

# Convex Optimization

**Prof. Nati Srebro**

## Lecture 4: Newton's Method

Reading: Boyd and Vandenberghe 9.5-9.7

# Gradient Descent with Backtracking Linesearch

```
Init       $x^{(0)} \in \text{dom}(f)$ 
Iterate    $\Delta x^{(k)} = -\nabla f(x^{(k)})$ 
          Stop if  $\|\nabla f(x^{(k)})\|_2^2 \leq 2\mu\epsilon$ 
          Set  $t^{(k)}$  by backtracking linesearch with params  $\alpha, \beta$ 
           $x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$ 
```

If  $f$  is  $\mu$ -strongly-convex and  $M$ -smooth on  $\{f(x) \leq f(x^{(0)})\}$ , then total #func evals and runtime:

$$O\left(\left\lceil \frac{\log M}{\log \frac{1}{\beta}} \right\rceil \frac{1}{\alpha} \max\left(\frac{1}{M}, \frac{1}{\beta}\right) \kappa \log\left(\frac{f(x^{(0)}) - p^*}{\epsilon}\right)\right)$$

$$\kappa = M/\mu$$

# Is Strong convexity necessary?

- Can we use Gradient Descent if  $f(x)$  is  $M$ -smooth, convex, but **not** strongly convex (or with very small  $\mu$ ) ?
  - Solution 1: optimize  $f_{\lambda}(x) \stackrel{\text{def}}{=} f(x) + \frac{\lambda}{2} \|x\|^2$ 
    - $\nabla^2 f_{\lambda}(x) = \nabla^2 f(x) + \lambda I$
    - $f_{\lambda}$  is  $(M + \lambda)$ -smooth and  $\lambda$ -strongly convex
    - $\rightarrow$  GD on  $f_{\lambda}$  finds  $f_{\lambda}(x^{(k)}) \leq \inf_x f_{\lambda}(x) + \epsilon$  with  $k = O\left(\frac{M+\lambda}{\lambda} \log \frac{1}{\epsilon}\right)$
- $f(x^{(k)}) \leq f_{\lambda}(x^{(k)}) \leq \inf_x f_{\lambda}(x) + \epsilon \leq f_{\lambda}(x^*) + \epsilon = f(x^*) + \frac{\lambda}{2} \|x^*\|^2 + \epsilon \leq p^* + 2\epsilon$
- Overall complexity:

$$O\left(\frac{M\|x^*\|_2^2}{\epsilon} \log \frac{1}{\epsilon}\right)$$

$$\lambda = \frac{2\epsilon}{\|x^*\|^2}$$

# Is Strong convexity necessary?

- Can we use Gradient Descent if  $f(x)$  is  $M$ -smooth, convex, but **not** strongly convex (or with very small  $\mu$ ) ?
- Solution 1: optimize  $f_{\lambda}(x) \stackrel{\text{def}}{=} f(x) + \frac{\lambda}{2} \|x\|^2$ 
  - Using  $\lambda = \frac{2\epsilon}{\|x^*\|^2}$ , #iter to find  $\epsilon$ -subopt:  $O\left(\frac{M\|x^*\|_2^2}{\epsilon} \log \frac{1}{\epsilon}\right)$
  - Can avoid log-factor by gradually decreasing  $\lambda$
- Solution 2: Run Gradient Descent directly on  $f(x)$ , with backtracking linesearch
  - Can find  $\epsilon$ -suboptimal solution with #iter
$$O\left(\frac{M\|x^* - x^{(0)}\|_2^2}{\epsilon}\right)$$

# Dependence on Conditioning

- Classical analysis depends on smoothness and strong convexity:

$$O\left(\frac{M}{\mu} \log\left(\frac{f(x^{(0)}) - p^*}{\epsilon}\right)\right)$$

- Linear dependence on  $M/\mu$  is real, and necessary in order to obtain linear convergence ( $\#\text{iter} \propto \log 1/\epsilon$ )
- Can always replace strong convexity  $\mu$  with dependence on  $\|x^*\|_2^2/\epsilon$

- All quantities depend on choice of norm  $\|\cdot\|_2$  and so on choice of basis:  
smoothness  $M$ , strong convexity  $\mu$ ,  $\|x^*\|_2$

$$f(x) + \langle \nabla f(x), \Delta x \rangle + \frac{\mu}{2} \|\Delta x\|_2^2 \leq f(x + \Delta x) \leq f(x) + \langle \nabla f(x), \Delta x \rangle + \frac{M}{2} \|\Delta x\|_2^2$$

$$\mu = \inf_{\|v\|_2=1} v^\top (\nabla^2 f) v \quad M = \sup_{\|v\|_2=1} v^\top (\nabla^2 f) v$$

# Pre-Conditioned Gradient Descent

Init	$x^{(0)} \in \text{dom}(f)$
	$H = \nabla^2 f(x^{(0)})$
Iterate	$\Delta x^{(k)} = -H^{-1} \nabla f(x^{(k)})$
	Set $t^{(k)}$ by backtracking linesearch with params $\alpha, \beta$
	$x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$

Complexity depends on conditioning w.r.t. the norm  $\|x\|_H = \sqrt{x^\top H x}$

$$\mu = \inf_{\|\nu\|_H=1} \nu^\top (\nabla^2 f(x)) \nu \qquad \qquad M = \sup_{\|\nu\|_H=1} \nu^\top (\nabla^2 f(x)) \nu$$

If the Hessian is fixed,  $\nabla^2 f(x) = \nabla^2 f(x^{(0)}) = H$ ,

$$\mu = \inf_{\|\nu\|_H=1} \nu^\top H \nu = 1 \qquad \qquad M = \sup_{\|\nu\|_H=1} \nu^\top H \nu = 1$$

# Pre-Conditioning a Quadratic

$$f(x) = \frac{1}{2}x^\top Hx + bx \quad H = \nabla^2 f(x^{(0)})$$

- The optimum  $x^*$  is given by:

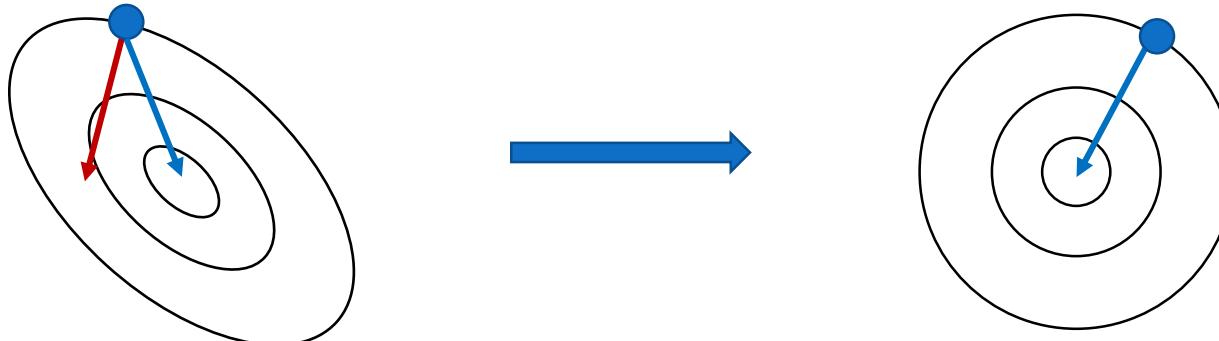
$$0 = \nabla f(x^*) = Hx^* + b \rightarrow x^* = -H^{-1}b$$

- Taking a single step of preconditioned GD:

$$\Delta x^{(0)} = -H^{-1}\nabla f(x^{(0)}) = -H^{-1}(Hx + b) = -(x + H^{-1}b)$$

- With a stepsize of  $t^{(0)} = 1$ :

$$x^{(1)} = x^{(0)} + \Delta x^{(0)} = x^{(0)} - (x^{(0)} + H^{-1}b) = -H^{-1}b = x^*$$



# Newton's Method

$$\min f(x)$$

Access to 1<sup>st</sup> and 2<sup>nd</sup> order oracles:  $x \mapsto f(x), \nabla f(x), \nabla^2 f(x)$

Init  $x^{(0)} \in \text{dom}(f)$

Iterate  $\Delta x^{(k)} = -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$

Set  $t^{(k)}$  by backtracking linesearch with params  $\alpha, \beta$

$x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$

- Affine invariant: consider  $\tilde{f}(\tilde{x}) = f(T\tilde{x} + b)$ , and running Newton on  $\tilde{f}$  and on  $f$ 
  - $\Delta x = T\Delta\tilde{x}$
  - If  $x^{(0)} = T\tilde{x}^{(0)}$  then  $x^{(k)} = T\tilde{x}^{(k)}$

# Newton's Method as Minimizing 2<sup>nd</sup> Order Approximation

$$f(x^{(k)} + \Delta x) \approx \underbrace{f(x^{(k)}) + \langle \nabla f(x^{(k)}), \Delta x \rangle + \frac{1}{2} \Delta x^\top \nabla^2 f(x^{(k)}) \Delta x}_{\hat{f}^{(k)}(x^{(k)} + \Delta x)}$$

$$\Delta x^{(k)} = \arg \min_{\Delta x} \hat{f}^{(k)}(x^{(k)} + \Delta x) = -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$$

With a step-size of  $t^{(k)} = 1$ :

$$x^{(k+1)} = \arg \min_x \hat{f}^{(k)}(x)$$

# Newton's Method

Init  $x^{(0)} \in \text{dom}(f)$

Iterate  $\Delta x^{(k)} = -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$

Set  $t^{(k)}$  by backtracking linesearch with params  $\alpha, \beta$

$x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$

Stopping condition?

Affine invariant?

# Newton's Decrement

$$\hat{f}^{(k)}(x^{(k)} + \Delta x) = f(x^{(k)}) + \langle \nabla f(x^{(k)}), \Delta x \rangle + \frac{1}{2} \Delta x^\top \nabla^2 f(x^{(k)}) \Delta x$$

- Suboptimality according to the quadratic approximation:

$$\begin{aligned} f^{(k)}(x^{(k)}) - \min_x \hat{f}^{(k)}(x) &= f(x^{(k)}) - \hat{f}^{(k)}(x^{(k)} + \Delta x^{(k)}) \\ &= f(x^{(k)}) - \left( f(x^{(k)}) + \langle \nabla f(x^{(k)}), -\nabla^2 f(x^{(k)}) \nabla f(x^{(k)}) \rangle + \frac{1}{2} \left( -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)}) \right)^\top \nabla^2 f(x^{(k)}) \left( -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)}) \right) \right) \\ &= \frac{1}{2} \nabla f(x^{(k)})^\top \nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)}) = \frac{1}{2} \lambda(x)^2 \end{aligned}$$

- Newton's Decrement:  $\lambda(x) \stackrel{\text{def}}{=} \sqrt{\nabla f(x)^\top \nabla^2 f(x)^{-1} \nabla f(x)}$

Stop if  $\frac{1}{2} \lambda(x)^2 \leq \epsilon$

- **Affine invariant!**

# Newton's Method

$$\min f(x)$$

Access to 1<sup>st</sup> and 2<sup>nd</sup> order oracles:  $x \mapsto f(x), \nabla f(x), \nabla^2 f(x)$

Init  $x^{(0)} \in \text{dom}(f)$

Iterate Stop if  $\frac{1}{2} \lambda(x^{(k)})^2 \leq \epsilon$

$$\Delta x^{(k)} = -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$$

Set  $t^{(k)}$  by backtracking linesearch with params  $\alpha, \beta$

$$x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$$

$$\lambda(x) \stackrel{\text{def}}{=} \sqrt{\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x)}$$

# Classical Analysis of Newton's Method

- Assumptions:
  - $\mu I \leq \nabla^2 f(x) \leq M I$
  - $\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L \|x - y\|_2$
- Part I (Damped Phase):  $\|\nabla f(x)\| \geq \eta \stackrel{\text{def}}{=} 3(1 - 2\alpha) \frac{\mu^2}{L}$ 
$$f(x^{(k)}) - f(x^{(k+1)}) \geq \gamma \stackrel{\text{def}}{=} \alpha \beta \eta^2 \frac{\mu}{M^2} \quad \rightarrow k \leq \frac{f(x^{(0)}) - p^*}{\gamma}$$
  - Backtracking is important! Might take short steps
- Part II (Quadratic Convergence):  $\|\nabla f(x)\| < \eta$ 
$$\|\nabla f(x^{(k+1)})\| \leq \frac{L}{2\mu^2} \|\nabla f(x^{(k)})\|^2$$
  - Stepsize  $t = 1$ 
$$\leq \frac{L}{2\mu^2} \eta = 3(1 - 2\alpha) < \frac{1}{2} \text{ when } \frac{1}{3} \leq \alpha \leq \frac{1}{2}$$
$$\rightarrow \|\nabla f(x^{(k+l)})\| \leq \underbrace{\left( \frac{L}{2\mu^2} \|\nabla f(x^{(k)})\| \right)}_{\leq \frac{L}{2\mu^2} \eta}^{(2^l)} \leq \frac{2\mu^2}{L} 2^{-2^l}$$
$$\rightarrow f(x^{(k+l)}) - p^* \leq \frac{1}{2\mu} \|\nabla f(x^{(k+l)})\|^2 < \frac{2\mu^3}{L^2} 2^{-2^l} \quad \rightarrow l < \log \log \frac{2\mu^3/L^2}{\epsilon}$$

# Classical Analysis of Newton's Method

- Assumptions:

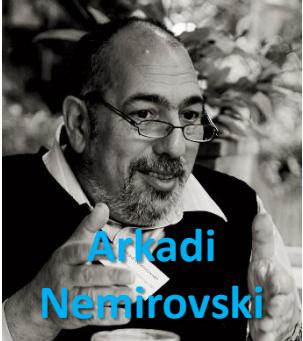
- $\mu I \leq \nabla^2 f(x) \leq M I$
- $\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L \|x - y\|_2$

Number of iterations to ensure  $f(x^{(k+l)}) \leq p^* + \epsilon$  using  $\alpha = 0.4$  and  $\beta = 0.9$ :

$$k + l \leq \frac{f(x^{(0)}) - p^*}{0.13 \frac{\mu^5}{L^2 M^2}} + \log \log \frac{2\mu^3/L^2}{\epsilon}$$

- #grad evals, #Hessian evals:  $k + l$
- #func evals:  $k \cdot (\#inner\ iter\ of\ backtracking\ linesearch) + l$
- Extra runtime:  $O((k + l)n^3)$

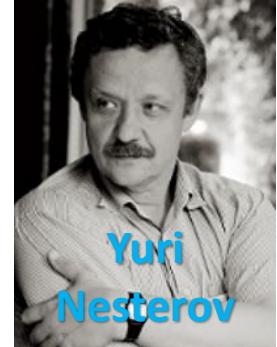
**Not affine invariant!**



Arkadi  
Nemirovski

# Self Concordant Analysis

## (Nemirovski and Nesterov 1994)



Yuri  
Nesterov

- Def:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is self concordant if

$$\forall x \forall \text{direction } v |f_v'''(x)| \leq 2f''(x)^{3/2}$$

- Examples:

- Linear ( $f_v'''(x) = 0$ )

- Quadratic ( $f_v'''(x) = 0$ )

- $f(x) = -\log x$

- $f'(x) = -\frac{1}{x}, f''(x) = \frac{1}{x^2}, f'''(x) = -\frac{2}{x^3}$

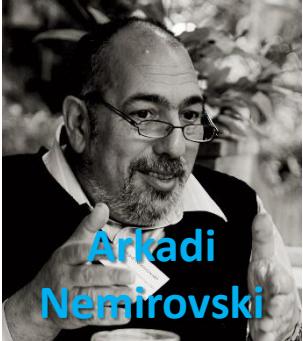
- What about  $f(x) = -\alpha \log x$  ?

- We need  $\frac{2\alpha}{x^3} \leq 2\frac{\alpha^{3/2}}{x^3} \rightarrow$  only if  $\alpha \geq 1$

- What about  $f(x) = e^x$ ?

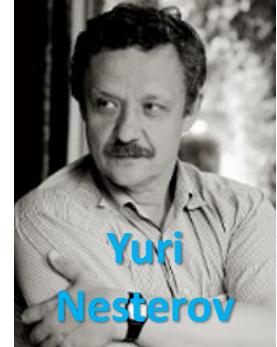
- $f'(x) = e^x, f''(x) = e^x, f'''(x) = e^x \rightarrow$  NOT self-concordant

- If  $f(x)$  is self-concordant, than so is  $\tilde{f}(x) = f(Tx + b)$ , for any affine transformation ( $T \in \mathbb{R}^{p \times n}, b \in \mathbb{R}^p$ )  
 $\rightarrow f(x) = \log(\langle a, x \rangle + a_0)$  is self concordant



Arkadi  
Nemirovski

# Self Concordant Analysis (Nemirovski and Nesterov 1994)



Yuri  
Nesterov

- Assumption:  $|f'''(x)| \leq 2f''(x)^{3/2}$
- Suboptimality **bounded** by Newton increment  $\lambda(x) = \sqrt{\nabla f(x)^\top \nabla^2 f(x)^{-1} \nabla f(x)}$   
$$f(x) \leq p^* + \lambda(x)^2$$
- Part I (Damped Phase):  $\lambda(x) \geq \eta \stackrel{\text{def}}{=} \frac{1}{4}(1 - 2\alpha)$   
$$f(x^{(k)}) - f(x^{(k+1)}) \geq \gamma \stackrel{\text{def}}{=} \alpha\beta \left(\frac{\eta^2}{1-\eta}\right)$$
  - Backtracking is important! Might take short steps
- Part II (Quadratic Convergence):  $\lambda(x) < \eta$   
$$\lambda(x^{(k+1)}) \leq 2\lambda(x^{(k)})^2 \quad \rightarrow f(x^{(k+l)}) \leq p^* + 2^{-2^l}$$
  - Stepsize  $t = 1$

# Newton's Method

$$\min f(x)$$

Access to 1<sup>st</sup> and 2<sup>nd</sup> order oracles:  $x \mapsto f(x), \nabla f(x), \nabla^2 f(x)$

Init  $x^{(0)} \in \text{dom}(f)$

Iterate Stop if  $\lambda(x^{(k)})^2 \leq \epsilon$

$$\Delta x^{(k)} = -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$$

Set  $t^{(k)}$  by backtracking linesearch with params  $\alpha, \beta$

$$x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$$

$$\lambda(x) \stackrel{\text{def}}{=} \sqrt{\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x)}$$

Theorem (NN'94): If  $f(x)$  is self-concordant, we have  $f(x^{(k+l)}) \leq p^* + \epsilon$  with

$$k + l \leq \frac{f(x^{(0)}) - p^*}{\gamma} + \log_2 \log_2 \frac{1}{\epsilon}$$

Depends only on linesearch params  
With  $\alpha = 0.15, \beta = 0.9$ :  $\gamma \approx 1/200$

# Newton for Non-Convex Functions

- Is the Newton direction  $\Delta x = -\nabla^2 f(x)^{-1} \nabla f(x)$  a descent direction?

$$\langle \nabla f(x), \Delta x \rangle = \langle \nabla f(x), -\nabla^2 f(x)^{-1} \nabla f(x) \rangle = -v^\top \nabla^2 f(x)^{-1} v < 0$$

$v = \nabla f(x)$

↑  
If  $f$  convex  
i.e.  $\nabla^2 f \geq 0$

- If  $f(x)$  is non-convex,  $\Delta x$  might not be a descent direction!

- Newton's method solves  $\nabla \hat{f}(x + \Delta x) = 0$ .
- $f(x)$  convex  $\rightarrow \nabla^2 f(x) \geq 0 \rightarrow \hat{f}(x)$  convex  $\rightarrow \nabla \hat{f} = 0$  is a minimum
- But if  $f(x)$  is not convex,  $x + \Delta x$  could be a saddle point or even maximum of  $\hat{f}$  !
- Newton (at least with  $t = 1$ ) converges to nearby critical point of  $f(x)$   
– local minima, local maxima or saddle point

# Newton vs Gradient Descent

## Gradient Descent

- Depends on choice of basis
- #iter  $\propto \kappa$
- Linear convergence,  
$$\propto \log \frac{1}{\epsilon}$$
- 1<sup>st</sup> order oracle
- Just vector ops—  
 $O(n)$  time per iter

## Newton

- Affine invariant
- Bad conditioning OK
- Quadratic convergence,  
$$\propto \log \log \frac{1}{\epsilon}$$
- 2<sup>nd</sup> order oracle
- Matrix inversion—  
 $O(n^3)$  per iter