

Algorithmic ML

Lecture 5 :

Last 2 lectures :

Gradient Descent (GD/NAGD)

Analysis of GD

This lecture :

Stochastic Gradient Descent (SGD)

GD : (parameter : "step size" or "learning rate" η)

1. Choose $w_0 \in \mathbb{R}^d$

2. for $i = 1, \dots, T$:

$$w_i = w_{i-1} - \eta \nabla f(w_{i-1})$$

Theorem : (GD convergence) $f: \mathbb{R}^d \rightarrow \mathbb{R}$ β -smooth & convex

$$\Rightarrow f(w_k) \leq f(w^*) + \frac{\beta \|w^* - w_0\|}{2k} \quad \text{when } \eta = \gamma_\beta .$$

↳ if $\beta, \|x^* - x_0\| = O(1)$

then GD returns w such that $f(w) \leq f(w^*) + \varepsilon$

in $O(\gamma_\varepsilon)$ iterations.

Pro. : Always converges in $O(\gamma_\varepsilon)$ iterations.

Con. : What is the time complexity per iteration?

→ Cost of computing $\nabla f(w)$.

Taking a step back: what f are we working with?

Recall ERM: (empirical risk minimization) Given labeled examples ("training data") $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \mathbb{R} \dots$

Empirical loss function $\underset{\text{define}}{\rightarrow} L(w) = \frac{1}{n} \sum_{i=1}^n l(h_w(x_i), y_i)$ where

Our "f" in GD. $h_w: \mathbb{R}^d \rightarrow \mathbb{R}$ the hypothesis corresponding to w (parametrized by)

Trying to find w^* that minimizes this loss. $l: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ a loss function.

$$l(a, b)$$

$$\begin{aligned} \text{Gradient: } \nabla L(w) &= \frac{1}{n} \sum_{i=1}^n \nabla l(h_w(x_i), y_i) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial a} l(h_w(x_i), y_i) \frac{\partial}{\partial w} h(x_i) \end{aligned}$$

Example: LSR (least squares regression)

$$\left. \begin{array}{l} \bullet h_w(x_i) = \langle w, x_i \rangle \\ \bullet l(a, b) = (a - b)^2 \end{array} \right\} \Rightarrow L(w) = \frac{1}{n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i)^2$$

$$\text{Gradient: } \nabla L(w) = \frac{2}{n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i) x_i$$

$$\begin{aligned} \nabla L(w)_j &= \frac{1}{n} \sum_{i=1}^n 2(\langle w, x_i \rangle - y_i) \underbrace{\frac{\partial}{\partial w_j} \langle w, x_i \rangle}_{= x_{ij}} \\ &= \frac{2}{n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i) x_{ij} = \frac{\partial}{\partial w_j} \sum_{i=1}^n w_i x_{ij} = x_{ij} \end{aligned}$$



Costs $O(n \cdot d)$ time to compute.

(n examples, d dimensions)

\Rightarrow Time complexity of GD $\approx O\left(\frac{n \cdot d}{\epsilon}\right)$

$n \cdot d$ very bad in many applications.

Stochastic Gradient Descent:

Idea: efficiently estimate $\nabla f(w)$.

- Let $G(w)$ be an estimator for $\nabla f(w)$

i.e. $E[G(w)] = \nabla f(w)$ (unbiased estimator)

$\rightarrow G(w)$ a random vector in \mathbb{R}^d

e.g. LSR: $\nabla L(w) = \frac{2}{n} \sum_{i=1}^n ((w, x_i) - y_i) x_i$

Could use: $G(w) = 2((w, x_{i^*}) - y_{i^*}) x_{i^*}$ where
 i^* sampled u.a.r.

or . for $k \ll n$:

$$G(w) = \frac{2}{k} \sum_{i \in S} ((w, x_i) - y_i) x_i \text{ where } S \subseteq \{1, \dots, n\}, |S| = k \text{ random.}$$

SGD: (using estimator G)

1) choose w_0

2) for $i=1, \dots, T$:

$$\text{let } w_{i+1} = w_i - \eta G(w_i)$$

How "good" is G_i ? Small variance is good.

$$\cdot \text{Var}(G_i) = \mathbb{E}[\|G_i(\omega) - \mathbb{E}[G_i(\omega)]\|_2^2]$$

$$\begin{aligned} & \mathbb{E}[\|G_i - \mathbb{E}[G_i]\|_2^2] \\ &= \mathbb{E}\left[\sum_{j=1}^d (G_{ij} - \mathbb{E}[G_{ij}])^2\right] = \sum_{j=1}^d \mathbb{E}[(G_{ij} - \mathbb{E}[G_{ij}])^2] \\ &= \sum_{j=1}^d \text{Var}(G_{ij}) \end{aligned}$$

Then (SGD convergence) :

f convex & β -smooth , $\eta \leq 1/\beta$,

$$\text{Var}(G_i(\omega)) \leq \sigma^2$$

$$\Rightarrow \mathbb{E}[f(\bar{\omega}_k)] \leq f(\omega^*) + \frac{\|\omega_0 - \omega^*\|_2^2}{2\eta \cdot k} + \eta \sigma^2$$

$$\bar{\omega}_k = \frac{1}{k} (\omega_0 + \dots + \omega_k)$$

error term added by
term varying in

Remarks :

- SGD is what is most commonly used in practice (variants).
- Can get "accelerated" version of SGD (as in NAGD) just by replacing $\nabla f(\omega_k)$ with $G_i(\omega_k)$ in NAGD.

Proof of SGD convergence :

useful things: $\forall u, v \in \mathbb{R}^d$:

- Convexity: $f(u) \leq f(v) + \langle \nabla f(v), u - v \rangle$

- Smoothness: $f(u) \leq f(v) + \langle \nabla f(v), u - v \rangle + \frac{\ell}{2} \|u - v\|_2^2$

- Cosine formula: $2 \langle u, v \rangle = \|u\|_2^2 + \|v\|_2^2 - \|u - v\|_2^2$

- SGD update: $w_{t+1} = w_t - \eta G(w_t)$

- Unbiased estimator: $\mathbb{E}[G(\omega)] = \nabla f(\omega)$

$$\Rightarrow \mathbb{E}\left[\|G(\omega)\|_2^2\right] - \|\nabla f(\omega)\|_2^2 \leq \sigma^2$$

$$\begin{aligned} \sigma^2 &\geq \mathbb{E}\left[\|G(\omega) - \mathbb{E}[G(\omega)]\|_2^2\right] = \mathbb{E}\left[\|G(\omega)\|_2^2\right] - \|\mathbb{E}[G(\omega)]\|_2^2 \\ &= \mathbb{E}\left[\|G(\omega)\|_2^2\right] - \|\nabla f(\omega)\|_2^2 \end{aligned}$$

Clm 1 : (SGD function decrease ineq) $\forall k$

(from smoothness)

$$\star \quad \mathbb{E}[f(w_{k+1})] \leq \mathbb{E}[f(w_k)] - \frac{\eta}{2} \mathbb{E}\left[\|\nabla f(w_k)\|_2^2\right] + \frac{\eta\sigma^2}{2}$$

pf:

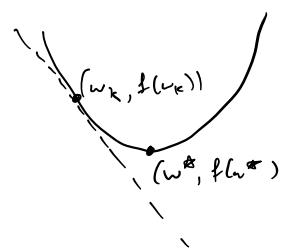
$$\begin{aligned} (\text{smoothness}) \quad f(w_{k+1}) &\leq f(w_k) + \langle \nabla f(w_k), w_{k+1} - w_k \rangle + \frac{\beta}{2} \|w_{k+1} - w_k\|_2^2 \\ &= f(w_k) + \langle \nabla f(w_k), -\eta G(w_k) \rangle + \frac{\beta\eta^2}{2} \|G(w)\|_2^2 \\ (\cancel{w_{k+1} = w_k}) \quad \sigma^2 &\geq \mathbb{E}\left[\|G(w) - \mathbb{E}[G(w)]\|_2^2\right] = \mathbb{E}\left[\|G(w)\|_2^2\right] - \|\mathbb{E}[G(w)]\|_2^2 \\ &= \mathbb{E}\left[\|G(w)\|_2^2\right] - \|\nabla f(w)\|_2^2 \end{aligned}$$

$$\begin{aligned} \mathbb{E}[f(w_{k+1})] &\leq \mathbb{E}[f(w_k)] - \eta \mathbb{E}\left[\|\nabla f(w_k)\|_2^2\right] + \frac{\beta\eta^2}{2} \mathbb{E}\left[\|\nabla f(w_k)\|_2^2\right] + \frac{\beta\eta^2\sigma^2}{2} \\ \left(\eta \leq \frac{1}{\beta}\right) \quad &\leq \mathbb{E}[f(w_k)] - \frac{\eta}{2} \mathbb{E}\left[\|\nabla f(w_k)\|_2^2\right] + \frac{\eta\sigma^2}{2}. \quad \square \end{aligned}$$

Now, similar to GD analysis, let's use convexity of f :

$$\star \star \quad \text{Convexity: } f(w_k) \leq f(w^*) + \langle \nabla f(w_k), w_k - w^* \rangle$$

$$(- \langle \nabla f(w_k), w_k - w^* \rangle \leq f(w^*) - f(w_k))$$



$$\star \star \star \quad \frac{1}{2\eta} \mathbb{E}\left[\|w_{k+1} - w^*\|_2^2 - \|w_k - w^*\|_2^2\right]$$

(Relates point distance
to the gradient)

$$= \frac{\eta}{2} \mathbb{E}\left[\|\nabla f(w_k)\|_2^2\right] - \mathbb{E}\left[\langle \nabla f(w_k), w_k - w^* \rangle\right] + \frac{\eta\sigma^2}{2}$$

Twice neg:

$$\text{Pf : } \|\omega_{k+1} - \omega^*\|_2^2 = \underbrace{\|\omega_k - \omega^* - \eta g(\omega_k)\|_2^2}_{\text{in}} = \|\omega_k - \omega^*\|_2^2 + \eta^2 \|g(\omega_k)\|_2^2 - 2\eta \langle \omega_k - \omega^*, g(\omega_k) \rangle$$

$$\Rightarrow \mathbb{E} \left[\|\omega_{k+1} - \omega^*\|_2^2 - \|\omega_k - \omega^*\|_2^2 \right] + \eta^2 \mathbb{E} \left[\|\nabla f(\omega_k)\|_2^2 + \eta^2 \sigma^2 \right] - 2\eta \mathbb{E} \left[\langle \nabla f(\omega_k), \omega_k - \omega^* \rangle \right]$$

multiply by $\frac{1}{2\eta} \dots \square$

Main proof continued...

$\star\star\star + \star\star \Rightarrow$

$$(-\mathbb{E}[\langle \nabla f(\omega_k), \omega_k - \omega^* \rangle]) \leq f(\omega^*) - \mathbb{E}[f(\omega_k)]$$

$$\begin{aligned} \frac{1}{2\eta} \mathbb{E} \left[\|\omega_{k+1} - \omega^*\|_2^2 - \|\omega_k - \omega^*\|_2^2 \right] &\leq \frac{\eta}{2} \mathbb{E} \left[\|\nabla f(\omega_k)\|_2^2 \right] + f(\omega^*) - \mathbb{E}[f(\omega_k)] + \frac{\eta\sigma^2}{2} \\ &\leq f(\omega^*) - \mathbb{E}[f(\omega_{k+1})] + \eta\sigma^2 \quad (\text{by } \star) \end{aligned}$$

$$\Rightarrow \forall i=1,\dots,k : \quad (\text{Telescoping sum trick})$$

$$\mathbb{E}[f(\omega_i)] \leq f(\omega^*) - \frac{1}{2\eta} \mathbb{E} \left[\|\omega_i - \omega^*\|_2^2 - \|\omega_{i-1} - \omega^*\|_2^2 \right] + \eta\sigma^2$$

$$\begin{aligned} \Rightarrow \sum_{i=1}^k \mathbb{E}[f(\omega_i)] &\leq k \cdot f(\omega^*) - \frac{1}{2\eta} (\|\omega_k - \omega^*\|_2^2 - \|\omega_0 - \omega^*\|_2^2) + k\eta\sigma^2 \\ &\leq k \cdot f(\omega^*) + \frac{\|\omega_0 - \omega^*\|_2^2}{2\eta} + k\eta\sigma^2 \end{aligned}$$

$$\Rightarrow \mathbb{E} \left[\frac{1}{k} \sum_{i=1}^k f(\omega_i) \right] \leq f(\omega^*) + \frac{\|\omega_0 - \omega^*\|_2^2}{2\eta} + \eta\sigma^2$$

$$\text{Now, } f \text{ convex} \Rightarrow f \left(\frac{1}{k} \sum_{i=1}^k \omega_i \right) \leq \frac{1}{k} \sum_{i=1}^k f(\omega_i)$$

and recall that SGD output is $\bar{\omega}_k = \frac{1}{k} \sum_{i=1}^k \omega_i$. \square

Concluding Remarks :

Then (SGD convergence) :

f convex & β -smooth , $\eta \leq \frac{1}{\beta}$,

$$\text{Var}(G(\omega)) \leq \sigma^2$$

$$\Rightarrow \mathbb{E}[f(\bar{\omega}_k)] \leq f(\omega^*) + \frac{\|\omega_0 - \omega^*\|_2^2}{2\eta \cdot k} + \eta \sigma^2$$

Choosing step-size :

in standard GD :

just want to pick η as large as possible (i.e. $\eta = \frac{1}{\beta}$)

large η is good:
allows us to converge faster
(as long as $\eta \leq \frac{1}{\beta}$)

small η is good:
mitigates error of estimator

in SGD : have to balance the 2 error terms

general guideline : choose $\eta \approx \frac{1}{\sqrt{k}}$

$$\eta = \frac{1}{\sqrt{k}} \Rightarrow \mathbb{E}[f(\bar{\omega}_k)] \leq f(\omega^*) + \frac{1}{\sqrt{k}} \left(\frac{\|\omega_0 - \omega^*\|_2^2}{2} + \sigma^2 \right)$$

$$\eta = \frac{\|\omega_0 - \omega^*\|_2}{\sigma \sqrt{2k}} \Rightarrow \mathbb{E}[f(\bar{\omega}_k)] \leq f(\omega^*) + \frac{\sqrt{2} \cdot \sigma \|\omega_0 - \omega^*\|_2}{\sqrt{k}}$$

~ in practice, step-size choice can be more nuanced
depending on the specific application.

... or use $\eta_t = \frac{1}{\sqrt{t}}$ (step-size decreases over time)

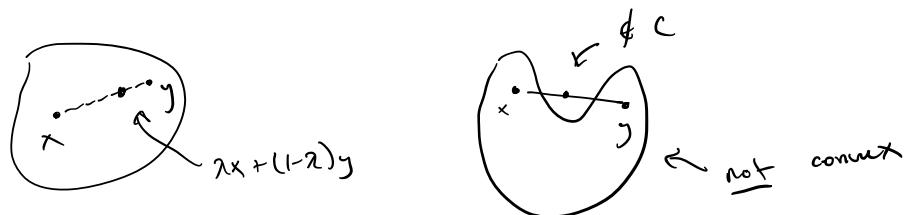
Constrained Optimization: So far we have only considered optimizing over all of \mathbb{R}^d .

$$\text{i.e. } \boxed{\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} L(w)}$$

In many situations we want w to satisfy some constraints.
i.e. want w to lie in some bounded region of \mathbb{R}^d .

It turns out we can adapt GD / SGD / NAGD to this situation
as long as the constraint region is a convex set.

Defn: $C \subseteq \mathbb{R}^d$ is convex set if
 $\forall x, y \in C, \frac{x+y}{2} \in C$ (midpoint also in C)
equiv: $\forall x, y \in C, \forall \lambda \in [0, 1], \lambda x + (1-\lambda)y \in C$



Goal: compute $\boxed{\underset{w \in C}{\operatorname{argmin}} L(w)}$

Projected Gradient Descent:

Projection: $\text{Proj}_C(y) = \underset{x \in C}{\operatorname{argmin}} \|x - y\|_2$ ($\begin{matrix} \text{closest point in} \\ C \text{ to } y \end{matrix}$)

PGD update : $w_{k+1} = \text{Proj}_C(w_k - \gamma \nabla f(w_k))$
 \uparrow
 or $G(w_k)$ in case of SGD

Remark : Everything known for GDO/SGO/NAGD

extends to PGD as long as

C is a convex set.

