

# Maximum Margin Classifiers: Support Vector Machines

Machine Learning and Pattern Recognition:  
Lecture 14

Sumit Chopra

# Outline of the Talk

- Quick Tutorial on Optimization
  - Basic idea behind Support Vector Machines
  - Optimization concepts and terminology
- Support Vector Machines in Detail
  - Given by **Fu Jie Huang**

# Binary Classification Problem

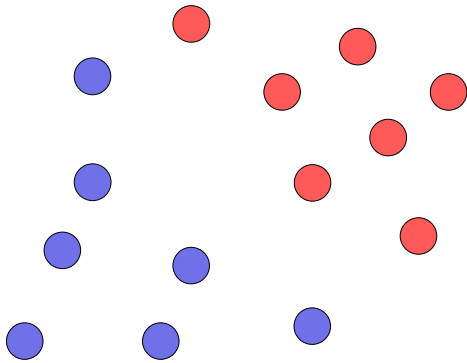
- **Given:** Training data generated according to the distribution  $D$

$$(x_1, y_1), \dots, (x_p, y_p) \in \mathcal{R}^n \times \{-1, 1\}$$

- **Problem:** Find a classifier (a function)  $h(x): \mathcal{R}^n \rightarrow \{-1, 1\}$  such that it generalizes well on the test set obtained from the same distribution  $D$
- **Solution:**
  - **Linear Approach:** linear classifiers - perceptron and many other.
  - **Non Linear Approach:** non-linear classifiers - neural nets and many other.

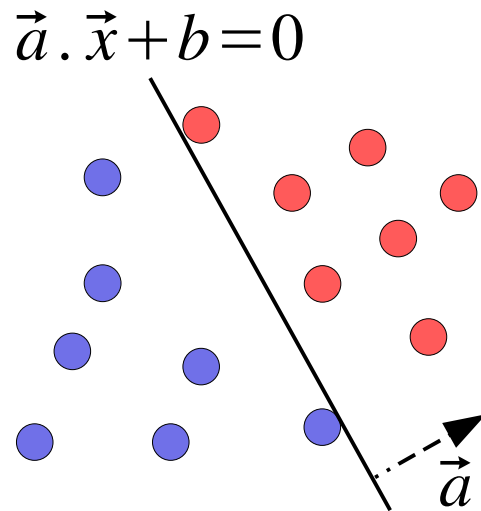
# Linearly Separable Data

- Assume that the training data is linearly separable



# Linearly Separable Data

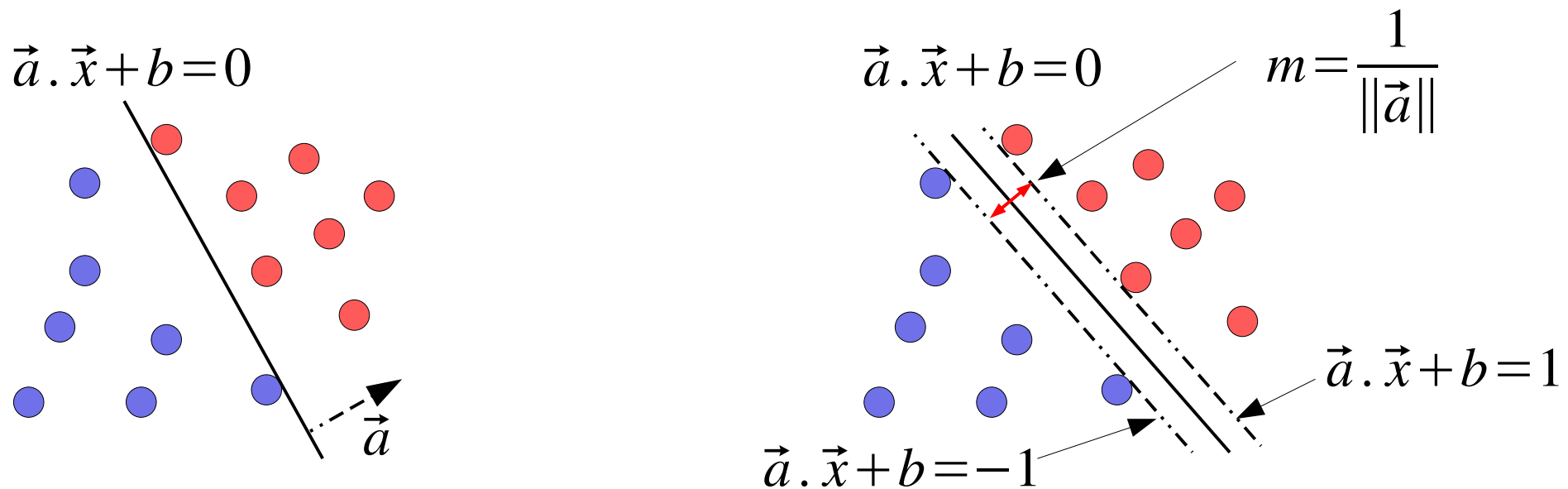
- Assume that the training data is linearly separable



- Then the classifier is:  $h(x) = \vec{a} \cdot \vec{x} + b$  where  $a \in \mathcal{R}^n, b \in \mathcal{R}$
- Inference:  $\text{sign}(h(x)) \in \{-1, 1\}$

# Linearly Separable Data

- Assume that the training data is linearly separable



- For the Closest Points:  $h(x) = \vec{a} \cdot \vec{x} + b \in -1, 1$
- Margin:  $m = \frac{1}{\|\vec{a}\|}$

# Optimization Problem

- Its a Constrained Optimization Problem

$$\min_x \frac{1}{2} \|\vec{a}\|^2$$

*s.t.:*

$$y_i (\vec{x}_i \cdot \vec{a} + b) \geq 1, \quad i=1, \dots, p$$

- A convex optimization problem
- Constraints are affine hence convex

# Optimization: Some Theory

- The problem:

$$\begin{aligned} \min_x f_0(x) & \longleftarrow \text{objective function} \\ \text{s.t. :} \\ f_i(x) \leq 0, \quad i = 1, \dots, m & \longleftarrow \text{inequality constraints} \\ h_i(x) = 0, \quad i = 1, \dots, p & \longleftarrow \text{equality constraints} \end{aligned}$$

- Solution of problem:  $x^o$

- Global Optimum – if the problem is convex
- Local Optimum – if the problem is not convex



# Optimization: Some Theory

- **Example:** Standard Linear Program (LP)

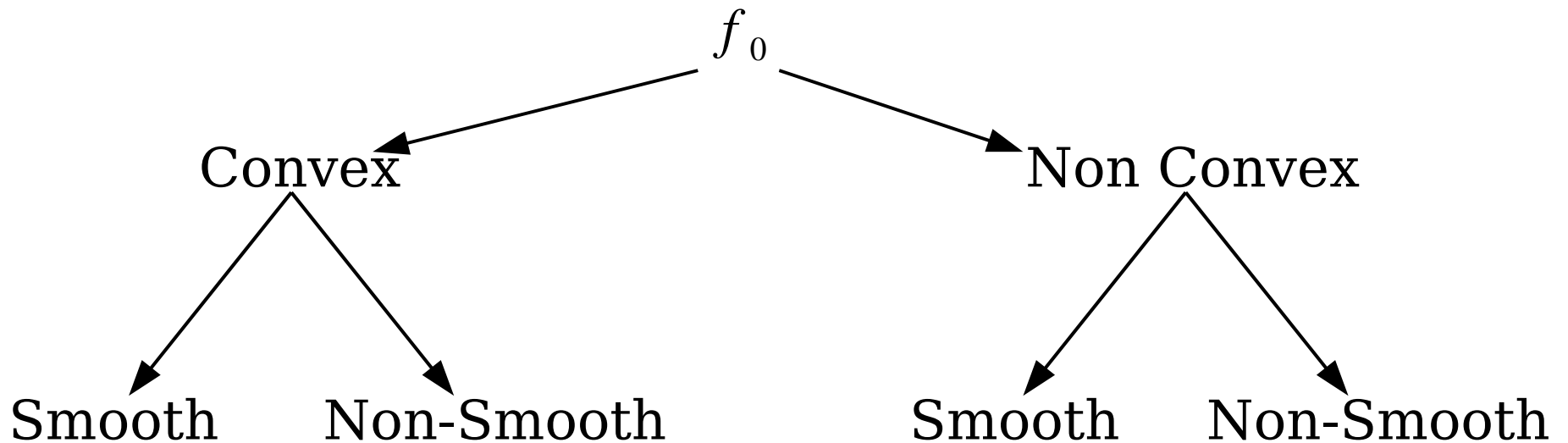
$$\begin{aligned} \min_x \quad & c^T x \\ \text{s.t.} \quad & \\ & Ax = b \\ & x \geq 0 \end{aligned}$$

- **Example:** Least Squares Solution of Linear Equations

$$\begin{aligned} \min_x \quad & x^T x \\ \text{s.t.} \quad & \\ & Ax = b \end{aligned}$$

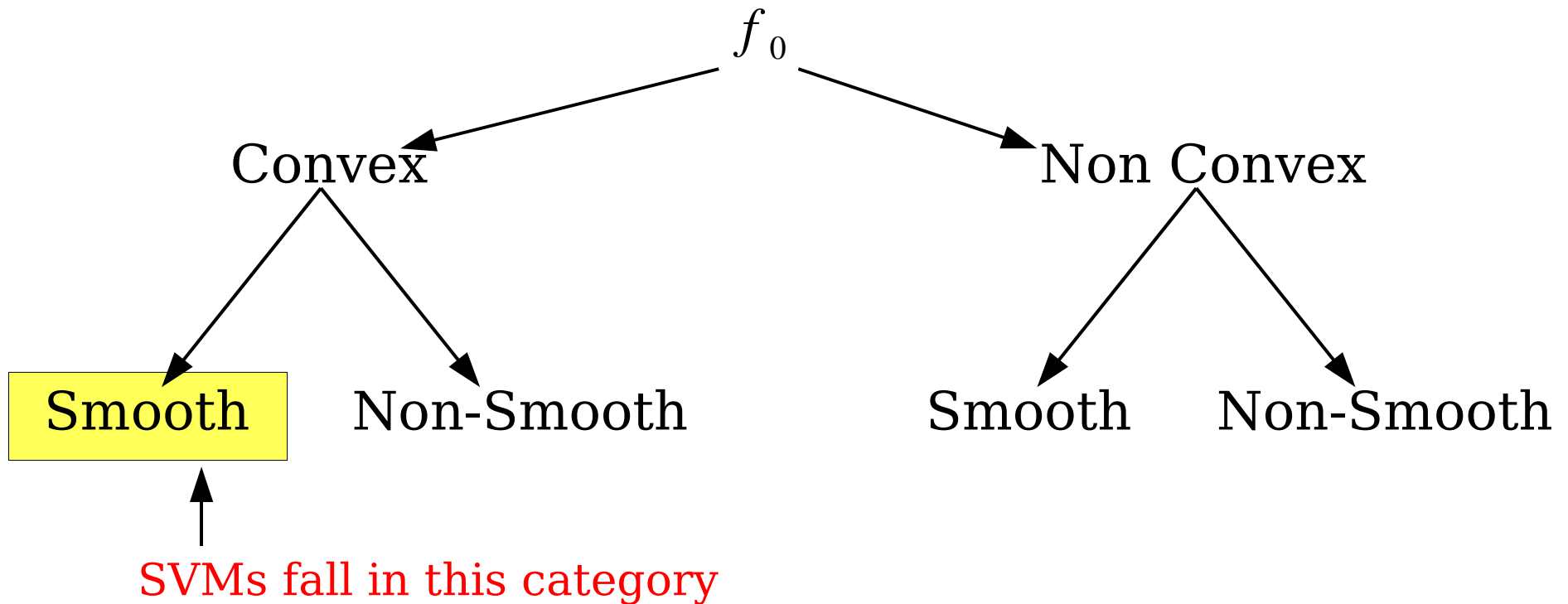
# The Big Picture

- Constrained / Unconstrained Optimization
- Hierarchy of object function



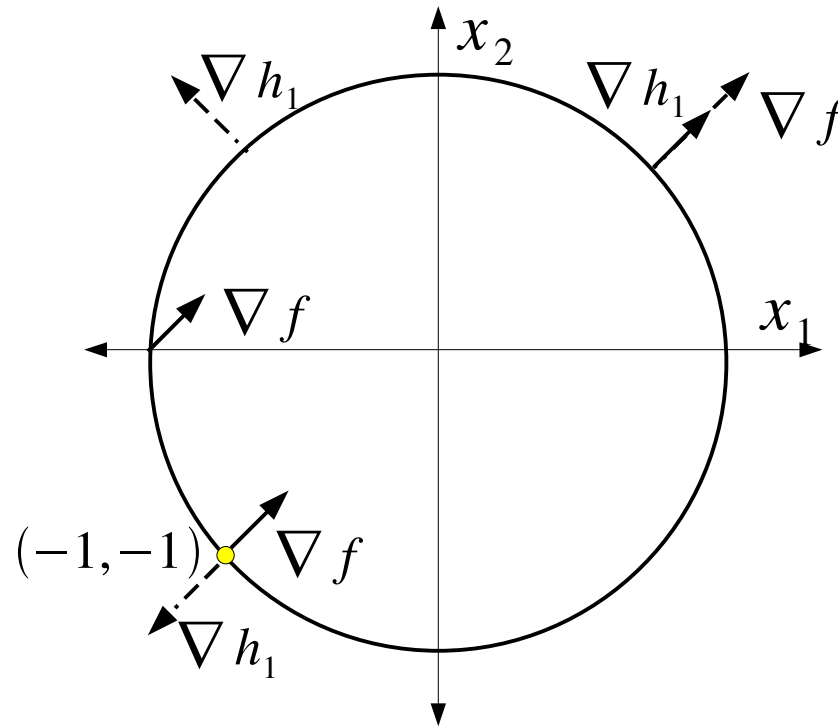
# The Big Picture

- Constrained / Unconstrained Optimization
- Hierarchy of object function



# A Toy Example: Equality Constraint

- Example 1:  $\min x_1 + x_2$   
 $s.t.: x_1^2 + x_2^2 - 2 = 0 \equiv h_1$



- At Optimal Solution:

$$\nabla f(x^o) = \lambda_1^o \nabla h_1(x^o)$$

# A Toy Example: Equality Constraint

- $x$  is not an optimal solution, if there exists an  $s \neq 0$  such that

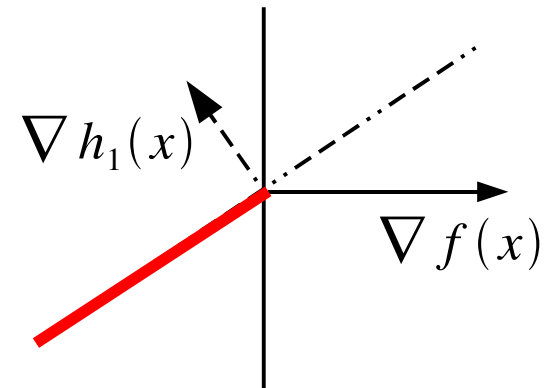
$$\begin{aligned}h_1(x+s) &= 0 \\ f(x+s) &< f(x)\end{aligned}$$

- Using first order Taylor's expansion

$$h_1(x+s) = h_1(x) + \nabla h_1(x)^T s = \nabla h_1(x)^T s = 0 \quad (1)$$

$$f(x+s) - f(x) = \nabla f(x)^T s < 0 \quad (2)$$

- Such an  $s$  can exist only when  $\nabla h_1(x)$  and  $\nabla f(x)$  are not parallel



# A Toy Example: Equality Constraint

- Thus we have

$$\nabla f(x^o) = \lambda_1^o \nabla h_1(x^o)$$

- The Lagrangian

$$L(x, \lambda_1) = f(x) - \lambda_1 h_1(x)$$

Lagrange multiplier  
or dual variable for  $h_1$



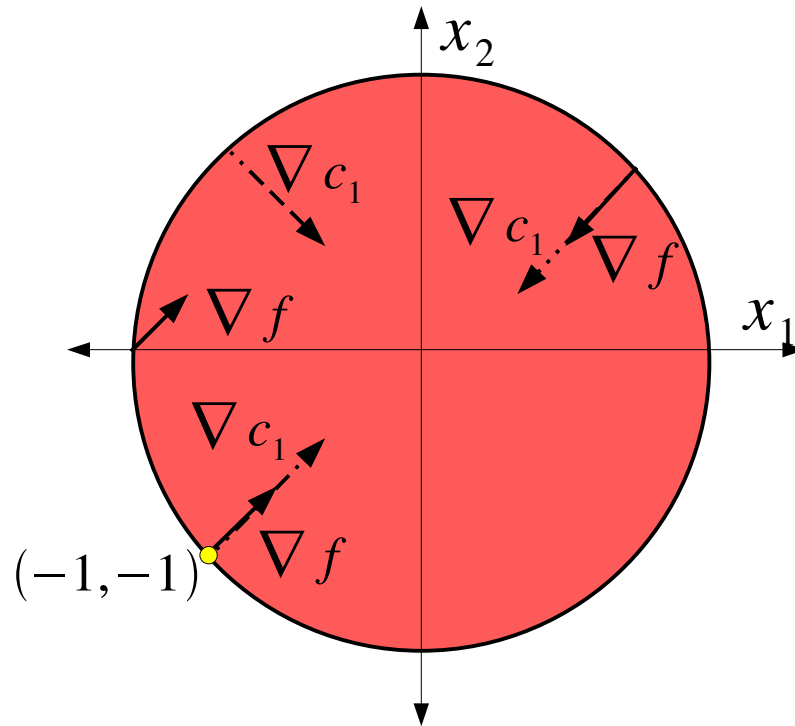
- Thus at the solution

$$\nabla_x L(x^o, \lambda_1^o) = \nabla f(x^o) - \lambda_1^o \nabla h_1(x^o) = 0$$

- This is just a necessary condition and not a sufficient condition.

# A Toy Example: Inequality Constraint

- Example 1:  $\min x_1 + x_2$   
 $s.t.: 2 - x_1^2 - x_2^2 \geq 0 \quad \equiv c_1$



# A Toy Example: Inequality Constraint

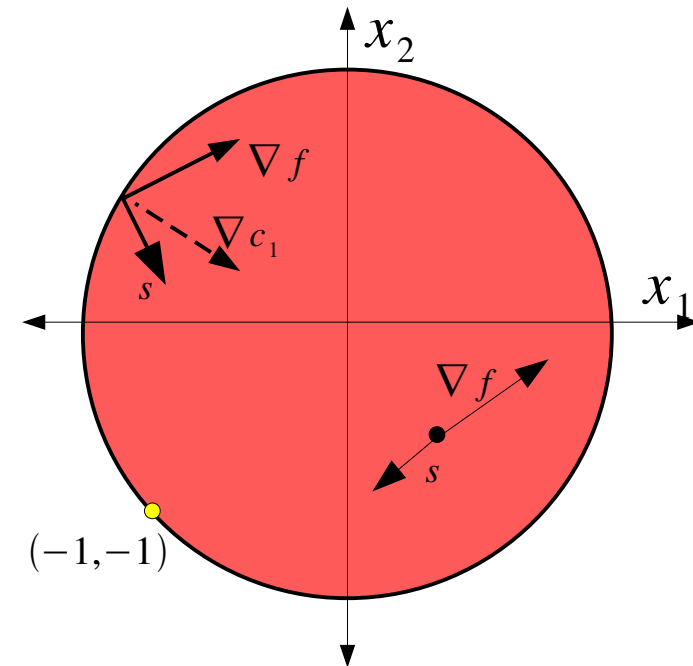
- $x$  is not an optimal solution, if there exists an  $s$  such that

$$\begin{aligned}c_1(x+s) &\geq 0 \\ f(x+s) &< f(x)\end{aligned}$$

- Using first order Taylor's expansion

$$c_1(x+s) = c_1(x) + \nabla c_1(x)^T s \geq 0 \quad (1)$$

$$f(x+s) - f(x) = \nabla f(x)^T s < 0 \quad (2)$$





# A Toy Example: Inequality Constraint

- **Case 1:** Inactive Constraint  $c_1(x) > 0$ 
  - Any sufficiently small  $s$  would do as long as

$$\nabla f_1(x) \neq 0$$

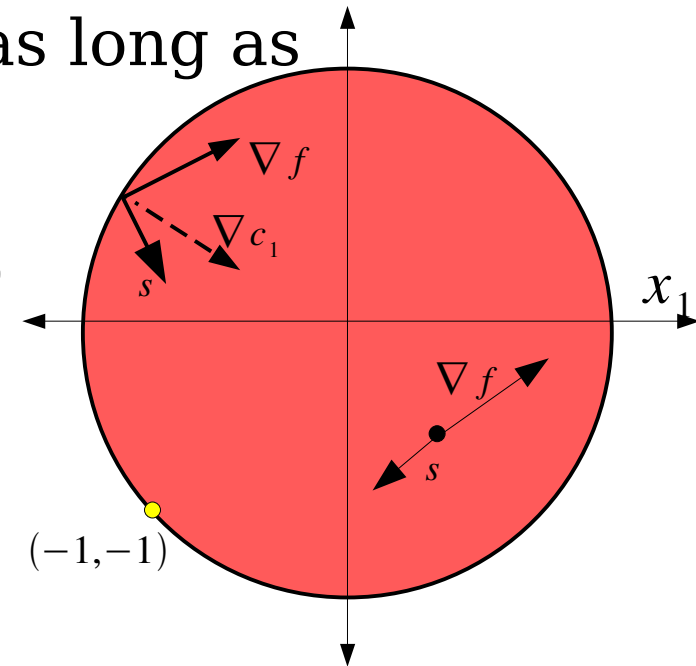
- Thus

$$s = -\alpha \nabla f(x) \quad \text{where } \alpha > 0$$

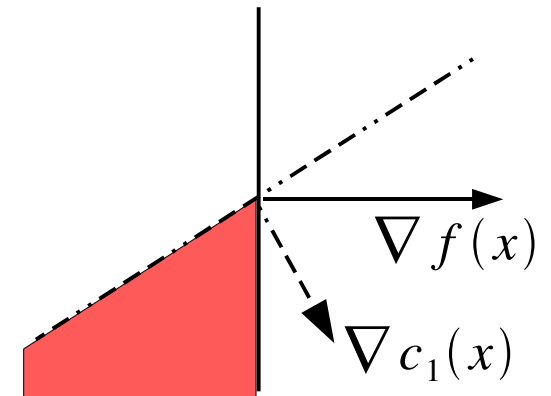
- **Case 2:** Active Constraint  $c_1(x) = 0$

$$\nabla c_1(x)^T s \geq 0 \quad (1)$$

$$\nabla f(x)^T s < 0 \quad (2)$$



$$\nabla f(x) = \lambda_1 \nabla c_1(x), \quad \text{where } \lambda_1 \geq 0$$



# A Toy Example: Inequality Constraint

- Thus we have the Lagrangian (as before)

$$L(x, \lambda_1) = f(x) - \lambda_1 c_1(x)$$

Lagrange multiplier  
or dual variable for  $c_1$

- The optimality conditions

$$\nabla_x L(x^o, \lambda_1^o) = \nabla f(x^o) - \lambda_1^o \nabla c_1(x^o) = 0 \quad \text{for some} \quad \lambda_1 \geq 0$$

and

$$\lambda_1^o c_1(x^o) = 0$$

Complementarity  
condition

# Same Concepts in a More General Setting

# The Lagrangian

- The Problem

$$\min_x f_0(x)$$

*s.t.:*

$$f_i(x) \leq 0, \quad i = 1, \dots, m$$

$$h_i(x) = 0, \quad i = 1, \dots, p$$

- The Lagrangian associated with the problem

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

dual variables or Lagrangian multipliers

# The Lagrange Dual Function

- Defined as the minimum value of the Lagrangian over  $x$

$$g: \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$$

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

# The Lagrange Dual Function

- Interpretation of Lagrange dual function:
  - Writing the original problem as unconstrained problem

$$\underset{x}{\text{minimize}} \left( f_0(x) + \sum_{i=1}^m I_0(f_i(x)) + \sum_{i=1}^p I_1(h_i(x)) \right)$$

where

$$I_0(u) = \begin{cases} 0 & u \leq 0 \\ \infty & u > 0 \end{cases}$$

$$I_1(u) = \begin{cases} 0 & u = 0 \\ \infty & u \neq 0 \end{cases}$$

indicator functions

# The Lagrange Dual Function

- Interpretation of Lagrange dual function:
  - The Lagrange multipliers in Lagrange dual function can be seen as “softer” version of indicator (**penalty**) function.

$$\text{minimize} \left( f_0(x) + \sum_{i=1}^m I_0(f_i(x)) + \sum_{i=1}^p I_1(h_i(x)) \right)$$

$$\inf_{x \in D} \left( f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x) \right)$$

# The Lagrange Dual Function

- Lagrange dual function gives a lower bound on optimal value of the problem.

$$g(\lambda, \nu) \leq p^o$$

- **Proof:** Let  $\hat{x}$  be a feasible point and let  $\lambda \geq 0$ .  
Then we have:

$$f_i(\hat{x}) \leq 0 \quad i=1, \dots, m$$

$$h_i(\hat{x}) = 0 \quad i=1, \dots, p$$



# The Lagrange Dual Function

- Lagrange dual function gives a lower bound on optimal value of the problem.

$$g(\lambda, \nu) \leq p^o$$

- **Proof:** Let  $\hat{x}$  be a feasible point and let  $\lambda \geq 0$ .  
Then we have:

$$\begin{aligned} f_i(\hat{x}) &\leq 0 & i=1, \dots, m \\ h_i(\hat{x}) &= 0 & i=1, \dots, p \end{aligned}$$

- Thus

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

# The Lagrange Dual Function

- Lagrange dual function gives a lower bound on optimal value of the problem.

$$g(\lambda, \nu) \leq p^o$$

- **Proof:** Let  $\hat{x}$  be a feasible point and let  $\lambda \geq 0$ . Then we have:

$$f_i(\hat{x}) \leq 0 \quad i=1, \dots, m$$

$$h_i(\hat{x}) = 0 \quad i=1, \dots, p$$

$$\leq 0$$

- Thus

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

# The Lagrange Dual Function

- Lagrange dual function gives a lower bound on optimal value of the problem.

$$g(\lambda, \nu) \leq p^o$$

- **Proof:** Let  $\hat{x}$  be a feasible point and let  $\lambda \geq 0$ . Then we have:

$$f_i(\hat{x}) \leq 0 \quad i=1, \dots, m$$

$$h_i(\hat{x}) = 0 \quad i=1, \dots, p$$

- Thus

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

- Hence

$$g(\lambda, \nu) = \inf_{x \in D} L(x, \lambda, \nu) \leq L(\hat{x}, \lambda, \nu) \leq f_0(\hat{x})$$

# The Lagrange Dual Problem

- Lagrange dual function gives a lower bound on optimal value of the problem.
- It is natural to seek the “best” lower bound.

$$\begin{array}{ll} \text{maximize} & g(\lambda, \nu) \\ \text{s.t.} & \lambda \geq 0 \end{array}$$

- Dual feasibility:

$$(\lambda, \nu): \quad \lambda \geq 0, \quad g(\lambda, \nu) \geq -\infty$$

- The dual optimal value and solution:

$$d^o = g(\lambda^o, \nu^o)$$

- The Lagrange dual problem is convex even if the original problem is not.

# Primal / Dual Problems

- Primal problem:

$$\begin{aligned} & \min_x f_0(x) \\ & \text{s.t.} : \\ & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_i(x) = 0, \quad i = 1, \dots, p \end{aligned} \quad p^\circ$$

- Dual problem:

$$\begin{aligned} & \max_{\lambda, \nu} g(\lambda, \nu) \\ & \text{s.t.} : \quad \lambda \geq 0 \end{aligned} \quad d^\circ$$

# Weak Duality

- Weak duality theorem:

$$d^o \leq p^o$$

- Optimal duality gap:

$$p^o - d^o \geq 0$$

- This bound is sometimes used to get an estimate on the optimal value of the original problem that is difficult to solve.

# Strong Duality

- Strong Duality:

$$d^o = p^o$$

- Strong duality does not hold in general.
- Slater's Condition: If  $x \in \text{relint } D$ , that it is **strictly feasible**.

$$f_i(x) < 0 \quad \text{for } i=1, \dots, m$$

$$h_i(x) = 0 \quad \text{for } i=1, \dots, p$$

- Strong duality theorem: Strong duality holds if Slater's condition holds.
- It also implies that the dual optimal value is attained.

$$\exists(\lambda^o, \nu^o) \quad \text{with} \quad g(\lambda^o, \nu^o) = d^o = p^o$$

# Optimality Conditions: First Order

- **Complementary slackness:** If strong duality holds, then at optimality

$$\lambda_i^o f_i(x^o) = 0 \quad i=1, \dots, m$$

- **Proof:** We have

$$f_0(x^o) = g(\lambda^o, \nu^o)$$

$$= \inf_x \left( f_0(x) + \sum_{i=1}^m \lambda_i^o f_i(x) + \sum_{i=1}^p \nu_i^o h_i(x) \right)$$

$$\leq f_0(x^o) + \sum_{i=1}^m \lambda_i^o f_i(x^o) + \sum_{i=1}^p \nu_i^o h_i(x^o) \leftarrow \text{less than 0}$$

$$\leq f_0(x^o)$$

- The result follows



# Optimality Conditions: First Order

- Karush-Kuhn-Tucker (KKT) Conditions: If the strong duality holds, then at optimality

$$f_i(x^o) \leq 0, \quad i=1, \dots, m$$

$$h_i(x^o) = 0, \quad i=1, \dots, p$$

$$\lambda_i^o \geq 0, \quad i=1, \dots, m$$

$$\lambda_i^o f_i(x^o) = 0, \quad i=1, \dots, m$$

$$\nabla f_0(x^o) + \sum_{i=1}^m \lambda_i^o \nabla f_i(x^o) + \sum_{i=1}^p \nu_i^o \nabla h_i(x^o) = 0$$

- KKT conditions are necessary in general and necessary and sufficient in case of convex problems.