

1. Explain what logistic regression is.

logistic regression یکی از مهم ترین روش های classification در یادگیری ماشین می باشد. برخلاف linear regression که عددی می تواند مقادیری پیوسته داشته باشد. در logistic regression عددی به یکی از چند مقدار محدود (class/category) تقسیم خواهد داشت. به عنوان مثال در یک دسته Binaray Classification ما دو مقدار 1 و 0 (True/false, Positive/negative) را خواهیم داشت. Sigmoid function / logistic function بهر صورت زیر تقریب می شود؟

$$f_{\vec{w},b}(\vec{x}) = g(z) = \frac{1}{1 + e^{-z}} \quad ; \quad z = \vec{w} \cdot \vec{x} + b$$

عددی تابع نرین عمده عددی بین 0 و 1 خواهد بود. در واقع این عدد احتمال تکراری داده درودی (مجموعه ویژگی های آن یعنی بردار \vec{x}) در کلاس 1 را نشان می دهد. برای این منظور threshold ای تعیین می شود. (فرضاً 0.5) که اگر احتمال مورد نظر از آن threshold بیشتر شد، داده درودی به کلاس 1 / Positive / True تعلق دارد.

$$f_{\vec{w},b}(\vec{x}) = P(y=1 | \vec{x}; \vec{w},b)$$

به عنوان مثال در سؤال مطرح شده که شابل تعدادی نقاط بارنگ های متفاوت است، ما به دنبال معنی حسیم را با کمترین مقدار خطا در بین نقاط بارنگ مختلف، امتحان کنیم. در حالت اول سؤال ما به دنبال یک خط به فرم $w_1x + b$ بودیم که چون یک خط نمی توانست این وجه تمایز دایره های با درجه بالاتر از 1 را نشان دهد، Accuracy پائینی بدست می آوردیم. در logistic regression، بهر از بدست آوردن معادله z و تکرار دادن آن در معادله sigmoid هر چه نقاط از معادله z دورتر باشد، احتمال بالاتری برای تعلق به کلاس وجود دارد. رابطه خط هم به صورت دیگری باشد؟

$$J(\vec{w},b) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{\vec{w},b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}^{(i)})) \right]$$

2. write its formula for this Problem --- 3.

$$Z = w_1x_1 + w_2x_2 + b = 0 \xrightarrow[\text{boundary}]{\text{decision}} x_2 = -1/1 x_1 + 2.0$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

2. Explain the following metrics:

$$\text{accuracy} = \frac{TP + TN}{P + N} = \frac{\text{No. correct Predictions}}{\text{No. All Predictions}}$$

از میان کل پیش بینی ها چه تعداد درست پیش بینی شده اند. یعنی مجموع تعداد داده هایی که به درستی به کلاس 0 و 1 تقسیم داده شده اند تقسیم بر کل داده هایی که بر آن ها Label تعلق گرفته است.

$$\text{Precision} = \frac{TP}{TP + FP}$$

از میان کل پیش بینی های Positive، یعنی تمامی داده هایی که به کلاس مثبت تقسیم داده شده اند، چه تعداد به درستی پیش بینی شده اند.

$$\text{Recall} = \frac{TP}{TP + FN}$$

از میان کل داده هایی که واقعاً Positive Label داشته، چه تعدادی را مدل ما پیدا کرده است.

specificity: True Negative Rate

$$\frac{TN}{TN + FP}$$

چنان Recall اما Label مثبت، Negative است.

در واقع از بین کل داده های Negative، چه تعداد به درستی پیش بینی شده اند.

sensitivity: True Positive Rate = Recall

* مقدار Metric حاضر قابل Colab *

3. Problems that Recall is more important than Precision?

طبیاً در مسائلی که پیدا کردن مقدار Positive برای ما بسیار حیاتی است و هزینه FN بسیار بیشتر از FP است. مثال: تشخیص این که یک تومور سرطانی است.

4. Problems that Precision is more important than Recall?

در مسائلی که FP بسیار پرهزینه تر از FN می باشد.

مثال: تشخیص اسپم.

Problem 2)

1. write the loss function for this formula.

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (\vec{w}_{w,b}(\vec{x}^{(i)}) - y^{(i)})^2$$

This is cost function, loss is simply this formula but for one sample, so we don't have $\frac{1}{2m}$ & Σ .

for this problem

we have $J(w_2, w_1) = \frac{1}{2m} \sum_{i=1}^m (w_2 x_2^{(i)} + w_1 x_1^{(i)} - y^{(i)})^2$

2. Drive the derivation of the loss function.

Gradient descent

repeat {

$$\frac{dJ}{dw_1} = \frac{1}{m} \sum_{i=1}^m (w_2 x_2^{(i)} + w_1 x_1^{(i)} - y^{(i)}) x_1^{(i)}$$

$$\frac{dJ}{dw_2} = \frac{1}{m} \sum_{i=1}^m (w_2 x_2^{(i)} + w_1 x_1^{(i)} - y^{(i)}) x_2^{(i)}$$

}

3. Explain the impact of input normalization ---

دستی درونی (input) نرمال سازی شده باشد: مثلاً feature اول در بازه $[0, 1]$ ، feature دوم در بازه $[0, 500]$ باشد، در این صورت روش gradient descent در مسیر رسیدن به بیشینه نوسان (oscillate) می کند. در نتیجه، به تدریج Converge می کند. علت این موضوع این است که تغییرات در استای یک feature بسیار بزرگتر از feature دیگر خواهد بود. همان مشق بر اساس یک feature

Problem 3)

Step 1.

1. Explain why does over fitting happen?

دلی از دلایلی که می تواند منجر به OverFitting شود؛ (۱) کم بودن تعداد Sample های Dataset.

(۲) بالا بودن تعداد feature ها.

(۳) بالا بودن پیچیدگی مدل و وابستگی بالا.

در این مثال به علت درجه بالا بودن، پیچیدگی مدل بالا بوده و مدل دارای وابستگی بالا است. به این معنی که با کوچکترین تغییر در یک داده مدل ما به شکل قابل توجهی تغییر می کند. همچنین تعداد داده های بسیار کم است. در نهایت مدل می تواند generalize شود و بسیار وابسته به داده های train می شود.

2. Discuss about the training error and generalization error of this model.

training error همان error روی training Dataset ما می باشد. از آنجایی که تعداد Data های ما ۲۱ بوده و مدل ما از مرتبه ۲۰ است. نکته داده های معنی افتاده و علاقه training error ما چیزی حدود ۰.۱ است. اما generalization error میزان دقت مدل روی داده های از قبل دیده شده را بررسی می کند. این مقدار می تواند با کوچک کردن overfitting افزایش می یابد. طبیعتاً overfitting باعث می شود تا بیشترین مدل که هوش پیدا کند و به اصطلاح مدل generalize نباشد. بنابراین generalization error افزایش می یابد.

Answer these questions:

1. همانطور که مشاهده می شود. با استفاده از regularization تا میانی خوبی از overfitting جلوگیری کردیم و به نسبت generalize بودن مدل پیش رفتیم. در واقع regularization با دارد کردن پارامترهای مدل به cost function سعی در کاهش وزن هایی که باعث می شود تا مدل overfit نشود. همانطور که مشاهده می شود L1 regularization بهترین خروجی را داشته و از همه بیشتر به تابع اصلی x^3 نزدیک است.

$$J(\vec{w}, b) = \frac{1}{2m} \left[\sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \rightarrow \text{Ridge}$$

$$\left[+ \frac{\lambda}{2m} \sum_{j=1}^n w_j \right] \rightarrow \text{Lasso}$$

2.

- Lasso \leftarrow برای Feature selection استفاده شد، ضریب برفی از Feature ها را صفر می کند.
- Ridge \leftarrow سعی کند تا ضرایب Feature ها را به شکل یک اندازه گیری که حس دهد.
- در واقع تفاوت Ridge و Lasso در Penalized term آن ها می باشد.

3. معادلش را گفتند L_1 (Lasso) در این مثال بهتر عمل کرده است. زیرا Lasso برفی از Feature ها را صفر کرده که باعث می شود تا با تعداد کمی Feature، generalization، بهتری داشته باشیم.

4. در histogram ها مشخص می شود که در L_2 توزیع یک اندازه گیری نسبت به histogram بدون Regularization داریم. در صورتی که در L_1 باعث Feature selection می شود برفی Feature ها صفر شده، مقدار آن ها نیز در محدوده صفر است که باعث می شود مدل از مرتبه کم انرژیک تر شود.

5. بی شک L_1 ، در واقع این موضوع مهم ترین برفی L_1 regularization است. چون باعث ایجاد Sparsity در ضرایب می شود و Feature های کم اهمیت تر و noise ها را حذف می کند. در نتیجه در برابر Outlier ها قدرتمندتر است.