

بخش چهارم:

A: جدول داده‌های کیفی (ناقص، تکراری، اشتباه) حذف می‌کنند و

Data cleaning کیفیت داده را بالا می‌برد و خطاهای داده‌ها را کم می‌کند و نتایج قابل اعتمادتری می‌دهد.

B:

حذف می‌کند ردیف‌ها یا ستون‌ها با داده‌های خراب (در صورت کمبود).

تجزیه و تحلیل با میانگین، میان و یا مقدار ثابت (مثلاً) می‌کند.

بیشترین استفاده از مدل‌های ساده برای درس‌های استفاده می‌شود.

C: داده‌هایی که از بقیه خیلی متفاوت اند (مثلاً خیلی بزرگ یا خیلی کوچک)

راه تشخیص ۱. IQR داده خارج از ۱.۵ برابر IQR - ۲ - score ۲ = داده‌هایی که

بیشتر از ۳ انحراف معیار از میانگین دور هستند. ۳. رسم Box Plot یا scatter هم‌گام می‌کند

۱، ۱، ۱۲ داده‌ها را به فرصت مناسب برای مدل تبدیل می‌کنند. ۲. مقیاس‌ها را یکسان می‌کنند

۳. روابط غیر خطی را ساده‌تر می‌کنند تا مدل بهتر یاد بگیرد.

:E

Label Encoding (دسته‌ها را به اعداد تبدیل می‌کند) (مثلاً قرمز = ۰، آبی = ۱)

برای داده‌های ترتیبی خوب است.

one-hot Encoding (برای هر دسته یک ستون با بیتی می‌سازد به طوری که فقط یکی از بیت‌ها ۱ باشد)

(قرمز = [۰، ۱] و آبی = [۱، ۰]) برای داده‌های غیر ترتیبی مناسب است.

به طور خلاصه: Label ترتیب فرض می‌کند ولی one-hot نه.

F: اکتیویشن‌های غیر ضروری را حذف و پیچیدگی مدل را کم می‌کند. ۲: سرعت یادگیری را

بالا می‌برد و از overfitting جلوگیری می‌کند. ۳: دقت سنجش را با تغییر روش داده‌های مهم‌تر می‌کند.

G: ۱- با دستورات SQL مثل Distinct یا Group By.

۲- در پایتون با `drop_duplicates()`.

۳- ابتدا داده‌های تک‌ارزشی شناسایی می‌شوند و بر اساس ستون و یا لیست خاص حذف می‌شوند.

H: ۱- مدل را لیج می‌کند و آنتروپی اشتباه یاد می‌گیرد. ۲- زمان پردازش و مصرف منابع

را به دلیل زیاد می‌کند. ۳- دقت بیش‌تر را کم می‌کند چون داده بی‌ربط نیز این‌ها را می‌کند.

Subject : \_\_\_\_\_

Date : \_\_\_\_\_

I: ۱- با برگردن منطق (مثل میانگین یا مدل پیش بینی) داده را کامل می کند. ۲- مثل های

ML به داده کامل نیاز دارند، هیچ Imputation فشرده ای می شود. ۳- حذف داده کم شده همیشه

مطلوب نیست

۵:

۱- رسم نمودار، هسته گرام یا Plot (۵-۵) (اگر رنگ های با خط صاف باشد، نرمال است)

۲- آزمون آماری، مثل P-value Kolmogorov-Smirnov, Shapiro-Wilk یا به معنی نرمال است

۳- بررسی skewness و kurtosis: اگر نزدیک صفر باشد، داده نرمال.