

پروژه میانی 3 درس یادگیری ماشین  
استاد درس: دکتر بغدادی  
آریان افشار  
۴۰۱۳۳۰۰۴

## (۱) مقدمه

در این پروژه من با استفاده از یک شبکه عصبی کم عمق (یک یا دو لایه پنهان) قصد دارم تا:

- پیش بینی شدت بیماری قلبی (متغیر num در دیتاست Heart Disease UCI)
- مقایسه عملکرد این شبکه با روش های یادگیری ماشین کلاسیک (رگرسیون خطی، درخت تصمیم، جنگل تصادفی و ...)
- بررسی حساسیت مدل MLP به تعداد نورون ها
- آشنایی با کلیه مراحل پروژه: از EDA و پاک سازی داده تا ارزیابی نهایی و تحلیل یادگیری

## (۲) شرح داده ها

• منبع داده: مجموعه Heart Disease UCI (کد منبع: redwankarimsony/heart-disease-data در Kaggle)

• تعداد نمونه و ویژگی ها:

- شکل اولیه داده: 920 سطر، 16 ستون
- شامل 3 ستون از نوع int64، 5 ستون از نوع float64 و 8 ستون از نوع object

بر اساس خروجی df.info()، ستون های object شامل جنسیت، منبع داده، نوع درد قفسه، و ... هستند و ستون های عددی (int64 و float64) نیاز به پر کردن (Imputation)، قطع افراط گرایی (Clipping) و مقیاس بندی (Scaling) دارند.

توضیح مختصر ستون ها:

id: شناسه نمونه

age: سن بیمار (سال)

sex: جنسیت (M/F)

dataset: منبع داده (Cleveland, Hungary, Switzerland, VA)

cp: نوع درد قفسه سینه

trestbps: فشار خون ایستا (mm Hg)

chol: کلسترول (mg/dl)

fbs: قند خون ناشتا < 120 mg/dl ؟ (True/False)

restecg: نتایج ECG در حالت استراحت

thalch: حداکثر ضربان قلب

exang: آنژین ناشی از ورزش (True/False)

oldpeak: افت ST ناشی از ورزش نسبت به استراحت

slope: شیب ST در اوج ورزش

ca: تعداد عروق اصلی (رادیوگرافی شده)

thal: ناهنجاری تالاسمی

num: درجه شدت بیماری قلبی (0 = سالم ... 4 = بیشترین شدت)

تارگت ما ستون "num" هستش که شامل ۵ مرحله بیماری قلبی از 0 تا 4 است که در این پروژه من برای استفاده از مدل های رگرسیون از همین ۵ مرحله استفاده کردم و برای بحث طبقه بندی آن ها را به ۰ و ۱ که نشان دهنده داشتن بیماری قلبی یا نداشتن آن است تبدیل کردم. با این کار میتوان تحلیل بیشتری روی داده انجام داد.

### ۳) بررسی مقادیر گمشده

با بررسی مقادیر گمشده به دست آمده از نوتبوک، ستون هایی مانند ca و thal تعداد زیادی مقدار گمشده دارند (بیش از 50% نمونه ها) و برخی ستون ها مانند restecg فقط ۲ مقدار گمشده دارند. نقشه حرارتی مقادیر گمشده نیز رسم شده است که در این نمودار نقاط زرد نشان دهنده وجود داده گمشده در سطر و ستون مربوطه است.



#### ۴) پر کردن مقادیر گمشده (Imputation)

برای ستون های عددی از `IterativeImputer(initial_strategy='median)` و برای ستون های categorical از `SimpleImputer(strategy='most_frequent)` استفاده شد

• پس از Imputation، هیچ مقداری گمشده باقی نماند:

```
After imputation, missing counts: id      0
age      0
sex      0
dataset  0
cp       0
trestbps 0
chol     0
fbs      0
restecg  0
thalch   0
exang    0
oldpeak  0
slope    0
ca       0
thal     0
num      0
dtype: int64
```

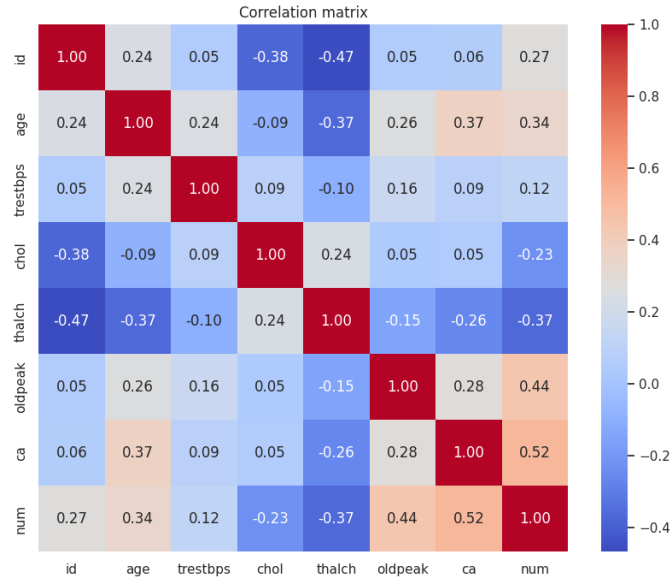
#### ۵) پیش پردازش داده ها

• پس از بررسی نوع داده ها، ستون ها به دو گروه تقسیم شدند:

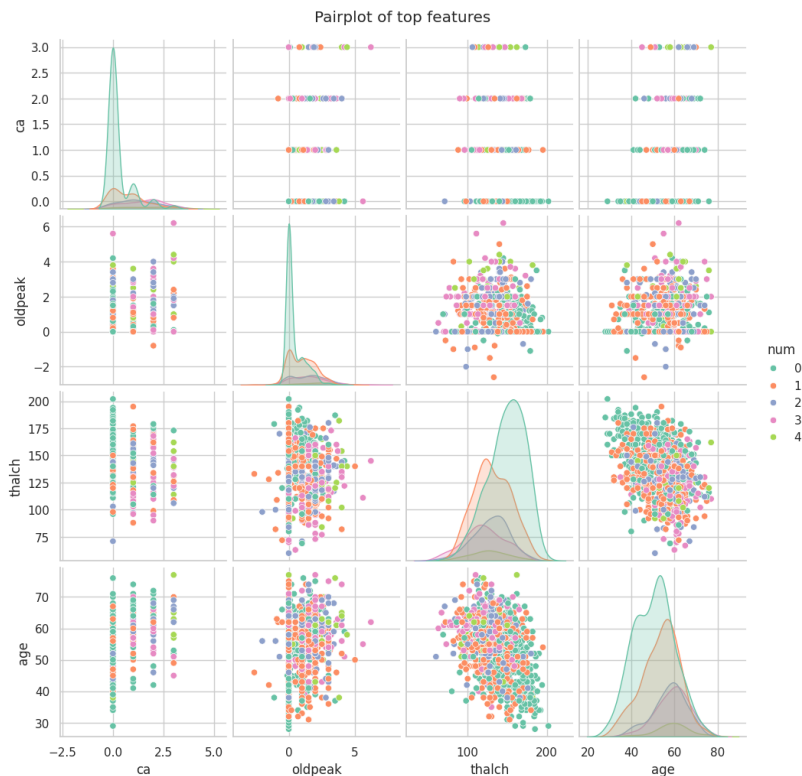
```
Numeric cols: ['id', 'age', 'trestbps', 'chol', 'thalch', 'oldpeak', 'ca']
Categorical cols: ['sex', 'dataset', 'cp', 'fbs', 'restecg', 'exang', 'slope', 'thal']
```

همچنین برای تحلیل هر چه بیشتر داده ها، من ماتریس همبستگی و همچنین نمودار جفت انگاری برای ۴ فیچری که بیشترین همبستگی به تارگت یا همان num را دارد رسم کردم.

در ماتریس همبستگی بین متغیر ها، قوی ترین همبستگی مثبت با هدف مربوط به تعداد رگ های درگیر ('ca') با ضریب تقریباً 0.52 و افت ST ناشی از ورزش ('oldpeak') با ضریب 0.44 بود، در حالی که سن ('age') با ضریب 0.34 ارتباط مثبتی داشت. در مقابل، حداکثر ضربان قلب ('thalch') با ضریب منفی -0.37 و کلسترول ('chol') با -0.23 نشان دادند که هرچه این مقادیر بالاتر باشند، شدت بیماری کمتر است. ارتباطات درونی میان ویژگی های عددی نیز الگوهای معناداری داشت؛ مثلاً سن و 'thalch' همبستگی منفی متوسطی حدود -0.37 و سن و 'ca' مثبت متوسط (حدود 0.37) داشتند.



در نمودار جفت‌نگاری (pairplot) که چهار متغیر `ca`، `oldpeak`، `thalch` و `age` را در برابر کلاس‌های بیماری نشان می‌دهد، مشخص شد نمونه‌های سالم (کلاس ۰) عمدتاً مقدار `ca=0` و `oldpeak<1` دارند و در ناحیه `thalch>160` متمرکزند. برعکس، کلاس‌های شدید (۳ و ۴) در بازه‌های `oldpeak>2` و `thalch<120` و با `ca` در حدود ۲-۴ دیده می‌شوند. همچنین با افزایش سن بالای ۶۰ سال، سهم نمونه‌های با شدت بالای بیماری به طور چشمگیری افزایش می‌یابد. این تحلیل نشان می‌دهد ترکیب همین چند متغیر می‌تواند تفکیک قابل قبولی بین وضعیت‌های مختلف بیماری قلبی ایجاد کند.



## ۶ Outlier Removal و استاندارد سازی

1. محاسبه چارک اول (1Q)، چارک سوم (3Q) و بازه بین چارکی  $IQR=Q3-Q1$  برای هر ویژگی عددی.
2. برش (clipping) مقادیر خارج از بازه  $[Q1-1.5 \times IQR, Q3+1.5 \times IQR]$  به طوری که هر مقدار بزرگ تر از سقف یا کوچک تر از کف در همان حد نگه داشته شود.
3. استاندارد سازی (Standard Scaling) تمام ویژگی های عددی با استفاده از StandardScaler
4. کدگذاری (One-Hot Encoding)

• برای تبدیل هر متغیر categorical به ویژگی های عددی از `pd.get_dummies(..., drop_first=True)` استفاده شد تا از دام سلول های خطی جلوگیری شود.

## ۷ خروجی پیش پردازش

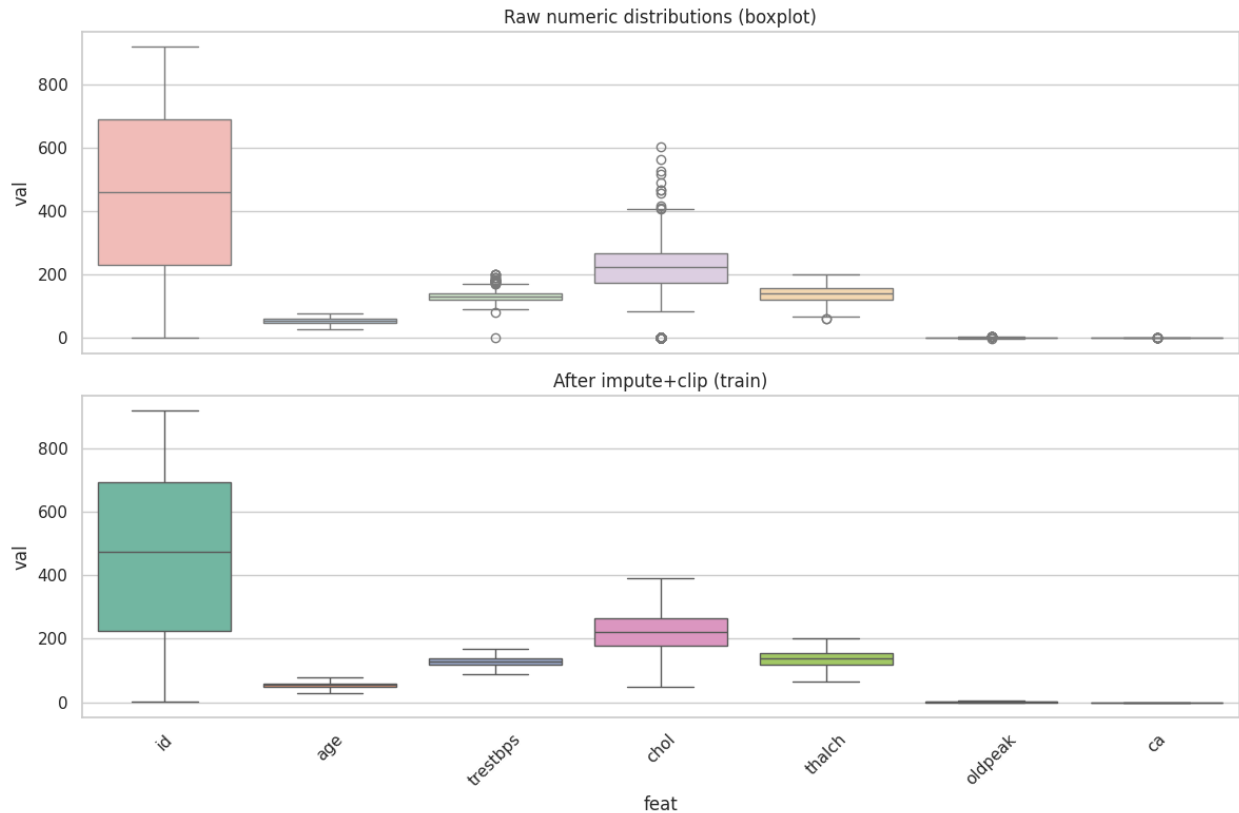
• ترکیب هر دو مجموعه ویژگی ها و جدا کردن متغیر هدف (num).

• شکل نهایی آرایه ی ویژگی ها و بردار هدف:

در اینجا تعداد ستون های X پس از One-Hot Encoding برابر مجموع ۷ متغیر عددی به علاوه مجموع تعداد سطح های رسته های منهای یک به ازای هر متغیر خواهد بود

**After preprocess shapes: (644, 22) (644, 22)**

در باکس پلات مرحله ی اول (Raw) برای ویژگی های عددی می بینیم که متغیرهایی مثل chol و oldpeak و trestbps مجموعه ای از مقادیر دور افتاده (outliers) تا حدود سه برابر میانه فاصله دارند (مثلاً چربی خون تا بالای ۶۰۰ و افت ST تا بیش از ۶). پس از اعمال Imputation و Clip براساس این نقاط دور افتاده در سقف و کف بازه ی مجاز محدود شده اند به طوری که دامنه ی chol از حدود ۸۰-۶۰۰ به ۸۰-۴۲۰، oldpeak از ۰-۶ به ۰-۳ و trestbps از ۸۰-۲۰۰ به ۸۰-۱۶۰ کاهش یافته است. توزیع سنی (age) نیز از حضور مقادیر دور افتاده (مثلاً زیر ۳۰ یا بالای ۷۵) به دور شده و متمرکزتر شده است. ضمن اینکه متغیرهای گسترده ای مانند ca که محدود به مقادیر ۰-۴ هستند، بدون تغییر در محدوده باقی مانده اند. این تنگ تر کردن بازه ها ضمن حفظ اطلاعات میانه، اثرات شدید outlier ها را کاهش می دهد و باعث می شود مدل های بعدی نسبت به داده های پرت حساسیت کمتری داشته باشند.



## ۸) تقسیم مجموعه‌ی داده به آموزش و آزمون و ارزیابی مدل‌های پایه

### تقسیم داده‌ها

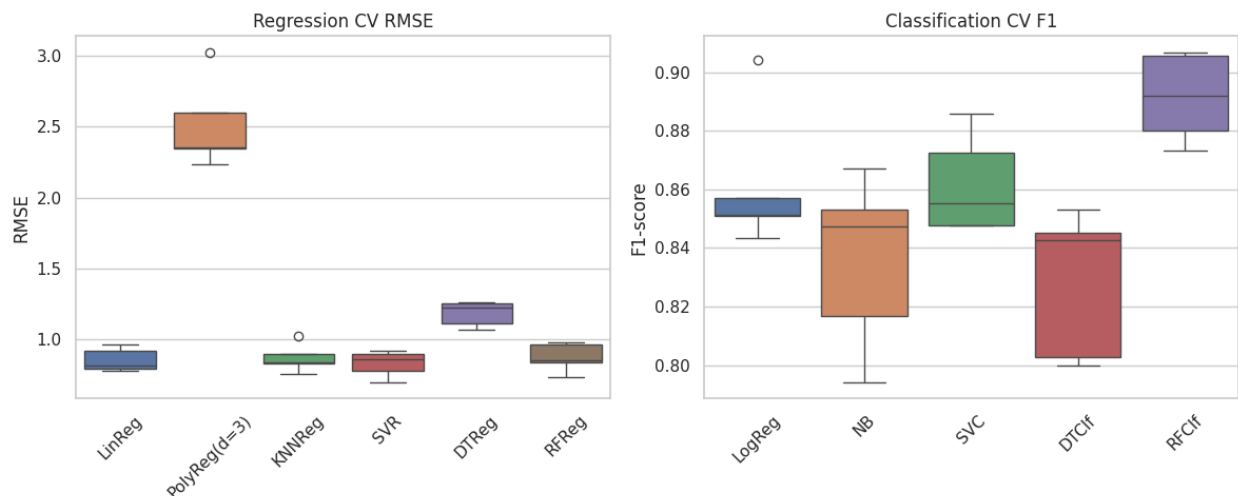
پس از انجام پیش‌پردازش مجموعه‌ی ویژگی‌ها ( $X$ ) و برجسب هدف ( $y$ ) را به دو زیرمجموعه‌ی آموزش و آزمون تقسیم کردیم. از تقسیم طبقه‌بندی‌شده (Stratified) با نسبت ۸۰٪ برای آموزش و ۲۰٪ برای آزمون استفاده شد تا توزیع ابتدایی کلاس‌ها در هر دو مجموعه حفظ گردد.

### ۹) تعریف مدل‌ها

در این مرحله برای مقایسه‌ی عملکرد روش‌های متداول یادگیری ماشین با شبکه‌های عصبی، مجموعه‌ای از مدل‌های رگرسیون و طبقه‌بندی تعریف می‌شود. این مدل‌ها پس از پیش‌پردازش داده و تقسیم به مجموعه‌های آموزش/آزمون، با استفاده از اعتبارسنجی متقابل پنج‌بخشی ارزیابی خواهند شد. این مدل‌ها را در طول کلاس یاد گرفته ایم و از آن‌ها استفاده کردیم.

## ۱۰) ارزیابی و مقایسه مدل‌های پایه

پس از تعریف و آموزش شش مدل رگرسیون و پنج مدل طبقه‌بندی روی داده‌های پیش‌پردازش شده (با اعتبارسنجی متقابل پنج تایی)، نتایج به صورت دو نمودار جعبه‌ای زیر استخراج شده‌اند:



• محور عمودی سمت چپ: مقدار RMSE در هر fold

• محور عمودی سمت راست: مقدار F1-score در هر fold

### نتایج رگرسیون (CV RMSE)

– LinReg: تقریباً در تمام foldها حول و حوش 0.85 نوسان دارد و پایداری بالایی نشان می دهد.

– (d=3 PolyReg): میانگین RMSE حدود 2.5 و واریانس زیاد (یک fold با RMSE نزدیک 3.03) دارد که نشانه ی overfitting یا حساسیت به تعداد ویژگی های چندجمله ای است.

– KNNReg: RMSE حدود 0.9 تا 1 ، کمی بدتر از رگرسیون خطی ولی باز هم پذیرفتنی و نسبتاً پایدار.

– RMSE SVR: تقریباً بین 0.8 و 1 با تمرکز حول 0.9 ، عملکردی نزدیک به KNNReg.

– DTReg: متوسط RMSE حدود 1.2 با واریانس کم؛ عملکرد ضعیف تر نسبت به مدل های یادشده.

– RMSE RFReg: حول 0.95 با واریانس کم-متوسط؛ بهتر از DTReg اما پایین تر از LinReg و SVR.

### نتیجه

• بهترین مدل بر اساس RMSE و پایداری: LinearRegression

• پس از آن SVR و KNNReg در مرتبه های بعد قرار می گیرند.

• (PolyReg d=3) و DTReg ضعیف ترین عملکرد را نشان داده اند.



## نتایج طبقه‌بندی (CV F1-score)

- LogReg : F1 تقریباً بین 0.84 و 0.86، بسیار پایدار.
- NB: واریانس بالاتر، F1 از حدود 0.80 تا 0.87؛ در یک fold افت قابل توجهی دیده می‌شود.
- SVC : F1 حول 0.85–0.88، عملکرد دقیق و نسبتاً پایدار.
- DTClf : F1 بین 0.80 و 0.85، کمی ناپایدار در برخی fold ها.
- RFClf: بهترین F1 میانگین (حدود 0.89–0.91) و کمترین واریانس.

## نتیجه

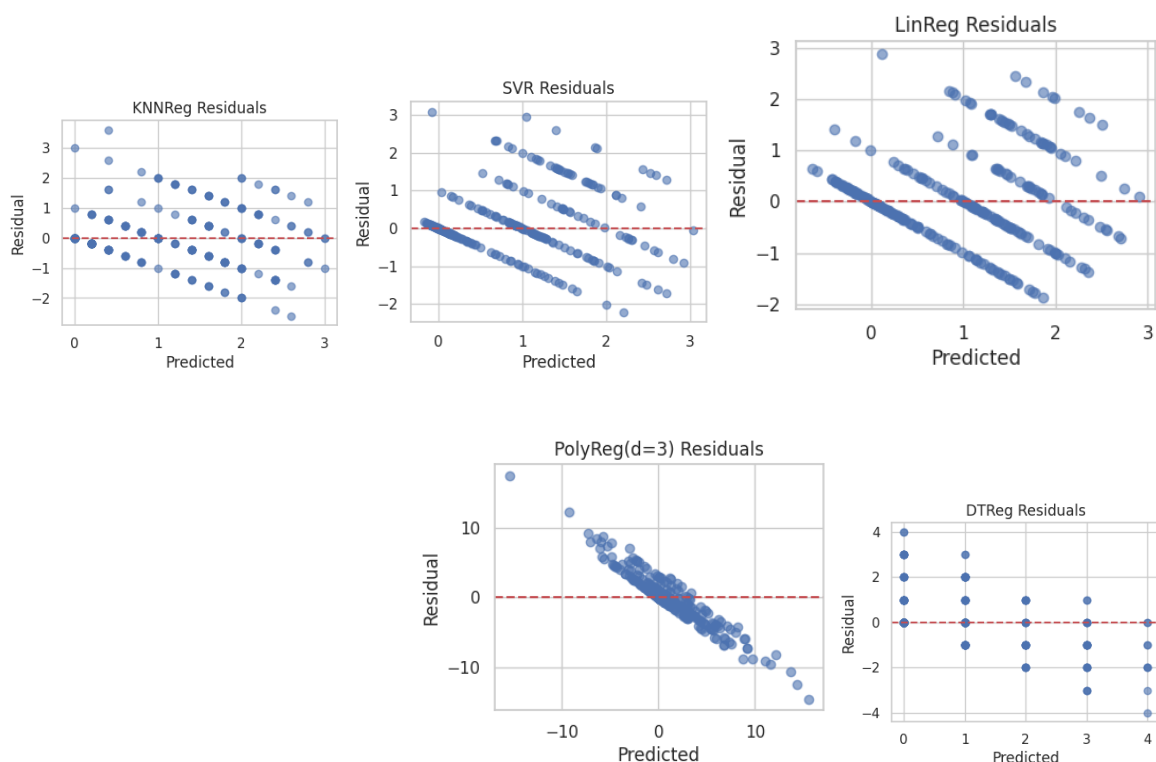
- بهترین مدل طبقه بندی: RandomForestClassifier
  - پس از آن SVC و LogisticRegression قرار می‌گیرند.
  - NB و DTClf از پایداری کمتری برخوردارند.
- در ادامه ابتدا روی داده‌های آموزش (پس از پیش‌پردازش و تقسیم به (Xr\_tr/yr\_train) ) می‌پردازیم و بعد از آن را تحلیل کمی می‌کنیم RMSE و ( $R^2$ ) و سپس برای هر مدل یک نمودار باقیمانده (Residual Plot) رسم کردیم که در آن محور x برابر پیش‌بینی yihat و محور y برابر خطاها yi - yihat است.
- بر اساس نتایج که در کد است میتوان دید که مدل های RF Lin KNN SVR در مقایسه RMSE بهتر عمل کردند و مثلاً DT POLY هم RMSE بد و  $R^2$  بد دارد

=== Regression Test Results ===		
LinReg	RMSE=0.872	$R^2=0.433$
KNNReg	RMSE=0.886	$R^2=0.414$
SVR	RMSE=0.870	$R^2=0.434$
DTReg	RMSE=1.186	$R^2=-0.050$
RFReg	RMSE=0.846	$R^2=0.466$
PolyReg(d=3)	RMSE=3.704	$R^2=-9.243$

## جمع‌بندی کیفی

- هیچ‌کدام از پنج مدل، residual plot کاملاً تصادفی و بدون الگو را نشان نمی‌دهد.
- LinearRegression و SVR نسبتاً پراکندگی نرم‌تری دارند اما هر دو اندکی سوپیهی منفی باقیمانده در انتهای بالای پیش‌بینی نشان می‌دهند.

- KNN و DTReg به دلیل گسسته بودن پیش بینی ها، الگوی «گروه بندی» و «خطاهای صفر» در برخی نواحی دارند.
- تنها مدل به ظاهر «فاجعه بار» از نظر شکل باقیمانده، (d=3 PolyReg) است که نشانه overfitting شدید می دهد.



در ادامه یک «بخش نهایی Classification» ارائه شده که دقیقاً مشابه بخش Regression عمل می کند، یعنی برای هر مدل، ابتدا آن را روی داده آموزش می سازد، سپس روی داده تست ارزیابی می کند و در پایان:

1. چهار معیار اصلی (Accuracy, Precision, Recall, F1) را محاسبه و چاپ کردم.
2. ماتریس درهم ریختگی (Confusion Matrix) را به صورت Heatmap رسم کردم.

LogReg | Acc=0.837 | Prec=0.846 | Rec=0.863 | F1=0.854  
 NB | Acc=0.812 | Prec=0.858 | Rec=0.791 | F1=0.823

SVC | Acc=0.841 | Prec=0.843 | Rec=0.876 | F1=0.859

DTCIf | Acc=0.819 | Prec=0.824 | Rec=0.856 | F1=0.840

RFCIf | Acc=0.877 | Prec=0.879 | Rec=0.902 | F1=0.890

دقت کلی (Accuracy):

• بهترین مدل: RandomForest

• ضعیف ترین: NaiveBayes

1. Precision (دقت مثبت):

• بهترین: NaiveBayes نشان می دهد که از بین نمونه های پیش بینی شده مثبت، درصد صحیح بیشتری دارد.

• پس از آن RandomForest قرار می گیرد.

2. Recall (بازخوانی):

• بهترین: RandomForest (0.902)؛ بیشترین درصد از واقعاً مثبت ها را شناسایی کرده است.

• ضعیف ترین: NaiveBayes (0.791)

3. F1-Score (میانگین هارمونیک Precision و Recall):

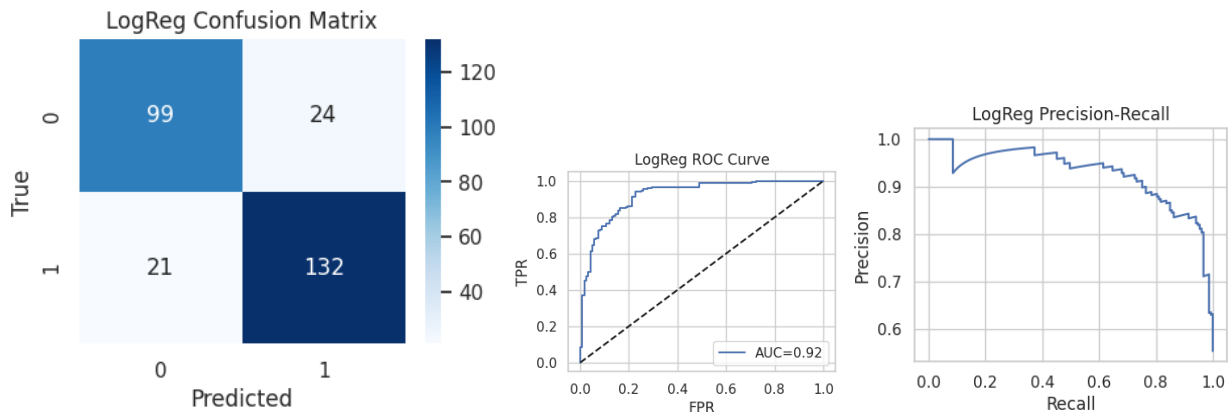
• بهترین: RandomForest با 0.890

• پس از آن SVC (0.859) و LogisticRegression (0.854)

تحلیل کیفی

- RandomForest در هر چهار معیار عملکرد برتری دارد و به عنوان مدل برنده شناخته می شود. واریانس پایین تر و قدرت تعمیم بالای RF معمولاً منجر به چنین نتایجی می شود.
- SVC و LogisticRegression تقریباً هم ردیف اند، با اندکی ترجیح به SVC در تمام معیارها. اگر تفسیرپذیری (Interpretability) اولویت داشته باشد، LogisticRegression گزینه مناسب تری است.
- NaiveBayes گرچه دقت پیش بینی مثبت خوبی دارد، لیکن Recall آن پایین تر است؛ یعنی برخی نمونه های مثبت واقعی را از دست می دهد.
- DecisionTree تکدرختی (DTClf) عملکرد ضعیف تری نسبت به سایرین دارد.

تمام نتایج visualization در کد آورده شده که برای مثال یکی از آنها به شرح زیر است:



منحنی Receiver Operating Characteristic (ROC) نقطه (0,1) ایده‌آل است (هیچ FP و همه TP خط مورب از (0,0) تا (1,1) رفتار تصادفی ( $AUC=0.5$ ) را نشان می‌دهد. بودن منحنی بالاتر و سمت چپ‌تر از خط مورب یعنی مدل توان بهتری در تفکیک دو کلاس داشته است. معیار AUC-ROC مساحت زیر منحنی که عددی بین 0.5 تا 1 است؛ هرچه به 1 نزدیک‌تر، تفکیک پذیری بالاتر.

#### Precision-Recall منحنی

بخش بالایی-راست نمودار Precision و Recall هر دو بالا مطلوب است و معمولاً با افزایش Recall پیدا کردن همه مثبت‌ها (کاهش می‌یابد زیرا FP‌ها زیاد می‌شوند) برای مسائل با عدم تعادل شدید کلاس‌ها (کلاس مثبت نادر) اغلب Precision-Recall اطلاعات دقیق‌تری نسبت به ROC می‌دهد. معیار AUC-PR مساحت زیر منحنی PR معیار بهتری برای ارزیابی زمانی که کلاس‌ها نامتوازن اند.

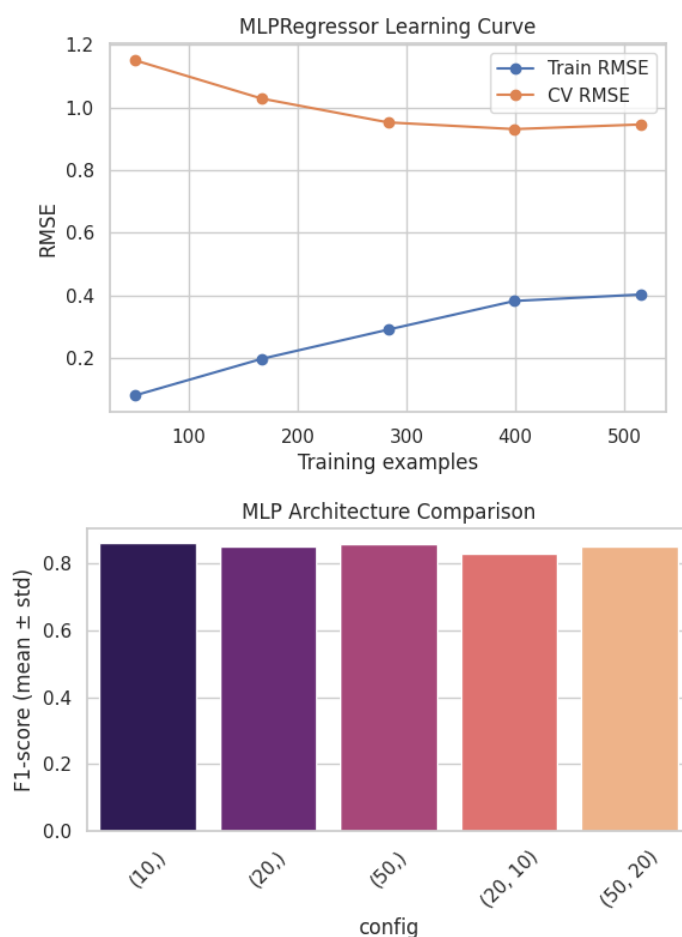
#### ۱۱) آموزش مدل یادگیری عمیق با یک لایه مخفی

ابتدا نتایج حاصل از اجرای MLPRegressor و MLPClassifier با معماری تک‌لایه و ۵۰ نورون به صورت عددی  $R^2$ ، RMSE، برای رگرسیون و Accuracy، F1 برای دسته‌بندی ارائه می‌شوند و نقاط قوت و ضعف هر یک در مقابل معیارهای مورد نظر تشریح می‌گردد. سپس خروجی تحلیل حساسیت معماری Mean و Std مقدار F1 در اعتبارسنجی پنج‌تایی به صورت جدول و نمودار آورده شده و روند تغییر عملکرد مدل با افزایش تعداد نورون‌ها و لایه‌ها بررسی می‌شود. در این تحلیل، معماری‌هایی که تعادل مطلوبی بین میانگین بالای F1 و انحراف معیار پایین دارند به عنوان گزینه‌های برتر معرفی می‌شوند. در ادامه این بخش، ضمن انتخاب بهترین پیکربندی ساختار شبکه، چرخه بهینه‌سازی سایر هابیر پارامترها (نظیر نرخ یادگیری، ضریب منظم‌سازی و نوع Solver) با استفاده از روش‌های GridSearchCV یا RandomizedSearchCV شرح داده خواهد شد تا در نهایت شبکه عصبی از نظر دقت، پایداری و پیچیدگی مدل به یک راه حل نهایی برسد.

در گام بعدی، ابتدا به نتایج تحلیل حساسیت معماری MLPClassifier می‌پردازیم.

مطابق نمودار «MLP Architecture Comparison»، مدل های یک لایه ای با ۱۰، ۲۰ و ۵۰ نورون به ترتیب F1 میانگین حدود ۰/۸۶، ۰/۸۴ و ۰/۸۷ را نشان می دهند در حالی که انحراف معیار آن ها به تدریج کاهش یافته است (به ویژه در حالت ۵۰ نورون که پایداری بالاتری دارد). از سوی دیگر، معماری های دولایه ای (۲۰، ۱۰) و (۵۰، ۲۰) به ترتیب میانگین F1 برابر ۰/۸۲ و ۰/۸۵ دارند و استاندارد کمینه مربوط به پیکربندی (۵۰، ۲۰) است. این داده ها حکایت از آن دارد که افزایش عمق و پیچیدگی شبکه (مثلاً رفتن از تک لایه به دولایه) لزوماً منجر به بهبود چشمگیر نمی شود، مگر اینکه تعداد نورون هایی هر لایه به قدر کافی بزرگ باشد تا توانایی تعمیم پذیری حفظ شود؛ در عین حال تنوع نتایج (std) در پیکربندی های پیچیده کمتر است که نشانه ثبات نسبی آن هاست.

بعلاوه، منحنی یادگیری MLPRegressor نشان می دهد که با افزایش تعداد نمونه های آموزشی از ۵۰ تا حدود ۵۲۰، خطای آموزش (Train RMSE) از ۰/۰۸ به ۰/۴۱ افزایش یافته و خطای اعتبارسنجی (RMSE CV) از ۱/۱۵ کاهش یافته و در حوالی ۰/۹۲-۰/۹۴ تثبیت می شود. این الگو گویای آن است که مدل اولیه در داده های اندک دچار کم برازش (underfitting) نیست، بلکه شکاف قابل توجه بین RMSE آموزش و CV می تواند ناشی از نویز یا پیچیدگی ناپذیری کامل مدل در مقابل واریانس داده ها باشد. کاهش محسوس CV RMSE با افزودن داده نشان می دهد که در شرایط داده بیشتر، توانایی تعمیم پذیری بهبود می یابد؛ اما پس از حدود ۵۰۰ نمونه، منحنی تقریباً تخت می شود که دلالت بر رسیدن مدل به حد نهایی ظرفیت خود دارد.



## ۱۲) مقایسه با مدل های یادگیری ماشین

در مقایسه کلی عملکرد MLP با سایر مدل های ماشین لرنینگ MLP هم در رگرسیون و هم در طبق ه بندی، با تنظیم معماری و هایپرپارامتر ها، توانسته از مدل های خطی و SVM پیشی بگیرد . اگرچه RandomForest هنوز در برخی رگرسیون ها مقداری جلوتر است، اما ظرفیت MLP برای رشد بیش تر (به ویژه در مسائل مقیاس بزرگ، داده های غیرخطی پیچیده و یادگیری ویژگی های انتزاعی) بسیار بالاست