

我的XGBoost学习经历及动手实践

原创 李祖贤 Datawhale 2020-06-21

收录于话题

#机器学习专题 19 #数据项目专栏 13

↑↑↑关注后"星标"Datawhale
每日干货 & 每月组队学习，不错过

Datawhale干货

作者：李祖贤 深圳大学，Datawhale高校群成员



知乎地址：<http://www.zhihu.com/people/meng-di-76-92>

我今天主要介绍机器学习集成学习方法中三巨头之一的XGBoost，这个算法在早些时候机器学习比赛内曾经大放异彩，是非常好用的一个机器学习集成算法。

XGBoost是一个优化的分布式梯度增强库，旨在实现高效，灵活和便携。它在Gradient Boosting框架下实现机器学习算法。XGBoost提供了并行树提升（也称为GBDT，GBM），可以快速准确地解决许多数据科学问题。

相同的代码在主要的分布式环境（Hadoop，SGE，MPI）上运行，并且可以解决超过数十亿个样例的问题。XGBoost利用了核外计算并且能够使数据科学家在一个主机上处理数亿的样本数据。最终，将这些技术进行结合来做一个端到端的系统以最少的集群系统来扩展到更大的数据集上。

XGBoost原理介绍

从0开始学习，经历过推导公式的波澜曲折，下面展示下我自己的推公式的手稿吧，希望能激励到大家能够对机器学习数据挖掘更加热爱！

Xgboost: GBDT的推广

1. 优化目标: $Obj = \sum_{i=1}^n l(y_i, \bar{y}_i) + \sum_{k=1}^K \Omega(f_k)$, 其中: $\bar{y} = \sum_{k=1}^K f_k(x)$, $f_k \in \Gamma$

(模型损失函数)
n为样本个数, y_i 为样本真实标签, \bar{y} 模型输出. K表示树的个数, f_k 表示第k棵树, 用于由样本至叶子节点值的映射($x \rightarrow R$), Ω 为模型复杂度函数.

⇒ 任务: 找到一组树 使得Obj最小.

2. 追加法训练 (Additive Training Boosting)

核心思想: 已经训练好的树 $T_1 \sim T_{t-1}$ 不再调整, 那么第 T_t 棵树:

$\bar{y}^{(0)} = 0$, $\bar{y}^{(1)} = f_1(x) = \bar{y}^{(0)} + f_1(x)$, 训练第一棵树.

$\bar{y}^{(2)} = f_1(x) + f_2(x) = \bar{y}^{(1)} + f_2(x)$, 训练第二棵树. 不再调整第一棵树.

$\dots \bar{y}^{(t)} = \sum_{k=1}^t f_k(x) = \bar{y}^{(t-1)} + f_t(x)$ 训练第t棵树. 不再调整第t-1棵树.

⇒ 假设此时对第t棵树训练, 则目标函数为:

$Obj^{(t)} = \sum_{i=1}^n l(y_i, \bar{y}_i^{(t)}) + \sum_{i=1}^t \Omega(f_i) = \sum_{i=1}^n l(y_i, \bar{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) + \underbrace{\sum_{i=1}^{t-1} \Omega(f_i)}_{\text{constant}}$

由Taylor公式=近似: $f(x+\Delta x) \approx f(x) + f'(x)\Delta x + \frac{1}{2}f''(x)(\Delta x)^2$

$Obj^{(t)} = \sum_{i=1}^n [l(y_i, \bar{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t) + \text{constant}$

*注: 将 $\bar{y}_i^{(t-1)}$ 视为 x , $f_t(x_i)$ 视为 Δx , $l(y_i, \bar{y}_i^{(t-1)})$ 视为 $f(x)$, $l(y_i, \bar{y}_i^{(t-1)} + f_t(x_i))$ 视为

$f(x+\Delta x)$, $g_i = \frac{\partial l(y_i, \bar{y}_i^{(t-1)})}{\partial \bar{y}_i^{(t-1)}}$, $h_i = \frac{\partial^2 l(y_i, \bar{y}_i^{(t-1)})}{\partial (\bar{y}_i^{(t-1)})^2}$

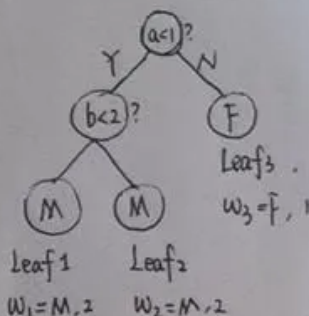
⇒ 当前面t-1棵树已知, 那么 $\sum_{i=1}^n [l(y_i, \bar{y}_i^{(t-1)})] = \text{constant}$

⇒ $Obj^{(t)} = \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \underbrace{\Omega(f_t)}_{\text{还未定义}}$

⇒ 定义 $\Omega(f_t)$

假设第t棵树有T个节点, 叶子节点用的输出向量 $[w_1, w_2, \dots, w_T]$.

那么 $f_t(x) = w_{q(x)}$, $w \in R^T$, $q(x): R^d \rightarrow \{1, 2, \dots, T\}$.



$q(x): R^d \rightarrow \{1, 2, 3\}$.

$f_t(x) = w_{q(x)} = M \text{ or } F$

知乎 @萌弟

XGBoost公式1

XGBoost 定义 $\Omega(f_t) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$

T 为叶子节点数, w_j 为叶子节点 j 的输出, γ 为系数.

$$\Omega = 3 \cdot \gamma + \frac{1}{2} \lambda (2^2 + 2^2 + 1^2) = 3\gamma + \frac{9}{2} \lambda$$

$$\begin{aligned} \Rightarrow \bar{Obj}^{(t)} &= \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t) \\ &= \sum_{i=1}^n [g_i w_{I(x_i)} + \frac{1}{2} h_i w_{I(x_i)}^2] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\ &= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) w_j + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) w_j^2] + \gamma T \end{aligned}$$

其中 $I_j = \{i | I(x_i) = j\}$

$$\text{令 } G_j = \sum_{i \in I_j} g_i, \quad H_j = \sum_{i \in I_j} h_i, \quad \text{则}$$

$$\bar{Obj}^{(t)} = \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + \gamma T. \quad (\text{参数为 } w_j) \Rightarrow \text{累加二次函数, 关于 } w_j, j=1, 2, \dots, T$$

$$\argmin(\bar{Obj}^{(t)}; \underbrace{w_1, \dots, w_T}_{\text{参数}}) = (-\frac{G_1}{H_1 + \lambda}, -\frac{G_2}{H_2 + \lambda}, \dots, -\frac{G_T}{H_T + \lambda})$$

$$\text{即: } w_j^* = -\frac{G_j}{H_j + \lambda}$$

$$\bar{Obj}^{(t)}_{\min} = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{(H_j + \lambda)} + \gamma T$$

~~总结流程:~~

~~① 迭代生成第 j 棵树 T_j~~

~~② 迭代初始, 对每个样本计算 g_i, h_i, \dots~~

3. 生成第 j 棵树的方法:

DT: ID3, C4.5, Gini, 左 右 分裂得分

$$\text{XGBoost: Gain} = \frac{1}{2} \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \lambda$$

Gain 越大, 说明分裂后能使目标函数值减小越多.

寻找最优节点: $\begin{cases} \text{Basic Exact Greedy Algorithm 精确贪心算法. CART.} \\ \text{Approximate Algorithm 近似算法. (V)} \end{cases}$

连续值 \rightarrow 离散化.

由 $1, \dots, m$ 平均值.

知乎 @萌弟

XGBoost公式2

现在我们对手稿的内容进行详细的讲解:

1. 优化目标:

$$Obj = \sum_{i=1}^n l(y_i, \bar{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

其中， n 为样本个数， y_i 为第 i 个样本真实标签， \bar{y}_i 表示模型的第 i 个样本输出值， K 表示树的个数， f_k 表示第 k 棵树，用于样本到叶子结点的映射($x \rightarrow R$)， Ω 为模型复杂度函数。

我们的任务是找到一组树使得Obj最小，很明显这个优化目标Obj可以看成是样本的损失和模型的复杂度惩罚相加组成。

2. 使用追加法训练 (Additive Training Boosting)

核心思想是：在已经训练好了 $T_1 \sim T_{t-1}$ 棵树后不再调整前 T_{t-1} 棵树，那么第 t 棵树可以表示为：

$$\bar{y}^{(t)} = \sum_{k=1}^t f_k(x) = \bar{y}^{(t-1)} + f_t(x)$$

(1). 那此时如果我们对第 t 棵树训练，则目标函数为：

$$Obj^{(t)} = \sum_{i=1}^n l(y_i, \bar{y}_i^{(t)}) + \sum_{i=1}^t \Omega(f_i) = \sum_{i=1}^n l(y_i, \bar{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) + \boxed{\sum_{i=1}^{t-1} \Omega(f_i)}$$

前 $t-1$ 棵树已知下为常数

对上式进行泰勒二阶展开：

$$Obj^{(t)} = \sum_{i=1}^n [l(y_i, \bar{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t) + constant$$

* 注：Taylor公式是 $f(x + \Delta x) \approx f(x) + f'(x)\Delta x + \frac{1}{2} f''(x)\Delta x^2$,

将 $\bar{y}_i^{(t-1)}$ 视为 x ， $f_t(x_i)$ 视为 Δx ， $l(y_i, \bar{y}_i^{(t-1)})$ 视为 $f(x)$ ， $l(y_i, \bar{y}_i^{(t-1)} + f_t(x_i))$ 视为 $f(x + \Delta x)$ ，

$$\text{其中 } g_i = \frac{\partial l(y_i, \bar{y}_i^{(t-1)})}{\partial \bar{y}_i^{(t-1)}}, h_i = \frac{\partial^2 l(y_i, \bar{y}_i^{(t-1)})}{\partial^2 \bar{y}_i^{(t-1)}}$$

由于前 $t-1$ 棵树已知，那么

$\sum_{i=1}^n l(y_i, \bar{y}_i^{(t-1)}) = constant$ 常数，常数不影响我们的优化故删除，得：

$$Obj^{(t)} = \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t)$$

(2). 我们已经对前半部分的损失函数做出了充分的讨论，但是后半部分的 $\Omega(f_t)$ 还只是个符号并未定义，那我们现在就来定义 $\Omega(f_t)$ ：假设我们待训练的第 t 棵树有 T 个叶子结点：叶子结点的输出向量表示如下：

$$[w_1, w_2, \dots, w_T]$$

假设 $q(x) : R^d \rightarrow \{1, 2, 3, \dots, T\}$ 表示样本到叶子结点的映射，那么 $f_t(x) = w_{q(x)}$, $w \in R^T$ 。

那么我们定义：

$\Omega(f_t) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$, 其中 T 为叶子结点数, w_j 为叶子结点 j 的输出, γ 为系数

(3). 我们的目标函数最终化简为:

$$\begin{aligned} OB_j^{(t)} &= \sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\ &= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) w_{q(x_i)} + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) w_j^2] + \gamma T \end{aligned}$$

其中, $I_j = \{i | q(x_i) = j\}$, 令 $G_j = \sum_{i \in I_j} g_i$, $H_j = \sum_{i \in I_j} h_i$, 则最终

$$OB_j^{(t)} = \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + \gamma T$$

我们找到了目标函数就需要对目标函数进行优化:

$$\operatorname{argmin}(OB_j^{(t)}; w_1, \dots, w_T) = \operatorname{argmin}(\sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + \gamma T)$$

我们可以看到, 这是一个累加独立的二次函数, 因此根据二次函数最小点 $x = -\frac{b}{2a}$ 得:

$$\text{最优的参数 } w_j^* = -\frac{G_j}{H_j + \lambda}, \text{ 最小值为: } OB_{j_{min}}^{(t)} = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + \gamma T$$

3. 生成树的策略:

我们刚刚的假设前提是已知前 $t-1$ 棵树, 因此我们现在来探讨怎么生成树。根据决策树的生成策略, 再每次分裂节点的时候我们需要考虑能使得损失函数减小最快的节点, 也就是分裂后损失函数减去分裂前损失函数我们称之为 Gain:

$$\text{某个节点的 } Gain = \frac{1}{2} \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \lambda$$

Gain 越大越能说明分裂后目标函数值减小越多。(因为从式子来看: $\frac{G_j^2}{H_j + \lambda}$ 越大, 反而 OB_j 越小)

4. 寻找最优节点:

- 精确贪心算法 (Basic Exact Greedy Algorithm)

▪ 近似算法 (Approximate Algorithm)

在决策树 (CART) 里面, 我们使用的是精确贪心算法 (Basic Exact Greedy Algorithm), 也就是将所有特征的所有取值排序 (耗时耗内存巨大), 然后比较每一个点的Gini, 找出变化最大的节点。当特征是连续特征时, 我们对连续值离散化, 取两点的平均值为分割节点。可以看到, 这里的排序算法需要花费大量的时间, 因为要遍历整个样本所有特征, 而且还要排序!!

Algorithm 1: Exact Greedy Algorithm for Split Finding

Input: I , instance set of current node

Input: d , feature dimension

$gain \leftarrow 0$

$G \leftarrow \sum_{i \in I} g_i, H \leftarrow \sum_{i \in I} h_i$

for $k = 1$ **to** m **do**

$G_L \leftarrow 0, H_L \leftarrow 0$

for j **in** $sorted(I, \text{by } \mathbf{x}_{jk})$ **do**

$G_L \leftarrow G_L + g_j, H_L \leftarrow H_L + h_j$

$G_R \leftarrow G - G_L, H_R \leftarrow H - H_L$

$score \leftarrow \max(score, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda})$

end

end

Output: Split with max score

知乎 @萌弟

论文的精确贪心算法的伪代码

因此在XGBoost里面我们使用的是近似算法 (Approximate Algorithm): 该算法首先根据特征分布的百分位数(percentiles)提出候选分裂点, 将连续特征映射到由这些候选点分割的桶中, 汇总统计信息并根据汇总的信息在提案中找到最佳解决方案。对于某个特征 k , 算法首先根据特征分布的分位数找到特征切割点的候选集合 $S_k = \{S_{k_1}, S_{k_2}, \dots, S_{k_l}\}$, 然后将特征 k 的值根据集合 S_k 划分到桶(bucket)中, 接着对每个桶内的样本统计值 G 、 H 进行累加, 最后在这些累计的统计量上寻找最佳分裂点。

Algorithm 2: Approximate Algorithm for Split Finding

```

for  $k = 1$  to  $m$  do
    | Propose  $S_k = \{s_{k1}, s_{k2}, \dots, s_{kl}\}$  by percentiles on feature  $k$ .
    | Proposal can be done per tree (global), or per split(local).
end
for  $k = 1$  to  $m$  do
    |  $G_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} \geq \mathbf{x}_{jk} > s_{k,v-1}\}} g_j$ 
    |  $H_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} \geq \mathbf{x}_{jk} > s_{k,v-1}\}} h_j$ 
end
Follow same step as in previous section to find max
score only among proposed splits.

```

知乎 @萌弟

论文的近似算法的伪代码

XGBoost动手实践：

1. 引入基本工具库：

```

1 # 引入基本工具库
2 import numpy as np
3 import pandas as pd
4 import xgboost as xgb
5 import matplotlib.pyplot as plt
6 plt.style.use("ggplot")
7 %matplotlib inline

```

2. XGBoost原生工具库的上手：

```

1 import xgboost as xgb # 引入工具库
2 # read in data
3 dtrain = xgb.DMatrix('demo/data/agaricus.txt.train') # XGBoost的专属数据格式,
4 dtest = xgb.DMatrix('demo/data/agaricus.txt.test') # XGBoost的专属数据格式,
5 # specify parameters via map
6 param = {'max_depth':2, 'eta':1, 'objective':'binary:logistic' } # 设置XGB的
7 num_round = 2 # 使用线程数
8 bst = xgb.train(param, dtrain, num_round) # 训练
9 # make prediction
10 preds = bst.predict(dtest) # 预测

```

3. XGBoost的参数设置(括号内的名称为sklearn接口对应的参数名字)

XGBoost的参数分为三种：

1. 通用参数

- `booster`:使用哪个弱学习器训练，默认`gbtree`，可选`gbtree`，`gblinear` 或`dart`
- `nthread`：用于运行XGBoost的并行线程数，默认为最大可用线程数
- `verbosity`：打印消息的详细程度。有效值为0（静默），1（警告），2（信息），3（调试）。
- Tree Booster的参数：
 - `eta (learning_rate)`：`learning_rate`，在更新中使用步长收缩以防止过度拟合，

通用参数有两种类型的`booster`，因为`tree`的性能比线性回归好得多，因此我们很少用线性回归。

2. 任务参数

- 差，排名的平均平均精度)，用户可以添加多个评估指标
- `rmse`，均方根误差；`rmsle`：均方根对数误差；`mae`：平均绝对误差；`mphe`：平均伪Huber错误；`logloss`：负对数似然；`error`：二进制分类错误率；
 - `merror`：多类分类错误率；`mlogloss`：多类`logloss`；`auc`：曲线下面积；`aucpr`：PR曲线下的面积；`ndcg`：归一化累计折扣；`map`：平均精度；
 - `seed`：随机数种子，[默认= 0]。

这个参数用来控制理想的优化目标和每一步结果的度量方法。

3. 命令行参数

这里不说了，因为很少用命令行控制台版本

4. XGBoost的调参说明：

参数调优的一般步骤：

- 1.确定（较大）学习速率和提升参数调优的初始值
- 2.`max_depth` 和 `min_child_weight` 参数调优
- 3.`gamma`参数调优
- 4.`subsample` 和 `colsample_bytree` 参数调优
- 5.正则化参数`alpha`调优
- 6.降低学习速率和使用更多的决策树

5. XGBoost详细攻略:

1). 安装XGBoost

方式1:

```
1 pip3 install xgboost
```

方式2:

```
1 pip install xgboost
```

2). 数据接口 (XGBoost可处理的数据格式DMatrix)

```
1 # 1.LibSVM文本格式文件
2 dtrain = xgb.DMatrix('train.svm.txt')
3 dtest = xgb.DMatrix('test.svm.buffer')
4 # 2.CSV文件(不能含类别文本变量, 如果存在文本变量请做特征处理如one-hot)
5 dtrain = xgb.DMatrix('train.csv?format=csv&label_column=0')
6 dtest = xgb.DMatrix('test.csv?format=csv&label_column=0')
7 # 3.NumPy 数组
8 data = np.random.rand(5, 10) # 5 entities, each contains 10 features
9 label = np.random.randint(2, size=5) # binary target
10 dtrain = xgb.DMatrix(data, label=label)
11 # 4.scipy.sparse 数组
12 csr = scipy.sparse.csr_matrix((data, (row, col)))
13 dtrain = xgb.DMatrix(csr)
14 # pandas数据框dataframe
15 data = pandas.DataFrame(np.arange(12).reshape((4,3)), columns=['a', 'b', 'c'])
16 label = pandas.DataFrame(np.random.randint(2, size=4))
17 dtrain = xgb.DMatrix(data, label=label)
```

笔者推荐: 先保存到XGBoost二进制文件中将使加载速度更快, 然后再加载进来

```
1 # 1.保存DMatrix到XGBoost二进制文件中
2 dtrain = xgb.DMatrix('train.svm.txt')
3 dtrain.save_binary('train.buffer')
4 # 2. 缺少的值可以用DMatrix构造函数中的默认值替换:
5 dtrain = xgb.DMatrix(data, label=label, missing=-999.0)
6 # 3. 可以在需要时设置权重:
7 w = np.random.rand(5, 1)
```

```
8 dtrain = xgb.DMatrix(data, label=label, missing=-999.0, weight=w)
```

3). 参数的设置方式:

```
1 # 加载并处理数据
2 df_wine = pd.read_csv('https://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data')
3 df_wine.columns = ['Class label', 'Alcohol', 'Malic acid', 'Ash', 'Alcalinity of ash',
4                   'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins', 'Color intensity',
5                   'Hue', 'Hue (modulated)', 'Hue (averaged)', 'Hue (blue-green)', 'Hue (blue-violet)',
6                   'Hue (red-green)', 'Hue (red-violet)', 'Hue (red-blue)', 'Hue (blue-red)', 'Hue (blue-green)',
7                   'Hue (blue-violet)', 'Hue (red-green)', 'Hue (red-violet)', 'Hue (red-blue)', 'Hue (blue-red)']
8 df_wine = df_wine[df_wine['Class label'] != 1] # drop 1 class
9 y = df_wine['Class label'].values
10 X = df_wine[['Alcohol', 'OD280/OD315 of diluted wines']].values
11 from sklearn.model_selection import train_test_split # 切分训练集与测试集
12 from sklearn.preprocessing import LabelEncoder # 标签化分类变量
13 le = LabelEncoder()
14 y = le.fit_transform(y)
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
16 dtrain = xgb.DMatrix(X_train, label=y_train)
17 dtest = xgb.DMatrix(X_test)
18 # 1.Booster 参数
19 params = {
20     'booster': 'gbtree',
21     'objective': 'multi:softmax', # 多分类的问题
22     'num_class': 10, # 类别数, 与 multisoftmax 并用
23     'gamma': 0.1, # 用于控制是否后剪枝的参数, 越大越保守, 一般0.1、
24     'max_depth': 12, # 构建树的深度, 越大越容易过拟合
25     '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, 模型越简单
26     'subsample': 0.7, # 随机采样训练样本
27     'colsample_bytree': 0.7, # 生成树时进行的列采样
28     'min_child_weight': 3,
29     'silent': 1, # 设置成1则没有运行信息输出, 最好是设置为0.
30     'eta': 0.007, # 如同学习率
31     'seed': 1000,
32     'nthread': 4, # cpu 线程数
33     'eval_metric': 'auc'
34 }
35 plst = params.items()
36 # evallist = [(dtest, 'eval'), (dtrain, 'train')] # 指定验证集
```

4). 训练

```
1 # 2. 训练
2 num_round = 10
3 bst = xgb.train( plst, dtrain, num_round)
4 #bst = xgb.train( plst, dtrain, num_round, evallist )
```

5). 保存模型

```
1 # 3. 保存模型
2 bst.save_model('0001.model')
3 # dump model
4 bst.dump_model('dump.raw.txt')
5 # dump model with feature map
6 #bst.dump_model('dump.raw.txt', 'featmap.txt')
```

6). 加载保存的模型

```
1 # 4. 加载保存的模型:
2 bst = xgb.Booster({'nthread': 4}) # init model
3 bst.load_model('0001.model') # Load data
```

7). 设置早停机制

```
1 # 5. 也可以设置早停机制（需要设置验证集）
2 train(..., evals=evals, early_stopping_rounds=10)
```

8). 预测

```
1 # 6. 预测
2 ypred = bst.predict(dtest)
```

9). 绘图

```
1 # 1. 绘制重要性
2 xgb.plot_importance(bst)
3 # 2. 绘制输出树
4 #xgb.plot_tree(bst, num_trees=2)
5 # 3. 使用xgboost.to_graphviz()将目标树转换为graphviz
6 #xgb.to_graphviz(bst, num_trees=2)
```

6. 实战案例:

1). 分类案例

```
1 from sklearn.datasets import load_iris
2 import xgboost as xgb
3 from xgboost import plot_importance
4 from matplotlib import pyplot as plt
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score # 准确率
7 # 加载样本数据集
8 iris = load_iris()
9 X,y = iris.data,iris.target
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
11
12 # 算法参数
13 params = {
14     'booster': 'gbtree',
15     'objective': 'multi:softmax',
16     'num_class': 3,
17     'gamma': 0.1,
18     'max_depth': 6,
19     'lambda': 2,
20     'subsample': 0.7,
21     'colsample_bytree': 0.75,
22     'min_child_weight': 3,
23     'silent': 0,
24     'eta': 0.1,
25     'seed': 1,
26     'nthread': 4,
27 }
28
29 plst = params.items()
30
31 dtrain = xgb.DMatrix(X_train, y_train) # 生成数据集格式
32 num_rounds = 500
33 model = xgb.train(plst, dtrain, num_rounds) # xgboost模型训练
34
35 # 对测试集进行预测
36 dtest = xgb.DMatrix(X_test)
37 y_pred = model.predict(dtest)
```



```
38
39 # 计算准确率
40 accuracy = accuracy_score(y_test,y_pred)
41 print("accuracy: %.2f%%" % (accuracy*100.0))
42
43 # 显示重要特征
44 plot_importance(model)
45 plt.show()
```

2). 回归案例

```
1 import xgboost as xgb
2 from xgboost import plot_importance
3 from matplotlib import pyplot as plt
4 from sklearn.model_selection import train_test_split
5 from sklearn.datasets import load_boston
6 from sklearn.metrics import mean_squared_error
7
8 # 加载数据集
9 boston = load_boston()
10 X,y = boston.data,boston.target
11
12 # XGBoost 训练过程
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
14
15 params = {
16     'booster': 'gbtree',
17     'objective': 'reg:squarederror',
18     'gamma': 0.1,
19     'max_depth': 5,
20     'lambda': 3,
21     'subsample': 0.7,
22     'colsample_bytree': 0.7,
23     'min_child_weight': 3,
24     'silent': 1,
25     'eta': 0.1,
26     'seed': 1000,
27     'nthread': 4,
28 }
```

```
29
30 dtrain = xgb.DMatrix(X_train, y_train)
31 num_rounds = 300
32 plst = params.items()
33 model = xgb.train(plst, dtrain, num_rounds)
34
35 # 对测试集进行预测
36 dtest = xgb.DMatrix(X_test)
37 ans = model.predict(dtest)
38
39 # 显示重要特征
40 plot_importance(model)
41 plt.show()
```

7. XGBoost调参 (结合sklearn网格搜索)

```
1 import xgboost as xgb
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.model_selection import GridSearchCV
5 from sklearn.metrics import roc_auc_score
6
7 iris = load_iris()
8 X,y = iris.data,iris.target
9 col = iris.target_names
10 train_x, valid_x, train_y, valid_y = train_test_split(X, y, test_size=0.3, ran
11 parameters = {
12     'max_depth': [5, 10, 15, 20, 25],
13     'learning_rate': [0.01, 0.02, 0.05, 0.1, 0.15],
14     'n_estimators': [500, 1000, 2000, 3000, 5000],
15     'min_child_weight': [0, 2, 5, 10, 20],
16     'max_delta_step': [0, 0.2, 0.6, 1, 2],
17     'subsample': [0.6, 0.7, 0.8, 0.85, 0.95],
18     'colsample_bytree': [0.5, 0.6, 0.7, 0.8, 0.9],
19     'reg_alpha': [0, 0.25, 0.5, 0.75, 1],
20     'reg_lambda': [0.2, 0.4, 0.6, 0.8, 1],
21     'scale_pos_weight': [0.2, 0.4, 0.6, 0.8, 1]
22
```

```
23 }
24
25 xlf = xgb.XGBClassifier(max_depth=10,
26                         learning_rate=0.01,
27                         n_estimators=2000,
28                         silent=True,
29                         objective='multi:softmax',
30                         num_class=3,
31                         nthread=-1,
32                         gamma=0,
33                         min_child_weight=1,
34                         max_delta_step=0,
35                         subsample=0.85,
36                         colsample_bytree=0.7,
37                         colsample_bylevel=1,
38                         reg_alpha=0,
39                         reg_lambda=1,
40                         scale_pos_weight=1,
41                         seed=0,
42                         missing=None)
43
44 gs = GridSearchCV(xlf, param_grid=parameters, scoring='accuracy', cv=3)
45 gs.fit(train_x, train_y)
46
47 print("Best score: %0.3f" % gs.best_score_)
48 print("Best parameters set: %s" % gs.best_params_ )
```

本文电子版 后台回复 **XGBoost** 获取

Datawhale

和学习者一起成长

一个专注于AI的开源组织，让学习不再孤独



长按扫码关注

🔗 点击[阅读原文](#)，关注作者