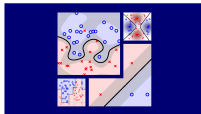


Machine Learning Techniques (機器學習技法)



Lecture 11: Gradient Boosted Decision Tree

Hsuan-Tien Lin (林軒田)

`htlin@csie.ntu.edu.tw`

Department of Computer Science
& Information Engineering

National Taiwan University
(國立台灣大學資訊工程系)



Roadmap

- 1 Embedding Numerous Features: Kernel Models
- 2 Combining Predictive Features: Aggregation Models

Lecture 10: Random Forest

bagging of **randomized C&RT** trees with
automatic validation and **feature selection**

Lecture 11: Gradient Boosted Decision Tree

- Adaptive Boosted Decision Tree
- Optimization View of AdaBoost
- Gradient Boosting
- Summary of Aggregation Models

- 3 Distilling Implicit Features: Extraction Models

From Random Forest to AdaBoost-DTree

需要让决策树处理带有不同权重的样本。
得到加权错误率最小的 g_t

function **RandomForest**(\mathcal{D})

For $t = 1, 2, \dots, T$

- ① request size- N' data $\tilde{\mathcal{D}}_t$ by bootstrapping with \mathcal{D}
- ② obtain tree g_t by Randomized-DTree($\tilde{\mathcal{D}}_t$)

return $G = \text{Uniform}(\{g_t\})$

function **AdaBoost-DTree**(\mathcal{D})

For $t = 1, 2, \dots, T$

- ① reweight data by $\mathbf{u}^{(t)}$
- ② obtain tree g_t by DTree($\mathcal{D}, \mathbf{u}^{(t)}$)
- ③ calculate 'vote' α_t of g_t

return $G = \text{LinearHypo}(\{(g_t, \alpha_t)\})$

need: **weighted** DTree($\mathcal{D}, \mathbf{u}^{(t)}$)

随机森林:

把决策树bagging起来。

不同 g_t 不同的随机化方式: bootstrap sample/随机特征子空间

adaboost提升树:

1每个 g_t 也是决策树

2每个 g_t 接收的数据带有不同的权重 u_t .根据权重选择最优的 g_t 并根据 g_t 结果,得到下一次 u_{t+1} 的权重,用来算 g_{t+1}

3算 g_t 同时得到 g_t 的权重 α_t

4最后按权重合并所有 g_t 的预测结果,得到 G

Weighted Decision Tree Algorithm

Weighted Algorithm

最小化
加权错误率

$$\text{minimize (regularized)} \quad E_{\text{in}}^{\mathbf{u}}(h) = \frac{1}{N} \sum_{n=1}^N u_n \cdot \text{err}(y_n, h(\mathbf{x}_n))$$

本来应该是在看算法计算 E_{in} 时，哪里能用到权重，做相应改变。但决策树改算法很麻烦。

因此根据 u_n 抽样，来得到对应这些权重的样本。

if using existing algorithm as **black box** (no modifications),
to get $E_{\text{in}}^{\mathbf{u}}$ approximately optimized.....

'Weighted' Algorithm in Bagging

weights \mathbf{u} expressed by
bootstrap-sampled copies
—request size- N' data $\tilde{\mathcal{D}}_t$
by bootstrapping with \mathcal{D}

A General Randomized Base Algorithm

weights \mathbf{u} expressed by
sampling proportional to u_n
—request size- N' data $\tilde{\mathcal{D}}_t$
by sampling $\propto \mathbf{u}$ on \mathcal{D}

bagging里，
样本被重复抽到 $\tilde{\mathcal{D}}_t$ 以后，共有几份，就可以看作该样本权重是多少

adaboost-Tree: 进入 g_t 前，先通过权重 u_n 抽样。得到的数据 $\tilde{\mathcal{D}}_t$ 就和 u_n 希望的权重一样。

AdaBoost-DTree: often via
AdaBoost + **sampling** $\propto \mathbf{u}^{(t)}$ + DTree($\tilde{\mathcal{D}}_t$)
without modifying DTree

之后直接用决策树拟合加权的的数据 $\tilde{\mathcal{D}}_t$ ~
算error + g_t 权重at
其中error用 \mathcal{D} , u_n 算
而不是用 $\tilde{\mathcal{D}}_t$ 算。
否则容易是0

Weak Decision Tree Algorithm

AdaBoost: **votes** $\alpha_t = \ln \diamond_t = \ln \sqrt{\frac{1-\epsilon_t}{\epsilon_t}}$ with **weighted error rate** ϵ_t

if **fully grown** tree trained on **all** \mathbf{x}_n

$\Rightarrow E_{in}(g_t) = 0$ if **all** \mathbf{x}_n different

$\Rightarrow E_{in}^u(g_t) = 0$

$\Rightarrow \epsilon_t = 0$

$\Rightarrow \alpha_t = \infty$ (**autocracy!!**)

如果一棵树 g_t 能完全拟合输入数据:

g_t 错误率是0。该模型 α_t 会是无穷大。只有这个树发挥作用, 别的 g_t 没有用了。

希望这个树犯一点错, $E_{in}(g_t) \neq 0$ 。不要完全拟合。

need: **pruned** tree trained on **some** \mathbf{x}_n to be **weak**

- **pruned**: usual pruning, or just **limiting tree height** 剪枝/限制树高
- **some**: **sampling** $\propto \mathbf{u}^{(t)}$ 用 \mathbf{u} 抽样, 总有抽不到的数据

使 g_t 不能在全部数据上完全拟合

但算error还是在整体 D 上用 \mathbf{u} 算

AdaBoost-DTree: often via AdaBoost +
sampling $\propto \mathbf{u}^{(t)}$ + **pruned** DTree(\tilde{D})
 符合 \mathbf{u} 分布 非完全的剪枝

AdaBoost with Extremely-Pruned Tree

what if DTree with **height** ≤ 1 (extremely pruned)?

DTree (C&RT) with **height** ≤ 1

learn **branching criteria** 最弱的树：树高只有一层。只切一次

$$b(\mathbf{x}) = \underset{\text{decision stumps } h(\mathbf{x})}{\operatorname{argmin}} \sum_{c=1}^2 |\mathcal{D}_c \text{ with } h| \cdot \text{impurity}(\mathcal{D}_c \text{ with } h)$$

—if **impurity** = **binary classification error**,

just a decision stump, remember? :-)

如果不纯净度用分类误差算，每个树就是决策桩。这时候树比较简单，再用adaboost时，un可以直接算误差时算进去，不用在进入gt前抽样得Dt~

AdaBoost-Stump
= **special case** of AdaBoost-DTree

Fun Time

When running AdaBoost-DTree with sampling and getting a decision tree g_t such that g_t achieves zero error on the sampled data set \tilde{D}_t . Which of the following is possible?

- 1 $\alpha_t < 0$
- 2 $\alpha_t = 0$
- 3 $\alpha_t > 0$
- 4 all of the above

Fun Time

When running AdaBoost-DTree with sampling and getting a decision tree g_t such that g_t achieves zero error on the sampled data set \tilde{D}_t . Which of the following is possible?

- ① $\alpha_t < 0$
- ② $\alpha_t = 0$
- ③ $\alpha_t > 0$
- ④ all of the above

Reference Answer: ④

While g_t achieves zero error on \tilde{D}_t , g_t may not achieve zero weighted error on $(\mathcal{D}, \mathbf{u}^{(t)})$ and hence ϵ_t can be anything, even $\geq \frac{1}{2}$. Then, α_t can be ≤ 0 .

Example Weights of AdaBoost

$$\begin{aligned}
 u_n^{(t+1)} &= \begin{cases} u_n^{(t)} \cdot \blacklozenge_t & \text{if incorrect} \\ u_n^{(t)} / \blacklozenge_t & \text{if correct} \end{cases} && \text{adaboost: 根据gt预测结果更新ut。使} \\
 & && \text{ut+1在gt上表现很差, gt+1,gt diverse} \\
 &= u_n^{(t)} \cdot \blacklozenge_t^{-y_n g_t(\mathbf{x}_n)} = u_n^{(t)} \cdot \text{exp}(-y_n \alpha_t g_t(\mathbf{x}_n)) \\
 & \text{简化写法 } \alpha_t = \ln(t) && \alpha_t: \text{gt的权重}
 \end{aligned}$$

$$u_n^{(T+1)} = u_n^{(1)} \cdot \prod_{t=1}^T \text{exp}(-y_n \alpha_t g_t(\mathbf{x}_n)) = \frac{1}{N} \cdot \text{exp}\left(-y_n \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n)\right)$$

- recall: $G(\mathbf{x}) = \text{sign}\left(\sum_{t=1}^T \alpha_t g_t(\mathbf{x})\right)$

y_n 是固定的

可以看作是T个模型对 \mathbf{x}_n 的加权投票分数

- $\sum_{t=1}^T \alpha_t g_t(\mathbf{x})$: **voting score** of $\{g_t\}$ on \mathbf{x}

adaboost每轮的新权重和加权score有关

$$\text{AdaBoost: } u_n^{(T+1)} \propto \text{exp}(-y_n (\text{voting score on } \mathbf{x}_n))$$

Voting Score and Margin

linear blending = **LinModel** + hypotheses as transform + ~~constraints~~

$$G(\mathbf{x}_n) = \text{sign} \left(\overbrace{\sum_{t=1}^T \underbrace{\alpha_t}_{w_i} \underbrace{g_t(\mathbf{x}_n)}_{\phi_i(\mathbf{x}_n)}}^{\text{voting score}} \right)$$

voting score各项
相当于linear model里的 $w_i, \phi_i(\mathbf{x})$
voting score可以看作是一种没有正规化的距离（类比SVM）

and hard-margin SVM **margin** = $\frac{y_n \cdot (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$, **remember? :-)**

$y_n(\text{voting score})$ = signed & unnormalized **margin**

want $y_n(\text{voting score})$ **positive & large**

$\Leftrightarrow \exp(-y_n(\text{voting score}))$ **small** 预测结果如果正确, $y_n * \text{voting_score}$ 会很大
 $\Leftrightarrow u_n^{(T+1)}$ **small** $y_n = 1/-1$, 得分相差很大。
 因此希望margin : $y_n * \text{voting_score}$ 较大

对应到adaboost,最后是希望权重的值 u_n 较小(都学会了)

claim: AdaBoost **decreases** $\sum_{n=1}^N u_n^{(t)}$ 所有点的总权重

AdaBoost Error Function

claim: AdaBoost **decreases** $\sum_{n=1}^N u_n^{(t)}$ and thus somewhat **minimizes**

$$\sum_{n=1}^N u_n^{(T+1)} = \frac{1}{N} \sum_{n=1}^N \exp \left(-y_n \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n) \right)$$

voting score: 不同类别的margin变大

linear score $s = \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n)$

- $\text{err}_{0/1}(s, y) = \mathbb{I}[ys \leq 0]$
- $\widehat{\text{err}}_{\text{ADA}}(s, y) = \exp(-ys)$:
upper bound of $\text{err}_{0/1}$
—called **exponential error measure**

$\widehat{\text{err}}_{\text{ADA}}$: **algorithmic error measure**
by **convex upper bound** of $\text{err}_{0/1}$

AdaBoost Error Function

claim: AdaBoost **decreases** $\sum_{n=1}^N u_n^{(t)}$ and thus somewhat **minimizes**

$$\sum_{n=1}^N u_n^{(T+1)} = \frac{1}{N} \sum_{n=1}^N \exp \left(-y_n \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n) \right)$$

linear score $s = \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n)$

- $\text{err}_{0/1}(s, y) = \mathbb{I}[ys \leq 0]$
- $\widehat{\text{err}}_{\text{ADA}}(s, y) = \exp(-ys)$:
upper bound of $\text{err}_{0/1}$
—called **exponential error measure**

$\widehat{\text{err}}_{\text{ADA}}$: **algorithmic error measure**
by **convex upper bound** of $\text{err}_{0/1}$

AdaBoost Error Function

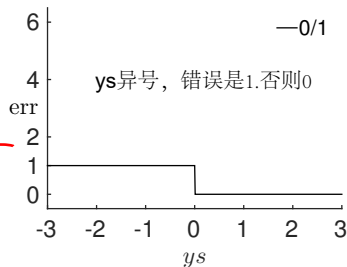
claim: AdaBoost **decreases** $\sum_{n=1}^N u_n^{(t)}$ and thus somewhat **minimizes**

$$\sum_{n=1}^N u_n^{(T+1)} = \frac{1}{N} \sum_{n=1}^N \exp \left(-y_n \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n) \right)$$

voting score: 不同类别的margin变大

linear score $s = \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n)$

- $\text{err}_{0/1}(s, y) = \mathbb{I}[ys \leq 0]$
- $\widehat{\text{err}}_{\text{ADA}}(s, y) = \exp(-ys)$:
upper bound of $\text{err}_{0/1}$
—called **exponential error measure**



$\widehat{\text{err}}_{\text{ADA}}$: **algorithmic error measure**
by **convex upper bound** of $\text{err}_{0/1}$

AdaBoost Error Function

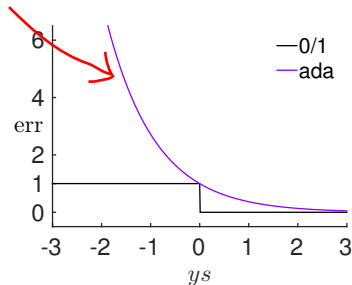
claim: AdaBoost decreases $\sum_{n=1}^N u_n^{(t)}$ and thus somewhat **minimizes**

$$\sum_{n=1}^N u_n^{(T+1)} = \frac{1}{N} \sum_{n=1}^N \exp \left(-y_n \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n) \right)$$

adaboost实际降低的是这个分类错误函数

linear score $s = \sum_{t=1}^T \alpha_t g_t(\mathbf{x}_n)$

- $\text{err}_{0/1}(s, y) = \mathbb{I}[ys \leq 0]$
- $\widehat{\text{err}}_{\text{ADA}}(s, y) = \exp(-ys)$:
upper bound of $\text{err}_{0/1}$ 是0/1错误的上限
—called **exponential error measure**



所以降低权重之和，相当于降低分类错误的凸上限。

$\widehat{\text{err}}_{\text{ADA}}$: **algorithmic error measure**
by **convex upper bound** of $\text{err}_{0/1}$

Gradient Descent on AdaBoost Error Function

recall: gradient descent (**remember? :-)**), at iteration t

$$\min_{\|\mathbf{v}\|=1} E_{\text{in}}(\mathbf{w}_t + \eta \mathbf{v}) \approx \underbrace{E_{\text{in}}(\mathbf{w}_t)}_{\text{known}} + \underbrace{\eta}_{\text{步长}} \underbrace{\mathbf{v}^T}_{\text{单位向量}} \underbrace{\nabla E_{\text{in}}(\mathbf{w}_t)}_{\text{known}}$$

复习一下梯度下降

希望选一个更新方向 \mathbf{v} , 使 $f(\mathbf{w}_t + 1) = f(\mathbf{w}_t + \eta \mathbf{v})$ 最小。 $\mathbf{w}_t + 1 = \mathbf{w}_t + \eta \mathbf{v}$

按照泰勒展开 单位向量 \mathbf{v} 和梯度反向(夹角180度)结果最小。因此 \mathbf{v} 选择负梯度方向

更新 \mathbf{w}_t : $\mathbf{w}_t + 1 = \mathbf{w}_t + \eta \mathbf{v}$

\mathbf{w}_t 处的梯度方向

以前是找一个使表现最好的向量 \mathbf{w} 。通过更新向量, 使表现更好。现在是找一个使表现最好的 \mathbf{g}_t 。通过更新 \mathbf{g}_t , 使整体表现更好

向量是给定 index, 输出对应值。函数 \mathbf{g}_t 输入是实数

at iteration t , to find \mathbf{g}_t , solve

$$\begin{aligned} \min_h \hat{E}_{\text{ADA}} &= \frac{1}{N} \sum_{n=1}^N \exp \left(-y_n \left(\sum_{\tau=1}^{t-1} \alpha_\tau g_\tau(\mathbf{x}_n) + \eta h(\mathbf{x}_n) \right) \right) \\ &= \sum_{n=1}^N u_n^{(t)} \exp(-y_n \eta h(\mathbf{x}_n)) \\ &\stackrel{\text{taylor}}{\approx} \sum_{n=1}^N \underbrace{u_n^{(t)}}_{\text{固定}} (1 - y_n \eta h(\mathbf{x}_n)) = \sum_{n=1}^N u_n^{(t)} - \eta \sum_{n=1}^N u_n^{(t)} y_n h(\mathbf{x}_n) \end{aligned}$$

当前的整体的函数

我们希望通过加一个新的 \mathbf{g}_{t+1} , 更新 \mathbf{G} , 使 \hat{E}_{Ada} 向减小的方向移动

对 $\exp(x)$ 在原点处展开

adaboost 相当于通过选择合适的 $n, h(\mathbf{x})$, 最小化这部分

选择合适的 h
最小化 error

good h : minimize $\sum_{n=1}^N u_n^{(t)} (-y_n h(\mathbf{x}_n))$

Learning Hypothesis as Optimization

finding good h (function direction) \Leftrightarrow minimize $\sum_{n=1}^N u_n^{(t)} (-y_n h(\mathbf{x}_n))$

for binary classification, where y_n and $h(\mathbf{x}_n)$ both $\in \{-1, +1\}$:

$$\sum_{n=1}^N u_n^{(t)} (-y_n h(\mathbf{x}_n)) = \sum_{n=1}^N u_n^{(t)} \begin{cases} -1 & \text{if } y_n = h(\mathbf{x}_n) \\ +1 & \text{if } y_n \neq h(\mathbf{x}_n) \end{cases}$$

抽出一个-un

$$= -\sum_{n=1}^N u_n^{(t)} + \sum_{n=1}^N u_n^{(t)} \begin{cases} 0 & \text{if } y_n = h(\mathbf{x}_n) \\ 2 & \text{if } y_n \neq h(\mathbf{x}_n) \end{cases}$$

如果是二分类
相当于选一个 h ,
最小化 h 在 N 个训练数据上的加权错误率

而adaboost里

gt就是使得 u_t 加权下的 E_{in} 最小的算法

ht=gt

—who minimizes $E_{in}^{u^{(t)}}(h)$? **A in AdaBoost! :-)**

$$= -\sum_{n=1}^N u_n^{(t)} + 2E_{in}^{u^{(t)}}(h) \cdot N$$

固定的 $h(x)$ 上的错误率

A: good $g_t = h$ for 'gradient descent'

adaboost选的是一个好的, G 的更新方向。使 G 朝着 $-h$ 的方向更新

Deciding Blending Weight as Optimization

AdaBoost finds g_t by approximately $\min_h \hat{E}_{\text{ADA}} = \sum_{n=1}^N u_n^{(t)} \exp(-y_n \eta h(\mathbf{x}_n))$
 adaboost每次找一个最好的更新方向h, 更新G

after finding g_t , how about $\min_{\eta} \hat{E}_{\text{ADA}} = \sum_{n=1}^N u_n^{(t)} \exp(-y_n \eta g_t(\mathbf{x}_n))$
 如何选步长呢? 在最优方向基础上, 选一个最优的步长 η

在当前 g_t 情况下, 选最优的方向和步长, 贪婪的更新G (尽管未必全局最优)

- optimal η_t somewhat ‘**greedily faster**’ than fixed (small) η —called **steepest** decent for optimization
- two cases inside summation: 考虑2分类, 可以把Eada错误写成:

- $y_n = g_t(\mathbf{x}_n) : u_n^{(t)} \exp(-\eta)$ (correct)

- $y_n \neq g_t(\mathbf{x}_n) : u_n^{(t)} \exp(+\eta)$ (incorrect)

加权错误率/归一化

- $\hat{E}_{\text{ADA}} = \left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left((1 - \epsilon_t) \exp(-\eta) + \epsilon_t \exp(+\eta) \right)$

最优化 η

by solving $\frac{\partial \hat{E}_{\text{ADA}}}{\partial \eta} = 0$, **steepest** $\eta_t = \ln \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} = \alpha_t$, **remember? :-)**
 —AdaBoost: **steepest** decent with **approximate functional gradient**

Fun Time

With $\hat{E}_{\text{ADA}} = \left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left((1 - \epsilon_t) \exp(-\eta) + \epsilon_t \exp(+\eta) \right)$, which of the following is $\frac{\partial \hat{E}_{\text{ADA}}}{\partial \eta}$ that can be used for solving the optimal η_t ?

- 1 $\left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left(+ (1 - \epsilon_t) \exp(-\eta) + \epsilon_t \exp(+\eta) \right)$
- 2 $\left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left(+ (1 - \epsilon_t) \exp(-\eta) - \epsilon_t \exp(+\eta) \right)$
- 3 $\left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left(- (1 - \epsilon_t) \exp(-\eta) + \epsilon_t \exp(+\eta) \right)$
- 4 $\left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left(- (1 - \epsilon_t) \exp(-\eta) - \epsilon_t \exp(+\eta) \right)$

Fun Time

With $\hat{E}_{\text{ADA}} = \left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left((1 - \epsilon_t) \exp(-\eta) + \epsilon_t \exp(+\eta) \right)$, which of the following is $\frac{\partial \hat{E}_{\text{ADA}}}{\partial \eta}$ that can be used for solving the optimal η_t ?

- ① $\left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left(+ (1 - \epsilon_t) \exp(-\eta) + \epsilon_t \exp(+\eta) \right)$
- ② $\left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left(+ (1 - \epsilon_t) \exp(-\eta) - \epsilon_t \exp(+\eta) \right)$
- ③ $\left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left(- (1 - \epsilon_t) \exp(-\eta) + \epsilon_t \exp(+\eta) \right)$
- ④ $\left(\sum_{n=1}^N u_n^{(t)} \right) \cdot \left(- (1 - \epsilon_t) \exp(-\eta) - \epsilon_t \exp(+\eta) \right)$

Reference Answer: ③

Differentiate $\exp(-\eta)$ and $\exp(+\eta)$ with respect to η and you should easily get the result.

Gradient Boosting for Arbitrary Error Function

AdaBoost

$$\min_{\eta} \min_h \frac{1}{N} \sum_{n=1}^N \exp \left(-y_n \left(\underbrace{\sum_{\tau=1}^{t-1} \alpha_{\tau} g_{\tau}(\mathbf{x}_n)}_s + \eta h(\mathbf{x}_n) \right) \right)$$

s是vote function

with **binary-output hypothesis** h

adaboost相当于最小化二分类误差上限 $\exp(-y_n * s)$

对于固定步长 n , 选一个最佳更新方向 g_t ,

对于固定更新方向, 选一个最佳更新步长 n

组合起来更新 G , 使 $G(x)$ 对应的error function降低

GradientBoost

同样
先最优方向 $h(x)$
再优化步长 n
作为我新的 g_{t+1}
at+1

$$\min_{\eta} \min_h \frac{1}{N} \sum_{n=1}^N \text{err} \left(\underbrace{\sum_{\tau=1}^{t-1} \alpha_{\tau} g_{\tau}(\mathbf{x}_n)}_{\text{已有的voting score}} + \eta h(\mathbf{x}_n), y_n \right)$$

沿着使error function下降最快的步长 n 和方向 $h(x)$ 更新 G

换成任意感兴趣的error function
with **any hypothesis** h (usually **real-output hypothesis**)

GradientBoost: allows **extension to different err** for regression/soft classification/etc.

GradientBoost for Regression

误差是回归误差。G(x)输出是实数

$$\min_{\eta} \min_h \frac{1}{N} \sum_{n=1}^N \text{err} \left(\underbrace{\sum_{\tau=1}^{t-1} \alpha_{\tau} g_{\tau}(\mathbf{x}_n)}_{s_n} + \eta h(\mathbf{x}_n), y_n \right) \quad \text{with } \text{err}(s, y) = (s - y)^2$$

先考虑回归：平方误差函数

已有的voting score
可能是实数

$$\begin{aligned} \min_h \dots &\approx \min_h \frac{1}{N} \sum_{n=1}^N \underbrace{\text{err}(s_n, y_n)}_{\text{constant}} + \frac{1}{N} \sum_{n=1}^N \eta h(\mathbf{x}_n) \quad \left. \frac{\partial \text{err}(s, y_n)}{\partial s} \right|_{s=s_n} \\ &= \min_h \text{constants} + \frac{\eta}{N} \sum_{n=1}^N h(\mathbf{x}_n) \cdot 2(s_n - y_n) \end{aligned}$$

error function在 s_n 处一阶泰勒展开。对每个点 (\mathbf{x}_n, y_n) 都这么做，最后加起来求平均。

原来的错误率，常数

会展出一个一阶梯度的项是error函数在 s_n 点处的梯度

最小化这部分。让 $h(\mathbf{x})$ 与 $(s_n - y_n)$ 反向

naïve solution $h(\mathbf{x}_n) = -\infty \cdot (s_n - y_n)$ if no constraint on h

只需要考虑 h 的方向。大小由 n 控制。但限制 $|h|=1$ 不好优化。用正则项制约下

Learning Hypothesis as Optimization

$$\min_h \quad \text{constants} + \frac{\eta}{N} \sum_{n=1}^N 2h(\mathbf{x}_n)(s_n - y_n)$$

- **magnitude** of h does not matter: because η will be optimized next
- **penalize large magnitude** to avoid naïve solution

$$\begin{aligned} \min_h \quad & \text{constants} + \frac{\eta}{N} \sum_{n=1}^N (2h(\mathbf{x}_n)(s_n - y_n) + (h(\mathbf{x}_n))^2) \\ & \text{防止} h(\mathbf{x}_n) \text{的绝对值太大。} \\ & \text{我们只在意} h \text{方向。大小主要由} \eta \text{控制} \\ = \quad & \text{constants} + \frac{\eta}{N} \sum_{n=1}^N \left(\text{constant} + \underbrace{(h(\mathbf{x}_n) - (y_n - s_n))^2}_{(\text{sn}-\text{yn})^{**2} \text{ 常数}} \right) \\ & \text{可以凑出这样一个形式} \end{aligned}$$

- solution of **penalized approximate functional gradient**:
squared-error regression on $\{(\mathbf{x}_n, \underbrace{y_n - s_n}_{\text{residual}})\}$

希望学到的gt, 平方误差上逼近 $y_n - s_n$
 给gt的数据是 $(\mathbf{x}_n, \text{residual})$
 希望拟合的数据是残差,
 是 $y_n - s_n$ (\mathbf{x}_n 对应的当前的 voting score)
- 这是一个回归问题。
 如果用决策树作为gt, 选择regression error作为不纯净度
 gt最终的训练数据不是 (\mathbf{x}_n, y_n) , 而是 $(\mathbf{x}_n, y_n - s_n)$

GradientBoost for regression:

find $g_t = h$ by **regression** with **residuals**

Deciding Blending Weight as Optimization

after finding $g_t = h$,

找到好的 g_t/h 后，需要找到使error最小的最优步长 η 。即新 g_t 的权重 η

$$\min_{\eta} \min_h \frac{1}{N} \sum_{n=1}^N \text{err} \left(\underbrace{\sum_{\tau=1}^{t-1} \alpha_{\tau} g_{\tau}(\mathbf{x}_n)}_{s_n} + \eta g_t(\mathbf{x}_n), y_n \right) \text{ with } \text{err}(s, y) = (s - y)^2$$

重写一下

$$\min_{\eta} \frac{1}{N} \sum_{n=1}^N (s_n + \eta g_t(\mathbf{x}_n) - y_n)^2 = \frac{1}{N} \sum_{n=1}^N (\underbrace{(y_n - s_n)}_{\text{还是 } \mathbf{x}_n \text{ 在这点的残差}} - \underbrace{\eta g_t(\mathbf{x}_n)}_{\text{固定的实数}})^2$$

—one-variable linear regression on $\{(g_t\text{-transformed input}, \text{residual})\}$

\mathbf{X} : N 个样本，每个样本只有一维特征 $g_t(\mathbf{x}_n)$

只需要一个变量 w_i

$\text{error} = \text{mean} (y_n - \sum_i (w_i * x_{n,i}))^2$

\mathbf{Y} : N 个样本，每个样本的输出是 $(y_n - s_n)$

相当于单维度的
线性回归问题。
求出单变量 w_i
即 g_t 的权重

GradientBoost for regression: $\alpha_t = \text{optimal } \eta$
by $g_t\text{-transformed linear regression}$

Putting Everything Together

Gradient Boosted Decision Tree (GBDT)

$s_1 = s_2 = \dots = s_N = 0$ $T=0$ 所有点 $s_n=0$ 。因为还没有voting score.

for $t = 1, 2, \dots, T$

① obtain g_t by $\mathcal{A}(\{(\mathbf{x}_n, y_n - s_n)\})$ where \mathcal{A} is a (squared-error) regression algorithm
 1 先解一个regression问题: $x_n \rightarrow y_n - s_n$ 得到最优的 g_t
 —**how about sampled and pruned C&RT?** 这里用决策树cart做regression

② compute $\alpha_t = \text{OneVarLinearRegression}(\{(g_t(\mathbf{x}_n), y_n - s_n)\})$

③ update $s_n \leftarrow s_n + \alpha_t g_t(\mathbf{x}_n)$ 2 之后解一个一维的线性regression问题 $g_t \rightarrow y_n - s_n$ 得到 α_t
 3 根据 g_t, α_t , 更新所有样本的voting score s_n

return $G(\mathbf{x}) = \sum_{t=1}^T \alpha_t g_t(\mathbf{x})$ 4 T棵树后, 输出总的 $G(\mathbf{x})$ 因为是regression, 直接输出加权实值

GBDT: 'regression sibling' of AdaBoost-DTree
 —**popular in practice**

是Adaboost树的回归版本。降低的是不同的error function

adaboost-决策树是最小化二分类损失的上限。GBDT是最小化回归问题的平方损失
 基本函数 A 都是决策树CART。用不同的加权函数组合起来, 得到 G

Fun Time

Which of the following is the optimal η for

$$\min_{\eta} \quad \frac{1}{N} \sum_{n=1}^N ((y_n - s_n) - \eta g_t(\mathbf{x}_n))^2$$

- ① $(\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)) \cdot (\sum_{n=1}^N g_t^2(\mathbf{x}_n))$
- ② $(\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)) / (\sum_{n=1}^N g_t^2(\mathbf{x}_n))$
- ③ $(\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)) + (\sum_{n=1}^N g_t^2(\mathbf{x}_n))$
- ④ $(\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)) - (\sum_{n=1}^N g_t^2(\mathbf{x}_n))$

第二步，算最优 η
可以得到公式解

$\eta = \frac{\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)}{\sum_{n=1}^N g_t^2(\mathbf{x}_n)}$

Fun Time

Which of the following is the optimal η for

$$\min_{\eta} \quad \frac{1}{N} \sum_{n=1}^N ((y_n - s_n) - \eta g_t(\mathbf{x}_n))^2$$

- ① $(\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)) \cdot (\sum_{n=1}^N g_t^2(\mathbf{x}_n))$
- ② $(\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)) / (\sum_{n=1}^N g_t^2(\mathbf{x}_n))$
- ③ $(\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)) + (\sum_{n=1}^N g_t^2(\mathbf{x}_n))$
- ④ $(\sum_{n=1}^N g_t(\mathbf{x}_n)(y_n - s_n)) - (\sum_{n=1}^N g_t^2(\mathbf{x}_n))$

Reference Answer: ②

Derived within Lecture 9 of ML Foundations,
remember? :-)

Map of Blending Models

blending: aggregate **after getting diverse** g_t 已经有 g_t 了

uniform

simple
voting/averaging of g_t

降低方差, 结果稳定

non-uniform

linear model on
 g_t -transformed inputs

线性组合, linear model学习权重at

conditional

nonlinear model on
 g_t -transformed inputs

g_t 的非linear 组合

uniform for 'stability';
non-uniform/conditional **carefully** for
'complexity'

相对容易过拟合

Map of Aggregation-Learning Models

learning: aggregate **as well as getting** **diverse** g_t

Bagging

diverse g_t by
bootstrapping;
 uniform vote
 by nothing :-)

$D_t \sim$
 +决策树 (不同维度特征)
 (完全长成 差异大)

AdaBoost

diverse g_t
 by **reweighting**;
 linear vote
 by steepest search

数据D
 g_t 降低
 加权分
 类误差

不断改变 u_t /贪心的学习最优的 g_t, a_t , 降低分类误差 exp error

GradientBoost

diverse g_t
 by **residual fitting**;
 linear vote
 by steepest search

与adaboost类似。但最小化的是regression error
 决策树 g_t 拟合的是残差。 a_t 最小化的也是残差

Decision Tree

diverse g_t
 by data splitting;
 conditional vote
 by branching

可以用surrogate解决缺失值问题
 可以回归/分类
 可以处理category类型的特征

输出是实数

boosting-like algorithms most popular

Map of Aggregation of Aggregation Models

Bagging

AdaBoost

Decision Tree

Random Forest

randomized bagging
+ 'strong' DTree

AdaBoost-DTree

AdaBoost
+ 'weak' DTree

一开始弱一点。否则 $et=0$, at 无穷

如果是决策树, ut 不好直接搞到
通过 ut sample 产生数据 D , 再送给 gt

最后在 D , ut 上, 计算 et , at

一般是非完全树, 提前剪枝。
限制阈值、树高

完全长成的 gt : 对输入更加敏感

加上 bootstrap + random d' 维特征

diversity

可以用 OOV 来做自我验证 + 特征选择

特征选择时, 用对应维度的特征取值
permutation 得到 D_p . 和原始输入 D 结果
作对比。

可以 G 不变, 只在 Eval 时做。

差距大, 该特征重要

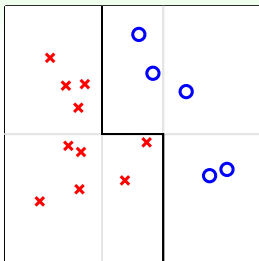
GradientBoost

GBDT

GradientBoost
+ 'weak' DTree

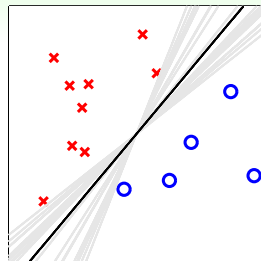
all three frequently used in practice

Specialty of Aggregation Models



cure underfitting

- $G(\mathbf{x})$ 'strong' 类似于特征变换
 - aggregation 通过特征组合使模型变强
有一定的解决欠拟合的效果
- ⇒ **feature transform**



cure overfitting

- $G(\mathbf{x})$ 'moderate' 有一定的解决过拟合的效果
 - aggregation 比如降低方差/增大margin
- ⇒ **regularization**

proper aggregation (a.k.a. 'ensemble')
⇒ **better performance**

Fun Time

Which of the following aggregation model learns diverse g_t by **reweighting** and calculates **linear vote** by **steepest search**?

- ① AdaBoost
- ② Random Forest
- ③ Decision Tree
- ④ Linear Blending

Fun Time

Which of the following aggregation model learns diverse g_t by **reweighting** and calculates **linear vote** by **steepest search**?

- ① AdaBoost
- ② Random Forest
- ③ Decision Tree
- ④ Linear Blending

Reference Answer: ①

Congratulations on being an **expert** in aggregation models! :-)

Summary

- 1 Embedding Numerous Features: Kernel Models
- 2 Combining Predictive Features: Aggregation Models

Lecture 11: Gradient Boosted Decision Tree

- Adaptive Boosted Decision Tree
 - sampling and pruning for 'weak' trees**
- Optimization View of AdaBoost
 - functional gradient descent on exponential error**
- Gradient Boosting
 - iterative steepest residual fitting**
- Summary of Aggregation Models
 - some cure underfitting; some cure overfitting**

- 3 Distilling Implicit Features: Extraction Models
 - **next: extract features other than hypotheses**