

تمرین اول

پیاده سازی درخت تصمیم برای دسته بندی و تقریب تابع

استاد درس: دکتر آرش عبدی هجران دوست

تدریس یار آموزشی: ایمان براتی

نکات تمرین :

- ❖ مهلت تمویل ۱۴۰۱/۱۲/۱۹ ساعت ۲۳:۵۵
- ❖ مهلت ارسال به هیچ وجه قابل تغییر نیست .
- ❖ مواردی که بعد از تاریخ فوق ارسال شوند قابل قبول نبوده و نمره ای نخواهند داشت .
- ❖ انجام تمرین تک نفره است. لطفا به تنهایی انجام شود، در غیر اینصورت نمره منفی در نظر گرفته خواهد شد .
- ❖ کل مشوای ارسالی زیپ شود و نام فایل زیپ ارسالی HW1_StudentNumber_YourName باشد .
- ❖ مشوای ارسالی دارای راهنما (read me) جهت تسهیل اجرا باشد .
- ❖ زبان برنامه نویسی دلفواه است. (پیشنهاد : پایتون)
- ❖ در صورت استفاده از زبان پایتون فایل کد ترجیفا به فرمت ipynb بوده و فایل کد حتما به صورت اجرا شده آپلود گردد و از وجود فروبی سلول ها اطمینان حاصل نمایید .
- ❖ موارد ارسال شده در تاریخی که بعدا مشخص خواهد شد و متعاقبا اعلام می گردد به صورت آنلاین نیز تمویل گرفته خواهند شد (صرفا آنچه در LMS طبق تاریخ فوق تمویل داده شده است بعدا به صورت حضوری تست شده و توضیح داده می شود).
- ❖ تنها تکالیفی که به LMS و قبل از مهلت ارسال، فرستاده می شوند بررسی خواهند شد .
- ❖ در صورت داشتن هرگونه سوال می توانید سوال خود را در گروه تلگرامی درس مطرح کنید .
- ❖ مراقب یک ساعت قبل از مهلت ارسال را احتیاطا هدف قرار دهید، تا مشکلات غیرقابل پیش بینی مانند موارد زیر باعث عدم آپلود پاسخ ها در LMS و ارسال آنها از طریق ایمیل نشوند :

(قطعی اینترنت، تنظیم نبودن دقیق ساعت سایت با ساعت گرینویچ، کمرش سیستم عامل و نیاز به فرمت، بارش زیبای شهاب سنگ از آسمان و ...)

بخش اول (پیاپی سازی اولیه)

در ابتدا، می‌خواهیم رده‌بندی^۱ درخت تصمیم را از بیخ و بن (برون استفاده از کتابخانه آماده)، برای داده‌های گسسته، مطابق شبه‌کد ارائه شده در اسلایدهای کلاس، خودمان پیاده‌سازی نماییم.

function DECISION-TREE-LEARNING(*examples*, *attributes*, *parent_examples*) **returns**
a tree

```

if examples is empty then return PLURALITY-VALUE(parent_examples)
else if all examples have the same classification then return the classification
else if attributes is empty then return PLURALITY-VALUE(examples)
else
     $A \leftarrow \operatorname{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$ 
    tree  $\leftarrow$  a new decision tree with root test A
    for each value  $v_k$  of A do
        exs  $\leftarrow \{e : e \in \text{examples} \text{ and } e.A = v_k\}$ 
        subtree  $\leftarrow$  DECISION-TREE-LEARNING(exs, attributes – A, examples)
        add a branch to tree with label ( $A = v_k$ ) and subtree subtree
    return tree

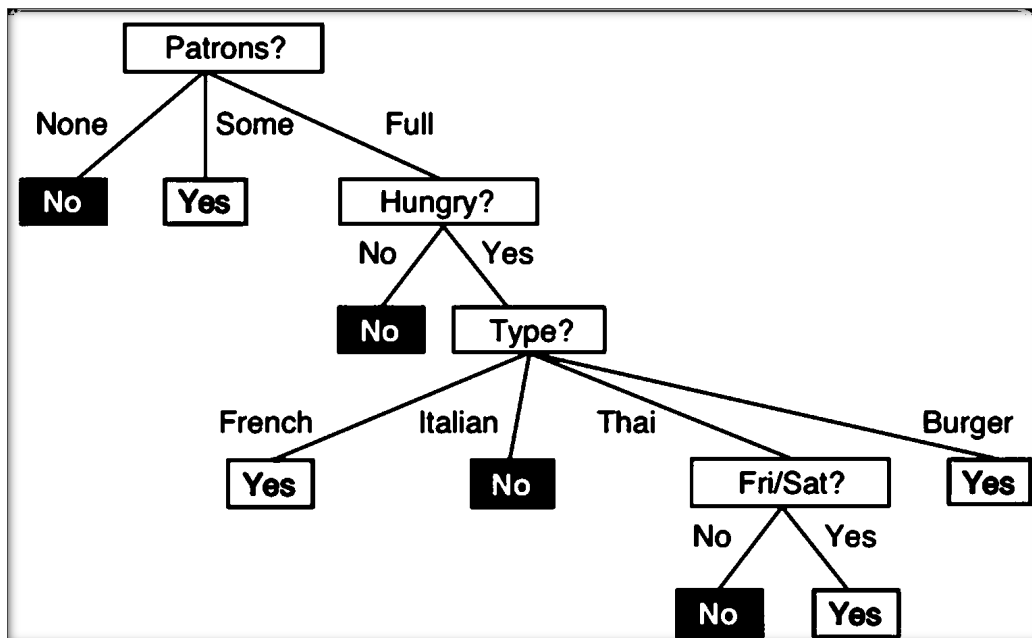
```

برای اطمینان از صحت پیاده‌سازی صورت گرفته، می‌توانید داده‌های ۱۲ تایی مثال رستوران (مطرح شده در کلاس درس) را مورد بررسی قرار دهید.

Example	Input Attributes										Goal
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
x₁	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0–10	$y_1 = \text{Yes}$
x₂	Yes	No	No	Yes	Full	\$	No	No	Thai	30–60	$y_2 = \text{No}$
x₃	No	Yes	No	No	Some	\$	No	No	Burger	0–10	$y_3 = \text{Yes}$
x₄	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10–30	$y_4 = \text{Yes}$
x₅	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	$y_5 = \text{No}$
x₆	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0–10	$y_6 = \text{Yes}$
x₇	No	Yes	No	No	None	\$	Yes	No	Burger	0–10	$y_7 = \text{No}$
x₈	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0–10	$y_8 = \text{Yes}$
x₉	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	$y_9 = \text{No}$
x₁₀	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10–30	$y_{10} = \text{No}$
x₁₁	No	No	No	No	None	\$	No	No	Thai	0–10	$y_{11} = \text{No}$
x₁₂	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30–60	$y_{12} = \text{Yes}$

¹ classifier

برای این کار تمام ۱۲ داده را به عنوان مجموعه آموزشی در نظر بگیرید (بدون مجموعه‌ی آزمایشی مجزا) و الگوریتم را برای این داده‌ها اجرا نمایید و سعی کنید سافت‌وار درخت آموزش دیده شده را نمایش دهید، خروجی صحیح مطابق تصویر زیر می‌باشد (البته ممکن است دو ویژگی آنتروپی یکسانی داشته باشند و سافت‌وار درخت شما در یک لایه متفاوت باشد).



بخش دوم (دسته‌بندی^۲)

در این مرحله می‌خواهیم با استفاده از رده‌بند درخت تصمیمی که در مرحله قبل پیاده‌سازی کرده‌اید، یک مسئله دسته‌بندی با داده‌های واقعی را حل کنیم. برای این کار از مجموعه داده‌گان این [لینک](#) می‌کنیم که در کنار این فایل با نام **part2.csv** نیز آمده است که یک مجموعه داده برای تشخیص تومور سرطانی با استفاده از تعدادی ویژگی می‌باشد.

در ابتدای کار می‌بایست داده‌ها را به دو بخش داده‌های آموزشی^۳ و داده‌های آزمایشی^۴ تقسیم کنید. نحوه بخش‌بندی داده‌ها به دو بخش آموزش و آزمایش به صورت کاملاً اختیاری و به دلخواه خودتان است (مثلاً ۹۰٪، ۱۰٪ - ۷۰٪، ۳۰٪ و ...)

همانطور که می‌دانید ورودی‌های درخت تصمیم باید به صورت گسسته باشد. برای گسسته‌سازی ورودی‌های از نوع پیوسته (مانند مجموعه داده فعلی) روش‌های مختلفی وجود دارد. ساده‌ترین ایده آن است که برای چنین ویژگی‌هایی، بازه مینیمم تا ماکزیمم اعداد در مجموعه آموزشی را به تعدادی بازه مساوی تقسیم کنید (چه تعداد؟ تعدادهای مختلف را آزمایش کنید) ایده‌های بهتر برای گسسته‌سازی مانند مرتب‌سازی و انتخاب نقاط برش در هر گره از درخت بر اساس نمونه‌هایی که در آن گره حاضرند را نیز امتحان کنید. همچنین می‌توانید

^۲ Classification

^۳ Training Set

^۴ Test Set

ایره های مطرح شده در کلاس یا ایره های جدید و فلاقانه خود نیز استفاده کنید و نتایج آن را با حالت های قبل (بازه های مساوی یا انتظاب نقاط برش بر حسب مرتب سازی) مقایسه کنید.

سپس با استفاده از الگوریتم نوشته شده درخت تصمیم در بخش قبلی آموزش مدل را بر روی داده های آموزشی انجام دهید (در نظر داشته باشید برای پیاده سازی درخت تصمیم نباید از توابع آماده استفاده کنید، لذا فرمول آنتروپی و ... را باید خودتان پیاده کنید. استفاده از توابع آماده برای قسمت های بعدی بلامانع است (و حتی توصیه میشود) مثلاً برای خواندن اکسل، نمایش گرافیکی فروبی درخت برای درک شهودی بهتر از فرآیند نحوه تصمیم گیری درخت تصمیم (که الزامی نیست)، نمایش دقت فروبی و ...

برای انتظاب بهترین ویژگی با توجه موارد تدریس شده از آنتروپی و **information gain** استفاده می شود، معیارهای دیگری مانند **gini index** نیز برای این کار وجود دارد. در ابتدا معیار **gini index** را تعریف کرده و چگونگی انتظاب بهترین ویژگی در الگوریتم درخت تصمیم را شرح دهید؛ سپس در کد نوشته شده این معیار را پیاده سازی کرده و آن را جایگزین **information gain** مبتنی بر آنتروپی کنید. درخت های حاصل بر اساس این دو معیار را با هم مقایسه کنید؛ دقت آنها، اندازه درخت، ترتیب ملاقات ویژگی های از ریشه تا برگ ها (مثلاً آیا ریشه درخت تغییر میکند؟ چند درصد گره های درخت متفاوت می شوند، و از این قبیل مقایسه ها).

سپس مجموعه داده گفته شده را با PCA^5 به ۱۰ بعد کاهش داده و با معیاری (**information gain** یا **gini index**) که در آموزش درخت تصمیم بر روی مجموعه داده با تمام ویژگی ها بهتر عمل کرده است، آموزش این مجموعه داده ی کاهش بعد داده شده را انجام دهید.

برای جلوگیری از بیش برازش^۶ از روش های توقف زودهنگام^۷ برای جلوگیری از برگ های با تعداد نمونه های بسیار کم (یا **gain** خیلی کوچک) قبل از سافت کامل درخت و روش های هرس کردن^۸ پس از آموزش کامل و سافت درخت، مورد استفاده قرار می گیرند. یک روش برای توقف زودهنگام و یک روش برای هرس کردن را پیاده سازی کنید و با مقادیر مختلف آزمایش کنید. ضمن توضیح مختصر نحوه پیاده سازی هر یک از روش ها، نتایج این دو روش را با یکدیگر مقایسه کنید و مزایا و معایب هر یک را بگویید؛

بعد از فرآیند آموزش درخت تصمیم، مقادیر زیر را برای داده های آموزشی و داده های آزمایشی همی حالت های گفته شده مناسبه کنید؛ (مورد ۶ و ۷ فقط برای داده های آزمایشی رسم و مناسبه شود)

⁵ Principal Component Analysis

⁶ Overfitting

⁷ early stopping or pre-pruning

⁸ Pruning or post-pruning

(۱) صحت^۹ (۲) دقت^{۱۰} (۳) فراخوانی^{۱۱} (۴) معیار^{۱۲} $f1$ (۵) ماتریس درهم‌ریختگی^{۱۳} (۶) نمودار^{۱۴} ROC (۷) مقدار^{۱۵} AUC

مقادیر ۱ تا ۵ با چند ضرب و تقسیم ساده به دست می‌آیند. مورد ۷ هم بر اساس مورد ۶ مناسبه میشود. برای نمودار ROC، لازم است یک پارامتر آزاد در نظر بگیریم و با تغییر آن، مقادیر مختلف برای TP و FP به دست آوریم (برای هر مقدار مشخص از پارامتر مربوطه، هر جفت این مقادیر مناسبه شده و یک نقطه در نمودار را تشکیل می‌دهد). پارامتر آزاد میتواند حد آستانه (threshold) برای دسته بندی باشد. فرض کنید دربرگها، به جای انتساب برپسب صفر یا یک (برای مساله دو کلاسه)، یک مقدار score برای کلاس مثبت در نظر گرفته شود، برین ترتیب که نسبت تعداد داده های با فروبی (برپسب) مثبت به کل داده ها در آن برگ را برابر با score مثبت بودن آن برگ در نظر بگیریم (عددی بین صفر (وقتی هیچ داده مثبتی نیست) و یک (وقتی همه داده های آن برگ مثبتند). حال با این درخت دارای Score در برگها، میتوان به هر داده تست، یک مقدار score نسبت داد (معادل احتمال مثبت بودن آن داده) و برای تصمیم گیری نهایی، با یک آستانه گذاری ساده، تصمیم نوعی مثبت یا منفی بودن داده گرفته شود. اینکه این آستانه چقدر باشد، میتواند پارامتری باشد که با تغییر آن، جفت ها FP, TP مناسبه شده و نمودار ROC ترسیم گردد. برای توضیح تکمیلی و نحوه پیاده سازی میتوانید لینک زیر را مطالعه بفرمایید که بیان دیگری از همین توضیحات است:

<https://stats.stackexchange.com/questions/105501/understanding-roc-curve/105577#105577>

دقت کنید که در آزمایشها یک مقدار حداقلی از مقادیر دقت مد نظر است (می‌توانید درخت تصمیم نوشته شده با کتابخانه آماده را به عنوان این مقدار مناسب در نظر بگیرید) چنانچه درخت شما مقادیر دقت قابل قبولی نداشت (در مقایسه با درخت آماده موجود در کتابخانه فیلی تفاوت زیاده بود) شما سعی کنید دقت را افزایش و نتایج را بهبود دهید. (با تغییر مجموعه هاپیر پارامترها چه در توزیع داده ها و چه در فود مدل و ...). در غیر این صورت منجر به کسر نمره فوادر شد.

بفش سوم (تقریب تابع)

در این بفش می‌خواهیم با استفاده از درخت تصمیم تقریب تابع (مسئله‌ی رگرسیون) انجام دهیم. برای این کار از مجموعه داده‌ی part3.csv موجود در کنار این فایل، استفاده نماییم که یک مجموعه داده برای تفمین قیمت یک ماشین با استفاده از تعدادی از ویژگی‌های آن ماشین می‌باشد. برای اطلاعات بیشتر درباره این مجموعه داده به این لینک مراجعه نمایید.

⁹ Accuracy Score

¹⁰ Precision

¹¹ Recall

¹² F1-Measure

¹³ Confusion Matrix

¹⁴ Receiver Operating Characteristic

¹⁵ Area Under Curve

در ابتدای کار نیاز است پس از تقسیم این مجموعه داده به مجموعه آموزشی و آزمایشی، یک پیش پردازش مناسب بر روی آن انجام داده و با ویژگی‌های آن بیشتر آشنا شوید و برای ویژگی‌های پیوسته، کسسته سازی به مانند بخش قبلی انجام دهید.

در کلام بعدی نیاز است تا درخت نوشته شده در بخش قبلی را با کمی تغییرات (که در کلاس درس به آن اشاره شده است)، به یک درخت تصمیم مسئله رگرسیون تغییر دهید.

در برگ‌ها نیز از روش میانگین‌گیری، میانه‌گیری و fit کردن تابع استفاده نمایید. (نتایج این سه روش را می‌خواهیم در ادامه با یکدیگر مقایسه نمایم)

می‌توان مجدداً از هرس کردن برای جلوگیری از بیش برازش (در صورت لزوم) استفاده کنید.

سعی کنید درخت تصمیم ساخته شده را رسم کنید (پیشنهادی) در غیر اینصورت سعی کنید سافت‌آر آن را متوجه شوید و به صورت کامل گزارش نمایید.

برای ارزیابی این مدل، معیارهای ارزیابی برای مسئله رگرسیون را مورد مطالعه قرار دهید و هم برای داده‌های آموزشی و هم برای داده‌های آزمایشی این معیارها را مناسبه و تحلیل نمایید.

آئپه تمویل داده می‌شود:

۱) کد اجرایی برنامه با توضیحات لازم برای اجرا

۲) درختی که برای مرحله دوم و سوم پیدا کرده اید (میتوانید گرافیکی نمایش دهید (به هر نوی که میتوانید) یا به صورت Text با پروتکلی که توضیح می‌دهید و قابل فهم باشد (بشود فهمید در هر کد کد نام ویژگی با چه مقداری فروبی تست شده اند و زیر شافه هایش کدامند و ...))

۳) گزارش کاملی از مسیر انجام کار، پالش‌هایی که مواجه شده‌اید، اجراهایی که گرفتید و نتایی که حاصل شده است. گزارش

کار از اهمیت بالایی برخوردار است، مهم آن و فرمت استاندارد آن اهمیت ندارد، اما باید نشان دهنده مسیر انجام

پروژه، پالشها، راه حلها و نتایج کار شما باشد.