

LECTURE 17: OVERVIEW OF OPTIMIZATION

STAT 598z: INTRODUCTION TO COMPUTING FOR STATISTICS

Vinayak Rao

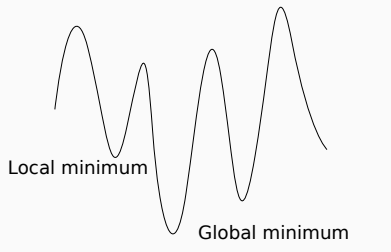
Department of Statistics, Purdue University

March 25, 2019

GLOBAL AND LOCAL MINIMUM

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

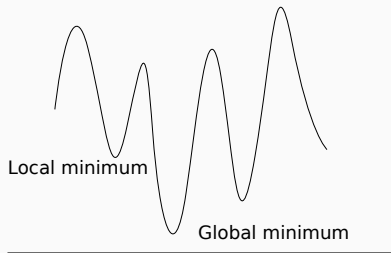
No global information (e.g. only function values, derivatives).



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 values, derivatives).



Finding a global minimum is hard!

We'll settle for a local minimum (maybe with multiple restarts).

ESTIMATING MLE

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

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

Maximum likelihood:

$$\theta_{MLE} = \operatorname{argmax} p(X|\theta) = \operatorname{argmax} \prod_{i=1}^N p(x_i|\theta)$$

More convenient to maximize the log-likelihood:

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

FIRST-ORDER CONDITIONS

The gradient $\nabla f = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right]^T$

FIRST-ORDER CONDITIONS

The gradient $\nabla f = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right]^T$

An infinitesimal step along $d\mathbf{x} = [dx_1, \dots, dx_N]$ gives a change

$$df = \nabla f \cdot d\mathbf{x}$$

FIRST-ORDER CONDITIONS

The gradient $\nabla f = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right]^\top$

An infinitesimal step along $d\mathbf{x} = [dx_1, \dots, dx_N]$ gives a change

$$df = \nabla f \cdot d\mathbf{x}$$

∇f : Direction of steepest ascent ($-\nabla f$ is steepest descent).

FIRST-ORDER CONDITIONS

The gradient $\nabla f = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right]^\top$

An infinitesimal step along $d\mathbf{x} = [dx_1, \dots, dx_N]$ gives a change

$$df = \nabla f \cdot d\mathbf{x}$$

∇f : Direction of steepest ascent ($-\nabla f$ is steepest descent).

At a local optimum $\nabla f = 0$.

FIRST-ORDER CONDITIONS

The gradient $\nabla f = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right]^\top$

An infinitesimal step along $d\mathbf{x} = [dx_1, \dots, dx_N]$ gives a change

$$df = \nabla f \cdot d\mathbf{x}$$

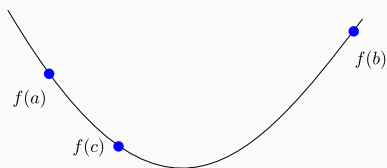
∇f : Direction of steepest ascent ($-\nabla f$ is steepest descent).

At a local optimum $\nabla f = 0$.

At a local minimum, Hessian $\nabla^2 f \succeq 0$ (positive semidefinite)

$$[\nabla^2 f]_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

ONE-DIMENSIONAL MINIMIZATION

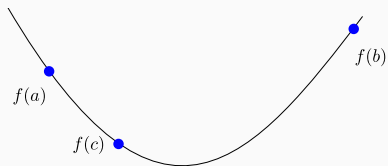


Find a bracket (a, b) with a third point $c \in (a, b)$, with

$$f(a) > f(c) < f(b)$$

Implies a local minimum lies in (a, b) .

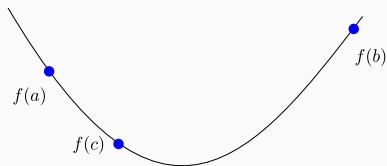
ONE-DIMENSIONAL MINIMIZATION



Finding an initial bracketing is not always easy. E.g.

- Pick two points l and r , $l < r$

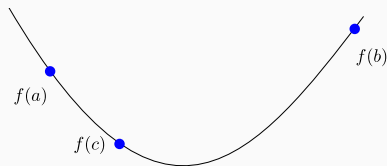
ONE-DIMENSIONAL MINIMIZATION



Finding an initial bracketing is not always easy. E.g.

- Pick two points l and r , $l < r$
- If $f(l) < f(r)$, $c = l$ and $b = r$, else $a = l$ and $c = r$.

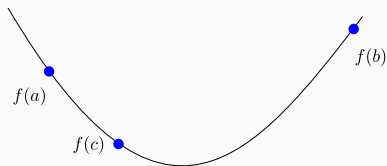
ONE-DIMENSIONAL MINIMIZATION



Finding an initial bracketing is not always easy. E.g.

- Pick two points l and r , $l < r$
- If $f(l) < f(r)$, $c = l$ and $b = r$, else $a = l$ and $c = r$.
- In the first case, choose $a < c$, and keep decreasing till $f(a) > f(c)$ (similarly with b for second case)

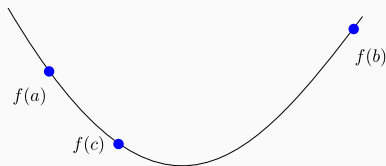
ONE-DIMENSIONAL MINIMIZATION



Having done this, successively refine (a, c) or (c, b) .

- Pick d in the longer interval (a, c) or (c, b) .

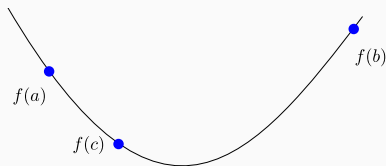
ONE-DIMENSIONAL MINIMIZATION



Having done this, successively refine (a, c) or (c, b) .

- Pick d in the longer interval (a, c) or (c, b) .
- Suppose it is (c, b) . Then either (a, c, d) and (c, d, b) forms a bracket.

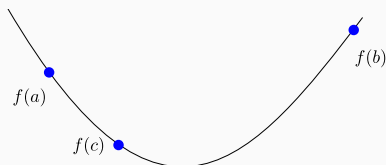
ONE-DIMENSIONAL MINIMIZATION



Having done this, successively refine (a, c) or (c, b) .

- Pick d in the longer interval (a, c) or (c, b) .
- Suppose it is (c, b) . Then either (a, c, d) and (c, d, b) forms a bracket.
- Choose and repeat

ONE-DIMENSIONAL MINIMIZATION



Having done this, successively refine (a, c) or (c, b) .

- Pick d in the longer interval (a, c) or (c, b) .
- Suppose it is (c, b) . Then either (a, c, d) and (c, d, b) forms a bracket.
- Choose and repeat

Doesn't extend easily to higher dimensions.

THE SIMPLEX ALGORITHM (NELDER & MEAD)

Find minimum of some function $f : \mathbb{R}^D \rightarrow \mathbb{R}$.

Requires only function evaluations.

Very general purpose, but not very efficient.

THE SIMPLEX ALGORITHM (NELDER & MEAD)

Find minimum of some function $f: \mathbb{R}^D \rightarrow \mathbb{R}$.

Requires only function evaluations.

Very general purpose, but not very efficient.

A 'simplex' in N -dimensions is the convex-hull of $N + 1$ points.
In 1-d: a line segment, 2-d: a triangle, 3-d: a tetrahedron etc.



THE SIMPLEX ALGORITHM (NELDER & MEAD)

Find minimum of some function $f: \mathbb{R}^D \rightarrow \mathbb{R}$.

Requires only function evaluations.

Very general purpose, but not very efficient.

A 'simplex' in N -dimensions is the convex-hull of $N + 1$ points.

In 1-d: a line segment, 2-d: a triangle, 3-d: a tetrahedron etc.



In 1-d we could bracket the minimum.

In higher dims, we must use other heuristics.

THE SIMPLEX ALGORITHM (NELDER & MEAD)

Start with an initial simplex.

Typically, pick an initial point \mathbf{P}_0 .

Also set $(\mathbf{P}_1, \dots, \mathbf{P}_{N+1})$ with $\mathbf{P}_i = \mathbf{P}_0 + \lambda_i \mathbf{e}_i$.

Here \mathbf{e}_i is the i th coordinate direction, and λ_i is the length-scale in that direction.

THE SIMPLEX ALGORITHM (NELDER & MEAD)

Start with an initial simplex.

Typically, pick an initial point \mathbf{P}_0 .

Also set $(\mathbf{P}_1, \dots, \mathbf{P}_{N+1})$ with $\mathbf{P}_i = \mathbf{P}_0 + \lambda_i \mathbf{e}_i$.

Here \mathbf{e}_i is the i th coordinate direction, and λ_i is the length-scale in that direction.

Assume $f(\mathbf{P}_0) \leq f(\mathbf{P}_1) \leq \dots \leq f(\mathbf{P}_{N+1})$.

THE SIMPLEX ALGORITHM (NELDER & MEAD)

Start with an initial simplex.

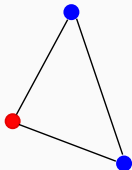
Typically, pick an initial point \mathbf{P}_0 .

Also set $(\mathbf{P}_1, \dots, \mathbf{P}_{N+1})$ with $\mathbf{P}_i = \mathbf{P}_0 + \lambda_i \mathbf{e}_i$.

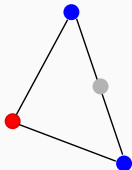
Here \mathbf{e}_i is the i th coordinate direction, and λ_i is the length-scale in that direction.

Assume $f(\mathbf{P}_0) \leq f(\mathbf{P}_1) \leq \dots \leq f(\mathbf{P}_{N+1})$.

At each step, try to improve the worst point \mathbf{P}_{N+1} using one of a sequence of moves.

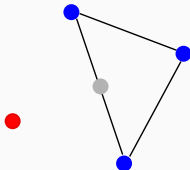


Get initial simplex



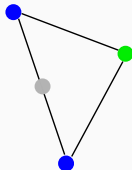
Find worst point, and find centroid of the remaining.

SIMPLEX ALGORITHM



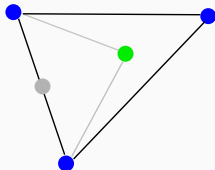
Reflect worst point.

This new point is now either worst, best or in the middle.

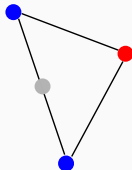


If this is the best point, extend.

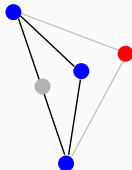
SIMPLEX ALGORITHM



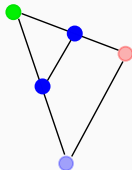
and go back to step one.



If this is the worst point, contract.



If this is the worst point, contract.



Else shrink all points except the best.

GRADIENT DESCENT

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

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

η : sometimes called the 'learning rate'

(terminology from the neural network literature)

GRADIENT DESCENT

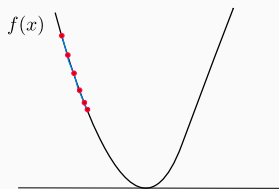
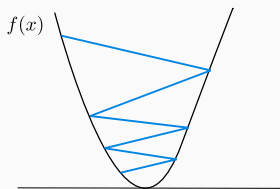
Let x_{old} be our current value

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

The steeper the slope, the bigger the move

η : sometimes called the 'learning rate'
(terminology from the neural network literature)

Choosing η is a dark art:



GRADIENT DESCENT

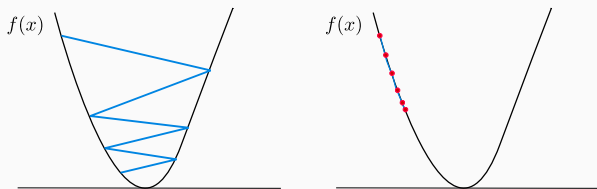
Let x_{old} be our current value

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

The steeper the slope, the bigger the move

η : sometimes called the 'learning rate'
(terminology from the neural network literature)

Choosing η is a dark art:

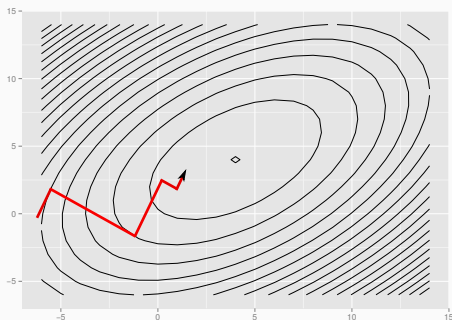


Better methods adapt step-size according to the curvature of f .

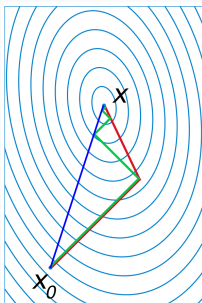
GRADIENT DESCENT IN HIGHER-DIMENSIONS

Gradient descent applies to higher dimensions too:

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



STEEPEST DESCENT



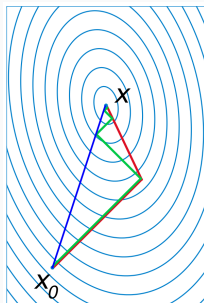
An any iteration, set \mathbf{p} to the direction of steepest descent.

$$\mathbf{p} = \nabla f(x_i)$$

Minimize along that direction: $\lambda_{min} = \operatorname{argmin}_{\lambda} f(\mathbf{x}_i + \lambda \mathbf{p})$

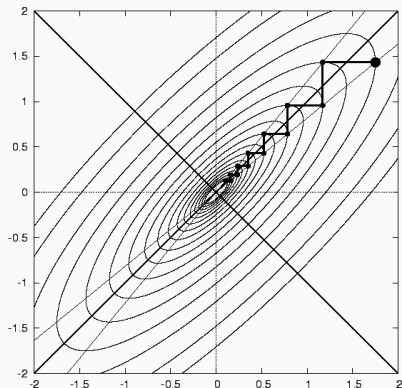
Set $\mathbf{x}_{i+1} = \mathbf{x}_i + \lambda_{min} \mathbf{p}$.

STEEPEST DESCENT



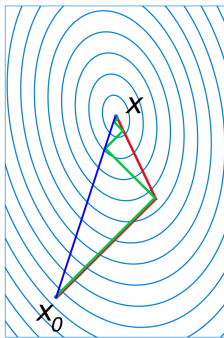
Can get trapped in long narrow valleys, where successive steps cancel each other.

STEEPEST DESCENT



Can get trapped in long narrow valleys, where successive steps cancel each other.

CONJUGATE DESCENT



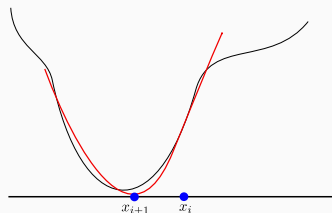
Conjugate gradient avoids moves along the same direction.

For a D -dim quadratic loss reaches minimum in D steps

Common default method

NEWTON'S METHOD

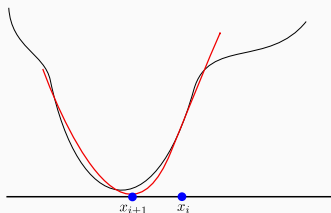
Uses the second derivative (curvature) to decide the step-size η .



At current point x_i , evaluate $f(x_i)$, $f'(x_i)$ and $f''(x_i)$.
Fit a parabola having these values and set x_{i+1} to its minimum.

NEWTON'S METHOD

Uses the second derivative (curvature) to decide the step-size η .

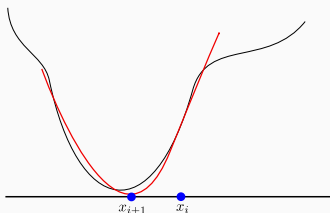


At current point x_i , evaluate $f(x_i)$, $f'(x_i)$ and $f''(x_i)$.
Fit a parabola having these values and set x_{i+1} to its minimum.
Easy to see that (show it!):

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

NEWTON'S METHOD

Uses the second derivative (curvature) to decide the step-size η .



At current point x_i , evaluate $f(x_i)$, $f'(x_i)$ and $f''(x_i)$.
Fit a parabola having these values and set x_{i+1} to its minimum.
Easy to see that (show it!):

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

If f'' is large, we're uncertain about f' , so take a small step.

Update rule:

$$\mathbf{x}_{i+1} = \mathbf{x}_i - [\nabla^2 f(\mathbf{x}_i)]^{-1} \nabla f(\mathbf{x}_i)$$

Update rule:

$$\mathbf{x}_{i+1} = \mathbf{x}_i - [\nabla^2 f(\mathbf{x}_i)]^{-1} \nabla f(\mathbf{x}_i)$$

Need to calculate the Hessian $\nabla^2 f$: N^2 elements.

Need to invert the Hessian: N^3 operations.

NEWTON'S METHOD IN HIGHER DIMENSIONS

Update rule:

$$\mathbf{x}_{i+1} = \mathbf{x}_i - [\nabla^2 f(\mathbf{x}_i)]^{-1} \nabla f(\mathbf{x}_i)$$

Need to calculate the Hessian $\nabla^2 f$: N^2 elements.

Need to invert the Hessian: N^3 operations.

Each iteration can be expensive.

Quasi-Newton methods try to alleviate this issue.

NEWTON'S METHOD IN HIGHER DIMENSIONS

Update rule:

$$\mathbf{x}_{i+1} = \mathbf{x}_i - [\nabla^2 f(\mathbf{x}_i)]^{-1} \nabla f(\mathbf{x}_i)$$

Need to calculate the Hessian $\nabla^2 f$: N^2 elements.

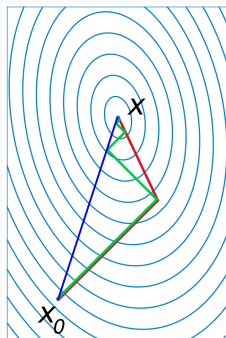
Need to invert the Hessian: N^3 operations.

Each iteration can be expensive.

Quasi-Newton methods try to alleviate this issue.

We have to be wary about taking wild steps.

NEWTON'S DESCENT



Set \mathbf{p} to the minimum of the local quadratic approximation.

$$\mathbf{p} = [\nabla^2 f(\mathbf{x})]^{-1} \nabla f(x_j)$$

Reaches minimum of quadratic loss in 1 step

QUASI-NEWTON METHODS

Newton's method: $\mathbf{p} = [\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x}_i)$

Steepest's descent: $\mathbf{p} = - \nabla f(\mathbf{x}_i)$

QUASI-NEWTON METHODS

Newton's method: $\mathbf{p} = [\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x}_i)$

Steepest's descent: $\mathbf{p} = - \nabla f(\mathbf{x}_i)$

Quasi-Newton methods use other matrices B :

$$\mathbf{p} = -B \nabla f(\mathbf{x}_i)$$

QUASI-NEWTON METHODS

Newton's method: $\mathbf{p} = [\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x}_i)$

Steepest's descent: $\mathbf{p} = - \nabla f(\mathbf{x}_i)$

Quasi-Newton methods use other matrices \mathbf{B} :

$$\mathbf{p} = \mathbf{B} \nabla f(\mathbf{x}_i)$$

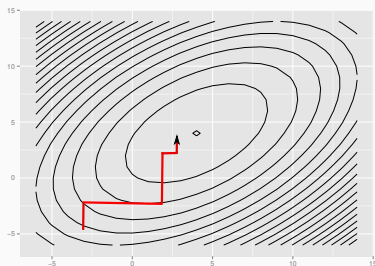
Usually, \mathbf{B} is allowed to vary from iteration to iteration, with

$$\mathbf{B}_i \rightarrow [\nabla^2 f(\mathbf{x})]^{-1}$$

Get benefits of Newton's method, without $O(N^3)$ computations.

E.g. BFGS

CO-ORDINATE DESCENT



Saw this last lecture

Simple, clean and inexpensive.

Often the 1-d problems can be solved exactly.

Convergence can be slow.

Exception: axis aligned ellipses need just D steps.

Use the `optim` function

Syntax:

```
optim(par, fn, gr = NULL, ...,  
      method = c('Nelder-Mead', 'BFGS', 'CG', 'L-BFGS-B', 'SANN',  
                'Brent'),  
      lower = -Inf, upper = Inf,  
      control = list(), hessian = FALSE)
```

`fn`: function to be optimized

`gr`: gradient function (calculate numerically if `NULL`)

`par`: initial value of parameter to be optimized (should be first argument of `fn`)