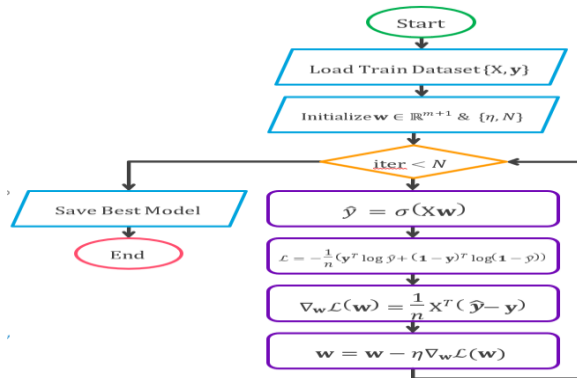


## پاسخ سوالات سری اول درس یادگیری ماشین

محمدرضا حاجی نیا

### سوال اول



فلوچارت کلی آموزش یک مدل طبقه بند خطی را در بالا مشاهده میکنیم.

مراحل فلوچارت به ترتیب شامل:

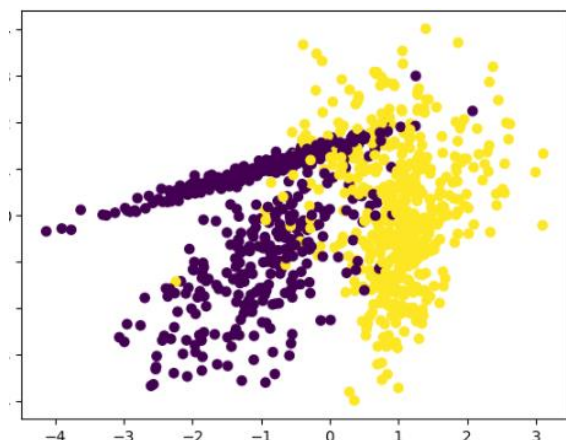
۱. بارگذاری داده‌های آموزش : در ابتدا، داده‌های آموزش (ورودی‌ها و برچسب‌ها) بارگذاری می‌شوند.
۲. مقداردهی اولیه وزن‌ها و بایاس : وزن‌ها و بایاس به صورت تصادفی مقداردهی اولیه می‌شوند.
۳. پیش‌بینی خروجی: در هر تکرار، خروجی پیش‌بینی شده با استفاده از تابع فعال‌ساز سیگموید محاسبه می‌شود.
۴. محاسبه خطا : خطای بین خروجی پیش‌بینی شده و برچسب واقعی با استفاده از تابع هزینه لگاریتمی محاسبه می‌شود.
۵. به‌روزرسانی وزن‌ها: گرادینان هزینه نسبت به وزن‌ها محاسبه می‌شود و سپس با استفاده از نرخ یادگیری، وزن‌ها به‌روزرسانی می‌شوند.
۶. تکرار فرآیند: باتوجه به تعداد دوره‌های آموزشی این فرایند تکرار می‌شود. و اگر به پایان برسد مدل ذخیره می‌شود.

تغییر حالت از دو کلاس به چند کلاس:

- ✓ تغییر تابع فعال‌ساز. تابع سیگموید دو حالت ۰ یا ۱ رو داره که باینری است میتون از توابع دیگه برای این موضوع استفاده کرد.
- ✓ همچنین باتغییر کلاس قطعا در مقداردهی وزن‌ها نیز تغییراتی ایجاد می‌شود.

## سوال دوم

یک دیتاست با ویژگی ها و کلاس ها خواسته شده ایجاد شد. در ابتدا این دیتاست چالش برانگیز و سخت بود زیرا داده ها بسیار به یکدیگر نزدیک کلاس بندی ها دارای دو خوشه بودند.



از بعضی پارامتر ها استفاده شد تا نوع کلاس بندی بهتر مشخص بشود. توضیح برخی از پارامتر های مهم را در این قسمت داریم.

پارامتر `n_clusters_per_class` برای تعیین تعداد خوشه در هر کلاس است. در اینجا مقدار ۱ است پس ، هر کلاس فقط یک خوشه دارد، یعنی داده های هر کلاس به یکدیگر نزدیک تر هستند.

پارامتر `class_sep` میزان تفکیک بین داده های دو کلاس را تعیین می کند. یعنی دوری و نزدیکی داده ها از هم.

پارامتر `n_redundant` تعداد ویژگی های اضافی و تکراری است که به داده ها اضافه می شود. با تنظیم این پارامتر به صفر، ویژگی های تکراری اضافه نمی شود.

برای چالش برانگیز تر کردن دیتاست میتوانیم مرز های تصمیم گیری را نزدیک تنظیم کنیم. یا خوشه بندی داده ها را بیشتر کنیم. مثلاً میتوانیم `class_sep` را کاهش بدهیم. یا `n_clusters_per_class` افزایش بدهیم. همچنین می توانیم بین داده ها نویز ایجاد کنیم که بسیار تاثیر گذار است.

## سوال سوم

برای ایجاد یک مدل طبقه بندی خطی از دو مدل

### LogisticRegression , SGDClassifier

استفاده شده است. داده ها به دو قسمت آموزش و تست تقسیم بندی می شوند. در ابتدا پارامتر مهمی که باید لحاظ بشه پارامتر stratify است. این پارامتر مشخص می کند که از هر کلاس به اندازه تقریباً برابر در مجموعه داده آموزش و تست قرار گیرد.

از مدل رگرسیون لجستیک LogisticRegression برای آموزش و ارزیابی عملکرد مدل استفاده می شود. سپس از داده های آموزشی برای آموزش مدل سپس ارزیابی انجام می شود. از ماتریس درهم ریختگی (confusion matrix) برای ارزیابی دقیق تر عملکرد مدل استفاده می شود.

بهینه ساز این مدل sag است یک روش تقریبی برای بهینه سازی وزن ها. که دو روش تصادفی و تکراری را انجام می دهد.

Max\_iter تعداد دوره های آموزشی می باشد.

random\_state برای یکسان بودن یا متفاوت بودن نتایج کاربرد بسیار زیادی دارد. (تکرار پذیری نتایج)

برای بهبود در این مدل تعداد دوره های آموزشی را تا حد مشخصی می شود بالا برد.

و مهمترین پارامتر solver است که باید با روش های مختلف و با توجه به داده ها انتخاب بشود.

مدل SGDClassifier از روش گرادیان کاهشی تصادفی برای آموزش استفاده می کند. پارامترهای مختلف مدل را استفاده کرده ایم و در این قسمت توضیح می دهیم.

تابع هزینه hinge: تابع هزینه که برای اندازه گیری خطا استفاده می شود. این پارامتر باید تست شود زیر توابع دیگری وجود دارند در این مدل.

max\_iter تعداد دوره های آموزشی

learning\_rate یا نرخ یادگیری یک پارامتر مهم است که مشخص می کند که میزان تغییرات وزن ها در هر مرحله از آموزش چقدر باشد. از این پارامتر جهت تنظیم سرعت یادگیری و همگرایی زودتر به مقدار بهینه استفاده می شود روش های مختلفی برای تعیین کردن دارد.

از پارامتر  $\alpha$  استفاده می کنیم تا بتوانیم بیش برآزش را کنترل کنیم. این پارامتر بازه تغییرات وزن را کنترل میکند تا مقادیر وزن در طول آموزش زیاد نشوند. افزایش یا کاهش این مقدار تاثیر زیادی در آموزش دارد. مقدار زیاد گاه باعث می شود مدل به سختی با داده ها هماهنگ شود و وزن های بسیار کوچکی داشته باشد. مقدار کم باعث می شود مدل به راحتی با داده ها هماهنگ بشه و وزن های بزرگتری بگیره. پس باید با امتحان کردن این مقدار بدست بیاد.

### سوال پنج

در این تمرین از کتابخانه `deawdata` استفاده می کنیم. ابتدا سه کلاس  $a, b, c$  ایجاد میکنیم با توجه به شکل مدنظر.

و ستون رنگ را حذف میکنیم. از انجایی که مقادیر  $x, y$  زیاد می باشند از نرمال سازی استفاده میکنیم. همچنین کلاس ها را به صورت عددی تغییر می دهیم یعنی  $a=0, b=1, c=2$ . سپس داده ها به مجموعه آموزشی و تست تقسیم می شوند.

از مدل های قبلی استفاده می شود تا ارزیابی بشود مدل. همچنین از ماتریس درهم ریختگی استفاده می شود تا تعداد داده های درست و غلط در هر کلاس مشخص بشود. در آخر هم مرز تصمیم گیری بین سه کلاس مشخص می شود.