

22. AdaBoost

Thursday, February 29, 2024

3:16 PM

→ a bit more complete, less vague, and see nice Fig 7.2
§10.2 in [SS] or §7 in Mohri

AdaBoost (= Adaptive Boosting) Binary classification, $Y = \{\pm 1\}$ ← vague in [SS]

Given dataset $S = ((x_1, y_1), \dots, (x_m, y_m))$

and a γ -weak-learner algorithm, proceeds in rounds t

Initialize $D_i^{(0)} = \frac{1}{m}$ // in general, $D^{(t)} \in \Delta = m\text{-dim. probability simplex}$
 $= \{D \in \mathbb{R}_+^m : \sum D_i = 1\}$

Iterate $t = 0, 1, \dots, T$

- The weak learner solves weighted ERM and returns h_t
using $D_i^{(t)}$

We (for now) treat the "learning problem" to be finding h

w/ low "true risk", $L_{D^{(t)}}(h) := \sum_{i=1}^m D_i^{(t)} \cdot \mathbb{1}_{h(x_i) \neq y_i}$

→
conceptual
leap: $D^{(t)}$

is empirical and
discrete but
treat it as
any other underlying
distribution

our weak learner either explicitly solves weighted ERM*
(w/ guarantee), or is more generic (i.e. bootstrap resample
 S according to $D^{(t)}$ possibly?)

* in fact, this
is directly solving
true risk (over D^t)
exactly

- Weak learning guarantee is that

$$\epsilon_t := L_{D^{(t)}}(h_t) \leq \frac{1}{2} - \gamma$$

with probability at least $1 - \delta$

Save h_t and give it a weight $w_t = \frac{1}{2} \log\left(\frac{1}{\epsilon_t} - 1\right)$

($w_t > 0$ if $\epsilon_t < \frac{1}{2}$) so lower error $\epsilon_t \Rightarrow$ higher weight

- adjust $D^{(t+1)}$ to give more weight to those samples we
misclassified: in fact, give equal mass to sets of
correct and incorrectly identified pts

$$\tilde{D}_i^{(t+1)} = D_i^{(t)} \cdot \exp\left(-w_t \underbrace{y_i h_t(x_i)}_{\substack{+1 \text{ if we got it right,} \\ -1 \text{ if we got it wrong}}}\right)$$

$$D_i^{(t+1)} = \frac{\tilde{D}_i^{(t+1)}}{\sum_{j=1}^m \tilde{D}_j^{(t+1)}} \quad \text{so that } D^{(t+1)} \in \Delta \quad \text{ie. is a probability}$$

call this normalization Z^t

22a. AdaBoost

Thursday, February 29, 2024

6:03 PM

After iterating, return $h(x) = \text{Sign} \left(\sum_{t=1}^T w_t h_t(x) \right)$ as strong learner
i.e. weighted majority vote

We need to assume realizability for weak learners

(i.e. wrt \mathcal{H} , not \mathcal{B}). We don't need to know the edge γ to run the algorithm (only for analysis / guarantees)

RECALL:

$$\varepsilon_t = L_{D^{(t)}}(h_t) \leq \frac{1}{2} - \gamma$$

$$w_t = \frac{1}{2} \log \left(\frac{1}{\varepsilon_t} - 1 \right) = \frac{1}{2} \log \left(\frac{1-\varepsilon_t}{\varepsilon_t} \right)$$

Before analysis, some lemmas

• Lemma $Z^{(t)} = 2 \sqrt{\varepsilon_t (1-\varepsilon_t)}$

proof: $Z^{(t)} := \sum_{i=1}^m D_i^{(t)} \exp \left(\underbrace{-\frac{1}{2} \log \left(\frac{1-\varepsilon_t}{\varepsilon_t} \right)}_{w_t} \underbrace{y_i h_t(x_i)}_{\varepsilon_t \pm 1} \right)$
 $\begin{matrix} +1 & \text{if correct} \\ -1 & \text{if incorrect} \end{matrix}$

$$= \sum_{i=1}^m D_i^{(t)} \left(\frac{1-\varepsilon_t}{\varepsilon_t} \right)^{-\frac{1}{2} \pm 1}$$

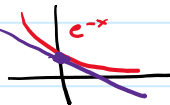
$$= \left(\sum_{i: \text{correct}} D_i^{(t)} \right) \cdot \sqrt{\frac{\varepsilon_t}{1-\varepsilon_t}} + \left(\sum_{i: \text{incorrect}} D_i^{(t)} \right) \cdot \sqrt{\frac{1-\varepsilon_t}{\varepsilon_t}}$$

$\underbrace{\qquad\qquad\qquad}_{=1-\varepsilon_t} \qquad\qquad\qquad \underbrace{\qquad\qquad\qquad}_{=\varepsilon_t}$

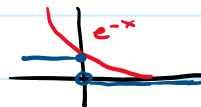
$$= 2 \sqrt{\varepsilon_t (1-\varepsilon_t)} \quad \square$$

• Fact $4\varepsilon(1-\varepsilon) = 1 - 4\left(\frac{1}{2} - \varepsilon\right)^2$

• Recall $1-x \leq e^{-x}$



• Fact $\mathbb{1}_{x \leq 0} \leq e^{-x}$



Thm (10.2 [SS] or 7.2 Mohri) "Adaboost solves ERM"

The training error decays exponentially fast: w.p. $\geq 1 - \delta T$,

after T rounds, $h_S = \text{sign} \left(\sum_{t=1}^T w_t h_t \right)$, $\hat{L}_S(h) := \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{h(x_i) \neq y_i}$,

then $\hat{L}_S(h_S) \leq e^{-\gamma^2 T}$ (using a γ -weak learner)

Proof

Let $f_T := \sum_{t=1}^T w_t h_t$ so $h_S = \text{sign} \circ f_T$. Observe:

$$(*) \quad \hat{L}_S(h_S) := \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i \cdot f_T(x_i) \leq 0} \leq \frac{1}{m} \sum_{i=1}^m e^{-y_i \cdot f_T(x_i)} \quad \text{via } (*)$$

22b. AdaBoost

Friday, March 1, 2024 10:10 AM

(proof continued)

Also observe $D_i^{(\tau+1)} = \frac{D_i^{(\tau)}}{Z^{(\tau)}} \cdot e^{-w_\tau y_i h_\tau(x_i)}$

apply recursively...

$$= \frac{D_i^{(\tau)}}{Z^{(\tau)}} \frac{D_i^{(\tau-1)}}{Z^{(\tau-1)}} \cdot \exp(-w_\tau y_i h_\tau(x_i) - w_{\tau-1} y_i h_{\tau-1}(x_i))$$

recall

$$D_i^{(0)} = \frac{1}{m}$$

$$= \frac{1}{m} \frac{e^{-y_i \sum_{t=1}^T w_t h_t(x_i)}}{\prod_{t=1}^T Z^{(t)}} = \frac{e^{-y_i f_T(x_i)}}{m \prod_{t=1}^T Z^{(t)}}$$

i.e. $e^{-y_i f_T(x_i)} = D_i^{(\tau+1)} \cdot m \cdot \prod_{t=1}^T Z^{(t)}$, plug into (*)

$$(*) \dots \leq \frac{1}{m} \sum_{i=1}^m e^{-y_i f_T(x_i)} = \frac{1}{m} \sum_{i=1}^m D_i^{(\tau+1)} \prod_{t=1}^T Z^{(t)}$$

$= 1$ Since a probability

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

now simplify...

$$= \prod_{t=1}^T \sqrt{1 - 4(\frac{1}{2} - \epsilon_t)^2}$$

$$\leq \prod_{t=1}^T \sqrt{e^{-4(\frac{1}{2} - \epsilon_t)^2}}$$

$1 - x \leq e^{-x}$

$$= \exp\left(-2 \sum_{t=1}^T \underbrace{(\frac{1}{2} - \epsilon_t)^2}_{\leq \gamma^2}\right)$$

$\epsilon_t \leq \frac{1}{2} - \gamma$

$$\leq e^{-2T\gamma^2} \quad \square$$

22c. AdaBoost

Friday, March 1, 2024 10:23 AM

Generalization error of AdaBoost

If each weak learner is $h_t \in \mathcal{B}$, then our final output h_S is in the space — [SS] uses different notation

$$\mathcal{H}_{\mathcal{B}, T} = \left\{ x \mapsto \text{sign} \left(\sum_{t=1}^T w_t h_t(x) \right) : w \in \mathbb{R}^T, h_t \in \mathcal{B} \right\}$$

Bounds:

(from Mohri) $\text{VCdim}(\mathcal{H}_{\mathcal{B}, T}) \leq 2(d+1) \cdot T \cdot \log_2((T+1) \cdot e)$

(from [SS]) " " $\leq (d+1) \cdot T \cdot (3 \ln(T(d+1)) + 2) \quad T, d \geq 3$

$$\approx O(d T \log(T))$$

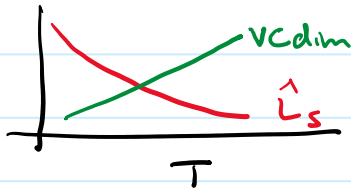
where $d = \text{VCdim}(\mathcal{B})$

For some classes (eg. $\mathcal{B} = L_D$, in which case shouldn't depend on T since L_D is closed under lin. comb.)

this isn't tight, but $\exists \mathcal{B}$ s.t. this is a nearly tight bound.

$$\text{i.e. } \text{VCdim}(\mathcal{H}_{\mathcal{B}, T}) \geq \Omega(d T)$$

So, finite VCdim \Rightarrow we can generalize



Tune T to find location in bias-variance tradeoff you want to be at

Misc

- Sometimes generalization still improves as $T \rightarrow \infty$, even after $\hat{L}_S = 0$
See Mohri for **margin** based analysis
(and coordinate descent interpretation, game theory interpretation, & connections to regularization)
- Doesn't work well w/ noise (i.e. \nexists oracle labeling)
- Another way to tune bias-variance is **early stopping** of SGD