

ABSTRACT

In this project, I built model to predict dropout in Massive Open Online Course(MOOC) platform, which is the topic in KDD cup 2015.

With my feature engineering result on this complicated three-dimensional dataset, I first explored different models and optimized parameters selection to reach best performance. With few models of good prediction on validation data, I ensemble them together using XGBoost Classifier to do a second level training prediction. In addition, I implemented Long Short Term Memory(LSTM) Recurrent Neural Networks with Keras on 30 days univariate Time series data and reached 0.857.

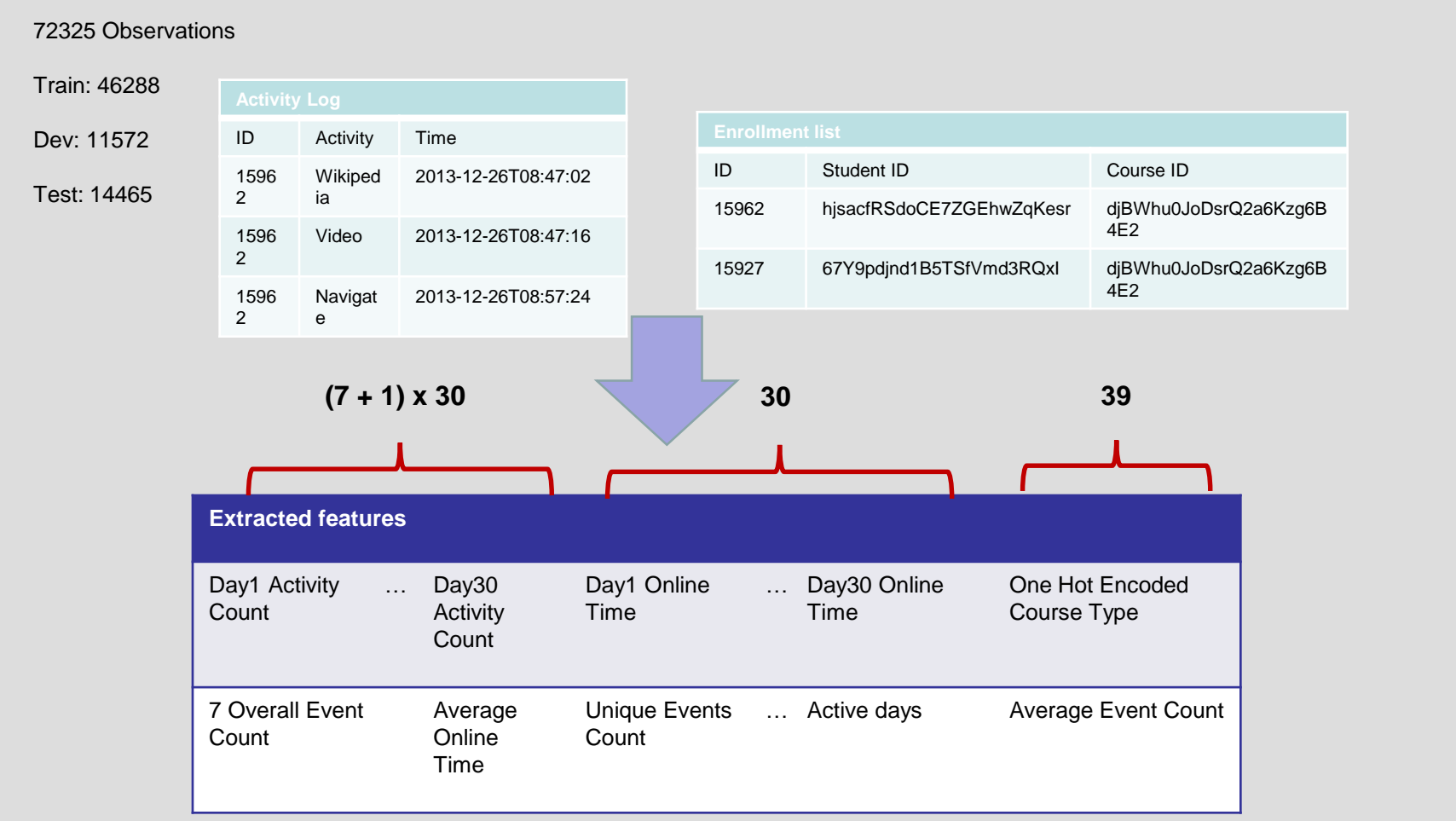
CONTACT

Zixun Yang
 Email: jasonyx@stanford.edu
 SUID# 06236719

GOAL

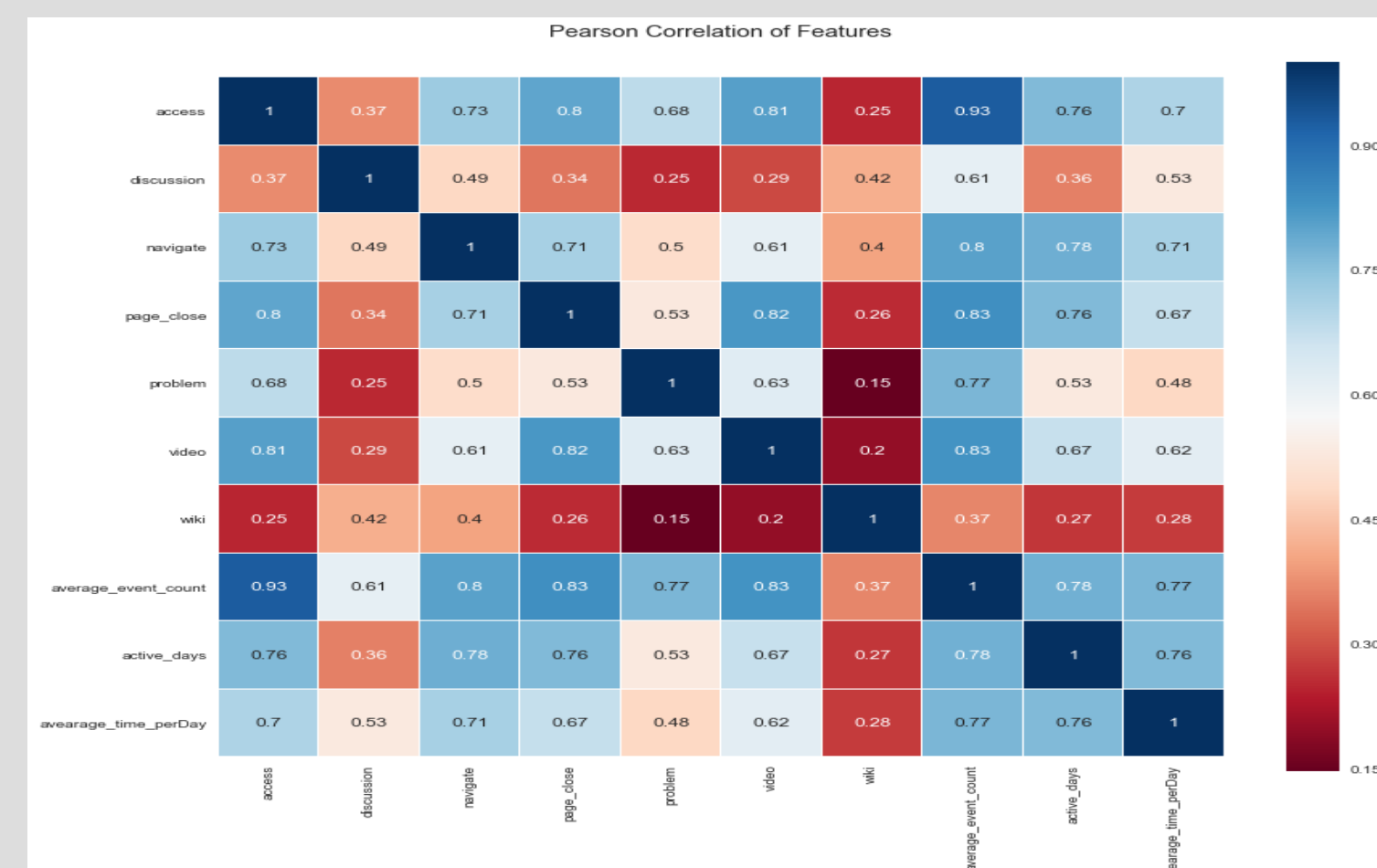
Predict whether or not a student will drop out MOOC.
 Massive Open Online Course (MOOC) has been revolutionizing the way people getting education. However, it also raises up concern that MOOC has very high dropout rate relative to traditional classes. An accurate prediction of dropout becomes very important because it can help MOOC course developers to adaptively tune web designs to students with high dropout rate.

Feature Engineering



Feature Selection

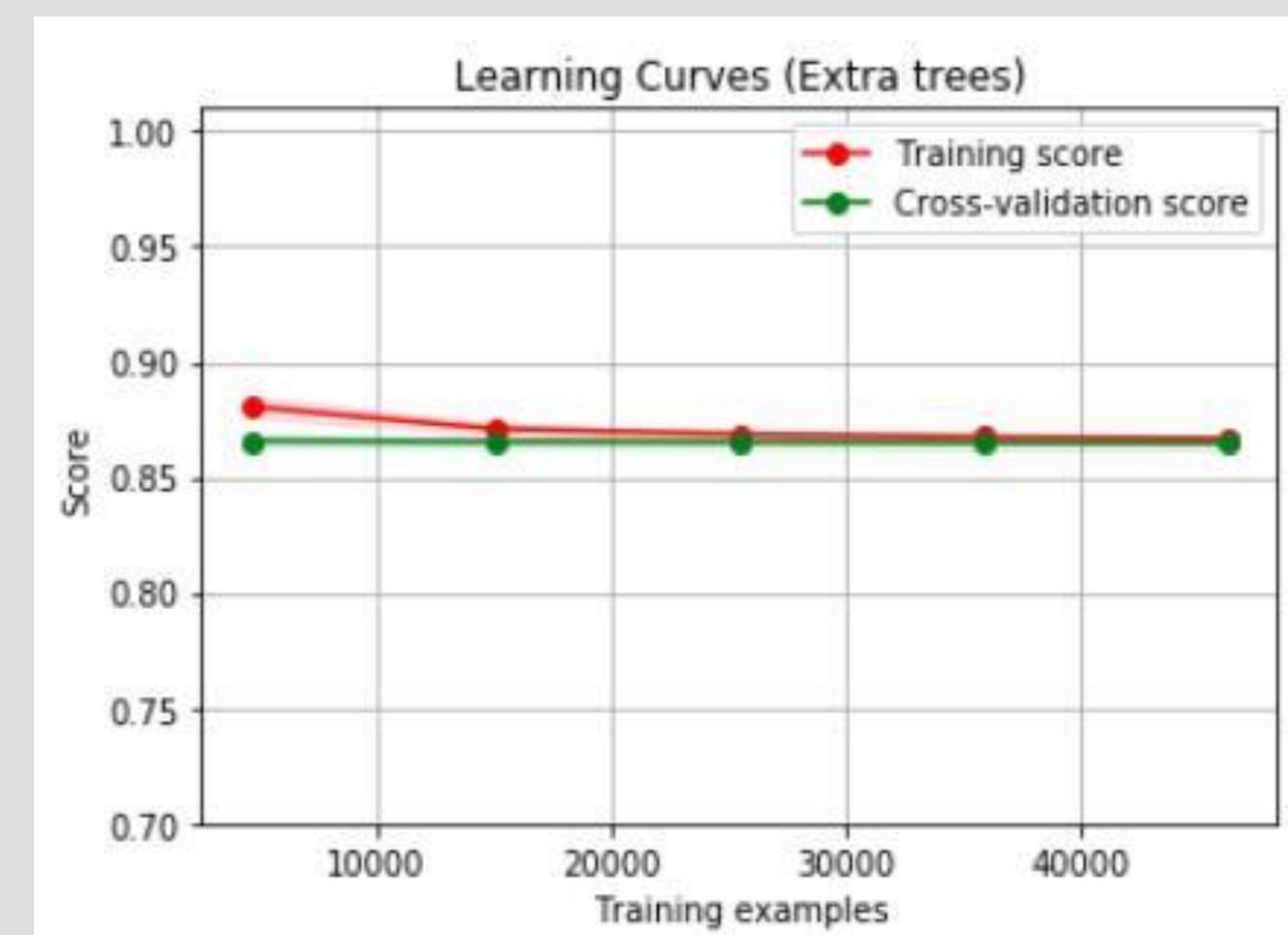
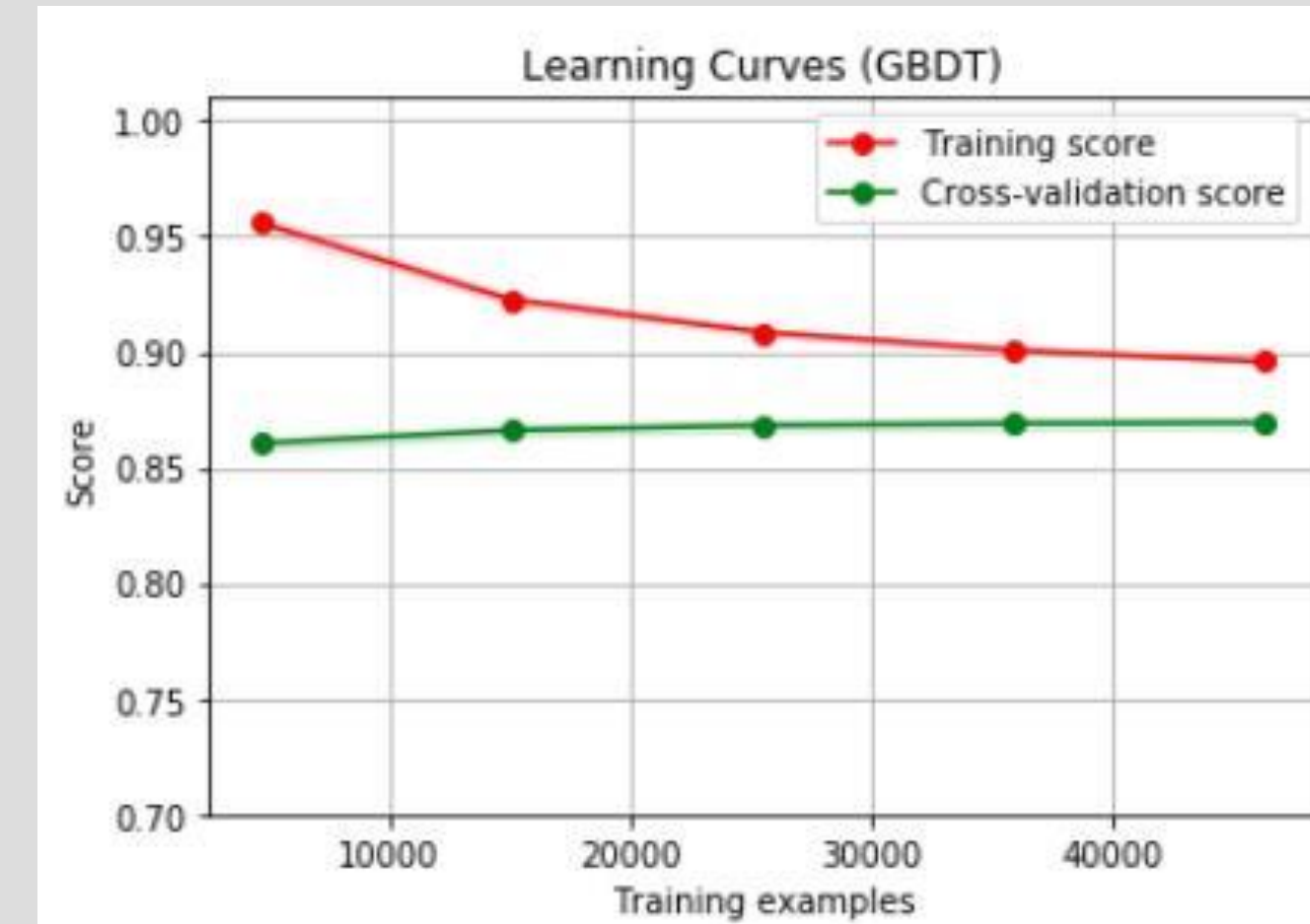
Pearson Correlation of Features



Model Ensemble

First Level of Model Ensemble

We fine-tuned SVC, Logistic Regression, Extra Trees and Gradient Boosting Decision Tree and Random Forest. Each of them have reached accuracy above 86%.



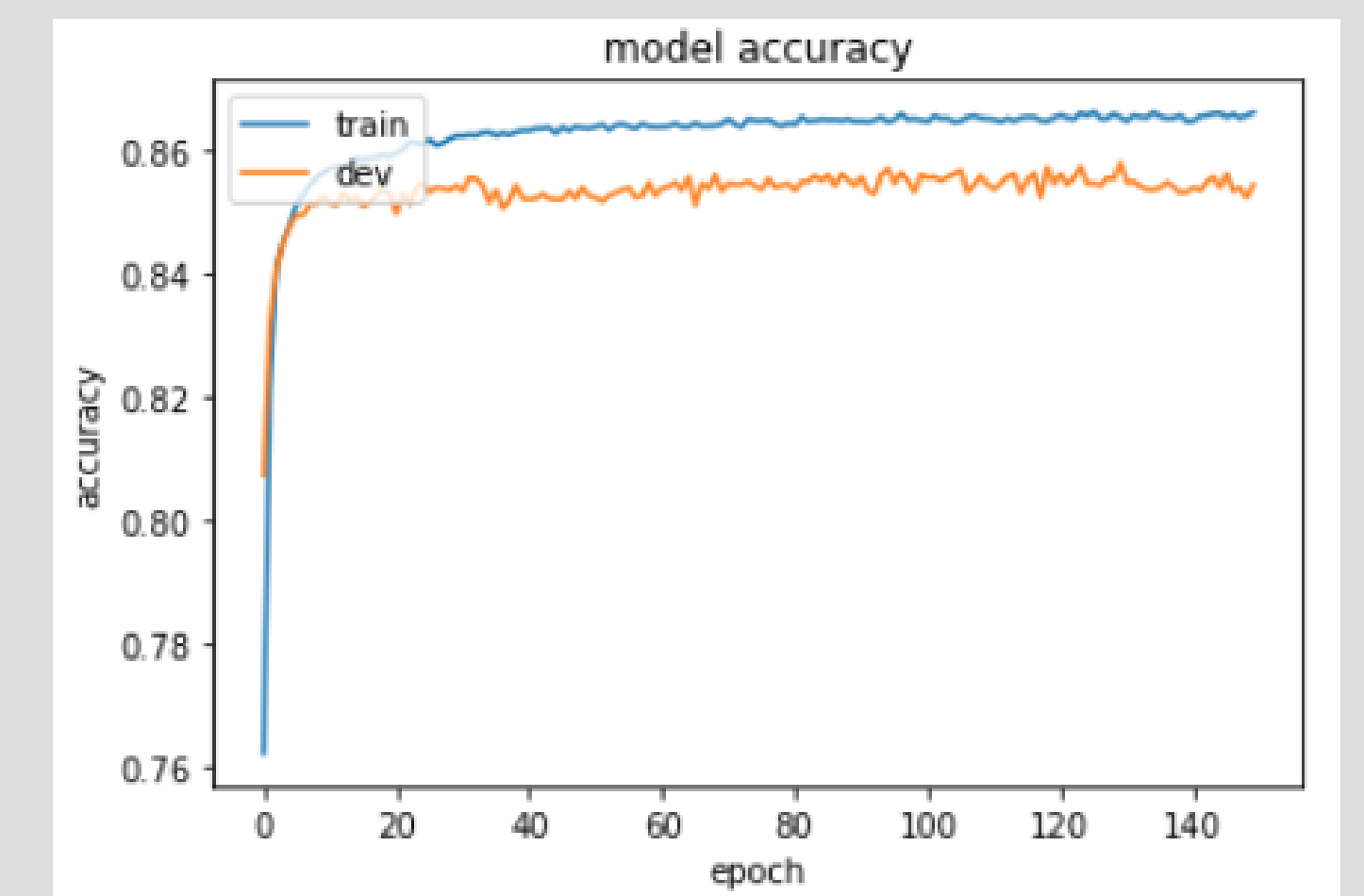
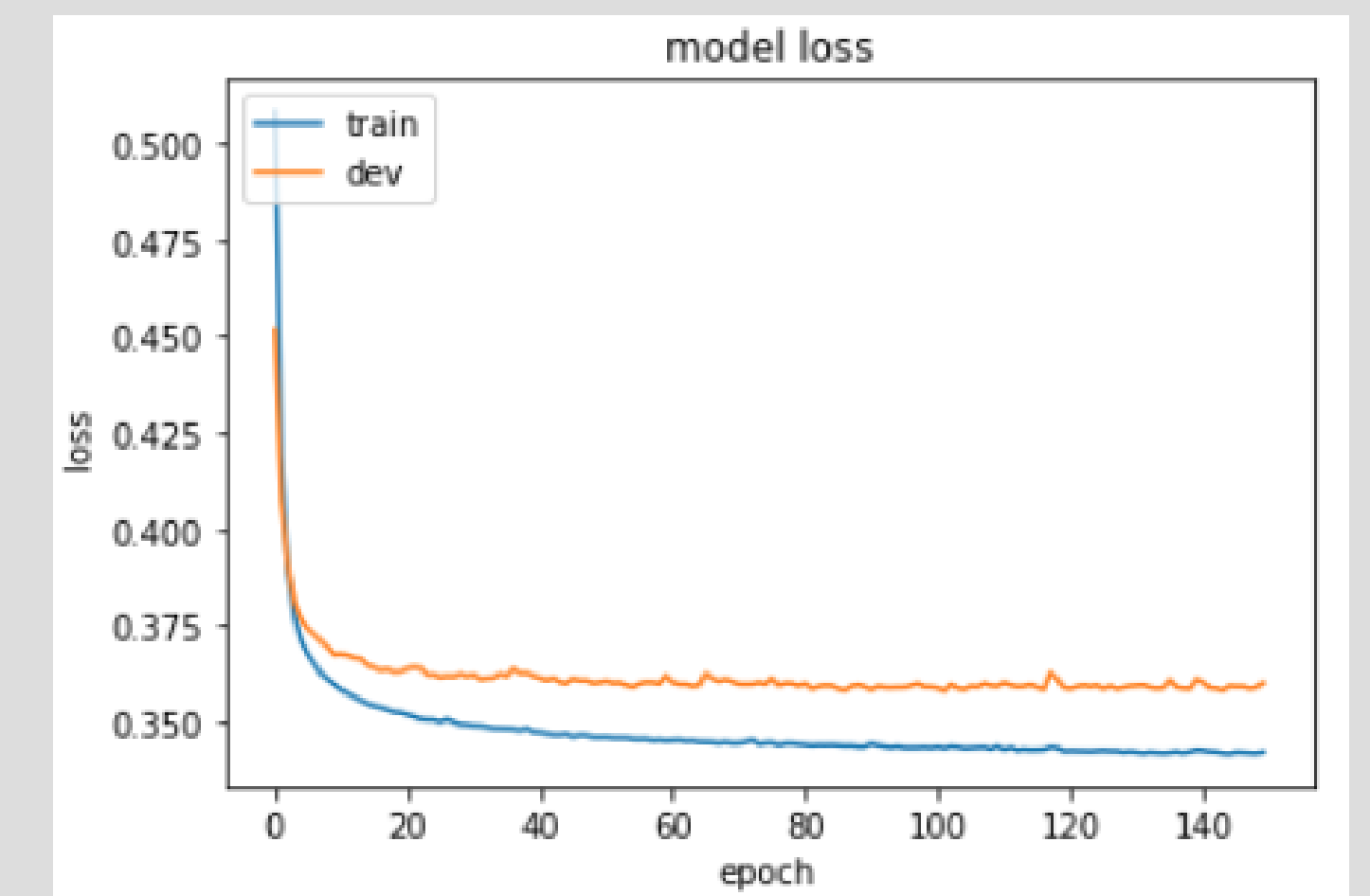
Conclusion

Second Level of Model Ensemble

After having these five classifiers, I ensemble/stacked the results of above five classifiers except LSTM and feed them to **XGboost** classifier for the second level training, reached accuracy of **87.5%**

LSTM

Since our data is time sequence data, and RNN is especially good at process sequence. We used all the day-based features, including seven types of event count, overall event count and online time, to train a stacked two layer LSTM model. Feature Size = 7+1+1 * 30. Below is the learning curve with batch size=200



Reference

- Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." Neural computation 9.8 (1997): 1735-1780.
- Liang, J., Li, C., & Zheng, L. (2016, August). Machine learning application in MOOCs: Dropout prediction. In Computer Science & Education (ICCSE), 2016 11th International Conference on (pp.52-57). IEEE.
- KDD Winner white paper: <http://www.conversionlogic.com/wp-content/uploads/2016/06/Whitepaper-KDD2015-JYL.pdf>