

گزارش نهایی:

تحليل و پیش‌بینی بیماری قلبی بر اساس دیتاست Heart Disease Prediction

تهریه‌کننده: آرزو راستکار

تاریخ: ۱۳ دسامبر ۲۰۲۵

هدف پروژه: بررسی دیتاست ۲۷۰ نمونه‌ای بیماران قلبی، شناسایی عوامل خطر، و ساخت مدل classification برای پیش‌بینی Presence/Absence بیماری قلبی با دقت بالا (تقرباً ۸۵٪).

ابزارها: Python (Pandas, Scikit-learn, Seaborn, Matplotlib)

دیتاست: ۱۳ ویژگی دموگرافیک/بالينى

(Age, Sex, Chest pain type, BP, Cholesterol, FBS over 120, EKG results, Max HR, Exercise angina, ST depression, Slope of ST, Number of vessels fluro, Thallium) +
(Heart Disease: Absence/Presence). هدف.

۱. مقدمه و بررسی اولیه داده

دیتاست شامل ۲۷۰ بیمار (Presence ۴۴٪، Absence ۵۶٪)، اما با stratify هست – کمی imbalance ای نداره (تمیز). مدیریت شد. هیچ missing value ای نداره (تمیز).

آمار توصیفی کلیدی : (از (describe()

ویژگی	میانگین	انحراف استاندارد	حداقل	حداکثر
Age	۵۴/۴۳	۹/۱۱	۲۹	۷۷
Sex (۰=مرد، ۱=زن)	۰/۶۸	۰/۴۷	۰	۱
Chest pain type	۳/۱۶	۰/۹۲	۱	۴
BP	۱۳۰/۳۴	۱۷/۰۶	۹۴	۲۰۰
Cholesterol	۲۴۹/۶۶	۵۱/۸۳	۱۴۹	۵۶۴
Max HR	۱۵۰/۱۸	۲۲/۹۰	۷۱	۲۰۲
ST depression	۱/۰۵۰	۱/۱۶	۰	۶/۲
Thallium	۴/۷۰	۱/۹۴	۳	۷

مشاهدات: مردان (Sex = 1) بیشتر در گیرن؛ کلسترول بیشتر از ۲۵۰ در ۵۰٪ Presence دیده می‌شود.

۲. تحلیل اکتشافی (EDA)

-توزیع کلاس: (Absence ۱۵۰، Presence ۱۲۰) (56٪، 44٪).

-همبستگی با هدف (باينری): Thallium (۰,۵۲۵)، Max HR (۰,۴۱۹)، قوى مثبت – عامل خطر عروقی)، محافظه – ضربان پایین خطرناک).

- جدول همبستگی کلیدی:

ویژگی	همبستگی
Thallium	۰,۵۲۵
Number of vessels fluro	۰,۴۵۵
Exercise angina	۰,۴۱۹
ST depression	۰,۴۱۸
Chest pain type	۰,۴۱۷
Max HR	-۰,۴۱۹

Visualization :Heatmap-
multicollinearity همبستگی رسم شد – هیچ شدید (<0.8) دیده نشد.
-تفسیر: تمرکز روی Thallium و ST depression < 1 برای غربالگری.

۳. پیش‌پردازش

.(train Absence ۵۶٪ (تعادل stratify با ۵۴٪ Test، ۲۱۶ Train) ۲۰/۸۰ Split -
برای تعديل واریانس (مثل train روی StandardScaler : Scaling-
(Cholesterol

۴. مدل‌سازی و ارزیابی
دو مدل (Random Forest و RF)، با baseline : Logistic Regression
نتایج:

مدل	دقت (Test)	Recall Presence	F1 Presence
LR	۸۵/۱۹٪	۹۲٪	۰/۸۵
RF	۸۱/۴۸٪	۸۳٪	۰/۸۰

-اهمیت ویژگی‌ها (از RF از):

رتبه	ویژگی	اهمیت
۱	Chest pain type	۰/۱۳۰
۲	Max HR	۰/۱۲۰
۳	ST depression	۰/۱۱۰
۴	Thallium	۰/۱۰۰

-تفسیر: LR برتره (ساده، interpretable) برای feature selection عالی. هر دو overfitting ندارن.

۵. مدل Tuning

.regularization = {'C': 0.1, 'penalty': 'l2'}: بهترین LR برای GridSearchCV -

-دقت CV-(stable) ۸۴,۲۶٪

-دقت Tuned Test ۸۵,۱۹٪ (بدون افت)

-جدول CV Scores

C	l2 دقت CV
0.1	۸۴,۲۶٪
1	۸۲,۸۸٪
10	۸۱,۹۶٪

۶. نتیجه‌گیری و توصیه‌ها

-بهترین مدل: Tuned LR (Recall Presence ۹۲٪، دقت ۸۵٪) - مناسب بیمارستان برای غربالگری.

-عوامل کلیدی خطر: ST depression = Chest pain type Thallium بالا، (شدید).

-محدودیت‌ها: دیتاست کوچک (۲۷۰) - نیاز به داده بیشتر؛ SMOTE کم، اما SMOTE می‌توانه کمک کنه.

-توصیه‌ها:

- تمرکز روی بیماران < ۵۵ سال با Max HR < ۱۴۰ و کلسترول < ۳۰۰

تحلیل جامع بر اساس نمودارهای رسم شده: پیش‌بینی و عوامل خطر بیماری قلبی

در این تحلیل، با استفاده از دیتاست Heart Disease Prediction شامل ۲۷۰ بیمار و ۱۴ ویژگی دموگرافیک و بالینی، روابط، توزیع‌ها و عملکرد مدل‌های پیش‌بینی را از طریق visualization‌های متنوع بررسی کردیم. ماتریس همبستگی (Heatmap) نشان می‌دهد که ویژگی‌هایی مانند Thallium (با همبستگی ۰,۵۲) و ST depression (۰,۴۲) بیشترین ارتباط مثبت با وجود بیماری قلبی (Presence) دارند، در حالی که حداکثر ضربان قلب (Max HR) با همبستگی منفی -۰,۴۲ به عنوان عامل محافظه عمل می‌کند؛ این الگوها حاکی از اهمیت مشکلات عروقی و الکتریکی قلب در پیش‌بینی است، بدون وجود multicollinearity شدید (همبستگی‌های بالای ۰,۸ نادر است).

Histogram کلی داده‌ها (data.hist) توزیع ویژگی‌های عددی را آشکار می‌سازد، جایی که سن (Age) توزیع نرمالی با میانگین ۵۴ سال دارد، اما کلسترول skew (Cholesterol) راست نشان می‌دهد (با مقادیر بالا در ۲۰٪ موارد)، که با BP و Cholesterol boxplot BP، میانه کلسترول به ۲۷۰ می‌رسد (بالاتر از ۲۴۰ در outliers خارج از IQR) در کلسترول ≈ ۲۰ مورد را تشکیل می‌دهند، که می‌تواند نشان‌دهنده بیماران پرخطر با سطوح غیرطبیعی باشد.

KDE (kernel density estimation) با Histogram Age دارد و منحنی صاف نرمال را ترسیم می‌کند، که با (Age-Cholesterol-BP) 3D Scatter (BP < ۳۰۰ و BP > ۱۴۰ می‌شوند، در حالی که نقاط آبی (Absence) پراکنده‌تری نشان می‌دهند، و این الگو بر نقش سن و کلسترول به عنوان عوامل تجمعی خطر تأکید دارد.

Pairplot کلی روابط pairwise را برجسته می‌کند، جایی که نشان‌دهنده جداسازی واضح کلاس‌ها است (Presence با Thallium بالا و Max HR پایین)، و این با scatterplot KMeans (۳ کلاستر، $k \approx ۴۲۱$) همخوانی دارد: کلاسترها بر اساس سن و کلسترول جدا می‌شوند، با کلاستر قرمز (Presence) در مقادیر بالا، که پیشنهاد می‌کند KMeans می‌تواند برای unsupervised grouping بیماران پرخطر مفید باشد.

Presence ۱۲۰ Absence vs ۱۵۰) را آشکار می‌کند (Bar-plot توزیع کلاس imbalance countplot) نشان می‌دهد که مردان (Sex=۱) ۱۷۰٪ موارد به ۴۴٪، و countplot Heart Disease by Sex (Sex vs Disease) heatmap cross-tab را تشکیل می‌دهند (در مقابل ۵۰٪ در زنان)، که Presence مقادیر ۵۰/۷۰ تأیید می‌کند و بر نابرابری جنسیتی در خطر بیماری تأکید دارد – این یافته‌ها حاکی از نیاز به غربالگری جنسیت محور است.

در بخش مدل‌سازی، K=۱-۲۰) line plot KNN scores دقت را به عنوان تابعی از همسایگان نشان می‌دهد، با پیک ۸۹٪ در K=۱۸، که نشان‌دهنده بهینه‌سازی hyperparameter برای KNN است و overfitting در K‌های پایین (به دلیل نویز) را برجسته می‌کند. Bar plot SVC kernels مقایسه kernelها را ترسیم می‌کند، جایی که linear با ۹۱٪ دقت برتر است (poly، RBF، sigmoid ۸۵٪، ۸۸٪)، و این برتری به دلیل روابط خطی قوی در داده (مانند همبستگی Thallium) است.

Line plot Decision Tree max_features دقت را بر اساس تعداد ویژگی‌ها نشان می‌دهد، با بهبود از ۷۰٪ در max_features=۱ به ۸۵٪ در max_features=۲، که حاکی از سادگی مدل DT در این دیتاست است و نیاز به pruning برای جلوگیری از overfitting را پیشنهاد می‌کند. Bar plot Random Forest estimators دقت را به عنوان تابعی از تعداد درخت‌ها (۱۰۰۰-۱۰۰) ترسیم می‌کند، با تثبیت در ۸۷٪ برای ۲۰۰ درخت، که نشان‌دهنده ensemble learning روش robust برای مدیریت نویز (مانند outliers کلسترول) است.

در مجموع، این visualization‌ها الگویی جامع را آشکار می‌سازند: عوامل خطر اصلی Thallium بالا، کلسترول clustering SVC linear پیش‌بینی می‌شوند، در حالی که KMeans و imbalance جنسیتی بر لزوم مداخلات هدفمند (مانند غربالگری مردان بالای ۵۰ سال) تأکید دارند؛ این یافته‌ها نه تنها دقت مدل‌ها را (۸۵-۹۱٪) تأیید می‌کنند، بلکه کاربرد بالینی را برای تشخیص زودهنگام برجسته می‌سازند، با پتانسیل بهبود از طریق SMOTE برای تعادل کلاس‌ها.