# LECTURE 19: ROOT-FINDING AND MINIMIZATION

STAT 545: INTRO. TO COMPUTATIONAL STATISTICS

Vinayak Rao

Purdue University

November 13, 2019

## Root-finding in one-dimension

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

$$f(x) = 0$$

Invariably need iterative methods.

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

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

Better: faster (asymptotic) convergence.

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

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

Finding an initial bracketing can be non-trivial.
Typically, start with an initial interval and expand or contract.

Below, we assume we have an initial bracketing.

# Root bracketing

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

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

Finding an initial bracketing can be non-trivial.
Typically, start with an initial interval and expand or contract.

Below, we assume we have an initial bracketing.

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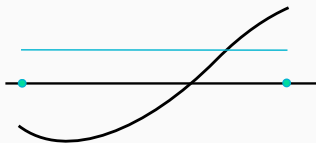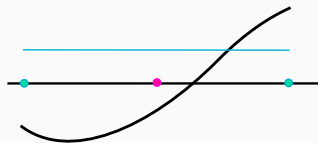
# Bisection method

Simplest root-finding algorithm.
Given an initial bracketing, cannot fail.
But is slower than other methods.

Successively halves the bracketing interval (binary search):



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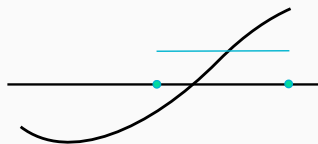  (whichever is a valid bracketing)

# Bisection method

Simplest root-finding algorithm.
Given an initial bracketing, cannot fail.
But is slower than other methods.

Successively halves the bracketing interval (binary search):



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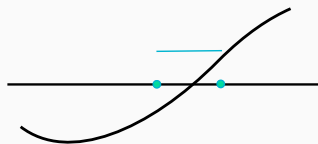  (whichever is a valid bracketing)

# Bisection method

Simplest root-finding algorithm.
Given an initial bracketing, cannot fail.
But is slower than other methods.

Successively halves the bracketing interval (binary search):



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

# Bisection method

Simplest root-finding algorithm.
Given an initial bracketing, cannot fail.
But is slower than other methods.

Successively halves the bracketing interval (binary search):



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

# Bisection method (contd)

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

Upperbounds error in root.

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

# Bisection method (contd)

Let $\epsilon_n$ be the interval length at iteration $n$.
Upperbounds error in root.

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

Linear convergence:

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

Let $\epsilon_n$ be the interval length at iteration *n*.
Upperbounds error in root.

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

Linear convergence:

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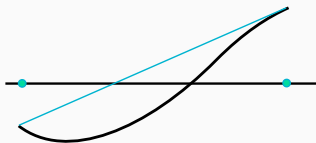
Superlinear convergence:

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

Quadratic convergence:
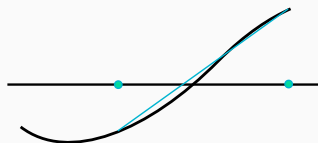Number of significant figures *doubles* every iteration.

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

## SECANT METHOD AND BISECTION METHOD

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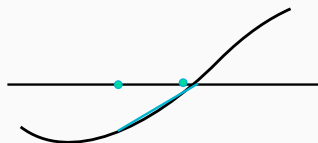


Secant method:

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

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

- Bracketing (and thus convergence) not guaranteed.

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

Secant method:
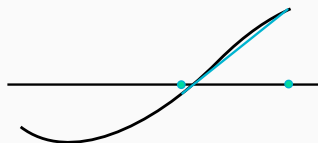


- always keep the newest point
- Superlinear convergence ($m = 1.618$, the golden ratio)

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

- Bracketing (and thus convergence) not guaranteed.

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



Secant method:

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

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

- Bracketing (and thus convergence) not guaranteed.

False position:

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

In practice, people use more sophiticated algorithms.

Most popular is Brent's method.

Maintains bracketing by combining bisection method with a quadratic approximation.

Lots of book-keeping.

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

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

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

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

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

Assume second- and higher-order terms are negligible.
Given $x_i$, choose $x_{i+1} = x_i + \delta$ so that $f(x_{i+1}) = 0$:

# Newton's method (a.k.a. Newton-Raphson)

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

Taylor expansion: $\quad f(x + \delta) = f(x) + \delta f'(x) + \frac{\delta^2}{2} f''(x) + \cdots$
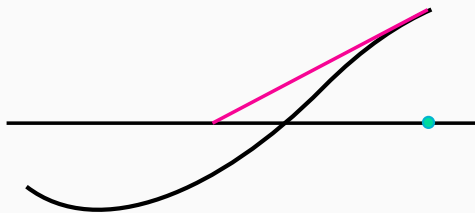
Assume second- and higher-order terms are negligible.
Given $x_i$, choose $x_{i+1} = x_i + \delta$ so that $f(x_{i+1}) = 0$:

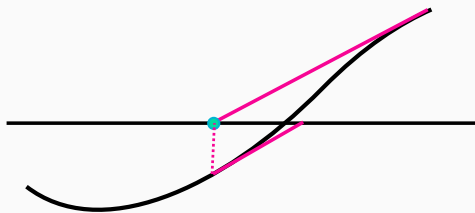$$0 = f(x_i) + \delta f'(x_i)$$

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

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

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

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

Letting $x^*$ be the root, we have

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

# Convergence of Newton's method

Letting $x^*$ be the root, we have

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

Also since $x_i = x^* + \epsilon_i$,

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

## Convergence of Newton's method

Letting $x^*$ be the root, we have

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

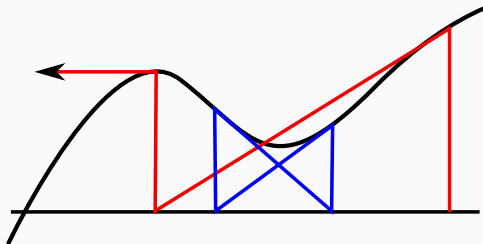Also since $x_i = x^* + \epsilon_i$,

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

This gives

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

## Convergence of Newton's method

Letting $x^*$ be the root, we have

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

Also since $x_i = x^* + \epsilon_i$,

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

This gives

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

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

Away from the root the linear approximation can be bad.

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

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

Often used in combination with some bracketing method.

# Root-finding for systems of nonlinear equations

Now have $N$ functions $F_1, F_2, \cdots, F_N$ of $N$ variables $x_1, x_2, \cdots, x_N$

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

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

Much harder than the 1-d case.

Much harder than optimization.

Again, consider a Taylor expansion:

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

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

Again, consider a Taylor expansion:

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

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

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

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

Solve $\delta\mathbf{x} = -\mathbf{J}(\mathbf{x})^{-1} \cdot \mathbf{F}(\mathbf{x})$ (e.g. by LU decomposition)

## Newton's method

Again, consider a Taylor expansion:

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

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

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

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

Solve $\delta\mathbf{x} = -\mathsf{J}(\mathbf{x})^{-1} \cdot \mathsf{F}(\mathbf{x})$ (e.g. by LU decomposition)

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

Can wildly careen through space if not careful.

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

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

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

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

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

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

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

This keeps our global objective in sight.

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

This keeps our global objective in sight.



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

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

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



Local minimum

Global minimum

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

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



Local minimum

Global minimum

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

Conceptually (deceptively?) simpler than EM.

# Gradient descent (iterative method)

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

Update $x_{new}$ as $\qquad x_{new} = x_{old} - \eta \left. \frac{\mathrm{d}f}{\mathrm{d}x} \right|_{x_{old}}$

The steeper the slope, the bigger the move.

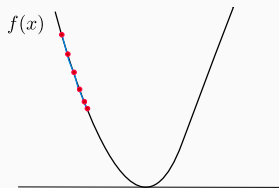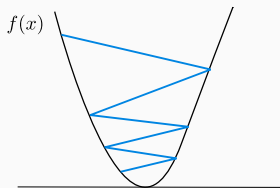# Gradient descent (iterative method)

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

Update $x_{new}$ as

$$x_{new} = x_{old} - \eta \left. \frac{\mathrm{d}f}{\mathrm{d}x} \right|_{x_{old}}$$

The steeper the slope, the bigger the move.

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

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

Update $x_{new}$ as $\qquad x_{new} = x_{old} - \eta \left. \dfrac{\mathrm{d}f}{\mathrm{d}x} \right|_{x_{old}}$

The steeper the slope, the bigger the move.

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

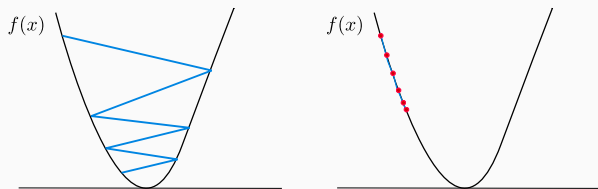Choosing $\eta$ is a dark art:

# GRADIENT DESCENT (ITERATIVE METHOD)

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

Update $x_{new}$ as $\qquad x_{new} = x_{old} - \eta \left. \frac{\mathrm{d}f}{\mathrm{d}x} \right|_{x_{old}}$

The steeper the slope, the bigger the move.

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

Choosing $\eta$ is a dark art:



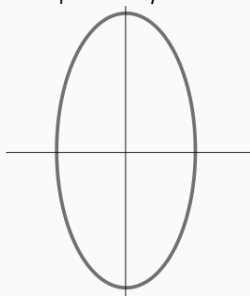Better methods adapt step-size according to the curvature of $f$.
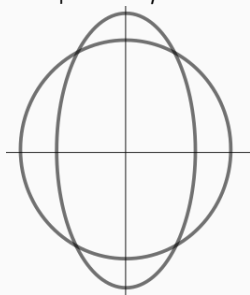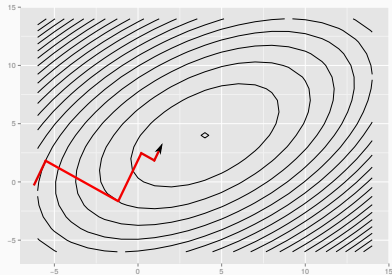
# Gradient descent in higher-dimensions

Steepest descent also applies to higher dimensions too:

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

At each step, solve a 1-d problem along the gradient

Now, even the optimal step-size $\eta$ can be inefficient:

# Gradient descent in higher-dimensions

Steepest descent also applies to higher dimensions too:
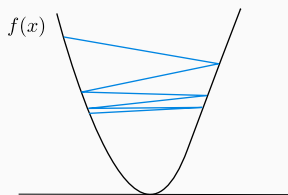
$$x_{new} = x_{old} - \eta \ \nabla f|_{x_{old}}$$

At each step, solve a 1-d problem along the gradient
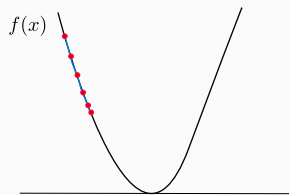
Now, even the optimal step-size $\eta$ can be inefficient:

Steepest descent also applies to higher dimensions too:

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

At each step, solve a 1-d problem along the gradient

Now, even the optimal step-size $\eta$ can be inefficient:

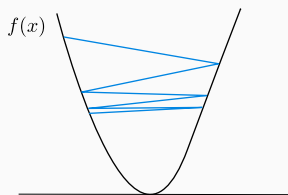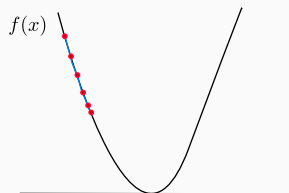Rather than the best step-size each step, find a decent solution



Big steps with little decrease



Small steps getting us nowhere

Rather than the best step-size each step, find a decent solution



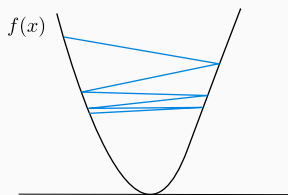Big steps with little decrease



Small steps getting us nowhere

Avg. decrease at least some fraction of initial rate:

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

Rather than the best step-size each step, find a decent solution



Big steps with little decrease



Small steps getting us nowhere

Avg. decrease at least some fraction of initial rate:

$$f(\mathbf{x} + \lambda \delta \mathbf{x}) \leq f(\mathbf{x}) + c_1 \lambda (\nabla f \cdot \delta \mathbf{x}), \qquad c_1 \in (0, 1) \; e.g. \; 0.9$$
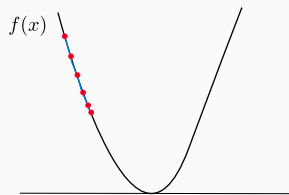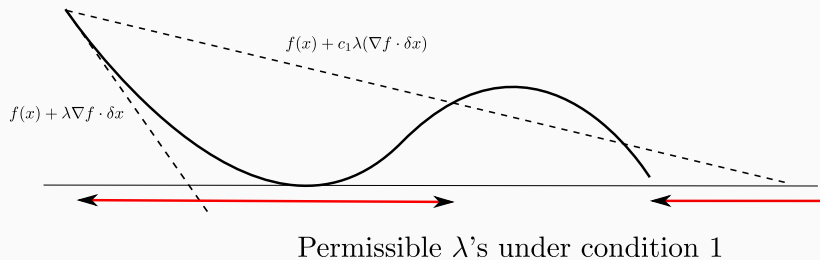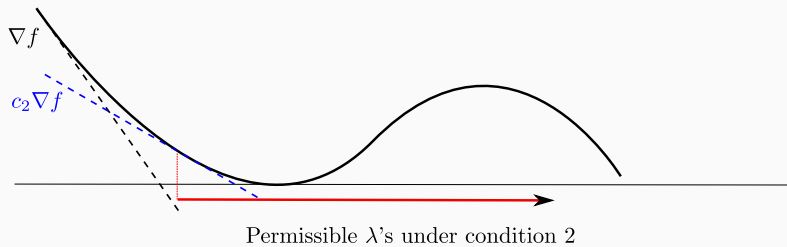
Final rate is greater than some fraction of initial rate:

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

$f(x) + c_1\lambda(\nabla f \cdot \delta x)$

$f(x) + \lambda \nabla f \cdot \delta x$

Permissible $\lambda$'s under condition 1
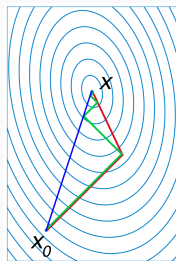
Permissible $\lambda$'s under condition 2

A simple way to satisfy Wolfe conditions:

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

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

# CONJUGATE GRADIENT DESCENT

Consider minimizing $\frac{1}{2}x^T A x - b^T x$:



Steepest descent can take many steps to get to the minimum
Problem: After minimizing along a direction, gradient is perpendicular to previous direction (why)
· Can 'cancel' out earlier gains

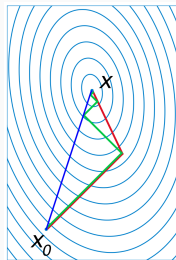A popular algorithm is conjugate gradient descent
Sequentially updates along directions $p_1, \cdots p_N$:

$$x_{t+1} = x_t + \lambda_{t+1} p_t, \text{ where } \lambda_{t+1} = \text{argmin}_\lambda \ f(x_t + \lambda p_t)$$
$$p_{t+1} = \nabla f(x_{t+1}) + \frac{\langle \nabla f(x_{t+1}), \nabla f(x_{t+1}) \rangle}{\langle \nabla f(x_t), \nabla f(x_t) \rangle} p_t$$

# CONJUGATE GRADIENT DESCENT

Consider minimizing $\frac{1}{2}x^T A x - b^T x$:



Steepest descent can take many steps to get to the minimum

Problem: After minimizing along a direction, gradient is perpendicular to previous direction (why)

· Can 'cancel' out earlier gains

If $f(x) = \frac{1}{2}x^T A x - b^T x$, $x \in \mathbb{R}^d$, CG takes max $d$ steps to converge

Can show the directions satisfy $\langle p_{t+1}, p_t \rangle_A := p_{t+1}^T A p_t = 0$

(this is unlike $p_{t+1}^T p_t = 0$ for steepest descent)