

Example Dependent cost-sensitive logistic regression

Team: Puneet Mangla (CS17BTECH11029)
Yash Khasbarge (CS17BTECH11044)
Rushikesh Tammewar (CS17BTECH11041)

Objective: Implement example dependent cost-sensitive logistic regression to classify fraudulent and non-fraudulent taxpayers.

Given a feature vector \mathbf{x} for each taxpayer and a label y which is 1 iff taxpayer is fraudulent else 0, train a Logistic regression classifier $\mathbf{h}_\theta(\mathbf{x})$ (outputs $P(y=1|\mathbf{x})$) by optimizing the following objective

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

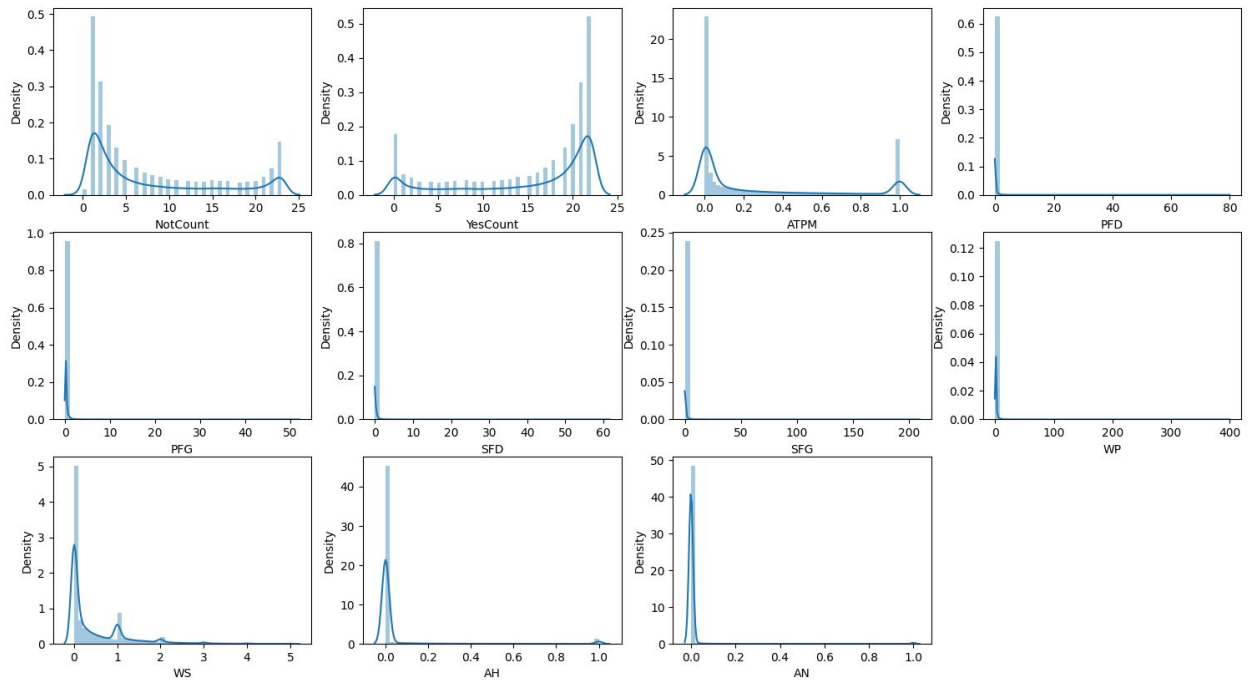
Here CTP, CFN, CFP, and CTN are costs for true positive, false negative, false positive, and true negative respectively.

At inference, we make predictions by obtaining probability, $P(y=1|\mathbf{x}) = \mathbf{h}_\theta(\mathbf{x})$. If $P(y=1|\mathbf{x}) > 0.5$ we classify it as positive else negative.

Dataset Summary:

- 147636 samples
- "Status" is the dependent variable. Class "0" has 103554 samples and Class "1" has 44082 samples.
- "FNC" gives us false negative cost for each example.
- Rest 11 fields are independent features used to predict.
- Description of each field:
 - Status (dependent), NotCount (not filed in time), YesCount (filed in time), ATPM (Average tax per month), PFD (purchase from fraudulent), PFG (purchase from genuine), SFD (Sales to fraudulent), SFG (sales to genuine), WP (purse value in waybill), WS (sales value in waybills), FNC (false negative cost), AH (One more independent variable), AN (One more independent variable)

- Distribution plots of 11 features



Attribute	Mean	Variance
NotCount	57.369688	7.7221
YesCount	57.713207	15.221
ATPM	0.134161	0.2532
PFD	0.112566	0.027062
PFG	0.139405	0.050789
SFD	0.088974	0.021648
SFG	1.289939	0.071709
WP	5.916689	0.271927
WS	0.582762	0.493587
AH	0.033224	0.045045
AN	0.010196	0.013474

Experimentation Results:

- Implementation Details:
 - FPC, TPC, and TNC = 150.0, 150.0, 0.0
 - Number of epochs: 20
 - Train: 100000 samples, Validation: 15000 samples, and Test: 32636 samples
 - Adam Optimizer use for optimizing
 - Checkpoint with best validation accuracy is used to evaluate on test data
- Evaluation Metrics:

- Accuracy: $100 * (\text{\#correct predictions}) / (\text{total predictions})$
- Precision: $(\text{\#True positives}) / (\text{\#True positives} + \text{\#False positives})$
- Recall : $(\text{\#True positives}) / (\text{\#True positives} + \text{\#False negatives})$

- Results:

- Test accuracy: 86.506 %
- Test precision: 0.779, Test Recall: 0.757, Test F1 Score: 0.768
- Train vs Val Accuracy/Loss Curves below

