

# Boosting

## Bias-Variance Tradeoff

$$E[(h_0(x) - y)^2] = \underbrace{E[(h_0(x) - \bar{h}(x))^2]}_{\text{variance}} + \underbrace{E[(\bar{h}(x) - \bar{y}(x))^2]}_{\text{bias}} + \underbrace{E[(\bar{y}(x) - y)^2]}_{\text{noise}}$$

$h \in \mathcal{H}$   
hypothesis class  
high bias  
weak learners

Kearns 88:  $\mathcal{H} \rightarrow \mathcal{H}^*$   
Schapire 90: Yes!

"strong learners"  
low bias classifier

$$H_T(x) = \sum_{t=1}^T \alpha_t h_t(x) = \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots$$

step size

$$\ell(H) = \frac{1}{n} \sum_{i=1}^n \ell(H(x_i), y_i)$$

imagine  
argmin  
 $h \in \mathcal{H}$

$$\ell(H_t + \alpha h_{t+1}(x)) \quad H_{t+1} = H_t + \alpha h_{t+1}$$

## Gradient Descent in Function Space

Goal: given  $\mathcal{H}$ : partial ensemble

find  $h$ : weak learner

$$\text{Taylor approx of } \ell(H + \alpha h) \approx \ell(H) + \alpha \langle \nabla \ell(H), h \rangle$$

Solution

$$\argmin_{h \in \mathcal{H}} \ell(H + \alpha h) \approx \argmin_{h \in \mathcal{H}} \ell(H) + \alpha \langle \nabla \ell(H), h \rangle$$

$$\langle g, f \rangle = \int g(x) f(x) dx$$

$$\langle g, f \rangle = \sum_{i=1}^n g(x_i) f(x_i)$$

$$= \sum_{i=1}^n \frac{\partial \ell}{\partial H}(x_i) h(x_i)$$

$$= \sum_{i=1}^n \frac{\partial \ell}{\partial [H(x_i)]} h(x_i)$$

$$h_{t+1} = \argmin_{h \in \mathcal{H}} \sum_{i=1}^n r_i h(x_i)$$

$$h_{t+1} = A(\{x_i, r_i\}_{i=1}^n) \rightarrow h_{t+1} \sum_{i=1}^n r_i h(x_i) < 0$$

## Generic Boosting

Input:  $\ell, \alpha, \{x_i, y_i\}_{i=1}^n, A$

$H_0 = 0$

for  $t = 0: T-1$  do

$$r_i, v_i = \frac{\partial \ell}{\partial [H(x_i)]}$$

$$h_{t+1} = A(\{x_i, r_i\}_{i=1}^n)$$

if  $\sum_{i=1}^n r_i h_{t+1}(x_i) < 0$  then

$$H_{t+1} = H_t + \alpha h_{t+1}$$

else

return  $f_t$

end

end

return  $H_T$

## #1 Gradient Boosted Regression Trees (GBRT)

Setting

- \*  $y \in \{-1, 1\}$   $y \in \mathbb{R}$
- \* fixed-depth (say  $k$ ) CART,  $h(x) \in \mathbb{R}$
- \*  $\alpha$  - fixed to small constant (hyperparameter)
- \*  $\ell$ : diff, convex, decomposes

$$\ell(H) = \sum_{i=1}^n \ell(H(x_i))$$

Goal

$$\text{find } h = \argmin_{h \in \mathcal{H}} \sum_{i=1}^n r_i h(x_i)$$

Assumptions

1.  $\sum_{i=1}^n h^2(x_i) = \text{constant}$  why?
2.  $\mathcal{H}$  is negation closed:  $\forall h \in \mathcal{H} \Rightarrow \exists -h \in \mathcal{H}$

$$\argmin_{h \in \mathcal{H}} \sum_{i=1}^n r_i h(x_i) \quad t_i = -r_i$$

$$= \argmin_{h \in \mathcal{H}} -2 \sum_{i=1}^n t_i h(x_i)$$

$$= \argmin_{h \in \mathcal{H}} \sum_{i=1}^n t_i^2 - 2 t_i h(x_i) + \underbrace{h(x_i)^2}_{\text{constant}}$$

$$= \argmin_{h \in \mathcal{H}} \sum_{i=1}^n (h(x_i) - t_i)^2$$

CART

GBRT

Input:  $\ell, \alpha, \{(x_i, y_i)\}$

$H_0 = 0$

for  $t = 1: T$  do

$$r_i = t_i = y_i - H(x_i) \quad \left( -\frac{\partial \ell}{\partial [H(x_i)]} \text{ in gen'l} \right)$$

$$h = \argmin_{h \in \mathcal{H}} \sum_{i=1}^n (h(x_i) - t_i)^2$$

$$H \leftarrow H + \alpha h$$

end

return  $H$

## #2 AdaBoost

Setting

- \*  $y \in \{-1, 1\}$
- \*  $h \in \mathcal{H}, h(x_i) \in \{-1, 1\}$
- \*  $\alpha$ : optimal  $\alpha$  every iteration
- \*  $\ell$ : exponential loss:  $\ell(H) = \sum_{i=1}^n e^{-y_i H(x_i)}$

Finding the best  $h$

$$r_i = \frac{\partial \ell}{\partial H(x_i)} = -y_i e^{-y_i H(x_i)}$$

$$w_i = \frac{1}{Z} e^{-y_i H(x_i)}$$

$$Z = \sum_{i=1}^n e^{-y_i H(x_i)}$$

$$\sum_{i=1}^n w_i = 1$$

$$h_{t+1} = \argmin_{h \in \mathcal{H}} \sum_{i=1}^n r_i h(x_i)$$

$$= \argmin_{h \in \mathcal{H}} - \sum_{i=1}^n y_i e^{-y_i H(x_i)} h(x_i)$$

$$= \argmin_{h \in \mathcal{H}} - \sum_{i=1}^n y_i w_i h(x_i)$$

$$= \argmin_{h \in \mathcal{H}} \sum_{i: h(x_i) \neq y_i} w_i - \sum_{i: h(x_i) = y_i} w_i$$

$$= \argmin_{h \in \mathcal{H}} \sum_{i: h(x_i) \neq y_i} w_i$$

$$\text{weighted classification error } \epsilon$$

$$\sum_{i: h(x_i) = y_i} w_i = 1 - \sum_{i: h(x_i) \neq y_i} w_i$$

Find the best  $\alpha$

$$\alpha = \argmin_{\alpha} \ell(H + \alpha h)$$

$$= \argmin_{\alpha} \sum_{i=1}^n e^{-y_i [H(x_i) + \alpha h(x_i)]}$$

$$\nabla_{\alpha} \ell = - \sum_{i=1}^n y_i h(x_i) e^{-[y_i H(x_i) + \alpha y_i h(x_i)]} \Rightarrow 0$$

$$= \sum_{i: h(x_i) y_i = 1} \frac{-y_i H(x_i) + \alpha y_i h(x_i)}{1} + \sum_{i: h(x_i) y_i = -1} \frac{-y_i H(x_i) + \alpha y_i h(x_i)}{-1} = 0$$

$$= - \sum_{i: h(x_i) y_i = 1} w_i e^{-\alpha} + \sum_{i: h(x_i) y_i = -1} w_i e^{\alpha} = 0$$

$$-(1 - \epsilon) e^{-\alpha} + \epsilon e^{\alpha} = 0$$

$$e^{2\alpha} = \frac{1 - \epsilon}{\epsilon}$$

$$\alpha = \frac{1}{2} \ln \frac{1 - \epsilon}{\epsilon}$$

Re-normalization

$$H_{t+1} = H_t + \alpha h$$

$$\hat{w}_i \leftarrow \hat{w}_i + e^{-\alpha h(x_i) y_i}$$

$$Z \leftarrow Z * 2 \sqrt{\epsilon(1 - \epsilon)}$$

$$w_i \leftarrow w_i \frac{e^{-\alpha h(x_i) y_i}}{2 \sqrt{\epsilon(1 - \epsilon)}}$$

Further analysis

weight update

$$h(x_i) y_i \in \{-1, 1\}$$

$$e^{\alpha} > 1 \text{ if } h(x_i) y_i = -1$$

$$e^{-\alpha} < 1 \text{ if } h(x_i) y_i = 1$$

normalization update

$$\ell(H) = Z = \prod_{t=1}^T 2 \sqrt{\epsilon_t(1 - \epsilon_t)}$$

$$C = \max_{\epsilon} \epsilon$$

$$\ell(H) \leq n [2 \sqrt{C(1 - C)}]^T$$

$$\frac{1}{2} = \argmax_{\epsilon} C(1 - \epsilon) \quad \epsilon_t < \frac{1}{2} \quad \forall t$$

$$C(1 - C) < \frac{1}{4}$$

$$C(1 - C) = \frac{1}{4} - \epsilon^2$$

$$\ell(H) \leq n [1 - 4\epsilon^2]^{\frac{T}{2}}$$

$< 1$  trivial error

margin error

$$\delta_{H(x) \neq y} < e^{-\alpha h(x) y}$$

$$\sum_{i=1}^n e^{-\alpha h(x_i) y_i} = \ell(H) \leq n [1 - 4\epsilon^2]^{\frac{T}{2}} < 1$$

Summarize

$h \in \mathcal{H} \rightarrow$  low bias

high bias

Final Points

\* Stochastic gradient boosting

$$h = \argmin_{h \in \mathcal{H}} \sum_{i=1}^n \dots$$

\* Test time

$$H(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

from error must be 0

$$O(\log(n))$$

$$T > \frac{2 \log(n)}{\log(1 - 4\epsilon^2)}$$