

تکلیف اول درس شبکه های عصبی

وحید ملکی

شماره دانشجویی: ۴۰۳۱۳۰۰۴

۲۶ آبان ۱۴۰۴

سوال ۲

مجموعه داده مسکن کالیفرنیا، که از سرشماری سال ۱۹۹۰ ایالات متحده استخراج شده است، یکی از مجموعه داده های معروف برای مسائل رگرسیون چندمتغیره است. این داده ها شامل ویژگی های اقتصادی، جمعیتی و جغرافیایی مناطق مسکونی در کالیفرنیا بوده و هدف، پیش بینی Median House Value است.

الف: رسم و تفسیر توزیع (هیستوگرام) هریک از ویژگی ها

با استفاده از کتابخانه sklearn، داده ها بارگذاری شده و هیستوگرام تمام ویژگی ها و متغیر هدف با ۵۰ ستون (bin) رسم گردید. شکل ۱ توزیع هریک از ویژگی ها را نشان می دهد.

تفسیر هیستوگرام ها

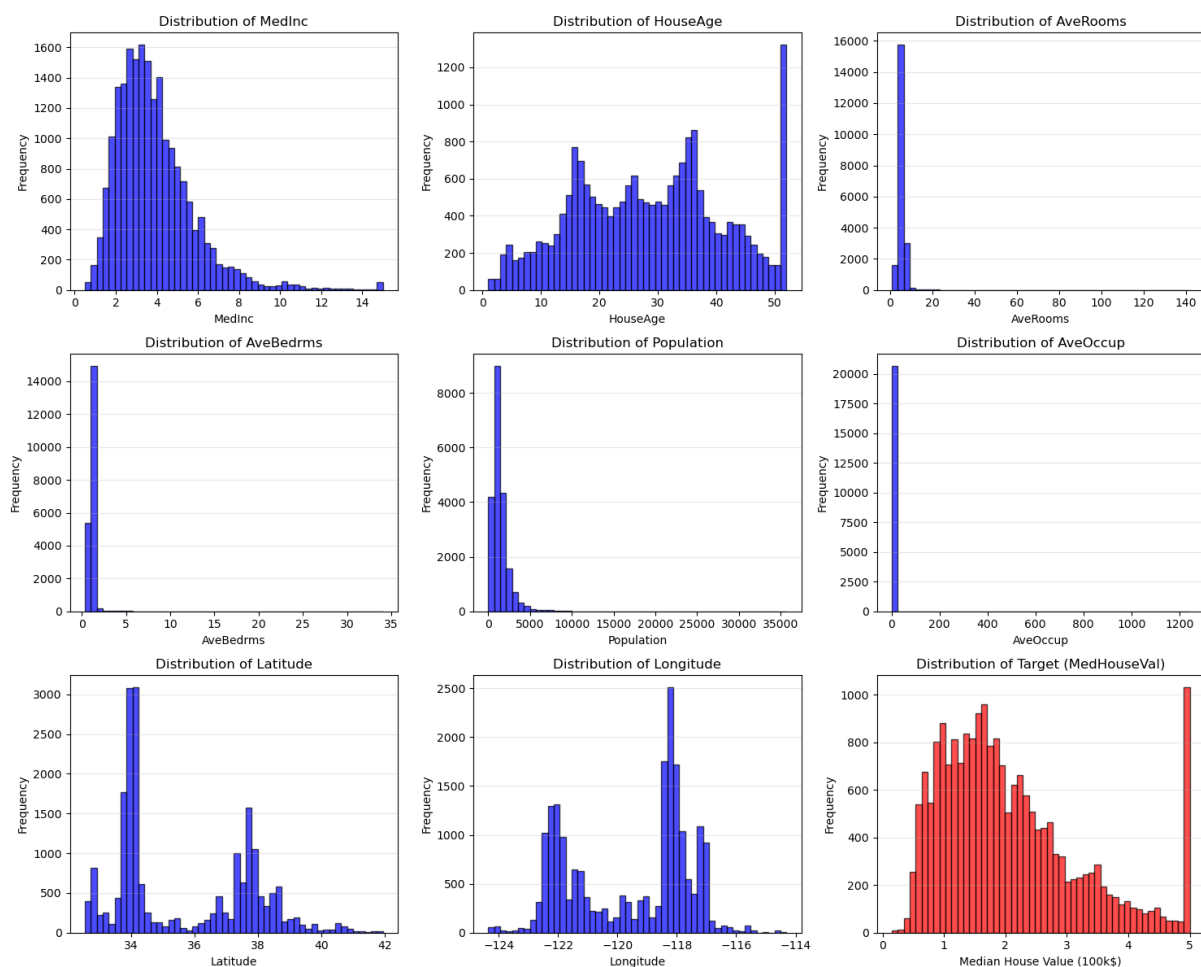
در ادامه، توزیع هر ویژگی به همراه تحلیل و اطلاعات استخراج شده از آن بررسی می شود:

۱. **MedInc (درآمد متوسط خانوار):** توزیع چپ کشیده (right-skewed) با پیک در حدود ۳ تا ۴ واحد (صد هزار دلار). بیشتر خانوارها درآمد متوسط پایینی دارند و تعداد کمی درآمد بسیار بالا (تا ۱۵) دارند. این نشان دهنده نابرابری درآمدی در مناطق است.

۲. **HouseAge (سن متوسط خانه):** توزیع نسبتاً یکنواخت با چند پیک در سنین ۱۵، ۳۰ و ۵۰ سال. این نشان می دهد خانه ها در سنین مختلف به طور نسبتاً یکسان توزیع شده اند، اما ساخت خانه در دهه های خاصی (مثلاً پس از جنگ جهانی دوم) بیشتر بوده است.

۳. **AveRooms (میانگین تعداد اتاق در هر خانه):** توزیع شدیداً چپ کشیده با پیک در حدود ۵ اتاق و دم بسیار طولانی تا بیش از ۱۲۰ اتاق. وجود مقادیر بسیار بالا (مثلاً خانه هایی با بیش از ۱۰۰ اتاق) نشان دهنده داده های پرت (outlier) یا احتمالاً خطا در داده هاست. این ویژگی نیاز به نرمال سازی یا حذف پرت دارد.

۴. **AveBedrms (میانگین تعداد اتاق خواب):** مشابه AveRooms، توزیع چپ کشیده با پیک در حدود ۱ تا ۵۰۱ و دم طولانی تا بیش از ۳۰. مقادیر بسیار بالا غیرواقعی به نظر می رسند و احتمالاً ناشی از خطا یا مناطق خاص (مثلاً خوابگاه ها) هستند.



شکل ۱: توزیع (هیستوگرام) ویژگی‌های مجموعه داده مسکن کالیفرنیا و متغیر هدف

۵. **Population (جمعیت منطقه):** توزیع شدیداً چپ کشیده با پیک در حدود ۱۰۰۰ نفر و دم تا بیش از ۳۵۰۰۰. بیشتر مناطق جمعیت کمی دارند و تعداد کمی منطقه پرجمعیت هستند. این ویژگی نیز دارای پرت‌های قابل توجه است.

۶. **AveOccup (میانگین تعداد افراد در هر خانه):** توزیع چپ کشیده با پیک در حدود ۵.۲ تا ۳ و دم بسیار طولانی تا بیش از ۱۲۰۰. مقادیر بسیار بالا کاملاً غیرعادی هستند و نشان‌دهنده داده‌های پرت یا خطا در ثبت اطلاعات است.

۷. **Latitude (عرض جغرافیایی):** توزیع چندمداله با پیک‌های اصلی در حدود ۳۴، ۳۷ و ۳۹ درجه. این الگو با توزیع جغرافیایی مناطق مسکونی در کالیفرنیا (از جنوب تا شمال) همخوانی دارد.

۸. **Longitude (طول جغرافیایی):** توزیع چندمداله با پیک‌های اصلی در حدود ۱۲۲- و ۱۱۸- درجه. این نشان‌دهنده تمرکز جمعیت در نواحی ساحلی (لس آنجلس، سان فرانسیسکو) و دره مرکزی است.

۹. **Median House Value (متغیر هدف):** توزیع چپ کشیده با پیک در حدود ۸۰.۱ تا ۲ واحد (۱۸۰ تا ۲۰۰ هزار دلار) و یک پیک غیرعادی در ۵ واحد (۵۰۰ هزار دلار). این پیک در انتهای بالایی به دلیل سقف‌گذاری (capping) در داده‌های اصلی است — در سرشماری ۱۹۹۰، قیمت‌های بالای ۵۰۰ هزار دلار به این مقدار گرد شده‌اند. این امر مدل‌سازی را پیچیده می‌کند و ممکن است نیاز به حذف یا اصلاح این نمونه‌ها باشد.

جمع‌بندی اطلاعات استخراج شده از هیستوگرام‌ها

- توزیع‌های نامتقارن و چپ کشیده: اکثر ویژگی‌ها (به‌ویژه MedInc، AveRooms، AveBedrms، Population، AveOccup) توزیع چپ کشیده دارند و نیاز به تبدیل لگاریتمی یا نرمال‌سازی دارند.
 - وجود پرت‌های شدید: ویژگی‌هایی مانند AveRooms، AveBedrms و AveOccup شامل مقادیر بسیار غیرعادی هستند که باید با روش‌های حذف پرت یا جایگزینی مدیریت شوند.
 - سقف‌گذاری در متغیر هدف: پیک در ۵۰۰ هزار دلار نشان‌دهنده censoring است و مدل‌های رگرسیون خطی ممکن است در این ناحیه خطای سیستماتیک داشته باشند.
 - اطلاعات جغرافیایی ارزشمند: توزیع Latitude و Longitude نشان‌دهنده الگوهای فضایی مشخص است که می‌تواند در مهندسی ویژگی (مثلاً محاسبه فاصله از ساحل) مفید باشد.
- این تحلیل اولیه نشان می‌دهد که پیش‌پردازش داده (نرمال‌سازی، حذف پرت، تبدیل متغیرها) برای بهبود عملکرد مدل‌های یادگیری ماشین ضروری است.

ب: مقایسه روش‌های نرمال‌سازی داده‌ها

در این بخش، چهار روش نرمال‌سازی زیر بر روی ویژگی‌های مجموعه داده مسکن کالیفرنیا اعمال شده و تأثیر هر یک بر توزیع داده‌ها، مقادیر پرت و آماره‌های توصیفی بررسی شده است:

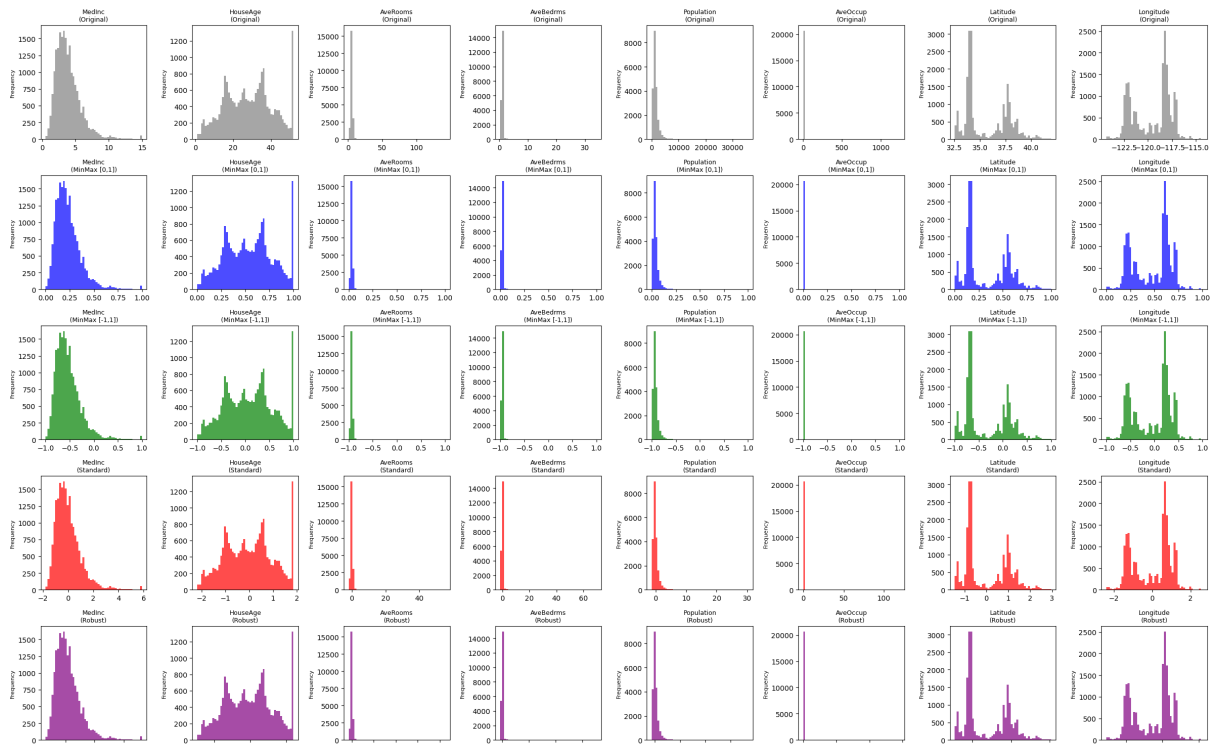
۱. نرمال‌سازی مین‌مکس به بازه $[0, 1]$

۲. نرمال‌سازی مین‌مکس به بازه $[-1, 1]$

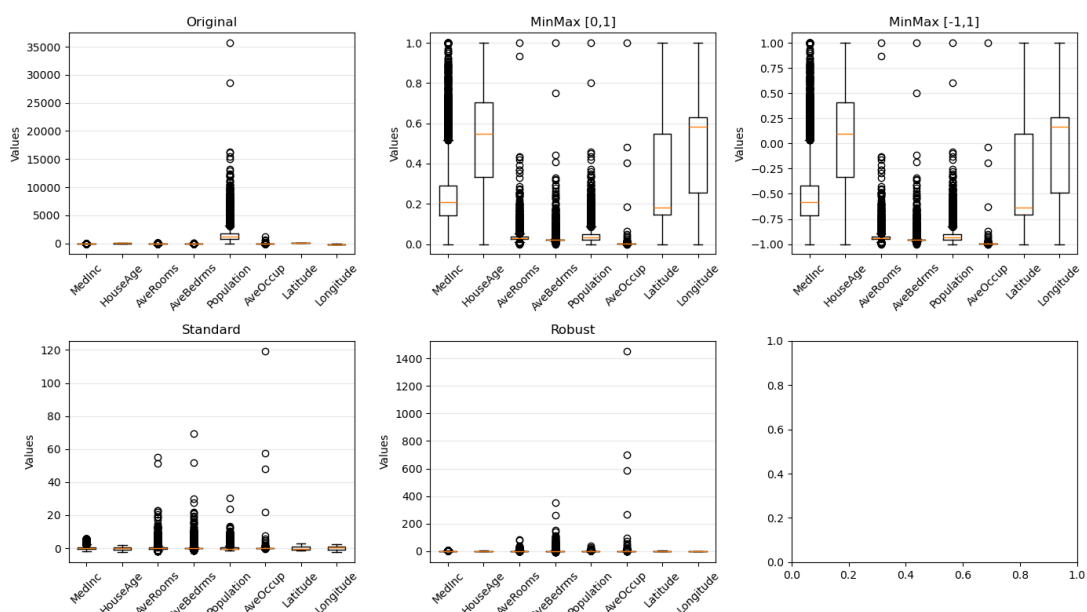
۳. نرمال‌سازی استاندارد (Z-score)

۴. نرمال‌سازی مقاوم (Robust Scaling)

شکل ۲ توزیع هیستوگرام هر ویژگی را قبل و بعد از اعمال هر روش نشان می‌دهد. همچنین شکل ۳ با استفاده از نمودار جعبه‌ای (boxplot)، تأثیر هر روش بر مدیریت داده‌های پرت را مقایسه می‌کند.



شکل ۲: مقایسه توزیع هیستوگرام ویژگی‌ها در حالت اصلی و پس از اعمال چهار روش نرمال‌سازی



شکل ۳: نمودار جعبه‌ای (boxplot) برای مقایسه تأثیر روش‌های نرمال‌سازی بر داده‌های پرت

فرمول ریاضی و تفسیر هر روش

در ادامه، برای هر روش، فرمول ریاضی و تفسیر عملی آن (اینکه دقیقاً چه کاری انجام می‌دهد و چه تأثیری بر داده‌ها دارد) ارائه شده است:

۱. نرمال‌سازی مین‌مکس به بازه $[0, 1]$

فرمول:

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

تفسیر: این روش کوچک ترین مقدار ویژگی را به ۰ و بزرگ ترین مقدار را به ۱ نگاشت می کند. تمام مقادیر دیگر به صورت خطی در این بازه توزیع می شوند. به عبارت دیگر، مقیاس ویژگی را به بازه استاندارد [۰, ۱] تبدیل می کند بدون تغییر در شکل نسبی توزیع. این روش برای الگوریتم هایی که به مقیاس مطلق حساس هستند (مثل شبکه های عصبی با تابع فعال سازی سیگموئید) مفید است، اما بسیار حساس به پرت است - یک مقدار پرت می تواند کل داده ها را فشرده کند.

۲. نرمال سازی مین مکس به بازه [-1, 1]
فرمول:

$$x' = 2 \times \frac{x - x_{\min}}{x_{\max} - x_{\min}} - 1$$

تفسیر: مشابه روش قبلی، اما بازه خروجی [-1, 1] است. ابتدا داده ها به [0, 1] نگاشت می شوند، سپس با ضرب در ۲ و کم کردن ۱، به بازه متقارن حول صفر تبدیل می شوند. این روش برای الگوریتم هایی که به مقیاس متقارن حول صفر نیاز دارند (مثل برخی روش های بهینه سازی یا شبکه های عصبی با تابع فعال سازی tanh) مناسب است. حساسیت به پرت همچنان بالاست.

۳. نرمال سازی استاندارد (Standard Scaling یا Z-score)
فرمول:

$$x' = \frac{x - \mu}{\sigma}$$

تفسیر: این روش میانگین (μ) ویژگی را به ۰ و انحراف معیار (σ) را به ۱ تبدیل می کند. به عبارت دیگر، داده ها را به توزیع نرمال استاندارد (میانگین ۰، واریانس ۱) نزدیک می کند. این روش فرض می کند داده ها تقریباً نرمال هستند و برای الگوریتم های مبتنی بر فاصله (مثل SVM، k-NN، رگرسیون خطی) بسیار مفید است. اما چون از میانگین و انحراف معیار استفاده می کند، تأثیر پرت ها بسیار زیاد است و می تواند توزیع را به شدت تحریف کند.

۴. نرمال سازی مقاوم (Robust Scaling)
فرمول:

$$x' = \frac{x - \text{median}}{Q_3 - Q_1} \quad (\text{IQR} = Q_3 - Q_1)$$

تفسیر: این روش از میانه (median) به جای میانگین و دامنه بین چارکی (IQR) به جای انحراف معیار استفاده می کند. میانه نماینده مرکزی مقاوم به پرت است و IQR فقط به ۵۰٪ داده های مرکزی وابسته است. بنابراین، این روش در حضور پرت ها بسیار مقاوم است و توزیع نسبی داده های اصلی را بدون تحریف حفظ می کند. برای داده های واقعی با نویز و پرت (مثل این دیتاست) بهترین انتخاب است.

تحلیل آماری کلی (از خروجی کد)

جدول ۱: مقایسه آماری روش‌های نرمال‌سازی (محدوده میانگین، انحراف معیار، حداقل و حداکثر)

روش	محدوده میانگین	محدوده انحراف معیار	حداقل	حداکثر
اصلی	$[-119.57, 1425.48]$	$[0.47, 1132.46]$	-124.35	35682.00
MinMax[۰،۱]	$[0.0019, 0.5420]$	$[0.0084, 0.2468]$	0.0000	1.0000
MinMax[-۱،۱]	$[-0.9962, 0.0839]$	$[0.0167, 0.4936]$	-1.0000	1.0000
Standard	$[\sim 10^{-15}, \sim 10^{-17}]$	$[1.0000, 1.0000]$	-2.39	119.42
Robust	$[-0.2849, 0.5125]$	$[0.5286, 12.1828]$	-7.66	1455.12

مشاهدات کلیدی از هیستوگرام‌ها و باکس‌پلات‌ها

- **MinMax [۰،۱] و [-۱،۱]:** این روش‌ها تمام مقادیر را به بازه محدود مین‌مکس می‌برند، اما حساس به پرت هستند. در ویژگی‌هایی مانند Population و AveRooms، پرت‌های شدید باعث فشرده شدن کل داده‌ها در نزدیکی ۰ یا ۱ می‌شوند.
- **Scaling: Standard** میانگین ≈ 0 و انحراف معیار $= 1$ برای همه ویژگی‌ها. اما به دلیل استفاده از میانگین و انحراف معیار، تأثیر پرت‌ها بسیار زیاد است و توزیع پس از نرمال‌سازی همچنان چپ کشیده باقی می‌ماند.
- **Scaling: Robust** از میانه و دامنه بین‌چارکی (IQR) استفاده می‌کند که مقاوم به پرت هستند. در باکس‌پلات، مشاهده می‌شود که پرت‌ها در این روش کمتر کشیده شده‌اند و توزیع متعادل‌تری دارند.

پ: کدام روش برای داده‌های دارای پرت بهتر است؟ چرا؟

Robust Scaling بهترین روش برای داده‌های دارای پرت است.

دلیل ریاضی

- **میانه (Median):** برخلاف میانگین، تحت تأثیر مقادیر بسیار بزرگ یا کوچک قرار نمی‌گیرد. اگر ۱۰٪ داده‌ها پرت باشند، میانه همچنان نماینده خوبی از مرکز داده‌هاست.
- **$IQR = Q_3 - Q_1$:** این معیار پراکندگی فقط به ۵۰٪ داده‌های مرکزی وابسته است و پرت‌های دور (بالای $Q_3 + 1.5 \times IQR$ یا زیر $Q_1 - 1.5 \times IQR$) در محاسبه آن دخیل نیستند.
- **نتیجه:** نرمال‌سازی با میانه و IQR مقیاس داده‌های اصلی را بدون تحریف توسط پرت‌ها حفظ می‌کند.

تأیید با نمونه از داده‌ها (ویژگی MedInc)

جدول ۲: مقایسه ۱۰ نمونه اول ویژگی MedInc پس از نرمال‌سازی

نمونه اصلی	MinMax[۰,۱]	MinMax[-۱,۱]	Standard	Robust
۱	۳۲۵۲.۸	۵۳۹۷.۰	۰.۷۹۳.۰	۳۴۴۸.۲
۲	۳۰۱۴.۸	۵۳۸۰.۰	۰.۷۶۱.۰	۱۸۶۷.۲
۳	۲۵۷۴.۷	۴۶۶۰.۰	۰.۶۷۹.۰	۷۰۷۷.۱
⋮	⋮	⋮	⋮	⋮
۱۰	۶۹۱۲.۳	۲۲۰۱.۰	۰.۹۴۵.۰	۰.۷۱۷.۰

- در MinMax، مقادیر بالا (مثل ۳۰۸) به نزدیکی ۵۰۰ فشرده می‌شوند، در حالی که بیشتر داده‌ها در بازه ۰ تا ۳۰۰ هستند □ از دست رفتن وضوح.
- در Standard، به دلیل وجود پرت‌های درآمد بالا، انحراف معیار بزرگ شده و مقادیر معمولی (مثل ۶۰۳) به نزدیکی صفر یا منفی می‌روند □ تحریف توزیع.
- در Robust، مقادیر منطقی (درآمد ۳ تا ۸) به بازه منطقی $[-1, 3]$ نگاشت می‌شوند و تفاوت نسبی بین نمونه‌ها حفظ می‌شود.

شواهد بصری از باکس‌پلات (شکل ۳)

- در روش‌های MinMax و Standard، پرت‌ها (نقاط سیاه) بسیار دور از جعبه هستند و توزیع نامتعادل است.
- در Robust، جعبه‌ها متعادل‌تر و پرت‌ها کمتر برجسته هستند □ نشان‌دهنده مدیریت بهتر نویز.

نتیجه‌گیری نهایی

Robust Scaling به دلیل استفاده از آماره‌های مقاوم (میان و IQR) بهترین عملکرد را در حضور داده‌های پرت دارد و توزیع نسبی داده‌ها را بدون تحریف حفظ می‌کند. این روش به‌ویژه برای مدل‌های حساس به مقیاس مانند شبکه‌های عصبی، SVM و k-NN توصیه می‌شود.