



# 262B-Lecture 9

Date created: 2021.02.16  
N. of Pages: 14

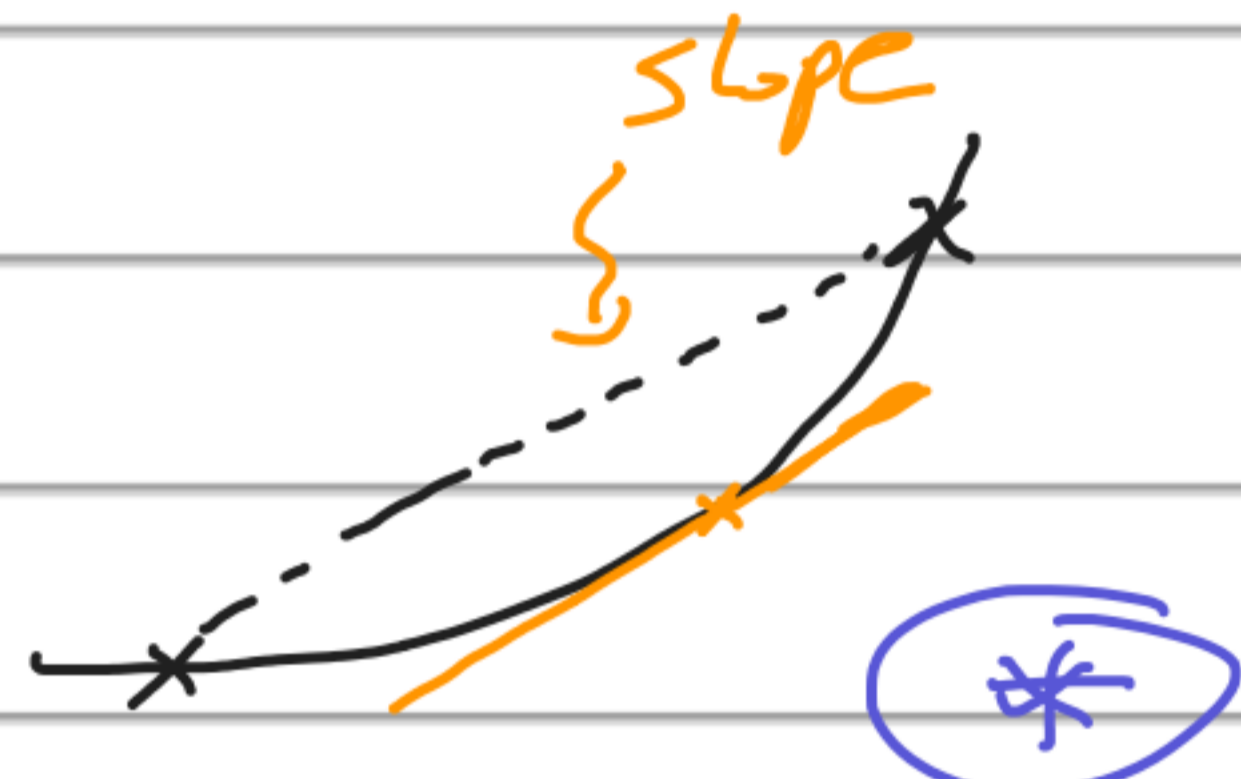
Newton's method:

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \underbrace{\nabla^2 f(x^{(k)})^{-1}}_{\text{need to have access to Hessian}} \nabla f(x^{(k)})$$

need to have access to Hessian

⇒ Can we approximate it?

Estimation method:



$$\nabla^2 f(x^{(k+1)}) (x^{(k+1)} - x^{(k)}) \approx \nabla f(x^{(k+1)}) - \nabla f(x^{(k)})$$

↓ new algorithm

$$f'(x) \sim \frac{f(x+\varepsilon) - f(x)}{(x+\varepsilon) - x}$$

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \underbrace{D^{(k)}}_{\text{Estimates inverse of Hessian}} \nabla f(x^{(k)})$$

Estimates inverse of Hessian

$$\underbrace{D^{(k+1)}}_{\text{matrix}} \underbrace{(\nabla f(x^{(k+1)}) - \nabla f(x^{(k)}))}_{\text{vector}} = \underbrace{x^{(k+1)} - x^{(k)}}_{\text{vector}}$$

If  $n=1 \Rightarrow$  we can solve for  $D^{(k+1)}$

If  $n > 1 \Rightarrow \underline{n}$  equations &  $\frac{n(n+1)}{2}$  variables in  $D^{(k+1)}$

So,  $D^{(k+1)}$  can't be uniquely found.

What if I find an optimal  $D^{(k+1)}$  ?

Fact :  $D^{(k+1)}$  can't be too far away

from  $D^{(k)}$   $\Rightarrow$  Smoothness in Hessian

$D^{(k+1)}$  : optimal solution to

$$\min_D \|D - D^{(k)}\|_Q \Rightarrow \text{weighted 2-norm}$$

$$\text{s.t. } D = D^T$$

$$D \left( \underbrace{\nabla f(x^{(k+1)}) - \nabla f(x^{(k)})}_{p^{(k)}} \right) = \underbrace{x^{(k+1)} - x^{(k)}}_{p^{(k)}}$$

$\Rightarrow$  This has a closed-form solution.

$\Rightarrow$   $D^{(k+1)}$  is a function of  $D^{(k)}$ .

$$\text{Definition: } q^{(k)} = \nabla f(x^{(k+1)}) - \nabla f(x^{(k)})$$

$$p^{(k)} = x^{(k+1)} - x^{(k)}$$

- Pick some arbitrary  $D^{(0)} > 0$

$$D^{(k+1)} = D^{(k)} + \left( 1 + \frac{(q^{(k)})^T D^{(k)} q^{(k)}}{(p^{(k)})^T q^{(k)}} \right) \frac{p^{(k)} (p^{(k)})^T}{(p^{(k)})^T p^{(k)}}$$

$$= \frac{D^{(k)} q^{(k)} (p^{(k)})^T + p^{(k)} (q^{(k)})^T D^{(k)}}{(p^{(k)})^T q^{(k)}}$$

BFGS rule : solution to previous

optimization for some  $\|\cdot\|$

Now :  $x^{(k+1)} = x^{(k)} - \alpha^{(k)} D^{(k)} \nabla f(x^{(k)})$

where  $D^{(k)}$  is obtained based

on BFGS  $\Rightarrow$  quasi Newton

Since finding  $D$  as an estimate of  $(\nabla^2 f)^{-1}$

wasn't unique, it is not the case that

$$\lim_{k \rightarrow \infty} D^{(k)} \rightarrow \lim_{k \rightarrow \infty} \nabla^2 f(x^{(k)})^{-1}$$

Since we started with  $x^{(k+1)} - x^{(k)}$  to estimate

Hessian, the method just worked through  
directional derivative.

$$\Rightarrow - \underbrace{D^{(k)}}_{\text{BFGS}} \nabla f(x^{(k)}) \longrightarrow \text{Newton's direction}$$

as  $k \rightarrow \infty$

(independent of the choice of  $D^{(0)}$ )

Convergence analysis if  $\nabla^2 f(x) \underset{>0}{\geq} \underbrace{m I}_{>0}$

Strongly convex case

gradient alg.  $\rightarrow$  linear convergence

- quasi-newton alg.  $\rightarrow$  superlinear  $\lim_{k \rightarrow \infty} \frac{e^{(k+1)}}{e^{(k)}} = 0$

newton alg.  $\rightarrow$  quadratic convergence



Conjugate gradient method:

so far, iteration complexity

↓

$$O(\log \log \frac{1}{\epsilon}), O(\log \frac{1}{\epsilon}), O(\frac{1}{\sqrt{\epsilon}}), O(\frac{1}{\epsilon}), O(\frac{1}{\epsilon^2})$$

This is about the number of iterations.

what is the role of  $n$ ?

Computational complexity = end-to-end complexity

$\sim$  number of iterations  $\times$  complexity per iteration  
e.g.  $O(\log \frac{1}{\epsilon})$  ?

Gradient alg. :  $\Delta x^{(k)} = -\nabla f(x^{(k)})$

Complexity =  $O(n)$  if complexity of evaluating each partial derivative is  $O(1)$ .

Newton's method:

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

Complexity =  $O(n^3)$  if you use

the naive method of taking the inverse

& doing the multiplication.

$\Rightarrow$  Can we improve this?

$$(*) \Rightarrow \underbrace{\nabla^2 f(x^{(k)})}_A \underbrace{\Delta x^{(k)}}_x = - \underbrace{\nabla f(x^{(k)})}_b$$

Find  $x$  satisfying  $Ax = b$

new method to solve equations.

Newton's method:  $A \succ 0$

General case:  $A \not\succeq 0$

$$Ax = b \Rightarrow \underbrace{A^T A}_{\text{new } A \succeq 0} x = \underbrace{A^T b}_{\text{new } b}$$

Let's focus on solving  $Qx_* = b$

This corresponds to solving

$$\min_x \frac{1}{2} x^T Q x - b^T x$$

Definition: A set of nonzero vectors

$d^{(1)}, d^{(2)}, \dots, d^{(k)}$  are called  $Q$ -conjugate

if  $(d^{(i)})^T Q (d^{(j)}) = 0 \quad \forall i, j \in \{1, \dots, k\}$   
 $i \neq j$

Thm:  $Q$ -conjugacy implies  $d^{(1)}, \dots, d^{(k)}$  are

linearly independent.

To prove by contradiction:

$$d^{(k)} = \underbrace{\alpha_1}_{\in \mathbb{R}} d^{(1)} + \underbrace{\alpha_2}_{\in \mathbb{R}} d^{(2)} + \dots + \underbrace{\alpha_{k-1}}_{\in \mathbb{R}} d^{(k-1)}$$

$$\Rightarrow \underbrace{(d^{(k)})^T Q d^{(k)}} = \sum_{i=1}^{k-1} \alpha_i \underbrace{(d^{(k)})^T Q d^{(i)}}_0$$

$$= 0 \Rightarrow \text{since } Q \succ 0 : d^{(k)} = 0 \quad \times$$



Devise a new algorithm:

- Consider  $n$  conjugate vectors

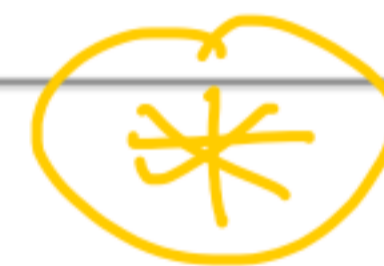
$d^{(0)}, d^{(1)}, \dots, d^{(n-1)} \Rightarrow$  treat them as directions

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)}, \quad k = 0, 1, \dots$$

where  $x^{(0)}$  is arbitrary and

$$\alpha^{(k)} \rightarrow \min_{\alpha} f(x^{(k)} + \alpha d^{(k)})$$

(positive or negative)



Focus on the quadratic case:

$$f(x) = \frac{1}{2} x^T Q x - b^T x$$

$$0 = \frac{\partial f(x^{(k)} + \alpha d^{(k)})}{\partial \alpha} \Big|_{\alpha = \alpha^{(k)}}$$

$$= (d^{(k)})^T \nabla f(x^{(k)} + \alpha^{(k)} d^{(k)})$$

$$\Rightarrow 0 = (d^{(k)})^T (Q(x^{(k)} + \alpha^{(k)} d^{(k)}) - b)$$

$$\Rightarrow \alpha^{(k)} = \frac{(d^{(k)})^T (b - Qx^{(k)})}{(d^{(k)})^T Q d^{(k)}}$$

> due to Q-conjugacy

Define:  $M^{(k)} = \{x \mid x = x^{(0)} + v \text{ where } v \in \text{subspace spanned by } d^{(0)}, \dots, d^{(k)}\}$

Make a space

by  $d^{(0)}, \dots, d^{(k)}$



shift it by  $x^{(0)}$

$$\Rightarrow \dim M^{(k)} = k + 1$$

$$\& \mathcal{M}^{(n-1)} = \mathbb{R}^n$$

$$\underbrace{M^{(0)}}_{1-d} \subset \underbrace{M^{(1)}}_{2-d} \subset \underbrace{M^{(2)}}_{3-d} \dots \subset M^{(n-1)}$$

Thm:  $x^{(k+1)}$  is solution to

$$\min_{x \in M^{(k)}} f(x)$$

$\forall k$

By product:  $k = n - 1$

$\Rightarrow x^{(n)}$  is solution to  $\min_{x \in M^{(n-1)}} f(x)$

$$= \min_{x \in \mathbb{R}^n} f(x)$$

$\Rightarrow$  This algorithm ( $x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)}$ )

stops after at most  $n$  steps.

Proof: 
$$\left. \frac{\partial f(x^{(i)} + \alpha d^{(i)})}{\partial \alpha} \right|_{\alpha = \alpha^{(i)}} = 0 \quad \forall i$$

$$\Rightarrow (d^{(i)})^T \nabla f(x^{(i+1)}) = 0 \quad \forall i \quad (*)$$

pick  $i \in \{0, \dots, k-1\}$ :

$$(d^{(i)})^T \nabla f(x^{(k+1)}) = (d^{(i)})^T (\underbrace{Q x^{(k+1)} - b}_{\text{...}})$$

$$\begin{aligned}
 x^{(k+1)} &\rightarrow \underbrace{x^{(k)} + \alpha^{(k)} d^{(k)}}_{\underbrace{x^{(k-1)} + \alpha^{(k-1)} d^{(k-1)}}} \\
 &\quad \vdots \\
 &\quad \underbrace{x^{(i+1)} + \alpha^{(i+1)} d^{(i+1)}}
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow (d^{(i)})^T \nabla f(x^{(k+1)}) &= (d^{(i)})^T \underbrace{(Qx^{(i+1)} - b)}_{\nabla f(x^{(i+1)})} \\
 + \sum_{j=i+1}^k (d^{(i)})^T Q d^{(j)} \alpha^{(j)} &= 0
 \end{aligned}$$

$$= (d^{(i)})^T \nabla f(x^{(i+1)}) = 0$$

due to  $\textcircled{*}$

$$\Rightarrow (d^{(i)})^T \nabla f(x^{(k+1)}) = 0 \quad i = 0, 1, \dots, k \quad \textcircled{**}$$

Note:  $x^{(k+1)} = x^{(0)} + \alpha^{(0)} d^{(0)} + \alpha^{(1)} d^{(1)} + \dots + \alpha^{(k)} d^{(k)}$

$$\Rightarrow \frac{\partial f(x^{(0)} + \beta_0 d^{(0)} + \beta_1 d^{(1)} + \dots + \beta_k d^{(k)})}{\partial \beta_i} \Big|_{\beta_i = \alpha^{(i)}} = 0$$

due to  $\textcircled{**}$   $\leftarrow i = 0, \dots, k$

$x^{(k+1)}$  is a solution to

$$\min_{\beta_0, \dots, \beta_k} f(x^{(0)} + \beta_0 d^{(0)} + \dots + \beta_k d^{(k)})$$

$$= \min_{x \in M^{(k)}} f(x)$$

$$x \in M^{(k)}$$

Conclusion: at every iteration the search space expands by  $\dim \underline{1}$ .

Generate Q-Conjugate directions:

Given linearly independent vectors  $v^{(0)}, \dots, v^{(n-1)}$

we can generate Q-Conjugate vectors  $d^{(0)}, \dots,$

$d^{(n-1)}$  such that

subspace spanned by  $d^{(0)}, \dots, d^{(i)}$   
= subspace spanned by  $v^{(0)}, \dots, v^{(i)}$  for  $i = 0, \dots, n-1$

$i=0 \Rightarrow$  pick  $d^{(0)} = v^{(0)}$

Have a constructive proof by induction.

Assume we have generated  $d^{(0)}, \dots, d^{(i)}$

$\rightarrow$  How to pick  $d^{(i+1)}$  ?

By induction step:

span of  $d^{(0)}, \dots, d^{(i)} = \text{span of } v^{(0)}, \dots, v^{(i)}$

aim to

span of  $d^{(0)}, \dots, d^{(i+1)} = \text{span of } v^{(0)}, \dots, v^{(i+1)}$

To achieve this, we write:

$$d^{(i+1)} = v^{(i+1)} + \beta_0 d^{(0)} + \beta_1 d^{(1)} + \dots + \beta_i d^{(i)}$$

For  $j=0, \dots, i$

$\Rightarrow$  to get  $Q$ -conjugacy

$$\Rightarrow (d^{(j)})^T Q d^{(i+1)} = (d^{(j)})^T Q v^{(i+1)} +$$

$$\beta_0 (d^{(j)})^T Q d^{(0)} + \dots + \beta_i (d^{(j)})^T Q d^{(i)}$$

$$= (d^{(j)})^T Q v^{(i+1)} + \beta_j (d^{(j)})^T Q d^{(j)}$$

$$\Rightarrow \beta_j = - \frac{(d^{(j)})^T \nabla v^{(i+1)}}{(d^{(j)})^T \nabla d^{(j)}}$$

$$\text{Now, } d^{(i+1)} = \nabla v^{(i+1)} + \beta_0 d^{(0)} + \dots + \beta_i d^{(i)}$$

with above coefficients will work.

Conjugate gradient method:

$$-\nabla f(x^{(0)}), -\nabla f(x^{(1)}), \dots, -\nabla f(x^{(n-1)})$$

⇓

generate Q-Conjugate directions

$$\Rightarrow x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)}$$