

بسمه تعالی

گزارش تمرین درخت تصمیم

درس یادگیری ماشین دکتر سیدین

تهیه کنندگان:

حمیدرضا دشت آبادی ۹۸۲۳۰۳۲

سید جلال طباطبایی ۹۸۲۳۰۵۹

در این پروژه یک کلاس برای تعریف گره ها (Node) که شامل ویژگی و شاخه ها (childs) تعریف شده است. در ادامه یک کلاس دیگر برای تعریف درخت تصمیم که ریشه آن شیئی از جنس کلاس Node می باشد و هر درخت تصمیم یک ریشه یا root دارد. تابع fit در هر مرحله ویژگی root درخت تصمیم فعلی را با محاسبه بهترین ویژگی به لحاظ information gain تعیین می کند. بهترین ویژگی با تابع getbestfeat بدست می آید که از تابع getIG برای محاسبه informationgain بهره می برد که آن هم از تابع getentropy یا getgini به عنوان مقیاس اندازه گیری استفاده می کند. همچنین در کلاس Node یک متغیر به اسم branch_childs نیز برای تعریف شاخه و گره مربوط به آن تعریف شده است. لذا برای هر شاخه آن (key)، root درخت سطح بعدی نسبت داده می شود که خود از جنس Node است (شیء گرایی). تابع fit در ادامه یک object موقت از نوع tree با دیتای مربوط به آن branch و عمق منهای یک و child آن شاخه می سازد. در ادامه به صورت بازگشتی تابع fit روی درخت موقتی ایجاد شده فراخوانده می شود تا root آن را بیابد و درون child آن شاخه بریزد. در ادامه هنگامی که درخت ساخته شد می توان تابع predict را با پاس دادن داده بدون لیبل فراخوانی کرد که در شاخه های درخت جلو می رود (به صورت بازگشتی با ورودی های بخشی از داده ای که مربوط به آن شاخه است و child آن شاخه فراخوانی می شود) تا به انتهای درخت برسد (یعنی leaf). حال ویژگی آن را به y_predict متناظر با سطرهای باقی مانده داده تست، نسبت (assign) می دهیم. تابع accuracy نیز داده y-predict با y-test مقایسه می کند.

بخش الف : در این بخش درخت با عمق ۸ آموزش داده شد و روش entropy روی ۸۰ درصد داده ها به صورت تصادفی که حدود ۳۰ ثانیه طول کشید و صحت آن تقریباً برابر ۱۰۰ درصد به دست آمد

بخش ب : در این بخش ۸ حالت ممکن بررسی شد. در حالت ۵۰ درصد عمق ۶ و ۸ و روش های آنتروپی و جینی صحت بدست آمده تقریباً یکسان است. در حالت ۷۵ درصد عمق ۸ صحت داده ها با روش جینی ۱۰۰ درصد و آنتروپی برخلاف انتظار ۹۵ درصد به دست آمد و درخت overfit شده است.