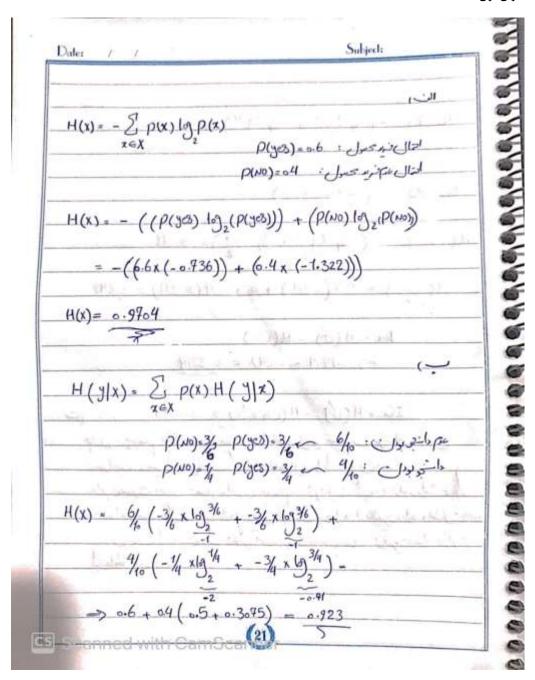
تمرین سری یکم مبانی یادگیری ماشین

معصومه پاسبانی 99243022

سوالات تئوري

سوال اول)



Dale: / /

H(c)9) = - (1/4 18) + 3/4 6) => - (-0.5-0.31125)

H(disto) = - (2 log 1 + 0) = 0

H()= - (3/4 lg 3/4 + 1/4 lg 1/4)= 0.811

H((1)) = 0.4 (0.811) + 0 + 0.4 (0.811) = 0.649

IG = H(b) - H(v) → 0.9704 - 0.649 = 0.3214

1. مسئله تعداد زياد ابعاد

در KNN ، فاصله بین نقاط (معمولاً فاصله اقلیدسی) محاسبه میشود. با افزایش تعداد ابعاد، فاصله بین نقاط دادهها نیز به شدت افزایش مییابد، و دادهها به نوعی پراکنده تر میشوند.به طوری که الگوریتم KNN به طور موثر نمی تواند بر روی دادهها عمل کند و این موضوع باعث میشود که تشخیص نزدیکی نقاط داده دشوار تر شود و در نتیجه، دقت مدل کاهش یابد.

در ابعاد زیاد، نقاط داده ممکن است از هم فاصله بگیرند و مدل نمی تواند تشخیص دقیقی داشته باشد، زیرا همه نقاط به نظر مشابه خواهند رسید، همچنین حاسبه فاصله بین نقاط نیاز به زمان بیشتری دارد.

برای رفع این مشکل از روش های کاهش ابعاد مانند PCA یا LDA میتوان استفاده کرد تا ویژگی های اصلی دادهها حفظ شوند و فضای جستجو کوچکتر شود؛ همچنین انتخاب ویژگیهای مهم و حذف ویژگیهای غیرضروری میتواند باعث کاهش ابعاد شود و تأثیر کمتری روی دقت مدل بگذارد.

2. مسئله مقياس بندى ويژگى ها

در KNN ، چون برای پیدا کردن نزدیکترین همسایگان از فاصلهها (مثلاً فاصله اقلیدسی) استفاده می شود، مقیاس ویژگیها بسیار مهم است. اگر ویژگیها مقیاسهای متفاوتی داشته باشند، ویژگیهایی با مقیاس بزرگتر به طور غیرعادلانه تأثیر بیشتری بر محاسبه فاصله خواهند داشت.

اگر ویژگیها مقیاسهای متفاوتی داشته باشند، ویژگیهایی که دارای مقیاس بزرگتر هستند، میتوانند تأثیر زیادی در انتخاب نزدیکترین همسایگان داشته باشند؛ همچنین تفاوت مقیاس بین ویژگیها باعث میشود که برخی از ویژگیها بیشتر از سایرین در تشخیص کلاسها تأثیر بگذارند.

برای رفع این مشکل میتوان از تکنیک های مقیاس بندی مانند standardization ،min-max scaling و normalization استفاده کرد که هر کدام به ترتیب به این صورت عمل میکنند: مقیاس داده ها را به بازه [0,1] تغییر داده، دادهها را به گونهای تغییر میدهد که میانگین آنها صفر و انحراف معیار یک باشد و دادهها را به صورت نسبی و در یک مقیاس خاص استاندارد کنید تا تأثیر ویژگیهای بزرگتر کاهش یابد.

سوال سوم)

Random Forest الگوریتمی است که از مجموعهای از درختهای تصمیم برای بهبود دقت و پایداری پیشبینیها استفاده می کند.

ابتدا برای ساخت هر درخت در جنگل تصادفی، نمونهای از دادههای آموزشی بهطور تصادفی با جایگزینی انتخاب می شود، یعنی امکان دارد دادهای چندین بار انتخاب شود. سپس برای هر درخت، در هر گره تنها از زیرمجموعهای از ویژگیها بهطور تصادفی استفاده می شود. این کار از ایجاد درختهایی که بیش از حد به یک ویژگی خاص وابسته هستند جلوگیری می کند و در نهایت در طبقه بندی با رای گیری label ای که بیشترین تعداد درختها آن را پیشبینی کردهاند به عنوان خروجی نهایی در نظر گرفته می شود.

مزایای Random Forest نسبت به درخت تصمیم:

Random Forest به دلیل ترکیب چندین درخت، دقت بالاتری نسبت به درخت تصمیم منفرد دارد. این ترکیب باعث می شود که ضعفهای درختهای منفرد کاهش پیدا کند و نتایج نهایی بهینه تر باشند. درختهای تصمیم به دلیل ساختارشان overfitting هستند، به ویژه زمانی که درخت عمیق و پیچیده می شود. در Random Forest، با ترکیب چندین درخت و ایجاد تنوع در آنها، این مشکل کاهش می یابد و مدل به دادههای آموزشی بیش از حد وابسته نمی شود. به دلیل اینکه پیش بینی براساس رأی گیری یا میانگین چندین درخت به دست می آید، Random Forest نسبت به نویز موجود در دادهها مقاوم تر است و عملکرد پایدار تری دارد.

Random Forest نه تنها در طبقه بندی بلکه در رگرسیون نیز به خوبی عمل می کند و یک روش چندمنظوره است.

سوال چهارم)

Weighted KNN: در الگوریتم استانداردKNN ، کلاسبندی بر اساس رأی گیری از نزدیک ترین همسایگان انجام می شود. هر همسایه به طور مساوی اهمیت دارد و تعداد همسایگان هر کلاس تعیین می کند که نمونه به کدام کلاس تخصیص داده شود. اما در Weighted KNN ، نزدیک ترین همسایگان با توجه به فاصله شان از نقطه موردنظر وزن دهی می شوند، به طوری که هرچه همسایه به نقطه نزدیک تر باشد، وزن بالاتری دارد.

این الگوریتم ابتدا فاصله بین نمونه جدید و تمام نمونه های مجود در مجموعه داده را محاسبه می کند، سپس نزدیک ترین k همسایه را انتخاب می کند. برای هر همسایه، weight = 1/distance محسابه می شود و در نهایت کلاس نمونه جدید با رای گیری وزنی تعیین می شود.

مزايا:

همسایههای نزدیکتر که اهمیت بیشتری در طبقهبندی دارند، تأثیر بیشتری روی پیشبینی نهایی دارند، که باعث افزایش دقت طبقهبندی میشود. این روش به دادههای توزیع غیریکسان و با فواصل مختلف حساس تر است و نسبت به KNN ساده که تمامی همسایگان را به یک اندازه در نظر می گیرد، انعطاف پذیری بهتری دارد.

NCA: الگوریتم NCAیک روش یادگیری بر اساس متریک است که به منظور افزایش دقت الگوریتم KNN طراحی شده است. در این روش، از بهینه سازی یک تابع هدف استفاده می شود که هدف آن یادگیری یک نگاشت خطی برای داده ها است. این نگاشت خطی، داده ها را به شکلی تغییر می دهد که همسایگان از نظر فاصله در یک فضای جدید بهینه سازی شوند.

NCA سعی می کند یک ماتریس تبدیل بیابد که دادهها را به یک فضای جدید منتقل کند به طوری که همسایگان مربوط به هر نمونه در این فضای جدید، طبقهبندی دقیقی داشته باشند. تابع هزینهای که در این روش بهینهسازی می شود، احتمال پیشبینی درست کلاس توسط نزدیک ترین همسایگان را اندازه گیری می کند. این احتمال بر اساس فاصلهها محاسبه می شود و NCA سعی دارد که فاصله نمونهها را در این فضای جدید به گونهای تغییر دهد که همسایگان نزدیک از همان کلاس باشند. با استفاده از روشهای بهینهسازی، این ماتریس نگاشت خطی به دست می آید، سپس دادهها به این فضای جدید تبدیل می شود.

مزايا:

با یادگیری نگاشتی که دادهها را به فضای جدیدی منتقل میکند، الگوریتم KNN در این فضا دقت بالاتری دارد. این روش به خصوص در مسائل پیچیده و دادههای با توزیع غیرهمگن عملکرد بهتری دارد، زیرا ساختار دادهها را بهینهتر تنظیم میکند.

سوالات عملي

سوال اول)

```
[55] X = data.iloc[:, :-1].values
    y = data.iloc[:, -1].values

    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

[56] depths = range(1, 70)
    depth_acc_dt, depth_prec_dt, depth_rec_dt = [], [], []
    depth_acc_rf, depth_prec_rf, depth_rec_rf = [], [], []

max_features = range(1, X.shape[1] + 1)
    feature_acc_dt, feature_prec_dt, feature_rec_dt = [], [], []

feature_acc_rf, feature_prec_rf, feature_rec_rf = [], [], []
```

ابتدا داده ها را به دو بخش train و test تقسیم کرده که این تقسیم بندی برای ارزیابی عملکرد مدل ها بر روی دادهای که در مرحله آموزش دیده نشدهاند، استفاده می شود.

```
for depth in depths:
57]
        dt = DecisionTreeClassifier(max depth=depth)
        dt.fit(X_train, y_train)
        y_pred_dt = dt.predict(X_test)
        depth_acc_dt.append(accuracy_score(y_test, y_pred_dt))
        rf = RandomForestClassifier(max_depth=depth)
        rf.fit(X train, y train)
        y_pred_rf = rf.predict(X_test)
        depth_acc_rf.append(accuracy_score(y_test, y_pred_rf))
58] for depth in depths:
        # Decision Tree
        dt = DecisionTreeClassifier(max_depth=depth)
        dt.fit(X train, y train)
        y pred_dt = dt.predict(X_test)
        depth_prec_dt.append(precision_score(y_test, y_pred_dt))
        # Random Forest
        rf = RandomForestClassifier(max_depth=depth)
        rf.fit(X_train, y_train)
        y_pred_rf = rf.predict(X_test)
        depth_prec_rf.append(precision_score(y_test, y_pred_rf))
```

سپس مدلهای درخت تصمیم و جنگل تصادفی با مقادیر مختلف max_depth آموزش داده می شوند. این پارامتر تعیین می کند که حداکثر عمق درخت چه مقدار باشد و تاثیر مستقیمی بر صحت، precisionو precisionمدلها دارد. برای هر مقدار از max_depth، سه متریک زیر محاسبه می شود:

(Accuracy): نسبت نمونههای درست پیش بینی شده به کل نمونهها.

:(Precision)نسبت نمونههای مثبت صحیح پیشبینی شده به کل نمونههای مثبت پیشبینی شده.

:(Recall)نسبت نمونههای مثبت صحیح پیشبینی شده به کل نمونههای مثبت واقعی.

```
for depth in depths:
   # Decision Tree
   dt = DecisionTreeClassifier(max depth=depth)
   dt.fit(X_train, y_train)
   y pred dt = dt.predict(X test)
   depth rec dt.append(recall score(y test, y pred dt))
   # Random Forest
   rf = RandomForestClassifier(max depth=depth)
   rf.fit(X train, y train)
   y pred rf = rf.predict(X test)
   depth rec rf.append(recall score(y test, y pred rf))
for feature in max features:
   # Decision Tree
   dt = DecisionTreeClassifier(max features=feature)
   dt.fit(X train, y train)
   y pred dt = dt.predict(X test)
   feature acc dt.append(accuracy score(y test, y pred dt))
   # Random Forest
   rf = RandomForestClassifier(max features=feature)
   rf.fit(X train, y_train)
   y pred rf = rf.predict(X test)
   feature acc rf.append(accuracy score(y test, y pred rf))
```

```
[61] for feature in max features:
         # Decision Tree
         dt = DecisionTreeClassifier(max features=feature)
         dt.fit(X train, y train)
         y pred dt = dt.predict(X test)
         feature_prec_dt.append(precision_score(y_test, y_pred_
         # Random Forest
         rf = RandomForestClassifier(max features=feature)
         rf.fit(X train, y train)
         y_pred_rf = rf.predict(X_test)
         feature_prec_rf.append(precision_score(y_test, y_pred_
[62] for feature in max features:
         # Decision Tree
         dt = DecisionTreeClassifier(max features=feature)
         dt.fit(X train, y train)
         y_pred_dt = dt.predict(X test)
         feature rec dt.append(recall score(y test, y pred dt))
         # Random Forest
         rf = RandomForestClassifier(max features=feature)
         rf.fit(X train, y train)
         y pred rf = rf.predict(X test)
         feature_rec_rf.append(recall score(y test, y pred rf))
```

پارامتر max_featuresتعداد ویژگیهایی که مدل در هر بار انشعاب درخت استفاده میکند را تعیین میکند. این بخش مشابه قسمت قبل است، اما پارامتر max_features به جای max_depth تغییر میکند. برای هر مقدار از precision، صحت max_featuresو المصابه میشوند.

سپس نمودارها را برای تحلیل تغییرات صحت، دقت و فراخوانی نسبت به max_depth و max_features رسم کرده. این نمودارها کمک میکنند که تاثیر این پارامترها بر عملکرد مدلها را بهتر مشاهده و مقایسه کنیم. هر نمودار شامل خطی برای مدل درخت تصمیم و خط دیگری برای مدل جنگل تصادفی است.

توازنسازی با استفاده از نمونهبرداری مجدد Oversampling وUndersampling

Oversampling :در این روش، دادههای کلاس کمنمونه چندین بار تکرار میشوند یا نمونههای مصنوعی برای این کلاس SMOTE ایجاد می شود تا اندازه کلاسهای مختلف به یکدیگر نزدیک شود. یکی از روشهای پرکاربرد در این دسته SMOTE ایجاد می شود تا اندازه کلاسهای (Synthetic Minority Over-sampling Technique) است که با تولید نمونههای مصنوعی جدید برای کلاسهای کم نمونه، توازن ایجاد می کند.

Undersampling: در این روش، تعداد نمونههای کلاس غالب کاهش مییابد تا اندازه کلاسها متعادل شود. به عنوان مثال، می توان تعدادی از نمونههای کلاس غالب را به صورت تصادفی حذف کرد.

Cost sensitive learning

این روش به جای متعادل کردن دادهها، الگوریتم یادگیری را اصلاح می کند تا طبقهبندی نادرست کلاس اقلیت را به شدت جریمه کند. در طبقهبندی کنندههایی مانند درختهای تصمیم، ماشینهای بردار پشتیبان یا شبکههای عصبی، می توان هزینه های بالاتری را برای طبقهبندی اشتباه نمونههای کلاس اقلیت، بهبود عملکرد برای دادههای نامتعادل تعیین کرد.

Ensemble techniques

تکنیکهایی مانند جنگل تصادفی متعادل یا EasyEnsemble روشهای مجموعه را با استراتژیهای نمونه گیری ترکیب می کنند. برای مثال EasyEnsemble چندین مدل را در مجموعه داده های متعادل با کم نمونه برداری از کلاس اکثریت برای هر مدل آموزش می دهد. سپس نتایج تجمیع می شوند و عملکرد در طبقات اقلیت را افزایش می دهند.

```
7] le = LabelEncoder()
  values = ['gender', 'ever_married', 'work_type', 'Residence_type', 'smoking_status']
  for col in values:
    data[col] = le.fit_transform(data[col])
  data.head()
```

ابتدا ستون های غیرعددی با استفاده از labelencoder به داده های عددی تبدیل میشوند. سپس با ایتفاده از simpleImputer داده های میس شده در ویژگی های عددی با میانگین گیری پر میشوند.

```
[8] from sklearn.impute import SimpleImputer
    x = data.drop('stroke', axis=1)
    imputer = SimpleImputer(strategy='mean')
    x = pd.DataFrame(imputer.fit_transform(x), columns=x.columns, index=x.index) #convert b
    y = data['stroke']

Double-click (or enter) to edit

[9] smote = SMOTE(random_state=42)
    x_balanced, y_balanced = smote.fit_resample(x, y)
    print("Balanced Data Shape:", x_balanced.shape)
    print("Class Distribution After SMOTE:", Counter(y_balanced))

$\frac{1}{2}$

Balanced Data Shape: (85234, 11)
    Class Distribution After SMOTE: Counter(\{0: 42617, 1: 42617\})

[10] x_train, x_test, y_train, y_test = train_test_split(x_balanced, y_balanced, test_size=0)

[11] print("Train Set Shape:", x_train.shape)
    print("Test Set Shape:", x_test.shape)
```

از تکنیک SMOTE برای متعادل سازی کلاس ها استفاده شده است. این تکنیک نمونه های جدیدی را برای کلاس های کم تعداد ایجاد می کند تا توزیع داده ها متعادل شود. پس از اعمال SMOTE ، توزیع کلاس ها چاپ می شود تا مطمئن شویم داده ها به درستی متعادل شده اند

دادههای متعادل به نسبت ٪۶۷ و ٪۳۳ به مجموعههای آموزش و تست تقسیم میشوند. این تقسیمبندی کمک میکند تا از بخشی از دادهها برای آموزش مدل استفاده کرده و عملکرد مدل را روی بخشی دیگر ارزیابی کنیم.

```
def knn(x train, y train, x test, k=9):
   x train np = x train.values
   x test np = x test.values
   y_train_np = y_train.values
   max labels = []
   for test row in x test np:
       distances = np.linalg.norm(x train np - test row, axis=1)
       k indices = np.argsort(distances)[:k]
       k nearest labels = y train np[k indices]
       most common label = Counter(k nearest labels).most common(1)[0][0]
       max labels.append(most common label)
   return np.array(max labels)
prediction = knn(x_train, y train, x test)
print(classification report(y test, prediction))
                precision
                               recall f1-score
                                                     support
                      0.94
             0
                                 0.77
                                             0.85
                                                       14074
             1
                      0.81
                                 0.95
                                             0.87
                                                       14054
                                             0.86
     accuracy
                                                       28128
                                             0.86
   macro avg
                      0.87
                                 0.86
                                                       28128
                                 0.86
weighted avg
                     0.87
                                             0.86
                                                       28128
accuracy = accuracy score(y test, prediction)
accuracy
0.860494880546075
```

سپس تابع knn را به صورت دستی پیاده سازی میکنیم. برای هر نمونه در مجموعه تست، فاصله اقلیدسی آن نسبت به تمام نمونههای مجموعه آموزش محاسبه شده و k نزدیک ترین همسایهها انتخاب می شوند. لیبل رایج ترین همسایهها به عنوان پیشبینی آن نمونه تست انتخاب می شود. پس از پیشبینی کلاسها، عملکرد مدل با استفاده از classification_report و accuracy_scoreرزیابی می شود.

```
from sklearn.preprocessing import StandardScaler, PolynomialFeatures

scaler = StandardScaler()
    x_train_scaled = scaler.fit_transform(x_train)
    x_test_scaled = scaler.transform(x_test)

poly = PolynomialFeatures(degree=2, interaction_only=True, include_bias=False)
    x_train_poly = poly.fit_transform(x_train_scaled)
    x_test_poly = poly.transform(x_test_scaled)
```

دادهها با استفاده از StandardScaler استانداردسازی میشوند تا ویژگیها میانگین صفر و واریانس یک داشته باشند.

ویژگیهای جدید چندجملهای (پلینومیال) با درجه ۲ با استفاده از PolynomialFeatures ایجاد میشوند، که شامل اثرات تعاملی بین ویژگیها میباشد.

```
metrics = ['euclidean', 'manhattan']
best k = None
best score = 0
best metric = None
for metric in metrics:
    for k in range(1,15):
        knn model = KNeighborsClassifier(n neighbors=k, metric=metric)
        knn_model.fit(x train, y train)
        prediction = knn model.predict(x test)
        score = accuracy score(y test, prediction)
        if score > best score:
            best score = score
            best k = k
            best metric = metric
        print(f"Metric: {metric}, k: {k}, Accuracy: {score:.4f}")
Metric: euclidean, k: 1, Accuracy: 0.9824
Metric: euclidean, k: 2, Accuracy: 0.9817
Metric: euclidean, k: 3, Accuracy: 0.9688
Metric: euclidean, k: 4, Accuracy: 0.9695
Metric: euclidean, k: 5, Accuracy: 0.9585
Metric: euclidean. k: 6. Accuracy: 0.9599
```

از KNeighborsClassifier برای پیادهسازی الگوریتم KNN استفاده کرده. این مدل از scikit-learn قابلیت استفاده از متریکهای فاصله مختلف (اقلیدسی و منهتن) و مقادیر متریکهای فاصله مختلف (اقلیدسی و منهتن) و مقادیر

مختلف k از ۱ تا ۱۴ آزمایش می شود تا بهترین ترکیب k و متریک فاصله پیدا شود. نتایج برای هر ترکیب k و متریک فاصله چاپ می شوند. چاپ می شوند و در نهایت بهترین مقدار k و متریک فاصله با بالاترین دقت گزارش می شوند.