



یادگیری عمیق

بهار ۱۴۰۰

مدرس: حمید بیگی

دانشکده مهندسی کامپیوتر

زمان تحویل: ۱۸ خرداد

یادگیری عمیق

تمرین سری چهارم

به نکات زیر توجه کنید:

فایل گزارش را به همراه تمامی کدها در یک فایل فشرده و با عنوان HW4_STD# بارگذاری نمایید. برای هر یک از سوالات عملی، پوشه ای مجزا در نظر بگیرید و کدها را درون آن قرار دهید و از شماره سوال برای نام پوشه استفاده کنید. در صورتی که در جواب از مقاله خاصی استفاده شده است لازم است که به آن اشاره شود. نمره کل تمرین ۷۴ است. ۲۲ نمره از کل تمرین امتیازی می باشد. برای بخش عملی دو نوتبوک در اختیار شما قرار خواهد گرفت و نمره عملی از آن محاسبه می شود که جمعا ۶۰ نمره است.

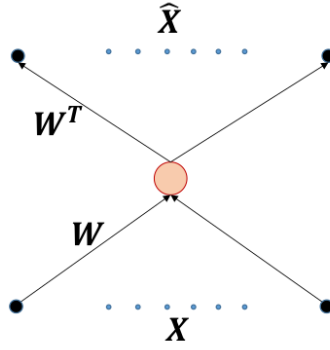
مسئله اول - KL-Divergence (۱۰ نمره)

(۱) kl-divergence بین دو توزیع گاوسی $N(\mu_1, \Sigma_1)$ و $N(\mu_2, \Sigma_2)$ را بدست آورید (ماتریس کواریانس قطری است).
(۲) توزیع $p(x)$ را یک توزیع ثابت دلخواه در نظر بگیرید که می خواهیم آن را با استفاده از توزیع $q(x) = N(x|\mu, I)$ تقریب بزنیم (ماتریس کواریانس همانی است). با نوشتن kl-divergence بین دو توزیع و مشتق گرفتن نسبت به μ نشان دهید که میانگین بهینه برابر است با:

$$\mu^* = \underset{\mu}{\operatorname{argmin}} \operatorname{KL}(p||q) = \mathbb{E}_p[x] \quad (۱)$$

مسئله دوم - autoencoders (۱۰ نمره)

(۱) مزایای استفاده از autoencoder چیست و چرا از آن استفاده می شود؟
(۲) نشان دهید که linear autoencoder شکل ۱ که تلاش به کمینه کردن L_2 divergence می کند همان PCA است.



شکل ۱: linear autoencoder

مسئله سوم- marginal likelihood estimation (۱۰ نمره)

یک latent variable model در نظر بگیرید که در آن x مقدار مشاهده شده و z مقدار نهان^۱ است. با استفاده از تابع بیشینه حاشیه‌ای^۲:

$$p_{\theta}(x) = \int_z p_{\theta}(x|z)p(z)dz \quad (۲)$$

با استفاده از توزیع پیشنهاد شده $q(z|x)$ یک importance sampling estimator را در نظر بگیرید:

$$\hat{L}(x) = \log\left(\frac{1}{M} \sum_{i=1}^M \frac{p_{\theta}(x|z_i)p(z)}{q(z_i|x)}\right) \quad (۳)$$

نشان دهید که \hat{L} یک biased estimator از $\log(p_{\theta}(x))$ است اما در صورتی که $M \rightarrow \infty$ آنگاه به طور تقریبی unbiased است.

$$\mathbb{E}_{z_i \sim q(\cdot|x)}[L(\hat{x})] \leq \log(p_{\theta}(x)) \quad (۴)$$

$$\lim_{M \rightarrow \infty} L(\hat{x}) = \log(p_{\theta}(x)) \quad (۵)$$

این ایده استفاده شده در مقاله importance weighted autoencoder است.

مسئله چهارم (امتیازی ۱۰ نمره)

(۱) در صورتی که latent variable در VAE ها به جای پیوسته گسسته باشد. چه تغییراتی لازم است داده شود.
(۲) در beta-VAE یک ضریب β پشت جمله KL اضافه شده است. با انجام چه محاسباتی این عمل انجام شده است و

^۱ latent
^۲ likelihood marginal

چه تاثیری در عملکرد VAE می گذارد؟

مسئله پنجم شبکه های تخصصی مولد^۳ (۱۲+۲۲)

یک) تابع هدف شبکه های تخصصی مولد در حالت پایه به صورت زیر تعریف می شود:

$$\mathcal{V}(G, D) = \mathbb{E}_{x \sim p_{Data}}[\log D(x)] + \mathbb{E}_{x \sim p_G}[\log(1 - D(x))] \quad (۶)$$

الف) (۲ نمره) اگر ظرفیت تمیزدهنده نامحدود باشد، نقاط بیشینه رابطه ۶ نسبت به تمیز دهنده را بیابید. (بر حسب p_G و p_{Data} بدست آورید)

ب) (۳ نمره) حال اگر در رابطه ۶، گام بهینه سازی نسبت به D به صورت بهینه انجام گیرد، نشان دهید کمینه کردن $\mathcal{V}(G, D^*)$ که D^* نقطه بهینه آن نسبت به D است، معادل کمینه کردن فاصله $JS(p_{Data} || p_G)$ (Shannon Jensen) خواهد بود.

پ) (۲ نمره) دیدگاه فوق به لحاظ نظری مطرح می شود. آنچه در عمل اتفاق می افتد، پارامترهای G و D در هر دسته داده^۴، با توجه به رابطه ۶ یک یا چند بار به روزرسانی می شوند. برای مثال در یک دسته داده می توان سه بار پارامترهای تمیزدهنده و سپس یکبار پارامترهای مولد را به روزرسانی کرد. حال تعیین نسبت به آموزش مولد به تمیز دهنده و سپس یکبار پارامترهای مولد را به روزرسانی کرد. حال تعیین نسبت آموزش مولد به تمیزدهنده یکی از مشکلات اساسی آموزش این خانواده از شبکه ها است. با توجه به دو بخش قبل، برای نزدیک شدن به دیدگاه نظری مطرح شده، چه نسبت آموزشی را بین مولد و تمیزدهنده پیشنهاد می کنید؟

۲) فرض کنید، دامنه توزیع داده اصلی و توزیع شبکه مولد همپوشانی نداشته باشند و همچنین D نزدیک به تمیز دهنده بهینه باشد.

الف) (۲ نمره) گرادیان $\log(1 - D(x))$ را نسبت به logit های شبکه D بدست آورید (اگر تمیز دهنده به صورت $\sigma(a)$ که $a = f(x)$ تعریف شده باشد، a ، logit خواهد بود.)

ب) (۲ نمره) در این حالت چه گرادیانی به شبکه G می رسد؟ چه مشکلی ایجاد می شود؟

پ) (۲ نمره) حال اگر از $-\log D(x)$ به عنوان تابع هزینه مولد شبکه مولد استفاده شود، چطور به حل این مشکل کمک می کند؟

۳) (امتیازی) رویکرد دیگری که در نحوه به روزرسانی پارامترهای دو شبکه G و D می توان متصور بود، آن است که تابع هدف ۶، به صورت همزمان نسبت به پارامترهای هر دو شبکه به روزرسانی شود. فرض کنید تابع هدف شما xy که نسبت به x بیشینه و y کمینه می گردد.

الف) (۱ نمره) نقاط زینی این رویه را بیابید.

ب) (۳ نمره) اگر از نقطه (۱، ۱) شروع کنیم و پارامترهای x و y را به صورت همزمان به روزرسانی کنیم، مسیر حرکت روی

Generative Adversarial Networks:^۳
Batch^۴

فضای xy را بررسی کنید (با استفاده از معادله دیفرانسیل و یا شبیه سازی). آیا این روش همگرا می شود؟
 (۴) (۴ نمره) از آنجا که در هر دو روش بیشینه درست‌نمایی و تخصمی به دنبال آموزش شبکه مولد هستیم، در این بخش می‌خواهیم روش تخصمی دیگری را بیابیم که در عمل با روش بیشینه درست‌نمایی یکسان باشد. می‌دانیم در روش بیشینه درست‌نمایی تابع هزینه زیر مورد استفاده قرار می‌گیرد.

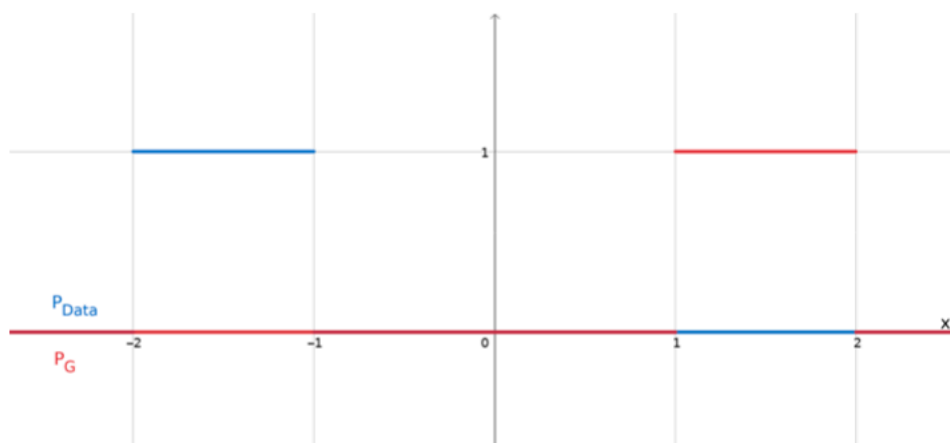
$$\mathcal{L}_{MLE}(\theta) = \mathbb{E}_{x \sim p_{Data}}[-\log p_G(x)] \quad (۷)$$

که θ پارامترهای شبکه مولد است.
 حال روش تخصمی را به این صورت تعریف می‌کنیم که تابع هزینه تمیزدهنده تغییری نکرده (رابطه ۶) اما تابع هزینه شبکه مولد به صورت زیر باشد.

$$\mathcal{L}_{MLE-GAN}(\theta) = \mathbb{E}_{x \sim p_G}[f(x)] \quad (۸)$$

تابع f چه باشد تا گرادیان $\mathcal{L}_{MLE-GAN}$ نسبت به θ ، با گرادیان \mathcal{L}_{MLE} نسبت به θ یکسان شود؟ این تابع را بر حسب D^* و logit آن به دست آورید. چند نکته:

- D^* تمیز دهنده بهینه است.
 - فرض کنید تابع f تابع مستقیمی از پارامترهای θ نباشد. به عبارت دیگر مشتق f نسبت به θ صفر است.
 - logit را با a نشان دهید.
- (۵) اگر ساختار تخصمی استاندارد (رابطه ۶) مورد استفاده باشد و توزیع داده آموزشی و شبکه مولد به صورت شکل زیر باشد، به سوالات زیر پاسخ دهید.



شکل ۲: توزیع داده اصلی (p_{Data}) و توزیع شبکه مولد (p_G)

الف) (۲ نمره) فواصل $JS(p_{Data}||p_G)$ ، $KL(p_{Data}||p_G)$ و $KL(p_G||p_{Data})$ را بیابید.

ب) (۱ نمره) در این حالت چه گرادسانی به شبکه مولد می‌رسد.

پ) (۲ نمره) یکی از راه‌های پایدارتر کردن آموزش شبکه‌های تخصصی اضافه کردن نویز به تصاویر آموزشی ورودی تمییز دهنده است. با توجه به جواب بخش الف (فاصله JS)، اضافه کردن نویز، چطور به بهینه سازی کمک خواهد کرد؟
 ۶) (امتیازی) فاصله دیگری به نام فاصله واسرشتاین^۵ وجود دارد که آشنایی با آن به دلیل محبوبیتش در حوزه شبکه‌های تخصصی مولد، خالی از لطف نیست. به زبان ساده، فاصله واسرشتاین بین دو توزیع کمترین هزینه جابجایی یک توزیع به توزیع دیگر است. فاصله (واسرشتاین ۱-) و به اختصار با W_1 نشان داده می‌شود، بین دو توزیع $p(x)$ و $q(x)$ را به صورت زیر تعریف می‌کنیم:

$$T(\gamma, c) = \int c(x, y) \gamma(x, y) dx dy \quad (۹)$$

$$W_1(p, q) = \inf_{\gamma \in \Gamma(p, q)} T(\gamma, c) \quad (۱۰)$$

که $c(x, y)$ هزینه جابجایی بین دو نقطه x و y (دارای خاصیت متر)، و $\Gamma(p, q)$ خانواده تمام روش‌های مختلف جابجایی^۶ بین این دو توزیع است. مقدار $\gamma(x, y)$ نشان‌دهنده مقداری است که می‌خواهید از نقطه x به y (از توزیع p به توزیع q) منتقل کنید. بنابراین $T(\gamma, c)$ هزینه جابجایی بین دو توزیع p و q تحت روش جابجایی $\gamma(x, y)$ خواهد بود.
 الف) (۳ نمره) ضمن مطالعه این فاصله و یا مقاله‌های مرتبط و بلاگ‌های مرتبط با روش شبکه تخصصی مولد واسرشتاین بیان کنید چه توابعی می‌توانند نشان‌دهنده یک روش جابجایی باشد؟ چه شروطی باید در این دسته توابع صدق کنند؟
 ۲ شرط باید ذکر کنید.

ب) (۳ نمره) اگر $c(x, y) = |x - y|$ باشد، ضمن ذکر یک روش جابجایی بین دو توزیع (لزومی به کمینه بودن هزینه روش جابجایی نیست) p_G و p_{Data} سوال قبل (شش)، هزینه جابجایی بین آن‌ها را محاسبه کنید.
 پ) (۲ نمره) با توجه به جواب قسمت ب، این فاصله چه مزیتی نسبت به فواصل JS و KL دارد؟

موفق باشید.