

Master-MIRI

Topics on Optimization and Machine Learning (TOML)

José M. Barceló Ordinas
Departament d'Arquitectura de Computadors
(UPC)

- **Descent Methods for unconstrained minimization**

- **Objective:** find the optimal point $\mathbf{x}^* \in \mathbb{R}^n$, that minimizes an objective function $f_0(\mathbf{x}) \rightarrow$ then the optimal value is $\mathbf{p}^* = \mathbf{f}_0(\mathbf{x}^*)$

The objective is to produce a minimizing sequence $\mathbf{x}^{(k)}$, $k=0, \dots$ such that:

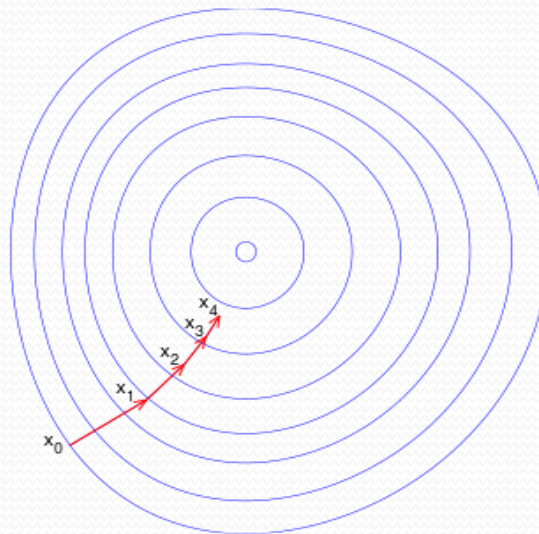
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t \Delta \mathbf{x},$$

Where:

$\Delta \mathbf{x} = \mathbf{d}^{(k)} \in \mathbb{R}^n$ is the **search direction** (vector) or **step**

$t \geq 0$ is the **step size**

k is the **iteration time**



- **Descent Methods for unconstrained minimization**

The objective is to produce a minimizing sequence $x^{(k)}$, $k=0, \dots$ such that:

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

Where:

$\Delta x = d^{(k)} \in \mathbb{R}^n$ is the **search direction** (vector) or **step**

$t \geq 0$ is the **step size**

k is the **iteration time**

- Any descent method should satisfy:

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

except when $x^{(k)}$ is optimal (with certain accuracy: $|x^{(k+1)} - x^{(k)}| \leq \varepsilon$), and then a descent method is characterized by:

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

i.e., must form an acute angle ($< 90^\circ$) with the negative gradient

- The rate of convergence tell us how fast the method approaches the optimal value. It can diverge.

- **Descent Methods for unconstrained minimization**

General Descent Method:

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

- Determine the descent direction $\Delta \mathbf{x} = \mathbf{d}^{(k)}$,
- Line search: Choose the step size t
- Update $k++$, calculate $\mathbf{x}^{(k+1)}$, and return to i)

Gradient Descent Method:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \rho^{(k)} \mathbf{d}^{(k)} = \mathbf{x}^{(k)} - t \nabla f(\mathbf{x}^{(k)})$$

- Choose an initial value $\mathbf{x}^{(0)}$
- The descent direction $\mathbf{d}^{(k)}$ is the **negative gradient**:
$$\Delta \mathbf{x} = \mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$$
- Line search: Choose the step size $t = \rho^{(k)}$
- Update $k++$, calculate $\mathbf{x}^{(k+1)}$, and return to i) if stop criterion is not fulfilled
- Stop Criterion: $\|\nabla f(\mathbf{x}^{(k)})\|_2 \leq \varepsilon$ with $\varepsilon > 0$

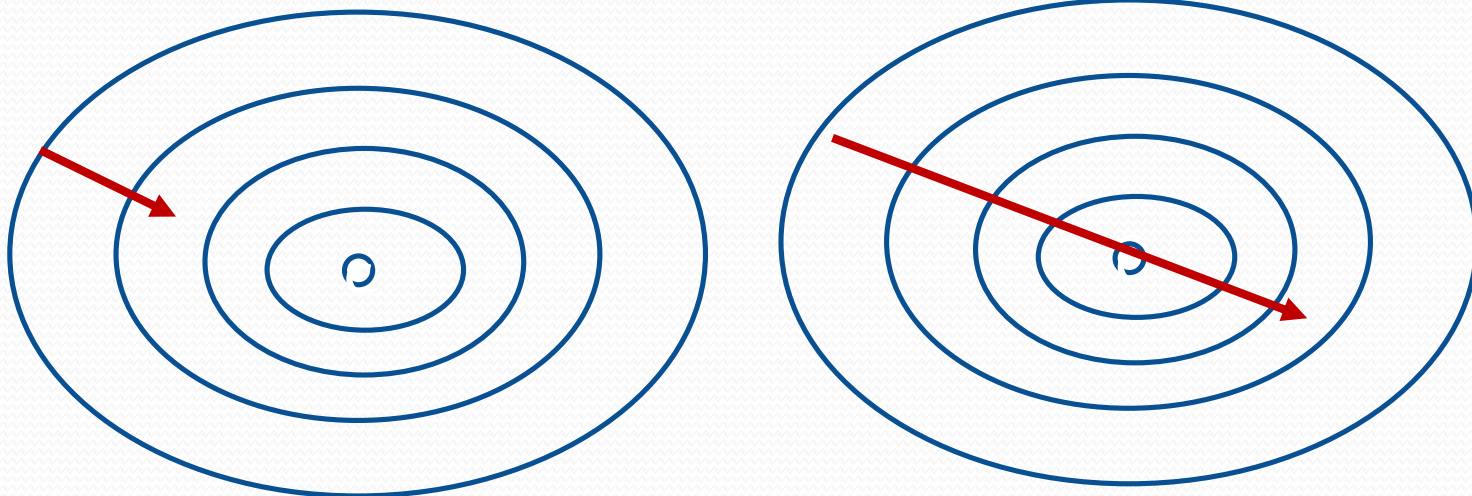
- **Descent Methods for unconstrained minimization**

Selection of the Step Size (Line Search):

- i. **Exact Line Search:** choose $t = \rho^{(k)}$ such that minimizes f along the ray $\{x + t\Delta x \mid t \geq 0\}$

$$t = \rho^{(k)} = \operatorname{argmin}_{s \geq 0} f(x^{(k)} + s \cdot d^{(k)})$$

However, sometimes the cost of this optimization problem is high and approximate methods are used (e.g. Backtracking)



• Descent Methods for unconstrained minimization

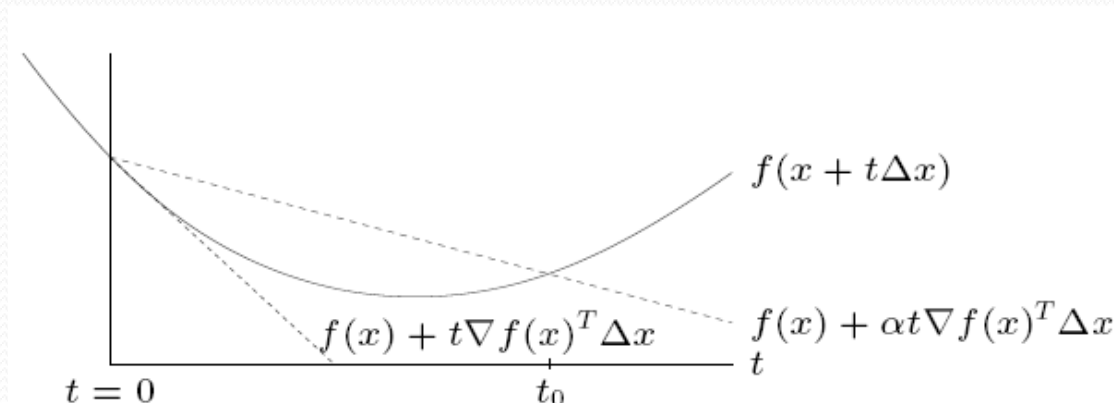
Selection of the Step Size (Line Search):

i. **Backtracking Line Search (Wolfe Condition):** choose $t = \rho^{(k)}$ such that approximately minimizes f along the ray $\{x + t\Delta x \mid t \geq 0\}$.

- given a descent direction $\Delta x = d^{(k)} = -\nabla f(x)$ for f at $x \in \text{dom } f$, $\alpha \in (0, 0.5)$ and $\beta \in (0, 1)$ and $t = 1$

- **while** $f(x^{(k)} + td^{(k)}) > f(x^{(k)}) + \alpha t \nabla f(x^{(k)})^T d^{(k)}$ (stop criterion) $\rightarrow t := \beta t$

- If α is chosen to be between 0.01-0.3, we want to decrease f between 1%-30% of the prediction based on linear interpolation. β is typically chosen 0.1 (very crude search) or 0.8 (less crude search)



• Descent Methods for unconstrained minimization

Steepest Descent Method:

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

- i. Choose an initial value $\mathbf{x}^{(0)}$
- ii. The descent direction $\Delta \mathbf{x}_{sd}$ makes the **directional derivative** $(\nabla f(\mathbf{x}^{(k)})^T \mathbf{v})$ as negative as possible
- ii. **Line search:** Choose the step size t using Exact line search or backtracking line search.
- iii. **Stop Criterion:** quit if $\|\nabla f(\mathbf{x}^{(k)})\|_2 \leq \varepsilon$ with $\varepsilon > 0$
- iv. **Update:** $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t \Delta \mathbf{x}$

$\Delta \mathbf{x}_{sd} = \operatorname{argmin} \{ \nabla f(\mathbf{x}^{(k)})^T \mathbf{v} \mid \|\mathbf{v}\| \leq 1 \}$ is the normalized steepest descent direction with respect norm $\|\cdot\|$.

- If the norm is the Euclidean norm $\rightarrow \Delta \mathbf{x}_{sd} = -\nabla f(\mathbf{x})$ and the steepest descent method is the gradient method.
- Other norms result in other steepest descent methods, e.g. $\|\cdot\|_1$ produces the **coordinate descent method**

• Descent Methods for unconstrained minimization

- The **condition number** measures how much a function f changes in proportion to small changes of the argument x . A problem with small condition number values is said **well-conditioned**, while a problem with large condition numbers is said to be **ill-conditioned**.

$K = x^T f'(x) / f(x)$ (if 1-dim), $K = \|x\| \|J\| / \|f\|$ if more than 1-dim

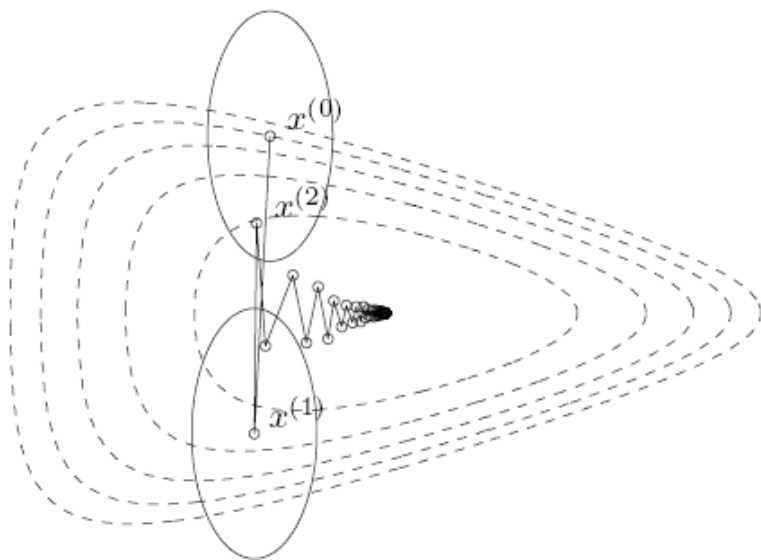


Figure 9.12 Steepest descent method, with quadratic norm $\|\cdot\|_{P_2}$.

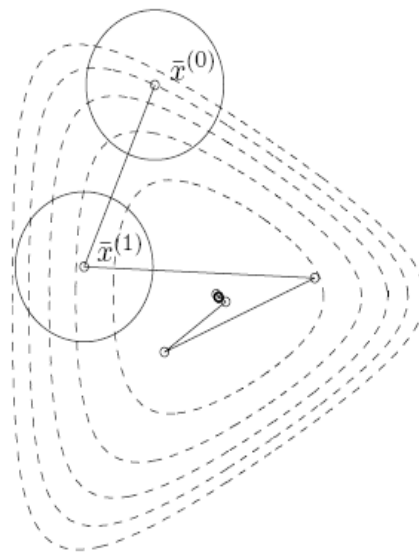


Figure 9.14 The iterates of steepest descent with norm $\|\cdot\|_{P_1}$, after the change of coordinates. This change of coordinates reduces the condition number of the sublevel sets, and so speeds up convergence.

- **Descent Methods for unconstrained minimization**

Newton's Method:

- i. Choose an initial value $x^{(0)}$
- ii. **Define the Newton decrement:** $\lambda^2 = \nabla f(x^{(k)})^T \nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$
- iii. **Define the Newton step:** $\Delta x = - \nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$
- iv. **Stopping Criterion:** quit if $\lambda^2/2 \leq \varepsilon$ with $\varepsilon > 0$
- ii. **Line search:** choose step size t using backtracking line search
- iii. **Update:** $x^{(k+1)} = x^{(k)} + t\Delta x$

Good estimate when x is near x^* , since the Newton step is a minimizer of second-order approximation:

$$f(x+v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$$

Before using Newton's method, the rate of convergence should be checked (it is quadratic) \rightarrow proof of quadratic convergence

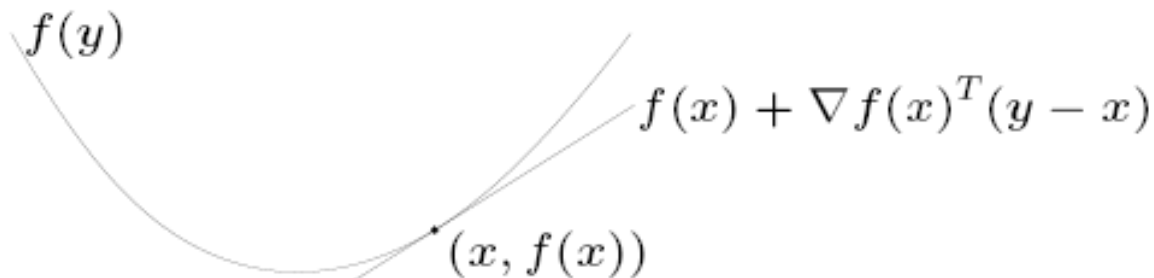
- **Subgradient methods**

We have always assumed that $f(x)$ is differentiable, that is that $\nabla f(x)$ exist for all $x \in X$. What happens if $f(x)$ is not differentiable ?

For example, $f(x) = |x|$, $x \in \mathbb{R}$ is not differentiable at $x=0$, however, the function is convex and the minimum is at $x=0$, How do we apply the Gradient Descent method if we can not obtain the gradient ?

Remember that the first order condition states for convex differentiable functions that:

$f(y) \geq f(x) + \nabla f(x)^T (y - x)$ (first order Taylor approximation of f near x) for all $x, y \in \text{dom } f$ (remember topic 12)



first-order approximation of f is global underestimator

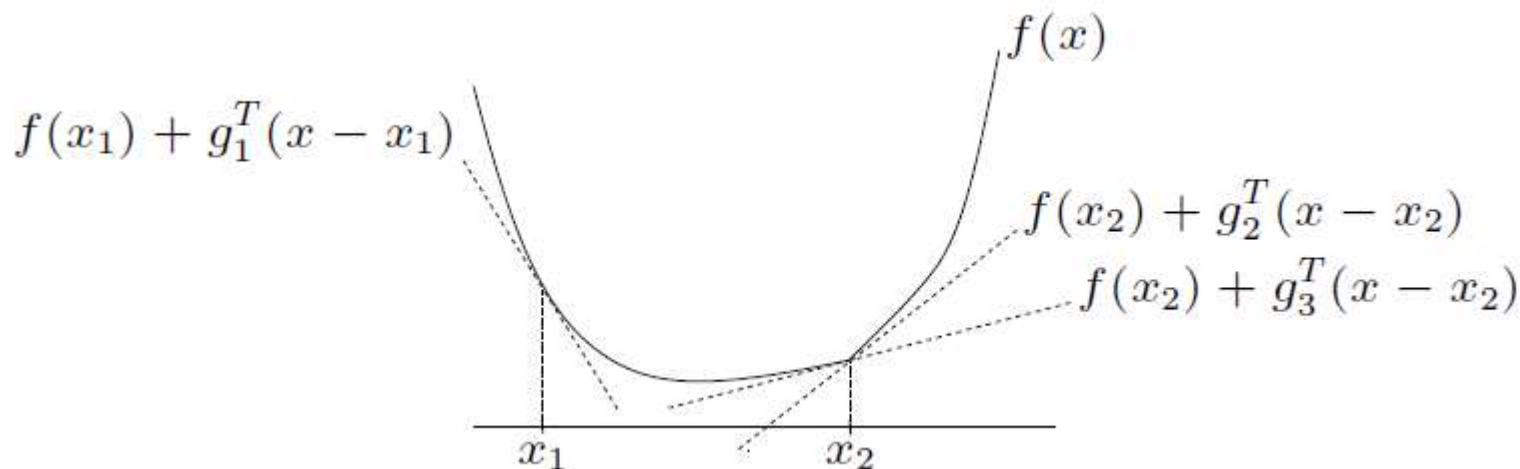
• Subgradient methods

The function g is a subgradient of f if:

$$f(y) \geq f(x) + g(x)^T(y - x) \text{ for all } y \in \text{dom } f$$

e.g. g_1 , g_2 and g_3 are subgradients of $f(x)$

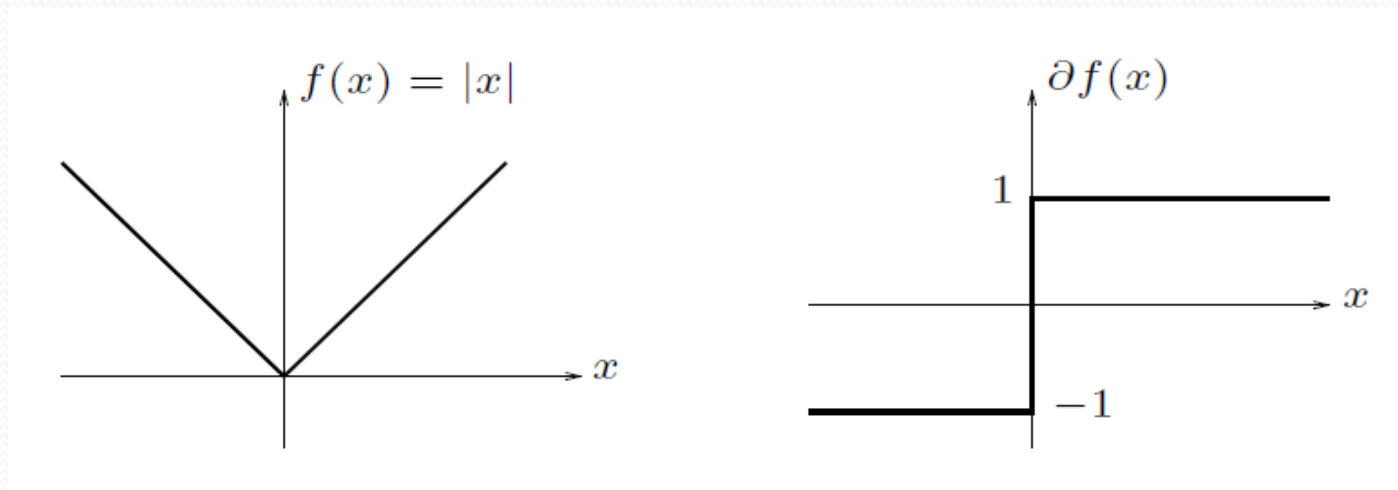
- if $f(x)$ is convex, it has at least one subgradient at every point in the relint of the domain
- if $f(x)$ is convex and differentiable then $\nabla f(x)$ is a subgradient of f at x .



- Subgradient methods

The set of all subgradients of $f(x)$ is called the subdifferential of f at x and is denoted as $\partial f(x)$.

e.g. $f(x) = |x| \quad \rightarrow \quad \begin{aligned} \partial f(x) &= -1 \text{ if } x < 0 \\ &= 1 \text{ if } x > 0 \end{aligned}$



There are a method called “subgradient method” to solve numerically

- Subgradient methods

There are descend methods called “subgradient method” to solve iteratively optimization problems in which the function is not differentiable but in which it is possible to find a subgradient.

In this case:

The objective is to produce a minimizing sequence $\mathbf{x}^{(k)}$, $k=1, \dots$ such that:
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - t \mathbf{g}^{(k)}(\mathbf{x}^{(k)}),$$

Where:

$\mathbf{g}^{(k)}(\mathbf{x}^{(k)})$ is a **subgradient** of \mathbf{f} at $\mathbf{x}^{(k)}$

$t \geq 0$ is the **step size**

k is the **iteration time**

If $-\mathbf{g}^{(k)}$ is not a descent direction of \mathbf{f} at $\mathbf{x}^{(k)}$, we maintain an \mathbf{f}_{best} that keeps track of the lowest objective function value found so far:

$$\mathbf{f}_{\text{best}}^{(k)} = \{ \mathbf{f}_{\text{best}}^{(k-1)}, \mathbf{f}(\mathbf{x}^{(k)}) \}$$

- **Descent Methods with equality constraints**

Let us assume a minimization problem with equality constraints:

$$\begin{array}{ll}\text{minimize} & f(x) \\ \text{subject to} & Ax = b\end{array}$$

and substitute the objective function by its Taylor second-order approximation near x :

$$\begin{array}{ll}\text{minimize} & f(x+v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v \\ \text{subject to} & A(x+v) = b \rightarrow Ax + Av = b \rightarrow Av = 0 \\ \text{var} & x, v\end{array}$$

quadratic COP that can be solved analytically.

Define Δx_{nt} , the **Newton Step at x** , as the solution of the former COP, it is to say the increment to x to solve the problem when the quadratic approximation is used in place of f .

- **Descent Methods with equality constraints**

The KKT conditions (remember last slide of topic 14) for quadratic problems with equality constraint is:

$$\begin{array}{ll} \text{minimize} & (1/2) \mathbf{x}^T \mathbf{P} \mathbf{x} + \mathbf{q}^T \mathbf{x} + r \rightarrow \mathbf{f}(\mathbf{x}) + \nabla \mathbf{f}(\mathbf{x})^T \mathbf{w} + (1/2) \mathbf{w}^T \nabla^2 \mathbf{f}(\mathbf{x}) \mathbf{w} \\ \text{subject to} & \mathbf{A} \mathbf{x} = \mathbf{b} \rightarrow \mathbf{A} \mathbf{x} + \mathbf{A} \mathbf{w} = \mathbf{b} \end{array}$$

where $\Delta \mathbf{x}_{\text{nt}}$ is the step and \mathbf{w} is the associated optimal variable of the dual problem

$$\begin{bmatrix} \mathbf{P} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}^* \\ \mathbf{v}^* \end{bmatrix} = \begin{bmatrix} -\mathbf{q} \\ \mathbf{b} \end{bmatrix} \rightarrow \begin{bmatrix} \nabla^2 \mathbf{f}(\mathbf{x}) & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x}_{\text{nt}} \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} -\nabla \mathbf{f}(\mathbf{x}) \\ \mathbf{0} \end{bmatrix}$$

Newton's Method with equality constraints:

- Choose an initial value $\mathbf{x}^{(0)}$, such that $\mathbf{x}^{(0)} \in \text{dom } \mathbf{f}$ with $\mathbf{A} \mathbf{x} = \mathbf{b}$ and choose tolerance $\varepsilon > 0$
- Compute the Newton $\Delta \mathbf{x}_{\text{nt}}$ step and decrement**
$$\lambda^2 = \nabla \mathbf{f}(\mathbf{x}^{(k)})^T \nabla^2 \mathbf{f}(\mathbf{x}^{(k)})^{-1} \nabla \mathbf{f}(\mathbf{x}^{(k)})$$
- Stopping Criterion:** quit if $\lambda^2/2 \leq \varepsilon$ with $\varepsilon > 0$
- Line search:** choose step size t using backtracking line search
- Update:** $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t \Delta \mathbf{x}_{\text{nt}}$

- Interior-Point Methods for COP

$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq 0 \quad i=1,\dots,m \\ & Ax = b\end{array}$$

That satisfies KKT conditions:

- i. **Primal constraints:** $f_i(x^*) \leq 0, \quad i=1,\dots,m$
- ii. **Primal constraints:** $Ax^* = b,$
- iii. **Dual constraints:** $\lambda_i^* \geq 0 \quad i=1,\dots,m$
- iv. **Complementary slackness:** $\lambda_i^* f_i(x^*) = 0 \quad i=1,\dots,m$
- v. **Gradient of Lagrangian vanishes:**

$$\nabla_x L(x, \lambda^*, v^*) = \nabla f_0(x^*) + \sum_{i=1,\dots,m} \lambda_i^* \nabla f_i(x^*) + A^T v^* = 0$$

Interior-point methods solve the COP problem (or KKT problem) by applying Newton's method to a sequence of equality constraint problems or sequence of modified versions of the KKT conditions.

- Interior-Point Methods – Logarithmic barrier

Objective: formulate the inequality constrained COP as an equality constrained problem to which Newton's method can be applied.

$$\begin{array}{ll}\text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0 \quad i=1,\dots,m \\ & \mathbf{Ax} = \mathbf{b}\end{array}$$

Can be re-written as:

$$\begin{array}{ll}\text{minimize} & f_0(\mathbf{x}) + \sum_{i=1,\dots,m} I_{-}(f_i(\mathbf{x})) \\ \text{subject to} & \mathbf{Ax} = \mathbf{b}\end{array}$$

Where:

$$I_{-}(u) = \begin{cases} 0 & u \leq 0 \\ \infty & u > 0 \end{cases}$$

Whose objective function is not differentiable (Newton's method can not be applied)

- Interior-Point Methods – Logarithmic barrier

Approximate the indicator function I_{-} as:

$$I_{-}(u) = -(1/t) \log(-u) \quad \text{with } \text{dom } I_{-} = -\mathbb{R}_{++}$$

where $t > 0$ sets the accuracy of the approximation and we note that I_{-} is convex and non-decreasing and takes value ∞ if $u > 0$.

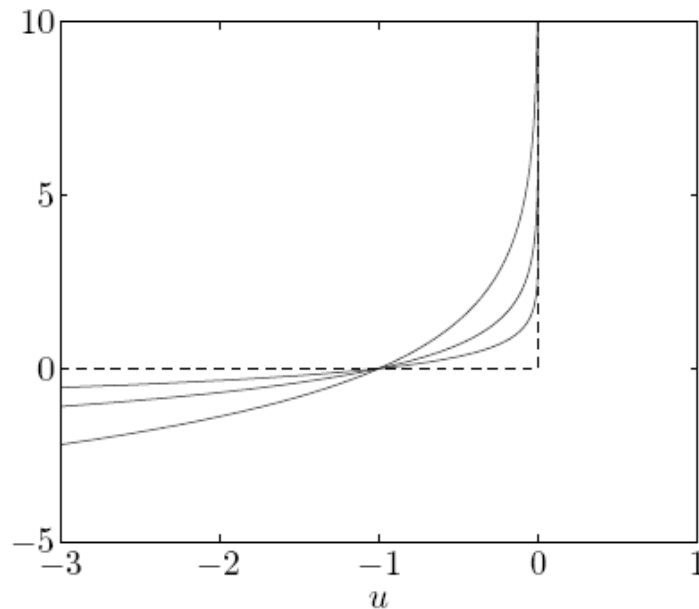


Figure 11.1 The dashed lines show the function $I_{-}(u)$, and the solid curves show $\hat{I}_{-}(u) = -(1/t) \log(-u)$, for $t = 0.5, 1, 2$. The curve for $t = 2$ gives the best approximation.

- Interior-Point Methods – Logarithmic barrier

Since:

$$l_{-}(u) = -(1/t) \log(-u) \quad \text{with } \text{dom } l_{-} = -\mathbb{R}_{++}$$

Then, now:

$$\begin{array}{ll} \text{minimize} & f_0(x) + \sum_{i=1, \dots, m} l_{-}(f_i(x)) \\ \text{subject to} & Ax = b \end{array}$$

can be re-written as:

$$\begin{array}{ll} \text{minimize} & f_0(x) + \sum_{i=1, \dots, m} -(1/t) \log(-f_i(x)) \\ \text{subject to} & Ax = b \end{array}$$

and Newton's method can be applied since the objective function is differentiable

• Interior-Point Methods – Logarithmic barrier

The function $\phi(\mathbf{x}) = -\sum_{i=1,\dots,m} \log(-f_i(\mathbf{x}))$ is called **Log barrier**.

- As t grows, the approximation improves
- As t grows, $f_0 + (1/t) \phi(\mathbf{x})$ is difficult to minimize using Newton's method

The gradient and Hessian of $\phi(\mathbf{x})$ are:

$$\nabla \phi(\mathbf{x}) = \sum_{i=1}^m \frac{-1}{f_i(\mathbf{x})} \nabla f_i(\mathbf{x})$$

$$\nabla^2 \phi(\mathbf{x}) = \sum_{i=1}^m \frac{1}{f_i(\mathbf{x})^2} \nabla f_i(\mathbf{x}) \nabla f_i(\mathbf{x})^T + \sum_{i=1}^m \frac{-1}{f_i(\mathbf{x})} \nabla^2 f_i(\mathbf{x})$$

• Interior-Point Methods – Central Path

The **central path** is defined as the set of **central points** $x^*(t)$, $t > 0$, that solve the minimization problem:

$$\begin{array}{ll} \text{minimize} & f_0(x) + (1/t) \phi(x) \\ \text{subject to} & Ax = b \end{array}$$

The **central point** $x^*(t)$ is characterized by being strictly feasible, it is to say:

$$Ax^*(t) = b$$

$$f_i(x^*(t)) < 0$$

The **Lagrangian** is $L(x, v) = f_0(x) + 1/t \phi(x) + v^T(Ax - b)$, and

The **Lagrange function** is $q(x, v) = \min_{x \in X} L(x, v) = \min_{x \in X} \{ f_0(x) + 1/t \phi(x) + v^T(Ax - b) \}$

We say that the **centrality condition** holds if (calculate the gradient of the Lagrangian with respect x):

$$\begin{aligned} 0 &= \nabla f_0(x^*(t)) + 1/t \nabla \phi(x^*(t)) + A^T v \\ &= \nabla f_0(x^*(t)) + \sum_{i=1, \dots, m} (-1)/(t f_i(x^*(t))) \nabla f_i(x^*(t)) + A^T v = \\ &= \nabla f_0(x^*(t)) + \sum_{i=1, \dots, m} \lambda_i \nabla f_i(x^*(t)) + A^T v \end{aligned}$$

• Interior-Point Methods – Central Path

Observe that we can interpret the former equation as coming from a Lagrangian such as: $L(x, \lambda, v) = f_0(x) + \sum_{i=1, \dots, m} \lambda_i f_i(x) + v^T(Ax - b)$ where $\lambda^*(t) = -1/(t f_i(x^*(t)))$

Then, $(\lambda^*(t), v^*(t))$ are dual feasible and the dual function $q(\lambda^*(t), v^*(t))$ is:

$$\begin{aligned} q(\lambda^*(t), v^*(t)) &= f_0(x^*(t)) + \sum_{i=1, \dots, m} \lambda_i^*(t) f_i(x^*(t)) + v^{*T}(t)(Ax^*(t) - b) \\ &= f_0(x^*(t)) + \sum_{i=1, \dots, m} 1/t = f_0(x^*(t)) - m/t \end{aligned}$$

where m is the number of inequalities, $i=1, \dots, m$

Finally, the duality gap tells us that:

$$q(\lambda^*(t), v^*(t)) = d^* \leq p^* \rightarrow f_0(x^*(t)) - m/t - p^* \leq 0 \rightarrow f_0(x^*(t)) - p^* \leq m/t$$

which tells us that $x^*(t)$ converges to p^* as $t \rightarrow \infty$

• Interior-Point Methods – Central Path

KKT interpretation: We can also interpret the central path conditions as a continuous deformation of the KKT optimality conditions: a point x is equal to $x^*(t)$ if and only if there exists λ, v such that:

- i. **Primal constraints:** $f_i(x) \leq 0, \quad i=1, \dots, m$
- ii. **Primal constraints:** $Ax = b$
- iii. **Dual constraints:** $\lambda_i \geq 0 \quad i=1, \dots, m$
- iv. **Complementary slackness:** $-\lambda_i f_i(x) = 1/t \quad i=1, \dots, m$
- v. **Gradient of Lagrangian vanishes:**

$$\nabla_x L(x, \lambda, v) = \nabla f_0(x) + \sum_{i=1, \dots, m} \lambda_i \nabla f_i(x) + A^T v = 0$$

Note that the only difference is in the slackness condition in which $-\lambda_i f_i(x) = 1/t$ instead of $\lambda_i f_i(x) = 0$.

In fact as $t \rightarrow \infty$, $\lambda_i f_i(x) \rightarrow 0$, for all $i=1, \dots, m$ and $\lambda(t)$ and $v(t)$ **almost satisfy** the KKT conditions.

- Interior-Point Methods – The Barrier Method

The Barrier Method SUMT (Sequential Unconstrained Minimization Technique)

Given strictly feasible x , $t=t^{(0)}$, $\mu>1$, tolerance $\varepsilon>0$

i. Centering step or outer iteration:

Compute $x^*(t)$ by minimizing $f_0 + 1/t \phi$, subject to $Ax=b$, starting at x .

ii. Update: $x=x(t)$

iii. Stopping criterion: quit if $m/t < \varepsilon$ (m/t is the duality gap)

iv. Increase t as $t=\mu t$

At each step, we compute the central point $x^*(t)$ starting from the previous computed central point. The algorithm also computes $\lambda^*(t)$ and $v^*(t)$

We refer to the Newton iterations or steps executed during the centering step as **inner iterations**. At each inner step, we have a primal feasible point; but we have a dual feasible point only at the end of each outer (centering) step.