

به نام خدا



درس: مقدمه‌ای بر یادگیری ماشین

استاد: دکتر صالح کلیبر

گزارش تمرین شماره ۱

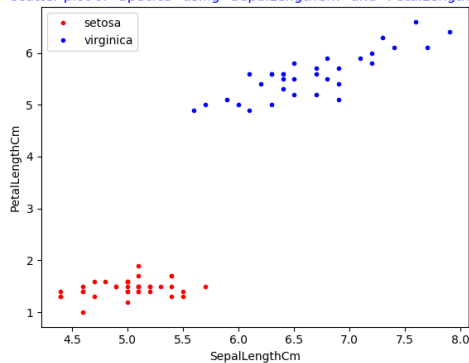
سید محمد امین منصوری طهرانی

۹۴۱۰۵۱۷۴

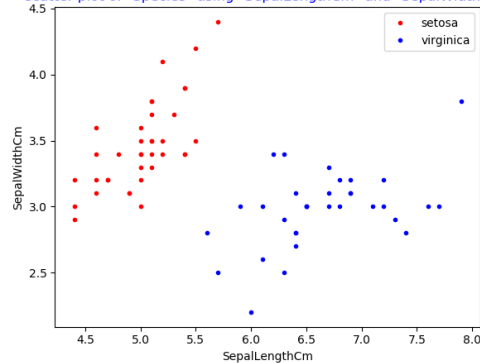
بخش اول – رسم نمودار

در این مسأله مطابق دستور کار، داده‌ها را به دو بخش `train` و `test` تقسیم می‌کنیم. با توجه این که از وجود ترتیب خاصی در داده‌ها اطلاع نداریم، انتخاب ۷۰ داده اول به عنوان `train` و ۳۰ داده آخر به عنوان `test` بدون اشکال به نظر می‌رسد و با سایر انتخاب‌ها از لحاظ تصادفی بودن تفاوتی ندارد (چون از ترتیب احتمالی در داده‌ها اطلاعی نداریم). برای بدست آوردن درک کلی از داده‌ها مطابق آن چه که خواسته شده نمودارهای زیر رسم می‌شوند. این تصاویر از اجرا کردن کد `jupyter notebook` بدست می‌آیند. هم‌چنین در کد `plotter.py` در فایل پیوست نیز آورده شده‌اند و با اجرای کد قابل مشاهده‌اند.

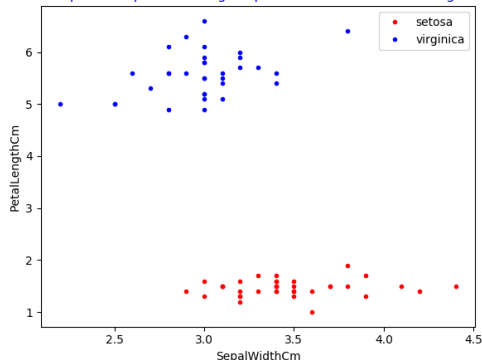
scatter plot of "Species" using "SepalLengthCm" and "PetalLengthCm"



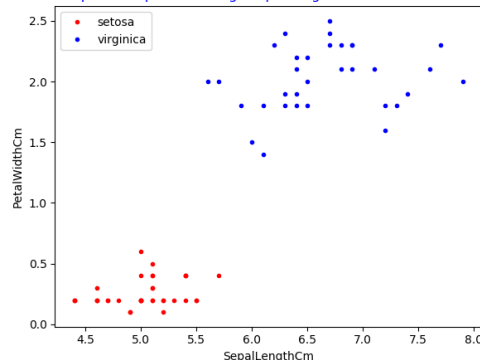
scatter plot of "Species" using "SepalLengthCm" and "SepalWidthCm"



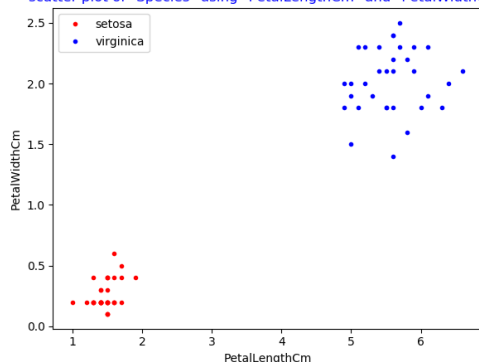
scatter plot of "Species" using "SepalWidthCm" and "PetalLengthCm"



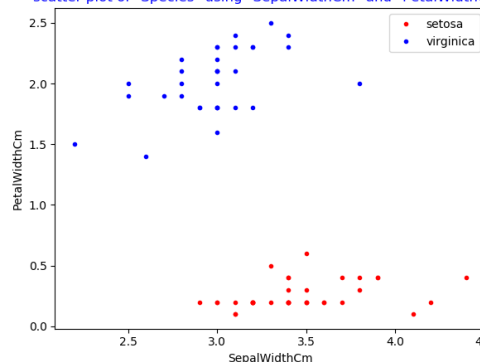
scatter plot of "Species" using "SepalLengthCm" and "PetalWidthCm"



scatter plot of "Species" using "PetalLengthCm" and "PetalWidthCm"



scatter plot of "Species" using "SepalWidthCm" and "PetalWidthCm"



در تصاویر بالا مشخص است همان‌طور که در توصیف داده‌ها گفته شده به شدت linear separable هستند و انتظار می‌رود عملیات classify کردن به سادگی صورت پذیرد. (با دقت کامل)

***توجه:** لطفاً کدهای py را در یک پوشه قرار دهید تا import ها دچار مشکل نشوند.

بخش دوم – یادگیری با مدل فاصله اقلیدسی

کد این بخش در قسمت اول فایل main.py که به پیوست تمرین قرار دارد نوشته شده‌است. برای این قسمت با توجه به این که مدل در دستور تمرین تعیین شده و کد نیز تماماً کامنت گذاری شده‌است، از آوردن توضیح اضافه در این جا خودداری می‌کنیم. لطفاً به کامنت‌ها مراجعه بفرمایید.

برای دقت مدل هم تعداد داده‌های تستی که نتیجه پیش‌بینی شده آن‌ها از مدل با برچسب آن‌ها یکی بوده را بر تعداد کل داده‌های تست تقسیم می‌کنیم.

نتیجه مطابق خواسته مسأله ۱۰۰ درصد درست پیش‌بینی شده‌است.

بخش سوم – یادگیری با مدل دلخواه

مدل انتخاب شده: Perceptron

از توضیح مدل پرسپترون با توجه به این که در کلاس درس تدریس شده‌است، صرف‌نظر می‌کنیم. توضیح کد:

برای استفاده از این روش یک class به نام Perceptron در فایل perceptron.py تعریف شده‌است که پس از پایان این بخش به توضیح آن می‌پردازیم.

پس از پایان کد قسمت قبل، دو آرایه label و training_inputs تعریف می‌شوند که برای آموزش طبقه‌بند استفاده می‌شوند. می‌خواهیم پرسپترون خطی استفاده کنیم که داده‌هایی با ۴ ویژگی را از هم جدا کند. (۴ ویژگی که در داده‌ها وجود دارد، طول و عرض گل‌برگ و کاس‌برگ) در حلقه اول برای داده‌های آموزش، مقادیر این ویژگی‌ها را با توجه به Perceptron ای که تعریف کرده‌ایم در متغیر training_inputs ذخیره می‌کنیم. سپس از آن جایی که پرسپترون ما به عنوان برچسب مقادیر ۰ و ۱ می‌پذیرد، گونه setosa را به عدد ۱ و گونه virginica را در همین داده‌های آموزش به عدد ۰ می‌نگاریم. سپس یک object از نوع Perceptron تعریف می‌کنیم با تعداد ویژگی ۴. (تنها آرگومانی که مشاهده می‌شود همین را نشان می‌دهد). سپس با صفتی که در کلاس برای آن تعریف کرده‌ایم (train) و ورودی داده‌های آموزش و برچسب صفر و یکی آن‌ها، آن را آموزش می‌دهیم. سپس داده‌های تست به فرمت قابل استفاده برای کلاس و صفت predict آن در می‌آیند و آرایه result در حلقه بعدی آن هر بار یک مقدار جدید به عنوان پیش‌بینی برچسب داده‌های تست به خود می‌پذیرد. (با صفت

perceptron.predict) در انتها نیز شمارش گری مانند کد قسمت قبل برای بررسی دقت مدل انتخاب می شود که اگر برجسب پیش بینی ۱ باشد و برجسب داده نیز setosa باشد یا اگر برجسب پیش بینی ۰ باشد و برجسب داده virginica باشد به معنی پیش بینی درست خواهد بود و این شمارش گر زیاد می شود. در غیر این صورت زیاد نمی شود. در نهایت نیز تعداد پیش بینی های درست به کل داده های تست تقسیم می شود و دقت چاپ می شود. نتیجه مطابق خواسته مسأله ۱۰۰ درصد درست پیش بینی شده است.

توضیح کد perceptron.py:

برای این کلاس چند متغیر ورودی در نظر گرفته ام. تعداد وزن هایی که لازم است یاد بگیریم، تعداد دفعات تکرار و آپدیت کردن وزن ها، نرخ یادگیری (پارامتر آپدیت کردن وزن ها). به این مقادیر اجازت دسترسی در زمان runtime نیز می دهیم تا در صورت لزوم بتوانیم آن ها را تغییر دهیم. هم چنین دقت می کنیم برای بدست آوردن شکل ماتریسی صحیح بدون حضور جمله بایاس، یک وزن دیگر هم اضافه شده و ورودی این بعد ۱ می باشد.

پس از آن متد پیش بینی را برای این کلاس تعریف می کنیم. مقدار متغیر summation همان طور که واضح است، حاصل ضرب داخلی ورودی در بردار وزن ها (بدون وزن مربوط به جمله بایاس) به اضافه جمله بایاس می باشد. با توجه به این که تابع فعالیت را پله در نظر می گیریم، اگر این مجموع (که همان پیش بینی با این وزن های فعلی می باشد) از صفر بیشتر بود ۱ و در غیر این صورت مقدار ۰ را باز می گردانیم.

پس از آن نیز متد یادگیری را برای این کلاس تعریف می کنم. این متد دو متغیر داده های آموزش و برجسب باینری آن ها را برای آموزش دیدن می خواهد. (که نحوه تبدیل برجسب ها به باینری که در بالا توضیح داده شد این جا به کار می آید). در ادامه این متد در صورتی که تعداد دفعات تکرار داده شده باشد از آن استفاده می کنیم و در غیر این صورت از مقدار default استفاده می شود. (آن را ۱۰۰ اختیار کردم). در ادامه یک شی که بتوان روی آن لوپ زد ایجاد می کنیم. این کار با زیپ کردن داده های آموزش و برجسب ها انجام می شود. در هر حلقه مقدار training_inputs در inputs و مقدار labels در label ذخیره می شود. سپس در ادامه همین حلقه، مقدار پیش بینی شده مدل با وزن های فعلی در prediction ذخیره می شود. پس از آن وزن ها همان طور که بلد هستیم تغییر می کنند و این حلقه به تعداد threshold تکرار می شود تا وزن ها بدست آیند. (البته وزن ها ممکن است قبل از رسیدن به تعداد تکرار تعیین شده نیز به مقدار صحیح رسیده باشند).

در واقع خطای فعلی در نرخ یادگیری و سپس در ورودی ها ضرب شده و با وزن های فعلی (بدون وزن ناشی از بایاس) جمع می شود. وزن بایاس هم که ناشی از جمع: خطا ضرب در نرخ یادگیری ضرب در ورودی ۱ است - با مقدار قبلی خود جمع می شود و آپدیت می شود و به این ترتیب همه وزن ها آپدیت می شوند و در تکرار بعدی همین حلقه دوباره پیش بینی و آپدیت داریم و این چرخه تا تعداد threshold تکرار می شود. (که امیدواریم تا این

تکرار وزن‌ها به مقدار صحیح رسیده باشند که در این مسأله به خاطر خوبی داده‌ها و تعداد نسبتاً کم آن‌ها این وضعیت مطلوب پیش می‌آید.)

در کد اصلی main با اجرا کردن کد، داده‌های تست تک تک توسط predict پیش‌بینی شده و دقت آن‌ها همان‌طور که گفته شد محاسبه و چاپ می‌شود.

خروجی کد پس از پایان اجرای کد main.py که خود از perceptron.py برای قسمت آخر استفاده می‌کند. دقت هر دو روش ۱۰۰ درصد است. با توجه به معروف بودن داده‌ها این انتظار می‌رفت. داده‌های این مشاهده به شدت linear separable هستند و این دقت را مدیون این ویژگی آن‌ها هستیم.

```
Python Console ×
...: score_method_2 = second_method_counter/len(test_data.Species)
...: print('Score using the second method (Linear Perceptron) is:')
...: print(score_method_2 * 100)
...:
Score using the first method (Euclidean Distance) is:
100.0
Score using the second method (Linear Perceptron) is:
100.0
```