

Отчет по Competition 3

Yandex SHAD & MIPT FIVT, ML, Spring 2015 [Kaggle.com]

Вихрева Мария, ВМК МГУ

13 мая 2015

Формулировка задачи

Yandex SHAD & MIPT FIVT, ML, Spring 2015 [Kaggle.com]

Дано: "важности" входящих в статью терминов

Задача: предсказать категории статьи

Функционал качества: $F1 = 2 \frac{pr}{p+r}$, где $p = \frac{tp}{tp+fp}$, $r = \frac{tp}{tp+fn}$

- 25640 безымянных признаков
- train – 10000 объектов, test – 10000 объектов
- матрица признаков – разреженная
- 83 категории

kNN (F1=0.30)

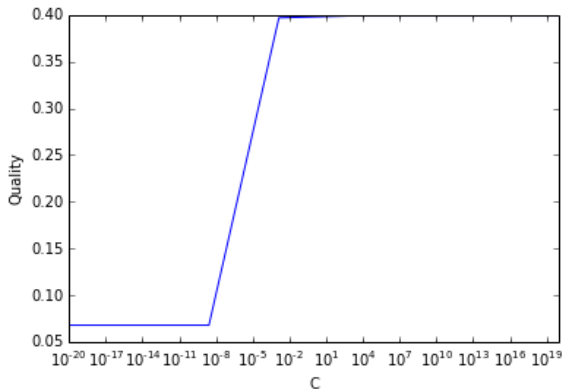
- обернут в OneVsRest классификатор
- веса соседей в соответствии с расстоянием до них

SVM (F1=0.39)

- обернут в OneVsRest классификатор
- на базе libsvm
- ядро - линейное

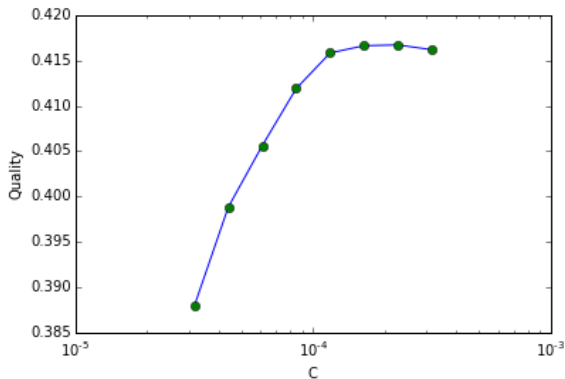
Настройка параметра C

Quality – средний F1-score на нескольких фолдах (фолдов 3)



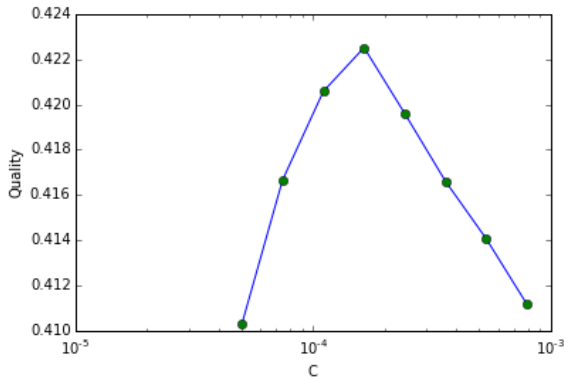
Настройка параметра C

Quality – средний F1-score на нескольких фолдах (фолдов 3)



Настройка параметра C

Quality – средний F1-score на нескольких фолдах (фолдов 7)

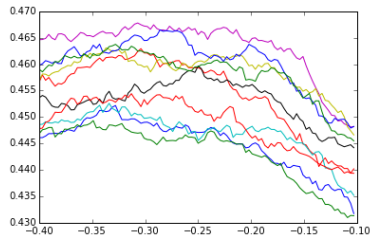
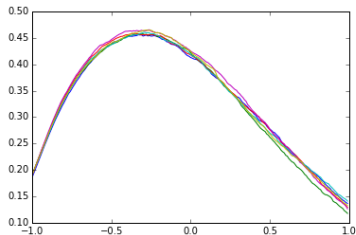


$C=0.0017$, CrossValidation=0.42, Leaderboard=0.43676

Настройка порога отсечения

Отступы объектов – SVM.decision_function

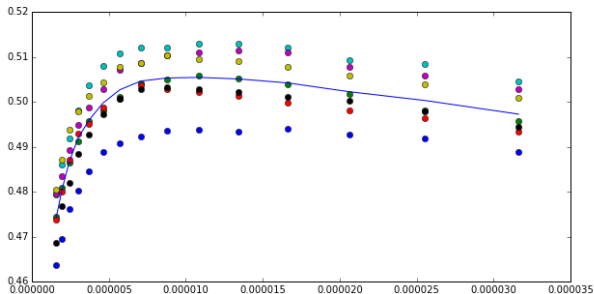
Quality – F1-score на нескольких фолдах (фолдов 7)



$C=0.0017$, $\text{threshold}=-0.32$, $\text{CrossValidation}=0.46$,
 $\text{Leaderboard}=0.46773$

Настройка C и порога одновременно

Quality – F1-score на нескольких фолдах (фолдов 7)



$C=1e-05$, threshold=-0.4, CrossValidation=0.505,
Leaderboard=0.51376

Привет!