



بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



# یادگیری و تعادل

کیانا هوشانفر  
مریم هارونی

رشته برق - کنترل  
تیرماه ۱۴۰۳

## تمرکز مقاله:

- مرور ادبیات اخیر در زمینه یادگیری غیرتعادلی در بازی‌ها.
- محدودیت‌های موجود به دلیل محدودیت فضا.
- تأکید بر مدل‌های یادگیری فردی.
- تمرکز کمتر بر مدل‌های تکاملی و مدل‌های تنظیم کوتاه‌مدت.



# مقدمه

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

## عقلانیت و تعادل

- عقلانیت  $\leftarrow$  تعادل نش
- آگاهی مشترک از عقلانیت بین بازیکنان  $\leftarrow$  تعادل نش
- هماهنگی در یک تعادل واحد

## یادگیری از تجربه

- بازی مکرر  $\leftarrow$  نتایجی نزدیک به پیش بینی های تعادلی
- ایده ایجاد تعادل از طریق یادگیری

## نظریه یادگیری در بازی ها

- فرآیند ایجاد تعادل از طریق یادگیری، تطبیق و/یا تقلید بلندمدت غیرتعادلی را رسمی می کند.



# مقدمه

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

## اهداف عاملان

تبادل نش  $\times$

حداکثر کردن سود خود  $\checkmark$

## سوال اصلی

چه زمانی یادگیری به رفتار تعادلی منجر می شود؟

## قابلیت قبول قوانین یادگیری

- نبود جایگزین های ساده و واضح بهتر.
- شرایط واقع گرایانه برای یادگیری.



## مقدمه

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

### مدل های خاص مورد بحث

#### تنظیمات ساده یادگیری

مشاهده استراتژی های عاملان در انتهای هر دور  
مواجه تصادفی عاملان با حریفان ناشناس

#### یادگیری در بازی های فرم گسترده

- ✓ فقط توالی اقدامات
- ✗ مشاهده استراتژی های حریفان
- بررسی محدودیت های یادگیری

# مقدمه

## بخش 2A بازی خیالی و بازی خیالی تصادفی

مدل‌ها نسبتاً ساده + تفکر عوامل بیزی + تفکر محیطی ثابت  
محیط ثابت یا تقریباً ثابت ← عملکرد خوب  
سادگی این مدل‌ها ← تحلیل با تکنیک‌های تقریب تصادفی

## بخش 2B عملکرد نامتقارن

همسانی هانان ← عملکرد خوب در محیط‌های ثابت  
شرایط قوی‌تر ← تضمین عملکرد خوب در محیط‌های عمومی‌تر  
تحت کالیبراسیون ← همگرایی بازی به مجموعه‌ای از تعادل‌های همبسته مدل‌های

## بخش 2C یادگیری پیچیده‌تر

مدل‌هایی که بازیکنان به‌طور فرضی ماتریس پرداخت را نمی‌دانند، شامل مدل‌های یادگیری تقویتی و مدل‌های تقلید.  
تفسیر بازی خیالی تصادفی به عنوان یادگیری تقویتی.

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# بازی خیالی

عاملان به طور مکرر یک بازی استراتژیک ثابت را بازی می کنند.

## فرضیات

عاملان فضای استراتژی و تابع سود خود را می دانند.  
عاملان استراتژی حریف خود را در هر دوره مشاهده می کنند.  
عاملان فرض می کنند که توزیع استراتژی های حریفان ثابت اما ناشناخته است.

## نزدیک بینی استراتژیک

فودنبرگ و کریس (۱۹۹۳) ← مدل "جمعیت بزرگ"

- هر دوره همه عاملان برای بازی روبهرو می شوند و فقط بازی مقابل خود را می بینند.
  - احتمال مواجهه عوامل با همان حریف در دوره های مکرر کم است.
- این تنظیم به عوامل اجازه می دهد تا بیشتر بر بازده های فوری تمرکز کنند تا بر تعاملات آینده با حریفان خاص.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری



# بازی خیالی

FP ← رویکرد بیزی به روز رسانی باورهای خود درباره استراتژی های حریفان

$$k_i^t(s^{-i}) = k_i^t(s^{-i}) + \begin{cases} 1 & \text{if } s^{-i} = s_{t-1}^{-i} \\ 0 & \text{if } s^{-i} \neq s_{t-1}^{-i} \end{cases}$$

$$\gamma_i^t(s^{-i}) = \frac{k_i^t(s^{-i})}{\sum_{\tilde{s}^{-i} \in S^{-i}} k_i^t(\tilde{s}^{-i})}$$

براساس این باورهای به روز شده بهترین اقدامات را برای به حداکثر رساندن بازده های مورد انتظار

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# بازی خیالی

مثالی از فودنبرگ و کریس برای نشان دادن محدودیت های احتمالی FP

i. وزن اولیه  $(1, \sqrt{2})$

ii. تخمین هردو از بازی رقیب B

iii. هردو A بازی میکنند

iv. وزن به روز شده  $(2, \sqrt{2})$

v. هردو B بازی میکنند

vi. ....

	A	B
A	0,0	1,1
B	1,1	0,0

همگرایی فرکانس های تجربی به تعادل نش مختلط ← بازده صفر به دلیل وجود چرخه



# بازی خیالی تصادفی

FP

بازیکنان به طور مکرر بازی را انجام می دهند و استراتژی های خود را بر اساس رفتار مشاهده شده حریفان تنظیم می کنند.

SFP

بازیکنان اقدامات را بر اساس یک تابع بهترین پاسخ تصادفی انتخاب می کنند.

تصادفی بودن : نوسانات غیرقابل پیش بینی در پاداش ها را منعکس می کند.

تصادفی بودن در تصمیم گیری ← SFP مدلی واقع بینانه تر از رفتار استراتژیک در بازی ها  
← همگرایی نرم تر به تعادل ها و ارائه سازگاری در پرداخت های بلندمدت

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# بازی خیالی تصادفی

توزیع بهترین پاسخ

هر بازیکن  $i$  شوک های تصادفی  $\eta^i$

$$\overline{BR}^i(\sigma^{-i})(s^i) = Prob[\eta^i \text{ s.t. } s^i \text{ is a best response to } \sigma^{-i}]$$

یک مثال خاص از بهترین پاسخ تصادفی = بهترین پاسخ تصادفی لجستیک

$$\overline{BR}^i(\sigma^{-i})(s^i) = \frac{\exp(\beta u(s^i, \sigma^{-i}))}{\sum_{\hat{s}^i} \exp(\beta u(\hat{s}^i, \sigma^{-i}))}$$

احتمال انتخاب  $s^i$

تقاطع توابع رو توزیع نش نامیدند.

با افزایش  $\beta$  به سمت بی نهایت، توزیع های نش به تعادل نش بازی با اطلاعات کامل همگرا می شوند.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# بازی خیالی تصادفی

## مزایای SFP

- همگرایی بهتر
- مقاومت در برابر نویز
- سازگاری جهانی

## تابع بهترین پاسخ صاف

یک بهترین پاسخ صاف عمومی می تواند از حداکثرسازی یک تابع پرداخت تعیینی مختل شده که اقدامات خالص را جریمه می کند باشد:

$$\arg \max_{\sigma^i} u^i(\sigma^i, \sigma^{-i}) + \beta^{-1} v^i(\sigma^i)$$

$$v^i(\sigma^i) = \sum_{s^i} -\sigma^i(s^i) \log \sigma^i(s^i)$$

مقدمه

بخش ۲

بخش ۳

نتیجه گیری



# بازی خیالی تصادفی

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

بررسی سیستم‌هایی از عوامل که همگی از بازی ساختگی تصادفی یا SFP استفاده می‌کنند.

استفاده از روش‌های "تقریب تصادفی" اعمال می‌شود.

نسخه Benaim برای SFP :

$$x_{n+1} - x_n = \frac{1}{n+1} (F(x_n) + U_n + b_n)$$
$$F(x(t)) = dx(t)/dt$$

تقریب تصادفی برای بازی دو نفره SFP :

$$\theta_{1,n+1} - \theta_{1,n} = (1/(n+1)) [\bar{B}\bar{R}_1(\theta_{2,n}) - \theta_{1,n} + U_{1,n} + b_{1,n}]$$
$$\theta_{2,n+1} - \theta_{2,n} = (1/(n+1)) [\bar{B}\bar{R}_2(\theta_{1,n}) - \theta_{2,n} + U_{2,n} + b_{2,n}]'$$

# بازی خیالی تصادفی

آن‌ها سپس از تقریب تصادفی برای ارتباط دادن رفتار حدی سیستم با رفتار سیستم قطعی زیر استفاده کردند:

$$\dot{\theta}_1 = BR_1(\theta_2) - \theta_1, \dot{\theta}_2 = BR_2(\theta_1) - \theta_2$$

توجه داشته باشید که نقاط تعادل این سیستم دقیقاً توزیع‌های تعادلی هستند. بنابراین تقریب تصادفی به طور تقریبی می‌گوید که SFP نمی‌تواند به یک توزیع نش ناپایدار خطی همگرا شود و باید به یکی از مجموعه‌های داخلی زنجیروار متعدی سیستم همگرا شود.

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# بازی خیالی تصادفی

Benaim

تکاملی از FP تا "FP تصادفی وزن دار" (نسبت هندسی وزن کمتری به مشاهدات قدیمی تر)

به طور خلاصه، FP صاف وزن دار با وزن هایی که به ۱ همگرا می شوند، مسیرها و مجموعه های حدی مشابه FP تصادفی دارد؛ اما تفاوت در سرعت حرکت است و از این روی در نحوه همگرایی توزیع تجربی متفاوت است.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری



# عملکرد مجانبی و همگرایی جهانی

مقدمه

چگونه SFP می‌تواند در محیط‌های قابل تعویض عملکرد خوبی داشته باشد — نیاز به قوانین یادگیری پیچیده‌تر — توسعه قوانین یادگیری بهتر و کارآمدتر

بخش ۲

## ضمانت

SFP تضمین می‌کند که بازیکنان حداقل به اندازه بیشینه سازی در برابر میانگین طول زمان بازی خوب عمل می‌کنند. (به شرط محیط ثابت)

بخش ۳

## محدودیت

SFP دارای محدودیت در بازی‌های دارای الگوها، روند یا چرخه

نتیجه‌گیری

# عملکرد مجانبی و همگرایی جهانی

## یادگیری بیزی

به روزرسانی باورهای بازیکنان بر اساس اطلاعات جدید

## محدودیت

رویدادهایی با احتمال پیشین صفر

## مسئله تعامل

تعامل پیچیده ← استراتژی هایی با پیش فرض غیرممکن ← شرایط بهینه غیر بیزی

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# سازگاری جهانی

## سازگاری جهانی

یک قانون یادگیری جهانی سازگار است اگر در طول زمان به اندازه ای که انگار فرکانس های اقدامات حریفان را از قبل می دانست، عمل کند.

**تئوری:** هانان و بلکول وجود چنین قوانینی را اثبات کردند. این قوانین تضمین می کنند که میانگین بازده طولانی مدت حداقل به اندازه مقدار مینیمکس است.

**مثال:** بازی تطابق سکه ها

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# سازگاری جهانی

## قابل دستیابی

### مفهوم

بلکول ایده قابل دستیابی را معرفی کرد که شامل ایجاد قوانینی است که دستیابی به مجموعه هدف را در طول زمان تضمین می کنند.

### کاربرد

ساخت قوانین یادگیری جهانی ← اطمینان از رسیدن استراتژی های بازیکنان به نتایج مطلوب در طول بازی های تکراری

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

## Probst و Lambson

دو بازیکن + طول الگوهای مساوی + همگرایی بازی دقیق یا FP  
← همگرایی توزیع تجمعی تجربی بازی به پوش محدب مجموعه تعادل های نش

ما انتظار داریم که تشخیص الگوهای طولانی تر یک مزیت باشد!

یک مثال جالب!

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

## Foster & Vohra

قواعد یادگیری مشتق شده از پیش‌بینی‌های کالبره‌شده

پیش‌بینی به خوبی کالبره شده یعنی چه؟!

کالبراسیون به نظر می‌رسد که ویژگی مطلوبی برای یک پیش‌بین باشد.

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# کالیبراسیون

سوال مربوط به کالیبراسیون : در تمام مواقعی که یک بازیکن اقدام خاصی را انجام داد، پاسخ آن چقدر خوب بود؟

هارت و ماس-کول (۲۰۰۰)

سنجش پشیمانی

اگر بدون توجه به بازی حریفان، بازیکن به طور نامحدود کالیبره شده باشد به این معنا که میانگین زمانی پشیمانی برای هر اقدام به صفر برسد، می‌گوییم که بازیکن به طور جهانی کالیبره شده است.

فاستر و ووهر (۱۹۹۷)

وجود روش‌های یادگیری دارای این ویژگی ← اگر همه بازیکنان از چنین قواعدی پیروی کنند، میانگین زمانی فرکانس بازی باید به مجموعه‌ای از تعادل‌های همبسته بازی همگرا شود.

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

الگوریتم اولیه فاستر و ووهره شامل روشی پیچیده برای یافتن ماتریس‌های تصادفی و بردارهای ویژه سوال:

- آیا فرض اینکه بازیکنان از قواعد جهانی کالیبره شده پیروی می‌کنند، تقریب خوبی است یا نه؟!
- از نظر مفهومی تشخیص اینکه آیا قواعد یادگیری در مسیر بازی کالیبره شده‌اند آسان است؟! اگر چنین باشد، میانگین زمانی بازی مشترک به مجموعه‌ای از تعادل‌های همبسته همگرا می‌شود.



# کالیبراسیون

تحقیقات بعدی

- گسترش قابل توجه مجموعه قواعدی به طور جهانی کالیبره شده
- ساده سازی بسیار الگوریتم‌ها و روش‌های

قواعد یادگیری به طور جهانی سازگار ← ساختن قواعد یادگیری به طور جهانی کالیبره شده ← حل یک مسئله نقطه ثابت

این مسئله نقطه ثابت یک مسئله خطی است که با وارونه‌سازی ماتریس حل می‌شود، همانطور که در فودنبرگ و لوین (۱۹۹۸) نشان داده شده است.

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# کالبراسیون

آیا وارونه‌سازی یک ماتریس ساده است؟!

تفاوت ظریف بین کالبره بودن و کالبره بودن جهانی!

Hart و Mas-Colell (2001) الگوریتم‌های ساده ← یادگیری کالبره

حل مسئله نقطه ثابت برای کالبراسیون جهانی

$$Rq = R^T q$$

$$q(b) = \begin{cases} \left(\frac{1}{\mu}\right) R(a, b) & \text{for } b \neq a \\ 1 - \sum_{c \neq a} q(c) & \text{for } b = a \end{cases}$$

$\mu$  یک عدد بزرگ

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

Hart & Mas-Colell (2000)

که اگر همه افراد از قوانین مشابهی استفاده کنند، کالبراسیون دارد.

Cahn (2001)

این قاعده کالبره است اگر همه افراد از قوانینی استفاده کنند که نرخ تغییر اقداماتشان مشابه باشد. به طور شهودی، اگر بازیکنان دیگر اقدامات خود را به سرعت زیادی تغییر ندهند، روش بالا به طور ضمنی ماتریس لازم برای حل معادله  $Rq = RTq$  را معکوس می کند.

# آزمون

یک تفسیر از مفهوم کالیبراسیون در روند یادگیری، آن است که "یادگیرنده" بتواند به درستی پاسخ‌های خود را با تکرار زمانی بهبود ببخشد، حتی اگر در ابتدا دانش زیادی نداشته باشد.

ساندرونی (۲۰۰۳)

دو ویژگی برای آزمون‌های یادگیری ارائه داده است:

۱. آزمون باید پس از یک تعداد محدود دوره‌ها اعلام کند که یادگیرنده از آزمون گذرانده است یا نه.
۲. آزمون باید با احتمال بالا واقعیت یادگیرنده را تأیید کند.

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# آزمون

دکل و فینبرگ (۲۰۰۶) و اولسیفسکی و ساندرونی (۲۰۰۶) شرطی را که آزمون باید در زمان محدود نتیجه قطعی دهد، رفع کردند.

فورتنو و وهر (۲۰۰۸)

یک الگوریتم نادانه که برخی آزمون‌ها را پشت سر گذاشته باشد، باید به ضرورت پیچیده‌ای محاسباتی باشد.

آل-نجار و وینشتاین (۲۰۰۷)

تمایز دادن بین دو یادگیرنده اطلاعاتی آسان‌تر است نسبت به ارزیابی یک یادگیرنده در تنهایی.

فینبرگ و استوارت (۲۰۰۷)

احتمال مقایسه بین چندین متخصص مختلف، برخی واقعی و برخی نادرست را در نظر گرفته و نشان دادند که تنها متخصصان واقعی تضمین می‌کنند که آزمون را با چه شرایط دیگری انجام دهند.

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# همگرایی به نش

چرا همگرایی به تعادل نش؟

استفاده گسترده‌ای از تعادل نش در نظریه بازی ها  
تعادل نش = توصیف شرایطی که در آن‌ها دیگر یادگیری امکان پذیر نیست

سوال:

یک کلاس از قوانین ← استفاده همه بازیکنان از قوانین یادگیری این کلاس استفاده  
← آیا به تعادل نش همگرا می شویم.

مثال

همه بازیکنان قوانین یادگیری کالیبره یکنواخت ← همگرایی بازی به مجموعه‌ی تعادل مشترک همگرا

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# همگرایی به نش

آیا کلاس‌هایی از قوانین یادگیری وجود دارند که زمانی که همه بازیکنان از آن‌ها استفاده می‌کنند، به همگرایی جهانی به یک تعادل نش منتهی می‌شوند؟  
نتایج اولیه در این راستا منفی بود.

مفهوم دیکوپله بودن

هارت و ماس-کلل (۲۰۰۳)

در سیستم‌های پیوسته و زمان حقیقی و بدون استفاده از اطلاعات گذشته، نمی‌توان به همگرایی به تعادل در بازی اطمینان داشت.

در ادامه، نشان دادند که در روش‌های گسسته و زمان تصادفی، همگرایی به تعادل نیز تضمین نمی‌شود در شرایط "۱-یادآوری" که در آن وضعیت سیستم بر اساس آخرین پروفایل بازی تعیین می‌شود ← بهبود نتایج همگرایی گذشته فوستر و یانگ تضمین همگرایی در شرایط استفاده از دو دوره گذشته

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# همگرایی به نش

مدل یادگیری تصادفی فوستر و یانگ

فرایند یادگیری از یک اقدام "وضع موجود" ← بازبینی دوره‌ای به صورت تصادفی

بازبینی‌ها در زمان‌های تصادفی :

وضع موجود ← زمان‌های بازبینی تصادفی ← مقایسه عادلانه‌ای بین نتایج ← بررسی  
عملکرد ← ... ← یک اقدام وضع موجود جدید به طور تصادفی ← ...

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری



# همگرایی به نش

برخی روش‌های یادگیری مانند بازی خیالی ← همگرا در بسیاری از محیط

نیاز به قوانین یادگیری پیچیده‌تر!

**مدل لوین (۱۹۹۱)**

بازیکنانی که فقط به حرکات فوری خود نگاه نمی‌کنند و کمی صبورتر هستند، از تعادل نش دور می‌شوند زیرا ارزش تعهدات بلندمدت خود را درک می‌کنند.

**شاما و ارسلان (۲۰۰۵)**

دینامیک‌های پیش‌بینی حرکت بعدی حریف بدون دید بلندمدت

چرا و چگونه به قوانین یادگیری شامل پیش‌بینی حرکات آینده حریفان نیاز داریم؟!

مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# همگرایی به نش

مدل شاما و ارسلان ← محیط بازی خیالی نرم با وزن دهی نمایی به مشاهدات

$$\sigma_i(t) = (1 - \lambda)z_i(t - 1) + \lambda\sigma_i(t - 1)$$

ثابت بین ۰ و ۱  
(نشان دهنده نرخ کاهش اهمیت مشاهدات)

یک بردار با یک درایه ۱ و ما بقی ۰  
(اقدام بازی شده در زمان  $t$  توسط بازیکن  $i$ )

فرکانس وزنی تجربی بازیکن  $i$  تا زمان  
 $t$  (میانگین وزنی اقدامات بازیکن  $i$  در طول زمان)

فرکانس وزنی تجربی همیشه بین مشاهدات جدید و گذشته تعادلی برقرار کند و به مرور زمان تغییرات را به طور نمایی در نظر بگیرد.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# همگرایی به نش

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

$\sigma_{-i}(t-1)$  نشان دهنده عمل بازیکن  $-i$  در گذشته نه **حال حاضر**  
← پیش بینی بازی فعلی بازیکن  $-i$  برای تخمین بهتری از بازی او در حال حاضر

شما و ارسال ← کنترل مشتق تناسبی ← معرفی  $r_{-i}$  برای استفاده در پیش بینی  
( $\sigma_{-i}$ ،  $r_{-i}$  را دنبال می کند.)

$$r_{-i}(t) = r_{-i}(t-1) + \lambda(\sigma_{-i}(t-1) - r_{-i}(t-1))$$

پیش بینی بازی بازیکن  $-i$

$$\sigma_{-i}(t-1) + \gamma(r_{-i}(t) - r_{-i}(t-1))$$

$\gamma$  پارامتر تعیین کننده وزن به تخمین  $r(t)$

$\lambda$  پارامتر تعیین کننده سرعت تنظیم متغیر کمکی

# همگرایی به نش

مدل شاما و ارسلان ← استفاده از روش‌های تقریب تصادفی ← بررسی همگرایی سیستم به نش

$$\dot{\sigma}_i = \phi[\beta_i(\sigma_{-i} + \gamma \dot{r}_{-i}) - \sigma_i]$$

$$\dot{r}_{-i} = \lambda(\sigma_{-i} - r_{-i})$$

$$\ddot{r}_{-i} = \lambda(\dot{\sigma}_{-i} - \dot{r}_{-i})$$

این مدل نشان می‌دهد که تحت شرایطی خاص، سیستم می‌تواند به تعادل نش همگرا شود.

اگر  $\ddot{r}_{-i}$  کوچک باشد:

○  $\dot{r}_{-i}$  تقریب خوبی از  $\dot{\sigma}_{-i}$  است

○ سیستم به طور جهانی به یک توزیع تعادل نش همگرا می‌شود

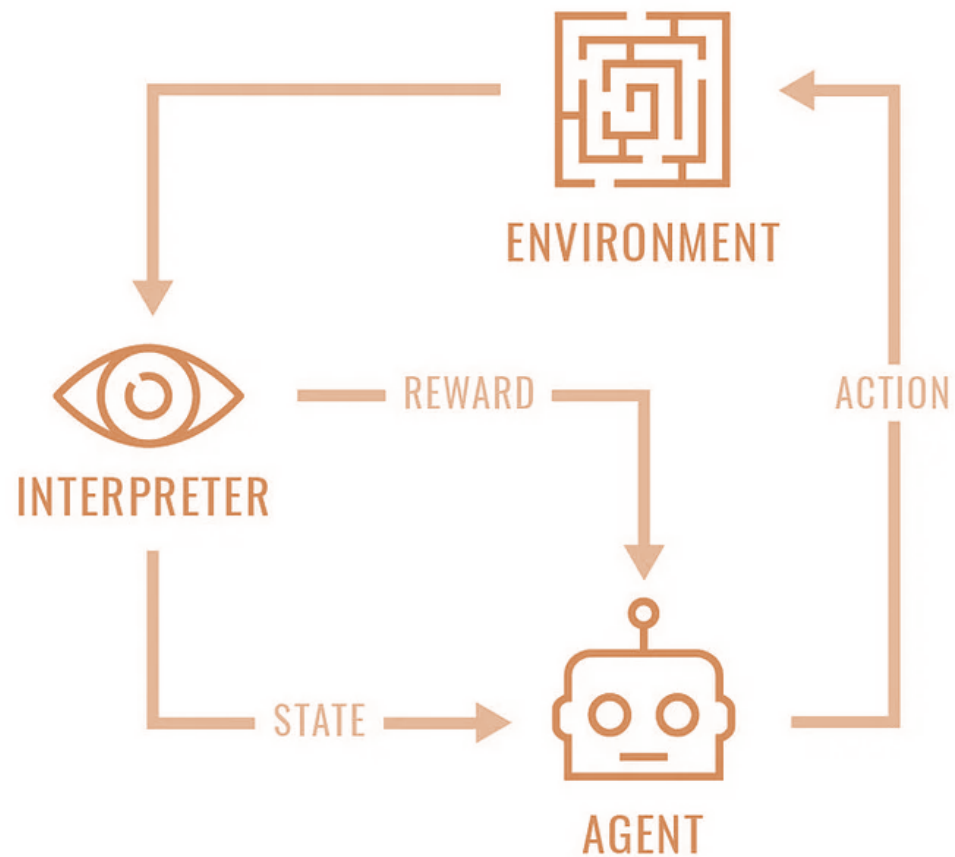
مقدمه

بخش ۲

بخش ۳

نتیجه‌گیری

# یادگیری تقویتی



مقدمه

بخش ۲

بخش ۳

نتیجه گیری



# پاداش تقویتی تجمعی (CPR)

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

پاداش تقویتی تجمعی **CPR**: یک مدل بنیادی از یادگیری تقویتی که توسط لاسلیه و همکاران مطالعه شده است. تنظیمات اولیه:

- عوامل با تخصیص وزنهای اولیه به اقدامات مختلف شروع می‌شوند.

مکانیسم به‌روزرسانی وزن:

- وزن‌ها بر اساس پاداش‌های دریافتی پس از هر اقدام به‌روزرسانی می‌شوند.

احتمال انتخاب اقدام:

- احتمال انتخاب یک اقدام خاص متناسب با پاداش تجمعی آن است.
- این احتمال با پاداش‌های کلی همه اقدامات مقایسه می‌شود.

# معادله احتمال عمل

احتمال یک عمل در CPR :

- بدست آمده توسط معادله:

$$P_k(t) = \frac{CU_k(t)}{\sum_j CU_j(t)}$$

جایی که  $CU_k(t)$  امتیاز تجمعی عمل  $k$  در زمان  $t$  است.

ویژگی های کلیدی:

- اطمینان حاصل می کند که هر عملی بی نهایت بارها امتحان شده است.
- به عوامل اجازه می دهد تا ارزش همه اقدامات را بیاموزند.

رویکرد احتمالی:

- اکتشاف و بهره برداری را متعادل می کند.
- به مرور زمان منجر به یادگیری موثرتر می شود.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری



# همگرایی و ثبات در CPR

## یافته های Laslier و همکاران:

- در سناریوهایی که یک بازیکن از یک استراتژی ثابت استفاده می کند، فرآیند یادگیری به بهترین اقدامات همگرا می شود.
- هر عمل بی نهایت و اغلب انجام می شود، تا اطمینان حاصل شود که عوامل ارزش همه اقدامات را یاد می گیرند.
- تجزیه و تحلیل با استفاده از تقریب تصادفی:

## مدل سازی شده با یک معادله دیفرانسیل زمان پیوست

$$\dot{x} = -x + r(x)$$

- X نشان دهنده کسری از زمان انتخاب هر عمل است.
- $r(x)$  احتمال هر عمل بر اساس وضعیت فعلی است.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری





# فرآیندهای تصادفی تقریبی

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

استفاده از فرآیندهای تصادفی تقریبی در RL :

- فرآیند RL زمان گسسته را با یک معادله دیفرانسیل زمان پیوسته تقریب می کند.

هدف:

- به محققان کمک می کند تا رفتار بلند مدت و ویژگی های همگرایی فرآیند یادگیری را درک کنند.

بینش های به دست آمده:

- درک چگونگی تنظیم استراتژی های خود توسط عوامل در طول زمان را فراهم می کند.

# یادگیری تقویتی آشفته

## مطالعه هاپکینز در مورد CPR:

- بررسی نسخه های CPR با اغتشاشات تصادفی کوچک.

## دینامیک تقویت آشفته:

آشفته ها تغییرات کوچک و تصادفی را در فرآیند یادگیری ایجاد می کنند.

## یافته های کلیدی:

- نشان دهید که چگونه استراتژی ها بر اساس موفقیت نسبی تکامل می یابند.
- استحکام مدل های RL را در برابر نوسانات و عدم قطعیت های جزئی توضیح دهید.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری



# مدل های تقویتی ثابت

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

## تحلیل بورگرز و سارین:

- یک مدل تقویتی با نرخ به روز رسانی ثابت مورد بررسی قرار گرفت.

## یافته های کلیدی:

- نشان داد که مدل ممکن است به اقدامات نادرست همگرا شود.
- نرخ به روز رسانی ثابت مانع از یادگیری مداوم بهترین اقدامات توسط بازیکنان می شود.

## اهمیت یافته ها:

- ضرورت مکانیسم های به روزرسانی تطبیقی را در مدل های RL برجسته می کند.

$$u_k(t+1) = u_k(t) + \gamma(R_k(t) - u_k(t))$$



# مدل های تقلیدی

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

## مدل های تقلیدی

عواملی را درگیر کنید که بر اساس موفقیت مشاهده شده از اقدامات دیگران تقلید می کنند.

## بیورنستاد و وایبول

عوامل مطالعه با اطلاعات پر سر و صدا در مورد پاداش دیگران. نشان داد که چگونه نمایندگان به استراتژی‌های با بهترین عملکرد تغییر می کنند.

## بینمور و ساموئلسون

تقلید کاوش شده با آرزوهای ثابت.

## نواک و همکاران

تقلید ترکیبی از استراتژی های موفق و محبوب. پویایی یادگیری متنوع را نشان داد.



# مدل های مبتنی بر آرزوها

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

## مدل های مبتنی بر آرزو

اگر پاداش های دریافتی کمتر از سطح هدف باشد، نمایندگان رفتار خود را تنظیم می کنند.

## آثار چو و ماتسووی

نشان داد که چگونه عوامل می توانند همکاری یا مخالفت را بر اساس رضایت خود از پاداش های فعلی یاد بگیرند.

## بینش های کلیدی

این مدل ها نشان می دهند که چگونه آرمان های کارگزاران بر فرآیندهای یادگیری و تصمیم گیری آنها تأثیر می گذارد.



# تعادل خود تأیید کننده SCE

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

## تعادل خود تأیید کننده SCE:

- مفهومی که در آن بازیکنان اعتقاداتی در مورد استراتژی های حریف دارند که بر اساس مشاهدات آنها صحیح است.

## عقاید و اعمال:

- این باورها اقدامات بازیکنان را هدایت می کنند، حتی اگر کاملاً دقیق نباشند.

## اهمیت SCE:

- توضیح می دهد که چگونه نتایج پایدار می تواند بدون آگاهی کامل رخ دهد.
- بازیکنان بر اساس واقعیت درک شده خود بهینه سازی می کنند.

# آزمایش و تعادل

## تعیین آزمایش کافی:

- برای دستیابی به تعادل robust ضروری است.

## یافته های فوندبرگ و کرپس:

- نتایج غیر نش نمی توانند پایدار باشند اگر هر اقدامی بی نهایت مرتب انجام شود.

## بینش کلیدی:

- اهمیت آزمایش کافی را برجسته می کند.
- اطمینان حاصل می کند که تمام استراتژی های بالقوه بررسی شده و تصمیمات بهینه گرفته شده است.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# مدل های مبتنی بر باور

## مدل های مبتنی بر اعتقاد:

- نیاز به باورهای بازیکنان برای همسویی با داده های مشاهده شده.
- آگاهی از عملکردهای بازده دیگران را فرض نکنید.

## مثال ها:

- تعادل حدسی منطقی روبینشتاین و وولینسکی.
- تعادل منطقی خود تأیید کننده دکل و همکاران.

## بینش های کلیدی:

- به درک اینکه چگونه باورها یادگیری و سازگاری استراتژی را شکل می دهند کمک کنید.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری



# یادگیری بیزی

## یادگیری بیزی:

- شامل بازیکنانی می شود که باورهای خود را بر اساس داده های مشاهده شده به روز می کنند.

## یافته های Acemoglu و همکاران:

- نشان داد که چگونه عوامل بیزی منطقی می توانند با وجود مواجهه با داده های مشابه، عقاید متفاوتی را حفظ کنند.
- پیچیدگی های یادگیری در زمینه های اقتصادی را برجسته می کند.

## اهمیت مدل های بیزی:

- چارچوبی برای درک به روزرسانی های باورها فراهم کنید.
- کمک به تصمیم گیری در شرایط عدم اطمینان.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# یادگیری استقرا معکوس

## استقرا معکوس:

- تکنیکی که در آن بازیکنان از پایان بازی به عقب کار می کنند تا استراتژی های بهینه را تعیین کنند.

## همگرایی به تعادل زیربازی کامل:

- با آزمایش کافی، بازیکنان می توانند به این همگرایی دست یابند.
- تصمیم گیری بهینه را در هر مرحله تضمین می کند.

## کاربرد:

- به ویژه در بازی های فرم گسترده با تصمیم گیری های متوالی مفید است.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری

# یادگیری غیرتعادلی در اقتصاد کلان

## نقش یادگیری در نظریه های اقتصاد کلان:

- برای درک پویایی های اقتصادی و سیاست گذاری بسیار مهم است.

## یادگیری غیر تعادلی:

- به خصوص تعادل خود تأیید کننده SCE

## برنامه های کاربردی:

- مطالعه چگونگی سازگاری اقتصادها با تغییرات
- درک چگونگی تثبیت اقتصادها در طول زمان.

مقدمه

بخش ۲

بخش ۳

نتیجه گیری



# کاربردهای عملی

## کاربردهای عملی مدل های یادگیری:

- در زمینه های مختلف از جمله تئوری حراج، کالاهای عمومی و پویایی بازار استفاده می شود.

## فواید:

- به سیاستگذاران کمک می کند تا مکانیسم های بهتری طراحی کنند.
- به پیش بینی رفتار در سیستم های اقتصادی و اجتماعی کمک می کند.

## مثال ها:

- طراحی حراج های کارآمد
- درک واکنش بازار به تغییرات سیاست

مقدمه

بخش ۲

بخش ۳

نتیجه گیری



## نتیجه گیری

مقدمه

Score-Based Equilibrium Learning in Multi-Player Finite Games with Imperfect Information

بخش ۲

Exploiting Hidden Structures in Non-Convex Games for Convergence to Nash Equilibrium

بخش ۳

Scalable and Independent Learning of Nash Equilibrium Policies in Stochastic Games with Unknown Independent Chains

نتیجه گیری

Paths to Equilibrium in Normal-Form Games

The Complexity of Approximate (Coarse) Correlated Equilibrium for Incomplete Information Games



# Thanks

