

Today: Convergence rates

Assignment 5 out this weekend
 Marked Assignment 3 due Mon Oct 20
 No tutorial this week

Last week: Artificial Neural Networks (no longer convex)

instead of $\sum_{i=1}^N f_i(\bar{w}^T \bar{x}_i)$, we have $\sum_{i=1}^N f_i(\bar{w}_1^T g(\bar{w}_2^T \bar{x}_i))$ like a set of features
 or $\sum_{i=1}^N f_i(\bar{w}_1^T g(\bar{w}_2^T h(\bar{w}_3^T \bar{x}_i)))$ add more layers

Subgradient Questions

- Where do subgradients exist?
 - everywhere for convex functions!
- Why focus on gradient methods? $O(nd)$. Stochastic: $O(d)$
 - Newton, linprog, quadprog, etc. cost: $O(d^2)$ or $O(n^2)$ or larger
 (great for small problems, but modern problems can't apply them)

Convergence Rate of Gradient Descent

Gradient Method: $x^{t+1} = x^t - \alpha_t \nabla f(x^t)$

How fast does this converge? ($O(nd)$, but how many iterations?)

- to answer this, you need assumptions
 (function or gradient can't change too quickly)

We'll assume $\mu I \preceq \nabla^2 f(x) \preceq L I$
 ($\mu > 0$) ($L < \infty$)
 "strongly convex" "strongly smooth"
 \circ 0th order \circ 1st order
 $f(x) - \frac{\mu}{2} \|x\|^2$ is convex (regularizer) $\nabla f(x)$ is "L-Lipschitz" continuous: $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$

$\forall x, v, \nabla^2 f(x) v \leq v^T (L I) v$
 $= L (v^T I v) = L v^T v$
 $= L \|v\|^2$

Proposition: With $\alpha_t = 1/L$, the gradient method satisfies

$f(x^t) - f(x_*) \leq (1 - \frac{\mu}{L})^t [f(x^0) - f(x_*)]$

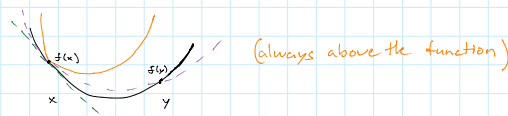


Proof: One version of the Taylor expansion

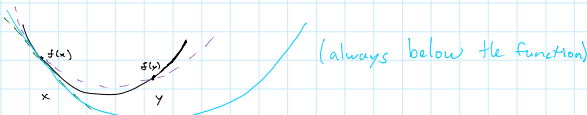
$\forall x, y, \exists z$ st. $f(y) = \underbrace{f(x) + \nabla f(x)^T (y-x)}_{\text{linearization}} + \underbrace{\frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x)}_{\text{Quadratic function of } y \text{ (Hessian in quadratic form)}}$



upper bound
 $f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} \|y-x\|^2$



lower bound
 $f(y) \geq f(x) + \nabla f(x)^T (y-x) - \frac{\mu}{2} \|y-x\|^2$



rearrange

Set $x = x^t, y = x^{t+1}$
 $f(x^{t+1}) \leq f(x^t) + \nabla f(x^t)^T (x^{t+1} - x^t) + \frac{L}{2} \|x^{t+1} - x^t\|^2$
 use $x^{t+1} - x^t = -\alpha_t \nabla f(x^t)$ $\alpha_t = 1/L$
 $f(x^{t+1}) \leq f(x^t) + \nabla f(x^t)^T (-\frac{1}{L} \nabla f(x^t)) + \frac{L}{2} \|\frac{1}{L} \nabla f(x^t)\|^2$
 $f(x^{t+1}) \leq f(x^t) - \frac{1}{2L} \|\nabla f(x^t)\|^2$ "guaranteed progress"

rearrange

Minimize in terms of y
 $(y = x - \frac{1}{L} \nabla f(x))$
 $f(x_*) \geq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2$
 rearrange (upper bound)
 $-\|\nabla f(x^t)\|^2 \leq -2L [f(x^t) - f(x_*)]$

Add $-f(x_*)$ to both sides of upper bound rearrange

$$f(x^{t+1}) - f(x_*) \leq f(x^t) - f(x_*) - \frac{1}{2L} \|\nabla f(x^t)\|^2$$

↓ using lower bound rearrange

$$f(x^{t+1}) - f(x_*) \leq f(x^t) - f(x_*) - \frac{2\mu}{2L} [f(x^t) - f(x_*)]$$

$$\leq \left(1 - \frac{\mu}{L}\right) [f(x^t) - f(x_*)]$$

$$f(x^t) - f(x_*) \leq \left(1 - \frac{\mu}{L}\right)^t [f(x^0) - f(x_*)]$$

$$\leq \left(1 - \frac{\mu}{L}\right) \left(1 - \frac{\mu}{L}\right) [f(x^{t-1}) - f(x_*)]$$

⋮ (keep doing this until we get x_0)

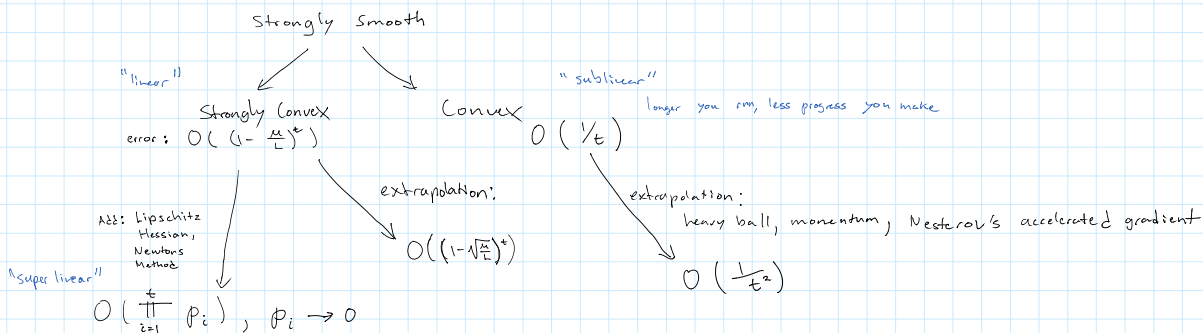
$$f(x^t) - f(x_*) \leq \left(1 - \frac{\mu}{L}\right)^t [f(x_0) - f(x_*)]$$

$$\left(1 - \frac{\mu}{L}\right)^t \leq \exp\left(-t \frac{\mu}{L}\right)$$

$$f(x^t) - f(x_*) \leq \epsilon \quad t = O(\log(\frac{1}{\epsilon}))$$

$$\text{total cost: } O(Nd \log(\frac{1}{\epsilon}))$$

Deterministic Gradient Complexity 200

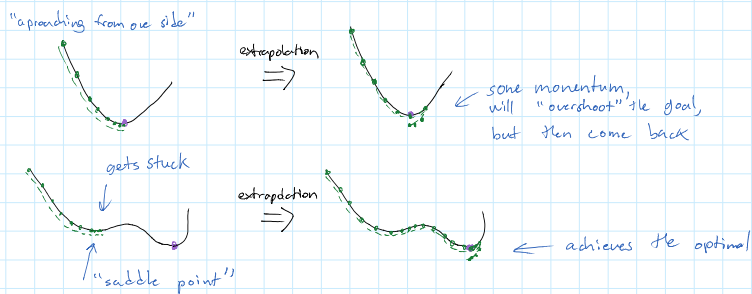


Approximate Newton:

- Barzilai - Borwein (limited memory)
- Quasi-Newton (L-BFGS)
- Hessian-free Newton

$$x^{t+1} = x^t - \alpha_t \nabla f(x_t) + \beta_t (x_t - x_{t-1})$$

- ⊗ Hic search
- important in practice
 - no effect on theory



Subgradient and/or Stochastic Gradient

x_* is not a fixed-point

$$x_* \neq x_* - \alpha g(x_*)$$

↖ subgradient

To get around this, need $\alpha_t \rightarrow 0$

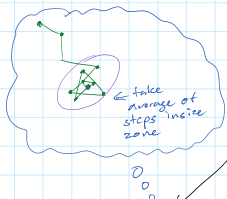
$$\text{Robbins - Monro } \sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} (\alpha_t)^2 < \infty$$

$$\text{Eg., } \alpha_t = \frac{1}{\sqrt{t}}$$

$$\text{gives } f(x^t) - f(x_*) = O\left(\frac{1}{\sqrt{t}}\right)$$

Subgradient / Stochastic Grad. Complexity Zoo



Lipschitz f (cant change too quickly)

Steady (max $O(\frac{\log(t)}{t})$)

Convex $O(\frac{1}{\sqrt{t}})$

finite, strongly smooth $O(\frac{1}{\sqrt{t}})$
SAG

Tail Averaging $O(\frac{1}{t})$

Strongly Smooth $O(\frac{1}{t})$

Smooth approximation $O(\frac{1}{\sqrt{t}})$

(Deterministic) extrapolation (exact gradient)

Constant Step-Size $O((1-2\alpha)^t) + O(\alpha)$

finite SAG $O(p^t)$

$O(\frac{1}{t})$

↑ additional error based off step size

