

# Support Vector Machine

Mengye Ren

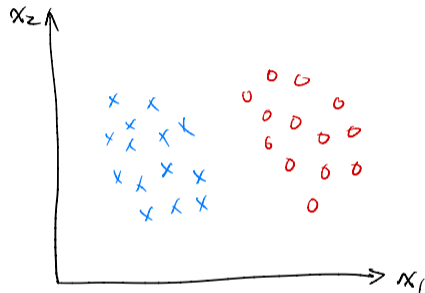
NYU

September 26, 2023

# Maximum Margin Classifier

# Linearly Separable Data

Consider a linearly separable dataset  $\mathcal{D}$ :



Find a separating hyperplane such that

- $w^T x_i > 0$  for all  $x_i$  where  $y_i = +1$
- $w^T x_i < 0$  for all  $x_i$  where  $y_i = -1$

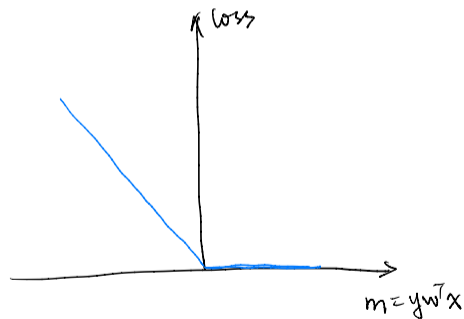
# The Perceptron Algorithm

- Initialize  $w \leftarrow 0$
- While not converged (exists misclassified examples)
  - For  $(x_i, y_i) \in \mathcal{D}$ 
    - If  $y_i w^T x_i < 0$  (wrong prediction)
    - Update  $w \leftarrow w + y_i x_i$
- Intuition: move towards misclassified positive examples and away from negative examples
- Guarantees to find a zero-error classifier (if one exists) in finite steps
- What is the loss function if we consider this as a SGD algorithm?

## Minimize the Hinge Loss

# Perceptron Loss

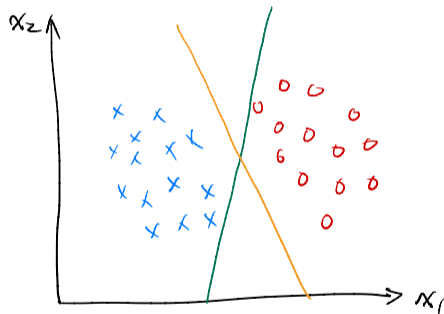
$$l(x, y, w) = \max(0, -yw^T x)$$



# Maximum-Margin Separating Hyperplane

For separable data, there are infinitely many zero-error classifiers.

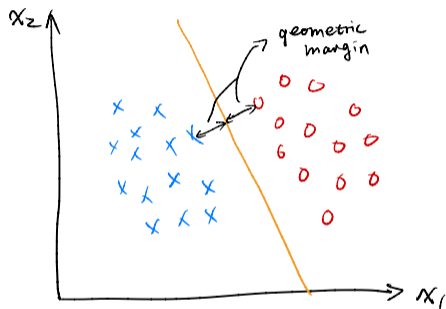
Which one do we pick?



(Perceptron does not return a unique solution.)

# Maximum-Margin Separating Hyperplane

We prefer the classifier that is farthest from both classes of points



- Geometric margin: smallest distance between the hyperplane and the points
- Maximum margin: *largest* distance to the closest points



## Geometric Margin

We want to maximize the distance between the **separating hyperplane** and the **closest** points.

Let's formalize the problem.

### Definition (separating hyperplane)

We say  $(x_i, y_i)$  for  $i = 1, \dots, n$  are **linearly separable** if there is a  $w \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  such that  $y_i(w^T x_i + b) > 0$  for all  $i$ . The set  $\{v \in \mathbb{R}^d \mid w^T v + b = 0\}$  is called a **separating hyperplane**.

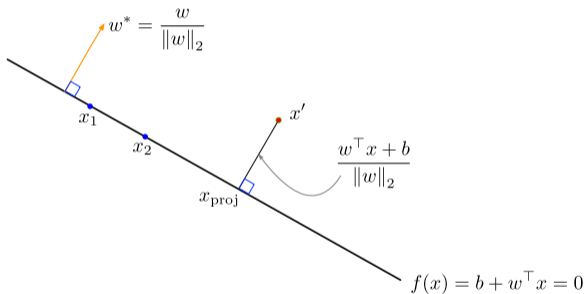
### Definition (geometric margin)

Let  $H$  be a hyperplane that separates the data  $(x_i, y_i)$  for  $i = 1, \dots, n$ . The **geometric margin** of this hyperplane is

$$\min_i d(x_i, H),$$

the distance from the hyperplane to the closest data point.

# Distance between a Point and a Hyperplane



- Any point on the plane  $p$ , and normal vector  $w/\|w\|_2$
- Projection of  $x$  onto the normal:  $\frac{(x'-p)^T w}{\|w\|_2}$
- $(x' - p)^T w = x'^T w - p^T w = x'^T w + b$  (since  $p^T w + b = 0$ )
- Signed distance between  $x'$  and Hyperplane  $H$ :  $\frac{w^T x' + b}{\|w\|_2}$
- Taking into account of the label  $y$ :  
$$d(x', H) = \frac{y(w^T x' + b)}{\|w\|_2}$$

# Maximize the Margin

We want to maximize the geometric margin:

$$\text{maximize } \min_i d(x_i, H).$$

Given separating hyperplane  $H = \{v \mid w^T v + b = 0\}$ , we have

$$\text{maximize } \min_i \frac{y_i(w^T x_i + b)}{\|w\|_2}.$$

Let's remove the inner minimization problem by

$$\begin{aligned} & \text{maximize } M \\ & \text{subject to } \frac{y_i(w^T x_i + b)}{\|w\|_2} \geq M \quad \text{for all } i \end{aligned}$$

Note that the solution is not unique (why?).

# Maximize the Margin

Let's fix the norm  $\|w\|_2$  to  $1/M$  to obtain:

$$\begin{aligned} & \text{maximize} && \frac{1}{\|w\|_2} \\ & \text{subject to} && y_i(w^T x_i + b) \geq 1 \quad \text{for all } i \end{aligned}$$

It's equivalent to solving the minimization problem

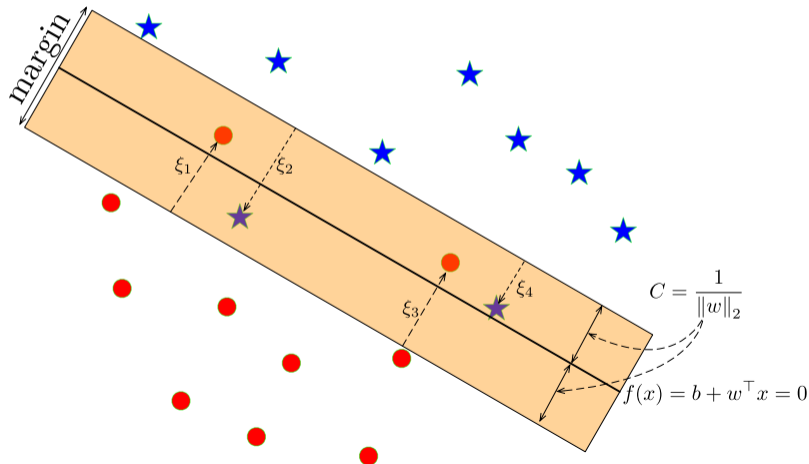
$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|w\|_2^2 \\ & \text{subject to} && y_i(w^T x_i + b) \geq 1 \quad \text{for all } i \end{aligned}$$

Note that  $y_i(w^T x_i + b)$  is the (functional) margin. The optimization finds the minimum norm solution which has a margin of at least 1 on all examples.

## Not linearly separable

What if the data is *not* linearly separable?

For any  $w$ , there will be points with a negative margin.



Introduce **slack variables**  $\xi$ 's to penalize small margin:

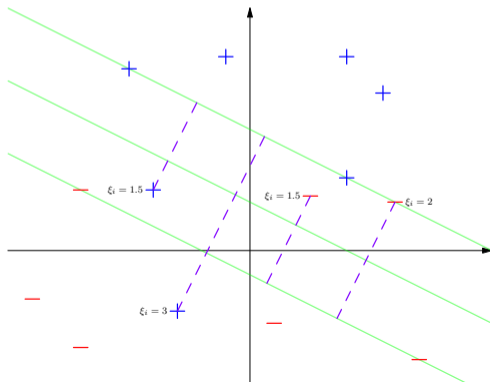
$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|w\|_2^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \\ & \text{subject to} && y_i (w^T x_i + b) \geq 1 - \xi_i \quad \text{for all } i \\ & && \xi_i \geq 0 \quad \text{for all } i \end{aligned}$$

- If  $\xi_i = 0 \forall i$ , it's reduced to hard SVM.
- What does  $\xi_i > 0$  mean?
- What does  $C$  control?

## Slack Variables

$d(x_i, H) = \frac{y_i(w^T x_i + b)}{\|w\|_2} \geq \frac{1 - \xi_i}{\|w\|_2}$ , thus  $\xi_i$  measures the violation by multiples of the geometric margin:

- $\xi_i = 1$ :  $x_i$  lies on the hyperplane
- $\xi_i = 3$ :  $x_i$  is past 2 margin width beyond the decision hyperplane

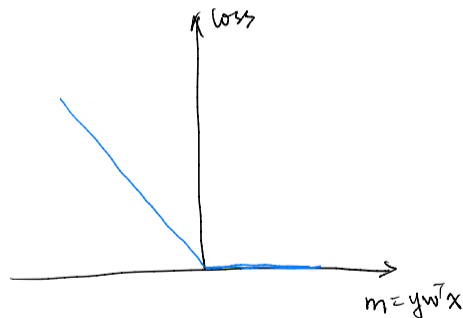


## Minimize the Hinge Loss



# Perceptron Loss

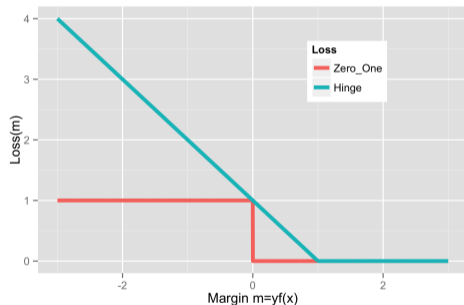
$$\ell(x, y, w) = \max(0, -yw^T x)$$



If we do ERM with this loss function, what happens?

# Hinge Loss

- SVM/Hinge loss:  $\ell_{\text{Hinge}} = \max\{1 - m, 0\} = (1 - m)_+$
- Margin  $m = yf(x)$ ; “Positive part”  $(x)_+ = x\mathbb{1}[x \geq 0]$ .



Hinge is a **convex, upper bound** on 0–1 loss. Not differentiable at  $m = 1$ .  
We have a “margin error” when  $m < 1$ .

# SVM as an Optimization Problem

- The SVM optimization problem is equivalent to

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & \xi_i \geq (1 - y_i [w^T x_i + b]) \text{ for } i = 1, \dots, n \\ & \xi_i \geq 0 \text{ for } i = 1, \dots, n \end{aligned}$$

which is equivalent to

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & \xi_i \geq \max(0, 1 - y_i [w^T x_i + b]) \text{ for } i = 1, \dots, n. \end{aligned}$$

# SVM as an Optimization Problem

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & \xi_i \geq \max(0, 1 - y_i [w^T x_i + b]) \text{ for } i = 1, \dots, n. \end{aligned}$$

Move the constraint into the objective:

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

- The first term is the L2 regularizer.
- The second term is the Hinge loss.

# Support Vector Machine

Using ERM:

- Hypothesis space  $\mathcal{F} = \{f(x) = w^T x + b \mid w \in \mathbb{R}^d, b \in \mathbb{R}\}$ .
- $\ell_2$  regularization (Tikhonov style)
- Hinge loss  $\ell(m) = \max\{1 - m, 0\} = (1 - m)_+$
- The SVM prediction function is the solution to

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i [w^T x_i + b]).$$

- **Not differentiable** because of the max

Two ways to derive the SVM optimization problem:

- Maximize the margin
- Minimize the hinge loss with  $\ell_2$  regularization

Both leads to the minimum norm solution satisfying certain margin constraints.

- **Hard-margin SVM:** all points must be correctly classified with the margin constraints
- **Soft-margin SVM:** allow for margin constraint violation with some penalty

# Subgradient Descent

Now that we have the objective, can we do SGD on it?

Subgradient: generalize gradient for non-differentiable convex functions

# SVM Optimization Problem (no intercept)

- SVM objective function:

$$J(w) = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i w^T x_i) + \lambda \|w\|^2.$$

- Not differentiable... but let's think about gradient descent anyway.
- Hinge loss:  $\ell(m) = \max(0, 1 - m)$

$$\begin{aligned} \nabla_w J(w) &= \nabla_w \left( \frac{1}{n} \sum_{i=1}^n \ell(y_i w^T x_i) + \lambda \|w\|^2 \right) \\ &= \frac{1}{n} \sum_{i=1}^n \nabla_w \ell(y_i w^T x_i) + 2\lambda w \end{aligned}$$



## “Gradient” of SVM Objective

- Derivative of hinge loss  $\ell(m) = \max(0, 1 - m)$ :

$$\ell'(m) = \begin{cases} 0 & m > 1 \\ -1 & m < 1 \\ \text{undefined} & m = 1 \end{cases}$$

- By chain rule, we have

$$\begin{aligned} \nabla_w \ell(y_i w^T x_i) &= \ell'(y_i w^T x_i) y_i x_i \\ &= \begin{cases} 0 & y_i w^T x_i > 1 \\ -y_i x_i & y_i w^T x_i < 1 \\ \text{undefined} & y_i w^T x_i = 1 \end{cases} \end{aligned}$$

## “Gradient” of SVM Objective

$$\nabla_w \ell(y_i w^T x_i) = \begin{cases} 0 & y_i w^T x_i > 1 \\ -y_i x_i & y_i w^T x_i < 1 \\ \text{undefined} & y_i w^T x_i = 1 \end{cases}$$

So

$$\begin{aligned} \nabla_w J(w) &= \nabla_w \left( \frac{1}{n} \sum_{i=1}^n \ell(y_i w^T x_i) + \lambda \|w\|^2 \right) \\ &= \frac{1}{n} \sum_{i=1}^n \nabla_w \ell(y_i w^T x_i) + 2\lambda w \\ &= \begin{cases} \frac{1}{n} \sum_{i: y_i w^T x_i < 1} (-y_i x_i) + 2\lambda w & \text{all } y_i w^T x_i \neq 1 \\ \text{undefined} & \text{otherwise} \end{cases} \end{aligned}$$

## Gradient Descent on SVM Objective?

- The gradient of the SVM objective is

$$\nabla_w J(w) = \frac{1}{n} \sum_{i: y_i w^T x_i < 1} (-y_i x_i) + 2\lambda w$$

when  $y_i w^T x_i \neq 1$  for all  $i$ , and otherwise is undefined.

Potential arguments for why we shouldn't care about the points of nondifferentiability:

- If we start with a random  $w$ , will we ever hit exactly  $y_i w^T x_i = 1$ ?
- If we did, could we perturb the step size by  $\varepsilon$  to miss such a point?
- Does it even make sense to check  $y_i w^T x_i = 1$  with floating point numbers?

# Subgradient

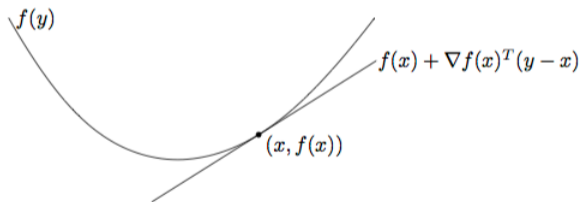
---

# First-Order Condition for Convex, Differentiable Function

- Suppose  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is **convex** and **differentiable**. Then for any  $x, y \in \mathbb{R}^d$

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

- The linear approximation to  $f$  at  $x$  is a **global underestimator** of  $f$ :



- This implies that if  $\nabla f(x) = 0$  then  $x$  is a global minimizer of  $f$ .

# Subgradient Descent

- Move along the negative subgradient:

$$x^{t+1} = x^t - \eta g \quad \text{where } g \in \partial f(x^t) \text{ and } \eta > 0$$

- This can **increase** the objective but gets us **closer to the minimizer** if  $f$  is convex and  $\eta$  is small enough:

$$\|x^{t+1} - x^*\| < \|x^t - x^*\|$$

- Subgradients don't necessarily converge to zero as we get closer to  $x^*$ , so we need **decreasing step sizes**.
- Subgradient methods are **slower** than gradient descent.

# Subgradient descent for SVM

SVM objective function:

$$J(w) = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i w^T x_i) + \lambda \|w\|^2.$$

Pegasos: stochastic subgradient descent with step size  $\eta_t = 1/(t\lambda)$

---

Input:  $\lambda > 0$ . Choose  $w_1 = 0, t = 0$

While termination condition not met

For  $j = 1, \dots, n$  (assumes data is randomly permuted)

$t = t + 1$

$\eta_t = 1/(t\lambda)$ ;

If  $y_j w_t^T x_j < 1$

$w_{t+1} = (1 - \eta_t \lambda) w_t + \eta_t y_j x_j$

Else

$w_{t+1} = (1 - \eta_t \lambda) w_t$

---

- Subgradient: generalize gradient for non-differentiable convex functions
- Subgradient “descent”:
  - General method for non-smooth functions
  - Simple to implement
  - Slow to converge



# The Dual Problem

- In addition to subgradient descent, we can directly solve the optimization problem using a QP solver.
- For convex optimization problem, we can also look into its **dual problem**.

# The Lagrangian

The general [inequality-constrained] optimization problem is:

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

## Definition

The **Lagrangian** for this optimization problem is

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

- $\lambda_i$ 's are called **Lagrange multipliers** (also called the **dual variables**).
- Weighted sum of the objective and constraint functions
- Hard constraints  $\rightarrow$  soft penalty (objective function)

# Lagrange Dual Function

## Definition

The **Lagrange dual function** is

$$g(\lambda) = \inf_x L(x, \lambda) = \inf_x \left( f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) \right)$$

- $g(\lambda)$  is **concave**
- **Lower bound property:** if  $\lambda \succeq 0$ ,  $g(\lambda) \leq p^*$  where  $p^*$  is the optimal value of the optimization problem.
- $g(\lambda)$  can be  $-\infty$  (uninformative lower bound)

# The Primal and the Dual

- For any **primal form** optimization problem,

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

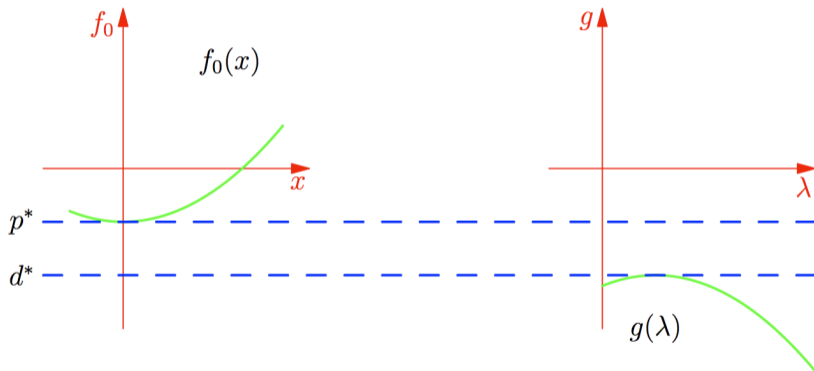
there is a recipe for constructing a corresponding **Lagrangian dual problem**:

$$\begin{array}{ll} \text{maximize} & g(\lambda) \\ \text{subject to} & \lambda_i \geq 0, \quad i = 1, \dots, m, \end{array}$$

- The dual problem is always a convex optimization problem.

# Weak Duality

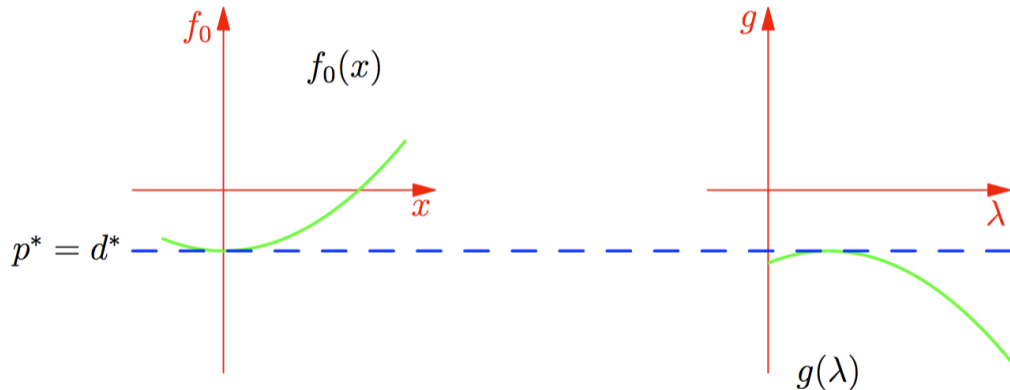
We always have **weak duality**:  $p^* \geq d^*$ .



Plot courtesy of Brett Bernstein.

# Strong Duality

For some problems, we have **strong duality**:  $p^* = d^*$ .



For convex problems, strong duality is fairly typical.

Plot courtesy of Brett Bernstein.

## Complementary Slackness

- **Assume strong duality.** Let  $x^*$  be primal optimal and  $\lambda^*$  be dual optimal. Then:

$$\begin{aligned} f_0(x^*) &= g(\lambda^*) = \inf_x L(x, \lambda^*) \quad (\text{strong duality and definition}) \\ &\leq L(x^*, \lambda^*) \\ &= f_0(x^*) + \sum_{i=1}^m \lambda_i^* f_i(x^*) \\ &\leq f_0(x^*). \end{aligned}$$

Each term in sum  $\sum_{i=1}^m \lambda_i^* f_i(x^*)$  must actually be 0. That is

$$\lambda_i > 0 \implies f_i(x^*) = 0 \quad \text{and} \quad f_i(x^*) < 0 \implies \lambda_i = 0 \quad \forall i$$

This condition is known as **complementary slackness**.

# The SVM Dual Problem



# SVM Lagrange Multipliers

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & -\xi_i \leq 0 \quad \text{for } i = 1, \dots, n \\ & (1 - y_i [w^T x_i + b]) - \xi_i \leq 0 \quad \text{for } i = 1, \dots, n \end{aligned}$$

Lagrange Multiplier	Constraint
$\lambda_i$	$-\xi_i \leq 0$
$\alpha_i$	$(1 - y_i [w^T x_i + b]) - \xi_i \leq 0$

$$L(w, b, \xi, \alpha, \lambda) = \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - y_i [w^T x_i + b] - \xi_i) + \sum_{i=1}^n \lambda_i (-\xi_i)$$

Dual optimum value:  $d^* = \sup_{\alpha, \lambda \geq 0} \inf_{w, b, \xi} L(w, b, \xi, \alpha, \lambda)$

# Strong Duality by Slater's Constraint Qualification

The SVM optimization problem:

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & -\xi_i \leq 0 \text{ for } i = 1, \dots, n \\ & (1 - y_i [w^T x_i + b]) - \xi_i \leq 0 \text{ for } i = 1, \dots, n \end{aligned}$$

Slater's constraint qualification:

- Convex problem + affine constraints  $\implies$  strong duality iff problem is feasible
- Do we have a feasible point?
- For SVM, we have **strong duality**.

# SVM Dual Function: First Order Conditions

Lagrange dual function is the inf over primal variables of  $L$ :

$$g(\alpha, \lambda) = \inf_{w, b, \xi} L(w, b, \xi, \alpha, \lambda)$$
$$= \inf_{w, b, \xi} \left[ \frac{1}{2} w^T w + \sum_{i=1}^n \xi_i \left( \frac{c}{n} - \alpha_i - \lambda_i \right) + \sum_{i=1}^n \alpha_i (1 - y_i [w^T x_i + b]) \right]$$

$$\partial_w L = 0 \iff w - \sum_{i=1}^n \alpha_i y_i x_i = 0 \iff w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\partial_b L = 0 \iff - \sum_{i=1}^n \alpha_i y_i = 0 \iff \sum_{i=1}^n \alpha_i y_i = 0$$

$$\partial_{\xi_i} L = 0 \iff \frac{c}{n} - \alpha_i - \lambda_i = 0 \iff \alpha_i + \lambda_i = \frac{c}{n}$$

# SVM Dual Function

- Substituting these conditions back into  $L$ , the second term disappears.
- First and third terms become

$$\frac{1}{2} w^T w = \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$
$$\sum_{i=1}^n \alpha_i (1 - y_i [w^T x_i + b]) = \sum_{i=1}^n \alpha_i - \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i - b \underbrace{\sum_{i=1}^n \alpha_i y_i}_{=0}$$

- Putting it together, the dual function is

$$g(\alpha, \lambda) = \begin{cases} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i & \begin{array}{l} \sum_{i=1}^n \alpha_i y_i = 0 \\ \alpha_i + \lambda_i = \frac{\epsilon}{n}, \text{ all } i \end{array} \\ -\infty & \text{otherwise.} \end{cases}$$

# SVM Dual Problem

- The dual function is

$$g(\alpha, \lambda) = \begin{cases} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i & \sum_{i=1}^n \alpha_i y_i = 0 \\ -\infty & \alpha_i + \lambda_i = \frac{c}{n}, \text{ all } i \\ & \text{otherwise.} \end{cases}$$

- The dual problem is  $\sup_{\alpha, \lambda \geq 0} g(\alpha, \lambda)$ :

$$\begin{aligned} \sup_{\alpha, \lambda} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \alpha_i + \lambda_i = \frac{c}{n} \quad \alpha_i, \lambda_i \geq 0, \quad i = 1, \dots, n \end{aligned}$$

## Insights from the Dual Problem

# KKT Conditions

For **convex** problems, if **Slater's condition** is satisfied, then **KKT conditions** provide **necessary and sufficient** conditions for the optimal solution.

- Primal feasibility:  $f_i(x) \leq 0 \quad \forall i$
- Dual feasibility:  $\lambda \succeq 0$
- Complementary slackness:  $\lambda_i f_i(x) = 0$
- First-order condition:

$$\frac{\partial}{\partial x} L(x, \lambda) = 0$$

# The SVM Dual Solution

- We found the SVM dual problem can be written as:

$$\begin{aligned} \sup_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_j^T x_i \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \alpha_i \in \left[0, \frac{c}{n}\right] \quad i = 1, \dots, n. \end{aligned}$$

- Given solution  $\alpha^*$  to dual, primal solution is  $w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$ .
- The solution is in the space spanned by the inputs.
- Note  $\alpha_i^* \in [0, \frac{c}{n}]$ . So  $c$  controls max weight on each example. (**Robustness!**)
  - What's the relation between  $c$  and regularization?



## Complementary Slackness Conditions

- Recall our primal constraints and Lagrange multipliers:

Lagrange Multiplier	Constraint
$\lambda_j$	$-\xi_j \leq 0$
$\alpha_j$	$(1 - y_j f(x_j)) - \xi_j \leq 0$

- Recall first order condition  $\nabla_{\xi_j} L = 0$  gave us  $\lambda_j^* = \frac{c}{n} - \alpha_j^*$ .
- By strong duality, we must have **complementary slackness**:

$$\alpha_j^* (1 - y_j f^*(x_j) - \xi_j^*) = 0$$

$$\lambda_j^* \xi_j^* = \left( \frac{c}{n} - \alpha_j^* \right) \xi_j^* = 0$$

## Consequences of Complementary Slackness

By strong duality, we must have **complementary slackness**.

$$\begin{aligned}\alpha_i^* (1 - y_i f^*(x_i) - \xi_i^*) &= 0 \\ \left(\frac{c}{n} - \alpha_i^*\right) \xi_i^* &= 0\end{aligned}$$

Recall “**slack variable**”  $\xi_i^* = \max(0, 1 - y_i f^*(x_i))$  is the hinge loss on  $(x_i, y_i)$ .

- If  $y_i f^*(x_i) > 1$  then the margin loss is  $\xi_i^* = 0$ , and we get  $\alpha_i^* = 0$ .
- If  $y_i f^*(x_i) < 1$  then the margin loss is  $\xi_i^* > 0$ , so  $\alpha_i^* = \frac{c}{n}$ .
- If  $\alpha_i^* = 0$ , then  $\xi_i^* = 0$ , which implies no loss, so  $y_i f^*(x) \geq 1$ .
- If  $\alpha_i^* \in (0, \frac{c}{n})$ , then  $\xi_i^* = 0$ , which implies  $1 - y_i f^*(x_i) = 0$ .

## Complementary Slackness Results: Summary

If  $\alpha^*$  is a solution to the dual problem, then primal solution is

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i \quad \text{where } \alpha_i^* \in [0, \frac{c}{n}].$$

Relation between margin and example weights ( $\alpha_i$ 's):

$$\alpha_i^* = 0 \implies y_i f^*(x_i) \geq 1$$

$$\alpha_i^* \in (0, \frac{c}{n}) \implies y_i f^*(x_i) = 1$$

$$\alpha_i^* = \frac{c}{n} \implies y_i f^*(x_i) \leq 1$$

$$y_i f^*(x_i) < 1 \implies \alpha_i^* = \frac{c}{n}$$

$$y_i f^*(x_i) = 1 \implies \alpha_i^* \in [0, \frac{c}{n}]$$

$$y_i f^*(x_i) > 1 \implies \alpha_i^* = 0$$

- If  $\alpha^*$  is a solution to the dual problem, then primal solution is

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

with  $\alpha_i^* \in [0, \frac{c}{n}]$ .

- The  $x_i$ 's corresponding to  $\alpha_i^* > 0$  are called **support vectors**.
- Few margin errors or “on the margin” examples  $\implies$  **sparsity in input examples**.