

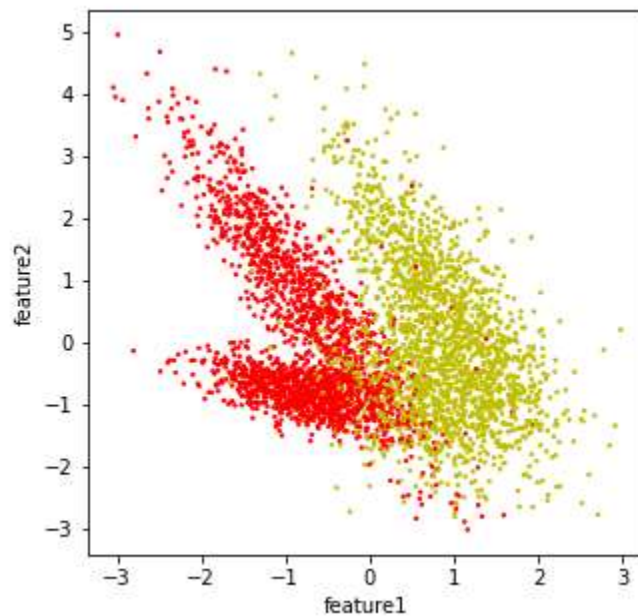
درس یادگیری ماشین

تکلیف برنامه نویسی: دسته بند پرسپترون

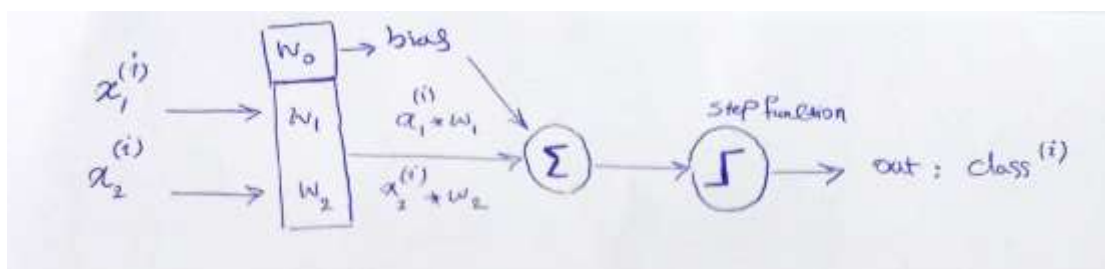
فایل ضمیمه : 02.ipynb

داده‌های آموزشی را از فایل `train.csv` به صورت یک دیتافریم کتابخانه `pandas` می‌خوانم. برای بخش کوچکی از محاسبه، کتابخانه `numpy` هم اضافه شد و از `matplotlib.pyplot` هم برای رسم نمودارها استفاده کردم. ماتریس داده‌های مساله شامل سه ستون است که دو ستون اول همان $x^{(i)}$ ها هستند و مقدار اسمی نظیر ستون سوم برچسب داده‌ها ست. ۴۰۰۰ داده‌ی آموزشی در اختیار هست. مساله دوکلاسه است و محتویات ستون `label` نظیر هر داده مقداری برابر ۰ یا ۱ دارد. برای راحت تر شدن روند کار این مقادیر را به ۱- و ۰- نظیر کردم.

نمودار نقطه‌ای داده‌های آموزشی را رسم می‌کنم. گروهی که برچسب ۱- دارند نقاط قرمز و نقاط زرد نظیر داده‌های با برچسب ۰- است. محور عمودی نظیر مقدار ستون دوم داده‌ها یا `feature2` یا همان $x_2^{(i)}$ روی محاسبات است و محور افقی هم نظیر مقدار ستون دوم داده‌ها یا `feature1` یا همان $x_1^{(i)}$ در جریان محاسبات.



مدل خواسته شده یک پرسپترون به صورت زیر است : (تابع پله‌ی استفاده شده همان تابع علامت است)

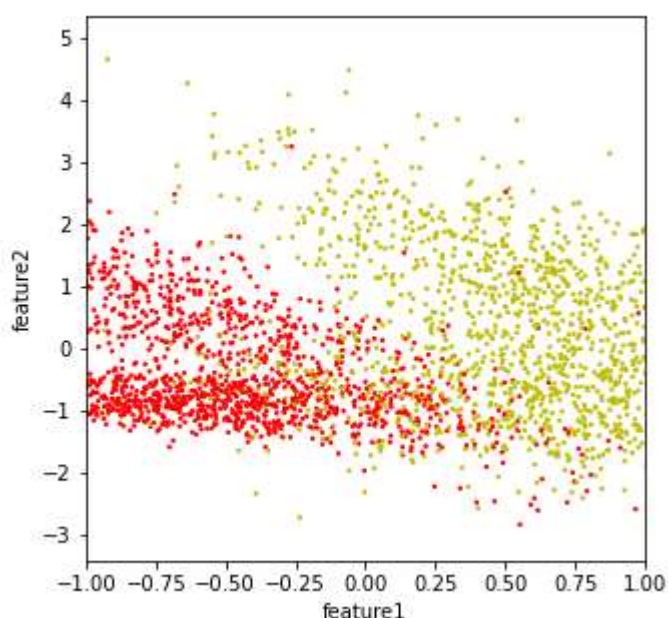


یعنی یک معادله ی خط به صورت زیر قرار است پیدا کنیم که اگر روی این نمودار رسم شود، دو کلاس داده ها را به صورت زیر از هم کاملاً تفکیک کند.

$$\omega_0 + \omega_1 x_1 + \omega_2 x_2 = 0$$

$$\forall x = [x_1, x_2] \quad \begin{cases} \text{if } \omega_0 + \omega_1 x_1 + \omega_2 x_2 \geq 0 & \text{class yellow : 1} \\ \text{else } \omega_0 + \omega_1 x_1 + \omega_2 x_2 < 0 & \text{class red : -1} \end{cases}$$

اگر به نمودار دوباره نگاه کنیم و مخصوصاً با دقت در برش زیر از دامنه ی feature1 به سادگی معلوم می‌شود که چنین خطی وجود ندارد. یعنی این داده ها به طور خطی تفکیک پذیر نیستند.



پس مدل خطی هرگز روی داده های آموزشی به دقت ۱۰۰ نمی‌رسد (همه را درست کلاس بندی کند) پس اگر الگوریتم تیوری پرسپترون را استفاده کنیم که شرط توقفش تا زمانی است که هیچ داده ای در دسته ی اشتباهی نباشد، الگوریتم هرگز متوقف نمی‌شود. برای اینکه جواب نزدیکی برای مساله با مدل خطی پیدا کنیم، قراردادی الگوریتم آپدیت وزن ها را تا مقدار مشخصی تکرار اجرا می‌کنیم. اینجا من ۱۰۰ انتخاب کردم.

برای راحتی بیشتر محاسبات بایاس مدل، یک ستون ۱ به ماتریس داده ها اضافه کردم. که روی محاسبه مقدار $x_0^{(i)}$ را می‌گیرد.

حالا با استفاده از داده هایی که آماده شد، الگوریتم یادگیری پرسپترون را اجرا می‌کنیم. حاصل کار، یک بردار سه مولفه ای w است. این همان خطی است که انتظار داشتیم داده ها را در دو طرف خود تفکیک کند.

یادآوری اینکه زمانی که خواستیم روی داده های تست عملکرد مدل را بسنجیم همین تغییرات باید روی ان ماتریس داده ها هم اعمال بشود.

حلقه ی اصلی آموزش در تابع `train_perceptron` پیاده شده که ماتریسی که آماده شد را می‌گیرد و تا زمانی که تمام داده ها درست طبقه بندی شوند یا تا حداکثر ۱۰۰ بار بردار وزن های پرسپترون را از رابطه ی زیر آپدیت می‌کند. که M مجموعه ی داده‌های اشتباه کلاس

بندی شده است و y هم همان مقدار label هست.

$$W_{new} = W_{old} + \eta \sum_{n \in M} x^{(n)} y^{(n)}$$

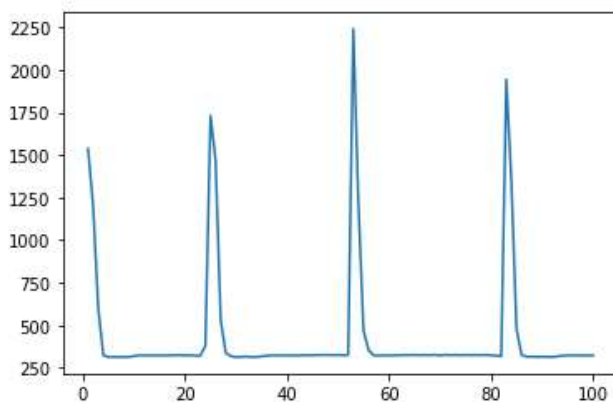
$$\begin{bmatrix} w0_{new} \\ w1_{new} \\ w2_{new} \end{bmatrix} = \begin{bmatrix} w0_{old} \\ w1_{old} \\ w2_{old} \end{bmatrix} + \eta \begin{bmatrix} \sum_{n \in M} x_0^{(n)} y^{(n)} \\ \sum_{n \in M} x_1^{(n)} y^{(n)} \\ \sum_{n \in M} x_2^{(n)} y^{(n)} \end{bmatrix}$$

مقدار اولیه وزن ها را تصادفی انتخاب کردم. هر بار با بردار وزن به دست آمده از مرحله ی قبل، با تابع **classify** روی کل داده های آموزشی کلاس بندی انجام می دهد و نتیجه ی کلاس به دست آمده از رابطه ی زیر، روی یک ستون به نام **assigned_label** از دیتافریم آموزشی ذخیره می شود.

در قدم بعدی لازم است این ستون را با ستون **label** مقایسه کرده و نقیض نتیجه را روی ستون دیگری به نام **missclass** قرار دهد. این کار در تابع **find_missed** انجام می شود. کل داده هایی که روی این ستون مقدار صحیح دارند، برای آپدیت وزن در تکرار فعلی حلقه استفاده می شوند. در واقع داریم از الگوریتم دسته ای (**batch mode**) برای الگوریتم آموزش پرسپترون استفاده می کنیم. برای ساختن مقداری که باید به مولفه های **w** اضافه و کم شود از یک دیتافریم موقتی استفاده می کنم و در نهایت با نرخ یادگیری که برابر **0.001** تعیین شده **w** را آپدیت می کنم. برای کنترل روند آموزشی تعداد داده هایی که اشتباه کلاس بندی شده اند در هر تکرار نگهداری می شود. وزن هایی که به دست آمدند :

W0	W1	W2
0.063877	0.0834247	0.226945

روی ۱۰۰ بار روند تکرار الگوریتم تعداد داده هایی که هربار اشتباه کلاس بندی شدند را رسم می کنم.



نمودار پرش ها و فرود های ناگهانی دارد. در تکرار های خاصی اگر آموزش متوقف شود خطای داده آموزش بسیار بالاست. این وضعیت به دلیل همان خطی تفکیک پذیر نبودن داده ها رخ می دهد. و در سایر نقاط نمودار با شیب خیلی کمی حول عرض کمتر از ۵۰۰ تا داده ی خطا می ماند. مدل در بهترین حالت روی داده های آموزشی این تعداد خطا دارد. این خطا ناشی از پیچیدگی کم مدل است و با این مدل قابل رفع نیست.

داده های تست را از روی فایل **test.csv** به صورت یک دیتا فریم خواندم و تغییرات لازم را انجام دادم.

برای ارزیابی مدل روی ماتریس داده های آموزش و تست، یک باز هر دو ماتریس را با مدل به دست آمده طبقه بندی می کنیم. بعد از ساخته شدن ستون های assigned_label و missclass نتیجه ی طبقه بندی را با دو تابع measurments و plot_result ارزیابی می کنیم. تابع measurments اطلاعاتی از عملکرد مدل به دست می دهد. این تابع تعداد کل داده هایی که اشتباه دسته بندی شده اند و نسبت آن ها به کل داده ها را و میزان دقت مدل روی آن داده های خواسته شده را گزارش می کند. اینجا دقت را برابر نسبت داد هایی که درست کلاس بندی شدند به کل داده ها گرفتیم (همان correct_ratio)

تابع plot_result برای نمایش داده های دو کلاس بر اساس رنگ کلاششان روی ماتریس داده ی ورودی تعریف شده. داده هایی که اشتباه کلاس بندی شده اند را به رنگ سیاه نمایش داده. روی همین نمودار یک خط سبز رنگ هم رسم می شود که نشان دهنده ی مقدار W نهایی به دست آمده ست. (این خط را از اعمال وزن های W روی دو نقطه ی ابتدا و انتهای محدوده ی محور عمودی رسم کردیم)

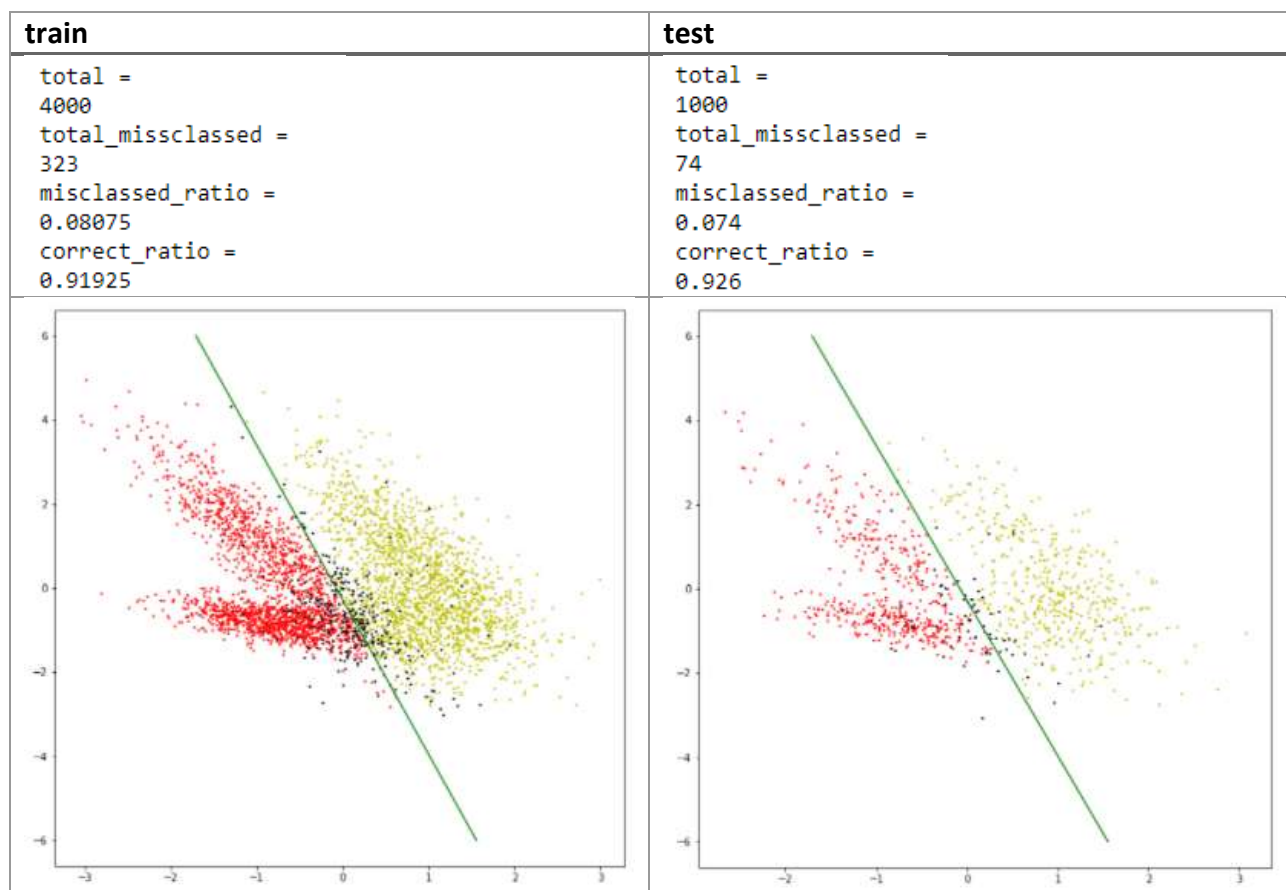
$$x_2 = \text{Feature 2} \in [-6, 6]$$

$$w_0 + w_1 x_1 + w_2 x_2 = f(x)$$

$$x_2 = -6 \rightarrow f(x) = ?$$

$$x_2 = +6 \rightarrow f(x) = ?$$

نتایج در ادامه آمده.



همانطور که از شکل نمودار اول گزارش انتظار می‌رفت بیشتر داده‌هایی که اشتباه کلاس بندی شده اند متعلق به همان بخش در هم رفتگی دو کلاس هستند که از رویش نتیجه گرفتیم مدل خطی نمی‌تواند جدا کننده‌ی کامل داده‌های آموزشی باشد. اگر به خاطر درهم رفتگی ذاتی این دو کلاس روی صفحه نبود، مدل می‌توانست طبقه بندی بهتری انجام دهد...

کل داده‌های تست و آموزش به ترتیب ۱۰۰۰ و ۴۰۰۰ تا هستند که روی هر کدام باز هم به ترتیب ۷ و ۸ درصد خطا داریم. که برای اسن داده‌ها عدد ثابت قبولی هم هست. همانطور که از نمودار آبی رنگ تعداد کلاس بندی‌های اشتباه در روند آموزش انتظار می‌رفت تعداد کل خطاهای آموزش عددی زیر ۵۰۰ برابر ۳۲۳ به دست آمده. اگر درصد کل‌های درست طبقه بندی شده روی تست را به عنوان دشت نهایی در نظر بگیریم، این مدل ۹۲ درصد داده‌های این مساله را می‌تواند درست طبقه بندی کند که عدد قابل قبولی است. اما اگر توزیع داده‌های تست طوری بود که بیشتر شامل نقاط درهم رفته‌ی دو کلاس میشد، عملکرد افت می‌کرد اما ما داده‌های تست و آموزش را همیشه از یک توزیع یکسان انتخاب می‌کنیم و علت نزدیک بودن نتیجه‌ی تست و آموزش هم در واقع همین است. اینجا مدل نهایی به لحاظ عددی هرچند عملکرد خوبی داشته ولی اگر به نمودارها نگاه کنیم واضح است که از مدل درستی استفاده نکردیم و من نتیجه‌ی عددی نسبتاً مطلوب را به این نسبت می‌دهم که هرچند دو کلاس درهم رفتگی دارند اما این درهم رفتگی روی توزیع داده‌ها عدد کوچکی است (نمودارهای اول گزارش) که ممکن است بتوان آن را حتی به زیاد بودن نویز داده‌ها نسبت داد اما چون نویز را یک فاکتور رندوم می‌گیریم، با حرف‌هایی که درباره تشابه توزیع (صرفاً به لحاظ ظاهری و نه بررسی آماری) این داده‌های درهم (سیاه رنگ روی نمودار) روی دو دسته‌ی داده‌ها زدم، نظر من این است که مساله‌ی اصلی کم بودن پیچیدگی مدل پرسپترون خطی برای طبقه بندی و یادگیری این داده‌هاست. با ساختن مدل‌های درجه بالاتر می‌توان این فرض را بررسی کرد.