

Basics of Numerical Optimization: Iterative Methods

Ju Sun

Computer Science & Engineering
University of Minnesota, Twin Cities

February 13, 2020

Find global minimum

1st-order necessary condition: Assume f is 1st-order differentiable at x_0 .
If x_0 is a local minimizer, then $\nabla f(x_0) = \mathbf{0}$.

x with $\nabla f(x) = \mathbf{0}$: **1st-order stationary point** (1OSP)

Find global minimum

1st-order necessary condition: Assume f is 1st-order differentiable at x_0 .
If x_0 is a local minimizer, then $\nabla f(x_0) = \mathbf{0}$.

x with $\nabla f(x) = \mathbf{0}$: **1st-order stationary point (1OSP)**

2nd-order necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$.

x with $\nabla f(x) = \mathbf{0}$ and $\nabla^2 f(x) \succeq \mathbf{0}$: **2nd-order stationary point (2OSP)**

Find global minimum

1st-order necessary condition: Assume f is 1st-order differentiable at x_0 .
If x_0 is a local minimizer, then $\nabla f(x_0) = \mathbf{0}$.

x with $\nabla f(x) = \mathbf{0}$: **1st-order stationary point (1OSP)**

2nd-order necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$.

x with $\nabla f(x) = \mathbf{0}$ and $\nabla^2 f(x) \succeq \mathbf{0}$: **2nd-order stationary point (2OSP)**

- **Analytic method:** find 1OSP's using gradient first, then study them using Hessian — for simple functions! e.g., $f(x) = \|\mathbf{y} - \mathbf{A}x\|_2^2$, or $f(x, y) = x^2y^2 - x^3y + y^2 - 1$

Find global minimum

1st-order necessary condition: Assume f is 1st-order differentiable at x_0 .
If x_0 is a local minimizer, then $\nabla f(x_0) = \mathbf{0}$.

x with $\nabla f(x) = \mathbf{0}$: **1st-order stationary point (1OSP)**

2nd-order necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$.

x with $\nabla f(x) = \mathbf{0}$ and $\nabla^2 f(x) \succeq \mathbf{0}$: **2nd-order stationary point (2OSP)**

- **Analytic method:** find 1OSP's using gradient first, then study them using Hessian — for simple functions! e.g., $f(x) = \|\mathbf{y} - \mathbf{A}x\|_2^2$, or $f(x, y) = x^2y^2 - x^3y + y^2 - 1$
- **Grid search:** incurs $O(\varepsilon^{-n})$ cost

Find global minimum

1st-order necessary condition: Assume f is 1st-order differentiable at x_0 .
If x_0 is a local minimizer, then $\nabla f(x_0) = \mathbf{0}$.

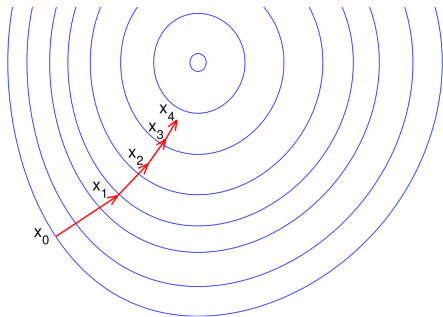
x with $\nabla f(x) = \mathbf{0}$: **1st-order stationary point** (1OSP)

2nd-order necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$.

x with $\nabla f(x) = \mathbf{0}$ and $\nabla^2 f(x) \succeq \mathbf{0}$: **2nd-order stationary point** (2OSP)

- **Analytic method:** find 1OSP's using gradient first, then study them using Hessian — for simple functions! e.g., $f(x) = \|y - Ax\|_2^2$, or $f(x, y) = x^2y^2 - x^3y + y^2 - 1$
- **Grid search:** incurs $O(\epsilon^{-n})$ cost
- **Iterative methods:** find 1OSP's/2OSP's by making consecutive small movements

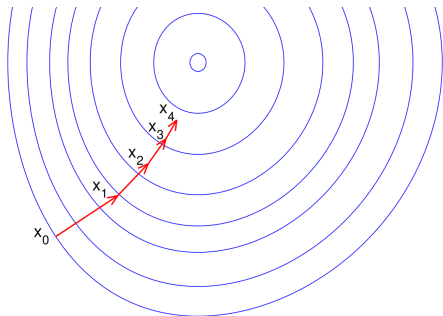
Iterative methods



Credit: aria42.com

Illustration of iterative methods on the contour/levelset plot (i.e., the function assumes the same value on each curve)

Iterative methods

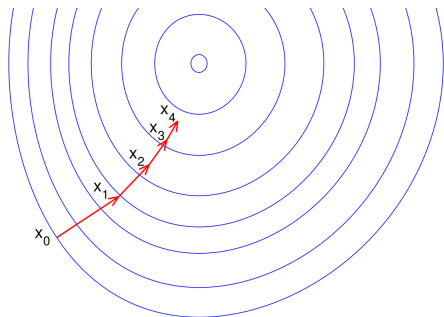


Credit: aria42.com

Illustration of iterative methods on the contour/levelset plot (i.e., the function assumes the same value on each curve)

Two questions: what direction to move, and how far to move

Iterative methods



Credit: aria42.com

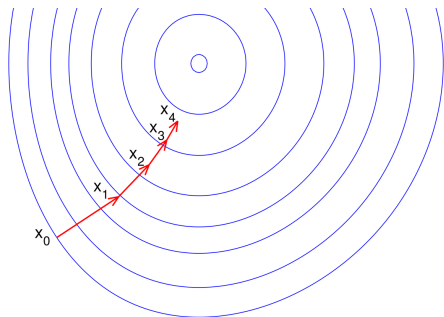
Illustration of iterative methods on the contour/levelset plot (i.e., the function assumes the same value on each curve)

Two questions: what direction to move, and how far to move

Two possibilities:

- **Line-search methods:** direction first, size second

Iterative methods



Credit: aria42.com

Illustration of iterative methods on the contour/levelset plot (i.e., the function assumes the same value on each curve)

Two questions: what direction to move, and how far to move

Two possibilities:

- **Line-search methods:** direction first, size second
- **Trust-region methods:** size first, direction second

Classic line-search methods

Advanced line-search methods

- Momentum methods

- Quasi-Newton methods

- Coordinate descent

- Conjugate gradient methods

Trust-region methods

Framework of line-search methods

A generic line search algorithm

Input: initialization \mathbf{x}_0 , stopping criterion (SC), $k = 1$

- 1: **while** SC not satisfied **do**
 - 2: choose a direction \mathbf{d}_k
 - 3: decide a step size t_k
 - 4: make a step: $\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$
 - 5: update counter: $k = k + 1$
 - 6: **end while**
-

Framework of line-search methods

A generic line search algorithm

Input: initialization \mathbf{x}_0 , stopping criterion (SC), $k = 1$

- 1: **while** SC not satisfied **do**
 - 2: choose a direction \mathbf{d}_k
 - 3: decide a step size t_k
 - 4: make a step: $\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$
 - 5: update counter: $k = k + 1$
 - 6: **end while**
-

Four questions:

- How to choose direction \mathbf{d}_k ?
- How to choose step size t_k ?
- Where to initialize?
- When to stop?

How to choose a search direction?

We want to decrease the function value toward global minimum...

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

for any fixed $t > 0$, using 1st order Taylor expansion

$$f(\mathbf{x}_k + t\mathbf{d}_{k+1}) - f(\mathbf{x}_k) \approx t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_{k+1} \rangle$$

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

for any fixed $t > 0$, using 1st order Taylor expansion

$$f(\mathbf{x}_k + t\mathbf{d}_{k+1}) - f(\mathbf{x}_k) \approx t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_{k+1} \rangle$$

$$\min_{\|\mathbf{v}\|_2=1} \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle$$

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

for any fixed $t > 0$, using 1st order Taylor expansion

$$f(\mathbf{x}_k + t\mathbf{d}_{k+1}) - f(\mathbf{x}_k) \approx t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_{k+1} \rangle$$

$$\min_{\|\mathbf{v}\|_2=1} \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle \implies \mathbf{v} = -\frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|_2}$$

How to choose a search direction?

We want to decrease the function value toward global minimum...

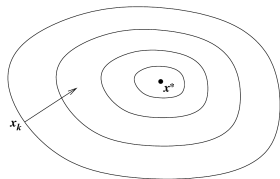
shortsighted answer: find a direction to decrease most rapidly

for any fixed $t > 0$, using 1st order Taylor expansion

$$f(\mathbf{x}_k + t\mathbf{d}_{k+1}) - f(\mathbf{x}_k) \approx t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_{k+1} \rangle$$

$$\min_{\|\mathbf{v}\|_2=1} \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle \implies \mathbf{v} = -\frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|_2}$$

Set $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$



gradient/steepest descent: $\mathbf{x}_{k+1} = \mathbf{x}_k - t\nabla f(\mathbf{x}_k)$

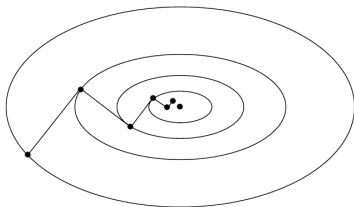
Gradient descent

$$\min_x \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{x}$$

Gradient descent

$$\min_x x^\top A x + b^\top x$$

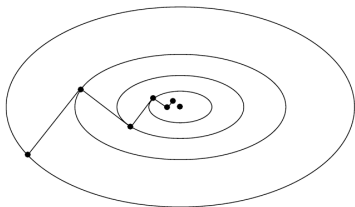
typical zig-zag path



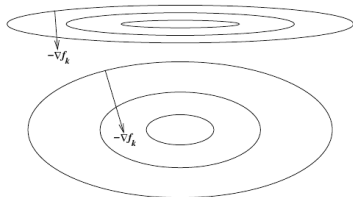
Gradient descent

$$\min_x x^T A x + b^T x$$

typical zig-zag path



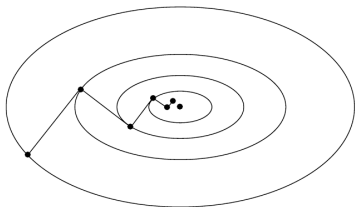
conditioning affects the path length



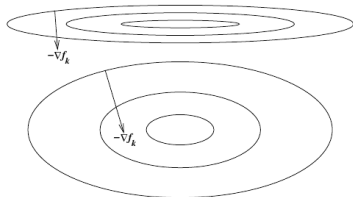
Gradient descent

$$\min_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}$$

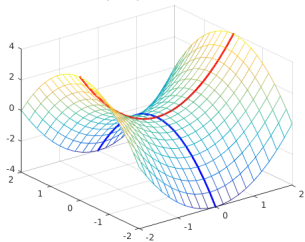
typical zig-zag path



conditioning affects the path length



$$f(x, y) = x^2 - y^2$$



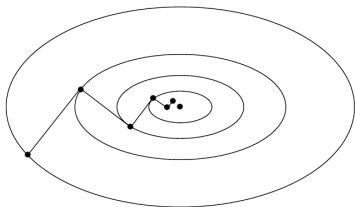
– remember direction curvature?

$$\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v} = \frac{d^2}{dt^2} f(\mathbf{x} + t\mathbf{v})$$

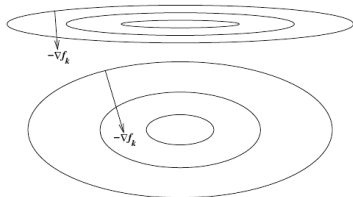
Gradient descent

$$\min_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}$$

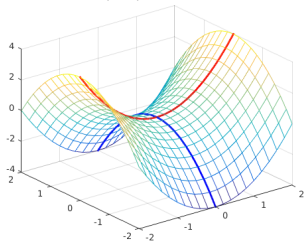
typical zig-zag path



conditioning affects the path length



$$f(x, y) = x^2 - y^2$$

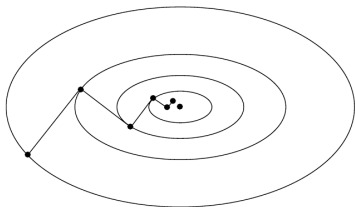


- remember direction curvature?
$$\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v} = \frac{d^2}{dt^2} f(\mathbf{x} + t\mathbf{v})$$
- large curvature \leftrightarrow narrow valley

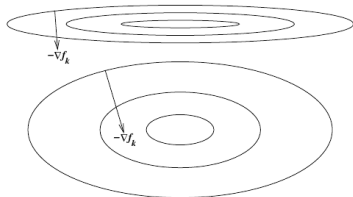
Gradient descent

$$\min_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}$$

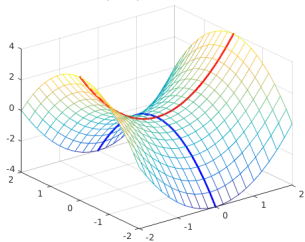
typical zig-zag path



conditioning affects the path length



$$f(x, y) = x^2 - y^2$$



- remember direction curvature?
$$\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v} = \frac{d^2}{dt^2} f(\mathbf{x} + t\mathbf{v})$$
- large curvature \leftrightarrow narrow valley
- directional curvatures encoded in the Hessian

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

farsighted answer: find a direction based on both gradient and Hessian

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

farsighted answer: find a direction based on both gradient and Hessian

for any fixed $t > 0$, using 2nd-order Taylor expansion

$$\begin{aligned} f(\mathbf{x}_k + t\mathbf{v}) - f(\mathbf{v}) &\approx t \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle \\ &\quad + \frac{1}{2} t^2 \langle \mathbf{v}, \nabla^2 f(\mathbf{x}_k) \mathbf{v} \rangle \end{aligned}$$

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

farsighted answer: find a direction based on both gradient and Hessian

for any fixed $t > 0$, using 2nd-order Taylor expansion

$$\begin{aligned} f(\mathbf{x}_k + t\mathbf{v}) - f(\mathbf{x}_k) &\approx t \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle \\ &\quad + \frac{1}{2} t^2 \langle \mathbf{v}, \nabla^2 f(\mathbf{x}_k) \mathbf{v} \rangle \end{aligned}$$

minimizing the right side

$$\implies \mathbf{v} = -t^{-1} [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

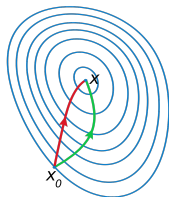
farsighted answer: find a direction based on both gradient and Hessian

for any fixed $t > 0$, using 2nd-order Taylor expansion

$$f(\mathbf{x}_k + t\mathbf{v}) - f(\mathbf{x}_k) \approx t \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle + \frac{1}{2} t^2 \langle \mathbf{v}, \nabla^2 f(\mathbf{x}_k) \mathbf{v} \rangle$$

minimizing the right side

$$\Rightarrow \mathbf{v} = -t^{-1} [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$



grad desc: green; Newton: red

How to choose a search direction?

We want to decrease the function value toward global minimum...

shortsighted answer: find a direction to decrease most rapidly

farsighted answer: find a direction based on both gradient and Hessian

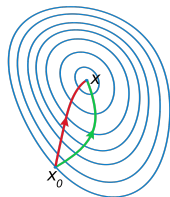
for any fixed $t > 0$, using 2nd-order Taylor expansion

$$f(\mathbf{x}_k + t\mathbf{v}) - f(\mathbf{x}_k) \approx t \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle + \frac{1}{2}t^2 \langle \mathbf{v}, \nabla^2 f(\mathbf{x}_k) \mathbf{v} \rangle$$

minimizing the right side

$$\Rightarrow \mathbf{v} = -t^{-1} [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

$$\text{Set } \mathbf{d}_k = [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$



grad desc: green; Newton: red

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

t can set to be 1.

Why called Newton's method?

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

Recall Newton's method for root-finding

$$x_{k+1} = x_k - f'(x_n) f(x_n)$$

Why called Newton's method?

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

Recall Newton's method for root-finding

$$x_{k+1} = x_k - f'(x_n) f(x_n)$$

Newton's method for solving nonlinear system $f(\mathbf{x}) = \mathbf{0}$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\mathbf{J}_f(\mathbf{x}_n)]^\dagger \mathbf{f}(\mathbf{x}_n)$$

Why called Newton's method?

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

Recall Newton's method for root-finding

$$x_{k+1} = x_k - f'(x_n) f(x_n)$$

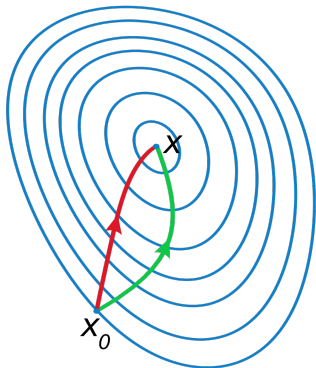
Newton's method for solving nonlinear system $f(\mathbf{x}) = \mathbf{0}$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\mathbf{J}_f(\mathbf{x}_n)]^\dagger \mathbf{f}(\mathbf{x}_n)$$

Newton's method for solving $\nabla f(\mathbf{x}) = \mathbf{0}$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_n)]^{-1} \mathbf{f}(\mathbf{x}_n)$$

How to choose a search direction?

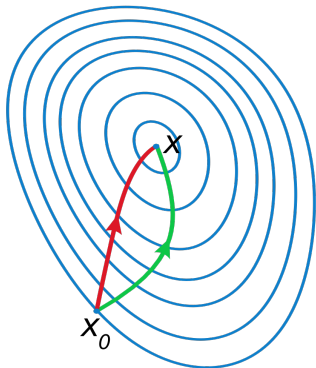


grad desc: green; Newton: red

Newton's method take fewer steps

How to choose a search direction?

nearsighted choice: cost $O(n)$ per step



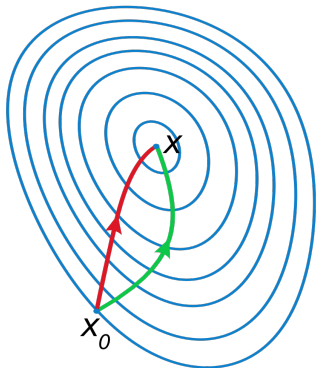
grad desc: green; Newton: red

Newton's method take fewer steps

gradient/steepest descent:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t \nabla f(\mathbf{x}_k)$$

How to choose a search direction?



grad desc: green; Newton: red

Newton's method take fewer steps

nearsighted choice: cost $O(n)$ per step

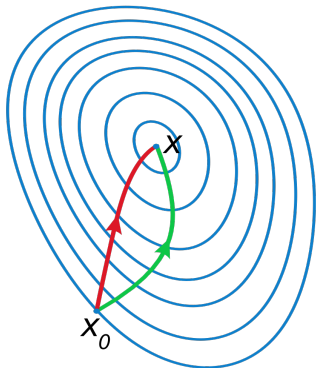
gradient/steepest descent:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t \nabla f(\mathbf{x}_k)$$

farsighted choice: cost $O(n^3)$ per step

Newton's method: $\mathbf{x}_{k+1} =$
 $\mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$

How to choose a search direction?



grad desc: green; Newton: red

Newton's method take fewer steps

nearsighted choice: cost $O(n)$ per step

gradient/steepest descent:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t \nabla f(\mathbf{x}_k)$$

farsighted choice: cost $O(n^3)$ per step

Newton's method: $\mathbf{x}_{k+1} =$
 $\mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$

Implication: The plain Newton never used for large-scale problems. More on this later ...

Problems with Newton's method

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

Problems with Newton's method

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

for any fixed $t > 0$, using 2nd-order Taylor expansion

$$\begin{aligned} f(\mathbf{x}_k + t\mathbf{v}) - f(\mathbf{v}) &\approx t \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle \\ &\quad + \frac{1}{2} t^2 \langle \mathbf{v}, \nabla^2 f(\mathbf{x}_k) \mathbf{v} \rangle \end{aligned}$$

minimizing the right side $\implies \mathbf{v} = -t^{-1} [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$

Problems with Newton's method

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

for any fixed $t > 0$, using 2nd-order Taylor expansion

$$\begin{aligned} f(\mathbf{x}_k + t\mathbf{v}) - f(\mathbf{v}) &\approx t \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle \\ &\quad + \frac{1}{2} t^2 \langle \mathbf{v}, \nabla^2 f(\mathbf{x}_k) \mathbf{v} \rangle \end{aligned}$$

minimizing the right side $\implies \mathbf{v} = -t^{-1} [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$

- $\nabla^2 f(\mathbf{x}_k)$ may be non-invertible

Problems with Newton's method

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

for any fixed $t > 0$, using 2nd-order Taylor expansion

$$\begin{aligned} f(\mathbf{x}_k + t\mathbf{v}) - f(\mathbf{v}) &\approx t \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle \\ &\quad + \frac{1}{2} t^2 \langle \mathbf{v}, \nabla^2 f(\mathbf{x}_k) \mathbf{v} \rangle \end{aligned}$$

minimizing the right side $\implies \mathbf{v} = -t^{-1} [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$

- $\nabla^2 f(\mathbf{x}_k)$ may be non-invertible
- the minimum value is $-\frac{1}{2} \langle \nabla f(\mathbf{x}_k), [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k) \rangle$. If $\nabla^2 f(\mathbf{x}_k)$ not positive definite, may be positive

Problems with Newton's method

$$\text{Newton's method: } \mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

for any fixed $t > 0$, using 2nd-order Taylor expansion

$$\begin{aligned} f(\mathbf{x}_k + t\mathbf{v}) - f(\mathbf{x}_k) &\approx t \langle \nabla f(\mathbf{x}_k), \mathbf{v} \rangle \\ &\quad + \frac{1}{2} t^2 \langle \mathbf{v}, \nabla^2 f(\mathbf{x}_k) \mathbf{v} \rangle \end{aligned}$$

minimizing the right side $\implies \mathbf{v} = -t^{-1} [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$

- $\nabla^2 f(\mathbf{x}_k)$ may be non-invertible
- the minimum value is $-\frac{1}{2} \langle \nabla f(\mathbf{x}_k), [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k) \rangle$. If $\nabla^2 f(\mathbf{x}_k)$ not positive definite, may be positive

solution: e.g., modify the Hessian $\nabla^2 f(\mathbf{x}_k) + \tau \mathbf{I}$ with τ sufficiently large

How to choose step size?

$$\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$$

How to choose step size?

$$\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$$

- Naive choice: sufficiently small constant t for all k

How to choose step size?

$$\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$$

- Naive choice: sufficiently small constant t for all k
- Robust and practical choice: back-tracking line search

How to choose step size?

$$\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$$

- Naive choice: sufficiently small constant t for all k
- Robust and practical choice: back-tracking line search

Intuition for back-tracking line search:

- By Taylor's theorem,
 $f(\mathbf{x}_k + t\mathbf{d}_k) = f(\mathbf{x}_k) + t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + o(t \|\mathbf{d}_k\|_2)$ when t sufficiently small — $t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle$ dictates the value decrease

How to choose step size?

$$\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$$

- Naive choice: sufficiently small constant t for all k
- Robust and practical choice: back-tracking line search

Intuition for back-tracking line search:

- By Taylor's theorem,
 $f(\mathbf{x}_k + t\mathbf{d}_k) = f(\mathbf{x}_k) + t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + o(t \|\mathbf{d}_k\|_2)$ when t sufficiently small — $t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle$ dictates the value decrease
- But we also want t large as possible to make rapid progress

How to choose step size?

$$\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$$

- Naive choice: sufficiently small constant t for all k
- Robust and practical choice: back-tracking line search

Intuition for back-tracking line search:

- By Taylor's theorem,
 $f(\mathbf{x}_k + t\mathbf{d}_k) = f(\mathbf{x}_k) + t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + o(t \|\mathbf{d}_k\|_2)$ when t sufficiently small — $t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle$ dictates the value decrease
- But we also want t large as possible to make rapid progress
- **idea**: find a large possible t^* to make sure
 $f(\mathbf{x}_k + t^* \mathbf{d}_k) - f(\mathbf{x}_k) \leq ct^* \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle$ (**key condition**) for a chosen parameter $c \in (0, 1)$, and no less

How to choose step size?

$$\mathbf{x}_k = \mathbf{x}_{k-1} + t_k \mathbf{d}_k$$

- Naive choice: sufficiently small constant t for all k
- Robust and practical choice: back-tracking line search

Intuition for back-tracking line search:

- By Taylor's theorem,

$f(\mathbf{x}_k + t\mathbf{d}_k) = f(\mathbf{x}_k) + t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + o(t \|\mathbf{d}_k\|_2)$ when t sufficiently small — $t \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle$ dictates the value decrease

- But we also want t large as possible to make rapid progress
- **idea:** find a large possible t^* to make sure $f(\mathbf{x}_k + t^* \mathbf{d}_k) - f(\mathbf{x}_k) \leq ct^* \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle$ (**key condition**) for a chosen parameter $c \in (0, 1)$, and no less
- **details:** start from $t = 1$. If the **key condition** not satisfied, $t = \rho t$ for a chosen parameter $\rho \in (0, 1)$.

Back-tracking line search

A widely implemented strategy in numerical optimization packages

Back-tracking line search

Input: initial $t > 0$, $\rho \in (0, 1)$, $c \in (0, 1)$

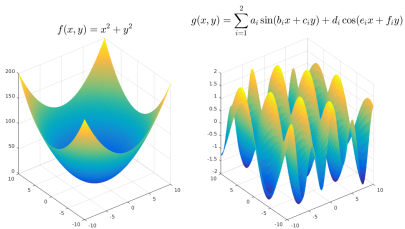
1: **while** $f(\mathbf{x}_k + t\mathbf{d}_k) - f(\mathbf{x}_k) \geq ct \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle$ **do**

2: $t = \rho t$

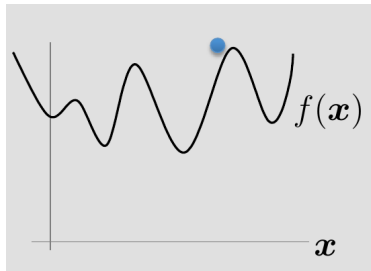
3: **end while**

Output: $t_k = t$.

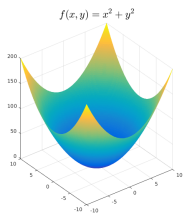
Where to initialize?



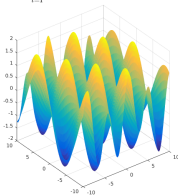
convex vs. nonconvex functions



Where to initialize?

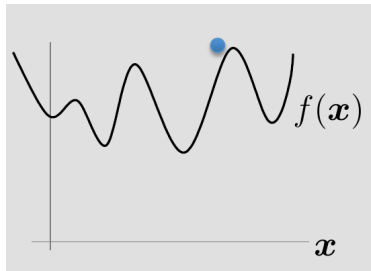


$$g(x, y) = \sum_{i=1}^2 a_i \sin(b_i x + c_i y) + d_i \cos(e_i x + f_i y)$$

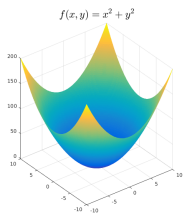


convex vs. nonconvex functions

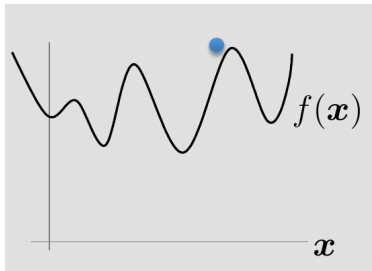
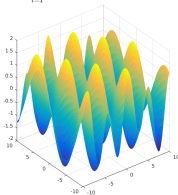
- **Convex**: most iterative methods converge to the global min no matter the initialization



Where to initialize?



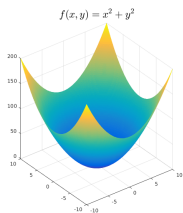
$$g(x, y) = \sum_{i=1}^2 a_i \sin(b_i x + c_i y) + d_i \cos(e_i x + f_i y)$$



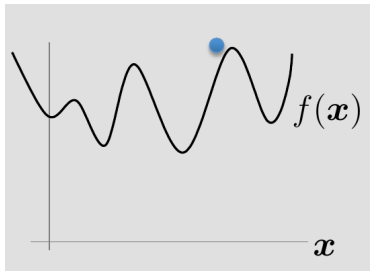
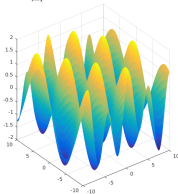
convex vs. nonconvex functions

- **Convex**: most iterative methods converge to the global min no matter the initialization
- **Nonconvex**: initialization matters a lot. Common heuristics: random initialization, multiple independent runs

Where to initialize?



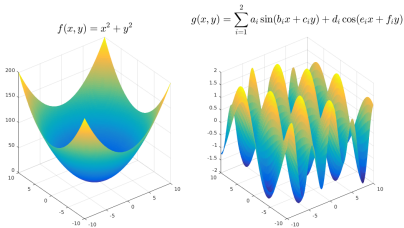
$$g(x, y) = \sum_{i=1}^2 a_i \sin(b_i x + c_i y) + d_i \cos(e_i x + f_i y)$$



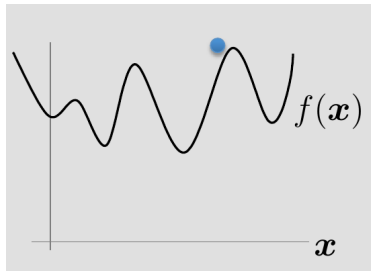
convex vs. nonconvex functions

- **Convex**: most iterative methods converge to the global min no matter the initialization
- **Nonconvex**: initialization matters a lot. Common heuristics: random initialization, multiple independent runs
- **Nonconvex**: clever initialization is possible with certain assumptions on the data:

Where to initialize?



convex vs. nonconvex functions



- **Convex**: most iterative methods converge to the global min no matter the initialization
- **Nonconvex**: initialization matters a lot. Common heuristics: random initialization, multiple independent runs
- **Nonconvex**: clever initialization is possible with certain assumptions on the data:

<https://sunju.org/research/nonconvex/>

and sometimes random initialization works!

When to stop?

1st-order necessary condition: Assume f is 1st-order differentiable at x_0 .

If x_0 is a local minimizer, then $\nabla f(x_0) = \mathbf{0}$.

2nd-order necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$.

When to stop?

1st-order necessary condition: Assume f is 1st-order differentiable at x_0 .

If x_0 is a local minimizer, then $\nabla f(x_0) = \mathbf{0}$.

2nd-order necessary condition: Assume $f(x)$ is 2-order differentiable at x_0 . If x_0 is a local min, $\nabla f(x_0) = \mathbf{0}$ and $\nabla^2 f(x_0) \succeq \mathbf{0}$.

Fix some positive tolerance values $\varepsilon_g, \varepsilon_H, \varepsilon_f, \varepsilon_v$. Possibilities:

When to stop?

1st-order necessary condition: Assume f is 1st-order differentiable at \mathbf{x}_0 .

If \mathbf{x}_0 is a local minimizer, then $\nabla f(\mathbf{x}_0) = \mathbf{0}$.

2nd-order necessary condition: Assume $f(\mathbf{x})$ is 2-order differentiable at \mathbf{x}_0 . If \mathbf{x}_0 is a local min, $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$.

Fix some positive tolerance values $\varepsilon_g, \varepsilon_H, \varepsilon_f, \varepsilon_v$. Possibilities:

- $\|\nabla f(\mathbf{x}_k)\|_2 \leq \varepsilon_g$

When to stop?

1st-order necessary condition: Assume f is 1st-order differentiable at \mathbf{x}_0 .

If \mathbf{x}_0 is a local minimizer, then $\nabla f(\mathbf{x}_0) = \mathbf{0}$.

2nd-order necessary condition: Assume $f(\mathbf{x})$ is 2-order differentiable at \mathbf{x}_0 . If \mathbf{x}_0 is a local min, $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$.

Fix some positive tolerance values ε_g , ε_H , ε_f , ε_v . Possibilities:

- $\|\nabla f(\mathbf{x}_k)\|_2 \leq \varepsilon_g$
- $\|\nabla f(\mathbf{x}_k)\|_2 \leq \varepsilon_g$ and $\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) \geq -\varepsilon_H$

When to stop?

1st-order necessary condition: Assume f is 1st-order differentiable at \mathbf{x}_0 .

If \mathbf{x}_0 is a local minimizer, then $\nabla f(\mathbf{x}_0) = \mathbf{0}$.

2nd-order necessary condition: Assume $f(\mathbf{x})$ is 2-order differentiable at \mathbf{x}_0 . If \mathbf{x}_0 is a local min, $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$.

Fix some positive tolerance values $\varepsilon_g, \varepsilon_H, \varepsilon_f, \varepsilon_v$. Possibilities:

- $\|\nabla f(\mathbf{x}_k)\|_2 \leq \varepsilon_g$
- $\|\nabla f(\mathbf{x}_k)\|_2 \leq \varepsilon_g$ and $\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) \geq -\varepsilon_H$
- $|f(\mathbf{x}_k) - f(\mathbf{x}_{k-1})| \leq \varepsilon_f$

When to stop?

1st-order necessary condition: Assume f is 1st-order differentiable at \mathbf{x}_0 .

If \mathbf{x}_0 is a local minimizer, then $\nabla f(\mathbf{x}_0) = \mathbf{0}$.

2nd-order necessary condition: Assume $f(\mathbf{x})$ is 2-order differentiable at \mathbf{x}_0 . If \mathbf{x}_0 is a local min, $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$.

Fix some positive tolerance values $\varepsilon_g, \varepsilon_H, \varepsilon_f, \varepsilon_v$. Possibilities:

- $\|\nabla f(\mathbf{x}_k)\|_2 \leq \varepsilon_g$
- $\|\nabla f(\mathbf{x}_k)\|_2 \leq \varepsilon_g$ and $\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) \geq -\varepsilon_H$
- $|f(\mathbf{x}_k) - f(\mathbf{x}_{k-1})| \leq \varepsilon_f$
- $\|\mathbf{x}_k - \mathbf{x}_{k-1}\|_2 \leq \varepsilon_v$

Nonconvex optimization is hard

Nonconvex: Even computing (verifying!) a local minimizer is NP-hard!
(see, e.g., [[Murty and Kabadi, 1987](#)])

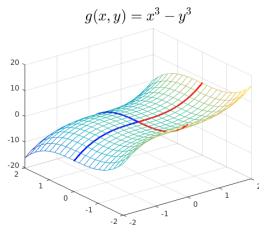
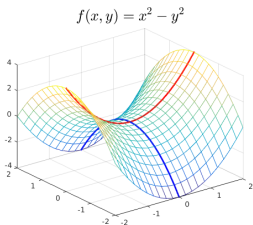
Nonconvex optimization is hard

Nonconvex: Even computing (verifying!) a local minimizer is NP-hard!

(see, e.g., [Murty and Kabadi, 1987])

2nd order sufficient: $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succ \mathbf{0}$

2nd order necessary: $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$



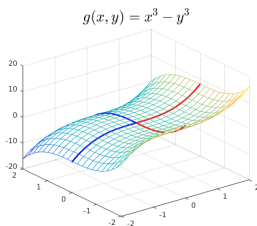
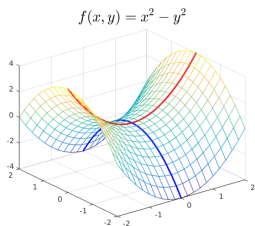
Nonconvex optimization is hard

Nonconvex: Even computing (verifying!) a local minimizer is NP-hard!

(see, e.g., [Murty and Kabadi, 1987])

2nd order sufficient: $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succ \mathbf{0}$

2nd order necessary: $\nabla f(\mathbf{x}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$



Cases in between: local shapes around SOSP determined by **spectral properties of higher-order derivative tensors**, calculating which is hard [Hillar and Lim, 2013]!

Classic line-search methods

Advanced line-search methods

- Momentum methods

- Quasi-Newton methods

- Coordinate descent

- Conjugate gradient methods

Trust-region methods

Classic line-search methods

Advanced line-search methods

- Momentum methods

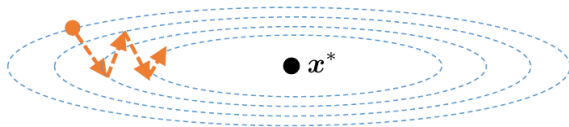
- Quasi-Newton methods

- Coordinate descent

- Conjugate gradient methods

Trust-region methods

Why momentum?

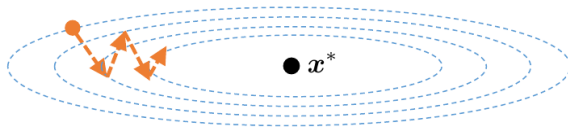


gradient descent

Credit: Princeton ELE522

- GD is cheap ($O(n)$ per step) but overall convergence sensitive to conditioning
- Newton's convergence is not sensitive to conditioning but expensive ($O(n^3)$ per step)

Why momentum?



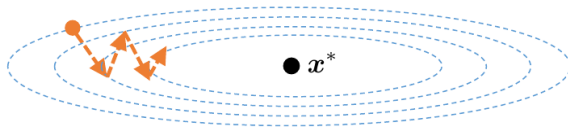
gradient descent

Credit: Princeton ELE522

- GD is cheap ($O(n)$ per step) but overall convergence sensitive to conditioning
- Newton's convergence is not sensitive to conditioning but expensive ($O(n^3)$ per step)

A cheap way to achieve faster convergence?

Why momentum?



gradient descent

Credit: Princeton ELE522

- GD is cheap ($O(n)$ per step) but overall convergence sensitive to conditioning
- Newton's convergence is not sensitive to conditioning but expensive ($O(n^3)$ per step)

A cheap way to achieve faster convergence? **Answer: using historic information**

Heavy ball method

In physics, a heavy object has a large inertia/momentum — resistance to change velocity.

Heavy ball method

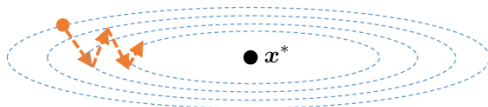
In physics, a heavy object has a large inertia/momentum — resistance to change velocity.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k \underbrace{(\mathbf{x}_k - \mathbf{x}_{k-1})}_{\text{momentum}} \quad \text{due to Polyak}$$

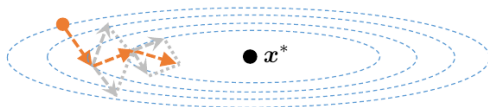
Heavy ball method

In physics, a heavy object has a large inertia/momentum — resistance to change velocity.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \underbrace{\beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})}_{\text{momentum}} \quad \text{due to Polyak}$$



gradient descent



heavy-ball method

Credit: Princeton ELE522

History helps to smooth out the zig-zag path!

Nesterov's accelerated gradient methods

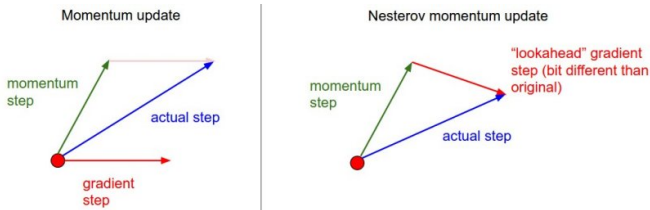
Another version, due to Y. Nesterov

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}) - \alpha_k \nabla f (\mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}))$$

Nesterov's accelerated gradient methods

Another version, due to Y. Nesterov

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}) - \alpha_k \nabla f (\mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}))$$



Credit: Stanford CS231N

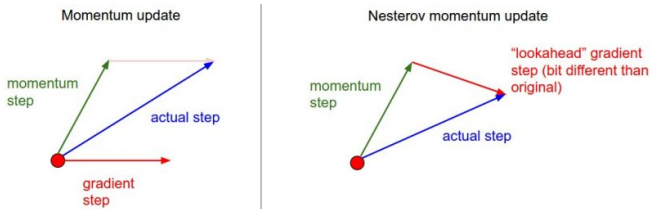
$$\text{HB} \begin{cases} x_{\text{ahead}} = x + \beta(x - x_{\text{old}}), \\ x_{\text{new}} = x_{\text{ahead}} - \alpha \nabla f(x). \end{cases}$$

$$\text{Nesterov} \begin{cases} x_{\text{ahead}} = x + \beta(x - x_{\text{old}}), \\ x_{\text{new}} = x_{\text{ahead}} - \alpha \nabla f(x_{\text{ahead}}). \end{cases}$$

Nesterov's accelerated gradient methods

Another version, due to Y. Nesterov

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}) - \alpha_k \nabla f (\mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}))$$



Credit: Stanford CS231N

$$\text{HB} \begin{cases} x_{\text{ahead}} = x + \beta(x - x_{\text{old}}), \\ x_{\text{new}} = x_{\text{ahead}} - \alpha \nabla f(x). \end{cases} \quad \text{Nesterov} \begin{cases} x_{\text{ahead}} = x + \beta(x - x_{\text{old}}), \\ x_{\text{new}} = x_{\text{ahead}} - \alpha \nabla f(x_{\text{ahead}}). \end{cases}$$

For more info, see Chap 10 of [Beck, 2017] and Chap 2 of [Nesterov, 2018].

Classic line-search methods

Advanced line-search methods

Momentum methods

Quasi-Newton methods

Coordinate descent

Conjugate gradient methods

Trust-region methods

Quasi-Newton methods

quasi-: seemingly; apparently but not really.

Newton's method: cost $O(n^2)$ storage and $O(n^3)$ computation per step

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

Quasi-Newton methods

quasi-: seemingly; apparently but not really.

Newton's method: cost $O(n^2)$ storage and $O(n^3)$ computation per step

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

Idea: approximate $\nabla^2 f(\mathbf{x}_k)$ or $[\nabla^2 f(\mathbf{x}_k)]^{-1}$ to allow efficient storage and computation — **Quasi-Newton Methods**

Quasi-Newton methods

quasi-: seemingly; apparently but not really.

Newton's method: cost $O(n^2)$ storage and $O(n^3)$ computation per step

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

Idea: approximate $\nabla^2 f(\mathbf{x}_k)$ or $[\nabla^2 f(\mathbf{x}_k)]^{-1}$ to allow efficient storage and computation — **Quasi-Newton Methods**

Choose \mathbf{H}_k to approximate $\nabla^2 f(\mathbf{x}_k)$ so that

- avoid calculation of second derivatives
- simplify matrix inversion, i.e., computing the search direction

Quasi-Newton methods

given: starting point $x_0 \in \text{dom } f$, $H_0 > 0$

for $k = 0, 1, \dots$

1. compute quasi-Newton direction $\Delta x_k = -H_k^{-1} \nabla f(x_k)$
 2. determine step size t_k (e.g., by backtracking line search)
 3. compute $x_{k+1} = x_k + t_k \Delta x_k$
 4. compute H_{k+1}
- Different variants differ on how to compute H_{k+1}
 - Normally H_k^{-1} or its factorized version stored to simplify calculation of Δx_k

Broyden–Fletcher–Goldfarb–Shanno (BFGS) method

BFGS update

$$H_{k+1} = H_k + \frac{yy^T}{y^T s} - \frac{H_k s s^T H_k}{s^T H_k s}$$

where

$$s = x_{k+1} - x_k, \quad y = \nabla f(x_{k+1}) - \nabla f(x_k)$$

Inverse update

$$H_{k+1}^{-1} = \left(I - \frac{sy^T}{y^T s} \right) H_k^{-1} \left(I - \frac{ys^T}{y^T s} \right) + \frac{ss^T}{y^T s}$$

Broyden–Fletcher–Goldfarb–Shanno (BFGS) method

BFGS update

$$H_{k+1} = H_k + \frac{yy^T}{y^T s} - \frac{H_k s s^T H_k}{s^T H_k s}$$

where

$$s = x_{k+1} - x_k, \quad y = \nabla f(x_{k+1}) - \nabla f(x_k)$$

Inverse update

$$H_{k+1}^{-1} = \left(I - \frac{sy^T}{y^T s} \right) H_k^{-1} \left(I - \frac{ys^T}{y^T s} \right) + \frac{ss^T}{y^T s}$$

Cost of update: $O(n^2)$ (vs. $O(n^3)$ in Newton's method), storage: $O(n^2)$

Broyden–Fletcher–Goldfarb–Shanno (BFGS) method

BFGS update

$$H_{k+1} = H_k + \frac{yy^T}{y^T s} - \frac{H_k s s^T H_k}{s^T H_k s}$$

where

$$s = x_{k+1} - x_k, \quad y = \nabla f(x_{k+1}) - \nabla f(x_k)$$

Inverse update

$$H_{k+1}^{-1} = \left(I - \frac{sy^T}{y^T s} \right) H_k^{-1} \left(I - \frac{ys^T}{y^T s} \right) + \frac{ss^T}{y^T s}$$

Cost of update: $O(n^2)$ (vs. $O(n^3)$ in Newton's method), storage: $O(n^2)$ To derive the update equations, three conditions are imposed:

- secant condition: $H_{k+1} s = y$ (think of 1st Taylor expansion to ∇f)
- Curvature condition: $s_k^T y_k > 0$ to ensure that $H_{k+1} \succ 0$ if $H_k \succ 0$
- H_{k+1} and H_k are close in an appropriate sense

See Chap 6 of [[Nocedal and Wright, 2006](#)] Credit: UCLA ECE236C

Limited-memory BFGS (L-BFGS)

Limited-memory BFGS (L-BFGS): do not store H_k^{-1} explicitly

- instead we store up to m (e.g., $m = 30$) values of

$$s_j = x_{j+1} - x_j, \quad y_j = \nabla f(x_{j+1}) - \nabla f(x_j)$$

- we evaluate $\Delta x_k = H_k^{-1} \nabla f(x_k)$ recursively, using

$$H_{j+1}^{-1} = \left(I - \frac{s_j y_j^T}{y_j^T s_j} \right) H_j^{-1} \left(I - \frac{y_j s_j^T}{y_j^T s_j} \right) + \frac{s_j s_j^T}{y_j^T s_j}$$

for $j = k - 1, \dots, k - m$, assuming, for example, $H_{k-m} = I$

- an alternative is to restart after m iterations

Limited-memory BFGS (L-BFGS)

Limited-memory BFGS (L-BFGS): do not store H_k^{-1} explicitly

- instead we store up to m (e.g., $m = 30$) values of

$$s_j = x_{j+1} - x_j, \quad y_j = \nabla f(x_{j+1}) - \nabla f(x_j)$$

- we evaluate $\Delta x_k = H_k^{-1} \nabla f(x_k)$ recursively, using

$$H_{j+1}^{-1} = \left(I - \frac{s_j y_j^T}{y_j^T s_j} \right) H_j^{-1} \left(I - \frac{y_j s_j^T}{y_j^T s_j} \right) + \frac{s_j s_j^T}{y_j^T s_j}$$

for $j = k - 1, \dots, k - m$, assuming, for example, $H_{k-m} = I$

- an alternative is to restart after m iterations

Cost of update: $O(mn)$ (vs. $O(n^2)$ in BFGS), storage: $O(mn)$ (vs. $O(n^2)$ in BFGS) — linear in dimension n ! recall the cost of GD?

See Chap 7 of [Nocedal and Wright, 2006] Credit: UCLA ECE236C

Classic line-search methods

Advanced line-search methods

Momentum methods

Quasi-Newton methods

Coordinate descent

Conjugate gradient methods

Trust-region methods

Block coordinate descent

Consider a function $f(\mathbf{x}_1, \dots, \mathbf{x}_p)$ with $\mathbf{x}_1 \in \mathbb{R}^{n_1}, \dots, \mathbf{x}_p \in \mathbb{R}^{n_p}$

Block coordinate descent

Consider a function $f(\mathbf{x}_1, \dots, \mathbf{x}_p)$ with $\mathbf{x}_1 \in \mathbb{R}^{n_1}, \dots, \mathbf{x}_p \in \mathbb{R}^{n_p}$

A generic block coordinate descent algorithm

Input: initialization $(\mathbf{x}_{1,0}, \dots, \mathbf{x}_{p,0})$ (the 2nd subscript indexes iteration number)

1: **for** $k = 1, 2, \dots$ **do**

2: Pick a block index $i \in \{1, \dots, p\}$

3: Minimize wrt the chosen block:

$$\mathbf{x}_{i,k} = \arg \min_{\boldsymbol{\xi} \in \mathbb{R}^{n_i}} f(\mathbf{x}_{1,k-1}, \dots, \mathbf{x}_{i-1,k-1}, \boldsymbol{\xi}, \mathbf{x}_{i+1,k-1}, \dots, \mathbf{x}_{p,k-1})$$

4: Leave other blocks unchanged: $\mathbf{x}_{j,k} = \mathbf{x}_{j,k-1} \forall j \neq i$

5: **end for**

Block coordinate descent

Consider a function $f(\mathbf{x}_1, \dots, \mathbf{x}_p)$ with $\mathbf{x}_1 \in \mathbb{R}^{n_1}, \dots, \mathbf{x}_p \in \mathbb{R}^{n_p}$

A generic block coordinate descent algorithm

Input: initialization $(\mathbf{x}_{1,0}, \dots, \mathbf{x}_{p,0})$ (the 2nd subscript indexes iteration number)

1: **for** $k = 1, 2, \dots$ **do**

2: Pick a block index $i \in \{1, \dots, p\}$

3: Minimize wrt the chosen block:

$$\mathbf{x}_{i,k} = \arg \min_{\boldsymbol{\xi} \in \mathbb{R}^{n_i}} f(\mathbf{x}_{1,k-1}, \dots, \mathbf{x}_{i-1,k-1}, \boldsymbol{\xi}, \mathbf{x}_{i+1,k-1}, \dots, \mathbf{x}_{p,k-1})$$

4: Leave other blocks unchanged: $\mathbf{x}_{j,k} = \mathbf{x}_{j,k-1} \forall j \neq i$

5: **end for**

– Also called **alternating direction/minimization methods**

Block coordinate descent

Consider a function $f(\mathbf{x}_1, \dots, \mathbf{x}_p)$ with $\mathbf{x}_1 \in \mathbb{R}^{n_1}, \dots, \mathbf{x}_p \in \mathbb{R}^{n_p}$

A generic block coordinate descent algorithm

Input: initialization $(\mathbf{x}_{1,0}, \dots, \mathbf{x}_{p,0})$ (the 2nd subscript indexes iteration number)

1: **for** $k = 1, 2, \dots$ **do**

2: Pick a block index $i \in \{1, \dots, p\}$

3: Minimize wrt the chosen block:

$$\mathbf{x}_{i,k} = \arg \min_{\boldsymbol{\xi} \in \mathbb{R}^{n_i}} f(\mathbf{x}_{1,k-1}, \dots, \mathbf{x}_{i-1,k-1}, \boldsymbol{\xi}, \mathbf{x}_{i+1,k-1}, \dots, \mathbf{x}_{p,k-1})$$

4: Leave other blocks unchanged: $\mathbf{x}_{j,k} = \mathbf{x}_{j,k-1} \forall j \neq i$

5: **end for**

- Also called **alternating direction/minimization methods**
- When $n_1 = n_2 = \dots = n_p = 1$, called **coordinate descent**

Block coordinate descent

Consider a function $f(\mathbf{x}_1, \dots, \mathbf{x}_p)$ with $\mathbf{x}_1 \in \mathbb{R}^{n_1}, \dots, \mathbf{x}_p \in \mathbb{R}^{n_p}$

A generic block coordinate descent algorithm

Input: initialization $(\mathbf{x}_{1,0}, \dots, \mathbf{x}_{p,0})$ (the 2nd subscript indexes iteration number)

1: **for** $k = 1, 2, \dots$ **do**

2: Pick a block index $i \in \{1, \dots, p\}$

3: Minimize wrt the chosen block:

$$\mathbf{x}_{i,k} = \arg \min_{\boldsymbol{\xi} \in \mathbb{R}^{n_i}} f(\mathbf{x}_{1,k-1}, \dots, \mathbf{x}_{i-1,k-1}, \boldsymbol{\xi}, \mathbf{x}_{i+1,k-1}, \dots, \mathbf{x}_{p,k-1})$$

4: Leave other blocks unchanged: $\mathbf{x}_{j,k} = \mathbf{x}_{j,k-1} \forall j \neq i$

5: **end for**

- Also called **alternating direction/minimization methods**
- When $n_1 = n_2 = \dots = n_p = 1$, called **coordinate descent**
- Minimization in Line 3 can be **inexact**: e.g.,

$$\mathbf{x}_{i,k} = \mathbf{x}_{i,k-1} - t_k \frac{\partial f}{\partial \boldsymbol{\xi}}(\mathbf{x}_{1,k-1}, \dots, \mathbf{x}_{i-1,k-1}, \mathbf{x}_{i,k-1}, \mathbf{x}_{i+1,k-1}, \dots, \mathbf{x}_{p,k-1})$$

Block coordinate descent

Consider a function $f(\mathbf{x}_1, \dots, \mathbf{x}_p)$ with $\mathbf{x}_1 \in \mathbb{R}^{n_1}, \dots, \mathbf{x}_p \in \mathbb{R}^{n_p}$

A generic block coordinate descent algorithm

Input: initialization $(\mathbf{x}_{1,0}, \dots, \mathbf{x}_{p,0})$ (the 2nd subscript indexes iteration number)

1: **for** $k = 1, 2, \dots$ **do**

2: Pick a block index $i \in \{1, \dots, p\}$

3: Minimize wrt the chosen block:

$$\mathbf{x}_{i,k} = \arg \min_{\boldsymbol{\xi} \in \mathbb{R}^{n_i}} f(\mathbf{x}_{1,k-1}, \dots, \mathbf{x}_{i-1,k-1}, \boldsymbol{\xi}, \mathbf{x}_{i+1,k-1}, \dots, \mathbf{x}_{p,k-1})$$

4: Leave other blocks unchanged: $\mathbf{x}_{j,k} = \mathbf{x}_{j,k-1} \forall j \neq i$

5: **end for**

– Also called **alternating direction/minimization methods**

– When $n_1 = n_2 = \dots = n_p = 1$, called **coordinate descent**

– Minimization in Line 3 can be **inexact**: e.g.,

$$\mathbf{x}_{i,k} = \mathbf{x}_{i,k-1} - t_k \frac{\partial f}{\partial \boldsymbol{\xi}}(\mathbf{x}_{1,k-1}, \dots, \mathbf{x}_{i-1,k-1}, \mathbf{x}_{i,k-1}, \mathbf{x}_{i+1,k-1}, \dots, \mathbf{x}_{p,k-1})$$

– In Line 2, many different ways of picking an index, e.g., cyclic, randomized, weighted sampling, etc

Block coordinate descent: examples

Least-squares $\min_{\mathbf{x}} f(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$

- $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 = \|\mathbf{y} - \mathbf{A}_{-i}\mathbf{x}_{-i} - \mathbf{a}_i x_i\|_2^2$

- coordinate descent: $\min_{\xi \in \mathbb{R}} \|\mathbf{y} - \mathbf{A}_{-i}\mathbf{x}_{-i} - \mathbf{a}_i \xi\|_2^2$

$$\implies x_{i,+} = \frac{\langle \mathbf{y} - \mathbf{A}_{-i}\mathbf{x}_{-i}, \mathbf{a}_i \rangle}{\|\mathbf{a}_i\|_2^2}$$

(\mathbf{A}_{-i} is \mathbf{A} with the i -th column removed; \mathbf{x}_{-i} is \mathbf{x} with the i -th coordinate removed)

Block coordinate descent: examples

Least-squares $\min_{\mathbf{x}} f(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$

- $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 = \|\mathbf{y} - \mathbf{A}_{-i}\mathbf{x}_{-i} - \mathbf{a}_i x_i\|_2^2$
- coordinate descent: $\min_{\xi \in \mathbb{R}} \|\mathbf{y} - \mathbf{A}_{-i}\mathbf{x}_{-i} - \mathbf{a}_i \xi\|_2^2$
 $\implies x_{i,+} = \frac{\langle \mathbf{y} - \mathbf{A}_{-i}\mathbf{x}_{-i}, \mathbf{a}_i \rangle}{\|\mathbf{a}_i\|_2^2}$

(\mathbf{A}_{-i} is \mathbf{A} with the i -th column removed; \mathbf{x}_{-i} is \mathbf{x} with the i -th coordinate removed)

Matrix factorization $\min_{\mathbf{A}, \mathbf{B}} \|\mathbf{Y} - \mathbf{A}\mathbf{B}\|_F^2$

- Two groups of variables, consider block coordinate descent
- Updates:

$$\begin{aligned} \mathbf{A}_+ &= \mathbf{Y}\mathbf{B}^\dagger, \\ \mathbf{B}_+ &= \mathbf{A}^\dagger\mathbf{Y}. \end{aligned}$$

(\cdot)[†] denotes the matrix pseudoinverse.)

Why block coordinate descent?

- may work with constrained problems and non-differentiable problems (e.g., $\min_{\mathbf{A}, \mathbf{B}} \|\mathbf{Y} - \mathbf{A}\mathbf{B}\|_F^2$, s. t. \mathbf{A} orthogonal, Lasso: $\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$)

Why block coordinate descent?

- may work with constrained problems and non-differentiable problems (e.g., $\min_{\mathbf{A}, \mathbf{B}} \|\mathbf{Y} - \mathbf{A}\mathbf{B}\|_F^2$, s. t. \mathbf{A} orthogonal, Lasso: $\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$)
- may be faster than gradient descent or Newton (next)

Why block coordinate descent?

- may work with constrained problems and non-differentiable problems (e.g., $\min_{\mathbf{A}, \mathbf{B}} \|\mathbf{Y} - \mathbf{A}\mathbf{B}\|_F^2$, s. t. \mathbf{A} orthogonal, Lasso: $\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$)
- may be faster than gradient descent or Newton (next)
- may be simple and cheap!

Why block coordinate descent?

- may work with constrained problems and non-differentiable problems (e.g., $\min_{\mathbf{A}, \mathbf{B}} \|\mathbf{Y} - \mathbf{A}\mathbf{B}\|_F^2$, s. t. \mathbf{A} orthogonal, Lasso: $\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$)
- may be faster than gradient descent or Newton (next)
- may be simple and cheap!

Some references:

- [\[Wright, 2015\]](#)
- [Lecture notes by Prof. Ruoyu Sun](#)

Classic line-search methods

Advanced line-search methods

Momentum methods

Quasi-Newton methods

Coordinate descent

Conjugate gradient methods

Trust-region methods

Conjugate direction methods

Solve linear equation $\mathbf{y} = \mathbf{A}\mathbf{x} \iff \min_{\mathbf{x}} \frac{1}{2}\mathbf{x}^T\mathbf{A}\mathbf{x} - \mathbf{b}^T\mathbf{x}$ with $\mathbf{A} \succ \mathbf{0}$

Conjugate direction methods

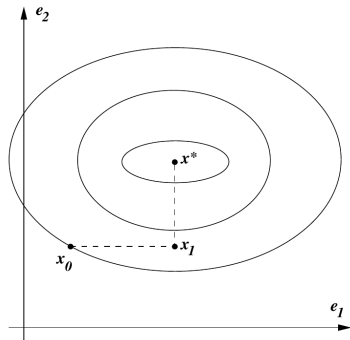
Solve linear equation $\mathbf{y} = \mathbf{A}\mathbf{x} \iff \min_{\mathbf{x}} \frac{1}{2}\mathbf{x}^T\mathbf{A}\mathbf{x} - \mathbf{b}^T\mathbf{x}$ with $\mathbf{A} \succ \mathbf{0}$

apply coordinate descent...

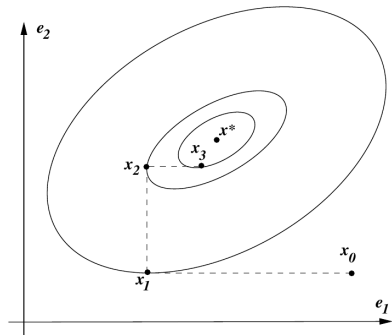
Conjugate direction methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2}x^T Ax - b^T x$ with $A \succ 0$

apply coordinate descent...



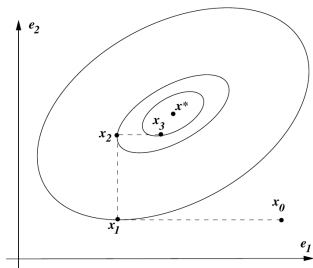
diagonal A : solve the problem in n steps



non-diagonal A : does not solve the problem in n steps

Conjugate direction methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2} x^T Ax - b^T x$ with $A \succ 0$

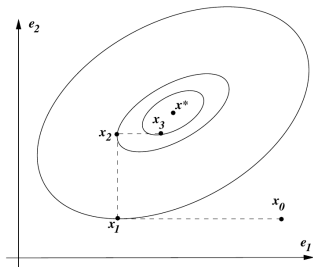


non-diagonal A : does not solve
the problem in n steps

Conjugate direction methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2}x^T Ax - b^T x$ with $A \succ 0$

Idea: define n “conjugate directions”
 $\{p_1, \dots, p_n\}$ so that $p_i^T A p_j = 0$ for all
 $i \neq j$ —conjugate as generalization of orthogonal



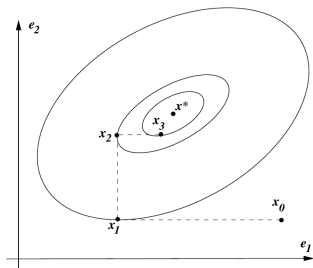
non-diagonal A : does not solve
the problem in n steps

Conjugate direction methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2}x^T Ax - b^T x$ with $A \succ 0$

Idea: define n “conjugate directions”
 $\{p_1, \dots, p_n\}$ so that $p_i^T A p_j = 0$ for all
 $i \neq j$ —conjugate as generalization of orthogonal

– Write $P = [p_1, \dots, p_n]$. Can verify that
 $P^T A P$ is diagonal and positive



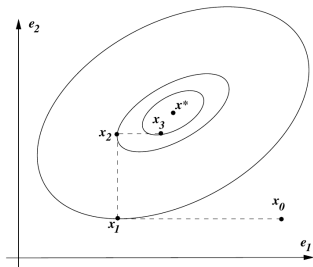
non-diagonal A : does not solve
the problem in n steps

Conjugate direction methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2}x^T Ax - b^T x$ with $A \succ 0$

Idea: define n “conjugate directions” $\{p_1, \dots, p_n\}$ so that $p_i^T A p_j = 0$ for all $i \neq j$ —conjugate as generalization of orthogonal

- Write $P = [p_1, \dots, p_n]$. Can verify that $P^T A P$ is diagonal and positive
- Write $x = P s$. Then $\frac{1}{2}x^T A x - b^T x = \frac{1}{2}s^T (P^T A P) s - (P^T b)^T s$ — quadratic with diagonal $P^T A P$



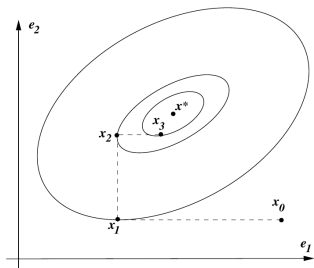
non-diagonal A : does not solve the problem in n steps

Conjugate direction methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2}x^T Ax - b^T x$ with $A \succ 0$

Idea: define n “conjugate directions” $\{p_1, \dots, p_n\}$ so that $p_i^T A p_j = 0$ for all $i \neq j$ —conjugate as generalization of orthogonal

- Write $P = [p_1, \dots, p_n]$. Can verify that $P^T A P$ is diagonal and positive
- Write $x = P s$. Then $\frac{1}{2}x^T A x - b^T x = \frac{1}{2}s^T (P^T A P) s - (P^T b)^T s$ — quadratic with diagonal $P^T A P$
- Perform updates in the s space, but write the equivalent form in x space



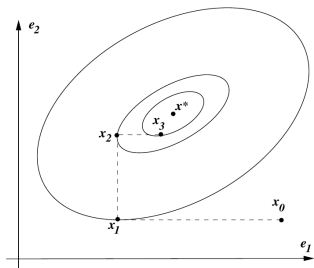
non-diagonal A : does not solve the problem in n steps

Conjugate direction methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2}x^T Ax - b^T x$ with $A \succ 0$

Idea: define n “conjugate directions” $\{p_1, \dots, p_n\}$ so that $p_i^T A p_j = 0$ for all $i \neq j$ —conjugate as generalization of orthogonal

- Write $P = [p_1, \dots, p_n]$. Can verify that $P^T A P$ is diagonal and positive
- Write $x = P s$. Then $\frac{1}{2}x^T A x - b^T x = \frac{1}{2}s^T (P^T A P) s - (P^T b)^T s$ — quadratic with diagonal $P^T A P$
- Perform updates in the s space, but write the equivalent form in x space
- The i -th coordinate direction in the s space is p_i in the x space



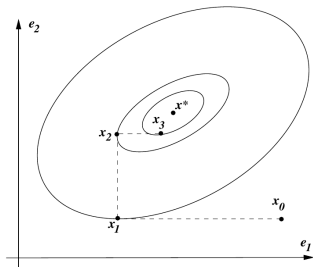
non-diagonal A : does not solve the problem in n steps

Conjugate direction methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2}x^T Ax - b^T x$ with $A \succ 0$

Idea: define n “conjugate directions” $\{p_1, \dots, p_n\}$ so that $p_i^T A p_j = 0$ for all $i \neq j$ —conjugate as generalization of orthogonal

- Write $P = [p_1, \dots, p_n]$. Can verify that $P^T A P$ is diagonal and positive
- Write $x = P s$. Then $\frac{1}{2}x^T A x - b^T x = \frac{1}{2}s^T (P^T A P) s - (P^T b)^T s$ — quadratic with diagonal $P^T A P$
- Perform updates in the s space, but write the equivalent form in x space
- The i -th coordinate direction in the s space is p_i in the x space



non-diagonal A : does not solve the problem in n steps

In short, change of variable trick!

Conjugate gradient methods

Solve linear equation $\mathbf{y} = \mathbf{A}\mathbf{x} \iff \min_{\mathbf{x}} \frac{1}{2}\mathbf{x}^\top \mathbf{A}\mathbf{x} - \mathbf{b}^\top \mathbf{x}$ with $\mathbf{A} \succ \mathbf{0}$

Idea: define n “conjugate directions” $\{\mathbf{p}_1, \dots, \mathbf{p}_n\}$ so that $\mathbf{p}_i^\top \mathbf{A}\mathbf{p}_j = 0$ for all $i \neq j$ —conjugate as generalization of orthogonal

Conjugate gradient methods

Solve linear equation $\mathbf{y} = \mathbf{A}\mathbf{x} \iff \min_{\mathbf{x}} \frac{1}{2}\mathbf{x}^\top \mathbf{A}\mathbf{x} - \mathbf{b}^\top \mathbf{x}$ with $\mathbf{A} \succ \mathbf{0}$

Idea: define n “conjugate directions” $\{\mathbf{p}_1, \dots, \mathbf{p}_n\}$ so that $\mathbf{p}_i^\top \mathbf{A}\mathbf{p}_j = 0$ for all $i \neq j$ —conjugate as generalization of orthogonal

Generally, many choices for $\{\mathbf{p}_1, \dots, \mathbf{p}_n\}$.

Conjugate gradient methods: choice based on ideas from steepest descent

Conjugate gradient methods

Solve linear equation $y = Ax \iff \min_x \frac{1}{2}x^T Ax - b^T x$ with $A \succ 0$

Idea: define n “conjugate directions” $\{p_1, \dots, p_n\}$ so that $p_i^T A p_j = 0$ for all $i \neq j$ —conjugate as generalization of orthogonal

Generally, many choices for $\{p_1, \dots, p_n\}$.

Conjugate gradient methods: choice based on ideas from steepest descent

Algorithm 5.2 (CG).

Given x_0 ;

Set $r_0 \leftarrow Ax_0 - b$, $p_0 \leftarrow -r_0$, $k \leftarrow 0$;

while $r_k \neq 0$

$$\alpha_k \leftarrow \frac{r_k^T r_k}{p_k^T A p_k}; \quad (5.24a)$$

$$x_{k+1} \leftarrow x_k + \alpha_k p_k; \quad (5.24b)$$

$$r_{k+1} \leftarrow r_k + \alpha_k A p_k; \quad (5.24c)$$

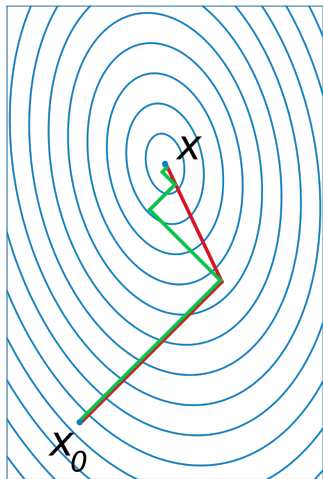
$$\beta_{k+1} \leftarrow \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}; \quad (5.24d)$$

$$p_{k+1} \leftarrow -r_{k+1} + \beta_{k+1} p_k; \quad (5.24e)$$

$$k \leftarrow k + 1; \quad (5.24f)$$

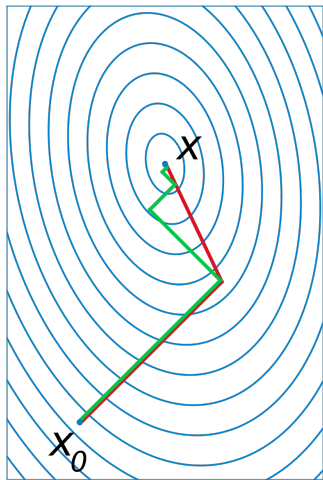
end (while)

Conjugate gradient methods



CG vs. GD (Green: GD,
Red: CG)

Conjugate gradient methods



CG vs. GD (Green: GD,
Red: CG)

- Can be extended to general non-quadratic functions
- Often used to solve subproblems of other iterative methods, e.g., truncated Newton method, the trust-region subproblem (later)

See Chap 5

of [\[Nocedal and Wright, 2006\]](#)

Classic line-search methods

Advanced line-search methods

- Momentum methods

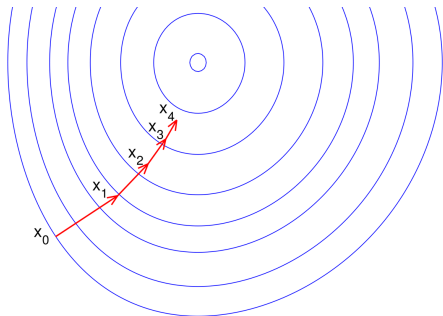
- Quasi-Newton methods

- Coordinate descent

- Conjugate gradient methods

Trust-region methods

Iterative methods



Credit: aria42.com

Illustration of iterative methods on the contour/levelset plot (i.e., the function assumes the same value on each curve)

Two questions: what direction to move, and how far to move

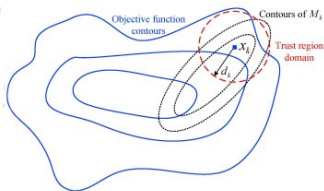
Two possibilities:

- **Line-search methods:** direction first, size second
- **Trust-region methods (TRM):** size first, direction second

Ideas behind TRM

Recall Taylor expansion $f(\mathbf{x} + \mathbf{d}) \approx f(\mathbf{x}) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle$

Start with \mathbf{x}_0 . Repeat the following:



Credit: [Arezki et al., 2018]

Ideas behind TRM

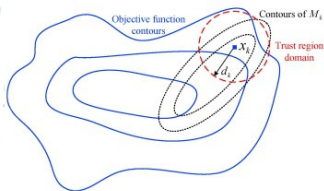
Recall Taylor expansion $f(\mathbf{x} + \mathbf{d}) \approx f(\mathbf{x}) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle$

Start with \mathbf{x}_0 . Repeat the following:

- At \mathbf{x}_k , approximate f by the quadratic function (called **model function** dotted black)

$$m_k(\mathbf{d}) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$$

i.e., $m_k(\mathbf{d}) \approx f(\mathbf{x}_k + \mathbf{d})$, and \mathbf{B}_k to approximate $\nabla^2 f(\mathbf{x}_k)$



Credit: [Arezki et al., 2018]

Ideas behind TRM

Recall Taylor expansion $f(\mathbf{x} + \mathbf{d}) \approx f(\mathbf{x}) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle$

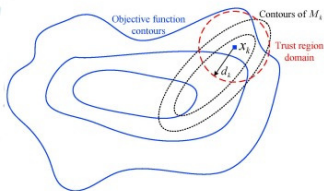
Start with \mathbf{x}_0 . Repeat the following:

- At \mathbf{x}_k , approximate f by the quadratic function (called **model function** dotted black)

$$m_k(\mathbf{d}) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$$

i.e., $m_k(\mathbf{d}) \approx f(\mathbf{x}_k + \mathbf{d})$, and \mathbf{B}_k to approximate $\nabla^2 f(\mathbf{x}_k)$

- Minimize $m_k(\mathbf{d})$ within a **trust region** $\{\mathbf{d} : \|\mathbf{d}\| \leq \Delta\}$, i.e., a norm ball (in red), to obtain \mathbf{d}_k



Credit: [Arezki et al., 2018]

Ideas behind TRM

Recall Taylor expansion $f(\mathbf{x} + \mathbf{d}) \approx f(\mathbf{x}) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle$

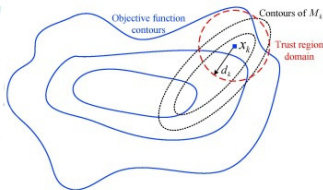
Start with \mathbf{x}_0 . Repeat the following:

- At \mathbf{x}_k , approximate f by the quadratic function (called **model function** dotted black)

$$m_k(\mathbf{d}) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$$

i.e., $m_k(\mathbf{d}) \approx f(\mathbf{x}_k + \mathbf{d})$, and \mathbf{B}_k to approximate $\nabla^2 f(\mathbf{x}_k)$

- Minimize $m_k(\mathbf{d})$ within a **trust region** $\{\mathbf{d} : \|\mathbf{d}\| \leq \Delta\}$, i.e., a norm ball (in red), to obtain \mathbf{d}_k
- If the approximation is inaccurate, decrease the region size; if the approximation is sufficiently accurate, increase the region size.



Credit: [Arezki et al., 2018]

Ideas behind TRM

Recall Taylor expansion $f(\mathbf{x} + \mathbf{d}) \approx f(\mathbf{x}) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle$

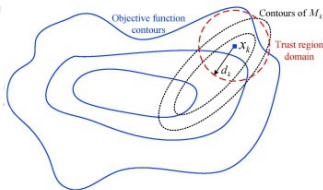
Start with \mathbf{x}_0 . Repeat the following:

- At \mathbf{x}_k , approximate f by the quadratic function (called **model function** dotted black)

$$m_k(\mathbf{d}) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$$

i.e., $m_k(\mathbf{d}) \approx f(\mathbf{x}_k + \mathbf{d})$, and \mathbf{B}_k to approximate $\nabla^2 f(\mathbf{x}_k)$

- Minimize $m_k(\mathbf{d})$ within a **trust region** $\{\mathbf{d} : \|\mathbf{d}\| \leq \Delta\}$, i.e., a norm ball (in red), to obtain \mathbf{d}_k
- If the approximation is inaccurate, decrease the region size; if the approximation is sufficiently accurate, increase the region size.
- If the approximation is reasonably accurate, update the iterate $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$.



Credit: [Arezki et al., 2018]

Framework of trust-region methods

To measure approximation quality: $\rho_k \doteq \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{d}_k)}{m_k(\mathbf{0}) - m_k(\mathbf{d}_k)} = \frac{\text{actual decrease}}{\text{model decrease}}$

Framework of trust-region methods

To measure approximation quality: $\rho_k \doteq \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{d}_k)}{m_k(\mathbf{0}) - m_k(\mathbf{d}_k)} = \frac{\text{actual decrease}}{\text{model decrease}}$

A generic trust-region algorithm

Input: \mathbf{x}_0 , radius cap $\widehat{\Delta} > 0$, initial radius Δ_0 , acceptance ratio $\eta \in [0, 1/4]$

```
1: for  $k = 0, 1, \dots$  do
2:    $\mathbf{d}_k = \arg \min_{\mathbf{d}} m_k(\mathbf{d})$ , s. t.  $\|\mathbf{d}\| \leq \Delta_k$  (TR Subproblem)
3:   if  $\rho_k < 1/4$  then
4:      $\Delta_{k+1} = \Delta_k/4$ 
5:   else
6:     if  $\rho_k > 3/4$  and  $\|\mathbf{d}_k\| = \Delta_k$  then
7:        $\Delta_{k+1} = \min(2\Delta_k, \widehat{\Delta})$ 
8:     else
9:        $\Delta_{k+1} = \Delta_k$ 
10:    end if
11:  end if
12:  if  $\rho_k > \eta$  then
13:     $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$ 
14:  else
15:     $\mathbf{x}_{k+1} = \mathbf{x}_k$ 
16:  end if
17: end for
```

Why TRM?

Recall the model function $m_k(\mathbf{d}) \doteq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$

Why TRM?

Recall the model function $m_k(\mathbf{d}) \doteq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$

- Take $\mathbf{B}_k = \nabla^2 f(\mathbf{x}_k)$

Why TRM?

Recall the model function $m_k(\mathbf{d}) \doteq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$

- Take $\mathbf{B}_k = \nabla^2 f(\mathbf{x}_k)$
- Gradient descent: stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$

Why TRM?

Recall the model function $m_k(\mathbf{d}) \doteq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$

- Take $\mathbf{B}_k = \nabla^2 f(\mathbf{x}_k)$
- Gradient descent: stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$
- Newton's method: $[\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$ may just stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$ or be ill-defined

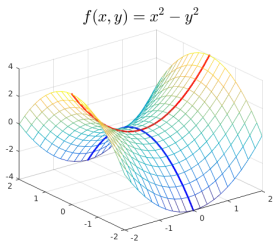
Why TRM?

Recall the model function $m_k(\mathbf{d}) \doteq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$

- Take $\mathbf{B}_k = \nabla^2 f(\mathbf{x}_k)$
- Gradient descent: stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$
- Newton's method: $[\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$ may just stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$ or be ill-defined
- Trust-region method: $\min_{\mathbf{d}} m_k(\mathbf{d}) \quad \text{s. t.} \quad \|\mathbf{d}\| \leq \Delta_k$

When $\nabla f(\mathbf{x}_k) = \mathbf{0}$,

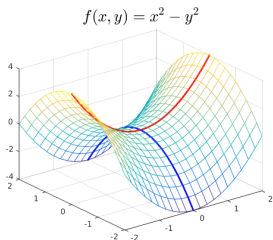
$$m_k(\mathbf{d}) - f(\mathbf{x}_k) = \frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle.$$



Why TRM?

Recall the model function $m_k(\mathbf{d}) \doteq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$

- Take $\mathbf{B}_k = \nabla^2 f(\mathbf{x}_k)$
- Gradient descent: stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$
- Newton's method: $[\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$ may just stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$ or be ill-defined
- Trust-region method: $\min_{\mathbf{d}} m_k(\mathbf{d}) \quad \text{s. t.} \quad \|\mathbf{d}\| \leq \Delta_k$



When $\nabla f(\mathbf{x}_k) = \mathbf{0}$,

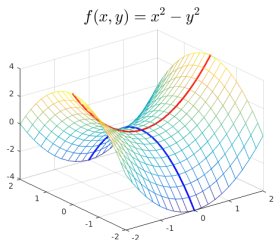
$$m_k(\mathbf{d}) - f(\mathbf{x}_k) = \frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle.$$

If $\nabla^2 f(\mathbf{x}_k)$ has **negative eigenvalues**, i.e., there are negative directional curvatures, $\frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle < 0$ for certain choices of \mathbf{d} (e.g., eigenvectors corresponding to the negative eigenvalues)

Why TRM?

Recall the model function $m_k(\mathbf{d}) \doteq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \langle \mathbf{d}, \mathbf{B}_k \mathbf{d} \rangle$

- Take $\mathbf{B}_k = \nabla^2 f(\mathbf{x}_k)$
- Gradient descent: stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$
- Newton's method: $[\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$ may just stop at $\nabla f(\mathbf{x}_k) = \mathbf{0}$ or be ill-defined
- Trust-region method: $\min_{\mathbf{d}} m_k(\mathbf{d}) \quad \text{s. t.} \quad \|\mathbf{d}\| \leq \Delta_k$



When $\nabla f(\mathbf{x}_k) = \mathbf{0}$,

$$m_k(\mathbf{d}) - f(\mathbf{x}_k) = \frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle.$$

If $\nabla^2 f(\mathbf{x}_k)$ has **negative eigenvalues**, i.e., there are negative directional curvatures, $\frac{1}{2} \langle \mathbf{d}, \nabla^2 f(\mathbf{x}_k) \mathbf{d} \rangle < 0$ for certain choices of \mathbf{d} (e.g., eigenvectors corresponding to the negative eigenvalues)

TRM can help to move away from “nice” saddle points!

To learn more about TRM

- A comprehensive reference [Conn et al., 2000]
- A closely-related alternative: cubic regularized second-order (CRSOM) method [Nesterov and Polyak, 2006, Agarwal et al., 2018]
- Example implementation of both TRM and CRSOM: Manopt (in Matlab) <https://www.manopt.org/> (choosing the Euclidean manifold)

- [Agarwal et al., 2018] Agarwal, N., Boumal, N., Bullins, B., and Cartis, C. (2018). **Adaptive regularization with cubics on manifolds.** *arXiv:1806.00065*.
- [Arezki et al., 2018] Arezki, Y., Nouira, H., Anwer, N., and Mehdi-Souzani, C. (2018). **A novel hybrid trust region minimax fitting algorithm for accurate dimensional metrology of aspherical shapes.** *Measurement*, 127:134–140.
- [Beck, 2017] Beck, A. (2017). **First-Order Methods in Optimization.** Society for Industrial and Applied Mathematics.
- [Conn et al., 2000] Conn, A. R., Gould, N. I. M., and Toint, P. L. (2000). **Trust Region Methods.** Society for Industrial and Applied Mathematics.
- [Hillar and Lim, 2013] Hillar, C. J. and Lim, L.-H. (2013). **Most tensor problems are NP-hard.** *Journal of the ACM*, 60(6):1–39.
- [Murty and Kabadi, 1987] Murty, K. G. and Kabadi, S. N. (1987). **Some NP-complete problems in quadratic and nonlinear programming.** *Mathematical Programming*, 39(2):117–129.
- [Nesterov, 2018] Nesterov, Y. (2018). **Lectures on Convex Optimization.** Springer International Publishing.

- [Nesterov and Polyak, 2006] Nesterov, Y. and Polyak, B. (2006). **Cubic regularization of newton method and its global performance.** *Mathematical Programming*, 108(1):177–205.
- [Nocedal and Wright, 2006] Nocedal, J. and Wright, S. J. (2006). **Numerical Optimization.** Springer New York.
- [Wright, 2015] Wright, S. J. (2015). **Coordinate descent algorithms.** *Mathematical Programming*, 151(1):3–34.