

ECEN 434: Optimization for ECEN

Lecture 10: Descent Methods for Unconstrained Problems

Tie Liu

Unconstrained minimization problems

- Consider the **unconstrained** optimization problem:

$$\text{minimize}_{\mathbf{x} \in \mathbb{X}} \quad f(\mathbf{x})$$

where f is **convex** and **twice continuously differentiable** and \mathbb{X} is a nonempty, open set

Unconstrained minimization problems

- Consider the **unconstrained** optimization problem:

$$\text{minimize}_{\mathbf{x} \in \mathbb{X}} \quad f(\mathbf{x})$$

where f is **convex** and **twice continuously differentiable** and \mathbb{X} is a nonempty, open set

- We will assume that the minimum value of f is **attainable**, i.e., there exists an optimal point \mathbf{x}^*

Unconstrained minimization problems

- Consider the **unconstrained** optimization problem:

$$\text{minimize}_{\mathbf{x} \in \mathbb{X}} \quad f(\mathbf{x})$$

where f is **convex** and **twice continuously differentiable** and \mathbb{X} is a nonempty, open set

- We will assume that the minimum value of f is **attainable**, i.e., there exists an optimal point \mathbf{x}^*
- Since f is convex and differentiable, a **necessary and sufficient** condition for a point \mathbf{x}^* to be optimal is

$$\nabla f(\mathbf{x}^*) = \mathbf{0}$$

Thus, solving the unconstrained minimization problem is the same as finding a solution to the above system of equations

- In a few special cases, we can find a solution to the unconstrained minimization problems by **analytically** solving the optimality equations, but usually the problem must be solved by an **iterative** algorithm

- In a few special cases, we can find a solution to the unconstrained minimization problems by **analytically** solving the optimality equations, but usually the problem must be solved by an **iterative** algorithm
- By this we mean an algorithm that computes a sequence of points $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots \in \mathbb{X}$ with $f(\mathbf{x}^{(k)}) \rightarrow p^*$ as $k \rightarrow \infty$. Such a sequence of points is called a **minimizing sequence** for the minimization problem

- In a few special cases, we can find a solution to the unconstrained minimization problems by **analytically** solving the optimality equations, but usually the problem must be solved by an **iterative** algorithm
- By this we mean an algorithm that computes a sequence of points $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots \in \mathbb{X}$ with $f(\mathbf{x}^{(k)}) \rightarrow p^*$ as $k \rightarrow \infty$. Such a sequence of points is called a **minimizing sequence** for the minimization problem
- The algorithm is terminated when $f(\mathbf{x}^{(k)}) - p^* \leq \epsilon$, where $\epsilon > 0$ is some specified **tolerance**

Descent methods

- We focus on the so-called **descent methods**, which start the algorithm by choosing a suitable starting point $\mathbf{x}^{(0)} \in \mathbb{X}$

Descent methods

- We focus on the so-called **descent methods**, which start the algorithm by choosing a suitable starting point $\mathbf{x}^{(0)} \in \mathbb{X}$
- The algorithm then produces a minimizing sequence through the following iterative update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$$

Descent methods

- We focus on the so-called **descent methods**, which start the algorithm by choosing a suitable starting point $\mathbf{x}^{(0)} \in \mathbb{X}$
- The algorithm then produces a minimizing sequence through the following iterative update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$$

- Here, $k = 0, 1, \dots$ denotes the **iteration number**, $\mathbf{x}^{(k)}$ is called the **step or search direction** (even though it need not have unit norm), and the scalar $t^{(k)} \geq 0$ is called the **step size or step length** at iteration k

Descent methods

- We focus on the so-called **descent methods**, which start the algorithm by choosing a suitable starting point $\mathbf{x}^{(0)} \in \mathbb{X}$
- The algorithm then produces a minimizing sequence through the following iterative update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$$

- Here, $k = 0, 1, \dots$ denotes the **iteration number**, $\mathbf{x}^{(k)}$ is called the **step or search direction** (even though it need not have unit norm), and the scalar $t^{(k)} \geq 0$ is called the **step size or step length** at iteration k
- The algorithm terminates when $\|\nabla f(\mathbf{x}^{(k)})\| \leq \eta$, where η is small and positive

Descent methods

- We focus on the so-called **descent methods**, which start the algorithm by choosing a suitable starting point $\mathbf{x}^{(0)} \in \mathbb{X}$
- The algorithm then produces a minimizing sequence through the following iterative update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}$$

- Here, $k = 0, 1, \dots$ denotes the **iteration number**, $\mathbf{x}^{(k)}$ is called the **step or search direction** (even though it need not have unit norm), and the scalar $t^{(k)} \geq 0$ is called the **step size or step length** at iteration k
- The algorithm terminates when $\|\nabla f(\mathbf{x}^{(k)})\| \leq \eta$, where η is small and positive
- The above iterative procedure is called a **descent** method if

$$f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$$

except when $\mathbf{x}^{(k)}$ is already optimal

Exact line search

- One method for determining the step sizes is **exact line search**, in which $t^{(k)}$ is chosen to minimize f along the ray $\{\mathbf{x}^{(k)} + t\Delta\mathbf{x}^{(k)} : t \geq 0\}$:

$$t^{(k)} = \arg \min_{t \geq 0} f(\mathbf{x}^{(k)} + t\Delta\mathbf{x}^{(k)})$$

Exact line search

- One method for determining the step sizes is **exact line search**, in which $t^{(k)}$ is chosen to minimize f along the ray $\{\mathbf{x}^{(k)} + t\Delta\mathbf{x}^{(k)} : t \geq 0\}$:

$$t^{(k)} = \arg \min_{t \geq 0} f(\mathbf{x}^{(k)} + t\Delta\mathbf{x}^{(k)})$$

- An exact line search is used when the cost of the minimization problem with one variable is low compared to the cost of computing the search direction itself. In some special cases, the minimizer along the ray can be found analytically and in others it can be computed efficiently

Gradient descent

- By **convexity** we have

$$f(\mathbf{y}) \geq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^t (\mathbf{y} - \mathbf{x}^{(k)})$$

Gradient descent

- By **convexity** we have

$$f(\mathbf{y}) \geq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^t (\mathbf{y} - \mathbf{x}^{(k)})$$

- Therefore, $\nabla f(\mathbf{x}^{(k)})^t (\mathbf{y} - \mathbf{x}^{(k)}) \geq 0$ implies $f(\mathbf{y}) \geq f(\mathbf{x}^{(k)})$

Gradient descent

- By **convexity** we have

$$f(\mathbf{y}) \geq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^t (\mathbf{y} - \mathbf{x}^{(k)})$$

- Therefore, $\nabla f(\mathbf{x}^{(k)})^t (\mathbf{y} - \mathbf{x}^{(k)}) \geq 0$ implies $f(\mathbf{y}) \geq f(\mathbf{x}^{(k)})$
- Consequently, the search direction in a **descent** method must satisfy

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

i.e., it must make an **acute** angle with the **negative** gradient

Gradient descent

- By **convexity** we have

$$f(\mathbf{y}) \geq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^t (\mathbf{y} - \mathbf{x}^{(k)})$$

- Therefore, $\nabla f(\mathbf{x}^{(k)})^t (\mathbf{y} - \mathbf{x}^{(k)}) \geq 0$ implies $f(\mathbf{y}) \geq f(\mathbf{x}^{(k)})$
- Consequently, the search direction in a **descent** method must satisfy

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

i.e., it must make an **acute** angle with the **negative** gradient

- A natural choice for the **search direction** is the negative gradient

$$\Delta \mathbf{x}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$$

The resulting algorithm is called the **gradient algorithm** or **gradient descent method**

Gradient descent for a quadratic function

- Consider the problem of minimizing

$$f(x_1, x_2) = \frac{1}{2} (\sigma_1 x_1^2 + \sigma_2 x_2^2)$$

over \mathbb{R}^2 , where $\sigma_1 \geq \sigma_2 > 0$ are two positive constants

Gradient descent for a quadratic function

- Consider the problem of minimizing

$$f(x_1, x_2) = \frac{1}{2} (\sigma_1 x_1^2 + \sigma_2 x_2^2)$$

over \mathbb{R}^2 , where $\sigma_1 \geq \sigma_2 > 0$ are two positive constants

- Apparently, $(x_1^*, x_2^*) = (0, 0)$ is the unique minimum point and the optimal value $p^* = 0$

Gradient descent for a quadratic function

- Consider the problem of minimizing

$$f(x_1, x_2) = \frac{1}{2} (\sigma_1 x_1^2 + \sigma_2 x_2^2)$$

over \mathbb{R}^2 , where $\sigma_1 \geq \sigma_2 > 0$ are two positive constants

- Apparently, $(x_1^*, x_2^*) = (0, 0)$ is the unique minimum point and the optimal value $p^* = 0$
- Consider gradient descent with exact line search:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - t \nabla f(\mathbf{x}^{(k)}) = \begin{bmatrix} (1 - \sigma_1 t) x_1^{(k)} \\ (1 - \sigma_2 t) x_2^{(k)} \end{bmatrix}$$

- Note that

$$f(\mathbf{x}^{(k+1)}) = \frac{1}{2} \left[\sigma_1 (1 - \sigma_1 t)^2 (x_1^{(k)})^2 + \sigma_2 (1 - \sigma_2 t)^2 (x_2^{(k)})^2 \right]$$

is a **quadratic** function of t

- Note that

$$f(\mathbf{x}^{(k+1)}) = \frac{1}{2} \left[\sigma_1 (1 - \sigma_1 t)^2 (x_1^{(k)})^2 + \sigma_2 (1 - \sigma_2 t)^2 (x_2^{(k)})^2 \right]$$

is a **quadratic** function of t

- Thus, an exact line search yields:

$$t^{(k)} = \frac{\sigma_1^2 (x_1^{(k)})^2 + \sigma_2^2 (x_2^{(k)})^2}{\sigma_1^3 (x_1^{(k)})^2 + \sigma_2^3 (x_2^{(k)})^2}$$

and hence

$$\mathbf{x}^{(k+1)} = \frac{x_1^{(k)} x_2^{(k)} (\sigma_1 - \sigma_2)}{\sigma_1^3 (x_1^{(k)})^2 + \sigma_2^3 (x_2^{(k)})^2} \begin{bmatrix} -\sigma_2^2 x_2^{(k)} \\ \sigma_1^2 x_1^{(k)} \end{bmatrix}$$

- For $k = 0$, we have

$$\mathbf{x}^{(1)} = \frac{x_1^{(0)} x_2^{(0)} (\sigma_1 - \sigma_2)}{\sigma_1^3 (x_1^{(0)})^2 + \sigma_2^3 (x_2^{(0)})^2} \begin{bmatrix} -\sigma_2^2 x_2^{(0)} \\ \sigma_1^2 x_1^{(0)} \end{bmatrix}$$

- For $k = 0$, we have

$$\mathbf{x}^{(1)} = \frac{x_1^{(0)} x_2^{(0)} (\sigma_1 - \sigma_2)}{\sigma_1^3 (x_1^{(0)})^2 + \sigma_2^3 (x_2^{(0)})^2} \begin{bmatrix} -\sigma_2^2 x_2^{(0)} \\ \sigma_1^2 x_1^{(0)} \end{bmatrix}$$

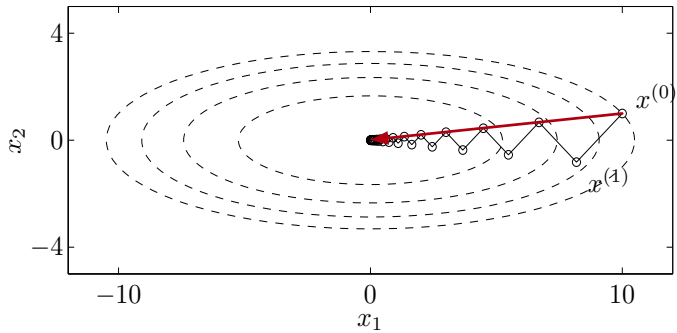
- For $k = 1$, we have

$$\mathbf{x}^{(2)} = \frac{x_1^{(1)} x_2^{(1)} (\sigma_1 - \sigma_2)}{\sigma_1^3 (x_1^{(1)})^2 + \sigma_2^3 (x_2^{(1)})^2} \begin{bmatrix} -\sigma_2^2 x_2^{(1)} \\ \sigma_1^2 x_1^{(1)} \end{bmatrix} = \rho \mathbf{x}^{(0)}$$

where

$$\rho = \frac{\sigma_1 \sigma_2 (\sigma_1 - \sigma_2)^2 (x_1^{(0)})^2 (x_2^{(0)})^2}{\left[\sigma_1 (x_1^{(0)})^2 + \sigma_2 (x_2^{(0)})^2 \right] \left[\sigma_1^3 (x_1^{(0)})^2 + \sigma_2^3 (x_2^{(0)})^2 \right]}$$

i.e., for every **two** iterations we move **directly in the line** toward the minimum point $\mathbf{x}^* = (0, 0)$



- Since f is a **quadratic** function

$$f(\mathbf{x}^{(2)}) - p^* = \rho^2 \left[f(\mathbf{x}^{(0)}) - p^* \right]$$

i.e., the optimality gap forms a **geometric series** with ratio $\approx \rho$ between two consecutive iterations

- Since f is a **quadratic** function

$$f(\mathbf{x}^{(2)}) - p^* = \rho^2 \left[f(\mathbf{x}^{(0)}) - p^* \right]$$

i.e., the optimality gap forms a **geometric series** with ratio $\approx \rho$ between two consecutive iterations

- In the context of iterative numerical methods, this is called **linear convergence** since the optimality gap of an iteration is a **linear function of the optimality gap** of the previous iteration

- Since f is a **quadratic** function

$$f(\mathbf{x}^{(2)}) - p^* = \rho^2 \left[f(\mathbf{x}^{(0)}) - p^* \right]$$

i.e., the optimality gap forms a **geometric series** with ratio $\approx \rho$ between two consecutive iterations

- In the context of iterative numerical methods, this is called **linear convergence** since the optimality gap of an iteration is a **linear function** of the optimality gap of the previous iteration
- Such a convergence behavior is considered to be **slow**, especially when the ratio ρ (known as **convergence rate**) is close to 1

- Since f is a **quadratic** function

$$f(\mathbf{x}^{(2)}) - p^* = \rho^2 \left[f(\mathbf{x}^{(0)}) - p^* \right]$$

i.e., the optimality gap forms a **geometric series** with ratio $\approx \rho$ between two consecutive iterations

- In the context of iterative numerical methods, this is called **linear convergence** since the optimality gap of an iteration is a **linear function** of the optimality gap of the previous iteration
- Such a convergence behavior is considered to be **slow**, especially when the ratio ρ (known as **convergence rate**) is close to 1
- Let $c := \sigma_1/\sigma_2 \geq 1$. The convergence rate

$$\rho = \frac{c(c-1)^2(x_1^{(0)})^2(x_2^{(0)})^2}{\left[c(x_1^{(0)})^2 + (x_2^{(0)})^2 \right] \left[c^3(x_1^{(0)})^2 + (x_2^{(0)})^2 \right]}$$

is apparently sensitive to the initial point $(x_1^{(0)}, x_2^{(0)})$

- The **worse** case occurs when $|x_2^{(0)}/x_1^{(0)}| = c$ and in this case, the convergence rate

$$\rho = \left(\frac{c-1}{c+1} \right)^2$$

which is very close to 1 when $c = \frac{\sigma_1}{\sigma_2} \gg 1$

- The **worse** case occurs when $|x_2^{(0)}/x_1^{(0)}| = c$ and in this case, the convergence rate

$$\rho = \left(\frac{c-1}{c+1} \right)^2$$

which is very close to 1 when $c = \frac{\sigma_1}{\sigma_2} \gg 1$

- Finally, note here that the **Hessian** for the quadratic function $f(x_1, x_2) = \frac{1}{2} (\sigma_1 x_1^2 + \sigma_2 x_2^2)$ is

$$\nabla^2 f = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}$$

so σ_1 and σ_2 are the **largest** and **smallest** eigenvalue of the Hessian respectively, and ρ is simply the condition number of the Hessian

Gradient descent for strongly convex functions

- For a general objective function f , gradient descent with exact line search is guaranteed to converge if f is **strongly convex** on its domain, i.e., there exists an $m > 0$ such that

$$\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}, \quad \forall \mathbf{x} \in \text{dom}(f)$$

Gradient descent for strongly convex functions

- For a general objective function f , gradient descent with exact line search is guaranteed to converge if f is **strongly convex** on its domain, i.e., there exists an $m > 0$ such that

$$\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}, \quad \forall \mathbf{x} \in \text{dom}(f)$$

- Under the strong convexity assumption, it can be shown that

$$f(\mathbf{x}^{(k+1)}) - p^* \leq \rho \left[f(\mathbf{x}^{(k)}) - p^* \right]$$

where

$$\rho = 1 - \frac{1}{\sup_{\mathbf{x} \in \text{dom}(f)} \frac{\lambda_1(\nabla^2 f(\mathbf{x}))}{\lambda_n(\nabla^2 f(\mathbf{x}))}}$$

Gradient descent for strongly convex functions

- For a general objective function f , gradient descent with exact line search is guaranteed to converge if f is **strongly convex** on its domain, i.e., there exists an $m > 0$ such that

$$\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}, \quad \forall \mathbf{x} \in \text{dom}(f)$$

- Under the strong convexity assumption, it can be shown that

$$f(\mathbf{x}^{(k+1)}) - p^* \leq \rho \left[f(\mathbf{x}^{(k)}) - p^* \right]$$

where

$$\rho = 1 - \frac{1}{\sup_{\mathbf{x} \in \text{dom}(f)} \frac{\lambda_1(\nabla^2 f(\mathbf{x}))}{\lambda_n(\nabla^2 f(\mathbf{x}))}}$$

- That is, the gradient method in general exhibits approximately **linear** convergence, i.e., the optimality gap $f(\mathbf{x}^{(k)}) - p^*$ converges to zero approximately as a **geometric** series

- The convergence rate depends greatly on the **condition number** of the Hessians. Convergence can be very slow, even for problems that are moderately well conditioned (say, with condition number in the 100s). When the condition number is larger (say, 1000 or more) the gradient method is so slow that it is useless in practice

Steepest descent

- The first-order Taylor approximation of $f(\mathbf{x} + \mathbf{v})$ around \mathbf{x} is

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^t \mathbf{v}$$

- The second term on the righthand side, $\nabla f(\mathbf{x})^t \mathbf{v}$, is the **directional derivative** of f at \mathbf{x} in the direction \mathbf{v} . It gives the approximate change in f for a small step \mathbf{v} . **The step \mathbf{v} is a descent direction if the directional derivative is negative**

Steepest descent

- The first-order Taylor approximation of $f(\mathbf{x} + \mathbf{v})$ around \mathbf{x} is

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^t \mathbf{v}$$

- The second term on the righthand side, $\nabla f(\mathbf{x})^t \mathbf{v}$, is the **directional derivative** of f at \mathbf{x} in the direction \mathbf{v} . It gives the approximate change in f for a small step \mathbf{v} . The step \mathbf{v} is a descent direction if the directional derivative is negative
- Consider the problem of how to choose \mathbf{v} to make the directional derivative as negative as possible. Since the directional derivative $\nabla f(\mathbf{x})^t \mathbf{v}$ is **linear** in \mathbf{v} , it can be made as negative as we like by taking \mathbf{v} large (provided \mathbf{v} is a descent direction). To make the question sensible we have to limit the **length** of \mathbf{v} or normalize by it

Steepest descent

- The first-order Taylor approximation of $f(\mathbf{x} + \mathbf{v})$ around \mathbf{x} is

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^t \mathbf{v}$$

- The second term on the righthand side, $\nabla f(\mathbf{x})^t \mathbf{v}$, is the **directional derivative** of f at \mathbf{x} in the direction \mathbf{v} . It gives the approximate change in f for a small step \mathbf{v} . The step \mathbf{v} is a descent direction if the directional derivative is negative
- Consider the problem of how to choose \mathbf{v} to make the directional derivative as negative as possible. Since the directional derivative $\nabla f(\mathbf{x})^t \mathbf{v}$ is **linear** in \mathbf{v} , it can be made as negative as we like by taking \mathbf{v} large (provided \mathbf{v} is a descent direction). To make the question sensible we have to limit the **length** of \mathbf{v} or normalize by it
- Let $\|\cdot\|$ be **any** norm on \mathbb{R}^n . We define a **normalized steepest descent direction** (with respect to the norm $\|\cdot\|$) as

$$\Delta \mathbf{x}_{nsd} = \arg \min \{ \nabla f(\mathbf{x})^t \mathbf{v} : \|\mathbf{v}\| = 1 \}$$

- It is also convenient to consider a steepest descent step $\Delta \mathbf{x}_{sd}$ that is **unnormalized**, by scaling the normalized steepest descent direction in a particular way:

$$\Delta \mathbf{x}_{sd} = \|\nabla f(\mathbf{x})\|_* \Delta \mathbf{x}_{nsd}$$

where $\|\cdot\|_*$ denotes the **dual** norm

- It is also convenient to consider a steepest descent step $\Delta \mathbf{x}_{sd}$ that is **unnormalized**, by scaling the normalized steepest descent direction in a particular way:

$$\Delta \mathbf{x}_{sd} = \|\nabla f(\mathbf{x})\|_* \Delta \mathbf{x}_{nsd}$$

where $\|\cdot\|_*$ denotes the **dual** norm

- It can be shown that for the steepest descent step, we have

$$\nabla f(\mathbf{x})^t \Delta \mathbf{x}_{sd} = \|\nabla f(\mathbf{x})\|_* \nabla f(\mathbf{x})^t \Delta \mathbf{x}_{nsd} = -\|\nabla f(\mathbf{x})\|_*^2$$

- It is also convenient to consider a steepest descent step $\Delta \mathbf{x}_{sd}$ that is **unnormalized**, by scaling the normalized steepest descent direction in a particular way:

$$\Delta \mathbf{x}_{sd} = \|\nabla f(\mathbf{x})\|_* \Delta \mathbf{x}_{nsd}$$

where $\|\cdot\|_*$ denotes the **dual** norm

- It can be shown that for the steepest descent step, we have

$$\nabla f(\mathbf{x})^t \Delta \mathbf{x}_{sd} = \|\nabla f(\mathbf{x})\|_* \nabla f(\mathbf{x})^t \Delta \mathbf{x}_{nsd} = -\|\nabla f(\mathbf{x})\|_*^2$$

- The **steepest descent method** uses the **steepest descent direction** as search direction (when exact line search is used, scale factors in the descent direction have no effect, so either normalized or unnormalized direction can be used)

Steepest descent for Euclidean norm

- If we take the norm $\| \cdot \|$ to be the **Euclidean** norm we find that the steepest descent direction is simply the negative gradient, i.e.,

$$\Delta \mathbf{x}_{sd} = -\nabla f(\mathbf{x})$$

Thus, the steepest descent method for the Euclidean norm **coincides** with the gradient descent method

Steepest descent for quadratic norms

- Consider the quadratic norm

$$\|z\|_P = \sqrt{z^t P z} = \|P^{1/2} z\|_2$$

where P is a positive definite matrix

Steepest descent for quadratic norms

- Consider the quadratic norm

$$\|z\|_P = \sqrt{z^t P z} = \|P^{1/2} z\|_2$$

where P is a **positive definite** matrix

- The normalized steepest descent direction is given by

$$\Delta x_{nsd} = - \left(\nabla f(x)^t P^{-1} \nabla f(x) \right)^{-1/2} P^{-1} \nabla f(x)$$

Steepest descent for quadratic norms

- Consider the quadratic norm

$$\|\mathbf{z}\|_{\mathbf{P}} = \sqrt{\mathbf{z}^t \mathbf{P} \mathbf{z}} = \|\mathbf{P}^{1/2} \mathbf{z}\|_2$$

where \mathbf{P} is a **positive definite** matrix

- The normalized steepest descent direction is given by

$$\Delta \mathbf{x}_{nsd} = - \left(\nabla f(\mathbf{x})^t \mathbf{P}^{-1} \nabla f(\mathbf{x}) \right)^{-1/2} \mathbf{P}^{-1} \nabla f(\mathbf{x})$$

- The dual norm is given by $\|\mathbf{z}\|_* = \|\mathbf{P}^{-1/2} \mathbf{z}\|_2$, so the steepest descent step with respect to $\|\cdot\|_{\mathbf{P}}$ is given by

$$\Delta \mathbf{x}_{sd} = -\mathbf{P}^{-1} \nabla f(\mathbf{x})$$

Steepest descent for quadratic norms

- Consider the quadratic norm

$$\|\mathbf{z}\|_{\mathbf{P}} = \sqrt{\mathbf{z}^t \mathbf{P} \mathbf{z}} = \|\mathbf{P}^{1/2} \mathbf{z}\|_2$$

where \mathbf{P} is a **positive definite** matrix

- The normalized steepest descent direction is given by

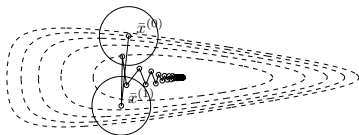
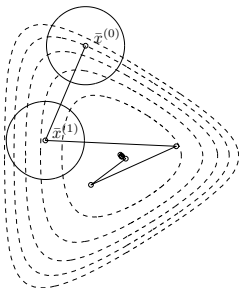
$$\Delta \mathbf{x}_{nsd} = - \left(\nabla f(\mathbf{x})^t \mathbf{P}^{-1} \nabla f(\mathbf{x}) \right)^{-1/2} \mathbf{P}^{-1} \nabla f(\mathbf{x})$$

- The dual norm is given by $\|\mathbf{z}\|_* = \|\mathbf{P}^{-1/2} \mathbf{z}\|_2$, so the steepest descent step with respect to $\|\cdot\|_{\mathbf{P}}$ is given by

$$\Delta \mathbf{x}_{sd} = -\mathbf{P}^{-1} \nabla f(\mathbf{x})$$

- Alternatively, the steepest descent method in the quadratic norm $\|\cdot\|_{\mathbf{P}}$ can be thought of as the **gradient method** applied to the problem after the **change of coordinates**

$$\tilde{\mathbf{x}} = \mathbf{P}^{1/2} \mathbf{x}$$



- This observation provides a prescription for choosing \mathbf{P} : It should be chosen so that the Hessians with respect to $\tilde{\mathbf{x}}$ are **well conditioned**

Steepest descent for ℓ_1 -norm

- Consider the steepest descent method for the ℓ_1 -norm. It can be shown that a normalized steepest descent direction is given by

$$\Delta \mathbf{x}_{nsd} = -\text{sgn} \left(\frac{\partial f(\mathbf{x})}{\partial x_{i^*}} \right) \mathbf{e}_{i^*}$$

where $i^* = \arg \max_i \left| \frac{\partial f(\mathbf{x})}{\partial x_i} \right|$, and \mathbf{e}_i is the i th **standard** basis vector

Steepest descent for ℓ_1 -norm

- Consider the steepest descent method for the ℓ_1 -norm. It can be shown that a normalized steepest descent direction is given by

$$\Delta \mathbf{x}_{nsd} = -\text{sgn} \left(\frac{\partial f(\mathbf{x})}{\partial x_{i^*}} \right) \mathbf{e}_{i^*}$$

where $i^* = \arg \max_i \left| \frac{\partial f(\mathbf{x})}{\partial x_i} \right|$, and \mathbf{e}_i is the i th **standard** basis vector

- An unnormalized steepest descent step is then

$$\Delta \mathbf{x}_{sd} = \|\nabla f(\mathbf{x})\|_{\infty} \Delta \mathbf{x}_{nsd} = -\frac{\partial f(\mathbf{x})}{\partial x_{i^*}} \mathbf{e}_{i^*}$$

Steepest descent for ℓ_1 -norm

- Consider the steepest descent method for the ℓ_1 -norm. It can be shown that a normalized steepest descent direction is given by

$$\Delta \mathbf{x}_{nsd} = -\text{sgn} \left(\frac{\partial f(\mathbf{x})}{\partial x_{i^*}} \right) \mathbf{e}_{i^*}$$

where $i^* = \arg \max_i \left| \frac{\partial f(\mathbf{x})}{\partial x_i} \right|$, and \mathbf{e}_i is the i th **standard** basis vector

- An unnormalized steepest descent step is then

$$\Delta \mathbf{x}_{sd} = \|\nabla f(\mathbf{x})\|_{\infty} \Delta \mathbf{x}_{nsd} = -\frac{\partial f(\mathbf{x})}{\partial x_{i^*}} \mathbf{e}_{i^*}$$

- Thus, the normalized steepest descent step in ℓ_1 -norm can always be chosen to be a standard basis vector (or a negative standard basis vector). It is the coordinate axis direction along which the approximate decrease in f is greatest

Coordinate descent

- The steepest descent algorithm in the ℓ_1 naturally motivates the following **coordinate descent algorithm**: At each iteration we select a component of $\nabla f(\mathbf{x})$ with maximum absolute value, and then decrease or increase the corresponding component of \mathbf{x} , according to the sign of $(\nabla f(\mathbf{x}))_{i^*}$

Coordinate descent

- The steepest descent algorithm in the ℓ_1 naturally motivates the following **coordinate descent algorithm**: At each iteration we select a component of $\nabla f(\mathbf{x})$ with **maximum absolute value**, and then decrease or increase the corresponding component of \mathbf{x} , according to the sign of $(\nabla f(\mathbf{x}))_{i^*}$
- Since only one component of the variable \mathbf{x} is updated at each iteration, this can greatly simplify, or even trivialize, the line search

Coordinate descent

- The steepest descent algorithm in the ℓ_1 naturally motivates the following **coordinate descent algorithm**: At each iteration we select a component of $\nabla f(\mathbf{x})$ with **maximum absolute value**, and then decrease or increase the corresponding component of \mathbf{x} , according to the sign of $(\nabla f(\mathbf{x}))_{i^*}$
- Since only one component of the variable \mathbf{x} is updated at each iteration, this can greatly simplify, or even trivialize, the line search
- Consider, for example, the problem:

$$\text{minimize } f(\mathbf{x}) = \ln \left(\sum_{i=1}^n \sum_{j=1}^n M_{i,j}^2 e^{x_i - x_j} \right)$$

- It is easy to minimize f one component at a time. Keeping all components except the k th fixed, we can write

$$f(\mathbf{x}) = \ln(\alpha_k + \beta_k e^{-x_k} + \gamma_k e^{x_k})$$

- It is easy to minimize f one component at a time. Keeping all components except the k th fixed, we can write

$$f(\mathbf{x}) = \ln(\alpha_k + \beta_k e^{-x_k} + \gamma_k e^{x_k})$$

- The minimum of $f(\mathbf{x})$, as a function of x_k , is obtained for

$$x_k = \frac{1}{2} \ln \left(\frac{\beta_k}{\gamma_k} \right)$$

So for this problem an exact line search can be carried out using a simple analytical formula

Convergence analysis

- Similar to gradient descent, it can be shown that for strongly convex functions, steepest descent with exact line search also converges linearly to the minimum point, and the linear convergence rate is again mainly controlled by the condition number of the Hessians

Newton's methods

- The **second-order** Taylor approximation of $f(\mathbf{x} + \mathbf{v})$ around \mathbf{x} is

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^t \mathbf{v} + \frac{1}{2} \mathbf{v}^t \nabla^2 f(\mathbf{x}) \mathbf{v}$$

which is a convex **quadratic** function of \mathbf{v} and minimized at

$$\mathbf{v} = -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x})$$

Newton's methods

- The **second-order** Taylor approximation of $f(\mathbf{x} + \mathbf{v})$ around \mathbf{x} is

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^t \mathbf{v} + \frac{1}{2} \mathbf{v}^t \nabla^2 f(\mathbf{x}) \mathbf{v}$$

which is a convex **quadratic** function of \mathbf{v} and minimized at

$$\mathbf{v} = -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x})$$

- This suggests the following **pure Newton method**:

$$\Delta \mathbf{x}_{nt} = -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}) \quad \text{and} \quad t = 1$$

Newton's methods

- The **second-order** Taylor approximation of $f(\mathbf{x} + \mathbf{v})$ around \mathbf{x} is

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^t \mathbf{v} + \frac{1}{2} \mathbf{v}^t \nabla^2 f(\mathbf{x}) \mathbf{v}$$

which is a convex **quadratic** function of \mathbf{v} and minimized at

$$\mathbf{v} = -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x})$$

- This suggests the following **pure Newton method**:

$$\Delta \mathbf{x}_{nt} = -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}) \quad \text{and} \quad t = 1$$

- It is straightforward to verify that

$$\nabla f(\mathbf{x})^t \Delta \mathbf{x}_{nt} = -\nabla f(\mathbf{x})^t \nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}) < 0$$

unless $\nabla f(\mathbf{x}) = \mathbf{0}$

- Unfortunately, pure Newton method may **not** be a descent method. That is, it is possible that $f(\mathbf{x}^{(k+1)}) > f(\mathbf{x}^{(k)})$, and this may happen if our starting point is **far away** from the minimum point

- Unfortunately, pure Newton method may **not** be a descent method. That is, it is possible that $f(\mathbf{x}^{(k+1)}) > f(\mathbf{x}^{(k)})$, and this may happen if our starting point is **far away** from the minimum point
- On the other hand, if the starting point is **sufficiently close** to the minimum point, it is observed that pure Newton method exhibits **superior** convergence properties: The optimality gap decreases as a **quadratic** function of the optimality gap from the previous iteration

- Unfortunately, pure Newton method may **not** be a descent method. That is, it is possible that $f(\mathbf{x}^{(k+1)}) > f(\mathbf{x}^{(k)})$, and this may happen if our starting point is **far away** from the minimum point
- On the other hand, if the starting point is **sufficiently close** to the minimum point, it is observed that pure Newton method exhibits **superior** convergence properties: The optimality gap decreases as a **quadratic** function of the optimality gap from the previous iteration
- To overcome the starting-point issue, one may consider using **exact line search** as opposed to a constant step size $t = 1$. This leads to the so-called **damped Newton method**

Convergence analysis

- The convergence behavior of damped Newton method can be divided into two stages

Convergence analysis

- The convergence behavior of damped Newton method can be divided into **two** stages
- In the first, the **damped Newton stage**, we are far from the minimum point (as measured by $\|\nabla f(\mathbf{x}^{(k)})\|_2$), but we make **constant progress**, i.e., $f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) - \gamma$ for some positive constant γ , (in the worst case) towards the optimal solution

Convergence analysis

- The convergence behavior of damped Newton method can be divided into **two** stages
- In the first, the **damped Newton stage**, we are far from the minimum point (as measured by $\|\nabla f(\mathbf{x}^{(k)})\|_2$), but we make **constant progress**, i.e., $f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) - \gamma$ for some positive constant γ , (in the worst case) towards the optimal solution
- When $\|\nabla f(\mathbf{x}^{(k)})\|_2$ is sufficiently small, we enter the **quadratically convergent phase**, where things are settled very quickly: At most **six or so** iterations are required to produce a solution of very high accuracy

Convergence analysis

- The convergence behavior of damped Newton method can be divided into **two** stages
- In the first, the **damped Newton stage**, we are far from the minimum point (as measured by $\|\nabla f(\mathbf{x}^{(k)})\|_2$), but we make **constant progress**, i.e., $f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) - \gamma$ for some positive constant γ , (in the worst case) towards the optimal solution
- When $\|\nabla f(\mathbf{x}^{(k)})\|_2$ is sufficiently small, we enter the **quadratically convergent phase**, where things are settled very quickly: At most **six or so** iterations are required to produce a solution of very high accuracy
- Finally, we mention that a very important feature of Newton's methods is that it is **independent of linear (or affine) changes of coordinates**. This is in stark contrast to the gradient (or steepest descent) method, which is strongly affected by changes of coordinates

Backtracking line search

- Most line searches used in practice are **inexact**: The step length is chosen to approximately minimize f along the ray $\{\mathbf{x} + t\Delta\mathbf{x} : t \geq 0\}$, or even to just reduce f “enough”

Backtracking line search

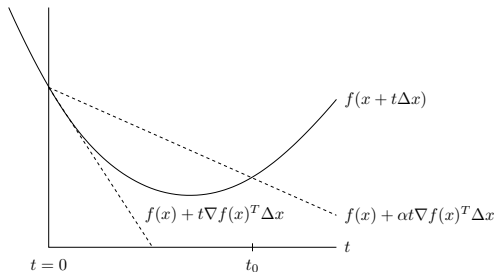
- Most line searches used in practice are **inexact**: The step length is chosen to approximately minimize f along the ray $\{\mathbf{x} + t\Delta\mathbf{x} : t \geq 0\}$, or even to just reduce f “enough”
- One inexact line search method that is very simple and quite effective is called **backtracking line search**. It depends on two constants α and β with $0 < \alpha < 0.5$ and $0 < \beta < 1$

Backtracking line search

- Most line searches used in practice are **inexact**: The step length is chosen to approximately minimize f along the ray $\{\mathbf{x} + t\Delta\mathbf{x} : t \geq 0\}$, or even to just reduce f “enough”
- One inexact line search method that is very simple and quite effective is called **backtracking line search**. It depends on two constants α and β with $0 < \alpha < 0.5$ and $0 < \beta < 1$
- The line search is called backtracking because it starts with a **unit** step size and then reduces it by the factor β until the stopping condition

$$f(\mathbf{x} + t\Delta\mathbf{x}) \leq f(\mathbf{x}) + \alpha t \nabla f(\mathbf{x})^t \Delta\mathbf{x}$$

holds



- Since $\Delta \mathbf{x}$ is a **descent** direction, we have $\nabla f(\mathbf{x})^T \Delta \mathbf{x} < 0$, so for small enough t we have

$$f(\mathbf{x} + t\Delta \mathbf{x}) \approx f(\mathbf{x}) + t\nabla f(\mathbf{x})^T \Delta \mathbf{x} < f(\mathbf{x}) + \alpha t\nabla f(\mathbf{x})^T \Delta \mathbf{x}$$

which shows that the backtracking line search eventually terminates. The constant α can be interpreted as the fraction of the decrease in f predicted by **linear extrapolation** that we will accept