

15D012 Advanced Computational Methods

Denitsa Panova, Felix Gutmann, Thomas Vicent

17.03.2016

1 Introduction

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

2 Review of Models and Predictions

Table 1 gives a broad overview of all the applied methods with the corresponding results on the leaderboard and on cross validation (if available). As depicted we tried an extensive amount of different models. In order to be concise we will briefly line out some basic approaches and in a little more detail the approach that worked best for us (both in performance and applicability).

2.1 Basic "off the shelf" approaches and benchmark

One of the first models we applied were some penalized (multinomial-) GLM's using Ridge, Lasso and Elastic Net regularization. This should give us another internal benchmark for further procedure. It turned out that all three didn't perform well isolated (leaderboard performance ≈ 0.50).¹ (Hmmmmmm)

SVM's are a widely used method. The `e1071` package provides a flexible setting for SVM's in case of tuning. It provides several kernel types and classification methods (`c` and `nu`) and arguments for class imbalances. A grid search was performed to tune both `C` (cost for misclassification) and `gamma` (margin). However, we stopped it after approximately 10 hours without a result (grid search with only two parameters for each argument). A submission with radial kernel and basic config obtained a score of 0.506 on the leaderboard. Mainly due to the lag of speed, but also because of tuning expenses and off the shelf performance SVM wasn't considered anymore for the following process.

The best "off the shelf" methods turned out to be random forest and boosting. Especially random forest performed extremely well. Due to the speed of the model it is very easy to tune. Table 1 shows some of the results for random forest in different settings. At the end it turned out that the overall second best performance was based on a random forest. We found this result quite early in the process. Based on the interim evaluation results this was adapted benchmark model.

2.2 Manual ensemble approach.

We briefly describe our so far most promising approach in more detail. Instead of running isolated models, we used different models and average prediction results.

The first try in that direction uses the `superLearner` package of R, which allows to run several model types combined. Furthermore, it provides internal cross validation and parallelized implementation for some models, which increases performance. Since it doesn't support multiple classification, we transformed the outcome variable to a binary matrix and trained a full model using a mix of `gbm`, random forest and a lasso glm. Despite the rather good performance on the leaderboard (0.543), it couldn't beat the best random forest specification so far. Despite parallelizing there is still a speed issue. A full training of the model and prediction lasts approximately six hours.

We wanted to exploit this more diverse prediction power in a more flexible setting regarding both model specification and speed performance. We use the `H2O` package, which provides an optimized performance of several different models. Furthermore, the supported methods are also parallelized, which brings again computational gains. Concerning the data we dropped the minority classes 4 and 5 to let the models focus more on the core classes. Furthermore we dropped observations, which in our opinion were questionable (as discussed in the introduction). In this setting we applied the following models for our ensembling:

- random forests
- `gbm`

¹Note, that we used it again in one ensemble approach, see later on

- (deep-) neural networks
- extremely randomized trees
- xgboosts
- rule-based models

Using such an extensive number of models a grid search of optimal parameters didn't seem to be feasible in reasonable amount of time. Therefore, we re-run all models several times and randomly choose parameters out of a reasonable parameter range ². To monitor our performance we randomly subset roughly 66% of the training data in each iteration as a input data set and use the other part as a test set to validate our models. This procedure turned out to be a solid working basis, because results could be produced fast within approximately less than an hour. Instead of directly predicting classes we used the probabilities for each classes produced by each model. Since we average at the end over all predictions of each model, we believe this provides a more granular outcome. This procedure gave us the best result on the public leaderboard with an accuracy of 0.557.

Experience with stacking

Conclusion and possible improvements

So far we saw that the best approach so far was the manual ensemble approach. Despite the fact it performed quite good, it also seem to have an "upper bound". Concerning the model setting the only thing would probably be to include SVM's using randomized parameter search.³ So far in our opinion the only thing that can be done to increase the prediction performance is construct new features, which add more predictive power.

²A "brute force" approach was considered using an AWS instance. Note that, while using the superLearner we set up and configured an instance on AWS to outsource some computational work and tuning efforts. However, it turned out that the cost benefit ratio here was not appropriate mainly due to stability and monetary issues.

³In the late process we experimented with some more "unorthodox" approaches. This involved tranforming the data using unsupervised techniques (principle components, clustering etc.) and perform prediction on that data set. The idea was try to condense and rearrange informations in a different setting. However using a random forest gave an (5-fold) cv-average accuracy of only around 0.51.[OPTIONAL]

Appendix

Table 1: Overview of Methods and results

Method	Specification	CV	LB	Package
Penalized multonomial GLM	-	-		glmnet
... Lasso	-	-	0.506	
... Elastic Net	-	-		
... Ridge	-	-		
Support Vector Machines	Kernel: Radial		0.496	e1071
Random Forest				randomForest
... Off the shelf	-	-		randomForest
... Balanced DS	-	-	0.506	randomForest, Smooth
Boosting				randomForest
Super Learner	RandomForest, Boosting,		0.543	SuperLearner
Manual Ensemble	RandomForest, LassoGLM,		0.543	H ₂ O
Unsupervised		0.516		Several