

Convex Optimization for Machine Learning

1. Optimization problems in ML

Minimize expected risk:

$$\min_{w \in \mathbb{R}^d} \left\{ R(w) := \mathbb{E}_{z \sim \mathcal{D}} [\ell(w, z)] \right\}$$

\mathcal{D} : data distribution over z (ex: $z = (x, y)$)

w : parameters

ℓ : loss function

e.g.:

- linear models

- $\ell(w, z) = \ell(w, (x, y)) = \tilde{\ell}(y, \langle w, \phi(x) \rangle)$

- $\tilde{\ell}$: square loss: $\ell(w, z) = \frac{1}{2} (y - \langle w, \phi(x) \rangle)^2$
logistic/hinge loss

- NN $w = (W_1, \dots, W_L)$, $\ell(w, z)$ is non-convex in w .

Empirical risk minimization (ERM)

Samples $z_i \sim \mathcal{D}$ i.i.d. $i=1, \dots, m$

$$\hat{R}(w) = \frac{1}{m} \sum_{i=1}^m \ell(w, z_i)$$

(Regularized) ERM:

$$\min_w \hat{R}(w) + \frac{\lambda}{2} \|w\|^2$$

$$\text{or} \quad \min_{\|w\| \leq B} \hat{R}(w)$$

Empirical vs Expected Risk:

$$R(w) = \underbrace{R(w) - \hat{R}(w)}_{\text{estimation}} + \underbrace{\hat{R}(w)}_{\text{optimization}}$$

\Rightarrow No need to optimize \hat{R} below the estimation error! (Bottou & Bousquet '08)

Ex: linear model, $\|\phi(x)\| \leq R$
loss function G -Lipschitz

$$\mathbb{E} \left[\sup_{\|w\| \leq B} |R(w) - \hat{R}(w)| \right] \leq \frac{G \cdot R \cdot B}{\sqrt{m}}$$

In some cases, "faster rates" are possible $\sim \frac{1}{m}$

2. Convexity, smoothness, gradient descent

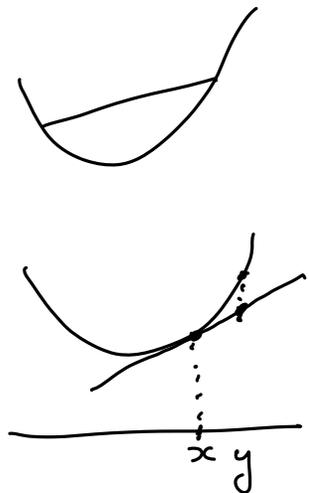
Refs: Nesterov '18 "Lectures on convex optimization"
Bubeck '15 "Conv. Opt.: algorithms & complexity"

Def: (convexity)

f is convex if for any x, y :

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$$

linear approximation at x

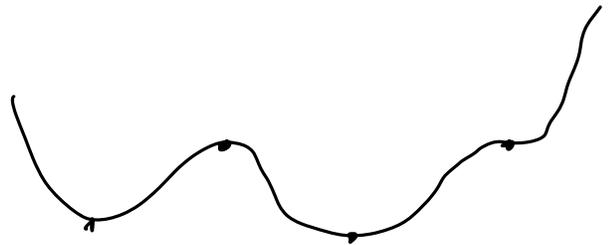


Why convexity?

→ local info \Rightarrow global info

[Fact: x is a global minimum
 \Leftrightarrow
 $\nabla f(x) = 0$

→ not true for non-convex

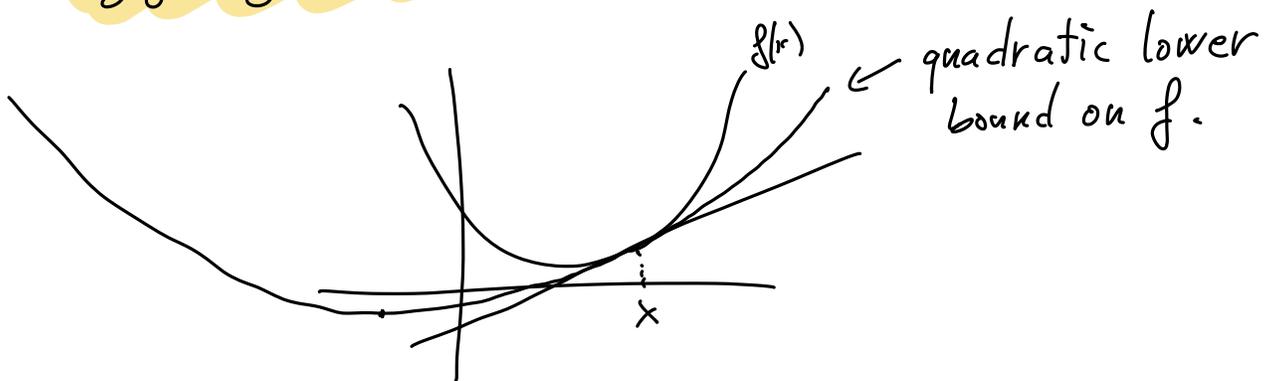


→ easy to bound $f(x) - f(x^*)$ if x^* is global min
suboptimality gap.

Strong convexity

[Def: f is μ -strongly convex if for all x, y .

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2$$



Ex: $\|Ax - b\|^2 = x^T \underbrace{A^T A}_{\text{need } \lambda_{\min}(A^T A) \geq \mu} x$

In the case of least squares, $\lambda_{\min}(\Sigma) \geq \mu$

where $\Sigma = \mathbb{E}[\phi(x)\phi(x)^T]$

- l^2 regularization $\frac{\mu}{2} \|w\|^2 \Rightarrow \mu$ -strong conv.

Fact: Equivalent definition: (see Nesterov '15)

$$f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{1}{2\mu} \|\nabla f(x) - \nabla f(y)\|^2$$

In particular, if $x = x^*$

$$\|\nabla f(y)\|^2 \geq 2\mu (f(y) - f(x^*)) \quad (\text{Łojasiewicz inequality})$$

(small gradient \Rightarrow small suboptimality)

Fact: Convergence of gradient flow on μ -strongly conv f :

$$\frac{d}{dt} x_t = -\nabla f(x_t) \quad x_0$$

$$(x_t = x_{t-1} - \gamma \nabla f(x_{t-1}))$$

We have:

$$f(x_t) - f(x^*) \leq \exp(-2\mu t) (f(x_0) - f(x^*))$$

proof:

$$\frac{d}{dt} (f(x_t) - f(x^*)) = \frac{d}{dt} (f(x_t))$$

$$= \langle \nabla f(x_t), \frac{d}{dt} x_t \rangle$$

$$= -\|\nabla f(x_t)\|^2$$

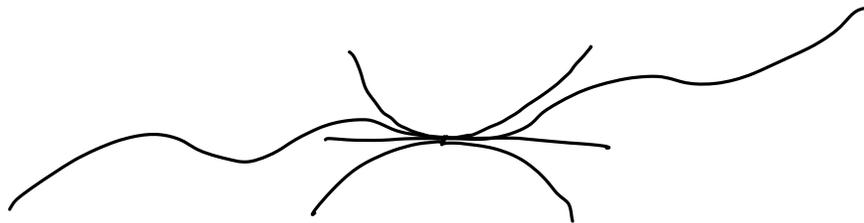
$$\leq -2\mu (f(x_t) - f(x^*))$$

(Gronwall's Lemma) integrate this leads to the result \square

Remark: What about discrete time?
 \Rightarrow need smoothness!

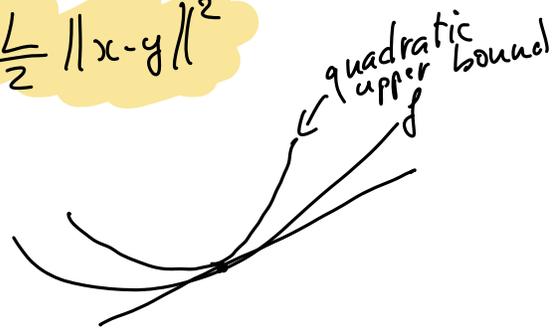
Smoothness

Def: f is L -smooth if ∇f is L -Lipschitz, i.e.
 $\|\nabla f(x) - \nabla f(y)\| \leq L \cdot \|x - y\|$ for any x, y .



Fact: If f is L -smooth, we have

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|x - y\|^2$$



proof: exercise (fund. thm. of calculus)

Ex: • $\|Ax - b\|^2 \Rightarrow L = \lambda_{\max}(A^T A)$

in least squares, $L = \lambda_{\max}(\epsilon)$

• If loss is G -Lipschitz $\ell(\hat{y}, y)$
 $\|\phi(x)\| \in \mathbb{R}$ " "
 $\langle w, \phi(x) \rangle$
 $L \leq G \cdot R^2$

Thm: Convergence of G.D. for
 L -smooth, μ -strongly convex f .

$$x_t = x_{t-1} - \gamma \nabla f(x_{t-1})$$

For $\gamma = \frac{1}{L}$:

$$f(x_t) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^t (f(x_0) - f(x^*))$$

Proof:

$$\begin{aligned} f(x_t) &\leq f(x_{t-1}) + \langle \nabla f(x_{t-1}), x_t - x_{t-1} \rangle + \frac{L}{2} \|x_t - x_{t-1}\|^2 \\ &= f(x_{t-1}) - \frac{1}{L} \|\nabla f(x_{t-1})\|^2 + \frac{1}{2L} \|\nabla f(x_{t-1})\|^2 \\ &= f(x_{t-1}) - \frac{1}{2L} \|\nabla f(x_{t-1})\|^2 \end{aligned}$$

$$\begin{aligned} f(x_t) - f(x^*) &= f(x_{t-1}) - f(x^*) - \frac{1}{2L} \|\nabla f(x_{t-1})\|^2 \\ &\leq f(x_{t-1}) - f(x^*) - \frac{\mu}{L} (f(x_{t-1}) - f(x^*)) \\ &= \left(1 - \frac{\mu}{L}\right) (f(x_{t-1}) - f(x^*)) \leq \dots \leq \left(1 - \frac{\mu}{L}\right)^t (f(x_0) - f(x^*)) \quad \square \end{aligned}$$

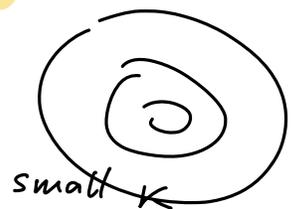
Remark: • $\left(1 - \frac{\mu}{L}\right)^t \leq e^{-\frac{\mu}{L}t}$ \rightarrow exponential convergence
 linear

• $\kappa = \frac{L}{\mu} \geq 1$ is the condition number

small $\kappa \Rightarrow$ fast convergence

$$\left(\kappa = \frac{\lambda_{\max}(\text{Hessian})}{\lambda_{\min}(\text{Hessian})} \right)$$

for quadratics



- In practice, κ may often be very large -

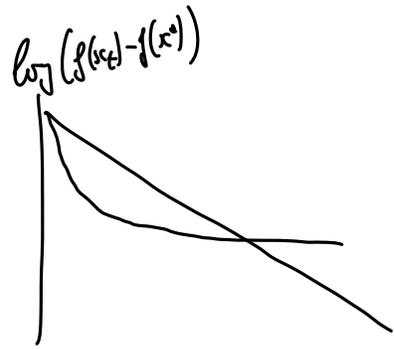
often $\mu \hat{=} \frac{1}{\sigma_m}$ or $\frac{1}{m}$

\Rightarrow can we obtain convergence rates independent of κ ?

Thm: Convergence of G.D. under L -smoothness and convexity

With $\gamma = \frac{1}{L}$, we have

$$f(x_t) - f(x^*) \leq \frac{L}{2t} \|x_0 - x^*\|^2$$



proof .. (Bansal & Gupta '2019)

control: $V_t(x_t) = t(f(x_t) - f(x^*)) + \frac{L}{2} \|x_t - x^*\|^2$

Other methods:

- Nesterov's acceleration can achieve faster convergence rates

$\rightarrow \exp(-\frac{t}{\sqrt{\kappa}})$ instead of $\exp(-\frac{t}{\kappa})$ (str. conv)

or $\frac{1}{t^2}$ instead of $\frac{1}{t}$ (conv)

This is optimal (cannot do better under these assumptions)

for "first-order" methods

(i.e. use only $\nabla f(x_t)$ at each t)

- **Newton's method** can converge much faster by computing Hessians

$$x_t = x_{t-1} - \gamma \text{Hess}(x_{t-1})^{-1} \nabla f(x_{t-1})$$

pros: break dependence on cond. number κ

cons: more costly (invert $d \times d$ matrix)

(Boyd & Vandenberghe book)

- **Proximal methods**

$$f(w) = \hat{R}(w) + \lambda \|w\|_1 \quad \checkmark$$

- assume access to prox-oracle

$$\text{prox}_{\|\cdot\|_1}(x) = \underset{z}{\text{argmin}} \|x-z\|^2 + \lambda \|z\|_1$$

- preserve similar convergence rates despite non-smooth term -

3. Stochastic Gradient Descent

Until now: optimize $f(w) = \hat{R}(w) \left[+ \frac{\lambda}{2} \|w\|^2 \right]$

↑
treated as black-box, access gradients $\nabla f(w)$

- pro: fast convergence rates (linear when $\kappa < \infty$)
- cons: computing a single gradient $\nabla f(w)$

requires an entire pass over data (EXPENSIVE!)

Q: Can we bypass optimizing $\hat{R}(w)$ and directly optimize $R(w)$?

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

optimization of $R(w) \leftrightarrow$ generalization

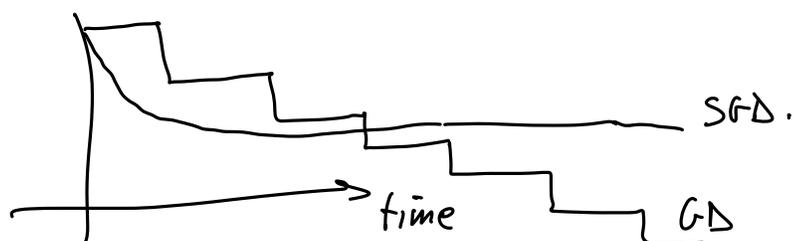
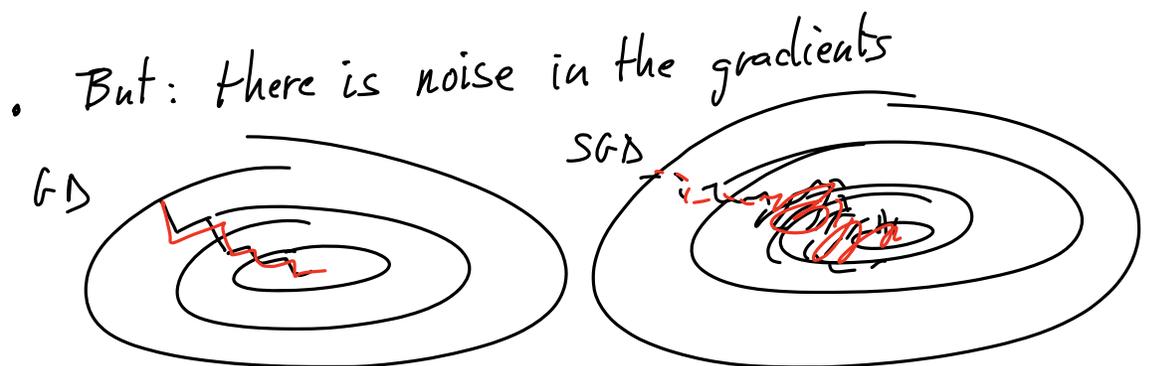
convergence rate
on $R(\hat{w}) - R(w^*) \leftrightarrow$ generalization bound

SGD: GD on $R(w)$, but replace $\nabla R(w)$ by stochastic gradients $\nabla_w l(w, z)$

$$\nabla R(w) = \nabla \mathbb{E}_z [l(w, z)] = \mathbb{E}_z [\nabla l(w, z)]$$

Remarks:

- each step is much cheaper (gradient on a single datapoint)



Theorem: Consider SGD updates:

$$w_t = w_{t-1} - \gamma \nabla \ell(w_{t-1}, z_t) \quad \text{with } z_t \sim \mathcal{D} \\ \text{(fresh sample)}$$

Assume:

$$\|\nabla \ell(w, z)\| \leq G \cdot R$$

($\tilde{\ell}$ is G -Lipschitz)

$$\|\phi(z)\| \leq R$$

$\Rightarrow \ell(\cdot, z)$ is GR -Lipschitz

$$\|w^*\| \leq \mathcal{B}$$

\uparrow

$$w^* \in \operatorname{argmin} R(w)$$

Then, with $\gamma = \frac{\mathcal{B}}{GR\sqrt{m}}$, SGD satisfies

$$\mathbb{E} [f(\bar{w}_m) - f(w^*)] \leq \frac{\mathcal{B}GR}{\sqrt{m}} \quad (f(w) = \mathbb{E}_z [\ell(w, z)])$$

$$\text{where } \bar{w}_m = \frac{1}{m} \sum_{t=0}^{m-1} w_t \quad (= \frac{m-1}{m} \bar{w}_{m-1} + \frac{1}{m} w_{m-1})$$

proof: We have:

$$\|w_t - w^*\|^2 = \|w_{t-1} - \gamma \nabla \ell(w_{t-1}, z_t) - w^*\|^2$$

$$= \|w_{t-1} - w^*\|^2 - 2 \langle w_{t-1} - w^*, \gamma \nabla \ell(w_{t-1}, z_t) \rangle + \underbrace{\|\gamma \nabla \ell(w_{t-1}, z_t)\|^2}_{\leq \gamma^2 G^2 R^2}$$

We have

$$\mathbb{E} [\|w_t - w^*\|^2 | w_{t-1}] \leq \|w_{t-1} - w^*\|^2 - 2\gamma \mathbb{E} [\langle w_{t-1} - w^*, \nabla \ell(w_{t-1}, z_t) \rangle | w_{t-1}] + \gamma^2 G^2 R^2$$

$$= \|w_{t-1} - w^*\|^2 - 2\gamma \underbrace{\langle w_{t-1} - w^*, \underbrace{\mathbb{E}_{z_t} [\nabla \ell(w_{t-1}, z_t)]}_{\nabla f(w_{t-1})} \rangle}_{\nabla f(w_{t-1})} + \gamma^2 G^2 R^2$$

use convexity

$$f(w^*) \geq f(w_{t-1}) + \langle \nabla f(w_{t-1}), w^* - w_{t-1} \rangle$$

$$- \langle \nabla f(w_{t-1}), w_{t-1} - w^* \rangle \leq - (f(w_{t-1}) - f(w^*))$$

This yields:

$$\mathbb{E}[\|w_t - w^*\|^2 | w_{t-1}] \leq \|w_{t-1} - w^*\|^2 - 2\gamma (f(w_{t-1}) - f(w^*)) + \gamma^2 G^2 R^2$$

Take expectation w.r.t. w_{t-1}

$$\mathbb{E}[\|w_t - w^*\|^2] \leq \mathbb{E}[\|w_{t-1} - w^*\|^2] - 2\gamma \mathbb{E}[f(w_{t-1}) - f(w^*)] + \gamma^2 G^2 R^2$$

Re-arrange:

$$2\gamma \mathbb{E}[f(w_{t-1}) - f(w^*)] \leq \mathbb{E}[\|w_{t-1} - w^*\|^2] - \mathbb{E}[\|w_t - w^*\|^2] + \gamma^2 G^2 R^2$$

Summing from $t=1$ to $t=n$,

$$\begin{aligned} 2\gamma \mathbb{E}\left[\sum_{t=1}^n f(w_{t-1}) - f(w^*)\right] &\leq \mathbb{E}[\|w_0 - w^*\|^2] - \mathbb{E}[\|w_n - w^*\|^2] + n \cdot \gamma^2 G^2 R^2 \\ &\leq B^2 + n \gamma^2 G^2 R^2 \end{aligned}$$

$$\frac{1}{n} \mathbb{E}\left[\sum_{t=1}^n f(w_{t-1}) - f(w^*)\right] \leq \frac{B^2}{2\gamma n} + \frac{\gamma G^2 R^2}{2} = \frac{BGR}{\gamma n}$$

$$\mathbb{E} f(\bar{w}_n) - f(w^*) \leq \left(\frac{1}{n} \mathbb{E} \sum_{t=0}^{n-1} f(w_t) \right) - f(w^*) \leq \frac{BGR}{\gamma n}$$

Jensen's ineq.

Remarks: • matches estim. error w/ Rademacher compl. !

⇒ one-pass SGD can suffice
(not always! e.g. fast rates)

- step-size depends on n , can use $\frac{1}{\sqrt{t}}$ decreasing step-size instead
- Smoothness does not help, but it can help with mini-batches
- strong convexity can improve rate to $\frac{1}{\mu n}$ instead of $\frac{1}{\sqrt{n}}$
- Can also use SGD to optimize

$$(*) \quad \hat{R}(w) = \frac{1}{m} \sum_{i=1}^m \ell(w, z_i) = \mathbb{E}_{i_t} [\ell(w, z_{i_t})]$$

by sampling $i_t \sim \text{Unif}(\{1, \dots, m\})$

→ but: no more generalization guarantees!

- For finite-sum problems (*), there are faster algorithms! "Variance-reduction"
(SAG, SAGA, SVRG, SDCA, MRSD, ...)

similar rate as GD, at a fraction of cost (similar to SGD)

$$\text{mini-batch : gradient: } g_t = \frac{1}{b} \sum_{i=1}^b \nabla \ell(w_{t-1}, z_t^i) \quad \begin{matrix} z_t^i \sim \mathcal{D} \\ i=1, \dots, b \end{matrix}$$

$$w_t = w_{t-1} - \gamma g_t \quad \text{instead of } g_t = \nabla \ell(w_{t-1}, z_t)$$

replace $\|\nabla \ell(w, z)\| \leq GR$

by $\mathbb{E} \|g_t - \mathbb{E} g_t\|^2 \leq \sigma^2$ condition

b times smaller with mini-batch.
 $\sigma^2 = \frac{\sigma_0^2}{b}$

different bound:

$$\mathbb{E} [f(\bar{w}_T) - f(w^*)] \leq \frac{L \|w^*\|^2}{T} + \frac{\sigma_0 \|w^*\|}{\sqrt{T \cdot b}} \approx \frac{\sigma_0 \|w^*\|}{\sqrt{T \cdot m}}$$