



گزارش تکلیف اول

داریوش حسن پور آده

۹۳۰۸۱۶۴

۱ قسمت اول تکلیف

- داده های جدول ۱ را به فرمت مناسب برای ورود به نرم افزار وکا آماده نمایید.

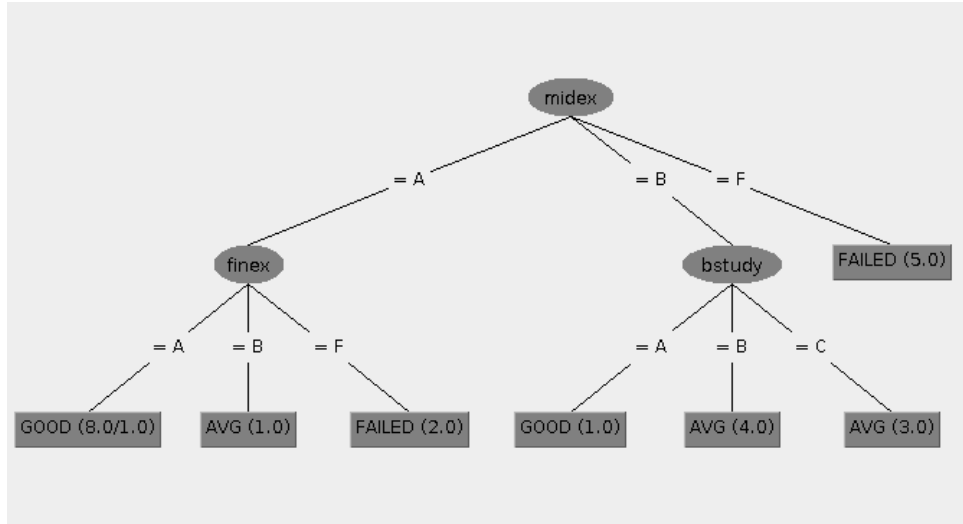
این داده ها را بصورت زیر به فرمت فایل قابل پردازش برای وکا بدست آورده شد - که به صورت پیوست در سامانه بارگذاری شده است. در داده های ایجاد شده اسامی ویژگی ها به صورت جدول ۱ آمده شده است.

نام فارسی ویژگی	نام انگلیسی ویژگی	نمایش «خوب»	نمایش «متوسط»	نمایش «ضعیف»	نمایش «زیر-۵»
حضور فعال در کلاس	apresent	A	B	C	-
مطالعه ی هفتگی کتاب	bstudy	A	B	C	-
مطالعه ی از روی جزوه	hstudy	A	B	C	-
میان نیمسال	midex	A	B	-	F
میان نیمسال	finex	A	B	-	F
تکالیف	asgnmnt	A	B	C	-
تحقیق	resrch	A	B	C	-
پروژه	project	A	B	C	-
نمره ی نهایی	final	GOOD	AVG	-	FAILED

جدول ۱: نمایش ویژگی ها در فایل به فرمت arff و معادل نمایش مقادیر خصیبه های آن ها با نمادین جدول ارائه شده.

- بدون هرس کردن درخت تصمیم گیری مرتبط با این داده های آموزشی را یاد بگیرید.

داده ها بعد از بارگذاری شدن در نرم افزار وکا با الگوریتم J.48 که همان معادل الگوریتم C4.5 می باشد به اجرا آورده شد، نتیجه ی درخت حاصله در شکل ۱ آمده است.



شکل ۱: درخت تصمیم ساخته شده از داده های جدول ۱ - بدون اعمال هرس

- درخت تصمیم گیری یادگرفته شده را بر روی دو نمونه آزمون جدول ۲ آزمایش نمایید.

نمره ی نهایی پیش بینی شده	project	resrch	asgnmnt	finex	midex	hstudy	bstudy	apresent
افتاده \equiv FAILED	خوب	خوب	خوب	خوب	زیر-۵۰	خوب	خوب	خوب
متوسط \equiv AVG	خوب	متوسط	خوب	متوسط	خوب	خوب	متوسط	متوسط

جدول ۲: نتایج آزمون داده های تست با درخت شکل ۱

- درخت را با استفاده از قوانین نمایش دهید.

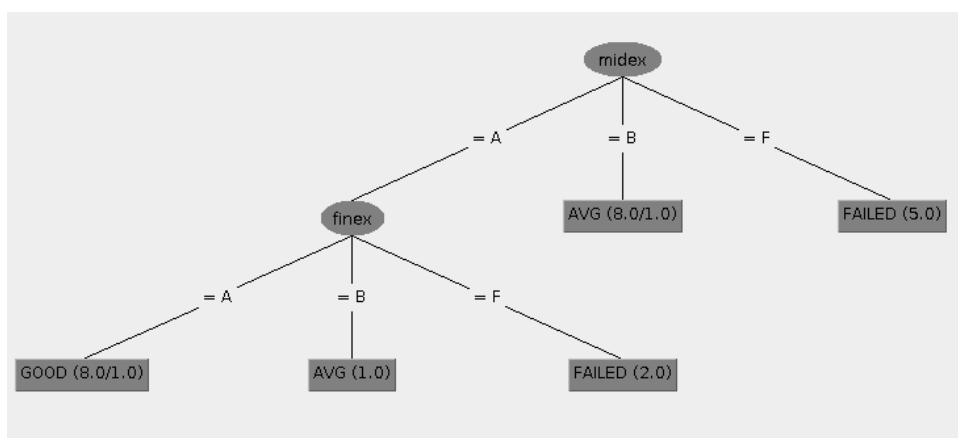
در جدول قوانین استخراج شده از درخت شکل ۱ در جدول ۳ به صورت اجزای مقدم و تالی مشخص گردیده اند که این قوانین را می توان به صورت مستقیم از درخت تصمیم استخراج کرد.

- درخت تصمیم گیری را با استفاده از هرس کردن بیاموزید و نتیجه آن را روی داده های آزمون بررسی کنید.

درخت تصمیم حاصل با اعمال هرس به صورت شکل ۲ بدست آمد. همانطور که مشاهده می شود گره «مطالعه ی هفتگی کتاب (bstudy)» حذف گردیده است. در این درخت نیز اگر داده ها تست را بیازماییم به همان نتایج جدول ۲ می رسیم که در جدول ۴ آمده است.

مقدم	تالی
میان نیمسال = زیر-۵۰	افتاده
میان نیمسال = متوسط و مطالعه‌ی هفتگی کتاب = خوب	خوب
میان نیمسال = متوسط و مطالعه‌ی هفتگی کتاب = متوسط	متوسط
میان نیمسال = متوسط و مطالعه‌ی هفتگی کتاب = ضعیف	متوسط
میان نیمسال = خوب و پایان نیمسال = خوب	خوب
میان نیمسال = خوب و پایان نیمسال = متوسط	متوسط
میان نیمسال = خوب و پایان نیمسال = زیر-۵۰	افتاده

جدول ۳: نمایش درخت تصمیم شکل ۱ به صورت قوانین



شکل ۲: درخت تصمیم ساخته شده از داده‌های جدول ۱ - با اعمال هرس

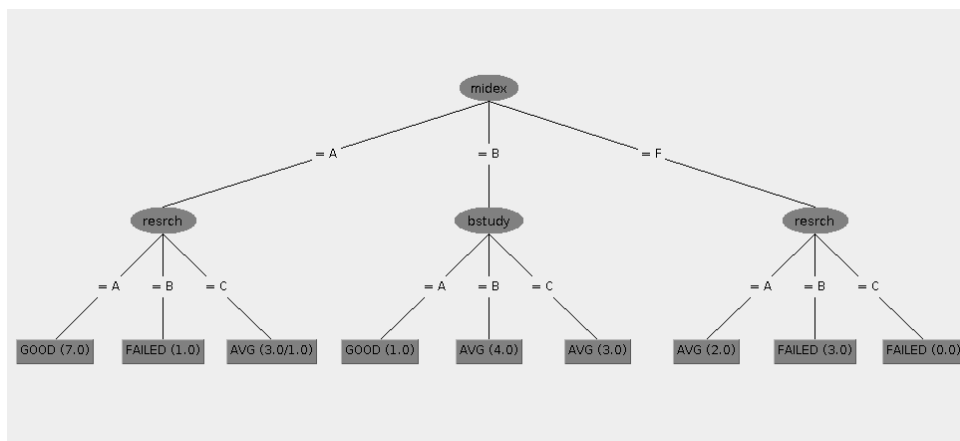
نمره‌ی نهایی پیش‌بینی شده	project	resrch	asgnmnt	finex	midex	hstudy	bstudy	apresent
افتاده \equiv FAILED	خوب	خوب	خوب	خوب	زیر-۵۰	خوب	خوب	خوب
متوسط \equiv AVG	خوب	متوسط	خوب	متوسط	خوب	خوب	متوسط	متوسط

جدول ۴: نتایج آزمون داده‌های تست با درخت شکل ۲

- خودتان با تغییراتی که لازم می‌دانید اثربیش پوشش را بررسی کرده و نشان دهید که چرا بیش پوشش انجام گرفته است و چگونه می‌توان جلوی آن را گرفت.

در درخت تصمیم یکی از مواقعی که بیش پوشش رخ می‌دهد داده‌های خطا دار به گونه‌ای باشند که در طی یادگیری متغیر هدف علاوه بر اینکه از فرضیه‌ی هدف منحرف شده‌ایم، عمق درخت نیز زیاد شود. به عنوان مثال اگر داده‌ها را بدین گونه تغییر دهیم که مقدار ۳ فرد اول که «افتاده» برچسب خورده‌اند را تغییر دهیم درخت حاصل از این تغییرات به صورت شکل ۳ بدست می‌آید. همان طور که در شکل ۳ مشاهده می‌شود شکل درخت فقط با دست کاری ۳ رکورد به کل تغییر پیدا کرد و چندین شاخه جدید بوجود آمد. و علت این بیش پوشش این است که چون داده‌ها دارای اغتشاش هستند مسیر رسیدن به فرضیه‌ی هدف گم می‌شود و الگوریتم سعی می‌کند به هر ترتیب که شده درختی با بهترین برازش برای این مجموعه داده بیابد بنابراین در نهایت به درختی می‌رسد که با داده‌های خطا دار حداکثر همخوانی را دارد که باعث ایجاد بیش پوشش شده است. برای هرس کردن چند روش وجود دارد که یکی از آن‌ها این است که یکی از گره‌ها را با زیردرخت آن را حذف کنیم و مقدار محتمل ترین برچسب در

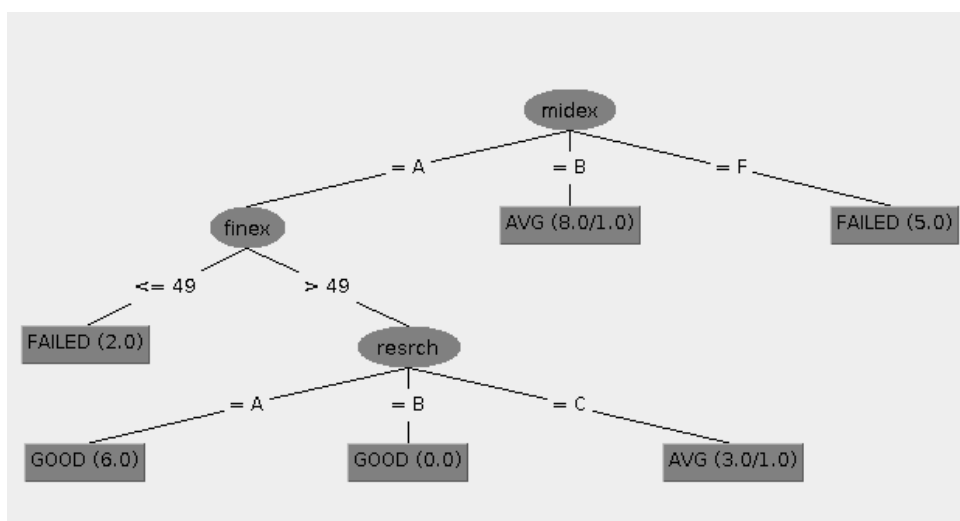
آن گره و زیردرختش را به عنوان برچسب گره در نظر بگیریم. و روش دیگر برای اجتناب از بیش‌پوشش حذف یکی از قوانین حاصله از درخت‌تصمیم که در واقع حذف یک مسیر از درخت‌تصمیم می‌باشد. در هر کدام از روش‌ها تا زمانی که میزان خطا افزایش پیدا نکرده است اقدام به حذف می‌کنیم و زمانی که خطای حاصله از آزمون درخت بعد از حذف گره یا مسیر افزایش پیدا کرد عمل هرس کردن را متوقف می‌کنیم.



شکل ۳: درخت تصمیم ساخته شده از داده‌های جدول ۱ - بدون اعمال هرس و با دست‌کاری در داده‌ها

- به دلخواه خود یکی از خصیصه‌ها را تبدیل به خصیصه پیوسته نمائید و دوباره درخت تصمیم‌گیری را یاد بگیرید.

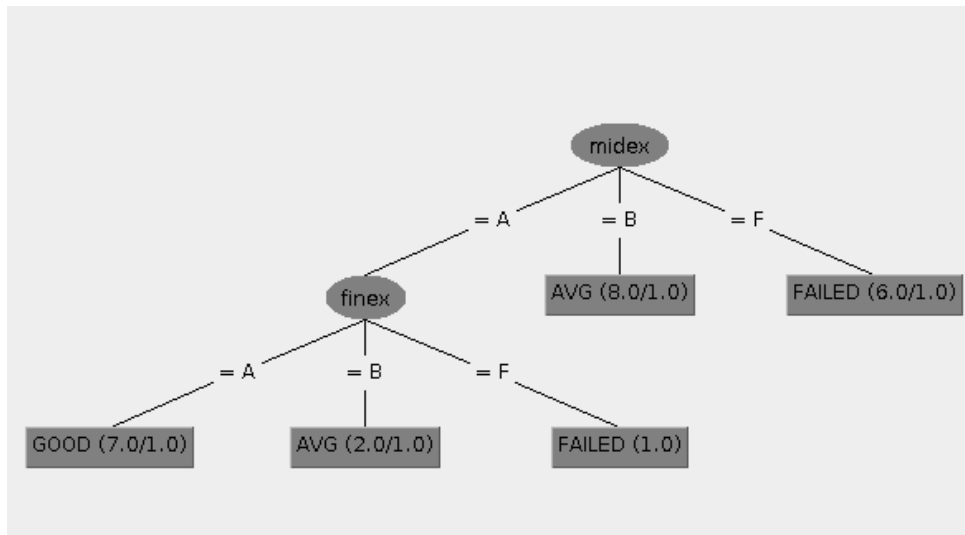
بنده ویژگی نمره‌ی پایان نمیسال رو پیوسته کردم و درخت بدست آمده بعد از پیوسته کردم این ویژگی در شکل ۴ آمده است، مقدار این ویژگی از ۰ الی ۱۰۰ براساس درجه‌ای که قبلاً داشته نمره گذاری شده است. همان‌طور که می‌بینیم الگوریتم بخوبی توانسته است که بر اساس این ویژگی پیوسته درخت را دوباره تشکیل دهد، که البته تا حدودی متفاوت‌تر از درخت‌های پیشین بدست آمد.



شکل ۴: درخت تصمیم ساخته شده از داده‌های جدول ۱ - با اعمال هرس و با ویژگی «پایان نمیسال» پیوسته

- اثر داده های خطادار و یا خصیصه های بدون مقدار را بررسی و گزارش کنید.

بنده مقدار خصیصه های «میان نیمسال : داده ی ۶ام» و «پایان نیمسال : داده ی ۹ام» را به ترتیب از $F \rightarrow B$ و $A \rightarrow F$ تغییر دادم و درخت حاصل از داده های مغشوش شده را در شکل ۵ آورده شده است. درخت بدست آمده با داده های خطادار مشابه درخت با داده های بدون خطا شکل ۲ می باشد و این یعنی الگوریتم توانسته با داده های خطادار بخوبی کنار بیاید و مفهوم هدف را دست ندهد ولی اگر به میزان خلوص گره های برگ توجه کنیم می بینیم میزان خلوص گره ها در درخت داده های خطا متفاوت تر از میزان خلوص با داده ها دست نخورده می باشد.



شکل ۵: درخت تصمیم ساخته شده از داده های جدول ۱ - با اعمال هرس و با داشتن داده های خطادار

۲ قسمت دوم تکلیف

- الف) چه تغییراتی الزم است که در داده های آموزشی داده شود تا بتوان فرضیه ای عطفی بوسیله الگوریتمهای Find-S یا حذف نامزد یادگرفت؟

الگوریتمهای Find-S یا حذف نامزد فقط توانایی تفکیک نمونه های مثبت و منفی را دارد در نتیجه چون خصیصه ی هدف داده ها دارای ۳ مقدار { خوب، متوسط و افتاده } می باشد این دو الگوریتم توانایی یادگیری براساس این ۳ خصیصه هدف ندارد پس ما برای اینکه بتوانیم این داده ها را با این دو الگوریتم آموزش دهیم باید نمونه ها را به صورت ۲ مقدار قبول شده (+) و افتاده (-) دسته بندی کنیم. بدین منظور در داده ها مقدار هر خصیصه ای که «افتاده» علامت گذاری نشده اند را + و مابقی را - علامت گذاری می کنیم.

- ب) پس از انجام تغییرات داده شده، الگوریتم Find-S چه فرضیه ای را یاد می گیرد؟

بعد از اعمال تغییرات لازم جهت یادگیری داده ها توسط الگوریتم Find-S داده هایی که با مقدار + علامت گذاری شده اند (کسانی که قبول شده اند) را به الگوریتم Find-S می دهیم که روند بدست آمدن فرضیه در جدول ۵ آمده است، همان طور که می بینیم الگوریتم Find-S نتوانست فرضیه ی درستی به دست بدهد و درست بعد از دریافت ۴ مثال + به عام ترین فرضیه

ممکن رسید که در این مثال خاص به معنی شکست کامل الگوریتم می‌باشد.

project	resrch	asgnmnt	finex	midex	hstudy	bstudy	apresent	فرد
∅	∅	∅	∅	∅	∅	∅	∅	–
خوب	خوب	خوب	خوب	خوب	خوب	خوب	خوب	۱
؟	؟	؟	خوب	خوب	خوب	خوب	خوب	۲
؟	؟	؟	؟	؟	؟	؟	خوب	۴
؟	؟	؟	؟	؟	؟	؟	؟	۶

جدول ۵: روند ایجاد خاص‌ترین فرضیه به ازای هر نمونه با مقدار هدف +

• (ج) این فرضیه داده‌های آزمون را چگونه دسته‌بندی می‌نماید؟

باتوجه به فرضیه‌ی بدست آمده که در سطر آخر جدول ۵ نشان داده شده است، فرضیه‌ی بدست آمده هردوی داده‌های تست را به عنوان مثال + دسته‌بندی خواهد کرد - که طبق دسته‌بندی‌ای که درخت تصمیم ارائه داد درست یکی از این نتایج درست نمی‌باشد.