

Convex Optimization

Lecture 12: Support Vector Machine

Lecturer: *Dr.* Wan-Lei Zhao

Autumn Semester 2025

Outline

- 1 Support Vector Machine
- 2 Solve SVM with Langrange Multiplier
- 3 References

Overview of discriminative classifier (1)

- Given a training set $(x_i, y_i)_{i=1 \dots m}$
- $x_i \in R^n$ is the observation, $y_i \in [-1, 1]$ is the class label
- A classifier is trained with this set
- Given a new instance $u \in R^n$
- The classifier makes the prediction

Overview of discriminative classifier (2)

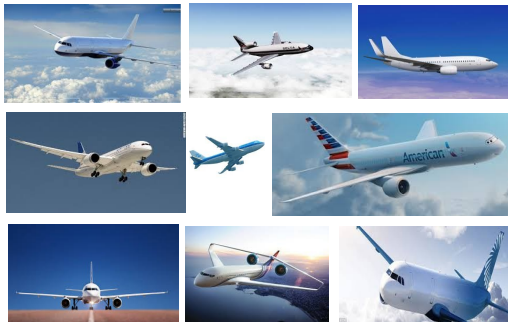
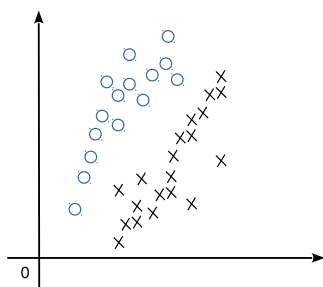


Figure: Training examples for plane

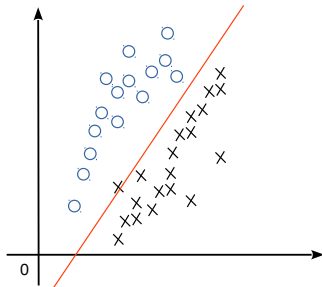


Figure: Samples to be predicted, **plane** or **not**

Opening discussion: linear regression (1)



(a) points groups



(b) linear regression on points groups

Figure: Two groups of points

- Linear regression helps to address this problem

Opening discussion: linear regression (2)

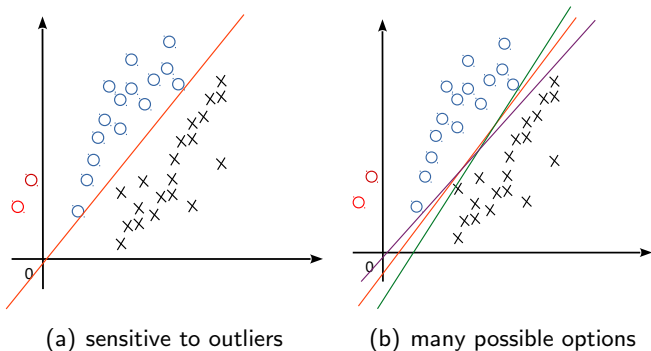


Figure: Two groups of points

- The results are not satisfactory

Discriminative classifier: the binary classification problem

- Let's now start with a simple problem
- Given any instance x , the classifier is asked to make a binary prediction
- **Yes** or **No**
- From another point of view, we want to train a function $g(\cdot)$

$$g(x) = \begin{cases} 1 & c1 \\ -1 & otherwise \end{cases} \quad (1)$$

- To simplify the problem
- We say $g(\cdot)$ is a linear function, that is

$$g(x) = w^T x + b \quad (2)$$

- w and b are parameters to be trained, $g(w^T x) \in [-1, 1]$
- **By default**, we are dealing with **column** vectors

SVM: the functional margin

- Now, it is clear that we want to maximize the margin
- Given $\gamma^{(i)}$, it is defined as

$$\gamma^{(i)} = y^{(i)}(w^T x^{(i)} + b) \quad (3)$$

- if $y^{(i)}=1$, $\gamma^{(i)}$ is going to be large only if $(w^T x^{(i)} + b)$ is large positive
- if $y^{(i)}=-1$, $\gamma^{(i)}$ is going to be large only if $(w^T x^{(i)} + b)$ is large negative
- However, this functional margin is not fully meaningful
- We can make it larger by $(2w^T x^{(i)} + 2b)$

SVM: the minimum functional margin (1)

- With w being normalized, the functional margin $\gamma^{(i)}$ is given as

$$\gamma^{(i)} = y^{(i)} \left(\frac{w^T}{\|w\|_2} x^{(i)} + b \right) \quad (4)$$

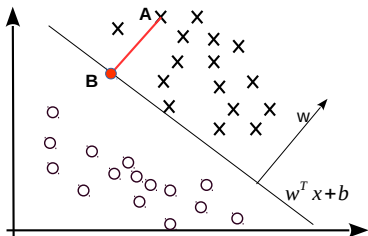


Figure: $\gamma^{(i)}$ is the distance from A to line $w^T x + b$

- It is equal to the distance between **A** and **B**
- Here, we consider the case $y_i=1$, similar analysis applies to $y_i=-1$

SVM: the minimum functional margin (2)

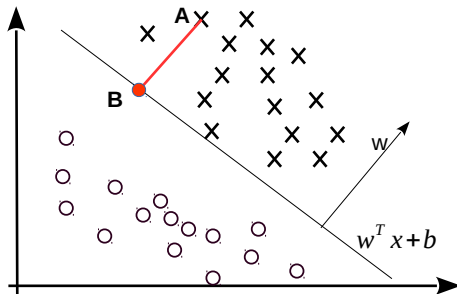


Figure: $\gamma^{(i)}$ is the distance from A to line $w^T x + b$

- Given $A = x^{(i)}$, **B** is on the line $w^T x + b$
- According to the theorem of distance between a point and a line/plain

$$r^{(i)} = \frac{w^T x^{(i)} + b}{\|w\|_2} \quad (5)$$

SVM: the minimum functional margin (3)

- So we have

$$r^{(i)} = \frac{w^T x^{(i)}}{\|w\|_2} + \frac{b}{\|w\|_2} \quad (6)$$

- Above equation is for positive training example
- For all the training examples, we have

$$r^{(i)} = y^{(i)} \left(\frac{w^T}{\|w\|_2} x^{(i)} + \frac{b}{\|w\|_2} \right) \quad (7)$$

SVM: the minimum functional margin (4)

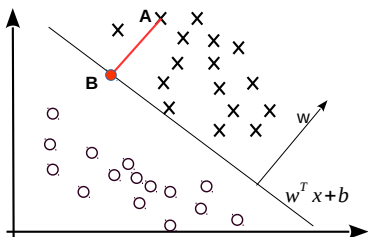


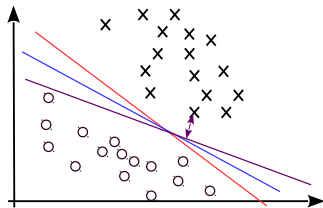
Figure: $\gamma^{(i)}$ is the distance from A to line $\frac{w^T}{\|w\|_2}x + b$

$$r^{(i)} = y^{(i)} \left(\frac{w^T}{\|w\|_2} x^{(i)} + \frac{b}{\|w\|_2} \right) \quad (10)$$

- We want to find the minimum margin, and maximize it

$$\gamma = \underset{i=1, \dots, m}{\operatorname{argmin}} \gamma^{(i)} \quad (11)$$

SVM: the model (1)

Figure: Searching for maximum γ

$$\gamma = \underset{i=1, \dots, m}{\operatorname{argmin}} \gamma^{(i)} \quad (11)$$

- Searching for w, b that maximizes γ

$$\begin{aligned} & \underset{w, b, \gamma}{\operatorname{Max.}} && \gamma \\ & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq \gamma, i = 1, \dots, m \\ & && \|w\|_2 = 1. \end{aligned}$$

SVM: the model (2)

- Searching for w, b that maximizes γ

$$\begin{aligned}
 & \underset{w, b, \gamma}{\text{Max.}} && \gamma && (11) \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq \gamma, i = 1, \dots, m \\
 & && \|w\|_2 = 1.
 \end{aligned}$$

- Unfortunately, above problem is not solvable
- Constraint $\|w\|_2 = 1$ is not convex

$$\begin{aligned}
 & \underset{w, b, \hat{\gamma}}{\text{Max.}} && \hat{\gamma} / \|w\|_2 && (12) \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, i = 1, \dots, m.
 \end{aligned}$$

- Unfortunately, above problem is not solvable either
- $\hat{\gamma}$ is functional margin, it is valid to scale it to $\hat{\gamma} = 1$

SVM: the model (3)

$$\begin{aligned}
 & \underset{w, b, \hat{\gamma}}{\text{Max.}} && \hat{\gamma} / \|w\|_2 \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, i = 1, \dots, m.
 \end{aligned} \tag{12}$$

- $\hat{\gamma}$ is functional margin, it is valid to scale it to $\hat{\gamma} = 1$

$$\begin{aligned}
 & \underset{w, b}{\text{Max.}} && \frac{1}{\|w\|_2} \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq 1, i = 1, \dots, m.
 \end{aligned} \tag{13}$$

- This is equivalent to solving following **quadratic optimization** problem

$$\begin{aligned}
 & \underset{w, b}{\text{Min.}} && \frac{1}{2} \|w\|^2 \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq 1, i = 1, \dots, m.
 \end{aligned} \tag{14}$$

SVM: the model (4)

- Binary classification problem is now modeled as **quadratic optimization** problem

$$\begin{aligned} \text{Min.}_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s. t.} \quad & 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, i = 1, \dots, m. \end{aligned} \tag{14}$$

- The unknowns are w and b
- The target function is quadratic
- m constraints are linear
- We are going to solve it with **Lagrange multiplier** method

$$\begin{aligned}
 & \underset{w, b, \gamma}{\text{Max.}} && \gamma \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq \gamma, i = 1, \dots, m \\
 & && \|w\|_2 = 1.
 \end{aligned}$$

$$\Downarrow$$

$$\begin{aligned}
 & \underset{w, b, \hat{\gamma}}{\text{Max.}} && \hat{\gamma} / \|w\|_2 \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, i = 1, \dots, m.
 \end{aligned}$$

$$\Downarrow$$

$$\begin{aligned}
 & \underset{w, b}{\text{Max.}} && \frac{1}{\|w\|_2} \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq 1, i = 1, \dots, m.
 \end{aligned}$$

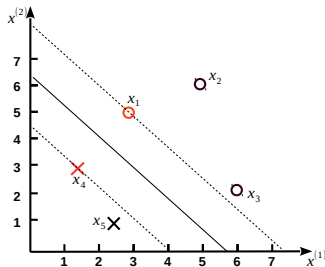
$$\Downarrow$$

$$\begin{aligned}
 & \underset{w, b}{\text{Min.}} && \frac{1}{2} \|w\|^2 \\
 & \text{s. t.} && y^{(i)}(w^T x^{(i)} + b) \geq 1, i = 1, \dots, m.
 \end{aligned}$$

Outline

- 1 Support Vector Machine
- 2 Solve SVM with Langrange Multiplier
- 3 References

Solve it as a QP problem (1)

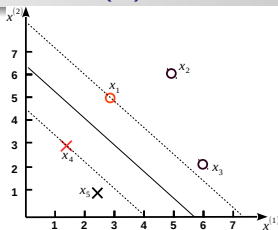


$$\text{Min.}_{w,b} \quad \frac{1}{2} \|w\|^2 \quad (1)$$

$$\text{sub. to} \quad 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, i = 1, \dots, m.$$

- ① The training samples (x_i, y_i) s are all known
- ② $w \in R^d$ and $b \in R$ are unknown
- ③ All constraints are linear inequation, it is QP problem!!

Solve it as a QP problem (2)



$$\text{Min.}_{w,b} \quad \frac{1}{2} \|w\|^2 \quad (2)$$

$$\text{sub. to} \quad 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, i = 1, \dots, m.$$

$$\Downarrow$$

$$\text{Min.}_{w,b} \quad \frac{1}{2} \|w\|^2 \quad (3)$$

$$\text{s. t.} \quad [-y^{(i)} x^{(i)T} - y^{(i)}] \begin{bmatrix} w \\ b \end{bmatrix} \cdot \leq -1, i = 1, \dots, m.$$

Solve it as a QP problem (3)

$$\text{Min.}_{w,b} \quad \frac{1}{2} \|w\|^2 \quad (4)$$

$$\text{s. t.} \quad 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, i = 1, \dots, m.$$

$$\Downarrow$$

$$\text{Min.}_{w,b} \quad \frac{1}{2} \|w\|^2 \quad (5)$$

$$\text{s. t.} \quad [-y^{(i)} x^{(i)T} \quad -y^{(i)}] \begin{bmatrix} w \\ b \end{bmatrix} \leq -1, i = 1, \dots, m.$$

$$\text{s.t.} \quad \begin{bmatrix} -y^{(1)} x_1^{(1)} & -y^{(1)} x_2^{(1)} & -y^{(1)} \\ -y^{(2)} x_1^{(2)} & -y^{(2)} x_2^{(2)} & -y^{(2)} \\ \vdots & \vdots & \vdots \\ -y^{(m)} x_1^{(m)} & -y^{(m)} x_2^{(m)} & -y^{(m)} \end{bmatrix} \cdot \begin{bmatrix} w \\ b \end{bmatrix} \preceq \begin{bmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{bmatrix}$$

Lagrangian duality of SVM (1)

- Given binary classification problem is modeled as

$$\begin{aligned} \text{Min.}_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s. t} \quad & 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, i = 1, \dots, m. \end{aligned} \quad (1)$$

- Given $\alpha_i \geq 0, i = 1, \dots, m$, we define

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1] \quad (2)$$

- Take derivative on w , we have

$$\nabla_w L(w, b, \alpha) = w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \quad (3)$$

Lagrangian duality of SVM (2)

- Lagrangian:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1] \quad (4)$$

- Take derivative on w , we have

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \quad (5)$$

- Take derivative on b , we have

$$\frac{\partial}{\partial b} L(w, b, \alpha) = \sum_{i=1}^m \alpha_i y^{(i)} = 0 \quad (6)$$

- Plug Eqn. 5 and Eqn. 6 into Eqn. 4, we have

$$L(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \quad (7)$$

Lagrangian duality of SVM (3)

- Lagrangian:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1] \quad (4)$$

- Notice that we have new form for $L(\cdot)$

$$L(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \quad (7)$$

- We say this is the dual problem of original $L(w, b, \alpha)$

$$\begin{aligned} \text{Max.}_{\alpha} \quad & G(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & \alpha_i \geq 0, i = 1 \cdots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0. \end{aligned} \quad (8)$$

Lagrangian duality of SVM (4)

- Dual form of Lagrangian:

$$\begin{aligned}
 \text{Max.}_{\alpha} \quad & G(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j < x^{(i)}, x^{(j)} > \\
 \text{s.t.} \quad & \alpha_i \geq 0, i = 1 \cdots, m \\
 & \sum_{i=1}^m \alpha_i y^{(i)} = 0.
 \end{aligned} \tag{8}$$

- This problem is easier to solve, only α to be considered
- Conventionally, this problem is solved by **coordinate ascent**¹
- Remember that once we get α , we have

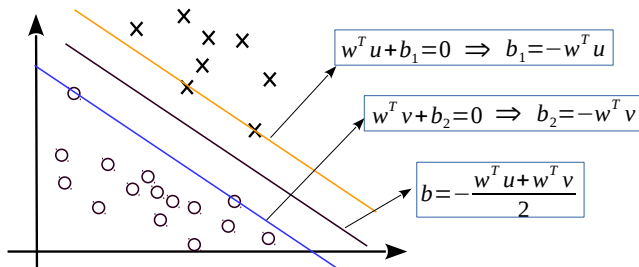
$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \quad b = - \frac{\max_{i:y^{(i)}=-1} w^T x^{(i)} + \min_{i:y^{(i)}=1} w^T x^{(i)}}{2}$$

¹Lecture notes on Machine Learning, Andrew Ng. <http://cs229.stanford.edu/>

Lagrangian duality of SVM (5)

- Solution to the Dual form of Lagrangian:

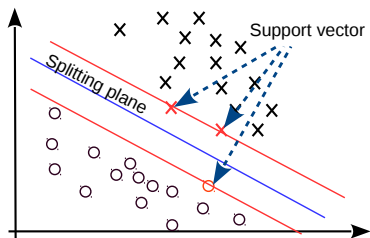
$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \quad b = -\frac{\max_{i:y(i)=-1} w^T x^{(i)} + \min_{i:y(i)=1} w^T x^{(i)}}{2}$$



SVM: brief procedure of the optimization

$$\begin{aligned}
 \text{Max.}_{\alpha} \quad & G(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j < x^{(i)}, x^{(j)} > \\
 \text{s.t.} \quad & \alpha_i \geq 0, i = 1 \cdots, m \\
 & \sum_{i=1}^m \alpha_i y^{(i)} = 0.
 \end{aligned} \tag{8}$$

- Coordinate ascent: the optimization process takes α_i and α_j out
- Check whether change them will lead to larger $G(\cdot)$

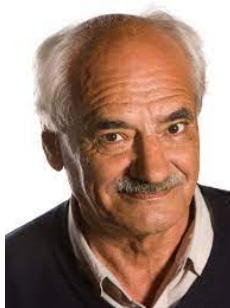


- The splitting plane (line) is determined by the support vectors

Comments about SVM



(a) Vladimir Vapnik (1936 -)



(b) A. Chervonenkis
(1938 - 2014)

Figure: The scientists invented SVM.

- The most popular classifier before the advent of deep learning
- The theory is still shining

- ① Top 10 algorithms in data mining, X. Wu, V. Kumar and et al. Knowledge and Information Systems, 2008, 14(1): 1-37
- ② The Nature of Statistical Learning Theory, Vladimir N. Vapnik , Springer-Verlag, 1995.
- ③ Lecture notes on Machine Learning, Andrew Ng., <http://cs229.stanford.edu/>