

تکلیف ششم درس شناسایی الگو

وحید ملکی
شماره دانشجویی: ۴۰۳۱۳۰۰۴

۲۰۲۵ نوامبر ۷

سوال ۷: نفرین ابعاد (Curse of Dimensionality)

الف- نسبت مساحت دایره به مربع در فضای دو بعدی

ما یک دایره به شعاع r محاط در یک مربع به ضلع $2r$ داریم، مساحت دایره (V_c) و مساحت مربع (V_s) به صورت زیر محاسبه می شوند:

$$V_c = \pi r^2$$

$$V_s = (2r)^2 = 4r^2$$

بنابراین نسبت مساحت ها برابر است با:

$$\frac{V_c}{V_s} = \frac{\pi r^2}{4r^2} = \frac{\pi}{4} \approx 0.785$$

این بدان معناست که در فضای دو بعدی، حدود 78.5% از ثوانه ها (با توزیع یکنواخت) درون دایره قرار می گیرند.

ب- نسبت فاصله میانگین تا ضلع و گوشه مربع

فاصله مرکز (مبدأ) تا نزدیکترین نقطه روی ضلع مربع برابر با شعاع دایره محاطی است:

$$d_{side} = r$$

فاصله مرکز تا گوشه مربع (وتر مثلث قائم الزاویه با اضلاع r) برابر است با:

$$d_{corner} = \sqrt{r^2 + r^2} = r\sqrt{2}$$

نسبت فاصله مرکز از ضلع به فاصله آن از گوشه برابر است با:

$$\frac{d_{side}}{d_{corner}} = \frac{r}{r\sqrt{2}} = \frac{1}{\sqrt{2}} \approx 0.707$$

ج- نسبت حجم فوق کره به فوق مکعب در فضای n بعدی

رابطه حجم فوق کره به شعاع r در فضای n بعدی داده شده است:

$$V_c(n) = \frac{\pi^{n/2} r^n}{\Gamma(\frac{n}{2} + 1)}$$

حجم فوق مکعب به ضلع $2r$ در فضای n بعدی برابر است با:

$$V_s(n) = (2r)^n = 2^n r^n$$

نسبت این دو حجم برابر است با:

$$\frac{V_c(n)}{V_s(n)} = \frac{\frac{\pi^{n/2} r^n}{\Gamma(\frac{n}{2} + 1)}}{\frac{2^n r^n}{2^n \Gamma(\frac{n}{2} + 1)}} = \frac{\pi^{n/2}}{2^n \Gamma(\frac{n}{2} + 1)} = \frac{(\frac{\pi}{4})^{n/2}}{\Gamma(\frac{n}{2} + 1)}$$

د- بررسی حد نسبت حجم‌ها وقتی $n \rightarrow \infty$

با استفاده از تقریب استرلینگ داده شده برایتابع گاما:

$$\Gamma(x+1) \approx (2\pi)^{0.5} e^{-x} x^{x+0.5}$$

با فرض $x = \frac{n}{2}$ ، مخرج کسر قسمت (ج) را تقریب می‌زنیم:

$$\Gamma(\frac{n}{2} + 1) \approx \sqrt{2\pi} e^{-n/2} (\frac{n}{2})^{\frac{n+1}{2}}$$

حال نسبت حجم‌ها را بازنویسی می‌کنیم:

$$\text{Ratio} \approx \frac{\pi^{n/2}}{2^n \left[\sqrt{2\pi} e^{-n/2} (\frac{n}{2})^{\frac{n}{2}} (\frac{n}{2})^{0.5} \right]}$$

با ساده‌سازی عبارت:

$$\text{Ratio} \approx \frac{1}{\sqrt{\pi n}} \left(\frac{\pi e}{2n} \right)^{n/2}$$

همانطور که مشاهده می‌شود، در عبارت داخل پرانتز، مخرج شامل ترم n است. وقتی $n \rightarrow \infty$ ، عبارت $\frac{\pi e}{2n}$ به سمت صفر میل می‌کند و به توان $n/2$ نیز می‌رسد که با سرعت بسیار زیادی به صفر نزدیک می‌شود.

$$\lim_{n \rightarrow \infty} \frac{V_c}{V_s} = 0$$

این نشان می‌دهد که در ابعاد بالا، حجم فوق کره (ناحیه مرکزی) نسبت به حجم فوق مکعب ناچیز است.

ه- نسبت فاصله مرکز تا وجه و گوشه در فضای n بعدی

فاصله مرکز تا نزدیکترین وجه فوق مکعب همچنان برابر است با:

$$d_{face} = r$$

فاصله مرکز تا دورترین نقطه (گوشه فوق مکعب) با استفاده از تعمیم قضیه فیثاغورس در فضای n بعدی برابر است با:

$$d_{corner} = \sqrt{\sum_{i=1}^n r^2} = \sqrt{nr^2} = r\sqrt{n}$$

نسبت این فواصل:

$$\frac{d_{face}}{d_{corner}} = \frac{r}{r\sqrt{n}} = \frac{1}{\sqrt{n}}$$

حد این نسبت وقتی $\infty \rightarrow n$

$$\lim_{n \rightarrow \infty} \frac{1}{\sqrt{n}} = 0$$

و- نتیجه‌گیری و مفهوم نفرین ابعاد

نتیجه‌گیری هندسی: بر اساس بند (د)، حجم در دسترس در فضای ابعاد بالا به سرعت به سمت صفر می‌کند، به این معنی که تقریباً تمام حجم یک فوق مکعب در "گوشه‌ها" (Corners) و پوسته‌های بیرونی آن متتمرکز شده است و مرکز آن تقریباً خالی است. بر اساس بند (ه)، فاصله گوشه‌ها از مرکز بسیار بیشتر از فاصله وجوده از مرکز می‌شود، به طوری که شکل به نوعی "تیغ‌دار" می‌شود. چرا این یک آفت (Curse) است؟

۱. پراکندگی داده‌ها (Data Sparsity): برای اینکه بتوانیم فضای ویژگی را به درستی با نمونه‌های آموزشی پوشش دهیم، با افزایش ابعاد، تعداد نمونه‌های مورد نیاز به صورت نمایی افزایش می‌یابد. از آنجا که تعداد داده‌ها محدود است، فضا بسیار خالی می‌ماند و تخمین چگالی دشوار می‌شود.

۲. بی‌معنی شدن فاصله اقلیدسی: همانطور که نشان داده شد، فاصله دورترین و نزدیک‌ترین نقاط به هم نزدیک می‌شوند (به صورت نسبی). این باعث می‌شود که معیارهای شباهت مبتنی بر فاصله (مانند الگوریتم KNN) کارایی خود را از دست بدنهند، زیرا همه نمونه‌ها تقریباً در فاصله یکسانی از هم قرار می‌گیرند و مفهوم "همسایگی" از بین می‌رود.

۳. خطر بیش‌برازش (Overfitting): به دلیل خالی بودن فضا، دسته‌بندها می‌توانند به راحتی مرزهای تصمیم‌گیری پیچیده‌ای پیدا کنند که نویزها را مدل کرده و روی داده‌های جدید تعمیم‌پذیری (Generalization) ضعیفی داشته باشند.

سوال ۱۱: دسته‌بند KNN با وزن‌دهی ویژگی و فاصله

الف- انتخاب وزن مناسب برای ویژگی‌ها

هدف از وزن‌دهی به ویژگی‌ها در اینجا، کاهش اثر مقیاس یا افزایش اهمیت ویژگی‌های تفکیک کننده است. با نگاه به نودار داده‌ها، مشاهده می‌کنیم:

- کلاس‌ها (علامت‌های + و -) عمدتاً در امتداد محور افقی (x_1) از هم جدا شده‌اند (کلاس مثبت در سمت چپ و کلاس منفی در سمت راست).

در امتداد محور عمودی (x_2), هر دو کلاس پراکنده‌گی زیادی دارند و همپوشانی کامل وجود دارد. این یعنی محور عمودی اطلاعات تفکیک کننده کمتری دارد و بیشتر شبیه به نویز عمل می‌کند.

بنابراین، برای افزایش دقت دسته‌بندی و برجسته کردن ویژگی مهم‌تر، باید وزن بیشتری به x_1 و وزن کمتری به x_2 بدهیم تا فاصله‌ها در راستای افقی (که نشان‌دهنده تغییر کلاس است) تاثیر بیشتری در محاسبه فاصله اقلیدسی داشته باشند.
نتیجه:

$$w_1 = \frac{2}{3} \quad w_2 = \frac{1}{3} \quad (\text{برای محور } x_2)$$

ب- تعیین کلاس نقطه مجهول با KNN ($k = 5$)

نقطه مجهول (تست) را با x_u نشان می‌دهیم. با مشاهده چشمی نودار، مختصات تقریبی نقطه مجهول برابر است با:

$$x_u \approx (5, 6.2)$$

فرمول فاصله اقلیدسی وزن‌دار برای دو نقطه x و y عبارت است از:

$$d(x, y) = \sqrt{w_1(x_1 - y_1)^2 + w_2(x_2 - y_2)^2} = \sqrt{\frac{2}{3}(x_1 - y_1)^2 + \frac{1}{3}(x_2 - y_2)^2}$$

حال ۵ نزدیک‌ترین همسایه را به صورت چشمی انتخاب کرده و فاصله آن‌ها را محاسبه می‌کنیم. (فرض می‌کنیم نزدیک‌ترین همسایه‌ها شامل ۴ نمونه از کلاس منفی و ۱ نمونه از کلاس مثبت که نزدیک‌تر است باشند):

| کلاس | مختصات تقریبی | اختلاف مختصات | فاصله وزن‌دار (d) | وزن رأی (d) |
|------|---------------|---------------|---|-------------|
| - | (5, 5) | (0, 1.2) | $\sqrt{0 + \frac{1}{3}(1.44)} \approx 0.69$ | 1.45 |
| - | (6, 6.2) | (1, 0) | $\sqrt{\frac{2}{3}(1) + 0} \approx 0.81$ | 1.23 |
| - | (5, 7.2) | (0, 1) | $\sqrt{0 + \frac{1}{3}(1)} \approx 0.58$ | 1.72 |
| + | (4, 6) | (1, 0.2) | $\sqrt{0.66 + 0.013} \approx 0.82$ | 1.21 |
| - | (6, 7.2) | (1, 1) | $\sqrt{0.66 + 0.33} \approx 1.00$ | 1.00 |

جدول ۱: محاسبه فاصله و وزن برای ۵ همسایه نزدیک

محاسبه مجموع امتیازات (Weights) برای هر کلاس:

- کلاس منفی (-): مجموع وزن‌ها = $1.45 + 1.23 + 1.72 + 1.00 = 5.40$

- کلاس مثبت (+): مجموع وزن‌ها = 1.21

نتیجه‌گیری: با توجه به اینکه مجموع وزن‌های کلاس منفی بسیار بیشتر از کلاس مثبت است، نقطه x_u به کلاس منفی (-) تعلق می‌گیرد.

سوال ۱۲: رفتار دسته‌بند ۱NN

پاسخ

بله، این امکان وجود دارد که تمامی نمونه‌های تست به کلاس مثبت اختصاص یابند، حتی اگر در داده‌های آموزشی کلاس منفی نیز وجود داشته باشد.

تشریک و مثال

در دسته‌بند «یک نزدیک‌ترین همسایه»، (1NN) فضای ویرگی بر اساس موقعیت داده‌های آموزشی به نواحی مختلفی تقسیم می‌شود که به آن‌ها «سلول‌های ورونوی» (Voronoi Cells) می‌گویند. هر نقطه در فضا برچسب نمونه‌ای را می‌گیرد که به آن نزدیک‌تر است. اگر موقعیت نمونه‌های تست به گونه‌ای باشد که هیگی در سلول‌های مربوط به نمونه‌های کلاس «مثبت» قرار گیرند (یعنی فاصله اقلیدسی آن‌ها تا یک نمونه مثبت کمتر از فاصله تا هر نمونه منفی باشد)، همگی مثبت برچسب گذاری می‌شوند. این حالت معمولاً زمانی رخ می‌دهد که نمونه‌های کلاس منفی «داده‌های پرت» (Outlier) باشند یا از ناحیه مرکز داده‌های تست بسیار دور باشند.

مثال عددی (در فضای یک بعدی):

فرض کنید داده‌های آموزشی (D_{train}) شامل دو نمونه باشد:

• یک نمونه از کلاس مثبت (+) در موقعیت 2 $x = 2$

• یک نمونه از کلاس منفی (-) در موقعیت 10 $x = 10$

مرز تصمیم گیری دقیقاً وسط این دو نقطه، یعنی در $x = 6$ قرار دارد. هر عددی کمتر از 6 به کلاس مثبت و هر عددی بیشتر از 6 به کلاس منفی تعلق می‌گیرد. حال فرض کنید داده‌های تست (D_{test}) شامل نقاط زیر باشند:

$$\{1, 3, 4, 5.5\}$$

از آنجا که تمام این اعداد از 6 کوچک‌تر هستند، فاصله همه آن‌ها تا نمونه مثبت (2) کمتر از فاصله‌شان تا نمونه منفی (10) است. بنابراین خروجی دسته‌بند برای تمام نمونه‌های تست، کلاس «مثبت» خواهد بود.

سوال ۱۳: محاسبات و مفاهیم KNN

الف- دسته‌بندی نمونه (3, 6) با فاصله منهتن (1NN)

$d = |x_1 - p_1| + |x_2 - p_2|$ در نظر می‌گیریم. فاصله منهتن (City Block) با فرمول Block) است را تا ۵ نمونه آموزشی موجود در جدول محاسبه می‌کنیم:

| فاصله نهایی (d) | محاسبه فاصله ($ 3 - x_1 + 6 - x_2 $) | برچسب (Y) | نمونه آموزشی (P_i) |
|---------------------|--|---------------|------------------------|
| 3 | $ 3 - 1 + 6 - 7 = 2 + 1$ | +1 | (1, 7) |
| 3 | $ 3 - 3 + 6 - 3 = 0 + 3$ | +1 | (3, 3) |
| 7 | $ 3 - 1 + 6 - 1 = 2 + 5$ | -1 | (1, 1) |
| 4 | $ 3 - 5 + 6 - 4 = 2 + 2$ | -1 | (5, 4) |
| 2 | $ 3 - 2 + 6 - 5 = 1 + 1$ | -1 | (2, 5) |

جدول ۲: محاسبه فاصله منهتن نمونه تست تا نمونه‌های آموزشی

نتیجه: کمترین فاصله برابر با 2 است که مربوط به نمونه (5, 2) می‌باشد. چون برچسب این نمونه 1- است، طبق قاعده 1NN، نمونه تست به کلاس منفی (-1) تعلق می‌گیرد.

ب- بررسی شکل مرز تصمیم‌گیری در فاصله منهتن

پاسخ: نادرست

دلیل: در فاصله منهتن (L_1)، مکان هندسی نقاطی که از دو نقطه خاص فاصله یکسانی دارند، لزوماً خطوط موازی با محورهای مختصات نیستند. این مرزها می‌توانند شامل پاره‌خط‌هایی با زاویه ۴۵ درجه (شیب ± 1) نیز باشند. بنابراین مرز تصمیم‌گیری مجموعه‌ای از قطعات خطی است که می‌تواند شامل خطوط افقی، عمودی و یا مورب باشد.

ج- تاثیر ضرب مختصات در اسکالر بر خروجی KNN

پاسخ: نادرست

دلیل: فاصله اقلیدسی نسبت به ضرب در یک عدد مثبت (Scaling) همگن است. اگر تمام مختصات در عدد 0.5 ضرب شوند، تمام فاصله‌های بین نقاط نیز دقیقاً در 0.5 ضرب می‌شوند:

$$d_{new}(x, y) = \sqrt{\sum (0.5x_i - 0.5y_i)^2} = 0.5 \times d_{old}(x, y)$$

چون تمام فاصله‌ها به یک نسبت کوچک می‌شوند، «ترتیب» (Ordering) همسایه‌ها تغییر نمی‌کند. یعنی نزدیک‌ترین همسایه قبل از تغییر مقیاس، همچنان نزدیک‌ترین همسایه باقی می‌ماند و نتیجه دسته‌بندی تغییر نخواهد کرد.

د- رابطه افزایش k و صحت دسته‌بندی

پاسخ: نادرست

دلیل: افزایش k همیشه باعث بهبود دقت نمی‌شود.

• وقتی k خیلی کوچک است (مثلاً $k = 1$)، مدل دچار بیش‌برازش (Overfitting) می‌شود و نسبت به نویز حساس است.

• وقتی k خیلی بزرگ می‌شود (به سمت N میل می‌کند)، مدل دچار کمپارازش (Underfitting) می‌شود و مرزهای تصمیم‌گیری بیش از حد هموار می‌شوند. در حالت $k = N$ ، مدل صرفاً کلاسی را انتخاب می‌کند که بیشترین تعداد کل نمونه‌ها را دارد (کلاس اکثریت)، صرف نظر از اینکه ورودی تست چگاست. بنابراین معمولاً یک k بهینه در میانه وجود دارد و غودار خطای U شکل است.