

- Data loading and preprocessing
  - o I load the data, handle missing values and split the data as features (x) and targets (y) with pandas
  - o The data includes categorical and numerical data so i used OneHotEncoder for categorical data and StandardScaler for numerical data to encode them. This way i can be sure that all the features are in comparable state.
- Dataset / Dataloader
  - o Since you asked me to write the code in OOP standards i choose to create a class for dataset(CustomDataset). This way the PyTorch can use the processed data as in being able to use DataLoader for batch processing.
- Since we have more than one target class i choose to create a multi-task NN
  - o This way i get a separate output heads for each task or target class(seller\_number, customer\_number, main\_account)
- Training and evaluation
  - o Training loop works in a way that it can handle all tasks simultaneously. And i used accuracy for each target class to evaluate the model performance
- I split the data as in Test and Training with train\_test\_split
  - o The size that you want me to use is 0.25 so i defined the test\_size=0.25 to ensure that 25% of the data(50 of 200) is held for testing.
  - o I used LabelEncoder to include all the classes from dataset
  - o I used normalization for "amount" class in a way that the mean of it is to be 0 and the deviation of 1. With this i aim to achieve to stabilize the convergence of the network while training.
- So i used tuples for data processing but i had some issues with it:
 

In CustomDataset class \_\_getitem\_\_ returns a tuple (x[idx], y[idx]). This is the PyTorch standard format. Where X[idx] is the feature vector for the one data point and Y[idx] contains the targets where defines the model to be a multi-task NN
- Model architecture
  - o The main part of the model consist of 2 shared layers but there is task specific output layers for each target class: output\_seller, output\_customer, output\_main\_account. With task specific models, i achieved to make the classifications on one model and not 3 separate models. The two layers are seemd fine but i might end up increasing it to 3

or 4 because i kina seem to approve the layers's output size to decrease to 1 but its not necessary or might be worse. And since i only have 200 documents to use both training and testing it seemd alright for to choose a less deeper model and just increase the epoch number and batch size or the learning rate.

- I had some input size issues while designning the model layers: The input sizes were not the same with the data size the layers accept. I fixed it with the line `input_size = x_training_final.shape[1]`. With this line the model dynamically adjust to the correct number of features.
- I used `acc_score` to quantify the model performance
- I aslo had some problems with as in how should i encode testing and training data, do i choose to encode them all together before splitting them was a hard choice.

here are some preformance values of the model:

Seller Accuracy: 60.00% Customer Accuracy: 75.00% Main Account Accuracy: 85.00%

when 500 epochs for training. Other epoch numbers was not very effective or different than this output so i chose not to include them. I also reviewed the code to chatgpt. You can read it from my [github](#). I used google colab. I could explain more of the problems that i encountered especially about encoding, but to be hones i am a bit tired. So pls excuse my laziness.