

پروژه میانی ۱ درس یادگیری ماشین
استاد درس: دکتر بغدادی
آریان افشار
۴۰۱۳۳۰۰۴

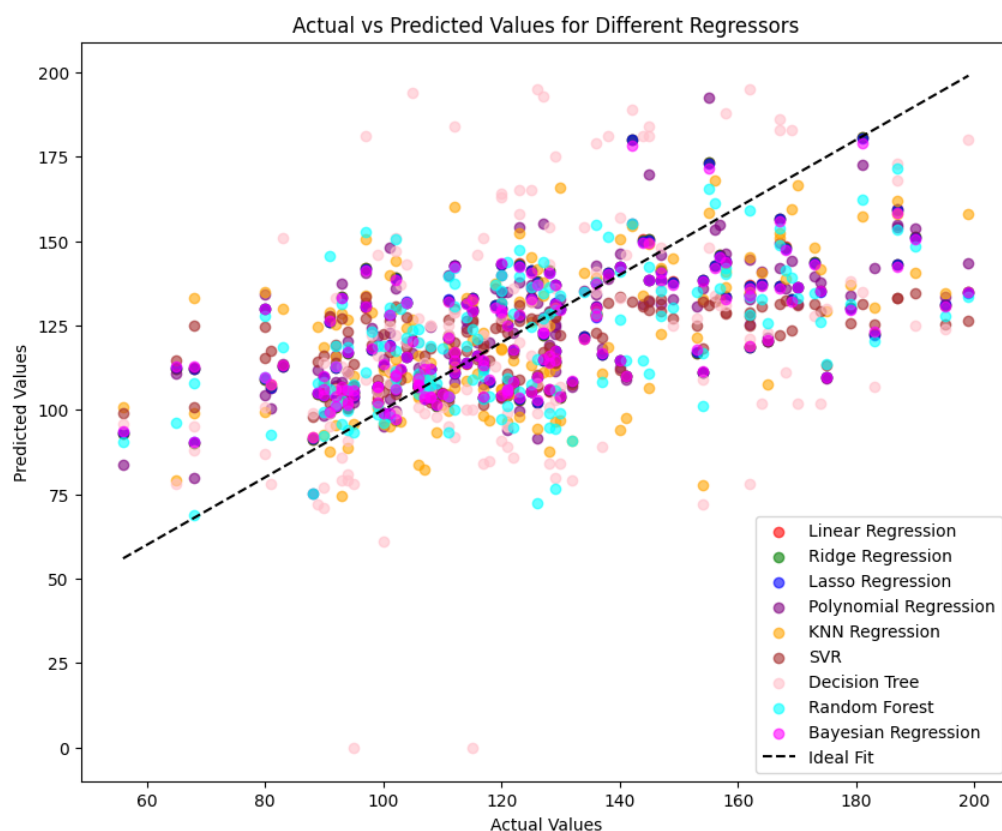
مقدمه:

در این پروژه، با استفاده از دیتاست Pima Indians Diabetes Database، که شامل ویژگی‌های پزشکی مختلفی از جمله سن، فشار خون، و سطح گلوکز است، به بررسی و مقایسه عملکرد مدل‌های مختلف رگرسیونی برای پیش‌بینی مقدار پیوسته‌ی سطح گلوکز خون پرداختیم. هدف اصلی، ارزیابی توانایی مدل‌های مختلف رگرسیون در یادگیری الگوی ارتباط بین ویژگی‌های ورودی و سطح گلوکز و انتخاب بهترین مدل برای این منظور است.

توضیحات مربوط به نوشتن کد در خود کد به صورت تکست آورده شده و در این گزارش کار من فقط به بخش تحلیل نمودارها و انتظارات خواسته شده پرداختم. همچنین با اجازه استاد از کدهای تمرین ۳ برای رسم نمودارها و پیاده‌سازی مدل‌ها استفاده کردم.

۱. مقایسه عملکرد مدل‌های مختلف رگرسیونی:

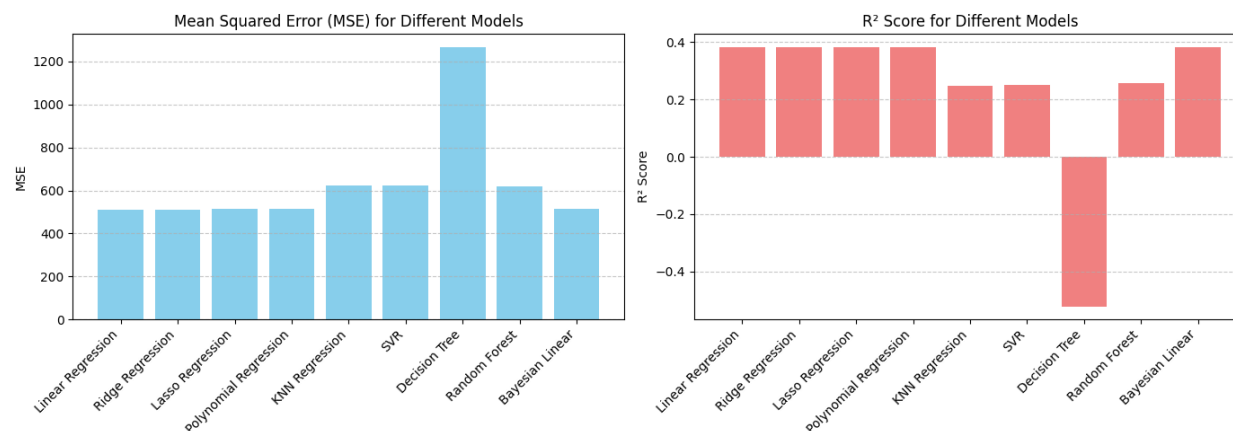
برای ارزیابی عملکرد مدل‌های مختلف رگرسیونی، از معیارهای رایج مانند Mean Squared Error (MSE) و R^2 (Squared Score) استفاده شده است. همچنین، نمودار مقادیر واقعی در برابر مقادیر پیش‌بینی شده برای هر مدل ترسیم شده است تا به صورت بصری عملکرد آن‌ها را مقایسه کنیم.



شکل ۱: نمودار Actual vs Predicted Values برای مدل‌های مختلف رگرسیونی

این نمودار مقادیر واقعی (Actual Values) را در برابر مقادیر پیش بینی شده (Predicted Values) برای هر یک از مدل های رگرسیونی نشان می دهد. خط چین سیاه نیز "Ideal Fit" را نشان می دهد که در آن مقادیر واقعی و پیش بینی شده برابر هستند.

با بررسی شکل ۱، می توان مشاهده کرد که مدل هایی که نقاط داده ها آن ها بیشتر در نزدیکی خط چین سیاه تجمع جمع شده اند، عملکرد بهتری در پیش بینی سطح گلوکز داشته اند. در نگاه اول، به نظر می رسد مدل های Linear Regression، Ridge Regression، Lasso Regression، Polynomial Regression و Bayesian Regression نقاطی دارند که به خط ایده آل نزدیک تر هستند نسبت به مدل هایی مانند Decision Tree و SVR. Random Forest نیز پراکندگی قابل قبولی از نقاط را نشان می دهد. مدل Decision Tree بیشترین پراکندگی نقاط را دارد که نشان دهنده دقت پایین تر آن در پیش بینی است.



شکل ۲: معیار های عملکرد مدل های مختلف رگرسیونی MSE و R² Score

این نمودار مقدار ضریب تعیین (R² Score) را برای هر یک از مدل ها نشان می دهد. R² Score نشان می دهد که چه نسبتی از واریانس متغیر هدف (در اینجا سطح گلوکز خون) توسط مدل توضیح داده می شود. مقادیر بالاتر R² (نزدیک به ۱) نشان دهنده عملکرد بهتر است. مقادیر منفی R² نشان دهنده این است که مدل بدتر از پیش بینی میانگین عمل کرده است. همچنین مقدار خطای میانگین مربعات (MSE) را برای هر یک از مدل ها نیز نشان می دهد. MSE یک معیار رایج برای ارزیابی عملکرد مدل های رگرسیونی است و مقادیر کمتر MSE نشان دهنده عملکرد بهتر است.

تحلیل Mean Squared Error (MSE)

بر اساس نمودار MSE در شکل ۲، مقادیر MSE برای مدل های مختلف به شرح زیر است:

Linear Regression تقریباً ۵۰۰

Ridge Regression تقریباً ۵۰۰

Lasso Regression تقریباً ۵۰۰

Polynomial Regression تقریباً ۵۰۰

KNN Regression تقریباً ۶۰۰

SVR تقریباً ۶۰۰

Decision Tree تقریباً ۱۲۵۰ (بالاترین MSE)

Random Forest تقریباً ۶۰۰

Bayesian Regression تقریباً ۵۰۰ (کمترین MSE)

البته که مقدار دقیق آن ها در کد چاپ شده است:

```
Linear Regression: MSE=511.9708, R2=0.3832
Ridge Regression: MSE=512.0028, R2=0.3831
Lasso Regression: MSE=512.6320, R2=0.3824
Polynomial Regression: MSE=513.6550, R2=0.3812
KNN Regression: MSE=624.5678, R2=0.2475
SVR: MSE=622.5677, R2=0.2499
Decision Tree Regression: MSE=1263.9156, R2=-0.5227
Random Forest Regression: MSE=617.0863, R2=0.2565
Bayesian Linear Regression: MSE=512.9068, R2=0.3821
```

همانطور که مشاهده می‌شود، مدل Linear Regression با کمترین مقدار MSE تقریباً (511) بهترین عملکرد را از نظر حداقل کردن خطای پیش بینی نشان می‌دهد. مدل‌های bayesian Regression ، Ridge Regression ، Lasso Regression و Polynomial Regression نیز MSE مشابه و پایینی دارند. در مقابل، مدل Decision Tree با MSE بسیار بالا، عملکرد ضعیف تری دارد.

تحلیل: R-squared Score (R^2)

بر اساس نمودار R^2 Score در شکل ۲، مقادیر R^2 برای مدل‌های مختلف به شرح زیر است:

Linear Regression تقریباً ۰.۴

Ridge Regression تقریباً ۰.۴

Lasso Regression تقریباً ۰.۴

Polynomial Regression تقریباً ۰.۴

KNN Regression تقریباً ۰.۲۵

SVR تقریباً ۰.۲۵

Decision Tree تقریباً -۰.۵ (مقدار منفی)

Random Forest تقریباً ۰.۲۵

Bayesian Linear تقریباً ۰.۴

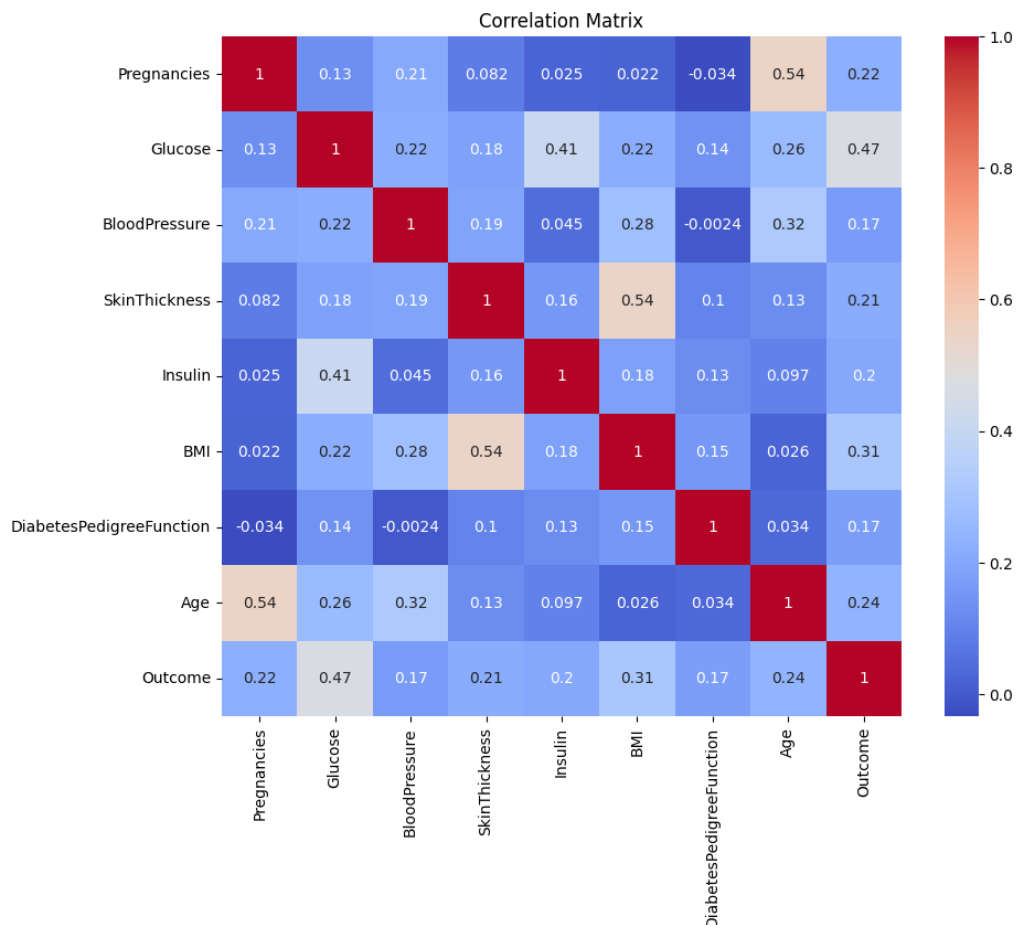
مدل‌های Linear Regression ، Ridge Regression ، Lasso Regression ، Polynomial Regression و Bayesian Linear با R^2 Score تقریباً ۰.۴، بیشترین نسبت از واریانس سطح گلوکز خون را نشان می‌دهند. این نشان دهنده توانایی نسبتاً خوب این مدل‌ها در کپچر کردن رابطه بین ویژگی‌های ورودی و متغیر هدف است. مدل‌های KNN ، Regression ، SVR و Random Forest R^2 Score کمتری دارند (تقریباً ۰.۲۵). قابل توجه است که مدل Decision Tree دارای R^2 Score منفی است که نشان دهنده عملکرد ضعیف‌تر آن نسبت به پیش‌بینی صرفاً میانگین سطح گلوکز خون است.

خلاصه مقایسه عملکرد:

با توجه به تحلیل MSE و R^2 Score ، مدل‌های Linear Regression ، Ridge Regression ، Lasso Regression ، Polynomial Regression ، Regression و Bayesian Regression عملکرد بهتری در پیش‌بینی سطح گلوکز خون در این دیتاست از خود نشان دادند. در میان آن‌ها، **linear Regression** با کمترین MSE، عملکرد اندکی بهتر دارد. مدل **Decision Tree ضعیف‌ترین** عملکرد را داشته است.

۲. تحلیل تاثیر ویژگی‌های مختلف بر پیش‌بینی سطح گلوکز خون:

تحلیل تاثیر هر ویژگی بر پیش‌بینی سطح گلوکز خون می‌تواند با بررسی ضرایب مدل‌های خطی (Linear, Ridge, Lasso, Bayesian) یا با استفاده از روش‌های تحلیل همبستگی انجام شود. با توجه به ماهیت دیتاست Pima Indians و همچنین رسم **correlation matrix** که ماتریس همبستگی مقدار سطح همبستگی میان هر جفت ویژگی‌ها را نشان می‌دهد که این مقادیر در بازه $[-1, 1]$ هستند. مقادیر نزدیک به 1 نشان دهنده همبستگی مثبت قوی، مقادیر نزدیک به -1 نشان دهنده همبستگی منفی قوی، و مقادیر نزدیک به 0 نشان دهنده عدم وجود ارتباط مشخص هستند که این کار را در آخر کد انجام داده‌ام.



تحلیل ویژگی‌ها:

1. Glucose (میزان گلوکز):

با ضریب همبستگی ۱ (مقدار هدف)، این ویژگی به صورت مستقیم بیشترین ارتباط را با مقدار هدف (Glucose) دارد و انتظار می‌رود بیشترین تأثیر را بر پیش‌بینی داشته باشد. این ویژگی اصلی‌ترین شاخص برای پیش‌بینی سطح گلوکز خون است.

2. Insulin (میزان انسولین سرم):

با ضریب همبستگی 0.41، میزان انسولین دومین ویژگی مهم در پیش‌بینی سطح گلوکز است. انسولین نقش کلیدی در تنظیم سطح قند خون دارد و رابطه مثبت با سطح گلوکز نشان‌دهنده اهمیت آن است.

3. BMI (شاخص توده بدنی):

ضریب همبستگی 0.22 نشان‌دهنده تأثیر مثبت BMI بر سطح گلوکز خون است. با توجه به ارتباط چاقی و اضافه وزن با دیابت، این ارتباط طبیعی به نظر می‌رسد.

4. Age (سن):

با ضریب همبستگی **0.26**، افزایش سن ارتباط مثبت با سطح گلوکز دارد. این ارتباط نشان دهنده این موضوع است که افراد مسن تر ممکن است مستعد افزایش سطح گلوکز و خطر ابتلا به دیابت باشند.

5. BloodPressure (فشار خون):

با ضریب همبستگی **0.22**، فشار خون نیز تأثیر نسبی بر سطح گلوکز دارد.

6. SkinThickness (ضخامت پوست):

ارتباط مثبت و اندک (ضریب همبستگی **0.18**) نشان دهنده تأثیر محدود این ویژگی در پیش‌بینی سطح گلوکز است.

7. DiabetesPedigreeFunction (تابع سابقه دیابت در خانواده):

با ضریب همبستگی **0.14**، این ویژگی نشان دهنده اثرات سابقه دیابت خانوادگی بر سطح گلوکز خون است. این ویژگی ممکن است در ترکیب با سایر ویژگی‌ها، اهمیت بیشتری پیدا کند.

8. Pregnancies (تعداد بارداری‌ها):

ضریب همبستگی **0.13** نشان دهنده ارتباط کم و مثبت تعداد بارداری‌ها با سطح گلوکز است. این ویژگی احتمالاً اثرات غیرمستقیم بر سطح گلوکز دارد.

تحلیل همبستگی نشان داد که ویژگی‌های **Glucose, Insulin, Age** و **BMI** بالاترین همبستگی را با هدف پیش‌بینی دارند و از اهمیت بیشتری برخوردارند. ویژگی‌هایی مانند **BloodPressure** و **SkinThickness** تأثیرات کمتری دارند، اما همچنان ممکن است در مدل‌سازی ترکیبی نقش ایفا کنند. ویژگی‌هایی مانند **Pregnancies** و **DiabetesPedigreeFunction** در این ماتریس ارتباط کمتری دارند و ممکن است با تنظیم دقیق مدل اهمیت بیشتری پیدا کنند.

۳. انتخاب بهترین مدل و بررسی اثرات تغییرات در پیش‌پردازش داده‌ها:

انتخاب بهترین مدل:

بر اساس نتایج تحلیل عملکرد **MSE** و **R² Score** مدل **Linear Regression** به عنوان بهترین مدل برای پیش‌بینی سطح گلوکز خون در این پروژه انتخاب می‌شود. این مدل با کمترین **MSE** و **R² Score** مشابه با سایر مدل‌های خطی برتر، عملکرد مناسبی را از خود نشان داده است. مدل‌های **Bayesian Regression**، **Ridge Regression**، **Lasso Regression** و **Polynomial Regression** نیز عملکرد خوبی دارند و می‌توانند به عنوان گزینه‌های جایگزین در نظر گرفته شوند.

بررسی اثرات تغییرات در پیش پردازش داده ها:

من در کد نتایج بالا را با استفاده از پیش پردازش **standard scaler** و حذف ۰ ها و **feature selection** بر اساس تحلیل ویژگی های مرحله قبل انجام داده ام و در پایین تاثیر آن ها را ذکر کردم , پس تمام کد و نتایج با استفاده از این روش ها بود ولی راه های دیگری نیز وجود دارد که در پایین ذکر کردم.

۱. استاندارد سازی (StandardScaler)

استاندارد سازی باعث می شود تمام ویژگی ها در یک مقیاس قرار گیرند (میانگین صفر و انحراف معیار یک). اگر یک ویژگی مانند سن مقادیر کوچکی داشته باشد ولی ویژگی دیگری مانند انسولین مقادیر بزرگی داشته باشد، مدل ممکن است به طور **unbalanced** وزن بیشتری به ویژگی بزرگ تر بدهد. با استاندارد سازی، مدل بین ویژگی ها تعادل برقرار می کند و در نتیجه، عملکرد و پایداری آن افزایش می یابد.

۲. حذف مقادیر صفر (Zero Removal)

در دیتاست Pima ، برخی ویژگی ها مثل فشار خون یا BMI دارای مقادیر صفر هستند، در حالی که در واقع این مقادیر نمی توانند صفر باشند (مثلاً هیچ فرد زنده ای نمی تواند فشار خون صفر داشته باشد). بنابراین، این صفرها به عنوان مقادیر گم شده تلقی می شوند و می توانند مدل را گمراه کنند. با حذف ردیف هایی که شامل این داده ها هستند، داده های واقعی تری در اختیار مدل قرار می گیرند. این کار معمولاً باعث افزایش دقت مدل می شود، هرچند که کاهش حجم داده ها نیز می تواند ریسک کمبود اطلاعات را به همراه داشته باشد.

۳. انتخاب ویژگی (Feature Selection)

تمام ویژگی های موجود در داده ها الزاماً برای پیش بینی هدف مفید نیستند. برخی از آن ها ممکن است نویز ایجاد کنند یا همبستگی کمی با خروجی داشته باشند. استفاده از روش های **feature selection** کمک می کند تنها ویژگی های مهم و مؤثر برای پیش بینی انتخاب شوند. این کار موجب ساده تر شدن مدل، کاهش احتمال بیش برآزش (overfitting) ، و در بسیاری از موارد بهبود دقت پیش بینی می شود.

۴. مقیاس بندی ویژگی ها:

مدل هایی مانند SVR و KNN به مقیاس ویژگی ها بسیار حساس هستند. اعمال روش های مقیاس بندی مانند Standard Scaler یا Min-Max Scaler می تواند دامنه مقادیر ویژگی ها را یکسان کرده و عملکرد این مدل ها را به طور قابل توجهی بهبود بخشد. در این پروژه، مقیاس بندی انجام شده است که عملکرد مدل های SVR و KNN را به سطحی قابل قبول رسانده است، هرچند نسبت به مدل های خطی ضعیف تر هستند.

مدیریت داده‌های پرت (Outliers):

وجود داده‌های پرت می‌تواند مدل‌های خطی را تحت تاثیر قرار داده و باعث افزایش خطای آن‌ها شود. روش‌های شناسایی و حذف یا تبدیل داده‌های پرت می‌توانند عملکرد این مدل‌ها را بهبود بخشند.

مهندسی ویژگی (Feature Engineering):

ایجاد ویژگی‌های جدید از ویژگی‌های موجود (مثلاً تعامل بین دو ویژگی یا تبدیل‌های غیرخطی) می‌تواند اطلاعات جدیدی را برای مدل فراهم کند و عملکرد آن را بهبود بخشد.

در این پروژه، با توجه به تفاوت عملکرد مدل‌های مختلف، می‌توان نتیجه گرفت که انتخاب روش‌های پیش‌پردازش به کار گرفته شده تاثیر مشخصی بر نتایج داشته است. برای مثال، اگر standard scaling انجام نمی‌شد، احتمالاً عملکرد مدل‌های SVR و KNN ضعیف‌تر از آنچه مشاهده می‌شود، بود.

۴. نتیجه‌گیری:

در این پروژه، با هدف پیش‌بینی سطح گلوکز خون در دیابت Pima Indians Diabetes، عملکرد چندین مدل رگرسیونی مورد ارزیابی قرار گرفت. نتایج نشان داد که مدل‌های Ridge Regression، Linear Regression، Lasso Regression، Polynomial Regression و Bayesian Regression عملکرد بهتری نسبت به مدل‌های KNN Regression، SVR و Decision Tree دارند. مدل Linear Regression با کمترین MSE، به عنوان بهترین مدل انتخاب شد. تحلیل تاثیر ویژگی‌ها نشان می‌دهد که ویژگی‌هایی مانند میزان گلوکز، BMI، و سن نقش مهمی در پیش‌بینی سطح گلوکز خون دارند. همچنین، اهمیت پیش‌پردازش داده‌ها در دستیابی به عملکرد مطلوب مدل‌ها مورد بررسی قرار گرفت.