

Artificial Intelligence II

Nefeli Tavoulari

Fall Semester 2021

1

The Cross Entropy loss function for one instance is defined as:

$$\mathcal{J}(y, \hat{y}) = - \sum_{k=1}^K (y_k \log(\hat{y}_k))$$

where K is the number of classes and \hat{y} is the softmax function.

Calculating the derivative of Cross Entropy with respect to the logit j:

$$\frac{\partial \mathcal{J}}{\partial \theta_j} = - \frac{\partial}{\partial \theta_j} \sum_{k=1}^K (y_k \log(\hat{y}_k))$$

$$\frac{\partial \mathcal{J}}{\partial \theta_j} = - \sum_{k=1}^K y_k \frac{\partial}{\partial \theta_j} \log \hat{y}_k, \text{ where:}$$

$$\log \hat{y}_k = \log\left(\frac{e^{\theta_i}}{\sum_{k=1}^K e^{\theta_k}}\right) = \theta_i - \log(\sum_{k=1}^K e^{\theta_k})$$

Now:

$$\frac{\partial(\theta_i - \log(\sum_{k=1}^K e^{\theta_k}))}{\partial \theta_j} = \frac{\partial \theta_i}{\partial \theta_j} - \frac{1}{\sum_{k=1}^K e^{\theta_k}} * \frac{\partial \sum_{k=1}^K e^{\theta_k}}{\partial \theta_j} =$$

$$- \frac{e_i^\theta}{\sum_{k=1}^K e_k^\theta} + \frac{\partial \theta_i}{\partial \theta_j} = -\hat{y} + \frac{\partial \theta_i}{\partial \theta_j}, \text{ where,}$$

if $j=i$ then $\frac{\partial \theta_i}{\partial \theta_j} = 1$, otherwise $\frac{\partial \theta_i}{\partial \theta_j} = 0$

$$\text{So, } \frac{\partial \mathcal{J}}{\partial \theta_j} = - \sum_{k=1}^K y_k (\{i == k\} - \hat{y})$$

$$= - \sum_{k=1}^K y_k \{i == k\} + \sum_{k=1}^K y_k \hat{y}$$

$$= -y_j + \sum_{k=1}^K y_k \hat{y}$$

$$= -y_j + \hat{y} \sum_{k=1}^K y_k$$

$$= -y_j + \hat{y} \sum_{k=1}^K y_k$$

$$= -y_j + \hat{y}$$

since the one-hot encoded vector y has a sum equal to 1.

So, $\frac{\partial \mathcal{J}}{\partial \theta} = \hat{y} - y$ and the same applies if we have m instances. In that case $\hat{y} - y$ is a vector.

2

Following the stages of the computational graph we get:

$$a = x * m + b$$

$$y = \text{RELU}(a) = \max(0, a) = \max(0, x * m + b)$$

$$f = \text{LOSS} = (y - y^*)^2 = (\max(0, a) - y^*)^2 = (\max(0, x * m + b) - y^*)^2$$

So, now I have to compute the gradients with respect to x, m, b, y^* :

$$\frac{\partial x}{\partial x} = \frac{\partial m}{\partial m} = \frac{\partial b}{\partial b} = \frac{\partial y^*}{\partial y^*} = 1$$

$$\frac{\partial a}{\partial x} = m \frac{\partial a}{\partial m} = x \frac{\partial a}{\partial b} = 1 \frac{\partial a}{\partial y^*} = 0$$

Using the chain rule:

$$\frac{\partial y}{\partial x} = \frac{\partial \max(0, a)}{\partial x} = 0 \text{ or } m \left(= \frac{\partial a}{\partial x} \right)$$

$$\frac{\partial y}{\partial m} = 0 \text{ or } x \left(= \frac{\partial a}{\partial m} \right)$$

$$\frac{\partial y}{\partial b} = 0 \text{ or } 1 \left(= \frac{\partial a}{\partial b} \right)$$

$$\frac{\partial y}{\partial y^*} = 0$$

$$\frac{\partial f}{\partial x} = \frac{\partial (\max(0, a) - y^*)^2}{\partial x} = \frac{2 * \partial (\max(0, a) - y^*)}{\partial x} * \frac{\partial y}{\partial x} =$$

$$0 \text{ or } 2 * (x * m + b - y^*) * \frac{\partial a}{\partial x} = 2 * (x * m + b - y^*) * m$$

$$\frac{\partial f}{\partial m} = 0 \text{ or } 2 * (x * m + b - y^*) * \frac{\partial a}{\partial m} = 2 * (x * m + b - y^*) * x$$

$$\frac{\partial f}{\partial b} = 0 \text{ or } 2 * (x * m + b - y^*) * \frac{\partial a}{\partial b} = 2 * (x * m + b - y^*)$$

$$\frac{\partial f}{\partial y^*} = 0 \text{ or } 2 * (x * m + b - y^*) * \frac{-\partial y^*}{\partial y^*} = -2(x * m + b - y^*)$$

(0 if $y = 0$)

3 Vaccine Sentiment Classification

In this notebook I trained a neural network, which classifies tweets as Neutral, Pro-vax or Anti-vax.

3.1 Preprocessing

First of all, I performed some preprocessing and feature extraction on the data, similar to the one I had done in the previous assignment. I removed all empty or duplicate tweets, as they offer no useful information to the model, I lowercased tweets, performed lemmatization, removed stopwords, special characters, emojis, words written in a non-latin alphabet and urls, since I noticed that the results are better this way. Then, I used the GloVe pretrained Word Embeddings, to take advantage of the relationships between the words. More precisely, I used the ones that are extracted from tweets, since we are dealing with tweets too and I experimented with all the dimensions. I started by finding every word of each tweet in the embeddings dictionary and then I find the average for every tweet (I sum each tweet's words' embeddings and I divide the sum by the dimension of the embeddings), in order to end up with a matrix of *numberOfInstances * dimension*. Then, I split the dataset into batches.

3.2 Models

The best model I created consists of 3 linear layers, the input size is equal to the dimension of the embeddings and the output size is equal to the number of classes (3). After experimenting with many different combinations of activation functions I found the best results using ReLU and Sigmoid in-between the hidden layers. I did not use an activation function in the output, since I use Cross Entropy as loss function, which combines `nn.LogSoftmax()` and `nn.NLLLoss()`. Moreover, I used Batch Normalization, which seemed to ameliorate my results. As far as the optimizer is concerned I used Adam, since it gave much better results than SGD or Adagrad.

In the following tables I will present some comparisons I made in order to find my most efficient model. The accuracy/precision/recall/F1 results I got were the following, with learning rate=1e-4 and SGD with momentum=0.2 and weight decay=0.01 with 3 layers(186,64 units):

	<i>SGD</i>	<i>Adagrad</i>	<i>Adam</i>
<i>Train</i>	0.5014	0.5529	0.6682
<i>Validation</i>	0.5157	0.5543	0.6498

Here we observe that there is also something wrong with the SGD and Adagrad results, since the accuracy on the validation set is higher than that of the train set.

Some other results I had to compare, concerning the number of units, hidden layers and activation functions:

	<i>4 – layers – no – activation</i>	<i>4 – layers – activation</i>
<i>Train</i>	0.6076	0.6952
<i>Validation</i>	0.6143	0.6498
	<i>3 – layers – no – activation</i>	<i>3 – layers – activation</i>
<i>Train</i>	0.60581	0.6682
<i>Validation</i>	0.61042	0.6498
	<i>2 – layers – no – activation</i>	<i>2 – layers – activation</i>
<i>Train</i>	0.6064	0.6563
<i>Validation</i>	0.6152	0.6428

It is obvious that the activation functions (ReLU, Sigmoid) and the Batch Normalization help the model become more expressive, in comparison with using just linear layers. Also, we notice that the model with the 3 layers has the best results for the validation set, while the one with the 4 layers seems to be overfitting. The number of units of the hidden layers of the 4-layer model is 256, 128, 64, of the 3-layer model is 128, 64 and of the 2-layer one is 64.

Using different number of units (in hidden layers) for the 3-layer model:

	[128, 64]	[256, 128]	[64, 32]
<i>Train</i>	0.6682	0.6815	0.6643
<i>Validation</i>	0.6498	0.6419	0.6520

So, the one one with the least units is the best one. Also, I experimented with some more activation functions, eg

softmax, logsoftmax, and higher or lower learning rate but still the best results were the ones I mentioned before.

The results with or without Dropout:

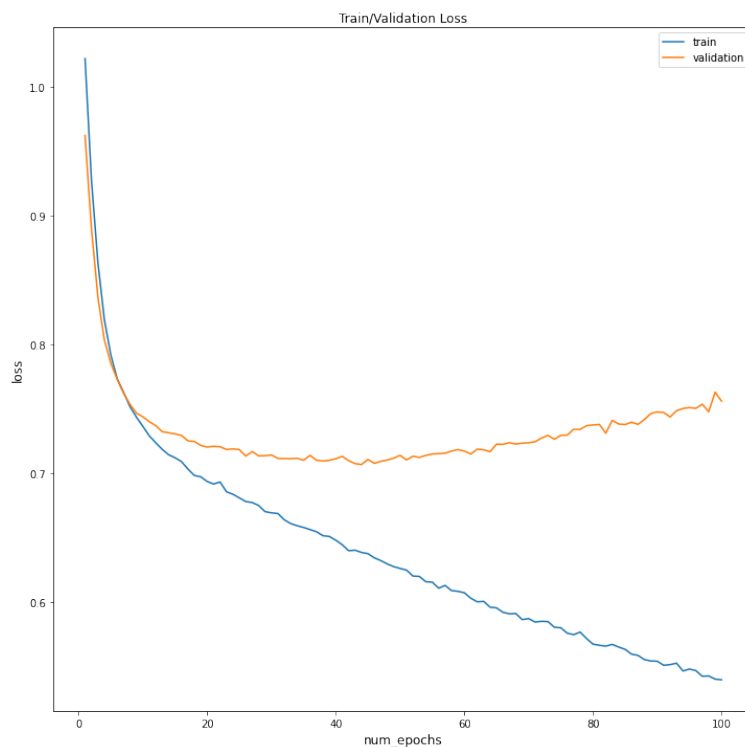
	<i>without</i>	<i>after – relu</i>	<i>after – sigmoid</i>	<i>both</i>
<i>Train</i>	0.6643	0.6436	0.6504	0.6415
<i>Validation</i>	0.6520	0.6538	0.6411	0.6419

Here, we notice that even though the model that uses Dropout after ReLU has the best score, it is overfitting, since the validation score is higher than the training score, which might not be necessarily a negative thing, but the first one is almost equally good, so it is preferred.

Testing the best model with different dimensions of word embeddings:

	<i>25d</i>	<i>50d</i>	<i>100d</i>	<i>200d</i>
<i>Train</i>	0.6643	0.7109	0.7818	0.8681
<i>Validation</i>	0.6520	0.6660	0.6822	0.6770

Even though, the model that uses 100d embeddings seems to be the best one, the epochs vs loss and epochs vs accuracy curves indicate that it is overfitting.



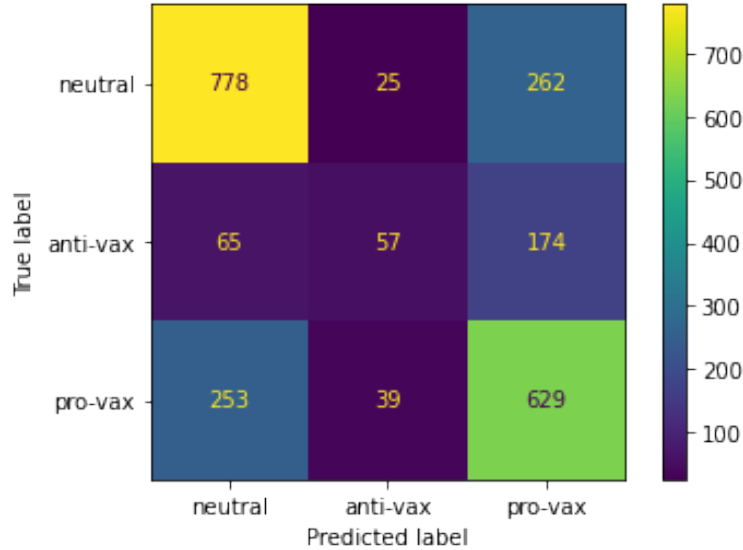
So, my best model uses 25d word embeddings and does not overfit.

In the notebook I have also included another model. The main difference this model has from the first one is that it contains one more layer, it uses Dropout and Log-Softmax on the output, because the used loss function is NLLLoss. The results are similar to the ones of my best model:

	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1</i>
<i>Train</i>	0.6644	0.6644	0.6644	0.6644
<i>Validation</i>	0.6560	0.6560	0.6560	0.6560

3.3 Evaluation - Comparison

Using different metrics, precision, recall, f1 score, accuracy, confusion matrix, roc auc, I was able to determine the quality of my model. From the confusion matrix, I understand that my model makes good enough predictions, since the diagonal has quite bigger numbers. I observe also that it does not predict that well the anti-vaxers, but that is expected, because the anti-vaxers in the data are very few, in comparison to the other classes. The confusion matrix ,using the aforementioned model :



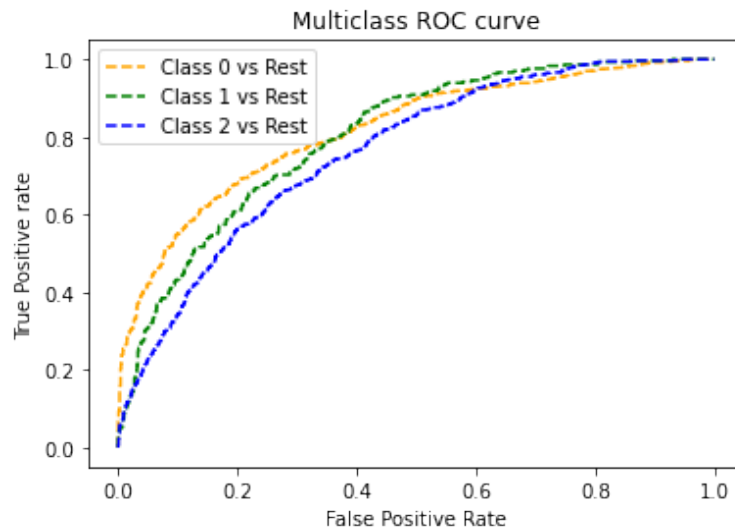
The results as I mentioned them before:

	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1</i>
<i>Train</i>	0.6643	0.6643	0.6643	0.6643
<i>Validation</i>	0.6520	0.6520	0.6520	0.6520

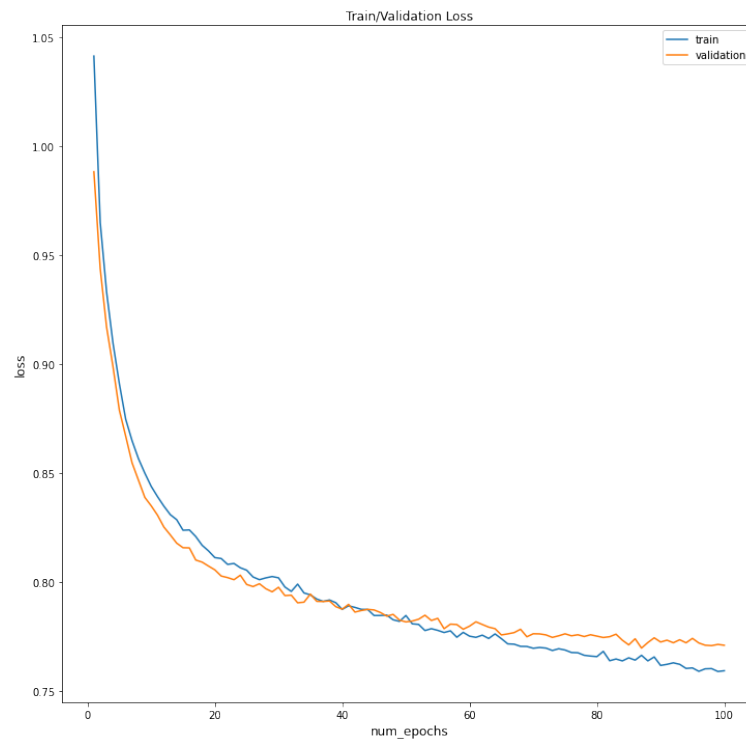
One-vs-One ROC AUC scores:
0.778982 (macro),
0.785546 (weighted by prevalence)

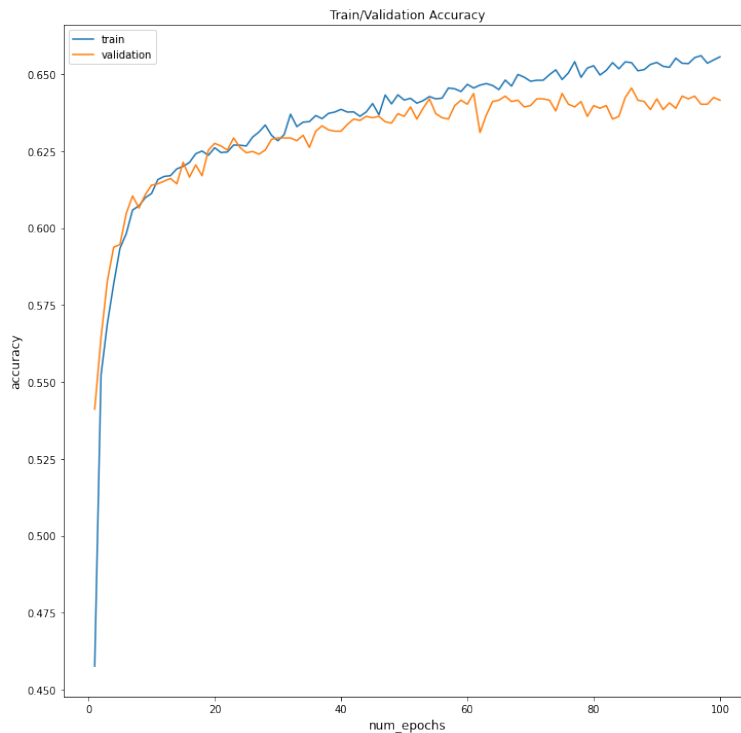
One-vs-Rest ROC AUC scores:
0.793916 (macro),
0.793266 (weighted by prevalence)

And some plots:



We observe that the anti-vax class is not that well predicted as the other two classes. However, the curves in general show that the model distinguishes the classes, since they lean towards the top left corner (TPR).





We observe that after approximately the 40th epoch, the model does not learn that much, the loss increases and the accuracy decreases, on the validation set. In general, however, the behavior of the curves is the expected one. That means that the accuracy increases and the loss decreases as the model learns.

Finally, I want to note that the results are slightly worse than the ones obtained by the Logistic Regression model of the previous assignment. One reason behind this could be the fact that some words related to Covid-19 are probably not present in the Embeddings. Another reason could be that a more clever neural network is needed (RNN, LSTM, GRU, etc).

References

- [1] [PyTorch Documentation](#)
- [2] [Word Embeddings](#)
- [3] [Optimizers](#)
- [4] [Loss Functions](#)