



# Mathematic Foundations for Computer Science

Probability and Optimization

## Chapter 4: Convex Optimization

Spring 2021

Instructor: Xiaodong Gu





# Recall: What is Optimization?

## Mathematical Optimization

The task of selecting the “**best**” configuration of a set of variables from a “**feasible**” set of configurations.

$$\begin{array}{ll} \text{minimize (or maximize)} & f(x) \\ \text{subject to} & x \in \chi \end{array}$$



## Continuous Optimization

Optimization problems where feasible set  $\chi$  is a connected subset of Euclidean space, and  $f$  is a **continuous** function.

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & g_i(x) \leq b_i, \text{ for } i \in C. \end{array}$$



# Standard Form

---

Finding the minimizer of a function subject to constraints:

$$\begin{aligned} & \underset{x}{\text{minimize}} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0 \quad i = 1, \dots, m \\ & && h_i(x) = 0 \quad i = 1, \dots, p \end{aligned}$$

where

$x = (x_1, \dots, x_n)$  is the optimization variable

$f_0 : \mathbf{R}^n \rightarrow \mathbf{R}$  is the objective function

$f_i : \mathbf{R}^n \rightarrow \mathbf{R}, \quad i = 1, \dots, m$  are inequality constraint functions

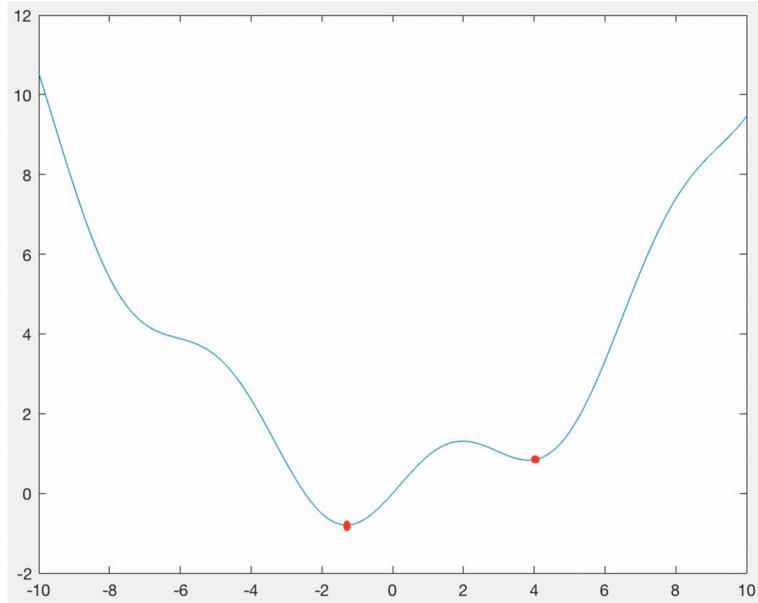
$h_i : \mathbf{R}^n \rightarrow \mathbf{R}, \quad i = 1, \dots, p$  are equality constraint functions.

- **Goal:** find an optimal solution  $x^*$  that minimizes  $f_0$  while satisfying all the constraints.



# Local Minima and Global Minima

- local minima: a solution that is optimal within a neighboring set.
- global minima: the optimal solution among all possible solutions



What kind of functions have  
local minima == global minima?





# Convex Optimization

## Convex Optimization

A problem of minimizing a **convex function** (or maximizing a concave function) over a **convex set**.

Standard Form:

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && g_i(x) \leq 0, \text{ for } i \in C_1. \\ & && h_i(x) = 0, \text{ for } i \in C_2. \end{aligned}$$

where  $f_0, g_i, h_i$  are convex

Convex optimization problems have local minima == global minima



# Overview of Chapter 3

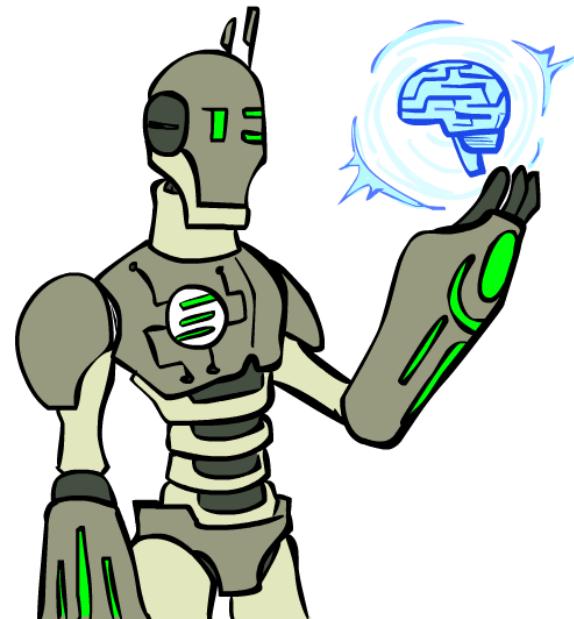
---

## Convex Problems

Convex set, convex function, quadratic programming, Lagrange duality

## Optimization Algorithms

Gradient descend, Newton's method, interior point method,

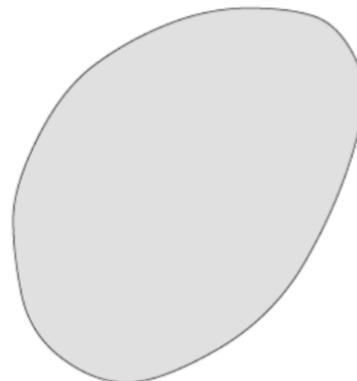




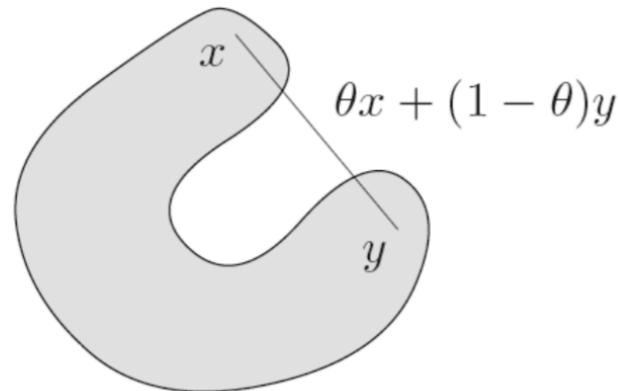
# Convex Set

- A set  $C \in \mathbf{R}^n$  is said to be **convex** if the line segment between any two points is in the set: for any  $x, y \in C$  and  $0 \leq \theta \leq 1$ ,

$$\theta x + (1 - \theta) y \in C.$$



convex



non-convex



# Examples of Convex Set

- Trivial:  $\emptyset$ , point, line, etc.

- Hyperplane:

$$C = \{x \mid a^T x = b\} \text{ where } a \in R^n, b \in R$$

- Halfplane:

$$C = \{x \mid a^T x \leq b\}$$

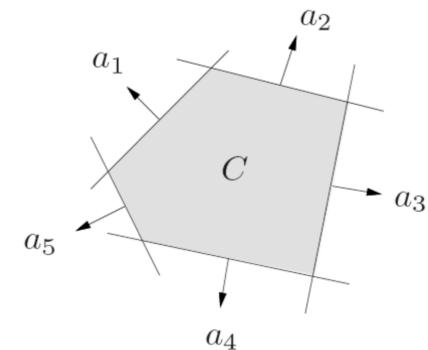
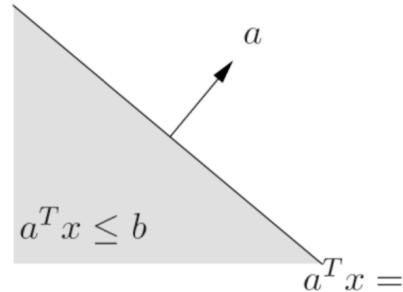
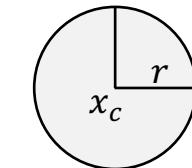
- Polyhedron:

$$C = \{x \mid Ax \leq b, Cx = d\}$$

where  $A \in R^{m \times n}, b \in R^m, C \in R^{p \times n}, d \in R^p$

- Euclidean ball:

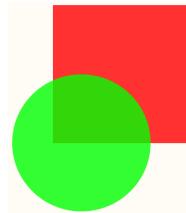
$$B(x_c, r) = \{x \mid \|x - x_c\|_2 \leq r\}$$





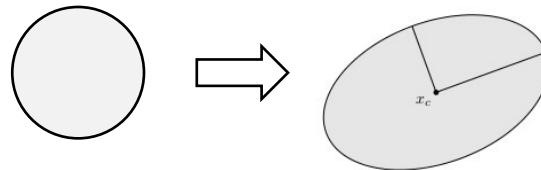
# Operations preserving convexity

- **Intersection:** the intersection of convex sets is convex.



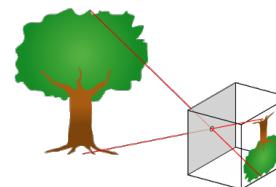
- **Affine Maps:** if  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is an affine function (i.e.,  $f(x) = Ax + b$ ) and  $C$  is convex, then  $f(C) = \{Ax + b : x \in C\}$  is convex.

**Example:** an ellipsoid is image of a unit ball after an affine map



- **Perspective Function:** Let  $P: \mathbb{R}^{n+1} \rightarrow \mathbb{R}^n$  be  $P(x, t) = x/t$ , and  $S \in \mathbb{R}^{n+1}$  is convex, then  $P(S)$  is convex.

The perspective function scales or normalizes vectors so the last component is one, and then drops the last component.



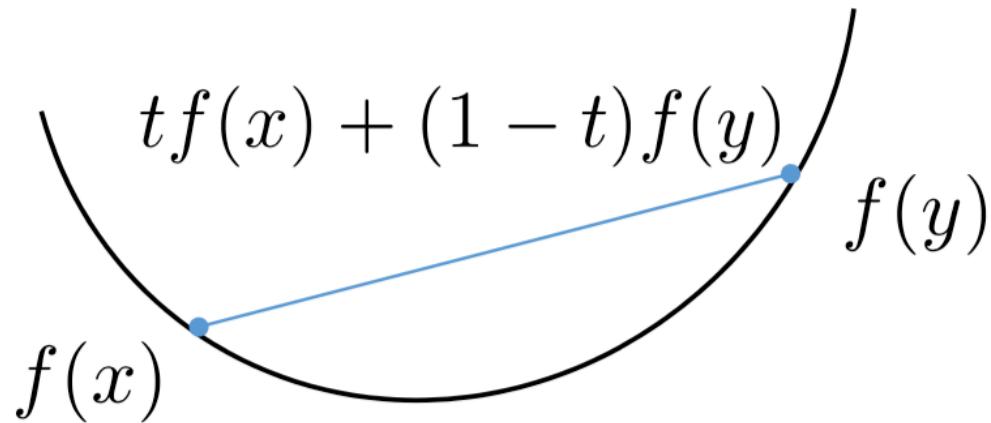


# Convex Functions

A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is **convex** if for  $x, y \in \text{dom } f$ ,

$$f(tx + (1-t)y) \leq t f(x) + (1-t) f(y),$$

for all  $t \in [0, 1]$ .





# Example of Convex Functions

---

- Exponential function:  $e^{ax}$
- logarithmic function  $\log(x)$  is concave
- Affine function:  $a^\top x + b$
- Quadratic function:  $x^\top Qx + b^\top x + c$  is convex if  $Q$  is positive semidefinite (PSD)
- Least squares loss:  $\|y - Ax\|_2^2$
- Norm:  $\|x\|$  is convex for any norm

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \quad \|x\|_1 = \sum_{i=1}^n |x_i|$$

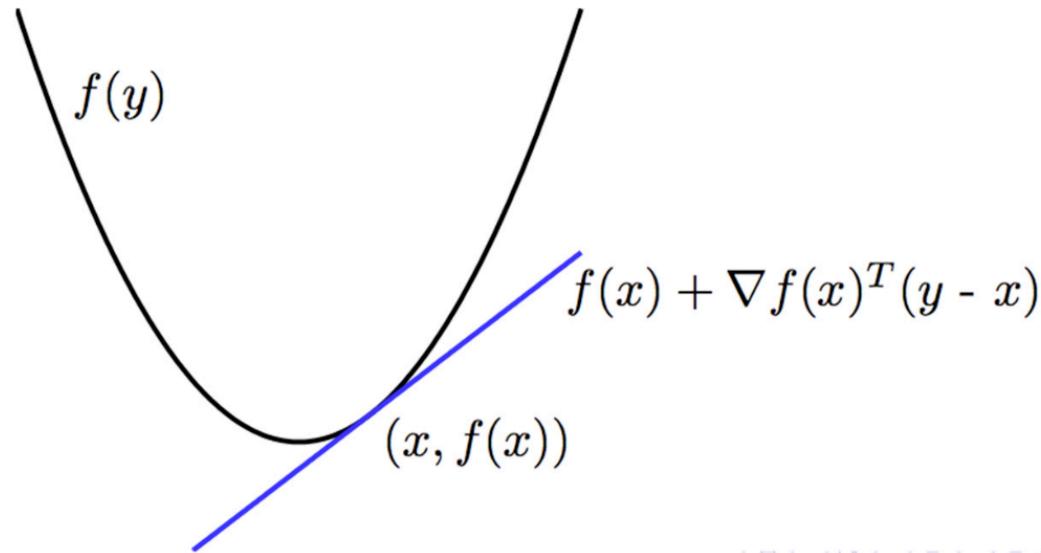


# First Order Convexity Conditions

## Theorem

Suppose  $f$  is differentiable. Then  $f$  is convex if and only if for all  $x, y \in \text{dom } f$

$$f(y) \geq f(x) + \nabla f(x)^T(y - x)$$





# Second Order Convexity Conditions

Suppose  $f$  is twice differentiable. Then  $f$  is **convex** if and only if for all  $x \in \text{dom } f$

$$\nabla^2 f(x) \succeq 0$$



# Properties of convex functions

- If  $x$  is a **local minimizer** of a convex function, it is a **global minimizer**.
- Suppose  $f$  is differentiable and convex. Then,  $x$  is a **global minimizer** of  $f(x)$  if and only if  $\nabla f(x) = 0$ .

- Proof:

- $\nabla f(x) = 0$ . We have

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) = f(x)$$

- $\nabla f(x) \neq 0$ . There is a direction of descent.



---

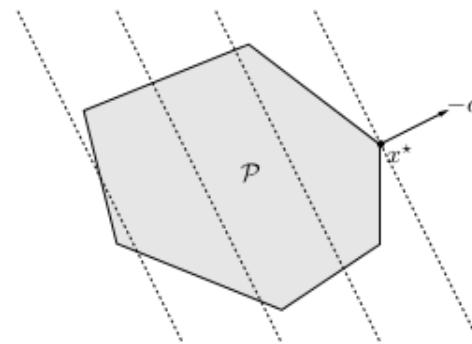
# Convex Optimization Problems



# Linear Programming

We have already seen linear programming

$$\begin{aligned} & \text{minimize} && c^T x \\ & \text{subject to} && Ax \leq b \\ & && x \geq 0 \end{aligned}$$





# Geometric Programming

## Definition

- A **monomial** (单项式) is a function  $f : \mathbb{R}_+^n \rightarrow \mathbb{R}_+$  of the form
$$f(x) = cx_1^{a_1}x_2^{a_2} \dots x_n^{a_n}$$
where  $c \geq 0$ ,  $a_i \in \mathbb{R}$ .
- A **posynomial** (正多项式) is a sum of monomials.

A **Geometric Program** is an optimization problem of the following form

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq b_i, \quad \text{for } i \in \mathcal{C}_1. \\ & && h_i(x) = b_i, \quad \text{for } i \in \mathcal{C}_2. \\ & && x \succeq 0 \end{aligned}$$

where  $f_i$ 's are posynomials,  $h_i$ 's are monomials, and  $b_i > 0$  (wlog 1).

## Interpretation

- GP model volume/area minimization problems, subject to constraints.



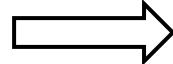
# Example: GP

A manufacturer is designing a suitcase (手提箱)

- Variables:  $h, w, d$
- Want to minimize surface area:  $2(hw + hd + wd)$  (i.e. amount of material used)
- Have a target volume:  $hwd \geq 5$
- Practical/aesthetic (美学) constraints limit aspect ratio:  
$$h/w \leq 2, h/d \leq 3$$
- Constrained by airline to  $h + w + d \leq 7$

$$\begin{array}{ll} \text{minimize} & 2hw + 2hd + 2wd \\ \text{subject to} & h^{-1}w^{-1}d^{-1} \leq \frac{1}{5} \\ & hw^{-1} \leq 2 \\ & hd^{-1} \leq 3 \\ & h + w + d \leq 7 \\ & h, w, d \geq 0 \end{array}$$

$$\begin{array}{l} \tilde{h} = \log h \\ \tilde{w} = \log w \\ \tilde{d} = \log d \end{array}$$



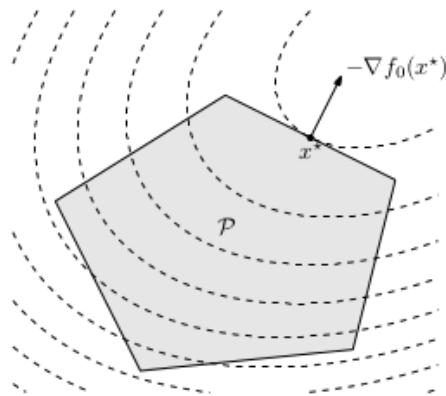
$$\begin{array}{ll} \text{minimize} & 2e^{\tilde{h}+\tilde{w}} + 2e^{\tilde{h}+\tilde{d}} + 2e^{\tilde{w}+\tilde{d}} \\ \text{subject to} & e^{-\tilde{h}-\tilde{w}-\tilde{d}} \leq \frac{1}{5} \\ & e^{\tilde{h}-\tilde{w}} \leq 2 \\ & e^{\tilde{h}-\tilde{d}} \leq 3 \\ & e^{\tilde{h}} + e^{\tilde{w}} + e^{\tilde{d}} \leq 7 \end{array}$$



# Quadratic Programming (二次规划)

Minimizing **convex** quadratic function (二次函数) over a polyhedron. Require  $P \geq 0$ .

$$\begin{aligned} & \text{minimize} && x^\top Px + c^\top x + d \\ & \text{subject to} && Ax \leq b \end{aligned}$$



$P \geq 0$  : positive semi-definite,  $P \in S^n_+$   
symmetric & all eigenvalues are nonnegative,  
 $x^\top Px \geq 0$  for all  $x$

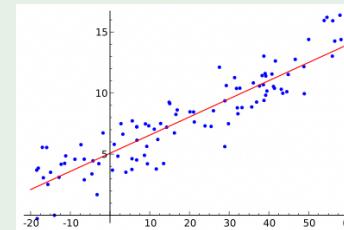


# Example: QP

## Constrained Least Squares

Given a set of measurements  $(x_1, y_1), \dots, (x_m, y_m)$ , where  $x_i \in \mathbb{R}^n$  is the  $i$ -th input and  $y_i \in \mathbb{R}$  is the  $i$ -th output, fit a linear function minimizing mean square error, subject to known bounds on the linear coefficients.

$$\begin{aligned} & \text{minimize } \|Ax - b\|_2^2 = x^T A^T A x - 2b^T A x + b^T b \\ & \text{subject to } l_i \leq x_i \leq u_i, \text{ for } i = 1, \dots, n. \end{aligned}$$





# Conic Optimization Problems (锥优化)

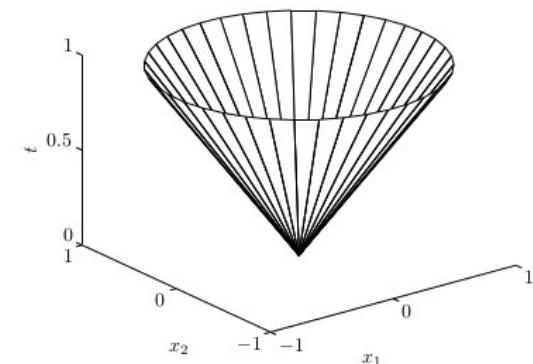
This is an umbrella term for problems of the following form

$$\begin{aligned} & \text{minimize} && c^\top x \\ & \text{subject to} && Ax + b \in K \end{aligned}$$

Where  $K$  is a convex cone (e.g.  $\mathbb{R}_{+}^n$ , positive semi-definite matrices, etc). Evidently, such optimization problems are convex.

**Example:** Second Order Cone Programming  
where  $K$  is a second order cone:

$$K = \{(x, t) : \|x\|_2 \leq t\}$$





# Semi-Definite Programming (SDP)

These are conic optimization problems where the cone in question is the set of positive semi-definite matrices.

$$\begin{aligned} & \text{minimize} && c^T x \\ & \text{subject to} && x_1 F_1 + x_2 F_2 + \dots + x_n F_n + G \succeq 0 \end{aligned}$$

where  $F_1, \dots, F_n, G$  are matrices, and  $\succeq$  refers to the positive semi-definite cone  $S^n_+$ .

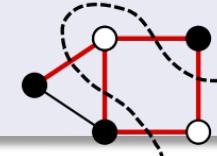


# Example: SDP

## Max Cut Problem

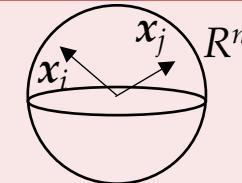
Given an undirected graph  $G = (V, E)$ , find a partition of  $V$  into  $(S, V \setminus S)$  maximizing number of edges with exactly one end in  $S$ .

$$\begin{aligned} & \text{maximize } \sum_{(i,j) \in E} (1 - x_i x_j)/2 \\ & \text{subject to } x_i \in \{-1, 1\}, \text{ for } i \in V. \end{aligned}$$



## Vector Program Relaxation

$$\begin{aligned} & \text{maximize } \sum_{(i,j) \in E} (1 - x_i x_j)/2 \\ & \text{subject to } \|x_i\|_2 = 1, \quad \text{for } i \in V. \\ & \quad x_i \in \mathbb{R}^n, \quad \text{for } i \in V. \end{aligned}$$



## SDP Relaxation

$$\begin{aligned} & \text{maximize } \sum_{(i,j) \in E} (1 - X_{ij})/2 \\ & \text{subject to } X_{ii} = 1, \quad \text{for } i \in V. \\ & \quad X \in \mathbb{S}_{++}^n, \quad \text{for } i \in V. \end{aligned}$$

$$X = \mathbf{x}\mathbf{x}^T = \begin{bmatrix} x_1 x_1 & \cdots & x_1 x_n \\ \vdots & \ddots & \vdots \\ x_n x_1 & \cdots & x_n x_n \end{bmatrix}$$

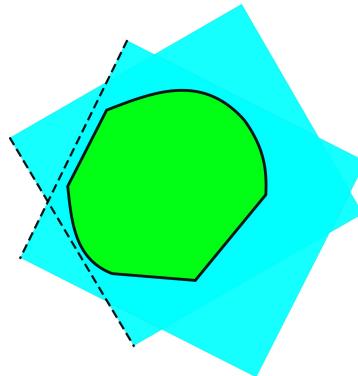


---

# Duality of Convex Optimization

# Geometric Duality of Convex Sets

---



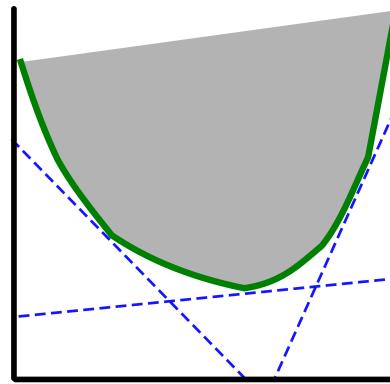
There are two equivalent ways to represent a convex set

- The family of points in the set (standard representation)
- The set of half-spaces containing the set ("dual" representation)

## Theorem

A closed convex set  $S$  is the intersection of all closed half spaces  $H$  containing it.

# Geometric Duality of Convex Functions



## Theorem

A convex function is the point-wise supremum of all affine functions under-estimating it everywhere.



# Recall: Duality in LP

## Primal LP

$$\begin{aligned} & \text{maximize} && c^T x \\ & \text{subject to} && Ax \leq b \end{aligned}$$

## Dual LP

$$\begin{aligned} & \text{minimize} && b^T y \\ & \text{subject to} && A^T y \geq c \\ & && y \geq 0 \end{aligned}$$

$$\text{maximize } 2x_1 + 3x_2$$

$$y_1 \times 4x_1 + 8x_2 \leq 12$$

$$y_2 \times 2x_1 + x_2 \leq 3$$

$$y_3 \times 3x_1 + 2x_2 \leq 4$$

$$\text{minimize } 12y_1 + 3y_2 + 4y_3$$

$$\text{subject to } 4y_1 + 2y_2 + 3y_3 \geq 2$$

$$8y_1 + y_2 + 2y_3 \geq 3$$

$$y_1, y_2, y_3 \geq 0$$



$$\begin{aligned} & (4y_1 + 2y_2 + 3y_3)x_1 + (8y_1 + y_2 + 2y_3)x_2 \leq 12y_1 + 3y_2 + 4y_3 \\ & z^* = 2x_1 + 3x_2 \leq (4y_1 + 2y_2 + 3y_3)x_1 + (8y_1 + y_2 + 2y_3)x_2 \end{aligned}$$

Dual LP **relax/soften** the constraints by replacing each with a linear “penalty term” or “cost” in the objective.



# The Lagrangian

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad \text{for } i = 1, \dots, m. \\ & && h_i(x) = 0, \quad \text{for } i = 1, \dots, k. \end{aligned}$$

The basic idea of Lagrangian duality is to **relax/soften** the constraints by replacing each with a linear “**penalty term**” or “**cost**” in the objective.

## The Lagrangian Function

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^k \nu_i h_i(x)$$

- $\lambda_i$  is **lagrange multiplier** for  $i$ -th inequality constraint
  - required to be **nonnegative**
- $\nu_i$  is **lagrange multiplier** for  $i$ -th equality constraint
  - allowed to be arbitrary sign



# Lagrange Dual Function

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

The lagrange dual function gives the optimal value of the primal problem subject to the **softened constraints**.

## The Lagrange Dual Function

$$g(\lambda, v) = \inf_{x \in \mathcal{D}} L(x, \lambda, v) = \inf_{x \in \mathcal{D}} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^k v_i h_i(x))$$

- Observe that:
  - The infimum is **unconstrained** (as opposed to the original constrained minimization problem)
  - $g$  is **concave regardless of original problem**
  - The Lagrange dual can be **unbounded** ( $-\infty$ ) for some  $\lambda$  and  $v$ .



# Interpretation: “Soft” Lower Bound

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

## The Lagrange Dual Function

$$g(\lambda, v) = \inf_{x \in \mathcal{D}} L(x, \lambda, v) = \inf_{x \in \mathcal{D}} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^k v_i h_i(x))$$

## Fact

$g(\lambda, v)$  is a **lower bound** on  $\text{OPT}(\text{primal})$  for every  $\lambda \geq 0$  and  $v \in \mathbb{R}^k$ .

## Proof

- Every primal feasible  $x$  incurs nonpositive penalty by  $L(x, \lambda, v)$
- Therefore,  $L(x^*, \lambda, v) \leq f_0(x^*)$
- So  $g(\lambda, v) \leq f_0(x^*) = \text{OPT}(\text{Primal})$



# Interpretation: “Soft” Lower Bound

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

## The Lagrange Dual Function

$$g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu) = \inf_{x \in \mathcal{D}} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^k \nu_i h_i(x))$$

## Interpretation

- A “hard” feasibility constraint can be thought of as imposing a penalty of  $+\infty$  if violated, and a penalty/reward of 0 if satisfied.
- Lagrangian imposes a “soft” linear penalty for violating a constraint, and a reward for slack.
- Lagrange dual finds the optimal subjects to these soft constraints.



# Lagrange Dual Problem

This is the problem of finding the **best lower bound** on OPT(primal) obtained from the Lagrange dual function

$$\begin{aligned} & \underset{\lambda, \nu}{\text{maximize}} && g(\lambda, \nu) \\ & \text{subject to} && \lambda \geq 0. \end{aligned}$$

- Note: this is a **convex** optimization problem, regardless of whether primal problem was convex
- By convention, sometimes we add “dual feasibility” constraints to impose “nontrivial” lower bounds (i.e.,  $g(\lambda, \nu) \geq -\infty$ )
- $(\lambda^*, \nu^*)$  solving the above are referred to as the **dual optimal solution**.



# Lagrange Dual Problem

## Example: Least-Norm Solution of Linear Equations

- Consider the problem

$$\begin{aligned} & \underset{x}{\text{minimize}} && x^T x \\ & \text{subject to} && Ax = b. \end{aligned}$$

- The Lagrangian is

$$L(x, \nu) = x^T x + \nu^T (Ax - b).$$

- To find the dual function, we need to solve an unconstrained minimization of the Lagrangian. We set the gradient equal to zero

$$\nabla_x L(x, \nu) = 2x + A^T \nu = 0 \implies x = -(1/2) A^T \nu$$



# Lagrange Dual Problem

## Example: Least-Norm Solution of Linear Equations

and we plug the solution in  $L$  to obtain  $g$ :

$$g(\nu) = L\left(-\left(1/2\right)A^T\nu, \nu\right) = -\frac{1}{4}\nu^T A A^T \nu - b^T \nu$$

- The function  $g$  is, as expected, a concave function of  $\nu$ .
- From the lower bound property, we have

$$\text{OPT(primal)} \geq -\frac{1}{4}\nu^T A A^T \nu - b^T \nu \text{ for all } \nu.$$

- The dual problem is the QP

$$\underset{\nu}{\text{maximize}} \quad -\frac{1}{4}\nu^T A A^T \nu - b^T \nu.$$

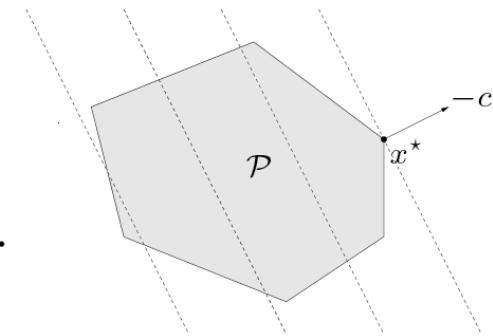


# Lagrange Dual Problem

## Example: Standard Form LP

- Consider the problem

$$\begin{aligned} & \underset{x}{\text{minimize}} && c^T x \\ & \text{subject to} && Ax = b, \quad x \geq 0. \end{aligned}$$



- The Lagrangian is

$$\begin{aligned} L(x, \lambda, \nu) &= c^T x + \nu^T (Ax - b) - \lambda^T x \\ &= (c + A^T \nu - \lambda)^T x - b^T \nu. \end{aligned}$$

- $L$  is a linear function of  $x$  and it is unbounded if the term multiplying  $x$  is nonzero.



# Lagrange Dual Problem

---

- Hence, the dual function is

$$g(\lambda, \nu) = \inf_x L(x, \lambda, \nu) = \begin{cases} -b^T \nu & \text{if } c + A^T \nu - \lambda = 0 \\ -\infty & \text{otherwise.} \end{cases}$$

- The function  $g$  is a concave function of  $(\lambda, \nu)$  as it is linear on an affine domain.
- From the lower bound property, we have

$$\text{OPT(primal)} \geq -b^T \nu \quad \text{if } c + A^T \nu \geq 0.$$

- The dual problem is the LP

$$\begin{aligned} & \underset{\nu}{\text{maximize}} && -b^T \nu \\ & \text{subject to} && c + A^T \nu \geq 0. \end{aligned}$$



# Weak duality

## Primal Problem

$$\begin{aligned} & \min f_0(x) \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad \forall i = 1, \dots, m. \\ & h_i(x) = 0, \quad \forall i = 1, \dots, k. \end{aligned}$$

## Dual Problem

$$\begin{aligned} & \max g(\lambda, \nu) \\ \text{s.t.} \quad & \lambda \succeq 0 \end{aligned}$$

### Weak Duality

$$\text{OPT(dual)} \leq \text{OPT(primal)}$$

Proof:  $g(\lambda, \nu) \leq \text{OPT(primal)}$  for feasible  $(\lambda, \nu)$

- Holds for every (even nonconvex) optimization problem
- **Duality Gap:** difference between optimal dual and primal values  
$$\text{OPT(primal)} - \text{OPT(dual)}$$
- Solving the dual problem may be used to find nontrivial lower bounds for difficult problems.



# Strong duality

Even more interesting is when **equality** is achieved in weak duality which is called **strong duality**:

## Strong Duality

$$\text{OPT}(\text{dual}) = \text{OPT}(\text{primal})$$

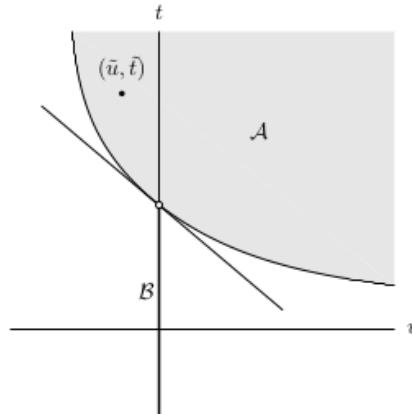
- Equivalently: there exists a setting of Lagrange multipliers so that  $g(\lambda, v)$  gives a tight lower bound on primal optimal value.
- Very desirable (we can solve a difficult problem by solving the dual)
- In general, **does not hold** for **non-convex** optimization problems
- Usually, but not always, holds for convex optimization problems.



# Slater's Condition for Strong Duality

## Slater's Condition

If the primal is a convex problem, and there exists at least one point  $x$  where all inequality constraints are strictly satisfied (i.e.,  $f_1(x) < 0, \dots, f_m(x) < 0$  and  $h_1(x) = \dots = h_n(x) = 0$ ). Then strong duality holds.



- A sufficient condition for strong duality.
- Forces supporting hyperplane to be non-vertical
- Can be weakened to requiring strict feasibility only of non-affine constraints



# Other Conditions

---

## KKT Conditions

[Optional for after-class reading]



---

# Optimization Algorithms

## Unconstrained Minimization

- Gradient descent, Newton's method

## Interior-point Methods



---

# Review of Calculus

## Gradient, Hessian

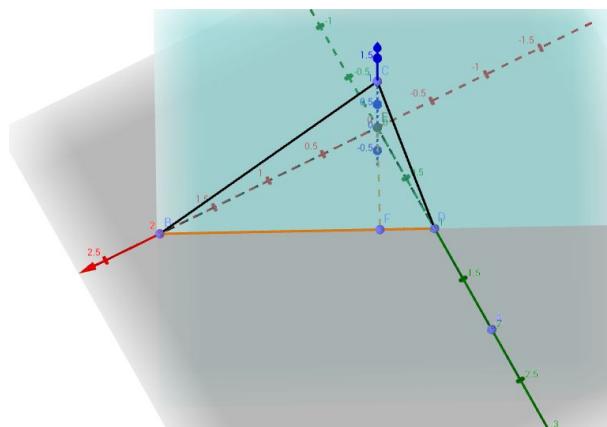


# Gradient

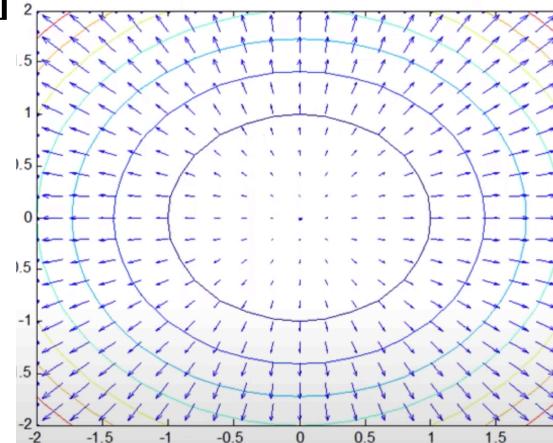
**gradient – “direction and rate of fastest increase”**

For  $f: R^n \rightarrow R$ , its gradient  $\nabla f: R^n \rightarrow R^n$  is defined at the point  $p = (x_1, \dots, x_n)$  as the vector

$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$$



<https://www.zhihu.com/question/36301367>

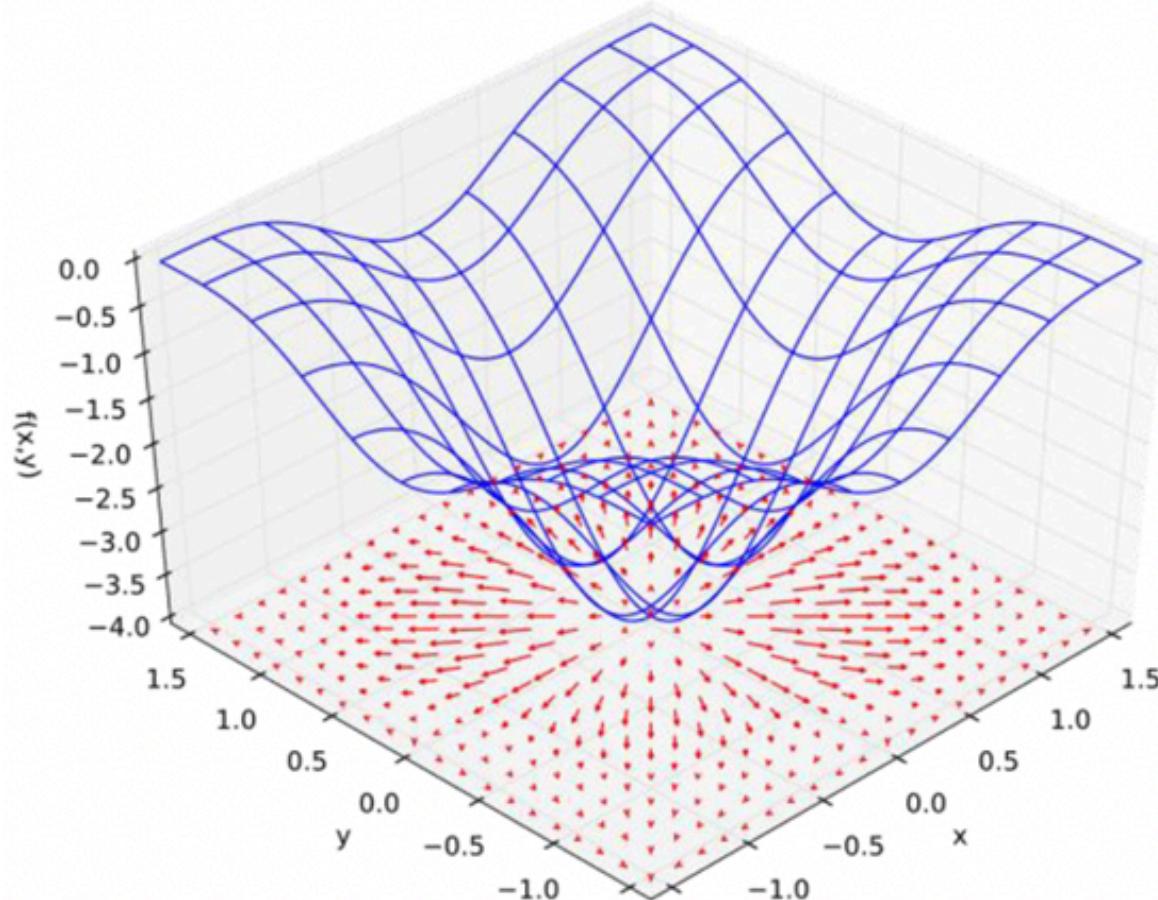


Contours for  $x^2+y^2$  with gradients



# Gradient

gradients for arbitrary function  $f(x,y)$





# Gradient Vector

## Example

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad f(x) = 2x_1^2 x_2 - x_1 x_3^3$$

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \frac{\partial}{\partial x_2} f(x) \\ \frac{\partial}{\partial x_3} f(x) \end{bmatrix} = \begin{bmatrix} 4x_1 x_2 - x_3^3 \\ 2x_1^2 \\ -3x_1 x_3^2 \end{bmatrix}$$



# Hessian

A square matrix of second-order partial derivatives of a scalar-valued function.

- The Hessian of  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is

$$\nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_d^2} \end{pmatrix}$$

Hessian describes the local curvature (曲率) of a function of many variables



# Unconstrained Minimization

---

$$\text{minimize} \quad f(x)$$

- $f$  convex, twice continuously differentiable (hence  $\text{dom } f$  open)
- we assume optimal value  $p^* = \inf_x f(x)$  is attained (and finite)

## unconstrained minimization methods

- produce sequence of points  $x^{(k)} \in \text{dom } f$ ,  $k = 0, 1, \dots$  with
$$f(x^{(k)}) \rightarrow p^*$$
- can be interpreted as iterative methods for solving optimality condition

$$\nabla f(x^*) = 0$$



# Unconstrained Minimization

---

## Algorithms

- Gradient Descend
- Steepest Descend
- Newton's Method
- ...



# Gradient Descend

## The simplest optimization algorithm

- Goal:

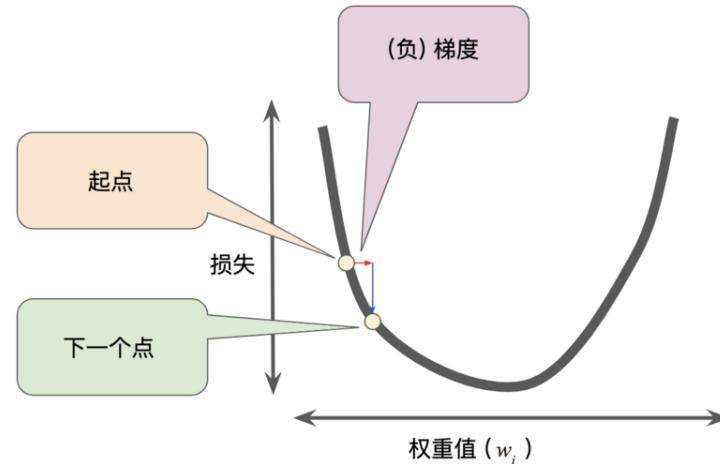
$$\min_x f(x)$$

Unconstrained  
Optimization

- Iteration:

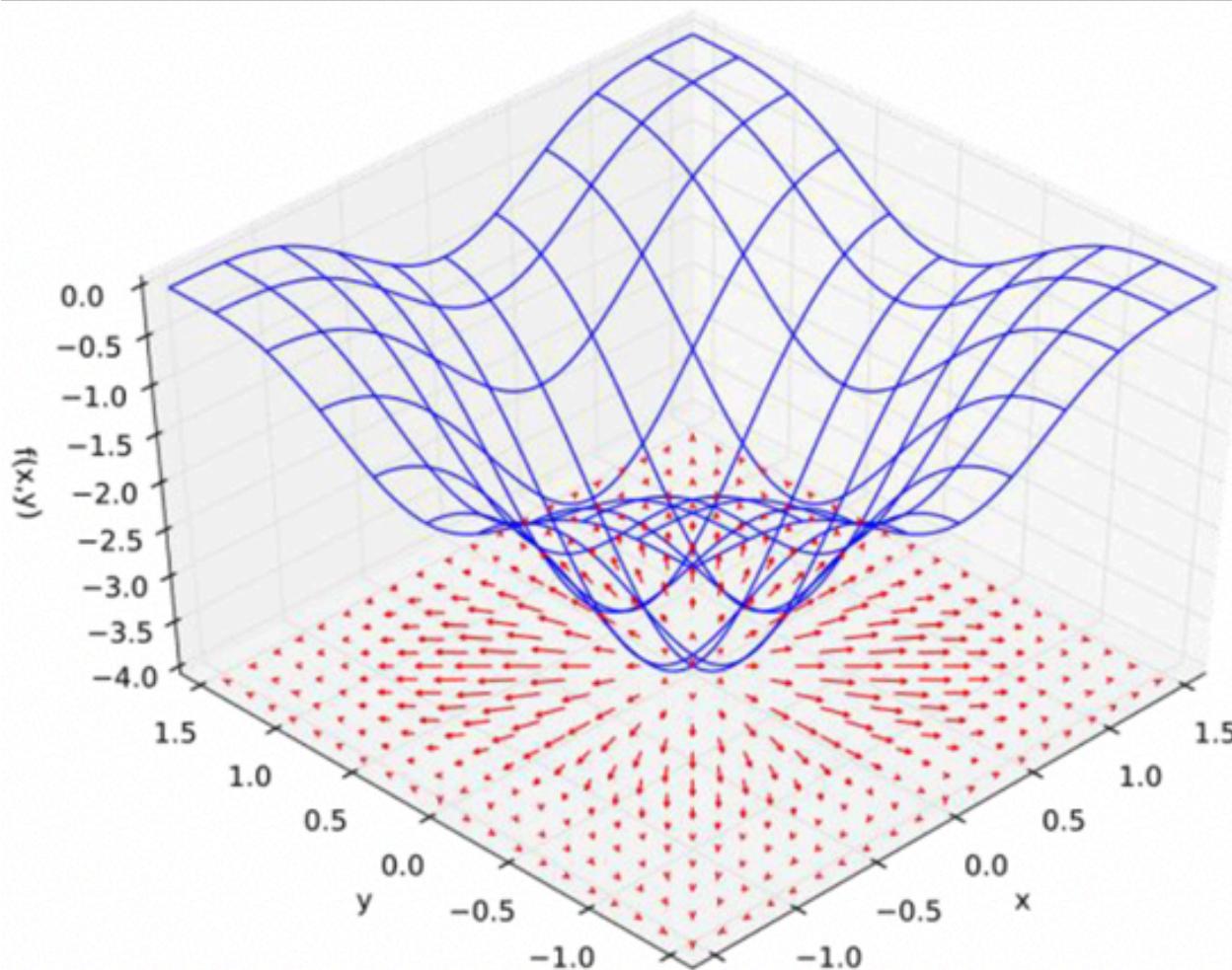
$$x_{t+1} = x_t - \eta_t \nabla f(x_t)$$

- $\eta_t$  is step size.





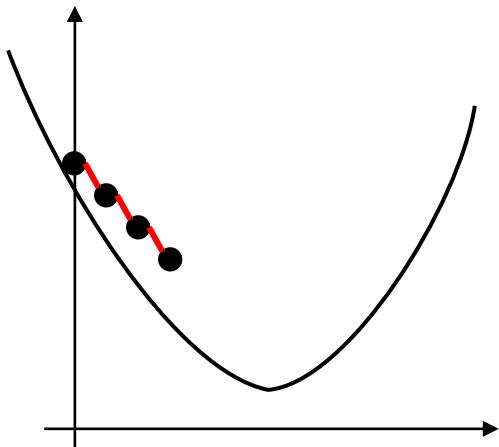
# Gradient Descend



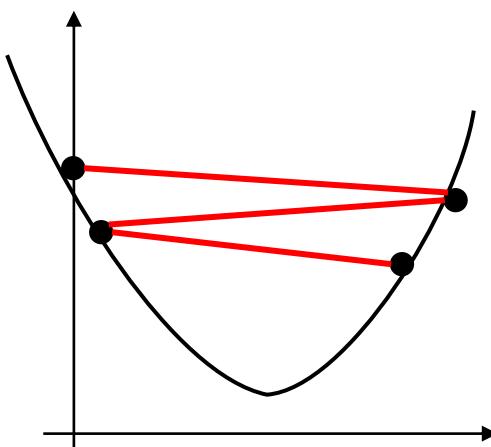


# How to Choose Step Size?

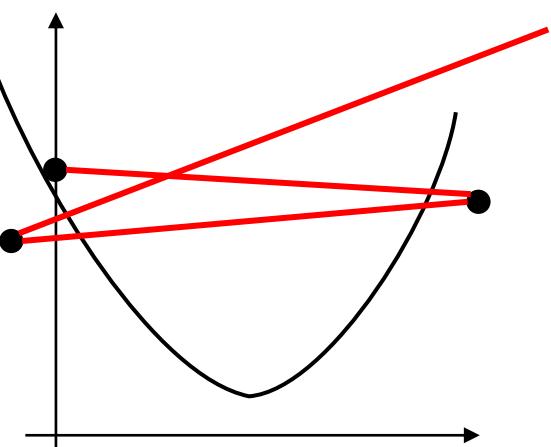
- If step size is too big, the value of function can diverge.
- If step size is too small, the convergence is very slow.



$\eta$  too small:  
slow progress



$\eta$  too large:  
oscillations



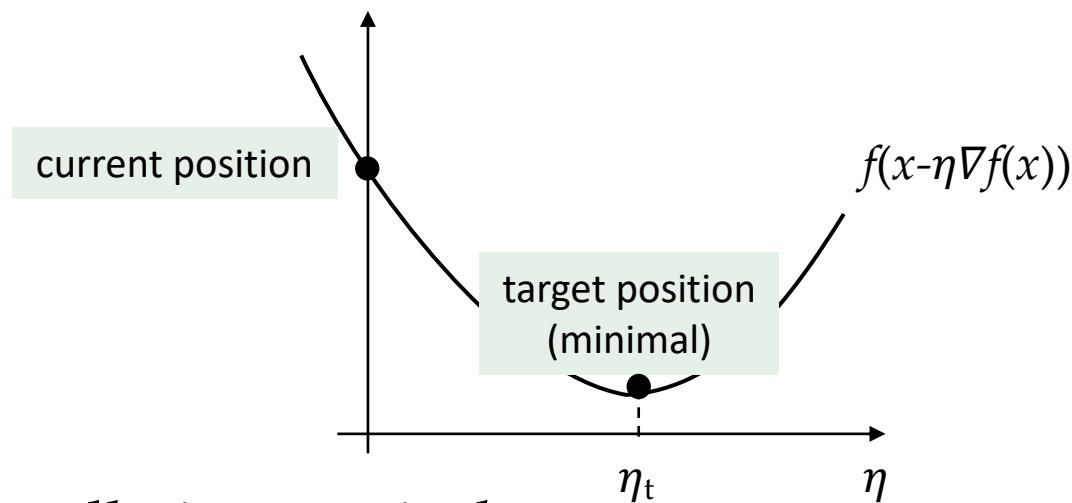
$\eta$  much too large:  
instability



# On the Step Size: Exact Line Search

- Idea: the best step directly leads the function value to its minimal

$$\eta_t = \arg \min_{\eta} f(x - \eta \nabla f(x))$$



- Usually impractical.

# On the Step Size: Backtracking Line Search

(回溯线搜索)



Using a relaxed tangent to approximate the upper bound of the minimal.

- Let  $\alpha \in (0, 1/2]$ ,  $\beta \in (0, 1)$ . Start with  $\eta = 1$  and multiply  $\eta = \beta\eta$  until

$$f(x - \eta \nabla f(x)) \leq f(x) - \alpha \eta \|\nabla f(x)\|^2$$

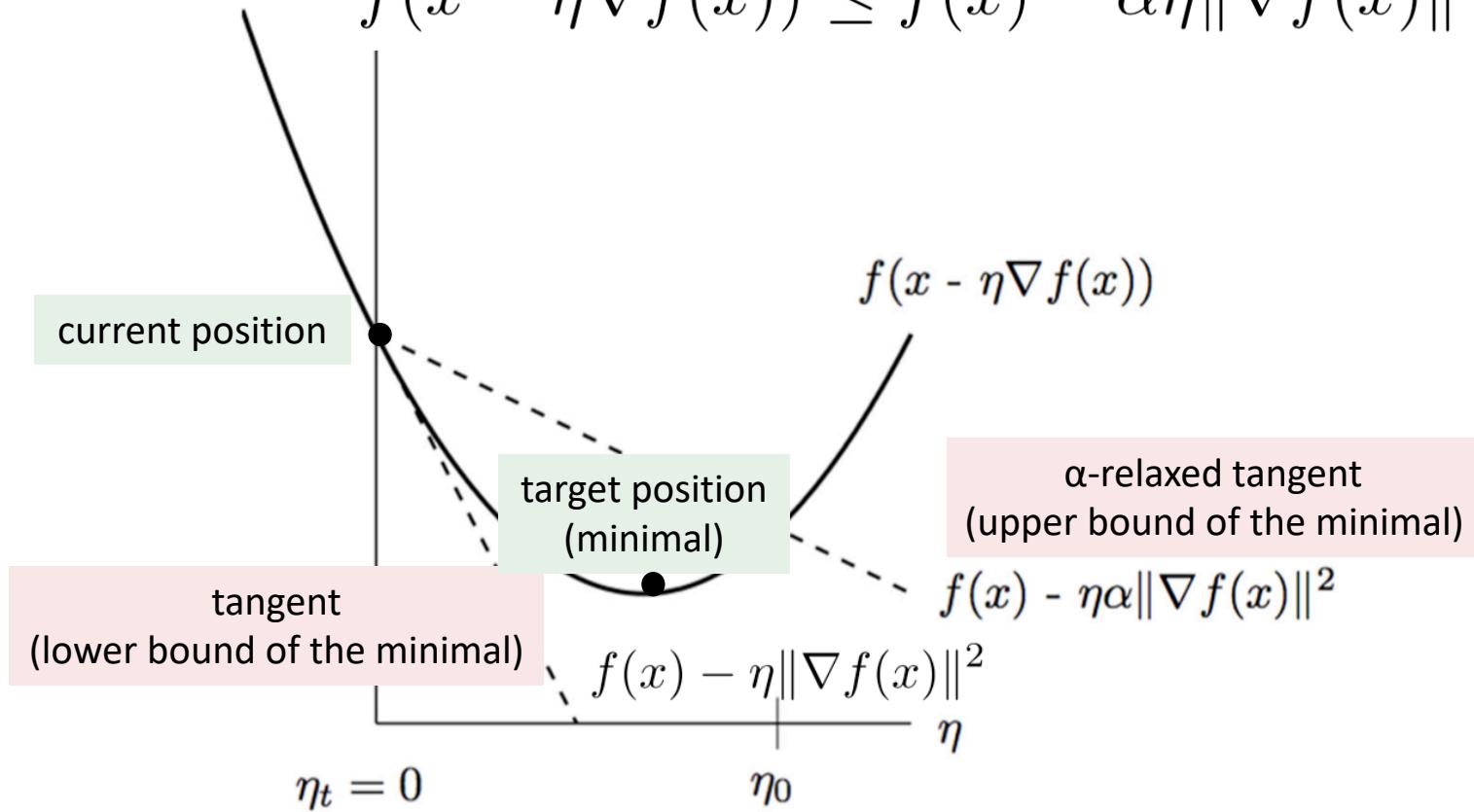
- Work well in practice.

# On the Step Size: Backtracking Line Search



The  $\alpha$ -relaxed tangent approximates an upper-bound of the minimal

$$f(x - \eta \nabla f(x)) \leq f(x) - \alpha \eta \|\nabla f(x)\|^2$$

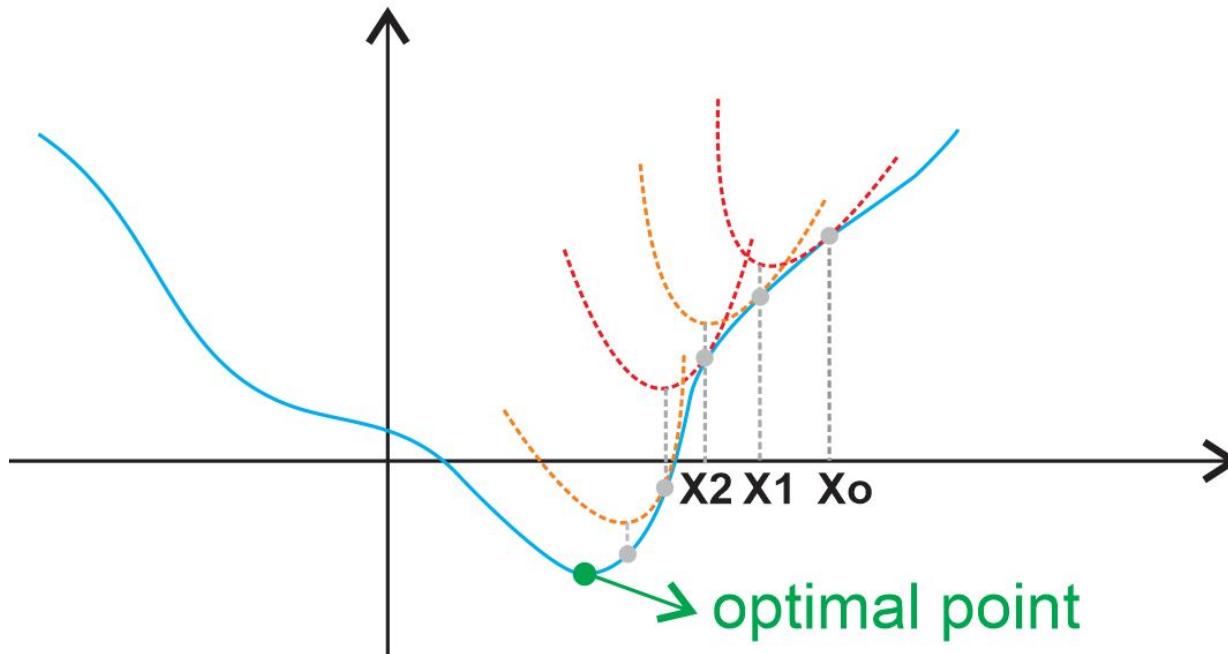




# Newton's Method

The next step moves to the minimal of the second order approximation of  $f$

$$x_{t+1} = x_t - [\nabla^2 f(x_t)]^{-1} \nabla f(x_t)$$





# Newton's Method

## interpretations

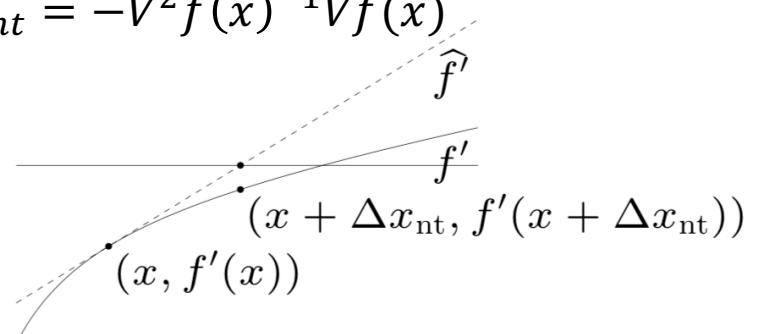
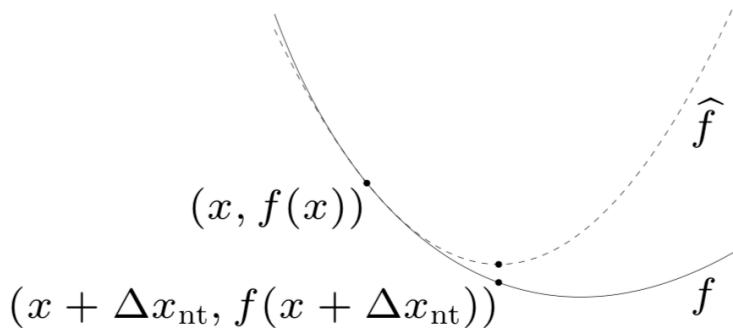
- $x + \Delta x_{nt}$  minimizes second order approximation

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

- $x + \Delta x_{nt}$  solves linearized optimality condition

$$\nabla f(x + v) \approx \nabla \hat{f}(x + v) = \nabla f(x) + \nabla^2 f(x)v = 0$$

$$\Rightarrow v^* = \nabla_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$$





# Interior-point Methods

Unconstrained minimization → Constrained minimization

- inequality constrained minimization
- logarithmic barrier function and central path
- barrier method



# Inequality constrained minimization

What about minimization with **inequality** constrains?

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

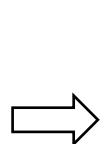
where  $f_i$ 's are assumed to be convex & twice continuously differentiable



# Trapping Points within the Feasible Set

The inequality constraints can be made implicit by rewriting (1) as

$$\begin{aligned} & \min f_0(x) \\ \text{s.t. } & Ax = b, \quad f_i(x) \leq 0 \end{aligned}$$



$$\begin{aligned} & \min f_0(x) + \sum_{i=1}^m I_-(f_i(x)) \\ \text{s.t. } & Ax = b \end{aligned}$$

where  $I_-(u) = 0$  for  $u \leq 0$ ,  $I_-(u) = \infty$  otherwise.

But  $I_-$  is not differentiable.

The basic idea: approximate  $I_-$  by some differentiable function.



# Logarithmic Barrier Function

Approximate  $I_-$  by

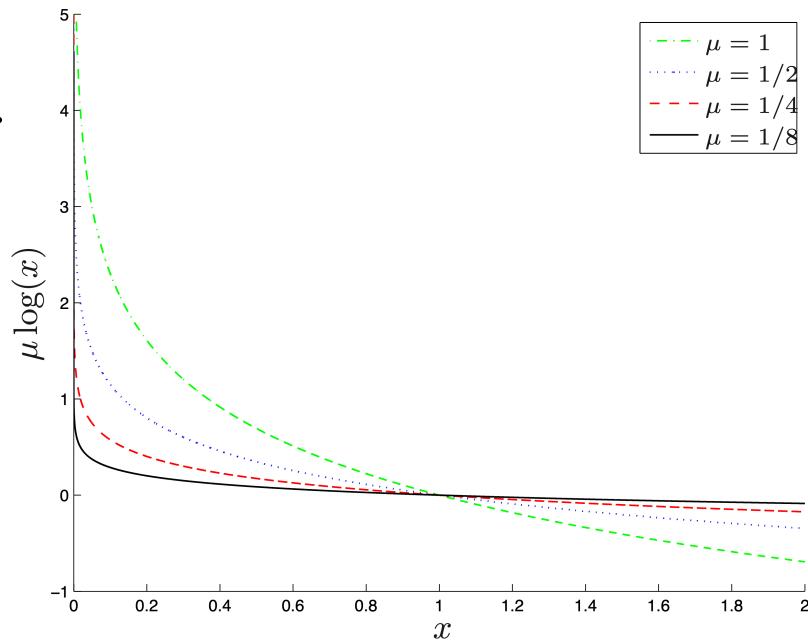
$$I_-(u) = -\mu \log(-\mu), \quad \text{dom } I_- = \{x \in \mathbb{R} \mid x < 0\}$$

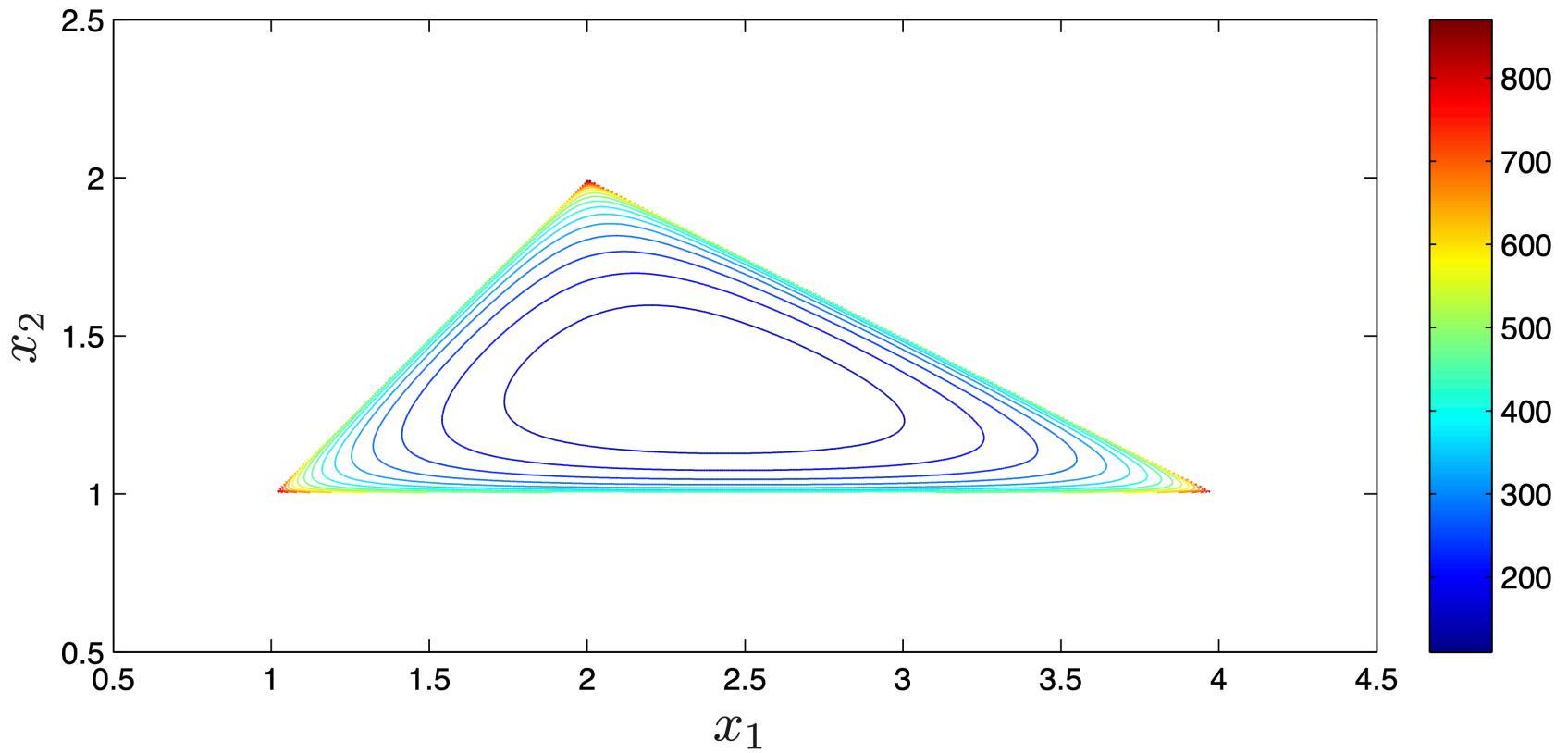
where  $\mu > 0$  is a parameter that controls the accuracy of the approx.

Example: a single inequality  $x \geq 0$ .

The corresponding approx. is

$$-\mu \log(x)$$





**Example:** A set of linear inequalities  $b_i - a_i^T x \leq 0$ ,  $i = 1, \dots, m$ .

The corresponding approx. is

$$-\mu \sum_{i=1}^m \log(a_i^T x - b_i)$$



# Logarithmic Barrier Function

The original problem can be approximated as:

$$\begin{aligned} & \min f_0(x) + \mu\phi(x) \\ \text{s.t. } & Ax = b \end{aligned}$$

$\min f(x)$  s.t.  $Ax=b$  is  
equivalent to  $\min_z f(Az+b+Fz)$

with  $\phi(x) = -\sum_{i=1}^m \log(-f_i(x))$ ,  $\text{dom } \phi = \{x \mid f_1(x) < 0, \dots, f_m(x) < 0\}$

$\varphi$  is called the **logarithmic barrier function**. Some nice properties:

- $\varphi$  is convex (by composition).
- $\varphi$  is twice differentiable:

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

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

This means that the objective fn. is convex and twice differentiable.



# Central Path

The log barrier approximation

$$\begin{aligned} \min \quad & f_0(x) + \mu \phi(x) \\ \text{s.t. } & Ax = b \end{aligned} \tag{*}$$

- is an accurate approximation of the original problem for  $\mu \rightarrow 0$  (with  $\mu > 0$ )
- becomes difficult to minimize as  $\mu \rightarrow 0$  (at least for Newton's method)

**The basic idea:** start with a large  $\mu$ , and iteratively reduce  $\mu$  until a desired solution accuracy is reached.

Define  $x^*(\mu)$  to be the solution of (\*).

**Central path**     $\{x \mid x=x^*(\mu), \mu > 0\}$

is the collection of optimal points for various  $\mu$ . A property:

$$f_0(x^*(\mu)) - p^* \leq m\mu$$



# Central Path

central path,  $\{x \mid x = x^*(\mu), \mu > 0\}$

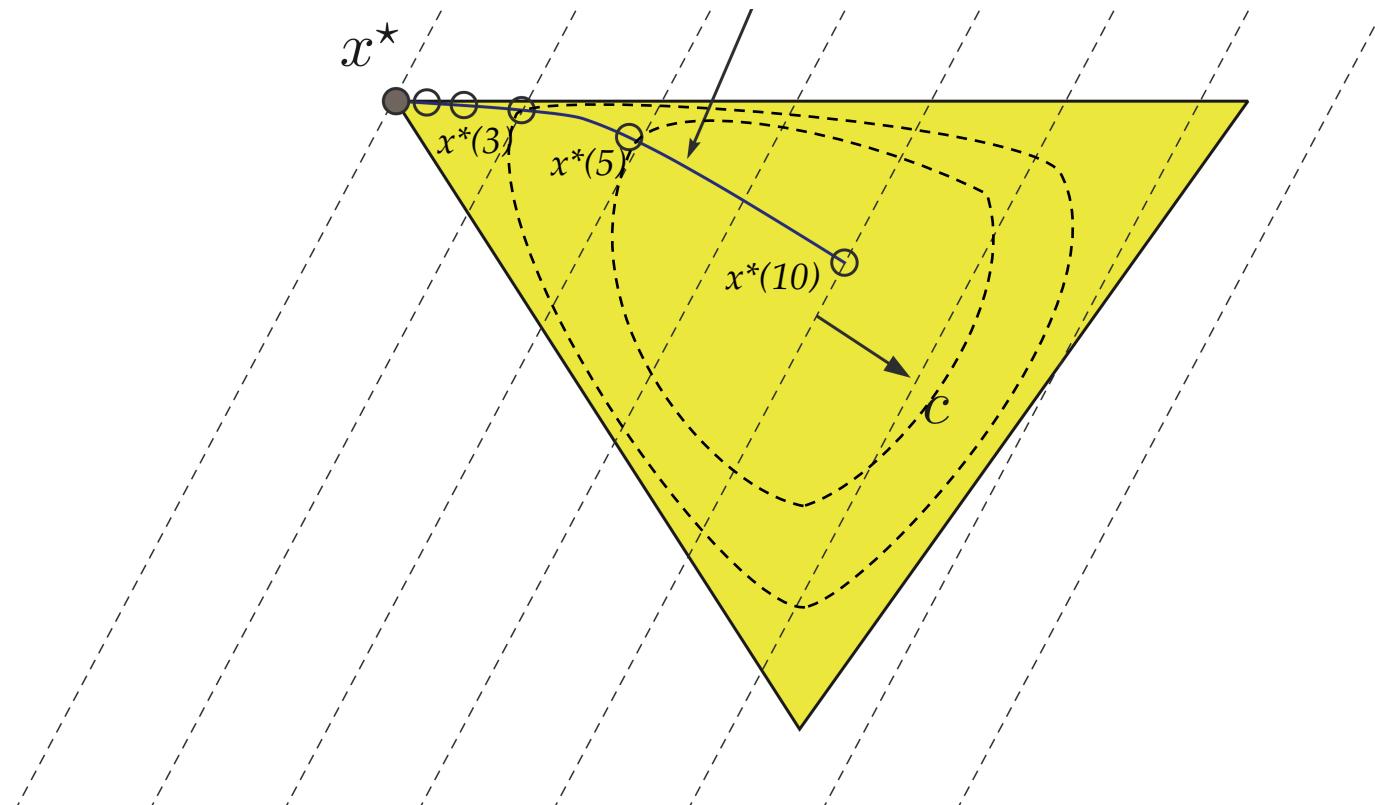


Illustration of the central path of LP. An interior-point method would follow the central path to iteratively approach the optimal solution.

# Barrier Method (Path Following Method)

---



## *Path Following Method*

**given** an initial strictly feasible point  $x, \mu, \epsilon > 0$ , &  $\beta < 1$ .

**repeat**

1. *Centering step.* Starting at  $x$ , use Newton's method to solve

$$x^*(\mu) = \min_x \mu f_0(x) + \phi(x) \text{ s.t. } Ax = b,$$

2. *Update.*  $x := x^*(\mu)$ .

3. *Stopping criterion.* **quit** if  $m\mu < \epsilon$ .

4. *Target shifting.*  $\mu := \beta\mu$ .
- 

- **Short-step path following:** choose  $\beta$  close to 1.
    - small number of Newton steps per outer iteration
    - large number of outer iterations
  - **Long-step path following:** choose a small  $\beta$ .
    - increased number of Newton steps per outer iteration
    - smaller number of outer iterations
-