

Convexity, Local and Global Optimality, etc.

Descent Algorithms

- A single iteration of the general descent algorithm consists of two main steps, *viz.*,
 - ① determining a good descent direction $\Delta \mathbf{x}^{(k)}$, which is typically forced to have unit norm and
 - ② determining the step size using some line search technique.
- If the function f is convex, from the necessary and sufficient condition for convexity restated here for reference:

& differentiable f

$$f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)}) + \nabla^T f(\mathbf{x}^{(k)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$

- We require that $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ and since $t^{(k)} > 0$, we must have

....as a necessary condition

Descent Algorithms

- A single iteration of the general descent algorithm consists of two main steps, *viz.*,
 - 1 determining a good descent direction $\Delta \mathbf{x}^{(k)}$, which is typically forced to have unit norm and
 - 2 determining the step size using some line search technique.
- If the function f is convex, from the necessary and sufficient condition for convexity restated here for reference:

$$f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)}) + \nabla^T f(\mathbf{x}^{(k)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$

- We require that $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ and since $t^{(k)} > 0$, we must have

$$\nabla^T f(\mathbf{x}^{(k)}) \Delta \mathbf{x}^{(k)} < 0$$

That is, the descent direction $\Delta \mathbf{x}^{(k)}$ must make an obtuse angle ($\theta \in \left(\frac{\pi}{2}, \frac{3\pi}{2}\right)$) with the gradient vector (in several settings, the dot product should be sufficiently negative to ensure descent)

Descent Algorithms (contd.)

Find a starting point $\mathbf{x}^{(0)} \in \mathcal{D}$

repeat

1. Determine $\Delta \mathbf{x}^{(k)}$.
2. Choose a step size $t^{(k)} > 0$ using ray^a search.
3. Obtain $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$.
4. Set $k = k + 1$.

until stopping criterion (such as $\|\nabla f(\mathbf{x}^{(k+1)})\| < \epsilon$) is satisfied

^aMany textbooks refer to this as line search, but we prefer to call it ray search, since the step must be positive.

Figure 7: The general descent algorithm.

There are many different empirical techniques for ray search, though it matters much less than the search for the descent direction. These techniques reduce the n -dimensional problem to a 1-dimensional problem, which can be easy to solve by use of plotting and eyeballing or even exact search.

Finding the step size t

- If t is too large, we get diverging updates of x
- If t is too small, we get a very slow descent
- We need to find a t that is *just right*
- We discuss two ways of finding t :
 - ① Exact ray search
 - ② Backtracking ray search

Exact ray search

$$\begin{aligned} t^{k+1} &= \operatorname{argmin}_t f\left(x^k - t\nabla f(x^k)\right) \\ &\quad + t \Delta x^k \\ &= \operatorname{argmin}_t \phi(t) \end{aligned}$$

- This method gives the most optimal step size in the given descent direction $\nabla f(x^k)$
- It ensures that $f(x^{k+1}) \leq f(x^k)$
- If f is itself quadratic, it gives an optimal solution to the minimization of f (since the quadratic approximation f_Q would become exact and no longer approximate)

H/W

Ray Search for Descent: Options

- ① **Exact ray search:** The exact ray search seeks a scaling factor t that satisfies

$$t = \underset{t>0}{\operatorname{argmin}} f(\mathbf{x} + t\Delta\mathbf{x}) \quad (26)$$

Not practical if f is not quadratic or some special form that yields itself to a closed form for this one dimensional optimization problem
Or if numerically solving this exact problem is expensive at every step

Ray Search for Descent: Options

- ① **Exact ray search:** The exact ray search seeks a scaling factor t that satisfies

$$t = \underset{t > 0}{\operatorname{argmin}} f(\mathbf{x} + t\Delta\mathbf{x}) \quad (26)$$

- ② **Backtracking ray search:** The exact line search may not be feasible or could be expensive to compute for complex non-linear functions. A relatively simpler ray search iterates over values of step size starting from 1 and scaling it down by a factor of $\beta \in (0, \frac{1}{2})$ after every iteration till the following condition, called the *Armijo condition* is satisfied for some $0 < c_1 < 1$.

$$f(\mathbf{x} + t\Delta\mathbf{x}) < f(\mathbf{x}) + c_1 t \nabla^T f(\mathbf{x}) \Delta\mathbf{x} \quad (27)$$

Based on first order convexity condition, it can be inferred that when $c_1 = 1$,

Ray Search for Descent: Options

- ① **Exact ray search:** The exact ray search seeks a scaling factor t that satisfies

$$t = \underset{t>0}{\operatorname{argmin}} f(\mathbf{x} + t\Delta\mathbf{x}) \quad (26)$$

- ② **Backtracking ray search:** The exact line search may not be feasible or could be expensive to compute for complex non-linear functions. A relatively simpler ray search iterates over values of step size starting from 1 and scaling it down by a factor of $\beta \in (0, \frac{1}{2})$ after every iteration till the following condition, called the *Armijo condition* is satisfied for some $0 < c_1 < 1$.

$$f(\mathbf{x} + t\Delta\mathbf{x}) < f(\mathbf{x}) + c_1 t \nabla^T f(\mathbf{x}) \Delta\mathbf{x} \quad (27)$$

Based on first order convexity condition, it can be inferred that when $c_1 = 1$, the right hand side of (27) gives a lower bound on the value of $f(\mathbf{x} + t\Delta\mathbf{x})$ and hence $c_1 < 1$

Ray Search for Descent: Options

- ① **Exact ray search:** The exact ray search seeks a scaling factor t that satisfies

$$t = \underset{t>0}{\operatorname{argmin}} f(\mathbf{x} + t\Delta\mathbf{x}) \quad (26)$$

- ② **Backtracking ray search:** The exact line search may not be feasible or could be expensive to compute for complex non-linear functions. A relatively simpler ray search iterates over values of step size starting from 1 and scaling it down by a factor of $\beta \in (0, \frac{1}{2})$ after every iteration till the following condition, called the *Armijo condition* is satisfied for some $0 < c_1 < 1$.

Keep scaling down t until you obtain an f that is upper bounded by the supposed lower bound were t to be further scaled down.

$$f(\mathbf{x} + t\Delta\mathbf{x}) < f(\mathbf{x}) + c_1 t \nabla^T f(\mathbf{x}) \Delta\mathbf{x} \quad (27)$$

Based on first order convexity condition, it can be inferred that when $c_1 = 1$, the right hand side of (27) gives a lower bound on the value of $f(\mathbf{x} + t\Delta\mathbf{x})$ and hence (27) can never hold. The Armijo condition simply ensures that t decreases f sufficiently.

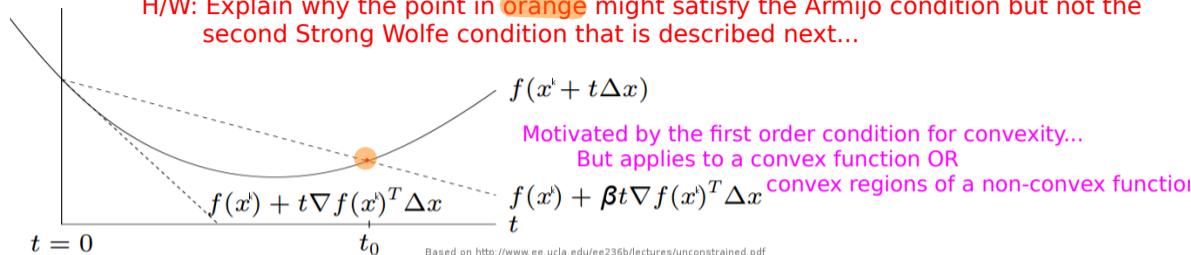
Backtracking ray search

- The algorithm

- ▶ Choose a $\beta \in (0, 1)$
- ▶ Start with $t = 1$
- ▶ While $f(\mathbf{x} + t\Delta\mathbf{x}) \geq f(\mathbf{x}) + c_1 t \nabla^T f(\mathbf{x}) \Delta\mathbf{x}$, do or Until the Armijo condition is met
 - ★ Update $t \leftarrow \beta t$

Interpretation of backtracking line search

H/W: Explain why the point in orange might satisfy the Armijo condition but not the second Strong Wolfe condition that is described next...



- $\Delta x =$ direction of descent $= -\nabla f(x^k)$ for gradient descent
- A different way of understanding the varying step size with β : Multiplying t by β causes the interpolation to tilt as indicated in the figure

Homeworks: Let $f(x) = x^2$ for $x \in \mathfrak{R}$. Let $x^0 = 2$, $\Delta x^k = -1$ for all k (since it is a valid descent direction of $x > 0$) and $x^k = 1 + 2^{-k}$. What is the step size t^k implicitly being used. Show that while t^k satisfies the Armijo condition (determine a c_1), it does not satisfy the second Strong Wolfe condition in the following slides. Why is the choice of step size bad?

Ray Search for First Order Descent: Strong Wolfe Conditions

Often, another condition is used for inexact line search in conjunction with the Armijo condition.

That is, the linear lower bound approximator has become sufficiently close to the x-y plane..

$$\left| \Delta \mathbf{x}^T \nabla f(\mathbf{x} + t\Delta \mathbf{x}) \right| \leq c_2 \left| \Delta \mathbf{x}^T \nabla f(\mathbf{x}) \right| \quad (28)$$

where $1 > c_2 > c_1 > 0$. This condition ensures that the slope of the function $f(\mathbf{x} + t\Delta \mathbf{x})$ at t is less than c_2 times that at $t = 0$.

- 1 The conditions in (27) and (28) are together called the strong Wolfe conditions. These conditions are particularly very important for non-convex problems.
- 2 While (11refeq:armijoCondition) ensures guaranteed decrease in $f(\mathbf{x} + \delta \mathbf{x})$, (28) provides guaranteed decrease in magnitude of slope and avoid too small steps.
- 3 Claim: If $1 > c_2 > c_1 > 0$ and the function $f(\mathbf{x})$ is convex and differentiable, there exists t such that (27) and (28) are both satisfied for any f . *Hint: Use the Mean Value Theorem*

convex

Convexity \Rightarrow Strong Wolfe Conditions

- Let $\phi(t) = f(\mathbf{x}^k + t\Delta\mathbf{x}^k) \geq f(\mathbf{x}^k) + t\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k$ (where the second inequality is by virtue of convexity). Remember that $\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k < 0$
- Since $0 < c_1 < 1$, the linear approximation $l(t) = f(\mathbf{x}^k) + tc_1\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k$ is unbounded below and it can be shown to intersect the convex function f for some t'

Convexity \Rightarrow Strong Wolfe Conditions

- Let $\phi(t) = f(\mathbf{x}^k + t\Delta\mathbf{x}^k) \geq f(\mathbf{x}^k) + t\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k$ (where the second inequality is by virtue of convexity). Remember that $\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k < 0$
- Since $0 < c_1 < 1$, the linear approximation $l(t) = f(\mathbf{x}^k) + tc_1\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k$ is unbounded below and it can be shown to intersect the graph of ϕ at least once.
- Let $t' > 0$ be the smallest intersecting value of t , that is:

$$f(\mathbf{x} + t'\Delta\mathbf{x}^k) = f(\mathbf{x}^k) + t'c_1\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k \quad (29)$$

- For all $t \in [0, t']$,

The curve must lie below this tangent/line

Convexity \Rightarrow Strong Wolfe Conditions $0 < c_2 < c_1 < 1$

- Let $\phi(t) = f(\mathbf{x}^k + t\Delta\mathbf{x}^k) \geq f(\mathbf{x}^k) + t\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k$ (where the second inequality is by virtue of convexity). Remember that $\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k < 0$
- Since $0 < c_1 < 1$, the linear approximation $l(t) = f(\mathbf{x}^k) + tc_1\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k$ is unbounded below and it can be shown to intersect the graph of ϕ at least once.
- Let $t' > 0$ be the smallest intersecting value of t , that is:

$$f(\mathbf{x} + t'\Delta\mathbf{x}^k) = f(\mathbf{x}^k) + t'c_1\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k \quad (29)$$

- For all $t \in [0, t']$,

$$f(\mathbf{x}^k + t\Delta\mathbf{x}^k) \leq f(\mathbf{x}^k) + tc_1\nabla^T f(\mathbf{x}^k)\Delta\mathbf{x}^k \quad (30)$$

That is, there exists a non-empty set of t such that the first Wolfe condition is met.

- By the mean value theorem, $\exists t'' \in (0, t')$ such that

$$f(\mathbf{x}^k + t'\Delta\mathbf{x}^k) - f(\mathbf{x}^k) = t'\nabla^T f(\mathbf{x}^k + t''\Delta\mathbf{x}^k)\Delta\mathbf{x}^k \quad (31)$$

Convexity \Rightarrow Strong Wolfe Conditions (contd.)

- Combining (29) and (31), and using $c_1 \succ c_2$, and $\nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k < 0$

$$\nabla^T f(\mathbf{x}^k + t'' \Delta \mathbf{x}^k) \Delta \mathbf{x}^k = c_1 \nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k > c_2 \nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k \quad (32)$$

- Again, since $\nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k < 0$, we get the $t^k = t''$ satisfying (28)

$$|\nabla^T f(\mathbf{x}^k + t'' \Delta \mathbf{x}^k) \Delta \mathbf{x}^k| < c_2 |\nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k| \quad (33)$$

- In fact, by continuity of $f(\cdot)$, there exists an interval around t'' for which Strong Wolfe conditions hold.

Empirical Observations on Ray Search

- A finding that is borne out of plenty of empirical evidence is that exact ray search does better than empirical ray search in a few cases only.
- Further, the exact choice of the value of β and α seems to have little effect on the convergence of the overall descent method.
- The trend of specific descent methods has been like a parabola - starting with simple steepest descent techniques, then accommodating the curvature hessian matrix through a more sophisticated Newton's method and finally, trying to simplify the Newton's method through approximations to the hessian inverse, culminating in conjugate gradient techniques, that do away with any curvature matrix whatsoever, and form the internal combustion engine of many sophisticated optimization techniques today.
- We start the thread by describing the steepest descent methods.

Algorithms: Steepest Descent

- The idea of steepest descent is to determine a descent direction such that for a unit step in that direction, the prediction of decrease in the objective is maximized

- However, consider $\Delta x = \operatorname{argmin}_v \begin{bmatrix} -5 & 10 & 15 \end{bmatrix} v$

$$\implies \Delta x = \begin{bmatrix} \infty \\ -\infty \\ -\infty \end{bmatrix}$$

which is unacceptable

- Thus, there is a necessity to restrict the norm of v
- The choice of the descent direction can be stated as:

$$\Delta x = \operatorname{argmin}_v \nabla^T f(x) v$$

$$\text{s.t. } \|v\| = 1$$

Algorithms: Steepest Descent

- Let $\mathbf{v} \in \mathbb{R}^n$ be a unit vector under some norm. By first order convexity condition for convex and differentiable f ,

$$f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k)} + \mathbf{v}) \leq -\nabla^T f(\mathbf{x}^{(k)})\mathbf{v}$$

- For small \mathbf{v} , the inequality turns into approximate equality. The term $-\nabla^T f(\mathbf{x}^{(k)})\mathbf{v}$ can be thought of as (an upper-bound on) the first order prediction of decrease.
- The idea in the steepest descent method is to choose a norm and then determine a descent direction such that for a unit step in that norm, the first order prediction of decrease is maximized. This choice of the descent direction can be stated as

$$\Delta\mathbf{x} = \operatorname{argmin} \left\{ \nabla^T f(\mathbf{x})\mathbf{v} \mid \|\mathbf{v}\| = 1 \right\}$$

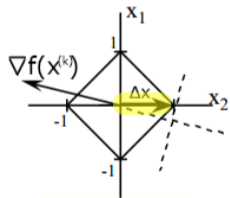
- Empirical observation:* If the norm chosen is aligned with the gross geometry of the sub-level sets³, the steepest descent method converges faster to the optimal solution. Else, it often amplifies the effect of oscillations.

³The alignment can be determined by fitting, for instance, a quadratic to a sample of the points.

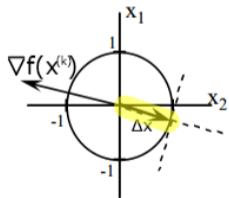
Various choices of the norm result in different solutions for Δx

- For 2-norm, $\Delta x = -\frac{\nabla f(x^{(k)})}{\|\nabla f(x^{(k)})\|_2}$
(gradient descent)
- For 1-norm, $\Delta x = -\text{sign}\left(\frac{\partial f(x^{(k)})}{\partial x_i^{(k)}}\right) e_i$, where e_i is the i th standard basis vector
(coordinate descent)
- For ∞ -norm, $\Delta x = -\text{sign}(\nabla f(x^{(k)}))$

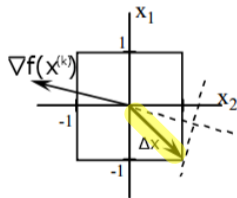
Where the gradient has maximum component along i^{th} dimension



SDD in L1-norm



SDD in L2-norm



SDD in L ∞ -norm

General Algorithm: Steepest Descent (contd)

Find a starting point $\mathbf{x}^{(0)} \in \mathcal{D}$.

repeat

1. Set $\Delta \mathbf{x}^{(k)} = \operatorname{argmin} \left\{ \nabla^T f(\mathbf{x}^{(k)}) \mathbf{v} \mid \|\mathbf{v}\| = 1 \right\}$.
2. Choose a step size $t^{(k)} > 0$ using exact or backtracking ray search.
3. Obtain $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$.
4. Set $k = k + 1$.

until stopping criterion (such as $\|\nabla f(\mathbf{x}^{(k+1)})\| \leq \epsilon$) is satisfied

Figure 8: The steepest descent algorithm.

Two examples of the steepest descent method are the [gradient descent method \(for the euclidian or \$L_2\$ norm\)](#) and the [coordinate-descent method \(for the \$L_1\$ norm\)](#). One fact however is that no two norms should give exactly opposite steepest descent directions, though they may point in different directions.

Convergence of Descent Algorithm

- Consider the general descent algorithm ($\nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k < 0$ for each k) with each step:
 $\mathbf{x}^{k+1} = \mathbf{x}^k + t^k \Delta \mathbf{x}^k$.

- ▶ Suppose f is bounded below in \mathbb{R}^n and
- ▶ is continuously differentiable in an open set \mathcal{N} containing the level set $\{\mathbf{x} | f(\mathbf{x}) \leq f(\mathbf{x}^0)\}$
- ▶ ∇f is Lipschitz continuous.

Then, $\sum_{k=1}^{\infty} \frac{(\nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k)^2}{\|\Delta \mathbf{x}^k\|^2} < \infty$ (that is, it is finite)

- Thus, $\lim_{k \rightarrow \infty} \frac{\nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k}{\|\Delta \mathbf{x}^k\|} = 0$. Dot product between gradient and unit step $\rightarrow 0$
- If we additionally assume that the descent direction is not orthogonal to the gradient, *i.e.*,
 $-\frac{\nabla^T f(\mathbf{x}^k) \Delta \mathbf{x}^k}{\|\Delta \mathbf{x}^k\| \|\nabla f(\mathbf{x}^k)\|} \geq \Gamma$ for some $\Gamma > 0$, then, we can show that $\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}^k)\| = 0$
- Before we try and prove this result, let us discuss Lipschitz continuity (recall from midsem).