

بسمه تعالی

جناب آقای دکتر فراهانی

نام و نام خانوادگی: سید مهرداد دستوری

شماره دانشجویی: ۹۹۴۲۲۰۸۳

تمرین سری ۳ واحد درسی داده کاوی

سوال ۱:

علت استفاده از کرنل ها :

علت استفاده از کرنل ها برای این است که داده ورودی را به یک داده دیگر و یک داده جدید در فضای جدید نگاشت کند، علت استفاده این است که چون svm در حالت کلی نمیتواند هر داده ای را جدا کند و بحث وکتورها برای همین از کرنل ها استفاده میشود تا آنها را به فضای جدید نگاشت کنیم و بتوانیم از آنها استفاده کنیم، پس در حالت کلی وظیفه کرنل این است که داده را به عنوان ورودی بگیرد و آن را به شکل مورد نظر تبدیل کند،

کاربردهای هر کرنلی عموماً با یک حالت و پس زمینه خاصی می آید، مثلاً کرنل چندجمله ای برای پردازش تصاویر استفاده میشود، کرنل RBF برای اهداف عمومی، کرنل تانژانت هیپربولیک و همچنین کرنل سیگموئید برای شبکه های عصبی استفاده میشوند.

برخی از کرنل های رایج مورد استفاده در SVM ها و کاربرد های آن ها :

۱- کرنل چند جمله ای

این کرنل در پردازش تصویر پرکاربرد است.

۲- کرنل گاوسی

این یک کرنل برای اهداف عمومی است. و هنگامی که هیچ دانش پیشینی در مورد داده ها وجود ندارد استفاده می شود.

۳- تابع پایه شعاعی گاوسی (RBF)

این کرنلی برای اهداف عمومی کاربرد دارد. و هنگامی که هیچ دانش پیشینی در مورد داده ها وجود نداشته باشد، مورد استفاده قرار می گیرد

۴- کرنل RBF لاپلاس

این هم یک کرنل برای اهداف عمومی است. و هنگامی که هیچ دانش پیشینی در مورد داده ها وجود ندارد استفاده می شود.

۵- کرنل تانژانت هیپربولیک (tanh)

می توانیم از آن در شبکه های عصبی استفاده کنیم.

۶- کرنل سیگموئید

می توان این کرنل را در شبکه های عصبی مورد استفاده قرار داد.

۷- کرنل تابع بسل (Bessel) از نوع اول

ما می توانیم از آن برای حذف مقطع عرضی در توابع ریاضی استفاده کنیم.

۸- کرنل پایه شعاعی ANOVA

ما می توانیم از آن در مسائل رگرسیون استفاده کنیم.

۹- کرنل spline خطی بصورت یک بعدی

این کرنل، هنگام کار با بردارهای بزرگ داده پراکنده، کاربرد زیادی دارد. این کرنل اغلب در دسته بندی متن مورد استفاده قرار می گیرد. کرنل spline همچنین در مسائل رگرسیون عملکرد خوبی دارد.

سوال ۲:

بر روی دیتاست کلاس بندی قیمت موبایل الگوریتم SVM را با تابع `svm_func_base` اجرا کردیم، نحوه عملکرد این تابع بدین گونه است که `svm` اجرا شده، روی داده ها `fit` شده و `predict` و در نهایت تابع `plot_results` برای نتایج فراخوانی میشود تا نتایج را نشان بدهد نتایج حاصل بعد از اجرا به این صورت میباشد که به صورت نمودار نمایش داده میشود و `confusion_matrix` نیز در خروجی چاپ میشود

`accuracy = ۰.۸۵۲`

`percision = ۰.۸۵۵`

`recall = ۰.۸۵۲`

`f_score = ۰.۸۵۳`

سوال ۳:

برای حالت های مختلف ما یک تابع `svm_func_۲` تعریف کردیم که در آن `kernel = linear` و خطی در نظر گرفته شد و `gamma = auto` مقداردهی شد که نتایج به صورت زیر است :

0.945accuracy =

0.944percision =

0.944recall =

0.944_score = 1f

حالت بعدی با تابع svm_func_۳ که در آن kernel خطی و gamma = scale در نظر گرفته شد
تعریف و اجرا کردیم

0.945accuracy =

0.944percision =

0.944recall =

0.944_score = 1f

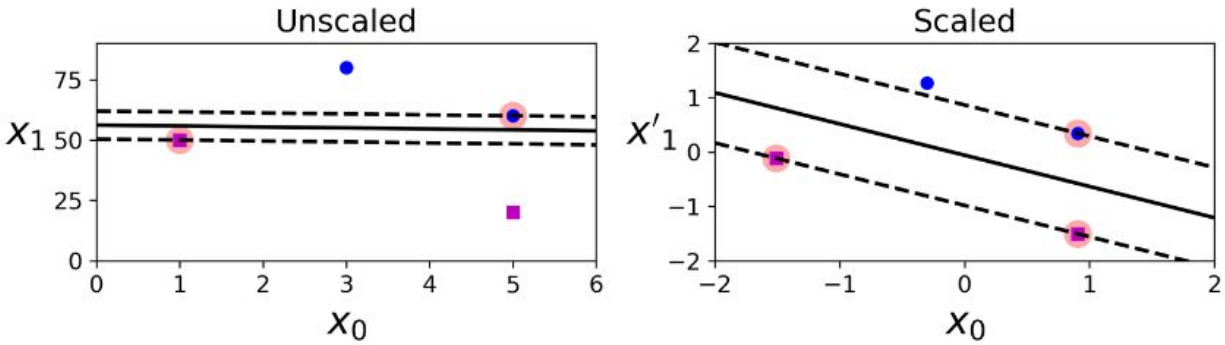
حالت بعدی با تابع svm_func_۴ که در آن kernel = poly و gamma = auto در نظر گرفته شد
تعریف و اجرا کردیم

حالت بعدی با تابع svm_func_۵ که در آن kernel = poly و gamma = scale در نظر گرفته شد
تعریف و اجرا کردیم

حالت بعدی با تابع svm_func_۶ که در آن پارامتر class_weight = balanced و
decision_function_shape = ovo در نظر گرفته شده تعریف و اجرا شد

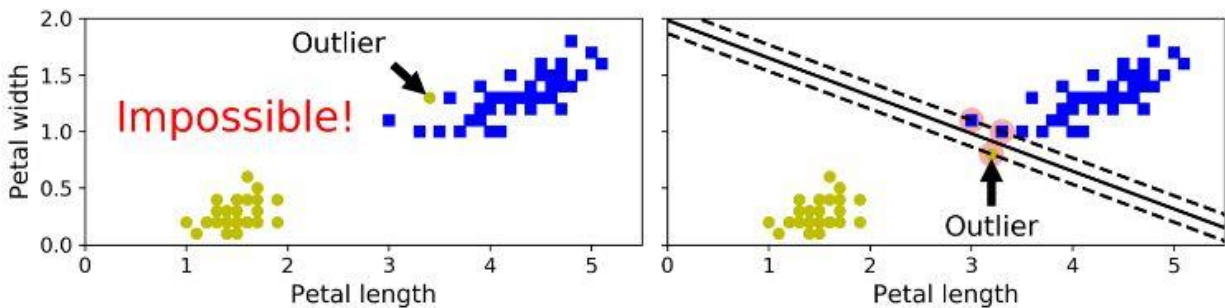
سوال ۴:

الگوریتم SVM به شدت به مقیاس داده ها حساس است. همون طور که در شکل سمت چپ می بینید
مقیاس عمودی بیشتر از افقی هست بخاطر همین فضای خالی نزدیک به افق هست. بعد از تغییر
مقیاس داده ها با استفاده از StandardScaler واقع در کتابخانه Scikit-Learn می بینیم که مرز های
تصمیم گیری شکل بهتری پیدا کردند.



طبقه بندی حاشیه نرم

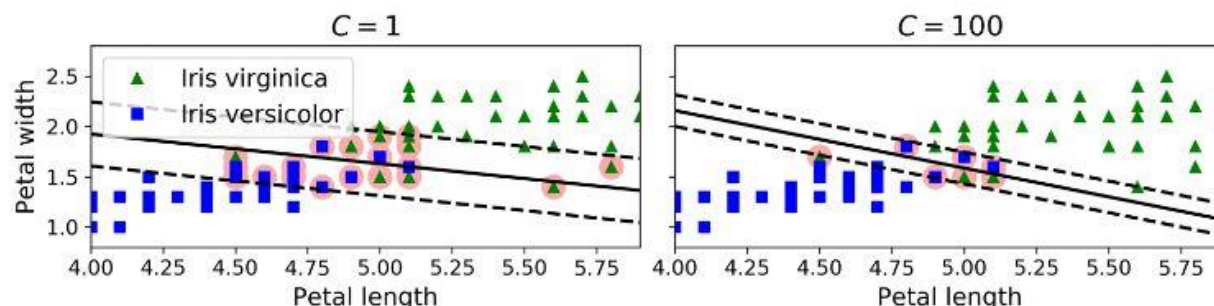
اگر تحمیل کنیم که تمام نمونه ها باید خارج از فضای خالی میانی باشند و حتما باید در سمت مناسب قرار بگیرند، در واقع روش **Hard Margin Classification** (طبقه بندی حاشیه سخت) نام دارد، این روش دو مشکل اساسی دارد، اول اینکه فقط برای داده‌هایی کار میکند که به صورت خطی جدا میشوند. دوم اینکه به داده‌های پرت به شدت حساس است. شکلی که می‌بینید دیتاست Iris رو نشان میدهد که در سمت چپ یک داده پرت وجود دارد و پیدا کردن یک **Hard Margin** برای آن غیر ممکن است. در سمت راست هم مرز تصمیم‌گیری با چیزی که در بالا دیدیم، خیلی متفاوت و احتمالا به خوبی نمیتواند عمومی سازی کند.



برای جلوگیری از این مشکلات، از یک مدل انعطاف پذیر تر استفاده می‌کنیم. هدف ما پیدا کردن یک تعادل مناسب بین بیشترین مقدار فضای خالی میانی و محدود کردن **Margin Violations** هست (نقض حاشیه یعنی قرار گرفتن نمونه‌ها در فضای میانی یا در سمت اشتباه) به این کار **Soft Margin Classification** می‌گویند

وقتی مدل‌های **SVM** رو با استفاده از **Scikit-Learn** می‌سازیم، یک سری هایپرپارامترها رو می‌توانیم تنظیم کنیم. یکی از این هایپرپارامترها **C** نام دارد. اگر مقدار **C** کم باشد، حاصل مدل

سمت چپ شکل زیر است، با مقدار زیاد به شکل سمت راست می‌رسیم. نقض حاشیه اتفاق خوبی نیست و بهتر است تا جایی که امکان دارد کمتر اتفاق بیفتد، اگرچه، در این مورد، درست است که مدل سمت چپ نقض حاشیه‌های زیادی دارد اما احتمالاً بهتر عمومی‌سازی می‌کند. در واقع هایپرپارامتر C مقدار مجازات مدل برای هر نمونه‌ای را که به صورت اشتباه طبقه‌بندی می‌کند تعیین میکند. هر چقدر C کمتر باشد، حاشیه نرم‌تر میشود



سوال ۴ :

ما برای این کار تابع `soft_margin_svm` را تعریف کردیم که مقدار $C = 0.1$ در نظر گرفتیم که بعد از فراخوانی نتیجه به صورت زیر حاصل شد :

`accuracy = 0.885`

`percision = 0.887`

`recall = 0.885`

`f_score = 0.885`

همچنین `hard_margin_svm` را تعریف کردیم مقدار $C = 1000$ برای آن در نظر گرفتیم که بعد از اجرا نتایج زیر حاصل شد :

`accuracy = 0.958`

`percision = 0.957`

`recall = 0.957`

$$0.957_score = 1f$$

پس هر چقدر مقدار c بالاتر باشد دقت بیشتر است

سوال ۵:

(آ)

بعد از انتخاب فیچر $battery_power$ مقدار $bin=10$ قرار دادیم و با دستور pd_cut این بین ها را اجرا کردیم نتیجه به این صورت شد که برای فیچر $battery_power$ تقسیم بندی بازه انجام شد و داده ها دسته بندی شدند بعد از نرمالایز کردن داده ها و فراخوانی svm بر روی آن نتایج زیر حاصل شد :

$$0.942accuracy =$$

$$0.942percision =$$

$$0.941recall =$$

$$0.941_score = 1f$$

یکبار دیگر با $bin=20$ اینکار را کردیم و نتایج زیر حاصل شد

$$0.920accuracy =$$

$$0.919percision =$$

$$0.919recall =$$

$$0.919_score = 1f$$

که در واقع دقت پایینتر شد

یکبار هم با $bin=30$ اینکار را انجام دادیم و نتایج زیر حاصل شد :

$$0.912accuracy =$$

0.912precision =

0.913recall =

0.912_score = 1f

که باز هم کمتر شد پس = 10bin نتایج بهتری دارد

(ج)

ما در کد از log transform استفاده کردیم و از دیتاها لگاریتم گرفتیم و نتایج به صورت زیر حاصل شد :

0.910accuracy =

0.909precision =

0.908recall =

0.908_score = 1f

داده هایی که نتیجه بی نهایت داشتند حذف شدند

علاوه بر log transform از تبدیل standard scaler transform نیز استفاده کردیم که نتایج به صورت زیر حاصل شد :

0.952accuracy =

0.952precision =

0.951recall =

0.951_score = 1f

که نتایج بهتری نسبت به log transform دارد

(د)

یک ویژگی به عنوان مساحت یا حجم گوشی ساختیم که `px_weight` و `px_height` را به آن دادیم و بعد از ضرب آن ها در هم یک ویژگی جدید مساحت ساخته شد که روی آن پیش بینی انجام میشود که نتایج به صورت زیر حاصل شد :

$$accuracy = 0.940$$

$$percision = 0.939$$

$$recall = 0.939$$

$$f_score = 0.939$$

سوال ۶:

الگوریتم `svm` را بر روی همه فیچرها اعمال کردیم و نتایج زیر حاصل شد :

$$accuracy = 0.825$$

$$percision = 0.827$$

$$recall = 0.823$$

$$f_score = 0.825$$

سوال ۷:

الگوریتم های مختلف درخت تصمیم:

الگوریتم ID۳ :

یکی از الگوریتم های بسیار ساده درخت تصمیم که در سال ۱۹۸۶ توسط Quinlan مطرح شده است. اطلاعات به دست آمده به عنوان معیار تفکیک به کار می رود. این الگوریتم هیچ فرایند هرس کردن را به کار نمی برد و مقادیر اسمی و مفقوده را مورد توجه قرار نمی دهد.

در این الگوریتم درخت تصمیم از بالا به پایین ساخته می شود. این الگوریتم با این سوال شروع می شود: کدام ویژگی باید در ریشه درخت مورد آزمایش، قرار بگیرد؟ برای یافتن جواب از معیار بهره اطلاعات استفاده می شود. با انتخاب این ویژگی، برای هر یک از مقادیر ممکن آن یک شاخه ایجاد شده و نمونه های آموزشی بر اساس ویژگی هر شاخه مرتب می شوند. سپس عملیات فوق برای نمونه های قرار گرفته در هر شاخه تکرار می شوند تا بهترین ویژگی برای گره بعدی انتخاب شود.

الگوریتم ۴.۵C :

این الگوریتم درخت تصمیم، تکامل یافته ID۳ است که در سال ۱۹۹۳ توسط Quinlan مطرح شده است که از معیار نسبت بهره (Gain ratio) استفاده می کند. الگوریتم هنگامی متوقف می شود که تعداد نمونه ها کمتر از مقدار مشخص شده ای باشد. این الگوریتم از تکنیک پس هرس استفاده می کند و همانند الگوریتم قبلی داده های عددی را نیز می پذیرد.

از نقاط ضعف الگوریتم ID۳ که در ۴.۵C رفع شده است می توان به موارد زیر اشاره کرد:

الگوریتم ۴.۵C می تواند مقادیر گسسته یا پیوسته را در ویژگی ها درک کند و الگوریتم ۴.۵C قادر است با وجود مقادیر گمشده نیز درخت تصمیم (decision tree) خود را بسازد، در حالی که الگوریتمی مانند ID۳ و بسیاری دیگر از الگوریتم های طبقه بندی نمی توانند با وجود مقادیر گمشده، مدل خود را بسازند. سومین موردی که باعث بهینه شدن الگوریتم ۴.۵C نسبت به ID۳ می شود، عملیات هرس کردن جهت جلوگیری از بیش برآزش می باشد. الگوریتم هایی مانند ID۳ به خاطر اینکه سعی دارند تا حد امکان شاخه و برگ داشته باشند (تا به نتیجه مورد نظر برسند) با احتمال بالاتری دارای پیچیدگی در ساخت مدل و این پیچیدگی در بسیاری از موارد الگوریتم را دچار بیش برآزش و خطای بالا می کند. اما با عملیات هرس کردن درخت که در الگوریتم ۵ انجام می شود، می توان

مدل را به یک نقطه بهینه رساند که زیاد پیچیده نباشد (و البته زیاد هم ساده نباشد) و بیش برآزش یا کم برآزش (Underfitting) رخ ندهد. الگوریتم $C4.5$ این قابلیت را دارد که وزن‌های مختلف و غیر یکسانی را به برخی از ویژگی‌ها بدهد.

الگوریتم CART :

برای برقراری درخت‌های رگرسیون و دسته‌بندی از این الگوریتم استفاده می‌شود. در سال ۱۹۸۴ توسط Breiman و همکارانش ارائه شده است. نکته حائز اهمیت این است که این الگوریتم درخت‌های باینری ایجاد می‌کند به طوری که از هر گره داخلی دو لبه از آن خارج می‌شود و درخت‌های بدست آمده توسط روش اثربخشی هزینه، هرس می‌شوند.

یکی از ویژگی‌های این الگوریتم، توانایی در تولید درخت‌های رگرسیون است. در این نوع از درخت‌ها برگ‌ها به جای کلاس مقدار واقعی را پیش‌بینی می‌کنند. الگوریتم برای تفکیک‌کننده‌ها، میزان مینیمم مربع خطا را جستجو می‌کند. در هر برگ، مقدار پیش‌بینی بر اساس میانگین خطای گره‌ها می‌باشد.

الگوریتم CHID :

این الگوریتم درخت تصمیم به جهت در نظر گرفتن مشخصه‌های اسمی در سال ۱۹۸۱ توسط Kass طراحی شده است. الگوریتم برای هر مشخصه ورودی یک جفت مقدار که حداقل تفاوت را با مشخصه هدف داشته باشد، پیدا می‌کند.

این الگوریتم با توجه به نوع برچسب کلاس از آزمون‌های مختلف آماری استفاده می‌کند. این الگوریتم هرگاه به حداکثر عمق تعریف شده‌ای برسد و یا تعداد نمونه‌ها در گره جاری از مقدار تعریف شده‌ای کمتر باشد، متوقف می‌شود. الگوریتم CHAID هیچگونه روش هرسی را اجرا نمی‌کند.

سوال ۸:

تابع `DT_func_base` را برای ساخت درخت تعریف کردیم که نتایج بعد از اجرا به صورت زیر حاصل شد:

0.840accuracy =

0.837percision =

0.839recall =

0.838_score = ۱f

سوال ۹:

تابع DT_func_۲ برای بررسی پارامترهای مختلف تعریف شد که در این تابع پارامتر $\text{criterion} = \text{gini}$ و پارامتر $\text{max_depth} = 5$ لحاظ شد و نتایج به صورت زیر حاصل شد :

0.828accuracy =

0.828percision =

0.825recall =

0.822_score = ۱f

تابع DT_func_۳ تعریف شد و پارامتر $\text{max_depth} = 20$ و پارامتر $\text{criterion} = \text{gini}$ قرار داده شد و نتایج به صورت زیر حاصل شد :

0.820accuracy =

0.819percision =

0.818recall =

0.818_score = ۱f

تابع DT_func_۴ تعریف شد و $\text{max_depth} = 50$ و پارامتر $\text{criterion} = \text{gini}$ قرار داده شد

تابع DT_func_۵ تعریف شد و پارامتر $\text{criterion} = \text{entropy}$ و $\text{max_depth} = 5$ قرار داده شد

تابع `DT_func_6` تعریف شد و پارامتر `criterion = entropy` و پارامتر `max_depth = 20` قرار داده شد

تابع `DT_func_7` تعریف شد و پارامتر `criterion = entropy` و `max_depth = 5` و پارامتر `min_samples_split = 5` قرار داده شد

تابع `DT_func_8` تعریف شد و پارامتر `criterion = entropy` و `max_depth = 5` و پارامتر `min_samples_split = 10` قرار داده شد

تابع `DT_func_9` تعریف شد و پارامتر `criterion = entropy` و `max_depth = 5` و پارامتر `max_features = 10` قرار داده شد

تابع `DT_func_10` تعریف شد و پارامتر `criterion = entropy` و `max_depth = 5` و پارامتر `max_features = 15` قرار داده شد و نتایج زیر حاصل شد :

`accuracy = 0.738`

`percision = 0.725`

`recall = 0.727`

`f_score = 0.725`

که دقت کمتر شده

سوال ۱۰:

هرس کردن درخت تصمیم

هرس کردن درخت تصمیم مقابل عمل تقسیم کردن است و با هرس کردن زیر گره هایی در درخت تصمیم حذف می گردد. زمانی که یک درخت تصمیم ساخته می شود، تعدادی از شاخه ها ناهنجاری هایی در داده های آموزش منعکس می کنند که ناشی از داده های پرت و یا نویز است.

در برخی الگوریتم های ایجاد درخت، هرس کردن جزئی از الگوریتم محسوب می شود. در حالی که در برخی دیگر، تنها برای رفع مشکل بیش برآزش از هرس کردن استفاده می شود.

چندین روش، معیارهای آماری را برای حذف کمتر شاخه های قابل اطمینان به کار می برند. درخت های هرس شده تمایل به کوچک تر بودن و پیچیدگی کم تر دارند و بنابراین به راحتی قابل فهم می باشند. آنها معمولاً در طبقه بندی صحیح داده های تست سریع تر و بهتر از درخت های هرس نشده عمل می کنند.

دو رویکرد رایج برای هرس درخت به شرح ذیل وجود دارد:

پیش هرس (Pre pruning) :

در این رویکرد یک درخت به وسیله توقف های مکرر در مراحل اولیه ساخت درخت، هرس می شود. به محض ایجاد یک توقف گره به برگ تبدیل می شود.

هرس پسین (Post pruning) :

رویکرد هرس پسین درخت تصمیم رایج تر است به این صورت که زیر درخت ها از یک درخت رشد یافته کامل را حذف می کند. یک زیر درخت در یک گره به وسیله حذف کردن شاخه ها و جایگزینی آنها با یک برگ، هرس می شود.

علت هرس کردن درخت : حذف داده های پرت یا ساده کردن مدل میباشد

سوال ۱۱:

نتیجه هرس کردن در کد

ما با معیار `ccp_alpha` در کد مشخص کردیم که چقدر هرس کردن اعمال شود، 0.01 یا 0.001 یا ... ، برای 0.01 نتیجه حدودا 0.75 میشود، با 0.001 که هرس کردیم دقت بالاتر رفت و نتیجه 0.82 را به ما داد پس تاثیرات مختلفی بر روی داده ها اعمال میشود

سوال ۱۲:

مجموعه ای از درختهای تصمیم `random forest` است که در انتها از آن درختها میانگین گیری میکند تا نتایج را پیش بینی کند، بعد از ساخت `random forest` و تخصیص مقادیر پارامترها نتایج مشخص میشود

سوال ۱۳:

از دیگر مزایای درخت تصمیم قابلیت تفسیرپذیری آن است. درخت های تصمیم ذات تفسیرپذیر و قابل درکی دارند که به استدلال آدمی شبیه است؛ .میتوان از درخت تشکیل شده یکسری قانون به صورت `if-then` استخراج کرد. از دیگر مزایای درخت های تصمیم میتوان به موارد زیر اشاره کرد:

- درخت تصمیم توانایی کار با داده های پیوسته و گسسته را دارد .
- درخت تصمیم از نواحی تصمیم گیری ساده استفاده می کند .
- مقایسه های غیرضروری در این ساختار حذف می شوند .
- از ویژگی های متفاوت برای نمونه های مختلف استفاده میشود .
- نیازی به تخمین تابع توزیع نیست .
- آماده سازی داده ها برای یک درخت تصمیم، ساده یا غیر ضروری است .
- درخت تصمیم یک مدل جعبه سفید است. توصیف شرایط در درختان تصمیم به آسانی و با منطق بولی امکان پذیر است. در حالی که شبکه های عصبی به دلیل پیچیدگی در توصیف

- نتایج آنها، مدل جعبه سیاه می باشند .
- تأیید یک مدل در درخت های تصمیم با استفاده از ارزیابی های آماری امکان پذیر است .
 - ساختارهای درخت تصمیم برای تحلیل داده های بزرگ در زمان کوتاه قدرتمند می باشند .
 - روابط غیرمنتظره یا نامعلوم را می یابند .
 - درختهای تصمیم قادر به شناسایی تفاوت های زیر گروه ها می باشند .
 - درختهای تصمیم قادر به سازگار کردن داده های فاقد مقدار هستند .
 - روش های درخت تصمیم به ویژه در آشکار کردن تراکنش های پیچیده بین متغیرها بسیار توانمند هستند. هر شاخه ای از درخت می تواند شامل ترکیبات مختلفی از متغیرها باشد و متغیرهای یکسان می توانند بیش از یک بار در قسمت های مختلف درخت ظاهر شوند. این امر می تواند مشخص کند که چگونه یک متغیر می تواند وابسته به متغیری دیگر باشد.

سوال ۱۴:

شماری از تکنیک های موفق و معروف یادگیری ماشین، نظیر طبقه بندی کننده های تقویتی، بردارهای پشتیبان و روش RIPPER در اصل برای مسائل دوکلاسه طراحی شده اند. البته لازم به ذکر است که روش RIPPER با هدف حل مسائل چندکلاسه تعریف شد اما این روش در واقع حاصل ترکیب دو روش AREP و IREP می باشد که هر دوی این روش ها در حوزه ی مسائل دوکلاسه تعریف شده اند. اما واقعیت این است که بسیاری از مسائل طبقه بندی در دنیای واقعی ابدا دوکلاسه نیستند بلکه متعلق به مسائل چندکلاسه می باشند. احتمال طبقه بندی نادرست در مسائل چندکلاسه بسیار بالاست و این احتمال با بالا رفتن تعداد کلاس ها، به سرعت افزایش می یابد، بنابراین واضح است که رسیدن به دقت بالا، در مسائل چندکلاسه، بسیار مشکل تر از مسائل دوکلاسه است.

تاکنون روش های موثر متفاوتی برای حل مسائل چندکلاسه ارائه شده است. این روش ها را می توان به دو دسته کلی تقسیم کرد:

روش های مبتنی بر تقسیم-و-تسخیر

روش‌های مبتنی بر تفکیک-و-تسخیر

دسته روش‌های تقسیم-و-تسخیر بر مبنای ایجاد درخت تصمیم‌گیری، بر روی داده‌های آموزشی استوار هستند. به این صورت که نقاط تصمیم‌گیری، در هر لایه از درخت تصمیم‌گیری، بر اساس آزمایش بر انواع ویژگی‌های داده، محاسبه می‌شود و داده‌ها بر اساس این آزمایش، طبقه‌بندی می‌شوند. این روند به صورت بازگشتی، آنقدر انجام می‌شود تا زمانی که دقت مورد نظر کسب شود. از جمله الگوریتم‌های متعلق به این دسته، می‌توان به ID۳ و CART و C۴.۵ و توسعه یافته آن C۵.۰ اشاره کرد که پیچیده‌ترین روش‌های القای درختی، در دو دهه اخیر هستند

در روش‌های مبتنی بر تفکیک-و-تسخیر، الگوریتم در هر زمان فقط بر یک کلاس تمرکز می‌کند. به بیان دیگر الگوریتم‌های این دسته سعی دارند که در هر زمان، قوانینی، را تولید کنند که بیشترین تعداد داده متعلق به کلاس جاری را به درستی پیدا کرده و تا حد ممکن داده‌های غیر از آن کلاس را نادیده بگیرند. پس از هر دور تولید قانون به روش مذکور، داده‌هایی که درست دسته‌بندی شده‌اند از داده‌های آموزشی حذف شده، سپس مابقی داده‌ها در دورهای آینده به صورت تکرار شونده، با روالی مشابه، مورد آموزش قرار می‌گیرند. الگوریتم‌هایی نظیر IREP و PRISM و RIPPER معروف‌ترین و موثرترین الگوریتم‌های این دسته می‌باشند.

سوال ۱۵:

بله میتوان از درخت تصمیم برای حل مسائل سری استفاده کرد ولی باید آنها را به دیتای سری زمانی تبدیل کرد که عموماً از Random forest برای این کار استفاده میکنند

سوال ۱۶:

برای تغییرات ستون date و حذف کاراکترهای اضافی از داده ها و تبدیل دیتاهای string به numeric باید بر روی آنها پیش پردازش انجام شود که ما اینکار را در کد با استفاده از تابع bit_data انجام دادیم

سوال ۱۷:

توابع مختلف تعریف شده :

Linear_regression_func

که ابتدا یک ابجکت از آن ها ساخته میشود، سپس رویدادهای **train** را فیت میکنیم ، رویدادهای **test** را پیش بینی میکنیم و مقادیر **mse** را پیش بینی میکنیم، بعد مقادیر **rmse** وقتی که **squared = false** باشد حساب میشود ، مقدار **mae** و **r** محاسبه میشوند، بعد مقدار دقت با تابع **get_acc** محاسبه شده و خروجی نهایتا در قالب اکسل ذخیره میشود

ridge_regression_func

lasso_regression_func

knn_regression_func

dt_regression_func

mlp_regression_func

random_frst_regression_func

bagging_regression_func

voting_regression_func

adaboost_regression_func

boosting_regression_func

نتایج بعد از اجرای این الگوریتم ها :

linear_regression_func

.....mse =

.....rmse =

0.007mae =

1.00 = 2r

99.794accuracy =

ridge_regression_func

0.000mse =

0.022rmse =

0.014mae =

0.99 = 2r

98.765accuracy =

dt_regression_func

0.059mse =

0.243rmse =

0.120mae =

0.19 = 2r

70.988accuracy =

votting_regression_func

0.15mse =

0.121rmse =

0.06mae =

0.8= 2r

72.634accuracy =

lasso_regression_func

0.155mse =

0.394rmse =

0.286mae =

1.11= - 2r

0accuracy =

mlp_regression_func

0.03mse =

0.058rmse =

0.029mae =

0.95= 2r

78.80accuracy =

random_frst_regression_func

0.051mse =

0.242rmse =

0.111mae =

0.19= 2r

71.399accuracy =

bagging_regression_func

0.051mse =

0.242rmse =

0.111mae =

0.19= 2r

71.399accuracy =

voting_regression_func

0.014mse =

0.120rmse =

0.059mae =

0.79= 2r

72.634accuracy =

adaboost_regression_func

$$0.060 \text{mse} =$$

$$0.246 \text{rmse} =$$

$$0.122 \text{mae} =$$

$$0.16 = 2r$$

$$71.399 \text{accuracy} =$$

ما بر روی داده های بیت کوین دو مدل یادگیری عمیق زدیم : CNN و LSTM

برای کدهای یادگیری عمیق از بستر **tensorflow** و کتابخانه **keras** استفاده کردیم ما ابتدا یک لایه **input** ساختیم، داده های **train** را به آن پاس دادیم، **shape** داده را مشخص کردیم که داده ها چند ویژگی دارند، بعد مدل یادگیری عمیق را ساختیم و لایه **input** را به آن پاس دادیم و بقیه لایه ها را به آن اضافه کردیم ، پارامترهای هر لایه را مشخص کرده و به آن ها اختصاص دادیم، مثلا برای لایه کانولوشنال از فیلتر ۱۶ استفاده کردیم، سپس لایه **D1MaxPool** و بعد لایه **flatten** را استفاده کردیم تا تبدیل به وکتورهای یک بعدی بشوند که بتوانیم لایه های **dense** را روی آنها اعمال کنیم ، بعد هم لایه های **dense** با نورون های مختلف استفاده کردیم و در انتها لایه خروجی با نام **out_layer** آورده شد که **activation** را برای آن **sigmoid** قرار دادیم چون رگرسیون داریم و تعداد نورون را ۱ قرار دادیم، بعد از اجرا یک مدل از **CNN** ساخته شد ، سپس مدل را کامپایل کردیم ، در کامپایل هم پارامترهای نوع **optimizer** و **loss function** و **metric** ها را مشخص کردیم و بعد از اجرا و کامپایل یک خروجی به صورت گرافیکی نمایش داده میشود که لایه ها و **shape** آنها و تعداد خروجی هایشان را شامل میشود

بعد از اینکار مدل را آموزش دادیم، برای اینکار ابتدا داده ها را **shape** میکنیم تا به صورت مناسب برای استفاده الگوریتم های یادگیری عمیق شود بعد مدل را **fit** میکنیم داده **train** و **label** را به آن پاس میدهیم و یک **epochs = 10** به آن پاس میدهیم که مشخص کننده تعداد اجرا میباشد و **batch_size** هم برای اینکه بداند چندتا چندتا داده ها را پاس بدهد ، بعد از اجرا الگوریتم شروع به

train کردن مدل میکند و مدل را آموزش میدهد تا نهایتا بتوانیم از مدل برای پیش بینی استفاده کنیم، بعد از اجرا میبینیم که مقدار **loss** و **rmse** تقریبا به صفر میرسد

بعد از آن برای تست مدل و پیش بینی دقت از تابع **get_acc** استفاده کردیم که بعد از فراخوانی **predict** نتایج زیر حاصل شد :

CNN_result

0.016mse =

0.127rmse =

0.075mae =

0.78= 2r

34.568accuracy =

بعد از آن یک مدل **LSTM** ساختیم ، مجددا مشابه مدل قبل لایه ورودی ، **input shape** و تعداد ویژگی ها را مشخص کردیم، لایه های بعدی و لایه آخر مجددا کامپایل و رسم شد ، سپس آموزش و **fit** کردن مدل را انجام دادیم و بعد با تابع **get_acc** تست را انجام دادیم و نتایج زیر حاصل شد :

LSTM_result

0.004mse =

0.060rmse =

0.028mae =

0.95= 2r

82.099accuracy =

ما برای این تست **100epochs** قرار دادیم که هر چه این عدد بالاتر باشد نتایج دقیقتر میشود و دقت بالاتری به دست می آید

سوال ۱۸:

استفاده از ایده **ensemble learning** بر روی ۵ تا از بهترین مدل ها و سپس اجزای تکنیک های **bagging** و **votting** و **boosting** نتایج به این صورت میباشد:

برای **bagging** از **mlp** به دلیل دقت بالا استفاده کردیم

برای **votting** از **linear_regression** و **ridge** استفاده کردیم و در **boosting** هم **GradientBoostingRegressor** استفاده شده

votting_model

$0.000 \text{mse} =$

$0.013 \text{rmse} =$

$0.008 \text{mae} =$

$1.00 = r^2$

$99.588 \text{accuracy} =$

که دقت بسیار بالایی دارد

bagging model

$0.006 \text{mse} =$

$0.074 \text{rmse} =$

$0.038 \text{mae} =$

$0.92 = r^2$

$75.309 \text{accuracy} =$

سوال ۱۹:

برای اجرای الگوریتم Adaboost بر روی دیتاست بیتکوین ما ۵ تابع مختلف با پارامترهای متفاوت با نام های `adaboost_۱` تا `adaboost_۵` تعریف کرده و اجرا کردیم که نتایج زیر حاصل شد:

`Adaboost_regression_result_on_bit`

`۰.۲۴۸mse =`

`۰.۱۵۷rmse =`

`۰.۷۸۶mae =`

`۰.۱۶= ۲r`

`۷۱.۴accuracy =`

`base_estimator_linear_regression`

`۰.۰۰۰mse =`

`۰.۸۰۰rmse =`

`۰.۴۰۰mae =`

`۰.۹۹= ۲r`

`۹۹.۸۰۰accuracy =`

`۲۰۰n_estimators_`

`۰.۶۱۰mse =`

`۰.۲۴۳rmse =`

`۰.۱۲۳mae =`

$$0.16 = 2r$$

$$71.4 \cdot \text{accuracy} =$$

$$\text{loss_square}$$

$$0.62 \cdot \text{mse} =$$

$$0.248 \cdot \text{rmse} =$$

$$0.124 \cdot \text{mae} =$$

$$0.16 = 2r$$

$$71.4 \cdot \text{accuracy} =$$

$$\text{learning_rate_}$$

$$0.247 \cdot \text{mse} =$$

$$0.157 \cdot \text{rmse} =$$

$$0.783 \cdot \text{mae} =$$

$$0.16 = 2r$$

$$71.4 \cdot \text{accuracy} =$$

$$0.001 \cdot \text{learning_rate_}$$

$$0.74 \cdot \text{mse} =$$

$$0.272 \cdot \text{rmse} =$$

$$0.144 \cdot \text{mae} =$$

$$r = 0.5$$

$$accuracy = 0.521$$

سوال ۲۰:

برای اجرای الگوریتم Random forest بر روی دیتاست بیتکوین ما ۵ تابع مختلف با پارامترهای متفاوت با نام های random_frst_۱ تا random_frst_۵ تعریف کرده و اجرا کردیم که نتایج زیر حاصل شد :

Random_forest_default_parameters

$$mse = 0.61$$

$$rmse = 0.247$$

$$mae = 0.123$$

$$r = 0.17$$

$$accuracy = 0.714$$

۲max_deep_

$$mse = 0.74$$

$$rmse = 0.271$$

$$mae = 0.142$$

$$r = 0.2$$

$$accuracy = 0.615$$

۹max_deep_

0.590mse =

0.243rmse =

0.119mae =

0.19= 2r

71.401accuracy =

max_features_sqrt

0.600mse =

0.246rmse =

0.121mae =

0.17= 2r

71.401accuracy =

300n_estimators_

0.590mse =

0.243rmse =

0.119mae =

0.19= 2r

71.401accuracy =

سوال ۲۱:

ما برای دسته بندی پیش بینی ها که اگر مقدار واقعی با مقدار پیش بینی شده ۰.۱ یا ۱۰ درصد اختلاف داشت ان را قبول کند و اگر نبود نه ، تابعی به نام `get_acc` تعریف کردیم که دو مقدار واقعی و پیش بینی شده را به عنوان ورودی میگیرد و دقت اختلاف را محاسبه میکند تا مقدار پیش بینی شده قبول یا رد شود

سوال ۲۲ :

برای اینکار ما **variation** ها را ساختیم به این صورت که اگر قیمت نسبت به روز قبل افزایشی باشد مقدار $\text{variation} = 1$ و اگر کاهش یافته باشد مقدار $\text{variation} = 0$ شود ، بعد `train` و `test` را جدا کردیم و مدل را روی داده `fit` کردیم تا نتایج حاصل شود

سوال ۲۳:

برای این کار ابتدا مسئله را به یک **Classification** تبدیل کردیم و اگر جدول داده ها را ببینیم یک ستون به نام **variation** اضافه کردیم که اگر قیمت افزایشی باشد ۱ و اگر کاهشی باشد ۰- مقدار میگیرد، در واقع از این فیچر جدید استفاده میکنیم تا نتایج را ببینیم ، بعد از نرمالسازی داده ها و اجرا و پاس کردن داده ها به **xgboost** نتایج حاصل شد و دقت بالاتر رفت

$\text{mse} = 0.000$

$\text{rmse} = 0.018$

$\text{mae} = 0.011$

$r^2 = 0.9$

$\text{accuracy} = 95.473$

سوال ۲۶:

حجم داده **tournament** خیلی زیاد بود و مجبور به اجرا در کولب شدیم، اولاً چون هر چقدر حجم داده زیاد باشد زمان فیت کردن مدل و یادگیری زیاد میشود، برای همین ما یک **chunk_size** در تابع **read_tournament_data** برای خواندن بخشی از داده ها تعریف کرده و بر روی آن اجرا کردیم، تابع **read_train_data** نیز برای خواندن داده های تَرین به کار بردیم، تابع **get_acc** برای محاسبه دقت و تابع **save_result** برای محاسبه نتایج به کار بردیم

ابتدا داده را خواندیم، **train** و **test** را جدا کردیم، یک **onehot encoding** بر روی داده ها اعمال کردیم به دلیل داشتن دو ستون غیر عددی تا به داده های عددی تبدیل شوند تا الگوریتم ها بتوانند بر روی آنها پردازش انجام دهند، سپس سه تا مدل مختلف **xgboost** و **Linear regression** و **ridge regression** را روی این داده ها فراخوانی کردیم، در نهایت مدل داده ها را فیت کرده و نتایج و خطای هر کدام را نشان میدهد و خروجی را ذخیره میکند

سوال ۲۷:

ما برای این کار از تکنیک **bagging** و با تابع **bagging_regression_func** استفاده کردیم، **votting** را دوبار فراخوانی کردیم یکبار با **linear_regression** یا **random_forest** و یکبار با **MLPRegressor** و **RandomForestRegressor** و همچنین **boosting** را دوبار فراخوانی کردیم یکبار با **GradientBoostingRegressor** و یکبار با **DecisionTreeRegressor** فراخوانی شدند و خروجی ها ذخیره شدند که بخشی از نتایج را مشاهده میکنید

votting

0.024mse =

0.155rmse =

0.112mae =

0.52= 2r

46.700accuracy =

votting

0.014mse =

0.119rmse =

0.090mae =

0.72= 2r

42.325accuracy =

boosting

0.048mse =

0.220rmse =

0.153mae =

0.03= 2r

49.920accuracy =

boosting

0.048mse =

0.219rmse =

0.156mae =

0.04= 2r

45.650accuracy =

سپس برای هر کدام **spearman_correlation** با استفاده از مقادیری که مدل های **ensemble** برگراندند محاسبه کردیم، برای این کار کتابخانه **spearman** را ایمپورت کردیم و مقادیر پیش بینی شده و لیبل های واقعی را به آن پاس دادیم تا آن را محاسبه کند، برای مشخص کردن کورولیت بودن سمپل ها از پارامتر $p = 0.05$ استفاده کردیم که اگر مقادیر بیشتر از این مقدار بود کورولیت نبوده و اگر کمتر بود کورولیت هستند که نتایج زیر حاصل شدند ، طبیعی است که مقادیر کورولیت هستند چون مقادیر پیش بینی شده به مقادیر واقعی نزدیک هستند

0.886votting by linear Spearman correlation coefficient :

0.000) p = 0Samples are correlated (reject H

0.871votting by mlp Spearman correlation coefficient :

0.000) p = 0Samples are correlated (reject H

0.235boosting by gradient descent Spearman correlation coefficient :

0.000) p = 0Samples are correlated (reject H

0.161boosting by dicision tree Spearman correlation coefficient :

0.000) p = 0Samples are correlated (reject H