

DSC 291: Stochastic Optimization

What this course is about:

Most tasks in supervised learning can be modeled as

$$\star \min_{\omega} \mathbb{E}_{z \sim P} l(\omega, z) + \lambda \cdot r(\omega)$$

Here

- z is usually a feature/label pair $z = (x, y) \in \mathbb{R}^d \times \mathbb{R}$
- $l(\cdot, z)$ is a random loss function

e.g.: $\hookrightarrow l(\omega, (x, y)) = (f_\omega(x) - y)^2$ MSE

$\hookrightarrow l(\omega, (x, y)) = \log(1 + e^{-y \cdot w^T \varphi(x)})$ Logistic

- $r: \mathbb{R}^d \rightarrow \{\mathbb{R}, V\} \cup \{\infty\}$ is a structure

inducing regularizer

e.g.: $\hookrightarrow r(\omega) = \|\omega\|_2^2$

$\hookrightarrow r(\omega) = \|\omega\|_1$

Ridge
 ℓ_1 -penalty

Assumption: only access to P is through sampling $z \sim P$.
Two approaches:

Strategy 1: Empirical Risk Minimization

$$\omega_n = \arg \min_{\omega} \frac{1}{n} \sum_{i=1}^n l(\omega, z_i) + \lambda \cdot r(\omega)$$

where $z_1, \dots, z_n \stackrel{iid}{\sim} P$.

Strategy 2: Stochastic Approximation

Algorithm that produces a sequence of points $\omega_1, \dots, \omega_t$, which in each iteration t draws a fresh batch

$$S_t = \{z_{t,1}, z_{t,2}, \dots, z_{t,b}\} \stackrel{iid}{\sim} P$$

and uses $\{f(\cdot, z)\}_{z \in S}$ to update S .

Main Example: **Stochastic Gradient Method**

$$\left\{ \begin{array}{l} \text{Draw batch } S_t \\ \omega_{t+1} = \omega_t - \frac{\eta}{b} \sum_{z \in S} \nabla l(\omega_t, z) \end{array} \right\}$$

When $\eta = 0$, $l(\cdot, z)$ is smooth.

Goal: Solve using as few samples as possible.

Background:

1.1 Inner products and linear maps.

- \mathbb{E} is an Euclidean space: finite dimensional real vector space with an inner-product $\langle \cdot, \cdot \rangle : \mathbb{E} \times \mathbb{E} \rightarrow \mathbb{R}$ satisfying
 - (Symmetry) $\langle x, y \rangle = \langle y, x \rangle$
 - (Bilinearity) $\langle ax + by, z \rangle = a\langle x, z \rangle + b\langle y, z \rangle$ $\forall x, y, z \in \mathbb{E}, a, b \in \mathbb{R}$

(Positive Semidefiniteness)

$$\langle x, x \rangle \geq 0 \text{ and } \langle x, x \rangle = 0 \iff x = 0.$$

Ex: $\mathbb{E} = \mathbb{R}^d$ are arrays of d real values with

$$\langle x, y \rangle = x^T y = \sum_{i=1}^d x_i y_i$$

Ex: $\mathbb{E} = \mathbb{R}^d$ with an inner product

$$\langle x, y \rangle_A = x^T A y = \sum_{i,j} A_{ij} x_i y_j,$$

where $A \in \mathbb{R}^{d \times d}$ is symmetric and positive definite (will come back to this in a second)

Ex: $\mathbb{E} = \mathbb{R}^{m \times n}$
with $\langle \mathbb{X}, \mathbb{Y} \rangle \leftarrow$ trace product

$$\begin{aligned}\langle \mathbb{X}, \mathbb{Y} \rangle &:= \langle \text{vec}(\mathbb{X}), \text{vec}(\mathbb{Y}) \rangle \\ &= \sum_{i,j} \mathbb{X}_{ij} \mathbb{Y}_{ij} \\ &= \text{tr}(\mathbb{X}^T \mathbb{Y})\end{aligned}$$

check
this

1.2 Norms:

A norm on a vector space \mathcal{V} is a function $\|\cdot\|: \mathcal{V} \rightarrow \mathbb{R}$ satisfying

(absolute homogeneity) $\|a\mathbf{x}\| = |a| \cdot \|\mathbf{x}\|$

(triangle inequality) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$

(positivity) $\|\mathbf{x}\| = 0 \iff \mathbf{x} = 0$

$\forall \mathbf{x}, \mathbf{y} \in \mathcal{V} \quad a \in \mathbb{R}.$

Ex: If $\langle \cdot, \cdot \rangle$ is an inner product on V , then $\|x\| := \sqrt{\langle x, x \rangle}$ is a norm
 e.g.: \mathbb{R}^n with dot-product $\rightarrow \|x\|_2 = \sqrt{\langle x, x \rangle} = \sqrt{x^T x} = \sqrt{\sum_{i=1}^n x_i^2}$
 induced norm by inner product.

Ex: $E = \mathbb{R}^{m \times n}$ with trace inner product.

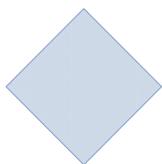
Then $\|\mathbf{X}\|_F := \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle} = \sqrt{\text{tr}(\mathbf{X}^T \mathbf{X})} = \|\text{vec}(\mathbf{X})\|_2 = \sqrt{\sum_{i,j} X_{ij}^2}$

Frobenius norm

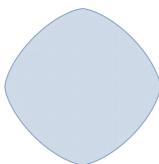
Lemma: (Cauchy-Schwartz)

$$|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$$

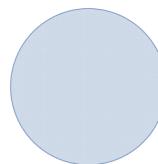
induced norms



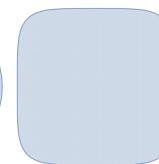
(a) $p = 1$



(b) $p = 1.5$



(c) $p = 2$



(d) $p = 5$



(e) $p = \infty$

Figure 1.1: Unit balls of ℓ_p -norms.

Ex: $\mathcal{V} = \mathbb{R}^d$. Define the ℓ_p -norm:

$$\|x\|_p = \begin{cases} \left(\sum_{i=1}^d |x_i|^p \right)^{1/p}, & \text{if } 1 \leq p < \infty \\ \max_{i=1,\dots,d} |x_i|, & \text{if } p = \infty \end{cases}$$

This is a norm. None of these are induced by any inner-product except $p=2$.

Thm: (von Neumann)

If $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}$ is sign-permutation

invariant norm on \mathbb{R}^d , then

$$\varphi(\Sigma) \triangleq \varphi(\sigma(\Sigma))$$

singular values

is a norm on $\mathbb{R}^{m \times n}$.

Eg:

$$\|\Sigma\|_1 = \|\sigma(\Sigma)\|_1$$

Nuclear norm

$$\|\Sigma\|_p = \|\sigma(\Sigma)\|_\infty$$

Operator norm

$$\|\Sigma\|_F = \sqrt{\text{tr}(\Sigma^\top \Sigma)} = \|\sigma(\Sigma)\|_2$$

Fact: all norms on \mathbb{E} are "equivalent".
For any two norms ρ_1 and ρ_2 on \mathbb{E} , there exists $\alpha, \beta > 0$ s.t.

$$\alpha \rho_1(x) \leq \rho_2(x) \leq \beta \rho_1(x) \quad \forall x \in \mathbb{E}$$

So ρ_1 and ρ_2 are within a constant multiplicative factor of each other.

Key caveat: α, β typically depend on the dimension $d = \dim(\mathbb{E})$

Ex: $\|x\|_2 \leq \|x\|_1 \leq \sqrt{d} \|x\|_2$

$$\|x\|_\infty \leq \|x\|_2 \leq \sqrt{d} \|x\|_\infty$$

$$\|x\|_\infty \leq \|x\|_1 \leq d \|x\|_\infty$$

We often think of d as being huge, so this is not super helpful.

1.3 Differentiability

- E and Y are Euclidean spaces.
 - $U \subseteq E$ be open set.
 - A map $F: U \rightarrow Y$
- is continuous at $x \in U$ if for any $x_i \xrightarrow{U} x$ we have $F(x_i) \rightarrow F(x)$

• F is L -Lipschitz if $\|F(y) - F(x)\| \leq L \|y - x\|$ Lipschitz constant $\forall x, y \in Q$.

If $L \in [0, 1)$ we say that F is a contraction. If $L = 1$,

we call F nonexpansive.

Defn: $f: E \rightarrow \mathbb{R}$ is differentiable at $x \in U$ if there exists $Df(x) \in E$

s.t.

$$\lim_{h \rightarrow 0} \frac{f(x+h) - [f(x) + \langle Df(x), h \rangle]}{\|h\|} = 0$$

We then call $Df(x)$ the gradient.

Defn.: $f: \mathbb{E} \rightarrow \mathbb{R}$ is twice differentiable at $x \in \mathbb{U}$ if there exists linear map $D^2f(x) : \mathbb{E} \times \mathbb{E} \rightarrow \mathbb{E}$ s.t.

$$\lim_{h \rightarrow 0} \frac{\|Df(x+h) - Df(x) - D^2f(x)h\|}{\|h\|} = 0$$

We then call $D^2f(x)$ the Hessian.

Ex: $\mathbb{E} = \mathbb{R}^d$. Then

$$Df(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_d} \end{bmatrix}$$

if f is differentiable

and $D^2f(x) = \left[\frac{\partial^2 f}{\partial x_i \partial x_j} \right]_{i,j}$

if f is twice differentiable.

H.W.: compute $Df(\underline{X})$, $D^2f(\underline{X})$ for
 $f(\underline{X}) = \log \det(\underline{X})$.

Def: f is C^1 -smooth if ∇f exists and is continuous.
 f is C^2 -smooth if $\nabla^2 f$ exists and is continuous.

Lemma (Accuracy in approximation)

If ∇f is L -Lipschitz, then

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

pt. in Homework. $= l_x(y)$ linearization.

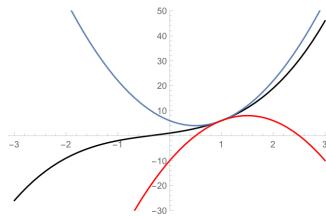


Figure 1.5: The black curve depicts the graph of a β -smooth function f ; the blue and red curves depict graphs of the quadratics $l_x(\cdot) + \frac{\beta}{2} \|\cdot - x\|^2$ and $l_x(\cdot) - \frac{\beta}{2} \|\cdot - x\|^2$, respectively.

Thus Consider C^2 -smooth $f: \mathbb{U} \rightarrow \mathbb{R}$ and $x \in \mathbb{U}$. Then the following are true.

1) (Necessary) If $x \in \mathbb{U}$ is a local m.m. then $\nabla f(x) = 0$ and $\nabla^2 f(x) \succ 0$.

2) (Sufficient) If $\nabla f(x) = 0$ and $\nabla^2 f(x) \succ 0$ then x is a local minimizer of f .

1.4 Convexity

Defn: $Q \subseteq \mathbb{E}$ is convex, if
 $\forall x + (1-\lambda)y \in Q \quad \forall x, y \in E$
 $\lambda \in [0, 1]$

Defn: $f: E \rightarrow \mathbb{R} \cup \{\pm\infty\}$ is convex, if
 $f(\lambda x + (1-\lambda)y) \leq \lambda \cdot f(x) + (1-\lambda)f(y) \quad \forall x, y \in E$
 $\lambda \in [0, 1]$

Thm: The following are equivalent for any C^1 -smooth function $f: U \rightarrow \mathbb{R}$ defined on a convex open set $U \subseteq E$.

① (convexity) f is convex:

② (gradient inequality)

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

If f is C^2 -smooth, then also

③ $\nabla^2 f(x) \succeq 0 \quad \forall x \in U$.

Note: For any C^1 -smooth convex f , TFAE

- $\nabla f(x) = 0$

- x is a global minimizer

Chapter I: (Stochastic) Gradient Descent for linear least squares.

Problem:

$$\text{min}_{x \in \mathbb{R}^d} f(x) = \frac{1}{2n} \|Ax - b\|_2^2 = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (a_i^T x - b_i)^2$$

where $A \in \mathbb{R}^{n \times d}$, $b \in \mathbb{R}^n$, and a_i are the rows of A .

Ex: Learning from the diabetes data set. Goal: predict diabetes disease progression one year after baseline.

Data: $\{(x_i, b_i)\}_{i=1}^n$

- Number of samples: $n = 442$ patients
- Training set (75%): 330
- Validation set (25%): 110
- Number of features: $d = 10$

You will experiment with this dataset in your homework.

$$f(x) = \frac{1}{2n} \|Ax - b\|^2$$

Observe:

$$\nabla f(x) = \frac{1}{n} A^T(Ax - b), \quad \nabla^2 f(x) = \frac{1}{n} A^T A \quad (\text{H.W})$$

Notice: $\frac{1}{n} A^T A \succeq 0$

[Reason
 $v^T A^T A v = \|Av\|_2^2 \geq 0 \forall v$]

So f is convex.

So x is optimal for \star

$$\Rightarrow \frac{1}{n} A^T(Ax - b) = 0$$

$$\Leftrightarrow \boxed{A^T A x = A^T b}$$

normal eqns.

Goal: Develop algorithms for solving normal eqns with cheap per iteration cost, and in particular that do not require forming $A^T A$.

$$\min_{\mathbf{x}} \frac{1}{n} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 =: f(\mathbf{x})$$

Gradient Descent ($\eta > 0$)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t) = \mathbf{x}_t - \frac{\eta}{n} (\mathbf{A}^T (\mathbf{A}\mathbf{x}_t - \mathbf{b}))$$

Let $\bar{\mathbf{x}}$ be a minimizer:

$$\mathbf{x}_{t+1} - \bar{\mathbf{x}} = \mathbf{x}_t - \bar{\mathbf{x}} - \frac{\eta}{n} (\mathbf{A}^T (\mathbf{A}\mathbf{x}_t - \mathbf{b}))$$

Recall $\mathbf{A}^T \mathbf{A} \bar{\mathbf{x}} = \mathbf{A}^T \mathbf{b}$

$$\begin{aligned} &= \mathbf{x}_t - \bar{\mathbf{x}} - \frac{\eta}{n} (\mathbf{A}^T \mathbf{A} \mathbf{x}_t - \mathbf{A}^T \mathbf{A} \bar{\mathbf{x}}) \\ &= \left(\mathbf{I} - \frac{\eta}{n} \mathbf{A}^T \mathbf{A} \right) (\mathbf{x}_t - \bar{\mathbf{x}}) \end{aligned}$$

Summary

$$\boxed{\mathbf{x}_{t+1} - \bar{\mathbf{x}} = \left(\mathbf{I} - \frac{\eta}{n} \mathbf{A}^T \mathbf{A} \right) (\mathbf{x}_t - \bar{\mathbf{x}})}$$

Define $H = \frac{1}{n} A^T A$ and $\Delta_t = x_t - \bar{x}$.

Then we have learned:

$$\Delta_t = (I - \gamma H) \Delta_{t-1} = (I - \gamma H)^t \Delta_0.$$

Let $\alpha = \lambda_{\min}(H)$, $\beta = \lambda_{\max}(H)$, $\kappa = \frac{\beta}{\alpha}$.

Setting $\gamma = \frac{1}{\beta}$ we see

$$\|I - \gamma H\|_{op} = \max_{\lambda \in [\alpha, \beta]} |1 - \gamma \lambda| = 1 - \kappa^{-1}$$

Therefore $\|\Delta_t\|_2 \leq (1 - \kappa^{-1})^t \|\Delta_0\|_2$.

In particular using $1 - a \leq e^{-a}$ for $a \in [0, 1]$,
we may write

$$\|\Delta_t\|_2 \leq \exp\left(-\frac{t}{\kappa}\right) \|\Delta_0\|_2 \leq \epsilon$$

Solving for ϵ yields

$$t \geq \kappa \log\left(\frac{\|\Delta_0\|_2}{\epsilon}\right)$$

Let's now look instead at $f(x_t) - f(\bar{x})$
 Since f is a pure quadratic,

$$f(x) = f(\bar{x}) + \underbrace{\langle \nabla f(\bar{x}), x - \bar{x} \rangle}_{\text{''}} + \frac{1}{2} \underbrace{\langle \nabla^2 f(\bar{x})(x - \bar{x}), x - \bar{x} \rangle}_{\text{''}}$$

Therefore

$$f(x_t) - f^* = \frac{1}{2} \Delta_0^T (I - \gamma H)^{2t} H \Delta_0$$

A very similar argument as before shows

$$f(x_t) - f^* \leq (1 - \gamma^{-1})^{2t} (f(x_0) - f^*)$$

for $\gamma = \frac{1}{B}$. In fact, we can obtain a rate that is insensitive to α :

$$f(x_t) - f^* \leq \frac{1}{2} \max_{\lambda \in [\alpha, B]} \lambda (1 - \frac{\lambda}{B})^{2t} \|\Delta_0\|_2^2$$

Observe

$$\begin{aligned} \lambda (1 - \frac{\lambda}{B})^{2t} &\leq \lambda \exp(-\frac{2t\lambda}{B}) = \frac{B}{2t} \cdot \frac{2t\lambda}{B} \exp(-\frac{2t\lambda}{B}) \\ &\leq \frac{B}{2t} \cdot \max_{s > 0} s e^{-s} \end{aligned}$$

$\frac{B}{2t} \cdot \max_{s > 0} s e^{-s}$

e^{-1}

$$\text{Thus } f(x_t) - f^* \leq \frac{\beta \|\Delta_0\|^2}{8\varepsilon} \leq \varepsilon$$

Solving for $\varepsilon > 0$ gives

$$t \geq \frac{\beta \|\Delta_0\|^2}{8\varepsilon}$$

Summary: Gradient Descent with $\gamma = \frac{1}{\beta}$ enjoys the guarantee:

$$f(x_t) - f^* \leq \min \left\{ (1-\gamma^{-1}) (f(x_0) - f^*), \frac{\beta \|x_0 - \bar{x}\|_2^2}{\varepsilon} \right\}$$

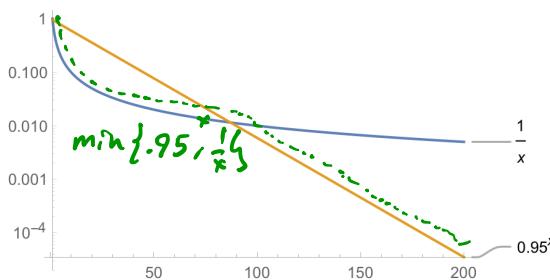


Figure 5.2: Sublinear vs. linear rates

Chapter 2:

Stochastic gradient descent:

Each iteration of GD requires $n \cdot d$ arithmetic operations. SGD has cost $O(d)$ per iteration instead.

$$\min_x f(x) := \mathbb{E}_{(a, b) \sim P} \frac{1}{2} (\langle a, x \rangle - b)^2$$

SGD: step $t = 1, \dots, T$:

$$\begin{cases} \text{Draw } (a_t, b_t) \sim P \\ x_{t+1} = x_t - \gamma_t (\langle a_t, x_t \rangle - b_t) \cdot a_t \end{cases}$$

Define

$$\Sigma := \mathbb{E}[aa^T], \quad v = \mathbb{E}[ba], \quad \bar{x} \in \arg \min f.$$

Note

$$0 = \nabla f(\bar{x}) = \mathbb{E}(\langle a, \bar{x} \rangle - b)a = \mathbb{E}[aa^T\bar{x} - ba]$$

$$\Rightarrow \boxed{\sum \bar{x} = v}$$

Set

$$\Delta_t = x_t - \bar{x}, \quad \varepsilon_t = b_t - \langle a_t, \bar{x} \rangle$$

$$\text{So } \langle a_t, x_t \rangle - b_t = \langle a_t, \Delta_t \rangle - \varepsilon_t$$

Therefore

$$\begin{aligned}\Delta_{t+1} &= \Delta_t - \gamma_t (\langle a_t, \Delta_t \rangle - \varepsilon_t) a_t \\ &= (I - \gamma_t a_t a_t^T) \Delta_t + \gamma_t \varepsilon_t a_t\end{aligned}$$

and as in the deterministic case

$$f(x) - f^* = \frac{1}{2} (x - \bar{x})^T \Sigma (x - \bar{x})$$

Let $E_t[\cdot] := E[\cdot | X_t]$

Assumptions: \nwarrow well-specified

$$\cdot E[\varepsilon | a] = 0, \quad E[\varepsilon^2 | a] \leq \sigma^2,$$

$$\alpha I \leq \Sigma, \quad E[\|a\|^2 a a^T] \leq R^2 \Sigma$$

e.g. true for $a \sim N(0, \Sigma)$

e.g. true for $\|a\| \leq R$

Write

$$\|\Delta_{t+1}\|_2^2 = \underbrace{\|(I - 2\gamma_t a_t a_t^\top) \Delta_t\|_2^2}_{A} - 2\gamma_t \langle \varepsilon_t a_t, (I - 2\gamma_t a_t a_t^\top) \Delta_t \rangle + \gamma_t^2 \|\varepsilon_t a_t\|_2^2$$

$$\begin{aligned} \mathbb{E}_t A &= \Delta_t^\top \mathbb{E} (I - 2\gamma_t a_t a_t^\top)^2 \Delta_t \\ &= \Delta_t^\top \left(I - 2\gamma_t \Sigma + \gamma_t^2 \mathbb{E}[\|a\|_2^2 a a^\top] \right) \Delta_t \\ &\leq \Delta_t^\top (I - 2\gamma_t (2 - \gamma_t R^2) \Sigma) \Delta_t \end{aligned}$$

$$\mathbb{E}_t B = 2\gamma_t \langle \underbrace{\mathbb{E}[\varepsilon_t | a_t]}_{0}, \dots \rangle = 0$$

$$\mathbb{E}_t C = \gamma_t^2 \mathbb{E}_{a_t} [\mathbb{E}[\varepsilon_t^2 | a_t] \cdot \|a_t\|^2] \leq \gamma_t^2 \delta^2 \underbrace{\mathbb{E}[\|a_t\|^2]}_{\text{tr}(\Sigma)}$$

So

$$\mathbb{E}_t \|\Delta_{t+1}\|_2^2 \leq (1 - \alpha \gamma_t (2 - \gamma_t R^2)) \|\Delta_t\|_2^2 + \gamma_t^2 \delta^2 \text{tr}(\Sigma)$$

Take expectation over x_1, \dots, x_t and use tower rule to get:

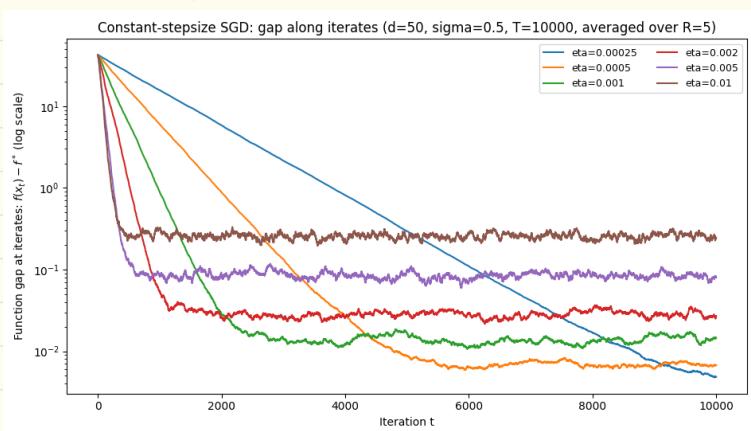
$$\mathbb{E} \|\Delta_{t+1}\|_2^2 \leq (1 - \alpha\gamma_t(2 - \gamma_t R^2)) \mathbb{E} \|\Delta_t\|_2^2$$

$\stackrel{\text{def}}{=} d_{t+1} + \gamma_t^2 \sigma^2 \cdot \text{tr}(\Sigma)$

Conclusion: With constant step size $\gamma \leq \frac{1}{R^2}$, get

$$\begin{aligned} d_{t+1} &\leq (1 - \alpha\gamma) d_t + \gamma^2 \sigma^2 \text{tr}(\Sigma) \\ &\stackrel{\dots}{\leq} (1 - \alpha\gamma)^t d_0 + \gamma^2 \sigma^2 \text{tr}(\Sigma) \cdot \sum_{i=0}^{t-1} (1 - \alpha\gamma)^i \\ &\leq (1 - \alpha\gamma)^t d_0 + \frac{\gamma^2 \sigma^2 \text{tr}(\Sigma)}{\alpha} \end{aligned}$$

initial condition decay exponentially fast noise floor



To get to zero error, there are two strategies

Thm: Set $\gamma_t = \frac{2}{\alpha t + 2R^2}$. Then x_t satisfy

$$d_t \leq \frac{\max \left\{ \alpha^2 \left(1 + \frac{2R^2}{\alpha} \right) d_0, 4\sigma^2 \cdot \text{tr}(\Sigma) \right\}}{\alpha^2 \left(t + \frac{2R^2}{\alpha} \right)}$$

[So roughly $d_t = O\left(\frac{\sigma^2 \cdot \text{tr}(\Sigma)}{\alpha^2 t}\right)$]

pf: From \star :

$$d_{t+1} \leq \left(1 - \frac{2}{t+2R^2/\alpha}\right) d_t + \frac{4\sigma^2 \cdot \text{tr}(\Sigma)/\alpha^2}{(t+2R^2/\alpha)^2}$$

The rest follows from the following lemma, which you will prove for hw.

Lemma: Consider sequence $D_t \geq 0$ and constants $t_0 \geq 0$, $a > 0$ satisfying

$$D_{t+1} \leq \left(1 - \frac{2}{t+t_0}\right) D_t + \frac{a}{(t+t_0)^2}$$

Then $D_t \leq \frac{\max \left\{ (1+t_0)D_0, a \right\}}{t+t_0} t_0^2$.



Second strategy is to run SGD in epochs
 Recall for fixed t and $\eta = \eta_0 = \frac{1}{R^2}$:

$$d_{t+1} \leq (1 - \alpha \eta) d_t + \eta_0 \frac{\Sigma}{\alpha}$$

$$\text{So after } t = \frac{1}{\alpha \eta_0} \log\left(\frac{d_0}{2 \eta_0 \sigma^2 \text{tr}(\Sigma) / \alpha}\right)$$

$$\text{can be sure } d_{t+1} \leq \frac{2 \eta_0 \sigma^2 \text{tr}(\Sigma)}{\alpha}$$

Now replace η_0 by $\eta_1 - \eta_0/2$ and
 repeat: In next epoch, can be
 sure that after

$$t = \frac{1}{\alpha \eta_1} \log\left(\frac{2 \eta_0 \sigma^2 \text{tr}(\Sigma) / \alpha}{2 \eta_1 \sigma^2 \text{tr}(\Sigma) / \alpha}\right)$$

$$= \frac{1}{\alpha \eta_1} \log(2)$$

$$\text{can be sure } d_{t+1} \leq \frac{2 \eta_1 \sigma^2 \text{tr}(\Sigma)}{\alpha}$$

You only need repeat this until

$$2^{-S} \eta_0 = \eta_S \leq \frac{\epsilon \alpha}{2 \sigma^2 \text{tr}(\Sigma)}$$

to get
 $d \leq \epsilon$.

$$\Rightarrow S \leq \log\left(\frac{1}{\epsilon}\right)$$

So total # iterations after careful accounting:

$$\frac{R^2}{d} \log\left(\frac{d_0}{\epsilon}\right) + \frac{25^2 \ln(\Sigma)}{\epsilon d^2}$$

+ obtain a point that is ϵ -close to \hat{x} .

Remarkable Fact: the dependence on d can be removed!!! when tracking the function values $f(\hat{x}_t) - f^*$ along the average iterate:

$$\hat{x}_t := \frac{1}{T+1} \sum_{i=0}^T x_i$$

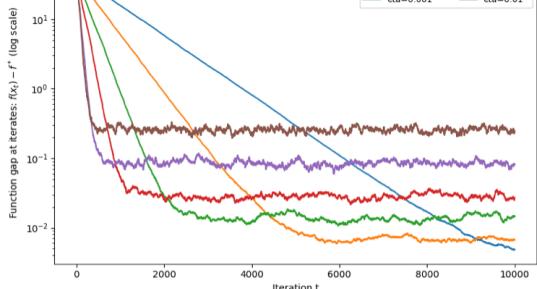
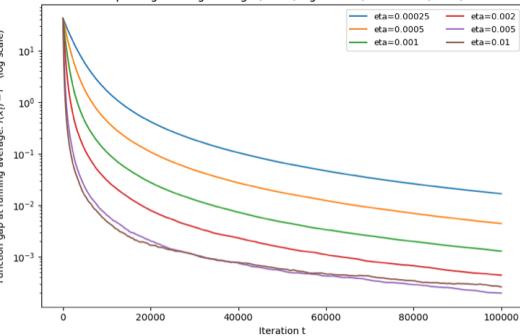
Then: Set $\alpha = \lambda_{\min}(\Sigma)$ and $\gamma = \frac{1}{\beta}$. constant!

Suppose $E[\varepsilon | a] = 0$. Then it holds:

$$E(f(\hat{x}) - f^*) \leq \frac{\beta^2 \cdot \Delta_0^2 \sum_{i=1}^{T-1} \Delta_0}{2T^2} + \frac{5^2 d}{2T} + O\left(\frac{\beta d_0}{T^2} + \frac{1}{\beta T}\right)$$

↳ See Bach-Montlives 2013

low-order

Gap along iterates ($d=50$, $\sigma=0.5$, $T=10000$, $R=5$)Gap along running average ($d=50$, $\sigma=0.5$, $S=100000$, $R=5$)

Pf sketch: The proof is quite technical and we will only sketch it. We have seen that Δ_t evolve according to

$$\Delta_{t+1} = (I - \gamma a_t a_t^T) \Delta_t + \gamma \epsilon_t a_t$$

$$= (I - \gamma \Sigma) \Delta_t + \gamma (\epsilon_t a_t - (a_t a_t^T \Sigma) \Delta_t)$$

Observe $\mathbb{E}_t Q^2 = O(\|\Delta_t\|_2^2) = O(\gamma)$

so Q^2 should be dominated by $\epsilon_t a_t$. Making this precise requires work. But taking it for granted let's analyze the evolution

$$\Delta_{t+1} = (I - \gamma \Sigma) \Delta_t + \gamma q_t \quad \text{with } q_t = \epsilon_t a_t$$

Let's expand the recursion:

$$\Delta_t = (I - \gamma \Sigma)^t \Delta_0 + \gamma \sum_{i=0}^t (I - \gamma \Sigma)^i q_{t-i}$$

$$\mathbb{E} \| \cdot \|_2^2 \approx \underline{\gamma^{-2} \text{tr}(\Sigma)}$$

(same computation
as before)

Now let's look at averaging:

$$\hat{x}_t - \bar{x} = \frac{1}{t+1} \sum_{i=0}^t (x_i - \bar{x}) = \frac{1}{t+1} \sum_{i=0}^t \Delta_i$$

Therefore

$$\begin{aligned} \hat{x}_{t-1} - \bar{x} &= \frac{1}{t} \sum_{i=0}^{t-1} \Delta_i \\ &= \underbrace{\frac{1}{t} \sum_{j=0}^{t-1} (I - \gamma \Sigma)^j \Delta_0}_{\text{bias term}} + \underbrace{\frac{1}{t} \sum_{j=0}^{t-1} \sum_{i=0}^j (I - \gamma \Sigma)^i q_{j-i}}_{\text{noise term}} \end{aligned}$$

Observe

$$\sum_{j=0}^{t-1} (I - \gamma \Sigma)^j = (\gamma \Sigma)^{-1} \cdot (I - (I - \gamma \Sigma)^t)$$

\uparrow
geometric series.

$$S_0 \quad B = \frac{1}{t} (\gamma \Sigma)^{-1} \cdot (I - (I - \gamma \Sigma)^t) \Delta.$$

For the noise term, swap seems:

$$\begin{aligned} N &= \frac{\gamma}{t} \sum_{j=0}^{t-1} \sum_{i=0}^j (I - \gamma \Sigma)^i q_{j-i} \\ &= \frac{\gamma}{t} \sum_{k=0}^{t-1} \left(\underbrace{\sum_{s=0}^{t-1-k} (I - \gamma \Sigma)^s}_{\text{''}} \right) q_k \\ &\quad \text{''} \\ &\quad (\gamma \Sigma)^{-1} (I - (I - \gamma \Sigma)^{t-k}) \end{aligned}$$

Now observe the algebraic identity

$$\begin{aligned} f(x) - f^* &= \mathbb{E} (\langle a, x \rangle - b)^2 - f^* \\ &= \mathbb{E} (\langle a, x \rangle - \langle a, \bar{x} \rangle - \varepsilon)^2 - \mathbb{E} \varepsilon^2 \\ &= \mathbb{E} \langle a, x - \bar{x} \rangle^2 \\ &= \mathbb{E} \operatorname{tr}(a a^T, (x - \bar{x})(x - \bar{x})^T) \\ &= \underbrace{(x - \bar{x})^T \sum (x - \bar{x})}_{\text{''}} \\ &\quad \text{''} \\ &\quad \|x - \bar{x}\|_2^2 \end{aligned}$$

Σ

$$f(\hat{x}_t) - f^* = \|\hat{x}_t - \bar{x}\|_{\Sigma}^2 = \|B\bar{x} + N\|_{\Sigma}^2$$

$$\leq 2\|B\|_{\Sigma}^2 + 2\|N\|_{\Sigma}^2$$

Now

$$\|B\|_{\Sigma}^2 = \frac{1}{\gamma^2 T^2} \Delta_0^T \underbrace{\left(I - (I - \gamma \Sigma)^T \right) \Sigma^{-1} \left(I - (I - \gamma \Sigma) \right) \Delta_0}_{\approx \Sigma^{-1}}$$

$$\leq \frac{\Delta_0^T \Sigma^{-1} \Delta_0}{\gamma^2 T^2}$$

Next define $A_{t,k} = I - (I - \gamma \Sigma)^{t-k}$

Then

$$\mathbb{E} \|N\|_{\Sigma}^2 = \mathbb{E} N^T \Sigma N$$

$$= \frac{\gamma^2}{T^2} \sum_{k=0}^{T-1} \mathbb{E} \left[q_k^T A_{t,k} \underbrace{\Sigma^{-1} A_{t,k}^T}_{\approx \Sigma^{-1}} q_k \right]$$

$$= \frac{\gamma^2}{T^2} \sum_{k=0}^{T-1} \underbrace{\text{tr}(\Sigma^{-1} \mathbb{E} q_k q_k^T)}_{\leq \sigma^2 \Sigma} + \left(\Sigma^{-1} \mathbb{E} q_k q_k^T \right) \leq \frac{\gamma^2 \sigma^2 d}{T}$$

The role of mini-batches. In practice, you run mini-batch SGD:

$$\left\{ \begin{array}{l} \text{Sample } S_t = \{(x_1, b_1), (x_2, b_2), \dots, (x_s, b_s)\} \text{ if} \\ x_{t+1} = x_t - \frac{\eta_t}{s} \sum_{(x, y) \in S_t} (a^T x - b) a; \\ g_t \end{array} \right.$$

Advantage:

- For reasonably sized s , the gradients for each data point in the batch can be computed in parallel.
- Does it help? Yes! Let's look at the last iterate convergence

$$\begin{aligned} \Delta_{t+1} &= \Delta_t - \eta_t g_t = \Delta_t - \eta_t \frac{1}{s} \sum_{(x, y)} (a^T x - b) a \\ &= \underbrace{\left(I - \eta_t \frac{1}{s} \sum_a a a^T \right)}_{E = \Sigma} \Delta_t + \eta_t \frac{1}{s} \sum_a \epsilon a \end{aligned}$$

same
as here

What is

$$\mathbb{E} \left\| \frac{1}{S} \sum_{i=1}^S \epsilon_i \alpha_i \right\|^2 = \frac{1}{S^2} \sum_{i=1}^S \mathbb{E} [\epsilon_i^2 \|\alpha_i\|_2^2] \\ \xrightarrow{\text{analogue of } \sigma^2 \text{tr}(\Sigma)} = \frac{1}{S^2} \sum_{i=1}^S \sigma^2 \text{tr}(\Sigma_i)$$

$$= \frac{1}{S} \cdot \sigma^2 \cdot \text{tr}(\Sigma_i)$$

All the results we have seen extend verbatim with σ^2 replaced by $\frac{\sigma^2}{S}$.

For example, for the last iterate with $y_t = y$, we get

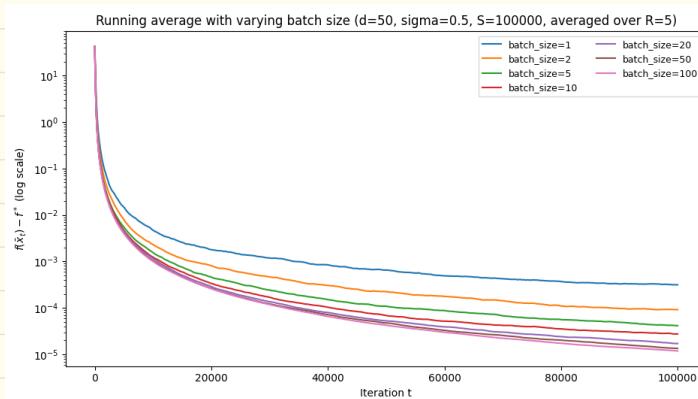
$$\|\Delta_t\|_2^2 \leq (1 - 2\alpha)^t \|\Delta_0\|_2 + \frac{\eta \sigma^2}{\alpha S} \text{tr}(\Sigma)$$

In particular, if you do step-decay, you need ~~bias~~ steps to ensure $\mathbb{E} \|x_t - \bar{x}\|^2 \leq \epsilon$.

$$t \geq \frac{R^2}{\alpha} \log\left(\frac{\Delta_0}{\epsilon}\right) + \frac{\sigma^2 \text{tr}(\Sigma)}{\frac{\sigma^2}{S} \alpha^2 \epsilon}$$

Batch saturation:

For sufficiently large s and finite time horizon T , the noise term is smaller than the bias so no improvement is seen with larger batches.

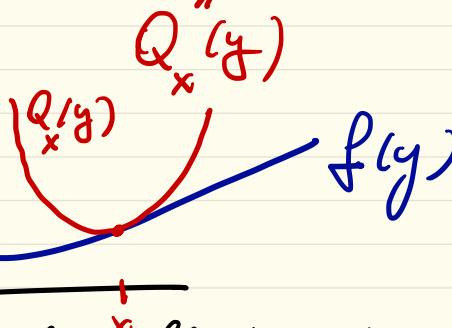


Chapter 2: General Convex Optimization

$$\min f(x)$$

where $f: \mathbb{E} \rightarrow \mathbb{R}$ is β -smooth:

$$f(y) \leq f(x) + \underbrace{\langle \nabla f(x), y-x \rangle}_{Q_x^*(y)} + \frac{\beta}{2} \|y-x\|^2$$



Remark: follows if $\|\nabla f(x) - \nabla f(y)\| \leq \beta \cdot \|x-y\|$

Gradient Descent:

$$x_{t+1} = x_t - \gamma \nabla f(x_t) = \arg \min_x f(x) - \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2} \|x - x_t\|^2$$

Lemma (Descent)

$$f(x_t) - f(x_{t+1}) \geq \gamma \left(1 - \frac{\gamma \beta}{2}\right) \|\nabla f(x_t)\|^2$$

p.t.: Plug in x_{t+1} into

Since $y \mapsto y(1 - \frac{\gamma\beta}{2})$ is maximized at $y = \frac{1}{\beta}$
we deduce that $x^* = x - \frac{1}{\beta} \nabla f(x)$ satisfies

$$f(x) - f(x^*) \geq \frac{1}{2\beta} \|\nabla f(x)\|^2$$

Then:

$$\min_{i=1,\dots,t} \|\nabla f(x_i)\|^2 \leq \frac{1}{t} \sum_{i=1}^t \|\nabla f(x_i)\|^2 \leq \frac{2\beta(f(x_0) - f^*)}{t}$$

Let's make stronger assumptions.

Then: Let $f: \mathbb{E} \rightarrow \mathbb{R}$ be convex and β -smooth.
Then with $\gamma = \frac{1}{\beta}$, GP iterates satisfy:

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

Pf: Letting $Q_x(y) = f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|^2$ get
 $f(x_{t+1}) \leq Q_{x_t}(x_{t+1})$
 $= f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$
 $\stackrel{cvx}{=} f(x_t) + \langle \nabla f(x_t), \bar{x} - x_t \rangle$
 $f(x_t) \stackrel{?}{\geq} - \langle \nabla f(x_t), \bar{x} - x_{t+1} \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$
 $\leq f^* + \frac{\beta}{2} (\|x_t - \bar{x}\|^2 - \|x_{t+1} - \bar{x}\|^2)$

Sum up these inequalities for $i=1 \dots t$.
 Right-side telescopes:

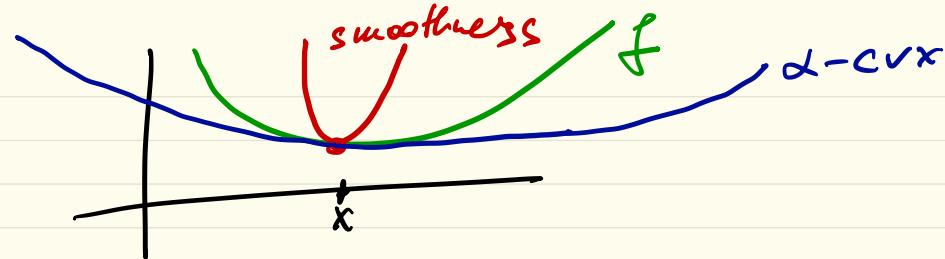
$$\frac{1}{t} \sum_{i=0}^{t-1} (f(x_i) - f^*) \leq \frac{\beta \|x_0 - \bar{x}\|^2}{t-1}$$

Defn: f is called α -strongly cvx
 if $x \mapsto f(x) - \frac{\alpha}{2} \|x\|^2$ is still convex.

Lemma: Suppose f is C^1 -smooth. Then
 f is α -strongly CVX iff

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

$\forall x, y$



Lemma: Suppose f is C^2 -smooth. Then f is α -cvx iff $D^2f(x) \geq \alpha \cdot I \quad \forall x$

For a β -smooth and α -cvx f define

Condition # to be $\lambda = \frac{\beta}{\alpha}$

Lemma: (Polyak-Łojasiewicz) not needed with a diff proof
Let f be C^1 -smooth and α -cvx.

Then

$$\frac{1}{2\alpha} \|Df(x)\|^2 \geq f(x) - f^* \quad \forall x$$

pf: Define $Q_x(y) = f(x) + \langle Df(x), y-x \rangle + \frac{\alpha}{2} \|y-x\|^2$

Then $f(\bar{x}) \geq Q_{\bar{x}}(\bar{x}) \geq \min_y Q_{\bar{x}}(y) = f(\bar{x}) - \frac{1}{2\alpha} \|Df(\bar{x})\|^2$

◻

Thm: Let f be β -smooth and α -cvx.
Then GD iterates with $\gamma = \frac{1}{\beta}$ satisfy.

$$\textcircled{1} \quad f(x_{t+1}) - f^* \leq \left(1 - \frac{1}{2\beta}\right) (f(x_t) - f^*)$$

$$\textcircled{2} \quad \|x_{t+1} - \bar{x}\|^2 \leq \left(\frac{\gamma t - 1}{\gamma + 1}\right) \|x_t - \bar{x}\|^2$$

Pf: We'll prove $\textcircled{1}$ and leave $\textcircled{2}$ for you to look at.

$$(f(x_t) - f^*) - (f(x_{t+1}) - f^*) \stackrel{\text{Descent lemma}}{\geq} \frac{1}{2\beta} \|Df(x_t)\|^2$$

$$\stackrel{\text{PL}}{\geq} \frac{1}{2\beta} (f(x_t) - f^*)$$

$$\rightarrow f(x_{t+1}) - f^* \leq \left(1 - \frac{1}{2\beta}\right) (f(x_t) - f^*)$$

Conclusion:

$$f(x_t) - f^* \leq \min \left\{ \frac{\beta \cdot \|x_0 - \bar{x}\|^2}{2t}, \left(1 - \frac{1}{2\beta}\right)^t (f(x_0) - f^*) \right\}$$

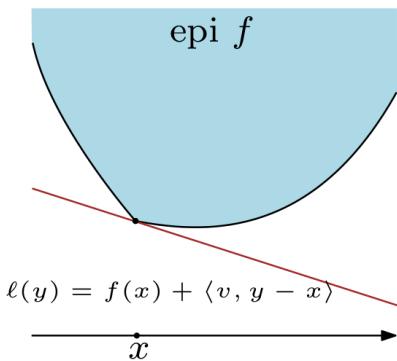
Stochastic algorithms (basic ones) are almost unaffected by smoothness or lack thereof of the convex f .

We will now focus on subgradients of nonsmooth CVX functions.

Defn: Let $f: \mathbb{E} \rightarrow [R \cup \infty]$ be CVX and let $x \in \mathbb{E}$ be s.t. $f(x) < \infty$. Then $v \in \mathbb{E}$ is a subgradient of f at x if

$$f(y) \geq f(x) + \langle v, y - x \rangle \quad \forall x, y$$

Then set of all such v is denoted by $\partial f(x)$ and is called the subdifferential.



(a) Subgradient of a convex function

Properties: Let f, g be cvx, $A \in \mathbb{R}^{m \times n}$

- If f is C^1 at x , then

$$\partial f(x) = \{\nabla f(x)\}$$

- If f is finite near x , then

(sum) $\partial(f+g)(x) = \partial f(x) + \partial g(x)$

where

$$A+B = \{a+b : a \in A, b \in B\}$$

- For $\varphi(x) = f(Ax - b)$, we have

$$\partial \varphi(x) = A^T \partial f(Ax - b)$$

provided f is finite near $Ax - b$.

Ex: $\varphi(x) = \|Ax - b\|_1 + \|x\|_\infty$

Then

$$\begin{aligned}\partial \varphi(x) &= \partial(x \mapsto \|Ax - b\|_1) + \partial\| \cdot \|_\infty(x) \\ &= A^T \partial\| \cdot \|_1(Ax - b) + \partial\| \cdot \|_\infty(x)\end{aligned}$$

More rules:

- If $\varphi(x) = \sum_{i=1}^d \varphi_i(x_i)$, then

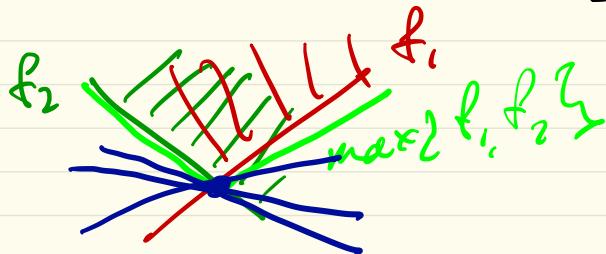
$$\partial \varphi(x) = \partial \varphi_1(x_1) \times \dots \times \partial \varphi_d(x_d)$$

- If f_1, \dots, f_d are ^{finite}CVX and $\varphi(x) = \max_{i=1,\dots,d} f_i(x)$. Then

$$\partial \varphi(x) = \text{conv} \left\{ \bigcup_{i \in I(x)} \partial f_i(x) \right\}$$

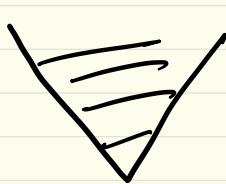
where

$$I(x) := \left\{ i : f_i(x) = \varphi(x) \right\}$$



So

$$\partial \|\cdot\|_I(x) = \partial I_1(x_1) \times \dots \times \partial I_d(x_d)$$



$$\partial I_i(s) = \begin{cases} \text{sign}(s), & s \neq 0 \\ [-1, 1], & s = 0 \end{cases}$$

And

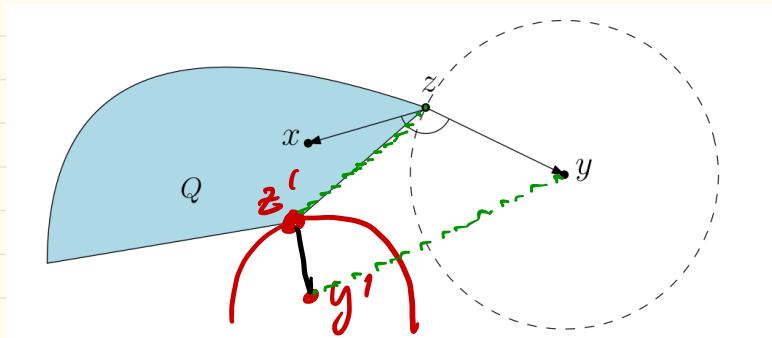
$$\partial \|\cdot\|_\infty(x) = \text{conv} \left\{ A_i : i \in I(x) \right\}$$

$$\text{where } I(x) = \left\{ i : \|x_i\|_\infty = \|x\|_\infty \right\}$$

$$\text{and } A_i = \left\{ \mathbb{R}^d \times \mathbb{R}^d_{x_{-i}} \times \underset{i^{\text{th}} \text{ entry}}{\underbrace{\partial I_i(x_i)}_{-\mathbb{R}^d}} \right\}$$

Last theory b.t:
For a set Q , we define

$$\text{proj}_Q(y) = \underset{x \in Q}{\operatorname{arg\,min}} \|x - y\|_2^2$$



Then: Suppose Q is closed CVX.
Then $\text{proj}_Q(y)$ is a singleton for
any y and satisfies

$$\|\text{proj}_Q(y) - \text{proj}_Q(y')\| \leq \|y - y'\|_{y,y'}$$

Problem:

$$\min_{\substack{x \in Q}} f(x)$$

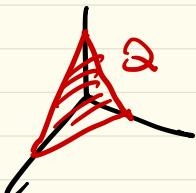
where Q is closed cvx and f is L -Lipschitz on a neighborhood of Q .

$$\text{Ex: } \min_z \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - y_i \langle \varphi(x_i), z \rangle\}$$

Kernel Support Vector Machines.

$$\begin{aligned} \text{Ex: } & \min \|Ax - b\|_1 \\ \text{s.t. } & \|x\|_\infty \leq \lambda \end{aligned}$$

$$\begin{aligned} \text{Ex: } & \min \|Ax - b\|_1 \\ \text{s.t. } & \sum_{i=1}^d x_i = 1, \quad x \geq 0 \end{aligned}$$



Stochastic oracle: we assume existence
of a map $(x, z) \mapsto g(x, z)$ with $z \sim P$ random s.t.

$$\mathbb{E}_z g(x, z) \in \partial f(x)$$

$$\mathbb{E}_z \|g(x, z)\|^2 \leq L^2 \quad \forall x \in Q$$

Ex: $g_l(x, z) = \nabla l(x, z)$. if $f(x) = \mathbb{E} l(x, z)$
if l is smooth.

Ex: $g(x, z) \in \partial l(x, z)$, if $f(x) = \mathbb{E} l(x, z)$
and l may be non-smooth.

Projected stochastic gradient method:

$$\left\{ \begin{array}{l} \text{Draw } z_t \sim P \\ \text{Set } x_{t+1} = \underset{Q}{\text{proj}}(x_t - \eta_t g(x_t, z_t)) \end{array} \right\}$$

Thm: Suppose f is cvx and L -Lip on a neighborhood of Q . Then

$$\mathbb{E} f\left(\frac{1}{\sum_{i=0}^t \gamma_i} \sum_{i=0}^t \gamma_i x_i\right) - f^* \leq \frac{\|x_0 - \bar{x}\|^2 + L^2 \sum_{i=0}^t \gamma_i^2}{2 \sum_{i=0}^t \gamma_i}$$

In particular if $\gamma_t = \frac{R}{L\sqrt{T+1}}$ where

$R \geq \|x_0 - \bar{x}\|$, we get

$$\mathbb{E} f\left(\frac{1}{T+1} \sum_{t=0}^T x_t\right) - f^* \leq \frac{RL}{\sqrt{T+1}}$$

Pf: Set $v_t = G(x_t, z_t)$. Then $\text{proj}_Q(\bar{x})$

$$\begin{aligned} \|x_{t+1} - \bar{x}\|^2 &= \|\text{proj}_Q(x_t - \gamma_t v_t) - \bar{x}\|^2 \\ &\leq \|x_t - \gamma_t v_t - \bar{x}\|^2 \\ &= \|x_t - \bar{x}\|^2 - 2\gamma_t \langle v_t, x_t - \bar{x} \rangle + \gamma_t^2 \|v_t\|^2 \end{aligned}$$

We take $\mathbb{E}_{t+1}[\cdot] = \mathbb{E}[\cdot | x_t]$ to get $\mathbb{E}_{t+1}[f(x_t)]$

$$\begin{aligned} \mathbb{E}_t \|x_{t+1} - \bar{x}\|^2 &\leq \|x_t - \bar{x}\|^2 - 2\gamma_t \langle \mathbb{E}[v_t], x_t - \bar{x} \rangle + \gamma_t^2 L^2 \\ &\leq \|x_t - \bar{x}\|^2 - 2\gamma_t (f(x_t) - f^*) + \gamma_t^2 L^2 \end{aligned}$$

Now take \mathbb{E} w.r.t x_t to get

$$\begin{aligned}\mathbb{E} \|x_{t+1} - \bar{x}\|^2 &\leq \mathbb{E} \|x_t - \bar{x}\|^2 - 2\eta_t (\mathbb{E} f(x_t) - f^*) + \eta_t^2 L^2 \\ &\leq \|x_0 - \bar{x}\|^2 - 2 \sum_{t=0}^T \eta_t (\mathbb{E} f(x_t) - f^*) + L^2 \sum_{t=0}^T \eta_t^2\end{aligned}$$

$$\Rightarrow \underbrace{\frac{\sum_{t=0}^T \eta_t (\mathbb{E} f(x_t) - f^*)}{\sum \eta_t}}_{\text{f(average)}} \leq \frac{\|x_0 - \bar{x}\|^2 + L^2 \sum_{t=0}^T \eta_t^2}{2 \sum_{t=0}^T \eta_t}$$

$$f(\text{average}) - f^* \quad \blacksquare$$

So to ensure $\mathbb{E} f(\bar{x}) - f^* \leq \epsilon$ suffices

$$T \geq \frac{R^2 L^2}{\epsilon^2}$$

Note this is a factor ϵ worse than
GD for smooth deterministic problems!

Then: Suppose f is α -cvx and L -Lip on a neighbourhood of closed cvx Q . Then with $\gamma_t = \frac{2}{\alpha(t+1)}$ get

$$\mathbb{E} f\left(\frac{2}{\alpha(t+1)} \sum_{i=1}^t i x_i\right) - f^* \leq \frac{2L^2}{\alpha(t+1)}$$

p.s.: Same as before:

$$\mathbb{E}_t \|x_{t+1} - \bar{x}\|^2 \leq \|x_t - \bar{x}\|^2 + 2\gamma_t (f^* - f(x_t)) - \frac{\alpha}{2} \|x_t - \bar{x}\|^2 + \gamma_t^2 L^2$$

Rearrange, take \mathbb{E} w.r.t. x_t , divide by γ_t :

$$\mathbb{E} f(x_t) - f^* \leq \left(\frac{1 - \alpha \gamma_t}{2\gamma_t}\right) \mathbb{E} \|x_t - \bar{x}\|^2 - \frac{1}{2\gamma_t} \mathbb{E} \|x_{t+1} - \bar{x}\|^2 + \frac{\gamma_t^2 L^2}{2}$$

Algebraic trick: $\gamma_t = \frac{2}{\alpha(t+1)}$ and multiply by t :

$$t \mathbb{E} f(x_t) - f^* \leq \frac{\alpha t(t+1)}{4} \mathbb{E} \|x_t - \bar{x}\|^2 - \frac{\alpha t(t+1)}{4} \mathbb{E} \|x_{t+1} - \bar{x}\|^2 + \frac{t}{\alpha(t+1)} L^2$$

Sum up and telescope:

$$\sum_{i=1}^t i (\mathbb{E} f(x_i) - f^*) \leq \sum_{i=1}^t \frac{i}{\alpha(i+1)} L^2 \leq \frac{tL^2}{\alpha}$$

Divide by $\sum_{i=1}^t i = \frac{t(t+1)}{2}$ and continue \square

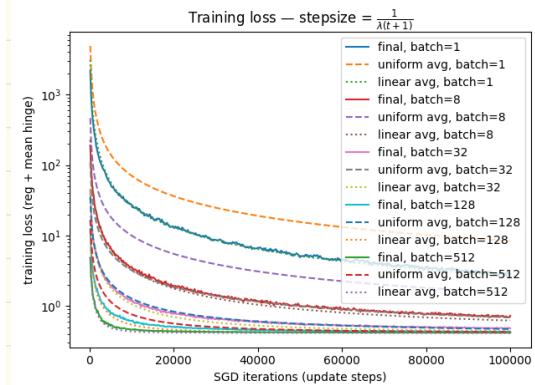
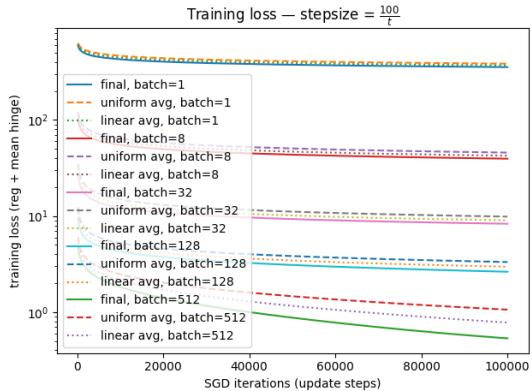
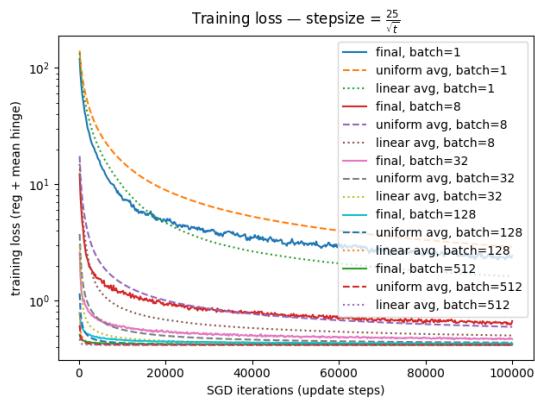
Regularized SVM

$$\min_{\mathbf{x}} \frac{1}{n} \sum_{i=1}^n \max\{0, 1 - y_i \cdot \mathbf{x}^T \mathbf{z}_i\}$$

$$+ \frac{\lambda}{2} \|\mathbf{x}\|^2$$

with data (\mathbf{z}_i, y_i)

with $\mathbf{z}_i \in \mathbb{R}^{2000}$, $n = 10,000$
 $\lambda = 10^{-3}$



Batch size does improve performance in experiments. In theory, this should not be the case because

$$\begin{aligned} \mathbb{E} \left\| \frac{1}{b} \sum_{i=1}^b \mathbf{v}_i \right\|_2^2 &= \mathbb{E} \left\| \frac{1}{b} \sum_{i=1}^b (\mathbf{v}_i - \mathbb{E} \mathbf{v}_i) + \frac{1}{b} \sum_{i=1}^b \mathbb{E} \mathbf{v}_i \right\|_2^2 \\ &= \mathbb{E} \left\| \frac{1}{b} \sum_{i=1}^b (\mathbf{v}_i - \mathbb{E} \mathbf{v}_i) \right\|_2^2 \approx O(\mathbb{E} \mathbf{f}(\mathbf{x}))^2 \\ \text{with } \mathbf{v}_i \in \partial f_i(\mathbf{x}_i) &\quad + \left\| \frac{1}{b} \sum_{i=1}^b \mathbb{E} \mathbf{v}_i \right\|_2^2 \approx \|\mathbb{E} f(\mathbf{x})\|^2 \end{aligned}$$

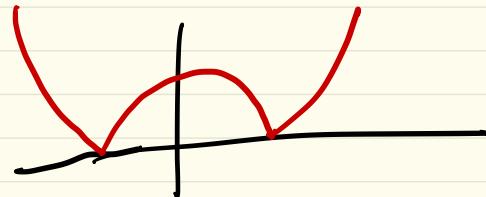
So

$$L^2 \approx O\left(\frac{\delta^2}{b}\right) + \|\nabla f(x)\|_2^2$$

If $\|\nabla f(x)\|$ is small (or gets smaller as you iterate) bigger batches help. This typical if $f(x) = \frac{1}{2} \sum l(x, z)$ is smooth.

If $\|\nabla f(x)\|$ is not small, then bigger batches don't help.

Eg: $f(x) = \underset{a \sim N(0, I)}{\mathbb{E}} |\langle x, a \rangle - \langle x_\#, a \rangle|^2$



We now replace Lipschitz continuity by smoothness.

Then: Let $f: E \rightarrow \mathbb{R}$ be β -smooth. and suppose that $g(x, z)$ satisfies

$$\mathbb{E}_z g(x, z) = Df(x) \text{ and}$$

$$\mathbb{E}_z \|g(x, z) - Df(x)\|_2^2 \leq \sigma^2 \quad \forall x.$$

*can be relaxed
in many ways.*

Then SGD with step-size $\gamma \leq \frac{1}{\beta}$ satisfies:

$$\mathbb{E} \|Df(x_{i^*})\|^2 \leq \frac{(f(x_0) - f^*) + (\sum_{j=1}^T \gamma_j)^2 \beta \sigma^2}{T}$$

where $i^* \in \{1, \dots, T\}$ is chosen according to

$$P(i^* = i) = \gamma_i / \sum_{j=1}^T \gamma_j.$$

pf: Compute for $V_t = g(x_t, z_t)$ the estimate

$$f(x_{t+1}) \leq f(x_t) + \langle Df(x_t), x_{t+1} - x_t \rangle + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$

$$= f(x_t) - \gamma_t \langle Df(x_t), V_t \rangle + \frac{\beta \gamma_t^2}{2} \|V_t\|^2$$

Now take $\mathbb{E}_t[f(x_{t+1})] = \mathbb{E}[\cdot | x_t]$ to get

$$\mathbb{E}_t f(x_{t+1}) \leq f(x_t) - \gamma_t \|Df(x_t)\|^2 + \frac{\gamma_t^2 \beta}{2} \mathbb{E} \|V_t\|^2$$

Note

$$\mathbb{E} \|g(x, z)\|^2 = \underbrace{\mathbb{E} \|g(x, z) - Df(x)\|^2}_{\leq \delta^2} + \|Df(x)\|_2^2$$

So now:

$$\begin{aligned} \mathbb{E}_t f(x_{t+1}) &\leq f(x_t) - \gamma_t \left(1 - \frac{\gamma_t \beta}{2}\right) \|Df(x_t)\|^2 \\ &\quad + \frac{\gamma_t^2 \beta \delta^2}{2} \end{aligned}$$

If $\gamma_t \leq \frac{1}{\beta}$, then $1 - \frac{\gamma_t \beta}{2} \geq \frac{1}{2}$

So

$$\begin{aligned} \mathbb{E}_t f(x_{t+1}) &\leq f(x_t) - \frac{\gamma_t}{2} \|Df(x_t)\|^2 \\ &\quad + \frac{\gamma_t^2 \delta^2 \beta}{2}. \end{aligned}$$

Take \mathbb{E} w.r.t. x_t to get

$$\mathbb{E} f(x_{t+1}) \leq \mathbb{E} f(x_t) - \frac{\gamma_t}{2} \mathbb{E} \|Df(x_t)\|^2 + \frac{\gamma_t^2 B^2 \sigma^2}{2}$$

So iterate to get

$$\mathbb{E} f(x_{t+1}) \leq f(x_0) - \sum_{i=1}^t \frac{\gamma_i}{2} \mathbb{E} \|Df(x_i)\|^2 + \left(\sum_{i=1}^t \gamma_i^2 \right) \cdot \frac{B\sigma^2}{2}$$

Rearrange:

$$\frac{\frac{1}{2} \sum_{i=1}^t \gamma_i \mathbb{E} \|Df(x_i)\|^2}{\sum_{i=1}^t \gamma_i} \leq \frac{\left(f(x_0) - f^* + \left(\sum_{i=1}^t \gamma_i^2 \right) \frac{B\sigma^2}{2} \right)}{\sum_{i=1}^t \gamma_i}$$

Pick i^* with probability $P(i^* = i) = \frac{n_i}{\sum_{j=1}^T n_j}$

and return $\|Df(x_{i^*})\|_2^2$. Then

$$\begin{aligned} \mathbb{E} \|Df(x_{i^*})\|_2^2 &= \frac{\sum_{i=1}^T \frac{n_i}{\sum_{j=1}^T n_j} \cdot \|Df(x_i)\|_2^2}{\sum_{i=1}^T n_i} \\ &\leq \underbrace{f(x_0) - f^* + \left(\sum_{i=1}^T n_i\right) \beta \gamma^2}_{\sum_{i=1}^T n_i} \end{aligned}$$

In particular, if $\gamma_t \equiv \gamma < \frac{1}{\beta}$.
Then

$$\mathbb{E} \|Df(x_{i^*})\|_2^2 \leq \frac{f(x_0) - f^*}{T\gamma} + \gamma \beta \gamma^2$$

Thm: Suppose $f: \mathbb{E} \rightarrow \mathbb{R}$ is β -smooth and cvx. Suppose

$$\mathbb{E}_z g(x; z) = \nabla f(x), \quad \mathbb{E}_z \|g(x; z) - \nabla f(x)\|^2 \leq \sigma^2$$

Then for $\gamma_t = \gamma < \frac{1}{2\beta}$ get

$$\mathbb{E} f(\hat{x}_T) - f^* \leq \frac{\|x_0 - \bar{x}\|^2}{2\gamma T} + \frac{\gamma}{2} \sigma^2.$$

where $\hat{x}_T = \frac{1}{T} \sum_{t=0}^T x_t$

Pf: Compute for $g_t = g(x_t, z_t)$ to get

$$\begin{aligned} \|x_{t+1} - \bar{x}\|^2 &= \|x_t - \gamma g_t - \bar{x}\|^2 \\ &= \|x_t - \bar{x}\|^2 - 2\gamma \langle g_t, x_t - \bar{x} \rangle + \gamma^2 \|g_t\|^2 \end{aligned}$$

Take \mathbb{E}_t to get

$$\begin{aligned} \mathbb{E}_t \|x_{t+1} - \bar{x}\|^2 &= \|x_t - \bar{x}\|^2 - 2\gamma \underbrace{\langle \nabla f(x_t), x_t - \bar{x} \rangle}_{\geq f(x_t) - f^*} + \gamma^2 \underbrace{\mathbb{E}_t \|g_t\|^2}_{= \mathbb{E} \|g_t - \nabla f(x_t)\|^2 + \|\nabla f(x_t)\|^2} \end{aligned}$$

$$\begin{aligned} &= \mathbb{E} \|g_t - \nabla f(x_t)\|^2 \\ &\quad + \|\nabla f(x_t)\|^2 \end{aligned}$$

$$\mathbb{E}_t \|\bar{x}_{t+1} - \bar{x}\|^2 \leq \|\bar{x}_t - \bar{x}\|^2 - 2\gamma (f(\bar{x}_t) - f^*)$$

$$+ \gamma^2 \sigma^2 + \gamma^2 \cdot \|\nabla f(\bar{x}_t)\|^2$$

Smoothness implies

$$\textcircled{*} \quad \|\nabla f(x)\|^2 \leq 2B(f(x) - f^*) \quad \text{we will check this.}$$

Then

$$\mathbb{E}_t \|\bar{x}_{t+1} - \bar{x}\|^2 \leq \|\bar{x}_t - \bar{x}\|^2 - 2\gamma \underbrace{(1 - \beta\gamma)}_{\geq \frac{1}{2}} (f(\bar{x}_t) - f^*) + \gamma^2 \sigma^2$$

Take \mathbb{E} and rearrange:

$$\gamma (\mathbb{E} f(\bar{x}_t) - f^*) \leq \mathbb{E} \|\bar{x}_t - \bar{x}\|^2 - \mathbb{E} \|\bar{x}_{t+1} - \bar{x}\|^2 + \gamma^2 \sigma^2$$

Sum up (right side telescopes), divide by $\gamma T \epsilon$ and lower bound left side by $f(\bar{x}_T) - f^*$ \square

If you do step-decay get complexity

$$O\left(\frac{\beta \|\bar{x}_0 - \bar{x}\|^2}{\epsilon} + \frac{\sigma^2}{\epsilon}\right)$$

Check at \star :

$$\frac{f(x) - f\left(x - \frac{1}{\beta} \nabla f(x)\right)}{\|f(x) - f^*\|} \geq \frac{1}{2\beta} \|\nabla f(x)\|^2$$

Consequences: If we optimize in y

we get: $\eta = \frac{\|x_0 - \bar{x}\|}{\sqrt{5T}}$

to get

$$\mathbb{E}(f(\bar{x}_T) - f^*)$$

$$\leq O\left(\frac{\sigma \cdot \|x_0 - \bar{x}\|}{\sqrt{T}}\right)$$

So if $f(x) = \mathbb{E}_{\mathcal{P}} l(x, z)$

and $g_t = \frac{1}{b} \sum_{i=1}^b \nabla l(x, z_i)$

for $z_1, \dots, z_b \stackrel{iid}{\sim} p$, then get

$$\frac{\sigma \|x_0 - \bar{x}\|}{b \cdot \sqrt{T}}$$

Then: Suppose f is β -smooth and α -cvx.
Suppose

$$\mathbb{E} g(x, z) = \nabla f(x), \quad \mathbb{E} \|g(x, z) - \nabla f(x)\|^2 \leq \sigma^2$$

Then with $\gamma_t \leq \frac{1}{4\beta}$ get

$$\eta_t \mathbb{E}(f(x_t) - f^*) \leq (1 - \alpha \gamma_t) \mathbb{E}\|x_t - \bar{x}\|^2 - \mathbb{E}\|x_{t+1} - \bar{x}\|^2 + \gamma_t^2 \sigma^2$$

pf: Compute for $g_t = g(x_t, z_t)$ get

$$\begin{aligned} \|x_{t+1} - \bar{x}\|^2 &\leq \|x_t - \bar{x} - \gamma_t g_t\|^2 \\ &= \|x_t - \bar{x}\|^2 - 2\gamma_t \langle g_t, x_t - \bar{x} \rangle + \gamma_t^2 \|g_t\|^2 \end{aligned}$$

Take \mathbb{E}_t and α -cvx to get

$$\begin{aligned} \mathbb{E}_t \|x_{t+1} - \bar{x}\|^2 &\leq \|x_t - \bar{x}\|^2 - 2\gamma_t [f(x_t) - f^* + \frac{\alpha}{2} \|x_t - \bar{x}\|^2] \\ &\quad + \gamma_t^2 [\|\nabla f(x)\|^2 + \sigma^2] \\ &\leq (1 - \alpha \gamma_t) \|x_t - \bar{x}\|^2 - 2\gamma_t [f(x_t) - f^*] \\ &\quad + [2B\gamma_t^2 (f(x_t) - f^*) + \gamma_t^2 \sigma^2] \end{aligned}$$

$$\begin{aligned}
 &= (1-\alpha\gamma_t) \|x_t - \bar{x}\|^2 - 2\gamma_t (1-2\beta\gamma_t) \frac{(f(x_t) - f^*)}{\gamma_t} \\
 &\leq (1-\alpha\gamma_t) \|x_t - \bar{x}\|^2 - \gamma_t (f(x_t) - f^*) + \frac{\gamma_t^2 \epsilon^2}{2}
 \end{aligned}$$

Take E and rearrange:

$$\begin{aligned}
 \mathbb{E}(\gamma_t (f(x_t) - f^*)) &\leq (1-\alpha\gamma_t) \mathbb{E} \|x_t - \bar{x}\|^2 - \mathbb{E} \|x_{t+1} - \bar{x}\|^2 \\
 &\quad + \frac{\gamma_t^2 \epsilon^2}{2}
 \end{aligned}$$

as claimed. \square

So if $\gamma_t = \gamma \leq \frac{1}{4\beta}$, get

$$\gamma \mathbb{E}(f(x_t) - f^*) \leq (1-\alpha\gamma) \mathbb{E} \|x_t - \bar{x}\|^2 - \mathbb{E} \|x_{t+1} - \bar{x}\|^2$$

Define: $\Gamma_t = \prod_{i=0}^t (1-\alpha\gamma)^i$. Divide by Γ_t :

$$\frac{\gamma}{\Gamma_t} \mathbb{E}(f(x_t) - f^*) \leq \frac{1}{\Gamma_{t-1}} \mathbb{E} \|x_t - \bar{x}\|^2 - \frac{1}{\Gamma_t} \mathbb{E} \|x_{t+1} - \bar{x}\|^2 + \frac{\gamma^2 \epsilon^2}{\Gamma_t}$$

So sum up to get:

$$\sum_{i=1}^t \frac{2}{\Gamma_i} E(f(x_i) - f^*) \leq \|x_0 - \bar{x}\|^2 + \sum_{i=0}^t \frac{\eta^2 \zeta^2}{\Gamma_i}$$

Algebraic Identity $\Leftarrow \frac{n\cdot \zeta^2}{2\Gamma_t}$

$$1 + \sum_{i=0}^t \frac{1}{\Gamma_i} = \frac{1}{\Gamma_t}$$

[Check!]

So the average iterate

$$\hat{x}_t = \frac{1}{\sum_{i=0}^t \frac{1}{\Gamma_i}} \sum_{i=0}^t \frac{n}{\Gamma_i} x_i$$

satisfies

$$E f(\hat{x}) - f^* \leq \frac{1}{t} \left(\|x_0 - \bar{x}\|^2 + \frac{2\zeta^2}{\sum_{i=0}^t \frac{1}{\Gamma_i}} \right)$$

Finally:

$$E f(\hat{x}_t) - f^* \leq (1-\alpha\gamma)^t \|x_0 - \bar{x}\|^2 + \frac{\gamma \varsigma^2}{\lambda}.$$

Now do step-decay to get complexity

$$\boxed{\frac{\beta}{\alpha} \log\left(\frac{\|x_0 - \bar{x}\|^2}{\epsilon}\right) + \frac{\varsigma^2}{\alpha \epsilon}}$$

(Check!)

Variance Reduction:

$$\min_x f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$$

where each f_i is β -smooth, f is α -cvx

Two options:

- ① Do GD \Rightarrow complexity $O\left(\frac{\beta}{\alpha} \log\left(\frac{1}{\epsilon}\right)\right)$
- ② Do SGD \Rightarrow complexity $O\left(\frac{n^2}{\alpha \epsilon}\right)$ ignoring variance

In terms of individual grad evals:

$$\underbrace{\frac{n \beta}{\alpha} \log\left(\frac{1}{\epsilon}\right)}_{GD} \quad vs \quad \underbrace{\frac{n^2}{\alpha \epsilon}}_{SGD}$$

Typically $\alpha \approx \frac{1}{n}$. So $n^2 \log\left(\frac{1}{\epsilon}\right)$ vs $\frac{n \beta^2}{\epsilon}$

We'll get an algo with complexity

$$(n + \frac{\beta}{\alpha}) \log\left(\frac{1}{\epsilon}\right)$$

We'll talk about SVRG!

Idea is to reduce the variance of $\nabla f_i(x_i)$, where i is sample uniformly from $\{1, \dots, n\}$.
 Candidate: Given x_1, \dots, x_n Look at

$$\nabla f_i(x_i) - \nabla f_i(y) + \nabla f(y)$$

Remark: y should ideally satisfy $y \approx \bar{x}$.

Lemma:

$$\mathbb{E} \| \nabla f_i(x) - \nabla f_i(\bar{x}) \|_2^2 \leq 2\beta(f(x) - f^*)$$

Pf: Define $h(x) = f_i(x) - f_i(\bar{x}) - \langle \nabla f_i(\bar{x}), x - \bar{x} \rangle$

Then since h is β -smooth, get:

$$\| \nabla h(x) \|_2^2 \leq 2\beta(h(x) - \underbrace{\min_{\bar{x}} h}_{\text{"Cvx"}})$$

$$\therefore \| \nabla f_i(x) - \nabla f_i(\bar{x}) \|_2^2 \leq 2\beta h(x)$$

Note $\mathbb{E} h(x) = 2\beta(f(x) - f^*)$ B

Consequently

$$E \| Df_i(x) - Df_i(y) + Df(y) \|^2$$

$$= E \| Df_i(x) - Df_i(\bar{x}) + Df_i(\bar{x}) - Df_i(y) + Df(y) \|^2$$

$$\left[(a+b)^2 \leq 2a^2 + 2b^2 \right]$$

$$\leq 2 \underbrace{E \| Df_i(x) - Df_i(\bar{x}) \|^2}_{\text{Lemma}} + 2 \underbrace{E \| Df_i(\bar{x}) - Df_i(y) + Df(y) \|^2}$$

$$2B \| f(x) - f^* \|^2$$

$$\begin{aligned} & E \| Df_i(\bar{x}) - Df_i(y) \|^2 \\ & \underbrace{\leq 2B \| f(y) - f^* \|^2}_{\text{Lemma}} \end{aligned}$$



SVRG: Let $y^{(0)} \in E$. For $s=1, 2, \dots$, let

$$x_i^{(s)} = y^{(s)}$$

For $t=1, \dots, k$ let

$$x_{t+1}^{(s)} = x_t^{(s)} - \eta \left(Df_{i_t^{(s)}}(x_t^{(s)}) - Df_{i_t^{(s)}}(y^{(s)}) \right)$$

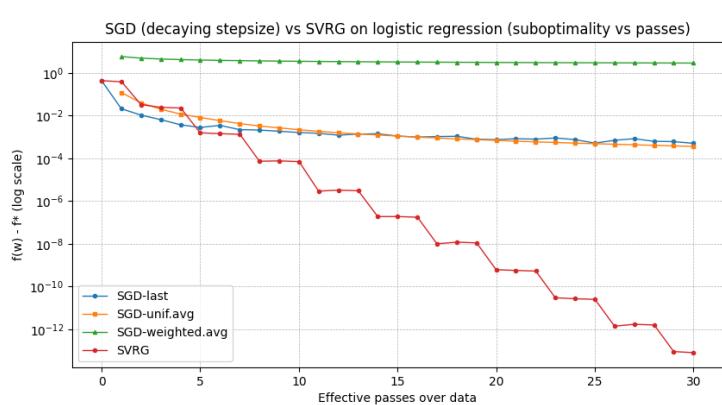
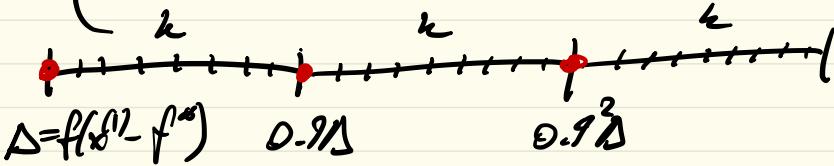
where $i_t^{(s)}$ is iid $\{1, \dots, n\}$. Let $y^{(s+1)} = \frac{1}{K} \sum_{t=1}^K x_t^{(s)}$

Thm: SVRG with $\gamma = \frac{1}{10B}$ and $k = 20 \frac{B}{\alpha}$ satisfies

$$\mathbb{E} f(y^{(s+1)}) - f^* \leq 0.9^s (f(y^{(1)}) - f^*)$$

So after $\frac{1}{0.9} \log\left(\frac{f(x^{(1)}) - f^*}{\epsilon}\right)$ epochs we have $\mathbb{E} f(y^{(s+1)}) - f^* \leq \epsilon$. Each epoch uses $\Theta\left(\frac{B}{\alpha}\right)$ individual gradients. Total # ind. grad evals is

$$O\left((n + \frac{B}{\alpha}) \log\left(\frac{f(x^{(1)}) - f^*}{\epsilon}\right)\right).$$



pt: Let's look at a single epoch initialized at anchor y .
 Set $V_t = \nabla f_{i_t^*}(x_t) - \nabla f_{i_t^*}(y) + \nabla f(y)$

Then

$$\|x_{t+1} - \bar{x}\|^2 = \|x_t - \bar{x}\|^2 - 2\gamma \langle V_t, x_t - \bar{x} \rangle$$

$$+ \gamma^2 \|V_t\|^2$$

Take E_t to get $\leq f(x_t) - f^*$

$$E_t \|x_{t+1} - \bar{x}\|^2 = \|x_t - \bar{x}\|^2 - 2\gamma \underbrace{\langle \nabla f(x_t), x_t - \bar{x} \rangle}_{+ \gamma^2 E_t \|V_t\|^2}$$

$$\cancel{+ \gamma^2 E_t \|V_t\|^2} \leq 2\beta(f(x_t) - f^*)$$

$$+ 2\beta(f(y) - f^*)$$

$$\leq \|x_t - \bar{x}\|^2 - 2\gamma(1-\beta\gamma)(f(x_t) - f^*)$$

$$+ 2\beta\gamma^2(f(y) - f^*)$$

Combining

Take \mathbb{E} and iterate:

$$\mathbb{E} \left[\frac{\|x_{t+1} - \bar{x}\|^2}{\gamma} \right] \leq \|y - \bar{x}\|^2 - 2\gamma(1 - 2\beta)\mathbb{E} \sum_{i=1}^k (f(x_i) - f^*) + 2\beta\gamma^2 k (f(y) - f^*)$$

Recall: $\frac{\|y - \bar{x}\|^2}{\gamma} \stackrel{\text{Cvx}}{\leq} (f(y) - f^*)$

So rearrange

$$\mathbb{E} \left[\frac{1}{\gamma} \sum_{i=1}^k (f(x_i) - f^*) \right] \leq \frac{(2 + 2\beta\gamma^2 k)(f(y) - f^*)}{2\gamma(1 - 2\beta) \cdot k}$$

$\stackrel{\text{Cvx}}{=} \left(\frac{1}{2(1 - 2\beta)k\gamma} + \frac{\beta k}{1 - 2\beta} \right) (f(y) - f^*) \leq 0.9 (f(y) - f^*)$

Done!

Note for some problems the variance actually shrinks to zero automatically without any variance reduction.

Recall:

$$\mathbb{E} \|Df_i(x) - Df_i(\bar{x})\|_2^2 \leq 2\beta(f(x) - f^*)$$

What if $Df_i(\bar{x}) = 0$?

↑ means \bar{x} is a minimizer of f_i .

Problems satisfying

such a condition are called interpolation problems