



## تمرین سوم – یادگیری تقویتی

دانشجو: سایه جارالهی

شماره دانشجویی: ۹۸۱۰۱۳۳۹

استاد: دکتر رهبان - آقای حسنی

بهار ۱۴۰۲

## سوال یک-

### بخش الف -

همانطور که میدانیم LQR یک روش کنترل است که در آن  $\text{cost function}$  به صورت یک تابع  $\text{quadratic}$  مدل میشود و دینامیک محیط نیز به صورت یک تابع  $\text{linear}$  مدل میشود. پس از آن  $\text{optimization}$  انجام میشود تا سری اکشن ها به دست آیند (درواقع ماتریس های  $K$  و  $k$  به دست آیند که بتوان اکشن را به دست آورد) باتوجه به تابع هایی که ارائه دادیم، فرض کردیم که دینامیک مدل با تابع خطی مدل می شود، در نتیجه اگر در واقعیت دینامیک مدل بسیار پیچیده تر باشد LQR به طور مستقیم همگرا نمیشود و باید ابتدا خطی سازی دینامیک را انجام دهیم. همچنین تابع درجه دو برای  $\text{cost function}$  مشخص شده و در صورتی که  $\text{cost function}$  بسیار پیچیده تر از این تابع باشد لزوماً همگرا به جواب بهینه نمیشویم.

یکی دیگر از موارد  $\text{fully observable}$  بودن محیط است. در صورتی که محیط کاملاً شناخته شده نباشد، مدل کردن دینامیک با تابع خطی کار نادرستی است و باعث میشود به جواب بهینه همگرا نشویم زیرا برخی اکشن ها و استیت ها ناشناخته هستند و در نتیجه با انجام یک اکشن ممکن است به جای کاملاً متفاوتی حرکت کنیم.

در روش LQR ساده، در نظر میگیریم که محیط  $\text{stable}$  است. در صورتی که رندومنس داشته باشیم و محیط  $\text{stable}$  نباشد در ابتدا باید محیط را  $\text{stable}$  کنیم.

همچنین برای اینکه LQR پاسخ داشته باشد، در برخی از مراحل باید واردن ماتریس  $C$  به دست آید. در نتیجه برای آنکه جواب داشته باشد باید ماتریس  $\text{cost}$  ها مثبت معین باشد.

### بخش ب-

طبق [مقاله](#) ای که قید شده است، یکی از راه های استفاده از LQR در محیط های  $\text{partially observable}$  آن است که با استفاده از Kalman Filter بخش هایی را که  $\text{unobserved}$  هستند تخمین بزنیم. در واقع Kalman Filter با استفاده از سری داده هایی که در طول زمان مشاهده کرده است، استیت یک پراسس را تخمین میزند، به صورت یک تابع بازگشتی پیاده سازی شده است و توسط دینامیکی که مشاهده میکند و مقادیری که دارد استیت را تخمین میزند.

درواقع میتوان گفت از LQG (Linear quadratic-Gaussian) استفاده میشود. در این روش از Kalman filter برای تخمین استفاده میشود. همچنین از یک  $\text{noise}$  گاوسی مشابه آنچه در کلاس درس دیدیم استفاده میشود تا در صورتی که Kalman filter تخمین بدون نویز را در نظر میگیرد، یک نویز نیز به آن اضافه کنیم. با این دو کار و به کمک روش LQR (پس از تخمین استیت) به خروجی نهایی میرسیم. [منبع](#)

### بخش ج-

همانطور که میدانیم LQR یک روش  $\text{model-based}$  است. میتوان مطابق مثالی که در اسلایدها برای ترکیب روش های  $\text{model-free}$  و  $\text{model-based}$  برای Half-cheetah زده شد (در آن مثال در ابتدا به روش بر پایه مدل ترین انجام میشد و سپس با روش  $\text{model-free}$ )، عمل کرد.

به این صورت که در ابتدا LQR به عنوان یک پایه در نظر گرفته شود و از خروجی آن خروجی بهینه به دست آید. در گام بعدی توسط یک روش  $\text{model-free}$  در RL یک  $\text{correction term}$  به دست آید تا در نقاطی که LQR ضعف دارد، خروجی و اکشن ها را اصلاح کند. یکی از مقالاتی که در این رابطه استفاده میشود در [لینک](#) آورده شده است.

همچنین یک روش دیگر آن است که خروجی LQR را با  $\text{model-free}$  با نح خاصی ترکیب کنیم. در واقع این دو جداگانه آموزش داده شوند و خروجی آن ها به دست آید و اکشن نهایی به صورت ضریبی از این دو باشد و در واقع  $\text{weighted sum}$  گرفته شود و یا با روش های دیگر ترکیب شوند.

### بخش د-

مشابه با آنچه در کلاس درس دیدیم، در روش iLQR با استفاده از اعمال یک توزیع گاوسی برای استیت بعدی میتوان  $\text{exploration}$  را هدایت کرد و یا  $\text{uncertainty}$  در سیستم را مدل کرد. برای اینکار میتوان یک نویز گاوسی را به  $\text{dynamic}$  مسئله اضافه کرد که

میانگین آن با استفاده از تابع های داینامیک تعریف شده باشد و ماتریس کواریانس نیز به اعمال نویز کمک میکند. با اینکار و با تغییر ماتریس کواریانس به مقادیر مناسب، میتوان exploration را تقویت کرد زیرا استتیت های بعدی که برمیگردند دارای کواریانس بیشتری هستند و به exploration کمک میکند. همچنین برای مدل کردن uncertainty محیط نیز مشابه با LQR میتواند مفید باشد زیرا به ازای استتیت و اکشن، استتیت بعدی از یک توزیع می آید و همیشه مقدار مشخصی ندارد.

$$r \sim \text{Gamma}(\alpha, \beta) \rightarrow p(r|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} r^{\alpha-1} e^{-\beta r} \quad (1)$$

$$t \sim \text{exp}(d) \rightarrow p(t|d) = d e^{-dt}$$

$$(t=10, r=100) \quad (t=5, r=200) : \text{data}$$

$$\begin{cases} \beta \sim \text{Gamma}(\varepsilon, \omega) \\ d \sim \text{Gamma}(\phi, \eta) \end{cases}$$

$$(t_1, r_1) \text{ data}$$

هناك خط ~ يعني اننا نستخدمه في توقع النتائج

$$P(R|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \times (100)^{\alpha-1} \times e^{-100\beta} \times \frac{\beta^\alpha}{\Gamma(\alpha)} \times (200)^{\alpha-1} \times e^{-200\beta}$$

$$= \left( \frac{\beta^\alpha}{\Gamma(\alpha)} \right)^2 \times (100 \times 200)^{\alpha-1} \times e^{-300\beta} = A$$

MLE, finding  $\beta$  by maximizing

$$\frac{dA}{d\beta} = \frac{(100 \times 200)^{\alpha-1}}{\Gamma(\alpha)^2} \left[ \alpha \beta^{\alpha-1} e^{-300\beta} - 300 \beta^\alpha e^{-300\beta} \right] = 0$$

$$\Rightarrow \alpha \beta^{\alpha-1} e^{-300\beta} = 300 \beta^\alpha e^{-300\beta}$$

$$\boxed{\frac{\alpha}{100} = \beta}$$

$$P(T|h) = h e^{-h\alpha} \times h e^{-\alpha h} = h^2 e^{-\gamma \cdot h} \quad \text{B}$$

$$\text{MLE: } \frac{dB}{dh} = \gamma h e^{-\gamma \cdot h} - \gamma_0 h^2 e^{-\gamma \cdot h} = 0$$

$$\cancel{\gamma h e^{-\gamma \cdot h}} = \gamma_0 h^2 \cancel{e^{-\gamma \cdot h}}$$

$$h = \frac{1}{\gamma_0}$$

$$P(B | r, \alpha, \epsilon, w)$$

posterior probability

احتمال وقوع  $B$  بشرط  $r, \alpha, \epsilon, w$  (posterior probability)

$$= \frac{P(B, r | \alpha, \epsilon, w)}{P(r | \alpha, \epsilon, w)} = \frac{P(r | \alpha, B, \epsilon, w) P(B | \alpha, \epsilon, w)}{P(r | \alpha, \epsilon, w)}$$

احتمال وقوع  $B$  بشرط  $r, \alpha, \epsilon, w$  (posterior probability)

$$= \frac{P(r | \alpha, B) P(B | \epsilon, w)}{P(r | \alpha, \epsilon, w)} = \frac{\beta^\alpha}{\Gamma(\alpha)} r^{\alpha-1} e^{-\beta r} \times \frac{w^\epsilon}{\Gamma(\epsilon)} \beta^{\epsilon-1} e^{-w\beta}$$

$$= \frac{\beta^{\alpha+\epsilon-1} \times r^{\alpha-1} \times e^{-(w+r)\beta} \times w^\epsilon}{\Gamma(\alpha)\Gamma(\epsilon) \times P(r | \alpha, \epsilon, w)} \quad (I)$$

$$(*) P(r | \alpha, \epsilon, w) = \int_B (P(B | \alpha, \epsilon, w) P(r | \alpha, B, \epsilon, w)) dB$$

$$= \int_B \frac{w^\epsilon}{\Gamma(\epsilon)} \beta^{\epsilon-1} e^{-w\beta} \frac{\beta^\alpha r^{\alpha-1} e^{-\beta r}}{\Gamma(\alpha)}$$

$$= \left[ \frac{w^\epsilon r^{\alpha-1}}{\Gamma(\epsilon)\Gamma(\alpha)} \int_B \beta^{\alpha+\epsilon-1} e^{-(w+r)\beta} \right]$$

در ترمینال (L) داریم:

$$\frac{B^{\alpha+\varepsilon-1} \times e^{-(w+r_1)B}}{\int_B B^{\alpha+\varepsilon-1} e^{-(w+r_1)B}}$$

(با کمک جدول لاپلاس می‌توانیم ثابت کنیم که این تابع توزیع گاما است، و در

درجه ۱،  $\frac{(r_1+w)^{\alpha+\varepsilon}}{\Gamma(\alpha+\varepsilon)}$  ضرب می‌کنیم. هیچ تغییری در  $\Gamma(\alpha+\varepsilon, r_1+w)$  ایجاد نمی‌شود.

چون این تابع در تمام مقادیر توزیع است و به هم می‌رسد. در نتیجه فقط صورت به هم می‌رسد.

$$\frac{(r_1+w)^{\alpha+\varepsilon}}{\Gamma(\alpha+\varepsilon)} B^{\alpha+\varepsilon-1} e^{-(w+r_1)B}$$

این تابع همان  $\Gamma(\alpha+\varepsilon, r_1+w)$  است.

$$\Rightarrow \begin{cases} \varepsilon' = \alpha + \varepsilon \\ w' = w + r_1 \end{cases}$$

$$t \sim \exp(\lambda) \rightarrow P(t|\lambda) \cdot d e^{-\lambda t}$$

$$\lambda \sim \text{Gamma}(\alpha, \eta)$$

$$P(\lambda | t_1, \alpha, \eta) = \frac{P(\lambda, t_1 | \alpha, \eta)}{P(t_1 | \alpha, \eta)} = \frac{P(\lambda | \alpha, \eta) P(t_1 | \alpha, \eta, \lambda)}{P(t_1 | \alpha, \eta)}$$

$$= \frac{\eta^\alpha \lambda^{\alpha-1} e^{-\lambda \eta} \times d e^{-\lambda t_1}}{\Gamma(\alpha) P(t_1 | \alpha, \eta)} = \frac{\eta^\alpha \lambda^\alpha e^{-\lambda(\eta+t_1)}}{\Gamma(\alpha) P(t_1 | \alpha, \eta)} \quad (\Pi)$$

$$(*) P(t_1 | \alpha, \eta) = \int_{\lambda} P(t_1, \lambda | \alpha, \eta) d\lambda = \int_{\lambda} P(\lambda | \alpha, \eta) P(t_1 | \alpha, \eta, \lambda) d\lambda$$

$$= \int_{\lambda} \frac{\eta^\alpha \lambda^\alpha e^{-\lambda(\eta+t_1)}}{\Gamma(\alpha)} d\lambda = \frac{\eta^\alpha}{\Gamma(\alpha)} \int_{\lambda} e^{-\lambda(\eta+t_1)} \lambda^\alpha d\lambda$$

$$(*) \Rightarrow (\Pi) = \frac{\lambda^\alpha e^{-\lambda(\eta+t_1)}}{\int_{\lambda} e^{-\lambda(\eta+t_1)} \lambda^\alpha d\lambda}$$

مطابق این متن باید متوجه شویم که در اینجا  $\frac{(n+t_1)^{\alpha+1}}{\Gamma(\alpha+1)}$  به دست می آید.

$$\lambda^\alpha e^{-\lambda(\eta+t_1)} \frac{(n+t_1)^{\alpha+1}}{\Gamma(\alpha+1)} = \text{Gamma}(\alpha+1, \eta+t_1)$$

معادله آمار

$$\Rightarrow \begin{cases} \alpha' = \alpha+1 \\ \eta' = \eta+t_1 \end{cases}$$



$$P(t_r | t_i) = \int_n P(t_r | t_i, d) \underbrace{P(d | t_i)}_{\text{prior}} dd \quad (1)$$

$$= \int_n \underbrace{P(t_r | t_i, d)}_{\text{prior}} \times \frac{n'^{\delta'} d^{\delta'-1} e^{-n'd}}{\Gamma(\delta')}$$

$$= \int_n d e^{-d t_r} \frac{n'^{\delta'} d^{\delta'-1} e^{-n'd}}{\Gamma(\delta')}$$

$$= \frac{n'^{\delta'}}{\Gamma(\delta')} \int_0^{\infty} d^{\delta'} e^{-d(t_r + n')} dd \quad (I)$$

$$= \frac{n'^{\delta'}}{\Gamma(\delta')} \left[ \frac{d^{\delta'} e^{-d(t_r + n')}}{-(t_r + n')} \Big|_0^{\infty} - \frac{(\delta')}{-(t_r + n')} \int_0^{\infty} d^{\delta'-1} e^{-d(t_r + n')} dd \right]$$

$$= \frac{n'^{\delta'}}{(\delta'-1)!} \times \frac{\delta'}{(t_r + n')} \int_0^{\infty} d^{\delta'-1} e^{-d(t_r + n')} dd$$

از معادله (I) و جدول 1 در پایین به مقدار متوالی رابطه بازگشتی را می توان نوشت و به

صورت دیگر می توان نوشت. (رابطه متوالی) است با مقدار احتمالی نویسنده جدول و استرال خیر - جدول دوم

$$= \frac{n'^{\delta'}}{(\delta'-1)!} \times \frac{\delta'!}{(t_r + n')^{\delta'+1}} = \frac{n'^{\delta'}}{(t_r + n')^{\delta'+1}} \times \delta' = \frac{n'^{\delta'}}{(n')^{\delta'+1}} \times \delta' = \frac{1}{n'} \times \delta'$$

$$= \frac{\delta'}{n'} \times \left( \frac{t_r}{n'} + 1 \right)^{-(\delta'+1)}$$

ج) اگر بهترین رفتار را معادله این سیستم به دست آوریم آنجا باید در حدی که بیشترین سود را بدهد

مراود به هیچ وقت بازماند و کارش را ادامه دهد.

اما وقتی رسید به این مرحله، چون مراد به با احتمال زیاد تر از حدی که بیش است حایز این

سود، پس احتمال (حرفه‌ای) را که او دارد است که وارد بازار شود.

اما وقتی رسید به این مرحله، باید بیند که در حایز حدی که بیشترین سود را بدهد.

وقتی رسید به این مرحله، با توجه به باور او به بازار (posterior) باید بیند که در حایز حدی که بیش

حایز حدی که بیش است پس در حایز حدی که بیش است  $E(t)$  در حایز حدی که بیش

$E(r)$  را نیز با توجه به posterior در حایز حدی که بیش است  $K \times E(t) < E(r)$  است

ترتیب به، در حایز حدی که بیش است با توجه به بازار دهد. توجه داریم که در حایز حدی که بیش است

بسیار است