

گزارش پروژه ۴

حسام رمضانیان

۸۱۰۱۰۰۲۴۸

سوال (۱)

درنوت بوک موجود است

سوال (۲)

درنوت بوک موجود است

سوال (۳)

باتوجه به نمودار هبستگی رسم شده در نوت بوک

MntCoffee,MntMeatProducts,Income, بیشترین همبستگی را دارند

سوال (۴)

در نوت بوک موجود است

سوال (۵)

نمودار ها در نوت بوک موجود است

سوال (۶)

فیچر هایی که مقادیر نامعتبر داشتند را با حذف ردیف ها اصلاح کردم و داده های پرت را از مجموعه داده ها حذف کردم در نوت بوک قابل مشاهده است

سوال ۷)

یکی از روش های مورد استفاده برای پر کردن داده های miss استفاده از یک مقدار ثابت میباشد که سریع و راحت است اما به علت اینکه این داده ها یکسان هستند میتواند یک bias بر روی داده ها ایجاد کند

روش دیگر استفاده از مقادیر ردیف های قبل و بعد است که با توجه به نوع داده های ما که وابستگی خاصی نسبت به هم ندارند این روش مناسب نیست

یک روش استفاده از linear interpolation است که بهتر از روش قبل عمل می کند اما مشابه روش قبل در صورتی که داده های مجاور وابستگی نداشته باشند این روش به درستی عمل نمیکند

سوال ۸)

Income = 223

Mntcoffee = 205

NumWebVisitsMonth = 200

MntGoldProds = 13

در نوت بوک miss ها برطرف شده است

سوال ۹)

Normalize : این روش مقادیر را به مقیاس ۰ تا ۱ در می آورد که باعث میشود مقیاس ها به یک

فرمت باشند تا در زمان train برای مدل اشکال ایجاد نشود

$$\text{Norm}(X) = (X - \min(X)) / (\max(X) - \min(X))$$

Standardizing : این روش مقیاس مقادیر را به گونه ای تغییر میدهد که میانگین ۰ و واریانس ۱

شود

$$\text{Standard}(X) = (X - \text{mean}(X)) / (\text{Std}(X))$$

علت استفاده از این روش ها این است که مدل های هوش مصنوعی برای train شدن و به دست آوردن

نتیجی بهتر نیاز دارند تا داده ها در یک مقیاس یکسان باشد تا روابط میان داده ها را به خوبی پیدا کنند

سوال (۱۰)

به طور معمول از دو روش برای تبدیل این داده ها استفاده میشود

(1) برای هر مقدار در فیچر یک فیچر جدید اضافه کنیم که مقادیر باینری ۰ و ۱ را میپذیرد

(2) آن ها را براساس یک ترتیب خاص متناظر با مقادیر عددی حسابی قرار دهیم

در صورتی که همه مقادیر موجود در فیچر معنی دار باشند و قصد حذف کردن فیچر را نداشته

باشیم برای تمامی مقادیر این پیش پردازش ضروری است

سوال (۱۱)

بله ستون هایی که همبستگی پایینی با ستون هدف دارند را میتوان از داده ها حذف کرد زیرا داده های مفیدی برای train کردن مدل نیستند

سوال (۱۲)

به طور معمول نسبت داده ها برای train بین ۷۰ تا ۸۰ درصد میباشد و سایر داده ها به منظور test مورد استفاده قرار میگیرند

یکی از روش ها به این صورت است که داده ها را به صورت رندوم بین train و test پخش کنیم روش دیگر به طور معمول برای داده هایی که بر اساس زمان مرتب شدن استفاده میشود به این صورت که بخش های ابتدایی را برای train و بخش های انتهایی را برای تست استفاده میکنیم

سوال (۱۳)

این دسته از داده ها برای ارزیابی هایپر پارامتر های مدل استفاده میشود که ایا مدل نسبت به قبل بهبود یافته یا خیر ، و بررسی اینکه ایا نیاز به توقف زود هنگام داریم یا نه با توجه با اینکه این داده ها متفاوت از داده های train و test است میتوان به خوبی با استفاده از ان ها این پیشرفت را بررسی کرد

سوال (۱۴)

در این روش داده ها را به k قسمت یا fold تقسیم میکنیم سپس k تکرار بر روی داده ها انجام میدهیم که در هر مرحله یک بخش به عنوان validation و سایر بخش ها را به عنوان train استفاده میکنیم این روش معمولاً برای دیتاست های کوچک استفاده میشود دقت شود که در هر تکرار یک بخش جدید به عنوان validation مورد استفاده قرار میگیرد پس از پایان این k تکرار میانگین ارزیابی این تکرار ها به عنوان نتیجه نهایی مورد استفاده قرار میگیرد

سوال ۱۵) در این بخش میخواهیم یک معادله خط پیدا کنیم که متغیر مستقل آن یکی از فیچرها و متغیر وابسته داده های ستون هدف باشد که بیشترین شباهت را با داده های ما داشته باشد به این منظور میخواهیم میزان خطا را کاهش دهیم پس باید معادله از RSS مشتق گرفته و با دو معادله ۲ مجهول به دست آمده ضرایب بهینه را به دست آوریم

سوال ۱۶)

MntCooffe بهترین فیچر است زیرا بیشترین همبستگی را با ستون هدف دارد. (در جدول رسم شده در نوت بوک قابل مشاهده است)

سوال ۱۷)

RSS: مربع خطا پیشبینی مدل نسبت به مقدار واقعی در ستون هدف میباشد که معیاری است که سعی در بهینه کردن آن داشتیم

$$RSS = \sum (y - \hat{y})^2$$

MSE: میانگین مقدار RSS میباشد (n: تعداد نمونه ها)

$$MSE = (\sum (y - \hat{y})^2) / n$$

RMSE: ریشه دوم MSE است که مقادیر کمتر بهترند و ۰ نشان دهنده عدم وجود خطاست.

$$RMSE = (\sum (y - \hat{y})^2 / n)^{(1/2)}$$

R2 Score: برخلاف RMSE که در مقادیر مثبت مرزی ندارد R2 score دارای مرز در مقادیر مثبت میباشد و بزرگتر از ۱ نمیشود اما در مقادیر منفی مرزی ندارد. R2 یا ضریب تعیین با مقایسه مجموع مربعات اختلاف بین مقادیر واقعی و پیش بینی شده متغیر وابسته با مجموع کل مربعات اختلاف بین مقادیر واقعی و میانگین متغیر وابسته محاسبه می شود.

$$R^2 = 1 - (SS_{res} / SS_{tot})$$

سوال ۱۸)

RMSE در صورتی که مدل عملکرد بهتری داشته باشد به صفر نزدیکتر است و R^2 به یک مقدار صفر R^2 یعنی پیشبینی تصادفی و مقادیر منفی به معنی پیشبینی بدتر از حالت تصادفی است برای ۴ فیچر که همبستگی بیشتری داشتند مدل را train کردم و مقدار این مقادیر در نوت بوک موجود است

سوال ۱۹)

طبق نتایج به دست آمده بیشترین دقت برای decisionTree پس از ان Logistic Regression و کمترین دقت برای k-nearest- Neighbours نتایج در نوت بوک موجود است

سوال ۲۰)

در نوت بوک موجود است

سوال ۲۱)

Overfitting: یعنی مدل بیش از حد به داده‌های آموزشی وابسته شده و نمی‌تواند روی داده‌های جدید عملکرد خوبی داشته باشد. این معمولاً زمانی رخ می‌دهد که مدل خیلی پیچیده باشد و سعی کند نویزها و جزئیات غیر ضروری داده‌های آموزشی را یاد بگیرد.

Underfitting: یعنی مدل نتوانسته روابط و الگوهای مهم در داده‌ها را یاد بگیرد. این معمولاً زمانی اتفاق می‌افتد که مدل خیلی ساده باشد و انعطاف‌پذیری کافی برای مدل‌سازی داده‌ها را نداشته باشد.

سوال ۲۲)

با حذف مرحله نرمالایز کردن از پیش پردازش نتیجه نهایی تمامی مدل‌ها دچار اندکی افت میشود

با حذف مرحله حذف داده های پرت نیز مدل نتیجه بدتری داشت زیرا این داده ها در مرحله train قدرت پیشبینی مدل را کاهش میدهند

با حذف مرحله حذف داده های بی معنی مثل kidhome با مقدار منفی نیز در صورتی که این داده ها در بخش test قرار بگیرند به شدت کاهش میابد

(سوال ۲۳)

در نوت بوک رسم شده است

(سوال ۲۴)

n_estimators : تعداد درخت هایی است که مدل train را با آن ها انجام میدهد که باعث طولانی شدن فرایند train میشود اما واریانس را کاهش میدهد

max_depth: عمق درخت هایی که ایجاد میشوند را مشخص میکند که با مقدار دهی مناسب ان میتوان از overfitting جلوگیری میکند

با توجه به نمودار های رسم شده در نوت بوک مقدار n_estimators تا حدود ۱۰۰-۲۰۰ باعث بهبود دقت مدل میشود اما پس از ان overfitting رخ میدهد

همچنین max_depth تا حدود عدد ۵-۱۰ باعث بهبود دقت مدل میشود اما پس از ان overfitting رخ میدهد

(سوال ۲۵)

bias: بررسی میکند که یک مدل چقدر به الگویی واقعی موجود در دیتا ها نزدیک شده است مقدار زیاد ان به معنی ساده سازی بیش از حد یا underfitting است

Variance: به معنی این است که مدل تا چه میزان به داده های train حساس شده است درواقع بالا بودن این مقدار به معنی overfitting میباشد

مدل Decision tree واریانس بالایی دارند و به راحتی دچار overfitting میشود در صورتی که در random forst با ایجاد درخت های متعدد سعی میشود تا بالانس خوبی میان بایاس و واریانس ایجاد شود

(سوال ۲۶)

اضافه کردن مقادیر کوچکی از نویز تصادفی به هر داده ی خاص می تواند شناسایی افراد را از روی داده ها دشوارتر کند، اما همچنان امکان train کردن مدل ها بر اساس این مجموعه ی داده ها ممکن است.

(سوال ۲۷)

نویز لاپلاس از توزیع لاپلاس جهت ایجاد نویز استفاده میکند که باعث میشود داده ها بیشتر تغییر کنند و حریم خصوصی بیشتر حفظ شود اما باعث میشود که داده دچار اشکال شوند در مقابل نویز نمایی از توزیع نمایی استفاده میکند که یعنی حریم خصوصی کمتر حفظ میشود اما داده ها سالم تر هستند

(سوال ۲۸)

از نویز لاپلاس استفاده شد که باعث مقداری کاهش در دقت مدل ها شد مقادیر در نوت بوک موجود است

(سوال ۲۹)

این یک روش یادگیری مجموعه ای است که پیش بینی های چندین مدل ضعیف را ترکیب می کند تا یک پیش بینی قوی تر تولید کند.

به این صورت عمل میکند که ابتدا یک مدل با یکی از الگوریتم ها مثل decisionTree میسازیم و آموزش میدهیم که به آن weak learner گفته میشود سپس خطای مدل را روی داده های train محاسبه میکنیم سپس یک weak learner جدید ایجاد میکنیم تا این خطا را پیشبینی کند و خطا مدل را بهبود ببخشد سپس مدل به مجموعه مدل ها اضافه میشود و نتیجه آن توسط میانگیری وزنی ترکیب میشود این کار را تا زمانی که به یک تعداد تکرار مشخص یا کم تر شدن تابع زیان از یک مقدار خاص ادامه میدهیم مدل نهایی میانگین وزنی تمام این weak learner ها را برمیگرداند

۳۰ سوال)

XGBoost مانند بوستینگ گرادیان از مجموعه ای از درختان تصمیم به عنوان مدل های پایه استفاده می کند. هر درخت بر روی باقی مانده ها یا خطاهای درخت های قبلی آموزش داده می شود. بنابراین اولین درخت بر روی داده های آموزش اصلی آموزش داده می شود. دومین درخت بر روی باقی مانده های درخت اول آموزش داده می شود و به همین ترتیب تا زمانی که متوقف شود ادامه دارد. هنگام آموزش یک درخت، XGBoost از الگوریتم greedy برای تقسیم داده ها در هر گره استفاده می کند و بهترین تقسیم که باعث بهبود حداکثری تابع هدف می شود را پیدا می کند. معیارهای متداول تقسیم شامل کاهش واریانس برای رگرسیون و افزایش اطلاعات برای طبقه بندی هستند. هر درخت از اشتباهات درخت های قبلاً آموزش داده شده می آموزد. پیش بینی های تمام درخت ها از طریق میانگین گیری وزنی با هم ترکیب می شوند تا پیش بینی نهایی انجام شود. بنابراین هر درخت متوالی سعی می کند درخت های قبلی را اصلاح و بهبود دهد. هر چه تعداد درخت ها بیشتر باشد، مدل دقیق تر می شود. XGBoost همچنین دارای پارامترهای منظم سازی است تا اورفیت شدن را با افزایش تعداد درخت ها کنترل کند.

سوال ۳۱) در نوت بوک موجود است