

تمرین اول داده کاوی

با توجه به دیتاست مشخص شده برای گروه تان، مراحل و گام های زیر را برای انجام تمرین انجام دهید.

۱. بررسی اولیه داده ها

- بارگذاری دیتاست: بارگذاری داده ها و بررسی ساختار کلی دیتاست (شامل تعداد سطرها و ستون ها، نوع داده ها).
- آمار توصیفی: استخراج انواع مختلف معیارهای آمار توصیفی از هر ستون.
- شناسایی مشکلات داده ها: شناسایی داده های از دست رفته (missing values)، داده های پرت (outliers)، داده های نویزی و دیگر مشکلات.
- کیفیت داده ها: بررسی کیفیت داده ها از جنبه های مختلف نظیر دقت (accuracy)، کامل بودن (completeness)، انسجام (consistency) و اعتبار (Validity).
- نمایش داده ها: نمایش داده های اولیه با استفاده از نمودارهای مختلف مانند هیستوگرام، باکس پلات، اسکتر پلات و اسکتر پلات متقابل برای شناسایی بهتر ویژگی های داده ها.
- توزیع داده ها: بررسی توزیع داده ها با استفاده از ابزارهای معرفی شده نظیر هیستوگرام یا روش های دیگر بر اساس دانش پیشین.
- چک لیست مشکلات: تهیه چک لیستی از مشکلات شناسایی شده و آماده کردن پیشنهاد روش های مناسب برای بهبود هر یک از مشکلات

۲. پیش پردازش داده ها

- جایگزینی مقادیر گمشده: شناسایی و جایگزینی مقادیر گمشده با حداقل دو روش مختلف مانند استفاده از معیارهای توصیفی آماری و توزیع داده ها.
- حذف داده های نویزی و پرت: شناسایی و حذف داده های نویزی و پرت با حداقل دو روش مختلف.
- اصلاح داده های نادرست و ناسازگار: شناسایی و اصلاح داده های نادرست و ناسازگار با استفاده از روش های مناسب.
- استانداردسازی/نرمال سازی: ایجاد یک دیتابیس جدید از دیتابیس قبلی که داده کاملاً استانداردسازی (یا نرمال سازی) شده باشد.
- حذف داده های تکراری: شناسایی و حذف داده های کاملاً تکراری و شبه تکراری (با درصد شباهت زیاد).
- مقایسه توزیع داده ها: مقایسه توزیع داده ها حداقل بین دو ویژگی با استفاده از دیورژانس کولبک-لایبر.
- نمایش داده های پاکسازی شده: نمایش داده های پس از پیش پردازش با استفاده از نمودارهای جدید برای ارزیابی تأثیر پیش پردازش بر داده ها.

- مقایسه قبل و بعد از پیش‌پردازش: مقایسه نتایج پیش‌پردازش با نتایج اولیه و تحلیل بهبودهای حاصل از پیش‌پردازش.

۳. نتیجه گیری و جمع بندی

- ارائه آمار بهبود: ارائه آماری از میزان بهبود در دیتاست پس از انجام مراحل پیش‌پردازش.
- بررسی میزان وابستگی بین داده‌های دیتاست با ابزارهای معرفی شده
- تهیه یک گزارش جامع شامل مراحل انجام‌شده، تحلیل‌های آماری و بصری، مشکلات شناسایی‌شده و راهکارهای پیشنهادی.

۴. قوانین و مقررات انجام تمرین:

- پروژه تنها به صورت انفرادی قابل انجام بوده و در صورت وجود هرگونه تشابه بین دو کد یا عدم تسلط به روند، نمره منفی به دانشجویان داده میشود.
- کپی و استفاده کورکورانه از منابع اینترنتی، یا ابزارهایی نظیر Chat-GPT ممنوع است.
- استفاده از روش‌های دقیق‌تر و خلاقانه در تحلیل و پیش‌پردازش داده با نمره اضافه همراه است.
- دانشجویان باید نتایج پیش‌پردازش را به صورت تصویری (با نمودارها) مقایسه کنند تا تاثیر هر مرحله به وضوح نمایش داده شود.
- فایل نهایی پروژه را با فرمت StudentNumber_Name_G#.zip (در قسمت G# به جای # شماره گروه خود را بگذارید.) زیپ کرده و در VU بارگذاری کند.