

# نظریه‌ی اطلاعات، آمار و یادگیری (۱-۲۵۱۱۰)



تمرین سری سوم

ترم بهار ۱۴۰۳-۰۴

دانشکده‌ی مهندسی برق

دانشگاه صنعتی شریف

استاد: دکتر محمدحسین یاسائی میبدی

مهلت تحویل: جمعه ۲ خرداد ۱۴۰۴ 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:

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

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  $X \sim \text{Ber}(1/2)$  باشد. قرار دهید  $P_0 = P_{Y|X=0}, P_1 = P_{Y|X=1}$ . ثابت کنید

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

راهنمایی: برای سمت چپ از نامساوی Pinsker استفاده کنید و برای سمت راست از نامساوی بین اطلاعات متقابل و  $\chi^2$  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 رؤس 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 رؤس 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 رؤس 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<sup>۱</sup> با پارامتر  $\delta$  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  $\rho$  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 رؤس در عمق  $k$  از این درخت باشد، ثابت کنید اگر  $d(1 - 2\delta)^2 < 1$  داریم:

$$d_{\text{TV}} \left( \mathbb{P}_{X_{\mathcal{S}_k} | X_{\rho}=+1}, \mathbb{P}_{X_{\mathcal{S}_k} | X_{\rho}=0} \right) \xrightarrow[k \rightarrow \infty]{} 0.$$

## ۳ تبخّر در اثبات نامساوی‌ها!

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

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

<sup>1</sup>Binary Symmetric Channel (BSC)

راهنمایی: از تکنیک تغییر اندازه و همچنین خواص تابع محدب 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:

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

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

## ۴ نامساوی قوی پردازش داده‌ها

هدف ما در این سوال، آشنایی بیشتر با نامساوی قوی پردازش داده‌ها<sup>۲</sup> است. فرض کنید تابع  $f: (\circ, \infty) \rightarrow \mathbb{R}$ ، تابعی محدب باشد که در  $x = 1$  محدب اکید است و همچنین داریم:  $f(1) = \circ$ . برای متغیرهای تصادفی  $X \in \mathcal{X}$  و  $Y \in \mathcal{Y}$  که در آن  $|\mathcal{X}|, |\mathcal{Y}| < \infty$  و کانال  $P_{Y|X}$  تعریف می‌کنیم:

$$\eta_f(P_{Y|X}, Q) \triangleq \sup_{P: \circ < D_f(P \| Q) < \infty} \frac{D_f(P_{Y|X} \circ P \| P_{Y|X} \circ Q)}{D_f(P \| Q)},$$

$$\eta_f(P_{Y|X}) \triangleq \sup_Q \eta_f(P_{Y|X}, Q).$$

که در آن منظور از  $P_{Y|X} \circ P$  توزیع القا شده بر روی متغیر تصادفی  $Y$  از روی  $X$  با توزیع  $P$  و تحت کانال  $P_{Y|X}$  است. ۱. با توجه به تعریف بالا، برای انحراف TV ثابت کنید:

$$\eta_{\text{TV}}(P_{Y|X}) = \sup_{x, x'} \{d_{\text{TV}}(P_{Y|X=x}, P_{Y|X=x'})\}.$$

۲. (\*) برای انحراف KL نشان دهید:

$$\eta_{\text{KL}}(P_{Y|X}, P_X) = \sup_{U: U \rightarrow X \rightarrow Y} \frac{I(U; Y)}{I(U; X)},$$

که در آن، سوپریموم بر روی تمام زنجیره‌مارکف‌های به شکل  $U \rightarrow X \rightarrow Y$  گرفته می‌شود که در آن توزیع توأم  $X, Y$  ثابت است.

۳. (\*) نشان دهید برای هر تابع  $f$  با ویژگی‌هایی که بیان شد، همواره داریم:

$$\eta_f(P_{Y|X}) \leq \eta_{\text{TV}}(P_{Y|X}).$$

## ۵ 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  $Z \sim \mathcal{N}(\circ, 1)$  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  $I(A) \triangleq I(X; \sqrt{A}X + Z)$  اثبات کنیم:

$$\frac{d}{dA} I(A) = \frac{1}{4} \mathcal{M}_E(A).$$

<sup>2</sup>Strong Data Processing Inequality

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

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

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

راهنمایی: از رابطه‌ی  $I(X; Y) = \mathbb{E}_X [D_{\text{KL}}(P_{Y|X} \| P_W)] - D_{\text{KL}}(P_Y \| P_W)$  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:

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

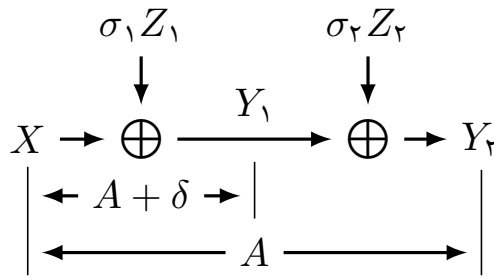
همین‌طور ثابت کنید:  $I(X; Y_1) - I(X; Y_2) = I(X; Y_1 | Y_2)$ .

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

## 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:

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

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

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

قضیه ۶-۱. در حالت  $SBM(\sigma, p, q)$  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:

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

که در آن:

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

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

همانطور که در درس دیدیم، انحراف TV دارای تعریف معادل زیر است:

$$d_{TV}(P||Q) = \sup_E \{P(E) - Q(E)\},$$

که در آن  $E$  عضو مجموعه‌ی پیشامدهای فضای احتمال است.

۱. با استفاده از این تعریف، قضیه‌ی Strassen را ثابت کنید:

$$d_{TV}(P||Q) = \inf_{P_{X\hat{X}} \in \Pi(P, Q)} P_{X\hat{X}}(X \neq \hat{X}),$$

که در آن

$$\Pi(P, Q) = \{P_{X\hat{X}} : \text{probability measure on } \mathcal{X}^2, P_X = P, P_{\hat{X}} = Q\}.$$

۲. فرض کنید  $X^n$  یک منبع i.i.d. باشد و سناریوی کدینگ منبع با اعوجاج زیر را در نظر بگیرید:



هدف ما در این سوال، پیدا کردن کران بالایی برای متوسط فاصله‌ی همینگ ورودی کدگذار و خروجی کدگشا برحسب انحراف KL است. ثابت کنید توزیع  $P_{X^n \hat{X}^n}$  مناسبی وجود دارد به گونه‌ای که

$$\frac{1}{n} \mathbb{E}[d_n(X^n, \hat{X}^n)] = \frac{1}{n} \sum_{i=1}^n P_{X^n \hat{X}^n}(X_i \neq \hat{X}_i) \leq \sqrt{\frac{1}{n} D_{KL}(P_{\hat{X}^n} || P_{X^n})}.$$

راهنمایی: از نامساوی Pinsker و قضیه‌ی Strassen استفاده کنید.

## ۸ انحراف Marton

۱. انحراف Marton به صورت زیر تعریف می‌شود:

$$D_M(P\|Q) = \int \left(1 - \frac{dP}{dQ}\right)_+ dQ.$$

ثابت کنید:

$$D_M(P\|Q) = \inf_{P_{XY} \in \Pi(P, Q)} \mathbb{E} [P_{XY}^*(X \neq Y|Y)].$$

۲. انحراف Marton متقارن به صورت زیر تعریف می‌شود:

$$D_{SM}(P\|Q) = D_M(P\|Q) + D_M(Q\|P).$$

ثابت کنید:

$$D_{SM}(P\|Q) = \inf_{P_{XY} \in \Pi(P, Q)} \{ \mathbb{E} [P_{XY}^*(X \neq Y|Y)] + \mathbb{E} [P_{XY}^*(X \neq Y|X)] \}.$$