سوال ١:

مراحل عملکرد کد به شرح زیر است:

1 .ابتدا، کتابخانههای مورد نیاز را وارد میکنیم، از جمله pandas برای کار با دادههای جدولی، sklearn برای استفاده از المانتهای مربوط به ماشین بردار پشتیبان و معیارهای ارزیابی، و matplotlib

2 .سپس، داده ها از یک فایل CSV با استفاده از تابع 'read_csv' در pandas خوانده می شوند. فایل CSV دو ستون دارد "label" :که برچسب همگونی یا هرزنامه را نشان می دهد و "message" که متن پیام ها را حاوی است.

3 .ستون "message" به عنوان متغیر ورودی (X) و ستون "label" به عنوان متغیر خروجی (y) در نظر گرفته می شوند. برچسبهای "ham" به مقدار 0 و برچسبهای "spam" به مقدار 1 نگاشت می شوند.

4 .سپس، بردار ویژگی TF-IDF بر روی متن پیامها اعمال میشود. بردار ویژگیهای TF-IDF اطلاعات مربوط به وزن و اهمیت کلمات در هر پیام را نشان میدهند.

5 .یک مدل SVM با کرنل خطی ایجاد می شود SVM .یک الگوریتم یادگیری ماشینی است که برای دسته بندی داده ها استفاده می شود و با استفاده از بردار های ویژگی، مرزهای تصمیمگیری بین دسته ها را تعیین میکند.

6 .مدل روی بردار ویژگیهای TF-IDF آموزش داده میشود تا بتواند پیشبینیهای همگونی/هرزنامه صحیحی را ارائه دهد.

7 با استفاده از مدل آموزش دیده، برچسبهای پیش بینی شده برای همه نمونهها محاسبه می شود.

8 .سپس، دقت مدل با استفاده از معیار دقت (accuracy) محاسبه میشود. دقت نسبت تعداد پیشبینیهای صحیح به کل نمونهها را نشان میدهد.

9 در ادامه، منحنی مشخصه عملکرد گیرنده (Receiver Operating Characteristic - ROC) برای ارزیابی کارایی مدل با استفاده از شاخصهای تشخیص همگونی و هرزنامه محاسبه می شود. این منحنی توزیع برچسبها در تمام حدود تصمیمگیری ممکن را نشان می دهد.

10 .در انتها، نمودار ROC با استفاده از کتابخانه 'matplotlib' رسم می شود. این نمودار نشان میدهد که با تغییر آستانه تصمیمگیری، نرخ اشتباه مثبت و نرخ صحیح مثبت چگونه تغییر میکنند.

سوال ٢:

مراحل عملکرد کد به شرح زیر است:

1 .ابتدا، کتابخانههای مورد نیاز را وارد میکنیم، از جمله pandas برای کار با دادههای جدولی، sklearn برای استفاده از المانتهای مربوط به ماشین بردار پشتیبان و اسکیلر استاندارد، و matplotlib

2 سپس، داده ها از یک فایل CSV با استفاده از تابع `read_csv` در pandas خوانده می شوند. فایل CSV شامل ویژگی های مربوط به بیماران دیابتی مانند سن، ضربان قلب، فشار خون و غیره است، و ستون "Outcome" برچسب دیابت را نشان می دهد.

3 .ستون "Outcome" به عنوان متغیر خروجی (y) در نظر گرفته می شود و سایر ستون ها به عنوان متغیر های ورودی.(X)

4 دادههای ورودی و خروجی به دو بخش آموزشی و آزمون تقسیم میشوند با استفاده از تابع 'train_test_split'از .sklearn در اینجا، 80٪ دادهها برای آموزش و 20٪ دادهها برای آزمون استفاده می شود.

5 اسکیلر استاندارد بر روی دادههای آموزشی و آزمون اعمال میشود تا ویژگیها را به مقادیر استاندارد مبدل کند. این مرحله باعث میشود ویژگیها بین دادهها قابل قیاس و مقایسه باشند.

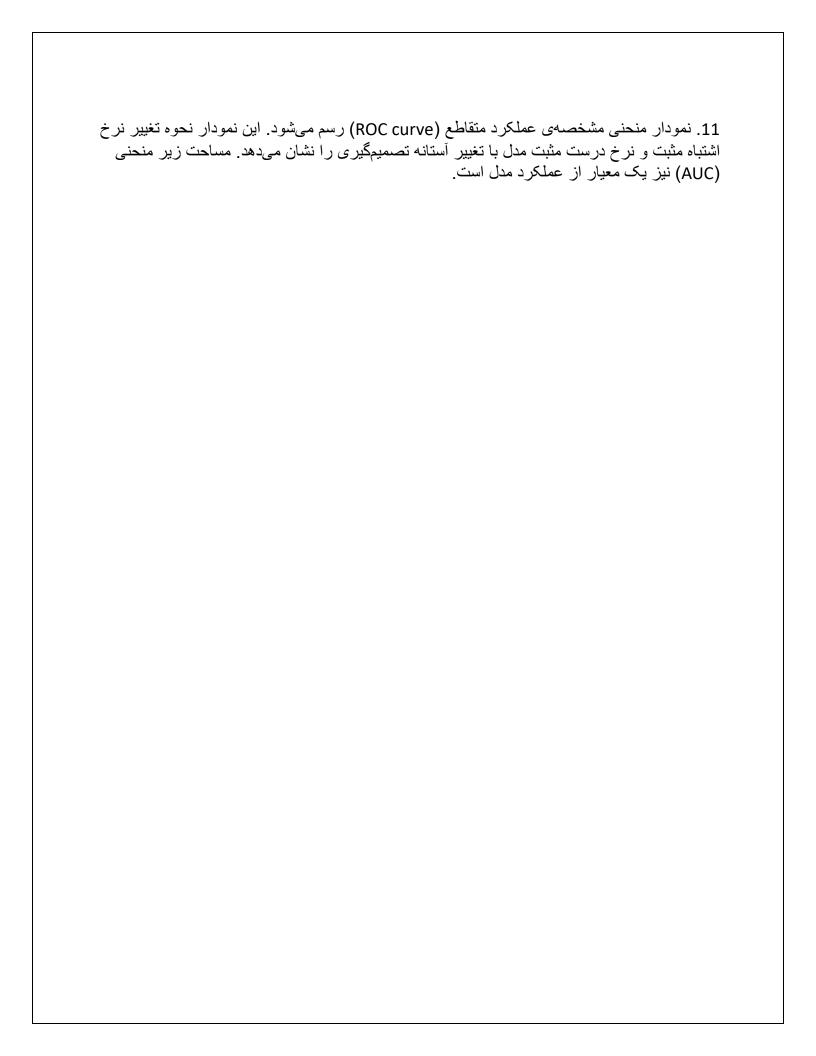
6 .یک مدل SVM با کرنل شعاعی (RBF) ایجاد می شود SVM .یک الگوریتم یادگیری ماشینی است که برای دسته بندی داده ها استفاده می شود و با استفاده از مرزهای تصمیمگیری بین دسته ها، داده ها را دسته بندی می کند.

7 مدل روی داده های استاندار د آموزش داده می شود تا بتواند پیشبینی دیابت صحیحی را ارائه دهد.

8 با استفاده از مدل آموزش دیده، برچسبهای پیش بینی شده برای دادههای آزمون محاسبه می شود.

9 دقت، بازیابی و امتیاز F1 مدل با استفاده از معیارهای مربوطه محاسبه میشوند. دقت (accuracy) نسبت تعداد نمونههای نسبت تعداد نمونههای نسبت تعداد نمونههای درست پیشبینی شده از یک کلاس به کل نمونههای واقعی این کلاس است و امتیاز (F1-score) F1 (F1-score) معیاری است که ترکیبی از دقت و بازیابی است.

10 .ماتریس در همریختگی (confusion matrix) بر اساس برچسبهای واقعی و برچسبهای پیشبینی شده برای دادههای آزمون رسم میشود. در اینجا، ماتریس در همریختگی نشان میدهد که چند داده به درستی به هر دسته تعلق میگیرد.



سو ال ٣:

این برنامه یک مدل پیشبینی قیمت مسکن با استفاده از روش پشتیبان وکتور رگرسیون Support) (Vector Regressionایجاد میکند. ابتدا دادههای مسکن از منبع مشخص شده در آدرس 'url' دریافت میشوند و در یک شیء DataFrame در Pandas ذخیره میشوند.

سپس دادههای ورودی (X) و خروجی (y) برای مدل ساخته می شود. دادههای ورودی شامل تمامی ستون های داده به جز ستون) "MEDV" قیمت مسکن) است و داده های خروجی شامل ستون "MEDV" هستند.

سپس داده ها به دو مجموعه آموزش و آزمون تقسیم می شوند (80 در صد برای آموزش و 20 در صد برای آرمون. (

سپس دادههای آموزش مقیاسبندی میشوند با استفاده از مقیاسبندی استاندارد (StandardScaler) تا تأثیر متغیر های ورودی را به یک مقیاس مشابه برساند. دادههای آزمون نیز با استفاده از همان مقیاسبندی استاندارد مقیاسبندی میشوند.

سپس یک مدل SVR با هسته RBF ایجاد می شود و روی داده های آموزش آموزش داده می شود.

مقادیر پیشبینی شده برای داده های آزمون به دست می آیند و خطای میانگین مربعات (MSE) بین مقادیر پیش بینی شده و واقعی محاسبه می شود و چاپ می شود.

در نهایت، نموداری از قیمتهای واقعی و پیشبینی شده رسم میشود تا بتوان تطابق مدل با دادههای واقعی را مشاهده کرد.

سوال ۴:

- ** آمادهساز ي دادهها: **
- این کد از مجموعه داده های 20 خبرگروه (20 (Newsgroups استفاده میکند که یک مجموعه داده رایج برای طبقه بندی متن است. این مجموعه داده به دو بخش آموزش و آزمایش تقسیم می شود.
 - ** .استخراج ویژگیها: **
 - کد از TfidfVectorizer استفاده میکند تا متن را به یک فرمت عددی تبدیل کند. این فرمت عددی می تواند به عنوان و رودی به مدل یادگیری ماشین داده شود.
 - ** .آمو ز ش مدل: **
- کد از یک طبقهبند SVM خطی (LinearSVC) استفاده میکند و آن را با استفاده از داده های آموزشی آموزشی آموزش میدهد.
 - ** .ارزیابی مدل: **
 - کد از مدل آموزش داده شده برای پیشبینی برچسبهای کلاس دادههای آزمایشی استفاده میکند.
 - دقت مدل با استفاده از تابع accuracy score محاسبه می شود.
 - ماتریس سردرگمی (confusion matrix) برای نشان دادن عملکرد مدل ترسیم می شود.
 - گزارش طبقهبندی (classification report) که شامل صحت، یادآوری، F1-score پشتیبانی برای هر کلاس است، چاپ میشود.
 - ** نمایش گزارش طبقهبندی: **
- کد یک نمودار جدید ایجاد میکند و گزارش طبقهبندی را به صورت متن روی آن نمایش میدهد. این نمایش خواناتر و بصری تر از چاپ ساده گزارش است.

سوال ۵:

1** بیشیر دازش دادهها: **

- اسکریپت فایل CSV را از URL ارائه شده میخواند و در یک DataFrame پاندا به نام 'df' ذخیره میکند.
- سپس 5 سطر اول DataFrame را چاپ میکند و همچنین ابعاد دیتاست (تعداد ردیفها و ستونها) را نمایش میدهد.
 - متغیر هدف به عنوان ''SMQ020' انتخاب می شود، که احتمالاً یک ستون دستهای است.
 - اسکریپت سپس هر ردیف با مقدار گمشده را حذف میکند.
 - متغیر هدف را با استفاده از `LabelEncoder از scikit-learn رمزگذاری میکند.

2** .تقسيم دادهها به آموزش و آزمايش: **

- اسکریپت داده ها را به مجموعه های آموزش و آزمایش با استفاده از تابع 'train_test_split' از scikit-learn' از scikit-learn
 - اندازه آزمایش به 20٪ از کل داده ها تنظیم شده است.

3** مقياسسازي ويژگيها: **

- اسکریپت از `StandardScaler از scikit-learn برای استاندارد کردن ویژگیهای ورودی (X) در مجموعههای آموزش و آزمایش استفاده میکند.

4** . آموزش مدل: **

- اسکریپت یک نمونه از طبقهبند ماشین بردار پشتیبان (SVM) از scikit-learn با هسته 'rbf' ایجاد میکند.
 - مدل SVM سپس بر روی دادههای آموزش برازش داده میشود.

5** .ارزيابي مدل: **

- اسکریپت با استفاده از مدل SVM آموزش دیده، پیشبینیهایی بر روی دادههای آزمایش انجام میدهد.
 - دقت، بازیابی و نمره F1 عملکرد مدل را محاسبه میکند.

6** .ترسيم تصوير: **

- اسکریپت یک نمودار ماتریس ابهام ایجاد میکند تا عملکرد طبقهبندی را نمایش دهد.
- همچنین یک منحنی عملکرد دریافت-عملیاتی (ROC) ایجاد میکند و معیار سطح زیر منحنی (AUC) را محاسبه میکند.

.