



## Lecture 24

- Boosting
- Nearest Neighbor classifier

**AdaBoost:** ensemble method for classification (or regression)  
 (adaptive boosting)

- trains multiple learners on weighted sample pts. (like bagging)
  - different weights for each learner
  - increased weights for misclassified pts.
  - higher weight to more accurate learners

Input:  $n \times d$  design matrix  $X$ , label vector  $\vec{y} \in \mathbb{R}^n$ ,  $y_i = \pm 1$

Main ideas:  "trees"

- train  $T$  classifiers,  $G_1, \dots, G_T$
- weight of sample pt.  $X_i$  in  $G_t$  grows according to how  $G_1, \dots, G_{t-1}$  misclassified  
 $\hookrightarrow$  if  $X_i$  misclassified by every learner, its weight grows even more if and vice versa
- train  $G_T$  to train harder to correctly classify sample pts. w/ larger weights
- **metaclassifier:** linear combo of all learners

$$M(\vec{z}) = \sum_{t=1}^T \beta_t G_t(\vec{z}) \quad - G_t(\vec{z}) = \pm 1, M \text{ continuous}$$

- return  $\text{Sign}(M(\vec{z}))$

Although boosting was originally designed for decision trees, it works well for most learning algorithms

Recall to weight pts. in trees, we used weighted entropy (proportion of weight in each class)

What  $\beta_T$  and  $G_T$  do we use?

- pick loss fn. and find  $G_T, \beta_T$  to minimize Risk =  $\frac{1}{n} \sum_{i=1}^n L(M(x_i), y_i)$

**Exponential loss fn:**  $L(p, l) = e^{-pl} = \begin{cases} e^{-p} & l = +1 \\ e^p & l = -1 \end{cases}$

$\hookrightarrow$  AdaBoost meta learner

uses this,  $G_t(\vec{z})$  uses another loss fn. if any

Important:  $\ell_t$ ,  $G_t$  binary but  $\ell = M(x_i)$  continuous

advantage of exponential loss fn: pushes hard against badly misclassified pts.

↳ better than squared error for metal learners

↳ same reason why in NNs we use cross-entropy instead of squared error

$$\begin{aligned} n \cdot \text{Risk} &= \sum_{i=1}^n \ell(M(x_i), y_i) \\ &= \sum_{i=1}^n e^{-y_i M(x_i)} \\ &= \sum_{i=1}^n e^{-y_i \sum_{t=1}^T \beta_t G_t(x_i)} \\ &= \sum_{i=1}^n \prod_{t=1}^T e^{-\beta_t y_i G_t(x_i)} \\ &= \sum_{i=1}^n (w_i)^T e^{-\beta_T y_i G_T(x_i)}, \quad w_i^T = \prod_{t=1}^{T-1} e^{-\beta_t y_i G_t(x_i)} \end{aligned}$$

Now, split sum in 2 parts depending on correct classification:  $y_i = G_T(x_i)$   
or incorrect:  $y_i \neq G_T(x_i)$

$$\begin{aligned} \rightarrow n \cdot \text{Risk} &= e^{-\beta_T} \sum_{y_i = G_T(x_i)} w_i^T + e^{\beta_T} \sum_{y_i \neq G_T(x_i)} w_i^T \\ &= e^{-\beta_T} \sum_{i=1}^n w_i^T + (e^{\beta_T} - e^{-\beta_T}) \sum_{y_i \neq G_T(x_i)} w_i^T \end{aligned}$$

What  $G_T$  minimizes the risk?

the teacher who minimizes  $w_i^T$  across the misclassified pts

Now we know what weight to assign to each pt.

- minimize risk  $\equiv$  minimize total weight of misclassified pts. for weight fn.  $w_i^T$

$$w_i^{T+1} = w_i^T e^{-\beta_T y_i G_T(x_i)} = \begin{cases} w_i^T e^{-\beta_T} & y_i = G_T(x_i) \\ w_i^T e^{\beta_T} & \text{o/w} \end{cases}$$

(recursion)

↳ benefit of choosing loss fn.

- Not always possible to pick learner that always classifies pts. correctly. Like SVMs and non-lin. separable cases
- NP-hard to find optimal linear classifier
- ↳ in practice  $G_T$  won't be true minimizer of training error

Optimal value of  $\beta_T$ :

$$\text{Set gradient of risk to 0: } 0 = -e^{-\beta_T} \sum_{i=1}^n w_i^{(T)} + (e^{\beta_T} + e^{-\beta_T}) \sum_{i=1}^n w_i T \\ y_i + G_T(x_i)$$

$$\text{divide both sides by 1st term: } 0 = -1 + (e^{2\beta_T} + 1) \text{err}_T, \text{err}_T = \frac{\sum_{i=1}^n y_i + G_T(x_i) w_i T}{\sum_{i=1}^n w_i T}$$

Optimal Metaclassifier

$$\Rightarrow \beta_T = \frac{1}{2} \ln \left( \frac{1 - \text{err}_T}{\text{err}_T} \right)$$

$\text{err}_T = 0 \Rightarrow \beta_T = \infty$  (perfect learner gets infinite vote/weight)

$\text{err}_T = 1/2 \Rightarrow \beta_T = 0$  (learner w/ 50% acc. gets no vote at all)

- More accurate learners get more votes in Metaclassifier

- acc < 50%  $\Rightarrow$  negative vote (learner w/ 40% acc = learner w/ 60% acc w/ vote sign flipped)

### AdaBoost Algorithm:

1) initialize weights  $w_i = \frac{1}{n}, i=1 \dots n$

2) for  $t=1 \dots T$ :

- compute weighted error rate  $\text{err} = \frac{\sum_{i \in \text{misclassified}} w_i}{\sum_{i \in \text{all}} w_i}$
  - compute weight  $\beta_t = \ln \left( \frac{1 - \text{err}}{\text{err}} \right)$
  - re-weight points:  $w_i = w_i \cdot \begin{cases} e^{\beta_t} & \text{if } y_i \text{ misclassified} \\ e^{-\beta_t} & \text{o/w} \end{cases}$
- $$= w_i \cdot \begin{cases} \sqrt{\frac{1 - \text{err}}{\text{err}}} & \text{if } y_i \text{ misclassified} \\ \sqrt{\frac{\text{err}}{1 - \text{err}}} & \text{o/w} \end{cases}$$

3) return Metaclassifier  $h(\vec{x}) = \text{sign} \left( \sum_{t=1}^T \beta_t g_t(x_i) \right)$

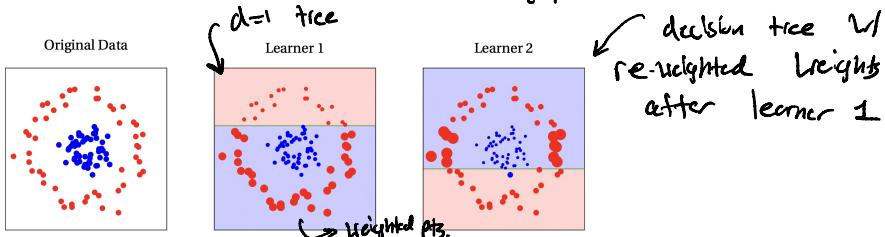


Figure 24.1: Example of AdaBoost with two learners

## Why boost decision trees? Why short trees?

- it's fast
  - training many learners and running many learners at classification time
  - short decision trees fast at training and testing
- no hyperparam search needed
- easy to make tree beat 55% training acc. (or any other threshold) consistently
- easy b/v control
  - boosting can overfit  $\rightarrow$  make AdaBoost trees short to avoid overfitting
    - as # learners  $\uparrow$ ,  $B \downarrow$
    - variance often decreases at first (since successive trees focus on different features)
    - increases later
- AdaBoost & short trees is a form of subset selection
  - features that don't improve metalearner aren't used at all
  - reduces overfitting & runtime especially if lots of useless features
- linear decision boundaries don't boost well
  - boosting linear classifiers gives another rough linear classifier
  - SVMs not good to use
  - nonlinear methods benefit more from boosting
    - allow boosting to reduce bias faster
  - sometimes AdaBoost uses w/  $d=1$  trees but this is linear (not ideal)
    - $d=2$  much better

## Facts about AdaBoost:

- posterior prob. can be approximated:  $P(Y=1|\vec{x}) \approx \frac{1}{1+e^{-2M(\vec{x})}}$
- if data vulnerable to outliers, use exponential loss
  - corrupted data  $\Rightarrow$  use different loss fn. (better loss functions usually have more complicated weights)
- if every learner beats acc.  $M = 50\%$ , then metalearner acc. will eventually be 100%
  - (proven in HW)

cont. below

## Nearst Neighbor Classification

given query pt.  $q$ , find  $k$  samples closest to  $q$

- use distance metric of choice

- regression can give average label of  $k$  pts, classification returns majority of  $k$  labels  
(neither are what we want)

↓  
1 nearest neighbor

Thm (Cover & Hart): as  $n \rightarrow \infty$ ,  $1\text{-NN}$  error rate  $\leq 2B - B^2$ ,  $B = \text{Bayes risk}$   
- if only 2 classes,  $1\text{-NN}$  error rate is  $\leq 2B - 2B^2$

→ lots of data → NN works quite well

Requirement of Thm: training & test pts. must all be drawn independently from same probability distr.

- applies for any dist metric, not just Euclidean

Thm (Fix & Hodges):  $n \rightarrow \infty, k \rightarrow \infty, \frac{k}{n} \rightarrow 0$ , then  $k\text{-nn}$  error converges to  $B$

→  $k\text{-nn}$  converges to be Bayes optimal