



لطفا نکات زیر را رعایت کنید:

- سوالات خود را از طریق پست مربوط به تمرین در Quera مطرح کنید.
- پاسخ ارسالی واضح و خوانا باشد.
- در هر کدام از سوالات، اگر از منابع خاصی استفاده کرده‌اید باید آن را ذکر کنید.
- اگر با افرادی همفکری کرده‌اید، نام ایشان را ذکر کنید.
- پاسخ ارسالی باید توسط خود شما نوشته شده باشد.
- تمام پاسخ‌های خود را در یک فایل با فرمت RL_HW#[SID]_[Fullname].zip روی کوئرا قرار دهید.
- برای ارسال هر تمرین تا ساعت ۲۳:۵۹ روز ددلاین فرصت دارید. علاوه بر آن، در هر تمرین می‌توانید تا سقف ۲ روز از تأخیر مجاز باقیمانده‌ی خود استفاده کنید و در مجموع ۵ روز تأخیر مجاز برای تمارین در اختیار دارید.

سوال ۱: (نظری) نظریه اطلاعات (۲۰ نمره)

برای یک توزیع احتمال متغیری به اسم آنروپی به شکل زیر تعریف می‌شود که برای بررسی عدم قطعیت یک توزیع از آن استفاده می‌شود.

$$H(x) = - \sum_x P(x) \log P(x)$$

همچنین یکی از ابزارها برای مقایسه میزان اطلاعات بین دو توزیع آنروپی نسبی است و برای دو توزیع P, Q به صورت زیر تعریف می‌شود:

$$D_{KL}(P, Q) = - \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

اطلاعات متقابل دو متغیر تصادفی X و Y نیز به صورت زیر تعریف می‌شود:

$$I(X; Y) = H(X) - H(X|Y)$$

$$I(X; Y|Z) = H(X|Z) - H(X|Y, Z)$$

حال با توجه به اطلاعات داده‌شده به سوال‌های زیر پاسخ دهید:

الف) برای متغیرهای تصادفی X, Y, Z مثالی بیاورید که هر یک از نامساوی‌های زیر برقرار باشد:

$$I(X; Y|Z) < I(X; Y) \quad (۱)$$

$$I(X; Y|Z) > I(X; Y) \quad (۲)$$

ب) هر یک از نامساوی زیر را ثابت و شرایطی که نامساوی به مساوی تبدیل می‌شود را توضیح دهید

$$H(X, Y, Z) - H(X, Y) \leq H(X, Z) - H(X) \quad (۱)$$

$$I(X; Z|Y) \geq I(Z; Y|X) - I(Z; Y) + I(X; Z) \quad (۲)$$

ج) فرض کنید شما یک ماجراجو هستید که در جستجوی گنجی افسانه‌ای هستید که در یک جزیره دورافتاده قرار دارد. شما به تنهایی در این ماجراجویی شرکت می‌کنید و باید از بین مسیرهای مختلفی که به گنج منتهی می‌شوند، مسیر بهینه را انتخاب کنید چهار مسیر A, B, C, D برای رسیدن به گنج وجود دارد. در ابتدا هیچ اطلاعاتی در خصوص اکشن بهینه نداریم. پس از انتخاب مسیر a و دریافت پاداش حال احتمال اینکه هر کدام از مسیرها مسیر بهینه باشد به صورت زیر تغییر پیدا می‌کند:

$$P(Z = A|Y = y) = \frac{1}{1 + e^{-0.5(y-5)}}$$

$$P(Z = B|Y = y) = \frac{1 - P(Z = A|Y = y)}{2}$$

$$P(Z = C|Y = y) = \frac{1 - P(Z = A|Y = y)}{6}$$

$$P(Z = D|Y = y) = \frac{1 - P(Z = A|Y = y)}{3}$$

(۱) فرض کنید منظور از متغیر Z اکشن بهینه باشد و منظور از Y میزان پاداشی است که پس از انجام آن عمل a دریافت می شود. مقدار $I(Z; Y)$ را محاسبه کنید می توانید تابع توزیع y دلخواه را در نظر بگیرید

(۲) $I(Y; Z)$ در این مساله نشان دهنده چیست؟ در صورتی که مقدار آن بزرگ یا کوچک باشد چه مفهومی دارد؟

(د) فرض کنید X, Y, Z سه متغیر تصادفی باشند حال تابع توزیع توام این سه متغیر را به صورت $p(x, y, z)$ در نظر بگیرید و ثابت کنید رابطه زیر برقرار است. همچنین بررسی کنید در چه شرایطی مقدار $DKL(p(x, y, z) || p(x)p(y)p(z))$ برابر صفر می گردد

$$DKL(p(x, y, z) || p(x)p(y)p(z)) = -H(X, Y, Z) + H(X) + H(Y) + H(Z)$$

(ه) فرض کنید $X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow X_4$ یک زنجیره مارکوف باشند ثابت کنید:

$$I(X_1; X_3) + I(X_2; X_4) \leq I(X_1; X_4) + I(X_2; X_3)$$

پاسخ:

حالتی را در نظر بگیرید که در آن X نیز یک متغیر تصادفی باینری یکنواخت باشد. همچنین فرض کنید $Y = X$ و $Z = Y$ باشد آنگاه:

$$I(X; Y) = H(X) - H(X|Y) = H(X) = 1$$

$$I(X; Y|Z) = H(X|Z) - H(X|Y, Z) = 0$$

(b) فرض کنید X, Y هر دو یک متغیر تصادفی باینری یکنواخت و مستقل از یکدیگر باشند و $Z = X + Y$ آنگاه:

$$I(X; Y) = 0$$

$$I(X; Y|Z) = H(X|Z) = 0.5$$

(ب)

(۱) می توانیم رابطه زیر را بنویسیم:

$$H(X, Y, Z) - H(X, Y) = H(Z|X, Y) = H(Z|X) - I(Y; Z|X) \leq H(Z|X) = H(X, Z) - H(X)$$

و مشخص است در حالتی نامساوی به مساوی تبدیل می شود که $I(Y; Z|X) = 0$ باشد.

(۲) در این حالت نیز رابطه زیر برقرار است

$$I(X; Z|Y) + I(Z; Y) = I(X, Y; Z) = I(Z; Y|X) + I(X; Z)$$

$$I(X; Z|Y) = I(Z; Y|X) - I(Z; Y) + I(X; Z)$$

(ج) برای راحتی کار، Y را شامل سه مجموعه به صورت $\{0, 1, 5\}$ که هر کدام با احتمال برابر رخ می دهد، قرار دهیم.

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

حال ابتدا $H(Z)$ را محاسبه می کنیم:

$$H(Z) = - \sum_z p(z) \log p(z)$$

لازم است $P(Z)$ را محاسبه کنیم

$$P(Z = A|Y = 0) = 0.075$$

$$P(Z = A|Y = 1) = 0.1192$$

$$P(Z = A|Y = 5) = 0.5$$

$$P(Z = A) = (0.5 + 0.075 + 0.1192) \left(\frac{1}{3} \right) = 0.2317$$

به همین ترتیب سایر احتمالات را نیز محاسبه می کنیم:

$$P(Z = B) = 0.38416$$

$$P(Z = C) = 0.1280$$

$$P(Z = D) = 0.2561$$

حال $H(Z)$ برابر مقدار زیر می شود:

$$H(Z) = 0.9020$$

حال لازم است $H(Z|Y)$ را محاسبه کنیم

$$H(Z|Y) = \sum_y P(y)H(Z|Y=y) = 1.7592$$

$$I(Z;Y) = 0.1427$$

$I(Z;Y)$ نیز میزان اطلاعات بدست آمده نسبت به مسیر بهینه براساس پاداش دریافتی را نشان می دهد. اگر $I(Z;Y)$ بزرگ باشد، دانستن پاداش دریافتی اطلاعات زیادی در مورد مسیر بهینه به ما می دهد و بالعکس.
(د) رابطه را به صورت زیر می نویسیم

$$\begin{aligned} D(p(x,y,z)||p(x)p(y)p(z)) &= E \left[\log \frac{p(x,y,z)}{p(x)p(y)p(z)} \right] \\ &= E [\log p(x,y,z)] - E [\log p(x)] - E [\log p(y)] - E [\log p(z)] \\ &= -H(X,Y,Z) + H(X) + H(Y) + H(Z) \end{aligned}$$

و همچنین تنها در شرایطی برابر صفر می شود که $p(x,y,z)=p(x)p(y)p(z)$
(ه) روابط را به صورت زیر بسط می دهیم

$$\begin{aligned} &I(X_1;X_4) + I(X_2;X_3) - I(X_1;X_3) - I(X_2;X_4) \\ &= H(X_1) - H(X_1|X_4) + H(X_2) - H(X_2|X_3) \\ &\quad - (H(X_1) - H(X_1|X_3)) - (H(X_2) - H(X_2|X_4)) \\ &= H(X_1|X_3) - H(X_1|X_4) + H(X_2|X_4) - H(X_2|X_3) \\ &= H(X_1,X_2|X_3) - H(X_2|X_1,X_3) - H(X_1,X_2|X_4) + H(X_2|X_1,X_4) \\ &\quad + H(X_1,X_2|X_4) - H(X_1|X_2,X_4) - H(X_1,X_2|X_3) + H(X_1|X_2,X_3) \\ &= -H(X_2|X_1,X_3) + H(X_2|X_1,X_4) - H(X_2|X_1,X_4) - H(X_2|X_1,X_3,X_4) \\ &= I(X_2;X_3|X_1,X_4) \geq 0 \end{aligned}$$

سوال ۲: (نظری) بهینه سازی (۱۵ نمره)

هدف از بهینه سازی پیدا کردن نقطه یا نقاطی است که مقدار یک تابع را روی فضای مشخصی کمینه یا بیشینه کند. فرض کنیم می خواهیم مقدار تابع f را کمینه کنیم

$$x^* = \operatorname{argmin}_{x \in U} f(x) \quad (۱)$$

با فرض اینکه تابع f کران پایین دارد و مشتق پذیر است و U نیز کران دار است، $x^* \in \partial(U) \cup \operatorname{Cr}(f,U)$ که $\partial(U)$ نقاط مرزی U می باشد و $\operatorname{Cr}(f,U)$ نقاط بحرانی f در U می باشند که مشتق f در این نقاط صفر می شود.

از بین تمام نقاط بحرانی، نقاطی که مقدار ماتریس هسیان^۱ در آن ها مثبت معین می شود کمینه محلی می باشند. اما همچنان کمینه بودن آن ها در کل فضای U مشخص نیست. این موضوع برای توابع محدب ساده تر است. زیرا در صورت محدب بودن f هر نقطه بحرانی یک نقطه کمینه برای f است.

برای مسائل بهینه سازی با قید نیز می توان از روش ضرایب لاگرانژ استفاده کرد. فرض کنید با فرض $g(x) = 0$ می خواهیم تابع $f(x)$ را کمینه کنیم. کافیهست معادله $\nabla f(x) = \lambda \nabla g(x)$ با فرض $g(x) = 0$ را حل کنیم. جواب مسئله اصلی داخل جواب های معادله جدید می باشد. در واقع مسئله مقید اولیه به مسئله بدون قید $f(x) + \lambda g(x)$ تبدیل می شود.

روش ضرایب لاگرانژ حالت خاص روشی برای حل مسائل بهینه سازی است که بر پایه مفهومی به نام مسئله دوگان استوار است. فرض کنید علاوه بر قیدهای برابر، قیدهایی به صورت نابرابری نیز وجود دارند

$$\min_x f(x) \quad \text{s.t.} \quad g_i(x) \leq 0, \quad h_j(x) = 0 \quad (۲)$$

دوگان این مسئله به صورت زیر تعریف می شود،

$$\max_{\mu, \lambda} L(\mu, \lambda) \quad \text{s.t.} \quad \mu \geq 0 \quad (۳)$$

$$L(\mu, \lambda) = \inf_x \{f(x) + \mu^T g(x) + \lambda^T h(x)\} \quad (۴)$$

^۱ Hessian Matrix

عبارت دوگان یک کران پایین برای تابع اصلی بهینه سازی ارائه می دهد و حل مسئله دوگان لزوماً منجر به حل مسئله اصلی نمی شود. در صورت برقراری شرایط KKT جواب مسئله دوگان معادل با جواب مسئله اصلی می شود. یکی از ویژگی های فرم دوگان این است که تابع L یک تابع مقعر است و بیشینه سازی آن همانند کمینه سازی یک تابع محدب ساده تر می باشد. گاهی در بهینه سازی به جای پیدا کردن نقاط بهینه به دنبال توابع بهینه هستیم. حل چنین مسائلی موضوع حوزه حساب تغییرات^۲ می باشد. برابری لاگرانژ اویلر یک شرط لازم برای بهینه بودن تابع معرفی می کند. فرض کنیم $y = y(x)$ تابعی از x است و هدف کمینه یا بیشینه کردن $\int F(x, y, y') dx$ است که F یک تابع دلخواه مشتق پذیر است. برای $y(x)$ بهینه داریم،

$$\frac{d}{dx} \frac{\partial F}{\partial y'} = \frac{\partial F}{\partial y} \quad (5)$$

(آ) نشان دهید در صورت محدب بودن f ، هر نقطه بحرانی، یک نقطه کمینه می باشد.

(ب) کمینه و بیشینه تابع $f(x, y, z) = x^2 + y^2 + z^2$ را با شرایط $z^2 = x^2 + y^2$ و $z = x + y + 1$ بیابید.

(ج) ثابت کنید

$$L(\mu, \lambda) \leq \min_x f(x)$$

و مثالی ارائه دهید که برای هیچ μ, λ ای رخ ندهد.

(د) مسئله بهینه سازی زیر را حل کنید،

$$\min_x \left\{ - \sum_{i=1}^n \log(\alpha_i + x_i) \right\} \quad s.t. \quad x \geq 0, 1^T x = 1 \quad (6)$$

(ه) توزیع احتمال با بیشینه آنتروپی پیوسته ای را بیابید که میانگین و واریانس مشخصی داشته باشد،

$$\operatorname{argmax}_P \left\{ - \int P(x) \log P(x) dx \right\} \quad \mathbb{E}_P[X] = \mu, \quad \mathbb{V}_P(X) = \sigma^2 \quad (7)$$

سوال ۳: (نظری) زنجیره مارکوف (۱۵ نمره)

یکی از کاربردهای زنجیره مارکوف استفاده از آن‌ها در مسائلی است که نمونه برداری با یک توزیع احتمال مشخص در آن‌ها دشوار می باشد در این تمرین یکی از این مسائل را بررسی می کنیم.

وقتی می گوئیم یک زنجیره مارکوف غیر قابل کاهش است منظور این است که تمامی وضعیت‌ها به یکدیگر دسترسی داشته باشند به عبارتی دیگر رابطه زیر برقرار باشد:

$$\forall i, j \in S \exists n \geq 0 \text{ such that } P(X_n = j | X_0 = i) > 0$$

حالت i دارای دوره تناوب k است اگر هر مسیر بازگشت به حالت i به طول مضارب k باشد. به زبان دیگر، دوره تناوب یک حالت برابر است با

$$k = \gcd\{n : \Pr(X_n = i | X_0 = i) > 0\}$$

حال در صورتی که دوره تناوب تمام حالت‌ها برابر یک باشد می گوئیم زنجیره مارکوف بدون تناوب است
الف) زنجیره مارکوفی را در نظر بگیرید که در آن به ازای هر وضعیت x و y اگر $x \neq y$ و $y \in N(x)$ باشد آنگاه:

$$P_{(x,y)} = 1/M \text{ و در غیر این صورت } P_{(x,y)} = 0 \text{ همچنین } P_{(x,x)} = 1 - \frac{N(x)}{M}$$

$N(x)$: set of neighbors of x

$$M \geq \max_{x \in \Omega} |N(x)|.$$

حال اگر این زنجیره مارکوف غیر قابل کاهش و بدون تناوب باشد، ثابت کنید توزیع مانا^۳ این زنجیره مارکوف یک توزیع یکنواخت است.

ب) مسئله شماردن تعداد حالت‌های مختلف در یک کوله پشتی را در نظر بگیرید. در این مسئله تعدادی شی در اختیار داریم که وزن هر کدام به صورت $a_1, a_2, a_3, \dots, a_n$ می باشد و وزن کل کوله پشتی نیز b می باشد. هدف تخمین تعداد بردارهایی به صورت $(x_1, x_2, x_3, \dots, x_n) \in \{0, 1\}^n$ طوری که $\sum_{i=1}^n a_i x_i < b$ است.

حال فرض کنید برای حل این مسئله به این صورت عمل می کنیم که به صورت یکنواخت یک بردار

^۲Variational Calculus
^۳stationary distribution

$\{0, 1\}^n \ni (x_1, x_2, x_3, \dots, x_n)$ را انتخاب کرده و این کار را به تعداد زیادی تکرار (N) می‌کنیم. تعداد حالت‌هایی که در آن شرط کوله پشتی نقض نمی‌شود را می‌شماریم و نام این مجموعه را p می‌نامیم. در نهایت مقدار $\frac{|P| \times 2^n}{N}$ به عنوان جواب نهایی برمی‌گردانیم این روش چه اشکالی دارد؟
(ج) حال فرض کنید برای مسئله بالا زنجیره مارکوفی را به صورت زیر می‌سازیم که در آن هر وضعیت X_j را به این صورت $(x_1, x_2, x_3, \dots, x_n)$ تعریف می‌کنیم. حال در هر گام زنجیره مارکوف یک متغیر $i \in [1, n]$ را به صورت یکنواخت انتخاب می‌کند: اگر $x_i = 1$ بود آن را صفر می‌کند و به حالت X_{j+1} می‌رود و در صورتی که $x_i = 0$ باشد آن را یک می‌کند در این حالت در صورتی که با یک کردن x_i شرط کوله پشتی نقض نشود به حالت X_{j+1} می‌رود و در غیر این صورت $X_j = X_{(j+1)}$

(۱) نشان دهید در صورتی که $\sum_{i=1}^n a_i < b$ آنگاه این زنجیره مارکوف دارای یک توزیع مانا به صورت توزیعی یکنواخت است.

(۲) توضیح دهید چگونه می‌توان از این زنجیره مارکوف جهت تخمین تعداد حالات برای مسئله کوله پشتی استفاده کرد.

پاسخ:

(الف) اولاً نشان می‌دهیم در یک زنجیره مارکوف محدود، غیرقابل کاهش و بدون تناوب با n وضعیت و ماتریس انتقالی P ، اگر $\pi = (\pi_0, \dots, \pi_n)$ وجود داشته باشد که $\sum_{i=0}^n \pi_i = 1$ و برای هر i, j رابطه زیر برقرار باشد:

$$\pi_i P_{(i,j)} = \pi_j P_{(j,i)}$$

آنگاه π توزیع پایا و یکتا این زنجیره مارکوف می‌باشد.
اثبات:

$$\sum_{i=0}^n \pi_i P_{i,j} = \sum_{i=0}^n \pi_j P_{j,i} = \pi_j$$

از این رو چون π در رابطه $\pi P = \pi$ و $\sum_{i=0}^n \pi_i = 1$ صدق می‌کند آنگاه π باید حتماً توزیع پایا و یکتا این زنجیره مارکوف باشد. حال در زنجیره مارکوف بیان شده طبق صورت مسئله برای $x \neq y$ اگر $\pi_y \pi_x = 0$ آنگاه:

$$\pi_x P_{(x,y)} = \pi_y P_{(y,x)}$$

حال در صورتی که π یک توزیع یکنواخت باشد طوری که $\sum_{i=0}^n \pi_i = 1$ با توجه به اینکه زنجیره مارکوف محدود، غیرقابل کاهش و بدون تناوب است آنگاه می‌توان گفت π توزیع مانا و یکتا این زنجیره می‌باشد.

(ب) مشکل این روش این است که توزیع حالت‌های مختلف این مسئله مشخص نیست و نمی‌توان به آسانی از یک توزیع مشخص جواب این مسئله را نمونه برداری کرد و در برخی حالات نیاز به تعداد نمایی نمونه برداری می‌باشد تا تخمین بدست آمده تخمین مناسبی باشد. به طور مثال، در حالتی که یک آیتم دارای سایز $1-b$ و سایر آیتم‌ها دارای سایز $2b$ باشد، تنها یک جواب ممکن برای مسئله وجود دارد و برای تخمین این عدد حدود 2^n نمونه برداری لازم است

(ج) با توجه به تعریف P واضح است که اگر گراف معادل این زنجیره مارکوف را رسم کنیم، تنها $X = (x_1, x_2, \dots, x_n)$ هایی در این گراف وجود دارند که شرط $\sum_{i=1}^n a_i x_i < b$ را ارضا می‌کنند. اگر دو گره X و Y که تنها در بیت i ام اختلاف دارند را در نظر بگیریم که در آن $x_i = 1$ و $y_i = 0$ ، مشخص است که یالی از Y به X وجود دارد. همچنین با توجه به اینکه یک کردن y_i منجر به جود آمدن X می‌شود و X نیز شرط $\sum_{i=1}^n a_i x_i < b$ را ارضا می‌کند، واضح است یالی از X به Y هم وجود دارد یا به عبارت دیگر یال‌های گراف دو طرفه هستند. حالت $z \in 0^n$ را در نظر بگیریم. از همه حالات X با تغییر x_i هایی که دارای مقدار ۱ هستند می‌توان به حالت z رسید و با توجه به اینکه گراف شرط $\sum_{i=1}^n a_i x_i < b$ را ارضا می‌کند، در نتیجه از z نیز به همه حالات Y می‌توان رسید. در نتیجه گراف همبند است. همچنین می‌توان نشان داد گراف حتماً دارای طوقه نیز می‌باشد زیرا اگر حالت بهینه کوله پشتی مانند w را در نظر بگیریم، با یک کردن هر آیتمی که برابر صفر می‌باشد (دقت شود با توجه به فرض $\sum_{i=1}^n a_i > b$) حتماً آیتمی وجود دارد که برداشته نشده باشد، حتماً طبق تعریف مسئله به خودش برمی‌گردد. حال با توجه به اینکه زنجیره مورد نیاز غیرقابل کاهش و بدون تناوب است و تابع انتقال مانند حالت الف می‌باشد توزیع مانا این زنجیره یکنواخت است. در نتیجه به استفاده از این زنجیره می‌توان به صورت یکنواخت از این مسئله نمونه برداری کرد و مشکل قسمت ب حل می‌شود زیرا می‌توان هر جواب valid را با احتمال برابر نمونه برداری کرد.

سوال ۴: (نظری) آمار بیزی (۱۵ نمره)

در یک دهه پرفراز و نشیب، جزیره ای به نام "آرکادیا" وجود دارد. مردم این جزیره تصمیم می‌گیرند زندگی را از سر بگیرند و یک شرکت به نام "Reminisce" راه اندازی که با استفاده از فناوری‌های جدید، یادگیری را بازیابی کند.

شرکت Reminisce یک فناوری به نام Mirror Memory را توسعه می‌دهد. این سیستم، به شما اجازه می‌دهد تا لحظات خوب و خاطره‌انگیز زندگی‌تان را بازگو کنید. اما این سیستم نیازمند شناسایی افراد است تا بتواند به درستی خاطرات را به هر فرد متصل کند.

سیستم Mirror Memory شامل یک دوربین و نرم‌افزاری است که ورودی اتاق را کنترل می‌کند. زمانی که شما می‌خواهید هویت خود را به سیستم ثبت کنید، ابتدا عکسی از خودتان گرفته می‌شود. سپس این عکس با گالری عکس‌های ذخیره شده در سیستم مقایسه می‌شود.

حال برای هر عکس موجود در گالری یک امتیاز بین صفر و یک در نظر گرفته می‌شود:

اگر عکس گرفته شده و عکس گالری نشان دهنده همان فرد باشد، امتیاز از توزیع زیر محاسبه می‌شود:

$$p(s|same) = \alpha_{ss} \exp(\lambda_s s)$$

اما اگر این عکس گالری نشان دهنده فردی متفاوت باشند، امتیاز از توزیع زیر محاسبه می شود

$$p(s|different) = \alpha_{ds} \exp(-\lambda_d s)$$

الف) توضیح دهید آیا توابع در نظر گرفته شده برای توزیع امتیازها منطقی می باشند یا خیر؟ علت استفاده از ضرایب α_{ds} و α_{ss} چیست؟
 ب) فرض کنید N عکس از N فرد مختلف در سیستم موجود باشد. حال با فرض داشتن N امتیاز مشخص به صورت $s_1, s_2, \dots, s_j, \dots, s_N$ احتمال اینکه تصویر j ، تصویر فرد مورد نظر باشد چقدر است؟ (فرض کنید s_i به صورت صعودی داده شده اند و احتمال $prior$ همه اشخاص مشابه می باشد.)

ج) حال فرض کنید هیچ امتیاز مشخصی در اختیار نداریم. سیستم صرفاً عکسی را برمی گرداند که بیشترین امتیاز را دارد. احتمال اینکه این عکس فرد درست تشخیص داده شود چقدر است؟ برای این موضوع یک فرمول عمومی ارائه دهید

پاسخ:

الف) بله، توابع در نظر گرفته شده برای امتیازها منطقی می باشند. زیرا همانطور که مشخص است، در صورتی که افراد مشابه باشند، شما با احتمال زیاد امتیاز بالایی می گیرید و در صورتی که افراد متفاوت باشند، با احتمال زیاد امتیاز خیلی کمی گرفته می شود. با توجه به اینکه جمع احتمال ها باید برابر یک باشد، جهت نرمال سازی مقادیر امتیازها نیز از ضرایب α_s و α_d استفاده شده است.

ب)

$$\begin{aligned} P(same_j, different_{i \neq j} | s_1, \dots, s_N) &= \frac{p(s_1, \dots, s_N | same_j, different_{i \neq j})}{p(s_1, \dots, s_N)} \\ &= \frac{p(s_j | same) \left(\prod_{i \neq j} p(s_i | different) \right) \left(\frac{1}{N} \right)}{p(s_1, \dots, s_N)} \\ &= \frac{\alpha_s \exp(\lambda_s s_j) \prod_{i \neq j} \alpha_d \exp(-\lambda_d s_i) \left(\frac{1}{N} \right)}{p(s_1, \dots, s_N)} \\ &= \frac{\alpha_s \alpha_d^{(N-1)} \exp(\lambda_s s_j) \exp(-\lambda_d (s_j - \sum_s s_i)) \left(\frac{1}{N} \right)}{p(s_1, \dots, s_N)} \\ &= \frac{\alpha_s \alpha_d^{(N-1)} \exp((\lambda_s + \lambda_d) s_j) \exp(-\lambda_d (\sum_s s_i)) \left(\frac{1}{N} \right)}{\sum_{j'} \alpha_s \alpha_d^{(N-1)} \exp((\lambda_s + \lambda_d) s_{j'}) \exp(-\lambda_d (\sum_s s_i)) \left(\frac{1}{N} \right)} \\ &= \frac{\exp((\lambda_s + \lambda_d) s_j)}{\sum_{j'} \exp((\lambda_s + \lambda_d) s_{j'})} \end{aligned}$$

ج)

$$\begin{aligned} p(correct\ recognition) &= \int_0^1 p(s' | same) P(s < s' | different)^{(N-1)} ds' \\ &= \int_0^1 \alpha_s \exp(\lambda_s s') \left(\frac{\exp(-\lambda_d s') - 1}{\exp(-\lambda_d) - 1} \right)^{(N-1)} ds' \\ &= \int_0^1 \frac{\lambda_s \exp(\lambda_s s')}{\exp(\lambda_s) - 1} \left(\frac{\exp(-\lambda_d s') - 1}{\exp(-\lambda_d) - 1} \right)^{(N-1)} ds' \\ &= \frac{1}{(\exp(\lambda_s) - 1) (\exp(-\lambda_d) - 1)^{(N-1)}} \int_0^1 \exp(\lambda_s s') (\exp(-\lambda_d s') - 1)^{(N-1)} ds' \end{aligned}$$

سوال ۵: (نظری) تئوری تخمین (۲۰ نمره)

فرض کنید متغیر تصادفی $X: \mathcal{X} \rightarrow \mathbb{R}$ از توزیع P_θ پیروی می کند که θ یک پارامتر نامعلوم است. هدف از تخمین پیدا کردن θ از روی نمونه های مشاهده شده از X است،

$$X = X_1, X_2, \dots, X_n \sim P_\theta(X) \quad (۸)$$

به هر تابعی مانند $W = W(X_1, X_2, \dots, X_n)$ یک تخمینگر می گوئیم. تخمینگرها بر اساس ویژگی های مختلفی که دارند مقایسه و بررسی می شوند. هدف اصلی یک تخمینگر، ارائه یک تخمین مناسب از پارامتر θ می باشد. با توجه به اینکه خروجی تخمینگر خود یک متغیر تصادفی است، برای مقایسه آن با پارامتر حقیقی از معیار میانگین مجذور خطا^۴ یا MSE استفاده می شود،

$$MSE(W(X), \theta) = \mathbb{E}[(W(X) - \theta)^2] = (\mathbb{E}[W(X)] - \theta)^2 + \mathbb{E}[(W(X) - \mathbb{E}[W(X)])^2] \quad (۹)$$

Mean Squared Error^۴

که بخش اول مجذور بایاس و بخش دوم واریانس تخمینگر است. به یک تخمینگر نااریب می گویند اگر بایاس آن صفر باشد. یکی از تخمینگرهای معروف تخمینگر بیشینه درست نمایی می باشد،

$$W(X) = \hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} P(X_1, \dots, X_n | \theta) = \operatorname{argmax}_{\theta} \mathcal{L}(\theta; X) \quad (10)$$

(آ) فرض کنید که $\theta = (\mu, \sigma^2)$ و $P_{\theta}(X) = \mathcal{N}(X; \mu, \sigma^2)$. تخمینگر بیشینه درست نمایی را برای θ پیدا کنید. آیا این تخمینگر نااریب است؟

(ب) برای تابع دلخواه τ نشان دهید اگر $\hat{\theta}$ تخمین بیشینه درست نمایی از θ باشد، $\tau(\hat{\theta})$ تخمین بیشینه درست نمایی برای $\tau(\theta)$ است. سپس تخمین بیشینه درست نمایی را برای $\frac{1}{\mu^2+1}$ و σ با توجه به تعریف های قسمت قبل بدست آورید.

(ج) یکی از کاربردهای تخمین، بدست آوردن امیدریاضی تابعی از یک متغیر تصادفی است. فرض کنید برای متغیر تصادفی $X \sim p(X)$ و تابع $f(x)$ درصد تخمین امیدریاضی این تابع تحت توزیع p هستیم،

$$\mu = \mathbb{E}_{X \sim P(X)}[f(X)] \quad (11)$$

فرض کنید نمونه های x_1, \dots, x_n از توزیع متفاوت $q(X)$ داریم. در روش Importance-Weighted Sampling (IS) برای تخمین μ از تخمینگر زیر استفاده می شود،

$$\hat{\mu}_{IS} = \frac{1}{n} \sum_{i=1}^n f(x_i) \frac{p(x_i)}{q(x_i)} = \frac{1}{n} \sum_{i=1}^n w_i f(x_i) \quad (12)$$

در واقع به جای میانگین ساده، میانگین وزن دار $f(x_i)$ ها محاسبه می شود.

۱. نشان دهید تخمینگر IS نااریب است.

۲. عبارتی برای واریانس تخمینگر IS پیدا کنید. توزیع q را برحسب p و f طوری انتخاب کنید که واریانس $\hat{\mu}_{IS}$ کمینه شود و مقدار کمینه را بیابید. همچنین نشان دهید p و q را می توان طوری تعیین کرد که واریانس این تخمینگر هر اندازه ای بزرگ شود.

۳. مشکل بزرگ شدن واریانس تخمینگر IS در نهایت باعث بزرگ شدن MSE می شود و صرف نااریب بودن، باعث مطلوب بودن این تخمینگر نمی شود. برای حل این مشکل، تخمینگر Normalized Importance-Weighted Sampling (N-IS) ارائه شده است. این تخمینگر به جای تعداد نمونه ها روی مجموع وزن ها نرمال می شود.

$$\hat{\mu}_{N-IS} = \frac{\sum_{i=1}^n w_i f(x_i)}{\sum_{i=1}^n w_i} \quad (13)$$

نشان دهید تخمینگر N-IS نااریب نیست اما

$$\lim_{n \rightarrow \infty} \hat{\mu}_{N-IS} = \mu \quad (14)$$

در نتیجه N-IS در بی نهایت نااریب است.

۴. با فرض کراندار بودن $f(x)$ نشان دهید،

$$\mathbb{V}(\hat{\mu}_{N-IS}) \leq \frac{(M-m)^2}{4} \quad (15)$$

که

$$M = \sup f(x) \quad m = \inf f(x) \quad (16)$$

در نتیجه واریانس تخمینگر N-IS کران بالایی مستقل از انتخاب p و q دارد.

۵. یکی از کاربردهای IS تخمین مقدار انتگرال است. فرض کنید می خواهیم انتگرال زیر را تخمین بزنیم:

$$H = \int_2^{\infty} e^{-\frac{x^2}{2}} dx \quad (17)$$

فرض کنید نمونه هایی از توزیع نرمال استاندارد به صورت زیر داریم:

$$S = \{0.26, -0.27, -1.56, 0.41, 0.40, -0.02, 0.10, -1.69, -0.28, -2.53\} \quad (18)$$

آ. ابتدا به روش monte-carlo تخمینی برای H بیابید.

ب. حال به روش IS و در نظر گرفتن توزیع $q = \mathcal{N}(3, 1)$ تخمین دیگری از H بیابید. مقدار این تخمین و تخمین بخش قبل را با مقدار واقعی H مقایسه کنید. آیا به تخمین بهتری برای H می‌رسیم؟ علت آن چیست؟

سوال ۶: (نظری) Variational Inference (۱۵ نمره)

یکی از چالش‌های موجود در استنباط آماری، تخمین توزیع پسین متغیرهای پنهان می‌باشد. فرض کنیم Z متغیر پنهان مسئله و X متغیر مشاهده شده باشد. توزیع پیشین Z ، $P(Z)$ و توزیع درست‌نمایی $P(X|Z)$ در مسئله داده شده‌اند. توزیع پسین Z با قاعده بیز به صورت زیر بدست می‌آید،

$$P(Z|X) = \frac{P(X|Z)P(Z)}{\int_{Z'} P(X|Z')P(Z')} \quad (۱۹)$$

محاسبه انتگرال مخرج در اکثر مواقع امکان‌پذیر نیست. از Variational Inference برای تخمین توزیع پسین استفاده می‌شود. برای این کار سعی می‌شود توزیع پسین با یک توزیع ساده شده $q(Z)$ تخمین زده شود. فاصله KL میان دو توزیع دلخواه p, q به صورت زیر تعریف می‌شود

$$D_{KL}(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx \quad (۲۰)$$

فاصله KL معیاری برای تفاوت دو توزیع است برای تخمین زدن توزیع پسین، می‌توان $D_{KL}(q(Z)||P(Z|X))$ را کمینه کرد که با بیشینه کردن $L(q)$ معادل است. حال فرض کنید توزیع q به تعدادی مولفه تجزیه شود،

$$q(Z) = q(Z_1, Z_2, \dots, Z_n) = q_1(Z_1)q_2(Z_2)\dots q_n(Z_n) \quad (۲۱)$$

(آ) نشان دهید فاصله KL در $p = q$ مقدار کمینه خود را که صفر است، می‌گیرد.

(ب) نشان دهید

$$D_{KL}(q(Z)||P(Z|X)) = \log P(X) - \mathbb{L}(q) = \log P(X) - (\mathbb{E}_q[\log P(X, Z)] - \mathbb{E}_q[\log q(Z)]) \quad (۲۲)$$

(ج) نشان دهید $\argmax_q \mathcal{L}(q)$ در معادلات زیر صدق می‌کند،

$$\forall 1 \leq i \leq n : \log q_i(Z_i) = \mathbb{E}_{q_j, j \neq i} [\log P(X, Z)] + const. \quad (۲۳)$$

(د) فرض کنید برای $1 \leq n \leq N$ ، $\mathbf{z}_n \in \mathbb{R}^K$ یک بردار K -of-1 دودویی تصادفی است و $\sum_{i=1}^K \pi_i = 1$ متغیرهای تصادفی با مجموع یک هستند. متغیر X مشاهده شده و متغیرهای Z, π, μ, Λ پنهان هستند. همچنین داریم:

$$\begin{aligned} P(X|Z, \mu, \Lambda) &= \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{x}_n | \mu_k, \Lambda_k^{-1})^{z_{nk}} \\ P(Z|\pi) &= \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \\ P(\pi) &\propto \sum_{k=1}^K \pi_k^{\alpha_0 - 1} \\ P(\mu, \Lambda) &= \prod_{k=1}^K \mathcal{N}(\mu_k | \mathbf{m}_0, (\beta_0 \Lambda_k)^{-1}) \mathcal{W}(\Lambda_k | \mathbf{W}_0, \nu_0) \end{aligned}$$

که \mathcal{W} توزیع ویشارت^۵ است. تخمین Variational را برای توزیع پسین متغیرهای پنهان با فرض زیر بدست آورید. الگوریتم EM برای انجام این تخمین را توصیف کنید.

$$q(\mathbf{Z}, \pi, \mu, \Lambda) = q(\mathbf{Z})q(\pi, \mu, \Lambda) \quad (۲۴)$$

^۵Wishart