

به نام خدا

تمرین سری اول
درس یادگیری ماشین
دکتر علی شریفی زارچی

فرزان رحمانی
۴۰۳۲۱۰۷۲۵

سوال اول

الف) تابع softmax معمولاً برای مسائل طبقه‌بندی، به‌ویژه طبقه‌بندی چند کلاسه، به چند دلیل استفاده می‌شود:

۱. تفسیر احتمال

تابع softmax خروجی خام (logits) یک مدل را به احتمالات تبدیل می‌کند. این کار را با تبدیل مقادیر خروجی به محدوده ای بین ۰ و ۱ انجام می‌دهد و اطمینان حاصل می‌کند که مجموع این مقادیر ۱ است و آنها را به عنوان احتمال قابل تفسیر می‌کند. برای یک مدل خروجی n کلاس، تابع softmax احتمال تعلق ورودی به هر کلاس را محاسبه می‌کند:

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

این به مدل اجازه می‌دهد تا عدم قطعیت را با اختصاص احتمالات مختلف به هر کلاس بیان کند.

۲. طبقه بندی چند کلاسه

Softmax به ویژه برای کارهای طبقه بندی چند کلاسه که هدف آن اختصاص یکی از چندین کلاس ممکن به یک ورودی است مفید است. برای مثال، اگر می‌خواهید یک تصویر را به یکی از ۱۰ دسته طبقه‌بندی کنید (مانند ارقام ۰ تا ۹)، softmax با مقایسه احتمالات نسبی برای هر کلاس، راهی برای تعیین محتمل‌ترین کلاس ارائه می‌کند.

۳. سازگاری با تابع ضرر آنتروپی متقابل (Cross-Entropy Loss)

هنگامی که در ترکیب با تابع ضرر آنتروپی متقابل استفاده می‌شود، softmax به خوبی کار می‌کند زیرا آنتروپی متقابل فاصله بین توزیع احتمال پیش بینی شده (از softmax) و توزیع واقعی (برچسب‌های کدگذاری شده) را اندازه می‌گیرد. خروجی softmax توزیع احتمالی را می‌دهد که می‌تواند با برچسب واقعی مقایسه شود و به حداقل رساندن آنتروپی متقابل مدل را تشویق می‌کند تا کلاس صحیح را با اطمینان بالاتر پیش بینی کند.

۴. مشتق پذیری

تابع softmax قابل مشتق گرفتن است، به این معنی که می‌توان از آن در فرآیند پس انتشار برای به روز رسانی پارامترهای مدل استفاده کرد. گرادینان های صاف آن باعث بهینه سازی کارآمد مدل در طول آموزش می‌شود و به gradient-based optimization ها کمک می‌کند.

۵. Score های normal

تابع softmax خروجی مدل را normal می‌کند تا امتیازهای بزرگ بر پیش‌بینی تسلط نداشته باشند، این مهم زمانی است که score خام (logits) می‌توانند دلخواه یا نامحدود باشند. بدون normal سازی، نمرات خام ممکن است خیلی شدید یا ناهموار باشد که می‌تواند مانع یادگیری شود.

(ب)

- **Variance:** Sensitivity of the model to training data

$$\text{Variance}(x) = \mathbb{E}[(h_w(x) - \mathbb{E}[h_w(x)])^2]$$

Explanation: Complex models tend to overfit

$$h_w(x) = w_0 + w_1x + w_2x^2 + \dots + w_mx^m$$

- Variance dominates when the model is too complex

$$\text{Variance} \gg \text{Bias}$$

- Fits noise, leading to high test error

در یادگیری ماشین، واریانس بالا به مدلی اطلاق می شود که بیش از حد داده های آموزشی را برازش می کند. (بر داده های آموزشی overfit می شود) برازش بیش از حد (overfitting) زمانی اتفاق می افتد که یک مدل نه تنها الگوهای اساسی در داده ها، بلکه نویز یا جزئیات نامربوط را نیز یاد می گیرد. در نتیجه، مدل بر روی داده های آموزشی به خوبی عمل می کند اما در داده های دیده نشده (test) ضعیف عمل می کند.

شاخص های واریانس بالا

- مدل دارای خطای بسیار کم در مجموعه آموزشی است اما خطای بالایی در اعتبارسنجی یا مجموعه تست دارد.
- مدل به تغییرات کوچک در داده های آموزشی بسیار حساس است، به این معنی که به خوبی تعمیم نمی دهد.

علل واریانس بالا (overfitting)

- مدل نسبت به مقدار داده های آموزشی بسیار پیچیده است (به عنوان مثال، استفاده از یک شبکه عصبی عمیق پیچیده برای یک مجموعه داده کوچک).
- مدل دارای پارامترهای بسیار زیادی است که به آن امکان می دهد به نویز موجود در داده های آموزشی fit شود.
- فقدان تکنیک های منظم سازی (L1, L2) برای محدود کردن یادگیری مدل.

راه هایی برای کاهش واریانس در مدل

۱. افزایش مقدار داده های آموزشی
 - با ارائه داده های بیشتر، مدل نمونه های بیشتری برای یادگیری دارد و احتمال برازش بیش از حد نویز را کاهش می دهد.
 - واریانس بالا اغلب زمانی اتفاق می افتد که مدل به دلیل داده های ناکافی، جزئیات بسیار خاصی را یاد می گیرد.
۲. تکنیک های منظم سازی
 - منظم سازی L1/L2: برای محدود کردن وزن مدل، یک عبارت جریمه به تابع ضرر اضافه کنید.
 - L1 (Lasso): sparsity را تشویق می کند، برخی از وزن ها را به صفر می رساند و منجر به مدل های ساده تر می شود.
 - L2 (Ridge): وزن های بزرگ را جریمه می کند و از پیچیده شدن مدل جلوگیری می کند.
 - Drop out (برای شبکه های عصبی): به طور تصادفی نورون ها را در حین آموزش «drop» می کند و شبکه را مجبور می کند تا با جلوگیری از سازگاری مشترک نورون ها، ویژگی های قوی تری را بیاموزد.
۳. ساده تر کردن مدل

- از یک مدل کمتر پیچیده با پارامترهای کمتر استفاده کنید، مانند کاهش تعداد لایه‌ها یا نورون‌ها در یک شبکه عصبی یا انتخاب یک الگوریتم ساده‌تر مانند رگرسیون لجستیک به جای شبکه عصبی عمیق در صورت لزوم.
- هرس کردن ویژگی‌های غیر ضروری یا استفاده از تکنیک‌های کاهش ابعاد مانند PCA نیز می‌تواند کمک کننده باشد.
- ۴. اعتبار سنجی متقابل (Cross-Validation)
 - از تکنیک‌هایی مانند اعتبار سنجی متقابل k-fold استفاده کنید تا اطمینان حاصل کنید که مدل بیش از حد به زیر مجموعه خاصی از داده‌ها تناسب ندارد.
 - با آموزش زیرمجموعه‌های مختلف داده و آزمایش بر روی fold باقی مانده، می‌توانید نحوه تعمیم مدل در پارتیشن‌های مختلف داده را بررسی کنید.
- ۵. توقف زودهنگام (Early Stopping)
 - هنگام آموزش مدل‌های یادگیری عمیق، از توقف زودهنگام استفاده کنید. این تکنیک عملکرد مدل را در یک مجموعه اعتبارسنجی نظارت می‌کند و پس از شروع بدتر شدن عملکرد، آموزش را متوقف می‌کند و از برازش بیش از حد مدل با داده‌های آموزشی جلوگیری می‌کند.
- ۶. داده افزایی (Data Augmentation)

برای داده‌های تصویر، می‌توانید از داده افزایی برای افزایش مصنوعی اندازه مجموعه داده با اعمال تبدیل‌ها (به عنوان مثال، flipping, rotating) به تصاویر آموزشی استفاده کنید. این به تعمیم بهتر مدل کمک می‌کند.

خلاصه

واریانس بالا به این معنی است که مدل شما بیش از حد با داده‌های آموزشی fit شده است و به خوبی به داده‌های دیده نشده تعمیم نمی‌دهد. برای کاهش واریانس، می‌توانید:

- داده‌های آموزشی را افزایش دهید.
- از تکنیک‌های منظم سازی (Drop out، L1/L2) استفاده کنید.
- مدل را ساده کنید.
- اعتبارسنجی متقابل را اعمال کنید.
- از توقف زودهنگام استفاده کنید.
- Data Augmentation را اعمال کنید.

پ) وقتی همه ویژگی‌ها تا حد خوبی با خروجی همبستگی داشته باشند رگرسیون Ridge بر رگرسیون Lasso ترجیح داده می‌شود زیرا Ridge تمایل دارد ضرایب همه ویژگی‌ها را به طور یکنواخت کوچک کند، در حالی که Lasso انتخاب ویژگی را با صفر کردن برخی از ضرایب انجام می‌دهد و به طور موثر ویژگی‌های کم اهمیت تر را حذف می‌کند. در اینجا نگاهی عمیق تر به این که چرا Ridge در این سناریو بهتر کار می‌کند می‌اندازیم:

تفاوت‌های کلیدی بین رگرسیون Ridge و Lasso

- رگرسیون Ridge (L2 regularization) جریمه‌ای متناسب با مجذور ضرایب اضافه می‌کند:

$$\text{Ridge Loss} = \text{MSE} + \lambda \sum_{j=1}^p \beta_j^2$$

این جریمه همه ضرایب را به سمت صفر کوچک می‌کند اما هرگز دقیقاً به صفر نمی‌رسد.

- رگرسیون Lasso (L1 regularization) جریمه‌ای متناسب با قدر مطلق ضرایب اضافه می‌کند:

$$\text{Lasso Loss} = \text{MSE} + \lambda \sum_{j=1}^p |\beta_j|$$

جریمه آن تمایل دارد برخی از ضرایب را دقیقاً به صفر کاهش دهد و انتخاب خودکار ویژگی (feature selection) را انجام دهد.

چرا Ridge در این حالت ترجیح داده می شود؟

۱. Feature Shrinkage vs. Feature Selection

- رگرسیون Ridge ضرایب ویژگی های همبسته را به طور متناسب کوچک می کند، به این معنی که هیچ ویژگی را به طور کامل حذف نمی کند بلکه تأثیر آنها را به طور یکنواخت کاهش می دهد. هنگامی که ویژگی ها به خوبی با خروجی همبستگی دارند، احتمالاً همه آنها حاوی اطلاعات ارزشمندی هستند و حذف هر یک از آنها ممکن است منجر به از دست دادن قدرت پیش بینی شود. (Feature Shrinkage)
- از سوی دیگر، رگرسیون Lasso می تواند برخی از ضرایب را دقیقاً به صفر برساند و برخی از ویژگی ها را به طور کامل از مدل حذف کند. اگر همه ویژگی ها با خروجی مرتبط باشند، حذف هر یک از آنها ممکن است اطلاعات مفید را از بین ببرد و منجر به کاهش عملکرد مدل شود. (Feature Selection)

۲. مدیریت Multicollinearity

- رگرسیون Ridge، Multicollinearity را (زمانی که ویژگی ها با یکدیگر همبستگی دارند) بهتر از Lasso کنترل می کند زیرا وزن ضرایب را در بین ویژگی های همبسته توزیع می کند. از آنجایی که Ridge ضرایب را کوچک می کند اما آنها را صفر نمی کند، به مدل اجازه می دهد تا قدرت پیش بینی همه ویژگی های مرتبط را حفظ کند.
- با این حال، Lasso ممکن است به طور تصادفی یک ویژگی را از گروهی از ویژگی های همبسته انتخاب کند و بقیه را روی صفر قرار دهد، که می تواند باعث بی ثباتی در مدل در هنگام برخورد با داده های بسیار همبسته شود.

۳. ثبات در برآورد ضرایب

- Ridge تخمین های پایدارتری را در حضور همبستگی بالا بین ویژگی ها ارائه می دهد. از آنجایی که ضرایب را بدون اجبار به صفر کوچک می کند، مدل قوی باقی می ماند، حتی اگر درجه بالایی از multicollinearity وجود داشته باشد.
- Lasso می تواند منجر به انتخاب ویژگی ناپایدار شود، به خصوص زمانی که همبستگی بین ویژگی ها بالا باشد، که بسته به اینکه کدام ویژگی حفظ یا حذف شود، منجر به نتایج متناقض می شود.

خلاصه

- Ridge زمانی ترجیح داده می شود که همه ویژگی ها به خوبی با خروجی همبستگی داشته باشند، زیرا ضرایب را بدون حذف هیچ ویژگی کوچک می کند و تضمین می کند که همه ویژگی های همبسته به مدل کمک می کنند.
- Lasso بیشتر برای انتخاب ویژگی مناسب است، و در موقعیت هایی با همبستگی ویژگی های بالا، ممکن است به طور خودسرانه برخی از ویژگی ها را حذف کند و به طور بالقوه منجر به مدلی با ثبات یا عملکرد پایین تر شود.

بنابراین، در مواردی که همه ویژگی ها آموزنده هستند و به خوبی با خروجی همبستگی دارند، رگرسیون Ridge به طور کلی انتخاب بهتری است زیرا در عین مدیریت Overfitting از طریق منظم سازی، تمام ویژگی ها را حفظ می کند.

(ت)

$$\mathbb{E}[(y - h_w(x))^2] = (\text{Bias})^2 + \text{Variance} + \text{Noise}$$

منظم سازی L2 (که در رگرسیون Ridge برای مدل های خطی استفاده می شود) نقش کلیدی در مدیریت تعادل بایاس واریانس در مدل های یادگیری ماشین بازی می کند. در اینجا آمده که چگونه منظم سازی L2 بر این تعادل تأثیر می گذارد:

۱. درک Bias-Variance Trade-Off

- سوگیری (Bias) به خطای ناشی از فرضیات ساده سازی شده در مدل اشاره دارد که باعث می شود مدل underfit شود و الگوهای مهم را از دست بدهد. سوگیری زیاد (High Bias) منجر به خطای زیاد آموزشی و تعمیم ضعیف می شود.
- واریانس به حساسیت مدل نسبت به نوسانات کوچک در داده های آموزشی اشاره دارد. مدلی با واریانس بالا به داده های آموزشی overfit می شود، نویز را یاد می گیرد و منجر به تعمیم ضعیف به داده های دیده نشده می شود.

هدف یافتن تعادل مناسب است که در آن مدل الگوهای اساسی (سوگیری کم) را بدون حساسیت بیش از حد به نویز (واریانس کم) ثبت می‌کند.

۲. اثر منظم سازی L2 بر مبادله بایاس-واریانس

منظم‌سازی L2 یک عبارت جریمه را معرفی می‌کند که با مجذور ضرایب مدل متناسب است:

$$\text{Loss with L2} = \text{MSE (or classification loss)} + \lambda \sum_{j=1}^p \beta_j^2$$

که در آن λ قدرت منظم سازی است (یک هاپیر پارامتر)، و β نشان دهنده ضرایب (وزن) مدل است.

چگونه منظم سازی L2 بر bias و واریانس تأثیر می‌گذارد:

- افزایش bias:
 - منظم‌سازی L2 ضرایب مدل را کوچک می‌کند، به این معنی که انعطاف‌پذیری مدل کاهش می‌یابد. با محدود کردن ضرایب مدل، مرز تصمیم‌گیری یا پیچیدگی مدل را ساده می‌کند. این امر بیش از حد برازش را کاهش می‌دهد، اما به قیمت ایجاد مقدار کمی سوگیری.
 - منظم‌سازی مدل را تشویق می‌کند تا به جای برازش کامل داده‌های آموزشی، بر الگوهای عمومی‌تر تکیه کند، که منجر به خطای آموزشی بالاتر و احتمالاً سوگیری کمی بالاتر می‌شود.
- کاهش واریانس:
 - با جریمه کردن ضرایب بزرگ، منظم سازی L2 از حساس شدن بیش از حد مدل به داده‌های آموزشی جلوگیری می‌کند و در نتیجه توانایی آن را برای تناسب با نویز کاهش می‌دهد. این باعث می‌شود مدل کمتر مستعد تغییرات شدید در پاسخ به تغییرات جزئی در داده‌های آموزشی باشد و واریانس را کاهش دهد.
 - مدلی با واریانس کمتر، بهتر به داده‌های جدید و دیده نشده تعمیم می‌دهد و عملکرد را در مجموعه‌های آزمایشی بهبود می‌بخشد.

۳. چگونه L2 Regularization تعادل مناسب را ایجاد می‌کند؟

- بدون منظم‌سازی، یک طبقه‌بندی‌کننده خطی ممکن است بایاس کم اما واریانس بالایی داشته باشد که منجر به بیش از حد برازش می‌شود. تمام جزئیات داده‌های آموزشی، از جمله نویز، را یاد می‌گیرد و روی داده‌های دیده نشده عملکرد ضعیفی دارد.
 - تنظیم L2 جریمه‌ای برای وزن‌های بزرگ اضافه می‌کند که با نرم‌تر کردن و کلی‌تر کردن مدل، به کاهش بیش از حد برازش (overfitting) کمک می‌کند. با کنترل پیچیدگی مدل، سوگیری را اندکی افزایش می‌دهد، اما واریانس را به میزان قابل توجهی کاهش می‌دهد، که معمولاً منجر به تعمیم بهتر می‌شود.
- با تنظیم قدرت تنظیم λ :
- λ کوچک: مدل دارای آزادی بیشتری است، که منجر به یک مرز تصمیم‌گیری پیچیده‌تر می‌شود، که ممکن است سوگیری را کاهش دهد اما واریانس را افزایش می‌دهد.
 - λ بزرگ: مدل محدودتر می‌شود، واریانس را کاهش می‌دهد اما سوگیری را افزایش می‌دهد. مرز تصمیم‌گیری مدل هموارتر و کلی‌تر می‌شود.

۴. نمایش گرافیکی

به طور کلی، با افزایش قدرت منظم سازی λ :

- سوگیری (bias) مدل به تدریج افزایش می‌یابد (از آنجایی که مدل انعطاف‌پذیری کمتری دارد).
- واریانس به شدت کاهش می‌یابد (از آنجایی که مدل نویز را در داده‌ها fit نمی‌کند).

منظم‌سازی بهینه تعادلی را پیدا می‌کند که در آن مجموع خطاهای بایاس و واریانس به حداقل برسد، که معمولاً منجر به تعمیم بهتر می‌شود.

خلاصه ای از اثر منظم سازی L2 بر تعادل سوگیری-واریانس:

- منظم سازی L2 با کاهش انعطاف پذیری مدل و کاهش توانایی آن در تناسب (fit) کامل با داده های آموزشی، سوگیری را کمی افزایش می دهد.
 - منظم سازی L2 با جلوگیری از fit شدن نويز در داده ها و قوی تر کردن مدل نسبت به داده های دیده نشده، واریانس را کاهش می دهد.
 - پارامتر منظم سازی λ قدرت جریمه را کنترل می کند و به مدل اجازه می دهد تا trade-off بایاس واریانس را هدایت کند و عملکرد تعمیم را بهبود بخشد.
- با معرفی منظم سازی L2، شما اساساً افزایش اندکی در سوگیری را برای کاهش قابل توجهی در واریانس معامله می کنید، که اغلب منجر به عملکرد بهتر در داده های دیده نشده می شود.

سوال دوم

الف) در رگرسیون خطی، هدف ما به حداقل رساندن مجذور خطا بین مقادیر پیش بینی شده و مقادیر واقعی است:

$$\min_w ||y - Xw||^2$$

هنگامی که رگرسیون را فقط روی ویژگی z انجام می دهیم، اساساً در حال حل مسئله ساده تر زیر هستیم، که X به ردیف z (ویژگی) ماتریس داده X اشاره دارد و بردار پیش بینی ها عبارت است از:

$$\hat{y} = w_j X_j$$

خطای مجذور تبدیل می شود به:

$$L(w_j) = \sum_{i=1}^N (y_i - w_j X_{ji})^2$$

ما باید $L(w_j)$ را با توجه به w_j کمینه کنیم. برای یافتن w_j بهینه، مشتق $L(w_j)$ را نسبت به w_j میگیریم:

$$\frac{\partial L(w_j)}{\partial w_j} = -2 \sum_{i=1}^N X_{ji} (y_i - w_j X_{ji}) = -2 \sum_{i=1}^N (X_{ji} y_i - w_j X_{ji}^2)$$

ساده سازی به شکل ماتریسی:

$$\frac{\partial L(w_j)}{\partial w_j} = -2(X_j y - w_j X_j X_j^T)$$

برای یافتن مینیمم خطا، مشتق را برابر صفر قرار دهید:

$$-2(X_j y - w_j X_j X_j^T) = 0 \quad \longrightarrow \quad X_j y - w_j X_j X_j^T = 0$$

حل برای w_j :

$$X_j y = w_j X_j X_j^T \quad \longrightarrow \quad w_j = \frac{X_j y}{X_j X_j^T}$$

بنابراین، نشان داده ایم که وزن w_j هنگام انجام رگرسیون فقط روی ویژگی z است:

$$w_j = \frac{X_j y}{X_j X_j^T}$$

ب) ایده کلیدی در اینجا این است که وقتی ویژگی‌ها مستقل هستند، term های متقابل در ماتریس کوواریانس $X^T X$ صفر می‌شوند و منجر به ساختار block-diagonal می‌شوند. این ویژگی مسئله را ساده می‌کند تا هر ویژگی به طور مستقل حل شود.

برای رگرسیون خطی با همه در نظر گرفتن همه ویژگی‌ها، ما عبارت زیر را به حداقل می‌رسانیم:

$$L(w) = ||y - Xw||^2$$

راه حل بهینه کلی برای رگرسیون خطی چند متغیره به شکل زیر است که در اسلاید های درس اثبات کردیم:

$$w = (X^T X)^{-1} X^T y$$

از آنجایی که ویژگی‌ها مستقل هستند، $(X^T X)_{ij} = 0$ برای $i \neq j$. بنابراین، ماتریس کوواریانس $X^T X$ به یک ماتریس diagonal تبدیل می‌شود:

$$X^T X = \text{diag}(X_1 X_1^T, X_2 X_2^T, \dots, X_L X_L^T)$$

همان طور که میدانیم معکوس آن نیز diagonal است که درایه های قطر اصلی آن به شکل زیر هستند:

$$X^T X^{-1} = \text{diag}\left(\frac{1}{X_1 X_1^T}, \frac{1}{X_2 X_2^T}, \dots, \frac{1}{X_L X_L^T}\right)$$

پس معادله بردار وزن ها $w = (X^T X)^{-1} X^T y$ به شکل زیر ساده می‌شود:

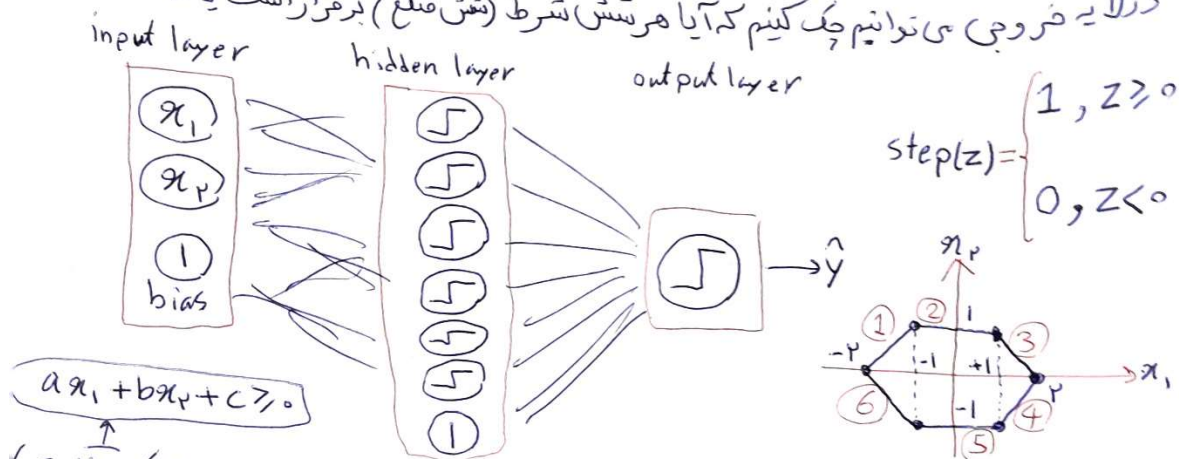
$$w = \text{diag}\left(\frac{1}{X_1 X_1^T}, \frac{1}{X_2 X_2^T}, \dots, \frac{1}{X_L X_L^T}\right) X^T y$$

بنابراین برای هر ویژگی j در بردار w داریم:

$$w_j = \frac{X_j^T y}{X_j X_j^T}$$

که این عبارت (پارامترهای بهینه) دقیقاً همان نتیجه ای است که هنگام انجام رگرسیون روی هر ویژگی به طور مستقل داریم، همانطور که در قسمت (الف) نشان داده شده است. از این رو، نتیجه می‌گیریم که وقتی ویژگی‌ها مستقل هستند، پارامترهای بهینه به‌دست آمده از آموزش رگرسیون روی همه ویژگی‌ها با پارامترهای به‌دست آمده از آموزش روی هر ویژگی به‌طور مستقل یکسان است.

۳- برای حل این مسئله نیاز است تا شش ضلعی را به ۶ ضلع (linear boundary) تقسیم کنیم و سپس هر ضلع با یک نورون قابل تشخیص است. در نهایت نیز با استفاده از یک نورون در لایه خروجی می‌توانیم چک کنیم که آیا هر شش شرط (شش ضلع) برقرار است یا نه.



حال باید معادله داخل و روی شش ضلعی را برای هر ضلع بنویسیم. (هر ضلع در واقع یک خط است).

① $\rightarrow x_2 = x_1 + 2$ پایین دردی خط $\rightarrow x_2 \leq x_1 + 2 \rightarrow x_1 - x_2 + 2 \geq 0$

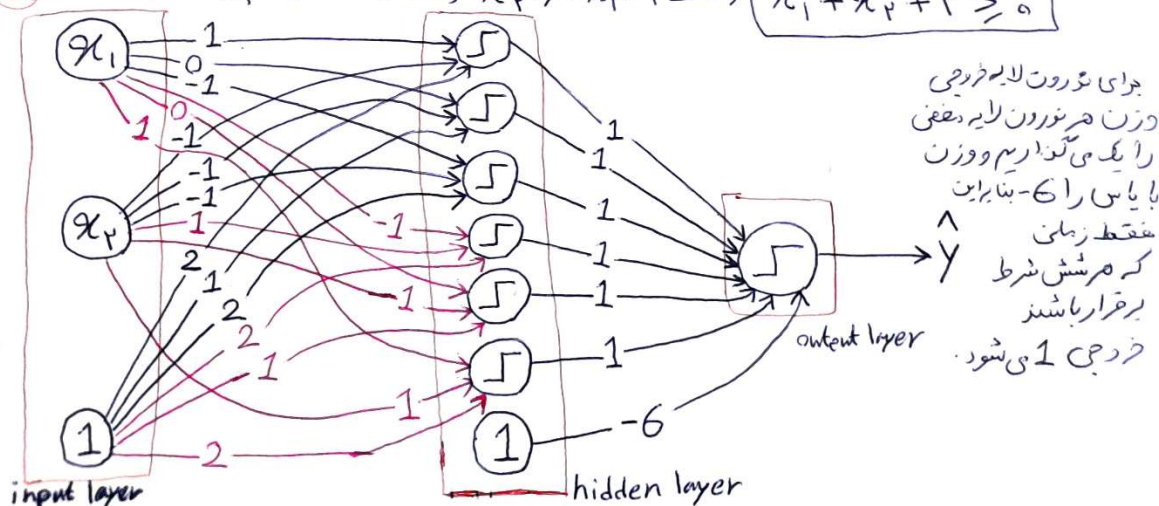
② $\rightarrow x_2 = 1$ پایین دردی خط $\rightarrow x_2 \leq 1 \rightarrow -x_2 + 1 \geq 0$

③ $\rightarrow x_2 = -x_1 + 2$ پایین و چپ دردی خط $\rightarrow x_2 \leq -x_1 + 2 \rightarrow -x_1 - x_2 + 2 \geq 0$

④ $\rightarrow x_2 = x_1 - 2$ بالا دردی خط $\rightarrow x_2 \geq x_1 - 2 \rightarrow x_2 - x_1 + 2 \geq 0$

⑤ $\rightarrow x_2 = -1$ بالا دردی خط $\rightarrow x_2 \geq -1 \rightarrow x_2 + 1 \geq 0$

⑥ $\rightarrow x_2 = -x_1 - 2$ بالا و چپ دردی خط $\rightarrow x_2 \geq -x_1 - 2 \rightarrow x_1 + x_2 + 2 \geq 0$



Case 1: $k = i$ $\frac{\partial \hat{y}_i}{\partial z_i} = \frac{\partial}{\partial z_i} \left(\frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \right)$ $\hat{y}_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$ (الف)

quotient rule $\rightarrow \frac{\partial \hat{y}_i}{\partial z_i} = \frac{e^{z_i} \left(\sum_{j=1}^n e^{z_j} \right) - e^{z_i} \cdot e^{z_i}}{\left(\sum_{j=1}^n e^{z_j} \right)^2} = \frac{e^{z_i} \left(\sum_{j=1}^n e^{z_j} - e^{z_i} \right)}{\left(\sum_{j=1}^n e^{z_j} \right)^2} = \frac{e^{z_i} \left(\sum_{j \neq i} e^{z_j} \right)}{\left(\sum_{j=1}^n e^{z_j} \right)^2}$

$\frac{\partial \hat{y}_i}{\partial z_i} = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \times \frac{\sum_{j \neq i} e^{z_j}}{\sum_{j=1}^n e^{z_j}} = \hat{y}_i \times (1 - \hat{y}_i) \rightarrow \boxed{\frac{\partial \hat{y}_i}{\partial z_i} = \hat{y}_i (1 - \hat{y}_i)}$

Case 2: $k \neq i$ $\frac{\partial \hat{y}_i}{\partial z_k} = \frac{\partial}{\partial z_k} \left(\frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \right) = \frac{0 - e^{z_i} e^{z_k}}{\left(\sum_{j=1}^n e^{z_j} \right)^2}$

$\frac{\partial \hat{y}_i}{\partial z_k} = - \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \times \frac{e^{z_k}}{\sum_{j=1}^n e^{z_j}} = - \hat{y}_i \hat{y}_k \rightarrow \boxed{\frac{\partial \hat{y}_i}{\partial z_k} = - \hat{y}_i \hat{y}_k}$

$k = i \rightarrow \boxed{\frac{\partial \hat{y}_i}{\partial z_i} = \hat{y}_i (1 - \hat{y}_i)}$, $k \neq i \rightarrow \boxed{\frac{\partial \hat{y}_i}{\partial z_k} = - \hat{y}_i \hat{y}_k}$ خلاصه:

ب) $L = - \sum_{i=1}^n y_i \log(\hat{y}_i)$ از قانون زنجیری (chain rule) استفاده می‌کنیم.

$\frac{\partial L}{\partial z_k} = \frac{\partial L}{\partial \hat{y}_i} \times \frac{\partial \hat{y}_i}{\partial z_k} = \frac{\partial}{\partial \hat{y}_i} \left(- \sum_{i=1}^n y_i \log(\hat{y}_i) \right) \left(\frac{\partial \hat{y}_i}{\partial z_k} \right) = - \sum_{i=1}^n \frac{\partial}{\partial \hat{y}_i} (y_i \log(\hat{y}_i)) \left(\frac{\partial \hat{y}_i}{\partial z_k} \right)$

$\frac{\partial L}{\partial z_k} = \left(- \sum_{i=1}^n \frac{y_i}{\hat{y}_i} \right) \left(\frac{\partial \hat{y}_i}{\partial z_k} \right)$ باتوجه به قسمت الف) و تقسیم می‌کنیم $\frac{\partial L}{\partial z_k} = - \frac{y_k}{\hat{y}_k} \times \hat{y}_k (1 - \hat{y}_k) - \sum_{i \neq k} \frac{y_i}{\hat{y}_i} \times (- \hat{y}_i \hat{y}_k)$

$\frac{\partial L}{\partial z_k} = - y_k (1 - \hat{y}_k) + \sum_{i \neq k} y_i \hat{y}_k = - y_k + y_k \hat{y}_k + \sum_{i \neq k} y_i \hat{y}_k = - y_k + \sum_{i=1}^n y_i \hat{y}_k$ یکی می‌کنیم

$\frac{\partial L}{\partial z_k} = - y_k + \hat{y}_k \sum_{i=1}^n y_i$ چون $\sum_{i=1}^n y_i = 1$ پس داریم $\sum_{i=1}^n y_i = 1$ و y_i های واقعی را همی است $\boxed{\frac{\partial L}{\partial z_k} = - y_k + \hat{y}_k = \hat{y}_k - y_k}$

این یک نتیم بسیار ساده است. نشان می‌دهد که گرادیان فقط تفاوت بین احتمال پیش‌بینی شده و بر صیب واقعی است که حس شهودی دارد. ما می‌خواهیم پیش‌بینی‌ها را متناسب با فاصله آنها با بر صیب‌های واقعی تنظیم کنیم.

بی دایم که $\hat{\beta}^{LS} = (X^T X)^{-1} X^T y$ و $\hat{\beta}^{Ridge}(\lambda) = (X^T X + \lambda I)^{-1} X^T y$

$$\hat{\beta}^{LS} = (X^T X)^{-1} X^T y = (X^T X)^{-1} X^T (X\beta + \varepsilon)$$

$$\hat{\beta}^{LS} = \beta + (X^T X)^{-1} X^T \varepsilon$$

error
ضرایب

حال فرض می کنیم error توزیع متقابل را دارد. $\varepsilon \sim N(0, \sigma^2)$

$$\text{Var}(\hat{\beta}^{LS}) = \text{Var}((X^T X)^{-1} X^T \varepsilon) = (X^T X)^{-1} X^T \text{Var}(\varepsilon) X (X^T X)^{-1}$$

$$\text{Var}(\varepsilon) = \sigma^2 I$$

$$\text{Var}(\hat{\beta}^{LS}) = \sigma^2 (X^T X)^{-1}$$

$$\hat{\beta}^{Ridge}(\lambda) = (X^T X + \lambda I)^{-1} X^T y = (X^T X + \lambda I)^{-1} X^T (X\beta + \varepsilon)$$

$$\hat{\beta}^{Ridge}(\lambda) = (X^T X + \lambda I)^{-1} X^T X \beta + (X^T X + \lambda I)^{-1} X^T \varepsilon$$

فرض می کنیم که $\varepsilon \sim N(0, \sigma^2)$

$$\text{Var}(\hat{\beta}^{Ridge}(\lambda)) = \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}$$

۱۳۹۹
۸ می خواهیم اثبات کنیم اگر $\lambda > 0$ آنگاه $Var(\hat{\beta}^{LS}) > Var(\hat{\beta}^{Ridge}(\lambda))$

۹ $X^T X = A \rightarrow Var(\hat{\beta}^{LS}) = \sigma^2 A^{-1}$

۱۰ $Var(\hat{\beta}^{Ridge}(\lambda)) = \sigma^2 (A + \lambda I)^{-1} A (A + \lambda I)^{-1}$

۱۱ فرض می کنیم تجزیه ماتریس متغایر درهای $Q = Q^T$ ویژه
۱۲ $A = Q \Lambda Q^T \rightarrow$ مقدارهای ویژه A به این شکل است. ماتریس قطری شامل مقادیر Λ ویژه A

۱۳ $Var(\hat{\beta}^{LS}) = \sigma^2 Q \Lambda^{-1} Q^T$

۱۴ $Var(\hat{\beta}^{Ridge}(\lambda)) = \sigma^2 Q (\Lambda + \lambda I)^{-1} \Lambda (\Lambda + \lambda I)^{-1} Q^T$

۱۵ برای هر مقدار ویژه a_i از A داریم
۱۶ $\rightarrow Var(\hat{\beta}^{LS}, i) = \frac{\sigma^2}{a_i}$ به این شکل است.

۱۷ $Var(\hat{\beta}^{Ridge}(\lambda), i) = \frac{\sigma^2 a_i}{(a_i + \lambda)^2}$ (به عناصر قطری توجه می کنیم.)
۱۸ \downarrow diagonal

۱۹ حال با این اثبات کنیم اگر $\lambda > 0$ باشد آنگاه $\frac{\sigma^2}{a_i} > \frac{\sigma^2 a_i}{(a_i + \lambda)^2}$

۲۰ $\frac{1}{a_i} > \frac{a_i}{(a_i + \lambda)^2} \xrightarrow{\times a_i (a_i + \lambda)^2} (a_i + \lambda)^2 > a_i^2 \rightarrow a_i + \lambda > a_i$

۱۳۹۹

$$a_i^2 + 2a_i\lambda + \lambda^2 > a_i^2 \rightarrow 2a_i\lambda + \lambda^2 > 0$$

چون $\lambda > 0$ و $a_i \geq 0$ پس این نامعادله همیشه درست است.

$$\text{یعنی } \frac{\sigma^2}{a_i} > \frac{\sigma^2 a_i}{(a_i + \lambda)^2} \quad \lambda > 0$$

$$\lambda > 0 \quad \text{آر } \text{Var}(\hat{\beta}^{LS}) > \text{Var}(\hat{\beta}^{Ridge}(\lambda))$$

12

(ب)

مرحله ۱: بیان واریانس پیش بینی ها

در اینجا، D_x یک ماتریس diagonal حاوی مقادیر تکن X است و $\hat{Y}(\lambda)$ مقدار پیش بینی شده برای رگرسیون Ridge است.

پیش بینی $\hat{Y}(\lambda)$ برای رگرسیون Ridge توسط فرم بسته زیر انجام میشود:

$$\hat{Y}(\lambda) = X\hat{\beta}^{Ridge}(\lambda)$$

$$\hat{\beta}^{Ridge}(\lambda) = (X^T X + \lambda I)^{-1} X^T y$$

$$\hat{Y}(\lambda) = X\hat{\beta}^{Ridge}(\lambda) = X(X^T X + \lambda I)^{-1} X^T y$$

با استفاده از عبارت واریانس برای $\hat{\beta}^{Ridge}(\lambda)$ ، می توانیم واریانس $\hat{Y}(\lambda)$ را محاسبه کنیم. واریانس $\hat{Y}(\lambda)$ برابر است با:

$$\text{Var}[\hat{Y}(\lambda)] = X \text{Var}[\hat{\beta}^{Ridge}(\lambda)] X^T$$

$$\text{Var}[\hat{\beta}^{Ridge}(\lambda)] = \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}$$

$$\text{Var}[\hat{Y}(\lambda)] = \sigma^2 X (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} X^T$$

مرحله ۲: استفاده از SVD

با استفاده از تجزیه مقدارهای منفرد X (SVD)، داریم:

$$X = U D_x V^T$$

که در آن U و V ماتریس های متعامد (orthogonal) هستند و D_x ماتریس diagonal حاوی مقادیر منفرد X است. می توانیم واریانس پیش بینی ها را به صورت زیر بازنویسی کنیم:

$$\text{Var}[\hat{Y}(\lambda)] = \sigma^2 U D_x V^T \cdot (V D_x^T U^T U D_x V^T + \lambda I)^{-1} \cdot V D_x^T U^T U D_x V^T \cdot (V D_x^T U^T U D_x V^T + \lambda I)^{-1} V D_x^T U^T$$

$$\text{Var}[\hat{Y}(\lambda)] = \sigma^2 U D_x V^T \cdot (V D_x^T D_x V^T + \lambda I)^{-1} \cdot V D_x^T D_x V^T \cdot (V D_x^T D_x V^T + \lambda I)^{-1} V D_x^T U^T$$

$$\text{Var}[\hat{Y}(\lambda)] = \sigma^2 U D_x (D_x^T D_x + \lambda I)^{-1} D_x^T D_x (D_x^T D_x + \lambda I)^{-1} D_x^T U^T$$

مرحله ۳: trace واریانس

trace ماتریس کوواریانس مجموع عناصر قطری آن است. از آنجایی که U یک ماتریس متعامد است، عملیات trace تحت تبدیل های متعامد ثابت است، بنابراین می توانیم روی ماتریس D_x diagonal تمرکز کنیم. trace می شود:

$$\begin{aligned}\text{tr}\{\text{Var}[\hat{Y}(\lambda)]\} &= \text{tr}\{\sigma^2 U D_x (D_x^\top D_x + \lambda I)^{-1} D_x^\top D_x (D_x^\top D_x + \lambda I)^{-1} D_x^\top U^\top\} \\ \text{tr}\{\text{Var}[\hat{Y}(\lambda)]\} &= \sigma^2 \text{tr}\{U D_x (D_x^\top D_x + \lambda I)^{-1} D_x^\top D_x (D_x^\top D_x + \lambda I)^{-1} D_x^\top U^\top\} \\ \text{tr}\{\text{Var}[\hat{Y}(\lambda)]\} &= \sigma^2 \text{tr}\{D_x (D_x^\top D_x + \lambda I)^{-1} D_x^\top D_x (D_x^\top D_x + \lambda I)^{-1} D_x^\top\} \\ \text{tr}\{\text{Var}[\hat{Y}(\lambda)]\} &= \sigma^2 \text{tr}\{D_x D_x^\top D_x D_x^\top (D_x^\top D_x + \lambda I)^{-2}\} \\ \text{tr}\{\text{Var}[\hat{Y}(\lambda)]\} &= \sigma^2 \text{tr}\{(D_x D_x^\top)^2 (D_x^\top D_x + \lambda I)^{-2}\} \\ \text{tr}\{\text{Var}[\hat{Y}(\lambda)]\} &= \sigma^2 \sum_{j=1}^p (D_x)_{jj}^4 [(D_x)_{jj}^2 + \lambda]^{-2}\end{aligned}$$

این رابطه مورد نیاز را ثابت می کند. نتیجه نشان می دهد که چگونه رگرسیون Ridge بر واریانس پیش بینی ها از طریق مقادیر منفرد X و پارامتر منظم سازی λ تأثیر می گذارد. با افزایش λ ، واریانس کاهش می یابد که خاصیت کاهش واریانس رگرسیون Ridge را نشان می دهد.

پایان