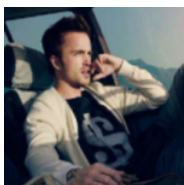


余音、未散的博客

目录视图

个人资料



余音、未散

关注

发私信



访问：64660次

积分：1726

等级：**BLOG > 4**

排名：千里之外

原创：108篇

转载：12篇

译文：0篇

评论：25条

github地址

<https://github.com/lytforgood>

文章搜索

文章分类

机器学习 (31)

R (41)

算法 (11)

天池比赛 (6)

机器学习算法MapReduce版 (3)

Hadoop相关 (25)

Spark (15)

linux (22)

Python (8)

爬虫 (6)

深入学习JVM笔记 (3)

Java (12)

神经网络 (11)

面试笔记 (4)

文章存档

xgboost使用调参

2016-12-20 15:14

664人阅读

评论

分类：**机器学习 (30)**

版权声明：本文为博主原创文章，未经博主允许不得转载。

目录(?)

[+]

github：<https://github.com/dmlc/xgboost>

论文参考：<http://www.kaggle.com/blobs/download/forum-message-attachment-files/4087/xgboost-paper.pdf>

基本思路及优点

<http://blog.csdn.net/q383700092/article/details/60954996>

参考<http://dataunion.org/15787.html>

<http://blog.csdn.net/china1000/article/details/51106856>

在有监督学习中，我们通常会构造一个目标函数和一个预测函数，使用训练样本对目标函数最小化学习到相关的参数，然后用预参数来对未知的样本进行分类的标注或者数值的预测。

1. Boosting Tree构造树来拟合残差，而Xgboost引入了二阶导来进行求解，并且引入了节点的数目、参数的L2正则来评估模型预测函数与目标函数。

2. 在分裂点选择的时候也以目标函数最小化为目标。

优点：

1. 显示的把树模型复杂度作为正则项添加到优化目标中。
2. 公式推导中用到了二阶导数，用了二阶泰勒展开。（GBDT用牛顿法貌似也是二阶信息）
3. 实现了分裂点寻找近似**算法**。
4. 利用了特征的稀疏性。
5. 数据事先排序并且以block形式存储，有利于并行计算。
6. 基于分布式通信框架rabit，可以运行在MPI和yarn上。（最新已经不基于rabit了）
7. 实现做了面向体系结构的优化，针对cache和内存做了性能优化。

原理推导及与GBDT区别

<http://blog.csdn.net/q383700092/article/details/60954996>

参考<http://dataunion.org/15787.html>

<https://www.zhihu.com/question/41354392>

余音、未散 :@zhangxiaozy123:参考新更新文档，已经改变用法<http://blog.csdn.net/...>

R可视化绘图三-recharts
zhangxiaozy123 :您好，下载echartR.R的链接已经失效，应该怎么才能运用echartR函数呢？现在无法下载到本地...

类别十分不平衡

scale_pos_weight = 1

```
1 from xgboost import XGBClassifier
2 xgb1 = XGBClassifier(
3     learning_rate =0.1,
4     n_estimators=1000,
5     max_depth=5,
6     min_child_weight=1,
7     gamma=0,
8     subsample=0.8,
9     colsample_bytree=0.8,
10    objective= 'binary:logistic',
11    nthread=4,
12    scale_pos_weight=1,
13    seed=27)
```

第二步：max_depth 和 min_weight 参数调优

grid_search参考

http://scikit-learn.org/stable/modules/generated/sklearn.grid_search.GridSearchCV.html

<http://blog.csdn.net/abcjennifer/article/details/23884761>

网格搜索scoring='roc_auc' 只支持二分类，多分类需要修改scoring(默认支持多分类)

```
1 param_test1 = {
2     'max_depth':range(3,10,2),
3     'min_child_weight':range(1,6,2)
4 }
5 #param_test2 = {
6     'max_depth':[4,5,6],
7     'min_child_weight':[4,5,6]
8 }
9 from sklearn import svm, grid_search, datasets
10 from sklearn import grid_search
11 gsearch1 = grid_search.GridSearchCV(
12     estimator = XGBClassifier(
13         learning_rate =0.1,
14         n_estimators=140, max_depth=5,
15         min_child_weight=1,
16         gamma=0,
17         subsample=0.8,
18         colsample_bytree=0.8,
19         objective= 'binary:logistic',
20         nthread=4,
21         scale_pos_weight=1,
22         seed=27),
23     param_grid = param_test1,
24     scoring='roc_auc',
25     n_jobs=4,
26     iid=False,
27     cv=5)
28 gsearch1.fit(train[predictors],train[target])
29 gsearch1.grid_scores_, gsearch1.best_params_, gsearch1.best_score_
30 #网格搜索scoring='roc_auc' 只支持二分类，多分类需要修改scoring(默认支持多分类)
```

第三步：gamma参数调优

```
1 param_test3 = {
2     'gamma':[i/10.0 for i in range(0,5)]
3 }
4 gsearch3 = GridSearchCV(
5     estimator = XGBClassifier(
6         learning_rate =0.1,
7         n_estimators=140,
8         max_depth=4,
9         min_child_weight=6,
10        gamma=0,
11        subsample=0.8,
12        colsample_bytree=0.8,
13        objective= 'binary:logistic',
14        nthread=4,
15        scale_pos_weight=1,
```

```

16 seed=27),
17 param_grid = param_test3,
18 scoring='roc_auc',
19 n_jobs=4,
20 iid=False,
21 cv=5)
22 gsearch3.fit(train[predictors], train[target])
23 gsearch3.grid_scores_, gsearch3.best_params_, gsearch3.best_score_

```

第四步：调整subsample 和 colsample_bytree 参数

```

1 #取0.6, 0.7, 0.8, 0.9作为起始值
2 param_test4 = {
3     'subsample':[i/10.0 for i in range(6, 10)],
4     'colsample_bytree':[i/10.0 for i in range(6, 10)]
5 }
6
7 gsearch4 = GridSearchCV(
8     estimator = XGBClassifier(
9         learning_rate =0.1,
10        n_estimators=177,
11        max_depth=3,
12        min_child_weight=4,
13        gamma=0.1,
14        subsample=0.8,
15        colsample_bytree=0.8,
16        objective= 'binary:logistic',
17        nthread=4,
18        scale_pos_weight=1,
19        seed=27),
20    param_grid = param_test4,
21    scoring='roc_auc',
22    n_jobs=4,
23    iid=False,
24    cv=5)
25 gsearch4.fit(train[predictors], train[target])
26 gsearch4.grid_scores_, gsearch4.best_params_, gsearch4.best_score_

```

第五步：正则化参数调优

```

1 param_test6 = {
2     'reg_alpha':[1e-5, 1e-2, 0.1, 1, 100]
3 }
4 gsearch6 = GridSearchCV(
5     estimator = XGBClassifier(
6         learning_rate =0.1,
7         n_estimators=177,
8         max_depth=4,
9         min_child_weight=6,
10        gamma=0.1,
11        subsample=0.8,
12        colsample_bytree=0.8,
13        objective= 'binary:logistic',
14        nthread=4,
15        scale_pos_weight=1,
16        seed=27),
17    param_grid = param_test6,
18    scoring='roc_auc',
19    n_jobs=4,
20    iid=False,
21    cv=5)
22 gsearch6.fit(train[predictors], train[target])
23 gsearch6.grid_scores_, gsearch6.best_params_, gsearch6.best_score_

```

第六步：降低学习速率

```

1 xgb4 = XGBClassifier(
2     learning_rate =0.01,
3     n_estimators=5000,
4     max_depth=4,
5     min_child_weight=6,
6     gamma=0,
7     subsample=0.8,
8     colsample_bytree=0.8,
9     reg_alpha=0.005,
10    objective= 'binary:logistic',

```

```

11 | nthread=4,
12 | scale_pos_weight=1,
13 | seed=27)
14 | modelfit(xgb4, train, predictors)

```

python示例

```

1 | import xgboost as xgb
2 | import pandas as pd
3 | #获取数据
4 | from sklearn import cross_validation
5 | from sklearn.datasets import load_iris
6 | iris = load_iris()
7 | #切分数据集
8 | X_train, X_test, y_train, y_test = cross_validation.train_test_split(iris.data, iris.target, test_size=
9 | #设置参数
10 | m_class = xgb.XGBClassifier(
11 |     learning_rate =0.1,
12 |     n_estimators=1000,
13 |     max_depth=5,
14 |     gamma=0,
15 |     subsample=0.8,
16 |     colsample_bytree=0.8,
17 |     objective= 'binary:logistic',
18 |     nthread=4,
19 |     seed=27)
20 | #训练
21 | m_class.fit(X_train, y_train)
22 | test_21 = m_class.predict(X_test)
23 | print "Accuracy : %.2f" % metrics.accuracy_score(y_test, test_21)
24 | #预测概率
25 | #test_2 = m_class.predict_proba(X_test)
26 | #查看AUC评价标准
27 | from sklearn import metrics
28 | print "Accuracy : %.2f" % metrics.accuracy_score(y_test, test_21)
29 | ##必须二分类才能计算
30 | ##print "AUC Score (Train): %f" % metrics.roc_auc_score(y_test, test_2)
31 | #查看重要程度
32 | feat_imp = pd.Series(m_class.booster().get_fscore()).sort_values(ascending=False)
33 | feat_imp.plot(kind='bar', title='Feature Importances')
34 | import matplotlib.pyplot as plt
35 | plt.show()
36 | #回归
37 | #m_regress = xgb.XGBRegressor(n_estimators=1000,seed=0)
38 | #m_regress.fit(X_train, y_train)
39 | #test_1 = m_regress.predict(X_test)

```

整理

xgb原始

```

1 | from sklearn.model_selection import train_test_split
2 | from sklearn import metrics
3 | from sklearn.datasets import make_hastie_10_2
4 | import xgboost as xgb
5 | #记录程序运行时间
6 | import time
7 | start_time = time.time()
8 | X, y = make_hastie_10_2(random_state=0)
9 | X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.5, random_state=0)##test_size测试
10 | #xgb矩阵赋值
11 | xgb_train = xgb.DMatrix(X_train, label=y_train)
12 | xgb_test = xgb.DMatrix(X_test, label=y_test)
13 | ##参数
14 | params={
15 |     'booster': 'gbtree',
16 |     'silent': 1, #设置成1则没有运行信息输出，最好是设置为0。
17 |     'nthread': 7, # cpu 线程数 默认最大
18 |     'eta': 0.007, # 如同学习率
19 |     'min_child_weight': 3,
20 |     # 这个参数默认是 1，是每个叶子里面 h 的和至少是多少，对正负样本不均衡时的 0-1 分类而言
21 |     #，假设 h 在 0.01 附近，min_child_weight 为 1 意味着叶子节点中最少需要包含 100 个样本。
22 |     #这个参数非常影响结果，控制叶子节点中二阶导的和的最小值，该参数值越小，越容易 overfitting。
23 |     'max_depth': 6, # 构建树的深度，越大越容易过拟合
24 |     'gamma': 0.1, # 树的叶子节点上作进一步分区所需的最小损失减少，越大越保守，一般 0.1、0.2 这样子。

```

```

25 'subsample':0.7, # 随机采样训练样本
26 'colsample_bytree':0.7, # 生成树时进行的列采样
27 'lambda':2, # 控制模型复杂度的权重值的L2正则化项参数，参数越大，模型越不容易过拟合。
28 '#alpha':0, # L1 正则项参数
29 '#scale_pos_weight':1, #如果取值大于0的话，在类别样本不平衡的情况下有助于快速收敛。
30 '#objective': 'multi:softmax', #多分类的问题
31 '#num_class':10, # 类别数，多分类与 multisoftmax 并用
32 'seed':1000, #随机种子
33 '#eval_metric': 'auc'
34 }
35 plst = list(params.items())
36 num_rounds = 100 # 迭代次数
37 watchlist = [(xgb_train, 'train'), (xgb_test, 'val')]
38
39 #训练模型并保存
40 # early_stopping_rounds 当设置的迭代次数较大时，early_stopping_rounds 可在一定的迭代次数内准确率没有提升
41 model = xgb.train(plst, xgb_train, num_rounds, watchlist, early_stopping_rounds=100, pred_margin=1)
42 #model.save_model('./model/xgb.model') # 用于存储训练出的模型
43 print "best_ntree_limit", model.best_ntree_limit
44 y_pred = model.predict(xgb_test, ntree_limit=model.best_ntree_limit)
45 print ('error=%f' % ( sum(1 for i in range(len(y_pred)) if int(y_pred[i]>0.5)!=y_test[i]) /float(len(y
46 #输出
47 cost_time = time.time()-start_time
48 print "Success!\n", "cost time:", cost_time, "(s)....."
49

```

xgb使用sklearn交叉验证

官方

会改变的函数名是：

eta -> learning_rate

lambda -> reg_lambda

alpha -> reg_alpha

```

1 from sklearn.model_selection import train_test_split
2 from sklearn import metrics
3 from sklearn.datasets import make_hastie_10_2
4 from xgboost.sklearn import XGBClassifier
5 X, y = make_hastie_10_2(random_state=0)
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.5, random_state=0)##test_size测试
7 clf = XGBClassifier(
8     silent=0, #设置成1则没有运行信息输出，最好是设置为0. 是否在运行升级时打印消息。
9     nthread=4, # cpu 线程数 默认最大
10    learning_rate= 0.3, # 如同学习率
11    min_child_weight=1,
12    # 这个参数默认是 1，是每个叶子里面 h 的和至少是多少，对正负样本不平衡时的 0-1 分类而言
13    #，假设 h 在 0.01 附近，min_child_weight 为 1 意味着叶子节点中最少需要包含 100 个样本。
14    #这个参数非常影响结果，控制叶子节点中二阶导的和的最小值，该参数值越小，越容易 overfitting。
15    max_depth=6, # 构建树的深度，越大越容易过拟合
16    gamma=0, # 树的叶子节点上作进一步分区所需的最小损失减少，越大越保守，一般0.1、0.2这样子。
17    subsample=1, # 随机采样训练样本 训练实例的子采样比
18    max_delta_step=0, #最大增量步长，我们允许每个树的权重估计。
19    colsample_bytree=1, # 生成树时进行的列采样
20    reg_lambda=1, # 控制模型复杂度的权重值的L2正则化项参数，参数越大，模型越不容易过拟合。
21    reg_alpha=0, # L1 正则项参数
22    #scale_pos_weight=1, #如果取值大于0的话，在类别样本不平衡的情况下有助于快速收敛。平衡正负权重
23    #objective= 'multi:softmax', #多分类的问题 指定学习任务和相应的学习目标
24    #num_class=10, # 类别数，多分类与 multisoftmax 并用
25    n_estimators=100, #树的个数
26    seed=1000 #随机种子
27    #eval_metric= 'auc'
28 )
29 clf.fit(X_train, y_train, eval_metric='auc')
30 #设置验证集合 verbose=False不打印过程
31 clf.fit(X_train, y_train, eval_set=[(X_train, y_train), (X_val, y_val)], eval_metric='auc', verbose=False)
32 #获取验证集合结果
33 evals_result = clf.evals_result()
34 y_true, y_pred = y_test, clf.predict(X_test)
35 print "Accuracy : %.4g" % metrics.accuracy_score(y_true, y_pred)
36 #回归
37 #m_regress = xgb.XGBRegressor(n_estimators=1000, seed=0)

```

网格搜索

可以先固定一个参数 最优化后继续调整

第一步：确定学习速率和tree_based 给个常见初始值 根据是否类别不平衡调节

max_depth,min_child_weight,gamma,subsample,scale_pos_weight

max_depth=3 起始值在4-6之间都是不错的选择。

min_child_weight比较小的值解决极不平衡的分类问题eg:1

subsample, colsample_bytree = 0.8: 这个是最常见的初始值了

scale_pos_weight = 1: 这个值是因为类别十分不平衡。

第二步：max_depth 和 min_weight 对最终结果有很大的影响

```
'max_depth':range(3,10,2),
```

```
'min_child_weight':range(1,6,2)
```

先大范围地粗调参数，然后再小范围地微调。

第三步：gamma参数调优

```
'gamma':[i/10.0 for i in range(0,5)]
```

第四步：调整subsample 和 colsample_bytree 参数

```
'subsample':[i/100.0 for i in range(75,90,5)],
```

```
'colsample_bytree':[i/100.0 for i in range(75,90,5)]
```

第五步：正则化参数调优

```
'reg_alpha':[1e-5, 1e-2, 0.1, 1, 100]
```

```
'reg_lambda'
```

第六步：降低学习速率

learning_rate =0.01,

```
1 from sklearn.model_selection import GridSearchCV
2 tuned_parameters= [{'n_estimators':[100,200,500],
3                       'max_depth':[3,5,7], ##range(3,10,2)
4                       'learning_rate':[0.5, 1.0],
5                       'subsample':[0.75,0.8,0.85,0.9]
6                       }]
7 tuned_parameters= [{'n_estimators':[100,200,500,1000]
8                       }]
9 clf = GridSearchCV(XGBClassifier(silent=0,nthread=4,learning_rate= 0.5,min_child_weight=1, max_depth=3,
10 clf.fit(X_train, y_train)
11 ##clf.grid_scores_, clf.best_params_, clf.best_score_
12 print(clf.best_params_)
13 y_true, y_pred = y_test, clf.predict(X_test)
14 print"Accuracy : %.4g" % metrics.accuracy_score(y_true, y_pred)
15 y_proba=clf.predict_proba(X_test)[:,:1]
16 print "AUC Score (Train): %f" % metrics.roc_auc_score(y_true, y_proba)
```

```
1 from sklearn.model_selection import GridSearchCV
2 parameters= [{'learning_rate':[0.01,0.1,0.3], 'n_estimators':[1000,1200,1500,2000,2500]}]
3 clf = GridSearchCV(XGBClassifier(
4     max_depth=3,
5     min_child_weight=1,
6     gamma=0.5,
7     subsample=0.6,
8     colsample_bytree=0.6,
9     objective= 'binary:logistic', #逻辑回归损失函数
10    scale_pos_weight=1,
11    reg_alpha=0,
12    reg_lambda=1,
13    seed=27
14 ),
15    param_grid=parameters,scoring='roc_auc')
16 clf.fit(X_train, y_train)
17 print(clf.best_params_)
18 y_pre= clf.predict(X_test)
19 y_pro= clf.predict_proba(X_test)[:,:1]
20 print "AUC Score : %f" % metrics.roc_auc_score(y_test, y_pro)
21 print"Accuracy : %.4g" % metrics.accuracy_score(y_test, y_pre)
```

输出特征重要性

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 feat_imp = pd.Series(clf.booster().get_fscore()).sort_values(ascending=False)
4 feat_imp.plot(kind='bar', title='Feature Importances')
5 plt.ylabel('Feature Importance Score')
6 plt.show()
```

顶 0 踩 0

- [上一篇](#) GBDT 原理与使用
- [下一篇](#) 并发和并行

我的同类文章

机器学习 (30)

- Word2vec原理与应用 2017-03-22 阅读 72
- XGBoost原理与应用 2017-03-09 阅读 182
- 机器学习面试准备(持续更新) 2017-02-20 阅读 414
- sklearn-数据预处理-特征变换 2017-01-16 阅读 162
- 数据处理与模型选择的一些注释 2016-12-26 阅读 198

- gcForest算法理解 2017-03-
- 机器学习面试问题汇总 2017-02-
- 聚类与距离学习笔记 2017-02-
- 结合Scikit-learn介绍几种常用的特征选... 2016-12-
- 对机器学习与数据竞赛的一些总结 2016-12-

更多文章

参考知识库



Python知识库
22793 关注 | 1607 收录



.NET知识库
3734 关注 | 833 收录



软件测试
22 关注



算法与数据结构知识库
6 关注 | 2320 收录

猜你在找

- Python数据分析实战：泰坦尼克...

Python算法实战视频课程--二叉树

Python算法实战视频课程--图

Python算法实战视频课程--队列...

使用决策树算法对测试数据进行...
- windows下在Java中使用xgboost ...

xgboost在windows下的安装与使用

在Python中使用XGBoost

python+xgboost在windows上的安...

使用g++编译C++程序链接时出现...

查看评论

暂无评论

您还没有登录,请[\[登录\]](#)或[\[注册\]](#)

* 以上用户言论只代表其个人观点，不代表CSDN网站的观点或立场

核心技术类目

全部主题	Hadoop	AWS	移动游戏	Java	Android	iOS	Swift	智能硬件	Docker	OpenStack	VP	
IE10	Eclipse	CRM	JavaScript	数据库	Ubuntu	NFC	WAP	jQuery	BI	HTML5	Spring	Apache
HTML	SDK	IIS	Fedora	XML	LBS	Unity	Splashtop	UML	components	Windows Mobile	Rails	
Cassandra	CloudStack	FTC	coremail	OPhone	CouchBase	云计算	iOS6	Rackspace	Web App	S		
Compuware	大数据	aptech	Perl	Tornado	Ruby	Hibernate	ThinkPHP	HBase	Pure	Solr	Android	
Cloud Foundry	Redis	Scala	Django	Bootstrap								

公司简介 | 招贤纳士 | 广告服务 | 联系方式 | 版权声明 | 法律顾问 | 问题报告 | 合作伙伴 | 论坛反馈

网站客服 杂志客服 微博客服 webmaster@csdn.net 400-600-2320 | 北京创新乐知信息技术有限公司 版权所有 | 江苏知之为计算机有限公司 | 江苏乐知网络技术有限公司
京 ICP 证 09002463 号 | Copyright © 1999-2016, CSDN.NET, All Rights Reserved

