




THEORY OF MACHINE LEARNING

LECTURE 15

STRONG CONVEXITY, REGULARIZATION, STABILITY

Homework 2 is out - (zip file submissions) .
ok.



SUMMARY OF GRADIENT DESCENT

- Argmin $f(x)$, over $x \in D$, where D is a convex domain
- Simple iterative algorithm ("first order")
- Fixed step-size

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

(project onto domain if needed.)

- low variance.
- Convergence with error $O(\frac{1}{\sqrt{T}})$ after T steps for any L -Lipschitz function
 - "Noisy gradient oracle" \rightarrow stochastic gradient descent

$$T = 10^8$$

$$\bar{x}_t \quad f(\bar{x}_t) - f(x^*) \leq O\left(\frac{1}{\sqrt{T}}\right)$$

- Error of $O(1/T)$ for "smooth" convex functions (derivative is M -Lipschitz), assuming step size $< \frac{1}{2M}$

$$10^4$$

gradient.

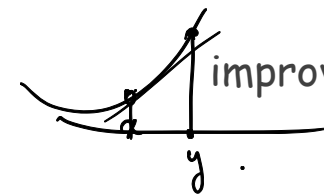
$$\text{target error} = 10^{-4}$$

- If function is also strongly convex with parameter μ , convergence bound

improves to roughly $\exp\left(-\left(\frac{\mu}{M}\right)T\right)$ (extends to Polyak-Lojasiewicz)

μ -strong convexity:

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



Parallelism & Complexity.

$$f: \mathbb{R}^d \rightarrow \mathbb{R}.$$

→ Optimization literature → d^2 → vector in \mathbb{R}^d gradients are "easy" but

$$\frac{\partial f}{\partial x_i}$$

Hessians are "complex".
→ $\mathbb{R}^{d \times d}$.

→ Parallelize grad. descent: → computing gradients can usually be parallelized.

$$\nabla f(w_{t+1}) = \nabla f(w_t - \eta \nabla f(w_t))$$

→ Second-order methods:

Beauzono - 2015.

[Mahoney..]

[Training] ^{imagined} $\log \log \left(\frac{1}{\epsilon} \right)$

$$f(x_{t+1}) - f(x^*) \leq \left(1 - \frac{1}{k}\right) \cdot (f(x_t) - f(x^*))$$

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

$n \times n$
 M
 $n \times k$.

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

$$\frac{1}{2} H^{-1} \nabla f(w_t)$$

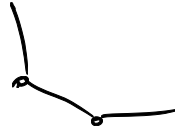
$$\rightarrow \begin{pmatrix} M^{-1} \end{pmatrix} \rightarrow n \times n$$

$$M^{-1} y$$

Strong convexity based results are good only if $\frac{M}{\mu}$ is "small".

$\hookrightarrow K$: Condition number.

"OPTIMAL" BOUNDS



- Turns out: under just the Lipschitz assumption, $\frac{1}{\sqrt{T}}$ cannot be improved, at least with "sub-gradient oracle"
- **Smoothness**: purely assuming smoothness, can get rate of $\frac{1}{T^2}$.
(Nesterov 1983), this is optimal for all "gradient based" methods
- Can we use information beyond the gradient?

10^{-2} iterations.
for err
 $= 10^{-4}$.

PRECONDITIONING

$$\nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_i \partial x_j} \end{pmatrix}$$

- Hessian plays informal role in most GD analyses (even M-smooth)

$$\|\nabla^2 f\| \leq M.$$

- "Directions" of Hessian can matter

$$\lambda^T \nabla^2 f \lambda \leq M \lambda^T \lambda$$

- Optimal movement using second order information

$$f(x_1, x_2, x_3) = \frac{1}{2} \left(x_1^2 + \frac{x_2^2}{4} + \frac{x_3^2}{9} \right)$$

$$\nabla f = \begin{pmatrix} x_1 \\ x_2/4 \\ x_3/9 \end{pmatrix}$$

Lipschitz const of ∇f is 2.

$$x_1^2 + (x_1 + x_2 - 2)^2 \dots$$

$$x^{(0)} = (1, 1, 1)$$

should do gradient descent with stepsize $< \frac{1}{4}$.

$$\nabla f(x^{(0)}) = \frac{1}{4} \begin{pmatrix} 2 \\ 1/2 \\ 2/9 \end{pmatrix}$$

$$x^{(1)} = x^{(0)} - \frac{1}{4} \cdot \nabla f(x^{(0)}) = \begin{pmatrix} 1/2 \\ 3/4 \\ 11/9 \end{pmatrix}$$

How much to move along diff. directions?



- what if we assume that function behaves like a quadratic in its neighborhood

$$f(x+\delta) \approx f(x) + \langle \delta, \nabla f(x) \rangle$$

$$f(x+\delta) = f(x) + \langle \delta, \nabla f(x) \rangle + \frac{1}{2} \delta^T \left(\nabla^2 f(x) \right) \cdot \delta + \dots$$

Hessian at x

$$f(x) + \delta(f'(x)) + \frac{\delta^2}{2!} (f''(x))$$

$$ax^2 + bx + c$$

$$\frac{-b}{2a}$$

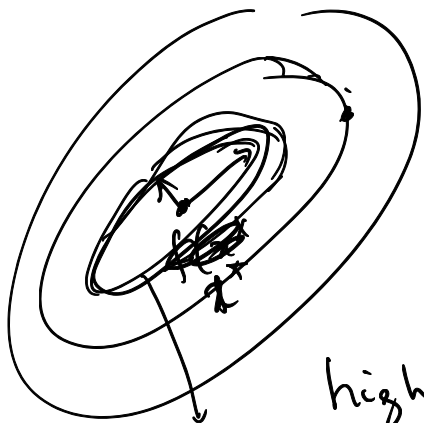
$$\delta = - \frac{f'(x)}{f''(x)}$$

$$\delta = - \left(\nabla^2 f(x) \right)^{-1} \nabla f(x)$$

(Newton's method)

high dim :

$$\{x : f(x) = f(x^*) + 0.1\}$$



IMPROVEMENTS, GENERALIZATIONS

for strongly convex f , instead of $e^{-\frac{1}{\mu}}$
 $e^{-1/\sqrt{\mu}}$.

- Polyak's "heavy ball" method (momentum)

- AdaGrad and related methods

- Second order (Newton) methods

- ...

$$\frac{1}{k} \left(\nabla f(x_1) \nabla f(x_1)^T + \dots + \dots \right)$$

"first order approx" to second order methods.

STRONG CONVEXITY, MOTIVATION

- Saw that strong convexity leads to "faster" optimization
- Additional benefit - "stability" to small perturbation
 - Example of quadratic

$f(x)$ is strongly convex $\rightarrow x^2$ 0

$g(x) = f(x) + \delta(x)$ $\rightarrow x^2 + \alpha x$ $-\frac{\alpha}{2}$

$\arg \min_x f(x)$ is "close" to $\arg \min_x g(x)$.

STABILITY OF LOSS MINIMIZATION

- Loss minimization with 'n' examples
- What happens if one example is "replaced"?

STABILITY IMPLIES GENERALIZATION

- Recall the notion of “generalization gap”
 - Can we phrase it in terms of stability?
- Stability versus “utility”!