

# LECTURE 2: FUNCTION MINIMIZATION

STAT 545: INTRO. TO COMPUTATIONAL STATISTICS

---

Vinayak Rao

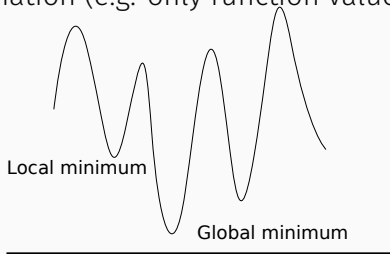
Purdue University

August 20, 2019

# GLOBAL AND LOCAL MINIMA

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

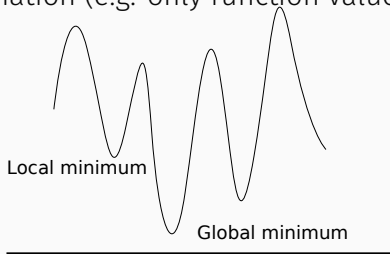
No global information (e.g. only function values, derivatives).



# GLOBAL AND LOCAL MINIMA

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

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



Finding global minima is hard! Usually settle for local minima.  
Even finding local minima is not easy. Usually need iterative algorithms. (Exceptions?)

## GRADIENT DESCENT (ITERATIVE METHOD)

Consider 1-d case. 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.

## GRADIENT DESCENT (ITERATIVE METHOD)

Consider 1-d case. 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.

$\eta$ : 'step-size' or 'learning rate'.

# GRADIENT DESCENT (ITERATIVE METHOD)

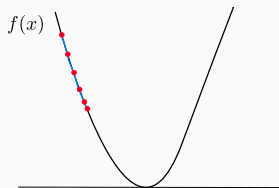
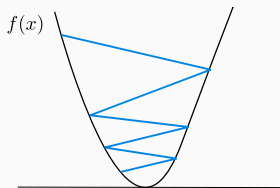
Consider 1-d case. Let  $x_{old}$  be our current value.

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

The steeper the slope, the bigger the move.

$\eta$ : 'step-size' or 'learning rate'.

Choosing  $\eta$  requires care (not too large or too small):



# GRADIENT DESCENT (ITERATIVE METHOD)

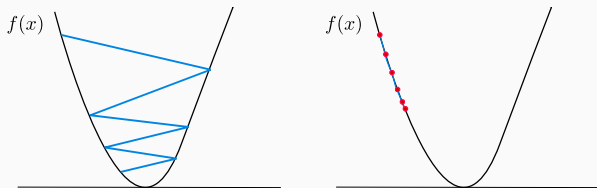
Consider 1-d case. Let  $x_{old}$  be our current value.

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

The steeper the slope, the bigger the move.

$\eta$ : 'step-size' or 'learning rate'.

Choosing  $\eta$  requires care (not too large or too small):



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

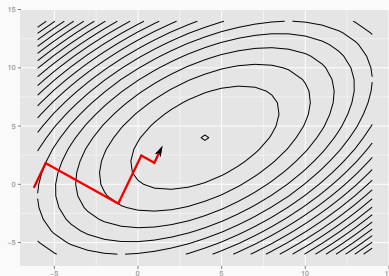
# GRADIENT DESCENT IN HIGHER-DIMENSIONS

Also applies to higher dimensions:  $x_{new} = x_{old} - \eta \nabla f|_{x_{old}}$

Again, need care choosing  $\eta$

Alternately, at each step, set  $\eta$  by minimizing along  $\nabla f$

• Note: even the optimal step-size  $\eta$  can be inefficient:

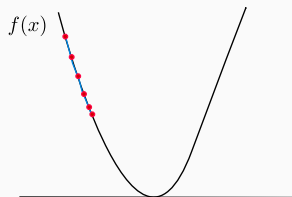
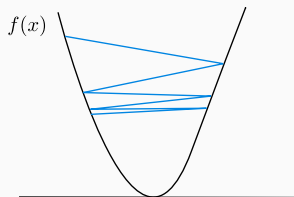


Rather than finding best step-size each step, save computation and find a decent solution



# WOLFE CONDITIONS TO DECIDE STEP-SIZE

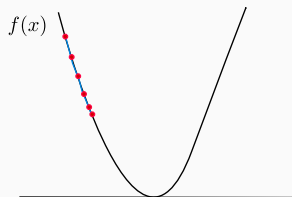
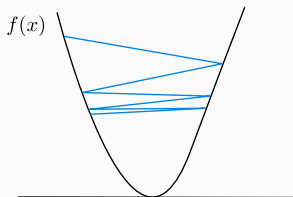
Bad step sizes along direction  $\mathbf{p}$  (for grad. descent,  $\mathbf{p} = -\nabla f$ ):



1) Big steps with little decrease 2) Small steps getting nowhere

# WOLFE CONDITIONS TO DECIDE STEP-SIZE

Bad step sizes along direction  $\mathbf{p}$  (for grad. descent,  $\mathbf{p} = -\nabla f$ ):



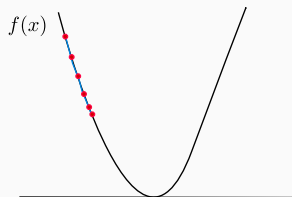
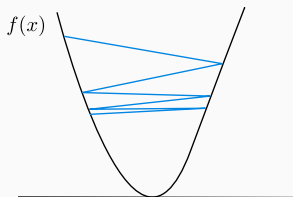
1) Big steps with little decrease 2) Small steps getting nowhere

Avoid (1): Avg. decrease at least some fraction of initial rate:

$$f(\mathbf{x} + \eta \mathbf{p}) \leq f(\mathbf{x}) + \eta c_1 (\nabla f \cdot \mathbf{p}), \quad c_1 \in (0, 1) \text{ e.g. } 0.1$$

# WOLFE CONDITIONS TO DECIDE STEP-SIZE

Bad step sizes along direction  $\mathbf{p}$  (for grad. descent,  $\mathbf{p} = -\nabla f$ ):



1) Big steps with little decrease 2) Small steps getting nowhere

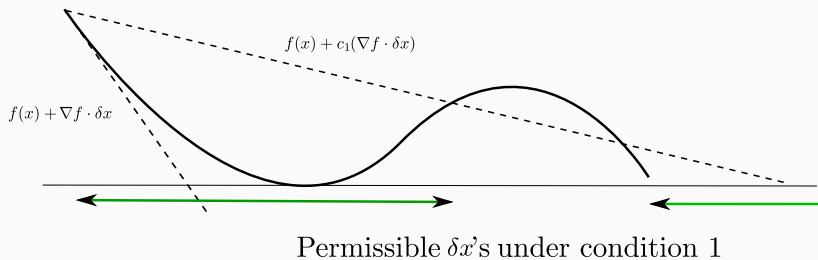
Avoid (1): Avg. decrease at least some fraction of initial rate:

$$f(\mathbf{x} + \eta \mathbf{p}) \leq f(\mathbf{x}) + \eta c_1 (\nabla f \cdot \mathbf{p}), \quad c_1 \in (0, 1) \text{ e.g. } 0.1$$

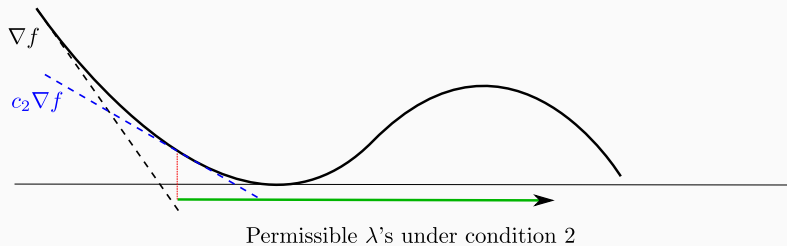
Avoid (2): Final rate is greater than some fraction of initial rate:

$$\nabla f(\mathbf{x} + \eta \mathbf{p}) \cdot \mathbf{p} \geq c_2 \nabla f(\mathbf{x}) \cdot \mathbf{p}, \quad c_2 \in (0, 1) \text{ e.g. } 0.9$$

# WOLFE CONDITIONS



# WOLFE CONDITIONS



A simple way to satisfy Wolfe conditions:

Set  $\mathbf{p} = -\nabla f$ ,  $c_1 = 0, 1$ ,  $c_2 = .9$

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

One way to understand/improve gradient descent is to view it as an approximation to 'gradient flow'.

Write  $\mathbf{x}_t$  for the position of a particle at time  $t$ , evolving as

$$\frac{d\mathbf{x}_t}{dt} = -\nabla f(\mathbf{x}_t), \quad \text{for some initialization at } t = 0.$$

$\mathbf{x}_t$  converges to  $\mathbf{x}^*$ , the minimum of  $f$  as  $t$  increases

- At minimum,  $\nabla f(\mathbf{x}^*) = 0$ .

Typically, not easy to solve the differential eq. for  $\mathbf{x}_t$

Different algs can be seen as approximations to this ideal

$$\frac{\Delta \mathbf{x}_t}{\Delta t} \approx \frac{d\mathbf{x}_t}{dt} = -\nabla f(\mathbf{x}_t) \quad (\text{forward Euler approximation})$$

$$\Rightarrow \mathbf{x}_{t+\Delta t} = \mathbf{x}_t - \Delta t \nabla f(\mathbf{x}_t)$$

This is just gradient descent with stepsize  $\eta = \Delta t$

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla f(\mathbf{x}_i)$$



## FORWARD AND BACKWARD METHODS

$$\frac{\Delta \mathbf{x}_t}{\Delta t} \approx \frac{d\mathbf{x}_t}{dt} = -\nabla f(\mathbf{x}_t) \quad (\text{forward Euler approximation})$$

$$\implies \mathbf{x}_{t+\Delta t} = \mathbf{x}_t - \Delta t \nabla f(\mathbf{x}_t)$$

This is just gradient descent with stepsize  $\eta = \Delta t$

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla f(\mathbf{x}_i)$$

$$\text{Backward Euler approx} \implies \mathbf{x}_{t+\Delta t} = \mathbf{x}_t - \Delta t \nabla f(\mathbf{x}_{t+\Delta t})$$

For a step size  $\eta$ , the iterates are:

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla f(\mathbf{x}_{i+1})$$

The updates are implicit ( $\mathbf{x}_{i+1}$  is on both LHS and RHS).

Why do we care?

## BACKWARD EULER METHOD

Backward method:  $\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla f(\mathbf{x}_{i+1})$

**Claim:** this is the same as solving

$$\mathbf{x}_{i+1} = \arg \min \mathbf{x} f(\mathbf{x}) + \frac{1}{2\eta} (\mathbf{x} - \mathbf{x}_i)^2$$

# BACKWARD EULER METHOD

Backward method:  $\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla f(\mathbf{x}_{i+1})$

**Claim:** this is the same as solving

$$\mathbf{x}_{i+1} = \arg \min \mathbf{x} f(\mathbf{x}) + \frac{1}{2\eta}(\mathbf{x} - \mathbf{x}_i)^2$$

Also called a proximal point method

# BACKWARD EULER METHOD

Backward method:  $\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla f(\mathbf{x}_{i+1})$

**Claim:** this is the same as solving

$$\mathbf{x}_{i+1} = \arg \min \mathbf{x} f(\mathbf{x}) + \frac{1}{2\eta}(\mathbf{x} - \mathbf{x}_i)^2$$

Also called a proximal point method

Now, we see that:

- $\mathbf{x}_{i+1} \leq \mathbf{x}_i$ , unlike gradient descent. Has faster convergence.
- This works even if  $\nabla f$  is not defined!
- Can be generalized to different distance functions:

$$\mathbf{x}_{i+1} = \arg \min \mathbf{x} f(\mathbf{x}) + \frac{1}{2\eta}d(\mathbf{x}, \mathbf{x}_i)$$

## OTHER APPROACHES TO IMPROVING CONVERGENCE

Newton's method: uses second derivatives/Hessians:

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

For vector-valued  $\mathbf{x}$ , writing  $\mathbf{H}f(\mathbf{x}_i)$  for the Hessian of  $f$  at  $\mathbf{x}_i$ ,

$$\mathbf{x}_{i+1} = \mathbf{x}_i - [\mathbf{H}f(\mathbf{x}_i)]^{-1} \nabla f(\mathbf{x}_i)$$

Intuition:

- Stepsize is small when gradient is changing rapidly
- Each iteration uses  $f(\mathbf{x}_i)$ ,  $\nabla f(\mathbf{x}_i)$  and  $\mathbf{H}f(\mathbf{x}_i)$  to construct a quadratic approximation to  $f$ , which is then minimized

MLE: maximum likelihood estimation

Consider a set of observations  $X = (x_1, \dots, x_N)$ .

Assume  $x_i \sim p(x|\theta)$

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

## BACK TO SIMPLE GRADIENT DESCENT

MLE: maximum likelihood estimation

Consider a set of observations  $X = (x_1, \dots, x_N)$ .

Assume  $x_i \sim p(x|\theta)$

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

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.)

## BACK TO SIMPLE GRADIENT DESCENT

MLE: maximum likelihood estimation

Consider a set of observations  $X = (x_1, \dots, x_N)$ .

Assume  $x_i \sim p(x|\theta)$

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

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 gradient requires evaluating likelihood  $N$  times. (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 gradient requires evaluating likelihood  $N$  times. (Each iteration must cycle through all datapoints.)
- *Lots* of redundancy, esp. for large  $N$ .

Pros:

- Convergence is better understood.

# 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:

$$\theta_{i+1} = \theta_i + \eta_i \widehat{\nabla} \ell(\theta_i) + \underbrace{\beta_i (\theta_i - \theta_{i-1})}_{\text{momentum}}$$

# STOCHASTIC GRADIENT DESCENT WITH MOMENTUM

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

$$\theta_{i+1} = \theta_i + \eta_i \widehat{\nabla} \ell(\theta_i) + \underbrace{\beta_i (\theta_i - \theta_{i-1})}_{\text{momentum}}$$

More generally,

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

Include many popular algorithms:

- Polyak's heavy ball method (HB):  $\gamma = 0$
- Nesterov's accelerated gradient (NAG):  $\gamma_i = \beta_i$

*Adaptive methods* accelerate convergence by using the entire history of iterates to determine step-sizes.



*Adaptive methods* accelerate convergence by using the entire history of iterates to determine step-sizes.

Often take the general form

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

where  $H_i$  is some combination of all previous gradients. E.g.

$$H_i = \text{diag} \left( \sum_{j=1}^i g_j \circ g_j \right),$$

with  $g_j = \widehat{\nabla} \ell(\theta_j + \gamma(\theta_j - \theta_{j-1}))$ , and  $\circ$  element-wise product.

Examples are AdaGrad, Adam etc.