چکیده

در این گزارش، ما با استفاده از دیتاست های در اختیار ، درخت تصمیم را می سازیم . مهم اینجا مرحله تمیز کردن دیتاست و پر کردن جا های خالی و در نهایت به تصویر کشیدن درخت است

۱ فاز ۱

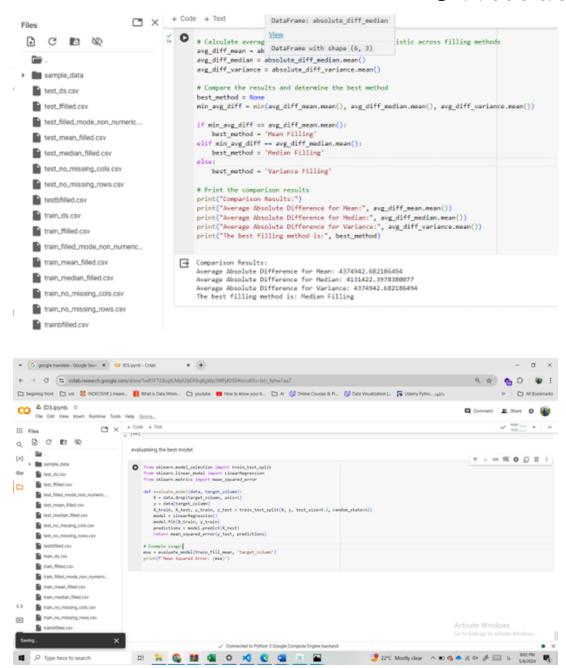
بهتر است در فاز ۱ برای پر کردن دادههای ناقص، از روشهای مختلفی استفاده کنیم:

- ۱. حذف دادههایی که مقادیر ناقص دارند: این روش ممکن است اطلاعات مهمی را از دست بدهد و باید با احتیاط استفاده شود.
- ۲. استفاده از میانگین یا مد: برای دادههای عددی، میتوانیم میانگین یا مد را به عنوان مقدار جایگزین استفاده کنیم. باید توجه داشت که این روش فقط برای دادههای عددی قابل استفاده است.
- ۳. استفاده از مد برای دادههای غیر عددی: در این حالت برای دادههایی که نوع عددی ندارند، از مد استفاده می کنیم. این نیاز به جدا کردن دادههای عددی و غیرعددی را دارد.
- ۱. استفاده از روش Backward | forward: این روشها بر اساس مقادیر مجاور یا پسین دادههای ناقص، مقادیر جایگزین را تخمین میزنند. در هر یک از این روشها، مدلهایی که برای تولید دادههای جایگزین استفاده می شوند، باید ذخیره شوند تا در صورت نیاز به آنها مراجعه شود.

٢ مقايسه الگوريتمها

در نهایت برای مقایسه این که کدام روش بهتر هست یک کد برای مقایسه نوشتم . هر خروجی که گرفتیم رو با دیتاست اصلی مقایسه میکنیم: تابع summarizedata را برای محاسبه آمار خلاصه (میانگین، میانه، واریانس) برای یک DataFrame اصلی تعریف می کند. محاسبه آمار اصلی بر روی دیتاست اصلی : آمار خلاصه (میانگین، میانه، واریانس) مجموعه داده اصلی را محاسبه می کند. summarizeFilledDatasets : تفاوت بین آمار هر مجموعه داده پر شده را با محاسبه آمار خلاصه آن خلاصه می کند. محاسبه تفاوت : تفاوت بین آمار خلاصه مجموعه داده اصلی و هر مجموعه داده پر شده را محاسبه می کند. محاسبه میانگین تفاوت : میانگین تفاوت را برای هر آمار (میانگین، میانه، واریانس) در تمام روش های پر کردن محاسبه می کند. تعیین بهترین روش: بهترین روش پر کردن را بر اساس کمترین میانگین تفاوت مطلق در تمام آمارها تعیین

می کند. نتایج مقایسه چاپ: نتایج مقایسه شامل میانگین تفاوت مطلق میانگین، میانه و واریانس و بهترین روش پر کردن را چاپ می کند.



٣ فاز ٢

۱-۳ توضیح کد

۱. تعریف کلاس گره: – در کلاس "Node" هر گره در درخت تصمیم را نشان می دهیم، که ویژگیهایی مانند ویژگی برای تقسیم (به عنوان "ویژگی")، مقدار آستانه برای ویژگیهای پیوسته (به عنوان

- "آستانه")، کلاس یا مقدار پیشبینی شده در یک گره برگ (به عنوان "مقدار")، و اشاره به گرههای فرزند چپ و راست (به عنوان "چپ" و "راست") را ذخیره میکند.
- ۲. توابع آنتروپی و IG: تابع آنتروپی یک مجموعه از برچسبها را محاسبه می کند که معیاری برای تصادفی بودن یا غیرقابل پیشبینی بودن دادهها است. آنتروپی کمتر به معنای عدم قطعیت کمتر است. تابع "informationgain": این تابع IG حاصل از تقسیم یک مجموعه داده به دو بخش را بر اساس یک ویژگی خاص محاسبه می کند. این اندازه گیری می کند که چقدر عدم اطمینان در برچسبها پس از تقسیم کاهش می یابد.
- ۳. عملکرد پیش پردازش داده: تابع "preprocessdata": این تابع مجموعه داده را برای مدلسازی آماده می کند. مقادیر از دست رفته در ستونهای عددی را با میانه آنها و در ستونهای طبقهبندی شده با حالت آنها پر می کند. همچنین متغیرهای طبقهبندی را با استفاده از "Label Encoder" در کدهای عددی رمزگذاری می کند و آنها را برای الگوریتم درخت تصمیم مناسب می سازد.
- ۴. بهترین عملکرد تقسیمبندی: تابع "findbestsplit": این تابع هر ویژگی در مجموعه داده را ارزیابی می کند تا مشخص کند کدام ویژگی و آستانه (برای ویژگیهای پیوسته) حداکثر IG را ارزیابی می دهد. روی تمام مقادیر ممکن (برای ویژگیهای طبقهبندی) تکرار می شود یا آستانههای بین مقادیر (برای ویژگیهای پیوسته) را محاسبه می کند تا بهترین تقسیم را پیدا کند.
- ۵. عملکرد ساختار درخت تصمیم: تابع " $build_decision_tree"$! این یک تابع بازگشتی است که درخت تصمیم را با انتخاب بهترین ویژگی برای تقسیم دادهها در هر گره (با استفاده از "findbestsplit") می سازد. برای هر نقطه تصمیم یک "گره" ایجاد می کند و تا زمانی که به حداکثر عمق مشخص شده برسد یا دسترسی به اطلاعات بیشتری امکان پذیر نباشد، به تقسیم شدن ادامه می دهد. برای گرههایی که همه نمونه ها به یک کلاس یا سایر موارد پایه تعلق دارند، یک گره برگ با رایج ترین بر چسب کلاس ایجاد می کند.

٣-٢ جواب سوالات

ویژگیهای پیوسته: کد ویژگیهای پیوسته را با محاسبه آستانه بهینه برای هر نقطه تقسیم بالقوه بر اساس IG کنترل می کند. این رویکرد به طور موثر داده های پیوسته را بر اساس اینکه آیا آنها زیر آستانه یا بالاتر از آستانه هستند، به دو دسته تقسیم می کند. – حداکثر مقادیر آنتروپی و gain: آنتروپی زمانی به حداکثر می رسد که داده ها بیشترین ترکیب را داشته باشند (همه کلاس ها به یک اندازه محتمل هستند)، و gain اطلاعات زمانی به حداکثر می رسد که یک تقسیم به طور کامل کلاس ها را به زیر مجموعه های خالص حدا کند.

1. ساختمان درختی بازگشتی (تابع builddecisiontree): - درخت تصمیم به صورت بازگشتی با انتخاب ویژگی و آستانه ای ساخته می شود که حداکثر اطلاعات را در هر گره به دست می آورد. این فرآیند بازگشتی از ریشه شروع می شود و تا زمانی ادامه می یابد که همه داده های یک گره برچسب کلاس یکسانی داشته باشند یا معیار توقف دیگری مانند رسیدن به حداکثر عمق ('madepth') برآورده شود. - در هر گره،

الگوریتم تمام ویژگی هایی را که هنوز در مسیر ریشه به گره استفاده نشده اند، بررسی می کند. این امر از استفاده مجدد از ویژگی ها در یک مسیر، جلوگیری از چرخه ها و اطمینان از استفاده از ویژگی های متنوع در قسمت های مختلف درخت جلوگیری می کند.

مدیریت ویژگی پیوسته، مقادیر منحصر به فرد مرتب شده و آستانه های بالقوه ارزیابی می شوند. این آستانه ها معمولاً به عنوان نقطه میانی بین مقادیر متوالی منحصر به فرد انتخاب می شوند. – این رویکرد گسستهسازی معمولاً به عنوان نقطه میانی بین مقادیر متوالی منحصر به فرد انتخاب می شوند. – این رویکرد گسستهسازی تضمین می کند که الگوریتم می تواند هر روش ممکن برای تقسیم داده ها به دو گروه را ارزیابی کند، که در آن یک گروه دارای مقادیر کمتر یا مساوی با آستانه و گروه دیگر دارای مقادیر بیشتر از آستانه است. – ارزیابی تقسیمات: – برای هر آستانه، مجموعه داده به دو زیر مجموعه تقسیم می شود. آنتروپی برای هر زیرمجموعه محاسبه شده و برای محاسبه سود کلی اطلاعات از ایجاد آن تقسیم استفاده می شود. – آستانه ای که منجر به بالاترین کسب اطلاعات می شود، به عنوان نقطه بهینه برای تقسیم داده های آن ویژگی در آن گره انتخاب می شود.

٣. حداكثر مقادير:

- حداکثر آنتروپی: - آنتروپی زمانی به حداکثر مقدار خود می رسد که داده های درون یک گره کاملاً مخلوط شوند، به این معنی که کلاس ها به یک اندازه محتمل هستند. برای طبقهبندی باینری، این حداکثر آنتروپی ۱ (بر حسب بیت) است که زمانی حاصل می شود که ۵۰ درصد داده ها به یک کلاس و ۵۰ درصد به کلاس دیگر تعلق داشته باشد.

حداکثر IG: IG زمانی به حداکثر می رسد که یک تقسیم منجر به کاهش قابل توجه آنتروپی شود. حداکثر IG زمانی اتفاق می افتد که زیر مجموعه های ایجاد شده توسط تقسیم کاملاً خالص باشند (یعنی تمام نمونه های هر زیر مجموعه به یک کلاس تعلق دارند). در این سناریو، IG اطلاعات برابر با آنتروپی مجموعه اولیه قبل از تقسیم است، زیرا آنتروپی هر زیر گروه پس از تقسیم \cdot است. در کد ، هر فراخوان بازگشتی به «builddecisiontree» همه ویژگی های استفاده نشده را در نظر می گیرد، همه تقسیم های ممکن (آستانه برای ویژگی های پیوسته) را ارزیابی می کند و IG اطلاعات را برای این تقسیم ها محاسبه می کند. IG و آستانه ای که IG را به حداکثر می رساند برای تقسیم داده ها در گره فعلی انتخاب می شود. اگر هیچ IG ای ممکن نباشد (یعنی حداکثر IG بی نهایت منفی باقی بماند)، گره به برگه ای با رایج ترین مقدار کلاس در میان نقاط داده باقیمانده تبدیل می شود.



۴ فاز ۳

وقوع مشکل: برازش بیش از حد در درختهای تصمیم زمانی اتفاق میافتد که مدل بیش از حد پیچیده میشود و شروع به گرفتن نویز در دادهها میکند نه فقط الگوی اصلی واقعی. این معمولا زمانی اتفاق می افتد که درخت اجازه دارد بدون محدودیت رشد کند تا زمانی که تمام نمونه های آموزشی را به طور کامل طبقه بندی کند. در اینجا دلیل این است که این می تواند به ویژه برای درخت های تصمیم مشکل ساز باشد:

۱. درختان عمیق: درختانی که در عمق رشد می کنند تمایل به یادگیری الگوهای بسیار نامنظم دارند
که تعمیم پذیری مدل را کاهش می دهد.

۲. اندازه برگ: اگر هر برگ درخت در نهایت تعداد بسیار کمی از نمونه های آموزشی را نشان دهد،
درخت بسیار مختص به داده های آموزشی می شود.

۳. تقسیمات پیچیده: داشتن تقسیمات زیاد ممکن است به این معنی باشد که مدل شروع به گرفتن نکات جزئی می کند، که فراتر از مجموعه داده آموزشی تعمیم نمی یابد.

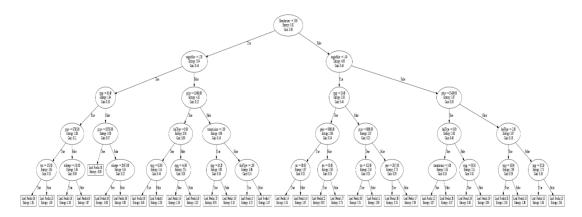
برای پرداختن به بیش از حد برازش در درختان تصمیم، چندین استراتژی را می توان به کار گرفت

- . **هرس**: پیش هرس (توقف زودهنگام): رشد درخت را قبل از طبقه بندی کامل داده های آموزشی متوقف کنید. این را می توان با تنظیم محدودیت هایی بر روی پارامترهایی مانند حداکثر عمق (' $min_samples_leaf$ ')، حداقل تعداد نمونه های مورد نیاز در یک گره برگ (' $min_samples_leaf$ ') یا حداقل افزایش مورد نیاز برای ایجاد یک تقسیم (' $min_impurity_aecrease'$) به دست آورد. پس از هرس: به درخت اجازه می دهد تا به عمق کامل خود رشد کند و سپس شاخه هایی را که قدرت کمی در پیش بینی متغیرهای هدف دارند حذف کند. این کار با ارزیابی بهبود خطای پیشبینی در زمانی که شاخهها هرس می شوند، انجام می شود.
- ۲. تنظیم محدودیت برای رشد درخت: محدود کردن حداکثر عمق درخت (max_depth'): یک راه ساده و موثر برای جلوگیری از پیچیده شدن بیش از حد درخت. درخت کم عمق کمتر رسا یک راه ساده و موثر برای جلوگیری از پیچیده شدن بیش از حد درخت. درخت کم عمق کمتر رسا است و در نتیجه کمتر به نویز تناسب دارد. افزایش $min_samples_split'$) MinimumSampleSplit است و حداقل تعداد نمونه هایی را که یک گره باید قبل از تقسیم داشته باشد را مشخص می کند. مقادیر بالاتر از یادگیری الگوهای بسیار ریز در مدل جلوگیری می کند، بنابراین قابل تعمیم تر است. افزایش حداقل نمونه ها در گره های برگ ($min_samples_leaf'$): این تضمین می کند که هر برگ بیش از تعداد معینی نمونه دارد، که با جلوگیری از تصمیم گیری هر برگ بر اساس نمونه های بسیار

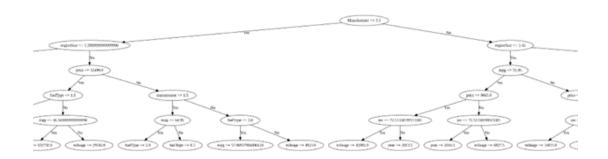
کمی، پیش بینی های مدل را هموارتر می کند.

• ۳. : - Bagging: ساخت چندین درخت به صورت موازی از نمونه های مختلف مجموعه داده آموزشی (نمونه های بوت استرپ) و میانگین گیری پیش بینی های آنها (مانند جنگل های تصادفی). این واریانس را کاهش می دهد و به جلوگیری از برازش بیش از حد کمک می کند. - تقویت: ساختن متوالی درختان، هر کدام بر پیش بینی صحیح مواردی که موارد قبلی بیشترین اشتباه را داشتند تمرکز می کنند. بنابراین، هر درخت به اصلاح پیشینیان خود کمک می کند، که به طور کلی منجر به مدل قوی تری می شود که بیشتر بر موارد سخت تر در مجموعه داده تمرکز می کند.

در نهایت در فاز τ ، هم با دیتاست تمیز شده برای train و هم از دیتاست test برای ساخت درخت استفاده کردم. در فاز τ با استفاده از کتابخانه گرافیکس خروجی ها را به صورت عکس ذخیره کردم.



test شکل ۱: برای دیتاست



شکل ۲: تکه از درخت تولید شده توسط روش mean