

集成学习学习笔记

Ensemble
Bagging (Random Forests)
Stacking
Boosting
AdaBoost
RegionBoost
GBDT(Gradient Boost Decision Tree)
XGBoost

GitHub <https://github.com/kongzz311/MachineLearningNotes>

假如有帮助欢迎加星

水平有限，若有错误欢迎指出：kozenzei@outlook.com

1. 集成学习学习笔记

1.1. Ensemble

Motivations

- 提高准确度
- 降低选择差模型的概率(Model Selection)

常见

- Bagging
- Boosting

前提分类器不同

- 不同的学习算法
- 不同的训练过程
 - 不同参数
 - 不同训练集
 - 不同特征集

弱分类器

- 容易产生不同的decision boundaries
- Stumps...(树桩)

怎样combine不同分类器

- 平均
- 投票

- Majority Voting
 - Random Forest
- Weighted Majority Voting(加权)
 - AdaBoost
- Learning Combiner
 - General Combiner
 - Stacking
 - Piecewise Combiner
 - RegionBoost
- No Free Lunch

1.2. Bagging (Random Forests)

Bootstrap Aggregation(Bagging)

- 有放回的采样(Resample with Replacement)

对D采样K次独立训练分类器然后voting

Random Forests

- 有放回的采样(Resample with Replacement)
- 使用大约 2/3 的原始数据 (why📥)

$$1 - \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n$$

- Majority Voting
- Number of Variables

$$\sqrt{k}$$

- k = available variables
- Number of Trees
 - 500 or more
- Self-Testing
 - Around one third of the original data are left out.
 - Out of Bag(OOB)
 - Similar to Cross-Validation

优点

- All data can be used in the training process
 - Data in **OOB** are used to evaluate the current tree
- High levels of predictive accuracy
 - 只要调少量的参数，如树的数量
 - 分类和回归都可以
- Resistant to overfitting

- No need for prior feature selection

1.3. Stacking

先生成不同的分类器再在输出上进行学习**权重**（模型）

Base Classifiers → Meta Classification

1.4. Boosting

Boosting

串行，先生成第一个C1，再在第一个的基础上生成后续模型

加入之前分错的数据进来

得到C2

加入分类器分类不一致的数据进来

得到C3

用来学习C1 和 C2 分类不同的

Boosting 和 bagging 的区别

Bagging aims at reducing **variance**, not **bias**.

In Boosting, classifiers are generated **sequentially**. Focuses on most informative data points.

Training samples are **weighted**. Outputs are combined via **weighted** voting. Can create arbitrarily **strong** classifiers. The base learners can be arbitrarily **weak**(比瞎猜好就可以). As long as they are better than random forests

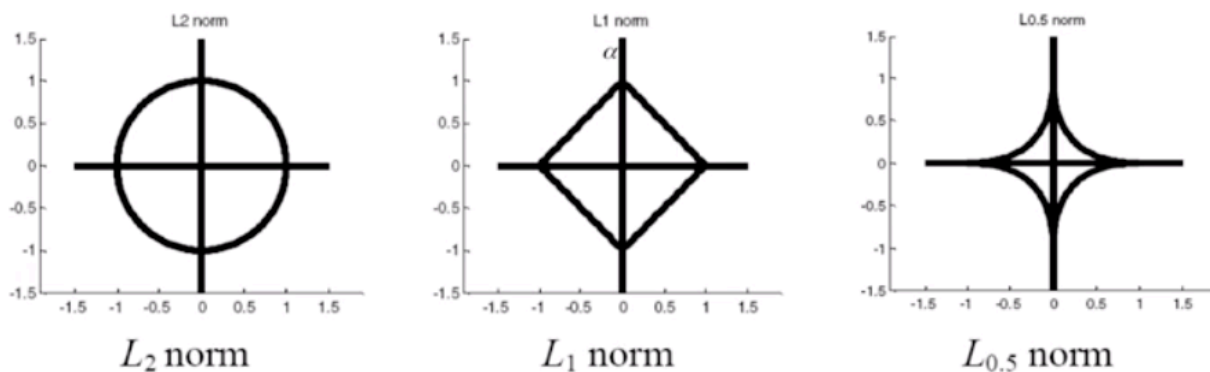
1.5. AdaBoost

1.6. RegionBoost

动态调整权重 (1998)

计算 $\alpha_j(x_i)$

- 通过KNN
- 计算正确率



1.7. GBDT(Gradient Boost Decision Tree)

与AdaBoost不同，**GBDT**每一次的计算是都为了减少上一次的残差，进而在残差减少（负梯度）的方向上建立一个新的模型。GBDT与Adaboost最主要的区别在于两者如何识别模型的问题。Adaboost用错分数据点来识别问题，通过调整错分数据点的权重来改进模型。GBDT通过负梯度来识别问题，通过计算负梯度来改进模型。

```
# Gradient Boosting
from sklearn.ensemble import GradientBoostingClassifier
clf = GradientBoostingClassifier()
# n_estimators = 100 (default)
# loss function = deviance(default) used in Logistic Regression
clf.fit(x_train,y_train)
clf.predict(x_test)
```

1.8. XGBoost

GBDT算法只利用了一阶的导数信息，xgboost对损失函数做了二阶的泰勒展开，并在目标函数之外加入了正则项对整体求最优解，用以权衡目标函数的下降和模型的复杂程度，避免过拟合。所以不考虑细节方面，两者最大的不同就是目标函数的定义。

```
import xgboost as xgb
# read in data
dtrain = xgb.DMatrix('demo/data/agaricus.txt.train')
dtest = xgb.DMatrix('demo/data/agaricus.txt.test')
# specify parameters via map
param = {'max_depth':2, 'eta':1, 'silent':1, 'objective':'binary:logistic' }
num_round = 2
bst = xgb.train(param, dtrain, num_round)
# make prediction
preds = bst.predict(dtest)
```

>