

Multi-Stage Ensemble and Feature Engineering for MOOC Dropout Prediction

Jeong-Yoon Lee
Conversion Logic

12300 Wilshire Blvd. Los Angeles,
CA 90025, USA

jeong@conversionlogic.com

Kohei Ozaki
AIG Japan Holdings KK

Kamiyacho MT Bldg 18F, 4-3-20
Toranomon, Minato-ku, Tokyo
105-0001, Japan

ozaki.kohei@aig.co.jp

Song Chen
American International Group,
Inc. (AIG)

175 Water Street, New York, NY
10038, USA

song.chen@aig.com

Andreas Toescher
Opera Solutions

Hauptplatz 12, 8580 Koeflach,
Austria

andreas.toescher@commendo.at

Mert Bay
Conversion Logic

12300 Wilshire Blvd. Los Angeles,
CA 90025, USA

mert@conversionlogic.com

Tam T. Nguyen
Institute for Infocomm
Research, A*STAR

1 Fusionopolis Way, #21-01
Connexis (South Tower), Singapore
138632

nguyentt@i2r.a-star.edu.sg

Michael Jahrer
Opera Solutions

Hauptplatz 12, 8580 Koeflach,
Austria

michael.jahrer@commendo.at

Peng Yan
NetEase Youdao

2nd Floor, Chuangye Building,
Tsinghua Science Park, Beijing,
100084, China

yanpeng@rd.netease.com

Xiaocong Zhou
Tsinghua University
Haidian District, Beijing, 100084,
China

infinitezxc@gmail.com

ABSTRACT

This paper describes the winning solution of KDD Cup 2015. The competition aims to predict dropouts in Massive Open Online Courses (MOOCs). Our approach begins with feature engineering to extract predictive features from activity logs of students and meta data. Then, we train sixty three individual classifiers with different subsets of features and seven algorithms. Lastly, we blend predictions of individual classifiers with the multi-stage ensemble framework. Our solution achieves AUC scores of 0.90918 and 0.90744 on the public and private leaderboards respectively.

Categories and Subject Descriptors

I.5.4 [Pattern Recognition]: Application

General Terms

Application

Keywords

KDD Cup, Feature Engineering, Ensemble Learning

1. INTRODUCTION

Our final solution is a joint work from 9 data scientists, distributed around the world. The pipeline from raw data to final solution is as follows:

- Hand crafted feature engineering (most of hard work)

- Automatic feature design (autoencoder)
- Individual models (gbm, nn, factor model,..)
- Stage-I ensemble (blends individual models)
- Stage-II ensemble (blends stage-I ensemble models)
- Stage-III ensemble (blends stage-II ensemble models)

2. DATASET

We got history from 200k enrollments, from 120k we know the labels. bla bla
bla bla

3. FEATURE ENGINEERING

All features.

3.1 Feature SK

Song and Kohei feature.

3.2 Feature RW

Peng and Xiaocong feature.

3.3 Feature TN

Tam feature.

3.4 Feature MJ

Michael feature.

3.5 Feature MB

Mert feature.

3.6 Feature JL

Jeong feature.

4. INDIVIDUAL MODELS

something like this.

4.1 Learning Algorithms

- Logistic Regression (LR)
- Kernel Ridge Regression (KRR)
- Factorization Machine/Field-aware Factorization Machine (FM/FFM)
- Neural Networks (NN)
- Extreme Trees (ET)
- Gradient Boosting Decision Trees (GBDT)

4.2 Individual Models

- Model 1: RandomForest(R). Dataset: X
- Model 2: Logistic Regression(scikit). Dataset: $\text{Log}(X+1)$
- Model 3: Extra Trees Classifier(scikit). Dataset: $\text{Log}(X+1)$ (but could be raw)
- Model 4: KNeighborsClassifier(scikit). Dataset: $\text{Scale}(\text{Log}(X+1))$
- Model 5: libfm. Dataset: Sparse(X). Each feature value is a unique level.

5. MULTI-STAGE ENSEMBLE

Stratified 5-fold cross validation (CV). We used xgboost [?], neural nets and linear regression for stage-II ensembling.

5.1 Stage-I Ensemble

We trained 20 stage-I ensemble classifiers with different subsets of CV predictions of 63 individual classifiers.

5.2 Stage-II Ensemble

We trained 2 stage-II ensemble classifiers with different subsets of CV predictions of 20 stage-I ensemble classifiers.

5.3 Stage-III Ensemble

We trained a stage-III ensemble classifier with CV predictions of 5 classifiers: 2 stage-II ensemble, 2 stage-I ensemble, and 1 individual classifiers.

Finally we ended up with linear ensembling with 39 courseID correction factors. These 39 factors improved the score from 0.90910 to 0.90918.

6. CONCLUSIONS

Our final AUC score of 0.90918 results from a complex pipeline from raw data to final score. Every part of that pipe needs to be (sub-)optimal implemented by our team to get the best score at the end. The first part “feature design” is the most important one and needs expertise, experience and of course a bit luck to capture all signals in the data.

7. ACKNOWLEDGEMENTS

Thanks to dropbox, github and skype to enable easy communication around the globe.