

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

موضوع: مینی پروژه - پاسخ سوال اول

درس مربوط: یادگیری ماشین لرنینگ

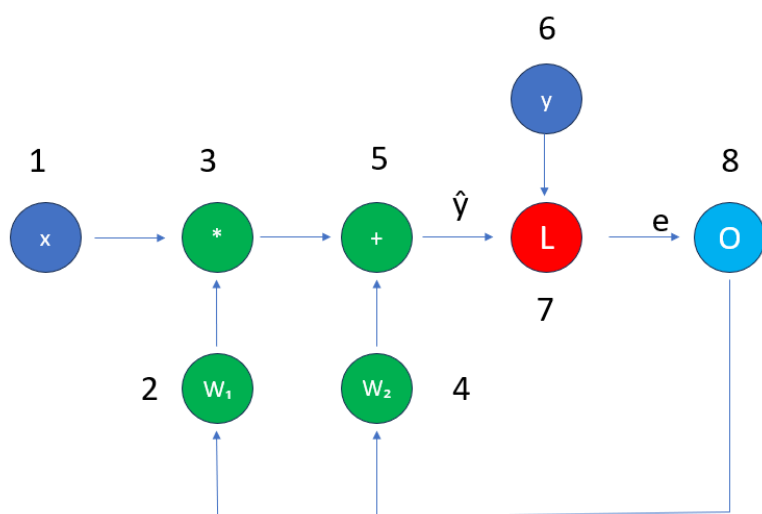
اساتید راهنما: آقایان دکتر علیاری و احمدی

دانشجو: ابوذر بختیاری برزیده

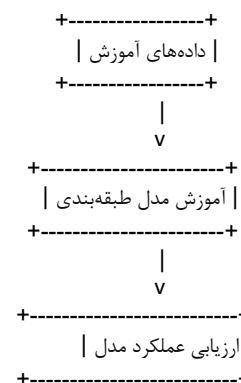
شماره دانشجویی: 4021634202

فروردین 1403

پاسخ سوال اول



1- یک دیاگرام بلوکی معمولاً برای نمایش فرآیند آموزش و ارزیابی یک مدل طبقه‌بندی خطی مفید است. این دیاگرام بلوکی می‌تواند به شکل زیر باشد:



قسمت 1: در این بلوک دیاگرام ما چندین قسمت مختلف داریم که هر کدام را به اختصار و در انتها توضیحات نهایی را می‌دهیم.

1- در این قسمت مقدار ورودی ما می‌باشد که به چه صورت وارد می‌شود که ممکن است یک کلاس و یا چند کلاس باشد که ما می‌توانیم از مقدار ورودی آن‌ها را شناسایی کنیم. این مرحله شامل داده‌هایی است که برای آموزش مدل استفاده می‌شود. این داده‌ها به دو بخش تقسیم می‌شوند: داده‌های ورودی (ویژگی‌ها) و برچسب‌ها (کلاس‌ها).

2- این قسمت معادله‌ای وارد بلوک دیاگرام ما می‌شود، این معادله می‌تواند معادله خط، معادله خطوط و یا معادله صفحه‌ای باشد.

3- ضرب مقدار ورودی و معادله موجود ما می‌باشد.

4- در این قسمت ما معادله از قبل را داریم که ممکن است تصادفی به دست آمده باشد و یا خیر در اطلاعات مسئله باشد.

5- در این قسمت ما جمع معادله تصادفی اولیه با حاصل ضرب مقداری ورودی و معادله موجود ما است در این زمان در حاصل جمع این دو جمع شونده، ما انتظار داریم که متغیر y بدست بیاید ولی در این قسمت y^{\wedge} بدست آمده است.

6- مقدار انتظار ما برای به وقوع پیوستن.

7- در این زمان باید اخلاف مقدار بدست آمده و مقدار مورد انتظار ما بدست آید که این کار بر عهده تابع اتلاف می باشد که توابع مختلفی برای این کار موجود می باشند.

8- در ادامه با استفاده از اختلاف بدست آمده ما می توانیم از روش های بهینه سازی برای مثال تابع گرادیان نزولی استفاده کنیم و معادلات در قسمت های دو و چهار را اپدیت کنیم تا در مرحله بعدی جواب بهتر ، نزدیک تر و حتی خود مقدار انتظار ما را برآورده کند و این روند را تا زمانی ادامه دهد که بهترین پاسخ ممکن را برای ما بدست آورد.

قسمت 2:

در کد زیر ما با استفاده از کتابخانه های پایتون به یک کتبخانه به نام (sklearn) استفاده می کنیم که بخشی از این کتبخانه دارای پایگاه داده های مختلفی می باشد از جمله breast cancer و make classification این پایگاه داده برای ساختن داده های طبقه بندی شده استفاده می شود.

- توضیح کوتاهی در مورد make classification (ساختن داده های طبقه بندی شده) : این تابع به ما اجازه می دهد تا داده های مصنوعی برای مسائل طبقه بندی شده تولید کنیم و این تابع در اکثر موارد در مسائل یادگیری ماشین کاربرد دارد که دارای بخش ها و آرگومان های مهمی می باشد که به اختصار توضیح میدهم :

- n_samples : تعداد نمونه هایی که باید تولید شوند
- n_features : تعداد ویژگی ها یا ویژگی های یک نمونه
- n_classes : تعداد کلاس های مسئله
- n_clusters_per_class : تعداد خوش های هر کلاس
- weights : وزن های هر کلاس
- class_sep : جدایی بین مراکز داده های مختلف
- and ect...

- ما با استفاده از ماژول inspect می باشیم امکاناتی برای انجام عملیات های با روش ها ، کلاس ها ، توابع و سایر اشیاء در زمان اجرا فراهم کند دستور from inspect import classify_class_attrs از کتبخانه دستور

classify_class_attrs را فراخوانی می کند و این دستور ویژگی های کلاس از نوع طبقه بندی شده را ساماندهی می کند.

- ما با استفاده از ماژول matplotlib که مخفف ان plt است ، استفاده می کنیم این کتابخانه به ما این امکان را می دهد تا تصاویر و نمودار تحت پایتون رسم کنیم.

ما داده های مصنوعی طبقه بندی ایجاد می کنیم که دارای 1000 نمونه است ، دارای 2 ویژگی می باشد برای مثال قد و وزن ، دما و رطوبت ، میزان مقدار مصرف سیگار و افزایش سرطان ریه ، در مرحله بعد میزان ویژگی های اضافی است که می تواند وابسته به ویژگی های اصلی باشد که در این مثال مقدار ان صفر است ، در این ارگومان برای تولید داده های تصادفی از عدد آنتالمپوس استفاده می شود و با تغییر مقدار این ارگومان داده های تصادفی تولید شده تغییر میکند ، در قسمت بعد تعداد خوشه های که برای هر کلاس ایجاد می شود مشخص می کند که در اینجا فقط یک خوشه داریم، در ارگومان بعدی جدایی بین مراکز داده های مختلف مقادیر کوچکتر یعنی که هر چه مقادیر کوچک تر نزدیک تر و هرچه مقادیر بزرگ تر دور تر می باشد در مرحله آخر هم تعداد کلاس ها تعیین شده است.

دستور scatter برای رسم نمودار پراکندگی استفاده می شود ، متغیر اول نشان دهنده این است که تمام مقادیر روی ردیف اول ماتریس و متغیر دوم تمام مقادیر روی ردیف دوم را نمایش بدهد ، متغیر آخر هم نمایش دهنده رنگ داده های ما می باشد.

از نظر چالشی بودن چالش خیلی زیاد چالش خاصی برای جدا کردن نداریم چون که ما میزان پراکندگی را 1.9 در مثال قبل قرار دادیم برای سخت تر کردن و چالش برانگیز کردن این مثال می توانیم این پراکندگی را کمتر کنیم که در مثال می بینیم.

توضیح مختصری دیگر: می توانید از کتابخانه Scikit-learn برای تولید دیتاست مورد نظر استفاده کنید. این کتابخانه دارای یک ماژول بنام datasets است که انواع مختلفی از دیتاست های آماده را ارائه می دهد. اما چون شما می خواهید یک دیتاست ساده با 1000 نمونه، 2 کلاس و 2 ویژگی ایجاد کنید، از تابع make_classification در این ماژول می توانید استفاده کنید.

در ادامه کدی آمده است که این دیتاست را ایجاد کرده و نمایش می دهد:

```
from sklearn.datasets import make_classification
```

```
import pandas as pd
```

تولید دیتاست با 1000 نمونه، 2 ویژگی و 2 کلاس

```
X, y = make_classification(n_samples=1000, n_features=2, n_classes=2,  
random_state=42)
```

#DataFrame تبدیل دیتاست به یک

```
df = pd.DataFrame(X, columns=['Feature 1', 'Feature 2'])
```

```
df['Class'] = y
```

نمایش اولیه از دیتاست

```
print("نمونه‌ها از دیتاست:")
```

```
print(df.head())
```

این کد یک دیتاست با 1000 نمونه، 2 ویژگی و 2 کلاس ایجاد می‌کند. سپس این دیتاست به صورت یک DataFrame پانداس نمایش داده می‌شود.

قسمت 3:

بررسی سرطان شایع پستان (breast_cancer) «علت انتخاب این است که اساتید محترم در ایم مورد مفصل توضیح داده اند»

load_breast_cancer این تابع از ماژول sklearn.datasets برای بارگیری مجموعه داده های سرطان

پستان استفاده می کند train_test_split این تابع از ماژول sklearn.model_selection برای تقسیم

داده‌ها به دو بخش آموزشی و آزمایشی استفاده می کند LogisticRegression و SGDClassifier: این دو

تابع از ماژول sklearn.linear_model برای ایجاد مدل های طبقه بندی لجستیک و SGD مورد استفاده

قرار می گیرند `accuracy_score`: این تابع از ماژول `sklearn.metrics` برای ارزیابی دقت مدل استفاده می کند.

این قسمت داده هایی که در دیتابیس سرطان پستان دارد را بارگزاری میکند و در دو متغیر قرار می دهد در این تابع یک آرگومان دارد که به معنای این است که داده ها را به صورت دو ارایه و جداگانه بر دو متغیر برگرداند در خط بعدی کد هم اندازه متغیر ها را نشان می دهیم.

این قسمت داده ها و متغیر ها را دو صورت تقسیم بندی می کنیم ، هم از متغیر ابتدایی دو نمونه آموزش و آزمایش استفاده می شود در تابع استفاده شده هم متغیر ابتدایی را ویژگی و متغیر دوم را بردار برچسب ها نام گذاری می شود و اندازه داده های تست را 20 درصد می گذارد بنابراین داده های آموزش باقی می ماند و در خط بعدی هم اندازه ان ها را نشان می دهد.

: Logistic Regression with prediction X_train

در خط اول در یک متغیر تابع رگرسیون لاجستیک را قرار می دهیم تا بتوانیم از ان استفاده کنیم این تابع یک نمونه جدید از مدل رگرسیون لجستیک را ایجاد می کند و در متغیر قرار می دهد و در خط دوم تابعی است که فرایند آموزش مدل را کنترل می کند و مدل های ما را آموزش میدهد و در خط سوم تابع پیش بینی کننده ای روی متغیر اولیه ما تعریف می شود که با استفاده داده ی آموزش داده شده ی ایکس پیش بینی را انجام می دهد و در خط چهارم هم با استفاده از تابع میزان دقت می توانیم دقت بین داده پیش بینی شده و داده آموزش داده شده بدست آوریم

: Logistic Regression with prediction X_test

در خط اول در یک متغیر تابع رگرسیون لاجستیک را قرار می دهیم تا بتوانیم از ان استفاده کنیم این تابع یک نمونه جدید از مدل رگرسیون لجستیک را ایجاد می کند و در متغیر قرار می دهد در خط دوم تابعی است که فرایند آموزش مدل را کنترل می کند و مدل های ما را آموزش میدهد در خط سوم تابع پیش بینی کننده ای روی متغیر تست و ازمون ما تعریف می شود و در

تابع یک مولفه خواهیم داشت در خط چهارم میزان دقت بین داده پیش بینی شده و تست ما صورت می پذیرد

: Stochastic gradient descent X_{train}

در خط اول ما در یک متغیر تابع گرادیان نزولی را قرار می دهیم که ویژگی هایی دارد در مرحله اول الفا که نرخ یادگیری می باشد و آن را کنترل می کند که مدل چقدر در هر بهینه سازی پارامتر ها به روز می شود و در اینجا خیلی کم می باشد در مولفه ی بعدی تکرار هایی که الگوریتم برای بهینه سازی انجام می دهد نشان می دهد و در اینجا حداکثر 1000 تکرار می باشد در خط دوم تابعی است که فرایند آموزش مدل را کنترل می کند و مدل های ما را آموزش میدهد در خط سوم تابع پیش بینی کننده ای روی متغیر آموزش ما تعریف می شود و در تابع یک مولفه خواهیم داشت در خط چهارم میزان دقت بین داده پیش بینی شده و آموزش ما صورت می پذیرد

: Stochastic gradient descent X_{test}

در خط اول ما در یک متغیر تابع گرادیان نزولی را قرار می دهیم که ویژگی هایی دارد در مرحله اول الفا که نرخ یادگیری می باشد و آن را کنترل می کند که مدل چقدر در هر بهینه سازی پارامتر ها به روز می شود و در اینجا خیلی کم می باشد در مولفه ی بعدی تکرار هایی که الگوریتم برای بهینه سازی انجام می دهد نشان می دهد و در اینجا حداکثر 1000 تکرار می باشد در خط دوم تابعی است که فرایند آموزش مدل را کنترل می کند و مدل های ما را آموزش میدهد در خط سوم تابع پیش بینی کننده ای روی متغیر تست ما تعریف می شود و در تابع یک مولفه خواهیم داشت در خط چهارم میزان دقت بین داده پیش بینی شده و تست ما صورت می پذیرد

: perceptron prediction with train data

در خط اول ابتدا تابع پرسپترون را در متغیری قرار می دهیم و دارای مولفه می باشد که الفا برای نرخ یادگیری و مولفه بعدی برای کنترل تکرار ها می باشد که در مثال قبل کامل توضیح دادم در خط دوم مدل یادگیری می باشد و آموزش مربوط به آن در خط سوم با استفاده از تابع پیش بینی

کننده با استفاده از داده های ایکس آموزش داده شده ی خودمون در خط چهارم با استفاده از تابعی میران دقت بین تابع پیش بینی شده و قسمت آموزش داده بررسی می کنیم

: perceptron prediction with test data

در خط اول ابتدا تابع پرسپترون را در متغییری قرار می دهیم و دارای مولفه می باشد که الفا برای نرخ یادگیری و مولفه بعدی برای کنترل تکرار ها می باشد که در مثال قبل کامل توضیح دادم در خط دوم مدل یادگیری می باشد و امزش مربوط به ان در خط سوم با استفاده از تابع پیش بینی کننده با استفاده از داده های ایکس تست خودمون در خط چهارم با استفاده از تابعی میران دقت بین تابع پیش بینی شده و قسمت تست بررسی می کنیم