# Introduction to Machine Learning (CSCI-UA 473): Fall 2021

## Lecture 7: Support Vector Machines - 2

**Sumit Chopra**
Courant Institute of Mathematical Sciences
Department of Radiology - Grossman School of Medicine
NYU

# Lecture Outline

Primal Formation of SVMs

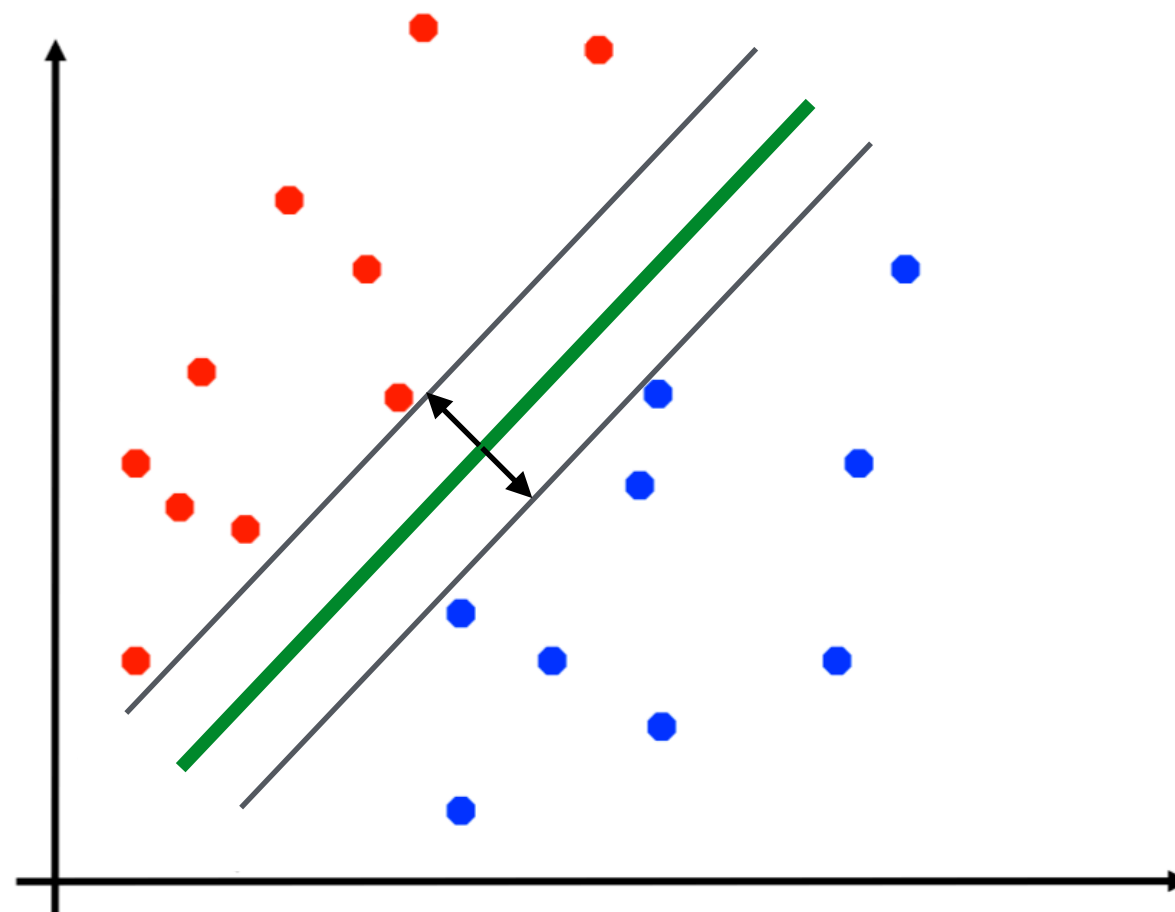Overview of Optimization

Dual Formulation of SVMs

SVMs for Non-Linearly Separable Case

# Maximum Margin Classifiers

SVMs solve the following optimization problem to compute a hyper-plane which has the maximum margin

Convex quadratic optimization with linear constraints

$$\min_{w,b} \quad \frac{1}{2}||w||^2$$

$$\mathbf{s.t.} \quad y^i(w^T x^i + b) \geq 1, \quad \forall i = 1,\ldots,n$$
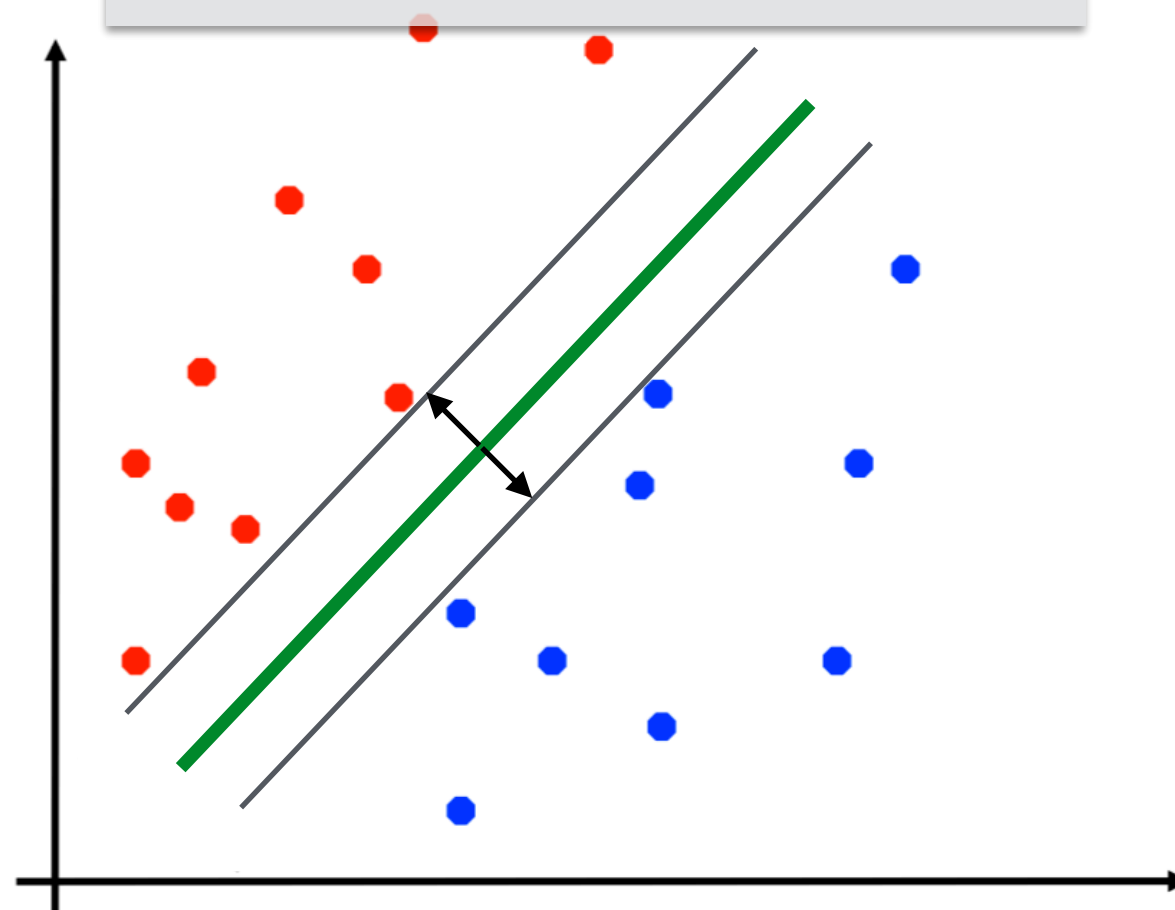
# Maximum Margin Classifiers

SVMs solve the following optimization problem to compute a hyper-plane which has the maximum margin

Convex quadratic optimization with linear constraints

$$\min_{w,b} \quad \frac{1}{2}||w||^2$$

$$\mathbf{s.t.} \quad y^i(w^T x^i + b) \geq 1, \quad \forall i = 1,\ldots,n$$
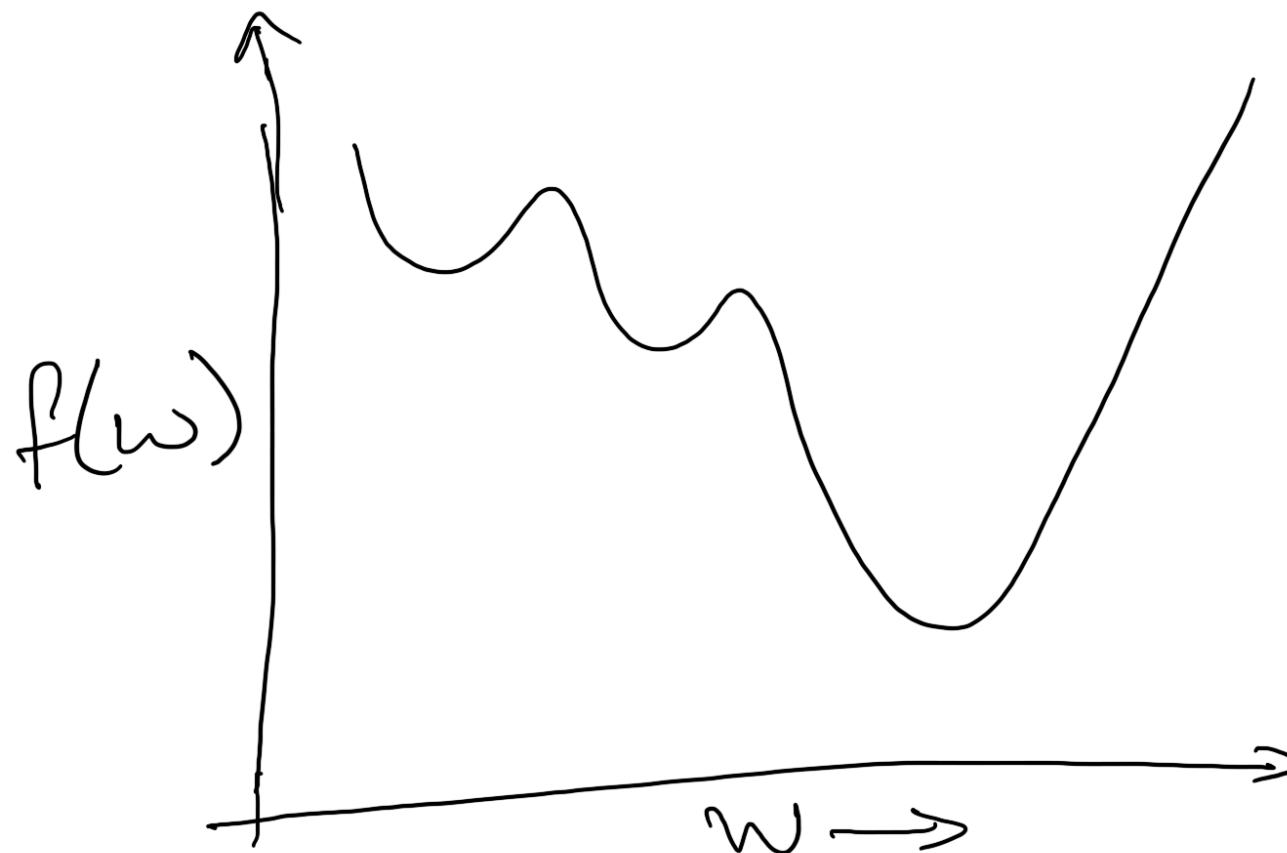
How do we solve this?

# 60,000 View of Function Optimization

# Overview of Optimization

$$\min(or \max) \quad f(w) \qquad\qquad w \in \mathfrak{R}^d$$

Unconstrained
Optimization Problem

$f(w)$

$w \rightarrow$

# Overview of Optimization

$$\min(or \max)_w \quad f(w) \qquad\qquad w \in \mathfrak{R}^d$$

$$g_1(w) \leq 0 \qquad\qquad\qquad h_1(w) = 0$$

$$g_1(w) \leq 0 \qquad\qquad\qquad h_1(w) = 0$$

$$\vdots \qquad\qquad\qquad\qquad \vdots$$

$$g_m(w) \leq 0 \qquad\qquad\qquad h_n(w) = 0$$

Inequality constrains $\qquad\qquad\qquad$ Equality constraints
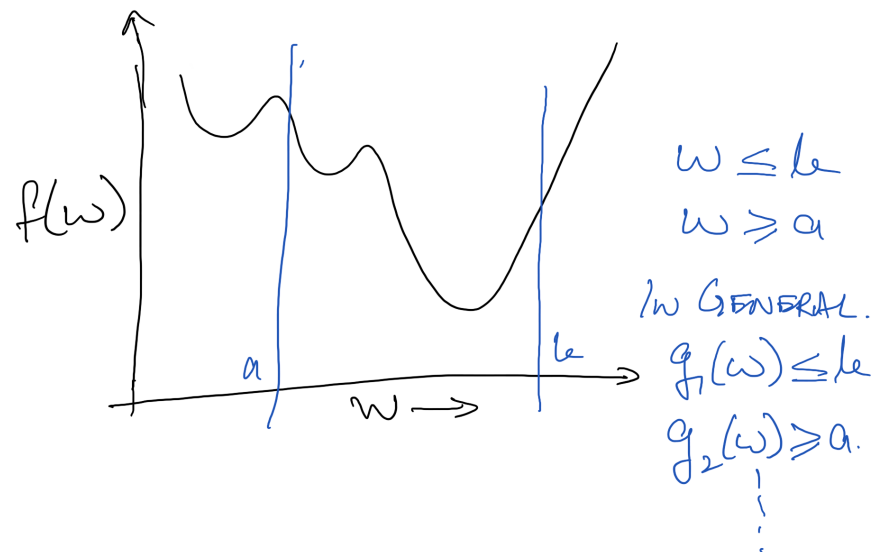
Any equality constraint can be represented as two inequality constraint

$$g(w) = 0 \qquad \rightarrow \qquad g(w) \leq 0 \quad -g(w) \leq 0$$



$$f(w)$$

$$w \leq b$$
$$w \geq a$$

In General.
$$g_1(w) \leq b$$
$$g_2(w) \geq a.$$

# Constrained Optimization Problem: Lagrangian Multipliers

<span style="color:red">Equality Constraints</span>

$$\min_{w} \quad f(w) \qquad \mathbf{s.t.} \qquad \begin{matrix} h_1(w) = 0 \\ h_1(w) = 0 \\ \vdots \\ h_e(w) = 0 \end{matrix} \qquad w \in \mathfrak{R}^d$$

Reduces the problem with $d$ variables with $e$ constraints into an unconstrained problem of $d + e$ variables

Introduces a scalar variable (called the Lagrangian multiplier) for each constraint and forms a linear combination involving the multipliers as coefficients

$$L(w, \beta) = f(w) + \sum_{i=1}^{e} \beta_i h_i(w)$$

$$\frac{\partial L}{\partial w_j} = 0, \qquad \frac{\partial L}{\partial \beta_i} = 0 \qquad \qquad \forall j \in \{1, \dots, d\} \qquad \forall i \in \{1, \dots, e\}$$

# Generalized Lagrangian

$$h_1(w) = 0 \qquad g_1(w) \leq 0$$

$$\min_{w} \quad f(w) \qquad \mathbf{s.t.} \qquad \begin{aligned} h_1(w) &= 0 & g_1(w) &\leq 0 \\ &\vdots & &\vdots \\ h_e(w) &= 0 & g_k(w) &\leq 0 \end{aligned} \qquad w \in \mathfrak{R}^d$$

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{e} \beta_i h_i(w)$$

# Generalized Lagrangian

$$O_p(w) = \max_{\alpha,\beta:\alpha_i \geq 0} L(w,\alpha,\beta) = \max_{\alpha,\beta:\alpha_i \geq 0} \left[ f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{e} \beta_i h_i(w) \right]$$

# Generalized Lagrangian

$$O_p(w) = \max_{\alpha,\beta:\alpha_i \geq 0} L(w,\alpha,\beta) = \max_{\alpha,\beta:\alpha_i \geq 0} \left[ f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{e} \beta_i h_i(w) \right]$$

If $w$ violates any of the constraints, either: $g_i(w) > 0$ or $h_i(w) \neq 0$ then

$$O_p(w) = \max_{\alpha=\{\alpha_1,\ldots,\alpha_k\},\beta=\{\beta_1,\ldots,\beta_e\}:\alpha_i \geq 0} \left[ f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{e} \beta_i h_i(w) \right]$$

$$= \infty$$

# Generalized Lagrangian

$$O_p(w) = \max_{\alpha,\beta:\alpha_i \geq 0} L(w,\alpha,\beta) = \max_{\alpha,\beta:\alpha_i \geq 0} \left[ f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{e} \beta_i h_i(w) \right]$$

If $w$ violates any of the constraints, either: $g_i(w) > 0$ or $h_i(w) \neq 0$ then

$$O_p(w) = \max_{\alpha=\{\alpha_1,\ldots,\alpha_k\},\beta=\{\beta_1,\ldots,\beta_e\}:\alpha_i \geq 0} \left[ f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{e} \beta_i h_i(w) \right]$$

$$= \infty$$

If $w$ satisfies all the constraints $\qquad g_i(w) \leq 0 \qquad h_i(w) = 0$

$$O_p(w) = \max_{\alpha=\{\alpha_1,\ldots,\alpha_k\},\beta=\{\beta_1,\ldots,\beta_e\}:\alpha_i \geq 0} \left[ f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{e} \beta_i h_i(w) \right]$$

$$= f(w)$$

12

# Generalized Lagrangian

If $w$ satisfies all the constraints

$$O_p(w) = \max_{\alpha=\{\alpha_1,\ldots,\alpha_k\},\beta=\{\beta_1,\ldots,\beta_e\}:\alpha_i\geq 0} \left[ f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{e} \beta_i h_i(w) \right]$$

$$= f(w)$$

So long as the constraints are satisfied solution to $O_p(w)$ is the same as $f(w)$

$$p* = \min_{w} O_p(w) = \min_{w} \max_{\alpha,\beta:\alpha_i\geq 0} L(w,\alpha,\beta) = \min_{w} f(x)$$

# Generalized Lagrangian

Define a Dual problem

$$\max_{\alpha,\beta:\alpha_i\geq 0} O_D(\alpha,\beta) = \max_{\alpha,\beta:\alpha_i\geq 0} \left[\min_w L(w,\alpha,\beta)\right]$$

$$d* = \max_{\alpha,\beta:\alpha_i\geq 0} \left[\min_w L(w,\alpha,\beta)\right]$$

# Generalized Lagrangian

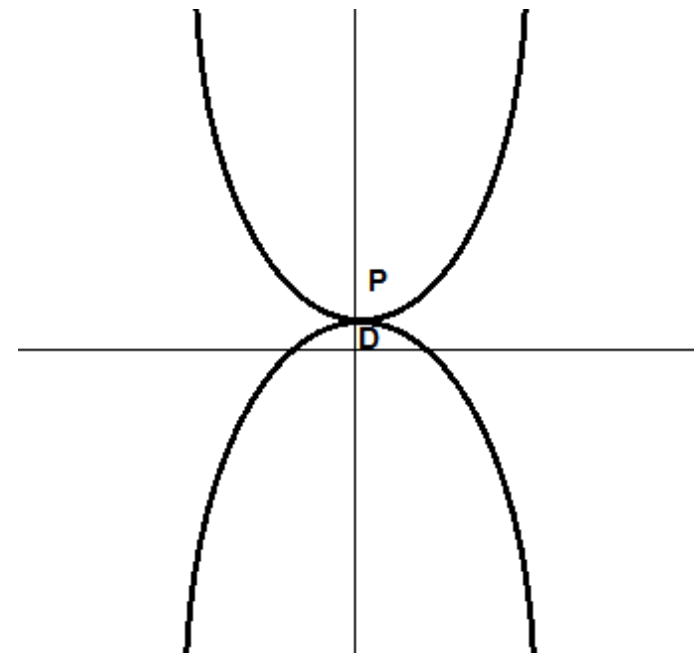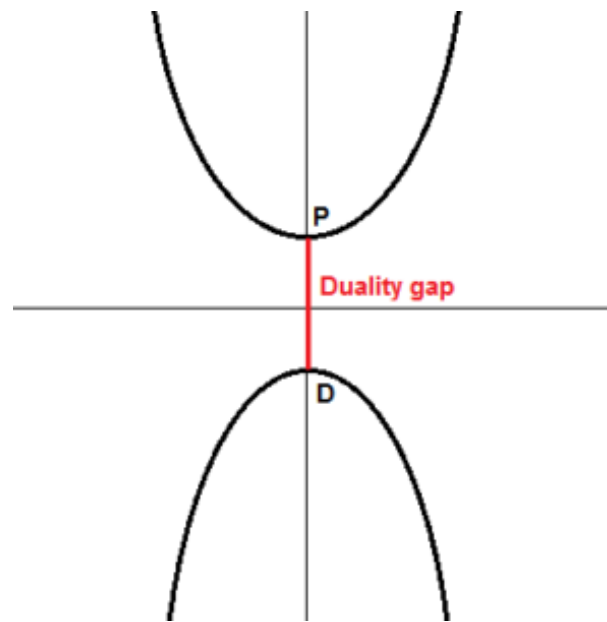Define a Dual problem

$$\max_{\alpha,\beta:\alpha_i \geq 0} O_D(\alpha,\beta) = \max_{\alpha,\beta:\alpha_i \geq 0} \left[ \min_w L(w,\alpha,\beta) \right]$$

$$d* = \max_{\alpha,\beta:\alpha_i \geq 0} \left[ \min_w L(w,\alpha,\beta) \right]$$

$d* \leq p*$: always true

$d* = p*$: sometimes true. Under certain condition

# Generalized Lagrangian

Define a Dual problem

$$\max_{\alpha,\beta:\alpha_i\geq 0} O_D(\alpha,\beta) = \max_{\alpha,\beta:\alpha_i\geq 0}\left[\min_w L(w,\alpha,\beta)\right]$$

$$d* = \max_{\alpha,\beta:\alpha_i\geq 0}\left[\min_w L(w,\alpha,\beta)\right]$$

$d* \leq p*$: always true

$d* = p*$: sometimes true. Under certain condition

$f, g_i$: are convex

$h_i$ :affine $(h_i(w) = a_i^T w + b_i)$

$g_i$ are strictly possible (i.e., there exists some $w$ such that $g_i(w) \leq 0$

# Generalized Lagrangian

Define a Dual problem

$$\max_{\alpha,\beta:\alpha_i\geq0} O_D(\alpha,\beta) = \max_{\alpha,\beta:\alpha_i\geq0}\left[\min_w L(w,\alpha,\beta)\right]$$

$$d* = \max_{\alpha,\beta:\alpha_i\geq0}\left[\min_w L(w,\alpha,\beta)\right]$$

$d* \leq p*$: always true
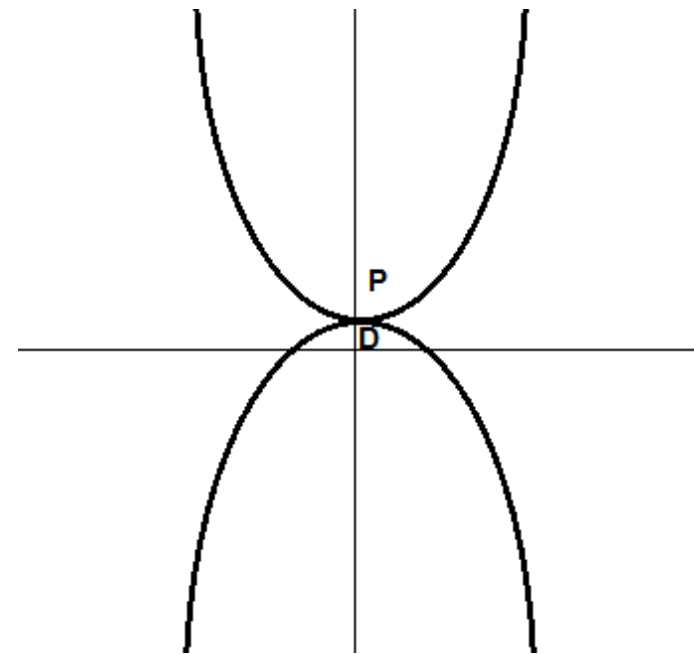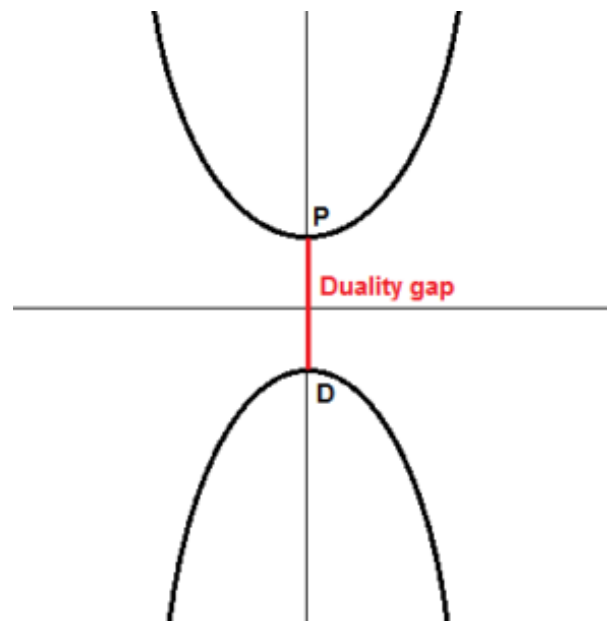
$d* = p*$: sometimes true. Under certain condition

$f, g_i$: are convex

Then there exists $w*, \alpha*, \beta*$, such that

$h_i$ :affine ($h_i(w) = a_i^T w + b_i$)

$w*$ solves primal problem

$g_i$ are strictly possible (i.e., there exists some $w$ such that $g_i(w) \leq 0$

$\alpha*, \beta*$ solves the dual problem

$$d* = p*$$

# Generalized Lagrangian

The solution $w*, \alpha*, \beta*$ satisfies the KKT (Karush-Kuhn-Tucker) Conditions

$$\frac{\partial}{\partial w_i} L(w*, \alpha*, \beta*) = 0 \qquad i = 1, \ldots, d$$

$$\frac{\partial}{\partial \beta_i} L(w*, \alpha*, \beta*) = 0 \qquad i = 1, \ldots, e$$

$$\alpha_i^* g_i(w*) = 0 \qquad i = 1, \ldots, k$$

$$g_i(w*) \leq 0 \qquad i = 1, \ldots, k$$

$$\alpha_i^* \geq 0 \qquad i = 1, \ldots, k$$

# Generalized Lagrangian

The solution $w*, \alpha*, \beta*$ satisfies the KKT (Karush-Kuhn-Tucker) Conditions

$$\frac{\partial}{\partial w_i} L(w*, \alpha*, \beta*) = 0 \qquad i = 1, \ldots, d$$

Complementarity
Condition

$$\frac{\partial}{\partial \beta_i} L(w*, \alpha*, \beta*) = 0 \qquad i = 1, \ldots, e$$

$$\alpha_i^* g_i(w*) = 0 \qquad i = 1, \ldots, k$$

$$g_i(w*) \leq 0 \qquad i = 1, \ldots, k$$

$$\alpha_i^* \geq 0 \qquad i = 1, \ldots, k$$

$$\alpha_i > 0 \implies g_i(w*) = 0$$

We say that constraint $g_i(w*) \leq 0$ is **active**

# Support Vector Machines

# Support Vector Machines

$$\min_{w,b} \quad \frac{1}{2}||w||^2$$

$$\mathbf{s.t.} \quad y^i(w^T x^i + b) \geq 1, \quad \forall i = 1,\ldots,n$$

The Lagrangian is given by

$$L(\alpha, w, b) = \frac{1}{2}||w||^2 - \sum_{i=1}^{N} \alpha_i \left( y^i(w^T x^i + b) - 1 \right)$$

The Dual of the problem is given by

$$O_D(\alpha) = \min_{w,b} L(\alpha, w, b)$$

Complementarity
Condition

$$\nabla_w L(\alpha, w, b) = w - \sum_{i=1}^{N} \alpha_i y^i x^i = 0 \quad \implies \quad w = \sum_{i=1}^{N} \alpha_i y^i x^i$$

$$\nabla_b L(\alpha, w, b) = - \sum_{i=1}^{N} \alpha_i y^i = 0$$

Plug the value of $w$ back into
the Lagrangian

$$\alpha_i[y^i(w^T x^i + b) - 1] = 0$$

# Support Vector Machines

The Lagrangian is given by

$$L(\alpha, w, b) = \frac{1}{2}||w||^2 - \sum_{i=1}^{N} \alpha_i \left( y^i(w^T x^i + b) - 1 \right)$$

$$L(\alpha, w, b) = \frac{1}{2}\left|\left|\sum_{i=1}^{N} \alpha_i y^i x^i w\right|\right|^2 - \sum_{i,j=1}^{N} \alpha_i \alpha_j y^i y^j (x^i \cdot x^j) - \sum_{i=1}^{N} \alpha_i y^i b + \sum_{i=1}^{N} \alpha_i$$

$$-\frac{1}{2}\sum_{i,j=1}^{N} \alpha_i \alpha_j y^i y^j (x^i \cdot x^j)$$

$$\nabla_w L(\alpha, w, b) = w - \sum_{i=1}^{N} \alpha_i y^i x^i = 0 \implies w = \sum_{i=1}^{N} \alpha_i y^i x^i$$

$$\nabla_b L(\alpha, w, b) = -\sum_{i=1}^{N} \alpha_i y^i = 0$$

$$\alpha_i[y^i(w^T x^i + b) - 1] = 0$$

# Support Vector Machines

$$\max_{\alpha} L(\alpha, w, b) = \max_{\alpha} \sum_{i=1}^{N} \alpha_i - \sum_{i,j=1}^{N} \alpha_i \alpha_j y^i y^j (x^i \cdot x^j)$$

$$\alpha_i \geq 0$$

$$\mathbf{s.t.} \quad \sum_{i=1}^{N} \alpha_i y^i = 0$$

This is Dual formulation of SVMs

We can solve the Dual problem and find $\alpha*$ instead of the Primal problem

$$\min_{w,b} \quad \frac{1}{2}||w||^2 \longrightarrow \quad w* = \sum_{i=1}^{N} \alpha_i* y^i x^i$$

$$\mathbf{s.t.} \quad y^i(w^T x^i + b) \geq 1, \qquad \forall i = 1,\ldots,n$$

$$b* = -\frac{\max_{i, y^i = -1}(w*)^T x^i + \min_{i, y^i = 1}(w*)^T x^i}{2}$$

# Support Vectors

KKT conditions at the optimal solution

$$\nabla_w L(\alpha, w, b) = w - \sum_{i=1}^{N} \alpha_i y^i x^i = 0 \qquad \Longrightarrow \qquad w = \sum_{i=1}^{N} \alpha_i y^i x^i$$

$$\nabla_b L(\alpha, w, b) = -\sum_{i=1}^{N} \alpha_i y^i = 0$$
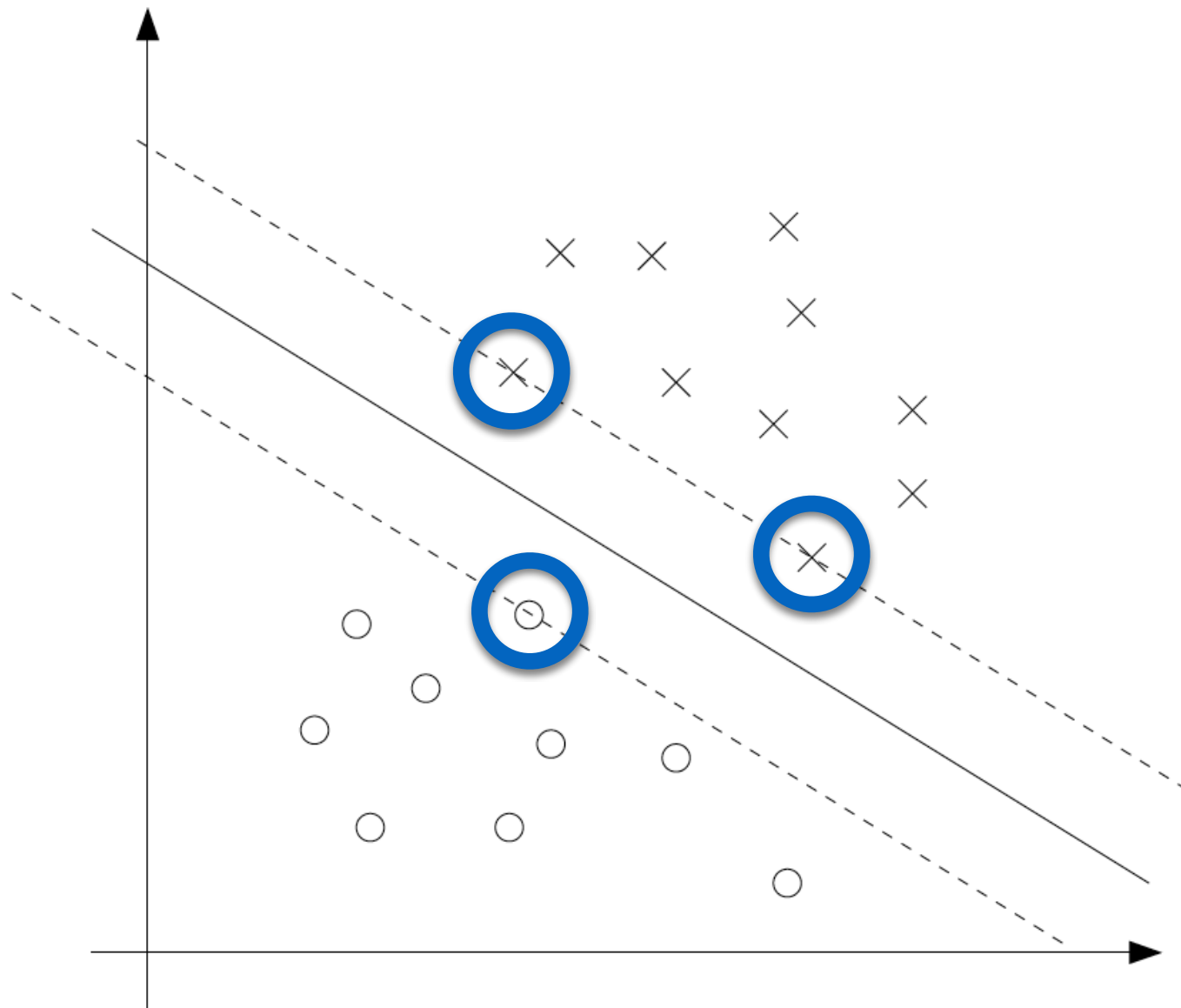
$$\alpha_i [y^i(w^T x^i + b) - 1] = 0$$

Either $\alpha_i = 0$ or $y^i(w^T x^i + b) - 1 = 0$

Thus for all $i$ for which $\alpha_i \neq 0$ we have $y^i(w^T x^i + b) = 1$

The constraint is active

These are the points for which the Geometric Margin is equal to 1

# Support Vectors



Constraint is active for only a few training points

These points are called the support vectors

# Predictions with SVMs

We have solved the optimization problem and found the solution $\alpha^*$ and hence $w^*$

What is the class for a new data point $x^0$?

Naive approach
Compute $(w^*)^T \cdot x^0 + b^*$ and assign class +1 if positive otherwise assign class -1

There is a better way

$$(w^*)^T x^0 + b^* = \left( \sum_{i=1}^{N} \alpha_i^* y^i x^i \right)^T \cdot x + b^*$$

$$= \sum_{i=1}^{N} \alpha_i^* y^i (x^{i^T} \cdot x^0) + b^*$$

You only need the inner products with the training samples!

Furthermore, $\alpha_i^*$ is non-zero for only a few training points (the support vectors)

# Summary

Solve

$$\max_{\alpha} \quad \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{N} y^{(i)} y^{(j)} \alpha_i \alpha_j < x^{(i)}, x^{(j)} >$$

s.t. $\quad \alpha_i \geq 0 \, , i = 1, \ldots, N$

$$\sum_{i=1}^{N} \alpha_i y^{(i)} = 0$$

Predict

$$\sum_{i=1}^{N} \alpha_i^* y^{(i)} < (x^{(i)})^T, x > + b$$
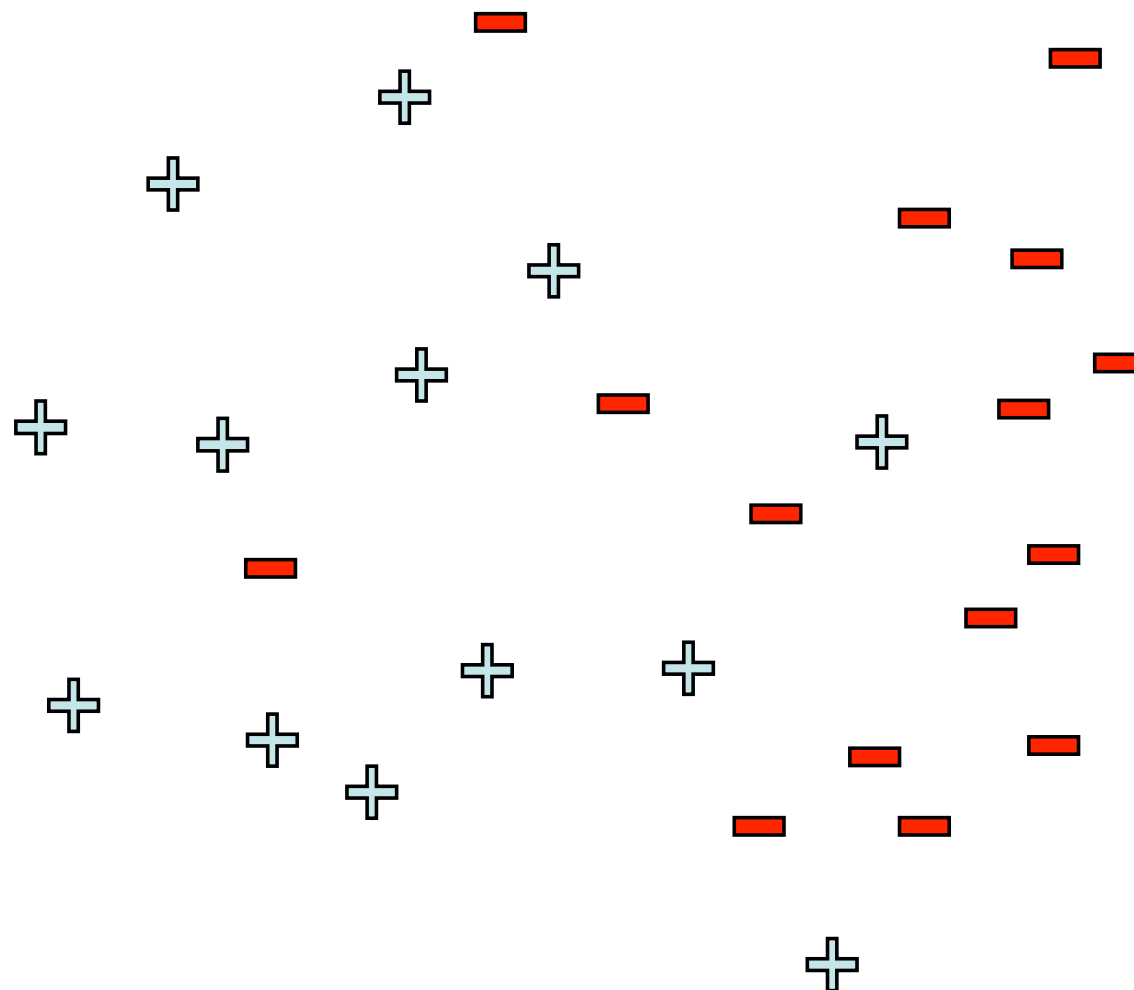
Inner products in SVMs is the key

Also forms the basis of the Kernels in SVMs

# When Linear Separability Does Not Exist

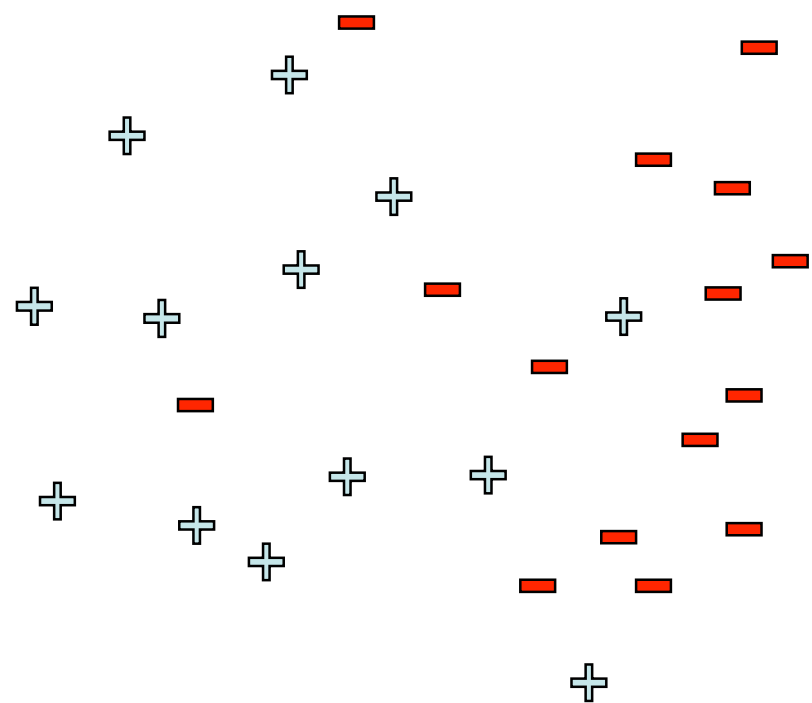What happens when there is no linear separability?

All the constraints get violated by $w, b, \gamma$

The geometric margin loses its meaning

$$\text{minimi}$$

$$\left( \mathbf{w}.\mathbf{x}_j + \right.$$

$$\min_{\mathbf{w},b}$$

# When Linear Separability Does Not Exist

One possible solution is to find $w, b$ such that the minimum number of constraints are violated

$$\min_{w,b} \quad \#mistakes$$

$$\underset{\mathbf{w},b}{\text{minimize}} \quad \mathbf{w}.\mathbf{w}$$

$$\text{s.t.} \quad (\mathbf{w}.\mathbf{x}_j + b) \, y (w^T x^i + b) y_i \geq 1 \quad , \forall i = 1,\ldots,N$$

$$\min_{\mathbf{w},b} \sum_j \ell_{0,1}(y_j, w.x_j + b) \quad \min_{\mathbf{w},b} \sum_{i=1}^{N} l_{0,1}(y^i, w^T x^i + b)$$

$$where \quad \ell_{0,1}(y, \hat{y}) = 1[y \neq \text{sign}(\hat{y})] \quad l_{0,1}(y, \hat{y}) = 1 \quad if \quad \hat{y} \neq y$$

This is an NP-Hard Problem

# When Linear Separability Does Not Exist
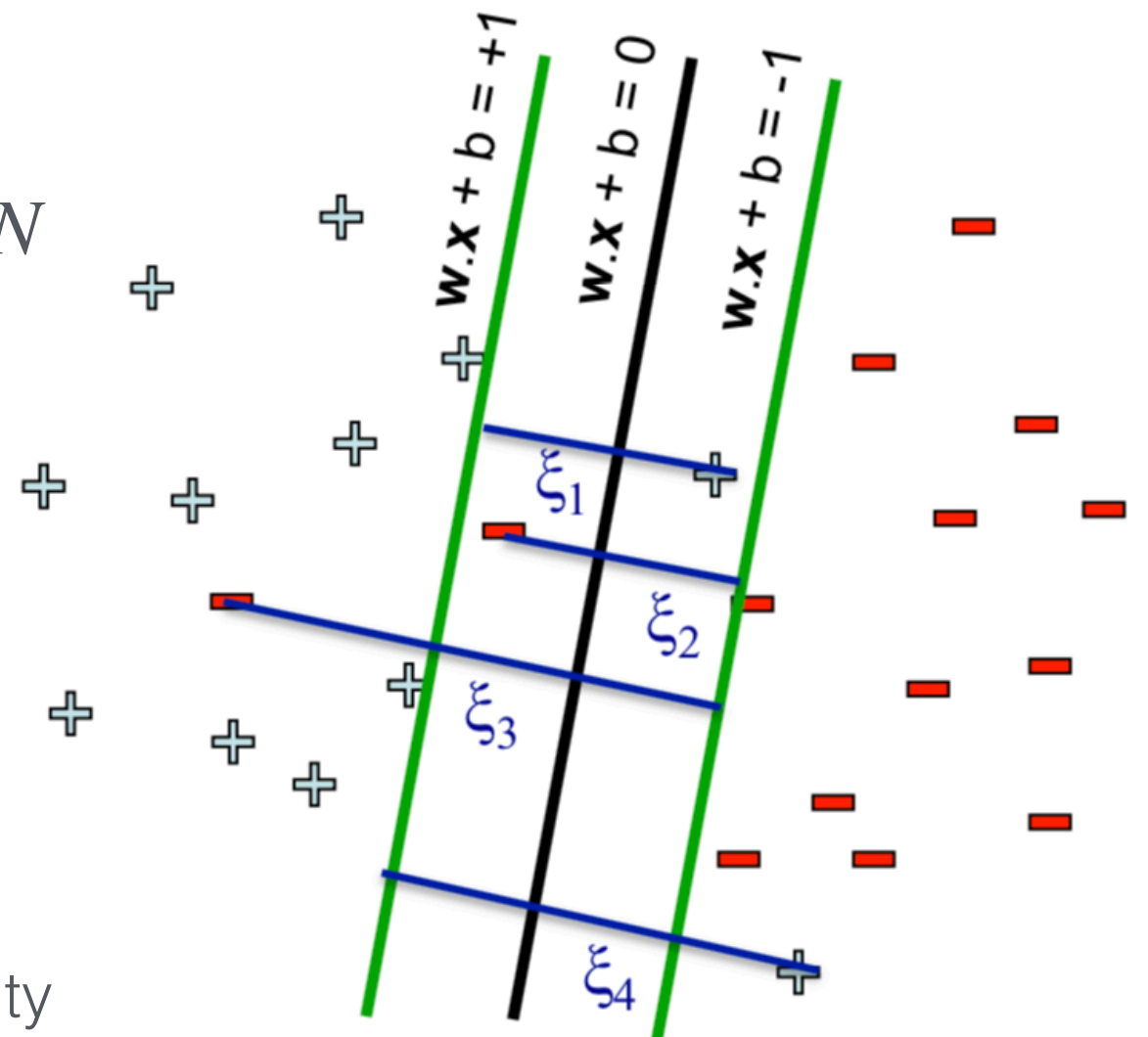
Another solution: Allow some slack!

Let's ignore that we are looking for a largest margin classifier

Instead look for a classifier with the minimal slack

$$\min_{w,b,\xi_i} \sum_{i=1}^{N} \xi_i$$

**s.t.** $\quad y^i(w^T x^i + b) \geq 1 - \xi_i \quad i = 1,\ldots,N$

$\qquad \xi_i \geq 0, i = 1,\ldots,N$



w.x + b = +1

w.x + b = 0

w.x + b = -1

$\xi_1$

$\xi_2$

$\xi_3$

$\xi_4$

If functional margin is $\geq 1$ the no penalty

If functional margin is $< 1$ then pay linear penalty

# When Linear Separability Does Not Exist

Another solution: Classifier with minimal slack

Optimal value of slack variables

$$\min_{w,b,\xi_i} \quad \sum_{i=1}^{N} \xi_i$$

**s.t.** $\quad y^i(w^T x^i + b) \geq 1 - \xi_i \quad i = 1,\ldots,N$

$\quad\quad\quad \xi_i \geq 0, i = 1,\ldots,N$
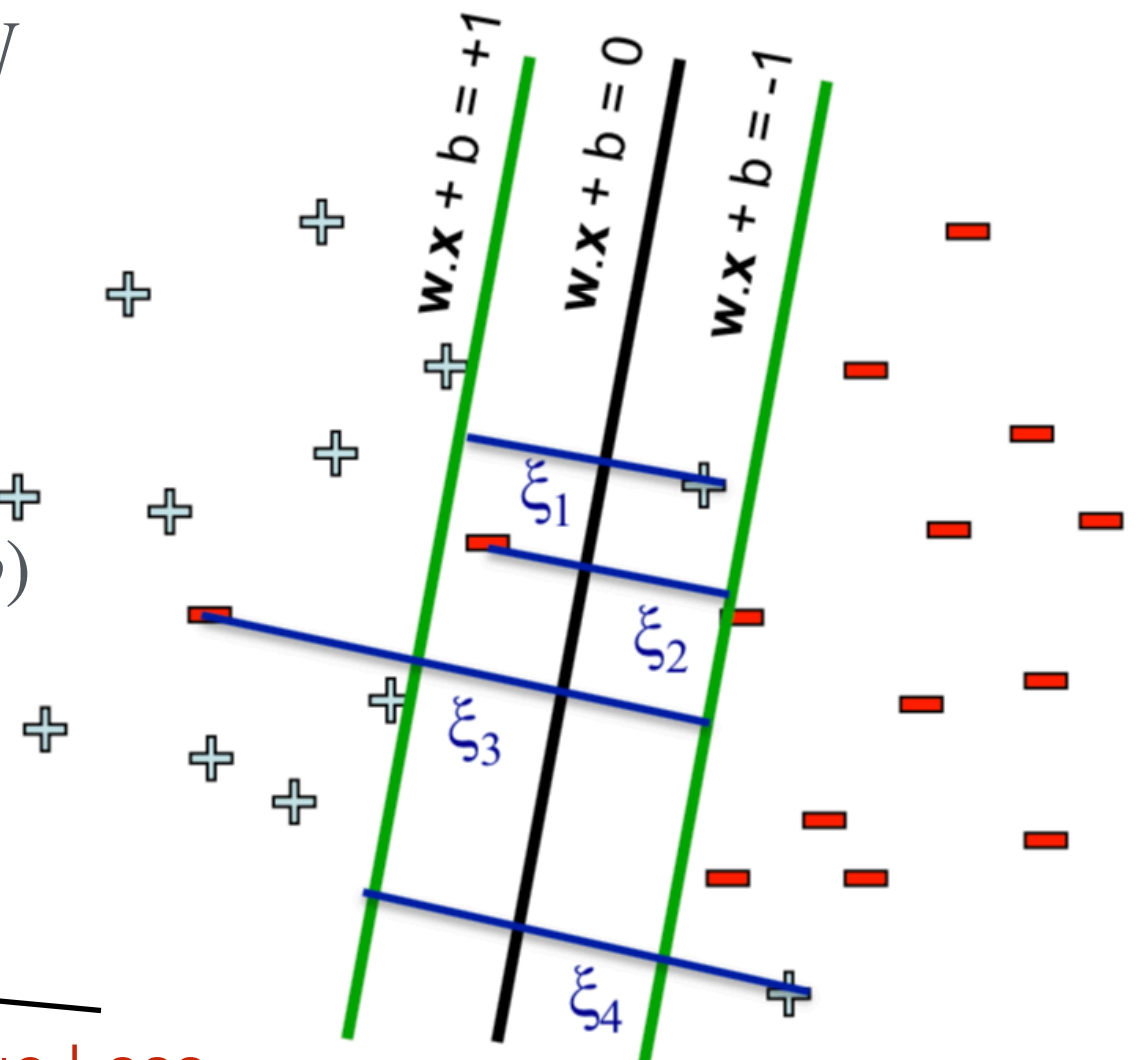
If $y^i(w^T x^i + b) \geq 1 \implies \xi_i = 0$

If $y^i(w^T x^i + b) < 1 \implies \xi_i = 1 - y^i(w^T x^i + b)$
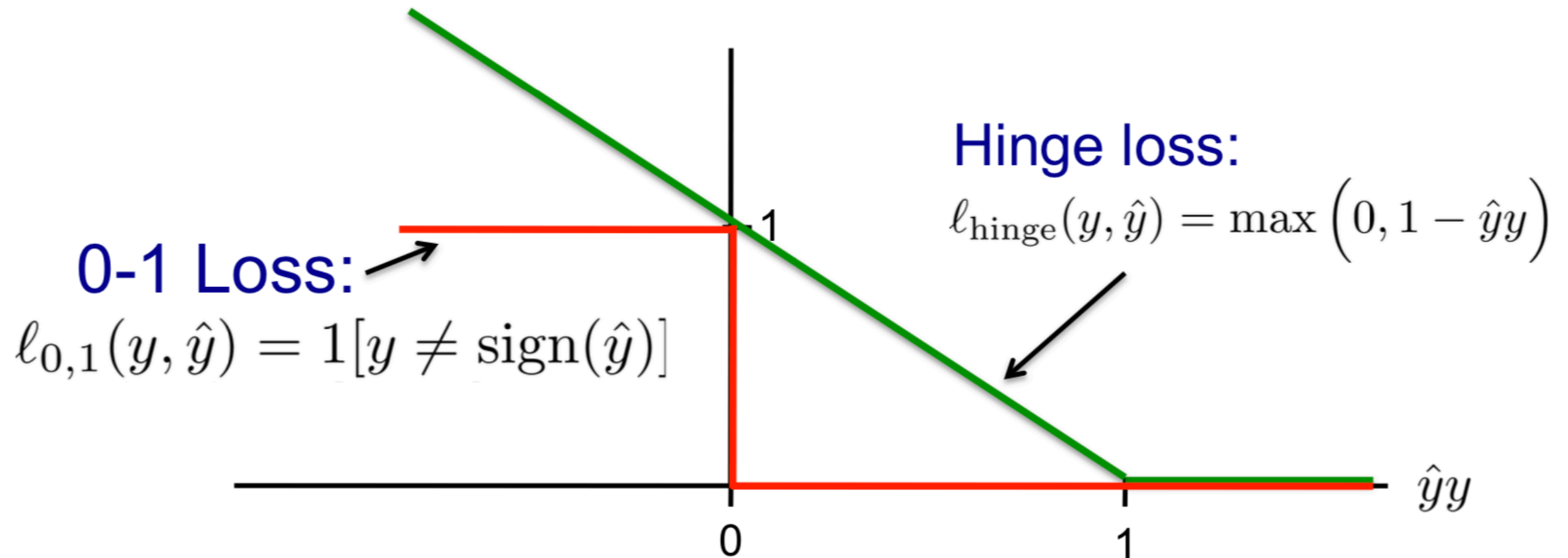
These two conditions can be written as

$$\xi_i = \max[0, 1 - y^i(w^T x^i + b)] \longleftarrow$$

Hinge Loss

$\mathbf{w.x} + b = +1$

$\mathbf{w.x} + b = 0$

$\mathbf{w.x} + b = -1$

$\xi_1$

$\xi_2$

$\xi_3$

$\xi_4$

# When Linear Separability Does Not Exist

This is the tightest convex upper bound of the intractable 0/1 loss



Hinge loss:

$\ell_{\text{hinge}}(y, \hat{y}) = \max\left(0, 1 - \hat{y}y\right)$

0-1 Loss:

$\ell_{0,1}(y, \hat{y}) = 1[y \neq \text{sign}(\hat{y})]$

$\hat{y}y$

# When Linear Separability Does Not Exist

With $\xi_i = \max(0, 1 - y^i(w^T x^i + b))$ we can write the optimization problem as

$$\min_{w,b} \sum_{i=1}^{N} \max(0, 1 - y^i \underbrace{(w^T x^i + b)}_{\hat{y}^i})$$

$$\min_{w,b} \sum_{i=1}^{N} L_{hinge}(y^i, w^T x^i + b)$$

# SVMs Under No Linear Separability

We find the largest margin classifier with some slack

$$\min_{w,b,\xi_i} \quad \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\xi_i$$

$$\textbf{s.t.} \quad y^i(w^T x^i + b) \geq 1 - \xi_i, i = 1,\ldots,N$$

$$\xi_i \geq 0, \qquad i = 1,\ldots,N$$

Thus there are two terms in the objective function that are balanced by the slack penalty $C$

If $C = \infty$ the you have to separate the data

If $C = 0$ then completely ignore the data

This also servers as another regularizer parameter

# SVMs Under No Linear Separability

Equivalent formulation via Hinge loss

$$\min_{w,b,\xi_i} \quad \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\xi_i$$

$$\textbf{s.t.} \quad y^i(w^T x^i + b) \geq 1 - \xi_i, \, i = 1,\ldots,N$$

$$\xi_i \geq 0, \qquad i = 1,\ldots,N$$

$$\min_{w,b} \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{N} L_{hinge}(y^i, w^T x^i + b)$$

Regularizer to prevent over fitting

# End of Lecture 07