



## یادگیری عمیق

پاییز ۱۴۰۱  
استاد: دکتر فاطمی زاده

گردآورندگان: -

تمرین اول      مفاهیم پایه      مهلت ارسال: سه شنبه ۱۸ آبان (با احتساب تاخیر)

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر پاسخ همه‌ی تمرین تا سقف ۶ روز و در مجموع ۲۰ روز، وجود دارد. پس از گذشت این مدت، پاسخ‌های ارسال شده پذیرفته نخواهند بود. همچنین، به ازای هر روز تأخیر غیر مجاز ۱۰ درصد از نمره تمرین به صورت ساعتی کسر خواهد شد.
- همکاری و هم‌فکری شما در انجام تمرین مانعی ندارد اما پاسخ ارسالی هر کس حتما باید توسط خود او نوشته شده باشد. (دقت کنید در صورت تشخیص مشابهت غیرعادی برخورد جدی صورت خواهد گرفت.)
- در صورت هم‌فکری و یا استفاده از هر منابع خارج درسی، نام هم‌فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفا تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.
- نتایج و پاسخ‌های خود را در یک فایل با فرمت zip به نام HW۱-Name-StudentNumber در سایت **Quera** قرار دهید. برای بخش عملی تمرین نیز لینک گیت‌هاب که تمرین و نتایج را در آن آپلود کرده‌اید قرار دهید. دقت کنید هر سه فایل نوتبوک تکمیل شده بخش عملی را در گیت‌هاب قرار دهید.
- لطفا تمامی سوالات خود را از طریق کوثرای درس مطرح بکنید (برای اینکه تمامی دانشجویان به پاسخ‌های مطرح شده به سوالات دسترسی داشته باشند و جلوی سوالات تکراری گرفته شود، به سوالات در بسترهای دیگر پاسخ داده نخواهد شد).
- دقت کنید کدهای شما باید قابلیت اجرای دوباره داشته باشند، در صورت دادن خطا هنگام اجرای کدتان، حتی اگر خطا بدلیل اشتباه تایپی باشد، نمره صفر به آن بخش تعلق خواهد گرفت.

### سوالات نظری (۳۰۰ نمره)

۱. (۲۵ نمره)

(آ) با توجه به تعریف نرم ماتریسی رابطه زیر را برای  $A \in R^{m \times n}$  ثابت کنید:

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{\text{rank}(A)} \|A\|_2$$

راهنمایی: از تجزیه SVD ماتریس  $A$  استفاده کنید.

(ب) i. با در نظر گرفتن متغیر تصادفی نامنفی  $X$ ، نامساوی زیر را اثبات کنید (نامساوی مارکوف)

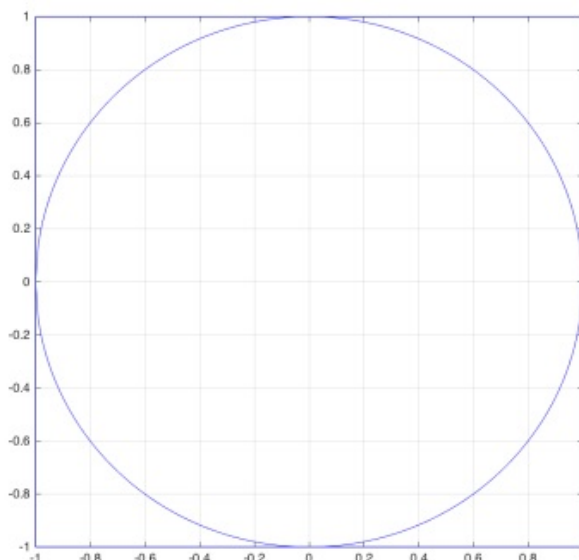
$$P(X \geq \alpha) \leq \frac{E(X)}{\alpha}$$

ii. با در نظر گرفتن نتیجه بخش الف، نشان دهید برای متغیر تصادفی دلخواه  $Z$  با امید ریاضی  $\mu$  و واریانس  $\sigma^2$  نامساوی زیر برقرار است (نامساوی چیشف):

$$P(|Z - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2}$$

iii. می خواهیم مقدار عدد  $\pi$  را تخمین بزنیم. برای این کار روی صفحه مختصات دوعدی، دایره ای به شعاع واحد و مربع محیطی آن را شبیه شکل زیر رسم می کنیم. مساحت این دایره  $\pi$  و مساحت مربع محیطی آن ۴ است. برای تخمین مقدار  $\pi$  تعدادی نقاط تصادفی داخل این مربع تولید کرده

و نسبت تعداد نقاطی که داخل دایره قرار می گیرند را به تعداد کل به عنوان مقدار عدد  $\pi$  در نظر می گیریم.



iv. با استفاده از نامساوی چبیشف تعداد عددهای تصادفی ای که باید تولید کنیم تا با قطعیت ۹۵ درصد بدانیم که خطای تخمین از ۱ درصد کمتر است را مشخص کنید.

۲. (۲۵ نمره)

• با فرض اینکه  $x \in \mathbb{R}$  و  $A$  ماتریسی مربعی به ابعاد  $n \times n$  باشد، دو مورد زیر را اثبات کنید.

$$\frac{\partial a^T x}{\partial x} = a^T$$

$$\frac{\partial x^T A x}{\partial x} = x^T (A + A^T)$$

• اگر درایه های ماتریس  $A$  تابع پارامتر اسکالر  $\beta$  باشند، رابطه ی زیر را اثبات کنید.

$$\frac{\partial A^{-1}}{\partial \beta} = -A^{-1} \frac{\partial A}{\partial \beta} A^{-1}$$

• نشان دهید

$$\nabla_A |A| = |A| A^{-T}$$

و به دنبال آن

$$\nabla_A \log |A| = A^{-1}$$

۳. (۱۰ نمره) با فرض اینکه  $\lambda_1, \lambda_2, \dots, \lambda_n$  مقادیر ویژه ماتریس  $A$  باشند، موارد خواسته شده را اثبات کنید.

$$\lambda_1 + \lambda_2 + \dots + \lambda_n = \text{trace}(A)$$

$$\lambda_1 \lambda_2 \dots \lambda_n = \det(A)$$

۴. (۲۰ نمره) فرض کنید ماتریس  $A$  تجزیه SVD به شکل  $A = U \Sigma V^T$  داشته باشد. برای این ماتریس عملگر Pseudo-Inverse به شکل  $A^\dagger = V \Sigma^{-1} U^T$  تعریف می شود. حال اثبات کنید:

$$A^\dagger = (A^T A)^{-1} A^T, \text{ اگر ماتریس } A \text{ رنک کامل ستونی داشته باشد,}$$

$$A^\dagger = A^T (A A^T)^{-1}, \text{ اگر ماتریس } A \text{ رنک کامل سطری داشته باشد,}$$

۵. (۴۰ نمره) فرض کنید ماتریس  $M$  را بصورت  $M = \begin{bmatrix} A_{n \times n} & B_{n \times k} \\ C_{k \times n} & D_{k \times k} \end{bmatrix}$  تعریف کنیم، برای این ماتریس اثبات کنید:

• اگر  $A$  وارون پذیر باشد آنگاه:

$$\det \begin{bmatrix} A_{n \times n} & B_{n \times k} \\ C_{k \times n} & D_{k \times k} \end{bmatrix} = \det(A) \det(D - CA^{-1}B)$$

• اگر  $D$  وارون پذیر باشد آنگاه:

$$\det \begin{bmatrix} A_{n \times n} & B_{n \times k} \\ C_{k \times n} & D_{k \times k} \end{bmatrix} = \det(D) \det(A - BD^{-1}C)$$

• با کمک دو بخش قبل اثبات کنید:

$$(A + BDC)^{-1} = A^{-1} - A^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}$$

• اگر  $A$  یک ماتریس وارون پذیر و  $u, v$  بردارهای ستونی باشند، با استفاده از تکنیک‌های استفاده شده در بخش‌های قبل ثابت کنید:

$$\det(A - uv^T) = \det(A)(1 - v^T A^{-1}u)$$

(نکته: مشخصاً هر کجا از ماتریسی وارون گرفته شده است آن ماتریس را وارون‌پذیر در نظر بگیرید.)

۶. (۷۰ نمره) در این سوال میخواهیم یک ویژگی جالب درباره ارتباط مقادیر ویژه دو ماتریس مرتبط با هم را کشف کنیم که شاید بعدها در رابطه با کاربرد آن بیشتر با هم صحبت کنیم، اما قبل اثبات حکم کلی چند لم و ابزار جری را با هم اثبات می‌کنیم.

• فرض کنید  $x$  و  $y$  دو بردار ستونی با درایه‌های در حالت کلی مختلط، با بعد  $n$  باشند، یعنی  $x, y \in \mathbb{C}^n$ ،  $a \in \mathbb{C}$  یک اسکالر مختلط باشد،  $B \in M_n$  ماتریس اولیه ما باشد و به کمک بردارها و اسکالر تعریف شده ماتریس جدیدی بصورت  $A = \begin{bmatrix} B & x \\ y^* & a \end{bmatrix} \in M_{n+1}$  تولید بکنیم، حال اثبات کنید که چندجمله‌ای مشخصه ماتریس جدید تولید شده بصورت زیر است:

$$p_A(t) = \det(tI - A) = (t - a)p_B(t) - y^*(adj(tI - B))x$$

• فرض کنید ماتریس‌های  $A, B \in M_n$  ماتریس‌های هرمیتی باشند و مقادیر ویژه این ماتریس‌ها به ترتیب از کوچک به بزرگ مرتب شده باشند یعنی رابطه زیروند مقادیر ویژه و بزرگی آنها بصورت:  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  باشد، حال برای این دو ماتریس و ماتریس حاصل جمع آنها رابطه زیر را اثبات بکنید:

$$\lambda_i(A + B) \leq \lambda_{i+j}(A) + \lambda_{n-j}(B) \quad j = 0, 1, \dots, n - i$$

برای هر  $i = 1, \dots, n$

• حال میخواهیم حکم کلی را اثبات کنیم، فرض کنید متغیرها و ماتریس‌هایی به شرح روبرو داشته باشیم:

$$B \in M_n, y \in \mathbb{C}^n, a \in \mathbb{R}, A = \begin{bmatrix} B & y \\ y^* & a \end{bmatrix}$$

یک ماتریس هرمیتی است، حال به کمک بخش‌های قبل اثبات بکنید:

$$\lambda_1(A) \leq \lambda_1(B) \leq \lambda_2(A) \leq \dots \leq \lambda_n(A) \leq \lambda_n(B) \leq \lambda_{n+1}(A)$$

😊 اگر خیلی در رابطه با کاربرد این نکته کنجکاو هستید، از این قضیه استفاده میشود تا اثبات کنیم افزایش پیچیدگی یک شبکه عصبی یا مدل یادگیری عمیق در صورت وجود رگولایزر مناسب در تابع هزینه شبکه، میتواند باعث از بین رفتن پدیده Double-Descent بشود. (با مفاهیم گفته شده در طول درس آشنا خواهید شد!)

۷. (۱۰ نمره) فرض کنید از توزیع با چگالی زیر، نمونه‌های  $X_1, X_2, \dots, X_n$  را دریافت کرده‌ایم، تخمینگر بیشینه likelihood را برای  $\theta$  بدست بیاورید.

$$f(x) = \frac{1}{\theta^2} x e^{-\frac{x}{\theta}}, \quad 0 < \theta < \infty$$

۸. (۲۰ نمره) توزیع نرمالی با واریانس معلوم  $\sigma^2$  و میانگین مجهول  $\mu$  داریم. با فرض اینکه  $n$  داده‌ی  $X_1, X_2, \dots, X_n$  از این توزیع بدست آورده‌ایم، ابتدا تخمینگر MLE و MAP را برای  $\mu$  بدست بیاورید، در نهایت بررسی کنید در صورت میل دادن تعداد سмпل‌ها به بی نهایت چه اتفاقی برای این تخمینگرها می‌افتد. (برای تخمینگر MAP تصور کنید که توزیع پیشینه خود یک توزیع نرمال با میانگین  $\gamma$  و واریانس  $\beta^2$  است.)

۹. (۲۰ نمره) فرض کنید یک توزیع گوسی چند متغیره بصورت  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  داریم که در آن متغیر تصادفی را میتوان بصورت زیر نوشت :

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$$

و به طبع آن خواهیم داشت :

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}$$

و برای ماتریس کوواریانس :

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix}$$

• حال اثبات کنید توزیع شرطی  $p(\mathbf{x}_a|\mathbf{x}_b)$  یک توزیع گوسی است که مشخصه‌هایی به شرح زیر دارد :

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}$$

• حال اثبات کنید توزیع حاشیه‌ای  $p(\mathbf{x}_a)$  توزیع گوسی‌ای است که مشخصه‌هایی به شرح زیر دارد:

$$E[\mathbf{x}_a] = \boldsymbol{\mu}_a$$

$$Cov[\mathbf{x}_a] = \boldsymbol{\Sigma}_{aa}$$

۱۰. (۶۰ نمره) در این سوال میخواهیم با الگوریتم بهینه‌سازی گرادیان کاهشی و ارتباط آن با عملگر  $Pseudo-inverse$  آشنا شویم. فرض کنید مسئله بهینه‌سازی بصورت  $\min_x \|Ax - b\|^2$  داریم،

- اثبات کنید اگر این مسئله چند جواب داشته باشد، جواب به فرم  $x^* = A^\dagger b$  کمترین نرم ۲ را در میان آنها دارد.

- فرض کنید میخواهیم این مسئله را با الگوریتم گرادیان کاهشی حل کنیم. بدین صورت که در مرحله اول مقدار متغیر بهینه‌سازی را با صفر مقداردهی میکنیم ( $x^{(0)} = 0$ ) و در هر مرحله متغیر را بصورت  $x^{(t+1)} = x^{(t)} - \nu A^T(Ax^{(t)} - b)$  آپدیت میکنیم که در آن  $\nu$  نرخ یادگیری است. اثبات کنید در این صورت جواب مسئله در صورت انتخاب پارامتر یادگیری بصورت مناسب پس از تعداد ایتريشن کافی به جواب  $x^* = A^\dagger b$  میل می‌کند.

- اثبات کنید یک کران مناسب برای نرخ یادگیری میتواند  $\frac{2}{\sigma_{max}^2(A)}$  باشد، بدین معنا که اگر نرخ یادگیری بصورت  $\nu \leq \frac{2}{\sigma_{max}^2(A)}$  انتخاب شود، در هر مرحله از پیاده‌سازی الگوریتم بهینه‌سازی، تابع هزینه نسبت به مرحله قبل کاهش خواهد داشت:  $\|Ax^{(t+1)} - b\| \leq \|Ax^{(t)} - b\|$  (تمامی نرم‌ها، نرم ۲ هستند).

---

### سوالات عملی (۳۰۰ نمره)

---

۱. (۱۰۰ نمره) فایل نوتبوکی در اختیار شما قرار داده شده است. راهنمایی‌های لازم برای نحوه انجام تمرین در فایل نوتبوک انجام شده است. در این تمرین با کار با کتابخانه pandas، نحوه تبدیل داده خام به داده مناسب برای ورودی مدل یادگیری و در نهایت پیاده‌سازی یک مدل یادگیری ماشین مبتنی بر یک روش بهینه‌سازی متناوب آشنا میشوید. برای بخش‌ترین کردن مدل به مقاله **Paper** مراجعه شده و تا بخش stochastic gradient descent خوانده شود.

۲. (۱۰۰ نمره) فایل نوتبوکی که در اختیارتان قرار داده شده است را کامل کنید. در این تمرین به کمک SVM به دسته‌بندی داده‌هایی که در فایل Heart Disease Dataset.csv قرار دارند میپردازیم. توصیفی از این داده‌ها در فایل Dataset Description.pdf آمده است. برای دریافت نمره کامل بخش MySVM باید دقت بهترین هسته پیاده‌سازی شده روی داده‌های آزمون حداقل ۹۰ درصد باشد. برای دقت‌های بالاتر محسوس، در صورت توضیح مناسب دلیل، نمره امتیازی در نظر گرفته خواهد شد.

۳. (۱۰۰ نمره) فایل نوتبوکی که در اختیارتان قرار داده شده است را کامل کنید. در این تمرین با رگرسیون خطی، رگرسیون ناپارامتری و نزدیک‌ترین همسایه و پیاده‌سازی آنها آشنا میشوید.