

Assignment 3: Example Dependent cost sensitive logistic regression

SAI HARSHA KOTTAPALLI - CS17BTECH11036

Tanmay R - CS17BTECH11042

OBJECTIVE

The objective for this assignment is to identify fraudulent tax payers from the given dataset containing various details regarding their transactions.

Generally, most models only try to minimize the so-called zero-one loss where all misclassifications have equal cost and the correct classifications having cost zero.

But in the real world, this tends to not be the case, the cost of misclassifications might significantly vary based on the application and its instance. Particularly in our case, depending on the type of transactions done by a fraudulent person classifying him into genuine has different costs. For example, the cost incurred on classifying a person who has embezzled a few lakhs as genuine is not the same as someone who has embezzled a few hundreds.

Logistic Regression estimates the posterior probability of the positive class, as the logistic sigmoid of a linear function of the feature vector.

For this assignment, we will use the logistic regression model and a cost sensitive loss function to predict if a person is fraudulent or genuine.

Here, the output is calculated as:

$$\hat{p}_i = P(y = 1|x_i) = h_{\theta}(x_i) = g\left(\sum_{j=1}^k \theta^j x_i^j\right),$$

where $g(x)$ is the sigmoid function.

For the purpose of cost sensitive linear regression, we use the following cost function:

$$J^c(\theta) = \frac{1}{N} \sum_{i=1}^N \left(y_i(h_{\theta}(x_i)C_{TP_i} + (1 - h_{\theta}(x_i))C_{FN_i}) \right. \\ \left. + (1 - y_i)(h_{\theta}(x_i)C_{FP_i} + (1 - h_{\theta}(x_i))C_{TN_i}) \right).$$

DATA SET SUMMARY

The dataset consists of 147636 people's transaction history comprising 11 features(independent variables) namely, NotCount, YesCount, Average tax per month(ATPM), purchase from fraudulent(PFD), purchase from genuine(PFG), Sales to fraudulent(SFD), sales to genuine(SFG), purse value in way bill(WP), sales value in way bills(WS), AH and AN.

It also provides us the label(Status) whether the person is fraudulent or genuine so that we can train the model accordingly and also the FNC associated with each person's classification, i.e., False Negative cost for misclassifying fraudulent person as genuine. Some of the cells are zero under FNC, this might be due to the fact that a person is associated with a very small business and there's not much impact in the real world by misclassifying him.

In the dataset, while training I tried to remove the data points which have extremely high FNC and the ones which have FNC as 0, but i didn't find any significant difference from taking random points , so i took a random split of 70% as training dataset and 30% for test dataset. In fact the accuracy dropped by a small amount, hence i decided to not remove the outliers.

Also since the data was already normalized but not exactly the method in which it was one, I decided to leave YesCount and NoCount as it is, so that the logistic regression can learn weights accordingly given that all of the features are independent.

EXPERIMENTATION RESULTS

In our case since it takes manpower to call the person identified as fraudulent to verify, we associate the cost of true positive and false positive as 150 each. There is no cost incurred for identifying genuine persons as genuine.

As for the false negative cost for misclassifying fraudulent people as genuine, as mentioned above is already provided in the dataset.

We use a seed at the beginning so that the results are reproducible from the code.

We then did do a grid search over learning rate from 1e-4 to 1e-8, with *Lecun initialization* for weights(the default in pytorch) and try out various other combinations for optimizations(Adam, SGD, etc) to find the optimal parameters and corresponding result:

Epochs	TP _(given)	TN _(given)	FP _(given)	FN _(given)	LR	Optimizer	Accuracy	F1 Score
2000	150	0	150	in dataset	1e-06	SGD	86.19%	76.92%

Where, TP = true positive cost, TN = true negative cot, FP = false positive cost, FN = false negative cost, LR = learning rate. SGD = Stochastic gradient descent.

Here, F1 score is calculated as :

$$F_1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

We also notice that as we increase epochs(on overfitting), the accuracy starts decreasing at some point after stabilizing.