

تکلیف پنجم درس شناسایی الگو

وحید ملکی
شماره دانشجویی: ۴۰۳۱۳۰۰۴

۱۴۰۴ آبان ۲۴

سؤال ۷

داده‌ها

مجموعه داده یک بعدی به صورت زیر است:

$$X = \{0, 1, 1, 1, 2, 2, 2, 2, 3, 4, 4, 4, 5\}, \quad N = 13.$$

تکرار هر مقدار به شرح زیر است:

- ۰: یک بار
- ۱: سه بار
- ۲: چهار بار
- ۳: یک بار
- ۴: سه بار
- ۵: یک بار

الف) رسم هیستوگرام با پهنهای بازه ۱

پهنهای بازه $1 = h$ و مرکز بازه‌ها در نقاط $\{0, 1, 2, 3, 4, 5\}$ قرار دارند. بنابراین بازه‌ها به صورت زیر تعریف می‌شوند:

- $[-0.5, 0.5)$ (مرکز ۰)
- $[0.5, 1.5)$ (مرکز ۱)
- $[1.5, 2.5)$ (مرکز ۲)
- $[2.5, 3.5)$ (مرکز ۳)
- $[3.5, 4.5)$ (مرکز ۴)
- $[4.5, 5.5)$ (مرکز ۵)

تعداد نمونه‌ها در هر بازه:

۰	: مرکز	$(x = 0)$
۱	: مرکز	$()$
۲	: مرکز	$()$
۳	: مرکز	$(x = 3)$
۴	: مرکز	$()$
۵	: مرکز	$(x = 5)$

تحمین چگالی هیستوگرام در هر بازه:

$$\hat{f}_{\text{hist}}(x) = \frac{\text{تعداد در بازه}}{N \cdot h}, \quad h = 1, \quad N = 13$$

بنابراین:

$$\begin{aligned} \hat{f}_{\text{hist}}(0) &= \frac{1}{13}, & \hat{f}_{\text{hist}}(1) &= \frac{3}{13}, & \hat{f}_{\text{hist}}(2) &= \frac{4}{13}, \\ \hat{f}_{\text{hist}}(3) &= \frac{1}{13}, & \hat{f}_{\text{hist}}(4) &= \frac{3}{13}, & \hat{f}_{\text{hist}}(5) &= \frac{1}{13}. \end{aligned}$$

هیستوگرام یک تابع پله‌ای است که در هر بازه مقدار ثابت دارد.

ب) رابطه کلی تخمین پارزن

تحمین چگالی پارزن با کرنل دخلواه K و پهنای پنجه h به صورت زیر است:

$$\hat{f}(x) = \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x - x_i}{h}\right).$$

ج) تخمین پارزن با کرنل مثلثی

کرنل مثلثی تعریف شده است:

$$K(u) = (1 - |u|) \cdot \delta(|u| \leq 1),$$

که (۱) نشانگر است ($\delta(|u| \leq 1)$ اگر $|u| \leq 1$ ، در غیر این صورت 0). همچنین:

$$u = \frac{x - x_i}{h}.$$

با انتخاب $h = 1$ (مانند هیستوگرام):

$$\hat{f}(x) = \frac{1}{N} \sum_{i=1}^N (1 - |x - x_i|) \cdot \mathbf{1}(|x - x_i| \leq 1).$$

حال تخمین در نقاط $x \in \{0, 1, 2, 3, 4, 5\}$ محاسبه می‌شود:

$$x = 0 : \text{نمونه‌های در فاصله } 1 - |0 - 1| = 0, 1 - |0 - 0| = 1 \text{ وزن‌ها: } x_i = 0, 1 : \leq 1 \text{ مجموع وزن = } 1 \Rightarrow \hat{f}(0) = \frac{1}{13} \Rightarrow 1$$

$\hat{f}(1) = \frac{3}{13} \Rightarrow 3$ نمونه‌های 2 وزن‌ها: 0 (سه بار)، 1 (سه بار)، 0 مجموع وزن = $x_i = 0, 1, 2$.

$\hat{f}(2) = \frac{4}{13} \Rightarrow 4$ نمونه‌های 3 وزن‌ها: 0 (سه بار)، 1 (چهار بار)، 0 مجموع وزن = $x_i = 1, 2, 3$.

$\hat{f}(3) = \frac{1}{13} \Rightarrow 1$ نمونه‌های 4 وزن‌ها: 0 (چهار بار)، 1 (سه بار)، 0 مجموع وزن = $x_i = 2, 3, 4$.

$\hat{f}(4) = \frac{3}{13} \Rightarrow 3$ نمونه‌های 5 وزن‌ها: 0 (سه بار)، 1 (سه بار)، 0 مجموع وزن = $x_i = 3, 4, 5$.

$\hat{f}(5) = \frac{1}{13} \Rightarrow 1$ نمونه‌های 0 وزن‌ها: 0 (سه بار)، 1 (سه بار)، 0 مجموع وزن = $x_i = 4, 5$.

بنابراین:

$$\hat{f}(0, 1, 2, 3, 4, 5) = \left(\frac{1}{13}, \frac{3}{13}, \frac{4}{13}, \frac{1}{13}, \frac{3}{13}, \frac{1}{13} \right).$$

این مقادیر دقیقاً با ارتفاع هیستوگرام در مرآکر بازه‌ها یکسان هستند.

د) تخمین MLE توزیع گوسی

فرض می‌کنیم داده‌ها از توزیع نرمال $X \sim \mathcal{N}(\mu, \sigma^2)$ هستند. تخمین‌های بیشینه درست نمایی:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2.$$

مجموع داده‌ها:

$$\sum x_i = 0 + 3 \times 1 + 4 \times 2 + 3 + 3 \times 4 + 5 = 31 \Rightarrow \hat{\mu} = \frac{31}{13}.$$

مجموع مربعات انحراف:

$$\sum (x_i - \hat{\mu})^2 = \frac{352}{13} \Rightarrow \hat{\sigma}^2 = \frac{1}{13} \cdot \frac{352}{13} = \frac{352}{169}.$$

پس توزیع گوسی تخمین‌زده شده:

$$X \sim \mathcal{N}\left(\frac{31}{13}, \frac{352}{169}\right).$$

۵) مقایسه روش‌ها

• هیستوگرام: ساده، مستقیم فرکانس‌های تجربی را نشان می‌دهد. ضعف: ناپوسته، وابستگی زیاد به عرض و موقعیت بازه‌ها، نویز در داده‌های کم. کاربرد: تحلیل اکتشافی سریع.

• پارزن با کونل مثالی: چگالی پوسته و خطی-تکه‌ای تولید می‌کند. بدون فرض پارامتری. ضعف: نیاز به انتخاب h ، بیش‌بازش با h کوچک، بیش‌صفاف کردن با h بزرگ. کاربرد: داده‌های کافی، عدم اطمینان به شکل پارامتری.

• MLE گوسی: مدل پارامتری بسیار صاف با دو پارامتر. ضعف: اگر توزیع واقعی چندوجهی یا غیرنرمال باشد، برآش ضعیف دارد. کاربرد: وقتی فرض نرمالی معتبر است و نمونه متوسط است.

در این داده، توزیع تجربی تک قله در اطراف ۲ دارد. نتیجه: هیستوگرام و پارزن دقیقاً الگوی داده را در نقاط شبکه بازتولید می‌کنند اما صاف نیستند. گوسی یک خلاصه صاف مناسب ارائه می‌دهد اما جزئیات محلی را از دست می‌دهد. در عمل، هیستوگرام برای ثمیش سریع، پارزن برای تخمین غیرپارامتری، و گوسی برای مدل‌سازی پارامتری مناسب است.