

Lecture 16

- * HW2 due on Gradescope on Wednesday
- If you're presenting in the 1st week, can get extension.
- Please mark pages on Gradescope when you submit.

RECAP

GD:

- ① Initialize w_0 ,
- ② For $t=1, \dots, T$
- ③ $w_{t+1} \leftarrow w_t - \eta \nabla f(w_t)$

Thm (convergence rate of GD) Let f be a convex, ℓ -Lipschitz function i.e.

$$|f(x) - f(y)| \leq \rho \|x - y\|_2 \quad \forall x, y \in \mathbb{R}^d.$$

Let $w^* = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} f(w)$. Let $\|w^* - w_0\| \leq B$, where

w_0 is initialization for GD. Suppose we run GD for T steps with $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ & let $\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$. Then \bar{w} satisfies

$$f(\bar{w}) - f(w^*) \leq \frac{B\rho}{\sqrt{T}}.$$

$$\text{If } T \geq \frac{B^2 \rho^2}{\varepsilon^2}, \quad f(\bar{w}) - f(w^*) \leq \varepsilon.$$

Stochastic gradient descent (SGD)

- ① Initialize w_0
- ② For $t=1 \dots T$
- ③ Choose any v_t s.t. $\mathbb{E}[v_t | w_t] = \nabla f(w_t)$.
- ④ $w_{t+1} \leftarrow w_t - \eta v_t$

Thm (convergence of SGD) Let f be a convex,

ρ -Lipschitz function. Let $w^* = \arg \min_{w \in \mathbb{R}^d} f(w)$ &

$\|w^* - w_0\| = B$ (where w_0 is the initialization). Suppose we run SGD for T steps with step size

$\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ & $\|v_t\| \leq \rho \forall t$. Let $\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$. Then,

$$\mathbb{E}[f(\bar{w})] - f(w^*) \leq \frac{B\rho}{\sqrt{T}}.$$

Proof

Let $v_{1:T}$ to denote the sequence v_1, \dots, v_T .

By convexity,

$$\mathbb{E}_{v_{1:T}} [f(\bar{w}) - f(w^*)] \leq \mathbb{E}_{v_{1:T}} \left[\frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*) \right]$$

Lemma (iterative update) Let v_1, \dots, v_T be an arbitrary sequence of vectors. Consider any algorithm with an update rule

$$w_{t+1} = w_t - \gamma v_t.$$

Then

$$\sum_{t=1}^T \langle w_t - w^*, v_t \rangle \leq \frac{\|w^* - w_1\|^2}{2\gamma} + \frac{\gamma}{2} \sum_{t=1}^T \|v_t\|^2$$

If $\|w^* - w_1\| \leq B$ & $\|v_t\| \leq \rho$ and we set

$$\gamma = \sqrt{\frac{B^2}{\rho^2 T}},$$

$$\frac{1}{T} \sum_{t=1}^T \langle w_t - w^*, v_t \rangle \leq \frac{B\rho}{\sqrt{T}}.$$

By above Lemma,

$$\mathbb{E}_{v_1, \dots, v_T} \left[\frac{1}{T} \sum_{t=1}^T \langle w_t - w^*, v_t \rangle \right] \leq \frac{B\rho}{\sqrt{T}}$$

∴ we only need to show

$$\mathbb{E}_{v_1, \dots, v_T} \left[\frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*) \right] \leq \mathbb{E}_{v_1, \dots, v_T} \left[\frac{1}{T} \sum_{t=1}^T \langle w_t - w^*, v_t \rangle \right]$$

$$\mathbb{E}_{v_1, \dots, v_T} \left[\frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*) \right] = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{v_1, \dots, v_T} [f(w_t) - f(w^*)]$$

$$= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{v_1:t-1} \left[f(w_t) - f(w^*) \right]$$

SGD requires $\mathbb{E}_{v_t} [v_t | w_t] = \nabla f(w_t)$

Since w_t only depends on $v_{1:t-1}$

$$\therefore \mathbb{E}_{v_t} [v_t | v_{1:t-1}] = \nabla f(w_t).$$

$$\begin{aligned} \mathbb{E}_{v_1:T} \left[\frac{1}{T} \sum_{t=1}^T f(w_t) - f(w^*) \right] &\leq \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{v_1:t-1} \left[\langle w_t - w^*, \mathbb{E}_{v_t} [v_t | v_{1:t-1}] \rangle \right] \\ &= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{v_1:t-1} \left[\mathbb{E}_{v_t} [\langle w_t - w^*, v_t \rangle | v_{1:t-1}] \right] \\ &= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{v_1:T} [\langle w_t - w^*, v_t \rangle] \end{aligned}$$

(because $\mathbb{E}_{v_1:t-1} [\mathbb{E}_{v_t} [\cdot | v_{1:t-1}]] = \mathbb{E}_{v_1:t} [\cdot]$)

$$= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{v_1:T} [\langle w_t - w^*, v_t \rangle]$$

$$= \mathbb{E}_{v_1:T} \left[\frac{1}{T} \sum_{t=1}^T \langle w_t - w^*, v_t \rangle \right]$$

and we are done.



Some more points about the convergence rate,

- Same convergence rate as GD!
- the convergence rate does not improve with the β -smoothness assumption.
- for strongly convex functions,
the convergence rate goes from $\frac{1}{\sqrt{T}} \rightarrow \frac{1}{T}$
we can't get the e^{-TK} rates.

So some tradeoffs in terms of cost per iteration
 k # iterations.

Learning with SGD

Suppose we are interested in minimizing the risk

$$R(w) = \mathbb{E}_{z \sim D} [l(w, z)]$$

Using ERM, we sample n datapoints & minimize training error. Use generalization bounds to ensure small test error.

SGD gives a different way to look at this,

$$\nabla R(w_t) = \nabla \mathbb{E}_{z \sim D} [l(w_t, z)] = \mathbb{E}_{z \sim D} [\nabla l(w_t, z)]$$

just the gradient at z !

\therefore If we set $v_t = \nabla l(w_t, z)$ (for $z \sim D$)
 $\mathbb{E}[v_t] = \nabla R(w_t)$

SGD for minimizing $R(w)$

- ① initialize w_0
- ② for $t = 1 \dots T$
- ③ $z \leftarrow E + (c, d)$
- ④ get $v_t = \nabla l(w_t, z)$
- ⑤ update $w_{t+1} = w_t - \eta v_t$
- ⑥ Output $\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$

Corollary Consider a convex, L -Lipschitz function $l(w, z)$. Let $w^* = \underset{w \in \mathcal{W}}{\operatorname{argmin}} R(w)$. Let $\|w^* - w_0\| \leq B$.

Then if we run SGD for T iterations with
 $\eta = \sqrt{\frac{B^2}{L^2 T}}$ where $T \geq \frac{B^2 L^2}{\epsilon^2}$, then the output \bar{w} satisfies,

$$\mathbb{E}[R(\bar{w})] \leq \min_{w \in \mathcal{W}} R(w) + \epsilon.$$

Note SGD is implementable in the SQ model
 (we need to extend SQ model to allow real-valued queries).

Online Learning

PAC learning | Statistical learning :

Learn under probabilistic assumptions on data
(train/test from some distribution)
Generalization theory to do well on unseen test points.

Online learning

No probabilistic assumptions, predict well on a datapoint as we see it.

Example : Weather forecasting

- We're interested in predicting rain/no rain.
- Every night, make prediction about next day.
- Next day, we see whether or not it rained.

There's no train/test split. Every example is both training & test example.

Formally,

At every time step t

- learner receives an input $x_t \in X$.
- makes prediction $p_t \in Y$.
- See true label $y_t \in Y$. suffer loss $l(p_t, y_t)$.

Think $Y = \{0, 1\}$

$$l(p_t, y_t) = \mathbb{I}\{p_t \neq y_t\}.$$

Realizability

Def (Mistake bound model) Let \mathcal{H} be a hypothesis class and A be a online learning algo. Given any sequence $S = (x_1, h^*(x_1)), \dots, (x_T, h^*(x_T))$ of T labelled datapoints where $h^* \in \mathcal{H}$, let $M_A(S)$ be the # mistakes A makes on the sequence S . We denote by $M_A(T)$ to the supremum of $M_A(S)$ over all possible S .

If there exists an algorithm A that satisfies a mistake bound of the form $M_A(T) \leq B < \infty$, we say \mathcal{H} is online learnable in the mistake bound model.

Note : B should be independent of length of sequence T ,
 \therefore as $T \rightarrow \infty$, avg # mistakes $(\leq \frac{B}{T}) \rightarrow 0$.

In PAC learning, ERM choose any consistent hypothesis over training set.

Alg: Consistent

① Initialize $V_1 = \mathcal{H}$

② for $t = 1, \dots, T$

receive x_t

choose any $h \in V_t$

predict $p_t = h(x_t)$

receive $y_t = h^*(x_t)$, loss $\mathbb{1}(p_t \neq y_t)$

update $V_{t+1} = \{ h \in V_t : h(x_t) = y_t \}$.

Proposition: Let \mathcal{H} be a finite hypothesis class.
The above algorithm gets a mistake bound

$$M_{\text{consistent}}(\mathcal{H}) \leq |\mathcal{H}| - 1.$$

Alg: Halving

① Initialize $V_1 = \mathcal{H}$

② for $t = 1, \dots, T$

receive x_t

predict $p_t = \operatorname{argmax}_{s \in \{0,1\}^t} |\{h \in V_t : h(x_t) = s\}|$

(if tie, $p_t = 1$)

receive $y_t = h^*(x_t)$, loss $\mathbb{1}(p_t \neq y_t)$

update $V_{t+1} = \{ h \in V_t : h(x_t) = y_t \}$.

Proposition Let \mathcal{H} be a finite hypothesis class.

Then halving algorithm satisfies the mistake bound

$$M_{\text{Halving}}(s) \leq \log_2 (|\mathcal{H}|).$$

Proof Whenever the algo. errors, $|V_{t+1}| \leq \frac{|V_t|}{2}$.

If M is # mistakes

$$|V_{t+1}| \leq |\mathcal{H}| 2^{-M}$$

$$\text{as } |V_{t+1}| \geq 1$$

$$\Rightarrow M \leq \log_2 (|\mathcal{H}|).$$

■