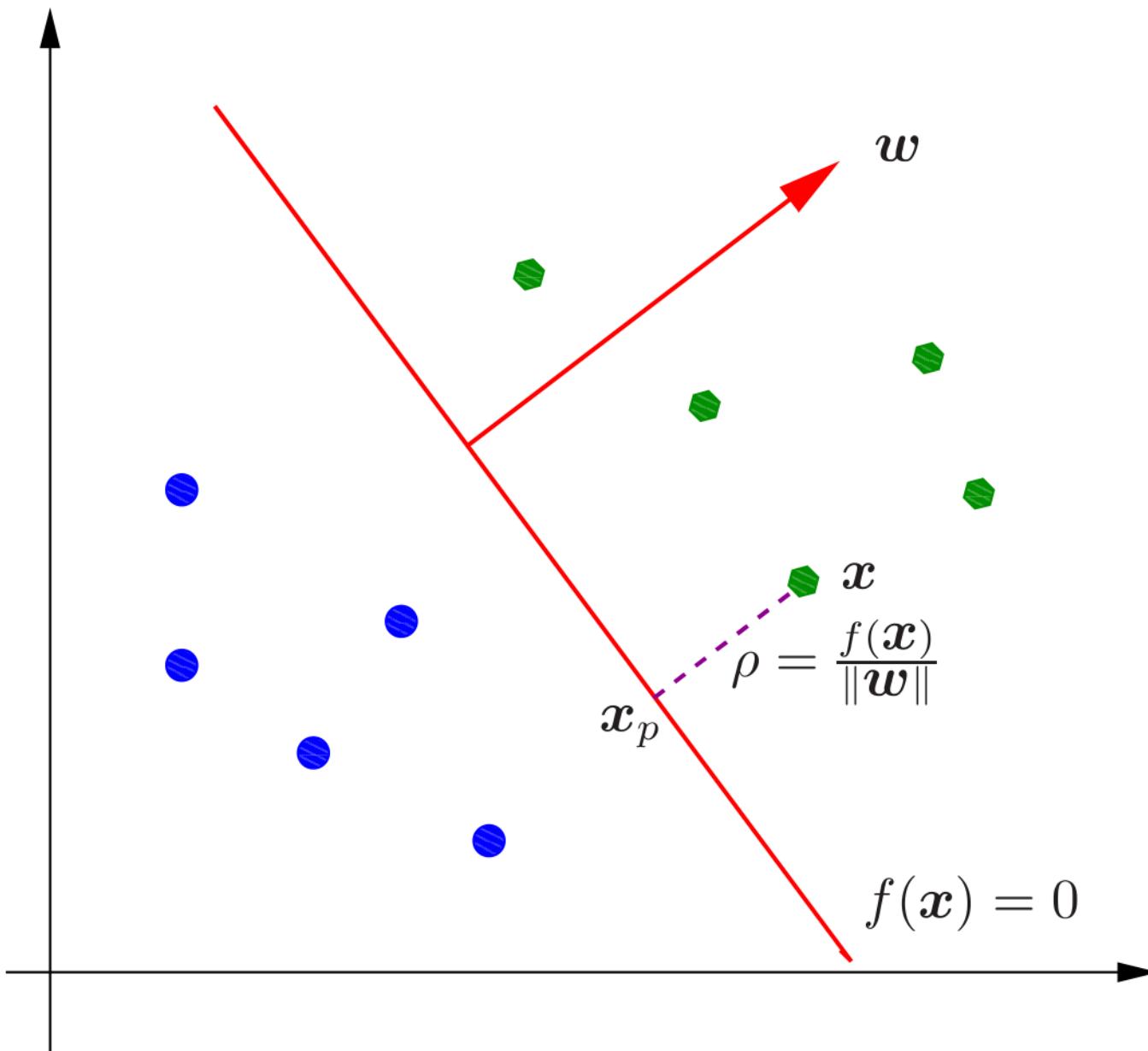


Lecture 10

Support Vector Machine (SVM)

- Separating Hyperplane and Binary Classification
- Margin and maximum margin classifier
- Support Vector
- Advanced topics
- Math foundations

Separating Hyperplane: Geometric View



Consider $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$, where a decision boundary is given by the hyperplane $f(\mathbf{x}) = 0$.

1. The weight vector \mathbf{w} is **normal** to the hyperplane.
2. $\rho = \frac{f(\mathbf{x})}{\|\mathbf{w}\|}$.

Why? \Rightarrow See next slide!

1. Suppose \mathbf{x}_1 and \mathbf{x}_2 are two points lying on the hyperplane, leading to

$$\mathbf{w}^\top \mathbf{x}_1 + b = 0,$$

$$\mathbf{w}^\top \mathbf{x}_2 + b = 0.$$

It follows from these equations that we have $\mathbf{w}^\top (\mathbf{x}_1 - \mathbf{x}_2) = 0$ which implies that \mathbf{w} is normal to the hyperplane.

2. Express $\mathbf{x} = \mathbf{x}_p + \rho \frac{\mathbf{w}}{\|\mathbf{w}\|}$. Note that $f(\mathbf{x}_p) = 0$. Then we have

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^\top \mathbf{x} + b \\ &= \mathbf{w}^\top \left(\mathbf{x}_p + \rho \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + b \\ &= \underbrace{\mathbf{w}^\top \mathbf{x}_p + b}_{0} + \rho \|\mathbf{w}\|. \end{aligned}$$

Binary Classification

- ▶ Training sample: $\mathcal{D} = \{(\mathbf{x}_t, y_t)\}_{t=1}^N$.
 - ▶ Instances: $\mathbf{x}_t \in \mathbb{R}^m$.
 - ▶ Labels: $y_t \in \{+1, -1\}$.
- ▶ Prediction function: $h : \mathcal{X} \mapsto \{+1, -1\}$, $h(\mathbf{x}) = \text{sign}(f(\mathbf{x}))$.
- ▶ Linear classification considers a linear discriminant function which has the form

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b,$$

where $(\mathbf{w}, b) \in \mathbb{R}^m \times \mathbb{R}$ ([weight vector](#), [bias](#)) are the parameters that control the function.

Margin

Definition

The **functional margin** of an example (\mathbf{x}_t, y_t) with respect to a hyperplane \mathbf{w} is defined as

$$\begin{aligned}\rho_t &= y_t f(\mathbf{x}_t) \\ &= y_t (\mathbf{w}^\top \mathbf{x}_t + b).\end{aligned}$$

Definition

The **geometric margin** is the functional margin for the normalized linear function $\frac{\mathbf{w}^\top}{\|\mathbf{w}\|} \mathbf{x} + \frac{b}{\|\mathbf{w}\|}$, i.e.,

$$\rho = y(f(\mathbf{x})/\|\mathbf{w}\|).$$

Canonical Form

For $\lambda \neq 0$, $(\lambda \mathbf{w}, \lambda b)$ describes the same hyperplane as (\mathbf{w}, b) , i.e.,

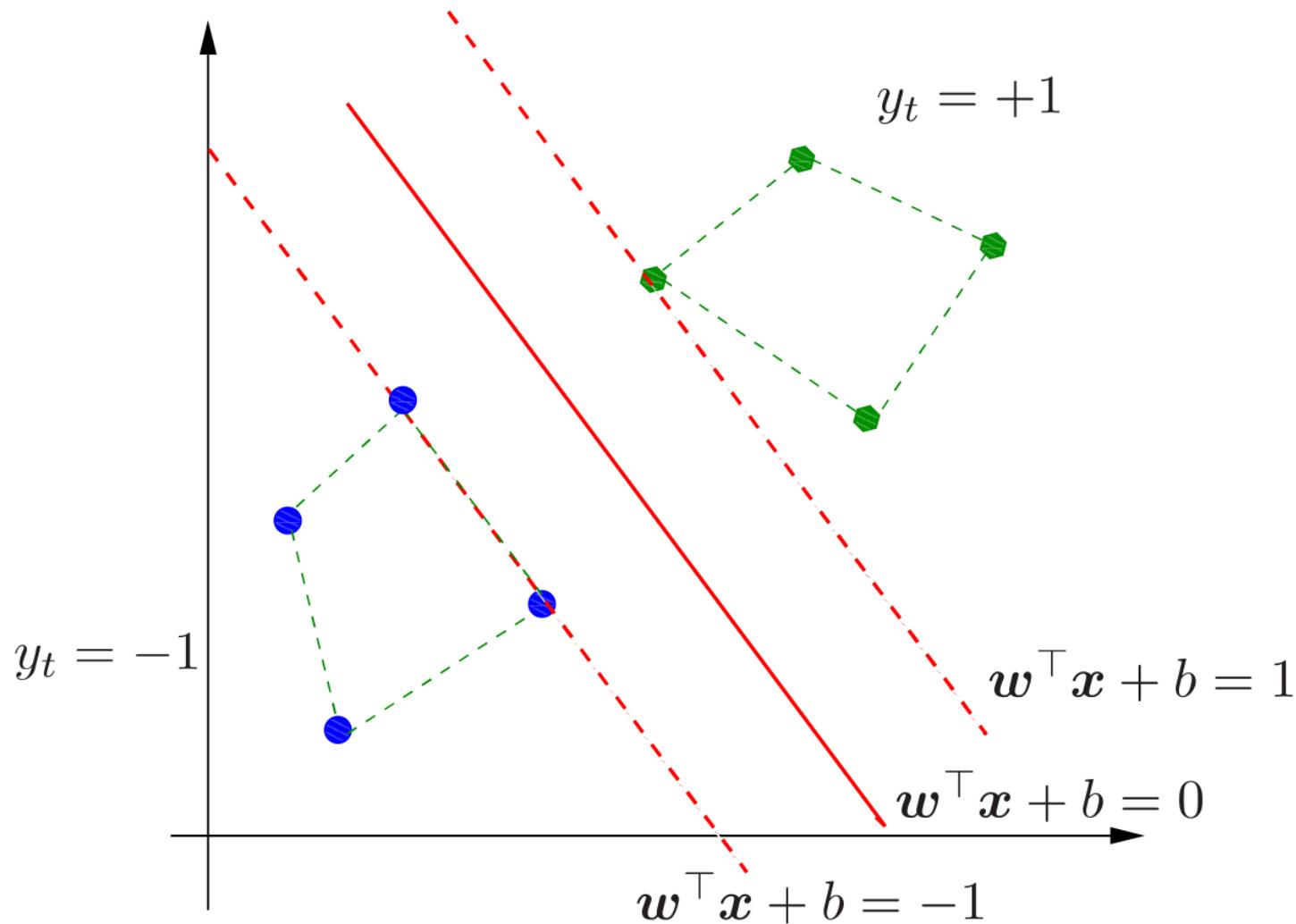
$$\{\mathbf{x} \mid \mathbf{w}^\top \mathbf{x} + b = 0\} = \{\mathbf{x} \mid \lambda (\mathbf{w}^\top \mathbf{x} + b) = 0\}.$$

Definition

The hyperplane is in [canonical form](#) if

$$\min_{\mathbf{x}_t \in \mathcal{X}} |\mathbf{w}^\top \mathbf{x}_t + b| = 1.$$

Canonical Optimal Hyperplane



$$\text{Margin} = \frac{1}{\|\mathbf{w}\|}$$

The geometric margin in previous slide is given by

$$\begin{aligned}\rho &= \frac{1}{2} \left\{ \frac{f(\mathbf{x}^+)}{\|\mathbf{w}\|} - \frac{f(\mathbf{x}^-)}{\|\mathbf{w}\|} \right\} \\ &= \frac{1}{2} \frac{1}{\|\mathbf{w}\|} \left\{ \underbrace{\mathbf{w}^\top \mathbf{x}^+ - \mathbf{w}^\top \mathbf{x}^-}_{2} \right\} \\ &= \frac{1}{\|\mathbf{w}\|}.\end{aligned}$$

Thus, maximizing margin is equivalent to minimizing the norm of the weight vector in the discriminant function.

Max Margin Classifier: Primal Form

Primal Problem

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \mathcal{J}(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} \quad & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, N \end{aligned}$$

Proposition

Given a linearly separable training sample $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, the hyperplane $\mathbf{w}^\top \mathbf{x}_i + b = 0$ that solves the above optimization problem realizes the maximal margin hyperplane with geometric margin $\rho = 1/\|\mathbf{w}\|$.

Primal Lagrangian

Primal Lagrangian $\mathcal{L}(\mathbf{w}, b, \alpha)$ is given by

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i (1 - y_i (\mathbf{w}^\top \mathbf{x}_i + b)).$$

$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0$ and $\frac{\partial \mathcal{L}}{\partial b} = 0$ yield respectively $\boxed{\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i}$ and

$\boxed{\sum_{i=1}^N \alpha_i y_i = 0}$. Substitute these relations into the primal form, leading to

$$\mathcal{L} = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^N \alpha_i.$$

Max Margin Classifier: Dual Form

Dual Problem

$$\begin{aligned} \max_{\alpha} \quad & \mathcal{G}(\alpha) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^N \alpha_i, \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0, \\ & \alpha_i \geq 0, \quad i = 1, \dots, N. \end{aligned}$$

Proposition

Suppose that α^* solves the dual problem, given a linearly separable training sample $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$. Then the weight vector $\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$ realizes the maximal margin hyperplane with geometric margin $\rho = 1/\|\mathbf{w}\|$.

Support Vectors

Note that $\boxed{\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i}$, i.e., \mathbf{w} is a linear combination of training data points \mathbf{x}_i .

It follows from KKT complimentary slackness condition that we have

$$\alpha_i [1 - y_i (\mathbf{w}^\top \mathbf{x}_i + b)] = 0.$$

- ▶ $y_i (\mathbf{w}^\top \mathbf{x}_i + b) \neq 1$ (x_i is not support vector) $\Rightarrow \alpha_i = 0$ (x_i is irrelevant).
- ▶ $\alpha_i \neq 0 \Rightarrow y_i (\mathbf{w}^\top \mathbf{x}_i + b) = 1$ (x_i is support vector).

Only support vectors influence the computation of \mathbf{w} .

Soft Margin Classifier: L_2 Norm

In cases where data are not linearly separable, the optimization problem cannot be solved as the primal has an empty feasible region and the dual an unbounded objective function.

To sidestep this problem, we introduce slack variables $\zeta_i \geq 0$ to allow the margin constraints to be violated:

$$y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \zeta_i, \quad i = 1, \dots, N.$$

Primal Problem

$$\min_{\mathbf{w}, b, \zeta} \quad \mathcal{J}(\mathbf{w}, \zeta) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \zeta_i^2,$$

$$\text{subject to} \quad y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \zeta_i, \quad i = 1, \dots, N.$$

- The primal Lagrangian is

$$\mathcal{L}(\mathbf{w}, b, \zeta) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{2} \sum_{i=1}^N \zeta_i^2 + \sum_{i=1}^N \alpha_i (1 - \zeta_i - y_i (\mathbf{w}^\top \mathbf{x}_i + b)).$$

- $\frac{\partial \mathcal{J}(\mathbf{w}, \zeta)}{\partial \mathbf{w}} = 0$ leads to $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$.
- $\frac{\partial \mathcal{J}(\mathbf{w}, \zeta)}{\partial b} = 0$ leads to $\sum_{i=1}^N \alpha_i y_i = 0$.
- $\frac{\partial \mathcal{J}(\mathbf{w}, \zeta)}{\partial \zeta} = 0$ leads to $\zeta = \frac{\alpha}{C}$.
- Substitute these relations into the primal Lagrangian to obtain the dual Lagrangian function:

$$\mathcal{G}(\alpha) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \frac{1}{2C} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N \alpha_i.$$

Soft Margin Classifier: Dual Form

Dual Problem

$$\begin{aligned} \max_{\alpha} \quad & \mathcal{G}(\alpha) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \left(\mathbf{x}_i^\top \mathbf{x}_j + \frac{1}{C} \delta_{ij} \right) + \sum_{i=1}^N \alpha_i, \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0, \\ & \alpha_i \geq 0, \quad i = 1, \dots, N. \end{aligned}$$

Soft Margin Classifier: L_1 Norm

Primal Problem

$$\begin{aligned} \min_{\mathbf{w}, b, \zeta} \quad & \mathcal{J}(\mathbf{w}, \zeta) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \zeta_i, \\ \text{subject to} \quad & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \zeta_i, \quad i = 1, \dots, N, \\ & \zeta_i \geq 0, \quad i = 1, \dots, N. \end{aligned}$$

The primal Lagrangian is

$$\mathcal{L}(\mathbf{w}, b, \zeta) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{2} \sum_{i=1}^N \zeta_i^2 + \sum_{i=1}^N \alpha_i (1 - \zeta_i - y_i (\mathbf{w}^\top \mathbf{x}_i + b)) - \sum_{i=1}^N \beta_i \zeta_i.$$

Soft Margin Classifier: Dual Form

Dual Problem

$$\max_{\alpha} \quad \mathcal{G}(\alpha) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^N \alpha_i,$$

subject to

$$\begin{aligned} & \sum_{i=1}^N \alpha_i y_i = 0, \\ & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, N. \end{aligned}$$

Regularization View

Consider a linear discriminant function

$$f(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i + b,$$

where labels are $y_i \in \{1, -1\}$

Regularized learning problem involves:

$$\min_f \quad \underbrace{\frac{1}{N} \sum_{i=1}^N V(y_i, f(\mathbf{x}_i))}_{\text{empirical risk}} + \lambda \|\mathbf{w}\|^2,$$

where $V(y_i, f(\mathbf{x}_i))$ is a loss function.

Why Regularization

- ▶ The problem of approximating a function from sparse data is **ill-posed** and a classical way to solve it is **regularization**.
- ▶ For a finite set of training examples the search for the best model or approximating function has to be constrained on an approximately **small hypothesis space** (which can be thought of as a space of model).
- ▶ If the space is too large, models can be found which will fit exactly the data but will have a poor generalization performance , that is poor predictive capability on new data.
- ▶ Vapnik's theory characterize these concepts in terms of **capacity** of a set of functions and **capacity control** depending on the training data.
- ▶ In the case of regularization, a form of capacity control leads to choosing an optimal regularization parameter λ .

Loss Functions

- ▶ L_2 Regularization Networks

$$V(y_i, f(\mathbf{x}_i)) = (y_i - f(\mathbf{x}_i))^2.$$

- ▶ Support Vector Regression

$$V(y_i, f(\mathbf{x}_i)) = |y_i - f(\mathbf{x}_i)|_\epsilon,$$

where $|\cdot|_\epsilon$ is epsilon-insensitive norm.

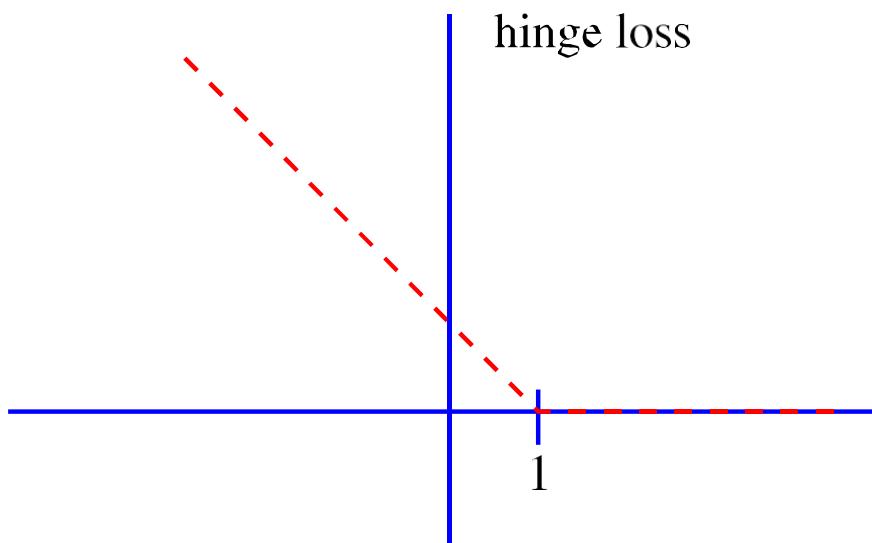
- ▶ Support Vector Machine

$$V(y_i, f(\mathbf{x}_i)) = [1 - yf(\mathbf{x})]_+,$$

where $[1 - yf(\mathbf{x})]_+$ is hinge loss (see the next slide).

SVM arises by considering **hinge loss** in the regularization framework:

$$V(y, f(\mathbf{x})) = [1 - yf(\mathbf{x})]_+ = \max(1 - yf(\mathbf{x}), 0).$$



The problem is not differentiable, so introduce slack variables ξ_i to make the problem easier to work with

$$\begin{aligned} \min_f \quad & \frac{1}{N} \xi_i + \lambda \|\mathbf{w}\|^2, \\ \text{subject to} \quad & y_i f(\mathbf{x}_i) \geq 1 - \xi_i, \quad i = 1, \dots, N, \\ & \xi_i \geq 0, \end{aligned}$$

which is equivalent to soft-margin SVM.

Regularization in RKHS

Regularized learning problem involves:

$$\min_{f \in \mathcal{H}} \underbrace{\frac{1}{N} \sum_{i=1}^N V(y_i, f(\mathbf{x}_i))}_{\text{empirical risk}} + \lambda \|f\|_K^2,$$

where $\|f\|_K^2$ is a norm in a Reproducing Kernel Hilbert Space \mathcal{H} define by the positive definite function κ .

Introducing slack variables ξ_i ,

$$\begin{aligned} \min_f \quad & \frac{1}{N} \xi_i + \lambda \|f\|_K^2, \\ \text{subject to} \quad & y_i f(\mathbf{x}_i) \geq 1 - \xi_i, \quad i = 1, \dots, N, \\ & \xi_i \geq 0. \end{aligned}$$

Apply the representer theorem

$$f^*(\mathbf{x}) = \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i),$$

leading to

$$\|f\|_K^2 = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha},$$

where $\mathbf{K} = [\kappa(\mathbf{x}_i, \mathbf{x}_j)] \in \mathbb{R}^{N \times N}$ and $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]^\top \in \mathbb{R}^N$.

Then we have the primal form for SVM

$$\begin{aligned} \min_{f \in \mathcal{H}} \quad & \frac{1}{N} \xi_i + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}, \\ \text{subject to} \quad & y_i \left(\sum_{j=1}^N \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) + b \right) \geq 1 - \xi_i, \quad i = 1, \dots, N, \\ & \xi_i \geq 0, \end{aligned}$$

where an unregularized bias term b is added.

Quick Overview of Constrained Optimization

Constrained Optimization

Consider a [primal form](#) of constrained optimization problem:

$$\begin{aligned} & \text{minimize} && \mathcal{J}(\mathbf{w}), \\ & \text{subject to} && g_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, M, \\ & && h_j(\mathbf{w}) = 0, \quad j = 1, \dots, L. \end{aligned}$$

Lagrangian is given by

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\nu}, \boldsymbol{\lambda}) = \mathcal{J}(\mathbf{w}) + \sum_{i=1}^M \nu_i g_i(\mathbf{w}) + \sum_{j=1}^L \lambda_j h_j(\mathbf{w}),$$

where $\nu_i > 0$ for $i = 1, \dots, M$ and λ_j for $j = 1, \dots, L$ are unrestricted in sign.

Karush-Kuhn-Tucker (KKT) Necessary Conditions

1. Optimality

$$\nabla \mathcal{L} = \nabla \mathcal{J}(\mathbf{w}) + \sum_{i=1}^M \nu_i \nabla g_i(\mathbf{w}) + \sum_{j=1}^L \lambda_j \nabla h_j(\mathbf{w}) = 0.$$

2. Feasibility

$$\begin{aligned} g_i(\mathbf{w}) &\leq 0, \quad i = 1, \dots, M \\ h_j(\mathbf{w}) &= 0, \quad j = 1, \dots, L. \end{aligned}$$

3. Complementary slackness

$$\nu_i g_i(\mathbf{w}) = 0, \quad i = 1, \dots, M \quad (\nu_i > 0).$$

Lagrangian Dual Problem

Lagrangian dual function is defined as the minimum value of the Lagrangian over \mathbf{w} :

$$\mathcal{G}(\boldsymbol{\nu}, \boldsymbol{\lambda}) = \inf_{\mathbf{w} \in \mathcal{W}} \mathcal{L}(\mathbf{w}, \boldsymbol{\nu}, \boldsymbol{\lambda}).$$

The Lagrangian dual problem has the form:

$$\begin{aligned} & \text{maximize} && \mathcal{G}(\boldsymbol{\nu}, \boldsymbol{\lambda}), \\ & \text{subject to} && \nu_i \geq 0, \quad i = 1, \dots, M. \end{aligned}$$

Weak duality: $\mathcal{J}(\mathbf{w}) \geq \mathcal{G}(\boldsymbol{\nu}, \boldsymbol{\lambda})$ if \mathbf{w} is a feasible solution of the primal problem and $(\boldsymbol{\nu}, \boldsymbol{\lambda})$ is a feasible solution of the dual problem.