

به نام خدا



یادگیری تقویتی

تمرین تئوری چهارم

یلدا شعبان زاده

۹۸۱۰۱۸۲۲

نیمسال دوم ۰۱-۰۲

## فهرست

- سوال ۱- کشف مهارت ..... ۳
- سوال ۲- یادگیری سلسله مراتبی ..... ۷
- سوال ۳- مروری بر یادگیری تقویتی معکوس ..... ۹
- سوال ۴- یادگیری برون خط ..... ۱۵

## سوال ۱- کشف مهارت

(آ)

۱.

$$I(S; Z) = H(Z) - H(Z|S)$$

در این عبارت ترم اول یا  $H(Z)$  با استفاده از uniform prior  $p(Z)$  ماکزیمم می‌شود و ترم دوم یا  $H(Z|S)$  با ماکزیمم کردن  $\log p(Z|S)$  مینیمم می‌شود و هرچه  $\log p(Z|S)$  بیشتر باشد یعنی مهارت‌ها بهتر متمایز شده‌اند. این عبارت اطلاعات متقابل بین حالت‌های محیط ( $S$ ) و متغیر پنهان نشان دهنده مهارت ( $Z$ ) را نشان می‌دهد. این مقدار اطلاعاتی را که متغیر مهارت در مورد وضعیت‌های محیط ارائه می‌دهد اندازه‌گیری می‌کند. مقدار بالاتر  $I(S; Z)$  نشان می‌دهد که مهارت‌ها آموزنده هستند و مقدار قابل توجهی از اطلاعات را در مورد حالت‌ها ارائه می‌دهند، به این معنی که مهارت‌های مختلف منجر به حالت‌های متمایز و قابل تشخیص محیط می‌شود.

۲.

$$I(A; Z|S) = H(Z|S) - H(Z|A, S) = H(A|S) - H(A|S, Z)$$

این عبارت اطلاعات متقابل مشروط بین اقدامات انجام شده توسط عامل ( $A$ ) و متغیر مهارت پنهان ( $Z$ ) را با توجه به حالات مشاهده شده محیط ( $S$ ) نشان می‌دهد. مقدار اطلاعاتی را که متغیر مهارت در مورد اقدامات ارائه می‌دهد، با در نظر گرفتن حالات مشاهده شده اندازه‌گیری می‌کند. مقدار بالاتر  $I(A; Z|S)$  نشان می‌دهد که مهارت‌ها آموزنده هستند و با توجه به حالت‌های مشاهده شده، مقدار قابل توجهی از اطلاعات را در مورد اقدامات ارائه می‌دهند. این نشان می‌دهد که مهارت‌های مختلف با در نظر گرفتن زمینه خاص ارائه شده توسط state‌ها منجر به الگوهای اقدام متمایز و قابل تشخیص می‌شود.

۳.

$$H(A|S)$$

این عبارت بیانگر آنتروپی مشروط اقدامات انجام شده توسط عامل ( $A$ ) با توجه به حالات مشاهده شده محیط ( $S$ ) است. عدم قطعیت یا تصادفی بودن اعمال عامل را با در نظر گرفتن حالات مشاهده شده اندازه‌گیری می‌کند. این مقدار اطلاعاتی را که با توجه به حالات مشاهده شده هنوز در مورد اقدامات نامشخص یا غیرقابل پیش‌بینی هستند، اندازه‌گیری می‌کند. مقدار کمتر  $H(A|S)$  نشان می‌دهد که اقدامات عامل با توجه به حالات مشاهده شده قطعی‌تر یا قابل پیش‌بینی‌تر هستند. این نشان می‌دهد که اقدامات عامل بیشتر محدود شده یا تحت تأثیر حالات مشاهده شده است و تنوع یا تصادفی کمتری در فرآیند انتخاب اکشن وجود دارد.

\*\*\*

(ب)

$$F(\theta) \triangleq I(S; Z) + H(A|S) - I(A; Z|S)$$

$I(S; Z)$ : این عبارت شرایطی را برآورده می‌کند که مهارت‌های مختلف باید منجر به مشاهده حالات مختلف محیط شود. مقدار بالاتر  $I(S; Z)$  نشان می‌دهد که مهارت‌ها اطلاعات بیشتری در مورد حالت‌ها ارائه می‌دهند، به این معنی که مهارت‌های مختلف منجر به حالت‌های متمایز و متنوع می‌شوند. (شرط اول) همچنین نشان می‌دهد برای تمییز بین مهارت‌های مختلف ما تنها به حالت‌های محیط نیاز داریم نه عمل‌هایی که توسط

عامل انجام می‌شود. درواقع  $H(Z|S)$  هرچقدر کمتر باشد یعنی با دانستن  $S$  می‌توان راحت‌تر  $Z$  را پیش‌بینی کرد و اینکه برای تمییز بین مهارت‌های مختلف ما تنها به حالت‌های محیط نیاز داریم نه عمل‌هایی که توسط عامل انجام می‌شود. (شرط دوم)

$H(A|S)$ : با به حداکثر رساندن  $H(A|S)$ ، این معیار عامل را تشویق می‌کند تا سطح بالاتری از تصادفی یا تغییرپذیری را در اعمال خود، با توجه به حالات مشاهده‌شده، نشان دهد. این کار باعث می‌شود کشف مهارت‌های متنوع برای عامل با اجازه دادن برای کشف گزینه‌های عمل مختلف در پاسخ به حالت‌های یکسان، بیشتر شود.

$I(A;Z|S)$ : این عبارت با در نظر گرفتن تأثیر حالت‌های مشاهده‌شده، با جریمه کردن وابستگی یا همبستگی بین اقدامات و مهارت‌ها، معیار را تکمیل می‌کند. با به حداقل رساندن  $I(A;Z|S)$ ، این معیار با توجه به حالات مشاهده‌شده، اقدامات را تشویق می‌کند تا کمتر به مهارت‌ها وابسته باشند. این امر کشف مهارت‌هایی را که از یکدیگر متمایز هستند، ترویج می‌کند، زیرا اقدامات تشویق می‌شوند تا اتکای کمتری به مهارت خاص مرتبط با حالت‌های مشاهده‌شده داشته باشند. (شرط سوم)

با ترکیب این عبارات در معادله (۱)، این معیار شرایط مشاهده حالات قابل تمایز، ترویج اقدامات متنوع با توجه به حالت‌ها و به حداقل رساندن وابستگی بین اقدامات و مهارت‌ها را در نظر می‌گیرد.

\*\*\*

(ج)

$$F(\theta) \triangleq I(S;Z) + H(A|S) - I(A;Z|S) = I(S;Z) + H(A|S) - H(A|S) + H(A|S,Z) \\ = H(Z) - H(Z|S) + H(A|S,Z)$$

عبارت اول توزیع prior را روی  $p(z)$  تشویق می‌کند تا آنتروپی بالایی داشته باشد.  $p(z)$  می‌تواند در این رویکردمان یکنواخت شود تا تضمین شود که حداکثر آنتروپی را دارد. (در DIAYN به این صورت است)

عبارت دوم نشان می‌دهد که استنتاج مهارت  $Z$  از وضعیت فعلی باید آسان باشد. درواقع  $H(Z|S)$  هرچقدر کمتر باشد یعنی با دانستن  $S$  می‌توان راحت‌تر  $Z$  را پیش‌بینی کرد و اینکه برای تمییز بین مهارت‌های مختلف ما تنها به حالت‌های محیط نیاز داریم نه عمل‌هایی که توسط عامل انجام می‌شود.

عبارت سوم نشان می‌دهد که هر مهارت باید تا حد امکان تصادفی عمل کند، که ما با استفاده از حداکثر سیاست آنتروپی برای نشان دادن هر مهارت به آن دست می‌یابیم. این کار باعث می‌شود کشف مهارت‌های متنوع برای عامل با اجازه دادن برای کشف گزینه‌های عمل مختلف در پاسخ به حالت‌های یکسان، بیشتر شود.

\*\*\*

(د)

$$\begin{aligned}
F(\theta) &\triangleq H(Z) - H(Z|S) + H(A|S, Z) = H(A|S, Z) + E_{z \sim p(z), s \sim \pi(z)}[\log p(z|s)] - E_{z \sim p(z)}[\log p(z)] \\
&= H(A|S, Z) + E_{z \sim p(z), s \sim \pi(z)}[\log p(z|s)] - E_{z \sim p(z), s \sim \pi(z)}[\log p(z)] \\
&= H(A|S, Z) + E_{z \sim p(z), s \sim \pi(z)}[\log p(z|s) - \log p(z)] \\
&= H(A|S, Z) + E_{z \sim p(z), s \sim \pi(z)} \left[ \log \frac{p(z|s)}{p(z)} \right] \\
&= H(A|S, Z) + E_{z \sim p(z), s \sim \pi(z)} \left[ \log \left( \frac{p(z|s)}{q_\phi(z|s)} \times \frac{q_\phi(z|s)}{p(z)} \right) \right] \\
&= H(A|S, Z) + E_{z \sim p(z), s \sim \pi(z)} \left[ \log \frac{p(z|s)}{q_\phi(z|s)} + \log \frac{q_\phi(z|s)}{p(z)} \right] \\
&= H(A|S, Z) + E_{z \sim p(z), s \sim \pi(z)} \left[ D_{KL} \left( p(z|s) || q_\phi(z|s) \right) \right] + E_{z \sim p(z), s \sim \pi(z)} \left[ \log \frac{q_\phi(z|s)}{p(z)} \right] \\
&\geq H(A|S, Z) + E_{z \sim p(z), s \sim \pi(z)} [\log q_\phi(z|s) - \log p(z)] = G(\theta, \phi)
\end{aligned}$$

\*\*\*

(ه) اولین ترم در  $G(\theta, \phi)$  نشان‌دهنده‌ی این است که هر مهارت باید تا حد امکان تصادفی عمل کند، یعنی آنتروپی بالایی از اکشن نسبت به دانستن استیت و مهارت می‌خواهیم. یعنی همان exploration؛ در دومین ترم در بخش  $E_{z \sim p(z), s \sim \pi(z)}[-\log p(z)]$  می‌خواهیم آنتروپی مهارت‌ها را زیاد کنیم. یعنی به احتمال یکسانی برای هر مهارت نزدیک شویم. در بخش  $E_{z \sim p(z), s \sim \pi(z)}[\log q_\phi(z|s)]$  هم می‌خواهیم احتمال discriminate کردن مهارت‌ها را به شرط دانستن استیت افزایش دهیم. به همین دلیل با افزایش ترم دوم؛ امکان تصادفی عمل کردن کاهش می‌یابد و برعکس. پس ما یک trade-off بین exploration و discrimination داریم. به همین دلیل در DIAYN با اضافه کردن ضریب  $\alpha$  به ترم اول این مصالحه را کنترل می‌کنند.

\*\*\*

(و) SAC آنتروپی سیاست را بر اقدامات به حداکثر می‌رساند، که از عبارت آنتروپی در هدف G ما مراقبت می‌کند. می‌توان ریوارد را با توجه به تابع هدف و discriminative بودن برابر زیر قرار داد:

$$r_z(s, a) = \log q_\phi(z|s) - \log p(z)$$

با این روش عامل برای بازدید از حالاتی که به راحتی قابل تشخیص است پاداش دریافت می‌کند، در حالی که discriminator آپدیت می‌شود تا مهارت Z را از حالت های بازدید شده بهتر استنتاج کند. تنظیم آنتروپی به عنوان بخشی از به روز رسانی SAC رخ می‌دهد.

\*\*\*

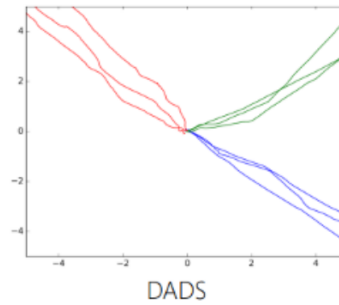
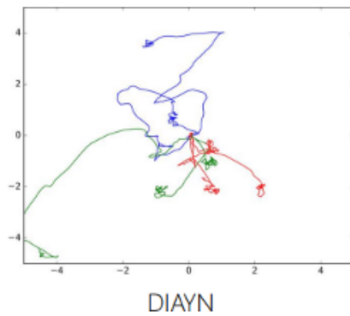
(ج)

$$\begin{aligned}
I(s'; z|s) &= H(z|s) - H(z|s', s) = H(s'|s) - H(s'|s, z) \\
&= -E_{s, s' \sim p}[\log p(s'|s)] + E_{z, s, s' \sim p}[\log p(s'|s, z)] = E_{z, s, s' \sim p} \left[ \log \frac{p(s'|s, z)}{p(s'|s)} \right] \\
&= E_{z, s, s' \sim p} \left[ \log \left( \frac{p(s'|s, z)}{q_\phi(s'|s, z)} \times \frac{q_\phi(s'|s, z)}{p(s'|s)} \right) \right] \\
&= E_{z, s, s' \sim p} \left[ \log \frac{p(s'|s, z)}{q_\phi(s'|s, z)} + \log \frac{q_\phi(s'|s, z)}{p(s'|s)} \right] \\
&= E_{s, z \sim p} \left[ D_{KL} \left( p(s'|s, z) || q_\phi(s'|s, z) \right) \right] + E_{z, s, s' \sim p} \left[ \log \frac{q_\phi(s'|s, z)}{p(s'|s)} \right]
\end{aligned}$$

اطلاعات متقابل در معادله اول می‌گوید که چقدر می‌توان در مورد  $s'$  با توجه به  $s, z$  یا به طور متقارن،  $z$  را با توجه به انتقال از  $s \rightarrow s'$  دانست. از معادله دوم، حداکثر کردن این هدف مربوط به به حداکثر رساندن تنوع انتقال‌های تولید شده در محیط است، که با آنتروپی  $H(s'|s)$  نشان داده می‌شود، در حالی که  $z$  را با به حداقل رساندن آنتروپی  $H(s'|s, z)$  در مورد وضعیت بعدی  $s'$  اطلاع رسانی می‌کند. به طور شهودی، مهارت‌های  $z$  را می‌توان به عنوان توالی‌های عمل انتزاعی تفسیر کرد که با انتقال‌های ایجاد شده در محیط (و نه فقط با وضعیت فعلی) قابل شناسایی هستند.

بنابراین، بهینه‌سازی این اطلاعات متقابل را می‌توان به عنوان رمزگذاری مجموعه‌ای از مهارت‌ها در فضای پنهان  $z$  درک کرد، در حالی که انتقال‌ها را برای  $z \in Z$  معین قابل پیش‌بینی می‌کند.

تفاوت این با تابع هدف در قسمت ب هم در این است که در این قسمت ترنیشن‌ها یا به عبارتی استیت‌های بعدی نیز در نظر گرفته شده‌اند تا از حرکات رندم در یک محدوده توسط عامل جلوگیری شود زیرا در اینجا  $H(s'|s)$  قرار است ماکزیمم شود و یعنی استیت‌های آینده احتمالات یکسانی بهتر است داشته باشند. اما با مینیمم کردن  $H(s'|s, z)$  یعنی شرطی شدن روی مهارت باعث افزایش احتمال استیت  $s'$  نسبت به  $s$  و  $z$  می‌شود. اما در ب وجود اکشن‌های متفاوت ممکن است باعث تمایز استیت نهایی در آن مهارت نشود. همچنین DIAYN بر این ایده استوار است که مهارت‌های متنوع را می‌توان برای حل وظایف مختلف مورد استفاده قرار داد و می‌تواند منجر به عملکرد بهتر از یادگیری یک پالیسی واحد برای همه کارها شود. این موجب افزایش پیچیدگی مهارت‌های یادگرفته شده می‌شود. اما DADS بر این ایده استوار است که مهارت‌ها را می‌توان با کشف دینامیک‌های زیربنایی محیط و سپس استفاده از این دانش برای یادگیری سیاست‌هایی که می‌تواند وظایف را حل کند، یاد گرفت.



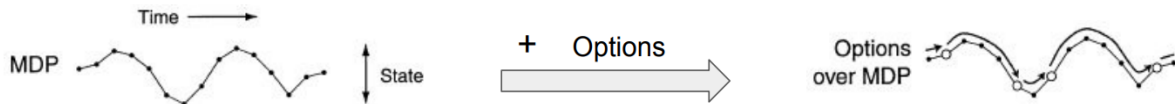
## سوال ۲- یادگیری سلسله مراتبی

(آ)

دلیل استفاده از semi-MDP زمانی که می‌خواهیم temporal abstraction داشته باشیم این است که اقدامات و برنامه‌ریزی‌های موقتی را با temporally abstract knowledge امکان‌پذیر می‌کند. علاوه بر این، انتقال دانش با استفاده از دانش دامنه برای تعریف گزینه‌ها امکان‌پذیر است و راه‌حل‌های اهداف فرعی را می‌توان مجدداً مورد استفاده قرار داد.

به طور خلاصه این فواید عبارتند از:

- Knowledge Transfer: استفاده از دانش دامنه برای تعریف option؛ راه حل های اهداف فرعی را می توان مجددا استفاده کرد.
- semi-MDP = option + MDP: یک تئوری معرفی شده است که در آن می‌گویند برای هر MDP و هر مجموعه‌ای از optionهای تعریف شده بر روی آن MDP، فرآیند تصمیم‌گیری که فقط از بین آن optionها انتخاب می‌شود و هر کدام را تا خاتمه اجرا می‌کند، یک نیمه MDP است.



- یادگیری و planning به صورت efficient

[منبع](#)

\*\*\*

(ب)

همانطور که می‌دانیم مقدار value حالت S زیر پالیسی  $\pi$  برابر است با: (state-value function  $\pi$ )

$$V^{\pi}(s) = E[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s, \pi] = E[r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_t = s, \pi]$$

$$= \sum_{a \in A_s} \pi(a|s) \left[ r_s^a + \gamma \sum_{s'} p_{ss'}^a V^{\pi}(s') \right]$$

حال optimal state-value function یا همان معادله بلمن برای S تحت  $\pi$  برابر است با:

$$V^*(s) = \max_{\pi} V^{\pi}(s) = \max_{a \in A_s} E[r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a] = \max_{a \in A_s} \left[ r_s^a + \gamma \sum_{s'} p_{ss'}^a V^*(s') \right]$$

حال اگر option داشته باشیم برای هر مارکف پالیسی  $\mu$  state-value function را می‌توان به صورت زیر نوشت:

$$V^{\mu}(s) = E[r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{k-1} r_{t+k} + \gamma^k V^{\mu}(s_{t+k}) | \varepsilon(\mu, s, t)] = \sum_{o \in O_s} \mu(s, o) \left[ r_s^o + \sum_{s'} p_{ss'}^o V^{\mu}(s') \right]$$

همچنین معادله بلمن برای  $s, a$  برابر است با:

$$\begin{aligned} Q^\mu(s, o) &= E[r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{k-1} r_{t+k} + \gamma^k V^\mu(s_{t+k}) | \varepsilon(o, s, t)] \\ &= E \left[ r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{k-1} r_{t+k} + \gamma^k \sum_{o \in O_s} \mu(s_{t+k}, o) Q^\mu(s_{t+k}, o') | \varepsilon(o, s, t) \right] \\ &= r_s^o + \sum_{s'} p_{ss'}^o \sum_{o \in O_{s'}} \mu(s', o') Q^\mu(s', o') \end{aligned}$$

معادله دوم بر اساس تعریف مقدار حالت بعدی به عنوان expectation مقادیر  $Q$  در وضعیت بعدی نسبت به سیاست فعلی به دست آمده است.

\*\*\*

(ج) ابتدا برای  $V_O^*(s)$  داریم:

$$\begin{aligned} V_O^*(s) &:= \max_{\mu \in \Pi(O)} V^\mu(s) = \max_{o \in O_s} E[r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{k-1} r_{t+k} + \gamma^k V_O^*(s_{t+k}) | \varepsilon(o, s, t)] \\ &= \max_{o \in O_s} \left[ r_s^o + \sum_{s'} p_{ss'}^o V_O^*(s') \right] = \max_{o \in O_s} E[r + \gamma^k V_O^*(s) | \varepsilon(o, s)] \end{aligned}$$

حال می‌دانیم که:  $V_O^*(s) = \max_{o \in O_s} Q_O^*(s, o')$

$$\begin{aligned} Q_O^*(s, o) &:= \max_{\mu \in \Pi(O)} Q^\mu(s, o) = E[r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{k-1} r_{t+k} + \gamma^k V_O^*(s_{t+k}) | \varepsilon(o, s, t)] \\ &= E \left[ r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{k-1} r_{t+k} + \gamma^k \max_{o \in O_{s_{t+k}}} Q_O^*(s_{t+k}, o') | \varepsilon(o, s, t) \right] \\ &= r_s^o + \sum_{s'} p_{ss'}^o \max_{o \in O_{s'}} Q_O^*(s', o') = E \left[ r + \gamma^k \max_{o \in O_{s'}} Q_O^*(s', o') | \varepsilon(o, s) \right] \end{aligned}$$



### سوال ۳- مروری بر یادگیری تقویتی معکوس

(آ) یادگیری تقویتی معکوس (IRL) تکنیکی است که به جای یادگیری مستقیم یک پالیسی، هدف آن یادگیری تابع پاداش از رفتار اکسپرت است. دلایل مختلفی وجود دارد که چرا ما ممکن است به جای یادگیری مستقیم یک پالیسی از رفتار اکسپرت در این زمینه، استفاده از یادگیری تابع پاداش از طریق IRL را انتخاب کنیم:

تابع پاداش آموخته شده می تواند بینشی در مورد ترجیحات و اهداف اساسی متخصص ارائه دهد که می تواند در درک کار و طراحی پالیسی های جدید مفید باشد. تابع پاداش آموخته شده می تواند وظیفه اساسی را به روشی قوی تر از یادگیری مستقیم یک پالیسی از رفتار اکسپرت نشان دهد، زیرا تابع پاداش می تواند رفتار متخصص را در طیف وسیعی از موقعیت ها توضیح دهد. درواقع می تواند generalization بهتری داشته باشد و در حالات OOD موجب خسارات برگشت ناپذیر نشود. یادگیری تابع پاداش، انعطاف پذیری بیشتری را در تعریف کار امکان پذیر می کند، زیرا می توانیم رفتار مورد نظر را از طریق تابع پاداش به جای تکیه بر مجموعه ثابتی از نمایش های اکسپرت مشخص کنیم.

پس با این کار می توان امید داشت که پالیسی یادگرفته ی ما بهتر از پالیسی اکسپرت بشود. درواقع ما intent را در اینجا می خواهیم یاد بگیریم. اما یک مشکل در IRL وجود دارد و آن هم under specification است. یعنی به ازای یک رفتار مشخص و ثابت از اکسپرت، ریوارد فانکشن های متفاوتی می توانند آن رفتار را نشان دهند. برای مثال اگر ما یک demonstration ای مثل این جهت های سبز رنگ در شکل زیر داشته باشیم، ریواردهای مختلفی می توانند با این رفتار متناسب باشند.



درواقع مسئله ill-posed است و ما باید بتوانیم بین این حالات یکی را به عنوان هدف انتخاب کنیم زیرا نمی توان مشخص کرد هدف اصلی کدام بوده است.

\*\*\*

(ب)

۱. روش های inverse RL یک مشکلی که دارند under specification است. پس باید یک سری constraint هایی داشته باشیم که این مشکل را حل کند. به همین دلیل فرض کردیم ریوارد فانکشن ما حالت پارامتریک دارد و از فیچرهای استیت-اکشن استفاده می کند و به شکل یک تابع خطی با پارامترهای ناشناخته  $\phi$  ریوارد تعریف می شود.

$$r_{\psi}(s, a) = \sum_{i=1}^n \psi_i f_i(s, a) = \psi^T f(s, a)$$

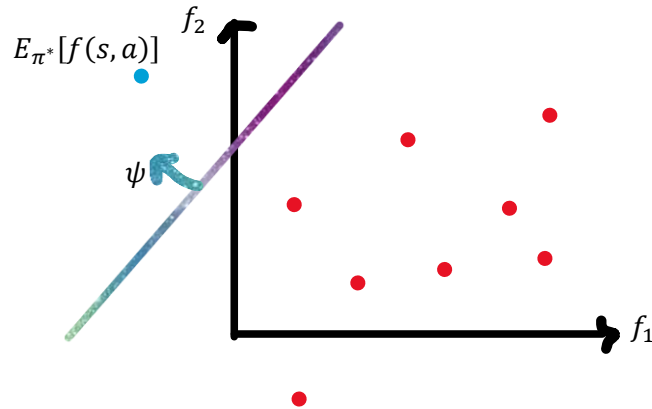
حال فرض کنید که  $\pi^{r_{\psi}}$  پالیسی آپتیمال برای  $r_{\psi}$  باشد. آنگاه داریم:

$$E_{\pi^r \psi}[f(s, a)] = E_{\pi^*}[f(s, a)]$$

این مفهوم feature matching است و درواقع می‌گوید این فیچرهایی که تعریف کردیم و ریوارد از آنها ساخته می‌شود می‌خواهیم با پالیسی‌ای که از اکسپرت می‌آید اکسپکتد یکسانی داشته باشد تا trajectoryهای بین اکسپرت با فیچرهایی که ما یاد گرفتیم مچ بشوند.

۲. این کار مشکل under specification را حل نمی‌کند. زیرا این تساوی می‌تواند با  $\psi$ های مختلف برقرار شود. مثلاً اگر در یک grid world فیچرها را به شکل باینری در نظر بگیریم که آیا در خانه  $i, j$  قرار گرفته‌ایم یا نه. یعنی در جدول  $m \times n, m \times n$  تا فیچر باینری تعریف به این صورت کنیم، در هر استیتی باشیم یکی ۱ و بقیه ۰ هستند. در این حالت پس باز هم مشکل under specification وجود دارد.

۳. برای حل این مشکل یک کار استفاده از ایده‌ی SVM و maximum margin است. به این صورت که در فضای فیچرها هرکدام از  $E[f(s, a)]$  یک نقطه هستند. اگر فیچر اسپیس ما به اندازه‌ی کافی قوی باشد نقطه‌ی متناظر با  $E_{\pi^*}[f(s, a)]$  از سایر نقاط که  $E_{\pi^r \psi}[f(s, a)]$  هستند (به ازای  $\psi$ های مختلف)، به خوبی تفکیک می‌شود.



این مسئله می‌تواند ill-posed باشد زیرا یک  $\psi$  واحد وجود ندارد. بنابراین از maximum margin استفاده می‌کنیم. پس می‌توانیم مسئله را به شکل زیر تعریف کنیم:

$$\begin{aligned} \max_{\psi, m} \quad & m \\ \text{s.t.} \quad & \psi^T \mathbb{E}_{\pi^*}[f(s, a)] \geq \max_{\pi} \psi^T \mathbb{E}_{\pi}[f(s, a)] + m \end{aligned}$$

که در اینجا  $m$  متغیر کمکی است. مشکلات این روش این است که ما در فضای همه‌ی پالیسی‌ها می‌خواهیم این کلسیفیکیشن را انجام دهیم. بنابراین ممکن است یکی از آن پالیسی‌ها همان  $\pi^*$  باشد. پس این  $m$  صفر می‌شود و نمی‌خواهیم این اتفاق بیفتد. کاری که برای حل این مشکل می‌کنند استفاده از KL-divergence است.

$$\min_{\psi} \frac{1}{2} \|\psi\|^2 \quad \text{such that } \psi^T E_{\pi^*}[\mathbf{f}(\mathbf{s}, \mathbf{a})] \geq \max_{\pi \in \Pi} \psi^T E_{\pi}[\mathbf{f}(\mathbf{s}, \mathbf{a})] + D(\pi, \pi^*)$$

e.g., difference in feature expectations!

در اینجا ما فاصله را برای سیاست‌هایی که متمایزتر از سیاست‌های اکسپرت هستند، بیشتر می‌کنیم، در حالی که سیاست‌های مشابه فاصله کمتری دارند و حاشیه کمتری دریافت می‌کنند.

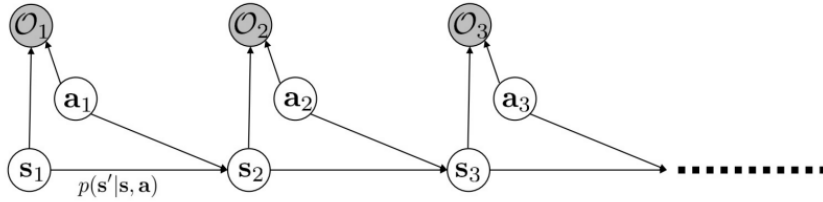
بنابراین مسئله به این آپتیمیزیشن می‌تواند تبدیل شود. مشکلاتی که باز هم دارد این است که فرض کردیم پالیسی اکسپرت آپتیمال است اما ممکن است sub optimal باشد. می‌توان slack variable را اضافه کرد اما کلا مسئله‌ی constraint optimization زیاد در دیپ لرنینگ خوش تعریف نیست.

\*\*\*

(ج)

۱.

$$\begin{aligned} p(O_t | s_t, a_t) &= \exp(r_{\psi}(s_t, a_t)) \\ \Rightarrow p(\tau | O_{1:T}) &= \frac{p(\tau, O_{1:T})}{p(O_{1:T})} \propto p(\tau) p(O_{1:T} | \tau) = p(\tau) \prod_{t=1}^T p(O_t | \tau) = p(\tau) \prod_{t=1}^T \exp(r_{\psi}(s_t, a_t)) \\ &= p(\tau) \exp\left(\sum_t r_{\psi}(s_t, a_t)\right) \end{aligned}$$



به طور دقیق‌تر می‌توان گفت که مساوی در حالت زیر است که در اینجا Z همان partition function است. درواقع:

$$p(\tau | O_{1:T}) = \frac{1}{Z} p(\tau) \exp\left(\sum_t r_{\psi}(s_t, a_t)\right), \quad Z = \int p(\tau) \exp\left(\sum_t r_{\psi}(s_t, a_t)\right) d\tau$$

در maximum likelihood learning داریم:

$$L_{\psi} := \frac{1}{N} \sum_{i=1}^N \log p(\tau_i | O_{1:T}, \psi)$$

همچنین:

$$p(\tau | O_{1:T}) \propto \underbrace{p(\tau)}_{\text{can ignore (independent of } \psi)} \exp\left(\sum_t r_{\psi}(s_t, a_t)\right)$$

بنابراین داریم:

$$\begin{aligned}
L_\psi &= \frac{1}{N} \sum_{i=1}^N \log p(\tau_i | O_{1:T}, \psi) = \frac{1}{N} \sum_{i=1}^N \log \frac{1}{Z} p(\tau_i) \exp \left( \sum_t r_\psi(s_t, a_t) \right) \\
&= \frac{1}{N} \sum_{i=1}^N \log p(\tau_i) + \log \exp \left( \sum_t r_\psi(s_t, a_t) \right) - \log Z \\
&= \frac{1}{N} \sum_{i=1}^N \log p(\tau_i) + \log \exp \left( \sum_t r_\psi(s_t, a_t) \right) - \log Z \\
\Rightarrow \max_\psi L_\psi &= \max_\psi \frac{1}{N} \sum_{i=1}^N \log p(\tau_i | O_{1:T}, \psi) = \max_\psi \frac{1}{N} \sum_{i=1}^N \log \exp \left( \sum_t r_\psi(s_t, a_t) \right) - \log Z \\
&= \max_\psi \frac{1}{N} \sum_{i=1}^N \sum_t r_\psi(s_{t,i}, a_{t,i}) - \log Z = \max_\psi \frac{1}{N} \sum_{i=1}^N r_\psi(\tau_i) - \log Z
\end{aligned}$$

بنابراین:

$$\max_\psi \frac{1}{N} \sum_{i=1}^N \log p(\tau_i | O_{1:T}, \psi) = \max_\psi \frac{1}{N} \sum_{i=1}^N r_\psi(\tau_i) - \log Z$$

نشان دادیم  $L_\psi$  در معادله‌ی بالا صدق می‌کند.

۲. اگر  $\log Z$  را نداشته باشیم؛ آنگاه انگار می‌خواهیم میانگین  $r_\psi(\tau_i)$  ها را حساب کنیم و این باعث می‌شود trajectoryهایی که اکسپرت دیده است نسبت به سایر trajectoryها محتمل‌تر بشوند و دیگر سایر آنها حتی اگر ريوارد بالایی داشته باشند محتمل نشوند.

۳.

$$\begin{aligned}
\nabla_\psi L_\psi &= \nabla_\psi \frac{1}{N} \sum_{i=1}^N r_\psi(\tau_i) - \log Z = \frac{1}{N} \sum_{i=1}^N \nabla_\psi r_\psi(\tau_i) - \nabla_\psi \log Z = \frac{1}{N} \sum_{i=1}^N \nabla_\psi r_\psi(\tau_i) - \frac{1}{Z} \\
&= E_{\tau \sim \pi^*(\tau)} [\nabla_\psi r_\psi(\tau)] - \frac{1}{Z} \int \underbrace{p(\tau) \exp(r_\psi(\tau))}_{p(\tau | O_{1:T}, \psi)} \nabla_\psi r_\psi(\tau) d\tau \\
&= E_{\tau \sim \pi^*(\tau)} [\nabla_\psi r_\psi(\tau)] - E_{p(\tau | O_{1:T}, \psi)} [\nabla_\psi r_\psi(\tau)]
\end{aligned}$$

حال برای ترم دوم می‌توان گفت:

$$E_{p(\tau | O_{1:T}, \psi)} [\nabla_\psi r_\psi(\tau)] = E_{p(\tau | O_{1:T}, \psi)} \left[ \nabla_\psi \sum_t r_\psi(s_t, a_t) \right] = \sum_{t=1}^T E_{(s_t, a_t) \sim p(s_t, a_t | O_{1:T}, \psi)} [\nabla_\psi r_\psi(s_t, a_t)]$$

همچنین می توان گفت:

$$p(s_t, a_t | O_{1:T}, \psi) = p(a_t | s_t, O_{1:T}, \psi) p(s_t | O_{1:T}, \psi)$$

که بنا بر تعریف بیزین از اسلایدها می دانیم:

$$\begin{aligned} p(a_t | s_t, O_{1:T}, \psi) &= \frac{\beta(s_t, a_t)}{\beta(s_t)}, & p(s_t | O_{1:T}, \psi) &\propto \beta(s_t) \alpha(s_t) \\ \Rightarrow p(s_t, a_t | O_{1:T}, \psi) &\propto \underbrace{\beta(s_t, a_t)}_{\text{backward message}} \underbrace{\alpha(s_t)}_{\text{forward message}} &= \mu_t(s_t, a_t) \end{aligned}$$

بنابراین داریم:

$$\begin{aligned} E_{p(\tau | O_{1:T}, \psi)} [\nabla_{\psi} r_{\psi}(\tau)] &= \sum_{t=1}^T E_{(s_t, a_t) \sim p(s_t, a_t | O_{1:T}, \psi)} [\nabla_{\psi} r_{\psi}(s_t, a_t)] \\ &= \sum_{t=1}^T \int \int \mu_t(s_t, a_t) \nabla_{\psi} r_{\psi}(s_t, a_t) ds_t da_t \end{aligned}$$

۴. در این عبارت جدید که بدست آوردیم ما نیاز داریم که روی فضای استیت-اکشن انتگرال بگیریم درحالی که در عبارت قبلی نیاز بود تا روی همه‌ی trajectoryها انتگرال بگیریم که این کار غیرممکن و بسیار سخت بود. بنابراین این روش یک روش عملی برای حساب کردن وزن‌های  $\psi$  است.

نکته‌ای که وجود دارد این است که در این روش نیز باید فضای استیت و اکشن محدود و کوچک باشد؛ اما در فضای محدود و کوچک نیز trajectoryها بزرگ می‌شوند.

۵. پس از هر آپدیت  $\psi$  باید هرسری بیزین RL بنویسیم و به همین دلیل خیلی scalable نیست و ایده‌اش قابل استفاده است. برای عملی شدن باید بتوانیم فضای بزرگ state action نیز هندل کنیم و اینکه state فقط از sampling بیایند و دینامیک ناشناخته باشد نیز مشکلات اپلای کردن این روش در فضاهای بزرگ است. حال فرض می‌کنیم  $\psi$  را در هر گام می‌دانیم. پس با این یک ریوارد تعریف می‌کنیم و یک MaxEnt-RL ران می‌کنیم و حال می‌توان از این RL سمپل گرفت. زیرا یک پالیسی رندم بدست می‌آید و با هر ران یک trajectory می‌گیریم.

MaxEnt RL به فرم زیر است:

$$J(\theta) = \sum_t E_{\pi(s_t, a_t)} [r_{\psi}(s_t, a_t)] + E_{\pi(s_t)} [H(\pi(a | s_t))]$$

حال داریم:

$$\nabla_{\psi} L_{\psi} \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\psi} r_{\psi}(\tau_i) - \frac{1}{M} \sum_{i=1}^M \nabla_{\psi} r_{\psi}(\tau_i)$$

که در آن عبارت اول از مسیرهای اکسپرت و عبارت دوم از ران کردن سیاست فعلی آمده است. این کار باعث می‌شود استیمیتور جمله‌ی دوم بایاس شود زیرا trajectoryها را از پالیسی‌ای گرفته است که الزاما بهینه نیست (یا حتی soft optimal) و ما صرفا یک گام  $p(s_t, a_t | o_{1:T}, \psi)$  را بهتر کردیم. بنابراین برای حل این مشکل باید از Importance sampling استفاده کنیم.

در صورت مقدار آپتیمالیتی را می‌گذاریم که exp جمع ریواردها را دارد و در مخرج نیز پالیسی ای که تا الان improve کردیم می‌گذاریم:

$$\nabla_{\psi} L_{\psi} \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\psi} r_{\psi}(\tau_i) - \frac{1}{\sum_j w_j} \sum_{i=1}^M w_j \nabla_{\psi} r_{\psi}(\tau_i), \quad w_j = \frac{p(\tau) \exp(r_{\psi}(\tau_j))}{\pi(\tau_j)}$$

پس از ساده کردن وزن‌ها داریم:

$$\begin{aligned} w_j &= \frac{p(\tau) \exp(r_{\psi}(\tau_j))}{\pi(\tau_j)} = \frac{p(s_1) \Pi_t p(s_{t+1} | s_t, a_t) \exp(r_{\psi}(s_t, a_t))}{p(s_1) \Pi_t p(s_{t+1} | s_t, a_t) \pi(a_t | s_t)} = \frac{\Pi_t \exp(r_{\psi}(s_t, a_t))}{\Pi_t \pi(a_t | s_t)} \\ &= \frac{\exp(\sum_t r_{\psi}(s_t, a_t))}{\Pi_t \pi(a_t | s_t)} = \frac{\exp(r_{\psi}(\tau_j))}{\Pi_t \pi(a_t | s_t)} \end{aligned}$$

## سوال ۴- یادگیری برون خط

(آ)

$$\mathcal{L}(\pi) = E_{\rho_{\pi}(\tau)} \left[ \sum_{t=0}^H \delta(a_t, a_t^*) \right], \quad \pi(a_t \neq a_t^* | s) \leq \epsilon$$

فرض کنید  $\epsilon_i$  برابر اکسپکتد 0-1 loss در زمان  $i$  از  $\pi$  باشد. یعنی داریم:

$$\epsilon_i = E_{s \sim d_{\pi^*}^i} [C_{\pi}(s)] = E_{s \sim d_{\pi^*}^i} \left[ E_{a \sim \pi(\cdot | s)} [\delta(a, \pi^*(s))] \right], \quad \text{for } i = 1, 2, \dots, H$$

آنگاه داریم:

$$\Rightarrow \epsilon \geq \frac{1}{H} \sum_{i=1}^H \epsilon_i$$

درواقع  $\epsilon_i$  احتمالی را نشان می‌دهد که  $\pi$  تحت توزیع  $d_{\pi^*}^i$  اشتباه عمل کند. حال فرض کنید  $p_t$  برابر احتمالی باشد که  $\pi$  در آن در  $t$  قدم اول اشتباه عمل نکرده باشد. (با در نظر گرفتن  $\pi^*$ )؛ و  $d_t$  توزیع حالات  $\pi$  روی زمان  $t$  به شرط اینکه تا الان اشتباهی عمل نکرده باشد. همچنین  $d_t'$  برابر توزیع حالات در زمان  $t$  که با پیروی از  $\pi^*$  به دست می‌آید، اما مشروط به این واقعیت است که  $\pi$  حداقل یک اشتباه در اولین  $t-1$  حالت‌های بازدید شده مرتکب شده است. پس داریم:

$$d_{\pi^*}^t = p_{t-1} d_t + (1 - p_{t-1}) d_t'$$

اکنون در زمان  $t$ ، هزینه مورد انتظار  $\pi$  اگر تا کنون اشتباه کرده است حداکثر ۱ است، یا اگر هنوز اشتباه نکرده است  $E_{s \sim d_t} [C_{\pi}(s)]$  است. پس داریم:

$$J(\pi) \leq \sum_{t=1}^H p_{t-1} E_{s \sim d_t} [C_{\pi}(s)] + (1 - p_{t-1}) \times 1$$

فرض کنید  $e_t$  و  $e_t'$  برابر احتمال اشتباه  $\pi$  در توزیع  $d_t$  و  $d_t'$  باشد. آنگاه:

$$E_{s \sim d_t} [C_{\pi}(s)] \leq E_{s \sim d_t} [C_{\pi^*}(s)] + e_t$$

همچنین داریم:

$$\epsilon_t = p_{t-1} e_t + (1 - p_{t-1}) e_t' \Rightarrow p_{t-1} e_t \leq \epsilon_t$$

$$p_t = (1 - e_t) p_{t-1} \Rightarrow p_t \geq p_{t-1} - \epsilon_t \geq 1 - \sum_{i=1}^t \epsilon_i \Rightarrow 1 - p_{t-1} + \epsilon_t \leq \sum_{i=1}^t \epsilon_i$$

و نیز داریم:

$$J(\pi^*) = \sum_{t=1}^H p_{t-1} E_{s \sim d_t} [C_{\pi^*}(s)] + (1 - p_{t-1}) E_{s \sim d_t'} [C_{\pi^*}(s)]$$

حال می‌توان گفت:

$$\begin{aligned}
J(\pi) &\leq \sum_{t=1}^H p_{t-1} E_{s \sim d_t} [C_{\pi}(s)] + (1 - p_{t-1}) \times 1 \leq \sum_{t=1}^H p_{t-1} (E_{s \sim d_t} [C_{\pi^*}(s)] + e_t) + (1 - p_{t-1}) \\
&= \sum_{t=1}^H p_{t-1} (E_{s \sim d_t} [C_{\pi^*}(s)]) + \underbrace{p_{t-1} e_t}_{\leq \epsilon_t} + (1 - p_{t-1}) \leq J(\pi^*) + \sum_{t=1}^H \epsilon_t + (1 - p_{t-1}) \\
&\leq J(\pi^*) + \sum_{t=1}^H \sum_{i=1}^t \epsilon_i \leq J(\pi^*) + H \sum_{t=1}^H \epsilon_t \leq J(\pi^*) + H^2 \epsilon \\
\Rightarrow J(\pi) - J(\pi^*) &\leq H^2 \epsilon
\end{aligned}$$

بنابراین خطای یادگیری دارای حد بالای  $O(H^2 \epsilon)$  است.

\*\*\*

(ب) در حالت on policy از آنجا که عامل تعامل با محیط دارد می‌توان گفت:

$$\epsilon_i = E_{s \sim d_{\pi}^i} [C_{\pi}(s)] = E_{s \sim d_{\pi}^i} \left[ E_{a \sim \pi(\cdot|s)} [\delta(a, \pi^*(s))] \right], \text{ for } i = 1, 2, \dots, H$$

و بنابراین به طور واضح داریم:

$$J(\pi) \leq J(\pi^*) + \sum_{i=1}^H \epsilon_i \leq J(\pi^*) + \sum_{i=1}^H \epsilon = J(\pi^*) + H \epsilon$$

بنابراین خطای یادگیری دارای حد بالای  $O(H \epsilon)$  است.