

# Support Vector Machines: Maximum Margin Classifiers

Machine Learning and Pattern Recognition:  
September 23, 2010

Piotr Mirowski

Based on slides by Sumit Chopra, Fu-Jie Huang and Mehryar Mohri

# Outline

- What is behind Support Vector Machines?
  - Constrained optimization
  - Lagrange constraints
  - “Dual” solution
- Support Vector Machines in detail
  - Kernel trick
  - LibSVM demo

# Binary Classification Problem

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

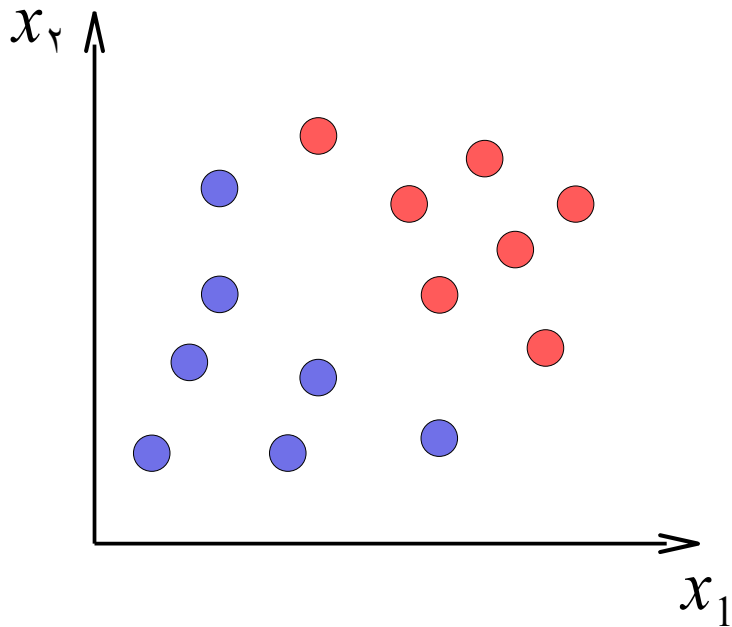
$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \mathcal{R}^n \times \{-1, 1\}$$

input label                      input label  
space space

- **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  
(e.g. logistic regression, Perceptron)
  - **Non Linear Approach:** non-linear classifiers  
(e.g. Neural Networks, SVM)

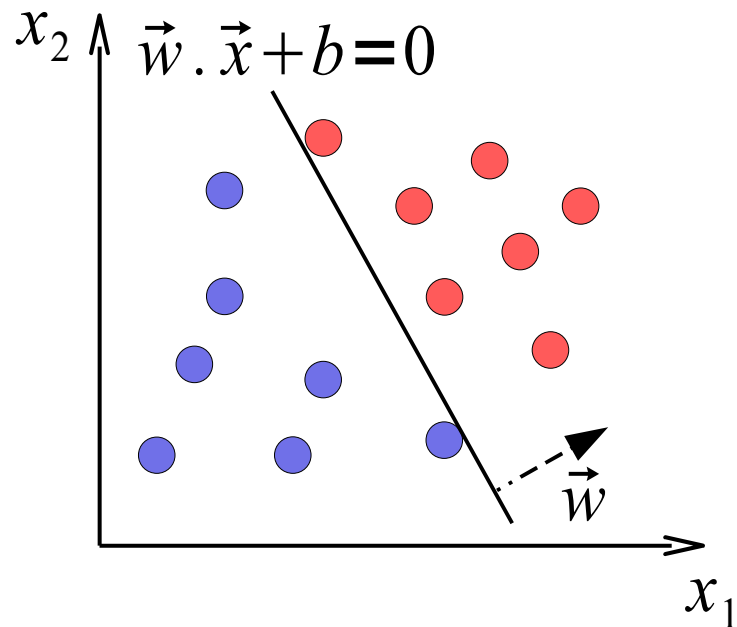
# Linearly Separable Data

- Assume that the training data is linearly separable



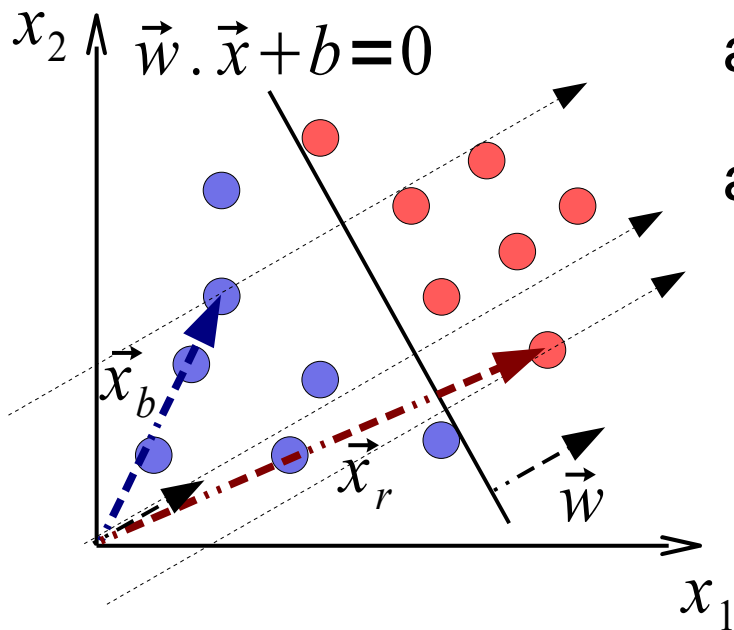
# Linearly Separable Data

- Assume that the training data is linearly separable



# Linearly Separable Data

- Assume that the training data is linearly separable



abscissa on axis parallel to  $\vec{w}$

abscissa of origin  $O$  is  $b$

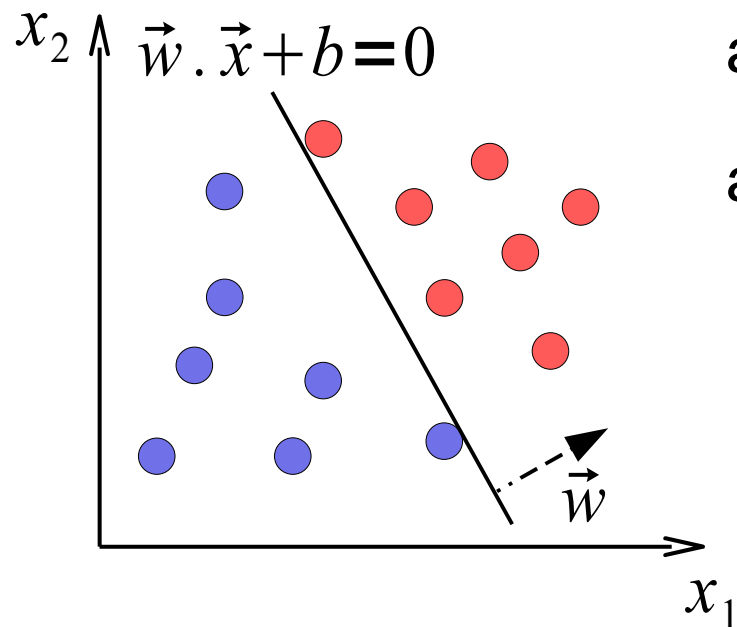
$$\vec{w} \cdot \vec{x}_b + b = \tilde{y}_b$$

$$\vec{w} \cdot \vec{x}_r + b = \tilde{y}_a$$

$$\vec{w} \cdot \vec{O} + b = b$$

# Linearly Separable Data

- Assume that the training data is linearly separable



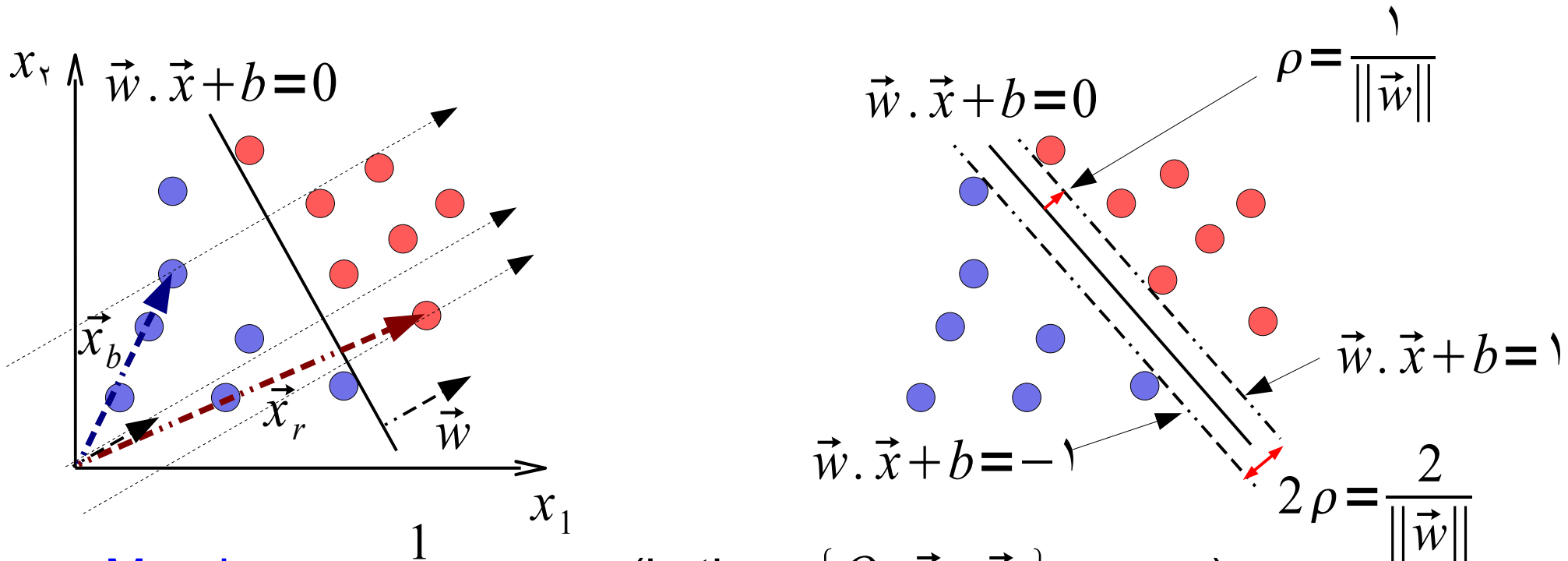
abscissa on axis parallel to  $\vec{w}$

abscissa of origin  $0$  is  $b$

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

# Linearly Separable Data

- Assume that the training data is linearly separable



- Margin**  $\rho = \frac{1}{\|\vec{w}\|}$  (in the  $\{O, \vec{x}_1, \vec{x}_2\}$  space)

- Maximize margin  $\rho$  (or  $2\rho$ ) so that:

For the closest points:  $h(x) = \vec{w} \cdot \vec{x} + b \in \{-1, 1\}$



# Optimization Problem

- A Constrained Optimization Problem

$$\min_w \frac{1}{2} \|\mathbf{w}\|^2$$

s.t.:

$$y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, m$$

label    input

- Equivalent to maximizing the margin  $\rho = \frac{1}{\|\mathbf{w}\|}$
- A convex optimization problem:
  - Objective is convex
  - Constraints are affine hence convex
  - Therefore, admits a unique optimum at  $\mathbf{w}_0$

# Optimization Problem

- Compare:

$$\min_w \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{objective}$$

*s.t.:*

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i=1, \dots, m$$

constraints

- With:

$$\min_w \left( \sum_{i=1}^m (-y_i(\mathbf{w} \cdot \mathbf{x}_i + b)) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right)$$

energy/errors regularization

# 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 (unique) optimum – if the problem is convex
- Local optimum – if the problem is not convex

(notation change: the parameters to optimize are noted  $x$ )

# Optimization: Some Theory

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

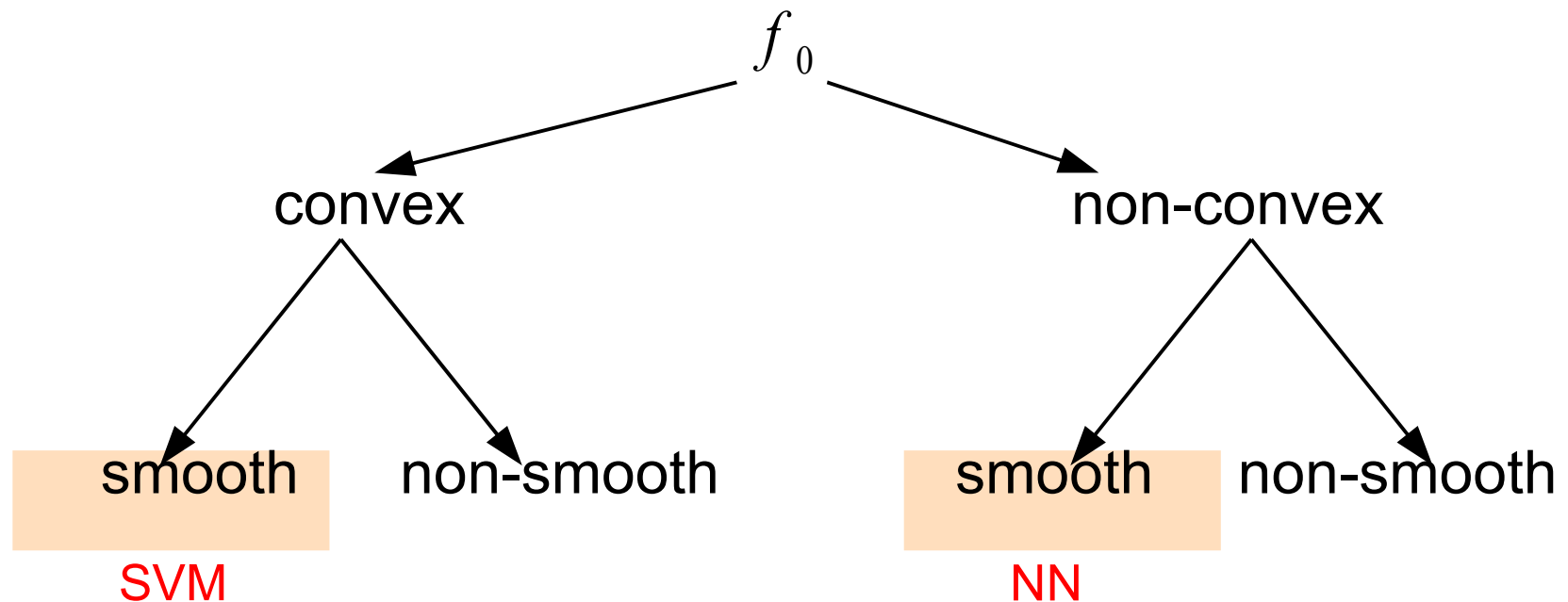
$$\begin{aligned} \min_x \quad & c^T x \\ \text{s.t.} \quad & \\ & Ax = b \\ & x \geq 0 \end{aligned}$$

- **Example:** Least Squares Solution of Linear Equations  
(with  $L_2$  norm regularization of the solution  $x$ )  
i.e. Ridge Regression

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

# Big Picture

- Constrained / unconstrained optimization
- Hierarchy of objective function:
  - smooth = infinitely derivable
  - convex = has a global optimum

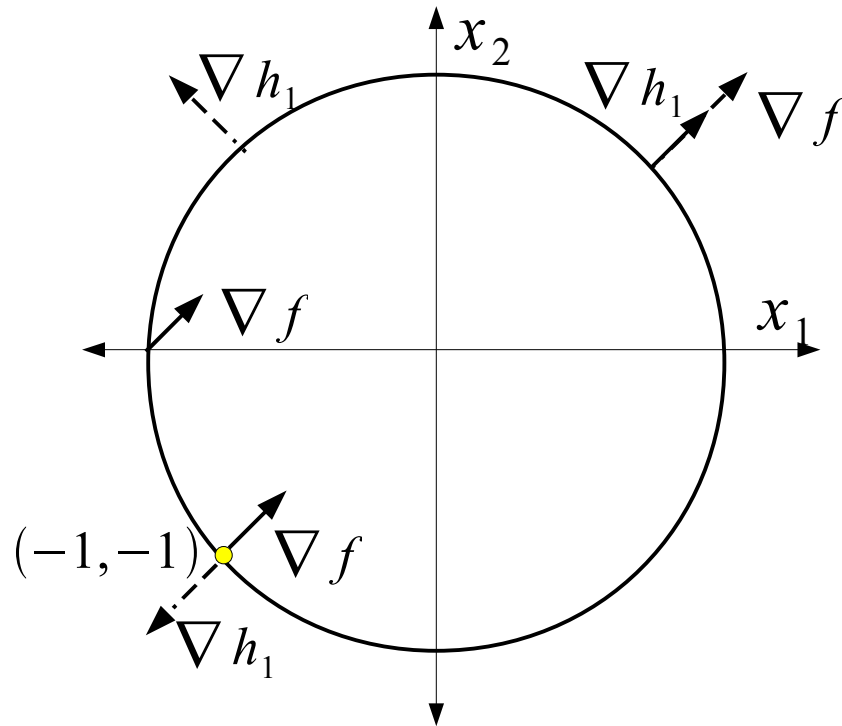


# Introducing the concept of Lagrange function on a toy example

# Toy Example: Equality Constraint

- Example 1:

$$\begin{aligned} \min \quad & x_1 + x_2 \quad \equiv f \\ \text{s.t.} \quad & x_1^2 + x_2^2 - 2 = 0 \quad \equiv h_1 \end{aligned}$$



$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{pmatrix}$$

$$\nabla h_1 = \begin{pmatrix} \frac{\partial h_1}{\partial x_1} \\ \frac{\partial h_1}{\partial x_2} \end{pmatrix}$$

- At Optimal Solution:

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

# Toy Example: Equality Constraint

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

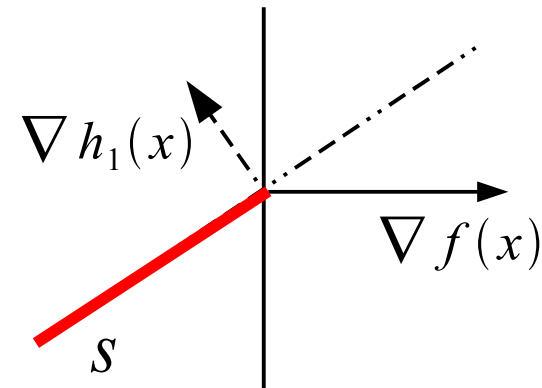
$$\begin{aligned}h_1(x+s) &= 0 \\f(x+s) &< f(x)\end{aligned}$$

- Using first order Taylor's expansion

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

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

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





# 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

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

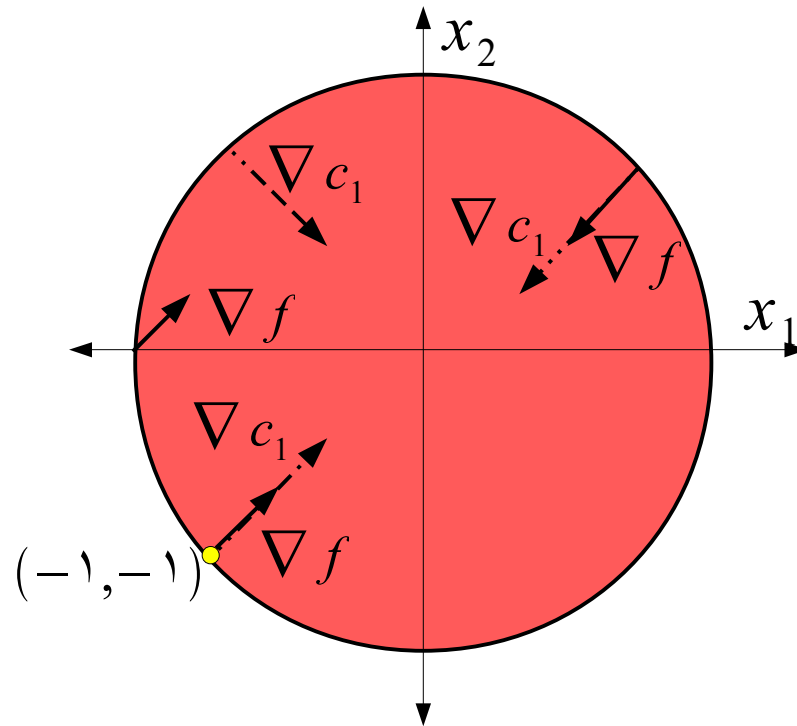
- This is just a necessary (not a sufficient) condition”  
 $x$  solution implies  $\nabla h_1(x) \parallel \nabla f(x)$

# Toy Example: Inequality Constraint

• Example 2:

$$\min x_1 + x_2 \quad \equiv f$$

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



$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{pmatrix}$$

$$\nabla c_1 = \begin{pmatrix} \frac{\partial c_1}{\partial x_1} \\ \frac{\partial c_1}{\partial x_2} \end{pmatrix}$$

# Toy Example: Inequality Constraint

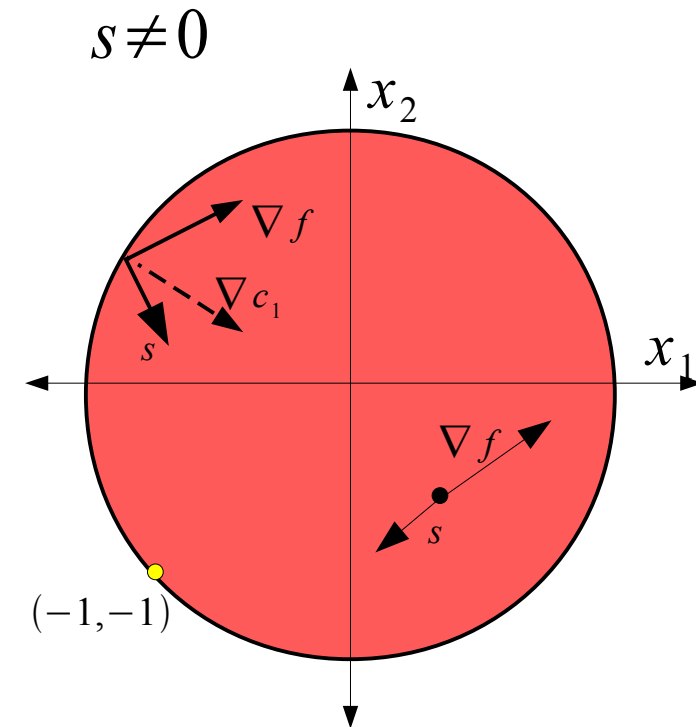
- $x$  is not an optimal solution, if there exists such that

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

- Using first order Taylor's expansion

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

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



# Toy Example: Inequality Constraint

- **Case 1: Inactive constraint**  $c_1(x) > 0$ 
  - Any sufficiently small  $s$  as long as  $\nabla f_1(x) \neq 0$

→ Thus  $s = -\alpha \nabla f(x)$  where  $\alpha > 0$

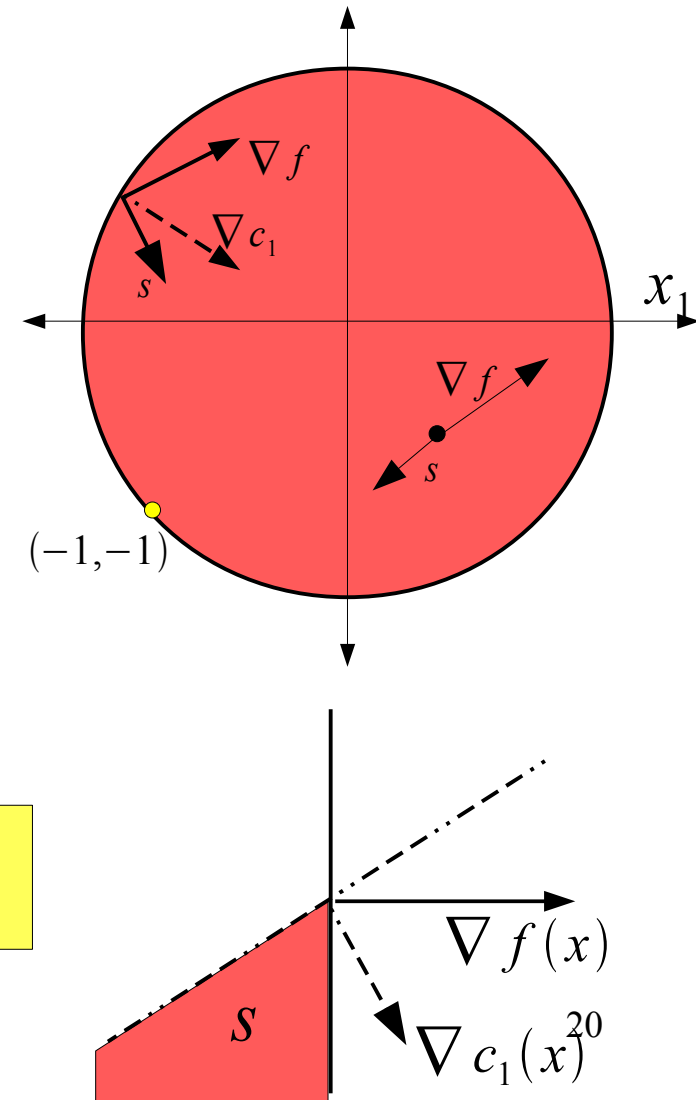
- **Case 2: Active constraint**  $c_1(x) = 0$

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

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

In that case,  $s = 0$  when:

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



# Toy Example: Inequality Constraint

- Thus we have the Lagrange function (as before)

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

Lagrange multiplier or  
dual variable for  $c_1$

- The optimality conditions

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

and

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

Complementarity  
condition

either  $c_1(x^o) = 0$  or  $\lambda_1^o = 0$   
(active) (inactive)

# Same Concepts in a More General Setting

# Lagrange Function

- The Problem

$$\min_x f_0(x)$$

objective function

s.t.:

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

$m$  inequality constraints

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

$p$  equality constraints

- Standard tool for constrained optimization:  
the Lagrange Function

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

dual variables or Lagrange multipliers

# Lagrange Dual Function

- Defined, for  $\lambda, \nu$  as the minimum value of the Lagrange function over  $x$

$m$  inequality constraints

$p$  equality constraints

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

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



# Lagrange Dual Function

- Interpretation of Lagrange dual function:
  - Writing the original problem as unconstrained problem but with hard indicators (penalties)

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

where

$$I_0(u) = \begin{cases} 0 & u \leq 0 \\ \infty & u > 0 \end{cases} \quad \begin{array}{l} \text{satisfied} \\ \text{unsatisfied} \end{array}$$
$$I_1(u) = \begin{cases} 0 & u = \cdot \\ \infty & u \neq \cdot \end{cases} \quad \begin{array}{l} \text{satisfied} \\ \text{unsatisfied} \end{array}$$

indicator functions

# Lagrange Dual Function

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

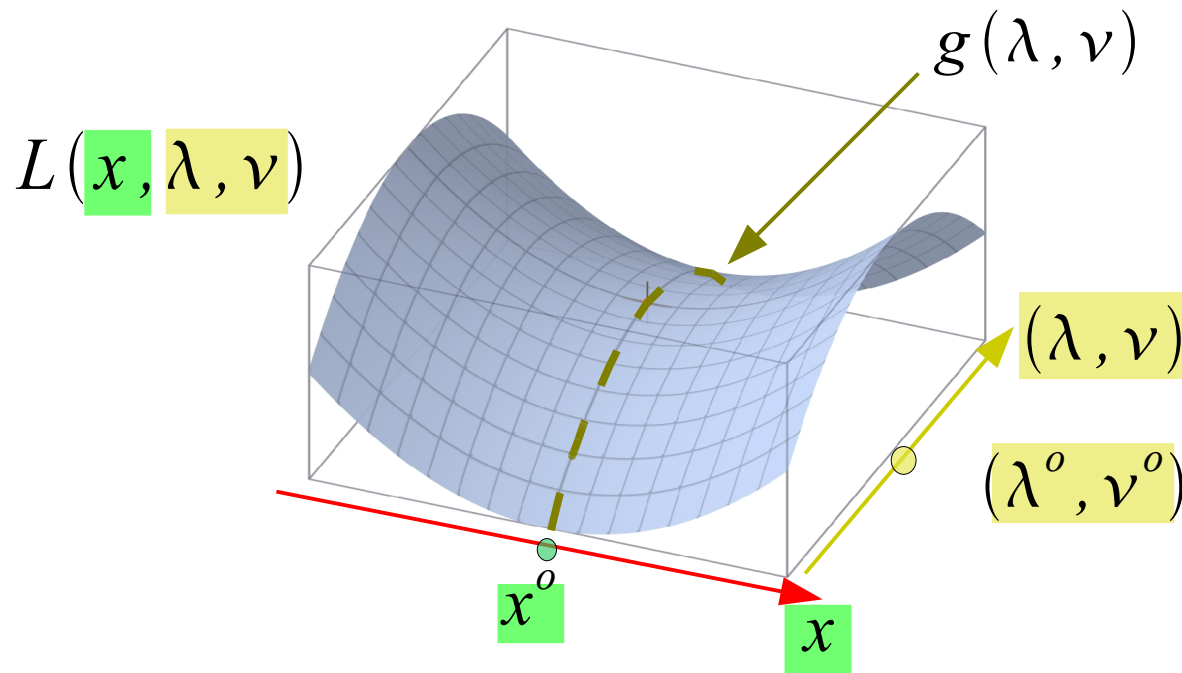
$$\underset{x}{\text{minimize}} \left( f_0(x) + \sum_{i=1}^m I_0(f_i(x)) + \sum_{i=1}^p I_1(h_i(x)) \right)$$

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

# Sufficient Condition

- If  $(x^o, \lambda^o, \nu^o)$  is a saddle point, i.e. if

$$\forall x \in \mathbb{R}^n, \quad \forall \lambda \geq 0, \quad L(x^o, \lambda, \nu) \leq L(x^o, \lambda^o, \nu^o) \leq L(x, \lambda^o, \nu^o)$$



- ... then  $(x^o, \lambda^o, \nu^o)$  is a solution of the primal problem  $p^o$

# Lagrange Dual Problem

- Lagrange dual function gives a **lower bound** on the **optimal value** of the problem.
- We seek the “**best**” **lower bound** to minimize the objective:

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

- 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:

$$p^o \quad \begin{array}{l} \min_{x \in D} f_0(x) \\ \text{s.t.}: \\ f_i(x) \leq 0, \quad i=1, \dots, m \\ h_i(x) = 0, \quad i=1, \dots, p \end{array}$$

- Dual problem:

$$d^o \quad \begin{array}{l} \max_{\lambda, \nu} g(\lambda, \nu) \\ \text{s.t.}: \quad \lambda \geq 0 \\ g(\lambda, \nu) = \inf_{x \in D} \left( f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x) \right) \end{array}$$

# 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 (local optimum)

- necessary and sufficient in case of convex problems (global optimum)