

# LECTURE 12: ROOT-FINDING AND MINIMIZATION

STAT 545: INTRO. TO COMPUTATIONAL STATISTICS

---

Vinayak Rao

Purdue University

October 16, 2018

Given some nonlinear function  $f: \mathbb{R} \rightarrow \mathbb{R}$ , solve

$$f(x) = 0$$

Invariably need iterative methods.

Assume  $f$  is continuous (else things are really messy).

More we know about  $f$  (e.g. gradients), better we can do.

Better: faster (asymptotic) convergence.

## ROOT BRACKETING

$f(a)$  and  $f(b)$  have opposite signs  $\rightarrow$  root lies in  $(a, b)$ .

$a$  and  $b$  *bracket* the root.

Finding an initial bracketing can be non-trivial.

Typically, start with an initial interval and expand or contract.

Below, we assume we have an initial bracketing.

## ROOT BRACKETING

$f(a)$  and  $f(b)$  have opposite signs  $\rightarrow$  root lies in  $(a, b)$ .

$a$  and  $b$  *bracket* the root.

Finding an initial bracketing can be non-trivial.

Typically, start with an initial interval and expand or contract.

Below, we assume we have an initial bracketing.

Not always possible e.g.  $f(x) = (x - a)^2$  (in general, multiple roots/nearby roots lead to trouble).

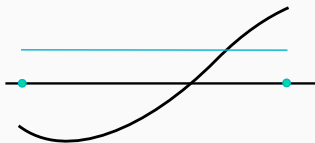
# BISECTION METHOD

Simplest root-finding algorithm.

Given an initial bracketing, cannot fail.

But is slower than other methods.

Successively halves the bracketing interval (binary search):



- Current interval =  $(a, b)$
- Set  $c = \frac{a+b}{2}$
- New interval =  $(a, c)$  or  $(c, b)$   
(whichever is a valid bracketing)

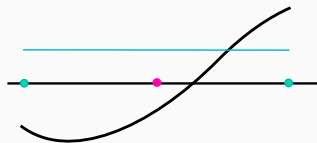
# BISECTION METHOD

Simplest root-finding algorithm.

Given an initial bracketing, cannot fail.

But is slower than other methods.

Successively halves the bracketing interval (binary search):



- Current interval =  $(a, b)$
- Set  $c = \frac{a+b}{2}$
- New interval =  $(a, c)$  or  $(c, b)$   
(whichever is a valid bracketing)

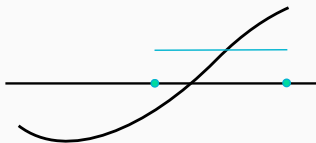
# BISECTION METHOD

Simplest root-finding algorithm.

Given an initial bracketing, cannot fail.

But is slower than other methods.

Successively halves the bracketing interval (binary search):



- Current interval =  $(a, b)$
- Set  $c = \frac{a+b}{2}$
- New interval =  $(a, c)$  or  $(c, b)$   
(whichever is a valid bracketing)

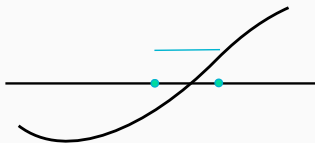
# BISECTION METHOD

Simplest root-finding algorithm.

Given an initial bracketing, cannot fail.

But is slower than other methods.

Successively halves the bracketing interval (binary search):



- Current interval =  $(a, b)$
- Set  $c = \frac{a+b}{2}$
- New interval =  $(a, c)$  or  $(c, b)$   
(whichever is a valid bracketing)



## BISECTION METHOD (CONTD)

Let  $\epsilon_n$  be the interval length at iteration  $n$ .

Upperbounds error in root.

$$\epsilon_{n+1} = 0.5 \epsilon_n \quad (\text{Linear convergence})$$

## BISECTION METHOD (CONTD)

Let  $\epsilon_n$  be the interval length at iteration  $n$ .

Upperbounds error in root.

$$\epsilon_{n+1} = 0.5 \epsilon_n \quad (\text{Linear convergence})$$

Linear convergence:

- each iteration reduces error by one significant figure.
- every (fixed)  $k$  iterations reduces error by one digit.
- error reduced exponentially with the number of iterations.

## BISECTION METHOD (CONTD)

Let  $\epsilon_n$  be the interval length at iteration  $n$ .

Upperbounds error in root.

$$\epsilon_{n+1} = 0.5 \epsilon_n \quad (\text{Linear convergence})$$

Linear convergence:

- each iteration reduces error by one significant figure.
- every (fixed)  $k$  iterations reduces error by one digit.
- error reduced exponentially with the number of iterations.

Superlinear convergence:

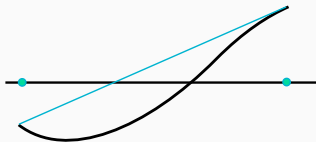
$$\lim_{n \rightarrow \infty} |\epsilon_{n+1}| = C \times |\epsilon_n|^m \quad (m > 1)$$

Quadratic convergence:

Number of significant figures *doubles* every iteration.

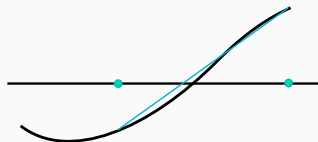
# SECANT METHOD AND BISECTION METHOD

Linearly approximate  $f$  to find new approximation to root.



# SECANT METHOD AND BISECTION METHOD

Linearly approximate  $f$  to find new approximation to root.



Secant method:

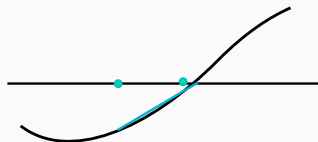
- always keep the newest point
- Superlinear convergence ( $m = 1.618$ , the golden ratio)

$$\lim_{n \rightarrow \infty} |\epsilon_{n+1}| = C \times |\epsilon_n|^{1.618}$$

- Bracketing (and thus convergence) not guaranteed.

# SECANT METHOD AND BISECTION METHOD

Linearly approximate  $f$  to find new approximation to root.



Secant method:

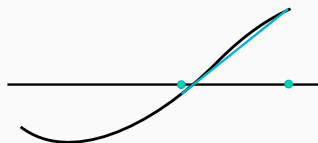
- always keep the newest point
- Superlinear convergence ( $m = 1.618$ , the golden ratio)

$$\lim_{n \rightarrow \infty} |\epsilon_{n+1}| = C \times |\epsilon_n|^{1.618}$$

- Bracketing (and thus convergence) not guaranteed.

# SECANT METHOD AND BISECTION METHOD

Linearly approximate  $f$  to find new approximation to root.



Secant method:

- always keep the newest point
- Superlinear convergence ( $m = 1.618$ , the golden ratio)

$$\lim_{n \rightarrow \infty} |\epsilon_{n+1}| = C \times |\epsilon_n|^{1.618}$$

- Bracketing (and thus convergence) not guaranteed.

False position:

- Can choose an old point that guarantees bracketing.
- Convergence analysis is harder.

In practice, people use more sophisticated algorithms.

Most popular is Brent's method.

Maintains bracketing by combining bisection method with a quadratic approximation.

Lots of book-keeping.



## NEWTON'S METHOD (A.K.A. NEWTON-RAPHSON)

At any point uses both function evaluation as well as derivative to form a linear approximation.

## NEWTON'S METHOD (A.K.A. NEWTON-RAPHSON)

At any point uses both function evaluation as well as derivative to form a linear approximation.

Taylor expansion:  $f(x + \delta) = f(x) + \delta f'(x) + \frac{\delta^2}{2} f''(x) + \dots$

# NEWTON'S METHOD (A.K.A. NEWTON-RAPHSON)

At any point uses both function evaluation as well as derivative to form a linear approximation.

Taylor expansion: 
$$f(x + \delta) = f(x) + \delta f'(x) + \frac{\delta^2}{2} f''(x) + \dots$$

Assume second- and higher-order terms are negligible.

Given  $x_i$ , choose  $x_{i+1} = x_i + \delta$  so that  $f(x_{i+1}) = 0$ :

# NEWTON'S METHOD (A.K.A. NEWTON-RAPHSON)

At any point uses both function evaluation as well as derivative to form a linear approximation.

Taylor expansion:  $f(x + \delta) = f(x) + \delta f'(x) + \frac{\delta^2}{2} f''(x) + \dots$

Assume second- and higher-order terms are negligible.

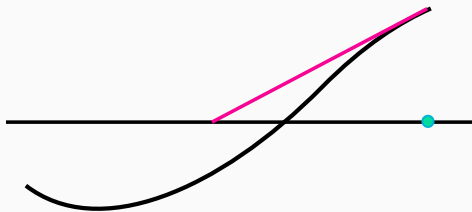
Given  $x_i$ , choose  $x_{i+1} = x_i + \delta$  so that  $f(x_{i+1}) = 0$ :

$$0 = f(x_i) + \delta f'(x_i)$$

$$x_{i+1} = x_i - f(x_i)/f'(x_i)$$

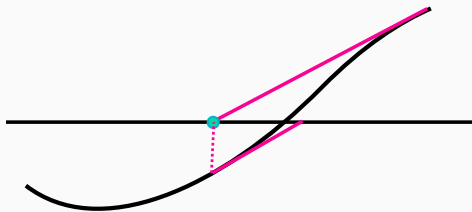
# NEWTON'S METHOD (A.K.A. NEWTON-RAPHSON)

$$x_{i+1} = x_i - f(x_i)/f'(x_i)$$



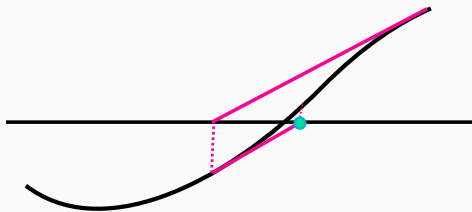
# NEWTON'S METHOD (A.K.A. NEWTON-RAPHSON)

$$x_{i+1} = x_i - f(x_i)/f'(x_i)$$



# NEWTON'S METHOD (A.K.A. NEWTON-RAPHSON)

$$x_{i+1} = x_i - f(x_i)/f'(x_i)$$



## CONVERGENCE OF NEWTON'S METHOD

Letting  $x$  be the root, we have

$$x_{i+1} - x = x_i - x - f(x_i)/f'(x_i)$$

$$\epsilon_{i+1} = \epsilon_i - f(x_i)/f'(x_i)$$



## CONVERGENCE OF NEWTON'S METHOD

Letting  $x$  be the root, we have

$$x_{i+1} - x = x_i - x - f(x_i)/f'(x_i)$$

$$\epsilon_{i+1} = \epsilon_i - f(x_i)/f'(x_i)$$

Also,

$$f(x_i) \approx f(x) + \epsilon_i f'(x) + \frac{\epsilon_i^2}{2} f''(x)$$

# CONVERGENCE OF NEWTON'S METHOD

Letting  $x$  be the root, we have

$$x_{i+1} - x = x_i - x - f(x_i)/f'(x_i)$$

$$\epsilon_{i+1} = \epsilon_i - f(x_i)/f'(x_i)$$

Also,

$$f(x_i) \approx f(x) + \epsilon_i f'(x) + \frac{\epsilon_i^2}{2} f''(x)$$

This gives

$$\epsilon_{i+1} = -\frac{f'(x_i)}{2f''(x_i)} \epsilon_i^2$$

# CONVERGENCE OF NEWTON'S METHOD

Letting  $x$  be the root, we have

$$x_{i+1} - x = x_i - x - f(x_i)/f'(x_i)$$

$$\epsilon_{i+1} = \epsilon_i - f(x_i)/f'(x_i)$$

Also,

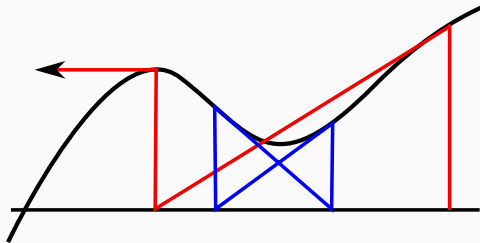
$$f(x_i) \approx f(x) + \epsilon_i f'(x) + \frac{\epsilon_i^2}{2} f''(x)$$

This gives

$$\epsilon_{i+1} = -\frac{f'(x_i)}{2f''(x_i)} \epsilon_i^2$$

Quadratic convergence (assuming  $f'(x)$  is non-zero at the root)

# PITFALLS OF NEWTON'S METHOD



Away from the root the linear approximation can be bad.

Can give crazy results (go off to infinity, cycles etc.)

However, once we have a decent solution can be used to rapidly 'polish the root'.

Often used in combination with some bracketing method.

Find  $(x_1, \dots, x_N)$  such that:

$$F_i(x_1, \dots, x_N) = 0 \quad i = 1 \text{ to } N$$

Much harder than the 1-d case.

Much harder than optimization.

Again, consider a Taylor expansion:

$$\mathbf{F}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{F}(\mathbf{x}) + \mathbf{J}(\mathbf{x}) \cdot \delta\mathbf{x} + O(\delta\mathbf{x}^2)$$

Here,  $\mathbf{J}(\mathbf{x})$  is the Jacobian matrix at  $\mathbf{x}$ , with  $J_{ij} = \frac{\partial F_i}{\partial x_j}$ .

## NEWTON'S METHOD

Again, consider a Taylor expansion:

$$\mathbf{F}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{F}(\mathbf{x}) + \mathbf{J}(\mathbf{x}) \cdot \delta\mathbf{x} + O(\delta\mathbf{x}^2)$$

Here,  $\mathbf{J}(\mathbf{x})$  is the Jacobian matrix at  $\mathbf{x}$ , with  $J_{ij} = \frac{\partial F_i}{\partial x_j}$ .

Again, Newton's method finds  $\delta\mathbf{x}$  by solving  $\mathbf{F}(\mathbf{x} + \delta\mathbf{x}) = 0$

$$\mathbf{J}(\mathbf{x}) \cdot \delta\mathbf{x} = -\mathbf{F}(\mathbf{x})$$

Solve e.g. by LU decomposition.

## NEWTON'S METHOD

Again, consider a Taylor expansion:

$$\mathbf{F}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{F}(\mathbf{x}) + \mathbf{J}(\mathbf{x}) \cdot \delta\mathbf{x} + O(\delta\mathbf{x}^2)$$

Here,  $\mathbf{J}(\mathbf{x})$  is the Jacobian matrix at  $\mathbf{x}$ , with  $J_{ij} = \frac{\partial F_i}{\partial x_j}$ .

Again, Newton's method finds  $\delta\mathbf{x}$  by solving  $\mathbf{F}(\mathbf{x} + \delta\mathbf{x}) = 0$

$$\mathbf{J}(\mathbf{x}) \cdot \delta\mathbf{x} = -\mathbf{F}(\mathbf{x})$$

Solve e.g. by LU decomposition.

Iterate  $\mathbf{x}_{new} = \mathbf{x}_{old} + \delta\mathbf{x}$  until convergence.

Can wildly careen through space if not careful.



Recall, we want to solve  $\mathbf{F}(\mathbf{x}) = 0$  ( $F_i(\mathbf{x}) = 0$ ,  $i = 1 \cdots N$ ).

Recall, we want to solve  $\mathbf{F}(\mathbf{x}) = 0$  ( $F_i(\mathbf{x}) = 0$ ,  $i = 1 \cdots N$ ).

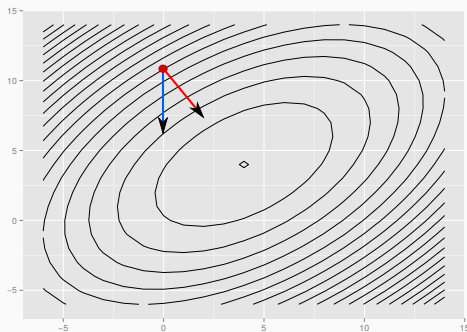
Minimize  $f(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^N |F_i(\mathbf{x})|^2 = \frac{1}{2} |\mathbf{F}(\mathbf{x})|^2 = \frac{1}{2} \mathbf{F}(\mathbf{x}) \cdot \mathbf{F}(\mathbf{x})$ .

Note: It is NOT sufficient to find a local minimum of  $f$ .

# GLOBAL METHODS VIA OPTIMIZATION)

We move along  $\delta \mathbf{x}$  instead of  $\nabla f = \mathbf{F}(\mathbf{x})\mathbf{J}(\mathbf{x})$ .

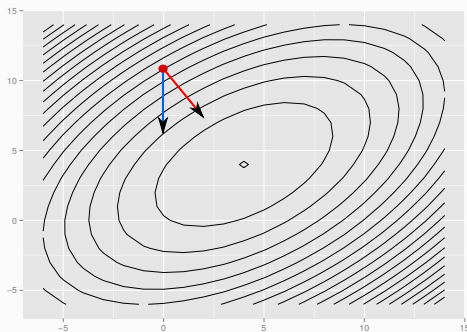
This keeps our global objective in sight.



# GLOBAL METHODS VIA OPTIMIZATION)

We move along  $\delta \mathbf{x}$  instead of  $\nabla f = \mathbf{F}(\mathbf{x})\mathbf{J}(\mathbf{x})$ .

This keeps our global objective in sight.



Note:  $\nabla f \cdot \delta \mathbf{x} = (\mathbf{F}(\mathbf{x})\mathbf{J}(\mathbf{x})) \cdot (-\mathbf{J}^{-1}(\mathbf{x})\mathbf{F}(\mathbf{x})) = -\mathbf{F}(\mathbf{x})\mathbf{F}(\mathbf{x}) < 0$

## NEWTON'S METHOD WITH BACKTRACKING

A full Newton step sets  $\mathbf{x}_{new} = \mathbf{x}_{old} + \delta\mathbf{x}$ .

This can cause  $f$  to increase i.e.  $f(\mathbf{x}_{new}) > f(\mathbf{x}_{old})$ .

In this case, *backtrack* and set  $\mathbf{x}_{new} = \mathbf{x}_{old} + \lambda\delta\mathbf{x}$ ,  $\lambda \in (0, 1)$ .

Since  $\delta\mathbf{x}$  is a descent direction, there exists a sufficiently small  $\lambda$  that causes  $f$  to decrease.

# NEWTON'S METHOD WITH BACKTRACKING

A full Newton step sets  $\mathbf{x}_{new} = \mathbf{x}_{old} + \delta\mathbf{x}$ .

This can cause  $f$  to increase i.e.  $f(\mathbf{x}_{new}) > f(\mathbf{x}_{old})$ .

In this case, *backtrack* and set  $\mathbf{x}_{new} = \mathbf{x}_{old} + \lambda\delta\mathbf{x}$ ,  $\lambda \in (0, 1)$ .

Since  $\delta\mathbf{x}$  is a descent direction, there exists a sufficiently small  $\lambda$  that causes  $f$  to decrease.

Finding best  $\lambda$ : too much work usually.

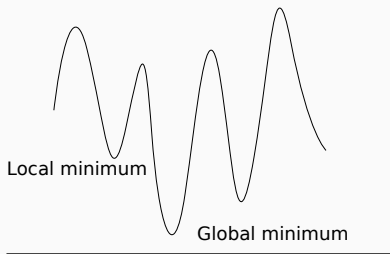
However, just causing  $f$  to decrease is not sufficient.

Can use Wolfe conditions (later)

# GLOBAL AND LOCAL MINIMUM

Find minimum of some function  $f: \mathbb{R}^D \rightarrow \mathbb{R}$ .  
(maximization is just minimizing  $-f$ ).

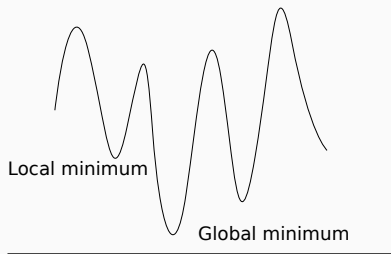
No global information (e.g. only function evaluations, derivatives).



# GLOBAL AND LOCAL MINIMUM

Find minimum of some function  $f: \mathbb{R}^D \rightarrow \mathbb{R}$ .  
(maximization is just minimizing  $-f$ ).

No global information (e.g. only function evaluations, derivatives).



Finding a global minimum is hard! Usually settle for finding a local minimum (like the EM algorithm).

Conceptually (deceptively?) simpler than EM.



## STEEPEST DESCENT (ITERATIVE METHOD)

Let  $x_{old}$  be our current value.

Update  $x_{new}$  as 
$$x_{new} = x_{old} - \eta \left. \frac{df}{dx} \right|_{x_{old}}$$

The steeper the slope, the bigger the move.

## STEEPEST DESCENT (ITERATIVE METHOD)

Let  $x_{old}$  be our current value.

Update  $x_{new}$  as 
$$x_{new} = x_{old} - \eta \left. \frac{df}{dx} \right|_{x_{old}}$$

The steeper the slope, the bigger the move.

$\eta$ : sometimes called the 'learning rate'  
(from neural network literature)

# STEEPEST DESCENT (ITERATIVE METHOD)

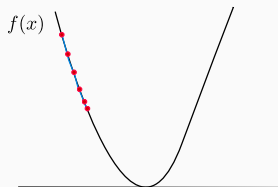
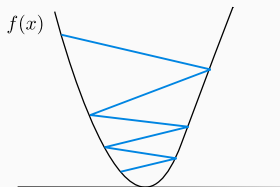
Let  $x_{old}$  be our current value.

Update  $x_{new}$  as 
$$x_{new} = x_{old} - \eta \left. \frac{df}{dx} \right|_{x_{old}}$$

The steeper the slope, the bigger the move.

$\eta$ : sometimes called the 'learning rate'  
(from neural network literature)

Choosing  $\eta$  is a dark art:



# STEEPEST DESCENT (ITERATIVE METHOD)

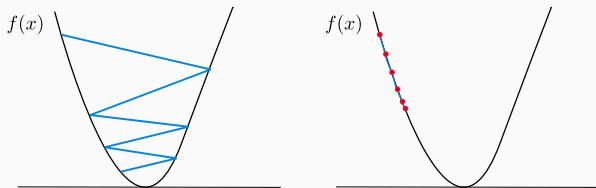
Let  $x_{old}$  be our current value.

Update  $x_{new}$  as 
$$x_{new} = x_{old} - \eta \left. \frac{df}{dx} \right|_{x_{old}}$$

The steeper the slope, the bigger the move.

$\eta$ : sometimes called the ‘learning rate’  
(from neural network literature)

Choosing  $\eta$  is a dark art:



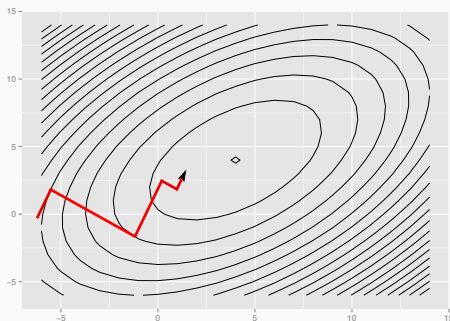
Better methods adapt step-size according to the curvature of  $f$ .

# STEEPEST DESCENT IN HIGHER-DIMENSIONS

Steepest descent also applies to higher dimensions too:

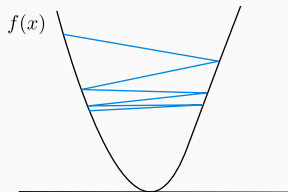
$$x_{new} = x_{old} - \eta \nabla f|_{x_{old}}$$

Now, even using the optimal  $\eta$  can be inefficient:

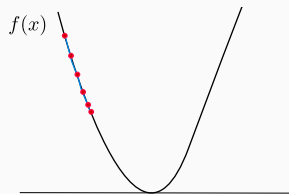


More on this later.

# WOLFE CONDITIONS

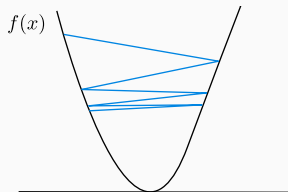


Big steps with little decrease

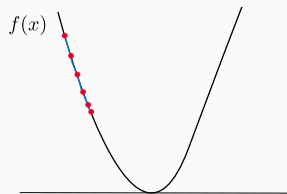


Small steps getting us  
nowhere

# WOLFE CONDITIONS



Big steps with little decrease

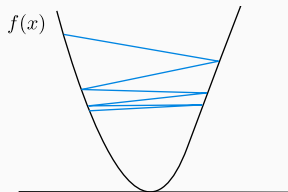


Small steps getting us  
nowhere

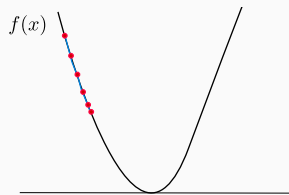
Avg. decrease at least some fraction of initial rate:

$$f(\mathbf{x} + \lambda \delta \mathbf{x}) \leq f(\mathbf{x}) + c_1 \lambda (\nabla f \cdot \delta \mathbf{x}), \quad c_1 \in (0, 1) \text{ e.g. } 0.9$$

# WOLFE CONDITIONS



Big steps with little decrease



Small steps getting us nowhere

Avg. decrease at least some fraction of initial rate:

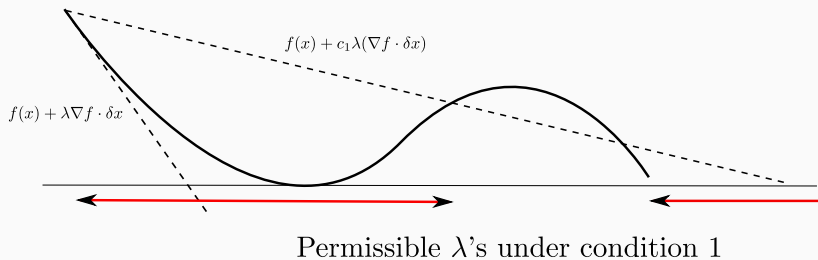
$$f(\mathbf{x} + \lambda \delta \mathbf{x}) \leq f(\mathbf{x}) + c_1 \lambda (\nabla f \cdot \delta \mathbf{x}), \quad c_1 \in (0, 1) \text{ e.g. } 0.9$$

Final rate is greater than some fraction of initial rate:

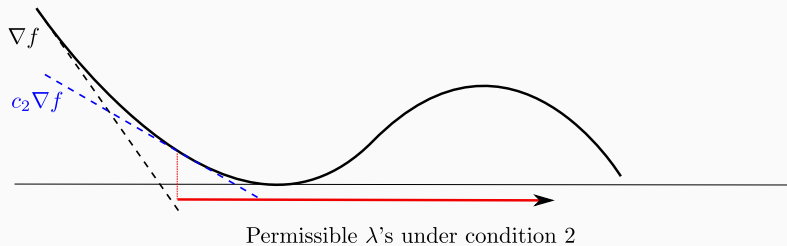
$$\nabla f(\mathbf{x} + \lambda \delta \mathbf{x}) \cdot \delta \mathbf{x} \geq c_2 \nabla f(\mathbf{x}) \delta \mathbf{x}, \quad c_2 \in (0, 1) \text{ e.g. } 0.1$$



# WOLFE CONDITIONS



# WOLFE CONDITIONS



A simple way to satisfy Wolfe conditions:

Set  $\delta x = -\nabla f$ ,  $c_1 = c_2 = .5$

Start with  $\lambda = 1$ , and while condition  $i$  is not satisfied, set  
 $\lambda = \beta_i t$  (for  $\beta_1 \in (0, 1)$ ,  $\beta_2 > 1$  and  $\beta_1 * \beta_2 < 1$ )

## CONVERGENCE ANALYSIS FOR FIXED STEP SIZE

Suppose  $f$  is convex and differentiable. Additionally, assume it is  $L$ -smooth.

$L$ -smoothness: gradients are Lipschitz continuous w. const.  $L$ :

$$\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\| \quad \forall x, y$$

Simple application of mean value theorem gives:

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{L}{2}\|y - x\|^2 \quad \forall x, y$$

(Upperbounded by a parabola for each  $x$ )

Theorem: Gradient descent with fixed step-size  $t \leq 1/L$  satisfies

$$f(x^{(k)}) - f(x^*) \leq \frac{\|x^{(0)} - x^*\|}{2tk}$$

This means it has convergence rate  $1/k$ , so that to get an error of  $\epsilon$  requires  $O(1/\epsilon)$  steps.

# STRONG CONVEXITY

For some  $d > 0$ ,  $\nabla^2 f(x) \geq dI$

Better lower bound (a parabola) than usual convexity.

Together with L-smoothness implies function lies between two quadratic functions.

Theorem: Gradient descent with fixed step-size  $t \leq 2/(d + L)$  satisfies

$$f(x^{(k)}) - f(x^*) \leq c^k \frac{L}{2} \|x^{(0)} - x^*\|^2, \quad c \in (0, 1)$$

This means it has convergence rate  $1/k$ , so that to get an error of  $\epsilon$  requires  $O(1/\epsilon)$  steps.

Error decays exponentially fast: to get an error of  $\epsilon$  requires  $O(\log(1/\epsilon))$  steps.

Called linear convergence

# ESTIMATING MLE

Consider a set of observations  $Y = (y_1, \dots, y_N)$ .

Assume  $y_i \sim p(y|\theta)$

$$\theta_{MLE} = \operatorname{argmax} \ell(\theta) = \operatorname{argmax} \sum_{i=1}^N \log p(x_i|\theta)$$

# ESTIMATING MLE

Consider a set of observations  $Y = (y_1, \dots, y_N)$ .

Assume  $y_i \sim p(y|\theta)$

$$\theta_{MLE} = \operatorname{argmax} \ell(\theta) = \operatorname{argmax} \sum_{i=1}^N \log p(x_i|\theta)$$

The gradient of the log-likelihood is  $\nabla \ell(\theta) = \sum_{i=1}^N \nabla \log p(x_i|\theta)$

The average of the gradients of each datapoint.



# ESTIMATING MLE

Consider a set of observations  $Y = (y_1, \dots, y_N)$ .

Assume  $y_i \sim p(y|\theta)$

$$\theta_{MLE} = \operatorname{argmax} \ell(\theta) = \operatorname{argmax} \sum_{i=1}^N \log p(x_i|\theta)$$

The gradient of the log-likelihood is  $\nabla \ell(\theta) = \sum_{i=1}^N \nabla \log p(x_i|\theta)$

The average of the gradients of each datapoint.

Starting with an initial  $\theta_0$ , iterate:

$$\theta_{i+1} = \theta_i + \eta_i \nabla \ell(\theta_i)$$

$$\nabla \ell(\theta) = \sum_{i=1}^N \nabla \log p(x_i|\theta)$$

Cons:

- Calculating the gradient is  $O(N)$ .  
(Each iteration must cycle through all datapoints.)
- *Lots* of redundancy, esp. for large  $N$ .

## GRADIENT DESCENT (CONTD.)

$$\nabla \ell(\theta) = \sum_{i=1}^N \nabla \log p(x_i|\theta)$$

Cons:

- Calculating the gradient is  $O(N)$ .  
(Each iteration must cycle through all datapoints.)
- *Lots* of redundancy, esp. for large  $N$ .

Pros:

- Convergence is better understood.
- Accelerated methods are available (e.g. Newton's method, conjugate gradient )

# STOCHASTIC GRADIENT DESCENT

Use a noisy gradient  $\widehat{\nabla}\ell$ .

Typically split data into  $N/B$  batches of size  $B$ .

Each iteration, calculate gradient on one of the batches  $B_i$ :

$$\widehat{\nabla}\ell(\theta) = \sum_{j \in B_i} \nabla \log p(x_j | \theta)$$

# STOCHASTIC GRADIENT DESCENT

Use a noisy gradient  $\widehat{\nabla}\ell$ .

Typically split data into  $N/B$  batches of size  $B$ .

Each iteration, calculate gradient on one of the batches  $B_i$ :

$$\widehat{\nabla}\ell(\theta) = \sum_{j \in B_i} \nabla \log p(x_j | \theta)$$

Pros:

- Calculating the gradient is  $O(B)$ .  
(Often, each batch is just a single datapoint)
- Much faster convergence (just one sweep through the data can get you a decent solution).
- Often, you get better solutions.
- Useful for online systems, tracking  $\theta$  that varies over time .

# STOCHASTIC GRADIENT DESCENT

Use a noisy gradient  $\widehat{\nabla}\ell$ .

Typically split data into  $N/B$  batches of size  $B$ .

Each iteration, calculate gradient on one of the batches  $B_i$ :

$$\widehat{\nabla}\ell(\theta) = \sum_{j \in B_i} \nabla \log p(x_j | \theta)$$

Cons:

- Convergence analysis is harder.
- Noisy gradients mean the algorithm will never converge.  
Typically need to reduce the step size every iteration.

We want

$$\eta_i \rightarrow 0, \quad \sum_{i=1}^{\infty} \eta_i = \infty$$

E.g.  $\eta_i = \frac{a}{b+i}$

One way to accelerate convergence is to include a momentum term:

$$\theta_{i+1} = \theta_i + \eta_i \nabla \ell(\theta_i) + \beta_i (\theta_i - \theta_{i-1})$$

More generally,

$$\theta_{i+1} = \theta_i + \eta_i \nabla \ell(\theta_i + \gamma(\theta_i - \theta_{i-1})) + \beta_i (\theta_i - \theta_{i-1})$$

Include many popular algorithms: Polyak's heavy ball method (HB):  $\gamma = 0$ , Nesterov's accelerated gradient (NAG):  $\gamma_i = \beta_i$

Another approach to accelerate convergence are adaptive methods that use the entire history of iterates to determine steps.