde no sólo a identificar las interacciones causales sino también a obtener el signo de las interacciones no cero.

Notemos que el resultado de inferir las interacciones es una matriz de adyacencia de la red no dirigida subyacente al sistema. Análogamente, inferir las interacciones con signo corresponde a inferir la matriz de adyacencia ponderada de la red dirigida subyacente al sistema, donde los pesos de las aristas indican la naturaleza de la interacción. Esta matriz de adyacencia es una matriz con elementos en $\{-1, 0, 1\}$. El método que proponemos para obtener estas interacciones se presenta con más detalle en el CAPÍTULO 2.

1.2 INFERENCIA DE REDES

Existen distintos enfoques para hacer inferencia sobre la red de interacciones subyacente de un sistema. Todos ellos asumen diferentes niveles de conocimiento del sistema y pueden requerir, o no, intervenir en el sistema usando acciones de control [30]. Usualmente los datos que se pueden observar corresponden a la actividad de cada uno de los agentes del sistema a lo largo del tiempo. En general el problema de inferencia de redes de interacción consiste en utilizar las observaciones que se puedan obtener sobre la actividad de los agentes de un sistema para inferir las interacciones causales (con o sin signo) subyacentes.

Los datos experimentales que pueden obtenerse del sistema casi siempre son mediciones de la actividad de los agentes del mismo. En otras palabras, lo que se obtiene son observaciones en distintos tiempos de una o varias trayectorias dadas del sistema dinámico. Las distintas estrategias parten de considerar estos datos como *series de tiempo*, algunas veces incluso en equilibrio. Los métodos que se han desarrollado para afrontar la inferencia pueden clasificarse según la maquinaria estadística en las que están basados [32]. Estos pueden ser métodos basados en *correlación*, ó métodos que buscan *causalidad estadística*, por ejemplo mediante información mutua o redes Bayesianas.

Los métodos de correlación son populares puesto que suelen ser "*simples, eficientes y generalmente realizables*"[18]. Sin embargo, dado que la correlación entre dos series de tiempo es una relación simétrica, no es posible asignar una dirección causal en las interacciones detectadas [29]. Aún más importante es el hecho de que los métodos basados en correlación se utilizan a pesar de que la correlación no es suficiente, ni necesaria, para establecer relaciones causales [15, 20, 21].

Utilizar la correlación para inferir causalidad es un error fundamental que conlleva a problemáticas y defectos en la inferencia. Estas

problemáticas ya se han detectado en distintos campos. Por ejemplo, se sabe que los métodos basados en correlación son propensos a generar redes demasiado densas o con presencia de correlaciones espurias [3]. Esto se debe en parte a que las redes de correlación se han utilizado incluso cuando la naturaleza de los datos es incompatible con el uso de la misma, y por lo tanto no debería utilizarse para detectar relaciones causales [10].

Estas limitantes se ven reflejadas en el desempeño de los algoritmos de inferencia existentes. Por ejemplo, incluso cuando se complementa con otro tipo de datos, los algoritmos de inferencia de redes existentes casi nunca producen información confiable sobre interacciones biológicas [5]. Incluso, se han reportado casos donde distintos métodos de inferencia de redes basados en correlación producen resultados diferentes de manera consistente [34].

Existen métodos específicos para buscar *causalidad estadística* que no están basados en correlación. Los principales métodos están respaldados con resultados teóricos de áreas como información mutua, redes Bayesianas, o causalidad en el sentido de Granger. Sin embargo, incluso estos métodos diseñados para datos y modelos específicos también presentan limitaciones importantes.

Por ejemplo, suele suceder que presentan dificultad para detectar interacciones indirectas entre agentes interpretandolas como interacciones directas. Es decir, si en un sistema las interacciones $i \rightarrow j$ and $j \rightarrow k$ existen, estos métodos suelen inferir que el sistema tiene la interacción $i \rightarrow k$. En particular las condiciones teóricas necesarias para utilizarlos son altamente específicas, y por lo tanto difícilmente las satisfacen los datos disponibles [18]. Los métodos bayesianos además son computacionalmente muy costosos, y sin garantía de converger a la red de interacciones auténtica [32].

Los métodos basados en causalidad de Granger suponen que en las interacciones causales a inferir la dinámica es lineal, requiere datos con una estructura temporal compartida y conocida, así como que el sistema tenga la propiedad de *separabilidad*, un requisito que se ha visto rara vez cumplen los sistemas dinámicos debido a no linealidades [29]. Además, los resultados son difíciles de interpretar y presentan problemáticas para inferir interacciones indirectas [16].

Las dificultades de adaptar estrategias teóricas a los datos experimentales hacen que la inferencia de las interacciones causales sea aún un problema abierto. Los métodos ya existentes llegan a presentar intersecciones tan pequeñas en la inferencia que se proponen estrategias que integren los distintos resultados [9]. Algunos de estos métodos presentan problemas de replicabilidad de resultados [19]. Aunado a lo anterior, la principal limitante de los métodos existentes es, en general, la poca eficiencia estadística a grado tal que incluso resultan ser estadísticamente equivalentes a hacer la inferencia al azar [22].

En síntesis, el problema de inferir interacciones causales es un problema de gran relevancia en distintas disciplinas. Con lo cual es necesario desarrollar herramientas específicas para inferir las interacciones que influyen la dinámica de un sistema. Aunque existen avances en esa dirección, la reconstrucción de redes de interacciones causales sigue siendo un problema abierto [25]. Así, inferir la dinámica de un sistema sigue siendo un desafío que requiere de nuevas bases teóricas y métodos novedosos.

Métodos de inferencia de redes causales

En su mayoría, los enfoques actuales para inferir las interacciones causales de un sistema requieren conocer *a priori* un modelo de su dinámica. Es decir, es necesario conocer las funciones matemáticas que modelan las interacciones, lo cual reduce el problema a una inferencia de parámetros multidimensional.

Los métodos que presuponen un modelo matemático para el cual sólo falta determinar los parámetros desconocidos no suelen adaptarse fácilmente a los datos que se obtienen en condiciones experimentales. Esto limita su aplicabilidad a sistemas cuya dinámica es bien conocida, dejando al lado una gran cantidad de sistemas biológicos, ecológicos y sociales complejos [7]. Es por ello que para realizar la inferencia en distintas áreas requiere adaptar, o crear, métodos específicos para cada situación, ver por ejemplo [11, 12, 17, 18, 24, 27, 37].

Suponer *a priori* un modelo dinámico puede generar errores fundamentales en la inferencia. Por un lado, presuponer una dinámica específica para las interacciones sobre un conjunto de datos puede producir errores sistemáticos en la inferencia si la dinámica real difiere de la supuesta [4]. Por otro lado, estos métodos requieren estructuras temporales en el muestreo que difícilmente se satisfacen en muchas condiciones, ó requieren tener la capacidad de intervenir en el sistema, lo cual no siempre es posible [36].

Estas limitaciones hacen evidente que es necesario un método que permita realizar la inferencia sin necesidad de conocer, o presuponer, el modelo de la dinámica del sistema. Entre los trabajos pioneros en el enfoque *libre de modelo* encontramos el trabajo de Nitzan, *et al.* [7]. Este enfoque no es completamente libre de modelo ya que busca estimar las funciones de interacción mediante aproximaciones sucesivas. Para ello, se requiere conocer la clase de funciones con la que se pueden aproximar la dinámica del sistema. Se desconocen las condiciones suficientes y necesarias que requiere tener un conjunto de datos para que el método efectivamente converja a la solución.

Una de las primeras estrategias de inferencia libre de modelo es el trabajo de Xiao *et al.* [36] y aunque también requiere aún de resultados teóricos que establezcan condiciones suficientes y necesarias que los datos deban satisfacer para poder realizar la inferencia, éste el primer

método que intenta manejar una de las principales carencias de todos los métodos anteriores ya que puede obtener más de una red de interacciones para los datos. Estos esfuerzos nos acercan un paso más al objetivo de obtener un método de inferencia de interacciones causales, pero presentan aún limitantes importantes.

Limitaciones fundamentales en los algoritmos de inferencia existentes

El proceso de inferencia está sujeto a limitaciones fundamentales inherentes a las características de los datos disponibles y al conocimiento de la dinámica del sistema, y que por tanto ningún método puede solventar. Independientemente de las limitaciones propias de cada método, la mayoría de los métodos proponen una *única* red de interacciones candidata a explicar las interacciones causales. Esto es un problema ya que se sabe que existen familias de sistemas dinámicos que generan exactamente las mismas trayectorias y, por lo tanto, cualquier conjunto de muestras que se obtenga de cualquiera de ellos no permitirá decidir sobre cuál de todas ellas provienen [1]. Es decir, si un conjunto de datos puede ser explicado por más de una red de interacciones, en general no es posible tener certeza de que las interacciones causales inferidas [6, 8, 13] por cualquiera de los métodos mencionados, puesto que no existe manera de decidir si la red inferida es la única red en la familia de sistemas dinámicos que pueden explicar los datos o si existe alguna otra red diferente que pudiera haber dado origen a los datos con propiedades distintas a la de la red obtenida.

Debido a que la utilidad de la inferencia radica fuertemente en la *confianza* que se puede tener en las interacciones inferidas, no es suficiente conocer *una* red de interacciones obtenida por un método. Es indispensable conocer *todas* las redes de interacciones que podrían haber dado origen a los datos con los que se cuenta. Esta es la única manera de garantizar que se tiene una interacción causal específica en el sistema.

A pesar de que Xiao *et al.* [36] son los primeros que tratan de afrontar esta limitación al proporcionar más de una red, aún no garantizan matemáticamente que obtienen todas las redes *admisibles* para el conjunto de datos dado. Debido a esto, la implementación de su método utiliza una técnica heurística. Otra limitante de la estrategia es que es un método que sólo admite datos en estado estable.

1.3 OBJETIVOS DE LA TESIS

Ante este panorama, en esta tesis se presentan los resultados teóricos que sustentan el enfoque de inferencia de interacciones causales libre de modelos presentado por Xiao *et al.* [36]. Adecuamos este método para obtener una inferencia *completa* de la familia de todas las

redes de interacciones que explican los datos, expandiendo el enfoque para utilizar tanto datos estables como datos dinámicos. Se describen también condiciones necesarias sobre el conjunto de datos para poder asociarle a los mismos una red de interacciones única. Y, aún en caso de que las limitantes en datos no permitan inferir la matriz completa, se muestra como podría ser posible recuperar parcialmente la red de interacciones.

El método de inferencia presentado es *completo* ya que infiere *exactamente todas* las redes de interacciones compatibles con los datos. Es decir, se infieren todas las redes de interacciones en la familia de dinámicas que podrían haber dado origen a los datos. Hasta donde tenemos conocimiento, es el primer método que logra hacerlo.

Notemos que, en el caso de que todas las redes inferidas compartan una misma arista, se tiene, entera confianza de que ésta interacción causal es parte de la dinámica del fenómeno estudiado. Por lo tanto, aún cuando los datos sean poco informativos, es posible que el método permita inferir algunas de las interacciones inequívocamente. De esta manera, es posible obtener *redes parciales* para las cuales se tiene confianza plena, aunque no se conozca ni la red completa ni los detalles de la dinámica de las interacciones ni los parámetros involucrados en las mismas.

El método es, también, *libre de modelos* ya que no es necesario suponer ningún modelo dinámico del sistema para llevar a cabo la inferencia. El método infiere todas las interacciones de las posibles dinámicas que pueden explicar el comportamiento de los datos. Las condiciones que se les pide a los sistemas dinámicos contemplados en este enfoque son mínimas. Éstas hipótesis son fácilmente satisfechas por las funciones más utilizadas para modelar una gran variedad de sistemas en distintas disciplinas.

Para describir de manera concreta las dinámicas a las que nos referimos consideraremos sistemas dinámicos que pueden ser descritos por

$$\frac{dx_i}{dt} = q_i(x_i) f_i(x_1, \dots, x_n),$$

para $i = 1, \dots, n$, donde $q_i : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$. Para esta clase de sistemas, la *red de interacciones* $S_f(t)$ corresponde a

$$S_f(t) = \text{sgn} \left(\frac{\partial f(x(t))}{\partial x} \right).$$

Esta clase de modelos de sistemas dinámicos es suficientemente amplia para modelar desde ecosistemas hasta epidemias [2]. Para ejemplos de dinámicas que pueden expresarse de esta manera ver TABLA 1.1. Esta familia de modelos son bastante útiles porque permiten separar la parte conocida de la dinámica intra-agentes $q_i(x_i)$, que suele ser más accesible de obtener, de la dinámica de interacción $f_i(x)$ que, usualmente, no es fácil de conocer. Notemos que la naturaleza de las

interacciones queda codificada explícitamente en el signo del Jacobiano de $f_i(x)$. En el CAPÍTULO 2 se explican a detalle los resultados teóricos que permiten realizar la inferencia partiendo de resultados teóricos concernientes a dicho Jacobiano.

SISTEMA	AGENTE	ACTIVIDAD DEL AGENTE	EJEMPLO DINÁMICA	$q_i(x_i)$
Ecosistema	especie	abundancia	$F_i(x) = r_i x_i + \sum_j a_{ij} \frac{x_i x_j}{1 + \theta_{ij} x_j}$	x_i
Regulación genética	genes	nivel de expresión	$F_i(x) = -b_i x_i + \sum_{j \neq i} a_{ij} \frac{x_j}{1 + \theta_{ij} x_j^h}$	$\mathbf{1}$
Neurociencia	neuronas	tasa de disparo	$F_i(x) = -x_i + \sum_j a_{ij} \frac{1}{1 + e^{-r(x_j - \mu)}}$	$\mathbf{1}$

Tabla 1.1: Ejemplos de sistemas dinámicos. En todos ellos, la red de interacciones subyacente queda codificada en la matriz de interacciones $A = (a_{ij})$. El modelo Lotka-Volterra generalizado [36], con respuesta funcional Holling II, es ampliamente utilizado para modelar ecosistemas. En el contexto de regulación genética, es usual considerar el modelo Michaelis-Menten [26]. Las poblaciones de neuronas suelen modelarse utilizando el modelo Wilson-Cowan [35].

INFERENCIA DE REDES CAUSALES

Conocer la matriz de interacciones de un sistema dinámico ayuda a entenderlo, predecirlo y controlarlo. Dado que no es posible inferir la dinámica específica a partir de observaciones puntuales de un sistema, resulta de interés poder conocer, por lo menos, la matriz de interacciones del sistema. Por otra parte, recordemos que es posible que un conjunto de datos puedan ser explicados por toda una familia de dinámicas entre las cuales no es posible decidir cual corresponde al fenómeno de interés. Luego, es necesario considerar las diferentes matrices de interacciones, ya que todas ellas son candidatas a explicar las relaciones causales entre los agentes.

Los resultados principales de esta tesis, que se presentan en este capítulo, constituyen un primer esfuerzo en este sentido. Esto es, proponemos un método para encontrar *todas* las matrices de interacciones compatibles con una muestra de datos provenientes de un sistema dinámico.

Para ello, la SECCIÓN 2.1 comienza dando una caracterización matemática de todas las redes admisibles para un conjunto de datos dado. Luego, utilizando esta caracterización, en la SECCIÓN 2.2 presentamos los resultados que permiten construir un algoritmo eficiente para inferir todas estas redes admisibles. En la SECCIÓN 2.3 se definen los indicadores que permiten medir que tan cercano está un conjunto de datos de permitir la inferencia exitosa de una *única* red de interés. Estos indicadores, introducidos por primera vez en este trabajo, permiten caracterizar en la SECCIÓN 2.4 las condiciones necesarias de un conjunto de datos que permitan inferir una *única* matriz de interacciones. Se ilustra una aplicación del algoritmo, y las limitaciones fundamentales, con un caso muy simple de inferencia de una red ecológica para un ecosistema de dos especies.

Es importante notar que las limitaciones en la inferencia dependen de las características de los datos y no del algoritmo de inferencia usado. Es por ello que el método que se presenta es el primero en inferir toda la información de la red causal que puede ser obtenida de la muestra de datos. Así, se cuenta por primera vez con un algoritmo de inferencia universal y completo.

Así, uno de los resultados principales es probar matemáticamente que para realizar una inferencia exitosa, en general, se requiere tanto de observaciones de puntos de equilibrio del sistema, muestras en

estado estable, como observaciones fuera de los mismos, *muestras en movimiento*.

El CAPÍTULO 3 consiste en validaciones del método mediante simulaciones de datos sintéticos de ecosistemas.

2.1 CARACTERIZACIÓN DE TODAS LAS REDES ADMISIBLES

El objetivo de la inferencia es conocer la matriz de interacciones del sistema del que proviene una muestra de datos. Es posible que a un conjunto de datos dados se le pueda asociar dos o más matrices de interacciones. En este caso, no hay elementos para decidir cual de estas matrices de interacciones es la asociada al proceso dinámico del que provienen los datos. Más aún, si se escoge una única matriz de interacciones puede suceder que el sistema dinámico del que provienen los datos tenga una matriz de interacciones diferente, y por lo tanto las conclusiones a las que se lleguen no sean de utilidad.

Es por ello que el proceso de inferencia de la matriz de interacciones requiere tener confianza en que la matriz asociada es efectivamente la *única* que es compatible con los datos. O bien, si esto no es posible, estudiar el sistema considerando todos los casos. Es decir, asociar *todas* las matrices de interacciones compatibles con una muestra de datos provenientes de un sistema dinámico y considerar este conjunto para estudiar y modelar el sistema de interés.

Dado que no se puede *observar* la totalidad del sistema, no siempre es posible determinar de manera única las características del mismo. Así, el objetivo del capítulo es caracterizar *todas* las redes de interacciones admisibles para un conjunto de datos. En particular, resulta de interés inferir las características comunes de *todas* las *redes admisibles* para el conjunto de datos dado.

Para ello es necesario empezar por establecer de manera formal el concepto de un *conjunto de datos*. Dado un sistema dinámico de n agentes, con la actividad de cada agente descrita como $\dot{x}_i(t) = q(x_i(t))f_i(x(t))$, con $q_i(x_i) \geq 0$. Sea $R \subset \mathbb{R}^n$ una región convexa del dominio de f para la cual $\text{sign}\left(\frac{\partial f_i}{\partial x_j}\right)$ es constante para toda $t \geq 0$ y para todos los agentes $i, j = 1, \dots, n$. Definimos un *conjunto de datos* como sigue:

DEFINICIÓN 2.1 Un *conjunto de datos* \mathcal{D} es una conjunto de observaciones del sistema dentro de la región R , junto con sus derivadas. Es decir

$$\mathcal{D} = \{(y, \dot{y})_k \mid y \in R\},$$

donde, $\dot{y}_i(t) = \frac{dy_i(t)}{dt}$. A los pares $(y_i(t), \dot{y}_i(t))$ les llamamos *muestras*.

Además, si $\dot{y}_i(t) = 0$, diremos que es una *muestra en equilibrio* y en otro caso es una *muestra dinámica*. Por último, denotamos mediante \mathcal{D}_i al subconjunto de muestras $y \in \mathcal{D}$ tales que $y_i > 0$.

Supongamos que tenemos un conjunto de datos \mathcal{D} que proviene de un sistema con dinámica $\dot{x}_i = q_i(x_i)f_i^*(x)$ donde las funciones q_i son conocidas pero las funciones f_i^* de interacciones no lo son. Supongamos también que la naturaleza de las interacciones no cambia en el tiempo, es decir que el $\text{sgn}(J(f^*))$ es constante. Así, inferir la naturaleza de las interacciones equivale a inferir $\text{sgn}\left(\frac{\partial f_i^*}{\partial x_j}\right)$.

Los resultados matemáticos sobre los que se basa la estrategia de inferencia resultan ser sorprendentemente elementales. En particular, los resultados en esta sección dependen fuertemente en el Teorema del Valor Medio para funciones de varias variables. Retomamos aquí el enunciado del mismo tal como aparece en [14].

TEOREMA 2.1 TEOREMA DEL VALOR MEDIO Sea U un abierto en \mathbb{R}^n , consideremos dos puntos $x, y \in U$ y una función $f : U \rightarrow \mathbb{R}^m$ clase C^1 . Supongamos que el segmento $x + ty$ para $t \in [0, 1]$ se encuentra contenido en U . Entonces,

$$f(x + y) - f(x) = \int_0^1 f'(x + ty)y dt = \int_0^1 f'(x + ty)dt \cdot y.$$

Donde, " \cdot " representa el producto interno de vectores en \mathbb{R}^n y $f'(x) = \partial f(x)/\partial x$.

Notemos que al considerar dos muestras arbitrarias $y, z \in \mathcal{D}_i$ podemos definir $x_1 = z$ y $x_2 = y - z$ y, así, aplicar el **TEOREMA 2.1** en la dinámica desconocida f_i^* del agente i . De esta manera se obtiene

$$f_i^*(y) - f_i^*(z) = \left(\int_0^1 f_i^{*\prime}(z + s(y - z))ds \right) \cdot (y - z). \quad (2.1)$$

Como además sabemos que $\dot{y}_i = q_i(y_i)f_i^*(y)$ y $\dot{z}_i = q_i(z_i)f_i^*(z)$, podemos sustituir

$$f_i^*(y) = \frac{\dot{y}_i}{q_i(y_i)} \text{ y } f_i^*(z) = \frac{\dot{z}_i}{q_i(z_i)}$$

en la ecuación 2.1 para obtener que

$$\frac{\dot{y}_i}{q_i(y_i)} - \frac{\dot{z}_i}{q_i(z_i)} = \left(\int_0^1 \frac{\partial f_i^*}{\partial x}(z + s(y - z))ds \right) \cdot (y - z). \quad (2.2)$$

Definiendo a $v_i(y, z) := \int_0^1 \frac{\partial f_i^*}{\partial x}(z + s(y - z))ds$ y $\beta_i(y, z) = \frac{\dot{y}_i}{q_i(y_i)} - \frac{\dot{z}_i}{q_i(z_i)}$, la ecuación puede reescribirse como

$$\beta_i(y, z) = v_i(y, z) \cdot (y - z), \quad (2.3)$$

lo que permite hacer una primera observación importante.

OBSERVACIÓN 2.1 Como por hipótesis el patrón de signos del Jacobiano de f_i^* es constante, y la integral en la Ecuación 2.2 preserva signos, entonces el patrón de signos del vector $v_i(y, z)$ es el mismo para

cualesquiera $y, z \in \mathcal{D}_i$. Más aún, este patrón de signos coincide con el patrón de interacciones $\text{sgn}(J(f_i^*))$.

Ahora bien, es posible que las muestras en \mathcal{D}_i puedan ser explicadas por alguna otra dinámica con un patrón de interacciones diferente a f_i^* , como se ilustra en [1]. Para considerar todas las dinámicas que pueden *explicar* los datos es necesario definir de manera formal cuando un vector de signos es *admisibile* para una muestra de datos.

DEFINICIÓN 2.2 Dado un conjunto de datos \mathcal{D} , se dice que un patrón de signos $s \in \{1, 0, -1\}^n$ es *admisibile* para el agente i , si existe una función continua $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ tal que

- i) $\text{sgn}\left(\frac{\partial f_i}{\partial x}\right) = s$ para todo x en el dominio de interés. Esto es, el patrón de signos del Jacobiano de f_i coincide con s .
- ii) $\dot{y}_i = q(y_i)f_i(y)$ para toda muestra $y \in \mathcal{D}$. Es decir, la dinámica f explica los datos.

Denotamos al conjunto de signos admisibles para el agente i , como

$$\mathcal{A}_i(\mathcal{D}) := \{s \in \{1, 0, -1\}^n \mid s \text{ es admisible para el agente } i \text{ dado } \mathcal{D}\}.$$

Por construcción, notemos que el *patrón de signos verdadero* s_i^* , correspondiente a la dinámica f_i^* de la que provienen las muestras, siempre pertenece a $\mathcal{A}_i(\mathcal{D})$. Sin embargo, es posible que no sea el único patrón de signos admisibles y que existan otros patrones de signos admisibles para el agente i para el conjunto de datos dado.

Supongamos que existe $s \in \mathcal{A}_i(\mathcal{D})$ diferente a s_i^* . Esto es, existe una dinámica g_i admisible para \mathcal{D}_i diferente a f_i^* . Notemos que para cada par de muestras $y, z \in \mathcal{D}_i$ definiendo $v'_i(y, z) = \int_0^1 \frac{\partial g_i}{\partial x}(z + s(y - z)) ds$, se tiene que

$$\beta_i(y, z) = v'_i(y, z) \cdot (y - z). \quad (2.4)$$

Por lo tanto, el vector de interacciones $\text{sgn}(J(g_i))$ coincide con $\text{sgn}(v'_i(y, z))$ para cada par de muestras $y, z \in \mathcal{D}_i$. Así, resulta que los vectores de interacción admisibles no sólo corresponden a una dinámica que explica los datos, sino que también están relacionados con las soluciones de la Ec. (2.3). A los vectores de signos correspondientes a la soluciones de dicha ecuación los llamaremos *compatibles*.

DEFINICIÓN 2.3 Dado un conjunto de datos \mathcal{D} , se dice que un vector de signos $s \in \{1, 0, -1\}^n$ es *compatible* para el agente i , si para cada par de muestras $y, z \in \mathcal{D}_i$ se tiene que

- i) Existe $v_i(y, z) \in \mathbb{R}^n$ tal que $\text{sgn}(v_i(y, z)) = s$, y además,
- ii) $\beta_i(y, z) = v_i(y, z) \cdot (y - z)$, donde $\beta_i(y, z) = \frac{\dot{y}_i}{q_i(y_i)} - \frac{\dot{z}_i}{q_i(z_i)}$.

Al conjunto de vectores compatibles con la especie i para un conjunto \mathcal{D} dado lo denotaremos por $\mathcal{S}_i(\mathcal{D})$.

Es importante observar que es posible decidir si un vector de signos s es *compatible* con un conjunto \mathcal{D}_i con base únicamente en los datos. No sólo eso, como se explica en la siguiente sección, es posible calcular $\mathcal{S}_i(\mathcal{D})$ a partir de los datos sin necesidad de construir los vectores v de manera explícita. Esta observación resulta crucial puesto que en la búsqueda de los vectores admisibles para un conjunto de datos, un primer criterio que deben de satisfacer es ser compatibles con la muestra y éstos pueden obtenerse de los datos. En resumen, partiendo de la definición de admisible y utilizando la **OBSERVACIÓN 2.1** se obtiene el siguiente lema.

LEMA 2.1 El conjunto de patrones de signos admisibles para un conjunto de datos $\mathcal{A}_i(\mathcal{D})$ está contenido en el conjunto de patrones de signos compatibles $\mathcal{S}_i(\mathcal{D})$.

Por otro lado, es necesario también considerar la definición equivalente a los vectores *admisibles* dada por el siguiente lema:

LEMA 2.2 Un vector $s \in \{1, 0, -1\}^n$ pertenece a $\mathcal{A}_i(\mathcal{D})$ si, y sólo si existe una función continua $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ tal que:

- i) $\text{sgn} \left(\frac{\partial f_i}{\partial x} \right) = s$ para todo x en el dominio de interés. Esto es, el patrón de signos del Jacobiano de f_i coincide con s .
- ii)' Para cada par de muestras $y, z \in \mathcal{D}_i$ se tiene que $\frac{y_i}{q_i(y_i)} - f_i(y) = \frac{z_i}{q_i(z_i)} - f_i(z)$.

DEMOSTRACIÓN

Probaremos primero que si un vector de signos $s \in \{1, 0, -1\}^n$ satisface estas dos nuevas condiciones, *i*) y *ii*)', entonces es admisible.

Supongamos pues que se satisface *ii*)' para cada par de muestras en \mathcal{D}_i . Entonces existe una constante $c \in \mathbb{R}$ tal que

$$c = \frac{y_i}{q_i(y_i)} - f_i(y) = \frac{z_i}{q_i(z_i)} - f_i(z).$$

Definamos ahora $\bar{f}_i = f_i + c$, claramente

$$s = \text{sgn} \left(\frac{\partial \bar{f}_i}{\partial x} \right) = \text{sgn} \left(\frac{\partial f_i}{\partial x} \right).$$

Más aún, observemos que

$$\frac{y_i}{q_i(y_i)} = f_i(y) + \frac{z_i}{q_i(z_i)} - f_i(z) = f_i(y) + c = \bar{f}_i(y).$$

Así pues, es precisamente \bar{f}_i la función que hace admisible a s para \mathcal{D}_i .

Por otro lado, supongamos que $s \in \mathcal{A}_i(\mathfrak{D})$. Tenemos ya, por hipótesis, que se satisface i). Así mismo, la segunda condición implica que

$$\frac{\dot{y}_i}{q(y_i)} - f_i(y) = 0,$$

y por lo tanto $\frac{\dot{y}_i}{q(y_i)} - f_i(y) = \frac{\dot{z}_i}{q_i(z_i)} - f_i(z)$ para cada par de muestras $y, z \in \mathfrak{D}_i$. □

Este último resultado da pie al teorema con el que cerramos esta sección y caracteriza *todos* los vectores de signos admisibles para una muestra.

TEOREMA 2.2 Un vector de signos s es admisible para el agente i dado un conjunto de datos \mathfrak{D} si, y sólo si, es compatible para el mismo.

DEMOSTRACIÓN

El **LEMA 2.1** establece ya que $\mathcal{A}_i \subset \mathcal{S}_i$. Por lo tanto basta probar que $\mathcal{S}_i \subset \mathcal{A}_i$ para concluir la demostración. Para ello, se procede por contradicción.

Supongamos que existe $s \in \mathcal{S}_i$ tal que $s \notin \mathcal{A}_i$. Tenemos pues, que para toda función $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ cuyo Jacobiano tiene signo s (i.e., $\text{sign}(\frac{\partial f_i}{\partial x}) = s$) al satisfacer la primera condición del **LEMA 2.2** no puede satisfacer la segunda. Esto es, si $\text{sign}(\frac{\partial f_i}{\partial x}) = s$, entonces existe un par de puntos $y^f, z^f \in \mathfrak{D}_i$ para los cuales

$$\frac{\dot{y}_i}{q_i(y_i)} - f_i(y) \neq \frac{\dot{z}_i}{q_i(z_i)} - f_i(z) \quad (2.5)$$

Así, al ser s compatible, por definición, tenemos que para cada par de muestras distintas $y, z \in \mathfrak{D}_i$ existe un vector $v_i(y, z)$ que satisface $s = \text{sgn}(v_i(y, z))$ y

$$\beta_i(y, z) = v_i(y, z) \cdot (y - z), \text{ donde } \beta_i(y, z) = \frac{\dot{y}_i}{q_i(y_i)} - \frac{\dot{z}_i}{q_i(z_i)}.$$

Notemos que esto se satisface si, y sólo si,

$$\begin{aligned} \frac{\dot{y}_i}{q_i(y_i)} - \frac{\dot{z}_i}{q_i(z_i)} &= v_i(y, z) \cdot y - v_i(y, z) \cdot z \\ \frac{\dot{y}_i}{q_i(y_i)} - v_i(y, z) \cdot y &= \frac{\dot{z}_i}{q_i(z_i)} - v_i(y, z) \cdot z \end{aligned}$$

Ya que esto se satisface para todas las muestras $y, z \in \mathfrak{D}_i$, entonces existe una constante $c = \frac{\dot{y}_i}{q_i(y_i)} - v_i(y, z) \cdot y$ para todo $y \in \mathfrak{D}_i$. Y así,

$$v_i(y, z) \cdot y = \frac{\dot{y}_i}{q_i(y_i)} - c. \quad (2.6)$$

Ahora bien, seleccionando $w \in \mathfrak{D}_i$ arbitrario y denotando por $y^*(x)$ al punto en \mathfrak{D}_i más cercano a x , se puede definir la función $j_i(x) =$

$v_i(y^*(x), w)$. Por construcción $j_i(x)$ es constante por pedazos y, por lo tanto, es integrable. Sea f_i la función que resulta de integrar $j_i(x)$. Por definición tenemos que $f_i(x) = v_i(y^*(x), w) \cdot x$. Más aún, por construcción se tiene que $s = \text{sgn}(v_i(y^*(x), w)) = \text{sgn}\left(\frac{\partial f_i}{\partial x}(x)\right)$.

Así pues, la hipótesis $s \notin \mathcal{A}_i$ implica que existen $y, z \in \mathcal{D}_i$ para los cuales se satisface la Ec. (2.5). Notemos que usando la definición de f_i y la Ec. (2.6) en estos dos puntos obtenemos que,

$$\begin{aligned} f_i(y) - f_i(z) &= v_i(y^*(y), w) \cdot y - v_i(y^*(z), w) \cdot z \\ &= v_i(y, w) \cdot y - v_i(z, w) \cdot z \\ &= \left(\frac{\dot{y}_i}{q_i(y_i)} - c \right) - \left(\frac{\dot{z}_i}{q_i(z_i)} - c \right) \\ &= \frac{\dot{y}_i}{q_i(y_i)} - \frac{\dot{z}_i}{q_i(z_i)}. \end{aligned}$$

Lo cual se contradice con la Ec. (2.5). Así, se concluye por contradicción que si $s \in \mathcal{S}_i$ entonces $s \in \mathcal{A}_i$. \square

Como consecuencia de este resultado, es posible calcular los vectores de signos admisibles, para el agente i , para un conjunto de datos \mathcal{D} . Para ello basta encontrar los vectores de signos compatibles para cada par de muestras en \mathcal{D}_i y quedarse con los vectores comunes a todas los pares de muestras.

OBSERVACIÓN 2.2 Es posible decidir si un vector $w \in \mathbb{R}^n$ satisface la ecuación $\beta_i(y, z) = w \cdot (y - z)$. Por lo tanto es posible calcular a partir de los datos cuales patrones de signos $s \in \{1, 0, -1\}^n$ son admisibles. En otras palabras, $\mathcal{A}_i(\mathcal{D}) = \mathcal{S}_i(\mathcal{D})$ es un *estadístico* de los datos pues puede ser construido utilizando los pares de muestras únicamente.

Más aún, como se muestra en la siguiente sección, es posible calcular todos los vectores admisibles a partir de los datos sin necesidad de construir los vectores $v_i(y, z)$ de manera explícita.

2.2 ALGORITMO PARA LA INFERENCIA DE REDES ADMISIBLES

Recordando que un vector $s \in \{1, 0, -1\}^n$ es admisible para un conjunto de datos si, y sólo si, es compatible para cada par de muestras en él, se concluye que para realizar la inferencia basta calcular todos los vectores admisibles para cada par $y, z \in \mathcal{D}_i$ y luego quedarse con los signos que son admisibles para todos los pares de muestras. Denotaremos por $S_i(y, z)$ a todos los vectores de signos admisibles para el agente i de un par de muestras $y, z \in \mathcal{D}_i$.

Esto es, el resultado principal de la sección anterior provee de un criterio constructivo para saber si un patrón de signos es admisible para una par de muestras. Por tanto, para realizar la inferencia, basta

Se estima que en el universo hay un total de 1×10^{24} estrellas ¹ Mientras que el número de átomos en el universo se estima entre 10^{70} y 10^{80} . ²

con aplicar este criterio para cada ortante del espacio. Sin embargo, este enfoque se vuelve rápidamente imposible de aplicar, puesto que el número de ortantes en \mathbb{R}^n es 3^n . Por ejemplo, para cada par de muestras de un sistema con 50 agentes habría que revisar un número cercano al estimado total de estrellas en el universo. Y para un sistema con 170 agentes, el número de ortantes 3^{170} es más grande ya que el número estimado de átomos en el universo.

Más aún, al considerar que para cada uno de los ortantes sería necesario probar que existe, o no, el vector $v_i(y, z)$ correspondiente, se vuelve evidente la necesidad de obtener un método más eficiente para realizar la inferencia. Por esta razón, esta sección se centra en construir resultados que permiten obtener todos los ortantes admisibles para un par de muestras $y, z \in \mathcal{D}_i$ sin necesidad de explorar el espacio de ortantes. Más aún, mostraremos como se puede obtener cada vector de signos s admisible sin necesidad de determinar los vectores $v_i(y, z)$ asociados a cada uno.

Para ello, recordemos que un vector de interacciones s es admisible para $y, z \in \mathcal{D}_i$ si, y sólo si existe un vector $v_i(y, z)$ tal que $s = \text{sgn}(v_i(y, z))$ y que satisfaga la ecuación lineal

$$\beta_i(y, z) = v_i(y, z) \cdot (y - z).$$

Así pues, la búsqueda de vectores de interacción se reduce a observar los ortantes que contienen al hiperplano de soluciones de la ecuación lineal correspondiente. En particular, si $\beta_i(y, z) \neq 0$ (es decir si las muestras no están en equilibrio y se encuentran en posición general), esta ecuación se puede reescribir como

$$1 = v_i(y, z) \cdot \frac{y - z}{\beta_i(y, z)}.$$

El siguiente resultado provee una caracterización constructiva de los patrones de signos asociados a las soluciones de esta familia de ecuaciones lineales. Denotaremos por $\mathcal{O}^n = \{1, 0, -1\}^n$ al espacio de ortantes de \mathbb{R}^n .

TEOREMA 2.3 Sea $w \in \mathbb{R}^n$ distinto de cero. Definimos al conjunto de signos $Q(w)$ asociados a w como

$$Q(w) = \{s \in \mathcal{O}^n \mid \exists v \in \mathbb{R}^n \text{ tal que } v \cdot w = 1 \text{ y } \text{sgn}(v) = s\}.$$

Entonces,

- a) $s \in Q(w)$ si, y sólo si, existe $w_i \neq 0$ tal que $\text{sgn}(w_i) = s_i$.
- b) Además, si w tiene $k > 0$ entradas no cero, entonces

$$|Q(w)| = \sum_{i=1}^k 3^{n-i} 2^{i-1}.$$

DEMOSTRACIÓN

Sea $s \in \mathcal{O}^n$ tal que $s_i \neq \text{sgn}(w_i)$ para todo $w_i \neq 0$ y sea $v \in \mathbb{R}^n$ tal que $s = \text{sgn}(v)$. Es fácil ver entonces que $v_i w_i \leq 0$ para todo $i \leq n$ y, así,

$$vw = \sum_{i=1}^n v_i w_i \leq 0.$$

Con lo cual $v \cdot w \neq 1$ para todo v tal que $s = \text{sgn}(v)$, y por lo tanto, $s \notin Q(w)$. De esta manera, si $s \in Q(w)$, entonces existe $w_i \neq 0$ tal que $s_i = \text{sgn}(w_i)$.

Por otro lado, consideremos $s \in \mathcal{O}^n$ tal que $s_i = \text{sgn}(w_i)$ para algún $w_i \neq 0$. Sean p el número de entradas no cero de w tales que $s_i = \text{sgn}(w_i)$, q el número de entradas no cero tales que $s_i = -\text{sgn}(w_i)$ y r el número de entradas cero de w . Definamos el siguiente vector v como sigue,

$$v_i = \begin{cases} 0 & \text{si } s_i = 0 \\ \frac{q+1}{p w_i} & \text{si } s_i = \text{sgn}(w_i) \\ -\frac{1}{w_i} & \text{si } s_i = -\text{sgn}(w_i) \end{cases}$$

Por construcción tenemos que $\text{sgn}(v) = s$, observemos además que,

$$\begin{aligned} v \cdot w &= \sum_{i=1}^n v_i w_i \\ &= \sum_{j=1}^p v_{i_j} w_{i_j} + \sum_{j=1}^q v_{i_j} w_{i_j} + \sum_{j=1}^r v_{j_j} w_{j_j} \\ &= \sum_{j=1}^p v_{i_j} w_{i_j} + \sum_{j=1}^q v_{i_j} w_{i_j} \\ &= \sum_{j=1}^p \frac{q+1}{p w_{i_j}} w_{i_j} + \sum_{j=1}^q -\frac{1}{w_{i_j}} w_{i_j} \\ &= \frac{q+1}{p} \sum_{j=1}^p \frac{w_{i_j}}{w_{i_j}} - \sum_{j=1}^q \frac{w_{i_j}}{w_{i_j}} \\ &= \left(\frac{q+1}{p} \right) p - q \\ &= 1, \end{aligned}$$

y por lo tanto, $s \in Q(w)$. De esta manera, se tiene que $s \in Q(w)$ si, y sólo si, $s = \text{sgn}(w_i)$ para algún $w_i \neq 0$.

Ahora bien, para calcular $|Q(w)|$ basta pues contar cuantos elementos de $s \in \mathcal{O}^n$ satisfacen que $s = \text{sgn}(w_i)$ para algún $w_i \neq 0$. Sea k el número de entradas no cero de w . Sin pérdida de generalidad, supongamos que son las primeras k entradas. Luego, para cada $i \leq k$ definamos

$$S_i = \{s \in \mathcal{O}^n \mid s_i = \text{sgn}(w_i) \text{ y } s_j \neq \text{sgn}(w_j) \text{ para toda } j < i\}.$$

Es fácil ver que $|Q(w)| = \sum_{i=1}^k |S_i|$. Ahora bien, observemos que

$$|S_1| = 3^{n-1},$$

ya que corresponde a contar los elementos $s \in \mathcal{O}^n$ tal que $s_1 = \text{sgn}(w_1)$. Análogamente,

$$|S_2| = 3^{n-1}2$$

puesto que hay únicamente dos opciones para elegir para la posición uno, y tres opciones para colocar en las entradas de la 3 a la n . En general tenemos que,

$$|S_i| = 3^{n-i}2^{i-1}.$$

Por lo tanto, se sigue que

$$|Q(w)| = \sum_{i=1}^k 3^{n-i}2^{i-1},$$

completando la prueba. □

Por otro lado, si $\beta_i(y, z) = 0$ la ecuación a resolver es

$$0 = v_i(y, z) \cdot (y - z).$$

Esto es, buscamos los signos asociados a los vectores ortogonales a $y - z$. Estos signos quedan caracterizados de manera constructiva en el siguiente resultado.

TEOREMA 2.4 Sea $w \in \mathbb{R}^n$ distinto de cero. Definimos el conjunto de *signos ortogonales a w* como,

$$P(w) = \{s \in \mathcal{O}^n \mid \exists v \in \mathbb{R}^n \text{ tal que } v \cdot w = 0 \text{ y } \text{sgn}(v) = s\}.$$

Entonces,

a) $s \in P(w)$ si, y sólo si, satisface una de las siguientes condiciones:

- i) $s_i = 0$ siempre que $w_i \neq 0$, ó
- ii) Existen índices i_+ e i_- tales que

$$s_{i_+} \text{sgn}(w_{i_+}) = 1 \text{ y } s_{i_-} \text{sgn}(w_{i_-}) = -1.$$

b) Además, si w tiene $k > 0$ entradas distintas de cero, se tiene que

$$|P(w)| = 3^{n-k} + 3^{n-k} \sum_{m=1}^{k-1} \binom{k}{m} \sum_{n=1}^{k-m} \binom{k-m}{n}$$

DEMOSTRACIÓN Supongamos, sin pérdida de generalidad, que únicamente las primeras k entradas de w son distintas de cero. Sea $s \in \mathcal{O}^n$

y supongamos que $s_i = 0$ para todo $i \leq k$. Claramente s es ortogonal a w , pues para todo vector v que satisface $\text{sgn}(v) = s$ se tiene que

$$v \cdot w = \sum_{i=1}^n v_i w_i = \sum_{i=1}^k v_i w_i + \sum_{i=k+1}^n v_i w_i = \sum_{i=1}^k 0 w_i + \sum_{i=k+1}^n v_i 0 = 0.$$

Por otra parte, si suponemos que $s_i \neq 0$ para algún $i \leq k$, tenemos tres casos:

- i) Supongamos, sin pérdida de generalidad, que $s_1 \text{sgn}(w_1) = 1$ y que $s_i w_i \geq 0$ para todo $i \in \{2, 3, \dots, k\}$. Luego, sea v cualquier vector tal que $\text{sgn}(v) = s$ tenemos que,

$$v \cdot w = v_1 w_1 + \sum_{i=2}^k v_i w_i > v_1 w_1 > 0.$$

Y por lo tanto, s no es ortogonal a w .

- ii) Supongamos ahora que $s_1 \text{sgn}(w_1) = -1$ y $s_i w_i \leq 0$ para todo $i \in \{2, 3, \dots, k\}$. De manera análoga se obtiene que para cualquier vector v tal que $\text{sgn}(v) = s$ se tiene que

$$v \cdot w < v_1 w_1 < 0,$$

y por lo tanto s no es ortogonal a w .

- iii) Por último supongamos, nuevamente sin pérdida de generalidad, que existen M, N tal que $M + N \leq k$ para los cuales se satisface que $s_i \text{sgn}(w_i) = 1$ para todo $i \leq M$, $s_i \text{sgn}(w_i) = -1$ para todo $M + 1 \leq i \leq N + M$ y $s_i = 0$ para $M + N + 1 \leq i \leq k$. Se puede construir un vector v tal que $\text{sgn}(v) = s$ como sigue:

$$v_i = \begin{cases} \frac{1}{w_i M} & \text{si } i \leq M \\ \frac{-1}{w_i N} & \text{si } M + 1 \leq i \leq N + M \\ 0 & \text{si } M + N + 1 \leq i \leq k \\ s_i & \text{si } i > k \end{cases}$$

Por construcción $s = \text{sgn}(v)$, ahora bien, observemos que

$$\begin{aligned}
v \cdot w &= \sum_{i=1}^n v_i w_i \\
&= \sum_{i=1}^M v_i w_i + \sum_{i=M+1}^{M+N} v_i w_i + \sum_{i=M+N+1}^k v_i w_i + \sum_{i=k+1}^n v_i w_i \\
&= \sum_{i=1}^M \frac{w_i}{w_i M} + \sum_{i=M+1}^{M+N} \frac{-w_i}{w_i N} \\
&= \sum_{i=1}^M \frac{1}{M} + \sum_{i=M+1}^{M+N} \frac{-1}{N} \\
&= M \left(\frac{1}{M} \right) + N \left(\frac{-1}{N} \right) \\
&= 1 - 1 \\
&= 0.
\end{aligned}$$

Con lo cual, se concluye que $s \in P(w)$.

De esta manera, queda demostrado el inciso a).

Ahora bien, para calcular $|P(w)|$ basta contar los vectores asociados para las dos clases dadas por el inciso anterior.

Notemos primero que los vectores de signos que satisfacen $s_i = 0$ para todo $i \leq k$, no tienen restricciones para las $n - k$ entradas restantes, por lo tanto cualquier combinación de $1, 0, -1$ en las últimas $n - k$ entradas resulta ser ortogonal a w , con lo cual tenemos un total de 3^{n-k} vectores de signo de esta clase.

Por otro lado, para contar cuantos vectores de signos pertenecen a la segunda clase es necesario observar que para obtener uno de ellos es equivalente hacer una partición de las entradas no cero de w en tres conjuntos, I_M, I_N, I_0 , tal que $|I_M|, |I_N| > 0$, que corresponden a los vectores con el mismo signo que las entradas de w , con el signo opuesto y las entradas cero respectivamente. Asimismo, para cada combinación de las entradas 1 a k tenemos 3^{n-k} maneras diferentes de completar el vector de signos correspondientes a las entradas cero de w . Así, para contar cuantos vectores de esta clase hay, es necesario primero ver todas las combinaciones posibles de I_M y de I_N , que resultan ser

$$\sum_{m=1}^{k-1} \sum_{n=1}^{k-m} \binom{k}{m} \binom{k-m}{n} 3^{n-k}$$

Así, se concluye que

$$|P(w)| = 3^{n-k} + 3^{n-k} \sum_{m=1}^{k-1} \binom{k}{m} \sum_{n=1}^{k-m} \binom{k-m}{n}.$$

□

Notemos que al sustituir w por $\frac{y-z}{\beta_i(y,z)}$ y $y - z$ en los Teoremas [TEOREMA 2.3](#) y [TEOREMA 2.4](#) respectivamente, la primera implicación

de cada uno de los teoremas permite obtener explícitamente $S(y, z)$ sin necesidad de revisar todo el espacio de ortantes ni construir los vectores $v(y, z)$ asociados en cada caso. De manera concreta, se tiene entonces el Algoritmo 1 para construir todas las redes de interacciones admisibles para cada uno de los agentes. Luego, las matrices de interacciones se obtienen tomando un vector de signos de $S_i(\mathcal{D}_i)$ para cada renglón i .

Algorithm 1 Algoritmo Base para calcular S_i

Input: \mathcal{D}_i : conjunto de muestras con actividad del agente i .

Output: S_i : Lista de vectores de signos admisibles para el agente i .

for cada par $y, z \in \mathcal{D}_i$ **do**

 Calcular $S_i(y, z)$ usando la caracterización dada por [TEOREMA 2.3](#) y [TEOREMA 2.4](#).

return $S_i = \bigcap_{y, z \in \mathcal{D}_i} S_i(y, z)$.

Es fácil ver que el tamaño de $S(y, z)$ depende de la cantidad de ceros en $y - z$ y del número de agentes del sistema. Como se observa en la Figura 2.1, el número de vectores de signos admisibles para una pareja de muestras es considerablemente alto incluso cuando la mayoría de las entradas de $y - z$ son cero. Esto puede hacer muy poco eficiente el proceso de inferencia.

Para hacer más eficiente este proceso denotemos por p_1, \dots, p_M a las parejas de elementos de \mathcal{D}_i . Observemos que el conjunto que buscamos inferir es la intersección

$$\bigcap_{k=1}^M S_i(p_k),$$

la cual está contenida en $S_i(p_1)$. Por lo cual no es necesario construir explícitamente el resto de los $S_i(p_k)$. Basta empezar con $S_i(p_1)$ y luego para cada pareja $p_k, k > 1$, quedarse con los elementos que también están en $S_i(p_k)$. Este proceso lo denotaremos como *Algoritmo secuencial para calcular S_i* y se presenta de manera explícita en el Algoritmo 2.

Algorithm 2 Algoritmo Secuencial para calcular S_i

Input: p_1, \dots, p_M parejas de elementos de \mathcal{D}_i

Output: S_i : Lista de vectores de signos admisibles para el agente i .

 Construir $S = S(p_1)$

for i desde 2 hasta M **do**

$S = \{s \in S \mid s \in S(p_i)\}$

return S

Notemos que el Algoritmo 2 reduce por $M - 1$ la cantidad de vectores de signos que es necesario calcular. Pero, como muestra la Figura 2.1, incluso construir los elementos de $S(p_1)$ de manera explícita puede representar un reto computacional considerable en el momento de la implementación.

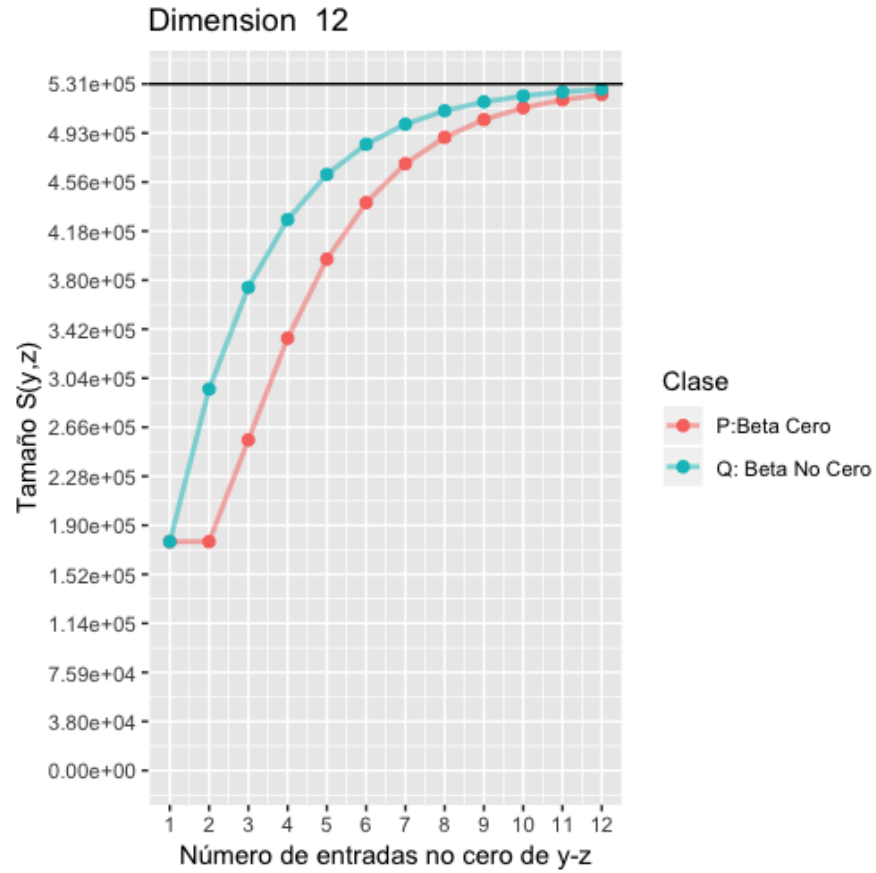


Figura 2.1: Conteo de vectores de signos contenidos en $S(y,z)$ en función de las entradas no cero de $w(y,z)$ para un sistema con $N = 12$ agentes.

Para resolver este problema, a continuación se presenta una estrategia para *representar* y *codificar* múltiples vectores de signos de manera simultánea. Esta codificación se basa en introducir la siguiente definición de un *patrón de signos*:

DEFINICIÓN 2.4 (PATRÓN DE SIGNOS) Un patrón de signos $pat \in \{1, 0, -1, 3\}^n$ representa todos los vectores de signos s que pueden obtenerse al sustituir todas las entradas $p_i = 3$ por cualquier combinación de elementos de $\{1, 0, -1\}$. A las entradas distintas de 3 de un patrón les llamaremos la parte *fija* del patrón.

Por ejemplo, el patrón de signos $p = (1, 0, -1, 3, 3, 3)$ representa un total de 18 vectores de signos, específicamente todos aquellos cuyas tres primeras entradas son $1, 0, -1$ respectivamente.

Decimos que pat_2 es una *extensión* de pat_1 si la parte fija de pat_1 está contenida en la parte fija de pat_2 . Esto es, ambos patrones coinciden en todas las entradas distintas de 3 de pat_1 .

Notemos que esta representación está inspirada en las condiciones necesarias para que los vectores de signos sean admisibles para un par de muestras dado y, z . Esto es debido a que es posible agrupar

todos los vectores de signo en $S(y, z)$ en patrones de signos cuya parte fija corresponde precisamente a las entradas no cero de $y - z$.

Esta representación resulta ser de utilidad debido a que todo vector de signos admisible para una pareja y, z está determinado por las entradas no cero de $y - z$. Por lo tanto, no es necesario considerar la información redundante contenida en las entradas correspondientes a las entradas cero de $y - z$ de los vectores de signos asociados.

De esta manera, definimos que un patrón pat es *admisible* para un par de muestras y, z si la parte fija de pat contiene a las entradas no cero de $y - z$ y todos los vectores de signos representados por este patrón son admisibles para y, z .

Denotaremos por $Pat(y, z)$ al conjunto de patrones admisibles que codifican a todos los vectores en $S(y, z)$. Esta representación de los vectores de signos permite además ampliar el proceso iterativo de quedarse sólo con los elementos de $Pat(y, z)$ que son admisibles para otro par de muestras.

Supongamos que se tienen los signos de $S(p_1)$ expresados en forma de patrones $Pat(p_1)$. Si lo que se busca es quedarse únicamente con los patrones de signos que representan a los vectores de signos que pertenecen a $S(p_2) \cap S(p_1)$ basta *adaptar* los patrones en $Pat(p_1)$ para que codifiquen a los vectores de signos admisibles para p_1 y para p_2 . Esto es, basta *adaptar* cada patrón de signos en $Pat(p_1)$. Con más precisión:

DEFINICIÓN 2.5 (ADAPTAR UN PATRÓN) Sea pat un patrón admisible para un par de muestras p_1 . Adaptar a pat para que también sea admisible para otro par p_2 corresponde a realizar una, y sólo una, de las siguientes acciones:

- a) Si el patrón pat ya es admisible para p_2 , adaptarlo significa conservarlo sin cambios.
- b) Si el patrón no es admisible para p_2 , hay dos opciones:
 - i) Es posible *extender* de una, o más maneras, el patrón pat para hacerlo admisible con p_2 . En ese caso, adaptarlo implica sustituir a pat por todas las extensiones de pat que sean admisibles para p_2 .
 - ii) Ninguna extensión de pat es admisible para p_2 . En este caso, no es posible adaptarlo y se descarta el patrón.

Notemos que es posible calcular los patrones admisibles para un conjunto de datos dado de manera análoga al Algoritmo 2. Definimos este proceso de manera formal en el Algoritmo 3.

La utilidad de recurrir a los patrones de signos para calcular los vectores de signos admisibles a un conjunto de datos queda de manifiesto en la Figura 2.2. En esta figura se puede observar que disminuye considerablemente los requerimientos computacionales con respecto al Algoritmo 2.

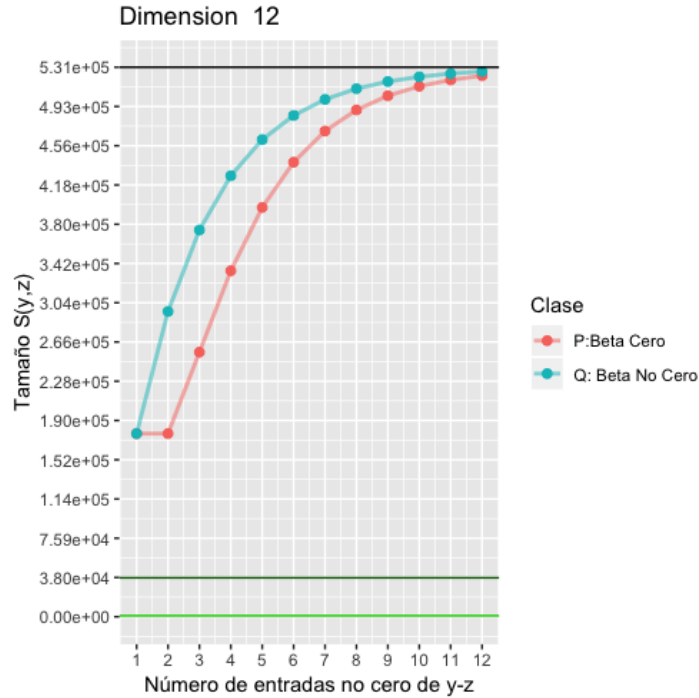
Algorithm 3 Calcular patrones admisibles**Input:** p_1, \dots, p_M parejas de elementos de \mathcal{D}_i **Output:** S_i : Lista de patrones admisibles para el agente i .Construir \mathcal{P} el conjunto de patrones admisibles para p_1 **for** i desde 2 hasta M **do**Adaptar los patrones en \mathcal{P} para que sean admisibles a p_i .**return** S 

Figura 2.2: Conteo de vectores de signos contenidos en $S(y, z)$ en función de las entradas no cero de $w(y, z)$ para un sistema con $N = 12$ agentes comparado con el número máximo de patrones (verde oscuro) y el tamaño promedio de los patrones (verde claro) que se alcanzaron al analizar muestras provenientes de un sistema de 12 agentes cuyos detalles se exponen en el siguiente capítulo.

Se realizó una implementación en R [23] de los algoritmos expuestos en esta sección. Se realizó también, como prueba de concepto, una serie de validaciones numéricas que se exponen a detalle en el CAPÍTULO 3.

Debido a la novedad del enfoque de la inferencia presentada, es necesario definir *indicadores ad hoc* que permitan cuantificar la calidad del proceso de inferencia, es decir de la información que es posible obtener de un conjunto de datos.

2.3 INDICADORES DE LA INFERENCIA

Una vez realizado el proceso de inferencia para cada agente, el proceso arroja n conjuntos $\{S_1, \dots, S_n\}$ correspondientes a los vectores de signos admisibles para cada agente. El número de vectores en cada uno de esos conjuntos depende de se presentalos datos \mathcal{D}_i . Así, las matrices de interacciones corresponden a todas las combinaciones que se pueden obtener al elegir para el renglón correspondiente los vectores de signo de S_i .

Dado que el tamaño y las características de cada S_i están determinados por el conjunto de datos, es importante enfatizar que el número de matrices de interacciones admisibles que se obtienen en el proceso de inferencia está determinado *solamente* por el conjunto de datos. En ese sentido, es posible que existan conjuntos de datos más informativos que otros. Es decir, datos a los cuales sólo es posible explicarlos mediante una única matriz de interacciones (*i.e.*, una única matriz de interacciones admisible), y datos que sean compatibles con una gran cantidad de matrices de interacciones (*i.e.*, existen más de dos matrices de interacciones admisibles). En general, dado un conjunto de datos \mathcal{D}_i hay tres casos:

1. $|S_i| > 0$ para todo $1 \leq i \leq n$ y existe j tal que $|S_j| > 1$. En este caso, hay más de una matriz de interacciones admisible para los datos.
2. $|S_i| = 1$ para todo $1 \leq i \leq n$. En este caso existe una única matriz de interacciones admisible para los datos.
3. Existe j tal que $|S_j| = 0$. Este caso puede ocurrir por dos razones: ya sea que $|\mathcal{D}_i| = 0$, ó que no se satisface alguna de las hipótesis necesarias (e.g. el Jacobiano cambia de signo). En este caso, no existe ninguna matriz de interacciones admisible para los datos.

Con el objetivo de medir la información que puede obtenerse de un conjunto de datos se proponen dos nuevos indicadores para la inferencia de redes. El primero de ellos, *informatividad*, cuantifica precisamente cuantas matrices de interacciones es posible asociar a un conjunto de datos.

DEFINICIÓN 2.6 (INFORMATIVIDAD) Sea \mathcal{D} un conjunto de datos y S_1, \dots, S_n los conjuntos de vectores de signos admisibles para cada agente. La *informatividad del i -ésimo agente*, $I_i(\mathcal{D}_i)$ se define como

$$I_i(\mathcal{D}_i) = \begin{cases} \frac{1}{|S_i|} & \text{si } |S_i| > 0, \\ 0 & \text{si } |S_i| = 0. \end{cases}$$

De manera análoga, la *informatividad de una muestra de datos*, $I(\mathcal{D})$, se define como

$$I(\mathcal{D}) = \prod_{i=1}^n I_i(\mathcal{D}_i).$$

Por definición se tiene que $I(\mathcal{D}) \in [0, 1]$. Además, observemos que si dos conjuntos de datos $\mathcal{D}_1, \mathcal{D}_2$ satisfacen que $\mathcal{D}_1 \subset \mathcal{D}_2$, entonces $I(\mathcal{D}_1) \leq I(\mathcal{D}_2)$. Notemos también que el algoritmo construido en la sección anterior permite, por primera vez, cuantificar la informatividad de un conjunto de datos dado, ya que se tiene una manera de calcular todas las matrices de interacciones admisibles para los datos. Por último, observemos que la matriz de interacciones puede ser determinada de manera única a partir de un conjunto de datos \mathcal{D} si, y sólo si, $I(\mathcal{D}) = 1$.

Ahora bien, incluso en el caso cuando $I(\mathcal{D}) < 1$, ya sea por que los datos son compatibles con más de una matriz de interacciones o bien porque no contienen información sobre uno o más agentes, puede darse el caso de que la matriz de interacciones asociada al sistema dinámico que dió origen a las muestras sea recuperada de manera *parcial*. Para ello es necesario observar que si hay entradas idénticas en todas las posibles matrices de interacciones compatibles, entonces las interacciones correspondientes a estas entradas han sido *unívocamente determinadas*. Así, independientemente de cual sea la matriz de interacciones asociada al sistema dinámico particular que dió origen a los datos, es posible reconstruir parcialmente la matriz de interés.

DEFINICIÓN 2.7 (INTERACCIONES UNÍVOCAMENTE INFERIDAS) Sea \mathcal{D} un conjunto de datos y S_1, \dots, S_n los conjuntos de vectores de signos admisibles para cada agente. Si S_i es no vacío, definimos al conjunto de *interacciones unívocamente determinadas* para el i -ésimo agente como:

$$u(S_i) = \{(s_j, k) | s_j = s'_j = k, \forall s, s' \in S_i\}.$$

Las entradas unívocamente inferidas permiten tener absoluta certeza sobre estas interacciones, pues estas son compartidas por todas las redes admisibles con los datos. De esta manera, al considerar las interacciones unívocamente inferidas de todos los agentes es posible reconstruir parcialmente la matriz de interacciones asociada al sistema dinámico del que provienen los datos. Para cuantificar la proporción de la matriz de interacciones auténtica que es posible inferir proponemos el indicador de *unicidad*.

DEFINICIÓN 2.8 (UNICIDAD) Sea \mathcal{D} un conjunto de datos y S_1, \dots, S_n los conjuntos de vectores de signos admisibles para cada agente. Definimos al indicador de *unicidad* $U(\mathcal{D})$ como

$$U(\mathcal{D}) = \frac{\sum_{i=1}^n |u(S_i)|}{n^2}.$$

De manera análoga a la informatividad, si dos conjuntos de datos $\mathcal{D}_1, \mathcal{D}_2$ satisfacen que $\mathcal{D}_1 \subset \mathcal{D}_2$, entonces $U(\mathcal{D}_1) \leq U(\mathcal{D}_2)$. Además,

observemos que al normalizar el número de entradas unívocamente determinadas mediante n^2 , el total de interacciones a inferir, este indicador sólo alcanza 1 cuando hay una sola matriz de interacciones admisible para la muestra de datos, esto es $U(\mathcal{D}) = 1$ si, y sólo si $I(\mathcal{D}) = 1$.

Es importante enfatizar que es la primera vez que es posible inferir con completa confiabilidad la naturaleza de las interacciones pues se están considerando todas las matrices de interacciones admisibles para un conjunto de datos dado.

2.4 LIMITACIONES FUNDAMENTALES

Con el objetivo de poder describir y estudiar con mayor detalle las características de un sistema dinámico, resulta de interés caracterizar los conjuntos de datos que tienen una única matriz de interacciones admisible. Esto permitiría buscar conjuntos de datos que permitan maximizar tanto la unicidad como la informatividad. Aunque describir las características necesarias para que esto suceda queda fuera del alcance de esta tesis, si daremos condiciones necesarias sobre un conjunto de datos para que sea posible asignarle una única matriz de interacciones. Estas condiciones ó *limitaciones fundamentales* no dependen del algoritmo de inferencia utilizado, ni de la dimensión del sistema. En ese sentido son limitaciones fundamentales intrínsecas de cada conjunto de datos.

Para ello, notemos primero que si el conjunto de datos contiene únicamente muestras en equilibrio, entonces todos los vectores de signos admisibles para un agente i se calculan bajo los criterios proporcionados por el [TEOREMA 2.4](#) y como consecuencia se tiene el siguiente corolario.

COROLARIO 2.1 Sean $y, z \in \mathcal{D}_i$ tales que $\beta_i(y, z) = 0$. Entonces, para cada $s \in S(y, z)$ se tiene que $-s \in S(y, z)$.

DEMOSTRACIÓN

Consideremos $s \in S(y, z)$. Dado que $\beta_i(y, z) = 0$, por el [TEOREMA 2.4](#), se tienen dos casos:

1. O bien $s_i = 0$ para todas las entradas i tales que $(y - z)_i \neq 0$. Y, en este caso es inmediato que $-s \cdot (y - z) = s \cdot (y - z)$ y por lo tanto $-s \in S(y, z)$.
2. O bien, existen índices i_+ e i_- tales que $s_{i_+} \text{sgn}((y - z)_{i_+}) = 1$ y $s_{i_-} \text{sgn}((y - z)_{i_-}) = -1$. Observando a $-s$ es claro que

$$-s_{i_-} \text{sgn}((y - z)_{i_-}) = 1 \text{ y } -s_{i_+} \text{sgn}((y - z)_{i_+}) = -1,$$

y por lo tanto $-s \in S(y, z)$.

□

Esta observación, al describir las limitaciones propias de cada par de muestras en equilibrio, permite caracterizar las limitaciones referentes tanto a la informatividad del mismo como las propiedades de las posibles entradas unívocamente determinadas de un conjunto de datos que consiste únicamente en muestras en equilibrio. Es decir:

PROPOSICIÓN 2.1 Consideremos un conjunto de datos en equilibrio \mathcal{D}_i y sea s el vector de interacciones de la dinámica de la que provienen los datos. Si $s \neq 0$ se tiene que $|S_i| < 1$ y por tanto es imposible inferir unívocamente todas las interacciones del i -ésimo agente.

DEMOSTRACIÓN

Sabemos ya que $s \in S_i(\mathcal{D}_i)$, ya que $s \neq 0$ entonces $s \neq -s$, y como consecuencia del **COROLARIO 2.1** se tiene que $s, -s \in S_i(\mathcal{D}_i)$. Así $|S_i| > 1$.

□

Este comportamiento se debe principalmente a que las únicas entradas de la matriz de interacciones que pueden ser unívocamente determinadas a partir de un conjunto de datos en equilibrio son las entradas cero de la matriz.

PROPOSICIÓN 2.2 Si todas las muestras en un conjunto \mathcal{D}_i están en equilibrio y s_j es una entrada unívocamente determinado, entonces $s_j = 0$.

DEMOSTRACIÓN

Supongamos que la entrada j está unívocamente determinada para el agente i . Tomemos $s \in S_i(\mathcal{D}_i)$ arbitrario, por el **COROLARIO 2.1** tenemos que $-s \in S_i(\mathcal{D}_i)$. Así, como la entrada j es unívocamente determinada, tenemos que $s_j = -s_j$ y por lo tanto concluimos que $s_j = 0$.

□

Tenemos pues que las muestras en equilibrio pueden distinguir la ausencia de interacciones, pero no permiten distinguir el signo de las entradas no cero. Observemos las implicaciones concretas en un ejemplo. Supongamos que tenemos un ecosistema con dos especies diferentes que interactúan de la siguiente manera

$$\begin{cases} \dot{x}_1 = x_1 \left(1.5 - x_1 + \frac{0.8x_2}{1+0.1x_2} \right), \\ \dot{x}_2 = x_2 (2.5 - x_2). \end{cases} \quad (2.7)$$

La matriz de interacciones asociada al sistema es

$$S = \begin{pmatrix} -1 & 1 \\ 0 & -1 \end{pmatrix}.$$

Para la inferencia, supondremos que la dinámica del sistema es completamente desconocida y que tenemos disponibles sólo algunas

muestras. Supongamos en un principio que tenemos únicamente tres muestras correspondientes, cada una, a los distintos estados estables del sistema (Fig. 2.3 a.). Notemos que sólo tenemos dos muestras donde estén presentes cada una de las especies (de hecho, para esta dinámica sólo existen estas dos muestras en estado estable correspondientes a los dos equilibrios del sistema). Así, los vectores admisibles para la primera especie resultan ser

$$S_1 = \{(-1, 1), (0, 0), (1, -1)\}.$$

De la misma manera, los vectores de signos admisibles para la especie dos resultan ser

$$S_2 = \{(1, 0), (0, 0), (-1, 0)\}.$$

De manera gráfica se observan ambos conjuntos en la Fig. 2.3 b. Claramente, la única entrada unívocamente determinada es la entrada $s_{2,1} = 0$. Así, utilizando únicamente datos en equilibrio no fue posible determinar la naturaleza exacta de ninguna interacción, pero sí fue posible inferir que la especie 1 no tiene ninguna influencia en la especie 2.

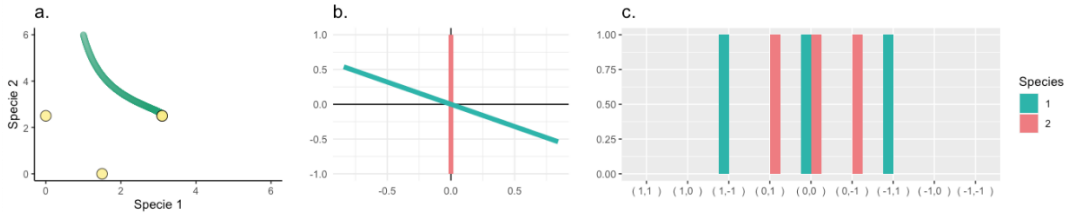


Figura 2.3: Ejemplo en dimensión $N = 2$ de que al considerar sólo muestras en equilibrio no es posible alcanzar unicidad 1. (a) Muestra de un sistema de dos dimensiones que consiste en los tres puntos en equilibrio no triviales. (b) Ortantes considerados para cada una de las especies. (c) Conteo de los ortantes compatibles para cada una de las especies.

Dado que para un conjunto formado únicamente por muestras en equilibrio es imposible asociarle una única matriz de interacciones, podría pensarse que los datos en movimiento son los más informativos y no presentan esta limitante. Sin embargo, resulta ser que si se consideran únicamente muestras dinámicas tampoco garantizan una informatividad de uno. Esto es una consecuencia de la caracterización que proporciona el **TEOREMA 2.3** para los vectores admisibles a las muestras en movimiento. Con más precisión:

COROLARIO 2.2 Sean $y, z \in \mathfrak{D}_i$ tales que $\beta_i(y, z) \neq 0$. Supongamos que existe $s \in S_i(y, z)$ tal que $s_j = 0$. Definimos

$$s^+ = (s_1, \dots, s_{j-1}, 1, s_{j+1}, \dots, s_n) \text{ y } s^- = (s_1, \dots, s_{j-1}, -1, s_{j+1}, \dots, s_n).$$

Entonces se tiene que $s^+, s^- \in S_i(y, z)$.

DEMOSTRACIÓN

Sea $s \in S_i(y, z)$ tal que $s_j = 0$. Luego, por el [TEOREMA 2.3](#) existe $s_l \neq 0$ tal que $\text{sgn}(\delta(y, z))_l = s_l$. Por construcción tenemos que $s_l^+ = s_l = s_l^-$ y por lo tanto concluimos que $s^+, s^- \in S(y, z)$. \square

Este corolario tiene consecuencias concretas en cuanto a la informatividad y las interacciones que pueden ser unívocamente determinadas a partir de un conjunto de muestras dinámicas. En particular, en el caso de que para un agente i de interés existan otros agentes con interacción nula sobre él, *i.e.*, no influyen en la dinámica de x_i , no será posible determinar esto a partir de un conjunto donde todas las muestras son dinámicas.

PROPOSICIÓN 2.3 Consideremos un conjunto de datos \mathfrak{D}_i conformado únicamente por muestras dinámicas. Sea s el vector de interacciones asociado a la dinámica del agente i , si existe $s_j = 0$ entonces $I_i(\mathfrak{D}_i) < 1$.

DEMOSTRACIÓN Por hipótesis, $s \in S_i(\mathfrak{D}_i)$, puesto que es admisible para cada par de muestras en \mathfrak{D}_i . Luego Sean s^+, s^- los vectores que define la j -ésima entrada con la notación utilizada en el [COROLARIO 2.2](#), como consecuencia del mismo corolario tenemos que s^+ y s^- son admisibles para cada par de muestras y por lo tanto $s^+, s^- \in S_i(\mathfrak{D}_i)$, de ésta manera se tiene que $|S_i(\mathfrak{D}_i)| > 1$. Con lo cual concluimos que $I_i(\mathfrak{D}_i) < 1$. \square

Este resultado implica que no es posible alcanzar máxima informatividad (*i.e.*, inferir la matriz de interacciones de manera única) utilizando únicamente muestras dinámicas. Más aún, esto está ligado a las características de las entradas que pueden ser unívocamente determinadas a partir de un conjunto de muestras dinámicas:

PROPOSICIÓN 2.4 Sea \mathfrak{D}_i un conjunto de muestras dinámicas. Supongamos que s_j es una entrada unívocamente determinada para el agente i . Entonces $s_j \neq 0$.

DEMOSTRACIÓN

Es consecuencia directa del [COROLARIO 2.2](#). \square

Las muestras dinámicas permiten distinguir las interacciones positivas de las negativas, pero no detectan la ausencia de interacciones. Volviendo al ejemplo del ecosistema de dos especies con dinámica dada por la Ecuación [2.7](#), si se tuviera un conjunto de cuatro muestras en movimiento, tal como se ve en la Figura [2.4 a](#). De la figura se puede observar que:

$$S_1 = \{(-1, 1), (-1, 0), (-1, -1)\},$$

$$S_2 = \{(1, -1), (0, -1), (-1, -1)\}.$$

En este ejemplo, las únicas entradas unívocamente determinadas resultan ser $s_{11} = -1$ y $s_{21} = -1$. Es posible que agregando más muestras dinámicas se pueda inferir de manera única la interacción s_{12} , pero no importa cuantas muestras dinámicas se tomen, para poder inferir la interacción $s_{21} = 0$ se requiere muestras en equilibrio.

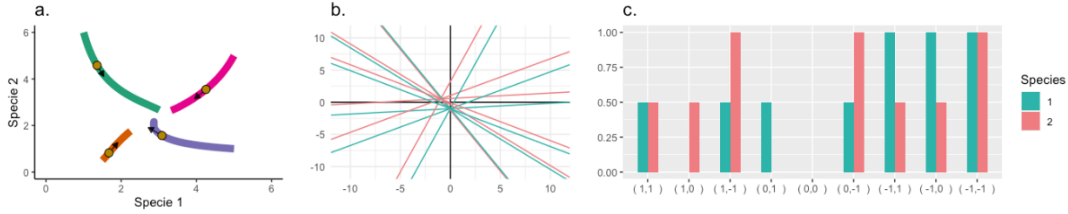


Figura 2.4: Ejemplo en dimensión $n = 2$ de que al considerar sólo muestras en movimiento no es posible alcanzar Unicidad 1. (a) Muestra de un sistema de dos dimensiones que consiste en los cuatro puntos en movimiento. (b) Ortantes compatibles para cada par de muestras y cada una de las especies. (c) Conteo de los ortantes compatibles para cada una de las especies. Se observa que las muestras en movimiento hacen que cada especie tenga tres ortantes compatibles.

Tenemos pues que para poder inferir de manera única la matriz de interacciones entera se necesitan tanto muestras en equilibrio como muestras dinámicas. Al considerar de manera conjunta la PROPOSICIÓN 2.2 y la PROPOSICIÓN 2.4 se tiene como consecuencia el siguiente resultado.

TEOREMA 2.5 Suponga que un conjunto de datos \mathcal{D} provienen de un sistema con al menos una interacción presente y al menos una interacción no presente. Para que \mathcal{D} tenga una única matriz de interacciones admisible es necesario que el conjunto contenga tanto muestras dinámicas como muestras en equilibrio.

Concluyendo también el ejemplo del ecosistema con dos especies, en la Figura 2.5 se observa que al considerar las muestras en equilibrio como las cuatro muestras en movimiento de la Figura 2.4 se logra inferir de manera única la matriz de interacciones correspondientes del sistema, puesto que los vectores de signos admisibles para cada especie resultan ser

$$S_1 = \{(-1, -1)\}, S_2 = \{(0, -1)\}.$$

2.5 PARÁMETRO DE TOLERANCIA

Por último, es necesario hacer notar que en la práctica el método puede ser muy sensible al ruido y a la precisión de las mediciones. Esto se debe principalmente a que para inferir los patrones admisibles para un agente usando el par de muestras y, z , el resultado de la inferencia

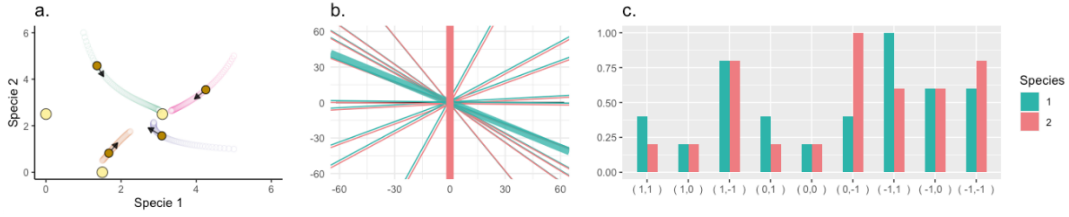


Figura 2.5: Ejemplo en dimensión $N = 2$ de que al considerar una conjunto de datos que contiene tanto muestras en equilibrio como muestras en movimiento. (a) Conjunto de datos de un sistema de dos dimensiones que consiste en los cuatro puntos en movimiento y los tres equilibrios no triviales. (b) Ortantes compatibles para cada par de muestras y cada una de las especies. (c) Conteo de los ortantes compatibles para cada una de las especies. Al observar la frecuencias se observa que si la muestra contiene datos en equilibrio y en movimiento, entonces si es posible alcanzar unicidad 1 para ambas especies.

depende fuertemente del signo de cada entrada del vector $y - z$. Así, el método depende de poder decidir con confiabilidad cuando la actividad de un agente cambia y con que signo es dicho cambio.

En particular, el resultado de la inferencia puede verse afectado cuando dos muestras presentan una actividad muy similar para uno de los agentes y no es posible asegurar que la diferencia sea consecuencia de la dinámica del sistema o simplemente sea debido a las limitaciones de la medición o por una fuente de ruido externo. Para manejar ésta situación proponemos incluir en la construcción de $\delta(y, z)$ el siguiente parámetro de tolerancia.

DEFINICIÓN 2.9 PARÁMETRO ϵ_T DE TOLERANCIA

Sean $y, z \in \mathfrak{D}_i$. Dado un parámetro de tolerancia $\epsilon_T > 0$, definimos la *diferencia con tolerancia* ϵ_T entre los vectores y, z de la siguiente manera:

$$(y - z)_j(\epsilon_T) = \begin{cases} 0 & \text{si } |(y - z)_j| < \epsilon_T, \\ y_j - z_j & \text{en otro caso.} \end{cases}$$

VALIDACIONES NUMÉRICAS

Para realizar una prueba de concepto del método se realizó una implementación en R [23] de los algoritmos presentados en el capítulo anterior. Como se verá a lo largo del presente capítulo, es todo un reto diseñar y realizar una validación exhaustiva. Sin embargo, mediante técnicas heurísticas, es posible realizar una validación inicial que arroja resultados prometedores. En particular, presentamos algunas conclusiones para orientar la interpretación de los resultados cuando los datos medidos contienen ruido.

Recordemos que las bases teóricas del método permiten utilizarlo en una amplia familia de sistemas dinámicos. Por lo cual, para acotar el problema, se eligió como punto de partida probar el método en modelos de ecosistemas microbianos.

Incluso en este contexto, explorar exhaustivamente todas las configuraciones y parámetros necesarios para simular ecosistemas con parámetros aleatorios es una tarea que queda fuera del alcance de esta tesis. Es por ello que se decidió simular datos sintéticos generados provenientes de un mismo ecosistema aleatorio, de tal manera que permita explorar el comportamiento del método en conjuntos de muestras similares a los muestreos experimentales propias de la disciplina. El objetivo principal de este capítulo es observar el desempeño del método al variar algunas condiciones sobre la cantidad, composición y calidad de las muestras disponibles. Asimismo, analizaremos como influye la elección del parámetro de tolerancia ϵ_T , incluyendo el caso cuando se le agregan distintos niveles de ruido a los datos.

Por último, mostramos el desempeño del método en un ecosistema del microbiota humano de doce especies. El sistema considera interacciones tipo Holling II, y los parámetros de la dinámica considerada se tomaron de resultados experimentales obtenidos en [31] por Venturelli *et al.* Se muestra como es posible recuperar parcialmente la matriz de interacciones al modelar la construcción de los conjuntos de datos siguiendo los experimentos realizados en dicho trabajo.

3.1 CASO DE ESTUDIO: ECOSISTEMA ALEATORIO

Muestreo de un ecosistema aleatorio

Se simularon muestras de un ecosistema aleatorio con $n = 5$ especies. Para simular la dinámica de la abundancia de cada especie se consideraron interacciones entre especies modeladas mediante respuestas funcionales tipo Holling II. Es decir, la dinámica de la abundancia de la i -especie esta determinada por la ecuación diferencial

$$\dot{x}_i = x_i \left(-b_i + \sum_{j \neq i}^n a_{ij} \frac{x_j}{1 + \theta_{ij} x_j} \right), \quad i = 1 \cdots n. \quad (3.1)$$

Arriba, los parámetros b_i, a_{ij}, θ_{ij} se determinaron mediante variables aleatorias. El proceso para generar los parámetros del ecosistema fue el siguiente:

1. La red de interacciones $A = (a_{ij})$ se obtuvo con una realización de un modelo de red aleatoria Erdős-Rényi $G(n, 0.5)$. Es decir, todas las interacciones a_{ij} tenían una probabilidad de 0.5 de aparecer en el modelo.
2. A cada una de las interacciones $a_{ij} \neq 0$ se le asignó un peso aleatorio con distribución Normal $\mathcal{N}(0, 0.05)$.
3. Las tasas de crecimiento b_i se asignaron mediante una variable uniforme $\mathcal{U}(1, 4)$.
4. Los valores θ_i de amortiguamiento de la respuesta funcional se asignaron mediante una variable uniforme $\mathcal{U}(0, 0.1)$.
5. Por último, se definió $a_{ii} = -1$ para toda $i \in \{1, 2, \dots, n\}$.

Los parámetros aleatorios resultantes pueden observarse en la Tabla

3.1.

Una vez fijos los parámetros del ecosistema aleatorio se simularon series de tiempo de las poblaciones de las especies para distintas configuraciones posibles. Aquí, nos referimos como *configuración* a la colección de especies presentes al inicio del experimento. Notemos que para un ecosistema de dimensión n , existen un total de $2^n - 1$ configuraciones diferentes. Ahora bien, dada una configuración de especies, se determinó una abundancia inicial para cada especie y se muestrearon puntos *en movimiento* y el *punto en equilibrio* correspondiente de la serie de tiempo obtenida de la siguiente manera:

1. Las muestras *en movimiento* se obtuvieron de la fase inicial de las series de tiempo. Se consideraron puntos para los cuales existiera cambio en al menos una de las especies. Esto es, se consideraron todos los tiempos para los cuales $\max(|\dot{x}_i|) > \epsilon_M$ para algún ϵ_M . Para nuestros resultados utilizamos un valor de $\epsilon_M = 0.01$.

Las configuraciones de los experimentos corresponden a los diferentes cocultivos que se pueden obtener. Es decir, a todas las posibles combinaciones de especies presentes en un experimento.

MATRIZ (a_{ij})				
-1.0000	0.0000	-0.0138	0.0000	0.0452
-0.0732	-1.0000	0.0000	0.0000	-0.0125
0.0000	0.0000	-1.0000	-0.0286	0.0661
0.0294	0.0000	0.0000	-1.0000	-0.0393
0.0127	0.0000	0.0064	-0.0629	-1.0000
VECTOR b_i				
3.4733	2.4902	2.2930	1.3355	3.6014
VALORES θ_{ij}				
0.0201	0.0674	0.0139	0.0545	0.0328

Tabla 3.1: Parámetros aleatorios para las interacciones del ecosistema simulado.

- De manera análoga, se consideró que una serie de tiempo se encontraba *en equilibrio* si el cambio en cada una de las especies presentes era mínimo, es decir si $\max(|\dot{x}_i|) < \epsilon_E$ para algún ϵ_E . En este caso se utilizó $\epsilon_E = 10^{-6}$.

Observemos que la elección de los parámetros de las simulaciones determina únicamente el universo muestral que se utilizará, pero no tienen influencia en el proceso de inferencia.

En la Fig. 3.1 se observa el comportamiento de las series de tiempo para todas las configuraciones diferentes de un ecosistema aleatorio simulado con los parámetros descritos en la Tabla 3.1. Dado que en este caso los parámetros de la simulación son conocidos, se calcularon las derivadas evaluando cada punto en movimiento en las funciones correspondientes. A los puntos en equilibrio se les asignó como derivada el vector cero.

De la misma manera, debido a que en las simulaciones se conoce el signo de las interacciones del sistema del que provienen los datos es posible determinar la *tasa de error* E_U de las entradas unívocamente determinadas como un indicador más del desempeño del método. Notemos que a diferencia de las limitaciones fundamentales, inherentes de un conjunto de datos, la tasa de error puede tener como origen distintos factores, entre ellos la aproximación numérica de las muestras en equilibrio y el parámetro de tolerancia ϵ_T .

El proceso para seleccionar las muestras para los distintos *conjuntos de datos* está inspirado en las posibles condiciones en las que se pueden obtener los datos reales al realizar experimentos en un laboratorio. Para ello, se seleccionan N_C configuraciones del ecosistema, y para cada configuración se obtienen N_M puntos en movimiento y la muestra en equilibrio correspondiente. De tal manera que, una vez elegidos

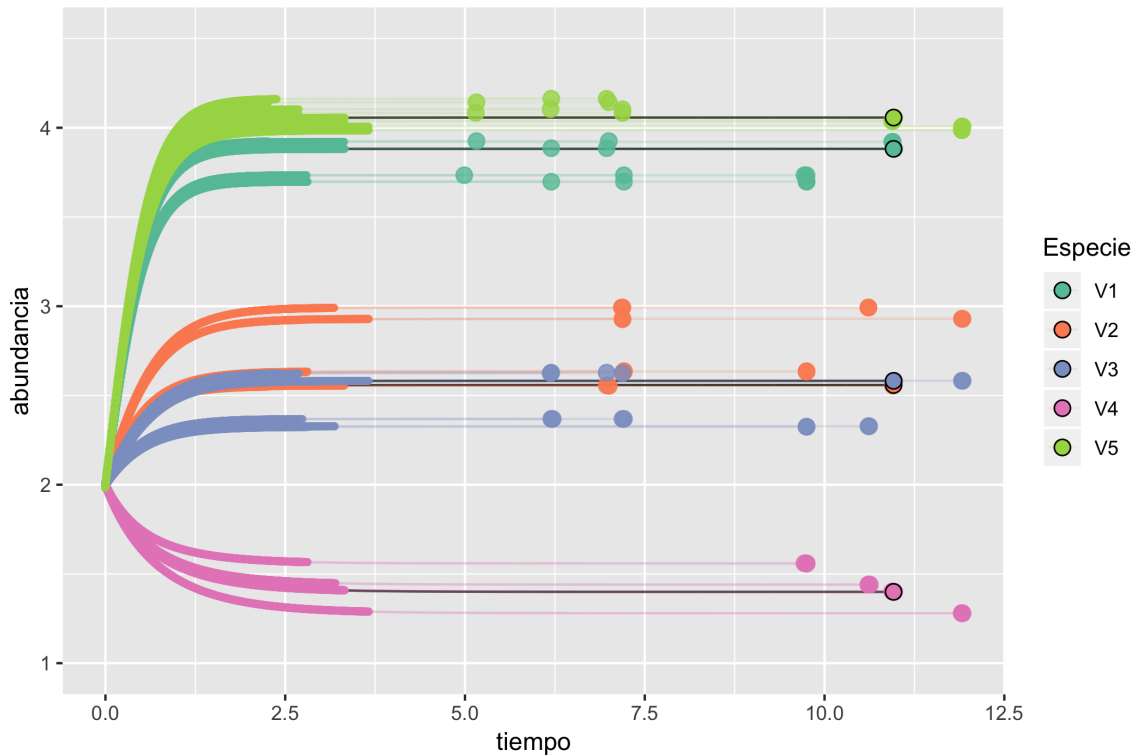


Figura 3.1: Simulación ecosistema aleatorio con cinco especies. Series de tiempo correspondientes a las abundancias de cada especie. El color de las líneas indica la especie. Las muestras en movimiento se tomaron en la primera parte del proceso (línea gruesa), y las *muestras en equilibrio* se tomaron cuando las series cambiaban mínimamente (Puntos a la derecha). Se observa como las diferentes configuraciones de especies determinan los valores de equilibrio a los que puede llegar cada especie. Resaltadas en negro se observan las abundancias correspondientes a la configuración con las cinco especies presentes.

N_C y N_M , el conjunto de datos tiene un total de $N_C(N_M + 1)$ muestras, N_C de las cuales están en equilibrio.

Como se observó en la Sección 2.3, se espera que los indicadores de la inferencia, tanto informatividad como unicidad, sean una función no decreciente con respecto al tamaño de muestra. Uno de los objetivos es explorar las diferentes combinaciones para N_C y N_M y observar los efectos de la composición de los conjunto de datos en el resultado de la inferencia. Por lo tanto, dado que para calcular las redes de interacciones admisibles para un conjunto dado es necesario hacerlo agregando una muestra a la vez, se calcularon primero las redes admisibles para todas las muestras en equilibrio presentes, y, en un segundo paso, se agregaron las muestras en movimiento.

Notemos que, por construcción, el método permite agregar las muestras una a una, y por lo tanto se puede obtener el cambio en la unicidad, la informatividad y la tasa de error conforme se van añadiendo las muestras en movimiento al análisis. En la Fig. 3.2 se puede observar el comportamiento esperado de cada uno de los indicadores al realizar el análisis.

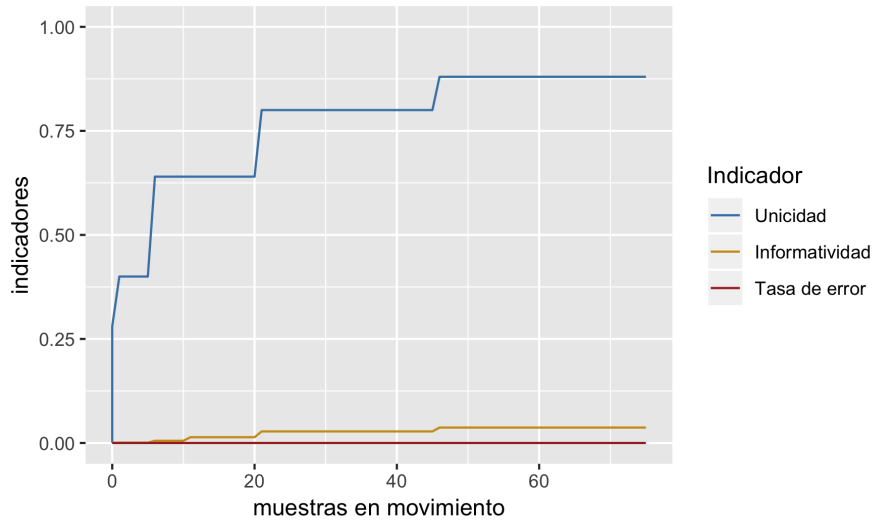


Figura 3.2: *Indicadores del proceso de inferencia.* Gráfica que muestra el cambio en los indicadores al agregar al análisis las muestras en movimiento. Suponiendo como punto de partida que ya se tomaron en cuenta todos las muestras en equilibrio disponibles.

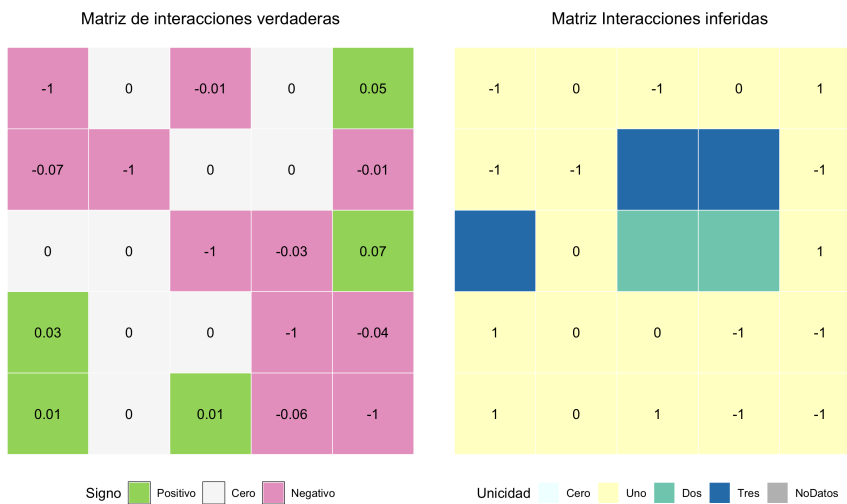


Figura 3.3: *Comparación de la matriz interacciones y la matriz de interacciones inferidas* Comparación de la matriz de parámetros originales (Izquierda) contra la matriz de interacciones inferidas (Derecha). La matriz de interacciones inferidas codifica en el color de cada entrada la cantidad de interacciones admisibles para cada entrada a_{ij} . En las entradas unívocamente determinadas se muestra el correspondiente valor inferido.

Por otro lado, el resultado del proceso de inferencia son los conjuntos S_1, \dots, S_d y se podrían enlistar a partir de ellos todas las redes de interacciones admisibles para los datos. En vez de ello, proponemos un resumen gráfico: la *Matriz de interacciones inferidas*. Esta matriz resume las entradas unívocamente inferidas, incluyendo el valor de las

mismas, así como la multiplicidad de las entradas que no es posible determinar de manera única. Esta representación gráfica, aunque no muestra todas las matrices de interacciones admisibles para los datos, resulta de utilidad puesto que al centrarse en las entradas unívocamente determinadas permite comparar los signos inferidos con la matriz de interacciones original tal como se muestra en la Fig 3.3. Denotaremos por s_{ij} a las interacciones verdaderas del modelo del que provienen los datos.

Parámetro de tolerancia ϵ_T

Como punto de partida, estamos interesados en observar el desempeño del método cuando los datos corresponden exactamente a las predicciones del modelo.

PARÁMETROS DE MUESTREO	VALORES CONSIDERADOS
Número configuraciones	6, 12, 18, 24, 31
Número muestras movimiento	0, 1, 5, 10
Parámetro tolerancia ϵ_T	0, 10^{-10} , 10^{-8} , 10^{-6} , 10^{-4} , 10^{-2}

Tabla 3.2: *Tabla con los valores utilizados para los parámetros de muestreo en la exploración inicial. Se realizaron 20 repeticiones para cada combinación de los parámetros, dando un total de 2400 realizaciones.*

En particular, queremos elegir primero el parámetro de tolerancia ϵ_T que tenga un mejor desempeño en el caso donde los datos no tienen ruido, es decir, corresponden a los valores obtenidos al simular numéricamente el sistema dinámico. En la Tabla 3.3 y en la Fig 3.4 se puede observar la media de cada uno de los indicadores según el valor de ϵ_T utilizado. Buscamos específicamente el valor ϵ_T que maximice la unicidad y minimice el error. Como se observa en la Fig 3.5, $\epsilon_T = 10^{-6}$ es el valor para el parámetro de tolerancia que obtiene un mejor desempeño. Dado que queremos explorar la influencia de cada parámetro de la simulación, se optó por fijar $\epsilon_T = 10^{-6}$ para las siguientes pruebas descriptivas.

Composición de una muestra con puntos estables y en movimiento

Para explorar como influyen en la informatividad de un conjunto de datos la composición del mismo, es decir el número de muestras en movimiento y de muestras en estado estable, se realizó una segunda ronda exploratoria de simulaciones para la cual se fijó el parámetro $\epsilon_T = 10^{-6}$ debido a que fue el valor que permitió a la inferencia tener el mejor desempeño. En particular, se eligió este valor puesto que en todas las simulaciones, la inferencia se obtuvo sin errores.

ÉPSILON	UNICIDAD	INFORMATIVIDAD	TASA DE ERROR
0	0.4444	0.03294712	0.20183654
10^{-10}	0.354	0.02262851	0.21173609
10^{-8}	0.4458	0.02338595	0.11707459
10^{-6}	0.6177	0.30334366	0.00000000
10^{-4}	0.5472	0.06723652	0.02464485
10^{-2}	0.2452	0.04074677	0.15615290

Tabla 3.3: *Análisis Exploratorio para seleccionar un parámetro de tolerancia* Valores promedio de los indicadores del proceso de inferencia para diferentes valores del parámetro de tolerancia ϵ_T .

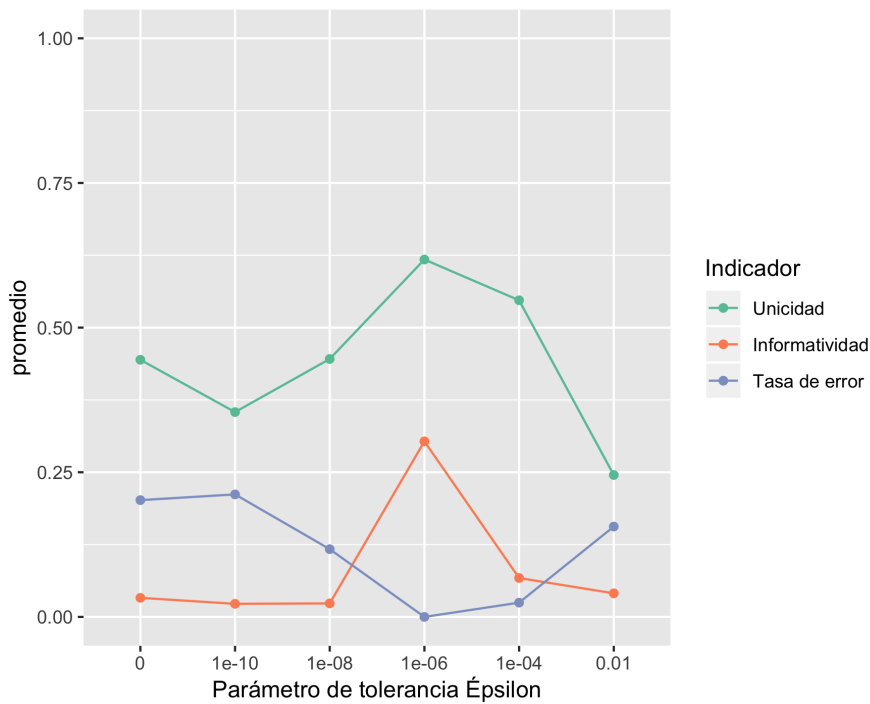


Figura 3.4: *Medias de los indicadores para distintos valores de ϵ_T .* Valores promedio de los indicadores del proceso de inferencia para diferentes valores del parámetro de tolerancia ϵ_T .

Se realizaron 30 repeticiones para cada combinación de los parámetros descritos en la Tabla 3.4. En cada una de ellas la selección de las configuraciones presentes en cada conjunto de datos así como el muestreo de los puntos en movimiento para cada configuración se realizaron de manera aleatoria. Como se observa en la Fig. 3.6, la unicidad alcanza valores mayores cuando el número de configuraciones incluidas es más mayor. En particular, cuando se consideran todas las configuraciones disponibles, $N_C = 31$, se alcanza $U(\mathcal{D}) = 1$ al agregar sólo dos muestras en movimiento.

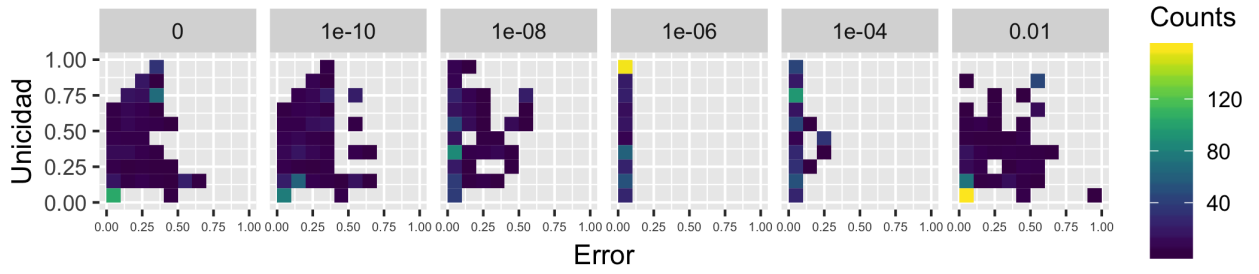


Figura 3.5: *Heatmap Error vs. Unicidad para distintos valores de ϵ_T* . Se muestra el resultado para distintos valores del parámetro de tolerancia ϵ_T . Como se corrobora en la Tabla 3.3, el único valor que obtuvo error cero es $\epsilon_T = 10^{-6}$. Es de destacar que, independientemente del valor ϵ_T , la gran mayoría de las inferencias que tienen unicidad no cero, presentan un error pequeño.

PARÁMETROS DE MUESTREO	VALORES CONSIDERADOS
Número configuraciones	6, 12, 18, 24, 31
Número muestras movimiento	0, 2, 4, 6, 8, 10
Parámetro tolerancia ϵ_T	10^{-6}

Tabla 3.4: *Tabla con los valores de los parámetros de muestreo utilizados para ϵ_T fijo*. Se realizaron 30 repeticiones para cada combinación de los parámetros, dando un total de 900 realizaciones.

Notemos además, que debido a que la matriz de interacciones tiene nueve entradas cero, se tiene que éste es el número máximo de entradas que pueden ser inferidas correctamente mediante conjuntos de muestras en equilibrio. Así, por la [PROPOSICIÓN 2.2](#), se esperaría que si la inferencia se realizó sin errores, entonces la unicidad esté acotada por 0.36 para los conjuntos que no contienen muestras en movimiento.

Ahora bien, como se observa en las Figs. 3.9 y 3.8, el desempeño de los indicadores parece estar más ligada a una composición equilibrada del conjunto de datos (*i.e.*, que se tengan tanto muestras en equilibrio como muestras en movimiento presentes) que al tamaño total de la muestra. Es decir, los resultados parecen indicar que para obtener resultados satisfactorios en la inferencia no es necesario tener una gran cantidad de muestras en movimiento para cada configuración, sino que es preferible considerar una amplia variedad de configuraciones y agregar pocas muestras en movimiento para cada una de ellas, lo cual sería suficiente para inferir de manera única la matriz entera. Notemos también que hay dos casos que parecen estar claramente acotados:

1. Cuando se tienen pocas configuraciones, pareciera que agregar muestras en movimiento no mejora el desempeño de los indicadores. Esto puede deberse a que cuando se tienen pocas configuraciones sólo se tiene acceso a información de un subconjunto de las especies.
2. Cuando no se cuenta con muestras en movimiento, no importa que se cuente con todas las configuraciones posibles, ambos indicadores permanecen acotados. Esto es consecuencia de que cuando sólo se tienen muestras en equilibrio, sólo es posible detectar las entradas cero de la matriz de interacciones.

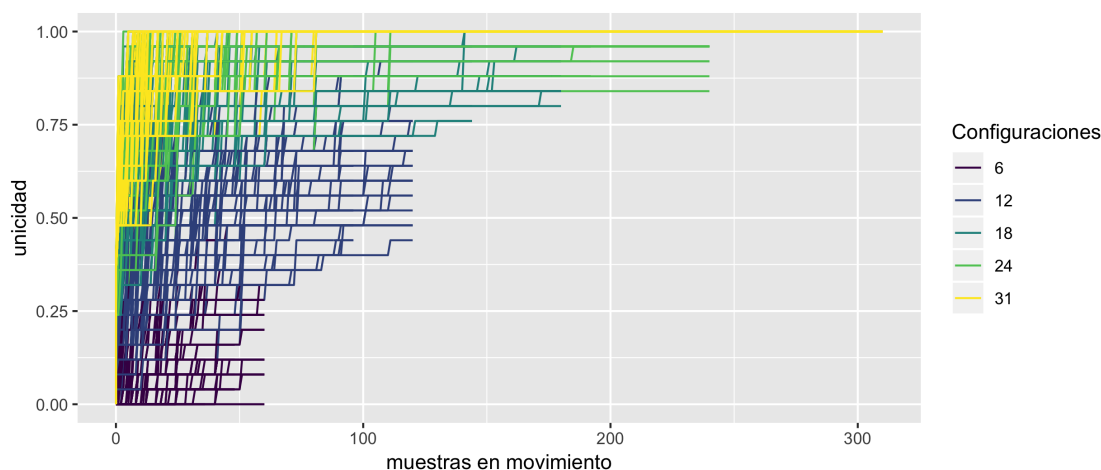


Figura 3.6: *Indicadores del proceso de inferencia.* Se muestra el cambio en los indicadores al agregar muestras en movimiento al análisis. En cada caso se tiene como punto de partida el resultado de analizar las muestras en equilibrio disponibles.

Estas observaciones son congruentes con las conclusiones del **TEOREMA 2.5**: para lograr un buen desempeño en la inferencia se necesita que el conjunto de datos contenga muestras en equilibrio y muestras en movimiento. En resumen, independientemente del número de muestras en movimiento, si se tienen *pocas* muestras en equilibrio, en este caso se corresponden a los casos con 6 o 12 equilibrios, es poco probable detectar la ausencia de interacciones (*i.e.*, *las entradas cero*). Análogamente, se observa que sin importar la cantidad de equilibrios presentes, si hay pocas muestras en movimiento difícilmente se infieren de manera única las interacciones no cero.

OBSERVACIÓN 3.1 IMPORTANCIA DE LA COMPOSICIÓN DE LOS CONJUNTOS DE DATOS Los resultados de la inferencia con simulaciones confirman las limitaciones fundamentales que presentan los conjuntos de datos cuando no contienen *suficientes* muestras en equilibrio, ó en movimiento.

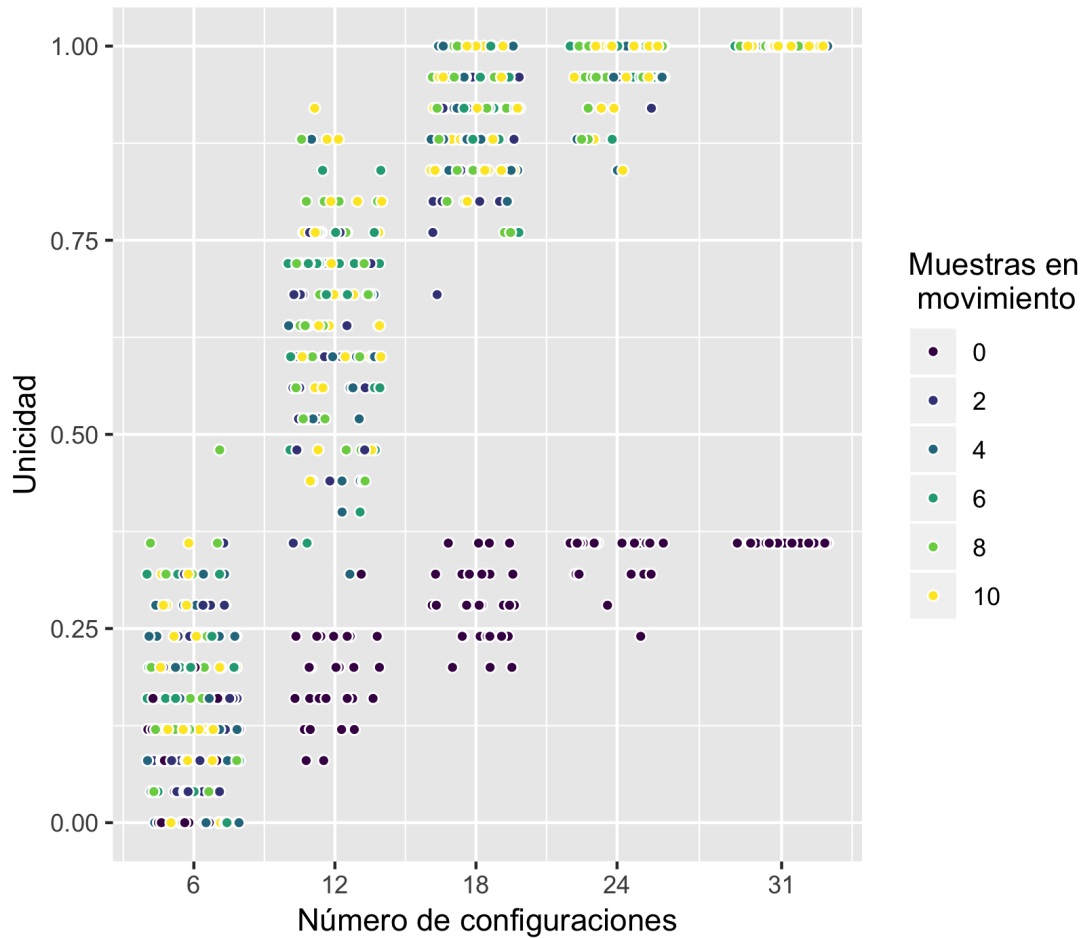


Figura 3.7: *Unicidad vs. Número de configuraciones*. Cada punto corresponde al resultado final de la inferencia para un conjunto de datos aleatorio. Como se observa, la composición del conjunto de datos, indicada en el eje x y en el color, influye fuertemente en la unicidad que se alcanza. Notemos que, cuando no hay muestras en movimiento, la unicidad queda acotada por el número de ceros en la matriz de interacciones, tal como se espera.

3.2 SENSIBILIDAD DEL METODO A MUESTRAS CON RUIDO

Por último, exploraremos el desempeño del método cuando los datos se obtienen con ruido. Es de interés observar si escoger adecuadamente el parámetro de tolerancia ϵ_T es suficiente para distinguir las interacciones en este caso, y bajo que condiciones de ruido esto resulta cierto.

Para ello, fijaremos ahora la *composición* de los conjuntos de datos, considerando tres casos: conjuntos que sabemos tienen poca informatividad (6 configuraciones, 2 muestras dinámicas por configuración), relativamente equilibrados (18 configuraciones, 2 muestras dinámicas por configuración) y conjuntos que esperamos aporten suficiente

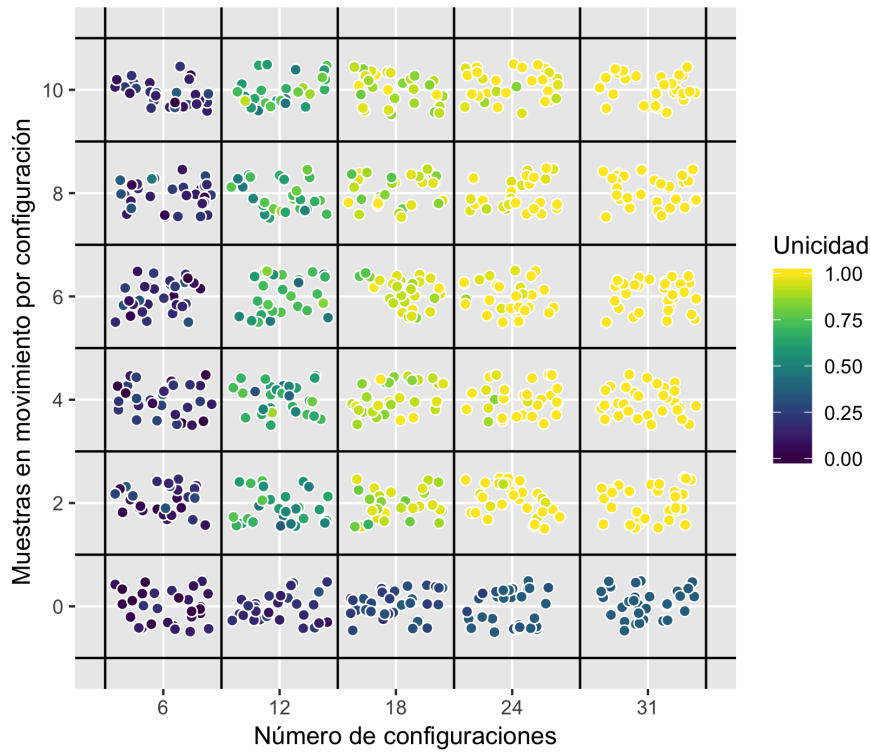


Figura 3.8: *Unicidad en distintas configuraciones de muestras.* Se observa como el valor de unicidad alcanzado está relacionado con el número de configuraciones y el número de muestras en movimiento consideradas. Para pocas configuraciones, la unicidad parece estancarse lejos del 1 independientemente del número de muestras en movimiento que se consideren. En cambio, cuando se consideran la mayoría de las configuraciones, incluso con muy pocas muestras en movimiento se alcanzan unicidades altas.

información (31 configuraciones, 2 muestras dinámicas por configuración). Para simular el ruido, para cada observación se consideró un ruido normal $N(0, \sigma_R)$ donde la abundancia x_i es reemplazada por $\max\{0, x_i + N(0, \sigma_R)\}$, para distintos valores de σ_R . Una vez agregado el ruido, se realizó el proceso de inferencia con distintos valores de ϵ_T . Los valores considerados en las simulaciones se pueden observar en la Tabla 3.5.

Es importante notar que además de las limitaciones fundamentales también existen otros factores que condicionan un resultado óptimo de la inferencia. En particular, es posible que los conjuntos de datos, inclusive datos sintéticos, no satisfagan las hipótesis necesarias y por lo tanto no se obtenga el resultado esperado. En particular, notemos que estamos asumiendo que todas las muestras son *compatibles*, es decir, todas ellas comparten al menos un vector admisible (el de los parámetros originales). Sin embargo, existen diversos motivos que pueden provocar que esta hipótesis no se siga satisfaciendo, entre

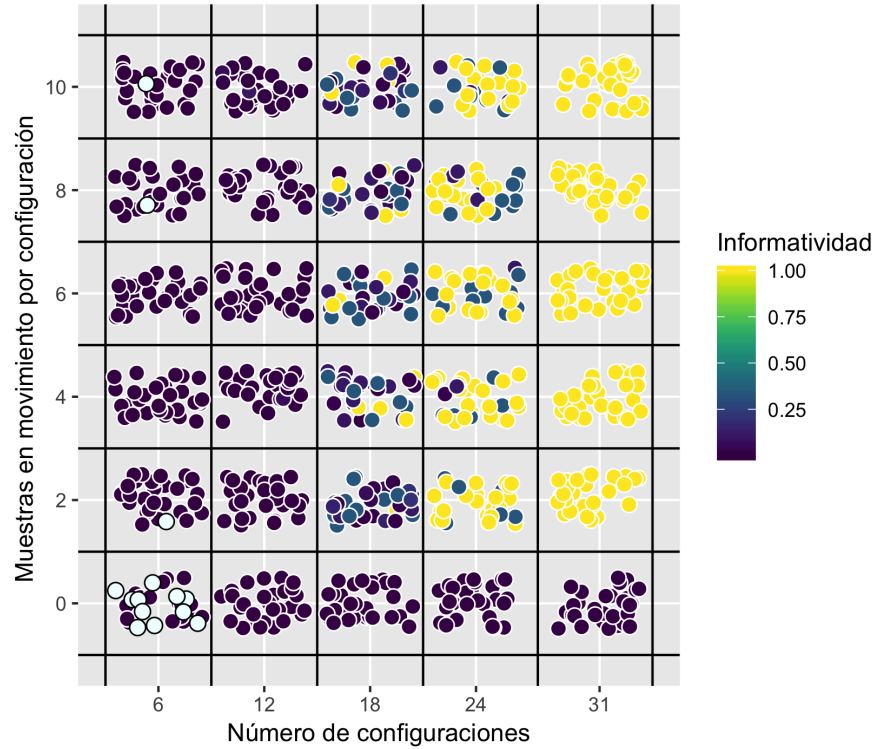


Figura 3.9: *Unicidad en distintas configuraciones de muestras.* Se observa como el valor de unicidad alcanzado está relacionado con el número de configuraciones y el número de muestras en movimiento consideradas. Para pocas configuraciones, la unicidad parece estancarse lejos del 1 independientemente del número de muestras en movimiento que se consideren. En cambio, cuando se consideran la mayoría de las configuraciones, incluso con muy pocas muestras en movimiento se alcanzan unicidades altas.

ellos, el valor ϵ_T seleccionado, las limitaciones numéricas o el ruido en los datos.

Una de las consecuencias de tener un conjunto de datos incompatibles para la especie i es que no se puede obtener un único vector de interacciones admisibles, más aún se obtiene un conjunto vacío. Así, existen sólo dos aspectos en la composición de un conjunto de datos que provocan que para alguna especie i se obtenga como resultado de la inferencia un conjunto vacío:

1. Si $|\mathcal{D}_i| \leq 1$, se sigue que los datos no contienen información sobre la dinámica del i -ésima especie y, por lo tanto, no es posible inferir vectores admisibles para la misma.
2. Si $|\mathcal{D}_i| > 1$, es decir, existe al menos un par de muestras para las contienen a la especie i , entonces para que $S_i(\mathcal{D}_i) = \emptyset$ es necesario que

$$\bigcap_{k=1}^M S_i(p_k) = \emptyset,$$

PARÁMETROS DE MUESTREO	VALORES CONSIDERADOS
Número configuraciones	6, 18, 31
Número muestras movimiento	2
Parámetro tolerancia ϵ_T	$0, 10^{-10}, 10^{-9}, 10^{-8}, 10^{-7}, 10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}$
Parámetro ruido σ_R	$0, 10^{-10}, 10^{-9}, 10^{-8}, 10^{-7}, 10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1$

Tabla 3.5: Tabla con los valores de los parámetros de muestreo utilizados para ϵ_T fijo . Se realizaron 10 repeticiones para cada combinación de los parámetros, dando un total de 7200 realizaciones.

para todas las parejas de elementos p_1, \dots, p_M de \mathcal{D}_i .

Que esta última intersección sea vacía basta que exista una pareja \hat{p} de elementos de \mathcal{D}_i tal que

$$S_i(\hat{p}) \cap \left(\bigcap_{p \neq \hat{p}} S_i(p) \right) = \emptyset.$$

Es decir, para que el conjunto \mathcal{D}_i sea incompatible, basta con que dos de sus elementos lo sean.

Esta observación resulta crucial, puesto que para una interacción a_{ij} sea inferida unívocamente de manera errónea es necesario primero que el vector de interacciones verdadero no sea admisible para todas las parejas de elementos de \mathcal{D}_i . En un segundo término, basta que exista un vector s_2 admisible a todas ellas tal que $s_{2j} = a_{ij}$, ya que en caso de que este segundo vector no exista no se obtiene el error, se obtiene un conjunto incompatible y por lo tanto se descarta el i -ésimo renglón completo.

De esta manera, si se tiene un conjunto \mathcal{D}_i para el cual un subconjunto del mismo $D_i \subset \mathcal{D}_i$ satisface que el único vector de signos admisible para el mismo es el vector de interacciones que dió origen a los datos, es imposible que las entradas del i -ésimo vector sean inferidas de manera única erróneamente. En todo caso, si existe otro subconjunto $D' \subset \mathcal{D}_i$ para el cual $s \notin S(D')$, se tiene que $S(\mathcal{D}_i) = \emptyset$. Luego, si la muestra tiene suficientes muestras compatibles con el vector de interacciones verdadero, al agregar muestras que el ruido vuelve incompatibles con el vector de interacciones original, en vez de obtener errores en la inferencia se descartan renglones completos, lo cual da inferencias con informatividad cero.

Este comportamiento se observa bajo distintas combinaciones de ϵ_T y σ_R , tal como se resume en la Fig 3.10. Notemos que conforme el nivel de ruido aumenta, se necesita un valor de ϵ_T mayor para obtener un mejor resultado de la inferencia. Sin embargo, si el ruido es suficientemente grande, la muestra se vuelve incompatible bajo todos los ϵ_T considerados.

Recordemos que el parámetro ϵ_T regula cuando observaciones distintas son consideradas *idénticas* o cuando es necesario diferenciarlas. Si el valor de ϵ_T es muy bajo se corre el riesgo de considerar como

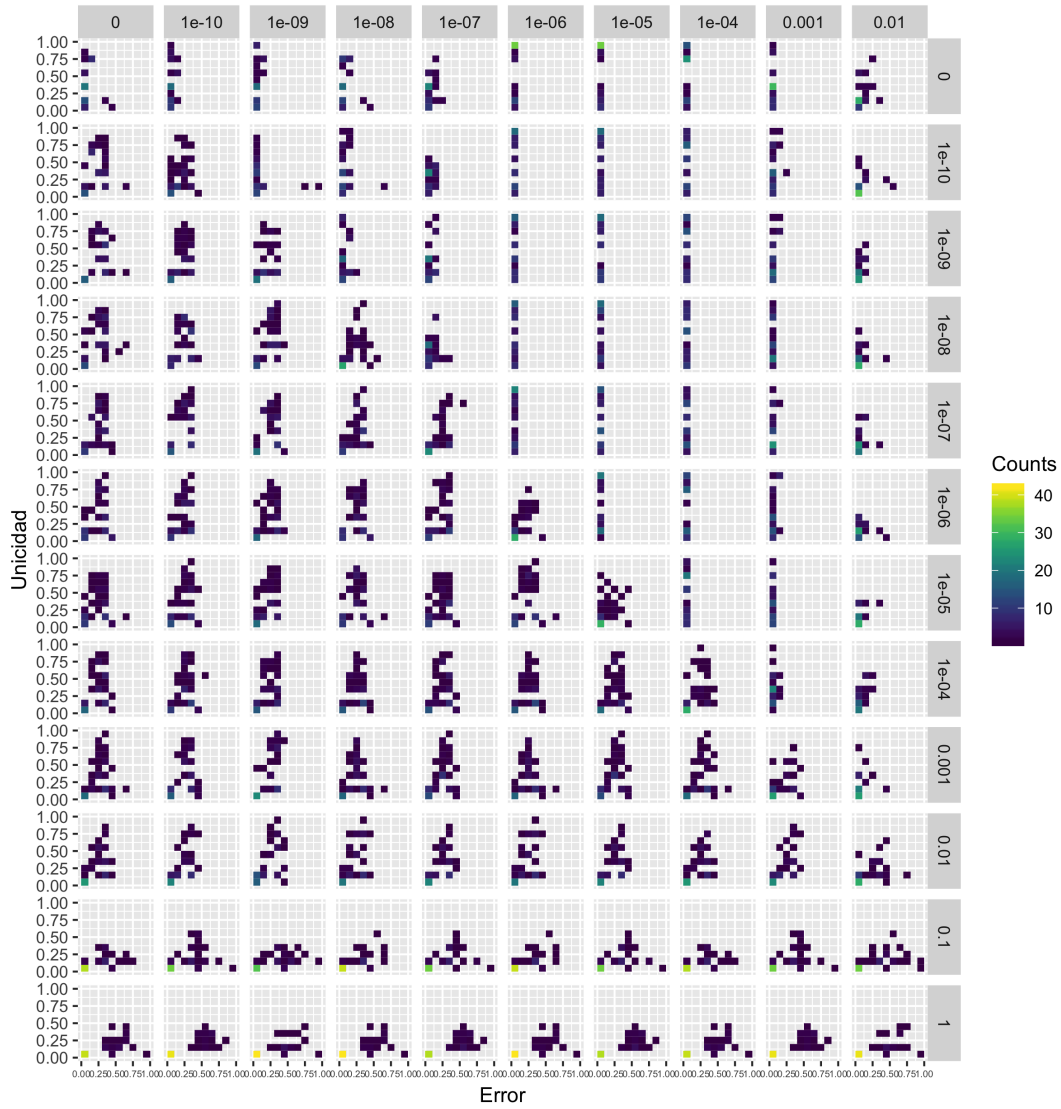


Figura 3.10: Error vs Unicidad para distintos parámetros de tolerancia para muestras con diferentes niveles de ruido. Mapas de calor de los conteos de los resultados de la inferencia. Los renglones corresponden a los distintos parámetros de ruido σ_R , mientras que los renglones corresponden a los parámetros de tolerancia ϵ_T considerados. En cada gráfica, el eje x corresponde a la tasa de error y el eje y a la unicidad obtenidas. Se busca que la inferencia maximice la unicidad minimizando la tasa de error. Por lo tanto, lo ideal es obtener la mayor concentración de resultados lo más cercano a la esquina superior izquierda. Por otro lado, si un conjunto se vuelve incompatible, o no tiene suficiente información, obtendrá unicidades bajas o cero.

diferentes valores que deberían ser el mismo, pero por resolución numérica o por ruido son diferentes. En cambio, si el valor de ϵ_T es muy grande, cabe la posibilidad de que consideramos como iguales observaciones que propiamente no lo son.

3.3 CASO DE ESTUDIO: MICROBIOTA INTESTINAL HUMANO

Para mostrar el desempeño de la implementación en un contexto más realista, se aplicó el método de inferencia en el modelo del microbiota intestinal humano de doce especies obtenido en [31]. En este trabajo, Ophelia S Venturelli *et al.* modelan las interacciones de 12 especies del microbiota intestinal humano mediante un modelo Lotka-Volterra generalizado. Las ecuaciones que modelan las abundancias de las especies están dadas por:

$$\dot{x}_i = x_i \left(-b_i + \sum_{j \neq i}^{12} a_{ij} x_j \right). \quad (3.2)$$

Los parámetros del modelo que proponen en dicho trabajo se pueden observar en la Tabla 3.6.

MATRIZ (a_{ij})												
	BH	CA	BU	PC	BO	BV	BT	EL	FP	CH	DP	ER
BH	-0.91	0.45	0.00	0.00	0.00	0.14	0.00	0.69	0.96	0.00	0.00	1.34
CA	-0.31	-0.83	0.00	-0.56	0.00	-0.66	0.00	-1.10	0.00	-0.24	0.04	0.00
BU	-0.23	-0.26	-0.88	-0.32	-0.63	-0.58	-0.75	-0.12	0.23	-0.15	-0.18	-0.06
PC	-0.53	-0.67	0.00	-0.62	0.00	0.00	0.00	-1.08	-0.40	-0.77	-0.43	0.00
BO	-0.21	-0.28	-0.92	-0.27	-0.73	-0.56	-0.82	-0.11	-0.10	-0.47	-0.20	-0.02
BV	-0.13	-0.17	-0.55	-0.20	-0.52	-0.66	-0.76	-0.05	0.76	0.04	-0.03	-0.03
BT	-0.27	-0.27	-0.82	-0.30	-0.62	-0.64	-0.91	-0.17	-0.07	-0.63	-0.20	0.00
EL	0.18	-0.45	3.38	-0.90	1.76	1.30	2.27	-2.44	-0.77	0.02	0.18	-0.14
FP	-0.23	-1.12	-0.78	-0.41	-0.21	-0.64	-0.70	-0.15	-1.04	-0.51	0.00	-0.17
CH	-0.35	0.31	0.07	0.27	-0.51	-0.05	-0.09	0.00	0.45	-1.45	-0.15	1.08
DP	-0.90	0.00	0.00	-0.98	0.00	-0.11	0.00	-0.41	1.01	-2.16	-1.25	0.00
ER	-0.55	0.00	0.00	-0.82	0.00	0.00	-0.74	0.00	0.00	-0.44	0.00	-1.27
VECTOR b												
	0.245	0.246	0.584	0.237	0.478	0.457	0.598	0.402	0.219	0.502	0.232	0.156

Tabla 3.6: Parámetros para el modelo Lotka-Volterra generalizado para doce especies del microbiota intestinal humano. Las doce especies consideradas son: *Prevotella copri* (PC), *Bacteroides vulgatus* (BV), *Bacteroides uniformis* (BU), *Bacteroides ovatus* (BO), *Bacteroides thetaiotaomicron* (BT), *Faecalibacterium prausnitzii* (FP), *Blautia hydrogenotrophica* (BH), *Eubacterium rectale* (ER), *Collinsella aerofaciens* (CA), *Eggerthella lenta* (EL), *Desulfovibrio piger* (DP) y *Clostridium hiranonis* (CH).

Para observar el comportamiento de la inferencia en este contexto, se simularon las trayectorias correspondientes al modelo propuesto.

Asimismo, los muestreos aleatorios se llevaron a cabo de acuerdo con los experimentos que se realizaron en este trabajo y los cuales sirvieron de base para inferir el modelo. En particular, consideramos la inferencia en el conjunto de datos denominado "T4". Este conjunto de datos consiste de muestras obtenidas de las siguientes configuraciones:

1. Los 12 experimentos monoespecie.
2. Los 66 experimentos posibles emparejando de dos en dos a las distintas especies.
3. El experimento con las 12 especies presentes.

De esta manera se tiene un total de 79 configuraciones presentes. Para cada experimento, en [31] tomaron 6 muestras, cada 12 horas, a la sexta muestra se consideró que el experimento ya había alcanzado un equilibrio.

PARÁMETROS DE MUESTREO	VALORES CONSIDERADOS
Número configuraciones	79 (fijas)
Número muestras movimiento	0, 1, 2, 3, 5
Parámetro tolerancia ϵ_T	0, 10^{-10} , 10^{-8} , 10^{-6} , 10^{-4} , 10^{-2}

Tabla 3.7: *Tabla con los valores utilizados para los parámetros de muestreo. Se realizaron 20 repeticiones para cada combinación de los parámetros, dando un total de 600 realizaciones.*

Como primer paso, y de manera análoga a la sección anterior, exploramos el rendimiento de la inferencia para distintos parámetros de tolerancia ϵ_T . Los parámetros considerados y los detalles de las simulaciones pueden observarse en la Tabla 3.7. El comportamiento de las medias de cada indicador se puede observar en la Tabla 3.8 y en la Figura 3.11. Asimismo, en la Figura 3.13 se observa como, de los valores considerados, el ϵ_T con mejor desempeño, tanto en unicidad como en la tasa del error, es $\epsilon_T = 10^{-4}$.

Una vez fijado el valor $\epsilon_T = 10^{-4}$, se procedió a explorar el comportamiento de la inferencia al considerar las mismas configuraciones, pero incluyendo distinto número de muestras en movimiento. Confirmando el comportamiento esperado, al considerar sólo muestras en equilibrio no es posible inferir de manera única las entradas no cero de la matriz de interacciones, ver Fig. 3.13. Asimismo, notemos que en la mayoría de los casos basta agregar una muestra en movimiento por experimento para que se logre inferir prácticamente la matriz entera.

Otra observación importante es que, a diferencia del ecosistema aleatorio de dimensión cinco donde utilizamos todas las configuraciones posibles de especies, los resultados son muy alentadores al considerar que se utilizan sólo 79 de las 4095 configuraciones posibles. Un ejemplo del resultado de la inferencia para un conjunto de datos

ÉPSILON	UNICIDAD	INFORMATIVIDAD	TASA DE ERROR
0	0.80000000	0.80000038	0.17222222
10^{-10}	0.80138889	0.80000038	0.16666667
10^{-8}	0.83055556	0.80000038	0.05000000
10^{-6}	0.79506944	0.19333346	0.06400080
10^{-4}	0.78333333	0.00299760	0.00007692
10^{-2}	0.33104167	0.00000000	0.16788924

Tabla 3.8: *Análisis Exploratorio para seleccionar un parámetro de tolerancia.* Valores promedio de los indicadores del proceso de inferencia para diferentes valores del parámetro de tolerancia ϵ_T .

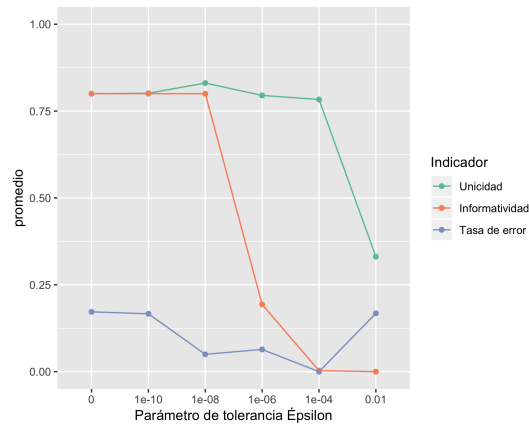


Figura 3.11: *Medias de los indicadores para distintos valores de ϵ_T .* Valores promedio de los indicadores del proceso de inferencia para diferentes valores del parámetro de tolerancia ϵ_T . Observemos que $\epsilon_T = 10^{-4}$ es el único valor que no presenta ningún error, asimismo, presenta una informatividad más baja que los ϵ_T menores. Esto sugiere que el alto valor de unicidad que alcanzan con los valores de ϵ_T menores a él se debe a la presencia de errores.

PARÁMETROS DE MUESTREO	VALORES CONSIDERADOS
Número de configuraciones	79 (fijas)
Número muestras movimiento	0, 1, 2, 3, 5
Parámetro tolerancia ϵ_T	10^{-4}

Tabla 3.9: *Tabla con los valores de los parámetros de muestreo utilizados para ϵ_T fijo.* Se realizaron 20 repeticiones para cada combinación de los parámetros, dando un total de 100 realizaciones.

que incluye una muestra en movimiento, y el equilibrio, para cada configuración, que logra una unicidad de 0.9444 se puede observar en la Fig.3.15.

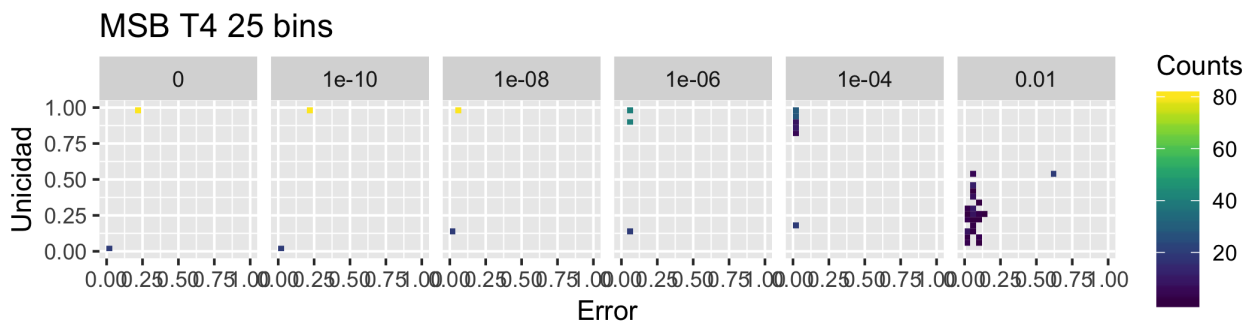


Figura 3.12: Heatmap Error vs. Unicidad para distintos valores de ϵ_T .

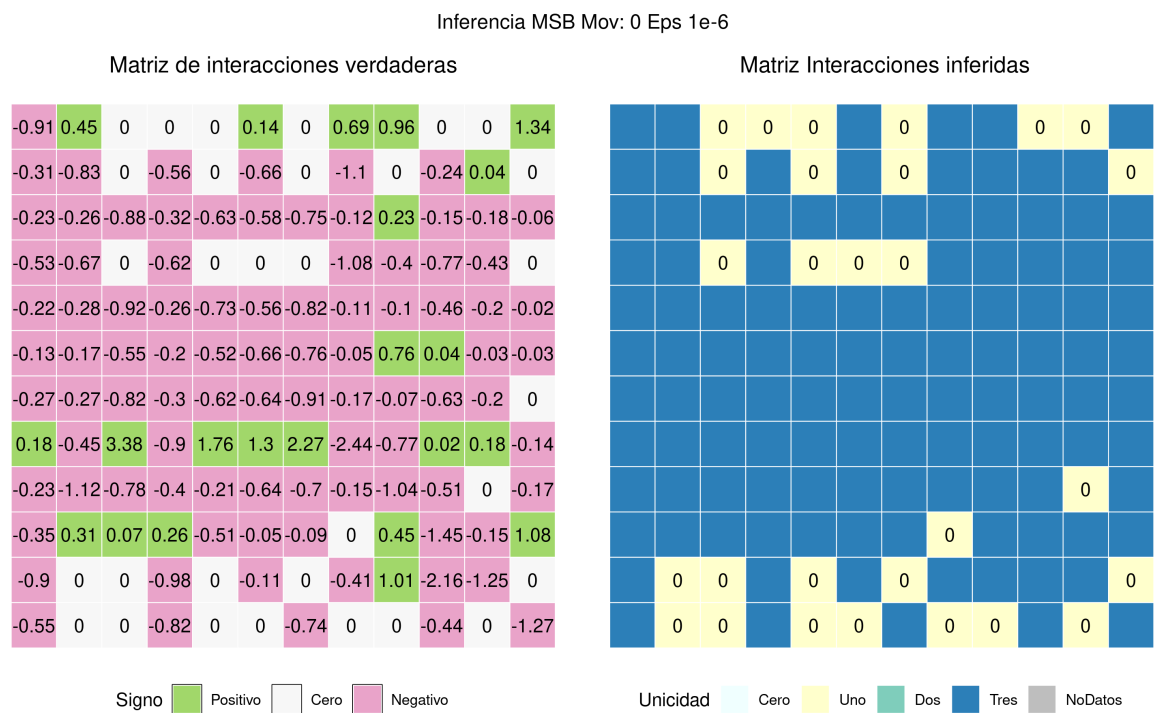


Figura 3.13: Matriz de interacciones inferidas. Resultado de la inferencia de un conjunto de datos con 79 muestras en movimiento y $\epsilon_T = 10^{-4}$. Este conjunto de datos muestra un comportamiento típico, alcanzó una unicidad de 0.1944 y cero errores.

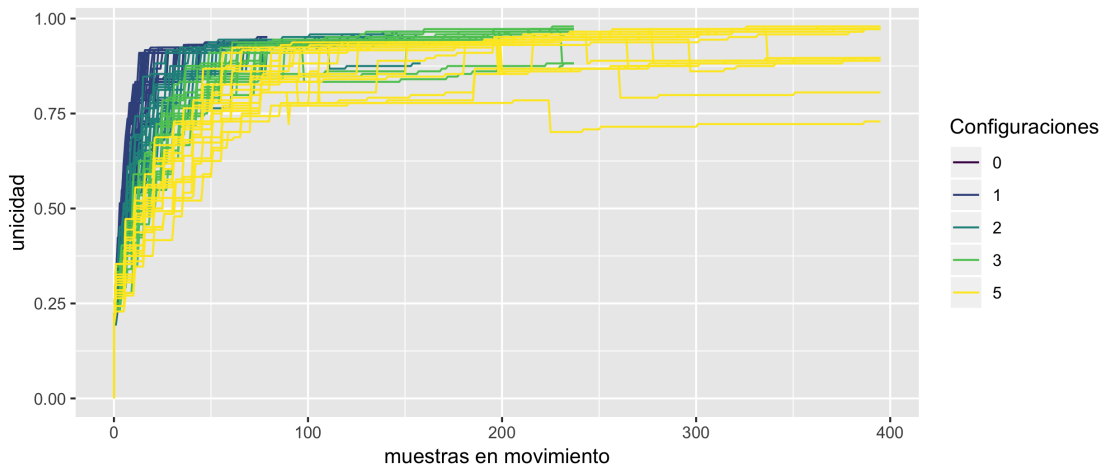


Figura 3.14: *Indicadores del proceso de inferencia.* Se muestra el cambio en los indicadores al agregar muestras en movimiento al análisis. En cada caso se tiene como punto de partida el resultado de analizar las muestras en equilibrio disponibles.

Inferencia MSB con 1 muestras en movimiento y Eps 1e-4

Matriz de interacciones verdaderas

-0.91	0.45	0	0	0	0.14	0	0.69	0.96	0	0	1.34
-0.31	-0.83	0	-0.56	0	-0.66	0	-1.1	0	-0.24	0.04	0
-0.23	-0.26	-0.88	-0.32	-0.63	-0.58	-0.75	-0.12	0.23	-0.15	-0.18	-0.06
-0.53	-0.67	0	-0.62	0	0	0	-1.08	-0.4	-0.77	-0.43	0
-0.22	-0.28	-0.92	-0.26	-0.73	-0.56	-0.82	-0.11	-0.1	-0.46	-0.2	-0.02
-0.13	-0.17	-0.55	-0.2	-0.52	-0.66	-0.76	-0.05	0.76	0.04	-0.03	-0.03
-0.27	-0.27	-0.82	-0.3	-0.62	-0.64	-0.91	-0.17	-0.07	-0.63	-0.2	0
0.18	-0.45	3.38	-0.9	1.76	1.3	2.27	-2.44	-0.77	0.02	0.18	-0.14
-0.23	-1.12	-0.78	-0.4	-0.21	-0.64	-0.7	-0.15	-1.04	-0.51	0	-0.17
-0.35	0.31	0.07	0.26	-0.51	-0.05	-0.09	0	0.45	-1.45	-0.15	1.08
-0.9	0	0	-0.98	0	-0.11	0	-0.41	1.01	-2.16	-1.25	0
-0.55	0	0	-0.82	0	0	-0.74	0	0	-0.44	0	-1.27

Matriz Interacciones inferidas

-1	1	0	0	0	1	0	1	1	0	0	1
-1	-1	0	-1	0	-1	0	-1	1	-1	1	0
-1	-1	-1	-1	-1	-1	-1	-1	1	-1	-1	-1
-1	-1	0	-1	0	0	0	-1	-1	-1	1	1
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	1	-1	-1	1
1	-1	1	-1	1	1	1	-1	-1	1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	0	-1
-1	1	1	1	-1	-1	-1	0	1	-1	1	1
-1	0	0	-1	0	-1	0	-1	1	-1	-1	0
-1	0	0	-1	0	0	-1	0	0	-1	0	-1

Signo ■ Positivo ■ Cero ■ Negativo

Unicidad ■ Cero ■ Uno ■ Dos ■ Tres ■ NoDatos

Figura 3.15: *Matriz de interacciones inferidas.* Resultado de la inferencia de un conjunto de datos con 158 muestras, 79 en equilibrio y 79 en movimiento. Este conjunto de datos muestra un comportamiento típico, alcanzó una unicidad de 0.94444 y cero errores.

DISCUSIÓN

En esta tesis, presentamos, por primera vez, los resultados teóricos que sirven de base para poder inferir de manera completa y libre de modelo todas las interacciones causales de un sistema dinámico a partir de observaciones del mismo.

Este método presenta un enfoque innovador y es, hasta donde tenemos conocimiento, el primer proceso de inferencia para el conjunto de *matrices de interacciones causales* con bases teóricas claras que delimitan el alcance del mismo. La inferencia de interacciones causales es esencial para entender, predecir y controlar una gran variedad de sistemas de diferentes disciplinas. Es así que al presentar las bases para dicha inferencia se da un avance significativo en el problema de reconstrucción de redes. Asimismo, se da un primer paso definitivo para entender diversos sistemas en una amplia variedad de disciplinas.

4.1 COMENTARIOS FINALES

La principal desventaja de los métodos existentes de reconstrucción de interacciones causales existentes es que el resultado de la inferencia es una *única* matriz de interacciones. Se conoce que hay conjuntos de datos que pueden ser igualmente explicados por más de una matriz de interacciones, es por ello que es necesario un enfoque *libre de modelos y completo*. Así, por primera vez, es posible tener total confianza en que las interacciones inferidas son las únicas admisibles *efectivamente* para los datos, ya que se están infiriendo las matrices de interacciones de todos los modelos que son admisibles a los datos.

Por otro lado, recordemos que las condiciones que se les piden a los modelos candidatos son mínimas y la capacidad de inferencia depende fuertemente de como sean los datos a los que quiere ajustarse el modelo. Por lo cual, es natural que el método resulte altamente sensible a los datos con los que se cuenta. Es por ello que el método también presenta algunas desventajas necesarias de mencionar. Resultan de interés los resultados teóricos que describen las limitaciones fundamentales inherentes a los conjuntos de datos, ya que la informatividad propia de un conjunto de datos determina el alcance de la inferencia. Así mismo, es necesario considerar la sensibilidad del método incluso a pequeños errores en los datos, ya que al tratar con datos experimentales hay diversos factores, desde la precisión de las mediciones hasta la presencia de ruido en el proceso, que podrían afectar significativamente el resultado de la inferencia.

Para manejar ésta última observación proponemos un *parámetro de tolerancia* que hace al método más robusto ante pequeñas variaciones en los datos y mediante las validaciones numéricas presentamos ejemplos de que es posible ajustarlo para detectar correctamente las interacciones causales.

Por otra parte, queda un largo camino para refinar y adecuar la implementación y las sutilezas de los procesos de inferencia en casos particulares y datos experimentales. A continuación, concluimos el capítulo con una serie de observaciones sobre distintas líneas de investigación que surgen a raíz de los resultados presentados en esta tesis.

1. ELECCIÓN DEL PARÁMETRO DE TOLERANCIA

Una de las adecuaciones inmediatas que se pueden hacer al método es considerar parámetros de tolerancia específicos para cada agente del sistema. Esto podría ser necesario cuando las series de tiempo de la actividad de distintos agentes presenten comportamientos en escalas diferentes. Un ejemplo inmediato son los ecosistemas microbianos donde los equilibrios de las especies pueden diferir en ordenes de magnitud. Por lo tanto, no sería apropiado utilizar un mismo parámetro de tolerancia para series de tiempo que cuantitativamente son muy distintas.

Por otra parte, para aplicar el método con un conjunto de datos para el cual no se conoce la matriz de interacciones asociado al sistema, es necesario elegir apropiadamente el parámetro de tolerancia, o un intervalo de los mismos, para los cuales se obtiene la inferencia óptima.

Las pruebas numéricas que se exploraron presentan un comportamiento similar, en niveles razonables de ruido, en el cual existe un valor de este parámetro ϵ_{T_0} para el cual los indicadores de desempeño son claramente superiores al resto de los valores probados para el parámetro de tolerancia. Más aún, pareciera existir un patrón en el cual, si el parámetro de tolerancia seleccionado es mucho menor que ϵ_{T_0} , es posible lograr una reconstrucción parcial con tasa de error baja, pero no es posible alcanzar niveles altos de unicidad. Asimismo, dicho patrón de comportamiento sugiere también que si el valor del parámetro de tolerancia es suficientemente más grande que ϵ_{T_0} , entonces se obtienen conjuntos de datos incompatibles y por lo tanto la inferencia se anula.

Un siguiente paso sería describir un proceso sistemático para la elección apropiada del parámetro de tolerancia para un conjunto de datos particulares que haga uso de estas observaciones y con resultados teóricos que lo sustenten.

2. PROCESO DE AUTOCORRECCIÓN DE ERRORES

Una de las conclusiones más interesantes de las validaciones numéricas se centra en observar que no fue posible obtener simultáneamente una unicidad alta con una tasa de errores alta. Esto es, cuando la tasa de error era significativa (mayor a 0.3, por ejemplo) rara vez se obtenían unicidades distintas de cero. Y en el caso de que así fuera, la unicidad no solía sobrepasar el valor de 0.5. En cambio, la mayoría de las veces que se presentaba un valor alto de unicidad venía acompañada de una baja tasa de error.

Este comportamiento, aunado a las observaciones que se hicieron en el Capítulo 3, parecían sugerir que cuando se tiene un conjunto de datos de *tamaño suficiente*, antes que inferir una interacción de manera errónea es más factible obtener un conjunto incompatible y por lo tanto descartar un renglón entero de la matriz de interacciones.

Sería interesante explorar estrategias que permitan diferenciar el comportamiento de las entradas correctamente inferidas de las que no y poder describir con mayor detalle el comportamiento de los errores. Por ejemplo ¿Qué entradas no cambian al considerar diferentes valores de ϵ_T ? ¿Cómo se comporta el método si se consideran bootstrap de las muestras?

3. ERROR POR INCOMPATIBILIDAD

En esta tesis se consideró únicamente como error cuando una entrada era unívocamente inferida de manera incorrecta. Sin embargo, cuando los datos cumplen todas las hipótesis, es de esperarse que tengan al menos un vector de signos admisible para cada uno de los agentes. A pesar de esto, se presentó el caso en las validaciones numéricas donde los datos eran incompatibles para uno o más agentes y de esta manera imposibilitaban inferir la matriz de interacciones completa.

Es claro que si *existe*, como se presupone, una matriz de interacciones para el conjunto de datos completo entonces ésta situación es un error. Es interesante cuantificar las situaciones que dan como resultado un error por incompatibilidad.

Además de explorar los resultados teóricos que describan condiciones necesarios y suficientes para evitar este error (o para que ocurra) hay otras opciones que vale la pena explorar. Por ejemplo, fijando el conjunto de datos asociado a un agente en particular, son de interés las siguientes preguntas:

- a) ¿Es posible obtener una reconstrucción parcial de un renglón cuando existe una entrada que hace al conjunto de datos incompatible? Además de descartar el renglón completo ¿que alternativas podrían considerarse para tratar de obtener más información del sistema?

- b) ¿Se puede cuantificar cuantas entradas hacen incompatible al conjunto de datos? Por ejemplo, no es lo mismo si existe una única entrada que no comparten todos los vectores a que más de la mitad de las entradas sean incompatibles.
- c) ¿Es posible detectar que los conjuntos incompatibles dentro de la muestra? Dicho de otra manera, ¿Es posible medir cuando difieren los conjuntos incompatibles? Sería de interés saber, por ejemplo, si existe una única muestra la que no es compatible con el resto de los datos o si en cambio existen dos, o más, subconjuntos de muestras compatibles entre ellas, pero incompatibles con las muestras de los demás subconjuntos.

4. CONJUNTOS INCOMPATIBLES Y DINÁMICAS INCOMPATIBLES

Continuando en el mismo sentido que la última pregunta del punto anterior. Si no se tiene la certeza de que exista un modelo único que explique todos los datos que se tienen y puede darse el caso de que tengamos muestras que provienen de sistemas con matrices de interacciones causales diferentes, esperaríamos ver que la inferencia no logre obtener entradas únicas sino que más bien el proceso *detecte* que el conjunto de datos es incompatible y se obtenga como resultado que no existen vectores de signos admisibles para los agentes que tienen dinámica diferente.

¿Es posible adaptar el método para diferenciar los conjuntos incompatibles que provienen de sistemas con matrices de interacciones diferentes? ¿Es posible utilizar este método para probar la hipótesis de que todas las muestras provienen de un mismo sistema? ¿Es posible utilizar el método para probar la hipótesis de que no efectivamente las muestras provienen de sistemas diferentes?

5. ESTRATEGIAS PARA DETERMINAR DIFERENCIAS EN LAS OBSERVACIONES

Uno de los puntos sutiles que pueden determinar si se infiere correctamente la matriz de interacciones es el proceso de decisión sobre si dos equilibrios son iguales o si su diferencia es significativa. Una de las estrategias para manejar esta situación fue la implementación del parámetro de tolerancia ϵ_T . Sin embargo, desde un enfoque experimental, puede haber más información que es necesario tomar en cuenta. Para determinar si el equilibrio que alcanza la actividad de un agente se ve afectado por la presencia de otro agente, si se tiene una sola observación de cada caso se debe proceder con una resta, módulo ϵ_T si se quiere. Pero si se cuenta con repeticiones múltiples de los experimentos, para decidir si una diferencia es significativa puede recurrirse a una prueba de hipótesis sobre diferencia de medias. El resultado

de la misma sirve para definir un análogo del vector diferencia que se necesita para inferir los vectores de signos admisibles.

¿Qué otras consideraciones pueden tomarse en cuenta para *pre-procesar* los datos al proceso de inferencia? ¿Qué otros aspectos del proceso se pueden adaptar para utilizar la maquinaria estadística propia de casos particulares?

BIBLIOGRAFÍA

- [1] Marco Tulio Angulo, Jaime A Moreno, Gabor Lippner, Albert-László Barabási y Yang-Yu Liu. «Fundamental limitations of network reconstruction from temporal data». En: *Journal of the Royal Society Interface* 14.127 (2017), pág. 20160966.
- [2] Baruch Barzel y Albert-László Barabási. «Universality in network dynamics». En: *Nature physics* 9.10 (2013), pág. 673.
- [3] Danielle S Bassett y Olaf Sporns. «Network neuroscience». En: *Nature neuroscience* 20.3 (2017), pág. 353.
- [4] Hong-Tai Cao, Travis E Gibson, Amir Bashan y Yang-Yu Liu. «Inferring human microbial dynamics from temporal metagenomics data: Pitfalls and lessons». En: *BioEssays* 39.2 (2017).
- [5] Alex Carr, Christian Diener, Nitin S Baliga y Sean M Gibbons. «Use and abuse of correlation analyses in microbial ecology». En: *The ISME journal* (2019), pág. 1.
- [6] Jose Casadiego y Marc Timme. «Network dynamics as an inverse problem». En: *Mathematical technology of networks*. Springer, 2015, págs. 39-48.
- [7] Jose Casadiego, Mor Nitzan, Sarah Hallerberg y Marc Timme. «Model-free inference of direct network interactions from nonlinear collective dynamics». En: *Nature communications* 8.1 (2017), pág. 2192.
- [8] Gheorghe Craciun y Casian Pantea. «Identifiability of chemical reaction networks». En: *Journal of Mathematical Chemistry* 44.1 (2008), págs. 244-259.
- [9] Riet De Smet y Kathleen Marchal. «Advantages and limitations of current network inference methods». En: *Nature Reviews Microbiology* 8.10 (2010), pág. 717.
- [10] Jonathan Friedman y Eric J Alm. «Inferring correlation networks from genomic survey data». En: *PLoS computational biology* 8.9 (2012), e1002687.
- [11] Timothy S Gardner, Diego Di Bernardo, David Lorenz y James J Collins. «Inferring genetic networks and identifying compound mode of action via expression profiling». En: *Science* 301.5629 (2003), págs. 102-105.
- [12] Ryota Kobayashi y Renaud Lambiotte. «Tideh: Time-dependent hawkes process for predicting retweet dynamics». En: *Tenth International AAAI Conference on Web and Social Media*. 2016.

- [13] Arun Krishnan, Alessandro Giuliani y Masaru Tomita. «Indeterminacy of reverse engineering of gene regulatory networks: the curse of gene elasticity». En: *PLoS One* 2.6 (2007), e562.
- [14] Serge Lang. *Undergraduate analysis*. Springer Science & Business Media, 2013.
- [15] B Lünsmann. «Reconstruction of physical interactions in stationary stochastic network dynamics». Tesis doct. thesis, University of Göttingen, Germany, 2015.
- [16] Kumar Mainali, Sharon Bewick, Briana Vecchio-Pagan, David Karig y William F Fagan. «Detecting interaction networks in the human microbiome with conditional Granger causality». En: *PLoS computational biology* 15.5 (2019), e1007037.
- [17] Valeri A Makarov, Fivos Panetsos y Oscar de Feo. «A method for determining neural connectivity and inferring the underlying network dynamics using extracellular spike recordings». En: *Journal of Neuroscience Methods* 144.2 (2005), págs. 265-279.
- [18] Mor Nitzan, Jose Casadiego y Marc Timme. «Revealing physical interaction networks from statistics of collective dynamics». En: *Science advances* 3.2 (2017), e1600396.
- [19] Charlie Pauvert, Jessica Vallance, Laurent Delière, Marc Buée y Corinne Vacher. «Microbial networks inferred from metabarcoding data lack replicability: consequences for next-generation biomonitoring». En: *bioRxiv* (2019), pág. 642199.
- [20] Judea Pearl. *Causality: models, reasoning and inference*. Vol. 29. Springer, 2000.
- [21] Judea Pearl y Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic Books, 2018.
- [22] Robert J Prill, Daniel Marbach, Julio Saez-Rodriguez, Peter K Sorger, Leonidas G Alexopoulos, Xiaowei Xue, Neil D Clarke, Gregoire Altan-Bonnet y Gustavo Stolovitzky. «Towards a rigorous assessment of systems biology models: the DREAM3 challenges». En: *PloS one* 5.2 (2010), e9202.
- [23] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria, 2019. URL: <https://www.R-project.org/>.
- [24] Michael Rosenblum y col. «Reconstructing networks of pulse-coupled oscillators from spike trains». En: *Physical Review E* 96.1 (2017), pág. 012209.
- [25] Lisa Röttjers y Karoline Faust. «From hairballs to hypotheses—biological insights from microbial networks». En: *FEMS microbiology reviews* 42.6 (2018), págs. 761-780.

- [26] Armindo Salvador. *Uri Alon, An Introduction to Systems Biology: Design Principles of Biological Circuits*, Chapman & Hall/CRC, London, ISBN 1584886420, GBP 30.99, 2007 (320 pp.). 2008.
- [27] Srinivas Gorur Shandilya y Marc Timme. «Inferring network topology from complex dynamics». En: *New Journal of Physics* 13.1 (2011), pág. 013004.
- [28] Abhijeet R Sonawane, Scott T Weiss, Kimberly Glass y Amitabh Sharma. «Network Medicine in the age of biomedical big data». En: *Frontiers in Genetics* 10 (2019).
- [29] George Sugihara, Robert May, Hao Ye, Chih-hao Hsieh, Ethan Deyle, Michael Fogarty y Stephan Munch. «Detecting causality in complex ecosystems». En: *science* 338.6106 (2012), págs. 496-500.
- [30] Marc Timme y Jose Casadiego. «Revealing networks from dynamics: an introduction». En: *Journal of Physics A: Mathematical and Theoretical* 47.34 (2014), pág. 343001.
- [31] Ophelia S Venturelli, Alex V Carr, Garth Fisher, Ryan H Hsu, Rebecca Lau, Benjamin P Bowen, Susan Hromada, Trent Northen y Adam P Arkin. «Deciphering microbial interactions in synthetic human gut microbiome communities». En: *Molecular systems biology* 14.6 (2018).
- [32] Alejandro F Villaverde y Julio R Banga. «Reverse engineering and identification in systems biology: strategies, perspectives and challenges». En: *Journal of the Royal Society Interface* 11.91 (2014), pág. 20130505.
- [33] EUGENE P WIGNER. «The Unreasonable Effectiveness of Mathematics in the Natural Sciences». En: *COMMUNICATIONS ON PURE AND APPLIED MATHEMATICS* 13 (1960), págs. 001-14.
- [34] Sophie Weiss, Will Van Treuren, Catherine Lozupone, Karoline Faust, Jonathan Friedman, Ye Deng, Li Charlie Xia, Zhenjiang Zech Xu, Luke Ursell, Eric J Alm y col. «Correlation detection strategies in microbial data sets vary widely in sensitivity and precision». En: *The ISME journal* 10.7 (2016), pág. 1669.
- [35] Hugh R Wilson y Jack D Cowan. «Excitatory and inhibitory interactions in localized populations of model neurons». En: *Biophysical journal* 12.1 (1972), págs. 1-24.
- [36] Yandong Xiao, Marco Tulio Angulo, Jonathan Friedman, Matthew K Waldor, Scott T Weiss y Yang-Yu Liu. «Mapping the ecological networks of microbial communities». En: *Nature communications* 8.1 (2017), pág. 2042.
- [37] Yury V Zaytsev, Abigail Morrison y Moritz Deger. «Reconstruction of recurrent synaptic connectivity of thousands of neurons from simulated spiking activity». En: *Journal of computational neuroscience* 39.1 (2015), págs. 77-103.