



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



یادگیری ماشین

تمرین شماره ۲

نام و نام خانوادگی

امیرحسین چمی

آذر ۱۴۰۴

فهرست

فهرست شکل‌ها.....	۵
چکیده.....	۶
پرسش ۱ - PARZEN WINDOWS.....	۷
۱-۱. محاسبه‌ی میانگین تخمین پارزن.....	۷
۲-۱. نمودار pnx به ازای hn های مختلف.....	۸
۳-۱. تنظیم hn بر حسب میزان بایاس.....	۹
۴-۱. نمودار pnx به ازای hn مشخص در بازه‌ی خاص.....	۹
پرسش ۲ - SHIFT OF THE EIGENVALUE SPECTRUM.....	۱۱
۱-۲. اثبات ویژگی بردار ویژه.....	۱۱
۲-۲. چرا استفاده از β_{bridge} مناسب تر است.....	۱۱
۳-۲. مقایسه‌ی رگرسیون خطی و ریج.....	۱۲
پرسش ۳ - OPTIMIZING A KERNEL.....	۱۴
۱-۳. بخش اول: محاسبه‌ی تابع Loss.....	۱۴
۲-۳. بخش دوم: محاسبه به ازای ۲ بعد.....	۱۴
پرسش ۴ - KNN AND THE CURSE OF DIMENSIONALITY.....	۱۶
۱-۴. بخش A: مقدار مشاهدات مورد استفاده برای پیش‌بینی با $p = 1$	۱۶
۲-۴. بخش B: مقدار مشاهدات مورد استفاده برای پیش‌بینی با $p = 2$	۱۶
۳-۴. بخش C: مقدار مشاهدات مورد استفاده برای پیش‌بینی با $p = 100$	۱۷
۴-۴. بخش D: نتیجه‌گیری از بخش‌های قبل.....	۱۷
پرسش ۵ - CHOOSING METRICS.....	۱۸
۱-۵. بخش اول: طبقه‌بندی سه کلاس نوشیدنی.....	۱۸
۲-۵. بخش دوم: طبقه‌بندی چهار کلاس نوشیدنی.....	۱۹
۳-۵. بخش سوم: طبقه‌بندی سه کلاس خودرو.....	۲۰
۴-۵. بخش چهارم: استفاده از درخت تصمیم.....	۲۰

پرسش ۶ - KERNEL RIDGE REGRESSION ۲۲

۱-۶. مینیمم تابع هزینه رگرسیون ریج ۲۲

۲-۶. استفاده از ترفند کرنل در رگرسیون ریج ۲۲

پرسش ۷ - DERIVING LINEAR REGRESSION ۲۵

۱-۷. محاسبه‌ی $\Sigma yx = Ex, y[yx]$ ۲۵

۲-۷. محاسبه‌ی $\Sigma yx = Ex, y[yx]$ با استفاده از نمونه‌ها ۲۵

۳-۷. محاسبه‌ی w * بهینه ۲۵

۴-۷. محاسبه‌ی w * بهینه با استفاده از توزیع جدید ۲۶

۵-۷. اهمیت توزیع داده‌های X ۲۶

پرسش ۸ - LOGISTIC REGRESSION WITH NEWTON'S METHOD ۲۸

۱-۸. محاسبه‌ی $\nabla wJ(w)$ ۲۸

۲-۸. محاسبه‌ی ماتریس هسین ۲۹

۳-۸. رابطه‌ی آپدیت روش نیوتن ۲۹

۴-۸. بررسی مثال عددی ۳۰

پرسش ۹ - A BAYESIAN INTERPRETATION OF LASSO ۳۱

۱-۹. بسط احتمال پسین با استفاده از رابطه‌ی بیز ۳۱

۲-۹. بررسی تخمین MAP ۳۲

پرسش ۱۰ - PARZEN WINDOW CLASSIFICATION WITH GAUSSIAN KERNEL ۳۴

..... ۳۴

۱-۱۰. بخش A: طبقه‌بندی به ازای $h = 1$ ۳۴

۲-۱۰. بخش B: طبقه‌بندی به ازای $h = 1$ ۳۵

پرسش ۱۱ - MAXIMUM LIKELIHOOD ESTIMATION ۳۷

۱-۱۱. بخش A: محاسبه‌ی عددی MLE و مقایسه با حل فرم بسته ۳۷

۲-۱۱. بخش B: بررسی مینیمم شدن منفی LOG LIKELIHOOD ۳۷

۳-۱۱. بخش C: محاسبه‌ی ماتریس کوواریانس تخمینگر در حالت حدی ۳۸

۴-۱۱. بخش D: بهبود ماتریس کوواریانس تخمینگر در حالت حدی ۳۹

۴۰	۵-۱۱. بخش E: فاصله اطمینان ۹۵٪
۴۰	۶-۱۱. بخش F: تست فرض
۴۱	۱-۶-۱۱. <i>Z-test</i>
۴۱	۲-۶-۱۱. <i>Likelihood ratio chi-squared test</i>
۴۲	۳-۶-۱۱. <i>Wald chi-squared test</i>
۴۲	۷-۱۱. بخش G: COEFFICIENT OF VARIATION
۴۲	۱-۷-۱۱. خاصیت <i>invariance</i> نسبت به یک ضریب ثابت مثبت
۴۳	۲-۷-۱۱. <i>coefficient of variation</i> تخمین
۴۳	۳-۷-۱۱. <i>۹۵% Confidence Interval via the Delta Method</i>
۴۵	پرسش ۱۲ - NON-LINEAR REGRESSION USING BASIS FUNCTIONS
۴۵	۳-۱۲. نمودارهای GAUSSIAN BASIS REGRESSION به ازای لاندارهای مختلف
۴۵	۱-۳-۱۲. به ازای $\lambda = 0.001$
۴۶	۲-۳-۱۲. به ازای $\lambda = 0.05$
۴۸	۴-۱۲. نمودار میانگین مربع خطاها بر حسب λ
۵۰	مراجع

فهرست شکل‌ها

- شکل ۱. نمودار $pnx - x$ ۸
- شکل ۲. نمودار pnx به ازای مقادیر مشخص در بازه‌ی $0, 0.05$ ۱۰
- شکل ۳. نمودار رگرسیون به ازای $\lambda = 0.001$ ۴۶
- شکل ۴. نمودار رگرسیون به ازای $\lambda = 0.05$ ۴۷
- شکل ۵. نمودار رگرسیون به ازای $\lambda = 5$ ۴۸
- شکل ۶. نمودار میانگین مربع خطا به ازای مقادیر مختلف λ ۴۹

این سری از تمرین‌ها ترکیبی از مباحث تئوری و تحلیلی به همراه استدلال ریاضی و شهودی بودند و هدف اصلی آن‌ها، درک عمیق‌تر روش‌های مختلف یادگیری، تخمین توزیع، رگرسیون، طبقه‌بندی و تحلیل رفتار الگوریتم‌ها در شرایط مختلف بود. در طول تمرین، هم با روش‌های کلاسیک مانند رگرسیون خطی و لجستیک و هم با روش‌های پیشرفته‌تری مثل کرنل‌ها، regularization و تفسیر بیزی مدل‌ها سروکار داشتیم.

حجم تمرین نسبتاً زیاد و شامل سوالات متنوعی بود که هر کدام به جنبه خاصی از یادگیری ماشین می‌پرداختند. انجام این تمرین‌ها نیازمند صرف زمان قابل توجه و دقت در جزئیات محاسباتی و مفهومی بود و در نهایت باعث شد ارتباط بین مباحث مختلف درس و تأثیر انتخاب پارامترها و مدل‌ها در عمل بهتر درک شود.

این سری از تمرین‌ها به قدری زیاد بودند که نمی‌توان در این چکیده به تک تک سوالات پرداخت و هدف آن‌ها را ریز به ریز بیان کرد اما صراحتاً می‌توان گفت که این تمرین از تمرین قبلی عمیق‌تر بود و سعی می‌کرد مفاهیم (چه آن‌ها که تدریس شده و چه آن‌ها که تازه معرفی می‌شوند) را به صورت عمیق و پایه‌ای جا بیندازد.

پیش ۱ – Parzen Windows

در این سوال داده هایی با توزیع یکنواخت $p(x) \sim U(0, a)$ و تابع کرنل زیر داریم.

$$\varphi(x) = \begin{cases} e^{-x} & x > 0 \\ 0 & x \leq 0 \end{cases} \quad (1)$$

می خواهیم میانگین تخمین پنجره پارزن محاسبه و تاثیر پارامتر smoothing یعنی h_n را بررسی کنیم.

۱-۱. محاسبه ی میانگین تخمین پارزن

از جزوه ی درس می دانیم، توزیع تخمین زده شده به شکل زیر است.

$$p_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} \varphi\left(\frac{x - x_i}{h_n}\right) \quad (2)$$

و میانگین آن به صورت زیر بدست می آید.

$$\bar{p}_n(x) = E[p_n(x)] = \frac{1}{n} \sum_{i=1}^n E\left[\frac{1}{V_n} \varphi\left(\frac{x - x_i}{h_n}\right)\right] \quad (3)$$

با فرض i.i.d بودن داده های x_i تمامی expectation ها یکسان می شوند. پس داریم:

$$\bar{p}_n(x) = \frac{1}{n} \cdot n E\left[\frac{1}{V_n} \varphi\left(\frac{x - y}{h_n}\right)\right] = \int \frac{1}{V_n} \varphi\left(\frac{x - y}{h_n}\right) p(y) dy \quad (4)$$

برای حل این انتگرال ابتدا باید تکلیف بازه های غیر صفر تابع کرنل و $p(y)$ را مشخص کنیم.

$$\varphi\left(\frac{x - y}{h_n}\right) \rightarrow \frac{x - y}{h_n} > 0 \rightarrow y < x \quad (5)$$

$$p(y) \rightarrow 0 < y < a \quad (6)$$

پس در کل داده های y می توانند بین 0 تا $\min(a, x)$ باشند تا توابع بالا مقدار غیر صفر داشته باشند. حال می توانیم انتگرال (۴) را بنویسیم.

$$\begin{aligned} \bar{p}_n(x) &= \frac{1}{ah_n} \int_0^{\min(a, x)} e^{-\frac{(x-y)}{h_n}} dy = \frac{1}{ah_n} e^{-\frac{x}{h_n}} \int_0^{\min(a, x)} e^{\frac{y}{h_n}} dy \\ &= \frac{1}{ah_n} e^{-\frac{x}{h_n}} h_n \left(e^{\frac{y}{h_n}} \Big|_0^{\min(a, x)} \right) \end{aligned} \quad (7)$$

پس در نهایت خواهیم داشت.

$$\bar{p}_n(x) = \frac{1}{a} e^{-\frac{x}{h_n}} \left(e^{\frac{\min(a, x)}{h_n}} - 1 \right) \quad (8)$$

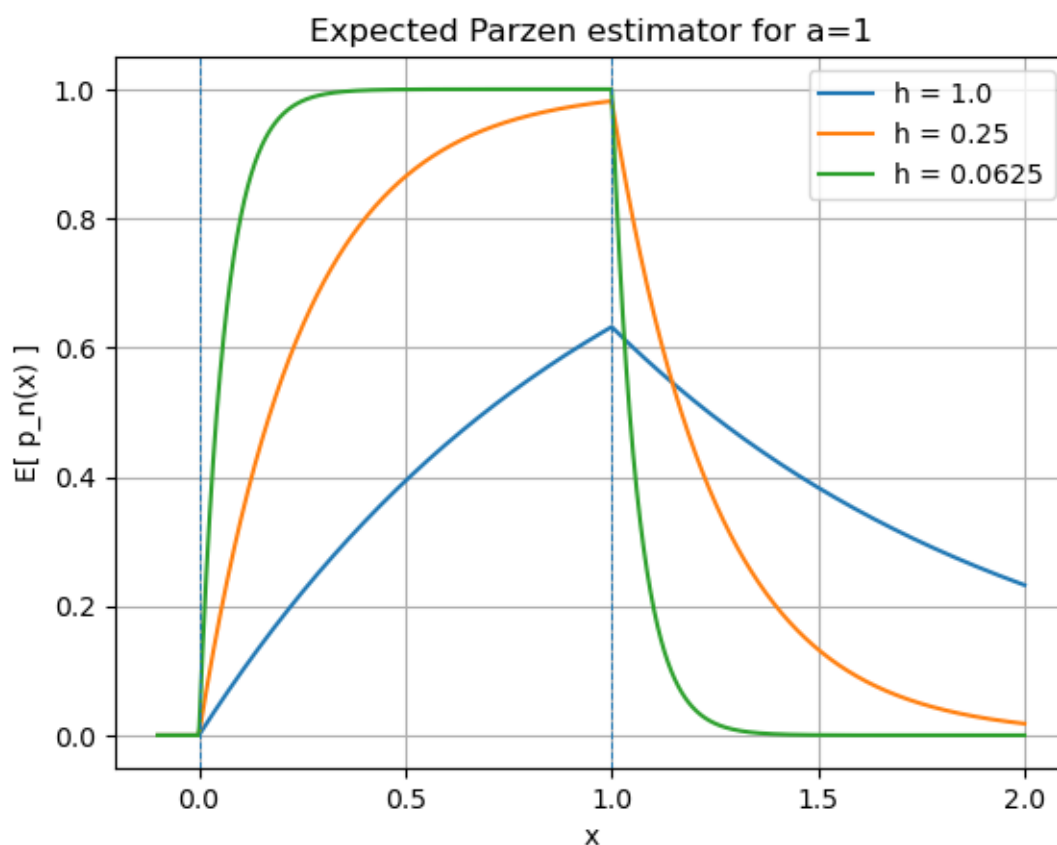
حال به ازای مقادیر مختلف x می توانیم (۸) را به یک تابع چند ضابطه ای تبدیل کنیم. اگر $x < 0$ آنگاه بازه‌ی $[0, x]$ اصلاً قابل تعریف نیست. اگر $0 \leq x \leq a$ آنگاه $\min(a, x) = x$ و اگر $x \geq a$ آنگاه $\min(a, x) = a$ می شود.

$$\bar{p}_n(x) = \begin{cases} 0 & x < 0 \\ \frac{1}{a} \left(1 - e^{-\frac{x}{h_n}} \right) & 0 \leq x \leq a \\ \frac{1}{a} \left(e^{\frac{a}{h_n}} - 1 \right) e^{-\frac{x}{h_n}} & x \geq a \end{cases} \quad (9)$$

پس به این شکل میانگین مورد نظر حساب شد.

۲-۱. نمودار $\bar{p}_n(x)$ به ازای h_n های مختلف

در این بخش نمودار $\bar{p}_n(x)$ به x را با فرض $h_n = 1, 0.25$ و 0.0625 و $a = 1$ رسم می کنیم. (برای رسم خط های عمودی به ازای $x = 0$ و $x = a$ از مرجع [۱] استفاده کردم).



شکل ۱. نمودار $\bar{p}_n(x) - x$

۳-۱. تنظیم h_n بر حسب میزان بایاس

مقدار $E[p_n(x)]$ در بازه $0 \leq x \leq a$ را در رابطه‌ی (۹) دیدیم. بر این اساس مقدار بایاس به صورت زیر است.

$$Bias = E[p_n(x)] - p(x) = \frac{1}{a} \left(1 - e^{-\frac{x}{h_n}} \right) - \frac{1}{a} = -\frac{1}{a} e^{-\frac{x}{h_n}} \quad (10)$$

پس در کل اندازه‌ی بایاس می‌شود: $|Bias| = \frac{1}{a} e^{-\frac{x}{h_n}}$

اندازه‌ی بایاس باید کمتر از یک درصد باشد. یعنی کمتر از یک درصد چگالی واقعی.

$$|Bias| < 0.01 p(x) = 0.01 \frac{1}{a} \quad (11)$$

$$\frac{1}{a} e^{-\frac{x}{h_n}} < 0.01 \frac{1}{a} \quad (12)$$

$$e^{-\frac{x}{h_n}} < 0.01$$

پس با لگاریتم گیری خواهیم داشت.

$$-\frac{x}{h_n} < \ln(0.01) = -4.605170186 \quad (13)$$

$$h_n < \frac{x}{4.605170186} \quad (14)$$

بازه $[0, a]$ طولی به اندازه‌ی a دارد. باید این بایاس برای هر بازه‌ی $[0.01a, a]$ برقرار باشد. کمترین مقدار این بازه $x_{min} = 0.01a$ است. پس رابطه‌ی (۱۴) به صورت زیر خواهد بود.

$$h_n < \frac{0.01a}{4.605170186} \quad (15)$$

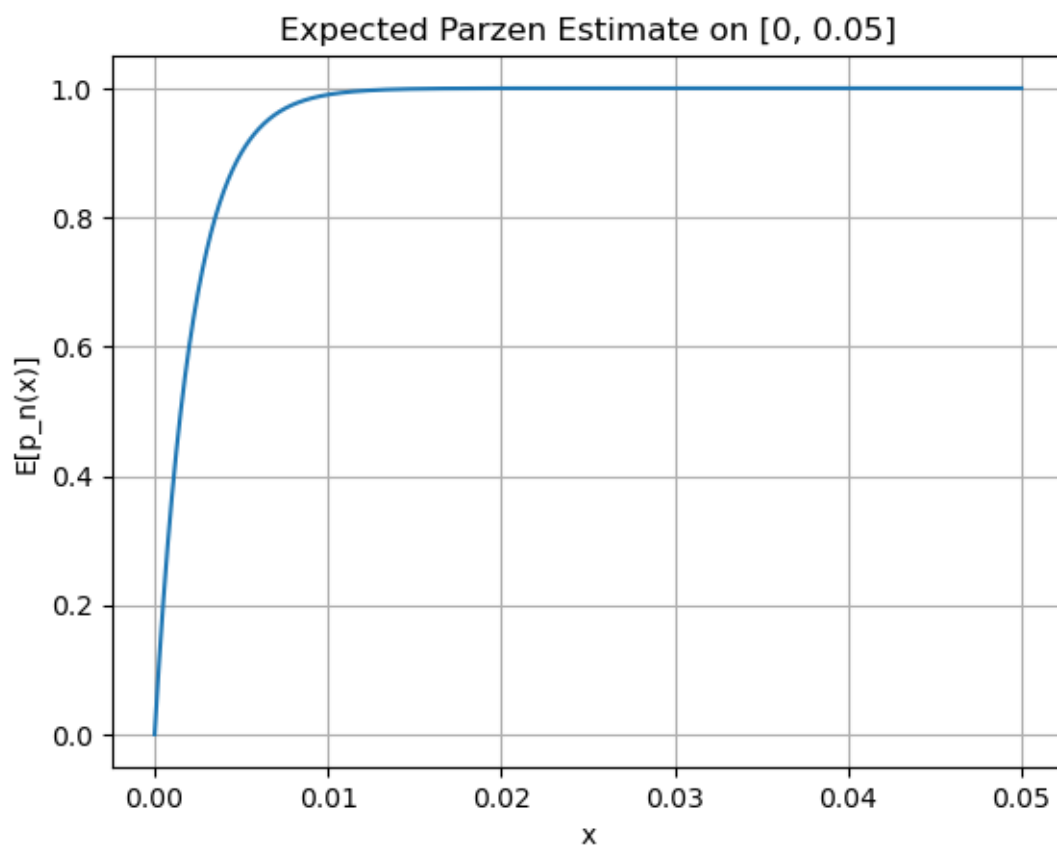
طبق رابطه‌ی (۱۵) حداکثر مقدار h_n بدست آمد. از طرفی طبق شرایط پارزن که در جزوه‌ی درس هم دیدیم، $\lim_{n \rightarrow \infty} h_n \rightarrow 0$ و از طرفی $\lim_{n \rightarrow \infty} nh_n \rightarrow \infty$ است. این به این معنی است که کمترین مقدار h_n صفر است ولی به صورت عملی ممکن نیست.

۴-۱. نمودار $\bar{p}_n(x)$ به ازای h_n مشخص در بازه‌ی خاص

در رابطه‌ی (۱۵) حداکثر مقدار h_n حساب شد. با فرض $a = 1$ خواهیم داشت.

$$h_n < \frac{0.01}{4.605170186} = 0.0021714724 \quad (16)$$

با استفاده از مقدار به دست آمده، $\bar{p}_n(x)$ در بازه‌ی $[0, 0.05]$ به صورت زیر است.



شکل ۲. نمودار $\bar{p}_n(x)$ به ازای مقادیر مشخص در بازه $[0, 0.05]$

پرسش ۲ – Shift of the Eigenvalue Spectrum

هدف این سوال همان طور که صراحتاً گفته شده، مروری بر جبر خطی است.

۱-۲. اثبات ویژگی بردار ویژه

می دانیم اگر ξ بردار ویژه برای ماتریس های A و B باشد باید روابط $A\xi = \lambda_A \xi$ و $B\xi = \lambda_B \xi$ برقرار باشد. پس برای اثبات باید به رابطه ای به همین شکل برسیم.

$$(A + B)\xi = A\xi + B\xi = \lambda_A \xi + \lambda_B \xi = (\lambda_A + \lambda_B)\xi \quad (17)$$

به همین سادگی اثبات شد که ξ بردار ویژهی ماتریس $A + B$ هم هست.

۲-۲. چرا استفاده از $\hat{\beta}^{ridge}$ مناسب تر است

می خواهیم نشان دهیم چرا چرا استفاده از $\hat{\beta}^{ridge}$ مناسب تر است. ابتدا با استفاده از تجزیه به مقادیر ویژه، ماتریس $X^T X$ که قطعاً متقارن است را تجزیه می کنیم.

$$X^T X = Q \Lambda Q^T \quad (18)$$

در این رابطه Q یک ماتریس ارتوگونال است که ستون های بردارهای ویژه ماتریس $X^T X$ هستند. و $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ یک قطری است.

حال از نتیجهی بخش ۱-۲ استفاده می کنیم. فرض می کنیم $A = X^T X$ و $B = \lambda I$ است. توجه کنید که هر برداری در \mathbb{R}^n یک بردار ویژه برای ماتریس I است. بر این اساس برای هر بردارویژهی q_i مقدار ویژهی متناظر با آن در ماتریس $X^T X + \lambda I$ برابر $\lambda_i + \lambda$ است.

$$(X^T X + \lambda I)q_i = X^T X q_i + \lambda I q_i = \lambda_i q_i + \lambda q_i = (\lambda_i + \lambda)q_i \quad (19)$$

چه فایده ای برای ما دارد؟

- اضافه شدن λI باعث افزایش مقادیرویژه می شود در حالی بردارویژه ها بدون تغییر می مانند.
- اگر ماتریس $X^T X$ تکین یا نزدیک به تکین باشد (یعنی بعضی از مقادیرویژه صفر یا نزدیک صفر باشند) باعث می شود که معکوس این ماتریس وجود نداشته باشد یا اینکه مقادیر نزدیک صفر باعث مقادیر بسیار بزرگ در معکوس شوند که همین باعث ناپایداری می شود. به همین دلیل باید λ مقدار مثبت و به قدر کافی بزرگ داشته باشد تا اطمینان داشته باشیم $X^T X + \lambda I$ همواره معکوس پذیر باشد.

- مورد سوم (که بی ربط به مورد دوم نیست) این است که مقدار condition number افزایش می یابد. طبق [۲] مقدار condition number معیاری است برای نشان دادن میزان حساسیت یک تابع به تغییرات در ورودی. برای ماتریس متقارن $X^T X$ مقدار condition number به صورت $K = \frac{\lambda_{max}}{\lambda_{min}}$ تعریف می شود که اگر $\lambda_{min} \approx 0$ در این صورت مقدار K بسیار بزرگ خواهد شد یعنی مقادیر کوچک می توانند باعث خطاهای بسیار بزرگ شوند و ناپایداری ایجاد کنند. در نتیجه استفاده از $X^T X + \lambda I$ مقدار condition number را کاهش داده و حساسیت به تغییر ورودی و در نتیجه خطا را کم می کند و از ناپایداری جلوگیری می کند.

۳-۲. مقایسه‌ی رگرسیون خطی و ریدج

با فرض اینکه مدل به صورت $y = X\beta + \epsilon$ بوده و $E\{\epsilon\} = 0$ و $var(\epsilon) = \sigma^2 I$ باشد. رگرسیون خطی از به صورت زیر است.

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (20)$$

امیدریاضی $\hat{\beta}$ به صورت زیر محاسبه می شود.

$$\begin{aligned} E\{\hat{\beta}\} &= E\{(X^T X)^{-1} X^T y\} = E\{(X^T X)^{-1} X^T (X\beta + \epsilon)\} \\ &= E\{(X^T X)^{-1} X^T X\beta\} + E\{(X^T X)^{-1} X^T \epsilon\} \\ &= E\{\beta\} + E\{(X^T X)^{-1} X^T\} E\{\epsilon\} = E\{\beta\} = \beta \end{aligned} \quad (21)$$

پس $bias = E\{\hat{\beta}\} - \beta = 0$. این روش بایاس ندارد. واریانس آن به صورت زیر است.

$$\begin{aligned} var(\hat{\beta}) &= var((X^T X)^{-1} X^T y) = var((X^T X)^{-1} X^T (X\beta + \epsilon)) \\ &= var((X^T X)^{-1} X^T \epsilon) \\ &= (X^T X)^{-1} X^T (\sigma^2 I) X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1} \end{aligned} \quad (22)$$

اگر $X^T X$ تکین یا نزدیک به تکین باشد، بعضی مقادیر ویژه‌ی آن خیلی کوچک بوده و در نتیجه $(X^T X)^{-1}$ مقادیری بسیار بزرگی خواهد داشت. پس این واریانس می تواند خیلی بزرگ باشد.

پس در کل با وجود اینکه رگرسیون خطی بایاس ندارد، اما می تواند واریانس خیلی زیادی داشته باشد.

اما در رگرسیون ریدج داریم:

$$\hat{\beta}_{ridge} = (X^T X + \lambda I)^{-1} X^T y \quad (23)$$

امیدریاضی آن به صورت زیر است.

$$\begin{aligned}
 E\{\hat{\beta}_{ridge}\} &= E\{(X^T X + \lambda I)^{-1} X^T y\} \\
 &= E\{(X^T X + \lambda I)^{-1} X^T (X\beta + \epsilon)\} \\
 &= E\{(X^T X + \lambda I)^{-1} X^T X\beta\} \neq \beta
 \end{aligned}
 \tag{24}$$

واضا این امیدریاضی β نیست. پس این روش بایاس غیر صفر دارد. این ویژگی مثبتی است. واریانس آن به صورت زیر است.

$$\begin{aligned}
 var(\hat{\beta}_{ridge}) &= var\left((X^T X + \lambda I)^{-1} X^T y\right) \\
 &= var\left((X^T X + \lambda I)^{-1} X^T (X\beta + \epsilon)\right) \\
 &= var\left((X^T X + \lambda I)^{-1} X^T \epsilon\right) \\
 &= (X^T X + \lambda I)^{-1} X^T (\sigma^2 I) X (X^T X + \lambda I)^{-1} \\
 &= \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}
 \end{aligned}
 \tag{25}$$

می توانیم این واریانس را با استفاده از تجزیه به مقادیر ویژه محاسبه کنیم. با فرض تجزیه ویژه $X^T X = U \Lambda U^T$ می توان واریانس رگرسیون ریدج را به صورت زیر نوشت.

$$var(\hat{\beta}_{ridge}) = \sigma^2 \frac{\lambda_i}{(\lambda_i + \lambda)^2} \tag{26}$$

در حالی که با استفاده از تجزیه ویژه برای واریانس رگرسیون خطی خواهیم داشت:

$$var(\hat{\beta}) = \frac{\sigma^2}{\lambda_i} \tag{27}$$

واضح است که واریانس رگرسیون ریدج کمتر از رگرسیون خطی است. پس در کل رگرسیون ریدج بایاس غیر صفر دارد اما واریانس کمتر از رگرسیون خطی است.

اگر بایاس صفر برای ما اهمیت داشته باشد به سراغ رگرسیون خطی می رویم و اگر واریانس کمتر مدل اهمیت داشته باشد باید به سراغ رگرسیون ریدج برویم.

پرسش ۳ – Optimizing a Kernel

این سوال درباره‌ی رگرسیون مبتنی بر کرنل است. هدف این است که تابع هزینه را برای این سوال بدست آورده و به ازای ۲ بعد هم تحلیل کنیم.

۳-۱. بخش اول: محاسبه‌ی تابع Loss

طبق صورت سوال تابع کرنل به صورت زیر است.

$$K(x_n, x_i) = \exp(-(x_n - x_i)W(x_n - x_i)^T) \quad (28)$$

خود سوال برای این بخش یک hint قرار داده که با توجه به آن باید جمع کرنل را روی تمام داده ها به جز خود داده x_i قرار دهیم. تابع پیش بینی به صورت زیر است.

$$f(x_i) = \frac{\sum_{n \neq i} K(x_n, x_i) y_n}{\sum_{n \neq i} K(x_n, x_i)} \quad (29)$$

با توجه به اینکه می خواهیم از residual sum of squares استفاده کنیم پس به صورت زیر loss را تعریف می کنیم.

$$\begin{aligned} L(W) &= \sum_{i=1}^N (y_i - f(x_i))^2 = \sum_{i=1}^N \left(y_i - \frac{\sum_{n \neq i} K(x_n, x_i) y_n}{\sum_{n \neq i} K(x_n, x_i)} \right)^2 \\ &= \sum_{i=1}^N \left(y_i - \frac{\sum_{n \neq i} \exp(-(x_n - x_i)W(x_n - x_i)^T) y_n}{\sum_{n \neq i} \exp(-(x_n - x_i)W(x_n - x_i)^T)} \right)^2 \end{aligned} \quad (30)$$

۳-۲. بخش دوم: محاسبه به ازای ۲ بعد

با توجه به اینکه بعد مسئله ۲ است پس هر داده به صورت $x_n = (x_n^{(1)}, x_n^{(2)})$ تعریف می شود. عبارتی که در توان exponential داریم به صورت زیر نوشته خواهد شد.

$$\begin{aligned} (x_n - x_i)W(x_n - x_i)^T &= W_{11} (x_n^{(1)} - x_i^{(1)})^2 + 2W_{12} (x_n^{(1)} - x_i^{(1)}) (x_n^{(2)} - x_i^{(2)}) \\ &\quad + W_{22} (x_n^{(2)} - x_i^{(2)})^2 \end{aligned} \quad (31)$$

پس بر این اساس تابع loss به صورت زیر است.

$$\begin{aligned}
L(W) &= \sum_{i=1}^N (y_i - f(x_i))^2 \quad (32) \\
&= \sum_{i=1}^N \left(y_i - \frac{\sum_{n \neq i} \exp \left(-W_{11} (x_n^{(1)} - x_i^{(1)})^r - r W_{1r} (x_n^{(1)} - x_i^{(1)}) (x_n^{(r)} - x_i^{(r)}) - W_{rr} (x_n^{(r)} - x_i^{(r)})^r \right) y_n}{\sum_{n \neq i} \exp \left(-W_{11} (x_n^{(1)} - x_i^{(1)})^r - r W_{1r} (x_n^{(1)} - x_i^{(1)}) (x_n^{(r)} - x_i^{(r)}) - W_{rr} (x_n^{(r)} - x_i^{(r)})^r \right)} \right)^2
\end{aligned}$$

(برای اینکه بتوان کسر را به صورت کامل نشان داد مجبور شدم سائز فونت را کوچک تر کنم.)

برای اینکه کوتاه تر بنویسیم می توان متغیر های $\Delta_{ni}^{(1)} = x_n^{(1)} - x_i^{(1)}$ و $\Delta_{ni}^{(r)} = x_n^{(r)} - x_i^{(r)}$ تعریف کرد. بر این اساس می توان کوتاه تر نوشت.

$$L(W) = \sum_{i=1}^N \left(y_i - \frac{\sum_{n \neq i} \exp \left(-W_{11} (\Delta_{ni}^{(1)})^r - r W_{1r} (\Delta_{ni}^{(1)}) (\Delta_{ni}^{(r)}) - W_{rr} (\Delta_{ni}^{(r)})^r \right) y_n}{\sum_{n \neq i} \exp \left(-W_{11} (\Delta_{ni}^{(1)})^r - r W_{1r} (\Delta_{ni}^{(1)}) (\Delta_{ni}^{(r)}) - W_{rr} (\Delta_{ni}^{(r)})^r \right)} \right)^2 \quad (33)$$

(برای اینکه بتوان کسر را به صورت کامل نشان داد مجبور شدم سائز فونت را کوچک تر کنم.)

پرسی ۴ – KNN and the Curse of Dimensionality

در این سوال مفهوم نفرین (نحسی) ابعاد را می فهمیم.

۴-۱. بخش a: مقدار مشاهدات مورد استفاده برای پیش‌بینی با $p = 1$

در این بخش با توجه به مثالی که در متن سوال زده شده، در مجموع می توانیم سه حالت را برای هر داده‌ی تست x در نظر بگیریم.

$$\begin{aligned} 0 &\leq x < 0.05 & (34) \\ 0.05 &\leq x \leq 0.95 \\ 0.95 &< x < 1 \end{aligned}$$

در واقع بازه ی ۰ الی ۱ به سه بخش تقسیم شده. بازه ی $0.05 \leq x \leq 0.95$ بازه ای است که اگر داده‌ی مشاهده در آن قرار گیرد می توانیم به اندازه‌ی ۰.۰۵ کمتر یا بیشتر از آن را بررسی کنیم. اما دو بازه‌ی دیگر جایی هستند که اگر مشاهده در آن قرار گیرد، با کم یا اضافه کردن ۰.۰۵ از بازه‌ی $[0, 1]$ که مسئله در آن تعریف شده است خارج می شویم. هر بازه‌ی بالا را می توان به ازای مشاهده‌ی x به صورت زیر بازنویسی کرد.

$$\begin{aligned} [0, x + 0.05] & \text{ with length } x + 0.05 & (35) \\ [x - 0.05, x + 0.05] & \text{ with length } 0.1 \\ [x - 0.05, 1] & \text{ with length } 1.05 - x \end{aligned}$$

از آنجا که توزیع x به صورت یکنواخت روی بازه ی ۰ الی ۱ است، میانگین داده های مشاهده‌ی مورد نظر با انگرال های زیر بدست می آید.

$$\begin{aligned} E[L] &= \int_0^{0.05} (x + 0.05) dx + \int_{0.05}^{0.95} 0.1 dx + \int_{0.95}^1 (1.05 - x) dx & (36) \\ &= \left(\frac{x^2}{2} + 0.05x \right) \Big|_0^{0.05} + (0.1x) \Big|_{0.05}^{0.95} + \left(1.05x - \frac{x^2}{2} \right) \Big|_{0.95}^1 \\ &= 0.00375 + 0.09 + 0.00375 = 0.0975 \end{aligned}$$

پس ۰.۰۹۷۵ یا ۹.۷۵٪ از داده های موجود برای پیش‌بینی استفاده می شوند. البته به صورت شهودی می توان گفت که با توجه به یکنواخت بودن توزیع داده ها و با توجه به اینکه بازه‌ی $0.05 \leq x \leq 0.95$ (که مقدار بیشتری داده نسبت به بازه‌ی های دیگر را شامل می شود) دقیقا ۹۰٪ فضای مسئله است، پس اگر از بازه های کوچک نزدیک *boundary* ها چشم پوشی کنیم می توان گفت که به تقریب به صورت میانگین از ۱۰٪ داده های موجود برای پیش‌بینی استفاده می کنیم. (منظورم این است که مقدار حدودی ۱۰٪ درصد قبل از محاسبات هم قابل حدس بود.)

۴-۲. بخش b: مقدار مشاهدات مورد استفاده برای پیش‌بینی با $p = 2$

در بخش قبلی به ازای یک ویژگی دیدیم که ۹.۷۵٪ داده ها مورد استفاده قرار می گیرد. حال با اضافه شدن ویژگی دوم و با فرض مستقل بودن ویژگی، میانگین داده های مورد استفاده به صورت زیر خواهد بود.

$$E[L_1 L_2] = E[L_1]E[L_2] = 0.0975 * 0.0975 = 0.00950625 \quad (37)$$

پس ۰.۰۰۹۵۰۶۲۵ یا ۰.۹۵۰۶۲۵٪ داده های موجود برای پیش‌بینی استفاده می شوند. خیلی عدد کمتری شد. شهود آن را به ازای یک ویژگی در بخش قبل دیدیم و بر اساس آن به صورت شهودی به ازای دو ویژگی تنها ۰.۰۱ یا ۱٪ از داده ها استفاده می شوند.

۳-۴. بخش c: مقدار مشاهدات مورد استفاده برای پیش‌بینی با $p = 100$

مشابه بخش های قبل به ازای ۱۰۰ تا ویژگی مستقل میانگین داده های مورد استفاده به صورت زیر خواهد بود.

$$average\ fraction = (E[L])^{100} = 0.0975^{100} \quad (38)$$

که برای همین حالت با استفاده از روش شهودی خواهد بود:

$$average\ fraction = 0.1^{100} \quad (39)$$

پس ۰.۱^{۱۰۰} یا ۱۰^{-۱۰۰}٪ داده های موجود برای پیش‌بینی استفاده می شوند. می بینیم که عدد به دست آمده مقدار بسیار کمی شده.

۴-۴. بخش d: نتیجه گیری از بخش های قبل

در بخش های قبل دیدیم که به ازای یک ویژگی ۱۰٪ از داده ها در نزدیکی نقطه‌ی تست هستند و این مقدار برای دو ویژگی تنها ۱٪ و برای ۱۰۰ ویژگی مقدار بسیار کم ۱۰^{-۱۰۰}٪ است. با اضافه شدن ویژگی ها، داده های نزدیک به داده‌ی تست به صورت نمایی کاهش پیدا می کنند.

فرض روش KNN این است که ما می توانیم تعداد مناسبی داده‌ی train در نزدیکی داده‌ی تست پیدا کنیم. اما وقتی تعداد ویژگی ها زیاد می شود، تقریباً هیچ داده ای در همسایگی نزدیک داده‌ی تست وجود ندارد. به عبارت دیگر نزدیک ترین همسایه در ابعاد بالا، بسیار دور است. پس باید تعداد داده های ما بی نهایت باشد که ممکن نیست یا اینکه برای پیدا کردن نزدیک ترین داده باید منطقه ی بسیار بزرگی از فضای ویژگی ها را بگردیم که این دیگر مفهوم نزدیکی و محلی بودن (locality) که در KNN مورد نظر بود را ندارد.

پس می توان گفت KNN بر اساس locality کار می کند ولی در ابعاد بالا اصلاً locality تعریف درستی ندارد. (تعریف نمی شود) و داده ای هم در نزدیکی داده ی تست وجود ندارد.

پرسی ۵ – Choosing Metrics

این سوال مستقیماً در مورد انتخاب معیار محاسبه‌ی فاصله از مباحث KNN است. در این سوال با معیارهای مختلف آشنا شده و در هر بخش بر حسب کاربرد بهترین را انتخاب می‌کنیم.

برای مرور فاصله‌ی اقلیدسی از [۳]، آشنایی با فاصله‌ی منهتن از [۴]، آشنایی با فاصله‌ی همینگ از [۵] و آشنایی با فاصله‌ی cosine از [۶] استفاده کردم. برای هر کدام از این معیارها تعریفی در نظر گرفته می‌شود که به صورت زیر است:

- فاصله‌ی اقلیدسی: فاصله‌ی کوتاه‌ترین فاصله‌ی مستقیم بین دو نقطه در فضای ویژگی‌ها است. این فاصله تفاوت دو نمونه را با در نظر گرفتن اندازه‌ی واقعی اختلاف‌ها در هر بعد اندازه‌گیری می‌کند.
- فاصله‌ی منهتن: فاصله دو نقطه را به جای کوتاه‌ترین خطی که آن‌ها را بهم متصل می‌کند، به صورت خطوط افقی و عمودی که آن دو نقطه را بهم می‌رساند نشان می‌دهد.
- فاصله‌ی همینگ: در تئوری اطلاعات فاصله‌ی همینگ برای دو رشته با طول مساوی، برابر تعداد مکان‌هایی است که سیمبل‌های متناظر متفاوت هستند. به عبارت دیگر فاصله‌ی همینگ تعداد ویژگی‌هایی را می‌شمارد که بین دو نمونه متفاوت هستند. این معیار فقط روی ویژگی‌های دسته‌ای یا باینری (۰/۱، True/False، .../A/B) قابل استفاده است، نه داده‌های پیوسته.
- فاصله‌ی Cosine: فاصله‌ی کسینوسی بر پایه‌ی زاویه بین دو بردار تعریف می‌شود (ضرب داخلی) و اندازه‌ی واقعی اهمیت ندارد. این معیار وقتی مفید است که الگوی توزیع ویژگی‌ها مهم‌تر از مقدار دقیق آن‌ها باشد؛ بنابراین فقط جهت بردار مهم است.

۵-۱. بخش اول: طبقه‌بندی سه کلاس نوشیدنی

معیار مناسب برای این طبقه‌بندی استفاده از فاصله‌ی اقلیدسی است. در ادامه هر معیار را تحلیل می‌کنم.
گزینه‌ی ۱: ویژگی‌ها در این مسئله به صورت عددی و پیوسته هستند. (مقدار کافئین و قند) این داده‌ها در فضای دو بعدی تعریف خواهند شد.

فاصله‌ی اقلیدسی، فاصله‌ی واقعی نقاط را در فضای ویژگی‌ها اندازه‌گیری می‌کند. برای مثال نوشیدنی که کمی کافئین و مقدار زیادی قند دارد واقعا از نوشیدنی با هر دو مقدار بالا فاصله‌ی بیشتری دارد.

گزینه‌ی ۲: استفاده از فاصله‌ی منهتن زمانی خوب است که داده‌ها ساختار شبکه‌ای دارند. یعنی دقیقاً بتوان تعریفی که از فاصله‌ی منهتن ارائه شد را روی آن پیاده‌سازی کرد.

ولی کافئین و قند چنین حالت شبکه ای را تشکیل نمی دهند. با این وجود، شاید با گسسته کردن مقدار کافئین و قند بتوان حالت شبکه ای را به صورت مصنوعی ایجاد کرد و از فاصله‌ی منتهن بهره برد ولی این کار بهینه نخواهد بود.

گزینه‌ی ۳: در این معیار فقط می‌گوییم که "این ویژگی فرق دارد یا خیر". به وضوح این مورد برای داده‌های Boolean تعریف می‌شود نه برای مسئله‌ی ما که ویژگی‌ها عددی و پیوسته هستند.

ما نیاز داریم بدانیم چقدر کافئین یا قند تفاوت دارد بنابراین همینگ در اینجا کاربردی ندارد.

گزینه‌ی ۴: از این معیار وقتی استفاده می‌شود که فقط نسبت ویژگی‌ها مهم است. همانطور که قبلاً اشاره شد در این مسئله مقدار واقعی کافئین و قند مهم است نه نسبت آنها.

برای مثال اگر نوشیدنی اول ۱۰ واحد کافئین و ۵ واحد قند داشته و نوشیدنی دوم ۲۰۰ واحد کافئین و ۱۰۰ واحد قند داشته باشد، طبق معیار cosine این دو نوشیدنی یکسان هستند زیرا نسبت‌ها یکسان است. پس این معیار برای این مسئله مناسب نیست.

نتیجه نهایی: بهترین معیار برای این مسئله استفاده از فاصله‌ی اقلیدسی (گزینه‌ی ۱) است.

۵-۲. بخش دوم: طبقه بندی چهار کلاس نوشیدنی

معیار مناسب برای این طبقه بندی استفاده از فاصله‌ی اقلیدسی است. در ادامه هر معیار را تحلیل می‌کنم.

گزینه‌ی ۱: توضیحی که خود سوال برای این گزینه ارائه کرده است تحلیل خوبی است. کلاس جدید آب با مقادیر (۰ و ۰) اضافه شده است. این مقدار باعث می‌شود که نسبت‌ها بی‌معنی شوند. به عبارت دیگر فاصله‌ی واقعی در فضای دو بعدی اهمیت بیشتری پیدا می‌کند و اقلیدسی معیاری است که مقدار واقعی ویژگی‌ها را درست در نظر می‌گیرد.

گزینه‌ی ۲: اضافه شدن آب همچنان هیچ کمکی نمی‌کند که کلاس‌های مسئله حالت شبکه ای را ایجاد کنند. همچنان تحلیلی که در بخش اول برای فاصله‌ی منتهن ارائه شد درست است یعنی؛ می‌توان از این معیار استفاده کرد ولی به وضوح معیار بهینه ای نیست.

گزینه‌ی ۳: همچنان با اضافه شدن کلاس آب، ویژگی‌ها عددی و پیوسته هستند و حالت Boolean ندارند. پس همینگ معیار خوبی نیست.

گزینه‌ی ۴: اضافه شدن کلاس آب مشکل این معیار را بزرگ‌تر هم می‌کند. آب یک بردار صفر (نقطه‌ی (۰ و ۰)) است. معیار cosine را با بردار صفر نمی‌توان تعریف کرد. نمی‌توان زاویه‌ی یک بردار را با نقطه‌ی (۰ و ۰) محاسبه کرد. پس این گزینه اصلاً قابل استفاده نیست.

نتیجه نهایی: پس همچنان معیار مناسب استفاده از فاصله‌ی اقلیدسی است. (گزینه ۱)

۵-۳. بخش سوم: طبقه بندی سه کلاس خودرو

معیار مناسب برای این طبقه بندی استفاده از فاصله‌ی همینگ است. در ادامه هر معیار را تحلیل می‌کنم.

گزینه‌ی ۱: داده‌های این مسئله عددی نیستند پس اصلاً استفاده از فاصله‌ی اقلیدسی بی‌معنی است.

گزینه‌ی ۲: این معیار هم باید روی داده‌های عددی اعمال شود که داده‌های ما عددی نیستند. پس قابل استفاده نیست.

گزینه‌ی ۳: داده‌های لوگو، نوع خودرو و مدل خودرو همگی حالت دسته‌ای (category) دارند و عددی و پیوسته نیستند. فاصله‌ی همینگ دقیقاً برای همین نوع داده‌هاست. هر بار که یک ویژگی بین دو نمونه متفاوت باشد، فاصله ۱ واحد اضافه می‌شود.

گزینه‌ی ۴: فاصله‌ی cosine نیازمند بردارهای عددی با اندازه و جهت است. لوگو و مدل و نوع خودرو مقادیر عددی نیستند. پس قابل استفاده نیست.

نتیجه نهایی: پس معیار همینگ بهترین و تنها گزینه مورد استفاده در اینجا است. (گزینه ۳)

۵-۴. بخش چهارم: استفاده از درخت تصمیم

گزینه سه جواب صحیح است. برای آشنایی با ID tree از [۷] استفاده کردم. در ادامه هر کدام از گزینه‌ها را بررسی می‌کنم.

گزینه ۱: توضیحات این گزینه کاملاً غلط است. روش KNN اصلاً نیازی به رسم داده‌ها ندارد. این روش صرفاً نیازمند یک معیار فاصله سنجی است. حتی وقتی داده‌ها غیر عددی هستند معیارهایی مثل فاصله‌ی همینگ که در بخش قبل دیدیم قابل استفاده هستند.

گزینه ۲: این دلیل هم درست نیست. اگر چند نمونه دقیقاً یکسان باشند همه‌ی آنها در یک فاصله قرار می‌گیرند. KNN بر اساس داده‌های همسایه کار می‌کند پس وقتی داده‌های تکراری در همسایگی یک نقطه داشته باشیم، احتمال اینکه کلاس غالب درست باشد افزایش می‌یابد. پس اگر تعداد آنها زیاد باشد باعث افزایش دقت می‌شود. پس وجود نمونه‌های تکراری اتفاقاً به نفع ماست.

گزینه ۳: همان‌طور که از متن سوال هم مشخص است، سه کلاس با استفاده از لوگو به تنهایی قابل طبقه بندی هستند. KNN تمامی ویژگی‌ها را به یک اندازه مهم در نظر می‌گیرد. ویژگی‌های کم‌اهمیت‌تر مانند مدل و نوع خودرو می‌تواند فاصله‌ها را خراب کند. در منبع [۷] می‌بینیم که درخت تصمیم چند مرحله دارد

که دومین مرحله مربوط به انتخاب ویژگی است که با استفاده از Attribute Selection Measure (ASM) بهترین ویژگی در مجموعه داده انتخاب می شود. همین مرحله باعث می شود ویژگی های نامرتبط کنار گذاشته شوند و اولین split را روی لوگو در گره ریشه اعمال می کنیم.

بنابراین درخت می تواند خیلی راحت تر و دقیق تر تاکسی ها را فقط بر اساس لوگو جدا کند.

گزینه ۴: ویژگی ها ثابت هستند (لوگو، مدل و نوع خودرو) و به حرکت تاکسی ها وابسته نیستند. این گزینه به وضوح غلط است.

نتیجه نهایی: پس بر این اساس می بینیم که ID tree با توجه به اینکه ویژگی های نامرتبط را کنار می گذارد، پس می تواند از KNN بهتر عمل کند. (گزینه ۳)

پیش ۶ – Kernel Ridge Regression

۶-۱. مینیمم تابع هزینه رگرسیون ریج

برای نشان دادن بردارها از حروف کوچک و bold استفاده می کنیم. تابع هزینه ی رگرسیون ریج به صورت زیر است.

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 + \frac{\lambda}{2} \|\theta\|^2 \quad (40)$$

نمایش برداری آن به صورت زیر است.

$$J(\theta) = \frac{1}{2} \|x\theta - y\|^2 + \frac{\lambda}{2} \theta^T \theta \quad (41)$$

ابتدا تابع هزینه را باز می کنیم و ریز تر می نویسیم.

$$J(\theta) = \frac{1}{2} \theta^T x^T x \theta - y^T x \theta + \frac{1}{2} y^T y + \frac{\lambda}{2} \theta^T \theta \quad (42)$$

حالا از این رابطه نسبت به θ مشتق می گیریم و مساوی صفر قرار می دهیم.

$$\nabla_{\theta} J(\theta) = x^T x \theta - x^T y + \lambda \theta = 0 \quad (43)$$

و

$$(x^T x + \lambda I) \theta = x^T y \quad (44)$$

پس با استفاده از θ زیر تابع هزینه مینیمم می شود.

$$\theta = (x^T x + \lambda I)^{-1} x^T y \quad (45)$$

۶-۲. استفاده از ترفند کرنل در رگرسیون ریج

با استفاده از یک تبدیل برای ویژگی ها، تابع هزینه به صورت زیر تبدیل شد.

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (\theta^T \phi(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2} \|\theta\|^2 \quad (46)$$

در این رابطه ماتریس ویژگی ها به صورت زیر است.

$$\Phi = \begin{bmatrix} \phi(x^{(1)})^T \\ \phi(x^{(2)})^T \\ \vdots \\ \phi(x^{(n)})^T \end{bmatrix} \quad (47)$$

پس تابع هزینه را به این صورت می توان بازنویسی کرد.

$$J(\theta) = \frac{1}{2} \|\Phi\theta - y\|^2 + \frac{\lambda}{2} \theta^T \theta \quad (48)$$

متناظرا از بخش اول می دانیم که با مساوی صفر قرار دادن گرادیان تابع هزینه به رابطه ی زیر خواهیم رسید.

$$\theta = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y \quad (49)$$

حال از رابطه ی معرفی شده در صورت سوال استفاده می کنیم. با فرض $A = \Phi$ و $B = \Phi^T$

$$(\Phi^T \Phi + \lambda I)^{-1} \Phi^T = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} \quad (50)$$

پس می توانیم (49) را به صورت زیر بازنویسی کنیم.

$$\theta = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y \quad (51)$$

از طرفی سوال گفته بردار θ را می توان به صورت ترکیب خطی از بردار ویژگی ها در نظر گرفت.

$$\theta = \sum_{i=1}^m \alpha_i \phi(x^{(i)}) = \Phi^T \alpha \quad (52)$$

با مقایسه ی (51) و (52) می توان دریافت:

$$\alpha = (\Phi \Phi^T + \lambda I)^{-1} y \quad (53)$$

از تعریف ماتریس کرنل به صورت $K = \Phi \Phi^T$ تعریف شده و هر درایه آن به صورت زیر است.

$$K_{ij} = \phi(x^{(i)})^T \phi(x^{(j)}) = k(x^{(i)}, x^{(j)}) \quad (54)$$

پس (53) را می توان به صورت زیر بازنویسی کرد.

$$\alpha = (K + \lambda I)^{-1} y \quad (55)$$

پس پیش بینی به ازای داده ی جدید x_{new} به صورت زیر است.

$$\hat{y}(x_{new}) = \theta^T \phi(x_{new}) = (\Phi^T \alpha)^T \phi(x_{new}) = \alpha^T \Phi \phi(x_{new}) \quad (56)$$

که در این رابطه:

$$\Phi\phi(x_{new}) = \begin{bmatrix} \phi(x^{(1)})^T \phi(x_{new}) \\ \vdots \\ \phi(x^{(m)})^T \phi(x_{new}) \end{bmatrix} = \begin{bmatrix} k(x^{(1)}, x_{new}) \\ \vdots \\ k(x^{(m)}, x_{new}) \end{bmatrix} \quad (57)$$

پس با توجه به (57) می توان رابطه ی پیش بینی را به صورت زیر هم نوشت.

$$\hat{y}(x_{new}) = \sum_{i=1}^m \alpha_i k(x^{(i)}, x_{new}) \quad (58)$$

پس به این ترتیب می توانیم در یک فضای ویژگی بزرگ (با ابعاد بالا) پیش بینی انجام دهیم بدون اینکه اصلا

نیازی به محاسبه ی $\phi(x)$ داشته باشیم.

پرسی ۷ – Deriving Linear Regression

در این مسئله می خواهیم فرمول بهینه w^* که خطای مربعات متوسط را کمینه می کند پیدا کنیم و نشان دهیم که این فرمول به شکل نسبتی بین کوواریانس و واریانس است.

۷-۱. محاسبه $\Sigma_{yx} = E_{x,y}[yx]$

مدل ما در این مسئله به صورت زیر است.

$$y = w^T x + \epsilon \quad (59)$$

امیدریاضی مورد نظر به صورت زیر است. توجه کنید که x مستقل از ϵ بوده و w اسکالر است.

$$\begin{aligned} \Sigma_{yx} = E_{x,y}[yx] &= E_{x,y}[(w^T x + \epsilon)x] = E[w^T x^2] + E[\epsilon x] \\ &= w^T E[x^2] + E[\epsilon]E[x] \end{aligned} \quad (60)$$

گفتیم w یک اسکالر است پس ترانپاده آن با خودش برابر است. با توجه به توزیع x و ϵ خواهیم داشت.

$$\Sigma_{yx} = w\Sigma_x + 0 = w\Sigma_x \quad (61)$$

۷-۲. محاسبه $\Sigma_{yx} = E_{x,y}[yx]$ با استفاده از نمونه ها

با فرض داده های $\{(x_n, y_n)\}_{n=1}^N$ مقدار امیدریاضی به صورت زیر است.

$$E_{x,y}[yx] = \frac{1}{N} \sum_{i=1}^N x_n y_n \quad (62)$$

۷-۳. محاسبه w^* بهینه

می خواهیم w را به گونه ای پیدا کنیم که امیدریاضی مربع خطا (residual loss) مینیمم شود. (توجه شود تمامی متغیر ها اسکالر است.)

$$A = E[(y - wx)^2] = E[y^2] - 2wE[yx] + w^2 E[x^2] \quad (63)$$

حال از A نسبت به w مشتق گرفته و مساوی صفر قرار می دهیم.

$$\frac{dA}{dw} = -2E[yx] + 2wE[x^2] = 0 \quad (64)$$

پس مقدار بهینه w به صورت زیر است.

$$w^* = \frac{E[yx]}{E[x^2]} \quad (65)$$

می دانیم کواریانس هر دو متغیر تصادفی X و Y به صورت زیر تعریف می شود.

$$\begin{aligned} cov(X, Y) &= \Sigma_{yx} = E[(X - E[X])(Y - E[Y])] \\ &= E[XY - XE[Y] - E[X]Y + E[X]E[Y]] \\ &= E[XY] - E[Y]E[X] - E[X]E[Y] + E[X]E[Y] \\ &= E[XY] - E[X]E[Y] \end{aligned} \quad (66)$$

پس خواهیم داشت:

$$E[XY] = \Sigma_{yx} + E[X]E[Y] = \Sigma_{yx} + \mu_x \mu_y \quad (67)$$

همچنین طبق تعریف کوواریانس برای $E[x^2]$ به طور مشابه خواهیم داشت:

$$E[x^2] = \Sigma_x + \mu_x^2 \quad (68)$$

حال با توجه به (67) و (68) می توانیم (65) را به صورت زیر بازنویسی کنیم. توجه داشته باشید که میانگین ها صفر هستند.

$$w^* = \frac{E[yx]}{E[x^2]} = \frac{\Sigma_{yx} + \mu_x \mu_y}{\Sigma_x + \mu_x^2} = \frac{\Sigma_{yx}}{\Sigma_x} \quad (69)$$

این بخش با رابطه‌ی (69) پایان می یابد. اما برای اطمینان از محدب بودن (63) مشتق دوم را هم حساب می کنیم تا مطمئن باشیم نقطه‌ی بدست آمده مینیمم تابع است.

$$\frac{d^2 A}{dw^2} = 2E[x^2] \quad (70)$$

از آنجا که $E[x^2] > 0$ پس مشتق دوم مثبت بوده و مقدار بدست آمده مینیمم است.

۴-۷. محاسبه‌ی w^* بهینه با استفاده از توزیع جدید

در این بخش باید با فرض توزیع $N(\mu_x, \Sigma_x)$ برای داده های x مقدار بهینه ی w^* را پیدا کنیم. این مورد را در بخش قبل در روابط (67)، (68) و (69) دیدیم. پس با توجه به آن ها مقدار بهینه به صورت زیر خواهد بود.

$$w^* = \frac{E[yx]}{E[x^2]} = \frac{\Sigma_{yx} + \mu_x \mu_y}{\Sigma_x + \mu_x^2} \quad (71)$$

۵-۷. اهمیت توزیع داده های x

در تمامی مراحل بخش ۳ و ۴ از گوسی بودن داده ها هیچ استفاده ای نکردیم. در تمامی روابط (63) الی (69) نهایتاً به ممان های اول و دوم x و y نیاز داریم. یعنی صرفاً به شرط داشتن مقادیر میانگین و کوواریانس، می توانیم به رابطه ی (71) برسیم. در واقع اصلاً هیچ فرضی در مورد ساختار و توزیع داده ها نیاز نداریم. تنها چیزی

که باید برقرار باشد (علاوه بر داشتن ممای اول و دوم) این است که $E[yx]$ و $E[x^2]$ مقادیر محدود و $E[x^2]$ مقدار غیر صفر داشته باشد تا تقسیم در (۷۱) معتبر باشد.

پرسی ۸ – Logistic Regression with Newton's Method

در این سوال تابع هزینه‌ی لجیستیک را در حالت ماتریسی بررسی می‌کنیم. قدم به قدم روابط را محاسبه می‌کنیم تا به فرمول آپدیت روش نیوتن برسیم.

توجه کنید همان طور که سوال هم اشاره کرده از حروف انگلیسی **bold** برای نشان دادن بردار استفاده کردم.

۸-۱. محاسبه‌ی $\nabla_w J(w)$

تابع هزینه به فرم برداری به صورت زیر نوشته می‌شود.

$$J(w) = -\mathbf{y} \cdot \ln \mathbf{s} - (\mathbf{1} - \mathbf{y}) \cdot \ln(\mathbf{1} - \mathbf{s}) \quad (72)$$

در این رابطه به شکل مستقیم از \mathbf{w} استفاده نشده بلکه \mathbf{w} درون \mathbf{s} تعریف شده است. به همین دلیل برای محاسبه‌ی گرادیان نسبت به \mathbf{w} باید از قاعده‌ی زنجیره‌ای استفاده کنیم.

بردار \mathbf{s} و هر المان آن به صورت زیر تعریف می‌شود.

$$\mathbf{s} = \begin{bmatrix} s(x_1^T \mathbf{w}) \\ s(x_2^T \mathbf{w}) \\ \vdots \\ s(x_n^T \mathbf{w}) \end{bmatrix}, \quad s(y) \triangleq \frac{1}{1 + e^{-y}}$$

ابتدا مشتق هر المان از بردار \mathbf{s} را مشخص می‌کنیم.

$$\frac{d}{dy} s(y) = \frac{e^{-y}}{(1 + e^{-y})^2} \quad (73)$$

می‌توانیم با استفاده از خود $s(y)$ مشتق آن را به صورت زیر نشان دهیم.

$$\frac{d}{dy} s(y) = s(y)(1 - s(y)) \quad (74)$$

پس مشتق نسبت به \mathbf{w} به صورت زیر تعریف می‌شود.

$$\left[\frac{\partial s_i}{\partial \mathbf{w}} \right]^T = s_i(1 - s_i) \mathbf{x}_i^T \quad (75)$$

و گرادیان به صورت زیر خواهد بود.

$$\nabla_w \mathbf{s} = \begin{bmatrix} s_1(1 - s_1) \mathbf{x}_1^T \\ s_2(1 - s_2) \mathbf{x}_2^T \\ \vdots \\ s_n(1 - s_n) \mathbf{x}_n^T \end{bmatrix} = \text{diag}(\mathbf{s} \odot (\mathbf{1} - \mathbf{s})) \mathbf{X} \quad (76)$$

منظور از \odot ضرب المان به المان (ضرب Hadamard) است.

حالا باید از $J(w)$ نسبت به s گرادیان بگیریم. برای هر المان مشتق به صورت زیر است.

$$\frac{\partial J}{\partial s_i} = -\frac{y_i}{s_i} + \frac{1 - y_i}{1 - s_i} \quad (77)$$

حالا همین مشتق های المانی را در قالب بردار می نویسیم و بردار بدست آمده را باز می کنیم تا بتوانیم بر حسب $\text{diag}()$ بازنویسی کنیم. به صورت زیر:

$$\begin{aligned} \nabla_s J &= \begin{bmatrix} \frac{\partial J}{\partial s_1} \\ \frac{\partial J}{\partial s_2} \\ \vdots \\ \frac{\partial J}{\partial s_n} \end{bmatrix} = \begin{bmatrix} -\frac{y_1}{s_1} + \frac{1 - y_1}{1 - s_1} \\ -\frac{y_2}{s_2} + \frac{1 - y_2}{1 - s_2} \\ \vdots \\ -\frac{y_n}{s_n} + \frac{1 - y_n}{1 - s_n} \end{bmatrix} = \begin{bmatrix} -\frac{y_1}{s_1} \\ -\frac{y_2}{s_2} \\ \vdots \\ -\frac{y_n}{s_n} \end{bmatrix} + \begin{bmatrix} \frac{1 - y_1}{1 - s_1} \\ \frac{1 - y_2}{1 - s_2} \\ \vdots \\ \frac{1 - y_n}{1 - s_n} \end{bmatrix} \quad (78) \\ &= -\text{diag}\left(\frac{1}{s}\right) y + \text{diag}\left(\frac{1}{1 - s}\right) (1 - y) \\ &= \frac{s - y}{s \odot (1 - s)} \end{aligned}$$

حال با استفاده از قاعده ی زنجیره ای و روابط (76) و (78) می توانیم گرادیان خواسته شده را محاسبه کنیم.

$$\begin{aligned} \nabla_w J(w) &= (\nabla_w s)^T (\nabla_s J) \quad (79) \\ &= X^T \text{diag}(s \odot (1 - s)) \frac{s - y}{s \odot (1 - s)} \\ &= X^T (s - y) \end{aligned}$$

به این شکل گرادیان تابع هزینه نسبت به w بدست آمد.

۸-۲. محاسبه ی ماتریس هسین

مقدار $\nabla_w J$ و $\nabla_w s$ را در بخش ۸-۱ حساب کردیم، با استفاده از آن می توانیم ماتریس هسین را هم بدست آوریم.

$$\nabla_w^2 J(w) = \frac{d}{dw} [X^T (s - y)] = X^T \frac{ds}{dw} = X^T \text{diag}(s \odot (1 - s)) X \quad (80)$$

۸-۳. رابطه ی آپدیت روش نیوتن

رابطه ی آپدیت روش نیوتن برای مینیمم کردن $J(w)$ به صورت زیر است.

$$w_{new} = w_{old} - [\nabla_w^2 J(w_{old})]^{-1} \nabla_w J(w_{old}) \quad (81)$$

گرادیان و هسین $J(w)$ در بخش های قبلی بدست آمد. صرفا آن ها را جاگذاری می کنیم.

$$w_{new} = w_{old} - [X^T \text{diag}(s \odot (1 - s)) X]^{-1} X^T (s - y) \quad (82)$$

۸-۴. بررسی مثال عددی

داده ها به صورت زیر تعریف می شوند.

$$X = \begin{bmatrix} 0.2 & 3.1 & 1 \\ 1.0 & 3.0 & 1 \\ -0.2 & 1.2 & 1 \\ 1.0 & 1.1 & 1 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad w^{(0)} = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \quad (83)$$

با استفاده از این مقادیر و نوشتن کد (کد نباید ارائه شود) optimization را تاجایی ادامه می دهیم که تغییرات w کمتر از 10^{-5} باشد. در این حالت به مقادیر زیر می رسیم.

$$w = \begin{bmatrix} 5.2146 \\ 7.8421 \\ -18.6193 \end{bmatrix} \quad (84)$$

پرسی ۹ – A Bayesian Interpretation of Lasso

این سوال در کل در مورد تخمین MAP می باشد.

۹-۱. بسط احتمال پسین با استفاده از رابطه‌ی بیز

با استفاده از قانون بیز می توانیم احتمال پسین را به صورت زیر بازنویسی کنیم. (در تساوی دوم از رابطه‌ی chain rule استفاده شده است).

$$f(w|(x_i, y_i)) = \frac{f((x_i, y_i)|w)f(w)}{f(x_i, y_i)} = \frac{f(y_i|x_i, w)f(x_i|w)f(w)}{f(x_i, y_i)} \quad (۸۵)$$

مخرج در این رابطه مستقل از w است پس کنار گذاشته می شود. از طرفی با فرض مستقل بودن x_i از w ها می توان نوشت $f(x_i|w) = f(x_i)$ و از طرفی این احتمال هم مستقل از w است پس کنار گذاشته می شود. پس بر این اساس این رابطه به صورت زیر تبدیل می شود.

$$f(w|(x_i, y_i)) \propto f(y_i|x_i, w)f(w) \quad (۸۶)$$

حال باید هر احتمال را محاسبه کنیم. با *likelihood* شروع می کنیم.

$$f(y_i|x_i, w) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y_i - w \cdot x_i)^2}{2\sigma^2}\right) \quad (۸۹)$$

با توجه به استقلال داده ها، برای احتمال کل داده ها می توان نوشت:

$$\begin{aligned} f(y|x, w) &= \prod_{i=1}^n f(y_i|x_i, w) \\ &= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w \cdot x_i)^2\right) \end{aligned} \quad (۹۰)$$

با توجه به استقلال تمامی المان های w از هم، احتمال پیشین هم به صورت زیر محاسبه می شود.

$$f(w_i) = \frac{1}{\sqrt{b}} \exp\left(-\frac{|w_i|}{b}\right) \quad (۹۱)$$

$$f(w) = \prod_{i=1}^d f(w_i) = \left(\frac{1}{\sqrt{b}}\right)^d \exp\left(-\frac{1}{b} \sum_{i=1}^d |w_i|\right) \quad (۹۲)$$

حال با استفاده از (۹۰) و (۹۲) می توانیم (۸۶) را به صورت زیر بازنویسی کنیم.

$$f(w|(x_i, y_i)) \propto \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^r \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w \cdot x_i)^2\right) \left(\frac{1}{rb}\right)^d \exp\left(-\frac{1}{b} \sum_{i=1}^d |w_i|\right) \quad (93)$$

برای اینکه رابطه در یک خطر قرار گیرد مجبور شدم سائز فونت را کوچک تر کنم.

خواسته‌ی سوال با رابطه‌ی (93) برآورده شد. اما همچنان در این رابطه ضرایبی وجود دارند که صرفاً نقش نرمالیزاسیون ایفا می کنند و می توان آنها را هم کنار گذاشته و به رابطه‌ی زیر رسید.

$$f(w|(x_i, y_i)) \propto \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w \cdot x_i)^2 - \frac{1}{b} \sum_{i=1}^d |w_i|\right) \quad (94)$$

۹-۲. بررسی تخمین MAP

طبق صورت سوال log-likelihood را به صورت زیر تعریف می کنیم.

$$\ell(w) \triangleq \ln f(w|(x_i, y_i)) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w \cdot x_i)^2 - \frac{1}{b} \sum_{i=1}^d |w_i| \quad (95)$$

تخمینگر MAP با استفاده از log-likelihood به صورت زیر محاسبه می شود.

$$\begin{aligned} \hat{w}_{MAP} &= \arg \max_w \ell(w) \\ &= \arg \max_w \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w \cdot x_i)^2 - \frac{1}{b} \sum_{i=1}^d |w_i| \right) \end{aligned} \quad (96)$$

ماکزیمم کردن یک تابع، برابر با مینیمم کردن قرینه‌ی آن تابع است. اگر این تابع را در ۱- ضرب کنیم تابع قرینه می شود و حالا در این تابع باید به جای ماکزیمم به دنبال مینیمم باشیم.

$$\hat{w}_{MAP} = \arg \min_w \left(\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w \cdot x_i)^2 + \frac{1}{b} \sum_{i=1}^d |w_i| \right) \quad (97)$$

ضرب کردن یک عدد مثبت ثابت در کل یک تابع هیچ تاثیری در نقطه‌ی ماکزیمم یا مینیمم نمی گذارد. پس می توانیم در $2\sigma^2$ ضرب کنیم و مطمئن باشیم که تخمینگر دقیقاً مساوی قبل عمل می کند.

$$\hat{w}_{MAP} = \arg \min_w \left(\sum_{i=1}^n (y_i - w \cdot x_i)^2 + \frac{2\sigma^2}{b} \sum_{i=1}^d |w_i| \right) \quad (98)$$

توجه کنید که $\sum_{i=1}^d |w_i| = \|w\|_1$

به همین سادگی نشان دادیم مینیمم کردن (۹۸) معادل ماکزیمم کردن (۹۶) است. پارامتر λ هم به صورت زیر تعریف می شود.

$$\lambda = \frac{2\sigma^2}{b} \quad (۹۹)$$

پرسش ۱۰ – Parzen Window Classification with Gaussian Kernel

در این سوال یک مسئله‌ی طبقه‌بندی داریم که در آن توزیع شرطی داده‌ها به شرط کلاس‌ها را با استفاده از روش پارزن محاسبه می‌کنیم. تابع کرنل در اینجا به صورت زیر است.

$$\varphi\left(\frac{x-x_i}{h}\right) \propto \exp\left[-\frac{(x-x_i)^T(x-x_i)}{2h^2}\right] \quad (100)$$

این توزیع یک ضریب هم پشت exp در رابطه‌ی (۱۰۰) دارد ولی چون صرفاً یک ضریب ثابت است که در تمامی کلاس‌ها تکرار می‌شود، می‌توان صرف نظر کرد.

هدف این سوال بررسی اثر مقدار پارامتر smoothing روی توزیع تخمین زده شده و طبقه‌بند است.

۱-۱۰. بخش a: طبقه‌بندی به ازای $h = 1$

مطابق آنچه در مورد تخمین توزیع با استفاده از روش پارزن در سوال ۱ دیدیم، توزیع شرطی به صورت زیر خواهد بود.

$$\begin{aligned} p_n(x|\omega_i) &= \frac{1}{n_i} \sum_{j=1}^{n_i} \frac{1}{h^d} \varphi\left(\frac{x-x_i}{h}\right) \\ &= \frac{1}{n_i h^d} A \sum_{j=1}^{n_i} \exp\left(-\frac{(x-x_i)^T(x-x_i)}{2h^2}\right) \end{aligned} \quad (101)$$

در رابطه‌ی (۱۰۱) منظور از A ضریب نرمالیزاسیونی است که در (۱۰۰) با ذکر دلیل صرف نظر شده بود (حذف این ضریب‌ها ممکن است باعث شود مقادیر بدست آمده برای احتمال بیشتر از ۱ شود. از آنجا که ماکزیمم این احتمال‌ها مورد نظر ماست پس مشکلی نخواهد بود.) و d ابعاد مسئله است که در اینجا برابر با ۳ می‌باشد. با استفاده از قانون بیز و رابطه‌ی (۱۰۱) احتمال پسین را محاسبه کرده و کلاس متناظر با بیشترین مقدار را به عنوان کلاس برای داده‌ی تست، انتخاب می‌کنیم.

$$p(\omega_i|x) = \frac{p_n(x|\omega_i)p(\omega_i)}{\sum_j p_n(x|\omega_j)p(\omega_j)} \quad (102)$$

صورت مسئله چیزی درباره‌ی احتمال‌های پیشین نگفته است. از طرفی اگر بخواهیم با استفاده از نمونه‌های هر کلاس برآوردی از احتمال پیشین داشته باشیم به علت کم بودن تعداد نمونه‌ها این برآورد دقیق نخواهد بود. پس فرض می‌شود احتمال‌های پیشین برابر هستند.

از طرفی مخرج در (۱۰۲) صرفاً ضریبی است که برای همه‌ی کلاس‌ها تکرار می‌شود پس قابل صرف نظر است. پس تصمیم‌گیری به صورت زیر خواهد بود.

$$\omega_i = \arg \max_i p(\omega_i|x) = \arg \max_i p_n(x|\omega_i) \quad (10.3)$$

خروجی کد به ازای $h = 1$ به صورت زیر است.

Test point ۱: [۰,۵ ۱. ۰.]

ω_1 : density = ۱,۲۵۹۱۹۰e-۰۱

ω_2 : density = ۴,۷۱۰۷۸۵e-۰۱

ω_3 : density = ۳,۷۷۴۶۵۶e-۰۱

Predicted class: ω_2

Test point ۲: [۰,۳۱ ۱,۵۱ -۰,۵]

ω_1 : density = ۱,۵۳۴۳۸۶e-۰۱

ω_2 : density = ۴,۸۲۷۹۶۳e-۰۱

ω_3 : density = ۲,۰۵۰۴۰۳e-۰۱

Predicted class: ω_2

Test point ۳: [-۰,۳ ۰,۴۴ -۰,۱]

ω_1 : density = ۱,۳۹۹۰۸۱e-۰۱

ω_2 : density = ۳,۷۸۲۵۴۳e-۰۱

ω_3 : density = ۱,۷۷۰۶۶۴e-۰۱

Predicted class: ω_2

می بینیم که طبقه بند برای هر سه داده‌ی تست کلاس ۲ را مشخص کرده است. به ازای هر داده‌ی تست مقدار احتمال شرطی به شرط هر کلاس حساب شده و کلاس متناظر با بیشترین مقدار انتخاب شده است. در اینجا بیش از این توضیح نمی دهیم زیرا تحلیل و نتیجه گیری بهتر در انتهای بخش بعدی خواهد بود.

۱۰-۲. بخش b: طبقه‌بندی به ازای $h = 1$

خروجی کد به ازای $h = 0.1$ به صورت زیر است.

Test point ۱: [۰,۵ ۱. ۰.]

ω_1 : density = ۸,۷۷۴۶۹۵e-۲۰

ω_2 : density = ۶,۷۶۹۰۲۰e-۰۵

ω_3 : density = ۸,۰۰۰۱۸۹e-۱۷

Predicted class: ω_2

Test point ۲: [۰,۳۱ ۱,۵۱ -۰,۵]

ω_1 : density = $۲,۸۷۱۱۱۶e-۲۰$

ω_2 : density = $۱,۲۰۰۱۸۴e-۰۵$

ω_3 : density = $۲,۱۵۹۳۲۹e-۲۵$

Predicted class: ω_2

Test point ۳: [-۰,۳ ۰,۴۴ -۰,۱]

ω_1 : density = $۳,۶۳۷۱۲۳e-۱۲$

ω_2 : density = $۳,۷۸۳۰۹۴e-۰۴$

ω_3 : density = $۱,۰۶۴۰۲۶e-۳۹$

Predicted class: ω_2

باز هم می بینیم که طبقه بند در هر سه نقطه‌ی تست کلاس دوم را انتخاب کرده است.

در بخش a که $h = ۱$ است توزیع های گوسی تخمین زده شده پهن تر هستند. یعنی توزیع ها نرم (smooth) و کم نوسان هستند. در این حالت چون دامنه ی گوسی ها بزرگ تر است، مقدار احتمال برای داده های تست از مرتبه ی ۰.۱ هستند.

اما در بخش b که $h = ۰.۱$ است توزیع های گوسی تخمین زده شده تیز تر و باریک تر می شوند. به این معنی که مقدار توزیع در دامنه ها بسیار کمتر خواهد بود. در واقع مقدار چگالی در اکثر نقاط به طور واضحی کوچک می شود. همین باعث می شود که مقدار احتمال برای داده های تست خیلی کمتر باشد.

این رفتار دو پیامد دارد. با کاهش h به نظر می رسد که مدل کلاس غالب را با اطمینان بیشتری تشخیص می دهد اما این یک مزیت نیست (پیامد اول). این تمایز بیش از حد ناشی از حساسیت شدید مدل به نمونه های آموزشی است و عملاً منجر به overfitting کاهش توانایی generalization می شود (پیامد دوم).

پرسش ۱۱ – Maximum likelihood estimation

۱-۱۱. بخش a: محاسبه‌ی عددی MLE و مقایسه با حل فرم بسته

برای داده های نرمال likelihood به صورت زیر است.

$$L(\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \quad (104)$$

لگاریتم آن به صورت زیر است.

$$\ell(\mu, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \quad (105)$$

با مشتق گرفتن نسبت به μ و σ^2 و مساوی صفر قرار دادن می توان تخمین MLE آن ها را به دست آورد.

$$\frac{\partial \ell}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \hat{\mu}_{MLE} = \bar{X} = \frac{1}{n} \sum_{i=1}^n x_i \quad (106)$$

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (x_i - \mu)^2 = 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2 \quad (107)$$

برای پیاده سازی این بخش از [۸] و [۹] استفاده کردم. برای مینیمم کردن تابع log likelihood باید از یک نقطه شروع می کردیم من از (۰ و ۰) استفاده کردم. مقادیر عددی و تئوری به صورت زیر هستند.

Numerical MLE:

mu_hat numerical=۱۰۰,۴۸۰۵۱۳۶۸۳۱۴۴۳۵ , sigma²_hat numerical=۲۱۶,۳۱۸۶۷۰۳۶۵۹۱۲۴۵

Closed-form MLE:

mu_hat closed form=۱۰۰,۴۸۰۴۹۹۹۹۹۹۹۹۹۹ , sigma²_hat closed form=۲۱۶,۳۱۸۴۶۹۷۵

می بینیم که مقادیر حل عددی با تئوری تا ۴ رقم اعشار کاملاً مشابه هستند.

۱۱-۲. بخش b: بررسی مینیمم شدن منفی log likelihood

این مورد را هم از دید تئوری و هم عددی می توان بررسی کرد. ابتدا با بررسی تئوری شروع می کنم. برای بررسی مینیمم بودن منفی log likelihood در نقاط بدست آمده در تحلیل تئوری می توان مشتق دوم این عبارت

را نسبت به هر کدام از متغیرها محاسبه کرد. اگر این مشتق دومها مثبت باشند یعنی نقاط بدست آمده همان مینیممها هستند.

$$\frac{\partial^2(-\ell)}{\partial \mu^2} = \frac{n}{\sigma^2} > 0. \quad (108)$$

پارامترهای n و σ^2 به ذاته مثبت هستند. پس مشتق دوم مثبت است و نقطه‌ی بدست آمده مینیمم است.

$$\frac{\partial^2(-\ell)}{\partial (\sigma^2)^2} = \frac{n}{2(\sigma^2)^2} + \frac{1}{(\sigma^2)^3} \sum_{i=1}^n (x_i - \mu)^2 > 0. \quad (109)$$

پارامترهای این مشتق دوم هم یا توان دو دارند یا به ذاته مثبت هستند پس نقطه بدست آمده مینیمم است.

برای اینکه به صورت عددی نشان دهیم، مقادیر بدست آمده مینیمم هستند ابتدا مقدار تابع negative log likelihood را در نقطه‌ی بدست آمده از بخش قبل حساب می‌کنیم و سپس مقدار تابع را به ازای افزایش یا کاهش به اندازه‌ی خیلی کم (epsilon) نقطه حساب می‌کنیم. اگر نقاط بدست آمد در بخش قبل واقعا مینیمم تابع باشند باید مقدار تابع در نقاط کمتر از بقیه‌ی نقاط باشد.

NLL at MLE: ۸۲۱,۴۶۲۸۷۸۵۲۲۸۱۹۹

NLL mu + eps: ۸۲۱,۴۶۲۹۲۴۷۵۰۹۵۷۵

NLL mu - eps: ۸۲۱,۴۶۲۹۲۴۷۵۰۹۵۷۵

NLL sigma² + eps: ۸۲۱,۴۶۳۱۴۴۸۳۱۷۸۲۶

NLL sigma² - eps: ۸۲۱,۴۶۳۱۴۶۴۷۸۳۱۶۷

خوب می‌بینیم که به ازای تغییر اندک، مقدار تابع بیشتر از مقدار تابع در نقطه‌ی مینیمم بدست آمده شده. پس به این شکل نشان دادیم نقطه‌ی بدست آمده مینیمم تابع است.

۱۱-۳. بخش c: محاسبه‌ی ماتریس کوواریانس تخمینگر در حالت حدی

در واقع این بخش به زبان ساده می‌گوید اگر ما بخواهیم آزمایش را بارها و بارها تکرار کنیم، تخمین‌های $\hat{\mu}$ و $\hat{\sigma}^2$ چگونه تغییر می‌کنند.

برای محاسبه‌ی این بخش نیاز به استفاده از ماتریس اطلاعات فشر داریم. من با ماتریس اطلاعات فشر برای محاسبه یک باند برای تخمینگر در درس تئوری آشکارسازی و تخمین آشنا شدم. در مورد ماتریس اطلاعات فشر می‌توانید در [۱۰] بخوانید.

برای تعداد زیادی نمونه، توزیع به صورت زیر خواهد شد.

$$\hat{\theta} = (\hat{\mu}, \hat{\sigma}^2) \sim N(\theta, I(\theta)^{-1}) \quad (110)$$

که در این رابطه، $I(\theta)$ همان ماتریس اطلاعات فیشر است. برای یک نمونه $x_i \sim N(\mu, \sigma^2)$ ماتریس اطلاعات فیشر به صورت زیر است.

$$I_1(\mu, \sigma^2) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix} \quad (111)$$

برای n تا نمونه خواهیم داشت.

$$I_n(\mu, \sigma^2) = n \cdot I_1(\mu, \sigma^2) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix} \quad (112)$$

بر این اساس ماتریس کوواریانس MLE برابر معکوس ماتریس اطلاعات فیشر بوده و به صورت زیر بدست می آید.

$$\text{asymtotic covriance matrix} = \Sigma = I_n^{-1} = \begin{pmatrix} \frac{\sigma^2}{n} & 0 \\ 0 & \frac{2\sigma^4}{n} \end{pmatrix} \quad (113)$$

با استفاده از کد برای این دیتاست به صورت زیر بدست می آید.

```
array([[ ۱,۰۸۱۵۹۲۳۵, ۰. ],
       [ ۰. , ۴۶۷,۹۳۶۸۰۳۵۵]])
```

۱۱-۴. بخش d: بهبود ماتریس کوواریانس تخمینگر در حالت حدی

در بخش قبل ماتریس کوواریانس را برای تخمینگر MLE حساب کردیم ولی یک مشکل وجود داشت و آن هم این است که σ^2 بایاس دارد یعنی می توان بعد استفاده از یک ضریب، بایاس آن را از بین برد و این ماتریس کوواریانس را بهبود بخشید.

$$\begin{aligned} E[\hat{\sigma}^2] &= \frac{1}{n} E \left[\sum_{i=1}^n (x_i - \bar{X})^2 \right] \\ &= \frac{1}{n} \left(\sum_{i=1}^n E[(x_i - \mu)^2] - n E[(\bar{X} - \mu)^2] \right) \\ &= \frac{1}{n} \left(n\sigma^2 - n \frac{\sigma^2}{n} \right) = \frac{1}{n} (n\sigma^2 - \sigma^2) = \frac{n-1}{n} \sigma^2 \end{aligned} \quad (114)$$

حال می توانیم با تعریف متغیر $s^2 = \frac{n}{n-1} \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2$ بایاس را صفر کنیم. حالا ماتریس کوواریانس در حالت حدی (asymptotic) به صورت زیر می شود.

$$\text{asymtotic covriance matrix} = \Sigma = \begin{pmatrix} \frac{s^2}{n} & . \\ . & \frac{2s^4}{n} \end{pmatrix} \quad (115)$$

با استفاده از کد برای این دیتاست به صورت زیر بدست می آید.

Better estimated asymptotic covariance matrix:

$$\begin{bmatrix} 1,087,2749 & 0. \\ 0. & 472,6515,229 \end{bmatrix}$$

۱۱-۵. بخش e: فاصله اطمینان ۹۵٪

با استفاده از بخش c می توان گفت $\hat{\sigma}^2 \xrightarrow{d} N\left(\sigma^2, \frac{2\sigma^4}{n}\right)$ و خطای استاندارد (standard error-SE) به صورت زیر است.

$$\widehat{SE}(\hat{\sigma}^2) = \sqrt{\frac{2\hat{\sigma}^4}{n}} \quad (116)$$

برای داشتن فاصله اطمینان ۹۵٪ باید داشته باشیم $\hat{\sigma}^2 \pm 1.96 \sqrt{\frac{2\hat{\sigma}^4}{n}}$

حال با استفاده از کد می توان این مقادیر را محاسبه کرد.

۹۵٪ large-sample CI for sigma^۲: (np.float64(۱۷۳.۹۲۰۰۴۹۶۷۹),
np.float64(۲۵۸,۷۱۶۸۸۹۸۲۱))

۱۱-۶. بخش f: تست فرض

در این بخش یک مسئله آشکار سازی داریم که با استفاده از روش های مختلف تست فرض آن را حل می کنیم.

$$H_0: \mu = 10.3, \quad H_1: \mu \neq 10.3 \quad (117)$$

۱۱-۶-۱. Z-test

برای آشنایی با این روش تست فرض، از [۱۱] استفاده کردم. در کل این روش وقتی استفاده می شود که توزیع آماری تحت فرضیه ی صفر می تواند به صورت یک توزیع نرمال تخمین زده شود. یعنی تحت فرضیه ی صفر داریم: $Z \xrightarrow{d} N(0,1)$

مقدار Z-statistic به صورت زیر محاسبه می شود. ($\mu. = ۱۰۳$)

$$Z = \frac{\hat{\mu} - \mu.}{\sqrt{\frac{\hat{\sigma}^2}{n}}} \quad (۱۱۸)$$

با استفاده از کد، مقادیر p-value و Z-statistic به صورت هستند.

Z-test statistic: -۲,۴۲۲۶۰۴۷۲۲۶۲۷۲۱۰۷

Z-test p-value: ۰,۰۱۵۴۰۹۶۸۳۵۹۱۶۵۴۴

مقدار بزرگ Z-statistic و مقدار اندک p-value، فرضیه ی صفر را رد می کنند.

۱۱-۶-۲. Likelihood ratio chi-squared test

برای آشنایی با این روش تست فرض از [۱۲] استفاده کردم. این روش وقتی مورد استفاده قرار میگیرد که بین درستی یک فرضیه و متمم آن تصمیم می گیریم. آماری آن به این صورت محاسبه می شود.

$$LR = 2(\ell(\hat{\mu}, \hat{\sigma}^2) - \ell(\mu., \tilde{\sigma}^2)) \quad (۱۱۹)$$

که $(\hat{\mu}, \hat{\sigma}^2)$ بدون هیچ محدودیتی و برای تمامی تخمین های MLE و $(\mu., \tilde{\sigma}^2)$ برای تخمینگر های MLE محدود به فرضیه ی صفر است.

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu.)^2 \quad (۱۲۰)$$

به صورت حدی خواهیم داشت: $LR \xrightarrow{d} \chi^2(1)$ و حل فرم بسته آن به صورت زیر است.

$$LR = n \log \left(\frac{\tilde{\sigma}^2}{\hat{\sigma}^2} \right) \quad (۱۲۱)$$

با استفاده از کد مقادیر به صورت زیر محاسبه شدند.

LR statistic: ۵,۷۸۴۵۴۸۷۷۸۵۹۶۵۷۲

LR test p-value: ۰,۰۱۶۱۶۷۶۴۷۵۱۶۰۷۷۵۷۷

۱۱-۶-۳. Wald chi-squared test

برای آشنایی با این روش تست فرض از [۱۳] استفاده کردم. این تست فرض برای وقتی است که می خواهیم بدانیم پارامتر تخمین زده شده توسط نمونه برابر با پارامتر مورد نظر است یا خیر. آماره‌ی آن به صورت زیر تعریف می شود.

$$W = \frac{(\hat{\mu} - \mu_*)^2}{\sqrt{\frac{\hat{\sigma}^2}{n}}} \quad (122)$$

و به صورت حدی داریم: $W \xrightarrow{d} \chi^2(1)$

با توجه به Z-test به راحتی می توان دریافت که $W = Z^2$. با استفاده از کد مقادیر به صورت زیر محاسبه شدند.

Wald statistic: ۵,۸۶۹۰۱۳۶۴۲۰۹۵۶۶۳

Wald test p-value: ۰,۰۱۵۴۰۹۶۸۳۵۹۱۶۵۴۲۸۸

با توجه به هر سه تست فرض می توان گفت، برای تعداد زیادی نمونه، هر سه تست فرض p-value های مشابهی می دهند ولی هر کدام از رابطه و مفهوم مختص خودشان می آیند.

۱۱-۷. بخش g: coefficient of variation

پارامتر coefficient of variation به صورت زیر محاسبه می شود.

$$CV = \frac{\sigma}{\mu} \quad (123)$$

۱۱-۷-۱. خاصیت invariance نسبت به یک ضریب ثابت مثبت

اگر X یک متغیر تصادفی با میانگین μ و واریانس σ^2 باشد. آنگاه متغیر تصادفی $Y = cX$ دارای میانگین $c\mu$ و انحراف معیار $c\sigma$ است.

بر این اساس coefficient of variation برای متغیر Y به صورت زیر است.

$$CV(Y) = \frac{c\sigma}{c\mu} = \frac{\sigma}{\mu} = CV(X) \quad (124)$$

۱۱-۷-۲. تخمین coefficient of variation

خاصیت invariance می گوید؛ اگر $\hat{\theta}$ تخمین MLE برای θ باشد آنگاه $g(\hat{\theta})$ تخمینگر MLE برای $g(\theta)$ است. پس بر این اساس اگر $\hat{\mu}$ تخمین MLE برای μ و $\hat{\sigma}$ تخمین MLE برای σ باشد، آنگاه تخمین MLE برای coefficient of variation به صورت زیر است.

$$\widehat{CV} = \frac{\hat{\sigma}}{\hat{\mu}} \quad (125)$$

با استفاده از کد می توان این مقدار را برای داده های مشخص شده محاسبه کرد.

MLE of coefficient of variation: ۰٫۱۴۶۳۷۴۳۶۱۴۶۲۶۳۴۷۴

۱۱-۷-۳. ۹۵% Confidence Interval via the Delta Method

بردار پارامتر $\theta = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}$ و تابع $g(\theta) = \frac{\sqrt{\sigma^2}}{\mu}$ است. مشتق این تابع نسبت به μ و σ^2 به صورت زیر است.

$$\frac{\partial g}{\partial \mu} = -\frac{\sqrt{\sigma^2}}{\mu^2}, \quad \frac{\partial g}{\partial \sigma^2} = \frac{1}{2\mu\sqrt{\sigma^2}}$$

پس گرادیان تابع به صورت زیر است.

$$\nabla_{\theta} g(\theta) = \begin{pmatrix} -\frac{\sqrt{\sigma^2}}{\mu^2} \\ \frac{1}{2\mu\sqrt{\sigma^2}} \end{pmatrix} \quad (126)$$

در بخش های قبلی ماتریس کوواریانس حالت حدی را حساب کردیم و با Σ نمایش دادیم. با استفاده از روش دلتا داریم:

$$var(\widehat{CV}) = \nabla_{\theta} g(\theta)^T \Sigma \nabla_{\theta} g(\theta) = \frac{1}{n} \left(\frac{\sigma^4}{\mu^4} + \frac{\sigma^2}{2\mu^2} \right) \quad (127)$$

با استفاده از این مقادیر، مقدار خطای استاندارد (standard error-SE) به صورت زیر است.

$$\widehat{SE}(\widehat{CV}) = \sqrt{\frac{1}{n} \left(\frac{\hat{\sigma}^4}{\hat{\mu}^4} + \frac{\hat{\sigma}^2}{2\hat{\mu}^2} \right)} \quad (128)$$

حال فاصله اطمینان ۹۵٪ به صورت زیر بدست می آید.

$$\widehat{CV} \pm 1.96 \times \widehat{SE}(\widehat{CV}) \quad (129)$$

با استفاده از کد و داده ها می توان به صورت زیر محاسبه کرد.

95% delta-method CI for CV: (np.float64(0.1317255563598098),
np.float64(0.1610231660604097))

پرسش ۱۲ – Non-Linear Regression using basis functions

این سوال در مورد پیاده سازی رگرسیون غیر خطی با استفاده از توابع پایه است که با استفاده از یک تبدیل هر داده را به یک بردار تبدیل می کند. تمامی مراحل مطابق آنچه سوال خواسته پیش می برم و در اینجا مستقیماً نتایج را بررسی خواهیم کرد.

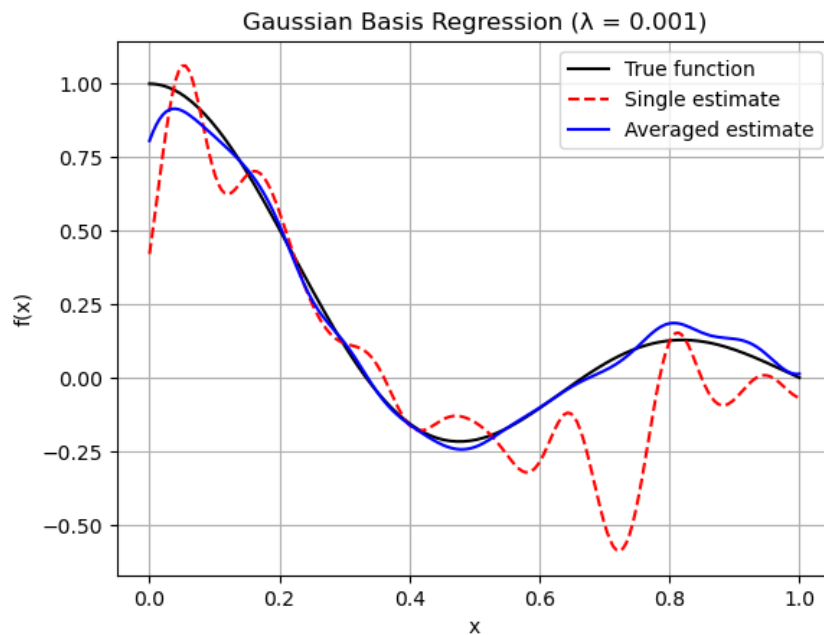
۱۲-۳. نمودارهای gaussian basis regression به ازای لاندارهای مختلف

در این بخش به ازای لاندارهای مختلف تابع $f(x)$ رسم شده و در ادامه تحلیل می شود.

۱۲-۳-۱. به ازای $\lambda = 0.001$:

در این حالت مدل تقریباً بدون regularization است. ضرایب β آزادی زیادی دارند. در تخمینی که به ازای یک مجموعه داده زدیم (نمودار خط چین قرمز) می بینیم که مدل بسیار نوسان دار است و به نویز حساس است. و در بعضی نواحی به شدت از تابع واقعی فاصله می گیرد. علت این اتفاق این است که مدل دچار overfitting شده است. تخمین میانگین گیری شده (نمودار آبی رنگ) نسبت به تخمین تک بار (یک مجموعه داده) صاف تر است. عملاً نمودار خیلی بهتر شده و به تابع سینک اصلی بسیار نزدیک تر شده.

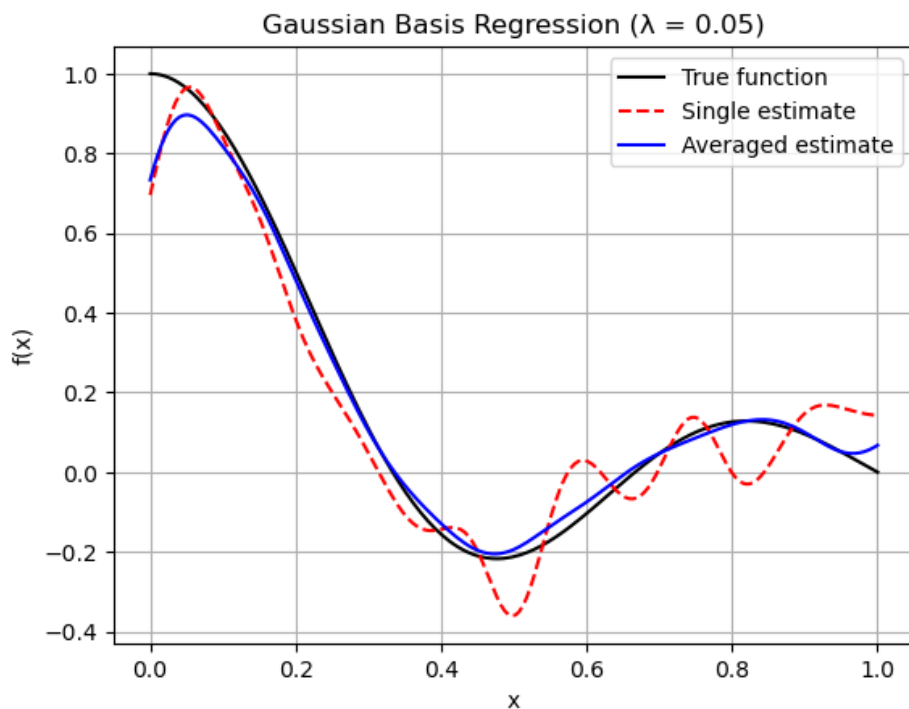
در کل در این حالت مدل وازیانس بالایی دارد. راه حل ارائه شده عملکرد چندان راضی کننده ای ندارد



شکل ۳. نمودار رگرسیون به ازای $\lambda = 0.001$

۱۲-۳-۲. به ازای $\lambda = 0.05$:

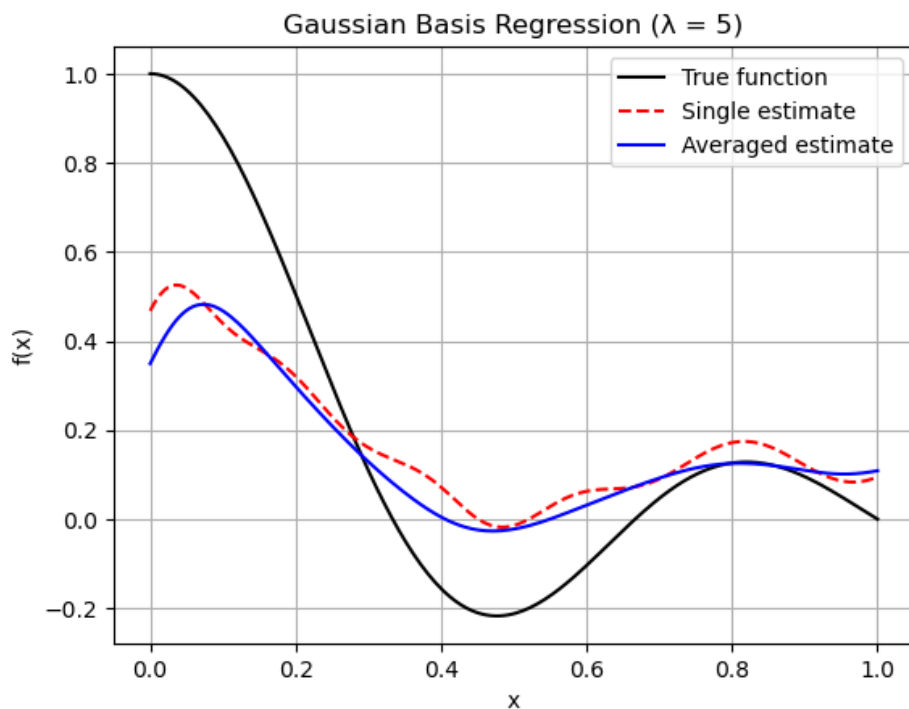
در این حالت مدل تعادل خوبی بین انعطاف پذیری و محدود سازی ضرایب ایجاد کرده. در تخمینی که به ازای یک مجموعه داده زدیم (نمودار خط چین قرمز) می بینیم که بسیار به نمودار واقعی نزدیک تر و نوسان آن کم تر شده (هر چند همچنان نوسان را می توان به وضوح دید). تخمین میانگین گیری شده (نمودار آبی رنگ) نسبت به تخمین یک مجموعه خیلی بهتر شده. تقریباً بر تابع سینک منطبق شده و صاف (smooth) شده است. بر حسب این نمودار می بینیم که $\lambda = 0.05$ بهترین انتخاب است. (مقدار $\lambda = 5$ را هم در ادامه خواهیم دید که وضع آن بدتر است).



شکل ۴. نمودار رگرسیون به ازای $\lambda = 0.05$

۱۲-۳-۳. به ازای $\lambda = 5$:

در این حالت ضرایب به شدت به صفر میل داده می شوند. یعنی مدل بیش از حد ساده می شود. در تخمین به ازای یک مجموعه داده (نمودار خط چین قرمز) می بینیم که دامنه تابع به شدت کاهش یافته و مدل نمی تواند ساختار تابع سینک را بیاموزد. در این حالت مدل دچار **underfitting** شده است. بایاس بالا دارد و کیفیت تخمین پایین است. به حدی که حتی تخمین میانگین گیری شده هم که در نمودار های قبل خوب عمل می کرد، نمی تواند به خوبی عمل کند.



شکل ۵. نمودار رگرسیون به ازای $\lambda = 5$

پس همان طور که گفته شد $\lambda = 0.05$ بهترین انتخاب است. علت اینکه تابع میانگین گیری شده عملکرد بهتری از خود نشان می دهد این است که واریانس مدل در ۱۰۰ بار تکرار و جمع شده و با میانگین گیری، واریانس آن بسیار کاهش می یابد. به همین دلیل نمودار تابع میانگین گیری شده بسیار به تابع واقعی نزدیک می شود.

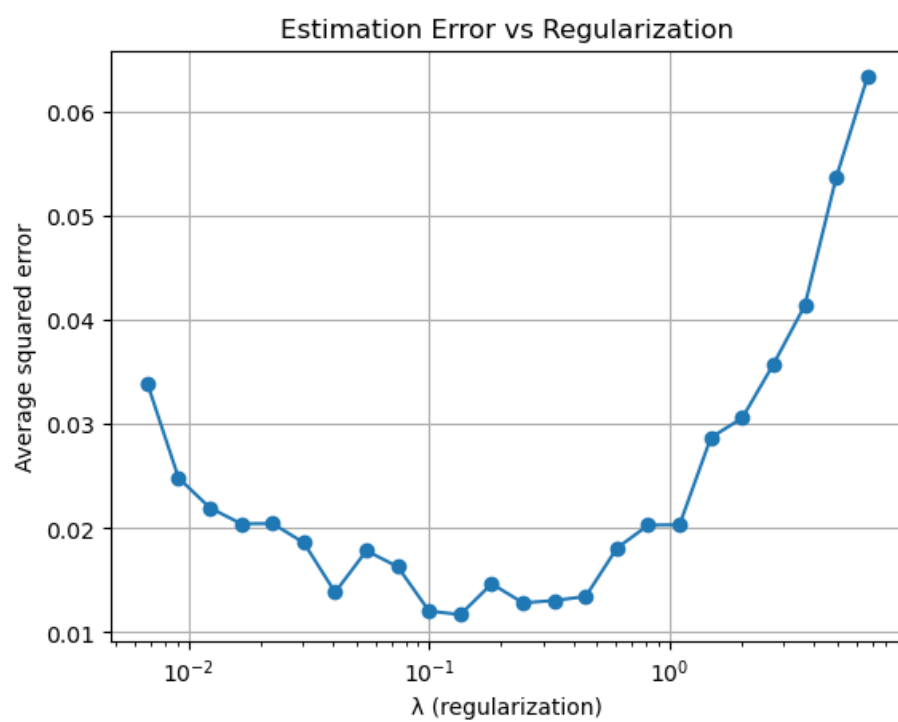
۱۲-۴. نمودار میانگین مربع خطاها بر حسب λ

در این بخش خواسته شده که دو سمت نمودار میانگین مربع خطا را تحلیل کنیم.

سمت چپ نمودار که مقادیر لاندار بسیار کم است در واقع regularization ناچیز است. مدل ها همگی به نوبت حساس هستند و ضرایب β بزرگ و ناپایدار می شوند. مدل ها به ازای این لاندا های کوچک واریانس زیادی داشته و دچار overfitting شده اند.

سمت راست نمودار که مقادیر لاندار بسیار بزرگ هستند، شاهد regularization زیادی هستیم. ضرایب به شدت به سمت صفر میل داده می شوند. مدل بیش از حد ساده می شود و خطا افزایش می یابد. در این حالت مدل بایاس زیادی پیدا می کند و دچار underfitting می شود.

در مورد نقطه‌ی میانی این نمودار هم می توان گفت که بهترین trade-off را بین بایاس و واریانس مدل برقرار کرده.



شکل ۶. نمودار میانگین مربع خطا به ازای مقادیر مختلف λ

- [۱] "axvline," [Online]. Available:
https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.axvline.html.
- [۲] Available: .[خطی ادرن], Condition number,” wikipedia“
.https://en.wikipedia.org/wiki/Condition_number
- [۳] Available: .[خطی ادرن], Euclidean distance,” wikipedia“
.https://en.wikipedia.org/wiki/Euclidean_distance
- [۴] Available: .[خطی ادرن], Taxicab geometry,” wikipedia“
.https://en.wikipedia.org/wiki/Taxicab_geometry
- [۵] Available: .[خطی ادرن], Hamming distance,” wikipedia“
.https://en.wikipedia.org/wiki/Hamming_distance
- [۶] Available: .[خطی ادرن], Cosine similarity,” wikipedia“
.https://en.wikipedia.org/wiki/Cosine_similarity
- [۷] “درخت تصمیم در یادگیری ماشین – به زبان ساده + نمونه کد”, Faradars , [ادرن خطی].
Available: https://blog.faradars.org/%D%A%AF%D%B%E%DAA-%D%.%A%.%B%.%C%.%D%.%E%.%F%.%G%.%H%.%I%.%J%.%K%.%L%.%M%.%N%.%O%.%P%.%Q%.%R%.%S%.%T%.%U%.%V%.%W%.%X%.%Y%.%Z%.%[%D%.%DB%AC%D^AV%D^AFDA%AF%DB%AC%D^BV^DB%AC-Dq/.8d/Dk/Av/Dk/Bf/Db/kc/Dq/.86/#/dk.af/dk/bv/dk/ae/dk/aa-.%dk/aa/dk/bd/dq/.8d/db/kc/dq/.8d-/dk.af/dk/bv-/db/kc/dk/av/dk/af-
%
- [۸] Available: .[خطی ادرن], minimize,” SciPy“
.https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.minimize.html
- [۹] Available: .[ادرن خطی], Intro to Scipy Optimization: Minimize Method,” youtube“
.https://www.youtube.com/watch?v=G-yP-TM-oag
- [۱۰] Available: .[خطی ادرن], Fisher information,” wikipedia“
.https://en.wikipedia.org/wiki/Fisher_information

- [۱۱] Available: <https://en.wikipedia.org/wiki/Z-test> .[ادرون خطی] ,Z-test,” wikipedia“
- [۱۲] Available: .[ادرون خطی] ,Likelihood-ratio test,” wikipedia“
https://en.wikipedia.org/wiki/Likelihood-ratio_test
- [۱۳] Available: .[ادرون خطی] ,Wald test,” wikipedia“
https://en.wikipedia.org/wiki/Wald_test