

GD basics $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$

$$\sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2)$$

Lipschitz convex: 1/eps^2

$$\gamma := \frac{R}{B\sqrt{T}} \quad \frac{1}{T} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{RB}{\sqrt{T}}$$

Smooth convex: 1/eps

$$\gamma := \frac{1}{L} \quad f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$
$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2$$

suff. descent:

Smooth convex accelerated: 1/sqrt(eps)

Nesterov's accelerated gradient descent ('83)

Smooth strongly-convex: (L/mu) log(1/eps)

$$\gamma := \frac{1}{L} \quad \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2$$
$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

PGD $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$ $\mathbf{x}_{t+1} := \Pi_X(\mathbf{y}_{t+1})$

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \leq \|\mathbf{y}_{t+1} - \mathbf{x}^*\|^2$$

(second term can be seen as noise, often cancels out)

Lipschitz convex: idem, 1/eps^2

Technically, only need bounded gradient \ni lipschitz (X closed)

Smooth convex: 1/eps

suff. "descent" $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$

Smooth strongly-convex: (L/mu) log(1/eps)

square distance to OPT still geom. decreasing, but

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \|\nabla f(\mathbf{x}^*)\| \left(1 - \frac{\mu}{L}\right)^{T/2} \|\mathbf{x}_0 - \mathbf{x}^*\| + \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T > 0.$$

Proximal grad: $\mathbf{f}=\mathbf{g}+\mathbf{h}$, $\mathbf{x}_{t+1} := \text{prox}_{h,\gamma}(\mathbf{x}_t - \gamma \nabla g(\mathbf{x}_t))$

$$\text{prox}_{h,\gamma}(\mathbf{z}) := \underset{\mathbf{y}}{\text{argmin}} \left\{ \frac{1}{2\gamma} \|\mathbf{y} - \mathbf{z}\|^2 + h(\mathbf{y}) \right\} \quad \text{non-expansive}$$

g smooth, g,h convex, $\gamma := \frac{1}{L} : f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$

Subgradient descent $\mathbf{g}_t \in \partial f(\mathbf{x}_t)$

Lipschitz convex: idem, 1/eps^2

Tame strongly-convex $B = \max_{t=1}^T \|\mathbf{g}_t\| : 1/\text{eps}$

$$\gamma_t := \frac{2}{\mu(t+1)} \quad f\left(\frac{2}{T(T+1)} \sum_{t=1}^T t \cdot \mathbf{x}_t\right) - f(\mathbf{x}^*) \leq \frac{2B^2}{\mu(T+1)}$$

Reason for step-size choice: must multiply by t before telescoping

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{B^2 \gamma_t}{2} + \frac{(\gamma_t^{-1} - \mu)}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \frac{\gamma_t^{-1}}{2} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2$$

Lower bound (Nesterov): \exists f B-lipschitz s.t

$$\text{for any subgradient method,} \quad f(\mathbf{x}_T) - f(\mathbf{x}^*) \geq \frac{RB}{2(1 + \sqrt{T+1})}$$

SGD

Smoothness never helps! Only bounded gradients ("tame")
(Smoothness may help when doing variance-reduced SGD)

Lipschitz convex: idem, 1/eps^2

Technically, only need bounded stoch. gradients $\mathbb{E}[\|\mathbf{g}_t\|^2] \leq B^2$

Tame strongly-convex $B = \max_{t=1}^T \mathbb{E}[\|\mathbf{g}_t\|] : 1/\text{eps}$

$$\gamma_t := \frac{2}{\mu(t+1)} \quad \mathbb{E}\left[f\left(\frac{2}{T(T+1)} \sum_{t=1}^T t \cdot \mathbf{x}_t\right) - f(\mathbf{x}^*)\right] \leq \frac{2B^2}{\mu(T+1)}$$

(same proof as in subgradient descent with expectations)

Mini-batch reduces variance: $\mathbb{E}[\|\tilde{\mathbf{g}}_t - \nabla f(\mathbf{x}_t)\|^2] \leq \frac{B^2}{m}$

where batch size=m and $\mathbb{E}[\|\mathbf{g}_t\|^2] \leq B^2$

Nonconvex functions

Def. of smooth is still only an upper bound!

TL;DR: $\|\nabla f(\mathbf{x}_t)\|^2 \rightarrow 0$ at same rate as $f(\mathbf{x}_t) - f(\mathbf{x}^*)$ for convex

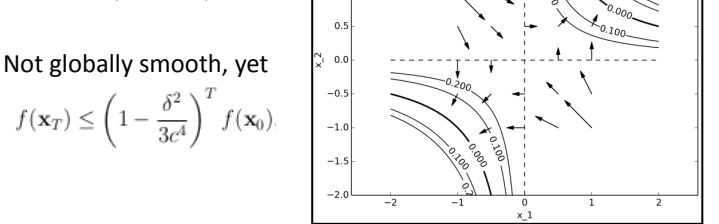
Smooth: 1/eps ON AVERAGE

$$\gamma := \frac{1}{L} \quad \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{2L}{T} (f(\mathbf{x}_0) - f(\mathbf{x}^*))$$

No-overshoot ppty: \nexists critical pt on segment $[\mathbf{x}_t, \mathbf{x}_{t+1}]$
that $\nabla f(\mathbf{x}) \neq \mathbf{0}$, i.e. \mathbf{x} is not a critical point. Suppose that f is smooth with parameter L over the line segment connecting \mathbf{x} and $\mathbf{x}' = \mathbf{x} - \gamma \nabla f(\mathbf{x})$, where $\gamma = 1/L' < 1/L$. Then \mathbf{x}' is also not a critical point.

The example

$$f(\mathbf{x}) := \frac{1}{2} \left(\prod_{k=1}^d x_k - 1 \right)^2$$



Not globally smooth, yet

$$f(\mathbf{x}_T) \leq \left(1 - \frac{\delta^2}{3c^4}\right)^T f(\mathbf{x}_0).$$

Newton's method $\mathbf{x}_{t+1} := \mathbf{x}_t - \nabla^2 f(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t)$

$$\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \mathbb{R}^d}{\text{argmin}} f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x} - \mathbf{x}_t) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_t)^\top \nabla^2 f(\mathbf{x}_t) (\mathbf{x} - \mathbf{x}_t)$$

(No step-size)

Rk. Newton's method is affine-invariant.

Thm. Suppose \exists ball around \mathbf{x}^* where (for spectral norm)

$$\|\nabla^2 f(\mathbf{x})^{-1}\| \leq \frac{1}{\mu} \quad \text{and} \quad \|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq B \|\mathbf{x} - \mathbf{y}\|$$

$$\text{Then if } \mathbf{x}_0 \in \text{ball}, \quad \|\mathbf{x}_{t+1} - \mathbf{x}^*\| \leq \frac{B}{2\mu} \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$

$$\text{Corr. If } \mathbf{x}_0 \in \text{ball and } \|\mathbf{x}_0 - \mathbf{x}^*\| \leq \frac{\mu}{B}, \quad \|\mathbf{x}_T - \mathbf{x}^*\| \leq \frac{\mu}{B} \left(\frac{1}{2}\right)^{2^{T-1}}$$

i.e to get $\|\mathbf{x}_T - \mathbf{x}^*\| < \varepsilon$, only need $T = \log \log(\frac{1}{\varepsilon})$

Local quadratic convergence ("double the number of correct digits in each iteration")

- affine invariant
- converge in 1 step for quadratics

Quasi-Newton $\mathbf{x}_{t+1} = \mathbf{x}_t - H_t^{-1} \nabla f(\mathbf{x}_t)$,

with H symmetric s.t secant condition:

$$\nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}_{t-1}) = H_t (\mathbf{x}_t - \mathbf{x}_{t-1})$$

$$x_{t+1} := x_t - f'(x_t) \frac{x_t - x_{t-1}}{f'(x_t) - f'(x_{t-1})}.$$

In 1D, only one secant method
Greenstadt family, of which (L)-BFGS
Newton \in Quasi-Newton \Leftrightarrow f nondegen. quadratic

Coordinate descent $\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_i \nabla_i f(\mathbf{x}_t) \mathbf{e}_i$

PL inequality: $\frac{1}{2} \|\nabla f(\mathbf{x})\|^2 \geq \mu(f(\mathbf{x}) - f(\mathbf{x}^*))$
 μ -strongly-convex $\Rightarrow \mu$ -PL
E.g $f(\mathbf{x}) = x_1^2$ is 1-PL but not SC
E.g $f(\mathbf{x}) := g(A\mathbf{x})$ for strongly convex g and arbitrary matrix A

GD on smooth + PL: $(L/\mu) \log(1/\epsilon)$
(Exact same proof as for smooth+SC)

Randomized CD $i \in [d]$ uniformly

If f is (L, \dots, L) -coord-wise-smooth and μ -PL, $\gamma_i = \frac{1}{L}$,
 $\mathbb{E}[f(\mathbf{x}_T) - f(\mathbf{x}^*)] \leq \left(1 - \frac{\mu}{dL}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*))$

Importance sampling CD $i \in [d]$ with probability $\frac{L_i}{\sum_{j=1}^d L_j}$

If f is (L_1, \dots, L_d) -coord-wise-smooth and μ -PL, stepsizes $\frac{1}{L_i}$,
 $\mathbb{E}[f(\mathbf{x}_T) - f(\mathbf{x}^*)] \leq \left(1 - \frac{\mu}{d\bar{L}}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*))$; $\bar{L} = \frac{1}{d} \sum_{i=1}^d L_i$
E.g $f(\mathbf{x}) = x_1^2$ is $(2, 0, \dots, 0)$ -smooth so $L=2$ and $\bar{L}=2/d$

Steepest CD aka Gauss-Southwell $i = \operatorname{argmax}_{i \in [d]} |\nabla_i f(\mathbf{x}_t)|$

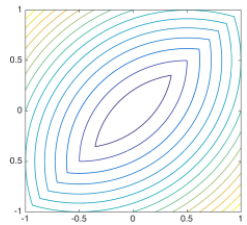
If f is (L, \dots, L) -coord-wise-smooth and μ -PL, same bound as for Randomized CD
 \rightarrow strictly worse bound, as per-iteration cost is $\sim d$

If f is (L, \dots, L) -coord-wise-smooth and μ_1 -PL w.r.t l_1 norm,
 $\gamma_i = \frac{1}{L}$, $f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \left(1 - \frac{\mu_1}{L}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*))$

Rk: μ_1 -SC w.r.t l_1 norm $\Rightarrow \mu_1$ -PL w.r.t l_1 norm
 $\frac{1}{2} \|\nabla f(\mathbf{x})\|_\infty^2 \geq \mu_1(f(\mathbf{x}) - f(\mathbf{x}^*))$

Greedy CD: line-search $\mathbf{x}_{t+1} := \operatorname{argmin}_{\lambda \in \mathbb{R}} f(\mathbf{x}_t + \lambda \mathbf{e}_i)$

May fail: $\|\mathbf{x}\|^2 + |x_1 - x_2|$
Thm. If $f = g+h$, g convex diffble and $h(\mathbf{x}) = \sum_i h_i(x_i)$ with h_i convex
then $\mathbf{x}_{t+1} = \mathbf{x}_t \Rightarrow \mathbf{x}_t$ global min of f



Frank-Wolfe aka conditional gradient

$$\mathbf{s} := \operatorname{LMO}_X(\nabla f(\mathbf{x}_t))$$
$$\mathbf{x}_{t+1} := (1 - \gamma_t)\mathbf{x}_t + \gamma_t \mathbf{s}$$

Lin. Min. Oracle $\operatorname{LMO}_X(\mathbf{g}) := \operatorname{argmin}_{\mathbf{z} \in X} \mathbf{g}^\top \mathbf{z}$

If $X = \operatorname{conv}(A)$, then $\operatorname{LMO}_X(\mathbf{g}) \in A$ ("atoms")

Examples	\mathcal{A}	$ \mathcal{A} $	dim.	$\operatorname{LMO}_X(\mathbf{g})$
L1-ball	$\{\pm \mathbf{e}_i\}$	$2d$	d	$\pm \mathbf{e}_i$ with $\operatorname{argmax}_i g_i $
Simplex	$\{\mathbf{e}_i\}$	d	d	\mathbf{e}_i with $\operatorname{argmax}_i g_i$
Spectahedron	$\{\mathbf{x} \mathbf{x}^\top, \ \mathbf{x}\ = 1\}$	∞	d^2	$\operatorname{argmin}_{\ \mathbf{x}\ =1} \mathbf{x}^\top G \mathbf{x}$
Norms	$\{\mathbf{x}, \ \mathbf{x}\ \leq 1\}$	∞	d	$\operatorname{argmin}_{\ \mathbf{s}\ \leq 1} \langle \mathbf{s}, \mathbf{g} \rangle$
Nuclear norm	$\{Y, \ Y\ _* \leq 1\}$	∞	d^2	..

(Spectrahedron: PSD matrices with trace=1;
 $\operatorname{LMO}_X(G) = \mathbf{s}_1 \mathbf{s}_1^\top$ via eigenvector)

Duality gap $g(\mathbf{x}) := \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{s})$

$g(\mathbf{x}) \geq f(\mathbf{x}) - f(\mathbf{x}^*)$ and $g(\mathbf{x}^*) = 0$

Curvature constant

$$C_{(f,X)} := \sup_{\substack{\mathbf{x}, \mathbf{s} \in X, \gamma \in (0,1] \\ \mathbf{y} = (1-\gamma)\mathbf{x} + \gamma \mathbf{s}}} \frac{1}{\gamma^2} (f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}))$$

If f is L -smooth, then $C_{(f,X)} \leq \frac{L}{2} \operatorname{diam}(X)^2$.
Allows to capture that Frank-Wolfe algo is **affine-invariant**.

Smooth convex: $1/\epsilon$

Analysis is quite different from (S/P/prox)GD, CD, Newton.

Thm. If X convex compact, f convex, $C_{(f,X)} < \infty$,

step-size $\gamma_t = 2/(t+2)$. (indep of params!)

$$\operatorname{argmin}_{\gamma \in [0,1]} f((1-\gamma)\mathbf{x}_t + \gamma \mathbf{s}) \quad \text{or} \quad \min \left(\frac{g(\mathbf{x}_t)}{L \|\mathbf{s} - \mathbf{x}_t\|^2}, 1 \right),$$

then $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top \gamma_t (\mathbf{s} - \mathbf{x}) + \gamma_t^2 C_{(f,X)}$

and so $h(\mathbf{x}_{t+1}) \leq (1 - \gamma_t)h(\mathbf{x}_t) + \gamma_t^2 C_{(f,X)}$

and so $f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{4C_{(f,X)}}{T+1}$

Thm ("cv" of duality gap). Under the same conditions,

$$g(\mathbf{x}_t) \leq \frac{27/2 \cdot C_{(f,X)}}{T+1}$$

there exists $1 \leq t \leq T$ s.t

Zero-th-order/gradient-free optim

pick a random direction $\mathbf{d}_t \in \mathbb{R}^d$
 $\gamma := \operatorname{argmin}_{\gamma \in \mathbb{R}} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$ (line-search)

Random search: $\mathbf{x}_{t+1} := \mathbf{x}_t + \gamma \mathbf{d}_t$
(step-size: no other choice than line-search!)

Convergence rates: same as GD with optimal step-size, with slow-down factor of d

Smooth convex: $T < dL/\epsilon$

Smooth strongly-convex: $T < dL/\mu \log(1/\epsilon)$

Misc

In finite dim, convex \Rightarrow cont. and difble almost everywhere.

If param e.g $\gamma := \frac{1}{L}$ unknown, use doubling trick.
Using line-searched step-size, can only do better than fixed step-size.

To prove SC \Rightarrow PL: min over \mathbf{y} in

$$f(\mathbf{y}) > f(\mathbf{x}) + g(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \mu/2 \|\mathbf{y} - \mathbf{x}\|^2$$

similarly for non- l_2 norms

similarly, can prove L -smooth \Rightarrow "lower-PL"

For any convex L -smooth f ,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \leq 2L (f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}))$$

(proof: L -smooth \Rightarrow "lower-PL" on tilted $h(\mathbf{x}) = f(\mathbf{x}) - g(\mathbf{y})^\top \mathbf{x}$)