

BS6207 Assignment 4

SONG YIFAN

Table 1 result:

(Steps=2000)

Autoencoder vs. without Autoencoder

UCC and UCC^{2+} models had autoencoder branch and they were optimized jointly over both autoencoder loss and UCC loss. $UCC_{\alpha=1}$ and $UCC_{\alpha=1}^{2+}$ don't have autoencoder branch and optimized only by UCC loss. So, I just change the weight of the loss of autoencoder as 0 in the model to make sure that don't include the autoencoder loss.

Different uccs

UCC and $UCC_{\alpha=1}$ models trained on bags with labels of ucc1 to ucc4. (ucc_start=1, ucc_end=4) UCC^{2+} and $UCC_{\alpha=1}^{2+}$ models trained on bags with labels of ucc2 and ucc4 (ucc_start=2, ucc_end=4).

	ucc acc.	clustering acc.
UCC	0.960	0.951
UCC^{2+}	0.954	0.949
$UCC_{\alpha=1}$	0.970	0.950
$UCC_{\alpha=1}^{2+}$	0.947	0.904

From the result we can see that the ucc accuracy of UCC and UCC^{2+} seems higher than $UCC_{\alpha=1}$ and $UCC_{\alpha=1}^{2+}$, which means autoencoder helps to generate more accurate result for same steps. The score of ucc2 to ucc4 is less than ucc1 to ucc4. That because of the ucc prediction task becomes easier at the absence of pure sets and models.

Try different feature extractors

Different instances:

	Instance=16	Instance =32	Instance =64
Ucc acc	0.950	0.96	1.00

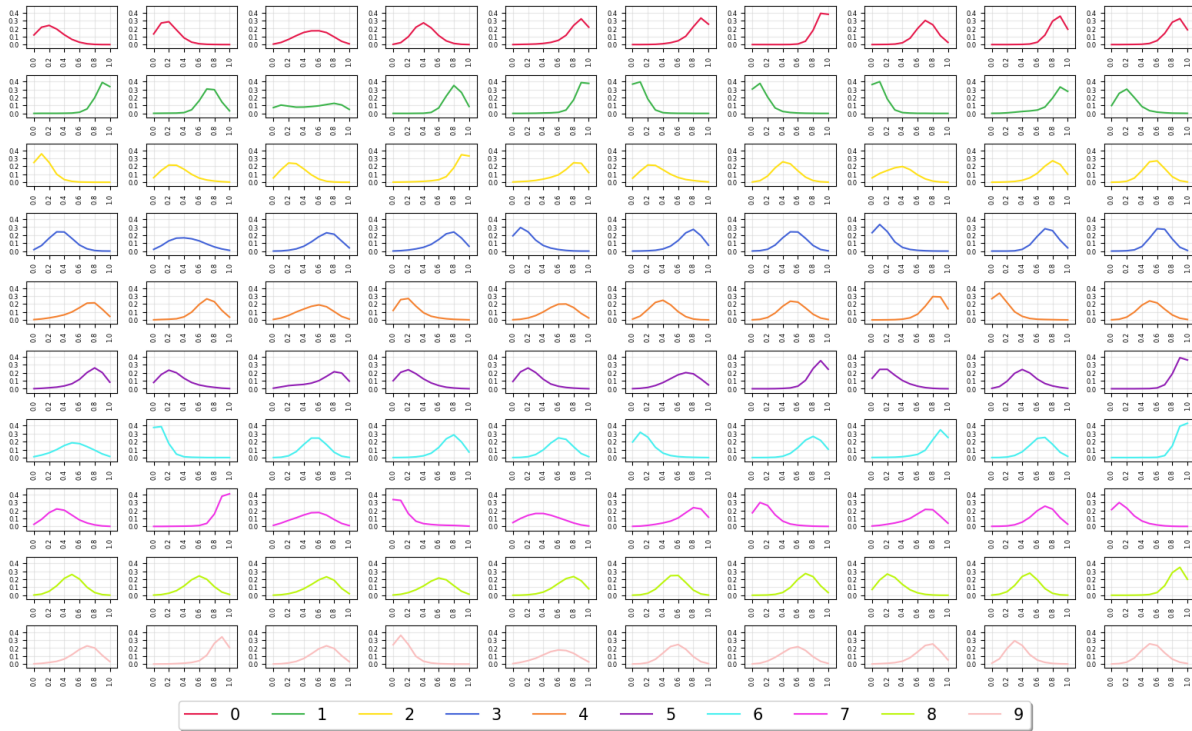
More instance means more complex the model is, so as expected, more instances, more running time, and more accurate. Instance means features in each bag, so more features will produce more accurate result.

Try different activation function (Relu to Tanh)

	ReLu	Tanh
Ucc_acc	0.960	0.610

ReLu is better than tanh. The biggest advantage of ReLu is the non-saturation of its gradient, which greatly speeds up the convergence of stochastic gradient descent compared to the tanh function.

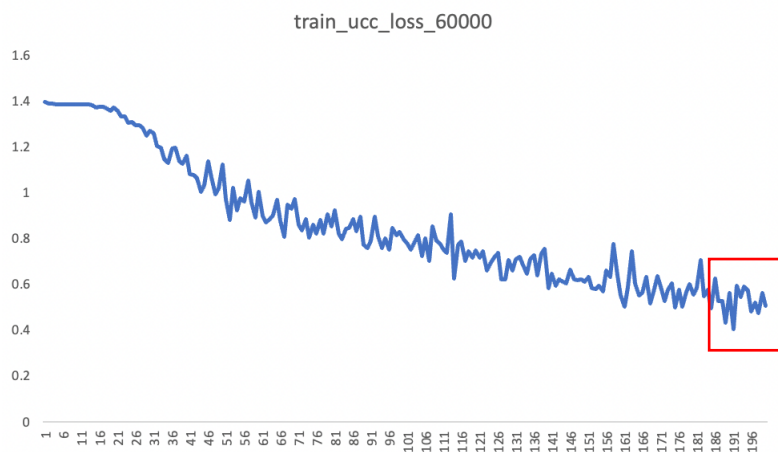
Distribution regression:



The kernel density estimation layer is used to construct the distribution of feature extracted by autoencoder. The figure shows above is the result of distribution regression. As we can see that the distribution of features in each class are not same. It means that KDE is a very efficient way to separate the distributions into different class.

Try different training number.

1. 60000 for training (default) ucc_auc=1.00



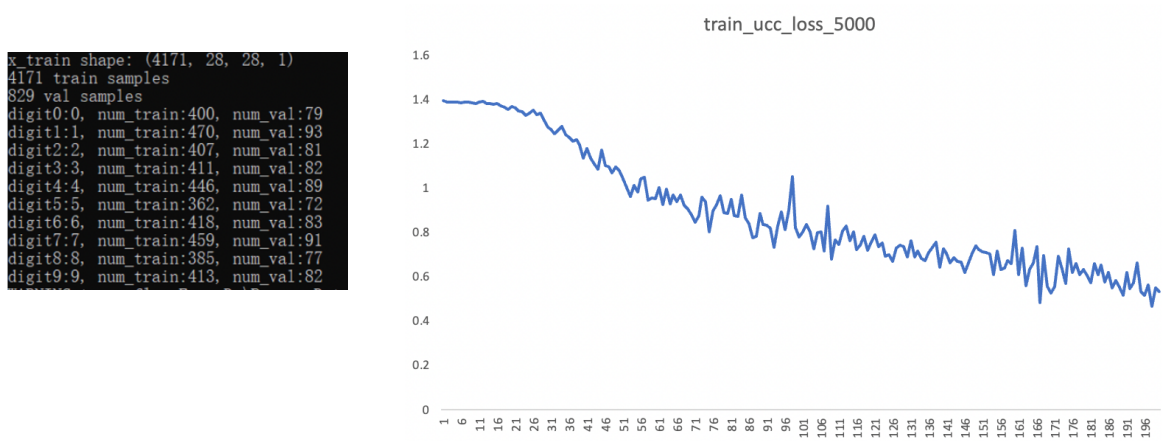
Seems like over fitting.

2. 20000 for training ucc_auc=0.950



Seems like over fitting.

3. 5000 for training ucc_acc=0.900



No overfitting until now, not sure for more steps.

4. 500 for training

Unable to get the result since small number of training data.

Hypothetically, if all hyperparameters were to be held constant, then more data means more steps along the gradient at the same learning rate, which could indeed overfit more easily. And observed the figure above, it's true that more training data means easier to overfit, but using more data makes the result more accurate.

In the proposed paper (Oner, Lee, & Sung, 2020), they implemented the model that can predict the number of unique classes in a package is learning patterns in the data that can be used to cluster the data into classes. I think the original model is accurate enough, however, it is time consuming because they default trained a million steps. To shorten the training time, I think some image augmentation can be applied. Using this method can both increase the accuracy and decrease the training time.

Reference

Oner, M. U., Lee, H., & Sung, K. W. (2020). WEAKLY SUPERVISED CLUSTERING BY EXPLOITING UNIQUE CLASS COUNT. *ICLR*.