

نظریه‌ی اطلاعات، آمار و یادگیری (۱-۲۵۱۱۰)



تمرین سری سوم

ترم بهار ۱۴۰۲-۰۳

دانشکده‌ی مهندسی برق

دانشگاه صنعتی شریف

استاد: دکتر محمدحسین یاسائی میبدی

مهلت تحویل: جمعه ۴ خرداد ۱۴۰۳ ساعت ۲۳:۵۹

(*) مسائلی که با ستاره مشخص شده‌اند امتیازی هستند و حل کردن آن‌ها نمره‌ی امتیازی خواهد داشت!

۱ TV و ارتباط آن با برخی انحراف‌ها

در این سوال به بررسی برخی از ویژگی‌های انحراف TV می‌پردازیم. می‌توانید به دلخواه به سه قسمت پاسخ دهید و مابقی نمره امتیازی خواهد داشت.

۱. فرض کنید P و Q دو توزیع احتمال بر روی $X_{1:n} = (X_1, \dots, X_n) \in \mathcal{X}^n$ باشند. همچنین فرض کنید $P_i(\cdot | x_{1:i-1})$ توزیع احتمال شرطی متغیر X_i به شرط $X_{1:i-1} = x_{1:i-1}$ باشد $Q_i(\cdot | x_{1:i-1})$ را نیز به طور مشابه در نظر بگیرید). نشان دهید:

$$\|P - Q\|_{TV} \leq \sum_{i=1}^n \mathbb{E}_{X_{1:i-1} \sim P} [\|P_i(\cdot | X_{1:i-1}) - Q_i(\cdot | X_{1:i-1})\|_{TV}],$$

که در آن امید ریاضی بر روی متغیر $X_{1:i-1}$ برحسب توزیع P گرفته می‌شود.

۲. نامساوی Bretagnolle-Huber: ثابت کنید برای هر دو توزیع P و Q داریم:

$$\|P - Q\|_{TV} \leq \sqrt{1 - \exp(-D_{KL}(P\|Q))} \leq 1 - \frac{1}{4} D_{KL}(P\|Q).$$

۳. برای هر دنباله از توزیع‌های P_n و Q_n نشان دهید هنگامی که $n \rightarrow \infty$ داریم:

$$d_{TV}(P_n^{\otimes n}, Q_n^{\otimes n}) \rightarrow 0 \quad \Leftrightarrow \quad D_{H^r}(P_n, Q_n) = o\left(\frac{1}{n}\right),$$

$$d_{TV}(P_n^{\otimes n}, Q_n^{\otimes n}) \rightarrow 1 \quad \Leftrightarrow \quad D_{H^r}(P_n, Q_n) = \omega\left(\frac{1}{n}\right),$$

که در آن $D_{H^r}(\cdot, \cdot)$ فاصله‌ی Hellinger است.

۴. فرم وردشی زیر را برای انحراف TV ثابت کنید:

$$d_{TV}(P_1, P_2) = \frac{1}{2} \inf_q \sqrt{\int_{x \in \mathcal{X}} \frac{(p_1(x) - p_2(x))^2}{q(x)} dx}.$$

راهنمایی: از نامساوی کوشی-شوارتز استفاده کنید.

۵. فرض کنید $P_{Y|X}$ یک کانال با ورودی باینری باشد. قرار دهید $P_0 = P_{Y|X=0}$, $P_1 = P_{Y|X=1}$. ثابت کنید

$$\frac{1}{4} d_{TV}^2(P_0, P_1) \leq I(X; Y) \leq d_{TV}(P_0, P_1).$$

راهنمایی: برای سمت چپ از نامساوی Pinsker استفاده کنید و برای سمت راست از نامساوی بین اطلاعات متقابل و χ^2 استفاده کنید.

۲ نشت اطلاعات

فرض کنید $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ یک گراف ساده‌ی بدون جهت متناهی باشد.

متغیرهای تصادفی $\{X_v : v \in \mathcal{V}\} \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(\frac{1}{2})$ را روی رئوس این گراف و متغیرهای تصادفی $\{Z_e : e \in \mathcal{E}\} \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(\delta)$ را بر روی یال‌های آن تعریف می‌کنیم. حال برای هر یال $e = (u, v) \in \mathcal{E}$ تعریف کنید $Y_e = X_u \oplus X_v \oplus Z_e$. مدل نشت روی گراف $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ یک اندازه‌ی احتمال روی تمام زیرگراف‌های \mathcal{G} است که در آن احتمال حضور هر کدام از یال‌های \mathcal{G} به طور مستقل برابر p باشد. چنین اندازه‌ی احتمالی را با $\mathbb{P}_{(\mathcal{G}, p)}$ نمایش می‌دهیم. همینطور پیشامد وجود مسیر بین دو زیر مجموعه‌ی $\mathcal{S}, \mathcal{S}'$ از رئوس را با $(\mathcal{S} \rightsquigarrow \mathcal{S}')$ نمایش می‌دهیم. در این سوال قصد داریم قضیه‌ی زیر را ثابت کنیم:

قضیه ۱-۲. برای هر زیرمجموعه از رئوس مانند $\mathcal{S} \subset \mathcal{V}$ و هر رأس مانند $v \in \mathcal{V}$ داریم:

$$I(X_v; X_{\mathcal{S}}, Y_{\mathcal{E}}) \leq \mathbb{P}_{(\mathcal{G}, \eta)}[v \rightsquigarrow \mathcal{S}] \log(2),$$

که در آن $\eta = (1 - 2\delta)^2$. همچنین منظور از $X_{\mathcal{S}}$ مجموعه‌ی $\{X_u : u \in \mathcal{S}\}$ است.

برای اثبات، ابتدا باید با نامساوی قوی پردازش داده‌ها آشنا شوید. اگر $U \rightarrow X \rightarrow Y$ یک زنجیره‌ی مارکف باشد، از نامساوی پردازش داده‌ها می‌دانیم: $I(U; Y) \leq I(U; X)$ حال اگر $P_{Y|X}$ ثابت باشد، می‌توانیم ضریب $\eta_{P_{Y|X}}$ را به صورت زیر تعریف کنیم:

$$\eta_{P_{Y|X}} = \sup_{P_{U,X}} \frac{I(U; Y)}{I(U; X)}$$

در این صورت برای این زنجیره‌ی مارکف همواره خواهیم داشت:

$$I(U; Y) \leq \eta_{P_{Y|X}} I(U; X).$$

می‌توان دید اگر $P_{Y|X}$ یک کانال دوتایی متقارن^۱ با پارامتر δ باشد، داریم: $\eta_{P_{Y|X}} = (1 - 2\delta)^2$.

۱. ثابت کنید: $I(X_v; Y_{\mathcal{E}}) = 0$.

۲. با استقرا روی $|\mathcal{E}|$ حکم مسئله را نتیجه بگیرید.

راهنمایی: برای گام استقرا از شرطی کردن اطلاعات متقابل استفاده کنید.

۳. (*) فرض کنید \mathcal{T} یک درخت منتظم با ریشه‌ی ρ باشد، که در آن درجه‌ی هر رأس $(d + 1)$ است. همین‌طور فرض کنید هر یال این درخت یک کانال $\text{BSC}(\delta)$ باشد. ابتدا یک بیت با توزیع $X_{\rho} \sim \text{Bernoulli}(\frac{1}{2})$ روی ریشه‌ی درخت تولید می‌شود. سپس این بیت از طریق کانال‌هایی که روی یال‌ها قرار دارند به سمت پایین انتشار می‌یابد. اگر \mathcal{S}_k مجموعه‌ی رئوس در عمق k از این درخت باشد، ثابت کنید اگر $d(1 - 2\delta)^2 < 1$ داریم:

$$d_{TV}(X_{\rho}, X_{\mathcal{S}_k}) \xrightarrow[k \rightarrow \infty]{} 0.$$

۳ تبخّر در اثبات نامساوی‌ها!

۱. فرض کنید $\infty \cup \mathbb{R}_{\geq 0} : [0, \infty] \mapsto f$ تابعی محدب باشد که $f(1) = f'(1) = 0$. همین‌طور فرض کنید μ, ν دو اندازه‌ی احتمال روی مجموعه‌ی \mathcal{X} باشند. ثابت کنید برای $M > 1$ داریم

$$\nu \left(\frac{d\nu}{d\mu} > M \right) \leq \frac{D_f(\nu || \mu)}{f'(M)}$$

¹Binary Symmetric Channel (BSC)

راهنمایی: از تکنیک تغییر اندازه و همچنین تغییر متغیر در انتگرال استفاده کنید.

۲. (*) فرض کنید P_{XY} توزیع مشترک (X, Y) باشد و \mathcal{E} واقعه‌ای مستقل از X باشد به طوری که $P[\mathcal{E}] = 1 - \delta$. برای توزیع دلخواه Q_Y که برای آن به صورت $P_X - a.s.$ داریم: $P_{Y|X} \ll Q_Y$ ، ثابت کنید:

$$D_{KL}(P_Y \| Q_Y) \leq \log(1 + D_{X^*}(P_{Y|\mathcal{E}} \| Q_Y)) + \delta \left(\log\left(\frac{1}{\delta}\right) + \mathbb{E}_X[D_{KL}(P_{Y|X} \| Q_Y)] \right) + \sqrt{\delta \text{Var} \left[\log \frac{dP_{Y|X}}{dQ_Y} \right]}.$$

راهنمایی: می‌توانید نامساوی $D_{KL}(P \| Q) \leq \log(1 + D_{X^*}(P \| Q))$ را دانسته فرض کنید. از تحدب D_{KL} و نامساوی کوشی-شوارتز استفاده کنید.

۴ اطلاعات متقابل و خطای تخمین

فرض کنید رابطه‌ی یک کانال با نویز گوسی به صورت $Y = \sqrt{A}X + Z$ باشد که در آن X ورودی کانال، Y خروجی کانال و $Z \sim \mathcal{N}(0, 1)$ است. فرض کنید می‌خواهیم با توجه به خروجی این کانال ورودی آن را با تابعی مانند $f(Y)$ تخمین بزنیم. خطای این تخمین را به صورت $\mathbb{E}[(X - f(Y))^2]$ در نظر می‌گیریم. همینطور خطای بهینه را برای این کانال به صورت $\mathcal{M}_E(A) = \min_f \{ \mathbb{E}[(X - f(Y))^2] \}$ تعریف می‌کنیم. در این سوال قصد داریم رابطه‌ی زیر را بین $\mathcal{M}_E(A)$ و $I(A)$ اثبات کنیم:

$$\frac{d}{dA} I(A) = \mathcal{M}_E(A).$$

۱. ابتدا ثابت کنید تابع تخمین بهینه همان $\mathbb{E}[X|Y]$ است.

۲. ثابت کنید اگر ورودی کانال $Y = \sqrt{\delta}X + Z$ توزیع گوسی داشته باشد، برای $\delta \rightarrow 0$ داریم:

$$I(X; Y) = \frac{\delta}{4} \mathbb{E}[(X - \mathbb{E}[X])^2] + o(\delta)$$

راهنمایی: از رابطه‌ی $I(X; Y) = \mathbb{E}_X [D_{KL}(P_{Y|X} \| P_W)] - D_{KL}(P_Y \| P_W)$ برای توزیع Z مناسب استفاده کنید.

۳. برای اثبات قضیه از ایده‌ی کانال با نویز افزایشی استفاده می‌کنیم. برای این کار از ترکیب دو کانال گوسی استفاده می‌کنیم. به این صورت که ابتدا مقداری نویز به ورودی اضافه می‌کنیم تا نسبت سیگنال به نویز برابر $A + \delta$ شود، و سپس نویز بیشتری اضافه می‌کنیم تا نسبت سیگنال به نویز به A کاهش یابد (شکل ۱ را ببینید). ثابت کنید که برای اثبات قضیه کافیت ثابت کنیم:

$$I(X; Y_1) - I(X; Y_2) = \frac{\delta}{4} \mathcal{M}_E(A) + o(\delta).$$

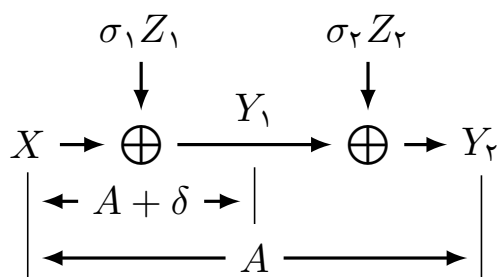
همین‌طور ثابت کنید: $I(X; Y_1) - I(X; Y_2) = I(X; Y_1 | Y_2)$.

۴. رابطه‌ی زیر را اثبات کنید:

$$(A + \delta)Y_1 = AY_2 + \delta X + \sqrt{\delta}Z$$

که در آن Z یک نرمال استاندارد و مستقل از X است.

۵. با توجه به قسمت‌های قبل حکم را نتیجه بگیرید.



شکل ۱: کانال با نویز افزایشی

۵ برخی خواص انحراف χ^2

فرض کنید یک خانواده‌ی پارامتری از توزیع‌ها به صورت $\{P_\theta : \theta \in \Theta\}$ داریم و π یک توزیع روی فضای Θ است. توزیع مخلوط P_π را به صورت زیر تعریف می‌کنیم:

$$P_\pi = \int P_\theta \pi(d\theta)$$

برای توزیع دلخواه Q تعریف می‌کنیم:

$$\mathcal{W}(\theta, \tilde{\theta}) = \mathbb{E}_Q \left[\frac{p_\theta p_{\tilde{\theta}}}{q^2} \right]$$

۱. ثابت کنید:

$$D_{\chi^2}(P_\pi \| Q) = \mathbb{E}_{\theta, \tilde{\theta} \sim \pi} [\mathcal{W}(\theta, \tilde{\theta})] - 1$$

۲. ثابت کنید اگر $D_{\chi^2}(P_n \| Q_n) = O(1)$ داریم:

$$\mathbb{P} \triangleleft \mathbb{Q},$$

که در آن $\mathbb{Q} = \{Q_n\}_{n=1}^\infty$ و $\mathbb{P} = \{P_n\}_{n=1}^\infty$ راهنمایی: به سوال ۶ از تمرین اول مراجعه کنید.

۶ مسئله‌ی تشخیص در SBM

Planted Partition Model یک مدل برای تولید گراف تصادفی است. فرض کنید $\sigma \in \{-1, 1\}^n$ باشد. در این صورت گراف تصادفی به صورت زیر تولید می‌شود:

$$A_{ij} \sim \begin{cases} \mathcal{P} & \sigma_i = \sigma_j \\ \mathcal{Q} & \sigma_i \neq \sigma_j \end{cases}$$

که در آن $\mathbf{A} = [A_{ij}]_{n \times n}$ ماتریس مجاورت وزن‌دار گراف است. این توزیع را با $G(\sigma, \mathcal{P}, \mathcal{Q})$ نمایش می‌دهیم. در حالتی که $\mathcal{P} \sim \text{Bernoulli}(p)$ و $\mathcal{Q} \sim \text{Bernoulli}(q)$ باشد، به این مدل Stochastic Block Model گفته می‌شود و آنرا با $\text{SBM}(\sigma, p, q)$ نمایش می‌دهیم. حال مسئله‌ی آزمون فرض دوتایی زیر را در نظر بگیرید:

$$H_0 : \mathcal{G} \stackrel{\text{i.i.d.}}{\sim} R_0 = G(n, \frac{\mathcal{P} + \mathcal{Q}}{2})$$

$$H_1 : \mathcal{G} \stackrel{\text{i.i.d.}}{\sim} R_1 = G(\mathcal{P}, \mathcal{Q}),$$

که در آن منظور از توزیع $G(\mathcal{P}, \mathcal{Q})$ این است که ابتدا بردار σ با توزیع Rademacher $\sigma_i \stackrel{\text{i.i.d.}}{\sim} (\pm 1)$ (با احتمال برابر) تولید می‌شود و سپس \mathcal{G} از توزیع $G(\sigma, \mathcal{P}, \mathcal{Q})$ نمونه‌برداری می‌شود. منظور از $G(n, \frac{\mathcal{P} + \mathcal{Q}}{2})$ نیز این است که وزن همه‌ی یال‌ها از توزیع $\frac{\mathcal{P} + \mathcal{Q}}{2}$ می‌آید. حال می‌خواهیم در چند گام قضیه‌ی زیر را ثابت کنیم:

قضیه ۶-۱. در حالت $\text{SBM}(\sigma, p, q)$ اگر $p = \frac{a}{n}, q = \frac{b}{n}$ باشد و داشته باشیم: $\frac{(a-b)^2}{2(a+b)} < 1$ ، در این صورت تشخیص بین دو فرض بالا غیرممکن می‌شود. یعنی خطای تشخیص نمی‌تواند به صفر همگرا شود، وقتی $n \rightarrow \infty$.

۱. اگر $P_{\sigma} = G(\sigma, P, Q)$ باشد، ثابت کنید:

$$\mathcal{W}(\sigma, \hat{\sigma}) = \mathbb{E}_{R_{\circ}} \left[\frac{p_{\sigma} p_{\hat{\sigma}}}{r_{\circ}^2} \right] \leq \exp\left(\frac{\rho}{\Upsilon} \langle \sigma, \hat{\sigma} \rangle^2\right),$$

که در آن:

$$\rho = \int_x \frac{(p(x) - q(x))^2}{\Upsilon(p(x) + q(x))} dx.$$

۲. ثابت کنید در حالت $\text{SBM}(\sigma, p, q)$ که $p = \frac{a}{n}, q = \frac{b}{n}$ ، با تعریف $\tau = \frac{(a-b)^2}{\Upsilon(a+b)}$ داریم:

$$\rho = \frac{\tau + o(1)}{n}.$$

۳. با استفاده از قضیه‌ی حد مرکزی حکم را ثابت کنید (فرض کنید در این جا همگرایی در توزیع همگرایی تابع مولد گشتاور را نتیجه می‌دهد، نیازی به اثبات این مورد نیست).

۷ تعیین ناحیه‌ی مشترک به وسیله‌ی توزیع‌های دوتایی

برای دو f -انحراف به شکل $D_f(P\|Q)$ و $D_g(P\|Q)$ ، ناحیه‌ی مشترک^۲ به صورت زیر تعریف می‌شود:

$$\mathcal{R} \triangleq \left\{ \left(D_f(P\|Q), D_g(P\|Q) \right) : P, Q \text{ are probability measures on some measurable space} \right\}.$$

همچنین ناحیه‌ی مشترک بر روی توزیع‌های k -تایی به صورت زیر تعریف می‌شود:

$$\mathcal{R}_k \triangleq \left\{ \left(D_f(P\|Q), D_g(P\|Q) \right) : P, Q \text{ are probability measures on } [k] \right\}.$$

در این پرسش قصد داریم نشان دهیم

$$\mathcal{R} = \text{co}(\mathcal{R}_{\Upsilon}),$$

که در آن منظور از $\text{co}(\mathcal{R}_{\Upsilon})$ پوش محدب^۳ مجموعه‌ی \mathcal{R}_{Υ} است. به بیان دیگر برای به دست آوردن ناحیه‌ی مشترک دو f -انحراف به جای بررسی تمامی توزیع‌ها، کافی است توزیع‌های دوتایی را در نظر گرفت.

۱. نشان دهید مجموعه‌ی \mathcal{R} محدب است.

۲. نشان دهید مجموعه‌های \mathcal{R}_k و \mathcal{R} را می‌توان به صورت زیر نمایش داد

$$\mathcal{R} = \left\{ \left(\begin{array}{c} \mathbb{E}[f(X)] + \tilde{f}(\circ)(1 - \mathbb{E}[X]) \\ \mathbb{E}[g(X)] + \tilde{g}(\circ)(1 - \mathbb{E}[X]) \end{array} \right)^{\top} : X \geq \circ, \mathbb{E}[X] \leq 1 \right\}$$

$$\mathcal{R}_k = \left\{ \left(\begin{array}{c} \mathbb{E}[f(X)] + \tilde{f}(\circ)(1 - \mathbb{E}[X]) \\ \mathbb{E}[g(X)] + \tilde{g}(\circ)(1 - \mathbb{E}[X]) \end{array} \right)^{\top} : \begin{array}{l} X \geq \circ, \mathbb{E}[X] \leq 1, X \text{ takes at most } k-1 \text{ values,} \\ \text{or,} \\ X \geq \circ, \mathbb{E}[X] = 1, X \text{ takes at most } k \text{ values} \end{array} \right\},$$

که در آن‌ها:

$$\tilde{f}(\circ) \triangleq \lim_{x \rightarrow \circ} x f\left(\frac{1}{x}\right),$$

$$\tilde{g}(\circ) \triangleq \lim_{x \rightarrow \circ} x g\left(\frac{1}{x}\right).$$

۳. (*) نشان دهید:

$$\mathcal{R} = \mathcal{R}_{\Upsilon}.$$

راهنمایی: از قضیه‌ی Fenchel-Eggleston-Caratheodory که در ادامه بیان می‌شود استفاده کنید.

²Joint range

³Convex hull

قضیه ۱-۷. فرض کنید $\mathcal{S} \subseteq \mathbb{R}^d$ و $\mathbf{x} \in \text{co}(\mathcal{S})$. آن گاه می توان مجموعه ای مانند \mathcal{S}' شامل $d + 1$ نقطه به صورت $\mathcal{S}' = \{\mathbf{x}_1, \dots, \mathbf{x}_{d+1}\} \subseteq \mathcal{S}$ یافت، به طوری که $\mathbf{x} \in \text{co}(\mathcal{S}')$. همچنین اگر \mathcal{S} همبند باشد آن گاه d نقطه کافی است.

۴. (*) برای $k \geq 2$ نشان دهید:

$$\mathcal{R}_{k+1} = \text{co}(\mathcal{R}_k).$$

۵. حکم سوال را اثبات کنید.