

Ch 21 Online Learning, part 2

Friday, March 27, 2020

3:05 PM

These notes:

- 1) doubling trick
- 2) online-to-batch conversion
- 3) online convex opt.
- 4) perceptron

Aside: the "Doubling Trick"

Given an algorithm that needs to know T in advance, and has an error $\sum_{i=1}^T \text{loss}_i \leq \alpha \cdot \sqrt{T}$ or regret

we can do the following: split T into $\log_2(T)$ groups of the following sizes:

1 2 4 8 16 ...

1 3 7 15 31 ... (running sum)

$$T \leq \sum_{m=0}^L 2^m \quad L = \lceil \log_2(T) \rceil \leq \log_2(T) + 1$$

i.e. $T=7 = 1+2+4 \quad L=2$
 $T=8 = 1+2+4 + \text{part of } 8 \quad L=3$
 $T=9 = 1+2+4 + \text{part of } 8 \quad L=3$

on each group of size 2^m ,

run algorithm, so error on this chunk

$$\text{is } \leq \alpha \cdot \sqrt{2^m} = \alpha \cdot \sqrt{2}^m$$

So, total error is $\leq \sum_{m=0}^L \alpha \cdot \sqrt{2}^m$

$$\begin{aligned} &= \alpha \frac{1 - \sqrt{2}^{L+1}}{1 - \sqrt{2}} \leftarrow \sqrt{2}^L = \sqrt{2}^L \leq \sqrt{T} \\ &\leq \alpha \frac{\sqrt{2}^{L+1} - 1}{\sqrt{2} - 1} \leq \alpha \frac{\sqrt{2}}{\sqrt{2} - 1} \sqrt{T} \approx 3.4 \sqrt{T} \end{aligned}$$

Recall if $S = \sum_{m=0}^L \rho^m = 1 + \rho + \dots + \rho^L$
 $\rho \cdot S = \rho + \dots + \rho^L + \rho^{L+1}$
 so $S - \rho S = 1 - \rho^{L+1}$
 so $\rho \neq 1 \Rightarrow S = \frac{1 - \rho^{L+1}}{1 - \rho}$

so if we knew T in advance, error bound is $\alpha \sqrt{T}$

if we don't know in advance, $3.4 \cdot \alpha \sqrt{T}$ not much worse

Used elsewhere in applied math + computer science

ex line searches in optimization methods \leftarrow mainly theoretical (not good practical performance)

ex growing memory $x = [3, 1, 4]$ // 3 units of memory allocated

$x.append(5)$ // now we need 4 units of memory, so find

better: allocate $2^{\lceil \log_2(n) \rceil}$ at first,

then when full, double memory.

Fewer copies, only at most $2 \times$ sub-optimal

4 free units, copy old memory to new location (slow!)
 since often want contiguous

Aside: Online-to-batch Conversion, converting an online learner to a PAC learner

Thm: Usual PAC learning problem (realizable), binary classification

and suppose we have an online algorithm A w/ a mistake bound $M_A(d) < \infty$

Run $m=T$ iid examples through algo A (labels via $h^* \in \mathcal{H}$), $S \sim \mathcal{D}^m$ aka \mathcal{D}^T

Let h_t be the hypothesis the online algo. used at round t

$$\text{Then } \mathbb{E}_{(x_t) \sim \mathcal{D}^m} L_0(h_r) \leq \frac{M_A(d)}{T} \quad (T=m)$$

$r \sim \text{Uniform}([1, T])$

not of PAC form,

but as we did in ch23 Reg.+Stability, convert to PAC via Markov's ineq (exer.18.1)

proof sketch

$$\mathbb{E}_{r \sim \text{Uni}([T])} L_D(h_r) = \frac{1}{T} \sum_{t=1}^T L_D(h_t)$$

so $\mathbb{E}_{\substack{(x_t) \sim D^m \\ r \sim \text{Uni}([T])}} L_D(h_r) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{(x_t) \sim D^m} L_D(h_t) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{(x_t) \sim D^m} \mathbb{E}_{x \sim D} \mathbb{1}_{h_t(x) \neq h^*(x)}$

h_t depends only on (x_1, \dots, x_{t-1})
 $x_t \sim D$

$$\text{so } \mathbb{E}_{x \sim D} \mathbb{1}_{h_t(x) \neq h^*(x)} = \mathbb{E}_{x_t \sim D} \mathbb{1}_{h_t(x_t) \neq y_t}$$

$$\dots = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{(x_t) \sim D^m} \mathbb{1}_{h_t(x_t) \neq y_t} \quad \text{i.e. } \mathbb{E}_{(x_t) \sim D^m} (\dots)$$

\leftarrow only 1^{st} t matters

$$\leq \frac{1}{T} \sum_{t=1}^T \mathbb{1}_{h_t(x_t) \neq y_t} \quad x_t \text{ adversarial} \quad = \mathbb{E}_{(x_t) \sim D^{t-1}} \mathbb{E}(\dots | (x_t) \sim D^{t-1})$$

\leftarrow law of total expectation

$$\leq \frac{1}{T} M_A(\mathcal{H}) \quad \square$$

Corollary

If the loss function is convex (eg, parameterize h by $w \in \mathcal{H}$, and $\forall z \in X \times Y, w \mapsto \ell(w, z)$ is convex)
 then can choose $h = \frac{1}{T} \sum h_t$ (instead of $h = h_r, r \sim \text{Uni}([T])$) and have the same bound

proof: Jensen's Ineq.

See "Online Learning and Online Convex Optimization", Shalev-Shwartz, Foundations and Trends in Machine Learning, 2011

② Online Convex Optimization

Setup is nearly the same, but instead of 0-1 loss $|h(x) - y|$, now allow a general loss $\ell: \mathcal{H} \times Z \rightarrow \mathbb{R}$, $Z = X \times Y$, parameterize h as $w \in \mathcal{H}$,

assume: 1) \mathcal{H} is a convex set (ex: $\mathcal{H} = \{w: \|w\| \leq B\}$)

2) $\forall z, w \mapsto \ell(w, z)$ is a convex function

i.e., regression or binary classifi w/ a convex surrogate loss

Analyze the regret:

Def: $\text{Regret}_A(w^*, T) := \sum_{t=1}^T \ell(w^{(t)}, z_t) - \sum_{t=1}^T \ell(w^*, z_t)$

Def: $\text{Regret}_A(\mathcal{H}, T) := \sup_{w^* \in \mathcal{H}} \text{Regret}_A(w^*, T)$
 $= \sum_{t=1}^T \ell(w^{(t)}, z_t) - \min_{w^* \in \mathcal{H}} \sum_{t=1}^T \ell(w^*, z_t)$

So, minimize $\sum_{t=1}^T \ell(w^{(t)}, z_t)$] online cvx opt (ch 21) \leftarrow training = testing
 w can change per iteration
 z_t adversarial, choose w_t before we see y_t

vs $\min \sum_{t=1}^T \ell(w, z_t)$] "batch" cvx opt (ch 14) \leftarrow this is training.
 Fix w, z_t iid, all labels known, Then later test

(Different rules than ch. 14)

focus on generalization

Algo: "Online ^{sub}Gradient Descent"

Stepsize $\eta > 0$

Init. $w^{(1)} = 0$ (if $0 \in H$)

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

Our "prediction" is $w^{(t)}$

Observe $z_t = (x_t, y_t)$

Loss function is $f_t(\cdot) = l(\cdot, z_t)$

so suffer loss $f_t(w^{(t)})$

Choose $v_t \in \partial f_t(w^{(t)})$ // Subgradient

Update

$$w^{(t+1)} = \text{Proj}_H(w^{(t)} - \eta v_t)$$

$$\text{Proj}_C(v) = \arg \min_{w \in C} \|w - v\| = \arg \min_{w \in C} \frac{1}{2} \|w - v\|^2$$

Then (online [projected sub]gradient descent) and assuming convexity

Choosing $w^{(t)}$ as in OGD, (i.e. $A = \text{OGD}$)

$$(1) \forall w^* \in H, \text{Regret}_{\text{OGD}}(w^*, T) \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$$

So...

$$(2) \text{ if } f_t \text{ is } \rho\text{-Lipschitz } \forall t \in [T] \text{ then } w_t, \eta = \frac{1}{\sqrt{T}} \text{ i.e. } \|v_t\| \leq \rho$$

$$\text{Regret}_{\text{OGD}}(w^*, T) \leq \frac{1}{2} \sqrt{T} \cdot (\|w^*\|^2 + \rho^2)$$

Sublinear regret!

$$(3) \text{ and if we also assume } H \text{ is } B\text{-bounded, then } w_t, \eta = \frac{B}{\sqrt{T}}$$

$$\text{Regret}_{\text{OGD}}(w^*, T) \leq B\rho\sqrt{T}$$

(if T unknown, then η unknown, so do doubling trick)

proof sketch see book, not very enlightening

compare to
Cor 14.12 in
batch PAC
learning case
Same η , same
bound, different
interpretation

③ Online Perceptron

For homogeneous half-spaces $H = \{x \mapsto \text{sign}(\langle w, x \rangle)\}$

+ binary classification

Ch 9: batch version

Ch 21: online version

$$X = \mathbb{R}^d, Y = \{-1, 1\}$$

... but we saw $\dim(H) = \infty$ if $d \geq 2$ (i.e. 1D + offset)

so a mistake bound not possible

... instead, use a surrogate convex loss

Because we're not trying to generalize, we can change the

surrogate loss to depend on the round t

(including on z_t and $w^{(t)}$), as long as it's a

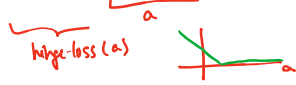
$$\text{Surrogate, i.e., } \underbrace{f_t(w^{(t)})}_{\text{surrogate loss}} \geq \underbrace{l(w^{(t)}, z_t)}_{\text{true loss, eg, 0-1 loss}}$$

So that we bound the risk.

$$\text{Use true loss is 0-1 loss, } l(w, (x, y)) = \frac{1}{2} [y \langle w, x \rangle \leq 0] = \begin{cases} 1 & y \neq \text{sign}(\langle w, x \rangle) \\ 0 & \text{else} \end{cases}$$

Surrogate is

$$f_t(w) = \begin{cases} 0 & \text{when we had correct prediction} \\ \max(0, 1 - y_t \langle w, x_t \rangle) & \text{if we got it wrong} \end{cases}$$



Algo: "Perceptron" Rosenblatt '58

$w_1 = 0$

no stepsize!

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

receive x_t

predict $p_t = \text{sign}(\langle w^{(t)}, x_t \rangle)$

if $y_t \langle w^{(t)}, x_t \rangle \leq 0$ // we got it wrong

$w^{(t+1)} = w^{(t)} + y_t x_t$

else

$w^{(t+1)} = w^{(t)}$

// we got it right
// no change

See Thm 21.16 for analysis