



UNIVERSIDAD NACIONAL DE ROSARIO

TESINA DE GRADO
PARA LA OBTENCIÓN DEL GRADO DE
LICENCIADO EN CIENCIAS DE LA COMPUTACIÓN

Vale por un título

Autor:

Tu persona

Directores:

XX

YY

Departamento de Ciencias de la Computación
Facultad de Ciencias Exactas, Ingeniería y Agrimensura
Av. Pellegrini 250, Rosario, Santa Fe, Argentina

22 de abril de 2025

May the force be with you.

*YODA (Y PARTICULARMENTE SEBA EN ESTE
MOMENTO)*

Agradecimientos

Quiero agradecer a quien quiera agradecer. (Por ejemplo, a la universidad pública).

Resumen

Versión cortita de tu tesina.

Índice general

Agradecimientos	III
Resumen	V
Índice general	VI
1 Introducción	1
1.1. Contribuciones	1
1.2. Organización del trabajo	1
2 Preliminares	3
2.1. sigma algebras	3
2.2. Rabin	3
3 Verificación de modelos con juegos	5
3.1. Sobre la verificación de modelos	5
3.2. Cadenas de Markov	5
3.3. Procesos de Decisión de Markov	7
3.4. Juegos estocásticos	11
3.4.1. Una breve historia sobre los juegos estocásticos	11
3.4.2. Draft	11
3.5. Juegos deterministas	12
3.6. Comparación de los distintos juegos	13
4 Juegos Estocásticos Politópicos - PSGs	15
4.1. Definiciones	15
4.1.1. Politopos	15
4.1.2. PSGs	15

4.2. Teoremas	18
5 Objetivos de Rabin en juegos estocásticos politópicos	19
5.1. Procesos de Decisión de Markov Politópicos - PMDPs	19
5.1.1. Draft	19
5.2. Juegos justos: la respuesta a la pregunta cualitativa	28
5.2.1. Juegos de adversario justo	28
5.2.2. Desrandomización de PSGs	30
5.2.3. Prueba de igualdad sobre los conjuntos ganadores	32
5.3. Transformar Rabin en alcanzabilidad: la respuesta a varias preguntas . .	37
5.3.1. Draft	37
6 Conclusiones	41
6.1. Trabajo Futuro	41
Referencias	43
A Título del Apendice	45
A.1. Título de la seccion	45

Capítulo 1

Introducción

1.1. Contribuciones

...

1.2. Organización del trabajo

- En el **Capítulo ??**, presentaré ...;
- En el **Capítulo 3**, extenderé y explicaré en detalle los métodos de esta tesina ...;
- En el **Capítulo 5**, presentaré mis resultados y conseguidos ...;
- Finalmente, en el **Capítulo 6**, concluyo con un resumen de los aportes realizados en esta tesina, menciono las implicancias de esta investigación y hallazgos realizados, y sugiero potenciales caminos para futuras investigaciones.

Capítulo 2

Preliminares

2.1. sigma algebras

2.2. Rabin

Capítulo 3

Verificación de modelos con juegos

3.1. Sobre la verificación de modelos

blabla

problema de síntesis de church

idea de la necesidad de juegos y nuestro entendimiento de ellos

¿Que preguntas nos hacemos cuando trabajamos con juegos? Mas que nada con juegos estocásticos tenemos las preguntas que se hace Kucera, las que presenta Chatterjee y se puede investigar mas Capaz esto iría mas en la introducción para plantear la idea de trabajo

3.2. Cadenas de Markov

Definición 3.2.1 (Cadena de Markov). *Una cadena de Markov es una tupla $M = (S, P)$ donde S es un conjunto (finito ?) no vacío de estados y $P : S \times S \rightarrow [0, 1]$ es una función*

tal que para todos los estados s vale que

$$\sum_{s' \in S} P(s, s') = 1$$

La función de probabilidad de transición P especifica para cada estado s la probabilidad $P(s, s')$ de moverse de s a s' en un solo paso. La restricción impuesta en P asegura que la función sea una distribución.

Una cadena de Markov induce un grafo subyacente, donde los estados actúan como vértices y hay una arista entre s y s' si y solo si $P(s, s') > 0$. Las cadenas de Markov se suelen representar por su grafo subyacente donde sus aristas estarán anotadas con las probabilidades en el intervalo $(0, 1]$.

Los caminos en una cadena de Markov son los caminos en el grafo subyacente. Son definidos como secuencias infinitas de estados $\omega = (s_0, s_1, s_2, \dots) \in S^\omega$ tales que $P(s_i, s_{i+1}) > 0$ para todo $i \geq 0$.

Veamos un pequeño ejemplo de cómo sería una cadena de Markov.

Ejemplo 3.2.1. *Supongamos que queremos observar el comportamiento de un pequeño robot llamado Roborto. Roborto se comporta probabilísticamente de la siguiente manera: desde su posición inicial de quietud tiene una probabilidad de ...*

Para poder asociar probabilidades a eventos en cadenas de Markov, la noción intuitiva de probabilidades en M es formalizada al asociarle un espacio de probabilidad (2.1). Los caminos infinitos de M juegan el rol de resultados. Esto es $Outc^M = Paths(M)$. La σ -álgebra asociada con M es generada por los conjuntos cilindro formados por los fragmentos de caminos finitos en M .

Me falta la noción de prefijo.

Definición 3.2.2 (Conjunto cilindro). *El conjunto cilindro de $\hat{\omega} = (s_0, \dots, s_n) \in Paths_{fin}(M)$ está definido como*

$$Cyl(\hat{\omega}) = \{\omega \in Paths(M) \mid \hat{\omega} \in pref(\omega)\}$$

Definición 3.2.3 (σ -álgebra de una cadena de Markov). La σ -álgebra \mathfrak{E}^M asociada a la cadena de Markov M es la σ -álgebra más pequeña que contiene todos los conjuntos cilindro $Cyl(\hat{\omega})$ donde $\hat{\omega} \in Paths_{fn}(M)$.

De conceptos clásicos de teoría de probabilidad se sigue que por cada estado s existe una única medida de probabilidad \mathbb{P}_s^M en la σ -álgebra \mathfrak{E}^M asociada a M , donde las probabilidades para los conjunto cilindros (es decir, los eventos) están dadas por:

$$\mathbb{P}_s^M(Cyl(s_0, \dots, s_n)) = \iota(s, s_0) \cdot \prod_{0 \leq i < n} P(s_i, s_{i+1}), \text{ donde } \iota(s, s_0) = \begin{cases} 1 & \text{si } s = s_0 \\ 0 & \text{en otro caso} \end{cases} \quad (3.1)$$

No estoy muy segura de cómo manejar lo de estado inicial.

Notación: en lo que sigue, usaremos notación LTL para describir ciertos eventos en cadenas de Markov. Por ejemplo, para un conjunto $B \subseteq S$ de estados, $\Diamond B$ denota el evento de llegar eventualmente a (algún estado en) B , mientras que $\text{siemprevent} B$ describe el evento en el que B es visitado infinitamente a menudo. Algo más? sí

EJEMPLO DE ROBOT y la pregunta de alcanzabilidad.

3.3. Procesos de Decisión de Markov

Un proceso de decisión de Markov (MDP, por sus siglas en inglés) es una generalización de una cadena de Markov donde un conjunto de acciones posibles es asociado a cada estado. A cada par estado-acción le corresponde una distribución de probabilidad en los estados, que es usada para seleccionar el próximo estado. A su vez, una cadena de Markov se corresponde a un MDP donde hay exactamente una acción asociada a cada estado. Asumiremos la existencia de un conjunto fijo de acciones Act . La definición de un proceso de decisión de Markov es como sigue:

Definición 3.3.1 (Proceso de Decisión de Markov). Un proceso de decisión de Markov $\mathcal{M} = (S, A, \theta)$ consiste de un conjunto finito de estados S y de dos componentes A y θ que especifican la estructura de transición:

- A es un conjunto de acciones. Para cada $s \in S$, $A(s) \subseteq \text{Acts}$ es el conjunto finito no vacío de acciones disponibles en s . Para cada estado $s \in S$ se requiere que $A(s) \neq \emptyset$.
- $\theta : S \times A \times S \rightarrow [0, 1]$ es una función de transición probabilística. Para todo estado $s \in S$, si $a \in A(s)$ tenemos que $\sum_{s' \in S} \theta(s, a, s') = 1$, mientras que si $a \notin A(s)$, $\sum_{s' \in S} \theta(s, a, s') = 0$. Para cada $s, t \in S$ y $a \in A(s)$, $\theta(s, a, t)$ es la probabilidad de transicionar de s a t cuando la acción a es seleccionada.

Agregar ejemplo de robot

Un comportamiento en un proceso de decisión de Markov es una secuencia alternante infinita de estados y acciones, construida iterativamente por un proceso de dos pasos. Primero, dado un estado s , una acción $a \in A(s)$ es seleccionada no-determinísticamente. Luego, el sucesor t de s es seleccionado de acuerdo a la distribución asociada a la acción a . La definición formal es como sigue:

Definición 3.3.2 (Comportamiento en un MDP). *Un comportamiento en un MDP \mathcal{M} es una secuencia infinita $\omega = (s_0, a_0, s_1, a_1, \dots)$ tal que $a_i \in A(s_i)$ y $a_i(s_{i+1}) > 0$ para todo $i \geq 0$.*

Dado un estado s , indicaremos con Ω_s el conjunto de todos los comportamientos que se originan en s , con $\Omega_{\mathcal{M}}$ el conjunto de todos los comportamientos en \mathcal{M} y con $\Omega_{\mathcal{M}}^{\text{fin}}$ el conjunto de todos los prefijos finitos de comportamientos en \mathcal{M} que terminan en un estado $s \in S$.

Para cadenas de Markov, el conjunto de caminos está equipado con una σ -álgebra y una medida de probabilidad que refleja la noción intuitiva de probabilidad para conjuntos (medibles) de caminos. Para los MDPs, esto es levemente distinto. Como no hay restricciones en la resolución de las elecciones no determinista, no hay una única medida de probabilidad asociada. (capaz acá podría ir algo tipo un ejemplo como el de la página 841)

Para poder razonar sobre probabilidades de conjuntos de comportamientos en un MDP necesitamos resolver de alguna manera el no determinismo, y para ello introduciremos el concepto de estrategia.

Definición 3.3.3 (Estrategia en un MDP). *Sea $\mathcal{M} = (S, A, \theta)$ un MDP. Una estrategia para \mathcal{M} es una función $\pi : \Omega^{\text{fin}} \rightarrow \text{Dist}(A)$ que asigna una distribución de probabilidad a cada prefijo finito de comportamiento tal que $\pi(\hat{\omega})(a) > 0$ solo si $a \in A(s)$.*

Comentar algo de que en la literatura se suelen bajar las acciones de la definición de estrategia? O capaz la idea sería dropear las acciones de la definición de estrategia acá?

Lo que está a continuación me hace un poco de ruido por la parte de la definición de la cadena de Markov. Es como lo presentan en el libro de Baier y Katoen y tiene esta idea de ir armando una continuidad entre MCs y MDPs en la explicación. La otra opción es ir poco más con el enfoque de Luca de Alfaro y presentar la σ -álgebra sin mencionar nada de MDPs.

Como una estrategia resuelve todas las elecciones no deterministas en un MDP, induce una cadena de Markov. Esto es, el funcionamiento de un MDP \mathcal{M} siguiendo las decisiones de una estrategia π puede ser formalizado por una cadena de Markov \mathcal{M}_π , donde los estados son los prefijos finitos de comportamientos en \mathcal{M} .

Definición 3.3.4 (Cadena de Markov de un MDP inducida por una estrategia). *Sea $\mathcal{M} = (S, A, \theta)$ un MDP y π una estrategia en \mathcal{M} . La cadena de Markov \mathcal{M}_π está dada por*

$$\mathcal{M}_\pi = (\Omega_{\mathcal{M}}^{\text{fin}}, P_\pi)$$

donde para $\hat{\omega} = (s_0, a_0, s_1, a_1, \dots, s_n)$:

$$P_\pi(\hat{\omega}, \hat{\omega}as_{n+1}) = \pi(\hat{\omega})(a) \cdot \theta(s_n, a, s_{n+1})$$

Nótese que \mathcal{M}_π cuenta con un espacio de estados infinito aun, aun cuando el MDP \mathcal{M} es finito. Entonces cuenta como una MC?.

Como \mathcal{M}_π es una cadena de Markov, uno ahora puede razonar sobre las probabilidades de los conjuntos medibles de caminos que siguen la estrategia π , simplemente usando las distintas medidas de probabilidad $\mathbb{P}_s^{\mathcal{M}_\pi}$ asociadas a la cadena de Markov \mathcal{M}_π (véase 3.1).

Responder preguntas de probabilidad de cosas en MDP robot

Intuitivamente, el estado (s_0, a_0, \dots, s_n) de \mathcal{M}_π representa la configuración donde el MDP \mathcal{M} está en el estado s_n y cuenta con la historia $(s_0, a_0, \dots, s_{n-1}, a_{n-1})$. Según la definición que vimos las estrategias pueden depender de la historia en su totalidad, produciendo resultados distintos si al menos una acción o estado en su historia cambia, pero es cierto que este caso no es lo usual. Veamos algunos tipos particulares de estrategias donde esto no sucede.

Definición 3.3.5 (Estrategias sin memoria). *Sea \mathcal{M} un MDP con espacio de estados S . Una estrategia π en \mathcal{M} es sin memoria sii para cada par de comportamientos (s_0, a_0, \dots, s_n) y (t_0, a'_0, \dots, t_m) con $s_n = t_m$ vale que:*

$$\pi(s_0, a_0, \dots, s_n) = \pi(t_0, a'_0, \dots, t_m)$$

En este caso, π puede ser vista como una función $\pi : S \rightarrow \text{Dist}(A)$.

Coloquialmente, una estrategia es sin memoria si no recuerda nada de la historia y solo elige probabilidades para las acciones basándose en el estado actual. Esto puede ser bastante extremo en ciertos casos, por eso existe una variante que busca reflejar la idea de finitud sin ser tan restrictiva: las estrategias de memoria finita. Una estrategia de memoria finita puede ser pensada intuitivamente como que solo puede guardar hasta una cantidad finita fija de información de la historia, por lo que no podrá ser distinta para **todo** prefijo finito de comportamiento. Formalmente, la definiremos a través de un autómata determinista finito (DFA). La distribución de probabilidad de las acciones será seleccionada a partir del estado actual en \mathcal{M} y el estado actual del autómata (al que llamaremos modo). Veamos su definición:

Definición 3.3.6 (Estrategias con memoria finita). *Sea \mathcal{M} un MDP con espacio de estados S y conjunto de acciones A . Una estrategia de memoria finita para \mathcal{M} es una tupla $\pi = (Q, f_\pi, \Gamma, \text{start})$ donde*

- Q es un conjunto finito de modos,
- $\Delta : Q \times A \times S \rightarrow Q$ es la función de transición del autómata,
- $\text{start} : S \rightarrow Q$ es la función que determina el modo en el que empieza el autómata para un estado inicial s ,

- $f_\pi : Q \times S \rightarrow \text{Dist}(A)$ es la función que asigna la distribución de probabilidad en las acciones desde un estado s , es decir, lo que veníamos entendiendo como estrategia en sí.

El funcionamiento del MDP bajo la estrategia de memoria finita sería como sigue. En principio, se inicializa el modo del DFA a $q_0 = \text{start}(s_0)$. Luego, desde cada estado s_i posterior el proceso será iterativo. Primero, se seleccionará la distribución de probabilidad en las acciones a partir del modo actual q_i del autómata con $f_\pi((q_i, s_i))$. Una vez tomada la decisión, se determina probabilísticamente la siguiente acción a_{i+1} , y, a partir de ella, se determina también probabilísticamente el siguiente estado s_{i+1} . Con la nueva acción y estado se seleccionará el próximo modo del DFA $q_{i+1} = \Delta(q_i, a_{i+1}, s_{i+1})$ y se repetirá el proceso.

Para $\hat{\omega} \in \Omega_{\mathcal{M}}^{\text{fin}}$, notaremos con $\pi(\hat{\omega})$ a la distribución obtenida al realizar el proceso explicado anteriormente con $\hat{\omega}$ pero qué pasa con la elección probabilística de las acciones? No habría que tenerlo? Y que hay de sin memoria? Ahí en realidad entonces no tendría que ser último estado y penúltima acción?

Además de su categorización en base a qué tanto dependen de su historia, existe otro tipo notable de estrategias: las estrategias deterministas. Una estrategia es determinista cuando En la literatura (véase [1] [2]) es usual encontrarse con el estudio de estrategias solo en su variante determinista.

Mencionar de que tipos son las estrategias vistas en el ejemplo de robot

3.4. Juegos estocásticos

3.4.1. Una breve historia sobre los juegos estocásticos

Enfoque mc ->mdp ->sg ->psg

Enfoque juegos deterministas ->juegos estocasticos (Shapley etc)

3.4.2. Draft

Cosas de Kucera y de la parte preliminar del paper de Pedro.

Me gustaría incluir cosas de Kucera/Banerjee y sus ideas de «1» val etc. (concepto de distintas maneras de ganar)

Capaz algo de omega regular stochastic games y qsy

Objetivos

distintos tipos y organizacion / ver de que manera presentarlos

3.5. Juegos deterministas

Agg transiciones de texto escrito entre definiciones

Definición 3.5.1 (juego de grafo de 2 jugadores). *Definimos a un juego de grafo de dos jugadores ($2G$) como una tupla $G = (V, V_\square, V_\diamond, E)$ donde $V = V_\square \uplus V_\diamond$ es un conjunto de vértices (o estados) particionado en V_\square y V_\diamond , y $E \subseteq (V \times V)$ es una relación que denota el conjunto de aristas (dirigidas) que representan transiciones de un estado a otro del juego.*

Los 2 jugadores son llamados \square y \diamond y controlan los vértices V_\square y V_\diamond , respectivamente.

Definición 3.5.2 (estrategia sobre un $2G$). *Una estrategia para un jugador i con $i \in \{\square, \diamond\}$ es una función $\sigma_i : V^*V_i \rightarrow V$ con la restricción de que $\sigma_i(wv) \in E(v)$ para todo $wv \in V^*V_i$.*

Definición 3.5.3 (jugada sobre un $2G$). *Una jugada en un juego de grafo de dos jugadores es una secuencia infinita de vértices $\rho = v_0v_1v_2\cdots \in V^\omega$, donde para todo $i \in \mathbb{N}_0$ tenemos que $v^i \in V$ y $(v^i, v^{i+1}) \in E$.*

Sean σ_\square y σ_\diamond un par de estrategias y sea v_0 un vértice inicial, la jugada que cumple con σ_\square y σ_\diamond es la única jugada $\rho = v_0v_1v_2\cdots$ para la cual cada $i \in \mathbb{N}_0$, si $v_i \in V_\square$ entonces $v_{i+1} = \sigma_\square(v_0 \dots v_i)$, y si $v_i \in V_\diamond$ entonces $v_{i+1} = \sigma_\diamond(v_0 \dots v_i)$.

Definición 3.5.4 (condición ganadora). *Una condición ganadora φ en un juego de grafo de dos jugadores es un conjunto de jugadas sobre el juego, i.e., $\varphi \subseteq V^\omega$. Usaremos notación LTL para describir conjuntos de jugadas específicos.*

Definición 3.5.5 (regiones ganadoras). *El jugador \square gana el juego de grafo de dos jugadores G para una condición ganadora φ desde un vértice v_0 si existe una estrategia*

π_{\square} tal que para cada π_{\diamond} , la jugada ρ que sigue π_{\square} y π_{\diamond} satisface φ , i.e., $\rho \in \varphi$. La región ganadora $\mathcal{W} \subseteq V$ para el jugador \square es el conjunto de vértices desde donde el jugador \square gana el juego.

Agg ejemplo

3.6. Comparación de los distintos juegos

Incluir cuadros comparativos entre juegos // ejemplos para todas las nociones nuevas explicadas

(Capaz agregar un apéndice con cuadritos de esto ?)

(En general, estaría bueno ver problemas y agregar un apéndice con ejemplos)

Capítulo 4

Juegos Estocásticos Politópicos - PSGs

4.1. Definiciones

4.1.1. Politopos

4.1.2. PSGs

Los juegos estocásticos politópicos fueron presentados en 2024 por Castro y D’Argenio [3]. La idea de su creación ...

Un juego estocástico politópico se caracteriza a través de una estructura que contiene un conjunto finito de estados divididos en dos conjuntos, cada uno perteneciente a un jugador diferente. Además, a cada estado se le asigna un conjunto finito de politopos convexos de distribuciones de probabilidad sobre los estados. La definición formal es como sigue:

Definición 4.1.1 (PSG). *Un juego estocástico politópico (abreviado PSG, por sus siglas en inglés) es una estructura $\mathcal{K} = (\mathcal{S}, (\mathcal{S}_\square, \mathcal{S}_\diamond), \Theta)$ tal que \mathcal{S} es un conjunto finito de estados particionado en $\mathcal{S} = \mathcal{S}_\square \uplus \mathcal{S}_\diamond$ y $\Theta : \mathcal{S} \rightarrow \mathcal{P}_f(\text{DPoly}(\mathcal{S}))$.*

La idea de un PSG es la esperada: en un estado $s \in \mathcal{S}_i$ (para $i \in \{\square, \diamond\}$), el jugador

i elige jugar un politopo $K \in \Theta(s)$ y una distribución $\mu \in K$. El siguiente estado s' se selecciona de acuerdo con la distribución μ , y el juego continúa desde s' repitiendo el mismo proceso.

El desarrollo de un juego estocástico politópico se interpreta en términos de un juego estocástico donde el número de transiciones salientes desde los estados de los jugadores puede ser no numerable. Formalmente, la interpretación de un PSG es la siguiente.

Definición 4.1.2 (Interpretación de un PSG). *La interpretación del juego estocástico politópico \mathcal{K} se define por el juego estocástico $\mathcal{G}_{\mathcal{K}} = (\mathcal{S}, (\mathcal{S}_{\square}, \mathcal{S}_{\diamond}), \mathcal{A}, \theta)$, donde $\mathcal{A} = \bigcup_{s \in \mathcal{S}} \Theta(s) \times \text{Dist}(\mathcal{S})$ y*

$$\theta(s, (K, \mu), s') = \begin{cases} \mu(s') & \text{si } K \in \Theta(s) \text{ y } \mu \in K \\ 0 & \text{en otro caso.} \end{cases}$$

Nótese que el conjunto de acciones \mathcal{A} puede ser no numerable, al igual que cada conjunto $\mathcal{A}(s) = \bigcup_{K \in \Theta(s)} \{K\} \times K$ de todas las acciones realizables en el estado s , identificado por el politopo elegido y la distribución seleccionada dentro del politopo. Por lo tanto, necesitamos extender las estrategias a este dominio no numerable, que debe estar dotado de una σ -álgebra adecuada.

Para esto, utilizamos una construcción estándar para darle a $\text{Dist}(S)$ una σ -álgebra: $\Sigma_{\text{Dist}(S)}$ se define como la σ -álgebra más pequeña que contiene los conjuntos $\{\mu \in \text{Dist}(S) \mid \mu(S) \geq p\}$ para todo $S \subseteq S$ y $p \in [0, 1]$. Ahora, dotamos a \mathcal{A} con la σ -álgebra producto $\Sigma_{\mathcal{A}} = \mathcal{P}(\bigcup_{s \in \mathcal{S}} \Theta(s)) \otimes \Sigma_{\text{Dist}(S)}$, y llamamos $\text{PMeas}(\mathcal{A})$ al conjunto de todas las medidas de probabilidad sobre $\Sigma_{\mathcal{A}}$. No es difícil comprobar que cada conjunto de acciones habilitadas $\mathcal{A}(s)$ es medible (es decir, $\mathcal{A}(s) \in \Sigma_{\mathcal{A}}$) y que la función $\theta(s, \cdot, s')$ es medible (es decir, $\{a \in \mathcal{A} \mid \theta(s, a, s') \leq p\} \in \Sigma_{\mathcal{A}}$ para todo $p \in [0, 1]$).

Ahora bien, para definir formalmente las estrategias, introduciremos el **concepto de comportamiento en PSGs**. Esto difiere de la formulación original hecha para PSGs en [3], pero no afecta a los resultados que se mostrarán en este trabajo.

Definición 4.1.3 (Comportamiento en un PSG). *Un comportamiento en un PSG es una secuencia infinita que alterna entre estados y conjuntos resultado. Formalmente un comportamiento es un $\omega = (s_0, (K_0, \mu_0), s_1, (K_1, \mu_1), \dots)$ tal que $s_i \in \mathcal{S}$, $(K_i, \mu_i) \in \mathcal{A}(s_i)$ y $\mu_i(s_{i+1}) > 0$ para todo $i \geq 0$.*

Notaremos al conjunto de todos los comportamientos de un PSG como Ω . Dado un estado s , indicaremos con Ω_s el conjunto de los comportamientos que se originan en s , y con Ω_s^n al conjunto de los comportamientos que se originan en s y tiene un total de n estados.

Definición 4.1.4 (Conjuntos soporte y acciones). Dada una acción (K, μ) definiremos su conjunto soporte, $\text{supp}((K, \mu))$ como el conjunto formado por los estados a los cuales (K, μ) les asigna una probabilidad positiva. Es decir,

$$\text{supp}((K, \mu)) = \{s \in \mathcal{S} \mid \mu(s) > 0\}$$

Dado un estado s y un conjunto de estados V , definiremos como $\text{acc}(s, V)$ a las acciones que parten desde s y tienen como conjunto soporte V . Es decir,

$$\text{acc}(s, V) = \{\alpha \in \mathcal{A}(s) \mid \text{supp}(\alpha) = V\} = \{(K, \mu) \in \mathcal{A}(s) \mid \mu(V) = 1 \wedge \forall v \in V, \mu(v) > 0\}$$

Por último, dado un estado s , definiremos como V_s al conjunto de todos los soportes que pueden tener las distintas acciones desde s . Es decir,

$$\begin{aligned} V_s &= \{\text{supp}(\mu) \mid \exists K \in \Theta(s) : \mu \in K\} = \\ &= \{V' \subseteq \mathcal{S} \mid \exists \mu \in K, K \in \Theta(s) \text{ tal que } \mu(V') = 1 \text{ y } \forall s' \in V', \mu(s') > 0\} \end{aligned}$$

Entonces, ahora sí podemos extender el concepto de estrategia en un PSG:

Definición 4.1.5 (Estrategia en un PSG). Una estrategia π_i para el jugador i ($i \in \{\square, \diamond\}$) en un PSG será una función $\pi_i : (\mathcal{S} \times \mathcal{A})\mathcal{S}_i \rightarrow \text{PMeas}(\mathcal{A})$ que asigna una medida de probabilidad a cada ωs tal que $\pi_i(\omega s)(\mathcal{A}(s)) = 1$.

Habría que agregar alguna justificación de esta idea de comportamientos?

Medida de probabilidad que no va a ir probablemente

Con la formalización del concepto de estrategia ahora podemos presentar formalmente la definición de $\mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond}}$, la medida de probabilidad definida por la cadena de Markov dada por $\mathcal{G}_{\mathcal{K}}$ y las estrategias π_{\square} y π_{\diamond} .

Para eso, primero para cada $n \geq 0$ y $s \in \mathcal{S}$ definimos $\mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond},n} : \Omega^{n+1} \rightarrow [0, 1]$ para todo $s' \in \mathcal{S}$ y $\hat{\omega} \in \Omega_s^{n+1}$ inductivamente como sigue:

$$\begin{aligned} \mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond},0}(s') &= \delta_s(s') \\ \mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond},n+1}(\hat{\omega}s_n V' s') &= \mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond},n}(\hat{\omega}s_n) \int_{\{a \in A(s_n) \mid \text{sop}(a)=V'\}} \theta(s_n, a, s') d(\pi_i(\hat{\omega}s_n)(a)) \\ &\quad \text{si } s_n \in \mathcal{S}_i \text{ con } i \in \{\square, \diamond\} \end{aligned}$$

($\text{sop}(a) = V'$ significa que el soporte de a es V')

(capaz restringirse sobre las acciones con soporte V' ? No sé bien cómo va a quedar esto)

y extendemos $\mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond},n} : \mathcal{P}(\Omega^{n+1}) \rightarrow [0, 1]$ a conjuntos como la suma de la medida sobre los elementos del conjunto.

Sea Σ_{Ω} la σ -álgebra discreta sobre Ω (pues tanto \mathcal{S} como los conjuntos soportes serán conjuntos finitos) y $\Sigma_{\Omega^{\omega}}$ la σ -álgebra producto usual sobre Ω^{ω} . Por el teorema de extensión de Carathéodory, $\mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond}} : \Sigma_{\Omega^{\omega}} \rightarrow [0, 1]$ se define como la única medida de probabilidad tal que para todo $n \geq 0$, y $SV_i \in \Sigma_{\Omega}$, $0 \leq i \leq n$,

$$\mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond}}(SV_0 \times \cdots \times SV_n \times \Omega^{\omega}) = \mathbb{P}_{\mathcal{G}_{\mathcal{K},s}}^{\pi_{\square},\pi_{\diamond},n}(SV_0 \times \cdots \times SV_n).$$

4.2. Teoremas

En esta sección presentamos los resultados desarrollados en [3] que nos servirán más adelante en el documento.

Algo más que teo1?

AGG RDOS DEL PAPER DE PEDRO

Capítulo 5

Objetivos de Rabin en juegos estocásticos politópicos

5.1. Procesos de Decisión de Markov Politópicos - PMDPs

5.1.1. Draft

Acá se está dando directamente la interpretación de un PMDP

Lo podría definir directamente como PSGs con un solo jugador.

Definición 5.1.1. *Un proceso de decisión de Markov politópico (PMDP por sus siglas en inglés) es una tupla $\mathcal{M} = (\mathcal{S}, Act, \theta')$ donde \mathcal{S} es un conjunto finito de estados, Act es un conjunto de pares (politopo, distribución) y $\theta' : \mathcal{S} \times Act \times \mathcal{S} \rightarrow [0, 1]$ es la función de transición entre estados. También podemos definir a un PDMP como un PSG donde $\mathcal{S}_{\square} = \emptyset$ o $\mathcal{S}_{\diamond} = \emptyset$.*

También introducimos el concepto de *comportamiento* en un PMDP. Un comportamiento será una secuencia alternante de estados y conjuntos de próximos estados posibles, que refleja un proceso de selección de dos pasos. Desde un estado s , primero, el jugador selecciona un politopo y una distribución cuyo conjunto soporte será un $V \in V_s$,

y luego el próximo estado s' será elegido probabilísticamente en base a la distribución seleccionada. Presentamos la definición formal de la siguiente manera:

Los comportamientos, las estrategias, y la medida de probabilidad en los PMDPs se definen tomando las mismas definiciones de los juegos estocásticos politópicos, mientras que las definiciones de conjuntos estado-acción, sub-MDPs y componentes finales de los procesos de decisión de Markov tradicionales se extienden a los PMDPs de la siguiente manera:

Definición 5.1.2 (conjuntos estado-resultado y sub-PMDPs). *Dado un PMDP $\mathcal{M} = (S, Act, \theta')$ un conjunto estado-resultado es un subconjunto $\chi \subseteq \{(s, V') \mid s \in S \wedge V' \in V_s\}$. Un sub-PMDP es un par (C, D) , donde $C \subseteq S$ y D es una función que asocia a cada $s \in C$ un conjunto $D(s) \subseteq V_s$ de subconjuntos de estados próximos posibles. Hay una relación uno-a-uno entre sub-PMDPs y conjuntos de estado-acción:*

- *dado un conjunto estado-resultado χ , denotamos $sub(\chi) = (C, D)$ al sub-PMDP definido por:*

$$C = \{s \mid \exists V'. (s, V') \in \chi\} \quad D(s) = \{V' \mid (s, V') \in \chi\}$$

- *dado un sub-PMDP (C, D) , denotamos por $er(C, D) = \{(s, V') \mid s \in C \wedge V' \in D(s)\}$ al conjunto estado-resultado correspondiente a (C, D) .*

Si definimos para cada V' , un vértice único nuevo $v_{V'}$ podemos ver que cada sub-PMDP (C, D) induce una *relación de aristas*: hay una arista $(s, v_{V'})$ de $s \in C$ a $v_{V'}$ para cada $V' \in D(s)$ y hay una arista de $(v_{V'}, t)$ de $v_{V'}$ con $V' \in D(s)$ a $t \in \mathcal{S}$ sii es posible ir de s a t en un paso con probabilidad positiva utilizando una acción cuyo conjunto resultado sea V' . La definición formal es como sigue:

Definición 5.1.3 (relación de aristas ρ). *Para un sub-PMDP, definimos la relación $\rho_{(C,D)}$ como*

$$\rho_{(C,D)} = \{(s, v_{V'}) \mid \exists V' \in D(s)\} \cup \{(v_{V'}, t) \mid t \in V'\}$$

Definición 5.1.4 (componente final). *Un sub-PMDP es una componente final si:*

- $V' \subseteq C$ para todo V' tal que existe un $s \in C$ donde $V' \in D(s)$

- el grafo $(C \cup \{v'_V \mid \exists s \in C : V' \in D(s)\}, \rho_{(C,D)})$ es fuertemente conexo.

Llamaremos $EC(\mathcal{M})$ al conjunto de todas las componentes finales en un PMDP \mathcal{M} .

Intuitivamente, una componente final representa un conjunto de pares estado-resultado que, una vez en ellos, es posible quedarse allí si la estrategia escoge las acciones de manera apropiada. Esta intuición se hará precisa con los siguientes teoremas.

Antes de enunciar estos teoremas, introducimos una abreviatura para el conjunto de estados-resultados que ocurren infinitamente a menudo en un comportamiento dado y una notación para el conjunto (infinito) de acciones con el mismo conjunto resultado.

Definición 5.1.5 (*inf*). Dado un comportamiento $\omega = (s_0, \alpha_0, s_1, \alpha_1, \dots)$ indicamos por

$$inf(\omega) = \{(s, V') \mid s_k = s \wedge \text{supp}(\alpha_k) = V' \text{ para infinitos } k \in \mathbb{N}_0\}$$

al conjunto de pares estado-resultado que ocurren infinitas veces en él.

Ahora sí, podemos pasar a presentar las primeras demostraciones:

Teorema 5.1.1 (estabilidad de componentes finales). Sea (C, D) una componente final. Entonces, para cada estrategia π existe una estrategia π' , que difiere de π solo en C , tal que:

$$\mathbb{P}_s^\pi(\diamond C) = \mathbb{P}_s^{\pi'}(\diamond C) = \mathbb{P}_s^{\pi'}(inf(\omega) = er(C, D)) \quad (5.1)$$

para todo $s \in S$.

Demostración. Considérese una estrategia π' definida como sigue para cada secuencia $s_0 \dots s_n$ con $n \geq 0$:

- Si $s_n \in C$, la estrategia asignará probabilidad positiva a una única acción $(K, \mu) \in A(s_n, V')$ para cada conjunto resultado V' (notaremos a esta acción particular con $(K, \mu)_{V'}$), y la probabilidad de elegir cada una de esas acciones se distribuirá de manera uniforme. Es decir,

$$\pi'(s_0 \dots s_n)(K, \mu) = \begin{cases} \frac{1}{|D(s_n)|} & \text{si } (K, \mu) = (K, \mu)_{V'} \text{ para algún } V' \in D(s); \\ 0 & \text{en otro caso} \end{cases}$$

- Si $s_n \notin C$, la estrategia π' coincide con π , i.e.

$$\pi(s_0 \dots s_n)(K, \mu) = \pi'(s_0 \dots s_n)(K, \mu) \quad \forall (K, \mu) \in Act$$

La primera igualdad en 5.1 es una consecuencia del hecho de que π y π' coinciden fuera de C .

Para la segunda igualdad, basta con ver que bajo la estrategia π' una vez que un comportamiento entra a C , nunca sale de C ni se elige una acción que no esté en D . Es más, una vez en C un comportamiento visitará todos los estados de C infinitamente a menudo con probabilidad 1. \square

El próximo resultado establece que, para cualquier estado inicial y cualquier estrategia (de memoria finita), un comportamiento terminará con probabilidad 1 en una componente final. Esta es la razón detrás del nombre “componente final”.

Teorema 5.1.2 (teorema fundamental de las componentes finales). *Sea \mathcal{M} un PMDP. Para todo $s \in S$, toda estrategia π de memoria finita,*

$$\mathbb{P}_{\mathcal{M},s}^{\pi}(\{\omega \in Paths(s) \mid sub(inft(\omega)) \text{ es una componente final}\}) = 1$$

Demostración. Consideremos un sub-PMDP (C, D) que no sea una componente final y sea $\Omega_s^{(C,D)} = \{\hat{\omega} \in \Omega_s \mid inft(\hat{\omega}) = er(C, D)\}$ el conjunto de comportamientos cuyo conjunto de pares estado-resultado que se repiten infinitas veces en él forman el sub-PMDP (C, D) .

Si podemos mostrar que

$$\mathbb{P}_{\mathcal{M},s}^{\pi}(\{\omega \in Paths(s) \mid \omega \in \Omega_s^{(C,D)}\}) = 0 \tag{5.2}$$

como (C, D) es un sub-PMDP cualquiera y como hay una cantidad finita de sub-PMDPs en \mathcal{M} , esto es lo mismo que mostrar que

$$\mathbb{P}_{\mathcal{M},s}^{\pi}(\{\omega \in Paths(s) \mid sub(inft(\omega)) \text{ es una componente final}\}) = 1$$

. Veamos que vale 5.2, dividiendo en casos según cuál es la condición de la definición 5.1.4 que no se cumple para (C, D) :

- Primero, asumamos que existe un $(t, V') \in er(C, D)$ tal que $V' \notin C$.

Sabemos que cada comportamiento en $\Omega_s^{(C,D)}$ toma el par estado-resultado (t, V') infinitas veces. Llamemos I al conjunto de índices infinito que representa los momentos en los que se visita (t, V') . Indiquemos con μ_i a la distribución elegida en el momento $i \in I$ y definamos como $r_i = \sum_{u \in C} \mu_i(u)$ a la probabilidad de quedarnos en C en el momento $i \in I$ (que en cada caso será menor a 1 porque $V' \notin C$).

Como π es de memoria finita sucederá que π solo puede elegir una cantidad finita de acciones (y, por lo tanto, distribuciones) distintas desde t . Esto hace que el conjunto $R = \{r_i \mid i \in I\}$ tenga un máximo, llamémoslo r .

Para que valga que ω esté en $\Omega_s^{(C,D)}$ tiene que valer que en infinitos momentos i nos quedemos en C . Entonces que vale que $\mathbb{P}_s^\pi(\omega \in \Omega_s^{(C,D)}) < r^k$ para todo $k > 0$ natural. Como sabemos que $r < 1$, tenemos que $\mathbb{P}_{\mathcal{M},s}^\pi(\{\omega \in Paths(s) \mid \omega \in \Omega_s^{(C,D)}\}) = 0$.

- Si no, asumamos que existen $t_1, t_2 \in C$ tales que no hay camino de t_1 a t_2 en $(C \cup \{v'_V \mid \exists s \in C : V' \in D(s)\}, \rho(C, D))$.

La falta de camino de t_1 a t_2 en $(C, \rho_{(C,D)})$ implica que para cada subsecuencia $s_m V_m s_{m+1} \dots s_n$ de comportamiento en $\Omega_s^{(C,D)}$ que vaya de $s_m = t_1$ a $s_n = t_2$, hay 2 opciones:

1. existe un $j \in [m+1, n-1]$ tal que $s_j \notin C$. Como cada comportamiento en $\Omega_s^{(C,D)}$ contiene infinitas subsecuencias de t_1 a t_2 y tenemos una cantidad finita de estados, sabemos que habrá una cantidad infinita de s_j iguales. Pero si infinitas veces se toma un estado s_j entonces, $s_j \in C$, lo que contradice la hipótesis anterior. Absurdo.
2. existe un $j \in [m, n-1]$ tal que $V_j \notin D(j)$. Como cada comportamiento en $\Omega_s^{(C,D)}$ contiene infinitas subsecuencias de t_1 a t_2 y tenemos una cantidad finita de conjuntos resultado, sabemos que habrá una cantidad infinita de V_j iguales, con lo que $V_j \in D(j)$, lo que contradice la hipótesis anterior. Absurdo

Con lo que arribamos a que $\mathbb{P}_{\mathcal{M},s}^\pi(\{\omega \in Paths(s) \mid \omega \in \Omega_s^{(C,D)}\}) = 0$

□

Esto nos permite definir el siguiente corolario útil para el análisis de condiciones de Rabin:

Corolario 5.1.2.1. *Sea \mathcal{M} un PMDP, s un estado en él y sea π una estrategia (lo hacemos con estrategias o planificadores?) en él. Una condición de Rabin $\tilde{R} = \{\langle G_1, R_1 \rangle, \dots, \langle G_k, R_k \rangle\}$ se satisface desde s , siguiendo la estrategia π con probabilidad 1 si y solo si para cada componente final U alcanzable desde s , existe un $j \in \{1, 2, \dots, k\}$ tal que $U \cap R_j = \emptyset$ y $U \cap G_j \neq \emptyset$.*

Ahora bien, una pregunta muy natural que surge es "¿por qué nos restringimos a estrategias de memoria finita en este último teorema?"

Nota sobre la limitación a estrategias finitas del Teorema 5.1.2

La prueba del teorema 5.1.2 está inspirada en las pruebas realizadas para MDPs (véase [1] [2]). En ellas, la prueba consiste en proponer un sub-PMDP (C, D) arbitrario que no cumple la definición 5.1.4, dividir el análisis por casos de cómo no se cumple la definición de componente final y mostrar que en cada caso la probabilidad de que un comportamiento que siga la estrategia del enunciado genere un sub-PMDP como el propuesto es 0.

Si nos remitimos a la prueba que mostramos en 5.1.2, vemos que el método de prueba se ajusta bien a nuestro caso hasta llegar al primer ítem donde el argumento está explícitamente respaldado en el hecho de que π es de memoria finita. De esta manera, se puede decir que el conjunto de las distintas probabilidades de quedarse en C en los momentos i tiene un máximo, r , lo que permite acotar la productoria $\prod_{i \in I} r_i$ por $\lim_{k \rightarrow \infty} r^k$, que sabemos que converge a 0, puesto que $r < 1$.

Si no nos restringimos a estrategias de memoria finita no podemos garantizar esto. Sabemos que la productoria va a estar acotada inferiormente por 0 y superiormente por 1, pero bien podría no converger a 0, sino a algún otro número. Por ejemplo, veamos lo siguiente:

Consideremos un juego estocástico politópico donde tenemos un vértice t y un vértice s . Desde t , existen acciones μ de la forma $\mu(s) = \frac{1}{k^2}$, $\mu(t) = 1 - \frac{1}{k^2}$, por lo que se puede definir una estrategia de memoria infinita π a partir de la cantidad de t que se encuentran

en el comportamiento de la siguiente manera:

$$\begin{aligned}\pi(\omega t)(s) &= \frac{1}{(\#_t(\omega))^2} \\ \pi(\omega t)(t) &= 1 - \frac{1}{(\#_t(\omega))^2}\end{aligned}$$

donde $\#_t$ se define inductivamente sobre comportamientos de la siguiente manera:

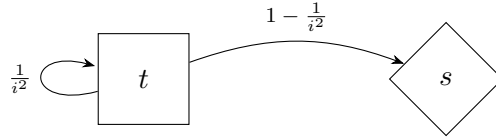


Figura 5.1: Interpretación del juego estocástico con la estrategia π fijada.

$$\begin{aligned}\#_t(\emptyset) &= 0 \\ \#_t(\omega V' t) &= \#_t(\omega) + 1 \\ \#_t(\omega V' s) &= \#_t(\omega)\end{aligned}$$

donde ω es un comportamiento, \emptyset representa el comportamiento vacío (no sé qué tan bien estaría esto) y $s \neq t$.

Si, remitiéndonos otra vez a la estrategia de prueba, suponemos que s no forma parte de C , tenemos que la probabilidad de quedarnos en C sería $\prod_{i \geq 0} 1 - (\frac{1}{i^2})$, que sabemos que converge a $\frac{1}{2}$.

No me gustan mucho estos siguientes párrafos, pero no sé muy bien cómo decir lo que quiero, y siento que vos habías dicho algo más interesante sobre esto.

Ahora bien, es importante aclarar que esto no constituye de ninguna manera un contraejemplo, además de que faltaría formalidad en su presentación, resulta que plantear $s \notin C$ siendo que tengo probabilidad positiva de visitarlo cada vez no tiene mucho sentido. (Esto no sé que tanto es así).

En cualquier caso, podría ser una muestra de que la aplicación de este fórmula de prueba no es útil *directamente* para probar el teorema para estrategias de memoria infinita. Pero resaltamos el directamente porque la productoria constituye una cota a la probabilidad de que $\omega \in \Omega_s^{(C,D)}$, no es directamente el valor de la probabilidad.

Esto de acá abajo capaz se va

A su vez, el teorema anterior nos permite empezar a considerar especialmente un tipo de propiedades, aquellas que solo dependen del comportamiento asintótico.

Definición 5.1.6 (propiedad límite). *Una propiedad de tiempo lineal se llama propiedad límite de tiempo lineal si para todos los comportamientos ω, ω' $\omega \in P \wedge \inf(\omega) = \inf(\omega') \implies \omega' \in P$.*

Si T es un subconjunto de estados para un PMDP \mathcal{M} , diremos que $T \models P$ para alguna propiedad límite P sii para todo comportamiento tal que $\text{sub}(\inf(\omega)) = (T, A)$ para algún A , vale que $\omega \in P$.

Por el teorema 5.1.2, solo los conjuntos T donde T es un conjunto de estados de una componente final son relevantes a la hora de analizar las probabilidades de una propiedad límite P . Denotemos U_P a la unión de todos los conjuntos T de todas las componentes finales (T, A) de \mathcal{M} tales que $T \models P$, y V_P la unión de los conjuntos T de todos los $(T, A) \in EC(\mathcal{M})$ tales que $\neg(T \models P)$. Esto nos permite presentar el siguiente interesante resultado:

Teorema 5.1.3. *Sea \mathcal{M} un PMDP y sea P una propiedad de límite. Entonces existe una estrategia de memoria finita π tal que para todo estado s de \mathcal{M} :*

- (a) $\sup_{\pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(P) = \sup_{\pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(\Diamond U_P)$
- (b) $\inf_{\pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(P) = 1 - \sup_{\pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(\Diamond V_P)$

Demostración. Primero, probemos el item (a).

Para cada estrategia tenemos que $\mathbb{P}_{\mathcal{M},s}^{\pi}(P) \leq \mathbb{P}_{\mathcal{M},s}^{\pi}(\Diamond U_P)$, puesto que por teorema 5.1.2 vale que:

$$\mathbb{P}_{\mathcal{M},s}^\pi(P) = \mathbb{P}_{\mathcal{M},s}^\pi\{\omega \in \text{Paths}(s) \mid \text{inf}(\omega) \in EC(\mathcal{M}) \wedge \text{inf}(\omega) \models P\}$$

y por definición de U_P , $\omega \models \Diamond U_P$ para cada camino ω tal que $\text{inf}(\omega)$ es una componente final y $\text{inf}(\omega) \models P$.

Si podemos probar que existe una estrategia π tal que $\mathbb{P}_{\mathcal{M},s}^\pi(P) = \sup_\pi \mathbb{P}_{\mathcal{M},s}^\pi(\Diamond U_P)$, podemos ver que vale la igualdad, por lo expuesto anteriormente.

Consideremos, entonces, una estrategia sin memoria π_0 que maximiza las probabilidades de llegar a U_P para cada s en \mathcal{M} (esta estrategia existe por resultado X de literatura y el teorema 1 de [3]). Además, para cada componente final (C, D) , definamos $\pi_{(C,D)}$: una estrategia que asegura quedarse siempre en C mientras visita infinitamente a menudo todos los estados $s \in C$ (esta estrategia existe por el teorema 5.1). Y, por último, para cada estado $u \in U_P$ seleccionamos una componente final $EC(u) = (C, D)$ tal que $u \in C$ y $C \models P$. Todo esto nos servirá para definir la estrategia π que nos asegura que $\mathbb{P}_{\mathcal{M},s}^\pi(P) = \sup_\pi \mathbb{P}_{\mathcal{M},s}^\pi(\Diamond U_P)$.

Definimos π como la estrategia que primero se comporta como π_0 hasta que llega a un estado $u \in U_P$. Desde ese momento, π se comporta como $\pi_{EC(u)}$. En consecuencia, para esta estrategia, los caminos que eventualmente entran a U_P visitarán todos los estados de una componente final (C, D) tal que $C \models P$ infinitamente a menudo con probabilidad 1. En particular, $\text{inf}(\pi) \models P$ vale para todo ω que sigue la estrategia π y eventualmente alcanza U_P . Esto lleva a que:

$$\begin{aligned} \mathbb{P}_{\mathcal{M},s}^\pi(P) &= \sum_{u \in U_P} \mathbb{P}_{\mathcal{M},s}^{\pi_0}((\neg U_P) \mathcal{U} u) \cdot \mathbb{P}_{\mathcal{M},u}^{\pi_{EC(u)}}(P) \\ &= \sum_{u \in U_P} \mathbb{P}_{\mathcal{M},s}^{\pi_0}((\neg U_P) \mathcal{U} u) \cdot 1 \\ &= \sum_{u \in U_P} \mathbb{P}_{\mathcal{M},s}^{\pi_0}((\neg U_P) \mathcal{U} u) \\ &= \sup_{\pi'} \mathbb{P}_{\mathcal{M},s}^{\pi'}(\Diamond U_P) \end{aligned}$$

(Revisar la manera en la que están escritas estas igualdades).

Hemos probado así el ítem (a).

Ahora, probemos el ítem (b). Este deriva del ítem anterior usando el hecho de que las propiedades límite son cerradas bajo negación (es decir, si P es una propiedad límite entonces \overline{P} también lo es) y que $\mathbb{P}_{\mathcal{M},s}^\pi(P) = 1 - \mathbb{P}_{\mathcal{M},s}^\pi(\overline{P})$ para toda estrategia π . Por lo tanto:

$$\inf_{\pi} \mathbb{P}_{\mathcal{M},s}^\pi(P) = 1 - \inf_{\pi} \mathbb{P}_{\mathcal{M},s}^\pi(\overline{P}) \quad (5.3)$$

y cualquier estrategia de memoria finita que maximiza la probabilidades para \overline{P} , minimiza las probabilidades para P . Para un subconjunto de estados T , tenemos que $T \models \overline{P}$ sii $\neg(T \models P)$. Entonces, el conjunto V_P de todos los estados t que están contenidos en una componente final (T, A) con $\neg(T \models P)$ coincide con el conjunto $U_{\overline{P}}$ que surge de la unión de todos las componentes finales (T, A) donde $T \models \overline{P}$. Y con eso, el ítem (a) aplicado a \overline{P} hace que valga:

$$\sup_{\pi} \mathbb{P}_{\mathcal{M},s}^\pi(\overline{P}) = \sup_{\pi} \mathbb{P}_{\mathcal{M},s}^\pi(\Diamond U_{\overline{P}}) = \sup_{\pi} \mathbb{P}_{\mathcal{M},s}^\pi(\Diamond V_P)$$

Y estas igualdades aplicadas a 5,3 nos dan lo que queríamos probar.

□

5.2. Juegos justos: la respuesta a la pregunta cualitativa

5.2.1. Juegos de adversario justo

Sea G un juego de grafo de dos jugadores y sea $E^l \subseteq (V_1 \times V) \cap E$ un conjunto dado de aristas vivas **ver traducción**. Sea $V^l := \text{dom}(E^l)$ el conjunto de vértices del jugador \Diamond que están en el dominio de E^l . Intuitivamente, las aristas en E^l representan suposiciones de justicia sobre el jugador \Diamond : para cada arista $(v, v') \in E^l$, si v es visitado infinitamente a menudo en una jugada, se espera que la arista (v, v') sea elegida también infinitamente a menudo por el jugador \Diamond . Es decir, si un vértice v es visitado infinitamente, se espera también que toda arista viva saliente de v sea tomada infinitamente a menudo.

Denotamos por $G^l = \langle G, E^l \rangle$ a un juego de grafo de dos jugadores con aristas vivas, y extendemos nociones como jugadas, estrategias, condiciones ganadoras, regiones ganadoras, etc., de juegos de grafo a juegos de grafos con aristas vivas. Una jugada π sobre G^l es fuertemente justa si satisface la siguiente fórmula en lógica temporal lineal (LTL):

$$\alpha := \bigwedge_{(v,v') \in E^l} (\Box \Diamond v \rightarrow \Box \Diamond (v \wedge \bigcirc v')) .$$

Dado G^l y una condición de victoria φ , el jugador \Box gana el juego de adversario justo sobre G^l con respecto a la condición de victoria φ desde un vértice $v_0 \in V$ si gana el juego sobre G^l para la condición de victoria $\alpha \rightarrow \varphi$ desde v_0 .

Hay dos observaciones interesantes sobre los juegos de adversario justo.

Primero, las aristas vivas permiten descartar ciertas estrategias del jugador \Diamond , facilitando que el jugador \Box gane en determinadas situaciones. Por ejemplo, consideremos un grafo de juego (figura X, parte superior) con dos vértices p y q . El vértice p pertenece al jugador \Diamond y el vértice q al jugador \Box . La arista (p, q) es una arista viva (representada con línea discontinua). Supongamos que la especificación para el jugador \Box es $\varphi = \Box \Diamond q$. Si la arista (p, q) no fuera viva, el jugador \Box no ganaría desde p , porque el jugador \Diamond podría mantener el juego atrapado en p eligiéndose a sí mismo como sucesor en cada turno. En contraste, el jugador \Box gana desde p en el juego adversarial justo, porque la suposición de liveness sobre la arista (p, q) fuerza al jugador \Diamond a elegir infinitamente a menudo la transición hacia q .

Segundo, las suposiciones de justicia modeladas por aristas vivas restringen las elecciones de estrategias del jugador \Diamond menos que lo que restringirían la suposición de que el jugador \Diamond elige probabilísticamente entre estas aristas. Consideremos, por ejemplo, un juego de adversario justo con un único vértice del jugador \Diamond , p (cuadrado) con dos aristas vivas salientes hacia los estados q y q' , como se muestra en la Figura 1 (parte inferior). Si el jugador \Diamond elige aleatoriamente entre las aristas (p, q) y (p, q') , toda secuencia finita de visitas a los estados q y q' ocurrirá infinitamente a menudo con probabilidad uno. Esto no es cierto en el juego de adversario justo. Aquí el jugador \Diamond puede elegir una secuencia particular de visitas a q y q' (por ejemplo, simplemente $qq'qq'qq' \dots$), siempre

que ambos sean visitados infinitamente a menudo.

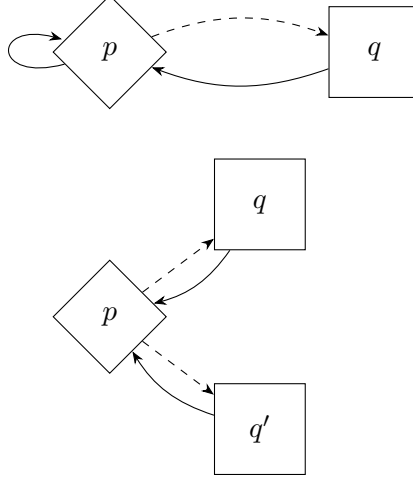


Figura 5.2: Dos juegos de adversario justo.

5.2.2. Desrandomización de PSGs

Dada la interpretación de un juego estocástico politópico $\mathcal{G}_{\mathcal{K}} = (\mathcal{S}, (\mathcal{S}_{\square}, \mathcal{S}_{\diamond}), \mathcal{A}, \theta)$, se define la *desrandomización* de $\mathcal{G}_{\mathcal{K}}$, $Derand(\mathcal{G}_{\mathcal{K}}) := ((V, V_{\square}, V_{\diamond}, E), E^l)$, como el (vamos con esta traducción?) juego de grafo de 2 jugadores extremadamente equitativo con:

$$\begin{aligned}\tilde{V} &= \bigcup_{\substack{s \in \mathcal{S} \\ V' \in V_s}} v_{V'} \\ V &= \mathcal{S} \cup \tilde{V} \\ V_{\square} &= \mathcal{S}_{\square} \\ V_{\diamond} &= \mathcal{S}_{\diamond} \cup \tilde{V} \\ E &= \{(s, v_V) : V \in V_s\} \\ E^l &= \{(v_V, s') : s' \in V\}\end{aligned}$$

y para el cual se cumple la siguiente condición de equidad:

$$\varphi^l := \bigwedge_{(v_{V'}, s') \in E^l} (\Box \Diamond v_{V'} \rightarrow \Box \Diamond (v_{V'} \wedge \bigcirc s'))$$

o lo que es lo mismo si todo camino cumple que para todo $v_V \in \tilde{V}$

$$\varphi^l := \bigwedge_{s' \in V'} (\Box \Diamond v_{V'} \rightarrow \Box \Diamond (v_{V'} \wedge \bigcirc s'))$$

Relación entre un PSG y su desrandomización

Nos será útil pensar en transformaciones entre los distintos juegos para las pruebas que tenemos más adelante, así que veremos algunas definiciones. Antes, fijemos la interpretación de un juego estocástico politópico \mathcal{G}_K y su desrandomización $Derand(\mathcal{G}_K)$.

Definición 5.2.1 (transformación de caminos). *Dado un comportamiento $\omega = (s_0, \alpha_0, s_1, \alpha_1, \dots)$ en \mathcal{G}_K , podemos obtener un único camino en $Derand(\mathcal{G}_K)$ al que llamaremos $derand(\omega)$ y tendrá la siguiente forma:*

$$derand(\omega) = (s_0, v_{\text{supp}(\alpha_0)}, s_1, v_{\text{supp}(\alpha_1)}, \dots)$$

Por otro lado, si tenemos un camino $\rho = (s_0, v_{V_0}, s_1, v_{V_1}, \dots)$ en $Derand(\mathcal{G}_K)$, existen varios comportamientos en \mathcal{G}_K que se corresponderían con él. Llamaremos a este conjunto de comportamientos $rand(\rho)$. Es decir,

$$rand(\rho) = \{(s_0, \alpha_0, s_1, \alpha_1, \dots) \in \Omega_{\mathcal{G}_K, s} \mid \forall i \geq 0, \text{supp}(\alpha_i) = V_i\}$$

Definición 5.2.2 (desrandomización de estrategias). *Dada una estrategia π_\Box en \mathcal{G}_K , vamos a definir una estrategia σ_\Box en $Derand(\mathcal{G}_K)$. π_\Box está definida para cada prefijo finito de comportamiento ω , debemos definir σ_\Box para cada prefijo finito $derand(\omega)$. Sin embargo, para cada prefijo finito $\rho = (s_0, v_{V_0}, s_1, v_{V_1}, \dots)$ en $Derand(\mathcal{G}_K)$, existe más de un $\hat{\omega}$ tal que $derand(\hat{\omega}) = \rho$. Por eso, lo primero que debemos hacer es elegir un camino ω particular tal que $derand(\omega) = \rho$ y definir $\sigma_\Box(derand(\omega))$ a partir de ese ω .*

Sabemos que $\pi_\Box(\omega s)$ es una función que asigna a cada acción $\hat{\alpha} \in \mathcal{A}(s)$ una probabilidad de ser seleccionada como próxima. Para cada $derand(\omega)s$, σ_\Box debemos asignar un próximo vértice. Lo que haremos para definir que vértice será seleccionado será elegir alguna acción $\hat{\alpha}^*$ a la que $\pi_\Box(\omega s)$ le haya asignado una probabilidad positiva, y entonces diremos que $\sigma_i(derand(\omega)s) = v_{\text{supp}(\hat{\alpha}^*)}$.

A esta nueva estrategia σ_\Box que hemos obtenido a partir de π_\Box la llamaremos $derand(\pi_\Box)$.

¿Capaz estaría bueno definirlo como tipo algoritmo? Y capaz dar alguna noción de corrección viendo que siempre existen las cosas que elegimos?

Definición 5.2.3 (randomización de estrategias). *Ahora supongamos que tenemos una estrategia σ_i en $\text{Derand}(\mathcal{G}_K)$ y queremos construir una estrategia π_i en \mathcal{G}_K que tenga un comportamiento similar. σ_i está definida para todos los prefijos finitos de caminos ρ en $\text{Derand}(\mathcal{G}_K)$ y nuestra idea es definir π_i para todos los comportamiento ω en \mathcal{G}_K . Nuestra idea será definir π_i de igual manera para todos los elementos del conjunto $\text{rand}(\rho)$.*

Solo nos va a interesar ver cómo σ_i se comporta en los prefijos finitos de caminos ρ que terminen en un estado $s \in \mathcal{S}$ ¹. Para cada uno de estos ρ , tenemos que $\sigma_i(\rho) = v_{\hat{V}}$ para algún v_V . Lo que haremos para la construcción de π_i es para cada $\omega \in \text{derand}(\rho)$ definir $\pi_i(\omega)(\hat{\alpha}) = 1$ para algún $\hat{\alpha}$ particular tal que $\text{supp}(\hat{\alpha}) = \hat{V}$.

Llamaremos a esta nueva estrategia π_i obtenida a partir de σ_i , $\text{rand}(\sigma_i)$.

5.2.3. Prueba de igualdad sobre los conjuntos ganadores

Teorema 5.2.1. *Sea $\mathcal{G}_K = (\mathcal{S}, (\mathcal{S}_{\square}, \mathcal{S}_{\diamond}), \mathcal{A}, \theta)$ la interpretación de un juego estocástico politópico, $\tilde{\mathcal{R}} = \{\langle G_1, R_1 \rangle, \dots, \langle G_k, R_k \rangle\}$ una condición de Rabin sobre \mathcal{S} , con su especificación LTL φ ,*

$$\varphi := \bigvee_{j \in [1, k]} \left(\diamond \square \overline{R_j} \wedge \diamond \square G_j \right)$$

y sea $\text{Derand}(\mathcal{G}_K)$ su desrandomización.

Sea $\mathcal{W} \subseteq \mathcal{S}$ el conjunto de todos los estados desde los cuales el jugador \square gana en $\text{Derand}(\mathcal{G}_K)$ y sea \mathcal{W}^{as} el conjunto de vértices desde los cuales el jugador \square gana casi seguramente con estrategias de memoria finita. Entonces, $\mathcal{W} = \mathcal{W}^{\text{as}}$.

Es más, a partir de una estrategia ganadora en $\text{Derand}(\mathcal{G}_K)$ se puede construir fácilmente una estrategia ganadora en \mathcal{G}_K , y viceversa.

¹Desde los prefijos finitos de caminos donde el último elemento es un vértice v_V , las decisiones son tomadas por el jugador \diamond y solo reflejan lo que sería la decisión probabilística en \mathcal{G}_K .

Demostración. Probaremos la doble contención:

Primera contención: $\mathcal{W} \subseteq \mathcal{W}^{as}$

Sea $s \in \mathcal{W}$. Entonces, sabemos que existe al menos una estrategia del jugador \square ganadora desde s en $Derand(\mathcal{G}_{\mathcal{K}})$. Llamemos a esta σ_{\square}^* .

Queremos ver que $s \in \mathcal{W}^{as}$, lo que requiere ver que existe una estrategia $\hat{\pi}_{\square}^*$ tal que:

$$\inf_{\pi_{\diamond} \in \Pi_{\diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}}(\varphi) = 1 \quad (5.4)$$

Proponemos $\hat{\pi}_{\square}^* = \text{rand}(\sigma_{\square}^*)$.

Ahora bien, ¿cómo podemos probar 5.4?

Por definición de σ_{\square}^* sabemos que si se está frente a una estrategia σ_{\diamond} que respeta la condición de equidad φ^l , entonces la jugada determinada por ambas estrategias cumple la especificación φ .

A esto de acá abajo es que le faltaría una justificación más precisa.

Pedir que valga φ^l en $Derand(\mathcal{G}_{\mathcal{K}})$ es pedir que valga $\hat{\varphi}^l$ en $\mathcal{G}_{\mathcal{K}}$ para todo comportamiento en el juego y acción α en el comportamiento, vale que:

$$\hat{\varphi}^l = \bigwedge_{s' \in \text{supp}(\alpha)} (\square \diamond \alpha \rightarrow \square \diamond (\alpha \wedge \bigcirc s'))$$

Capaz es notable el que por cada soporte V , por cómo está definido rand de estrategias, habrá una única acción α que se elije (por eso se sabe la equivalencia? por eso se sabe que habrá infinitas elecciones de α ?).

Entonces, por cómo $\hat{\pi}_{\square}^*$ está definida (a partir de σ_{\square}^*) sabemos que

$$\inf_{\pi_{\diamond} \in \Pi_{\diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}}(\hat{\varphi}^l \rightarrow \varphi) = 1 \quad (5.5)$$

Veamos ahora que la ecuación 5.4 valdrá porque la probabilidad de que no valga $\hat{\varphi}^l$ es 0. Veamos:

Esto no me está quedando bien con el tema de que esto es para cada α que aparece en el comportamiento, pero medio que entendiendo que es algo del camino capaz está bien y ni vale la pena hacerse mucha cabeza

$$\begin{aligned}
 & \mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\hat{\pi}_{\square}^*,\pi_{\diamond}}(\neg\hat{\varphi}^l) = \\
 & \mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\hat{\pi}_{\square}^*,\pi_{\diamond}} \left(\neg \left(\bigwedge_{s' \in \text{supp}(\alpha)} (\square\Diamond\alpha \rightarrow \square\Diamond(\alpha \wedge \bigcirc s')) \right) \right) = \\
 & \mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\hat{\pi}_{\square}^*,\pi_{\diamond}} \left(\bigvee_{s' \in \text{supp}(\alpha)} \neg (\square\Diamond\alpha \rightarrow \square\Diamond(\alpha \wedge \bigcirc s')) \right) = \\
 & \mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\hat{\pi}_{\square}^*,\pi_{\diamond}} \left(\bigvee_{s' \in \text{supp}(\alpha)} (\square\Diamond\alpha \wedge \Diamond\square\neg(\alpha \wedge \bigcirc s')) \right) \leq \\
 & \sum_{s' \in \text{supp}(\alpha)} \mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\hat{\pi}_{\square}^*,\pi_{\diamond}} (\square\Diamond\alpha \wedge \Diamond\square\neg(\alpha \wedge \bigcirc s'))
 \end{aligned}$$

Mostramos que este último término es igual a 0, viendo que para cada $s' \in \text{supp}(\alpha)$,

$$\mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\hat{\pi}_{\square}^*,\pi_{\diamond}} (\square\Diamond\alpha \wedge \Diamond\square\neg(\alpha \wedge \bigcirc s')) = 0$$

Sea ω una jugada. Notemos con μ a la distribución asociada a la acción α . Sea I el conjunto infinito de momentos donde se elige la acción α . La probabilidad de elegir s' como próximo vértice en cualquier momento $i \in I$ estará dada por $\mu(s')$.

Entonces, para cada momento i , la probabilidad de no visitar s' por los próximos k momentos i , está acotada superiormente por $(1 - \mu(s'))^k$, que converge a 0 cuando k tiende de a ∞ . Por lo tanto, tenemos que $\mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\hat{\pi}_{\square}^*,\pi_{\diamond}} (\square\Diamond\alpha \wedge \Diamond\square\neg(\alpha \wedge \bigcirc s')) = 0$.

Queda raro esto porque no hay un estado s de por sí.

Consideremos un $s' \in \text{supp}(\alpha)$ arbitrario. Sea ω una jugada y sea I el conjunto infinito de momentos donde se llega a s y se elige una acción $(K, \mu) \in \text{acc}(s, \text{supp}(\alpha))$. Llamemos μ_i a la distribución elegida en el momento $i \in I$. La probabilidad de elegir s' como próximo vértice en el momento $i \in I$ estará dada por $\mu_i(s')$.

Como estamos trabajando con estrategias de memoria finita, habrá solo una cantidad finita de distribuciones que π_\diamond puede elegir desde s , así que existirá una distribución, llamémosla μ_{\max} tal que da la probabilidad máxima de ir a s' .

No tengo porqué hablar de acciones de memoria finita mepa. Voy a elegir infinitas veces la misma acción. Solo elijo una acción por cada soporte por la manera en la que definimos rand.

Entonces, para cada momento i , la probabilidad de no visitar s' por los próximos k momentos i , está acotada superiormente por $(1 - \mu_{\max}(s'))^l$, que converge a 0 cuando l tiende de a ∞ . Por lo tanto, tenemos que $\mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\square \diamond \alpha \wedge \diamond \square \neg(\alpha \wedge \bigcirc s')) = 0$

En consecuencia, se sigue que $\sum_{s' \in \text{supp}(\alpha)} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\square \diamond \alpha \wedge \diamond \square \neg(\alpha \wedge \bigcirc s')) = 0$, lo que a su vez establece que $\mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\neg \hat{\varphi}^l) = 0$.

Luego vale que:

$$\begin{aligned} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\hat{\varphi}^l \rightarrow \varphi) &= \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\neg \hat{\varphi}^l \vee \varphi) \leq \\ \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\neg \hat{\varphi}^l) &+ \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\varphi) = 0 + \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\varphi) = \\ \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\varphi) \end{aligned}$$

Y por 5.5 entonces vale que $\inf_{\pi_{\diamond} \in \Pi_{\diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\hat{\pi}_{\square}^*, \pi_{\diamond}} (\varphi) = 1$ como queríamos probar.

Segunda contención: $\mathcal{W}^{as} \subseteq \mathcal{W}$

Sea $s \in \mathcal{W}^{as}$, veamos que $s \in \mathcal{W}$.

Como $s \in \mathcal{W}^{as}$, entonces existe π_{\square}^* en $\mathcal{G}_{\mathcal{K}}$ tal que $\inf_{\pi_{\diamond} \in \Pi_{\diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}^*, \pi_{\diamond}} (\varphi) = 1$.

Lo que queremos ver para probar que $s \in \mathcal{W}$ es que existe una estrategia σ_{\square}^* tal que

\square gana con ella desde s en $Derand(\mathcal{G}_K)$.

Definimos $\sigma_{\square}^* = \text{derand}(\pi_{\square}^*)$

Entonces ahora queremos ver que σ_{\square}^* es ganadora en $Derand(\mathcal{G}_K)$ desde s .

Una manera de ver esto es probar que para un comportamiento cualquiera ρ generado por σ_{\square}^* y una estrategia válida $\pi_{\diamond} \in \Pi_{\diamond}$ arbitraria, $\text{inf}(\rho)$ cumple con la especificación φ . Es decir, que vale la siguiente fórmula:

$$\varphi' := \bigvee_{j=1}^k ((\text{inf}(\rho) \cap R_j = \emptyset) \wedge (\text{inf}(\rho) \cap G_j \neq \emptyset)) \quad (5.6)$$

Si podemos probar que $\text{inf}(\rho)$ es un sub-PMDP alcanzable en el PMDP que se forma al fijar la estrategia de memoria finita π_{\square}^* en \mathcal{G}_K , por el teorema 5.1.2.1 podemos ver que vale la ecuación 5.6.

Primero, podemos ver que $\text{inf}(\rho)$ (entendiendo los vértices de la forma v_V como el conjunto resultado V) es una componente final en el polytopal markov decision process que se forma al fijar la estrategia π_{\square}^* en \mathcal{G}_K , llamémoslo (C, D) , viendo que cumple las dos propiedades:

- Para cada V^i que aparece como subíndice de vértices especiales en $\text{inf}(\rho)$, vale que $V^i \subseteq C$. En el caso de que existiese algún $s_x \in V^i$ tal que $s_x \notin C$, se estaría contradiciendo que π_{\diamond} sea una estrategia válida en $Derand(\mathcal{G}_K)$, puesto que visitaría infinitas veces v_{V^i} y solo finitas veces s_x , uno de sus sucesores.
- El grafo dirigido inducido por $\text{inf}(\rho)$ es fuertemente conexo. Si fuese de otra manera habría dos vértices $u, v \in \text{inf}(\rho)$ tales que v no sería alcanzable desde u , contradiciendo así que u y v son visitados infinitas veces por σ .

Luego, podemos ver que $\text{inf}(\rho)$ es alcanzable en \mathcal{G}_K viendo que existe una estrategia π_{\diamond}^* en \mathcal{G}_K que lo posibilita.

Definamos $\pi_{\diamond}^* = \text{rand}(\sigma_{\diamond})$.

Esta estrategia permite llegar con probabilidad positiva a $\inf(\rho)$, puesto que tanto π_{\diamond}^* como π_{\square}^* le da probabilidad positiva a los mismos vértices que π_{\diamond} asegura visitar infinitamente. **Capaz esto se podría explicar mejor.**

Con lo cual hemos probado que vale 5.6 y, por lo tanto, que σ_{\square}^* es ganadora en $Derand(\mathcal{G}_{\mathcal{K}})$ desde s , con lo que hemos probado que $s \in \mathcal{W}$.

□

¿Cómo represento gráficamente la idea de politopos? ¿Capaz dándoles forma a las acciones?

5.3. Transformar Rabin en alcanzabilidad: la respuesta a varias preguntas

5.3.1. Draft

La idea es intentar hacer algo al estilo

Teorema 5.3.1 (Reducción de Rabin). *Sea $\mathcal{G}_{\mathcal{K}}$ un juego estocástico politópico y sea $\mathcal{H}_{\mathcal{K}}$ su interpretación extrema tal como se presenta en [3]. Si tenemos una propiedad de Rabin R , y su respectivo conjunto de estados almost sure winning W_R entonces valen las siguientes ecuaciones:*

$$\begin{aligned}
 & \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(R) = \\
 & = \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R) = \\
 & = \inf_{\pi_{\diamond} \in \Pi_{\mathcal{H}_{\mathcal{K}}, \diamond}^{MD}} \sup_{\pi_{\square} \in \Pi_{\mathcal{H}_{\mathcal{K}}, \square}^{MD}} \mathbb{P}_{\mathcal{H}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R) = \\
 & = \sup_{\pi_{\square} \in \Pi_{\mathcal{H}_{\mathcal{K}}, \square}^{MD}} \inf_{\pi_{\diamond} \in \Pi_{\mathcal{H}_{\mathcal{K}}, \diamond}^{MD}} \mathbb{P}_{\mathcal{H}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R) = \\
 & = \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R) = \\
 & = \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(R)
 \end{aligned}$$

Demostración.

$$\begin{aligned}
& \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(R) \\
& \leq \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R) && \text{(por lema 5.3.2)} \\
& \leq \inf_{\pi_{\diamond} \in \Pi_{\mathcal{H}_{\mathcal{K}}, \diamond}^{MD}} \sup_{\pi_{\square} \in \Pi_{\mathcal{H}_{\mathcal{K}}, \square}^{MD}} \mathbb{P}_{\mathcal{H}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R) && \text{(por Teorema 1)} \\
& \leq \sup_{\pi_{\square} \in \Pi_{\mathcal{H}_{\mathcal{K}}, \square}^{MD}} \inf_{\pi_{\diamond} \in \Pi_{\mathcal{H}_{\mathcal{K}}, \diamond}^{MD}} \mathbb{P}_{\mathcal{H}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R) && \text{(por Teorema 1)} \\
& \leq \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R) && \text{(por Teorema 1)} \\
& \leq \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(R) && \text{(por lema 5.3.4)} \\
& \leq \inf_{\pi_{\diamond} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \diamond}} \sup_{\pi_{\square} \in \Pi_{\mathcal{G}_{\mathcal{K}}, \square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(R) && \text{(por prop de sup e inf)}
\end{aligned}$$

□

Lema 5.3.2. *Sea $\mathcal{G}_{\mathcal{K}}$ la interpretación (?) de un juego estocástico politópico y sea R una propiedad límite y W_R su correspondiente conjunto casi seguramente ganador. Entonces vale que*

$$\inf_{\pi_{\diamond} \in \Pi_{\diamond}} \sup_{\pi_{\square} \in \Pi_{\square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(R) \leq \inf_{\pi_{\diamond} \in \Pi_{\diamond}} \sup_{\pi_{\square} \in \Pi_{\square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_R)$$

Demostración. Nombremos π_{\diamond}^* a una estrategia tal que

$$\sup_{\pi_{\square} \in \Pi_{\square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}^*}(\diamond W_P) = \inf_{\pi_{\diamond} \in \Pi_{\diamond}} \sup_{\pi_{\square} \in \Pi_{\square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_P)$$

(explicar algo más hablado? decir algo sobre su existencia?). Ahora podemos hacer las siguientes deducciones:

$$\begin{aligned}
& \inf_{\pi_{\diamond} \in \Pi_{\diamond}} \sup_{\pi_{\square} \in \Pi_{\square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(P) \\
& \leq \sup_{\pi_{\square} \in \Pi_{\square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}^*}(P) && \text{(por def. de inf)} \\
& = \sup_{\pi_{\square} \in \Pi_{\square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}^*}(\diamond W_P) && \text{(por lema 5.3.3)} \\
& = \inf_{\pi_{\diamond} \in \Pi_{\diamond}} \sup_{\pi_{\square} \in \Pi_{\square}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}}, s}^{\pi_{\square}, \pi_{\diamond}}(\diamond W_P) && \text{(por def. de } \pi_{\diamond}^*)
\end{aligned}$$

y con esto hemos probado el enunciado. \square

Lema 5.3.3. *Sea \mathcal{M} un PMDP, π una estrategia en \mathcal{M} , s un estado en \mathcal{M} , R una propiedad de Rabin y W_R el conjunto casi seguramente ganador de R (en el juego estocástico debería ser, no?), entonces vale que:*

¿Cambiar w_r por conjunto desde el que existe una estrategia que da probabilidad 1 de ganar?

$$\sup_{\pi \in \Pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(R) = \sup_{\pi \in \Pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(\Diamond W_R)$$

Demostración. Por definición de conjunto casi seguramente ganador

$$\mathbb{P}_{\mathcal{M},s}^{\pi}(R) = \mathbb{P}_{\mathcal{M},s}(\{\omega \in Paths(s) \mid \text{inf}(\omega) \models R\})$$

Claramente los caminos que hacen que valga la propiedad también son caminos desde donde llego a algún estado desde donde existen estrategias para ganar con probabilidad 1. Por lo tanto, $\mathbb{P}_{\mathcal{M},s}^{\pi}(R) \leq \mathbb{P}_{\mathcal{M},s}^{\pi}(\Diamond W_R)$.

Para ver que vale la igualdad, veremos que existe una estrategia de memoria finita π que hace que $\mathbb{P}_{\mathcal{M},s}^{\pi}(R) = \sup_{\pi \in \Pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(\Diamond W_R)$. Para ello, consideremos la estrategia sin memoria π_0 que maximiza las probabilidades de alcanzar W_R desde todos los estados $s \in \mathcal{M}$ (sabemos que esta estrategia existe por [3, 4]). A su vez, sabemos que para cada estado $s \in W_R$ existe una estrategia π_s que asegura ganar con probabilidad 1 frente al objetivo R .

Sea entonces π la estrategia que primero se comporta como π_0 , hasta llegar a un estado t en W_R , y a partir de allí π se comporta como π_t . Con eso tenemos que:

$$\begin{aligned} \mathbb{P}_{\mathcal{M},s}^{\pi}(R) &= \sum_{t \in W_R} \mathbb{P}_{\mathcal{M},s}^{\pi_0}((\neg W_R \mathcal{U} t)) \cdot \underbrace{\mathbb{P}_{\mathcal{M},t}^{\pi_t}(R)}_{=1} \\ &= \sup_{\pi \in \Pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(W_R) \end{aligned}$$

Como $\sup_{\pi \in \Pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(W_R)$ es una cota superior las probabilidades para R bajo todas las estrategias, con esto podemos concluir la igualdad $\sup_{\pi \in \Pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(R) = \sup_{\pi \in \Pi} \mathbb{P}_{\mathcal{M},s}^{\pi}(\Diamond W_R)$. \square

Esto es lo importante que habría que probar y no está probado.

Lema 5.3.4. *Sea $\mathcal{G}_{\mathcal{K}}$ la interpretación (?) de un juego estocástico politópico y sea R una propiedad límite y W_R su correspondiente conjunto casi seguramente ganador. Entonces vale que*

$$\sup_{\pi_{\square} \in \Pi_{\square}} \inf_{\pi_{\diamond} \in \Pi_{\diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\pi_{\square},\pi_{\diamond}}(\Diamond W_R) \leq \sup_{\pi_{\square} \in \Pi_{\square}} \inf_{\pi_{\diamond} \in \Pi_{\diamond}} \mathbb{P}_{\mathcal{G}_{\mathcal{K}},s}^{\pi_{\square},\pi_{\diamond}}(R)$$

Capítulo 6

Conclusiones

6.1. Trabajo Futuro

Referencias

- [1] C. Baier y J.-P. Katoen. *Principles of Model Checking*. Vol. 26202649. The MIT Press, 2008. ISBN: 978-0-262-02649-9.
- [2] L. de Alfaro. *Formal Verification of Probabilistic Systems*. Inf. téc. Stanford, CA, USA, 1998.
- [3] P. F. Castro y P. R. D’Argenio. «Polytopal Stochastic Games». En: (2024). Enviado para publicación.
- [4] A. Condon. «The complexity of stochastic games». En: *Information and Computation* 96.2 (1992), págs. 203-224. ISSN: 0890-5401. DOI: [https://doi.org/10.1016/0890-5401\(92\)90048-K](https://doi.org/10.1016/0890-5401(92)90048-K). URL: <https://www.sciencedirect.com/science/article/pii/089054019290048K>.

Apéndice A

Titulo del Apendice

A.1. Titulo de la seccion