

Support Vector Machines: Maximum Margin Classifiers

Machine Learning and Pattern Recognition:
September 16, 2008

Piotr Mirowski

Based on slides by Sumit Chopra and Fu-Jie Huang

Outline

- What is behind Support Vector Machines?
 - Constrained Optimization
 - Lagrange Duality
- Support Vector Machines in Detail
 - Kernel Trick
 - LibSVM demo

Binary Classification Problem

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

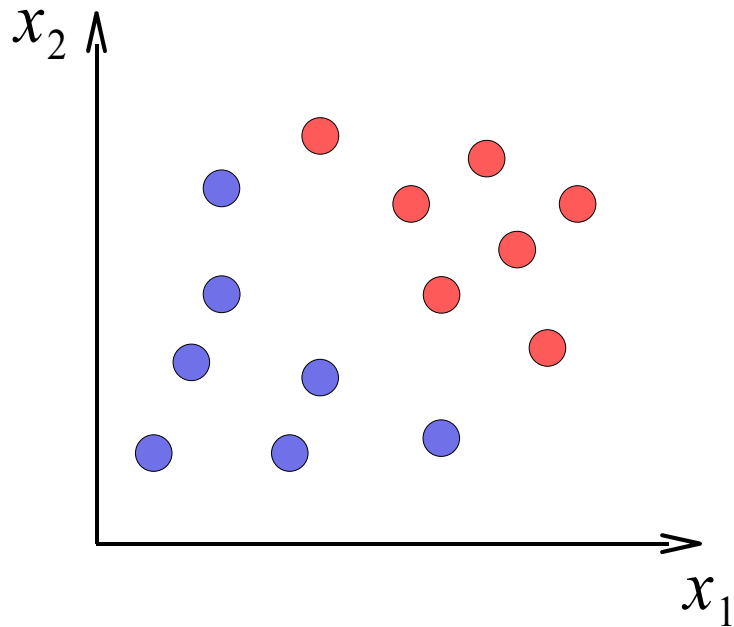
$$(x_1, y_1), \dots, (x_p, y_p) \in \mathcal{R}^n \times \{-1, 1\}$$

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

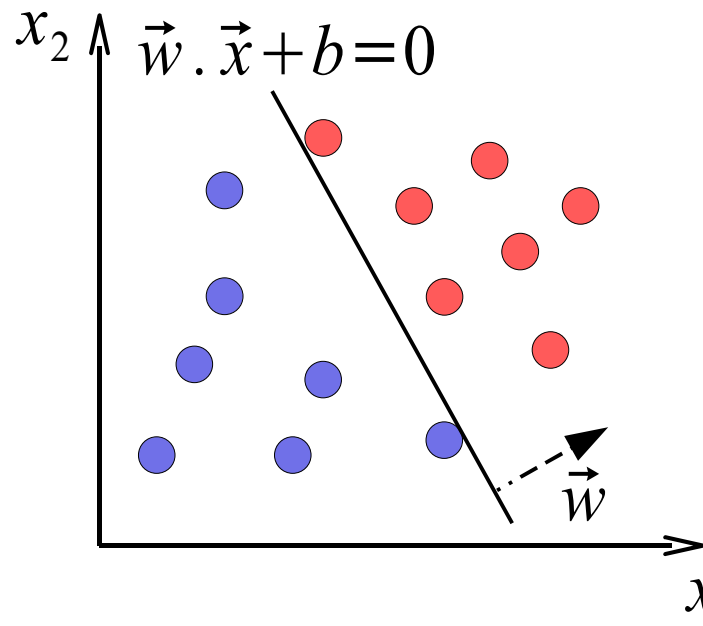
Linearly Separable Data

- Assume that the training data is linearly separable



Linearly Separable Data

- Assume that the training data is linearly separable



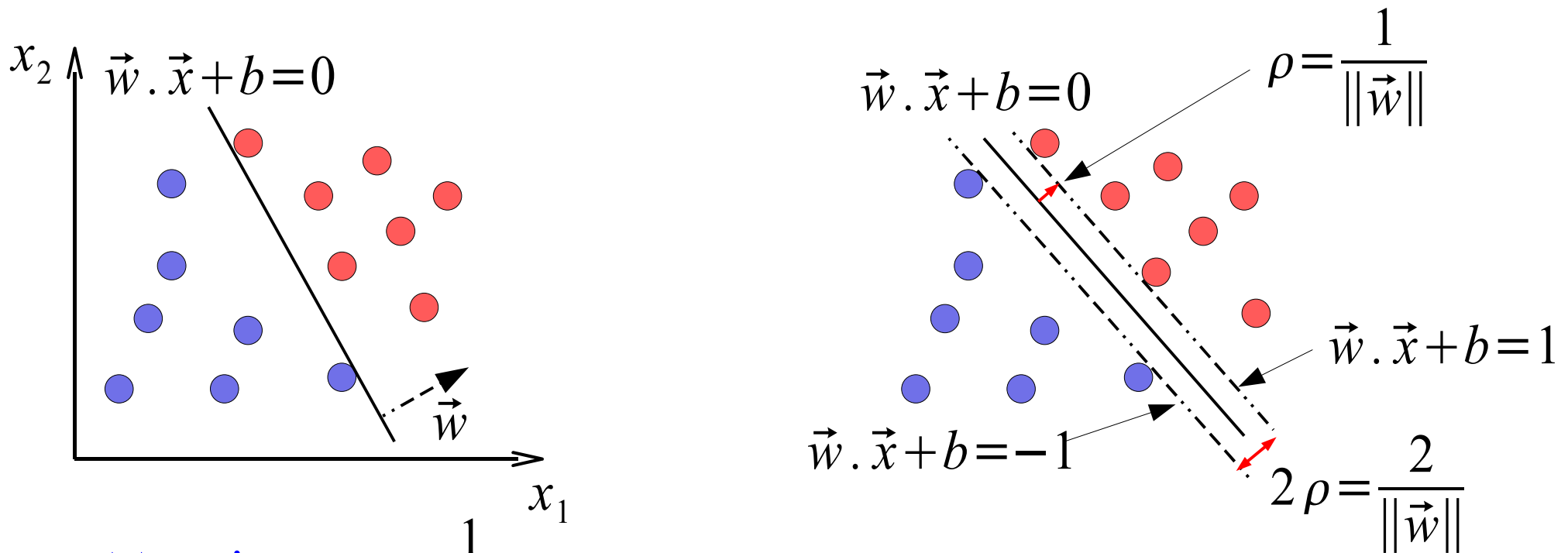
abscissa on axis parallel to \vec{w}

abscissa of origin 0 is b

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

Linearly Separable Data

- Assume that the training data is linearly separable



- Margin $\rho = \frac{1}{\|\vec{w}\|}$

- Maximize margin ρ (or 2ρ) so that:

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

Optimization Problem

- A Constrained Optimization Problem

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

s.t.:

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

label input

- Equivalent to maximizing the margin $\rho = \frac{1}{\|\vec{w}\|}$
- A convex optimization problem:
 - Objective is convex
 - Constraints are affine hence convex

Optimization Problem

- Compare:

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

s.t. :

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

constraints

- With:

$$\min_w \left(\sum_{i=1}^p \underbrace{(-y_i (\vec{w} \cdot \vec{x}_i + b))}_{\text{energy/errors}} + \underbrace{\frac{\lambda}{2} \|\vec{w}\|^2}_{\text{regularization}} \right)$$

Optimization: Some Theory

- The problem:

$$\min_x f_0(x) \quad \leftarrow \text{objective function}$$

s.t.:

$$f_i(x) \leq 0, \quad i = 1, \dots, m \quad \leftarrow \text{inequality constraints}$$

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

- Solution of problem: x^o

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

Optimization: Some Theory

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

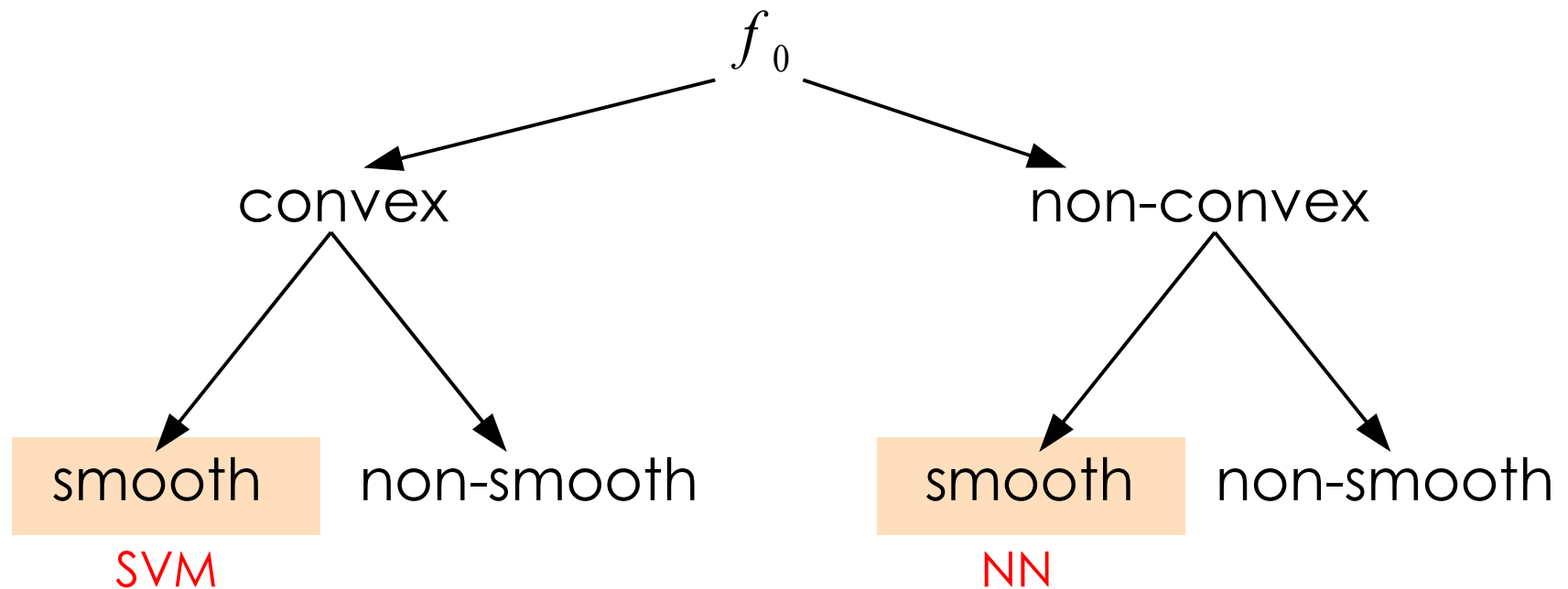
$$\begin{aligned} \min_x \quad & c^T x \\ \text{s.t.} \quad & \\ & Ax = b \\ & x \geq 0 \end{aligned}$$

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

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

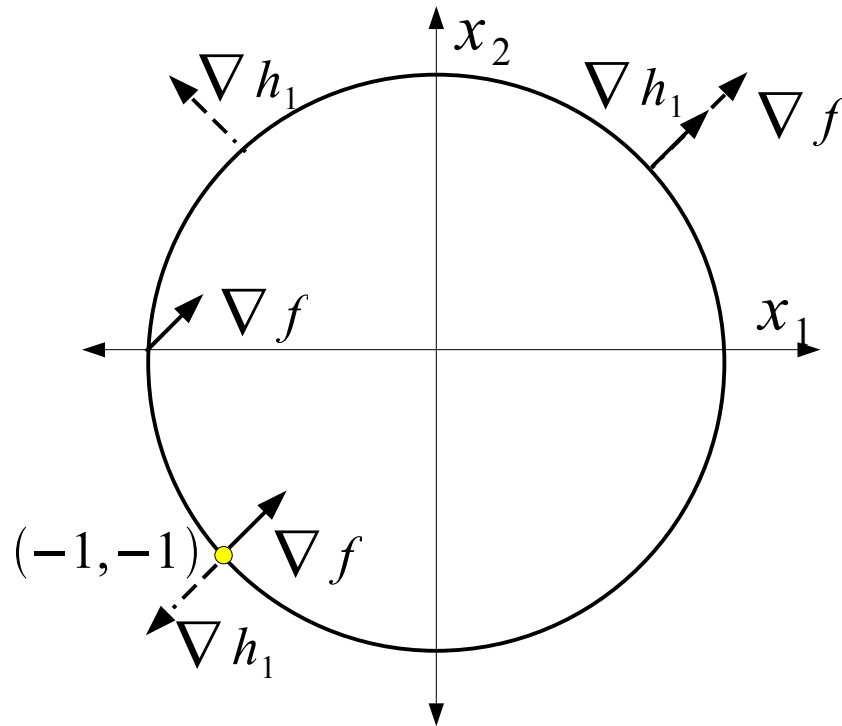
Big Picture

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



Toy Example: Equality Constraint

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



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

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

- 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

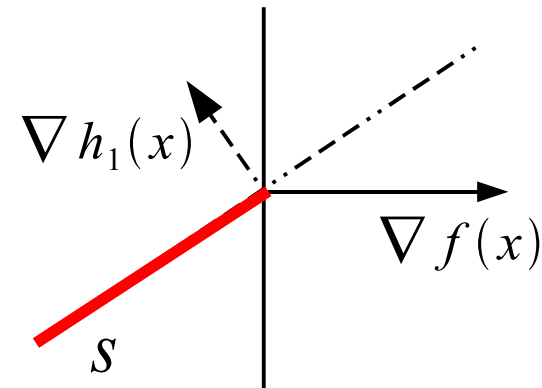
$$\begin{aligned}h_1(x+s) &= 0 \\f(x+s) &< f(x)\end{aligned}$$

- 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

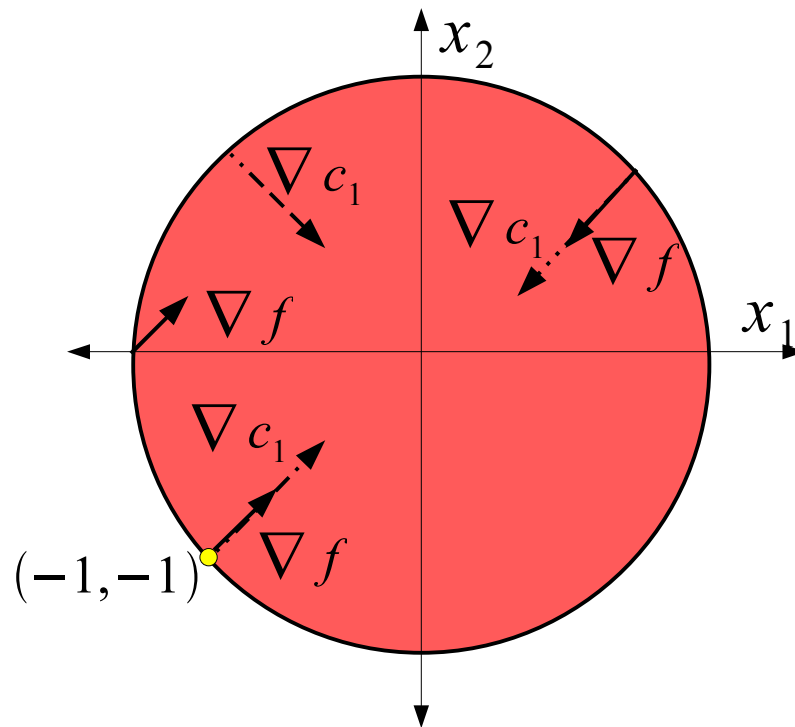
- Thus at the solution

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

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



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

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

Toy Example: Inequality Constraint

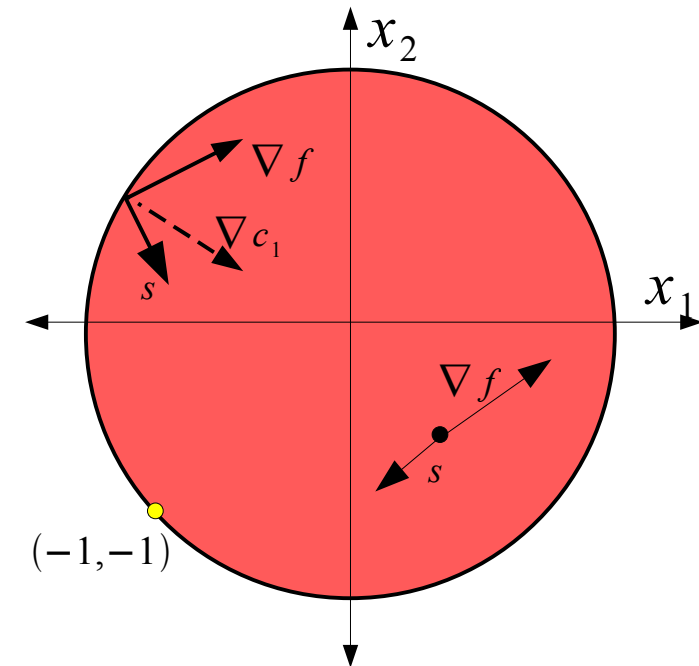
- x is not an optimal solution, if there exists $s \neq 0$ such that

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

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

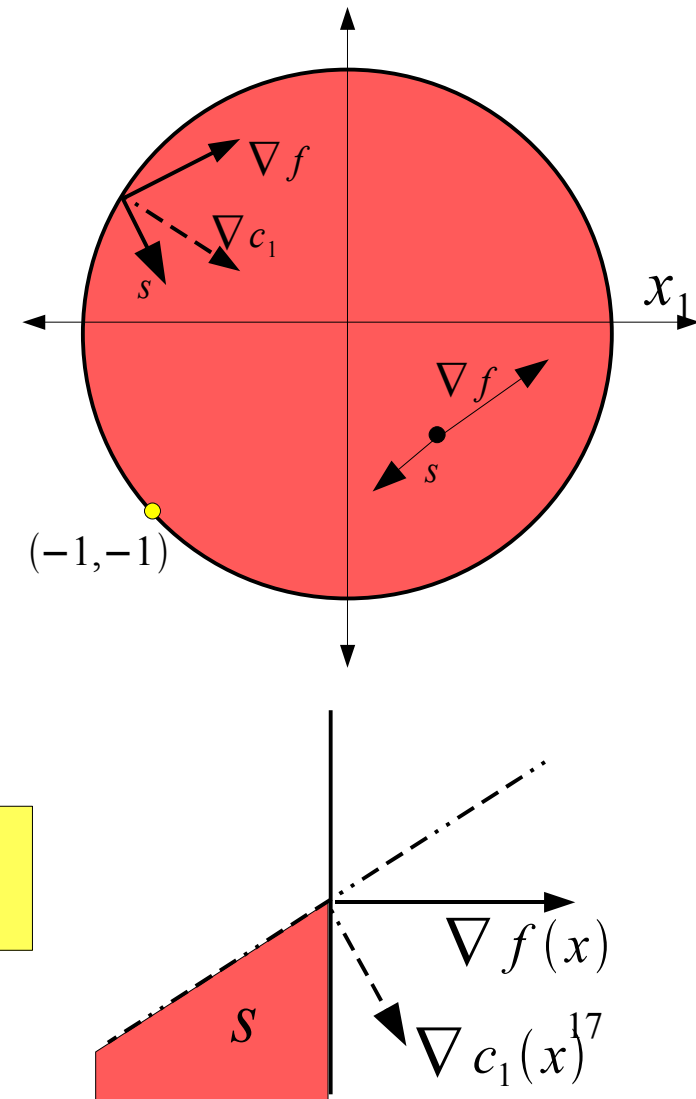
→ Thus $s = -\alpha \nabla f(x)$ where $\alpha > 0$

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

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

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

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



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

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 λ, ν 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

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

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.

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

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

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 optimal point and let $\lambda \geq 0$. Then we have:

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

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 optimal 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

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

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 optimal 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

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

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 optimal point and let $\lambda \geq 0$. Then we have:

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

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

$$\cdot \leq 0$$

- Thus

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

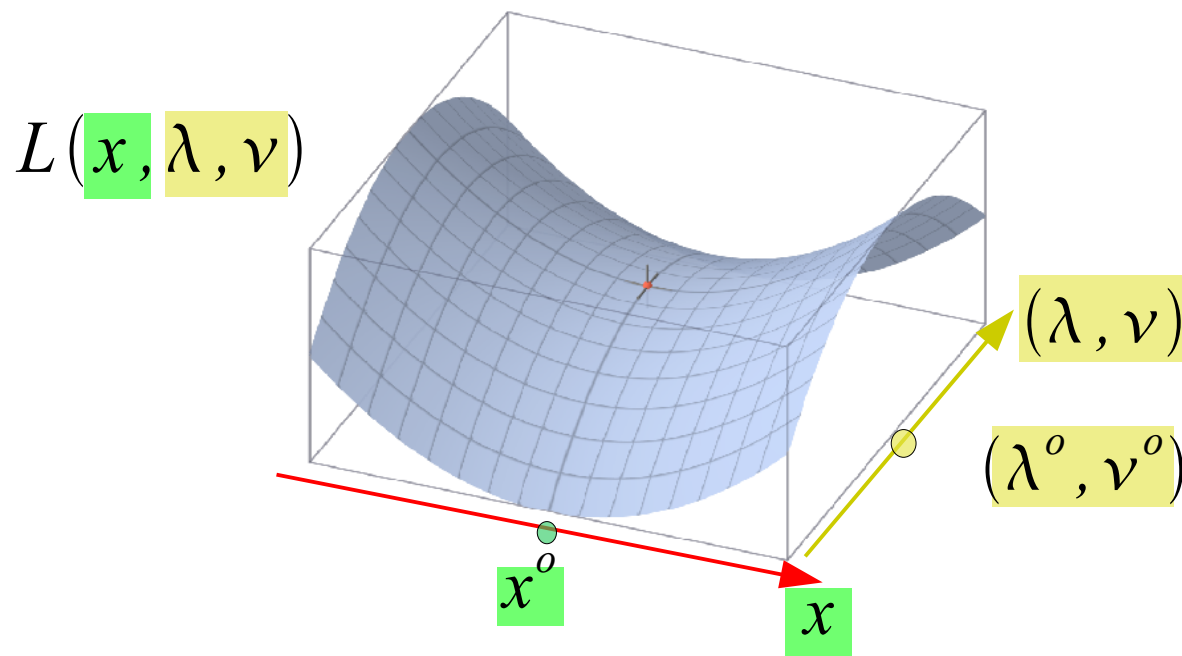
- Hence

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

Sufficient Condition

- If (x^o, λ^o, ν^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, λ^o, ν^o) is a solution of 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:

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

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 $x \in D$ and it is **strictly feasible**:

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

- Strong Duality theorem:

if Slater's condition holds and the problem is **convex**, then **strong duality** is attained:

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

Optimality Conditions: First Order

- Complementary slackness:
if **strong duality** holds, then at optimality (x^o, λ^o, ν^o)

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

- Proof:

$$f_0(x^o) = g(\lambda^o, \nu^o) \quad (\text{strong duality})$$

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

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

$$\leq f_0(x^o) \quad \forall i, h_i(x) = 0$$

$$\forall i, f_i(x) \leq 0, \lambda_i \geq 0$$

Optimality Conditions: First Order

- Karush-Kuhn-Tucker (KKT) conditions

If the strong duality holds, then at optimality:

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

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

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

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

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

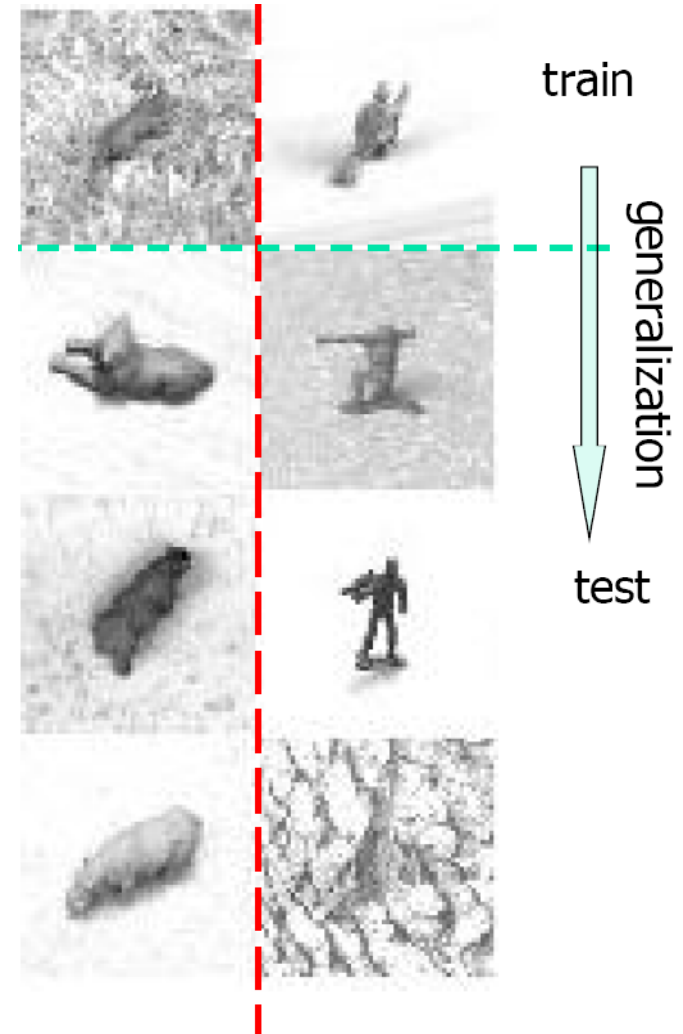
- KKT conditions are

- necessary in general (local optimum)

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

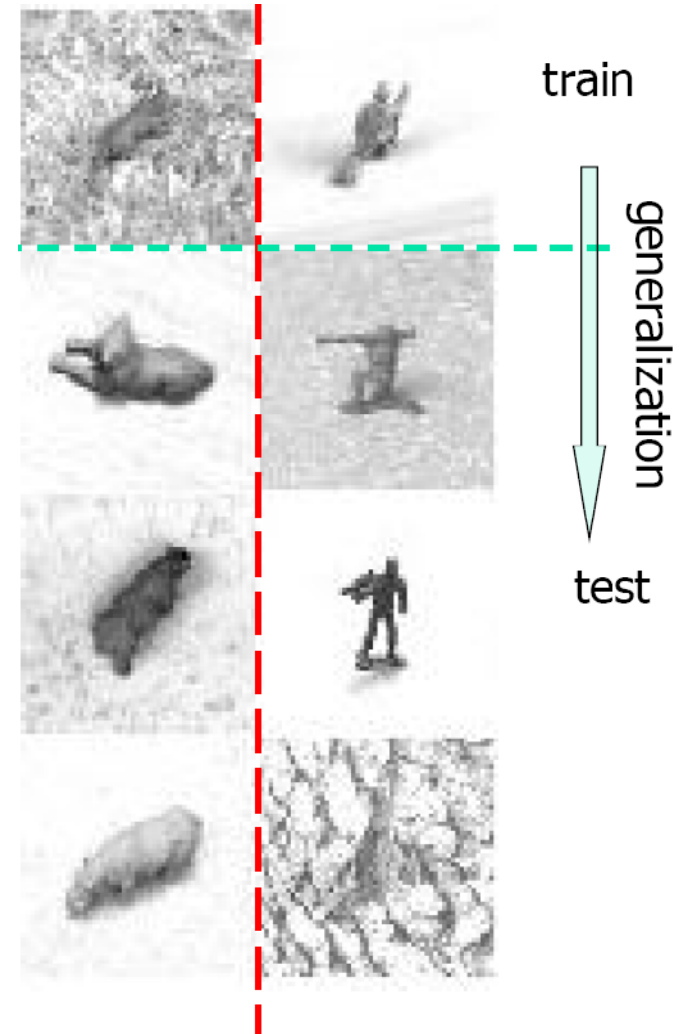
What are Support Vector Machines?

- Linear classifiers
- (Mostly) binary classifiers
- Supervised training
- Good generalization with explicit bounds



Main Ideas Behind Support Vector Machines

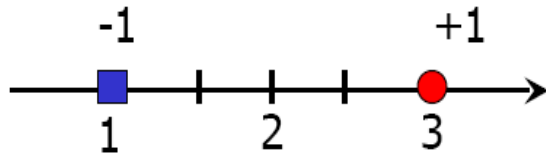
- Maximal margin
- Dual space
- Linear classifiers in high-dimensional space using non-linear mapping
- Kernel trick



Quadratic Programming

$$\max_{w,b} \min_i \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|} \xrightarrow{\min_i |\mathbf{w}^T \mathbf{x}_i + b| = 1} \min_{w,b} \frac{1}{2} \langle \mathbf{w}^T \cdot \mathbf{w} \rangle$$

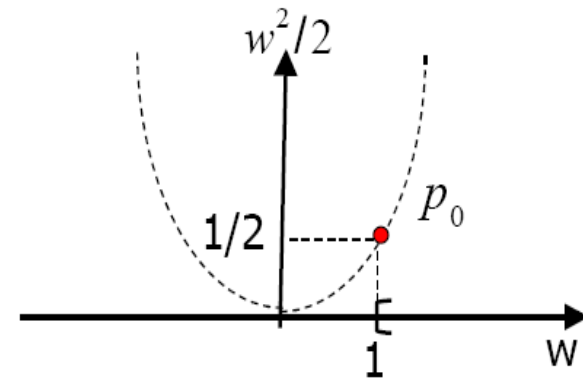
$$y_i (\langle \mathbf{w}^T \cdot \mathbf{x}_i \rangle + b) \geq 1$$



$$\min_w \frac{w^2}{2}$$

$$(+1)(w \cdot 3 + b) \geq 1$$

$$(-1)(w \cdot 1 + b) \geq 1$$



Using the Lagrangian

- Combine target and constraints
- Minimize over primal
- Maximize over dual

$$L(x, \boldsymbol{\lambda}) = f_0(x) - \sum \lambda_i f_i(x)$$

$$Q(\boldsymbol{\lambda}) = \min_x L(x, \boldsymbol{\lambda})$$

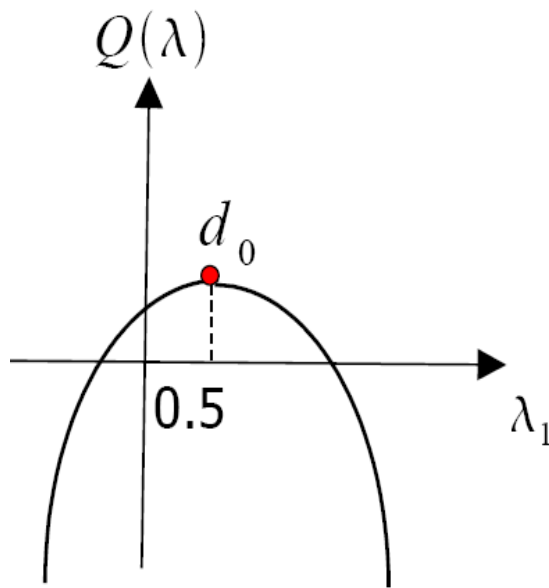
$$\max_{\boldsymbol{\lambda}} Q(\boldsymbol{\lambda}), \lambda > 0$$

Dual Space

$$\min_w \frac{w^2}{2}$$

$$(+1)(w \cdot 3 + b) \geq 1$$

$$(-1)(w \cdot 1 + b) \geq 1$$



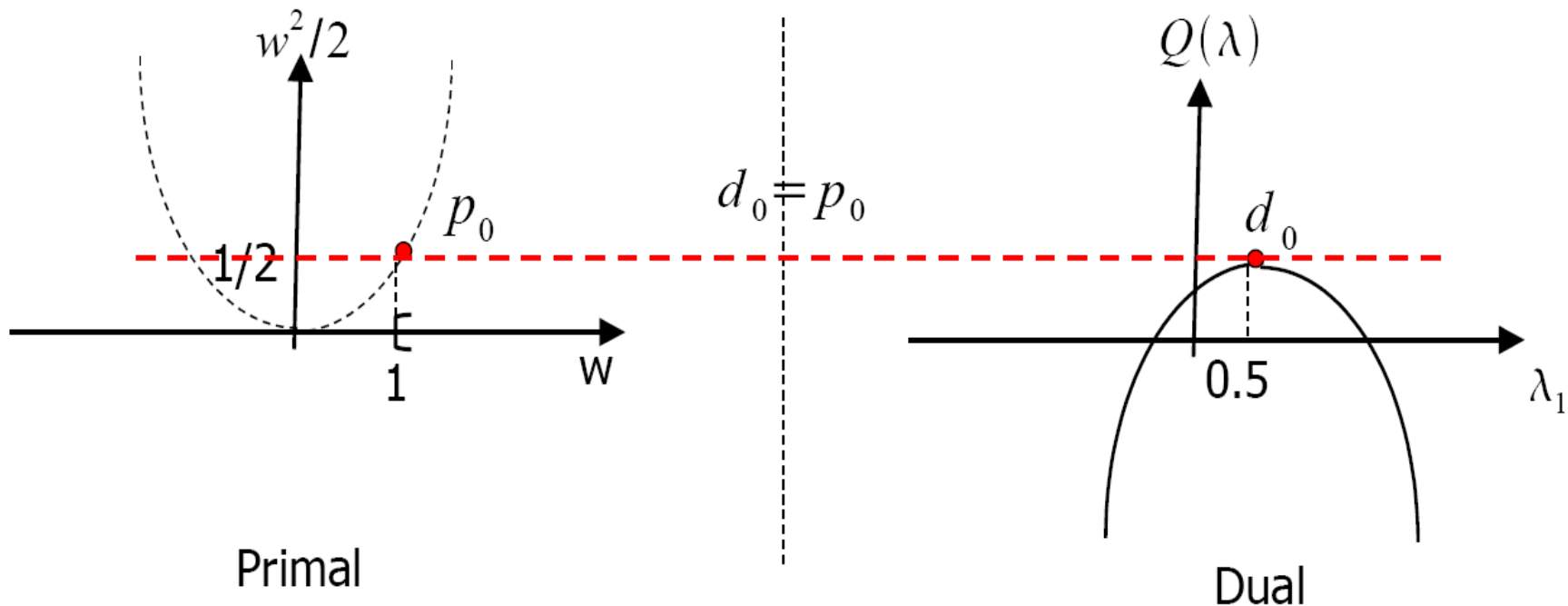
$$L(w, b, \boldsymbol{\lambda}) = w^2/2 - \lambda_1(3w + b - 1) - \lambda_2(-w - b - 1)$$

$$\min_{w, b} L(w, b, \boldsymbol{\lambda}) \Rightarrow \begin{cases} \lambda_1 = \lambda_2 \\ w = 3\lambda_1 - \lambda_2 = 2\lambda_1 \\ Q(\boldsymbol{\lambda}) = Q(\lambda_1) = -2\lambda_1^2 + 2\lambda_1 \end{cases}$$

$$\max_{\boldsymbol{\lambda}} Q(\boldsymbol{\lambda}) \Rightarrow \lambda_1 = \lambda_2 = 1/2, w = 1, b = 2$$

Strong Duality

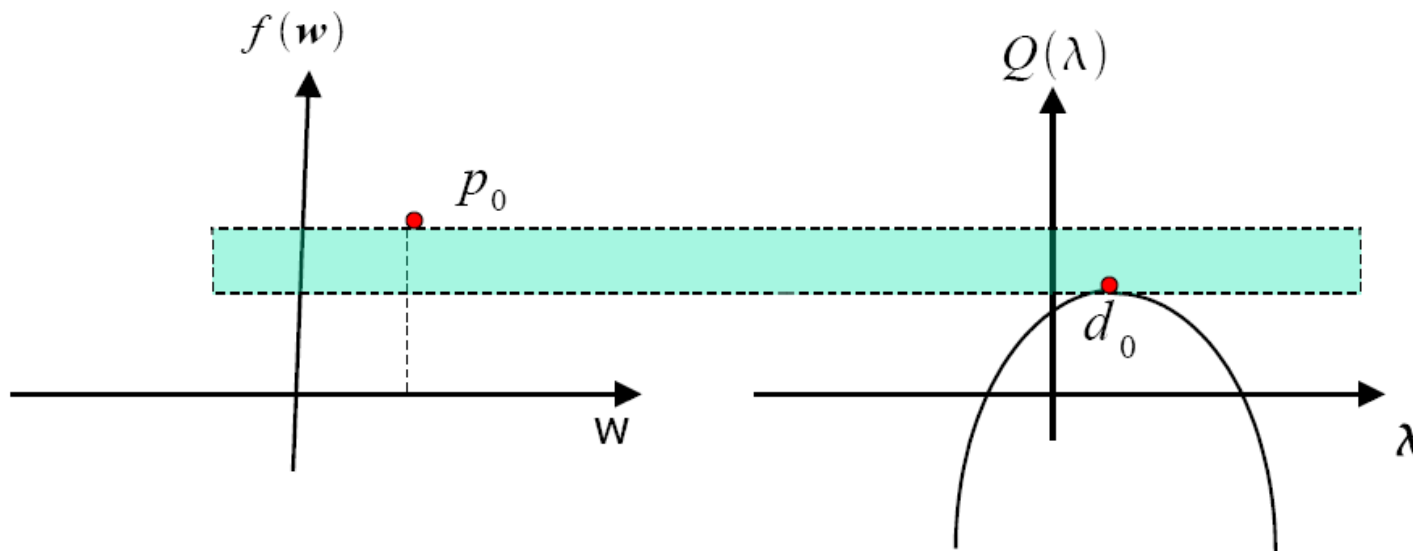
- Primal and dual space optimization:
 - Same result!



Duality Gap

$$d_0 < p_0$$

- In a general case
 - Strong duality is not true
 - “Step 1-2-3” a lower bound, not a solution



No Duality Gap Thanks to Convexity

- Convex function
 - Quadratic programming
- Convex set
 - Linear constraints
- No duality gap

$$\min_{w, b} \frac{1}{2} \langle \mathbf{w}^T \cdot \mathbf{w} \rangle$$

$$y_i (\langle \mathbf{w}^T \cdot \mathbf{x}_i \rangle + b) \geq 1$$

$$d_0 = p_0$$

Dual Form

- H
 - Hessian matrix
 - Gram matrix
- Lambda
 - Support vector
 - Sparse

$$\max_{\lambda} Q(\lambda) = -0.5 \lambda^T H \lambda + f^T \lambda$$

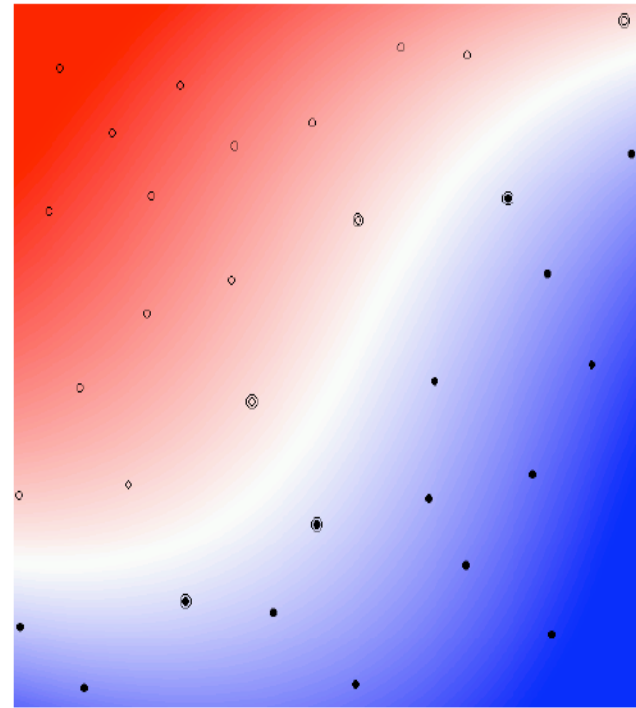
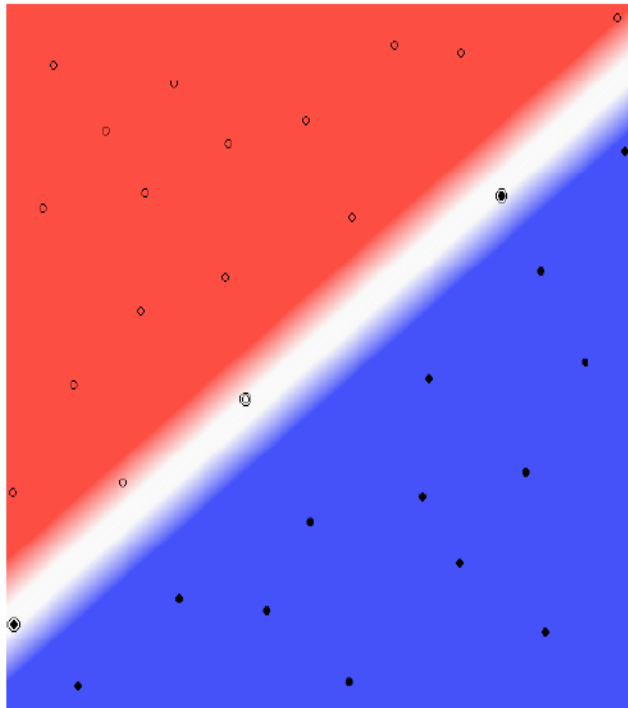
$$y^T \lambda = 0$$

$$\lambda \geq 0$$

$$\text{where, } H_{ij} = y_i y_j \langle \mathbf{x}_i^T \cdot \mathbf{x}_j \rangle$$

f is a unit vector

Non-linear separation of datasets



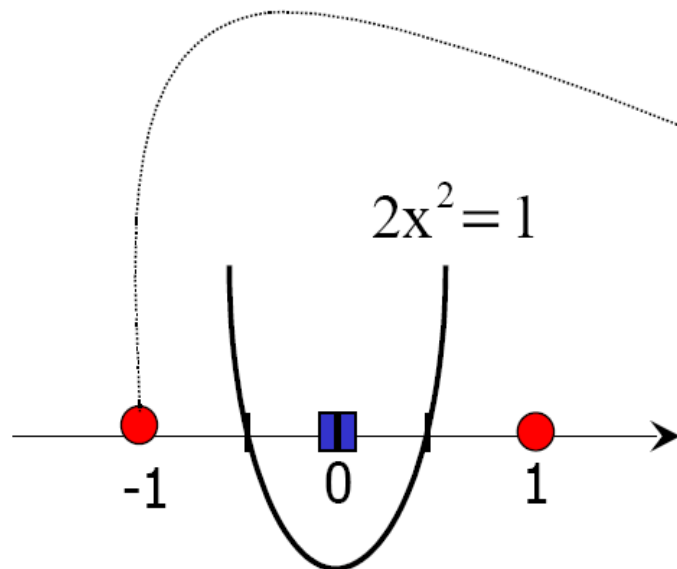
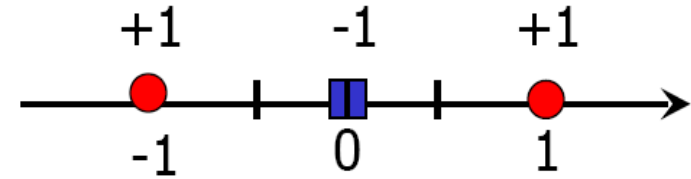
- Non-linear separation is **impossible** in **most problems**

Non-separable datasets

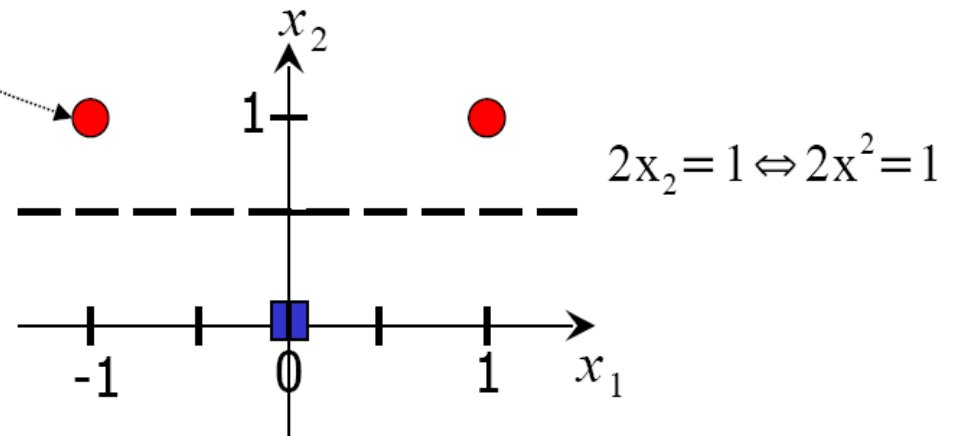
- Solutions:

1) Nonlinear classifiers

2) **Increase dimensionality** of dataset and add a **non-linear mapping Φ**

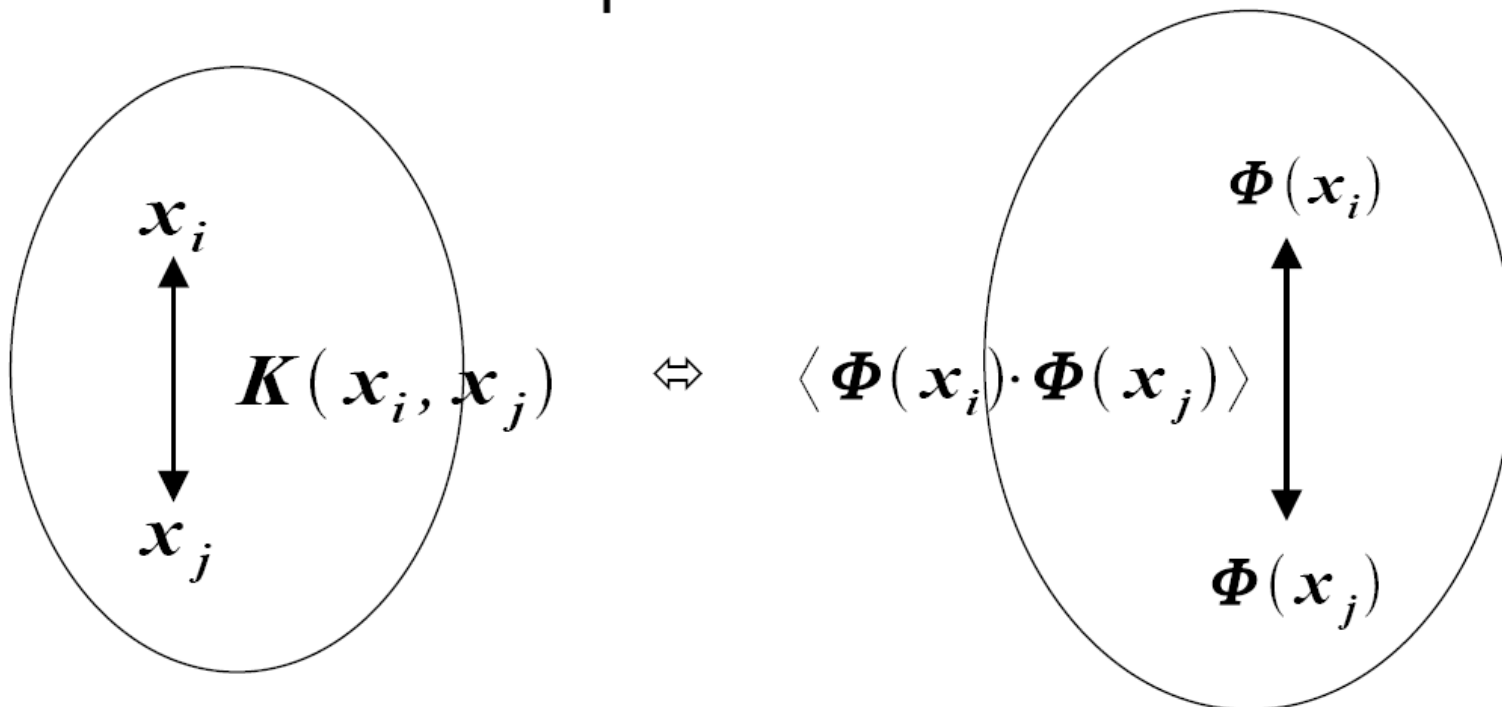


$$[x] \rightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$

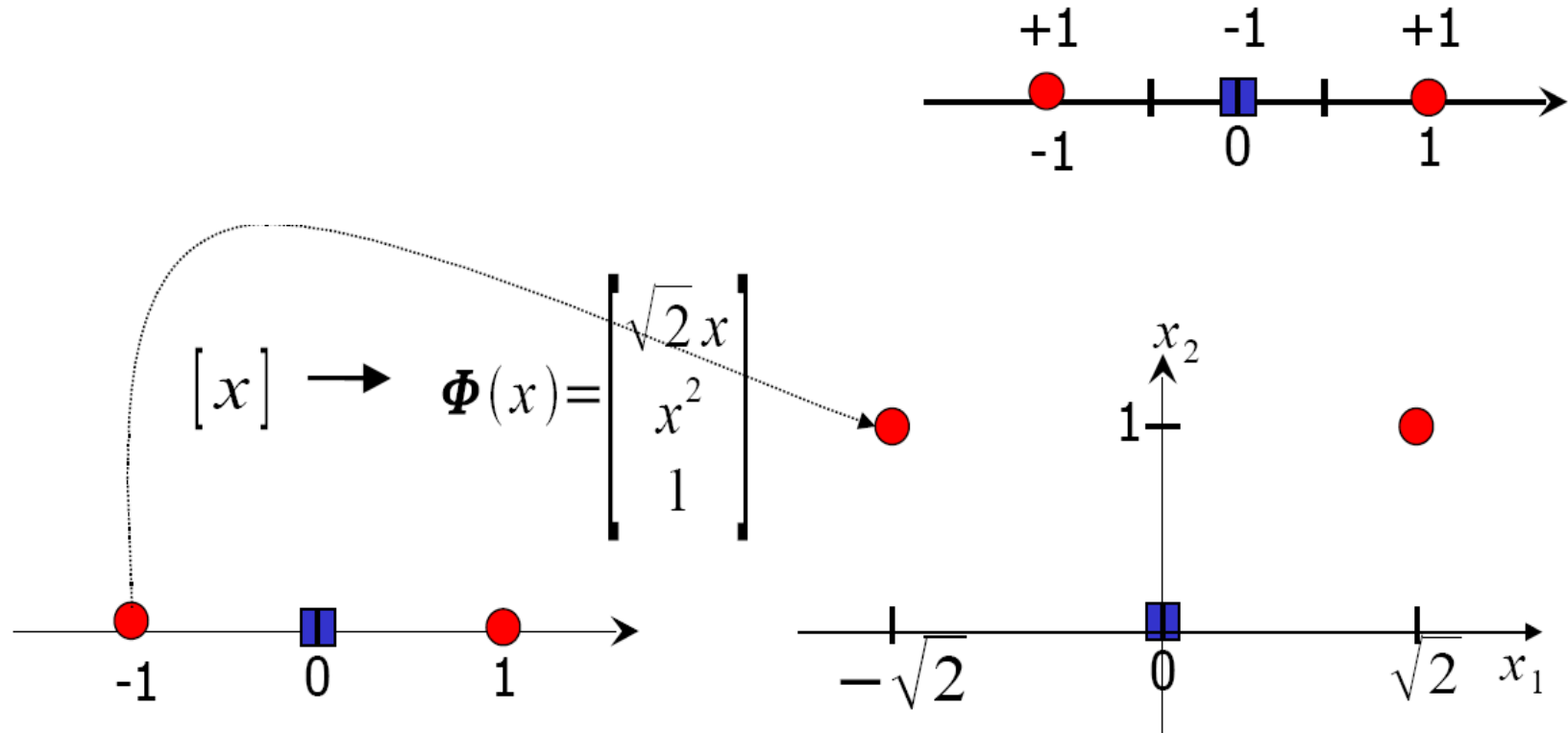


Kernel Trick

- Kernel function
 - in the original space
 - Inner product
 - In the feature space with increased dimension
- “similarity measure”
between 2 data samples



Kernel Trick Illustrated



$$K(x_i, x_j) = (x_i x_j + 1)^2$$

$$\langle \Phi(x_i) \cdot \Phi(x_j) \rangle = 2x_i x_j + x_i^2 x_j^2 + 1 = (x_i x_j + 1)^2 = K(x_i, x_j) \quad 47$$

Curse of Dimensionality Due to the Non-Linear Mapping

- Primal space
 - Makes optimization much harder
- Dual space
 - Can be avoided

$$\min_{w, b} \frac{1}{2} \langle \Phi^T(w) \cdot \Phi(w) \rangle$$
$$y_i (\langle \Phi^T(w) \cdot \Phi(x_i) \rangle + b) \geq 1$$

$$\max_{\lambda} Q(\lambda) = -0.5 \lambda^T H \lambda + f^T \lambda$$

$$y^T \lambda = 0$$

$$\lambda \geq 0$$

where, $H_{ij} = y_i y_j K(x_i, x_j)$

f is a unit vector

Positive Symmetric Definite Kernels (Mercer Condition)

- Dual form is convex
 - H is P.S.D.
 - Kernel must be P.S.D.

$$Q(\lambda) = -0.5 \lambda^T H \lambda + f^T \lambda$$

where, $H_{ij} = y_i y_j K(x_i, x_j)$

- Mercer kernels
 - Polynomial
 - Gaussian

$$K(\mathbf{x}, \mathbf{y}) = [\langle \mathbf{x}^T \mathbf{y} \rangle + 1]^p$$

$$K(\mathbf{x}, \mathbf{y}) = e^{-\frac{(\mathbf{x}-\mathbf{y})^T \Sigma^{-1} (\mathbf{x}-\mathbf{y})}{2}}$$

Advantages of SVM

- Work very well...
- Error bounds easy to obtain:
 - Generalization error **small** and **predictable**

$$E_{test} = E_{train} + E_{generalization} \leftarrow \frac{|SV|}{N}$$

- Fool-proof method:
 - (Mostly) three kernels to choose from:
 - Gaussian
 - Linear and Polynomial
 - Sigmoid
 - Very small number of parameters to optimize

Limitations of SVM

→ Size limitation:

- Size of kernel matrix is quadratic with the number of training vectors

→ Speed limitations:

- 1) During **training**:
very large quadratic programming problem solved numerically
 - Solutions:
 - **Chunking**
 - **Sequential Minimal Optimization (SMO)**
breaks QP problem into many small QP problems solved analytically
 - **Hardware** implementations
- 2) During **testing**:
number of support vectors
 - Solution: **Online SVM**