



مینی پروژه شماره یک

۱ پیش بینی بقا در کشتی تایتانیک^۱ با استفاده از رگرسیون لجستیک^۲

مقدمه^۳

کشتی تایتانیک در اولین سفر خود در تاریخ ۱۵ آوریل ۱۹۱۲ با برخورد به یک کوه یخی غرق شد. این حادثه، یکی از مرگبارترین فجایع دریایی تاریخ محسوب می‌شود که جان بیش از ۱۵۰۰ نفر را گرفت. بررسی داده‌های مسافران این کشتی نشان می‌دهد که احتمال زنده ماندن افراد به عوامل مختلفی مانند سن، جنسیت و کلاس سفر بستگی دارد. تحلیل این داده‌ها می‌تواند به ما کمک کند تا الگوهای مؤثر در بقا را شناسایی کنیم. یکی از روش‌های پرکاربرد در این زمینه، رگرسیون لجستیک است که به دلیل قابلیت آن در دسته‌بندی داده‌های دودویی^۴، روش مناسبی برای پیش‌بینی بقا یا فوت مسافران محسوب می‌شود. در این پروژه، از همین روش برای تخمین احتمال زنده ماندن هر مسافر بر اساس ویژگی‌های او استفاده خواهیم کرد.

بخش اول: آشنایی با مجموعه داده

مجموعه داده تایتانیک^۵ یکی از معروف‌ترین مجموعه داده‌ها در حوزه یادگیری ماشین^۶ است که اطلاعات مربوط به مسافران کشتی تایتانیک را شامل می‌شود. در این پروژه، قصد داریم با تحلیل ویژگی‌های مسافران، روابط بین متغیرها را بررسی کنیم و مدلی برای پیش‌بینی بقای مسافران بسازیم. این مجموعه داده شامل ویژگی‌های مختلفی از مسافران تایتانیک است که در عکس زیر قابل مشاهده است:

survival	Survival (0 = No; 1 = Yes)
pclass	Passenger Class (1 = 1st; 2 = 2nd; 3 = 3rd)
name	Name
sex	Sex
age	Age
sibsp	Number of Siblings/Spouses Aboard
parch	Number of Parents/Children Aboard
ticket	Ticket Number
fare	Passenger Fare
cabin	Cabin
embarked	Port of Embarkation (C = Cherbourg; Q = Queenstown; S = Southampton)

برای استفاده از این مجموعه داده‌ها، روش‌های مختلفی وجود دارد:

- می‌توانید آن را از طریق [این لینک](#) دریافت کنید،
- یا با استفاده از API Kaggle، مجموعه داده را مستقیماً در پروژه خود آپلود و استفاده نمایید.

^۱ Titanic

^۲ logistic regression

^۳ برای آشنایی بیشتر با logistic regression می‌توانید به [این لینک](#) مراجعه کنید.

^۴ binary

^۵ Titanic Dataset

^۶ machine learning

۱.۱ بررسی مجموعه داده:

وقتی یک پروژه یادگیری را شروع می‌کنیم، داده‌هایی که در ابتدا با آنها شروع می‌کنیم، داده‌های خام^۱ هستند. لذا نیاز داریم که آن‌ها را تجزیه و تحلیل کنیم و یک دید کلی نسبت به داده‌ها به دست آوریم و با ویژگی‌های آن‌ها آشنا شویم. به فاز اولیه تجزیه و تحلیل داده‌ها اصطلاحاً EDA^۲ می‌گویند. برای اجرا این فاز گام‌های زیر را انجام دهید:

- ساختار کلی داده‌ها را بدست آورید. (برای این کار می‌توانید از متدهای `info()` و `describe()` استفاده کنید.) مبهم
- نمودار همبستگی^۳ را رسم کنید و بررسی کنید که کدام ویژگی‌ها ارتباط قوی‌تری با متغیر هدف (Survived) دارند.
- نمودارهای پراکندگی^۴ و هگزین^۵ را برای بررسی رابطه بین ویژگی‌های عددی و متغیر هدف رسم کنید. توضیح دهید که این نمودارها چه اطلاعاتی را نمایش می‌دهند و چگونه می‌توان از آن‌ها برای تحلیل داده استفاده کرد.
- با استفاده از کتابخانه Plotly نمودار پراکندگی برای نمایش توزیع بازماندگان بر اساس سن و کرایه پرداختی رسم کنید. (محور افقی، سن و محور عمودی کرایه پرداختی است.) آیا افرادی که کرایه بیشتری پرداخته‌اند، شانس بقای بیشتری داشته‌اند؟ با استفاده از `countplot()` (Seaborn - بررسی کنید که چند درصد از مردان و چند درصد از زنان زنده ماندند. توزیع بازماندگان را بر اساس جنسیت مقایسه کنید (مثلاً چند درصد از زنان و چند درصد از مردان زنده ماندند).

۲.۱ تحلیل آماری مجموعه داده:

- تعداد اعضای خانواده‌ای که همراه هر مسافر سفر کرده‌اند را محاسبه کنید. بررسی کنید که آیا افراد دارای خانواده بزرگ‌تر یا کوچک‌تر، شانس بقای بیشتری داشته‌اند؟
- بررسی کنید که آیا یک مسافر به تنهایی سفر کرده است یا همراه با خانواده. تحلیل کنید که آیا تنهایی تأثیر منفی یا مثبتی بر نرخ بقا داشته است.
- مسافران را به گروه‌های سنی مختلف (مانند کودک، نوجوان، جوان، بزرگسال و سالمند) تقسیم کنید. بررسی کنید که آیا گروه‌های سنی مختلف، شانس بقای متفاوتی داشته‌اند.

بخش دوم: پیش پردازش مجموعه داده

مهم‌ترین فاز هر پروژه یادگیری ماشین، فاز پیش پردازش است. در این فاز فرمت داده‌ها را تغییر می‌دهیم، آن‌ها را اصلاح و خلاصه می‌کنیم، تا بتوانیم برای آموزش یک مدل یادگیری ماشین از آن استفاده کنیم. چرا که در دنیای واقعی، اطلاعات جمع آوری شده به راحتی کنترل نمی‌شوند و در نتیجه مقادیر خارج از محدوده، ناممکن، از دست رفته و به طور کلی گمراه کننده برای آموزش مدل در مجموعه داده‌ها وجود دارند. این بخش باعث می‌شود مدل کارا تری بتوانیم داشته باشیم و سرعت یادگیری بالاتر می‌رود.

- ممکن است برخی ستون‌های جدول دارای داده‌های از دست رفته باشند، تعداد و نسبت این داده‌ها را به دست بیاورید. روش‌های پر کردن داده‌های از دست رفته^۶ را توضیح دهید و حداقل سه روش را پیاده سازی کنید. دلیل استفاده از هر روش را مختصراً توضیح دهید.
- آیا امکان حذف برخی ستون‌ها وجود دارد؟ چرا؟ در صورتی که این امکان وجود دارد با ذکر دلیل ستون‌های لازم را حذف کنید.

^۱ raw data
^۲ exploratory data analysis
^۳ correlation matrix
^۴ scatter
^۵ hexbin
^۶ missing value

• کدام ویژگی‌ها را عددی و کدام‌ها را دسته‌ای می‌گویند؟ تفاوت این دو نوع از ویژگی‌ها در چیست؟ ویژگی‌های عددی^۱ و دسته‌ای^۲ را در این مجموعه دادگان مشخص کنید. برای ویژگی‌های دسته‌ای، که معمولاً بصورت یک string یا object در مجموعه داده ذخیره شده‌اند، در آموزش مدل چه پیش‌پردازش‌هایی مفید است؟ این موارد را در دیتاست اعمال کنید. در ویژگی‌های عددی نرمالایز کردن^۳ و استانداردسازی^۴ به چه منظور انجام می‌شود؟ تفاوت این دو روش در چیست؟ آیا در این پروژه نیاز به انجام این کار است؟

بخش سوم: انتخاب ویژگی^۵، آموزش^۶ و ارزیابی^۷

برای بهبود عملکرد مدل، ویژگی‌های جدیدی از داده‌های موجود استخراج می‌کنیم که می‌توانند به تشخیص الگوهای مؤثر بر بقا کمک کنند. همچنین در این مرحله، باید ویژگی‌های مهم را برای مدل یادگیری ماشین انتخاب کنیم تا دقت مدل افزایش یابد و از ویژگی‌های غیر ضروری اجتناب شود. پس از آموزش مدل‌های رگرسیون لجستیک، ضرایب مدل تحلیل شده و نسبت شانس^۸ برای هر ویژگی محاسبه می‌شود. این تحلیل به ما کمک می‌کند تا بفهمیم کدام ویژگی‌ها بیشترین تأثیر را بر احتمال بقا دارند.

• در مورد روش‌های مختلف انتخاب ویژگی تحقیق کنید. به طور خاص روش رگرسیون لاسو^۹ و حذف ویژگی‌ها به صورت بازگشتی^{۱۰} را پیاده‌سازی کنید (از توابع آماده موجود در کتابخانه‌ها می‌توانید استفاده کنید). و ویژگی‌های انتخاب شده توسط دو روش را مقایسه کنید. همچنین ویژگی‌های مورد نظر را انتخاب کنید (این ویژگی‌ها می‌تواند اجتماع ویژگی‌های دو روش باشند)

حال می‌خواهیم به کمک ویژگی‌های بدست آمده مدل‌های یادگیری ماشین برای پیش‌بینی بقای مسافران کشتی تایتانیک بسازیم. دو نوع رگرسیون لجستیک دودویی و رگرسیون لجستیک چندکلاسه را بررسی خواهیم کرد.

• مدل رگرسیون لجستیک دودویی را با ویژگی‌های منتخب آموزش دهید. بررسی کنید که دقت مدل روی داده‌های آموزشی و تستی چقدر است؟ ماتریس درهم‌ریختگی^{۱۱} را رسم کنید و تحلیل کنید که مدل کدام کلاس‌ها را بیشتر اشتباه پیش‌بینی کرده است؟ نمودار ROC^{۱۲} و مقدار AUC^{۱۳} را رسم کنید تا عملکرد مدل را ارزیابی کنید.

• دو مدل رگرسیون لجستیک چندکلاسه^{۱۴} و یکی در مقابل همه^{۱۵} را پیاده‌سازی کنید. (از توابع آماده موجود در کتابخانه sklearn می‌توانید استفاده کنید). برای این کار مسئله را اینگونه در نظر بگیرید:

0 → شانس بقا کم Low Chance

1 → شانس بقا متوسط Medium Chance

2 → شانس بقا زیاد High Chance

داده‌های بقا را بر اساس این سه کلاس دسته‌بندی کنید و بررسی کنید که چه تعداد نمونه در هر کلاس وجود دارد؟ دقت مدل را روی داده‌های آزمایش^{۱۶} بررسی کنید.

numerical^۱

categorical^۲

normalization^۳

standardization^۴

feature selection^۵

train^۶

evaluation^۷

odds ratio^۸

lasso regression^۹

recursive feature elimination^{۱۰}

confusion matrix^{۱۱}

receiver operating curve^{۱۲}

area under the ROC^{۱۳}

multinomial logistic regression^{۱۴}

one-vs-rest logistic regression^{۱۵}

test^{۱۶}

- ضرایب رگرسیون لجستیک دودویی و چندکلاسه را محاسبه کنید. لیستی از ویژگی‌های مدل همراه با ضرایب تخمین‌زده‌شده ارائه دهید. ویژگی‌هایی که ضرایب مثبت دارند، چگونه بر افزایش احتمال بقا تأثیر می‌گذارند؟ ویژگی‌هایی که ضرایب منفی دارند، چگونه احتمال بقا را کاهش می‌دهند؟ نسبت شانس را برای هر ویژگی محاسبه کنید. مقدار نسبت شانس را برای هر ویژگی تفسیر کنید.