

Introduction to Machine Learning

Lecture 12 Classification - Support Vector Machine

Hongteng Xu



中國人民大學
RENMIN UNIVERSITY OF CHINA

高瓴人工智能學院
Gaoling School of Artificial Intelligence

Outline

Review

- ▶ Classification, definition, evaluation
- ▶ Linear classifiers (Naïve Bayes classifier, Linear discriminant analysis, Logistic regression)

Outline

Review

- ▶ Classification, definition, evaluation
- ▶ Linear classifiers (Naïve Bayes classifier, Linear discriminant analysis, Logistic regression)

Today

- ▶ Support-vector machine (SVM)

Outline

Review

- ▶ Classification, definition, evaluation
- ▶ Linear classifiers (Naïve Bayes classifier, Linear discriminant analysis, Logistic regression)

Today

- ▶ Support-vector machine (SVM)
- ▶ Take it easy. Most of the following content can be found at Wikipedia:)
- ▶ https://en.wikipedia.org/wiki/Support-vector_machine

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)}$$

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} > T(=0)$$

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} > T(=0) \quad \Rightarrow \quad \mathbf{w}^T \mathbf{x} > c$$

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} > T(=0) \quad \Rightarrow \quad \mathbf{w}^T \mathbf{x} > c \quad (1)$$

(What is T for Naïve Bayes?)

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} > T(=0) \quad \Rightarrow \quad \mathbf{w}^T \mathbf{x} > c \quad (1)$$

(What is T for Naïve Bayes?)

Logistic regression

$$p(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^T \boldsymbol{\beta})}$$

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} > T(=0) \quad \Rightarrow \quad \mathbf{w}^T \mathbf{x} > c \quad (1)$$

(What is T for Naïve Bayes?)

Logistic regression

$$p(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^T \boldsymbol{\beta})} > 0.5$$

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} > T(=0) \quad \Rightarrow \quad \mathbf{w}^T \mathbf{x} > c \quad (1)$$

(What is T for Naïve Bayes?)

Logistic regression

$$p(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^T \boldsymbol{\beta})} > 0.5 \quad \Rightarrow \quad \mathbf{x}^T \boldsymbol{\beta} > 0$$

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} > T(=0) \quad \Rightarrow \quad \mathbf{w}^T \mathbf{x} > c \quad (1)$$

(What is T for Naïve Bayes?)

Logistic regression

$$p(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^T \boldsymbol{\beta})} > 0.5 \quad \Rightarrow \quad \mathbf{x}^T \boldsymbol{\beta} > 0 \quad (2)$$

Essentially, they aim at finding a boundary/hyperplane to separate the samples of different classes.

Revisit The Principle of Typical Classifiers

Linear determinant analysis

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} > T(=0) \quad \Rightarrow \quad \mathbf{w}^T \mathbf{x} > c \quad (1)$$

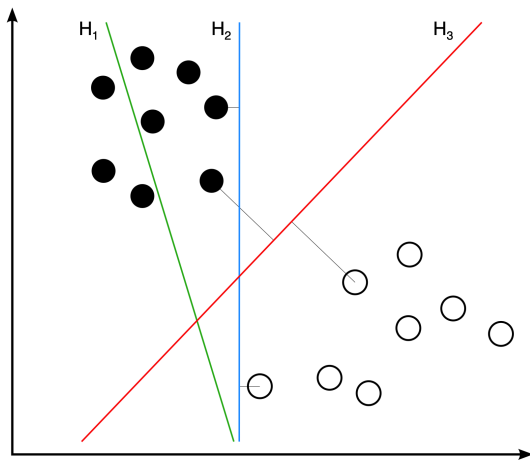
(What is T for Naïve Bayes?)

Logistic regression

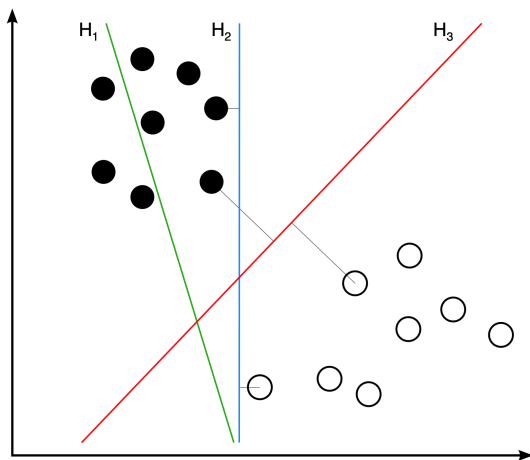
$$p(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^T \boldsymbol{\beta})} > 0.5 \quad \Rightarrow \quad \mathbf{x}^T \boldsymbol{\beta} > 0 \quad (2)$$

Essentially, they aim at finding a boundary/hyperplane to separate the samples of different classes. (Recall the analytic algebra you learned in your high school.)

Ambiguity of Decision Boundary/Hyperplane



Ambiguity of Decision Boundary/Hyperplane



How to find the best decision boundary? What is the criterion?

Support-Vector Machine (SVM)

Motivation:

- ▶ Find a hyperplane $\mathbf{x}^T \mathbf{w} = b$ with maximum-margin/largest separation,

Support-Vector Machine (SVM)

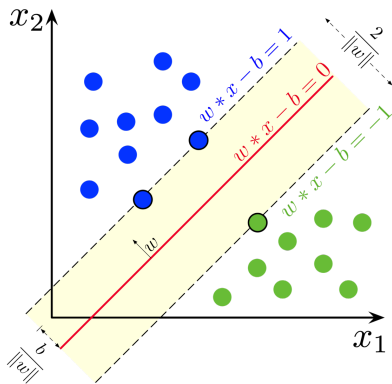
Motivation:

- Find a hyperplane $\mathbf{x}^T \mathbf{w} = b$ with maximum-margin/largest separation, such that the discriminative power is maximized, error risk is minimized.

Support-Vector Machine (SVM)

Motivation:

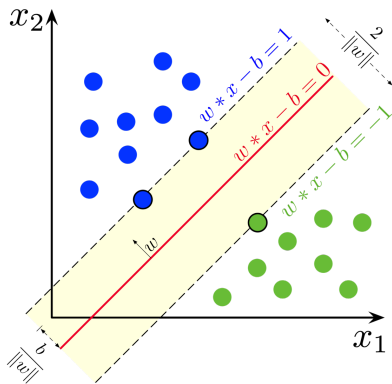
- Find a hyperplane $\mathbf{x}^T \mathbf{w} = b$ with maximum-margin/largest separation, such that the discriminative power is maximized, error risk is minimized.



Support-Vector Machine (SVM)

Principle:

- The desired hyperplane has the largest distance to the nearest training-data point of any class (so-called functional margin)



Implementation of SVM (Linearly Separable 2-Class)

- ▶ A set of training data $\{\mathbf{x}_n, y_n\}_{n=1}^N$, and each $y_n \in \{-1, 1\}$

Implementation of SVM (Linearly Separable 2-Class)

- ▶ A set of training data $\{\mathbf{x}_n, y_n\}_{n=1}^N$, and each $y_n \in \{-1, 1\}$
- ▶ Any hyperplane is the set of the points \mathbf{x} satisfying $\mathbf{w}^T \mathbf{x} - b = 0$.

Implementation of SVM (Linearly Separable 2-Class)

- ▶ A set of training data $\{\mathbf{x}_n, y_n\}_{n=1}^N$, and each $y_n \in \{-1, 1\}$
- ▶ Any hyperplane is the set of the points \mathbf{x} satisfying $\mathbf{w}^T \mathbf{x} - b = 0$.
- ▶ \mathbf{w} : the normal vector
- ▶ $\frac{b}{\|\mathbf{w}\|_2}$: the offset of the hyperplane shifting from the original along the normal vector. (Derive Them)

Implementation of SVM (Linearly Separable 2-Class)

- ▶ A set of training data $\{\mathbf{x}_n, y_n\}_{n=1}^N$, and each $y_n \in \{-1, 1\}$
- ▶ Any hyperplane is the set of the points \mathbf{x} satisfying $\mathbf{w}^T \mathbf{x} - b = 0$.
- ▶ \mathbf{w} : the normal vector
- ▶ $\frac{b}{\|\mathbf{w}\|_2}$: the offset of the hyperplane shifting from the original along the normal vector. (Derive Them)
- ▶ Shift the hyperplane along two opposite directions and define classification criteria (margins):

$$\mathbf{w}^T \mathbf{x} - b = 1 \quad \Rightarrow \quad \text{Any } \mathbf{x} \text{ on or above it is labeled by } 1$$

$$\mathbf{w}^T \mathbf{x} - b = -1 \quad \Rightarrow \quad \text{Any } \mathbf{x} \text{ on or below it is labeled by } -1.$$

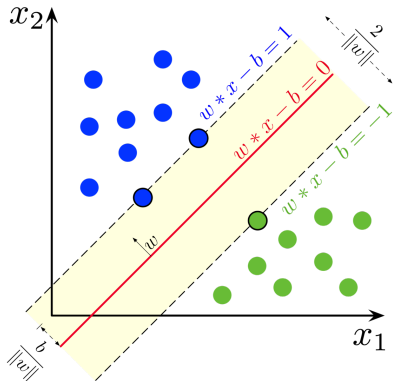
Implementation of SVM (Linearly Separable 2-Class)

- ▶ A set of training data $\{\mathbf{x}_n, y_n\}_{n=1}^N$, and each $y_n \in \{-1, 1\}$
- ▶ Any hyperplane is the set of the points \mathbf{x} satisfying $\mathbf{w}^T \mathbf{x} - b = 0$.
- ▶ \mathbf{w} : the normal vector
- ▶ $\frac{b}{\|\mathbf{w}\|_2}$: the offset of the hyperplane shifting from the original along the normal vector. (Derive Them)
- ▶ Shift the hyperplane along two opposite directions and define classification criteria (margins):

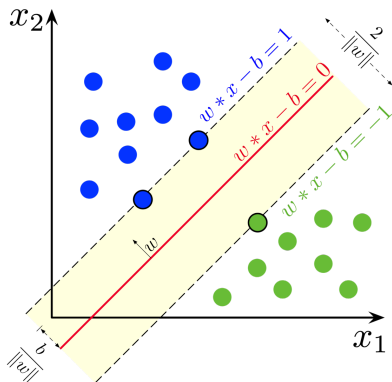
$$\begin{aligned}\mathbf{w}^T \mathbf{x} - b = 1 &\Rightarrow \text{Any } \mathbf{x} \text{ on or above it is labeled by } 1 \\ \mathbf{w}^T \mathbf{x} - b = -1 &\Rightarrow \text{Any } \mathbf{x} \text{ on or below it is labeled by } -1.\end{aligned}\tag{3}$$

- ▶ $\frac{2}{\|\mathbf{w}\|_2}$: The distance between these two margins. (Derive It)

Implementation of SVM (Linearly Separable 2-Class)



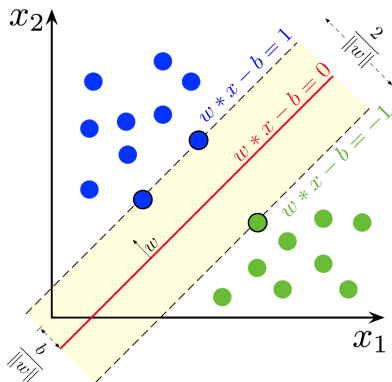
Implementation of SVM (Linearly Separable 2-Class)



For each \mathbf{x}_n, y_n , the data point \mathbf{x}_n must lie on the correct side of the margin.

$$\mathbf{w}^T \mathbf{x}_n - b \begin{cases} \geq 1 & \text{if } y_n = 1, \\ \leq -1 & \text{if } y_n = -1. \end{cases}$$

Implementation of SVM (Linearly Separable 2-Class)



For each \mathbf{x}_n, y_n , the data point \mathbf{x}_n must lie on the correct side of the margin.

$$\mathbf{w}^T \mathbf{x}_n - b \begin{cases} \geq 1 & \text{if } y_n = 1, \\ \leq -1 & \text{if } y_n = -1. \end{cases} \Rightarrow y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1, \forall n = 1, \dots, N. \quad (4)$$

Implementation of SVM (Linearly Separable 2-Class)

Principle:

- ▶ Maximize the distance between margins, i.e., $\max \frac{2}{\|\mathbf{w}\|_2}$

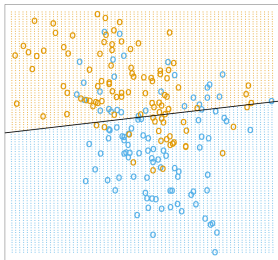
Implementation of SVM (Linearly Separable 2-Class)

Principle:

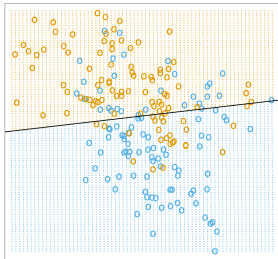
- ▶ Maximize the distance between margins, i.e., $\max \frac{2}{\|\mathbf{w}\|_2}$
- ▶ Equivalently, the problem becomes

$$\begin{aligned} \min_{w,b} \quad & \|\mathbf{w}\|_2 \\ \text{s.t.} \quad & y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1, \quad \forall n = 1, \dots, N \end{aligned} \tag{5}$$

How To Deal with Fuzzy Decision Boundary/Hyperplane?



How To Deal with Fuzzy Decision Boundary/Hyperplane?

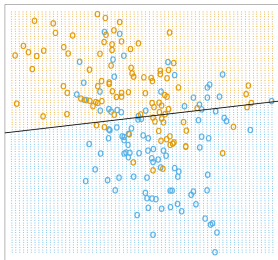


Soft-margin SVM: for noisy linearly-separable classes

$$\min_{\mathbf{w}, b, \xi} \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n$$

$$s.t. \ y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n, \ \xi_n \geq 0, \ \forall n = 1, \dots, N$$

How To Deal with Fuzzy Decision Boundary/Hyperplane?



Soft-margin SVM: for noisy linearly-separable classes

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n, \quad \xi_n \geq 0, \quad \forall n = 1, \dots, N \end{aligned} \tag{6}$$

where $\xi_n = \max\{0, 1 - y_n(\mathbf{w}^T \mathbf{x}_n - b)\}$, i.e., the smallest nonnegative number satisfying $y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n$.

Solve Soft-margin SVM (Primal Problem)

- The Primal Problem of Soft-margin SVM:

$$\begin{aligned} \min_{w,b,\xi} \quad & \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n, \quad \xi_n \geq 0, \quad \forall n = 1, \dots, N \end{aligned}$$

Solve Soft-margin SVM (Primal Problem)

- The Primal Problem of Soft-margin SVM:

$$\begin{aligned} \min_{w,b,\xi} \quad & \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n, \quad \xi_n \geq 0, \quad \forall n = 1, \dots, N \end{aligned} \tag{7}$$

- Plug the definition $\xi_n = \max\{0, 1 - y_n(\mathbf{w}^T \mathbf{x}_n - b)\}$ into (7):

$$\min_{w,b} \quad \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{Tikhonov reg.}} + \frac{1}{N} \sum_{n=1}^N \underbrace{\max\{0, 1 - y_n(\mathbf{w}^T \mathbf{x}_n - b)\}}_{\text{Hinge loss}}$$

Solve Soft-margin SVM (Primal Problem)

- ▶ The Primal Problem of Soft-margin SVM:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n, \quad \xi_n \geq 0, \quad \forall n = 1, \dots, N \end{aligned} \tag{7}$$

- ▶ Plug the definition $\xi_n = \max\{0, 1 - y_n(\mathbf{w}^T \mathbf{x}_n - b)\}$ into (7):

$$\min_{\mathbf{w}, b} \quad \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{Tikhonov reg.}} + \frac{1}{N} \sum_{n=1}^N \underbrace{\max\{0, 1 - y_n(\mathbf{w}^T \mathbf{x}_n - b)\}}_{\text{Hinge loss}} \tag{8}$$

- ▶ The objective function is convex (but non-smooth) w.r.t. \mathbf{w} and b , so GD or SGD can be applied.

Solve Soft-margin SVM (Primal Problem)

- ▶ The Primal Problem of Soft-margin SVM:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n, \quad \xi_n \geq 0, \quad \forall n = 1, \dots, N \end{aligned} \tag{7}$$

- ▶ Plug the definition $\xi_n = \max\{0, 1 - y_n(\mathbf{w}^T \mathbf{x}_n - b)\}$ into (7):

$$\min_{\mathbf{w}, b} \quad \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{Tikhonov reg.}} + \frac{1}{N} \sum_{n=1}^N \underbrace{\max\{0, 1 - y_n(\mathbf{w}^T \mathbf{x}_n - b)\}}_{\text{Hinge loss}} \tag{8}$$

- ▶ The objective function is convex (but non-smooth) w.r.t. \mathbf{w} and b , so GD or SGD can be applied.
- ▶ Derive the gradient of Hinge loss

Solve Soft-margin SVM (Dual Problem)

Because the Primal problem is convex optimization, we can solve it equivalently via solving its dual problem.

- Recall the Primal Problem of Soft-margin SVM:

$$\begin{aligned} \min_{w,b,\xi} \quad & \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & \underbrace{y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n}_{\text{Constraint 1}}, \quad \underbrace{\xi_n \geq 0}_{\text{Constraint 2}}, \quad \forall n = 1, \dots, N \end{aligned}$$

Solve Soft-margin SVM (Dual Problem)

Because the Primal problem is convex optimization, we can solve it equivalently via solving its dual problem.

- Recall the Primal Problem of Soft-margin SVM:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & \underbrace{y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n}_{\text{Constraint 1}}, \quad \underbrace{\xi_n \geq 0}_{\text{Constraint 2}}, \quad \forall n = 1, \dots, N \end{aligned} \tag{9}$$

- The Lagrangian dual of (9): Introduce dual variables $\{c_n\}_{n=1}^N$ for the first N constraints, and $\{\tau_n\}_{n=1}^N$ for the second N constraints

$$\begin{aligned} \max_{\{c_n, \tau_n\}_{n=1}^N} \quad & \left(\min_{\mathbf{w}, b, \xi} \lambda \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{N} \xi_n + c_n(1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n - b)) - \tau_n \xi_n \right) \right) \\ \text{s.t.} \quad & c_n \geq 0, \tau_n \geq 0, \quad \forall n = 1, \dots, N \end{aligned}$$

Solve Soft-margin SVM (Dual Problem)

Because the Primal problem is convex optimization, we can solve it equivalently via solving its dual problem.

- Recall the Primal Problem of Soft-margin SVM:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \lambda \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & \underbrace{y_n(\mathbf{w}^T \mathbf{x}_n - b) \geq 1 - \xi_n}_{\text{Constraint 1}}, \quad \underbrace{\xi_n \geq 0}_{\text{Constraint 2}}, \quad \forall n = 1, \dots, N \end{aligned} \quad (9)$$

- The Lagrangian dual of (9): Introduce dual variables $\{c_n\}_{n=1}^N$ for the first N constraints, and $\{\tau_n\}_{n=1}^N$ for the second N constraints

$$\max_{\{c_n, \tau_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \lambda \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{N} \xi_n + c_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n - b)) - \tau_n \xi_n \right) \right) \quad (10)$$

$$\text{s.t. } c_n \geq 0, \tau_n \geq 0, \forall n = 1, \dots, N$$

- Looks terrible... but simple actually;)

Solve Soft-margin SVM (Dual Problem)

Original dual problem:

$$\max_{\{c_n, \tau_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \lambda \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{N} \xi_n + \textcolor{red}{c}_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n - b)) - \textcolor{blue}{\tau}_n \xi_n \right) \right) \quad (11)$$

$$s.t. \ c_n \geq 0, \ \tau_n \geq 0, \ \forall n = 1, \dots, N$$

- **Firstly, we can ignore $\{\tau_n\}_{n=1}^N$ because $\tau_n^* = 0$ (Why?)**

Solve Soft-margin SVM (Dual Problem)

Original dual problem:

$$\begin{aligned} \max_{\{c_n, \tau_n\}_{n=1}^N} & \left(\min_{\mathbf{w}, b, \xi} \lambda \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{N} \xi_n + \mathbf{c}_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n - b)) - \tau_n \xi_n \right) \right) \\ \text{s.t. } & c_n \geq 0, \tau_n \geq 0, \forall n = 1, \dots, N \end{aligned} \quad (11)$$

- ▶ **Firstly, we can ignore $\{\tau_n\}_{n=1}^N$ because $\tau_n^* = 0$ (Why?)**
 - ▶ $\xi_n \geq 0$ and $\tau_n \geq 0$
 - ▶ Objective function: $\max_{\tau_n \geq 0} \min_{\xi_n \geq 0} -\tau_n \xi_n = 0$ when $\tau_n^* = 0$

Solve Soft-margin SVM (Dual Problem)

Original dual problem:

$$\begin{aligned} \max_{\{c_n, \tau_n\}_{n=1}^N} & \left(\min_{\mathbf{w}, b, \xi} \lambda \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{N} \xi_n + \mathbf{c}_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n - b)) - \tau_n \xi_n \right) \right) \\ \text{s.t. } & c_n \geq 0, \tau_n \geq 0, \forall n = 1, \dots, N \end{aligned} \quad (11)$$

► **Firstly, we can ignore $\{\tau_n\}_{n=1}^N$ because $\tau_n^* = 0$ (Why?)**

► $\xi_n \geq 0$ and $\tau_n \geq 0$

► Objective function: $\max_{\tau_n \geq 0} \min_{\xi_n \geq 0} -\tau_n \xi_n = 0$ when $\tau_n^* = 0$

► So, rewrite (11) as

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} & \left(\min_{\mathbf{w}, b, \xi} \lambda \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{N} \xi_n + \mathbf{c}_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n - b)) \right) \right) \\ \text{s.t. } & c_n \geq 0, \forall n = 1, \dots, N \end{aligned} \quad (12)$$

Solve Soft-margin SVM (Dual Problem)

Modified dual problem (Version 1):

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} & \left(\min_{\mathbf{w}, b, \xi} \lambda \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{N} \xi_n + \mathbf{c}_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n - b)) \right) \right) \\ \text{s.t. } & c_n \geq 0, \forall n = 1, \dots, N \end{aligned} \quad (13)$$

► **Secondly, the problem can be re-scale to**

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} & \left(\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{2N\lambda} \xi_n + \mathbf{c}_n (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{x}_n - b)) \right) \right) \\ \text{s.t. } & c_n \geq 0, \forall n = 1, \dots, N \end{aligned} \quad (14)$$

(Why?)

Solve Soft-margin SVM (Dual Problem)

Modified dual problem (Version 2):

$$\max_{\{c_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{2N\lambda} \xi_n + c_n (1 - \xi_n - y_n (\mathbf{w}^T \mathbf{x}_n - b)) \right)}_{L(\mathbf{w}, b, \xi, \mathbf{c})} \right) \quad (15)$$

$$s.t. \ c_n \geq 0, \ \forall n = 1, \dots, N$$

► **Thirdly, consider the optimality condition of inner problem:**

$$\frac{\partial L}{\partial \mathbf{w}} = 0$$

Solve Soft-margin SVM (Dual Problem)

Modified dual problem (Version 2):

$$\max_{\{c_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{2N\lambda} \xi_n + c_n (1 - \xi_n - y_n (\mathbf{w}^T \mathbf{x}_n - b)) \right)}_{L(\mathbf{w}, b, \xi, \mathbf{c})} \right) \quad (15)$$

$$s.t. \ c_n \geq 0, \ \forall n = 1, \dots, N$$

► **Thirdly, consider the optimality condition of inner problem:**

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w}^* = \sum_{n=1}^N c_n y_n \mathbf{x}_n \quad \text{Relation between opt. primal and dual}$$

Solve Soft-margin SVM (Dual Problem)

Modified dual problem (Version 2):

$$\max_{\{c_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{2N\lambda} \xi_n + c_n (1 - \xi_n - y_n (\mathbf{w}^T \mathbf{x}_n - b)) \right)}_{L(\mathbf{w}, b, \xi, \mathbf{c})} \right) \quad (15)$$

$$s.t. \ c_n \geq 0, \ \forall n = 1, \dots, N$$

► **Thirdly, consider the optimality condition of inner problem:**

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w}^* = \sum_{n=1}^N c_n y_n \mathbf{x}_n \quad \text{Relation between opt. primal and dual}$$
$$\frac{dL}{db} = 0$$

Solve Soft-margin SVM (Dual Problem)

Modified dual problem (Version 2):

$$\max_{\{c_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{2N\lambda} \xi_n + c_n (1 - \xi_n - y_n (\mathbf{w}^T \mathbf{x}_n - b)) \right)}_{L(\mathbf{w}, b, \xi, \mathbf{c})} \right) \quad (15)$$

$$s.t. \ c_n \geq 0, \ \forall n = 1, \dots, N$$

► **Thirdly, consider the optimality condition of inner problem:**

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w}^* = \sum_{n=1}^N c_n y_n \mathbf{x}_n \quad \text{Relation between opt. primal and dual}$$

$$\frac{dL}{db} = 0 \quad \Rightarrow \quad \sum_{n=1}^N c_n y_n = 0 \quad \text{Equality constraints of dual variables}$$

Solve Soft-margin SVM (Dual Problem)

Modified dual problem (Version 2):

$$\max_{\{c_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{2N\lambda} \xi_n + c_n (1 - \xi_n - y_n (\mathbf{w}^T \mathbf{x}_n - b)) \right)}_{L(\mathbf{w}, b, \xi, \mathbf{c})} \right) \quad (15)$$

$$s.t. \ c_n \geq 0, \ \forall n = 1, \dots, N$$

► **Thirdly, consider the optimality condition of inner problem:**

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w}^* = \sum_{n=1}^N c_n y_n \mathbf{x}_n \quad \text{Relation between opt. primal and dual}$$

$$\frac{dL}{db} = 0 \quad \Rightarrow \quad \sum_{n=1}^N c_n y_n = 0 \quad \text{Equality constraints of dual variables}$$

$$\frac{dL}{d\xi_n}$$

Solve Soft-margin SVM (Dual Problem)

Modified dual problem (Version 2):

$$\max_{\{c_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{2N\lambda} \xi_n + c_n (1 - \xi_n - y_n (\mathbf{w}^T \mathbf{x}_n - b)) \right)}_{L(\mathbf{w}, b, \xi, \mathbf{c})} \right) \quad (15)$$

$$\text{s.t. } c_n \geq 0, \forall n = 1, \dots, N$$

► **Thirdly, consider the optimality condition of inner problem:**

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \quad \Rightarrow \quad \mathbf{w}^* = \sum_{n=1}^N c_n y_n \mathbf{x}_n \quad \text{Relation between opt. primal and dual}$$

$$\frac{dL}{db} = 0 \quad \Rightarrow \quad \sum_{n=1}^N c_n y_n = 0 \quad \text{Equality constraints of dual variables}$$

$$\frac{dL}{d\xi_n} = \frac{1}{2N\lambda} - c_n$$

Solve Soft-margin SVM (Dual Problem)

Modified dual problem (Version 2):

$$\max_{\{c_n\}_{n=1}^N} \left(\min_{\mathbf{w}, b, \xi} \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{n=1}^N \left(\frac{1}{2N\lambda} \xi_n + c_n (1 - \xi_n - y_n (\mathbf{w}^T \mathbf{x}_n - b)) \right)}_{L(\mathbf{w}, b, \xi, \mathbf{c})} \right) \quad (15)$$

$$s.t. \ c_n \geq 0, \ \forall n = 1, \dots, N$$

► Thirdly, consider the optimality condition of inner problem:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} = 0 & \Rightarrow \mathbf{w}^* = \sum_{n=1}^N c_n y_n \mathbf{x}_n \quad \text{Relation between opt. primal and dual} \\ \frac{dL}{db} = 0 & \Rightarrow \sum_{n=1}^N c_n y_n = 0 \quad \text{Equality constraints of dual variables} \\ \frac{dL}{d\xi_n} = \frac{1}{2N\lambda} - c_n & \geq 0 \quad \text{(Inequality constraints of dual variables, why?)} \end{aligned} \quad (16)$$

Solve Soft-margin SVM (Dual Problem)

- Finally, plugging the relation and the constraints into (15), we obtain the final dual problem:

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} \quad & \sum_{n=1}^N c_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n c_n (\mathbf{x}_n^T \mathbf{x}_m) y_m c_m \\ \text{s.t.} \quad & \sum_{n=1}^N c_n y_n = 0, \quad 0 \leq c_n \leq \frac{1}{2N\lambda}, \quad \forall n = 1, \dots, N. \end{aligned}$$

Solve Soft-margin SVM (Dual Problem)

- Finally, plugging the relation and the constraints into (15), we obtain the final dual problem:

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} \quad & \sum_{n=1}^N c_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n c_n (\mathbf{x}_n^T \mathbf{x}_m) y_m c_m \\ \text{s.t.} \quad & \sum_{n=1}^N c_n y_n = 0, \quad 0 \leq c_n \leq \frac{1}{2N\lambda}, \quad \forall n = 1, \dots, N. \end{aligned} \tag{17}$$

- Derive it and write it in a matrix form.

Solve Soft-margin SVM (Dual Problem)

- ▶ Finally, plugging the relation and the constraints into (15), we obtain the final dual problem:

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} \quad & \sum_{n=1}^N c_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n c_n (\mathbf{x}_n^T \mathbf{x}_m) y_m c_m \\ \text{s.t.} \quad & \sum_{n=1}^N c_n y_n = 0, \quad 0 \leq c_n \leq \frac{1}{2N\lambda}, \quad \forall n = 1, \dots, N. \end{aligned} \tag{17}$$

- ▶ Derive it and write it in a matrix form.
- ▶ Why are the terms of ξ_n 's ignored?

Solve Soft-margin SVM (Dual Problem)

- ▶ Finally, plugging the relation and the constraints into (15), we obtain the final dual problem:

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} \quad & \sum_{n=1}^N c_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n c_n (\mathbf{x}_n^T \mathbf{x}_m) y_m c_m \\ \text{s.t.} \quad & \sum_{n=1}^N c_n y_n = 0, \quad 0 \leq c_n \leq \frac{1}{2N\lambda}, \quad \forall n = 1, \dots, N. \end{aligned} \tag{17}$$

- ▶ Derive it and write it in a matrix form.
- ▶ Why are the terms of ξ_n 's ignored?
- ▶ This problem can be solved by GD or Coordinate Gradient Descent.

Revisit Dual SVM as Kernel Methods

- ▶ Key relation:

$$\mathbf{w}^* = \sum_{n=1}^N c_n y_n \mathbf{x}_n$$

Revisit Dual SVM as Kernel Methods

- Key relation:

$$\mathbf{w}^* = \sum_{n=1}^N c_n y_n \mathbf{x}_n \quad (18)$$

- The dual SVM can be treated as applying a linear kernel function:

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} \quad & \sum_{n=1}^N c_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n c_n \underbrace{(\mathbf{x}_n^T \mathbf{x}_m)}_{K(\mathbf{x}_n, \mathbf{x}_m)} y_m c_m \\ \text{s.t.} \quad & \sum_{n=1}^N c_n y_n = 0, \quad 0 \leq c_n \leq \frac{1}{2N\lambda}, \quad \forall n = 1, \dots, N. \end{aligned} \quad (19)$$

Revisit Dual SVM as Kernel Methods

- ▶ Kernel SVM = Applying Linear Dual SVM in the feature space \mathcal{F} defined by the kernel function:

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{F}}, \quad \phi : \mathcal{X} \mapsto \mathcal{F}.$$

Revisit Dual SVM as Kernel Methods

- ▶ Kernel SVM = Applying Linear Dual SVM in the feature space \mathcal{F} defined by the kernel function:

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{F}}, \quad \phi : \mathcal{X} \mapsto \mathcal{F}. \quad (20)$$

- ▶ Accordingly, the key relation becomes

$$\mathbf{w}^* = \sum_{n=1}^N c_n y_n \phi(\mathbf{x}_n) \in \mathcal{F}$$

Revisit Dual SVM as Kernel Methods

- ▶ Kernel SVM = Applying Linear Dual SVM in the feature space \mathcal{F} defined by the kernel function:

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{F}}, \quad \phi : \mathcal{X} \mapsto \mathcal{F}. \quad (20)$$

- ▶ Accordingly, the key relation becomes

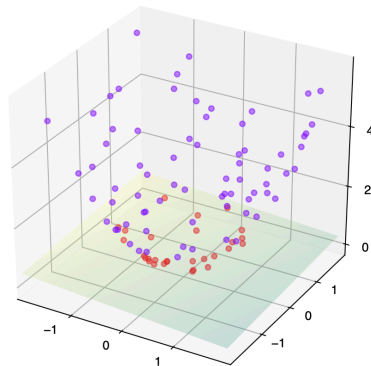
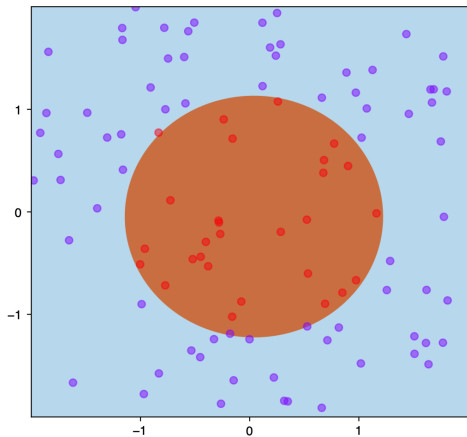
$$\mathbf{w}^* = \sum_{n=1}^N c_n y_n \phi(\mathbf{x}_n) \in \mathcal{F} \quad (21)$$

- ▶ Kernel SVM:

$$\begin{aligned} \max_{\{c_n\}_{n=1}^N} \quad & \sum_{n=1}^N c_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n c_n \underbrace{\langle \phi(\mathbf{x}_n), \phi(\mathbf{x}_m) \rangle_{\mathcal{F}}}_{K(\mathbf{x}_n, \mathbf{x}_m)} y_m c_m \\ \text{s.t.} \quad & \sum_{n=1}^N c_n y_n = 0, \quad 0 \leq c_n \leq \frac{1}{2N\lambda}, \quad \forall n = 1, \dots, N. \end{aligned} \quad (22)$$

Similar to other kernel method, we don't have to find ϕ explicitly — just define K .

Revisit Dual SVM as Kernel Methods



Extensions of SVM

Multi-class SVM

- ▶ Like LDA, “one against the rest” or pairwise classification

Extensions of SVM

Multi-class SVM

- ▶ Like LDA, “one against the rest” or pairwise classification

Support-Vector Regression (SVR)

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & |y_n - \mathbf{w}^T \mathbf{x}_n - b| \leq \epsilon \end{aligned} \tag{23}$$

In Summary

- ▶ Support-vector machine (SVM)
- ▶ Kernelized SVM

Next...

- ▶ Information theory in Statistic ML
- ▶ Decision tree model

HW 4: DDL June 3, 2022

Python Programming

- 1 Lab # 9 (4 Pts)
- 2 Lab # 10 (4 Pts)

Questions for Tech Report (6 Pts, ≤ 3 Pages)

- 1 Given $\mathbf{X} \in \mathbb{R}^{N \times D}$ and $\mathbf{Y} \in \mathbb{R}^{N \times C}$, where each row of \mathbf{X} is a D -dimensional feature, and each row of \mathbf{Y} is a one-hot vector indicating one of the C classes. **Can we solve classification as regression?** e.g.,

$$\min_{\mathbf{W}} \|\mathbf{Y} - \mathbf{XW}\|_F^2 \quad (24)$$

If it does not work well in general, when it works? (Hint: Find the answer in ESL)
(3 Pts)

- 2 Derive from (15) to (17) in details, and demonstrate that the terms related to ξ_n 's are ignorable (Hint: consider the objective function in (15)). (3 Pts)