

AI Descision Tree Project Report

ابتدا وبسایت Kaggle را برای پیدا کردن یک دیتاست مناسب جستجو کردم، بعد از جستجو، دیتاستی با عنوان "Flight Passenger Satisfaction" پیدا کردم.

با استفاده از Colab و github، دیتاست را به صورت یک DataFrame خواندم و آن را پاکسازی (clean) کردم. پس از پاکسازی، یک df بدون مقادیر غیر عدد صحیح (non-integer) داشتم و ستون های پیوسته (continuous columns) را دسته بندی کردم.

پس از پاکسازی df، داده ها را به سه بخش train، test، و valid تقسیم کردم و سپس برای هر کدام، ویژگی ها (X) و برچسب ها (y) را جدا نمودم.

اما دوباره به جستجو در Kaggle پرداختم و دیتاست دیگری برای پیش بینی دیابت بر اساس سن، BMI و غیره پیدا کردم. ابتدا آن را پاکسازی کردم. این دیتاست قبل از کلاس های عددی نگاشت شده بود، بنابر این نیاز به نگاشت دوباره نداشت. ستون های پیوسته را دسته بندی کردم.

برای پیش بینی دقیق تر، از نمودار دسته بندی BMI استفاده کردم.

پس از پاکسازی داده ها، کلاس مدل را برای ساخت مدل DecisionTree نوشتم.

پس از بررسی خروجی مدل (predictions)، متوجه شدم داده هایم بیش از حد ساده اند و نتوانستم دقت خوبی به دست بیاورم. (بود training و ۷۶٪ روی مجموعه validation بیشترین دقت فقط ۷۴٪ روی مجموعه)

تصمیم گرفتم دیتاستم را دوباره تغییر دهم، و پس از تغییر دیتاست، پاکسازی آن، و ساخت مدل، متوجه شدم دقت مدل بسیار بالاست. در نتیجه، به بررسی مجدد در Kaggle پرداختم و متوجه شدم این موضوع طبیعی است. در واقع، مدل من با دقت ۹۸٪ دچار underfitting بود.

پس از انجام grid search روی مدل، بهترین مدل را با دقت ۹۹.۶۰٪ پیدا کردم.

دلیل این دقت بسیار بالا این است که دیتاست من برای شناسایی نوع حمله یا درخواست عادی به یک وبسایت استفاده می شود.

ما پارامترهای زیادی داریم، و این پارامترها درخواست های عادی را از درخواست های حمله جدا می کنند.

به همین دلیل، درخواست های حمله همیشه نشانه هایی دارند که آن ها را قابل شناسایی می سازد.