

CS 599 Machine Learning

Lecture 9: Support Vector Machines

Hao Ji

Computer Science Department

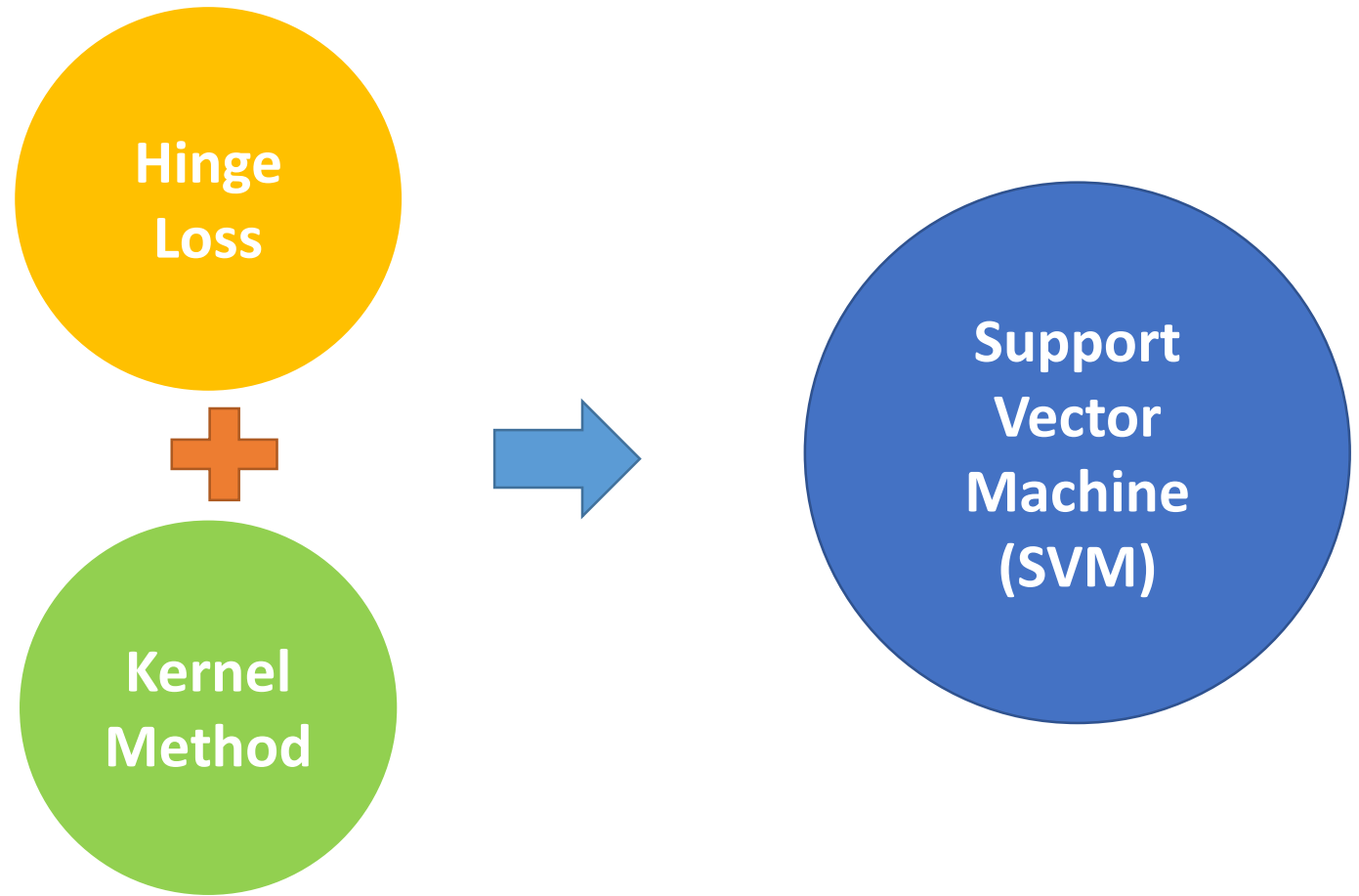
Cal Poly Pomona

SVM

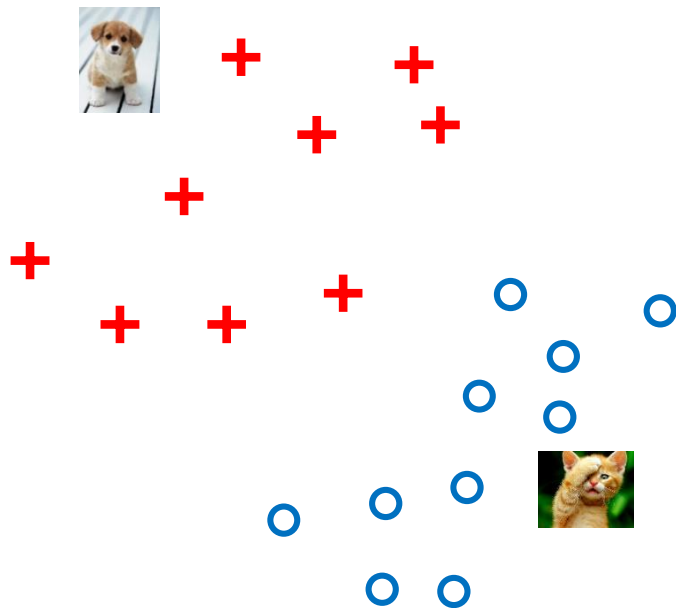
- SVM was inspired from Statistical Learning theory
- SVM was first introduced in 1992
 - Boser, Bernhard E., Isabelle M. Guyon, and Vladimir N. Vapnik. "A training algorithm for optimal margin classifiers." In *Proceedings of the fifth annual workshop on Computational learning theory*, pp. 144-152. ACM, 1992.
- SVM becomes popular because of its success in handwritten digit recognition
 - Bottou, Léon, Corinna Cortes, John S. Denker, Harris Drucker, Isabelle Guyon, Lawrence D. Jackel, Yann LeCun et al. "Comparison of classifier methods: a case study in handwritten digit recognition." In *Pattern Recognition, 1994. Vol. 2-Conference B: Computer Vision & Image Processing., Proceedings of the 12th IAPR International. Conference on*, vol. 2, pp. 77-82. IEEE, 1994.

SVM

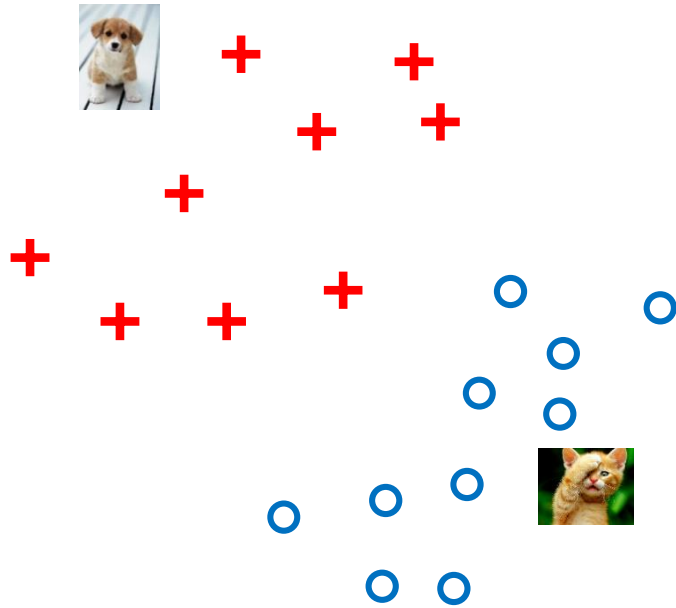
- Linear SVM
- Soft Margin
 - Hinge Loss
- Kernel Trick



Linear Classifier

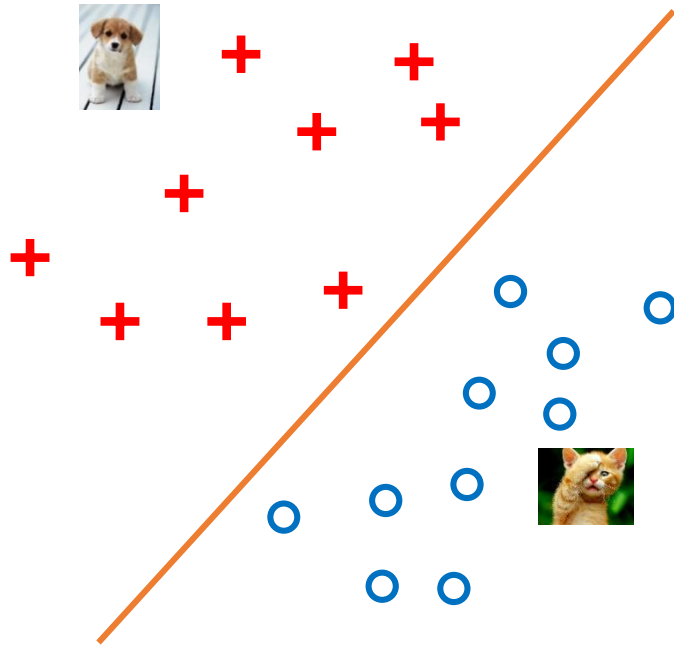


Linear Classifier



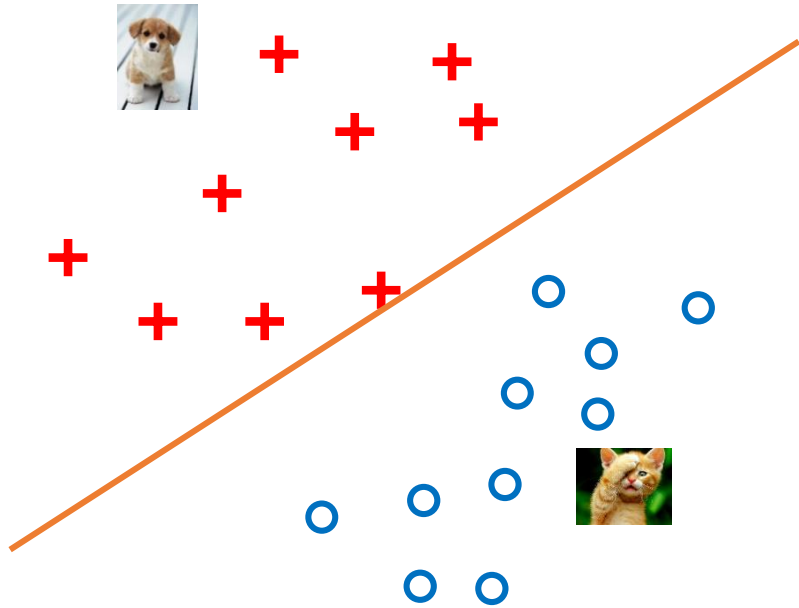
How would you classify this data?

Linear Classifier



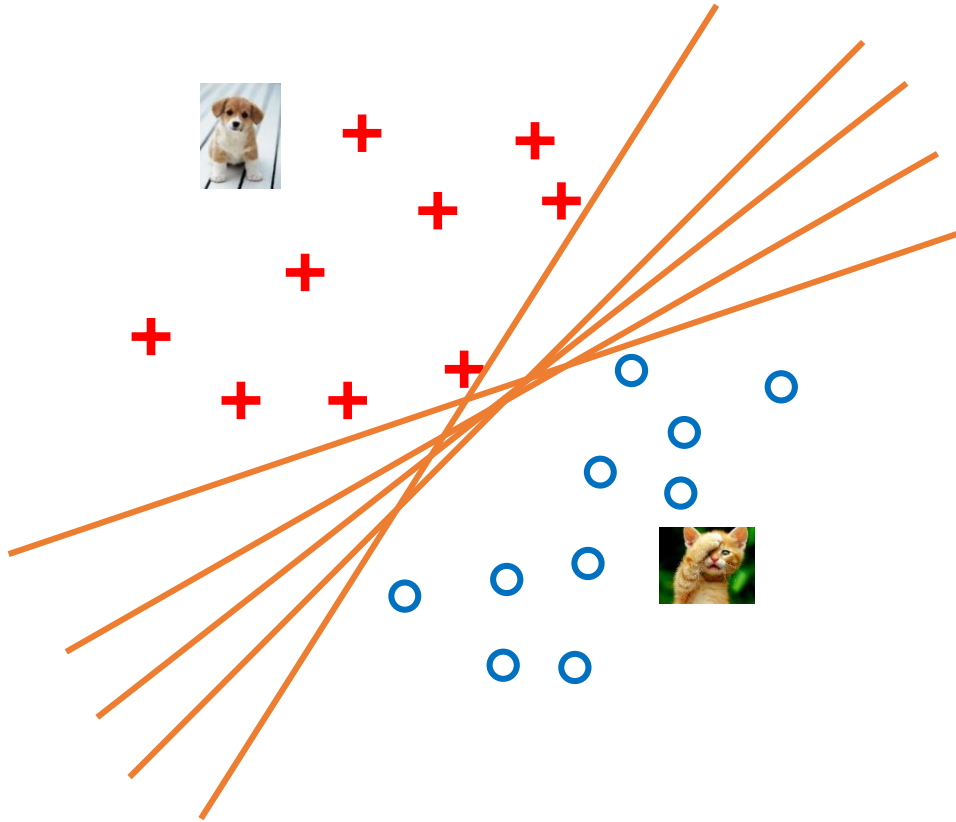
$$f(x) = \begin{cases} 1 & \text{if } w^T x + b \geq 0 \\ -1 & \text{if } w^T x + b < 0 \end{cases}$$

Linear Classifier



$$f(x) = \begin{cases} 1 & \text{if } w^T x + b \geq 0 \\ -1 & \text{if } w^T x + b < 0 \end{cases}$$

Linear Classifier

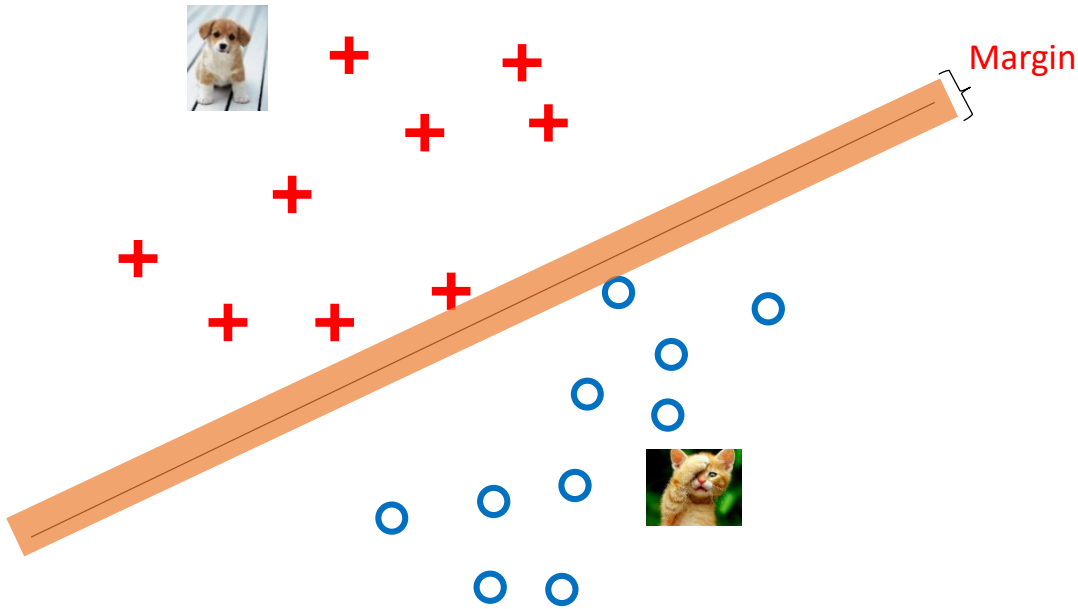


$$f(x) = \begin{cases} 1 & \text{if } w^T x + b \geq 0 \\ -1 & \text{if } w^T x + b < 0 \end{cases}$$

Any of these would be fine ...

But, which one is the best?

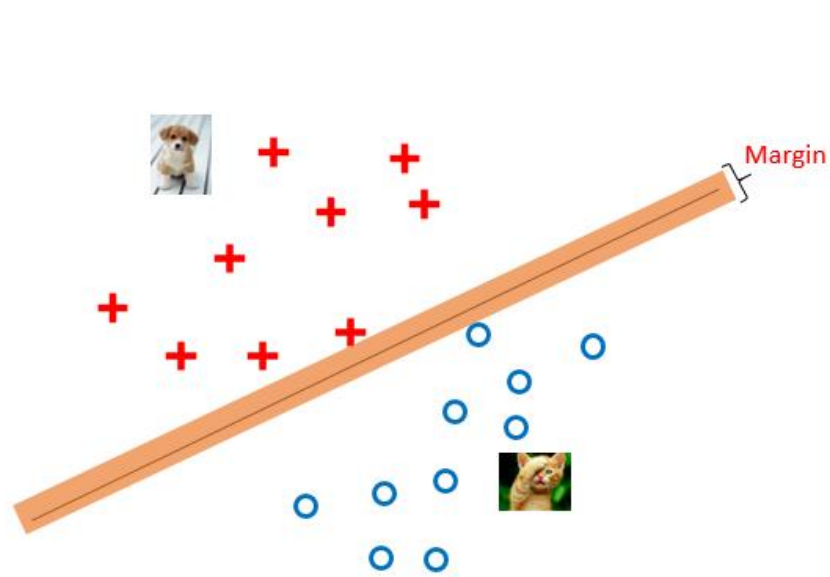
Linear Classifier



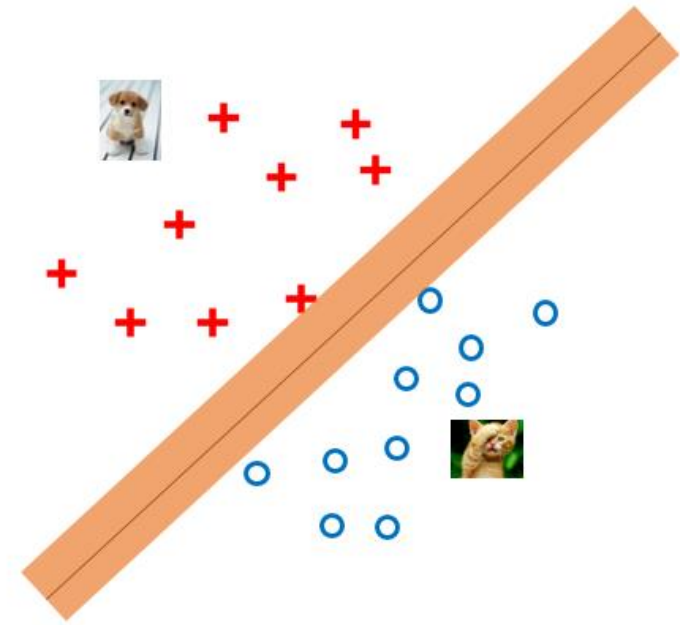
$$f(x) = \begin{cases} 1 & \text{if } w^T x + b \geq 0 \\ -1 & \text{if } w^T x + b < 0 \end{cases}$$

Margin: the width that the boundary could be increased by before hitting a datapoint

SVM



VS

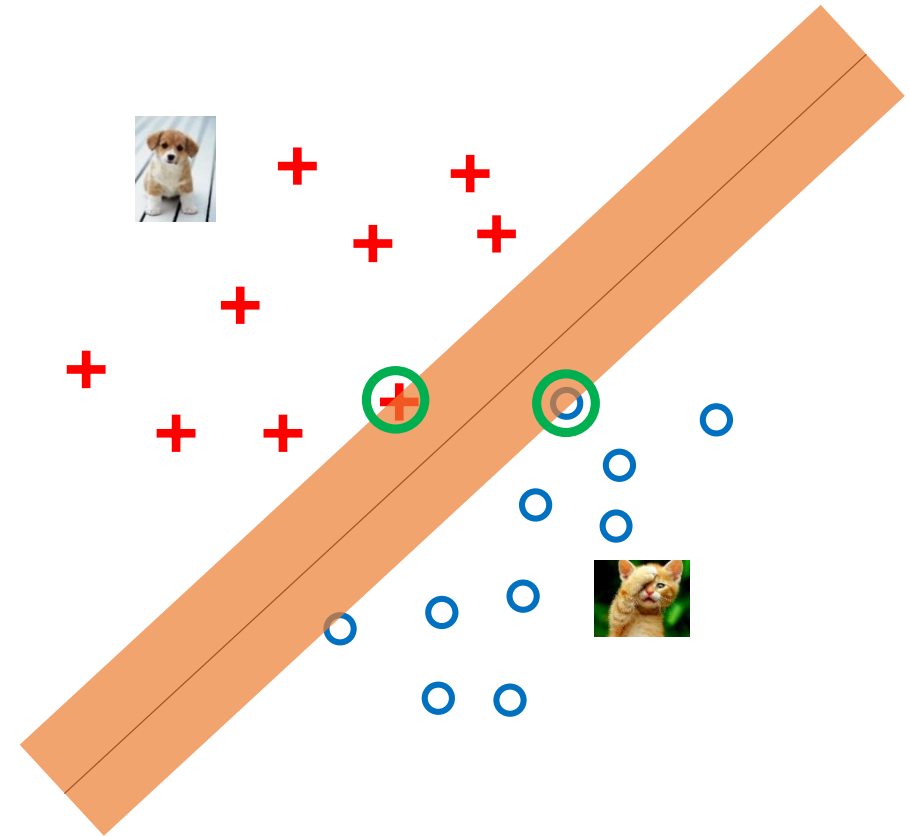


The maximum margin linear classifier: the linear classifier with the maximum margin

The simplest kind of SVM (called Linear SVM)

SVM

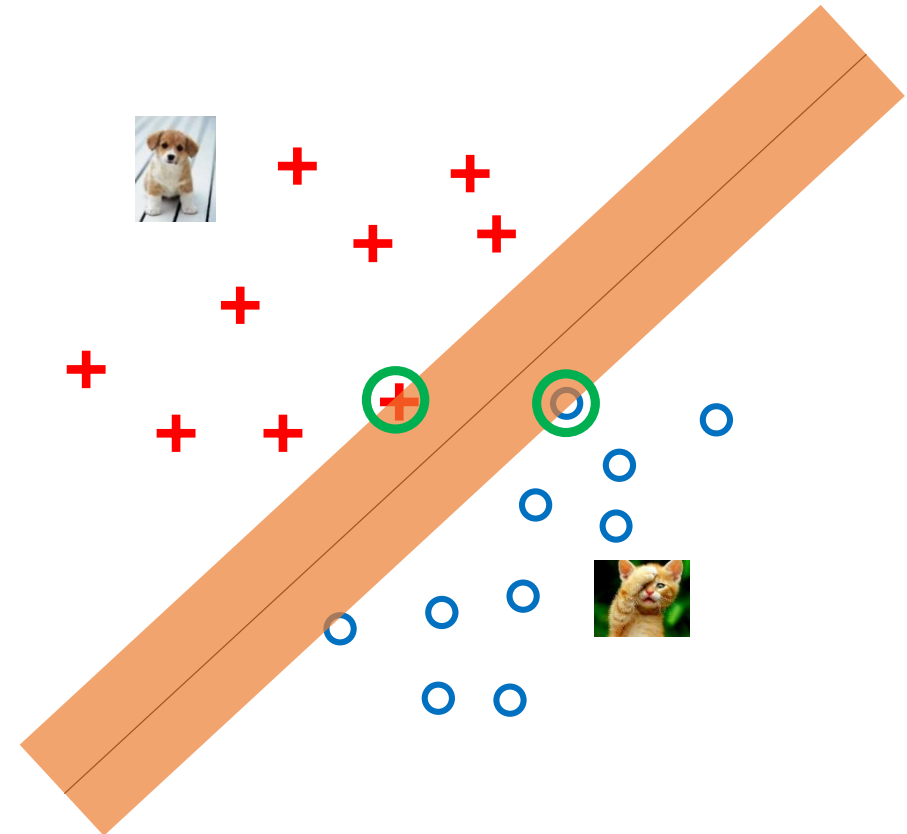
Support Vectors: the data points that the margin pushes up against



SVM

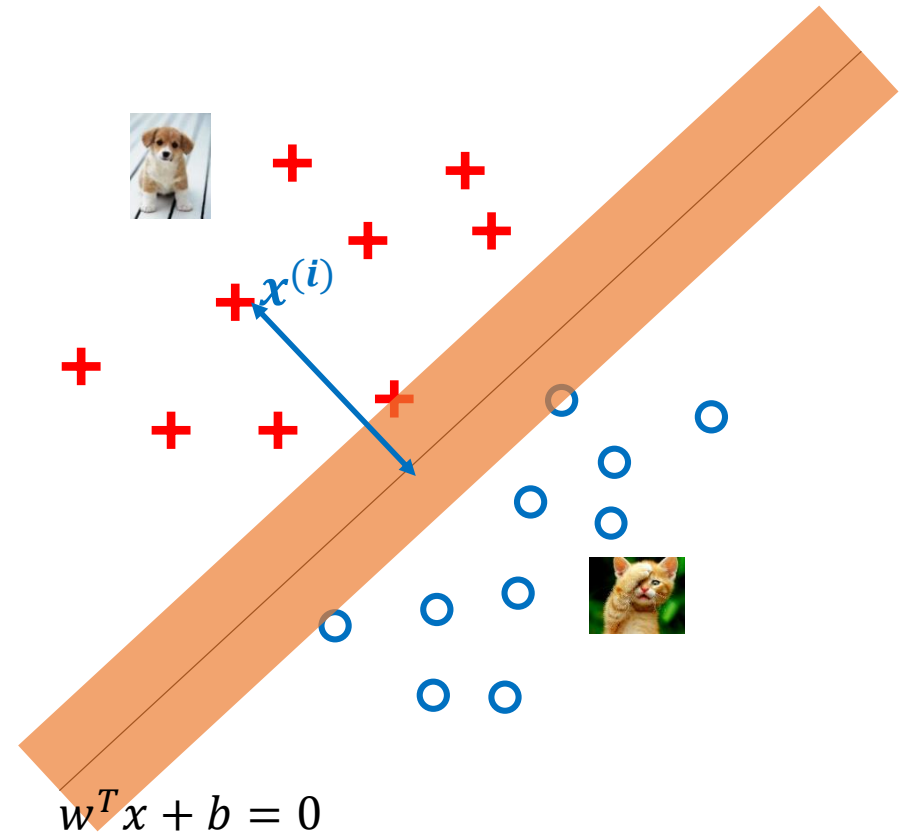
Support Vectors: the data points that the margin pushes up against

- Intuitively this feels safest
 - If we've made a small error in the location of the boundary, this gives us least chance of causing a misclassification
- The model is immune to removal of any non-support-vector datapoints
- There's some theory (using VC dimension) that is related to (but not the same as) the proposition that this is a good thing
- Empirically it works well.



SVM

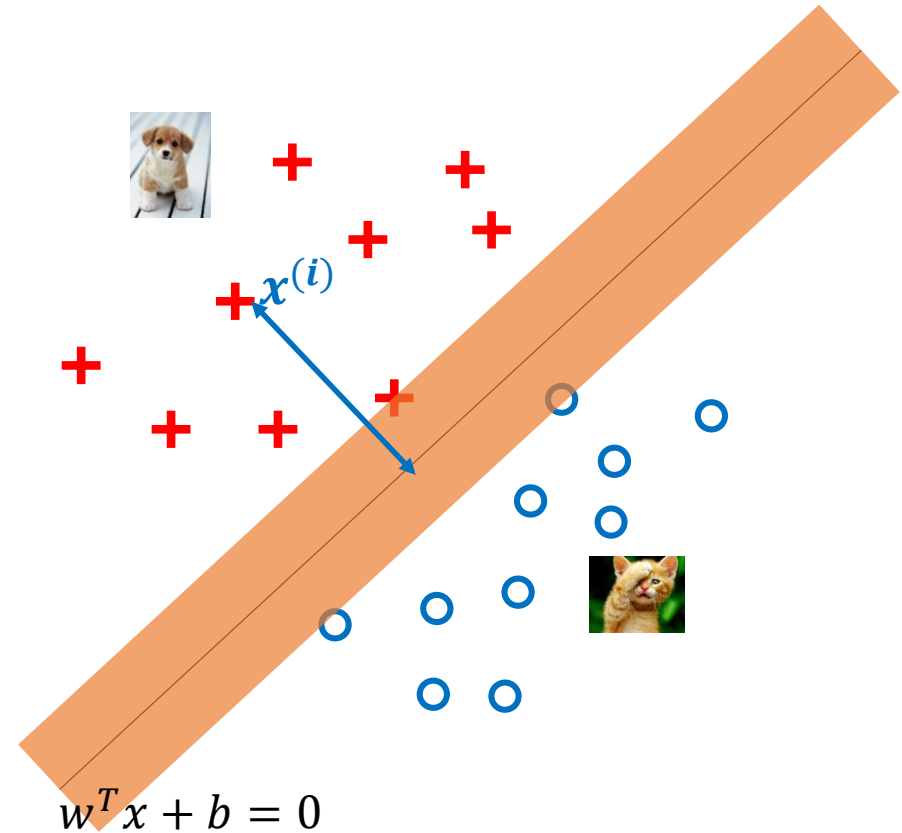
- What is the **distance** expression for a point $x^{(i)}$ to a line $w^T x + b = 0$?



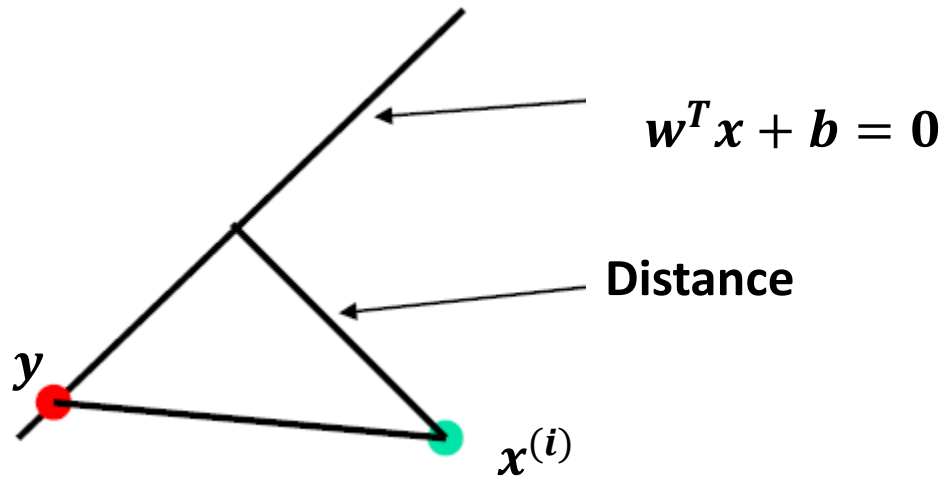
SVM

- What is the **distance** expression for a point $x^{(i)}$ to a line $w^T x + b = 0$?

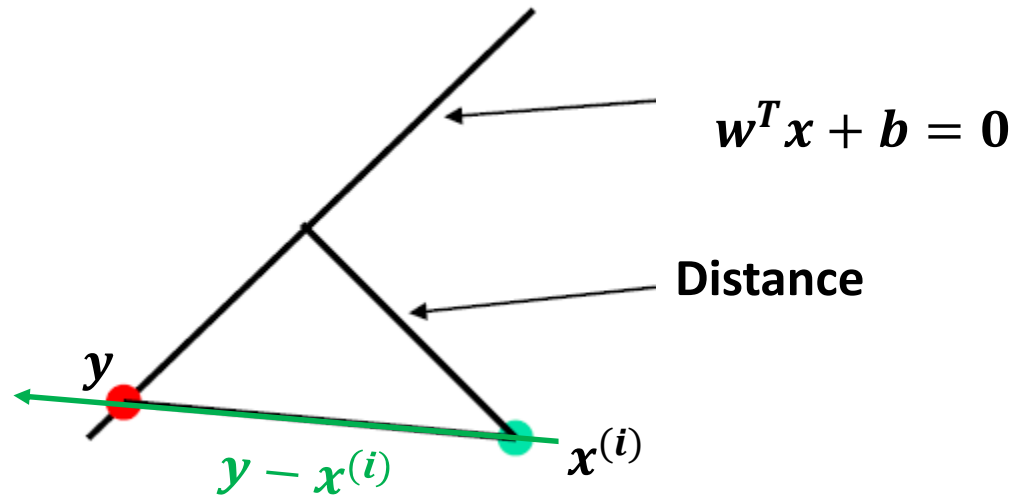
$$d(x^{(i)}) = \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$



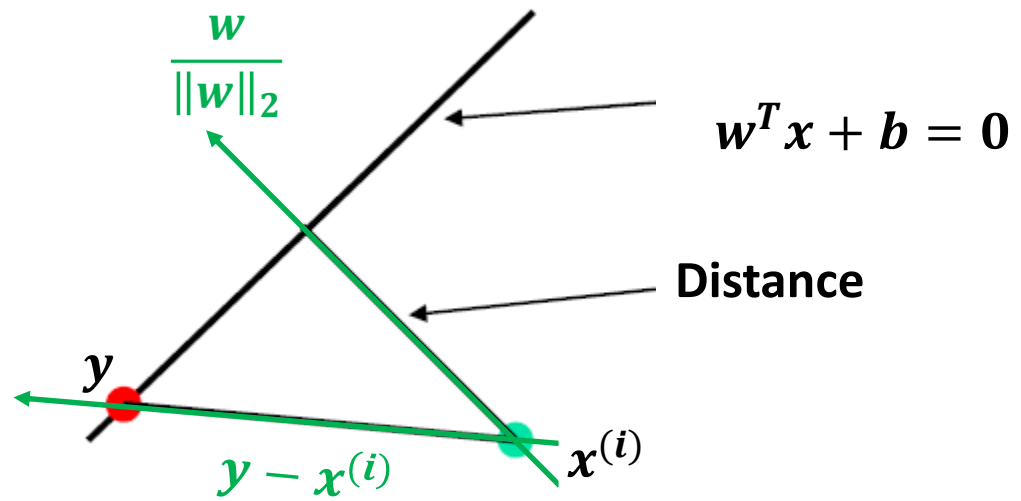
Distance Expression



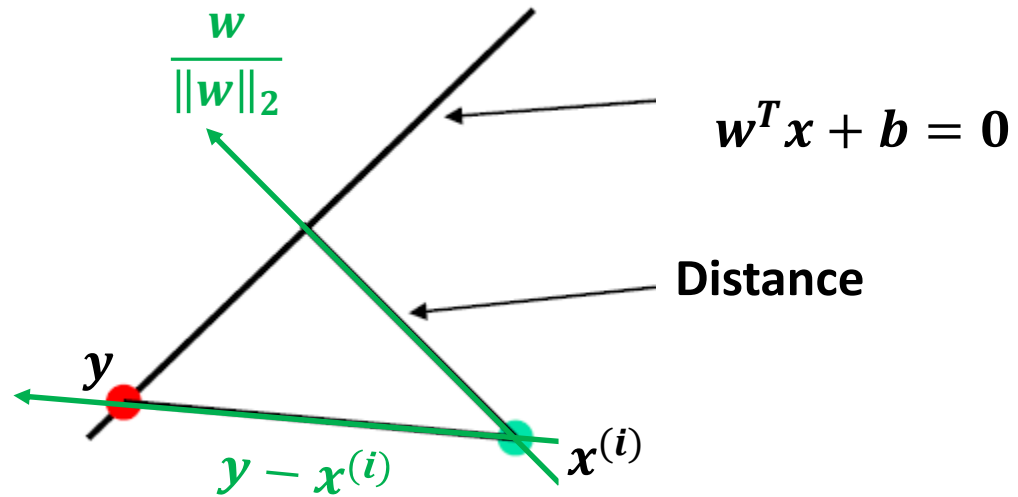
Distance Expression



Distance Expression

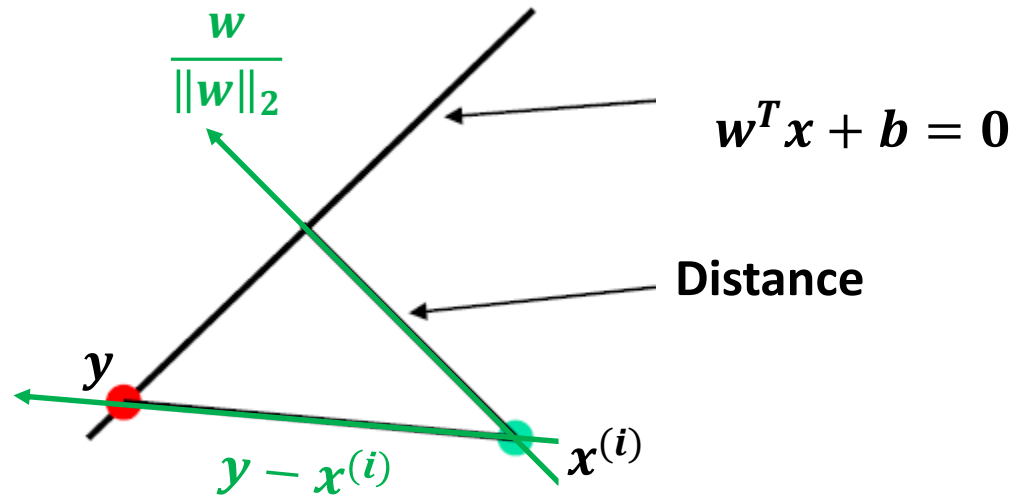


Distance Expression



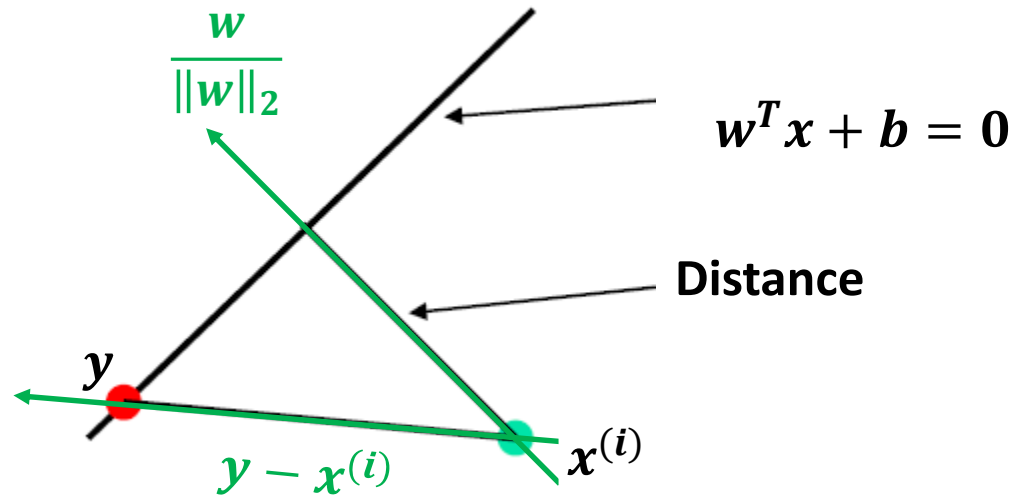
$$d(x^{(i)}) = \left| \frac{w^T}{\|w\|_2} (y - x^{(i)}) \right|$$

Distance Expression



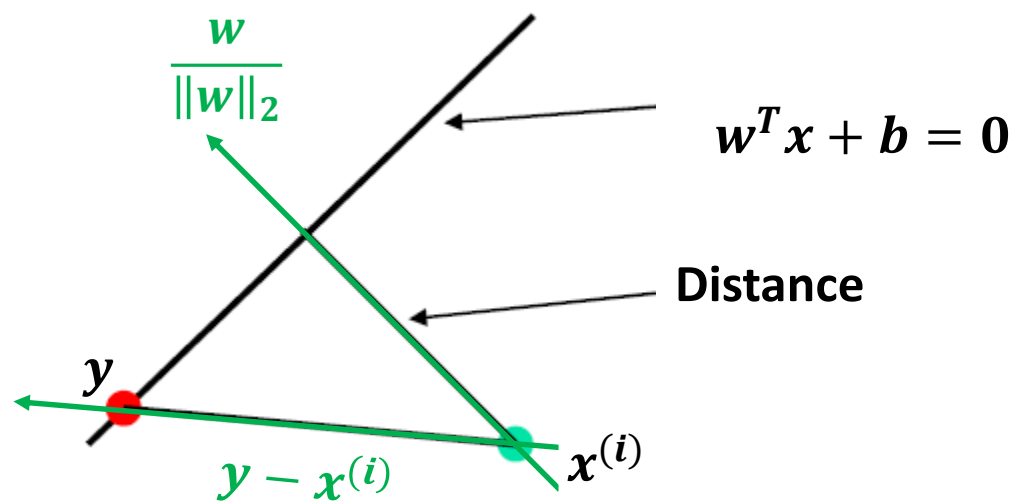
$$\begin{aligned} d(x^{(i)}) &= \left| \frac{w^T}{\|w\|_2} (y - x^{(i)}) \right| \\ &= \left| \frac{w^T y - w^T x^{(i)}}{\|w\|_2} \right| \end{aligned}$$

Distance Expression



$$\begin{aligned} d(x^{(i)}) &= \left| \frac{w^T}{\|w\|_2} (y - x^{(i)}) \right| \\ &= \left| \frac{w^T y - w^T x^{(i)}}{\|w\|_2} \right| \\ &= \left| \frac{-b - w^T x^{(i)}}{\|w\|_2} \right| \end{aligned}$$

Distance Expression

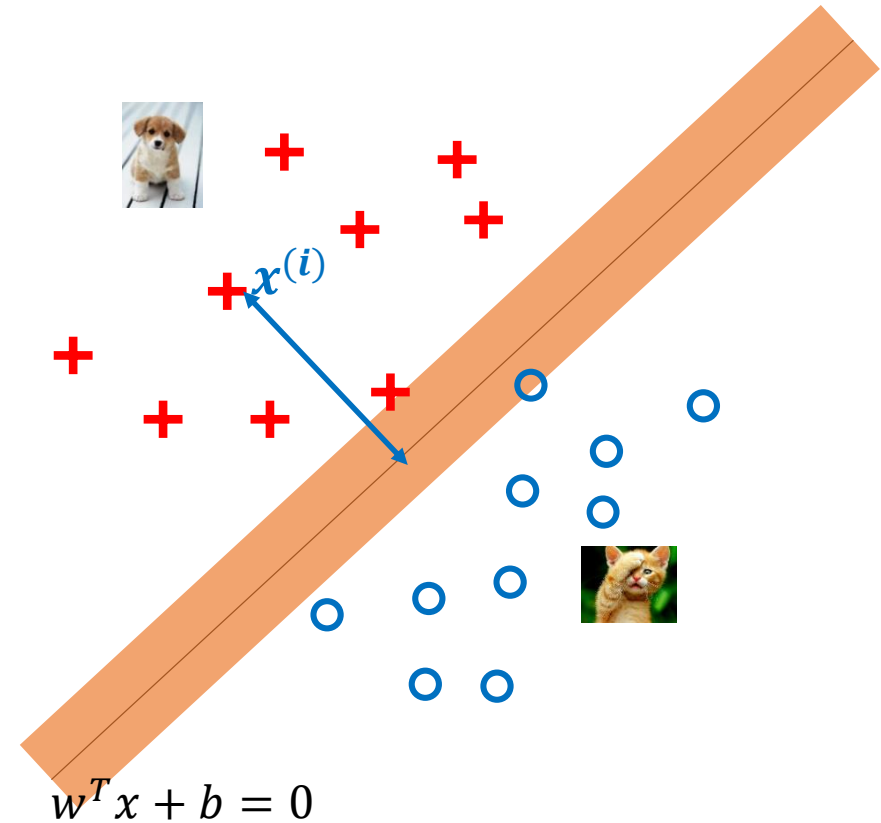


$$\begin{aligned} d(x^{(i)}) &= \left| \frac{w^T}{\|w\|_2} (y - x^{(i)}) \right| \\ &= \left| \frac{w^T y - w^T x^{(i)}}{\|w\|_2} \right| \\ &= \left| \frac{-b - w^T x^{(i)}}{\|w\|_2} \right| \\ &= \left| \frac{w^T x^{(i)} + b}{\|w\|_2} \right| \end{aligned}$$

SVM

- What is the **distance** expression for a point $x^{(i)}$ to a line $w^T x + b = 0$?

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



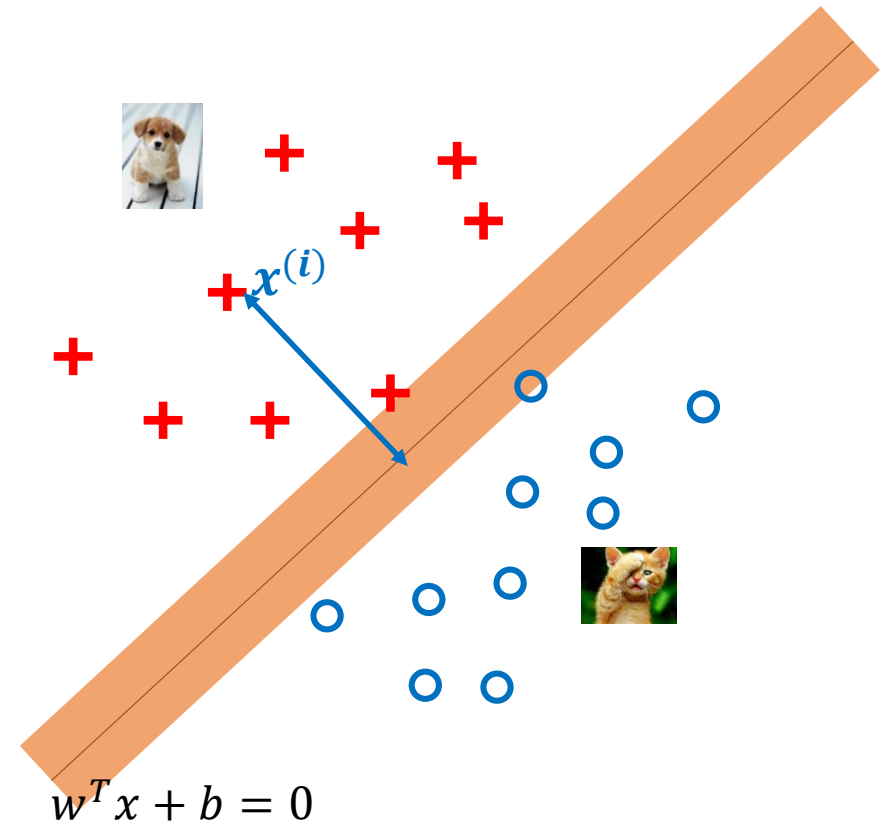
SVM

- What is the **distance** expression for a point $x^{(i)}$ to a line $w^T x + b = 0$?

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

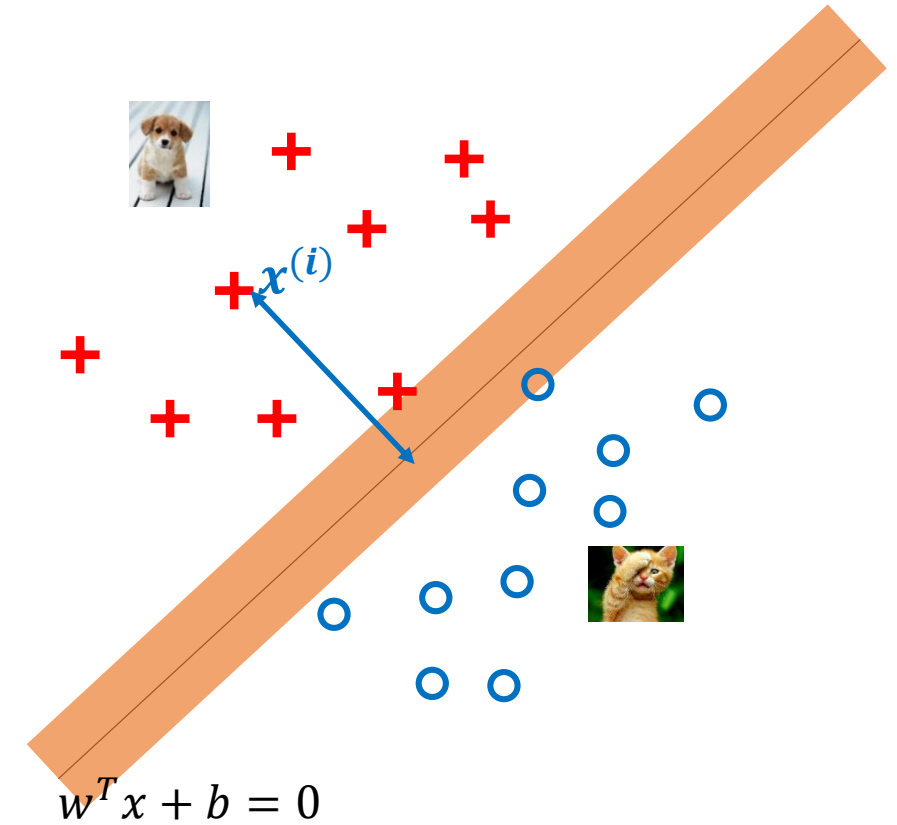
- Margin**: the width that the boundary could be increased by before hitting a datapoint

$$\begin{aligned} \text{margin} &\equiv 2 * \min_{x^{(i)} \in D} d(x^{(i)}) \\ &= 2 * \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2} \end{aligned}$$



SVM

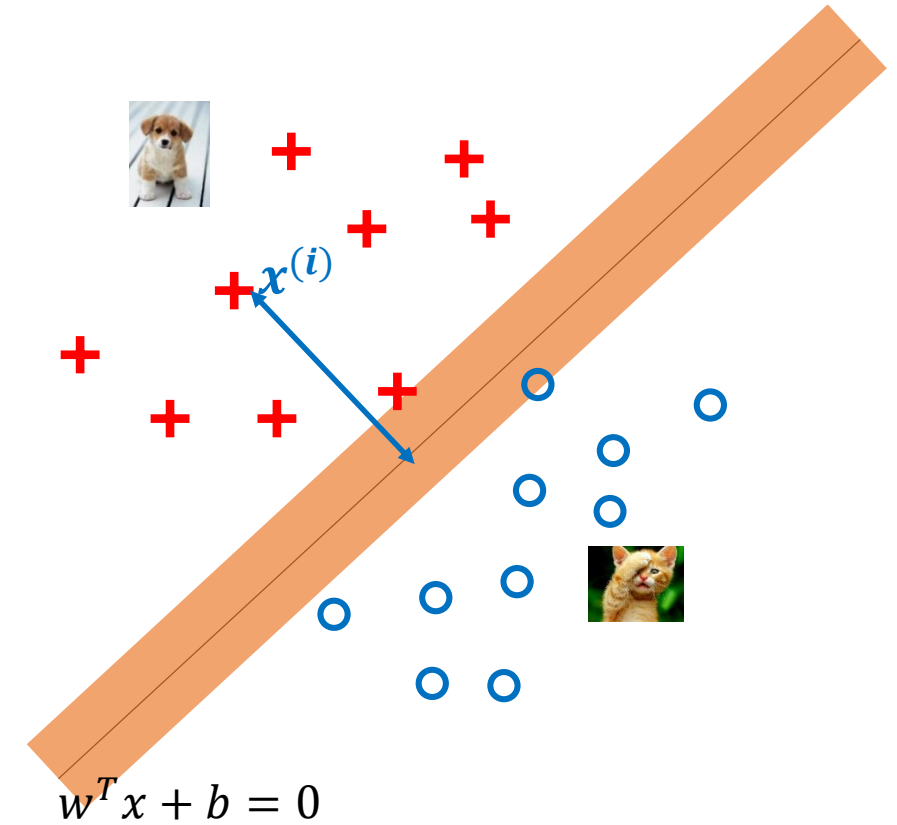
- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin



SVM

- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin

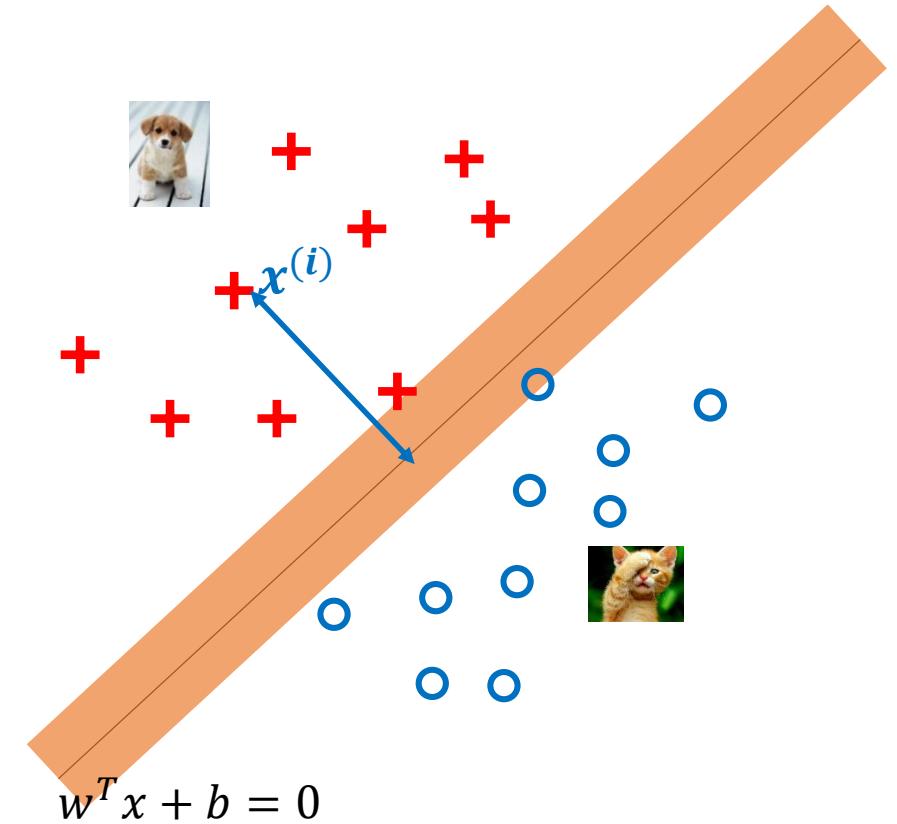
max margin
w, b



SVM

- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin

max margin
w, b

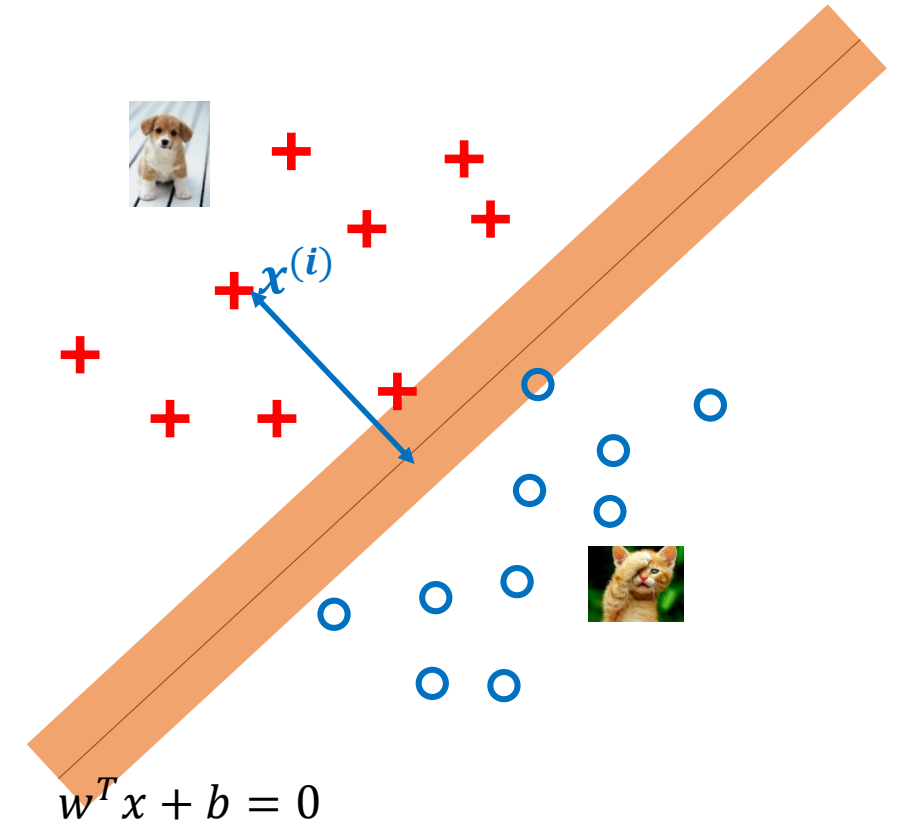


$$\text{margin} = \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

SVM

- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin

$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

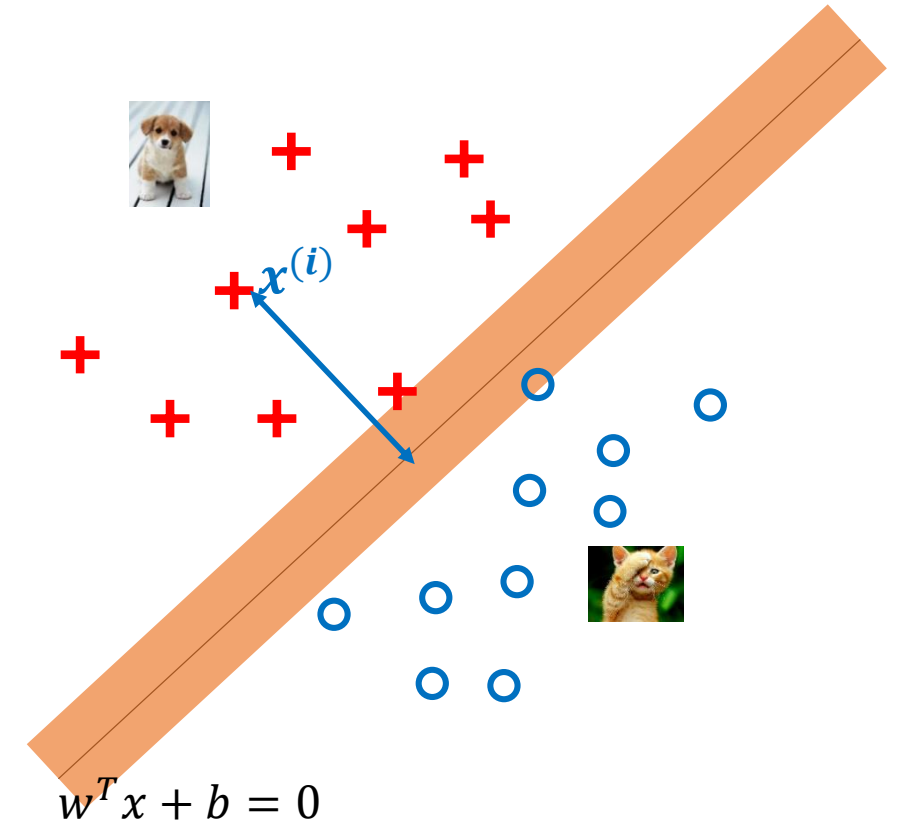


$$\text{margin} = \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

SVM

- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin

$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$



$$\text{margin} = \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

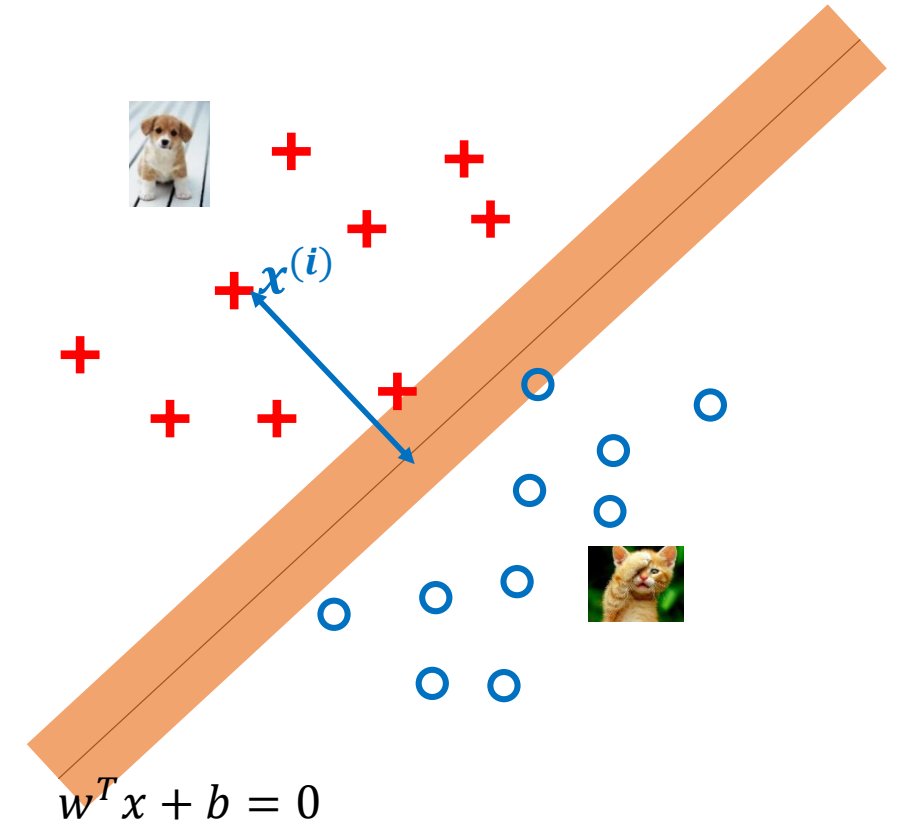
SVM

- **Linear SVM**

- A maximum margin classifier
- Preventing data points from falling into the margin

Hard margin

$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

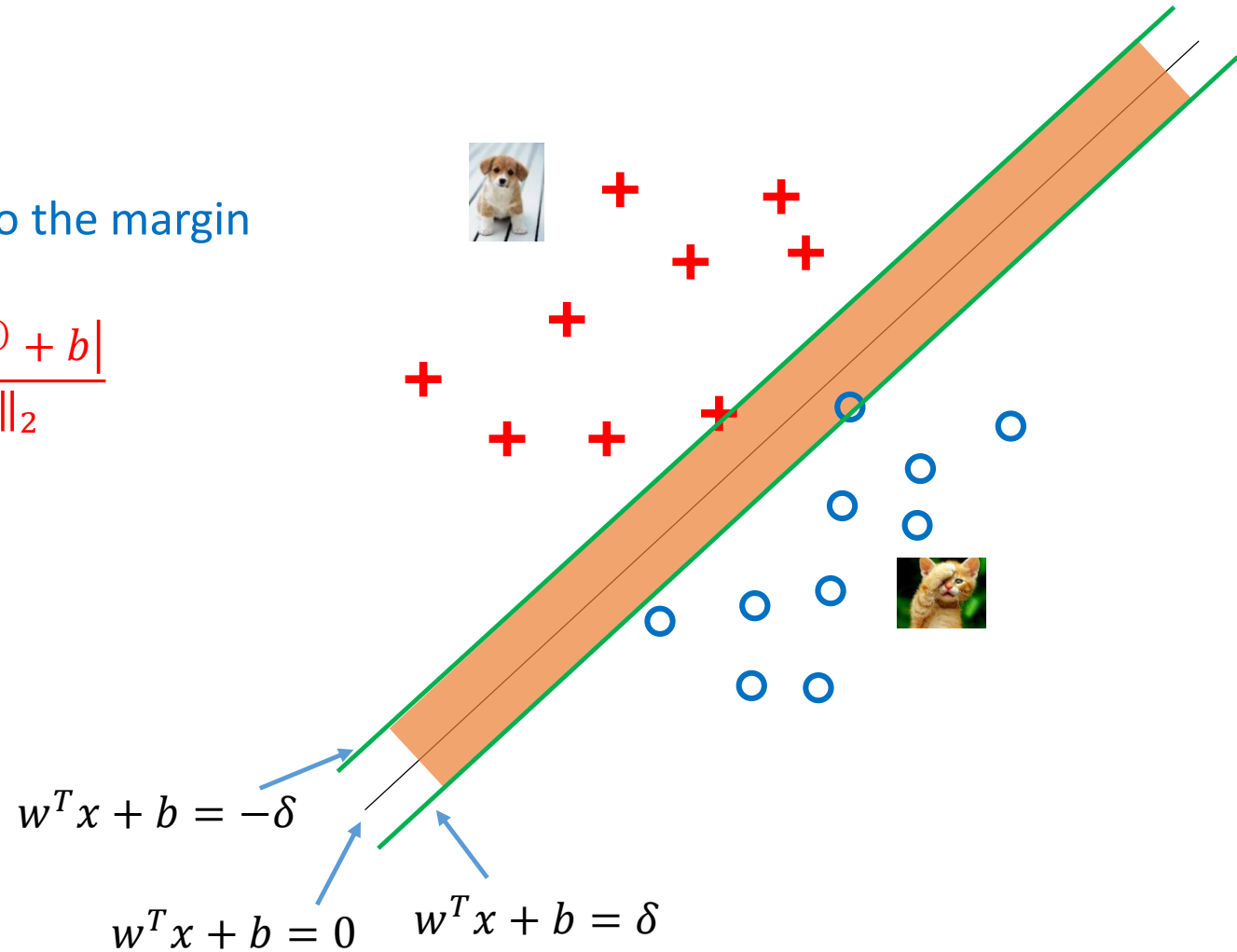


$$\text{margin} = \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

SVM

- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin

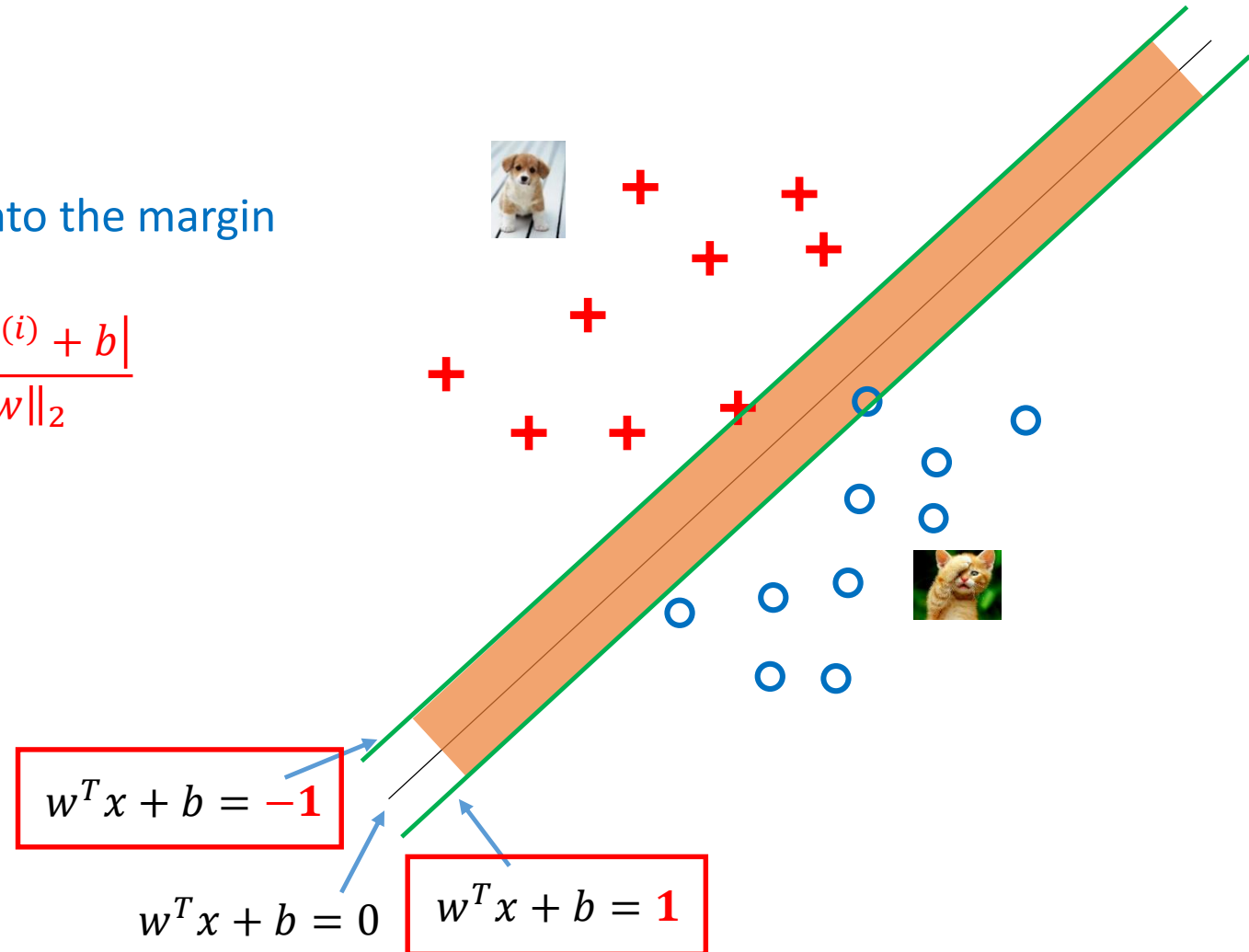
$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$



SVM

- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin

$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

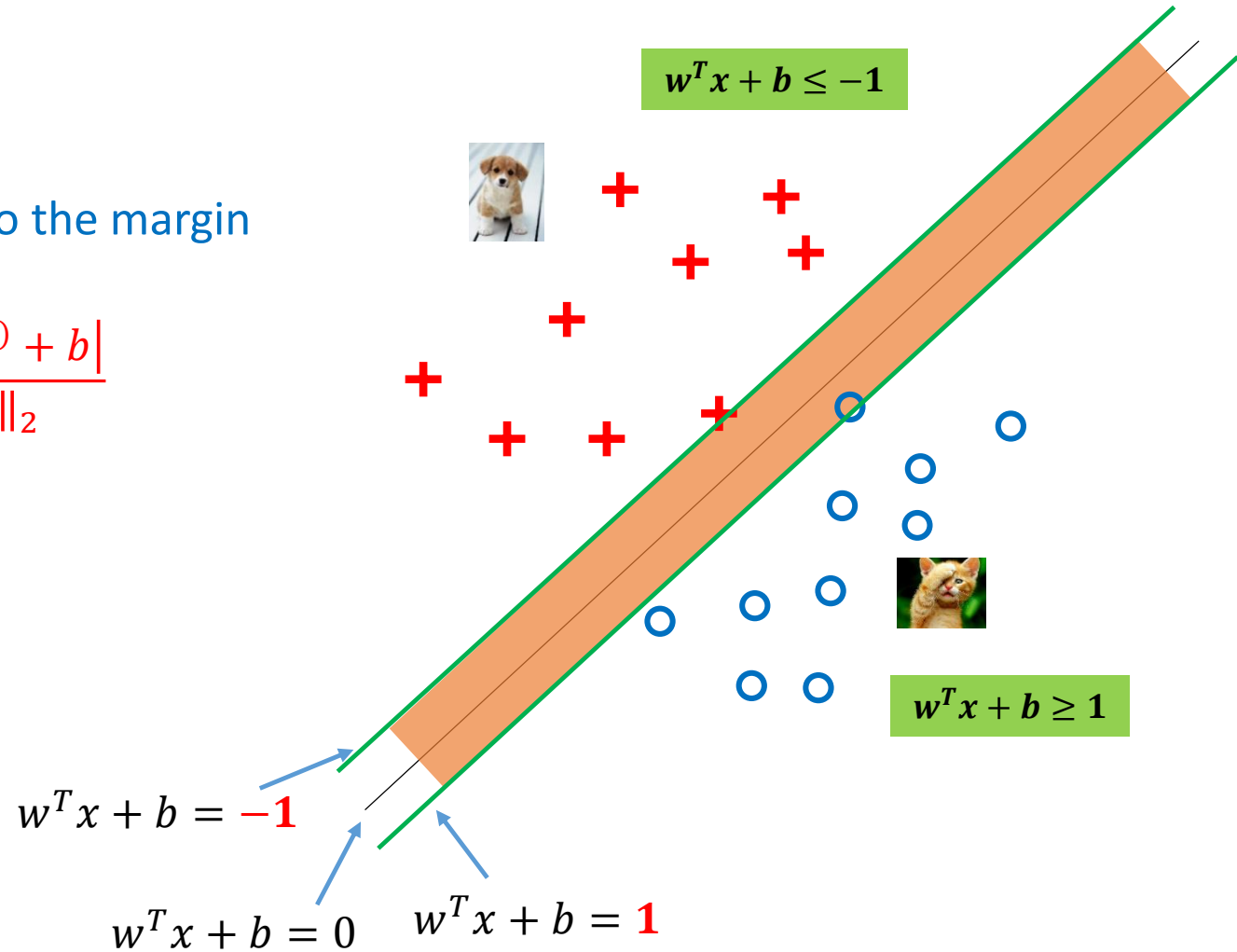


SVM

- **Linear SVM**

- A maximum margin classifier
- Preventing data points from falling into the margin

$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$



SVM

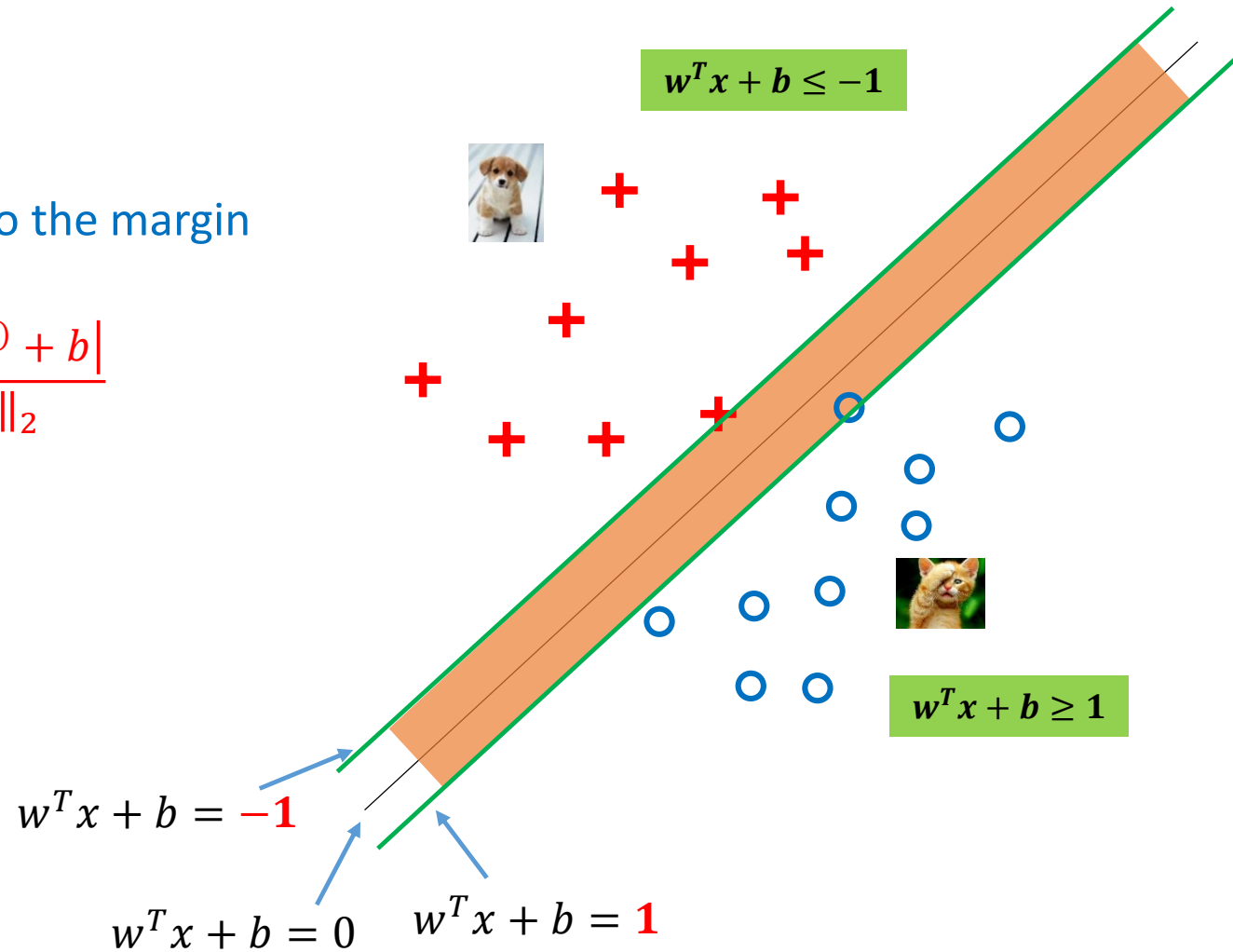
- **Linear SVM**

- A maximum margin classifier
- Preventing data points from falling into the margin

$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

For a data point $x^{(i)}$,

- if its target $t^{(i)} = 1$, $w^T x^{(i)} + b \geq 1$
- if its target $t^{(i)} = -1$, $w^T x^{(i)} + b \leq -1$



SVM

- **Linear SVM**

- A maximum margin classifier
- Preventing data points from falling into the margin

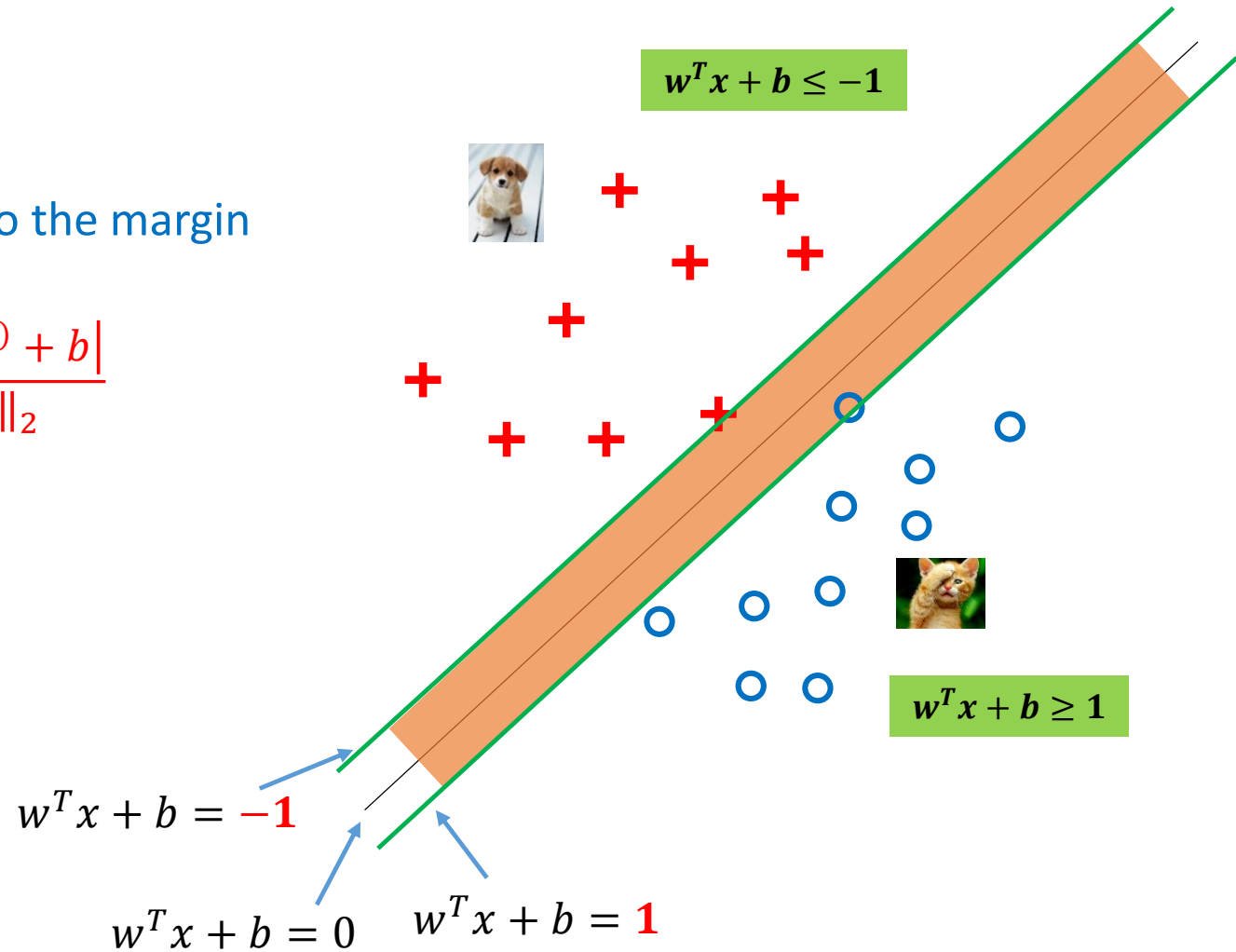
$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

For a data point $x^{(i)}$,

- if its target $t^{(i)} = 1$, $w^T x^{(i)} + b \geq 1$
- if its target $t^{(i)} = -1$, $w^T x^{(i)} + b \leq -1$

So we want

$$(w^T x^{(i)} + b) t^{(i)} \geq 1$$

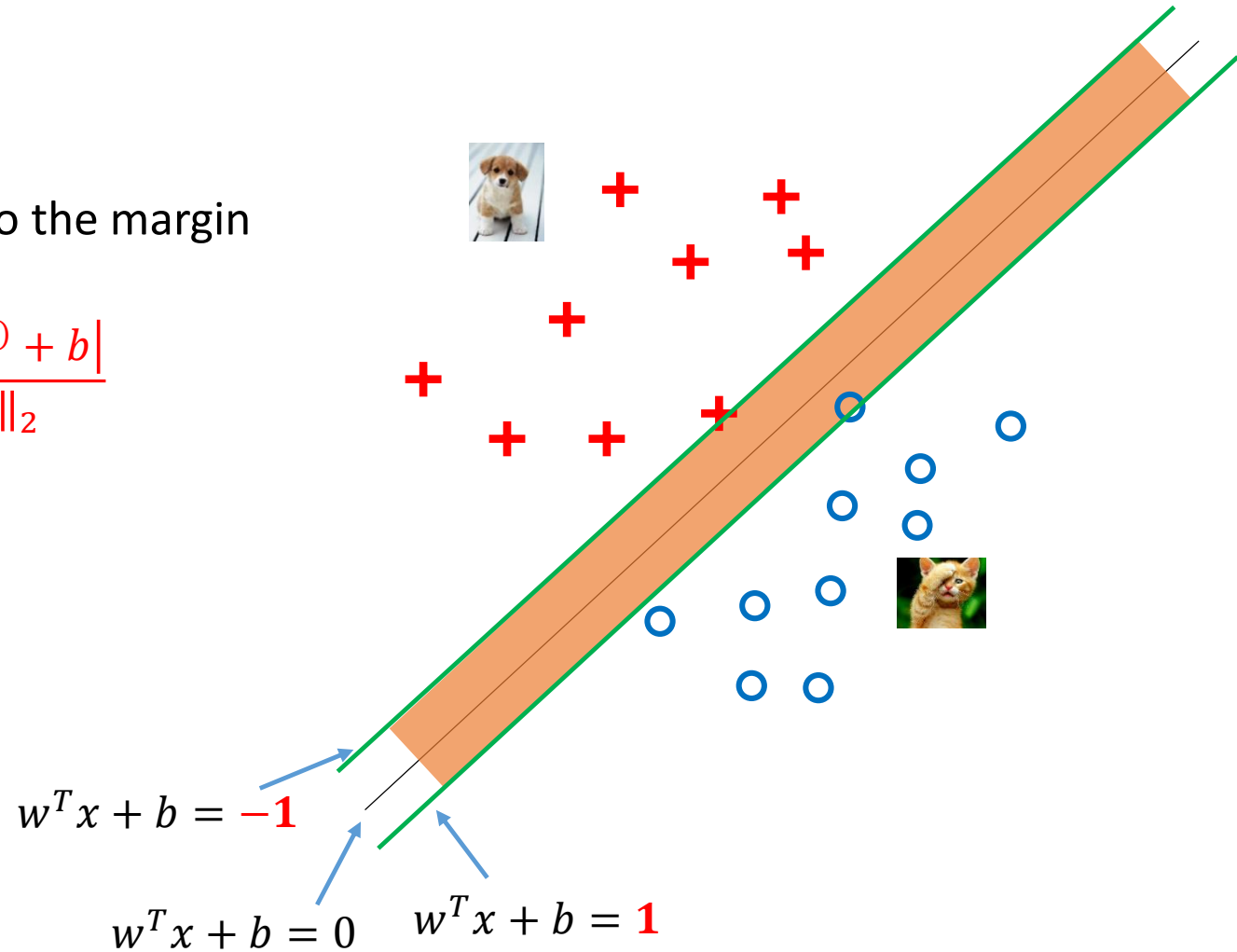


SVM

- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin

$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$



SVM

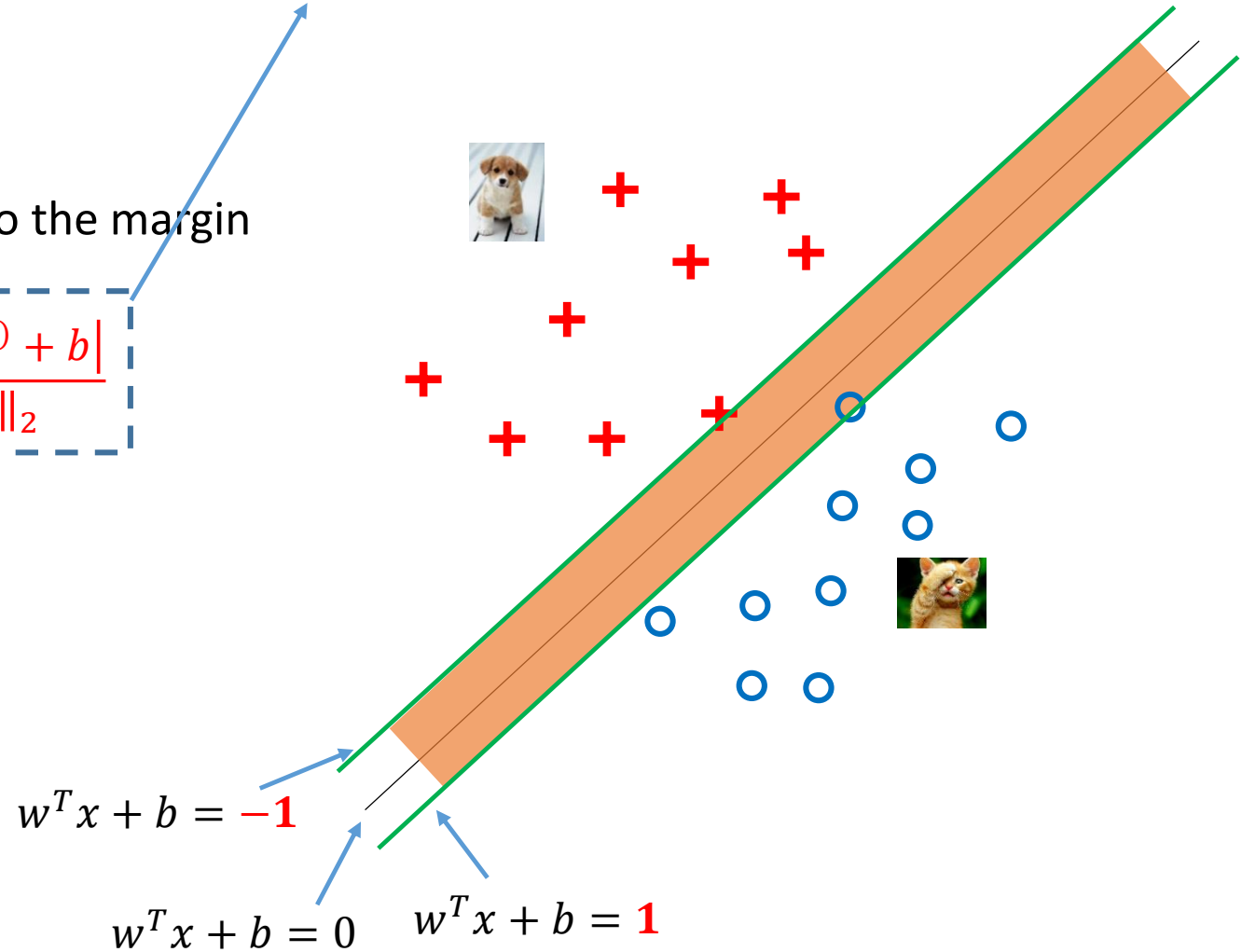
- **Linear SVM**

- A maximum margin classifier
- Preventing data points from falling into the margin

$$\max_{w,b} \text{margin} = 2 * \min_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2}$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$

$$\max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2} \rightarrow \max_{w,b} \frac{1}{\|w\|_2}$$



SVM

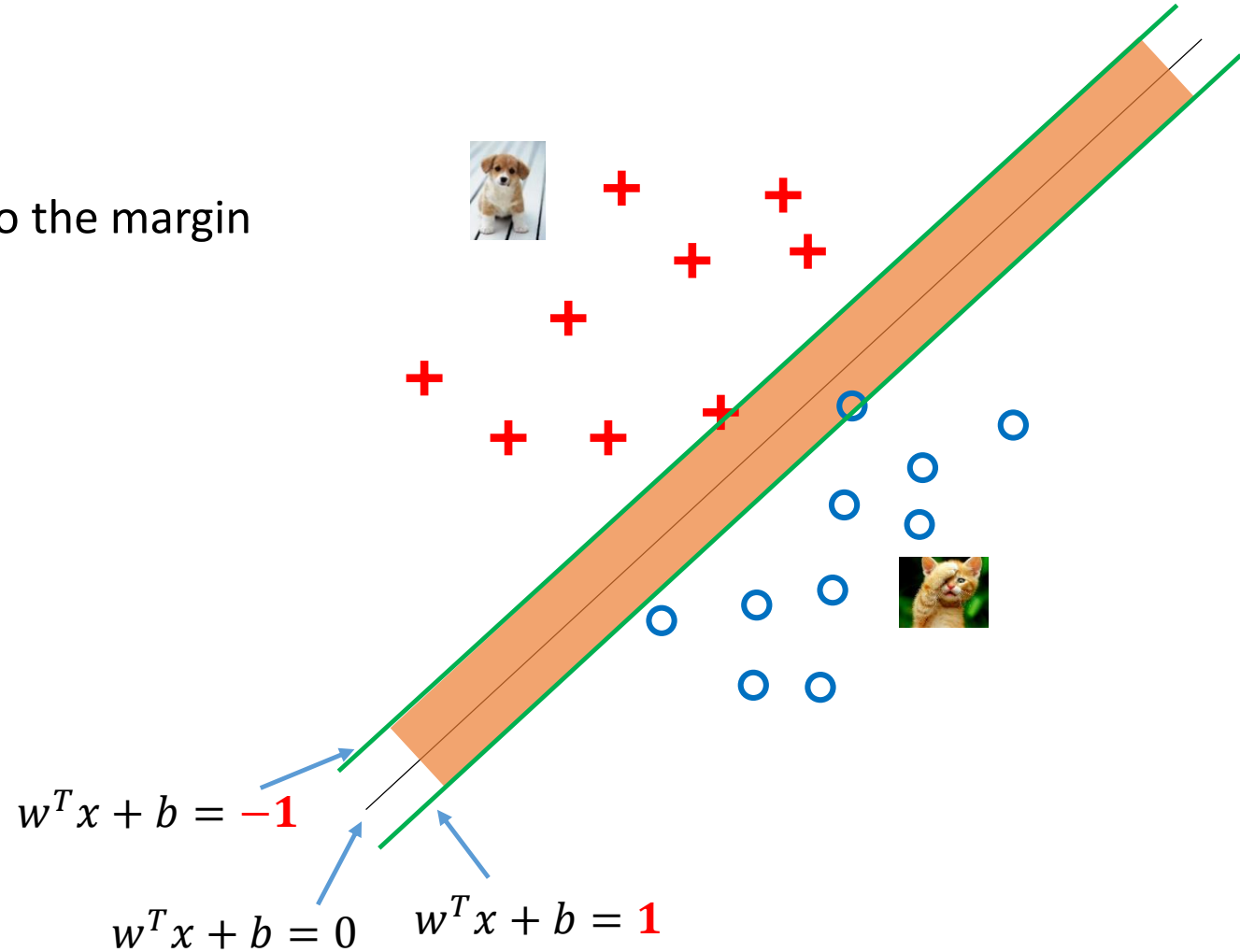
$$\max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2} \rightarrow \max_{w,b} \frac{1}{\|w\|_2}$$

- **Linear SVM**

- A maximum margin classifier
- Preventing data points from falling into the margin

$$\max_{w,b} \text{margin} = 2 * \max_{w,b} \frac{1}{\|w\|_2}$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$



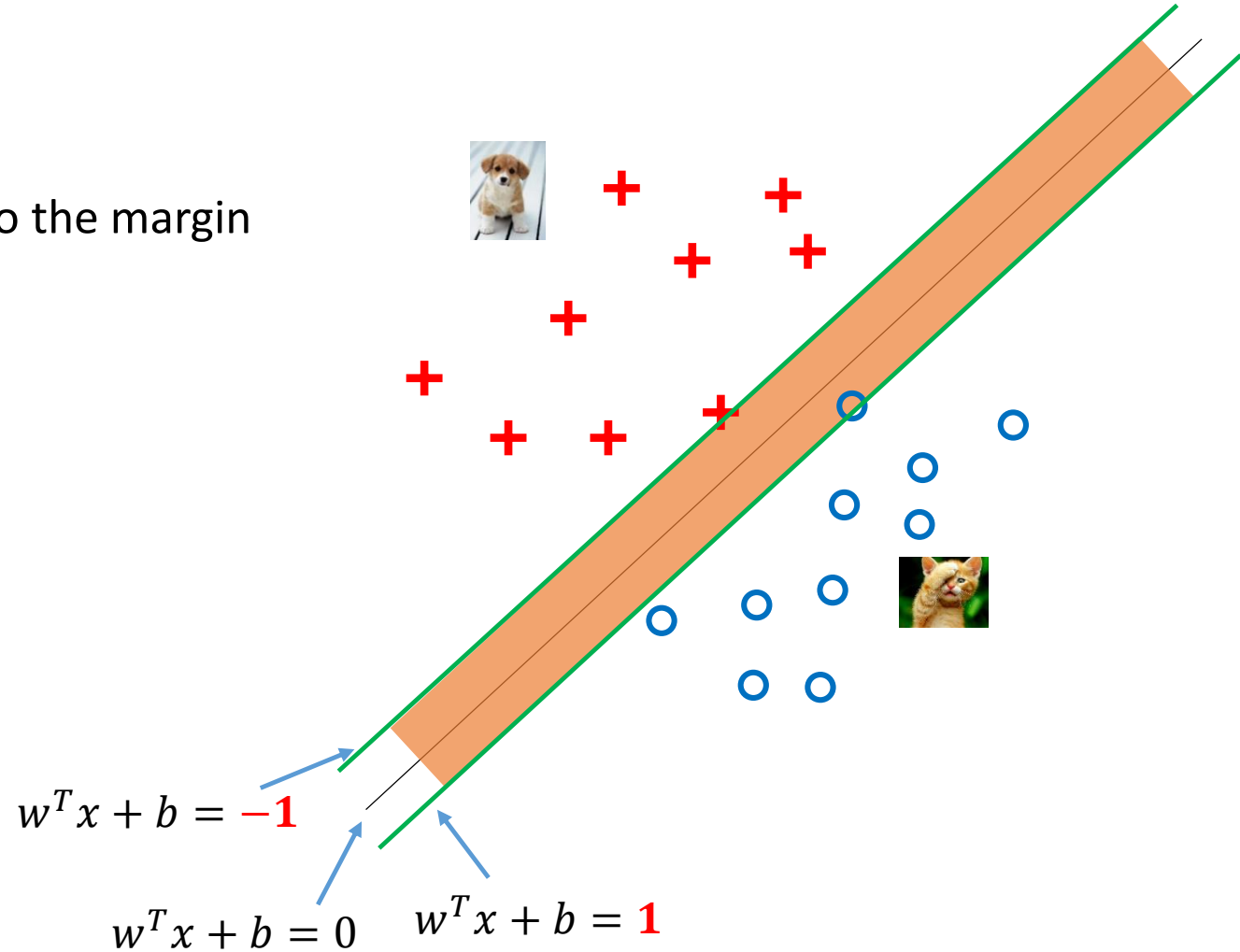
SVM

$$\max_{w,b} \min_{x^{(i)} \in D} \frac{|w^T x^{(i)} + b|}{\|w\|_2} \rightarrow \max_{w,b} \frac{1}{\|w\|_2} \rightarrow \min_{w,b} \frac{1}{2} \|w\|_2^2$$

- **Linear SVM**

- A maximum margin classifier
- Preventing data points from falling into the margin

$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$
$$s.t. \quad \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$

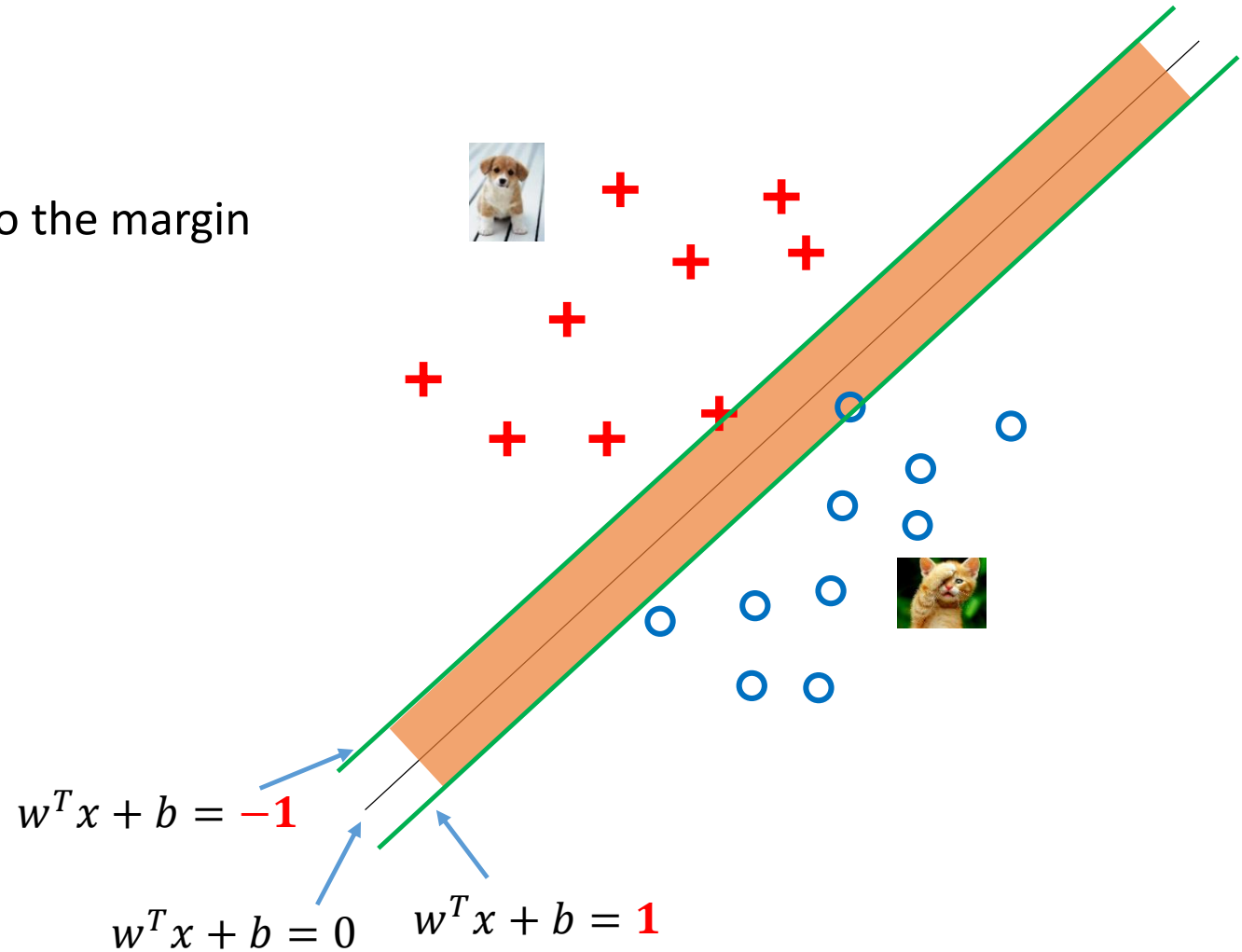


SVM

- **Linear SVM**
 - A maximum margin classifier
 - Preventing data points from falling into the margin

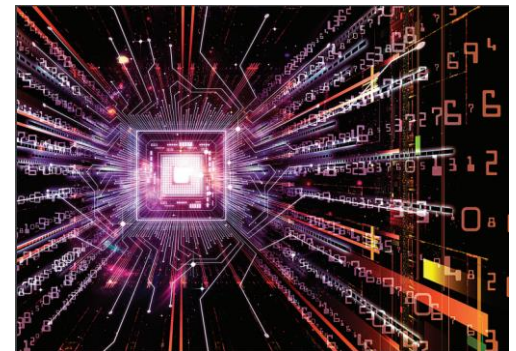
$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$
$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$

A constrained minimization



Machine Learning

- Learning

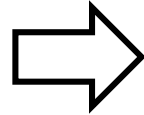


Three Steps for SVM

- Learning

Representation

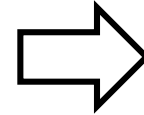
$$y(x) = w^T x + b$$



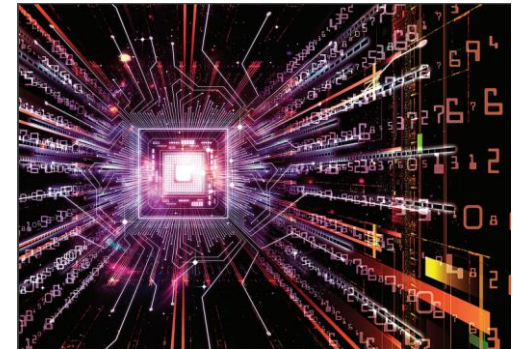
Evaluation

$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$



Optimization

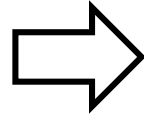


Three Steps for SVM

- Learning

Representation

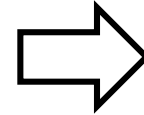
$$y(x) = w^T x + b$$



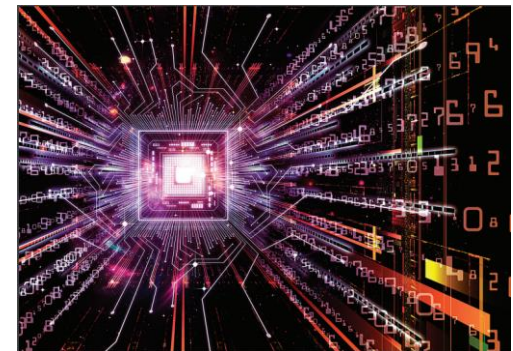
Evaluation

$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$



Optimization



$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$

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

$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$

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

Dual problem

Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - (w^T x^{(i)} + b) t^{(i)})$$

where $\alpha^{(i)}$'s are Lagrange multipliers

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

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

Dual problem

Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - (w^T x^{(i)} + b) t^{(i)})$$

where $\alpha^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. w, b for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

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

Dual problem

Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - (w^T x^{(i)} + b) t^{(i)})$$

where $\alpha^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. w, b for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

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

Dual problem

Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - (w^T x^{(i)} + b) t^{(i)})$$

where $\alpha^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. w, b for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

Then, substitute w back to function L

$$L(\alpha) = \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

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

Dual problem

Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - (w^T x^{(i)} + b) t^{(i)})$$

where $\alpha^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. w, b for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

Then, substitute w back to function L

$$L(\alpha) = \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

Next, we can obtain $\alpha^{(i)}$'s by solving the following optimization problem

$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

$$\text{s.t. } \alpha^{(i)} \geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$

Dual problem

Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - (w^T x^{(i)} + b) t^{(i)})$$

where $\alpha^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. w, b for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

Then, substitute w back to function L

$$L(\alpha) = \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

Next, we can obtain $\alpha^{(i)}$'s by solving the following optimization problem

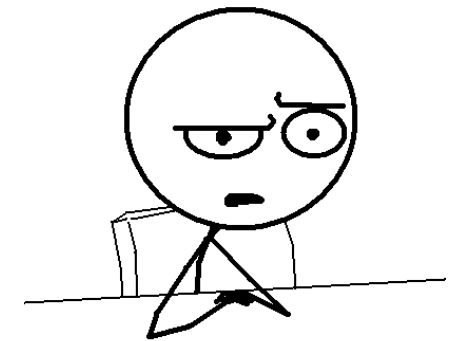
$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

$$\text{s.t. } \alpha^{(i)} \geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$

Dual problem



Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - (w^T x^{(i)} + b) t^{(i)})$$

where $\alpha^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. w, b for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

Then, substitute w back to function L

$$L(\alpha) = \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

Next, we can obtain $\alpha^{(i)}$'s by solving the following optimization problem

$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$


$$\text{s.t. } \alpha^{(i)} \geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$

Dual problem

Quadratic Programming

Find $\arg \max_{\mathbf{u}} \quad c + \mathbf{d}^T \mathbf{u} + \frac{\mathbf{u}^T R \mathbf{u}}{2}$  Quadratic criterion

Subject to

$$\left. \begin{aligned} a_{11}u_1 + a_{12}u_2 + \dots + a_{1m}u_m &\leq b_1 \\ a_{21}u_1 + a_{22}u_2 + \dots + a_{2m}u_m &\leq b_2 \\ &\vdots \\ a_{n1}u_1 + a_{n2}u_2 + \dots + a_{nm}u_m &\leq b_n \end{aligned} \right\} \begin{array}{l} n \text{ additional linear} \\ \text{inequality} \\ \text{constraints} \end{array}$$

And subject to

$$\left. \begin{aligned} a_{(n+1)1}u_1 + a_{(n+1)2}u_2 + \dots + a_{(n+1)m}u_m &= b_{(n+1)} \\ a_{(n+2)1}u_1 + a_{(n+2)2}u_2 + \dots + a_{(n+2)m}u_m &= b_{(n+2)} \\ &\vdots \\ a_{(n+e)1}u_1 + a_{(n+e)2}u_2 + \dots + a_{(n+e)m}u_m &= b_{(n+e)} \end{aligned} \right\} \begin{array}{l} e \text{ additional linear} \\ \text{equality} \\ \text{constraints} \end{array}$$

SVM

- Training a SVM model by solving

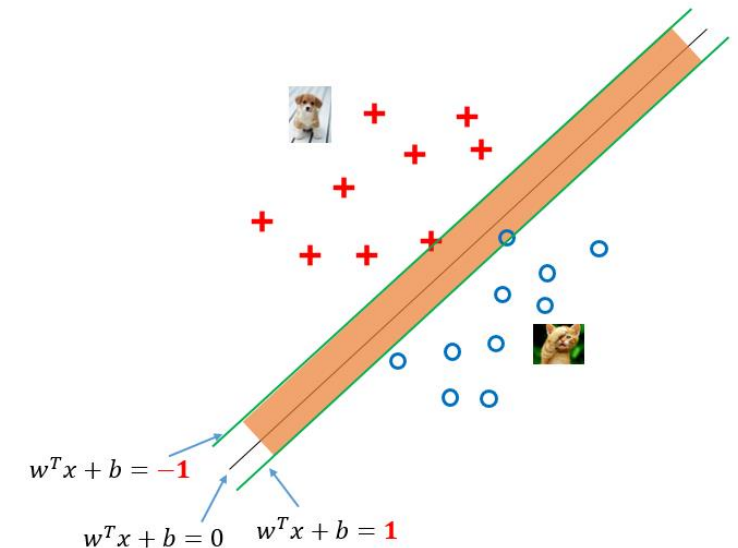
$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

s.t. $\alpha^{(i)} \geq 0$ and $\sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$

$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$

s.t. $\forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$

A constrained minimization



SVM

- Training a SVM model by solving

$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

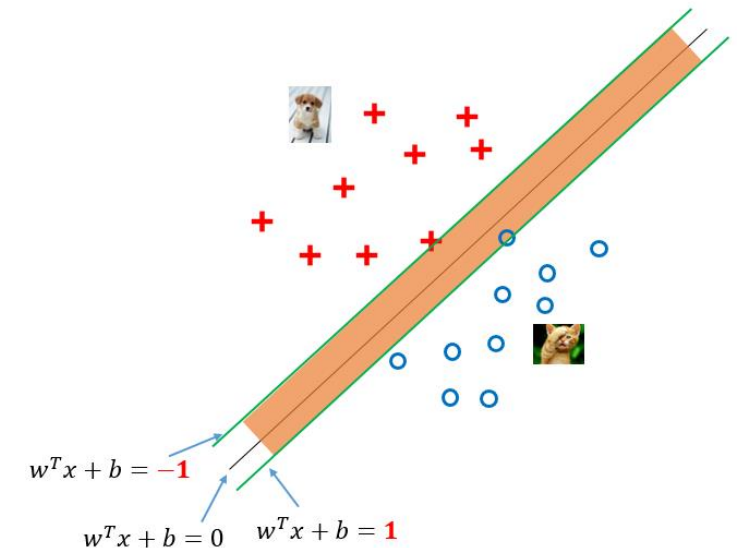
$$\text{s.t. } \alpha^{(i)} \geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

- Once $\alpha^{(i)}$ is obtained
 - The weights are

$$w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$
$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$

A constrained minimization



SVM

- Training a SVM model by solving

$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

$$\text{s.t. } \alpha^{(i)} \geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

- Once $\alpha^{(i)}$ is obtained

- The weights are

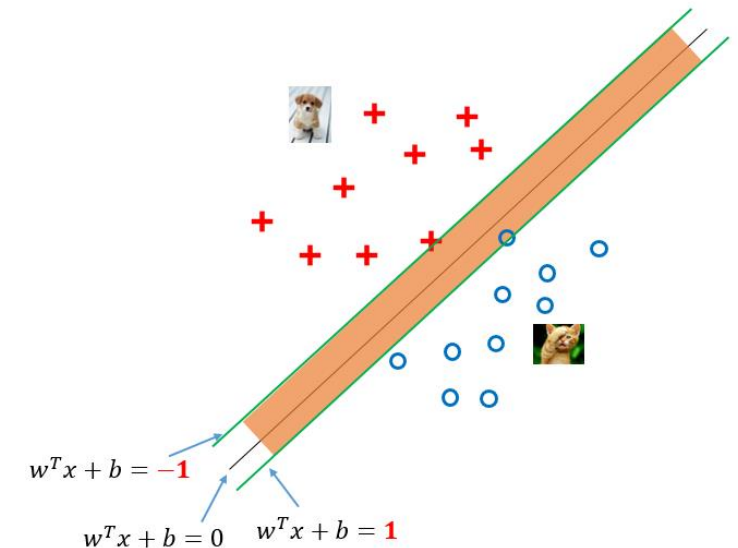
$$w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

- Prediction on a new example:

$$\begin{aligned} y^{(new)} &= \text{sgn}(w^T x^{(new)} + b) \\ &= \text{sgn}\left(\sum_{i=1}^N \alpha^{(i)} t^{(i)} (x^{(i)T} x^{(new)}) + b\right) \end{aligned}$$

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|_2^2 \\ \text{s.t.} \quad & \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 \end{aligned}$$

A constrained minimization



SVM

Only a small subset of $\alpha^{(i)}$'s will be nonzero, and the corresponding $x^{(i)}$'s are the **Support Vectors**.

- Training a SVM model by solving

$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

$$\text{s.t. } \alpha^{(i)} \geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

- Once $\alpha^{(i)}$ is obtained

- The weights are

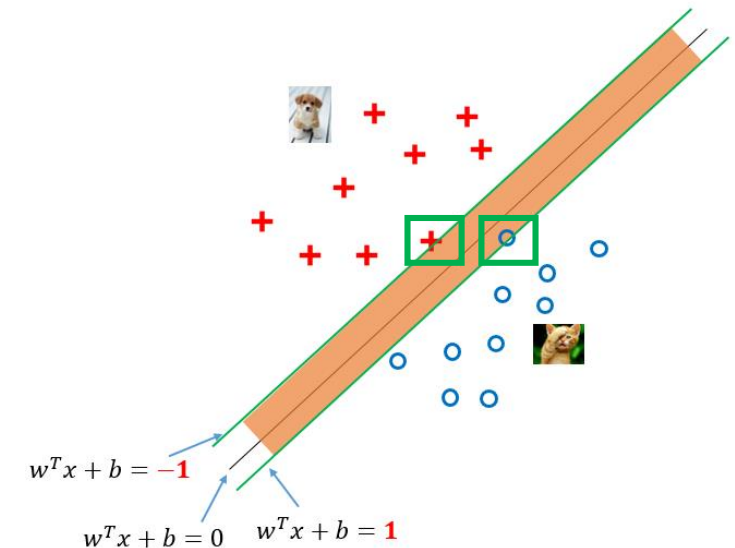
$$w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

- Prediction on a new example:

$$\begin{aligned} y^{(new)} &= \text{sgn}(w^T x^{(new)} + b) \\ &= \text{sgn}\left(\sum_{i=1}^N \alpha^{(i)} t^{(i)} (x^{(i)T} x^{(new)}) + b\right) \end{aligned}$$

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|_2^2 \\ \text{s.t.} \quad & \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 \end{aligned}$$

A constrained minimization



SVM

Only a small subset of $\alpha^{(i)}$'s will be nonzero, and the corresponding $x^{(i)}$'s are the **Support Vectors**.

Note that both the learning objective and the decision function depend only on dot products between datapoints

- Training a SVM model by solving

$$\begin{aligned} \max_{\alpha^{(i)}} L(\alpha) &= \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} \left(x^{(i)T} x^{(j)} \right) \\ \text{s.t. } \alpha^{(i)} &\geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0 \end{aligned}$$

- Once $\alpha^{(i)}$ is obtained

- The weights are

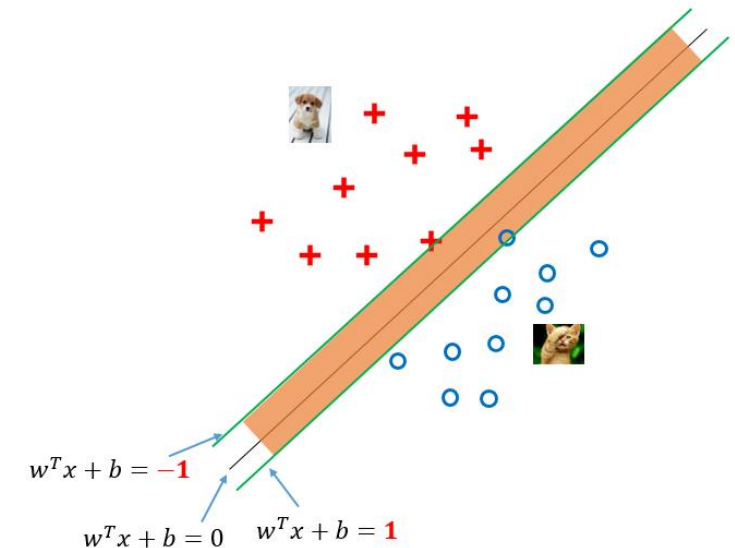
$$w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

- Prediction on a new example:

$$\begin{aligned} y^{(new)} &= \text{sgn}(w^T x^{(new)} + b) \\ &= \text{sgn} \left(\sum_{i=1}^N \alpha^{(i)} t^{(i)} \left(x^{(i)T} x^{(new)} \right) + b \right) \end{aligned}$$

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|_2^2 \\ \text{s.t.} \quad & \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 \end{aligned}$$

A constrained minimization



SVM

Only a small subset of $\alpha^{(i)}$'s will be nonzero, and the corresponding $x^{(i)}$'s are the **Support Vectors**.

Note that both the learning objective and the decision function depend only on dot products between datapoints

- Training a SVM model by solving

$$\begin{aligned} \max_{\alpha^{(i)}} L(\alpha) &= \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} \left(x^{(i)T} x^{(j)} \right) \\ \text{s.t. } \alpha^{(i)} &\geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0 \end{aligned}$$

`svm.fit()`

- Once $\alpha^{(i)}$ is obtained
 - The weights are

$$w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

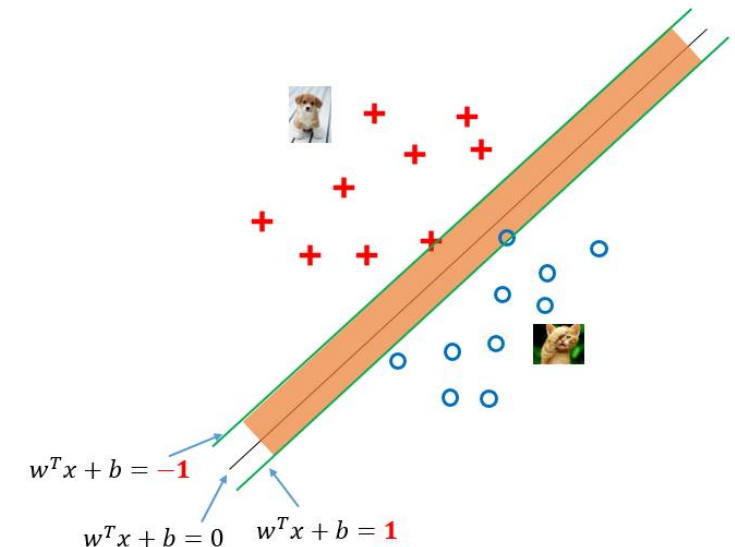
- Prediction on a new example:

$$\begin{aligned} y^{(new)} &= \text{sgn}(w^T x^{(new)} + b) \\ &= \text{sgn} \left(\sum_{i=1}^N \alpha^{(i)} t^{(i)} \left(x^{(i)T} x^{(new)} \right) + b \right) \end{aligned}$$

`svm.predict()`

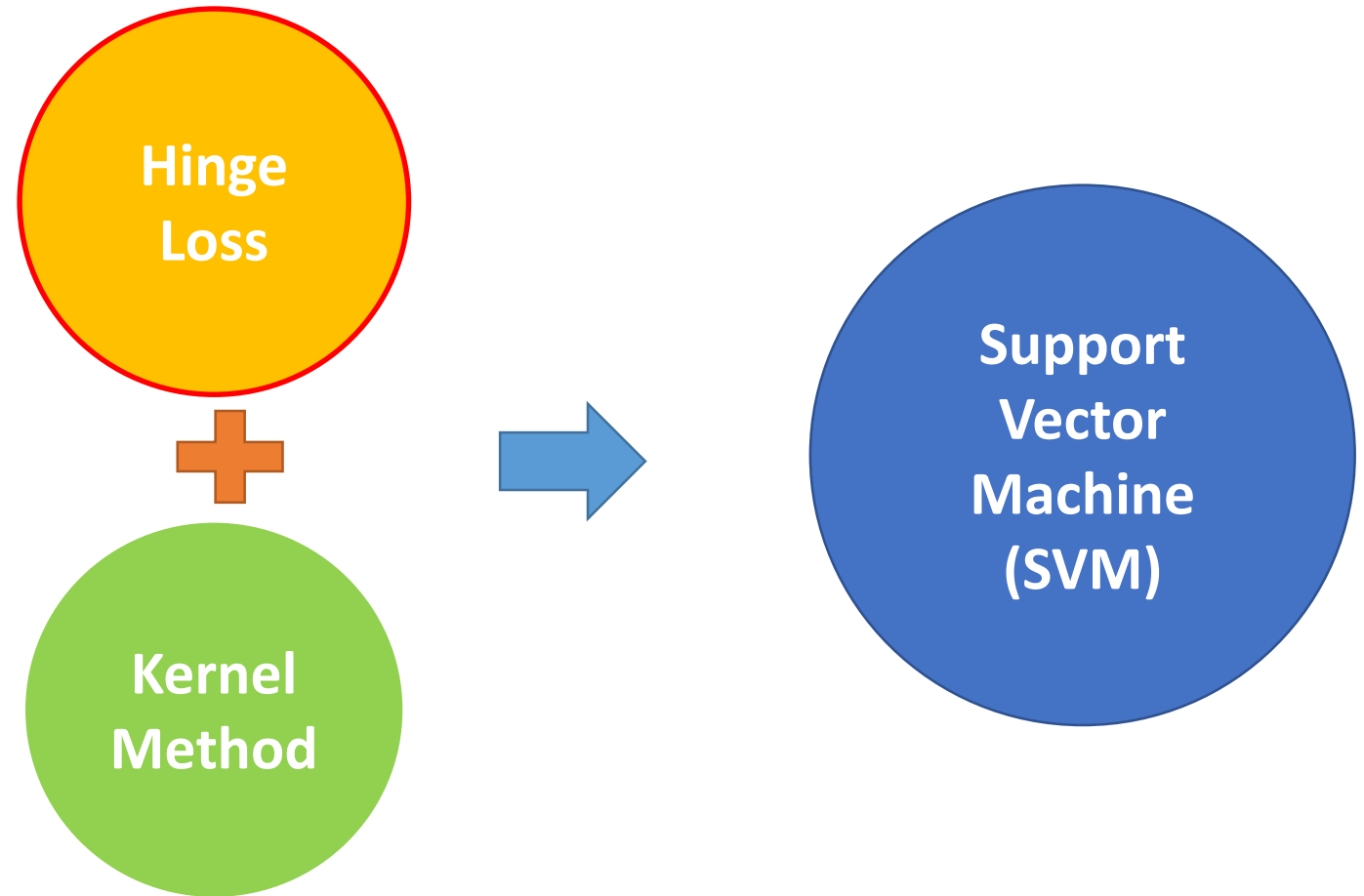
$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|_2^2 \\ \text{s.t.} \quad & \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 \end{aligned}$$

A constrained minimization

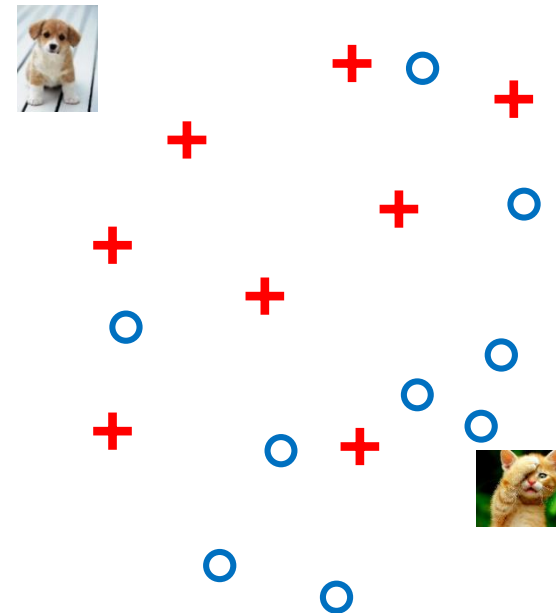


SVM

- Linear SVM
- Soft Margin (Non-linearly separable)
 - Hinge Loss
- Kernel Trick

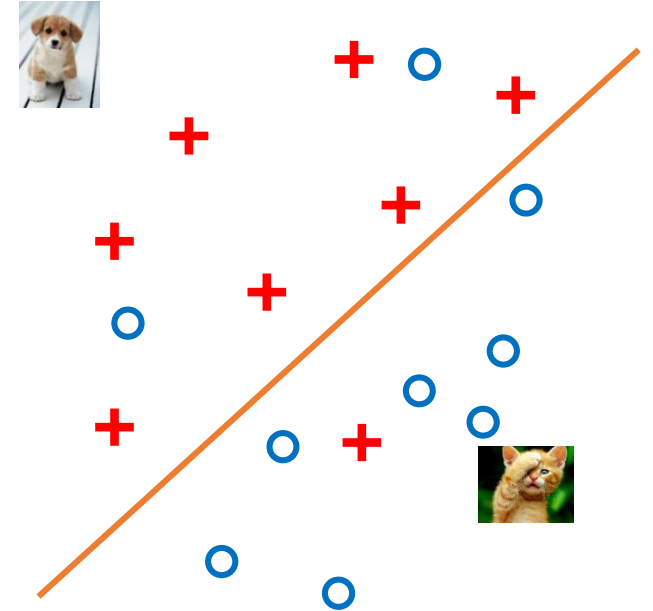


Non-Separable



Non-Separable

- What should we do?

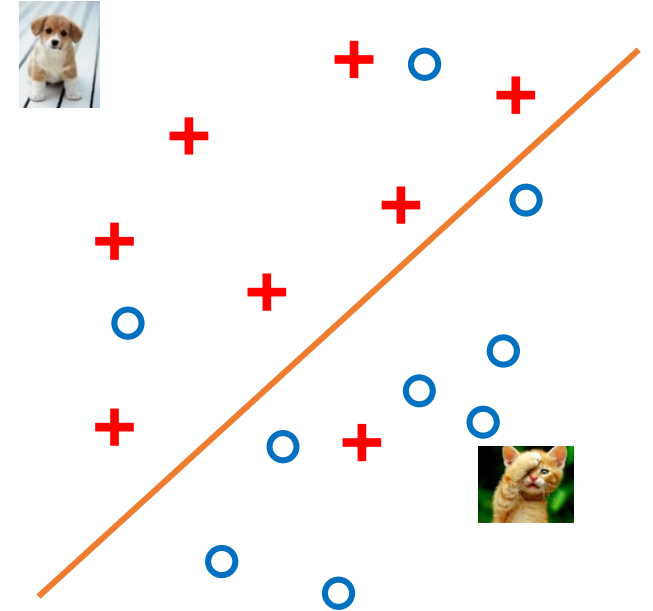


Non-Separable

- What should we do?

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + \textit{penalty}$$

Balance the tradeoff between margin and classification errors



Non-Separable

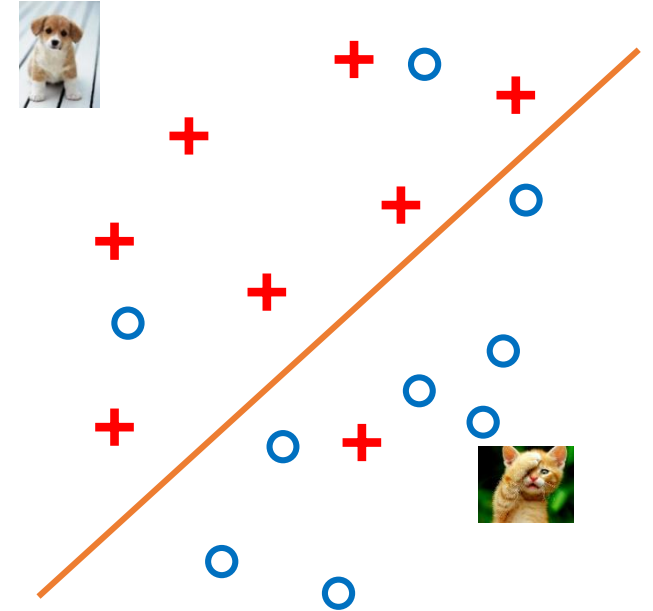
- What should we do?

- Idea 1:

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \times (\#train\ errors)$$

- Idea 2:

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \times (distance\ of\ error\ points\ to\ their\ correct\ place)$$



Non-Separable

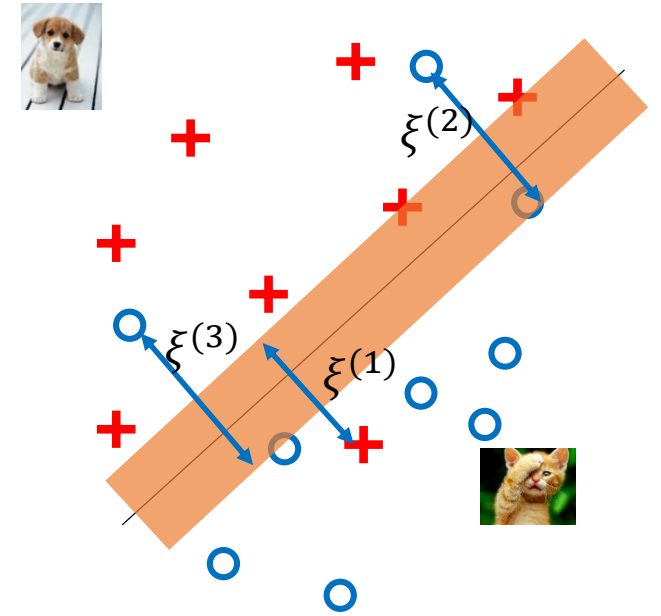
- What should we do?

- Idea 1:

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \times (\#train\ errors)$$

- Idea 2:

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \times (\text{distance of error points to their correct place})$$



Non-Separable

- What should we do?

- Idea 1:

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \times (\#train\ errors)$$

- Idea 2:

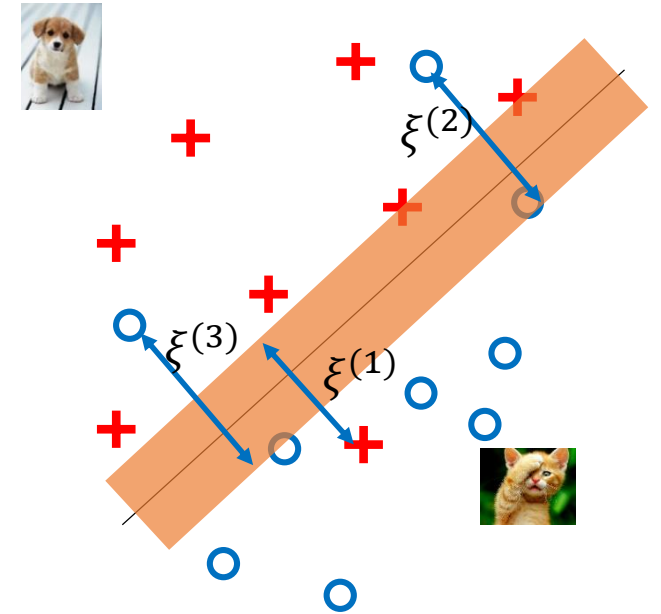
$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$s.t. (w^T x^{(1)} + b)t^{(1)} \geq 1 - \xi^{(1)}, \xi^{(1)} \geq 0$$

$$(w^T x^{(2)} + b)t^{(2)} \geq 1 - \xi^{(2)}, \xi^{(2)} \geq 0$$

$$\vdots$$

$$(w^T x^{(N)} + b)t^{(N)} \geq 1 - \xi^{(N)}, \xi^{(N)} \geq 0$$



Non-Separable

- What should we do?

- Idea 1:

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \times (\#train\ errors)$$

- Idea 2:

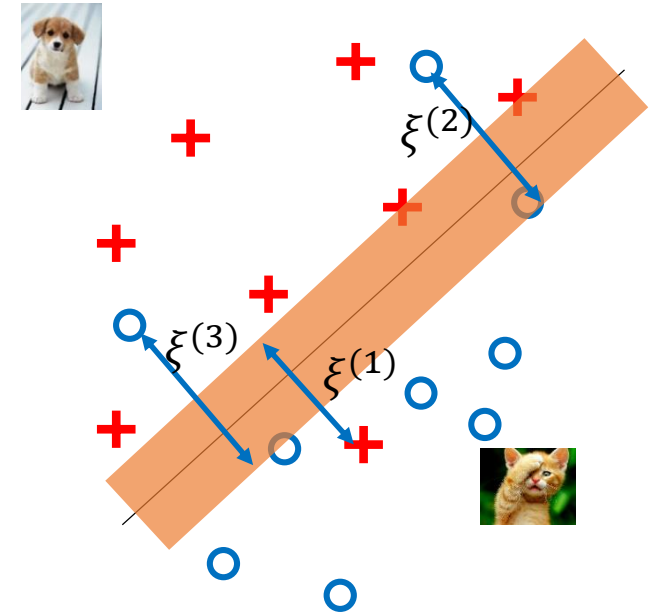
$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$s.t. (w^T x^{(1)} + b)t^{(1)} \geq 1 - \xi^{(1)}, \xi^{(1)} \geq 0$$

$$(w^T x^{(2)} + b)t^{(2)} \geq 1 - \xi^{(2)}, \xi^{(2)} \geq 0$$

\vdots

$$(w^T x^{(N)} + b)t^{(N)} \geq 1 - \xi^{(N)}, \xi^{(N)} \geq 0$$



$\xi^{(i)}$: slack variables

Non-Separable

- What should we do?

- Idea 1:

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \times (\#train\ errors)$$

- Idea 2:

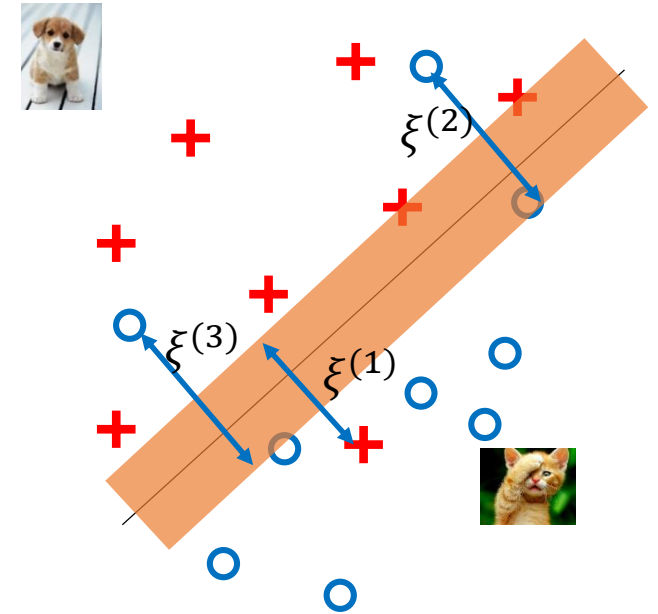
$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$s.t. (w^T x^{(1)} + b)t^{(1)} \geq 1 - \xi^{(1)}, \xi^{(1)} \geq 0$$

$$(w^T x^{(2)} + b)t^{(2)} \geq 1 - \xi^{(2)}, \xi^{(2)} \geq 0$$

\vdots

$$(w^T x^{(N)} + b)t^{(N)} \geq 1 - \xi^{(N)}, \xi^{(N)} \geq 0$$



Soft margin

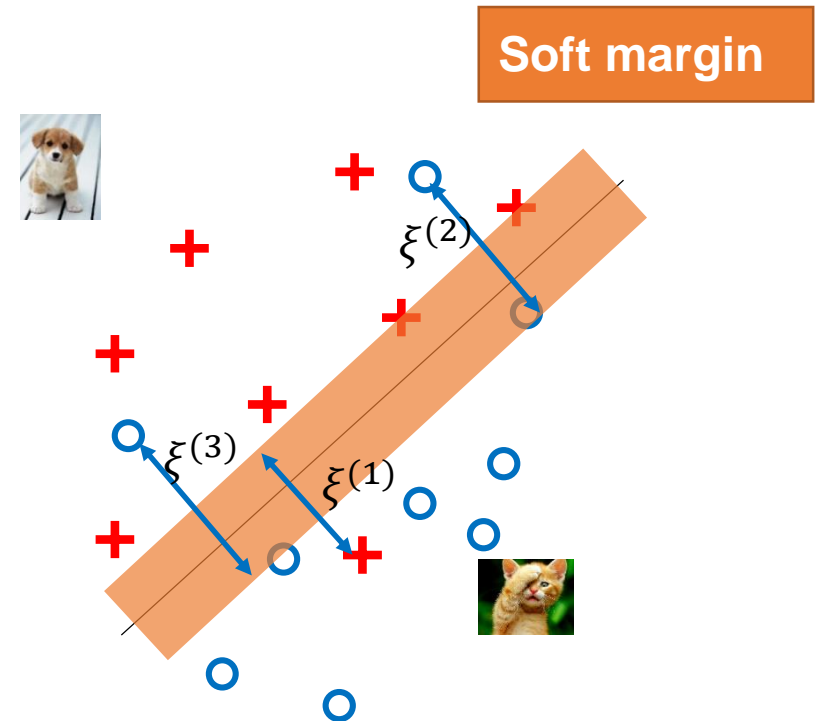
$\xi^{(i)}$: slack variables

SVM

- **Linear SVM with slack variables**
 - A maximum margin classifier
 - Considering the tradeoff between margin and classification errors

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

A constrained minimization

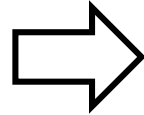


Three Steps for SVM

- Learning

Representation

$$y(x) = w^T x + b$$

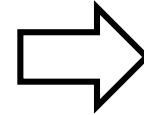


Evaluation

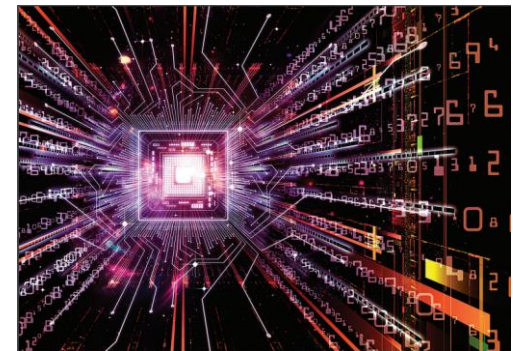
$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$s.t. \quad \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 - \xi^{(i)}$$

$$\forall i, \xi^{(i)} \geq 0$$



Optimization





Dual problem

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

A constrained minimization

Lagrange function

$$L(w, b, \xi, \alpha, \eta) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)} + \sum_{i=1}^N \alpha^{(i)} (1 - \xi^{(i)} - (w^T x^{(i)} + b) t^{(i)}) - \sum_{i=1}^N \eta^{(i)} \xi^{(i)}$$

where $\alpha^{(i)}$'s and $\eta^{(i)}$'s are Lagrange multipliers

Dual problem

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

A constrained minimization

Lagrange function

$$L(w, b, \xi, \alpha, \eta) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)} + \sum_{i=1}^N \alpha^{(i)} (1 - \xi^{(i)} - (w^T x^{(i)} + b) t^{(i)}) - \sum_{i=1}^N \eta^{(i)} \xi^{(i)}$$

where $\alpha^{(i)}$'s and $\eta^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. $w, b, \xi^{(i)}$ for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0$$

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial \xi^{(i)}} = C - \alpha^{(i)} - \eta^{(i)} = 0$$

Dual problem

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

A constrained minimization

Lagrange function

$$L(w, b, \xi, \alpha, \eta) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)} + \sum_{i=1}^N \alpha^{(i)} (1 - \xi^{(i)} - (w^T x^{(i)} + b) t^{(i)}) - \sum_{i=1}^N \eta^{(i)} \xi^{(i)}$$

where $\alpha^{(i)}$'s and $\eta^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. $w, b, \xi^{(i)}$ for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial \xi^{(i)}} = C - \alpha^{(i)} - \eta^{(i)} = 0$$

Dual problem

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

A constrained minimization

Lagrange function

$$L(w, b, \xi, \alpha, \eta) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)} + \sum_{i=1}^N \alpha^{(i)} (1 - \xi^{(i)} - (w^T x^{(i)} + b) t^{(i)}) - \sum_{i=1}^N \eta^{(i)} \xi^{(i)}$$

where $\alpha^{(i)}$'s and $\eta^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. $w, b, \xi^{(i)}$ for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0 \quad \frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial \xi^{(i)}} = C - \alpha^{(i)} - \eta^{(i)} = 0$$

Then, substitute w back to function L

$$L(\alpha) = \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

Dual problem

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

A constrained minimization

Lagrange function

$$L(w, b, \xi, \alpha, \eta) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)} + \sum_{i=1}^N \alpha^{(i)} (1 - \xi^{(i)} - (w^T x^{(i)} + b) t^{(i)}) - \sum_{i=1}^N \eta^{(i)} \xi^{(i)}$$

where $\alpha^{(i)}$'s and $\eta^{(i)}$'s are Lagrange multipliers

First, minimize function L w.r.t. $w, b, \xi^{(i)}$ for fixed Lagrange multipliers

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial w} = w - \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial b} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0 \quad \frac{\partial L(w, b, \xi, \alpha, \eta)}{\partial \xi^{(i)}} = C - \alpha^{(i)} - \eta^{(i)} = 0$$

Then, substitute w back to function L

$$L(\alpha) = \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

Next, we can obtain $\alpha^{(i)}$'s by solving the following optimization problem

$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (x^{(i)T} x^{(j)})$$

$$\text{s.t. } \alpha^{(i)} \in [0, C] \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

Dual problem

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

A constrained minimization

SVM

- Let's view SVM in another way

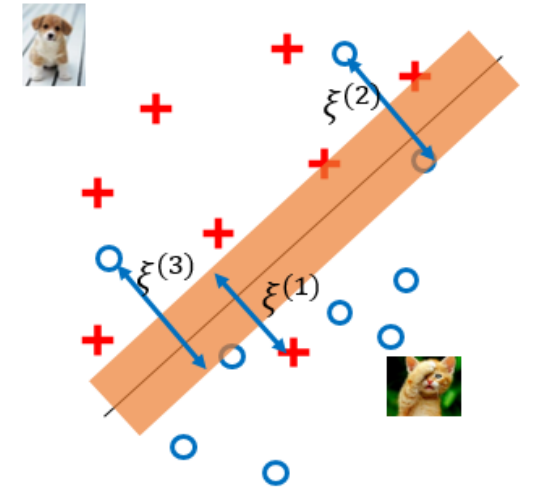
$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + \text{penalty}$$

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

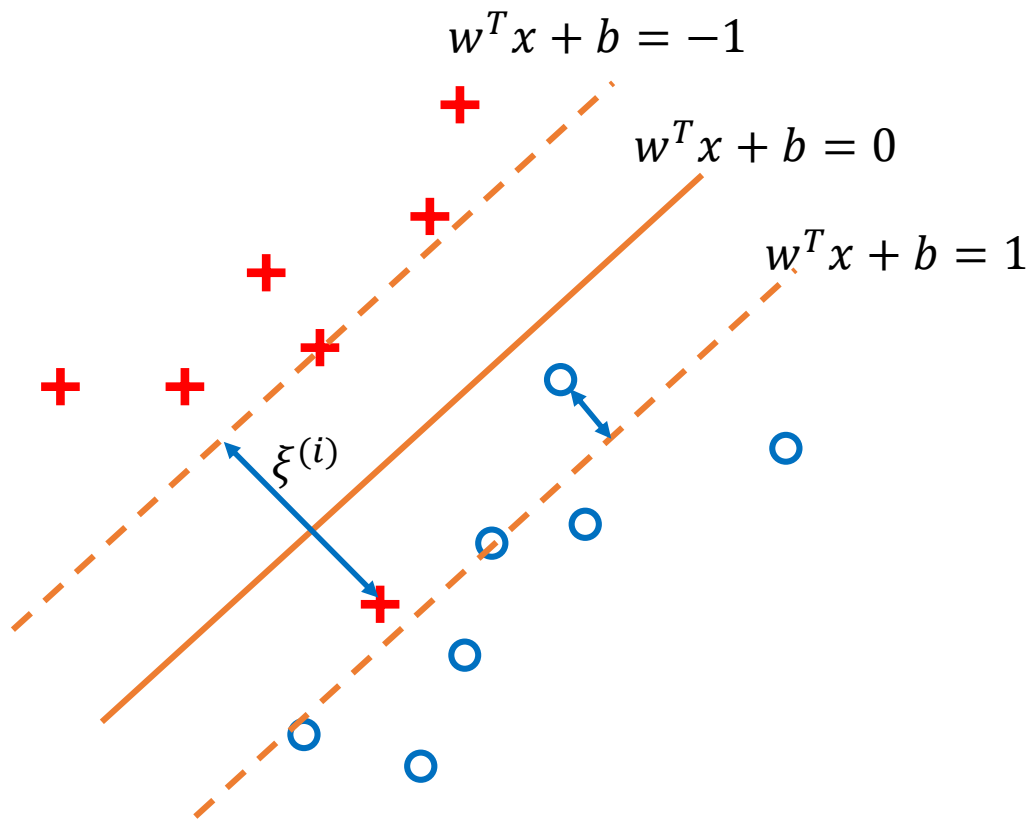
$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 - \xi^{(i)}$$

$$\forall i, \xi^{(i)} \geq 0$$

A constrained minimization

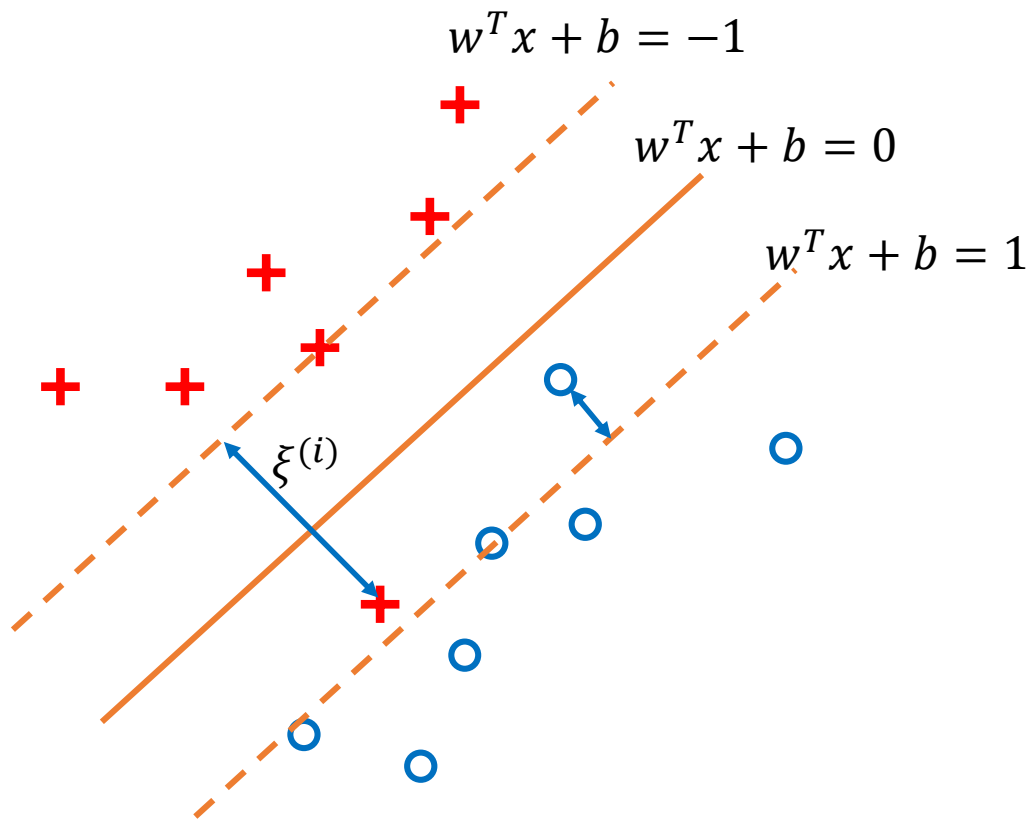


Non-Separable



- **Case 1:** $\xi^{(i)} \geq 1$
Misclassification
- **Case 2:** $0 < \xi^{(i)} < 1$
 $x^{(i)}$ is correctly classified, but lies inside the margin
- **Case 3:** $\xi^{(i)} = 0$
 $x^{(i)}$ is correctly classified, but lies outside the margin

Non-Separable



Hinge loss: $\max(0, 1 - t \cdot y)$

$$\varepsilon^{(i)} = \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

- **Case 1:** $\xi^{(i)} \geq 1$
Misclassification
- **Case 2:** $0 < \xi^{(i)} < 1$
 $x^{(i)}$ is correctly classified, but lies inside the margin
- **Case 3:** $\xi^{(i)} = 0$
 $x^{(i)}$ is correctly classified, but lies outside the margin

SVM

- Let's view SVM in another way

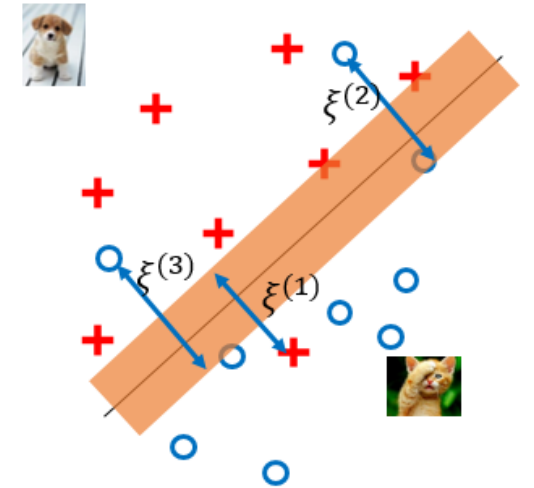
$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + \text{penalty}$$

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 - \xi^{(i)}$$

$$\forall i, \xi^{(i)} \geq 0$$

A constrained minimization



SVM

- Let's view SVM in another way

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + \text{penalty}$$

- For training data $\{(x^{(i)}, t^{(i)})\}_{i=1}^N$, use the following penalty

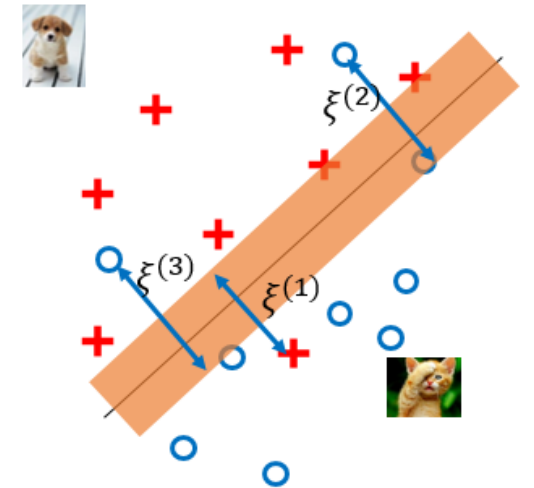
$$\max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b)t^{(i)} \geq 1 - \xi^{(i)}$$

$$\forall i, \xi^{(i)} \geq 0$$

A constrained minimization



SVM

- Let's view SVM in another way

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + \text{penalty}$$

- For training data $\{(x^{(i)}, t^{(i)})\}_{i=1}^N$, use the following penalty

$$\max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

- So we can define a loss function as

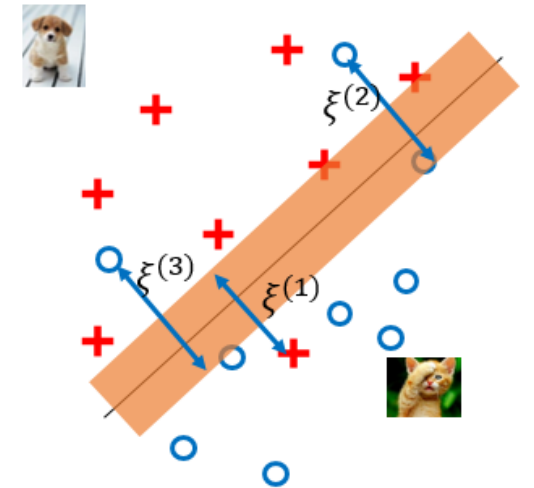
$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b)t^{(i)} \geq 1 - \xi^{(i)}$$

$$\forall i, \xi^{(i)} \geq 0$$

A constrained minimization



Logistic Regression

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: 0 or 1)

- Linear model

$$y(x) = \sigma(w_0 + w_1x_1 + w_2x_2 + \dots + w_Mx_M)$$

- Parameters

$$w_0, w_1, w_2, \dots, w_M$$

- Loss function

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

- Goal: minimize $\ell(w)$

Steps:

- Initialize w (e.g., randomly)
- Repeatedly update w based on the gradient

$$w = w - \epsilon \nabla_w \ell(w)$$

where ϵ is the learning rate.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Logistic Regression

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: 0 or 1)

- Linear model

$$y(x) = \sigma(w_0 + w_1x_1 + w_2x_2 + \dots + w_Mx_M)$$

- Parameters

$$w_0, w_1, w_2, \dots, w_M$$

- Loss function

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

- Goal: minimize $\ell(w)$

Steps:

- Initialize w (e.g., randomly)
- Repeatedly update w based on the gradient

$$w = w - \epsilon \nabla_w \ell(w)$$

where ϵ is the learning rate.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Linear SVM

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: 0 or 1)

- Linear model

$$y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$$

- Parameters

$$b, w_1, w_2, \dots, w_M$$

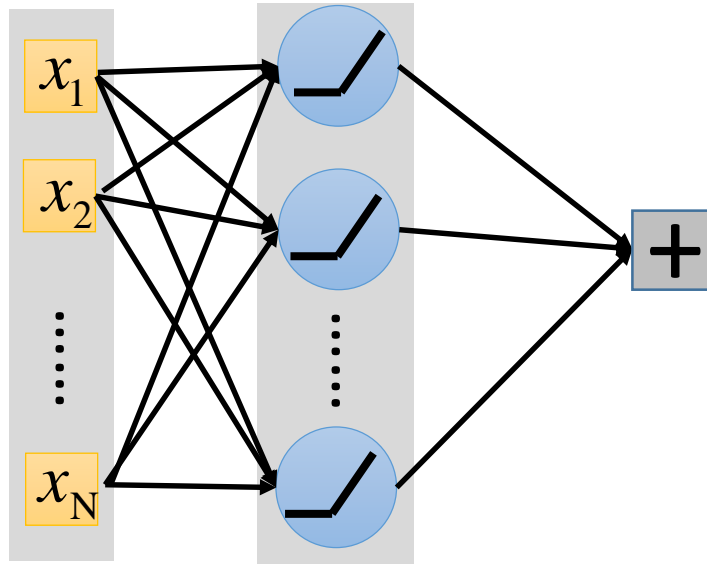
- Loss function

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b) t^{(i)})$$

- Goal: minimize $\ell(w)$

Linear SVM

Input



- Training datasets
 $\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$
(target $t^{(i)}$: 0 or 1)

- Linear model
$$y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$$

- Parameters
$$b, w_1, w_2, \dots, w_M$$

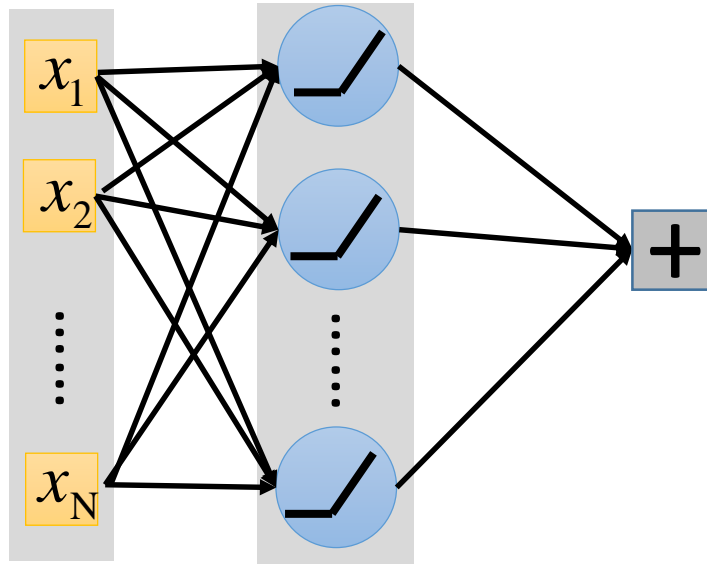
- Loss function

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

- Goal: minimize $\ell(w)$

Linear SVM

Input



Recall ReLU, Maxout Network

- Training datasets
 $\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$
(target $t^{(i)}$: 0 or 1)

- Linear model
$$y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$$

- Parameters
$$b, w_1, w_2, \dots, w_M$$

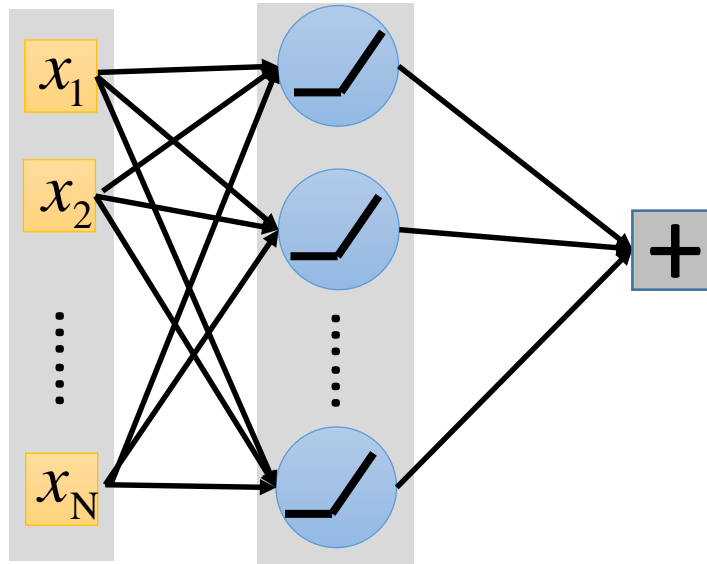
- Loss function

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

- Goal: minimize $\ell(w)$

Linear SVM

Input



Recall ReLU, Maxout Network



- Training datasets
 $\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$
(target $t^{(i)}$: 0 or 1)

- Linear model
 $y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$

- Parameters
 b, w_1, w_2, \dots, w_M

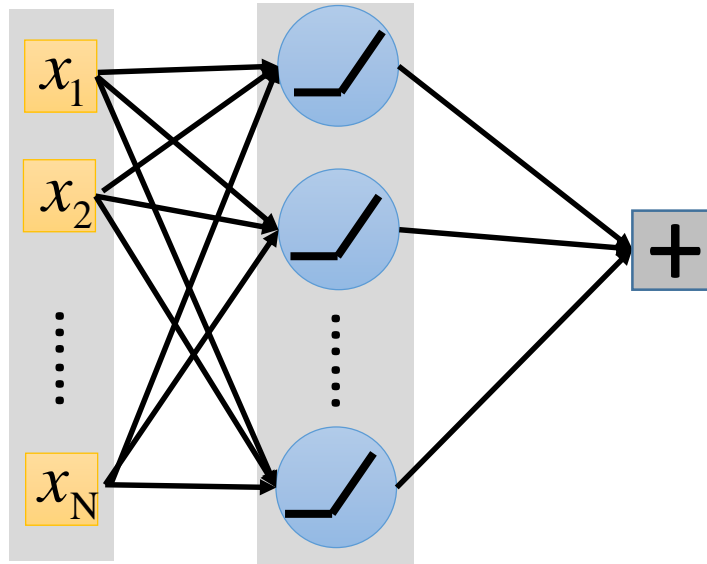
- Loss function

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

- Goal: minimize $\ell(w)$

Linear SVM

Input



Recall ReLU, Maxout Network



- Training datasets
 $\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$
(target $t^{(i)}$: 0 or 1)

- Linear model
 $y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$

- Parameters
 b, w_1, w_2, \dots, w_M

- Loss function

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

- Goal: minimize $\ell(w)$

Tang, Yichuan. "Deep learning using linear support vector machines." *arXiv preprint arXiv:1306.0239* (2013).

Logistic Regression

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: 0 or 1)

- Linear model

$$y(x) = \sigma(w_0 + w_1x_1 + w_2x_2 + \dots + w_Mx_M)$$

- Parameters

$$w_0, w_1, w_2, \dots, w_M$$

- Loss function

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

- Goal: minimize $\ell(w)$

Steps:

- Initialize w (e.g., randomly)
- Repeatedly update w based on the gradient

$$w = w - \epsilon \nabla_w \ell(w)$$

where ϵ is the learning rate.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Linear SVM

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: 0 or 1)

- Linear model

$$y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$$

- Parameters

$$b, w_1, w_2, \dots, w_M$$

- Loss function

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b) t^{(i)})$$

- Goal: minimize $\ell(w)$

Quadratic Programming

Logistic Regression

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: 0 or 1)

- Linear model

$$y(x) = \sigma(w_0 + w_1x_1 + w_2x_2 + \dots + w_Mx_M)$$

- Parameters

$$w_0, w_1, w_2, \dots, w_M$$

- Loss function

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

- Goal: minimize $\ell(w)$

Steps:

- Initialize w (e.g., randomly)
- Repeatedly update w based on the gradient

$$w = w - \epsilon \nabla_w \ell(w)$$

where ϵ is the learning rate.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Linear SVM

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: 0 or 1)

- Linear model

$$y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$$

- Parameters

$$b, w_1, w_2, \dots, w_M$$

- Loss function

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b) t^{(i)})$$

- Goal: minimize $\ell(w)$

But, we can also use gradient descent

- Initialize w, b (e.g., randomly)
- Repeatedly update w based on the gradient

$$w = w - \epsilon \nabla_w \ell(w, b)$$

$$b = b - \epsilon \nabla_b \ell(w, b)$$

where ϵ is the learning rate.

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + c \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

→ $\frac{\partial \ell(w, b)}{\partial w} = w - C \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$

where $\alpha^{(i)} = \begin{cases} -1 & \text{if } (w^T x^{(i)} + b)t^{(i)} < 1 \\ 0 & \text{otherwise} \end{cases}$

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

→ $\frac{\partial \ell(w, b)}{\partial w} = w - C \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$

where $\alpha^{(i)} = \begin{cases} -1 & \text{if } (w^T x^{(i)} + b)t^{(i)} < 1 \\ 0 & \text{otherwise} \end{cases}$

→ Let $\frac{\partial \ell(w, b)}{\partial w} = 0$

$$w^* = C \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

Only a small subset of $\alpha^{(i)}$'s will be nonzero, and the corresponding $x^{(i)}$'s are the **Support Vectors**.

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

→ $\frac{\partial \ell(w, b)}{\partial w} = w - C \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$

where $\alpha^{(i)} = \begin{cases} -1 & \text{if } (w^T x^{(i)} + b)t^{(i)} < 1 \\ 0 & \text{otherwise} \end{cases}$

→ Let $\frac{\partial \ell(w, b)}{\partial w} = 0$

$$w^* = C \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

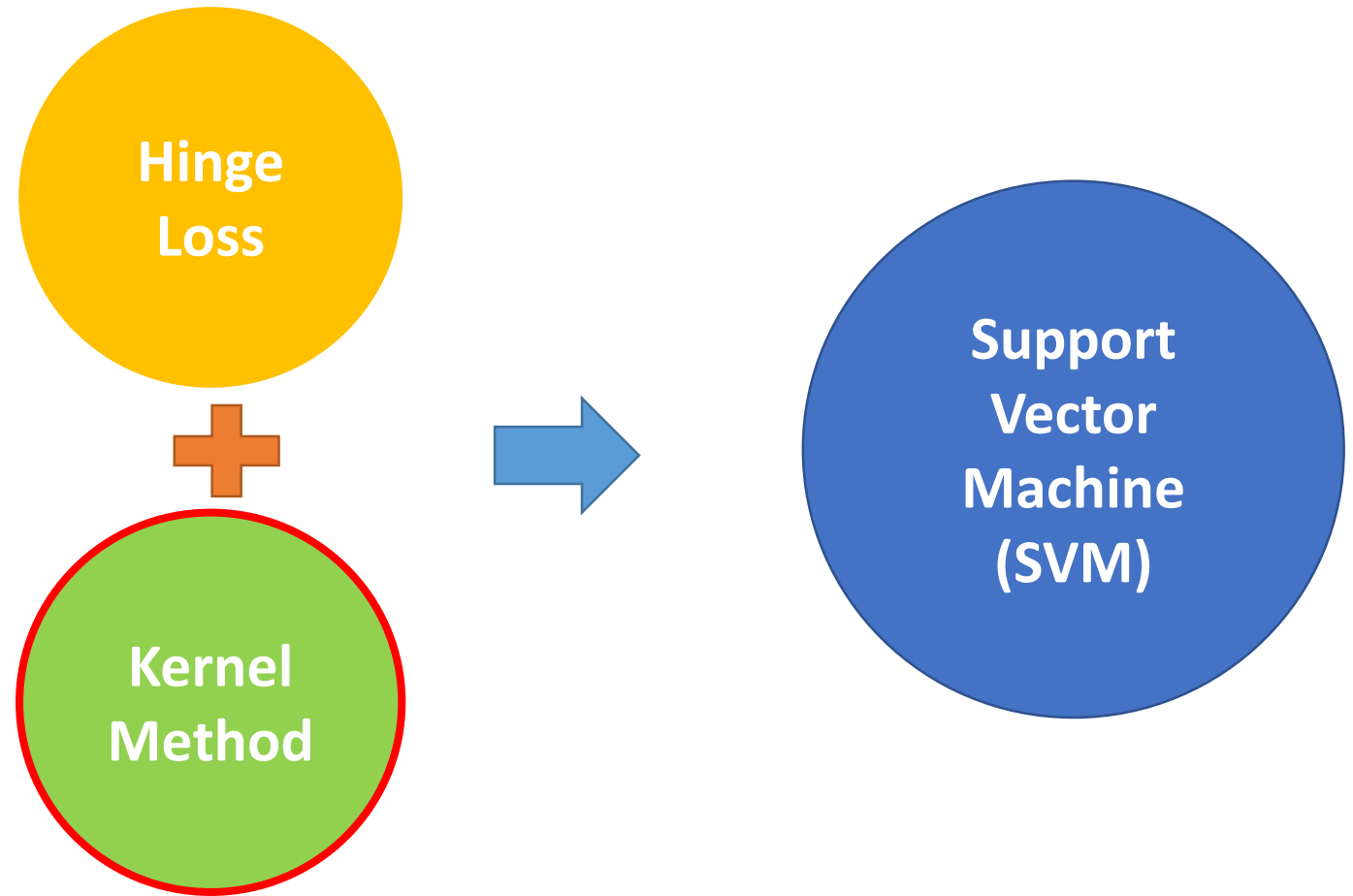
Only a small subset of $\alpha^{(i)}$'s will be nonzero, and the corresponding $x^{(i)}$'s are the **Support Vectors**.

Prediction on a new example:

$$\begin{aligned} y^{(new)} &= \text{sgn}(w^T x^{(new)} + b) \\ &= \text{sgn}\left(\sum_{i=1}^N \alpha^{(i)} t^{(i)} \left(x^{(i)T} x^{(new)}\right) + b\right) \end{aligned}$$

SVM

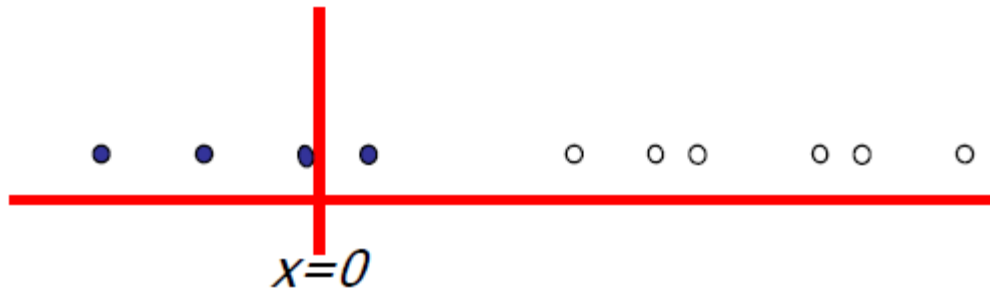
- Linear SVM
- Soft Margin
 - Hinge Loss
- Kernel Trick



SVM

- Suppose we're in 1-dimension

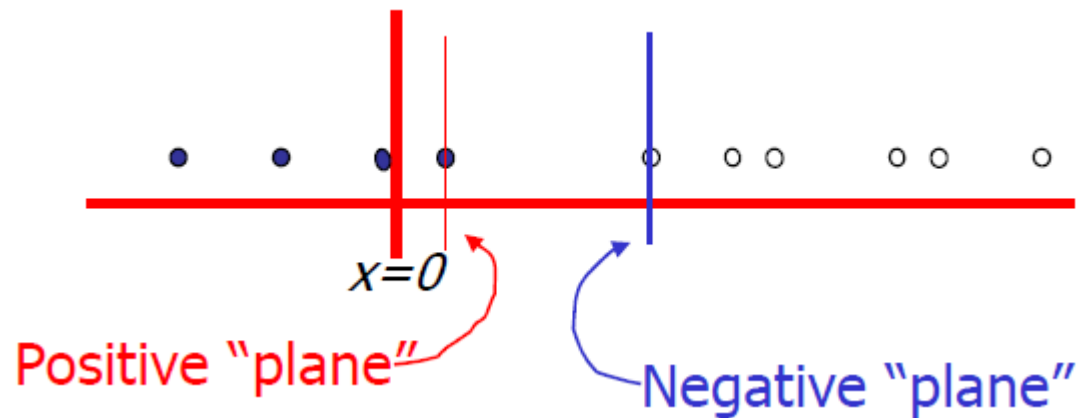
What would SVMs do with this data?



SVM

- Suppose we're in 1-dimension

Not a big surprise

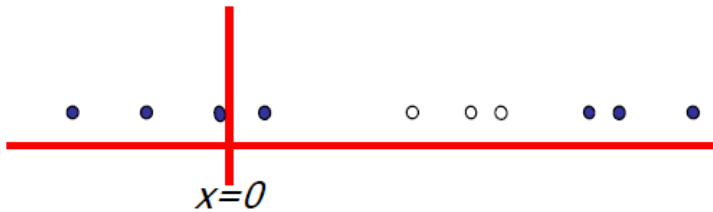


SVM

- Suppose we're in 1-dimension

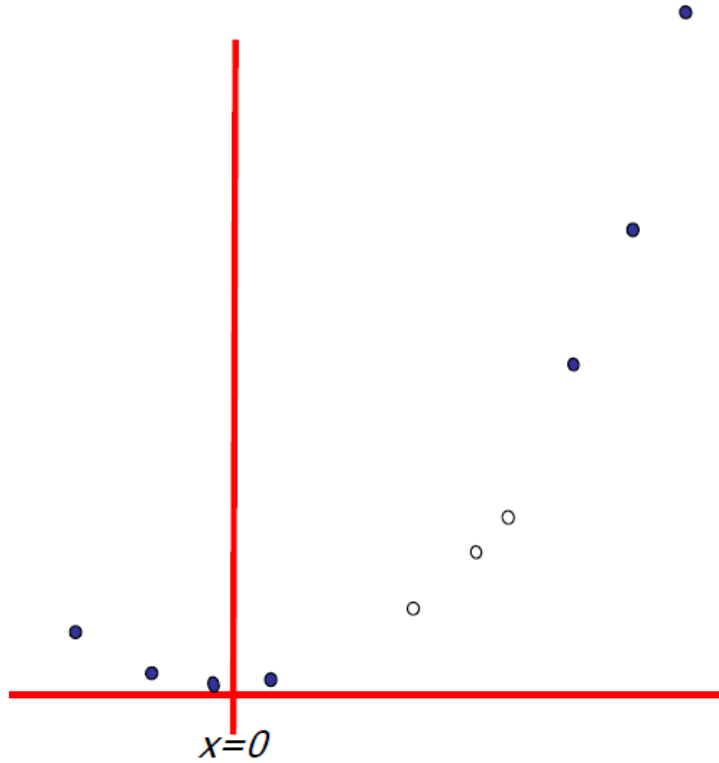
Apply the following map?

$$x_k \Rightarrow (x_k, x_k^2)$$



SVM

- Suppose we're in 1-dimension

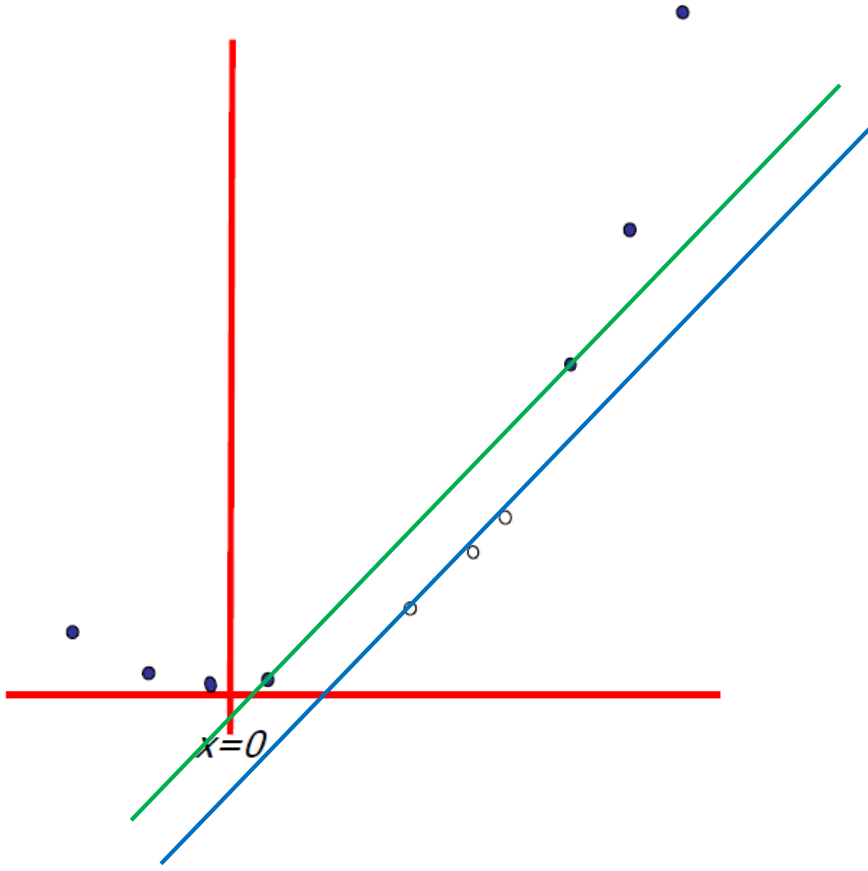


Apply the following map?

$$x_k \Rightarrow (x_k, x_k^2)$$

SVM

- Suppose we're in 1-dimension

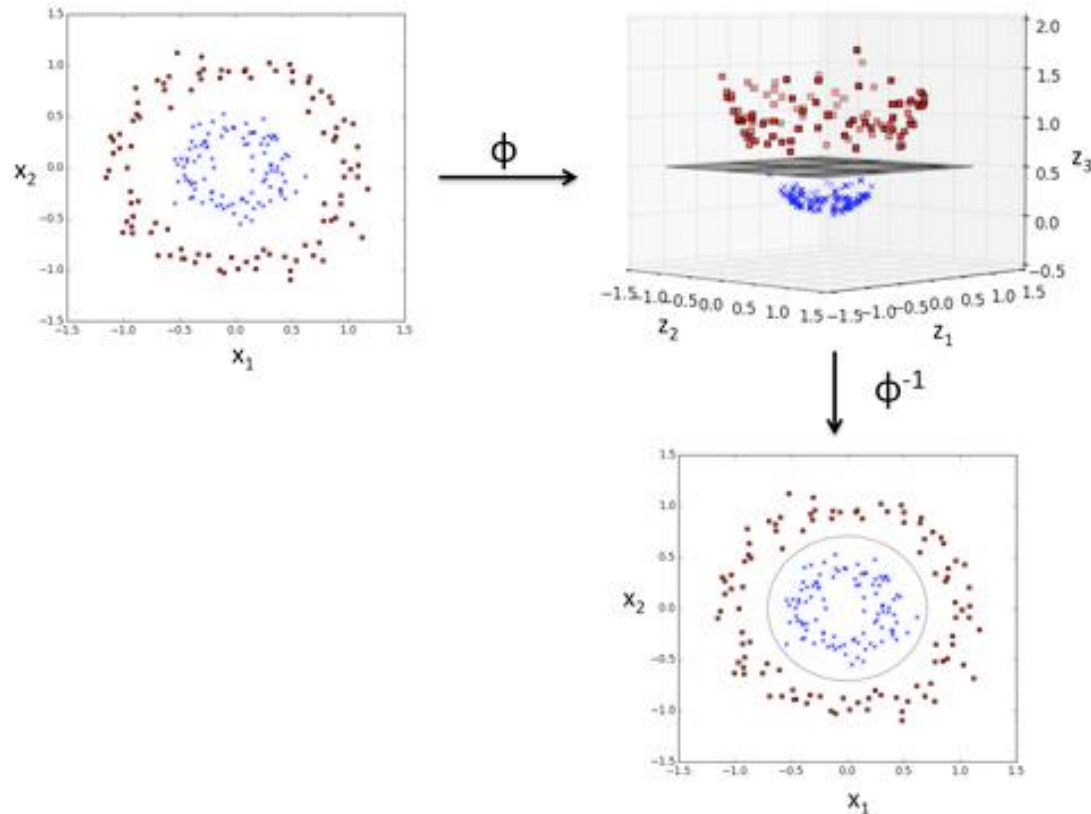


Apply the following map?

$$x_k \Rightarrow (x_k, x_k^2)$$

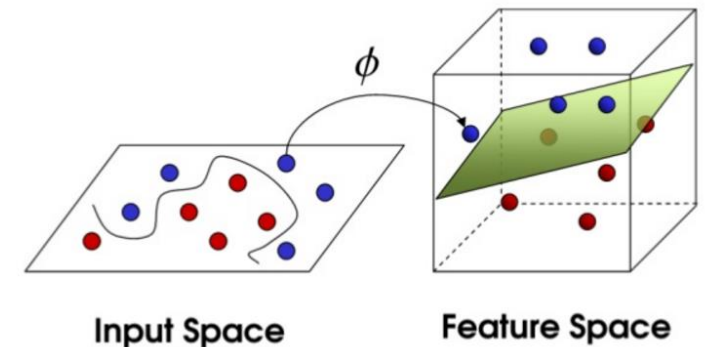
SVM

- Suppose we're in 2-dimension



Feature Transformation

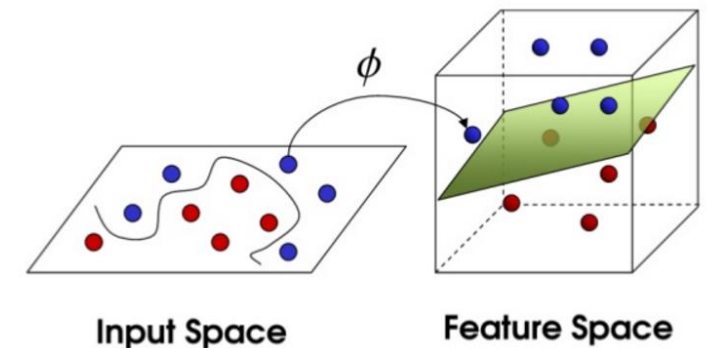
- To form non-linear decision boundaries in input space
 - Map data into feature space $x \rightarrow \phi(x)$
 - Replace dot products between inputs with feature points
$$x^{(i)T} x^{(j)} \rightarrow \phi(x^{(i)})^T \phi(x^{(j)})$$
 - Find a linear decision boundary in feature space



Feature Transformation

- To form non-linear decision boundaries in input space
 - Map data into feature space $x \rightarrow \phi(x)$
 - Replace dot products between inputs with feature points
$$x^{(i)T} x^{(j)} \rightarrow \phi(x^{(i)})^T \phi(x^{(j)})$$
 - Find a linear decision boundary in feature space

It is not easy to find a good transformation



SVM



Note that both the learning objective and the decision function depend only on dot products between datapoints

- Training a SVM model is to maximize

$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} (\mathbf{x}^{(i)T} \mathbf{x}^{(j)})$$

$$\text{s.t. } \alpha^{(i)} \geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

- Once $\alpha^{(i)}$ is obtained

- The weights are

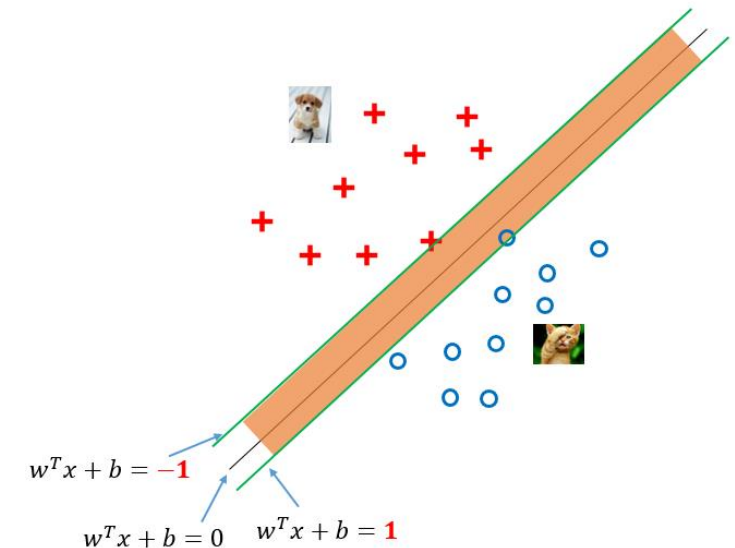
$$\mathbf{w} = \sum_{i=1}^N \alpha^{(i)} t^{(i)} \mathbf{x}^{(i)}$$

- Prediction on a new example:

$$\begin{aligned} y^{(new)} &= \text{sgn}(\mathbf{w}^T \mathbf{x}^{(new)} + b) \\ &= \text{sgn}\left(\sum_{i=1}^N \alpha^{(i)} t^{(i)} (\mathbf{x}^{(i)T} \mathbf{x}^{(new)}) + b\right) \end{aligned}$$

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \forall i, (\mathbf{w}^T \mathbf{x}^{(i)} + b) t^{(i)} \geq 1 \end{aligned}$$

A constrained minimization



SVM



Note that both the learning objective and the decision function depend only on dot products between datapoints

- Training a SVM model is to maximize

$$\max_{\alpha^{(i)}} L(\alpha) = \max_{\alpha^{(i)}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N t^{(i)} t^{(j)} \alpha^{(i)} \alpha^{(j)} \mathbf{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

$$\text{s.t. } \alpha^{(i)} \geq 0 \text{ and } \sum_{i=1}^N \alpha^{(i)} t^{(i)} = 0$$

$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1$$

A constrained minimization

- Once $\alpha^{(i)}$ is obtained

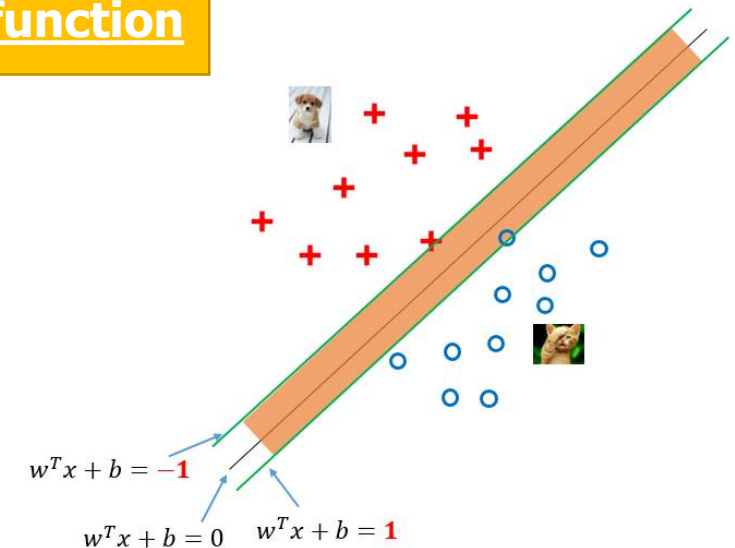
- The weights are

$$w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

- Prediction on a new example:

$$\begin{aligned} y^{(new)} &= \text{sgn}(w^T x^{(new)} + b) \\ &= \text{sgn}\left(\sum_{i=1}^N \alpha^{(i)} t^{(i)} \mathbf{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(new)}) + b\right) \end{aligned}$$

Replacing dot product with a kernel function



SVM

- Training datasets

$$\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$$

(target $t^{(i)}$: 0 or 1)

- Linear model

$$y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$$

- Parameters

$$b, w_1, w_2, \dots, w_M$$



Note that both the learning objective and the decision function depend only on dot products between datapoints

Replacing dot product with a kernel function

SVM

- Training datasets
 $\mathcal{D} = \{(x^{(1)}, t^{(1)}), \dots, (x^{(i)}, t^{(i)}), \dots, (x^{(N)}, t^{(N)})\}$
(target $t^{(i)}$: 0 or 1)
- Linear model
 $y(x) = b + w_1x_1 + w_2x_2 + \dots + w_Mx_M$
- Parameters
 b, w_1, w_2, \dots, w_M

$$w = \sum_{i=1}^N \alpha^{(i)} t^{(i)} x^{(i)}$$

$$\begin{aligned} y(x) &= w^T x + b \\ &= \sum_{i=1}^N \alpha^{(i)} t^{(i)} K(x^{(i)}, x) + b \end{aligned}$$

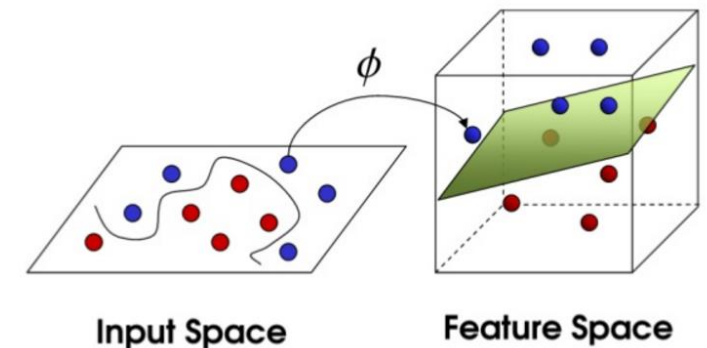


Note that both the learning objective and the decision function depend only on dot products between datapoints

Replacing dot product with a kernel function

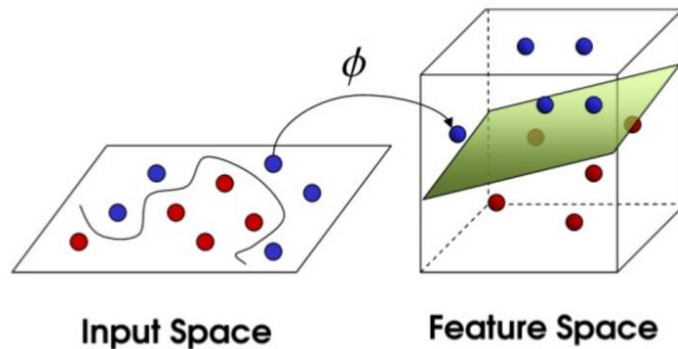
Kernel Trick

- To form non-linear decision boundaries in input space
 - Map data into feature space $x \rightarrow \phi(x)$
 - Replace dot products between inputs with feature points
$$x^{(i)T} x^{(j)} \rightarrow \phi(x^{(i)})^T \phi(x^{(j)})$$
 - Find a linear decision boundary in feature space
- In SVM, we use Kernel Tricks
 - (Pro) Introduce nonlinearity into the model
 - (Pro) Computational cheap
 - (Con) Potential overfitting problem



Kernel Trick

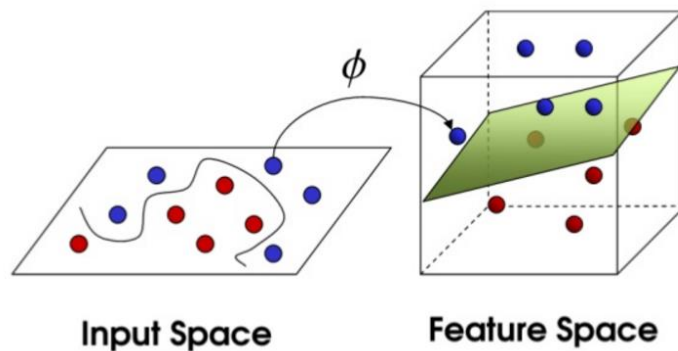
- Directly computing $K(x^{(i)}, x^{(j)})$ can be faster than “feature transformation + inner product” sometimes.



Kernel Trick

- Directly computing $K(x^{(i)}, x^{(j)})$ can be faster than “feature transformation + inner product” sometimes.

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow \phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}$$

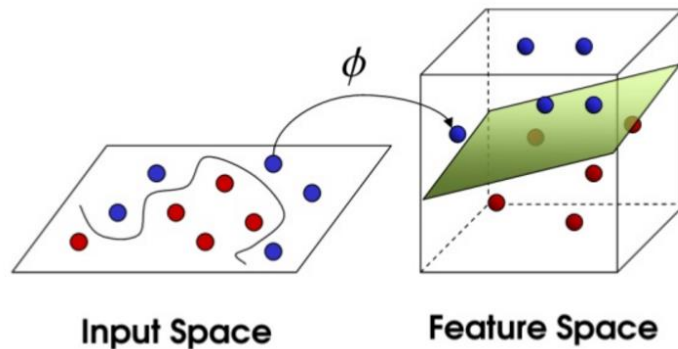


Kernel Trick

- Directly computing $K(x^{(i)}, x^{(j)})$ can be faster than “feature transformation + inner product” sometimes.

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow \phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}$$

$$K(x, z) = (x \cdot z)^2$$

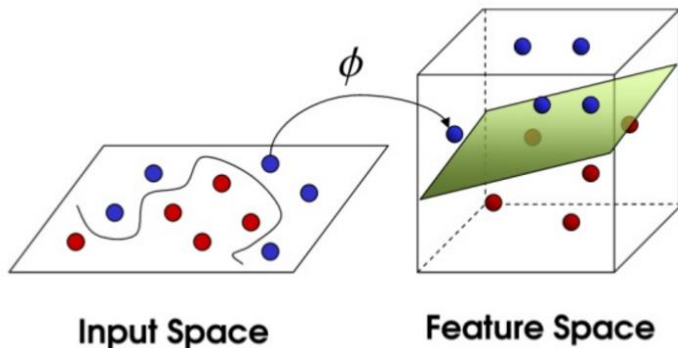


Kernel Trick

- Directly computing $K(x^{(i)}, x^{(j)})$ can be faster than “feature transformation + inner product” sometimes.

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow \phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}$$

$$K(x, z) = (x \cdot z)^2$$



$$\begin{aligned} K(x, z) &= \phi(x) \cdot \phi(z) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix} \\ &= x_1^2z_1^2 + 2x_1x_2z_1z_2 + x_2^2z_2^2 \\ &= (x_1z_1 + x_2z_2)^2 = \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \right)^2 \\ &= (x \cdot z)^2 \end{aligned}$$

Kernel Trick

- Directly computing $K(x^{(i)}, x^{(j)})$ can be faster than “feature transformation + inner product” sometimes.

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix} \quad Z = \begin{bmatrix} Z_1 \\ \vdots \\ Z_k \end{bmatrix}$$

$$\phi(x) = \begin{bmatrix} x_1^2 \\ \vdots \\ x_k^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \vdots \\ \sqrt{2}x_2x_3 \\ \vdots \end{bmatrix}$$

Kernel Trick

- Directly computing $K(x^{(i)}, x^{(j)})$ can be faster than “feature transformation + inner product” sometimes.

$$\begin{aligned} K(x, z) &= (x \cdot z)^2 \\ &= (x_1 z_1 + x_2 z_2 + \cdots + x_k z_k)^2 \\ &= \underline{x_1^2} \underline{z_1^2} + \underline{x_2^2} \underline{z_2^2} + \cdots + \underline{x_k^2} \underline{z_k^2} \\ &\quad + 2\underline{x_1 x_2} \underline{z_1 z_2} + 2\underline{x_1 x_3} \underline{z_1 z_3} + \cdots \\ &\quad + 2\underline{x_2 x_3} \underline{z_2 z_3} + 2\underline{x_2 x_4} \underline{z_2 z_4} + \cdots \\ &= \phi(x) \cdot \phi(z) \end{aligned}$$

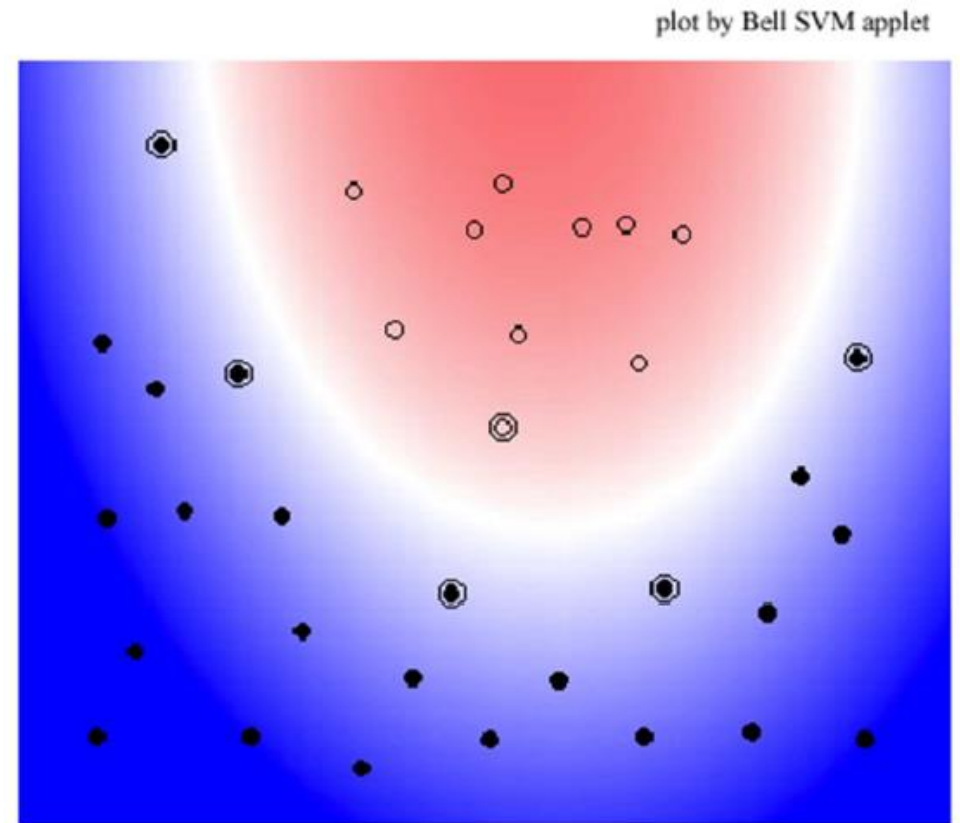
$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix} \quad z = \begin{bmatrix} z_1 \\ \vdots \\ z_k \end{bmatrix}$$

$$\phi(x) = \begin{bmatrix} x_1^2 \\ \vdots \\ x_k^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \vdots \\ \sqrt{2}x_2x_3 \\ \vdots \end{bmatrix}$$

Kernel Trick

- SVM with polynomial of degree 2

Kernel: $K(x^{(i)}, x^{(j)}) = (x^{(i)T} x^{(j)} + 1)^2$



Kernel Trick

- Polynomial Kernel

$$K(x, z) = (x^T z + 1)^d$$

- Gaussian Kernel

$$K(x, z) = e^{\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)}$$

- Sigmoid Kernel

$$K(x, z) = \tanh(\beta x^T z + a)$$

- ...

Kernel Trick

- Radial Basis Function Kernel

$$K(x, z) = \exp\left(-\frac{1}{2}\|x - z\|_2\right)$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \quad z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \end{bmatrix}$$

Kernel Trick

- Radial Basis Function Kernel

$$K(x, z) = \exp\left(-\frac{1}{2}\|x - z\|_2\right) = \phi(x) \cdot \phi(z)?$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \quad z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \end{bmatrix}$$

Kernel Trick

- Radial Basis Function Kernel

$$K(x, z) = \exp\left(-\frac{1}{2}\|x - z\|_2\right) = \phi(x) \cdot \phi(z)?$$
$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \quad z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \end{bmatrix}$$
$$= \exp\left(-\frac{1}{2}\|x\|_2 - \frac{1}{2}\|z\|_2 + x \cdot z\right)$$

Kernel Trick

- Radial Basis Function Kernel

$$\begin{aligned} K(x, z) &= \exp\left(-\frac{1}{2}\|x - z\|_2\right) = \phi(x) \cdot \phi(z)? & x &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} & z &= \begin{bmatrix} z_1 \\ z_2 \\ \vdots \end{bmatrix} \\ &= \exp\left(-\frac{1}{2}\|x\|_2 - \frac{1}{2}\|z\|_2 + x \cdot z\right) \\ &= \exp\left(-\frac{1}{2}\|x\|_2\right) \exp\left(-\frac{1}{2}\|z\|_2\right) \exp(x \cdot z) = C_x C_z \exp(x \cdot z) \end{aligned}$$

Kernel Trick

- Radial Basis Function Kernel

$$\begin{aligned} K(x, z) &= \exp\left(-\frac{1}{2}\|x - z\|_2\right) = \phi(x) \cdot \phi(z)? & x &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} & z &= \begin{bmatrix} z_1 \\ z_2 \\ \vdots \end{bmatrix} \\ &= \exp\left(-\frac{1}{2}\|x\|_2 - \frac{1}{2}\|z\|_2 + x \cdot z\right) \\ &= \exp\left(-\frac{1}{2}\|x\|_2\right) \exp\left(-\frac{1}{2}\|z\|_2\right) \exp(x \cdot z) = C_x C_z \exp(x \cdot z) \\ &= C_x C_z \sum_{i=0}^{\infty} \frac{(x \cdot z)^i}{i!} = C_x C_z + C_x C_z (x \cdot z) + C_x C_z \frac{1}{2} (x \cdot z)^2 \dots \end{aligned}$$

Kernel Trick

- Radial Basis Function Kernel

$$\begin{aligned}
 K(x, z) &= \exp\left(-\frac{1}{2}\|x - z\|_2\right) = \phi(x) \cdot \phi(z)? & x &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} & z &= \begin{bmatrix} z_1 \\ z_2 \\ \vdots \end{bmatrix} \\
 &= \exp\left(-\frac{1}{2}\|x\|_2 - \frac{1}{2}\|z\|_2 + x \cdot z\right) \\
 &= \exp\left(-\frac{1}{2}\|x\|_2\right) \exp\left(-\frac{1}{2}\|z\|_2\right) \exp(x \cdot z) = C_x C_z \exp(x \cdot z) \\
 &= C_x C_z \sum_{i=0}^{\infty} \frac{(x \cdot z)^i}{i!} = C_x C_z + C_x C_z (x \cdot z) + C_x C_z \frac{1}{2} (x \cdot z)^2 \dots
 \end{aligned}$$

Diagram illustrating the expansion of the RBF kernel into a dot product of feature vectors:

- $[C_x] \cdot [C_z]$ corresponds to the constant term $C_x C_z$.
- $\begin{bmatrix} C_x x_1 \\ C_x x_2 \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} C_z z_1 \\ C_z z_2 \\ \vdots \end{bmatrix}$ corresponds to the linear term $C_x C_z (x \cdot z)$.
- $\frac{1}{\sqrt{2}} \begin{bmatrix} C_x x_1^2 \\ \vdots \\ \sqrt{2} C_x x_1 x_2 \\ \vdots \end{bmatrix} \cdot \frac{1}{\sqrt{2}} \begin{bmatrix} C_z z_1^2 \\ \vdots \\ \sqrt{2} C_z z_1 z_2 \\ \vdots \end{bmatrix}$ corresponds to the quadratic term $C_x C_z \frac{1}{2} (x \cdot z)^2$.

Kernel Trick

- Radial Basis Function Kernel

$$K(x, z) = \exp\left(-\frac{1}{2}\|x - z\|_2\right) = \phi(x) \cdot \phi(z)? \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \quad z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \end{bmatrix}$$

$$= \exp\left(-\frac{1}{2}\|x\|_2 - \frac{1}{2}\|z\|_2 + x \cdot z\right) \quad \phi(*) \text{ has inf dim!!!}$$

$$= \exp\left(-\frac{1}{2}\|x\|_2\right) \exp\left(-\frac{1}{2}\|z\|_2\right) \exp(x \cdot z) = C_x C_z \exp(x \cdot z)$$

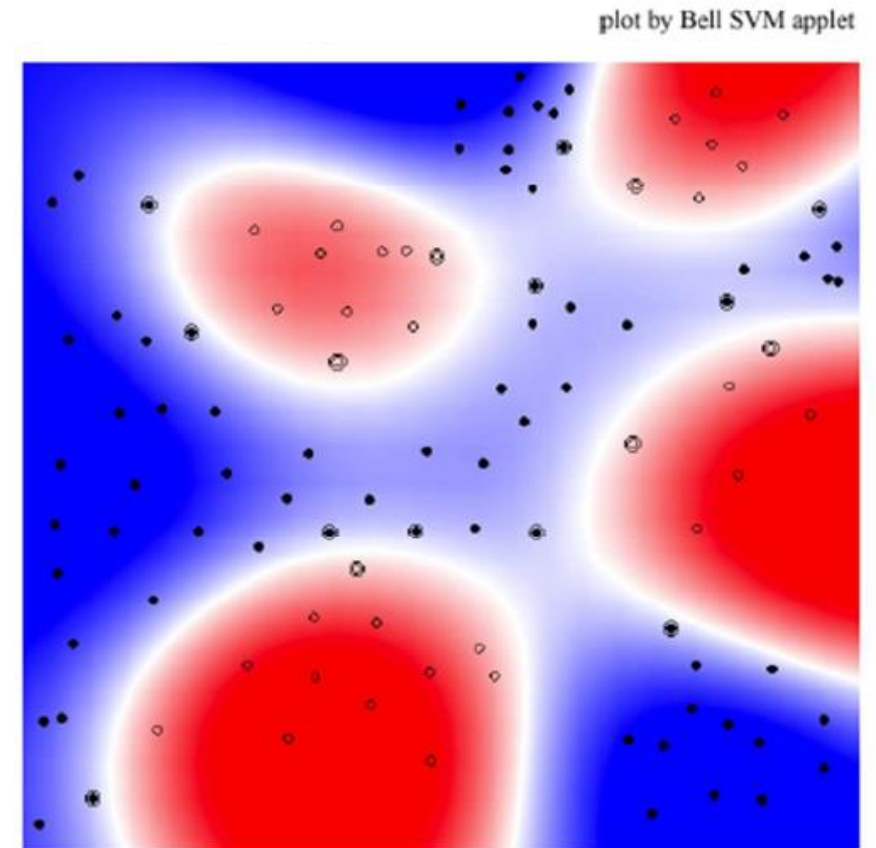
$$= C_x C_z \sum_{i=0}^{\infty} \frac{(x \cdot z)^i}{i!} = C_x C_z + C_x C_z (x \cdot z) + C_x C_z \frac{1}{2} (x \cdot z)^2 \dots$$

$$[C_x] \cdot [C_z] \quad \begin{bmatrix} C_x x_1 \\ C_x x_2 \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} C_z z_1 \\ C_z z_2 \\ \vdots \end{bmatrix} \quad \frac{1}{\sqrt{2}} \begin{bmatrix} C_x x_1^2 \\ \vdots \\ \sqrt{2} C_x x_1 x_2 \\ \vdots \end{bmatrix} \cdot \frac{1}{\sqrt{2}} \begin{bmatrix} C_z z_1^2 \\ \vdots \\ \sqrt{2} C_z z_1 z_2 \\ \vdots \end{bmatrix}$$

Kernel Trick

- SVM with Radial Basis Function Kernel (RBF)

Kernel: $K(x^{(i)}, x^{(j)}) = e^{\left(-\frac{\|x^{(i)} - x^{(j)}\|^2}{\sigma^2}\right)}$



Multi-Class Classification

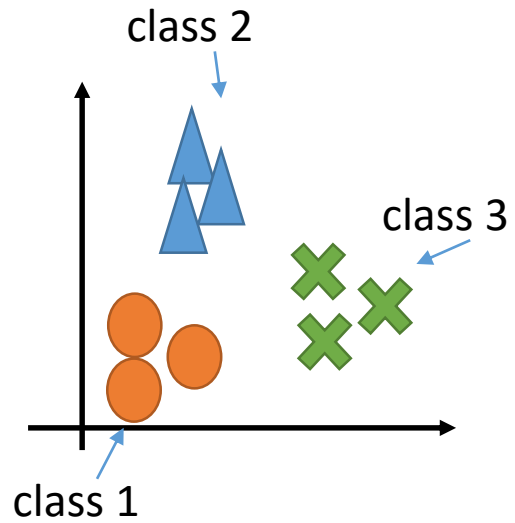
- SVMs are inherently two-class classifiers
- What can be done?

Multi-Class Classification

- SVMs are inherently two-class classifiers
- What can be done?
 - One vs All
 - Other approaches
 - Pair-wise SVM
 - Brunner, Carl, Andreas Fischer, Klaus Luig, and Thorsten Thies. "Pairwise support vector machines and their application to large scale problems." *Journal of Machine Learning Research* 13, no. Aug (2012): 2279-2292.
 - Multi-category SVM
 - Duan, Kai-Bo, and S. Sathiya Keerthi. "Which is the best multiclass SVM method? An empirical study." In *International Workshop on Multiple Classifier Systems*, pp. 278-285. Springer Berlin Heidelberg, 2005.

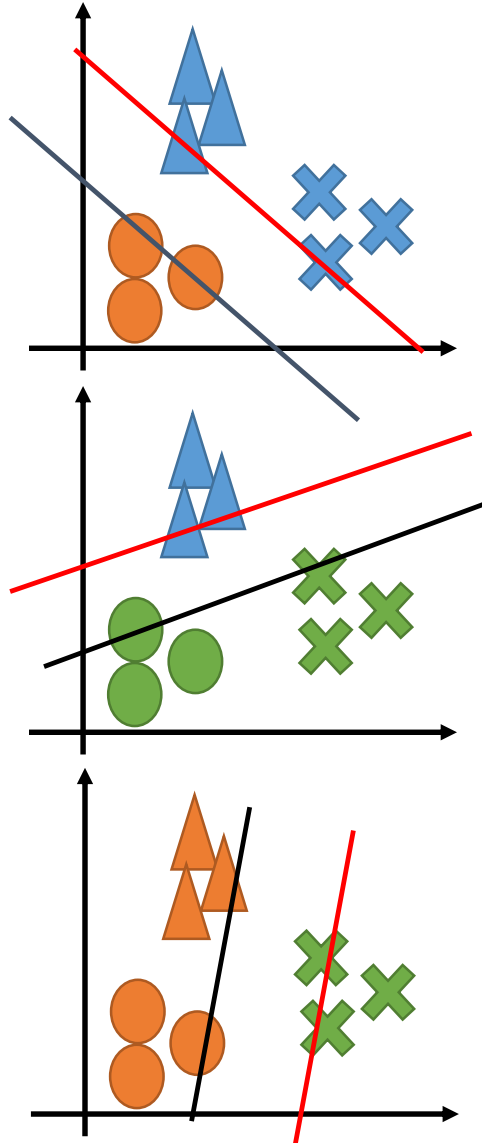
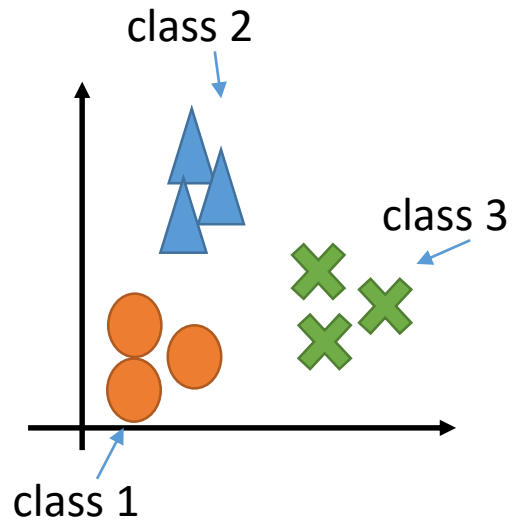
Multi-Class Classification

- One vs All



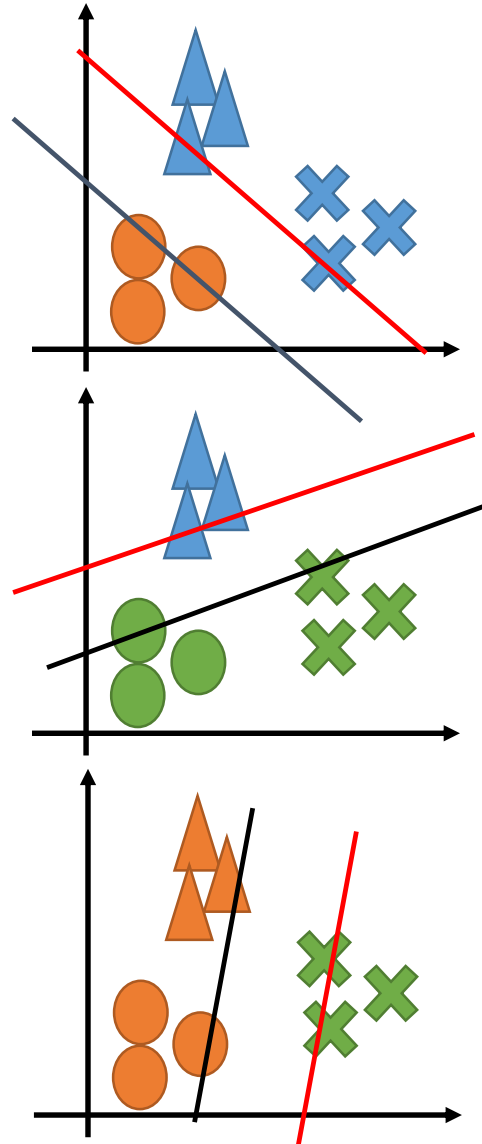
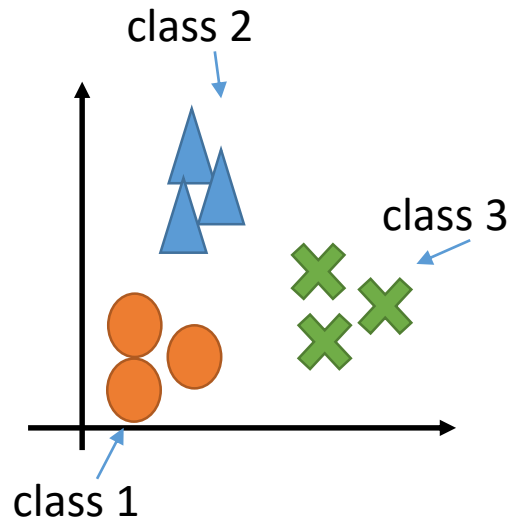
Multi-Class Classification

- One vs All



Multi-Class Classification

- One vs All

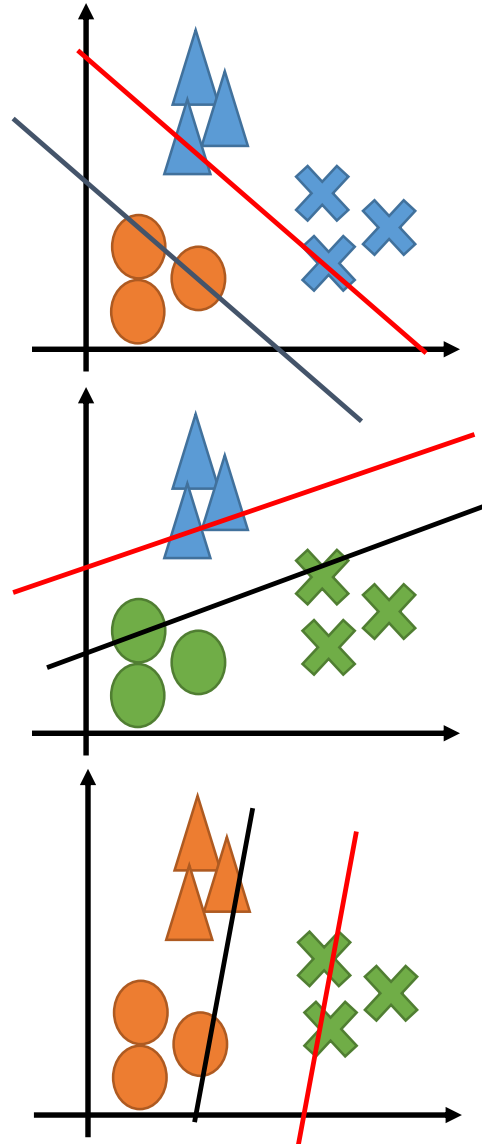
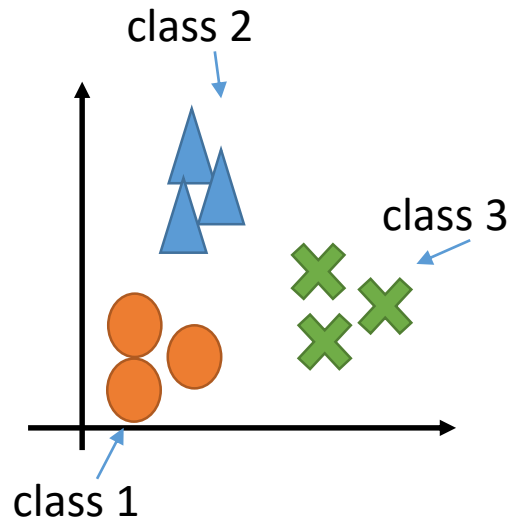


 New data



Multi-Class Classification

- One vs All



 New data



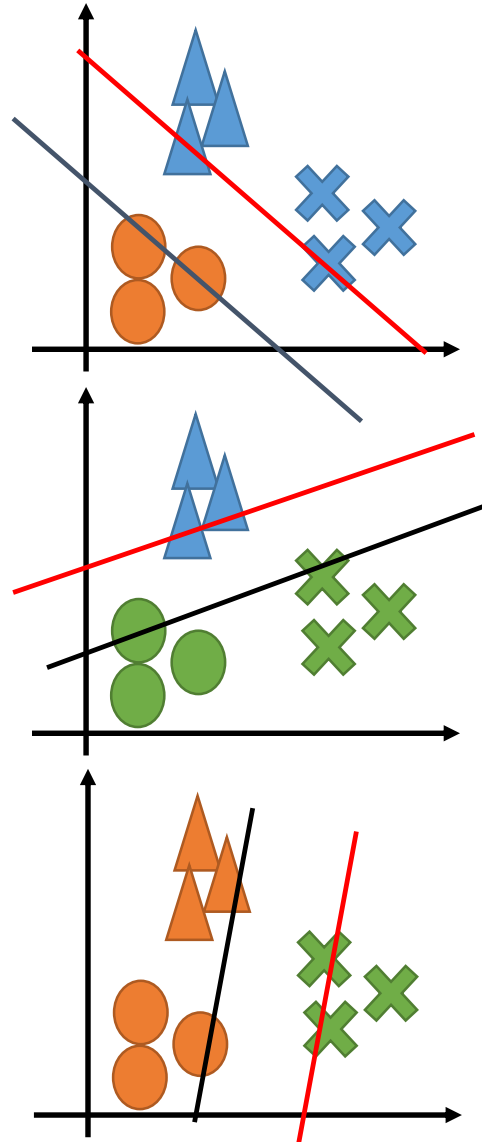
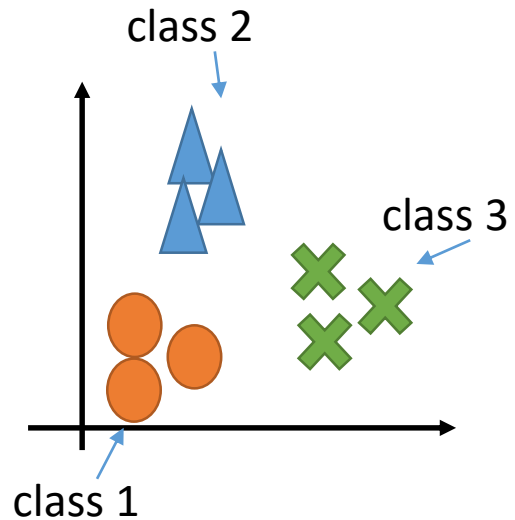
SVM 1: decision value in class 1

SVM 2: decision value in class 2

SVM 3: decision value in class 3

Multi-Class Classification

- One vs All



 New data



SVM 1: decision value in class 1

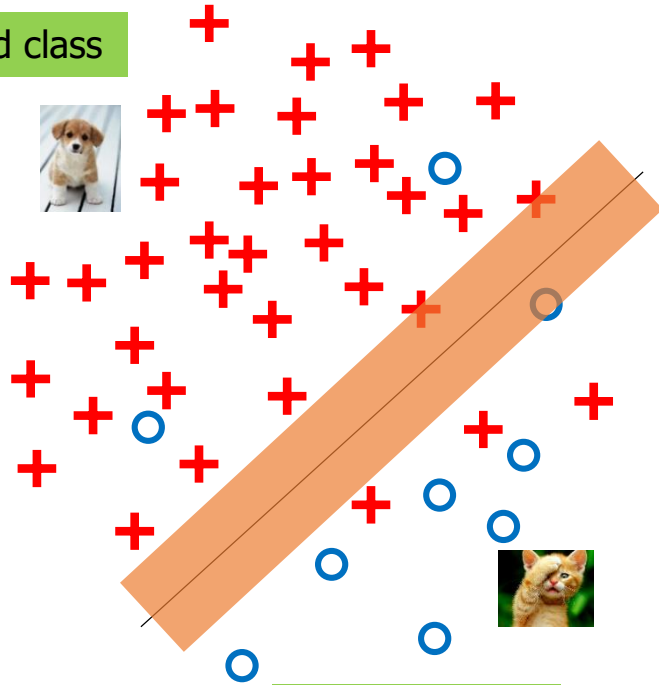
SVM 2: decision value in class 2

SVM 3: decision value in class 3

Choose a
class with
highest
decision
value

SVM for Unbalanced Data

The second class



The first class

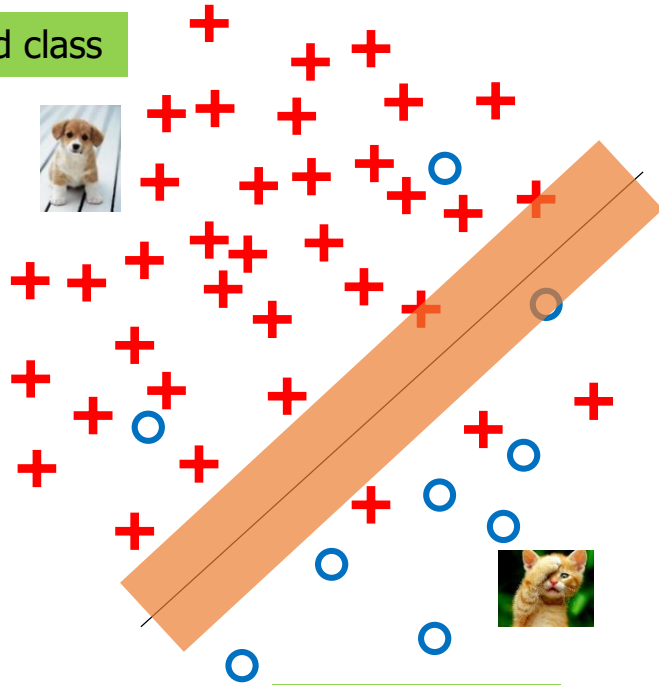
$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$s.t. \quad \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 - \xi^{(i)}$$

$$\forall i, \xi^{(i)} \geq 0$$

SVM for Unbalanced Data

The second class



The first class

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$s. t. \quad \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 - \xi^{(i)}$$

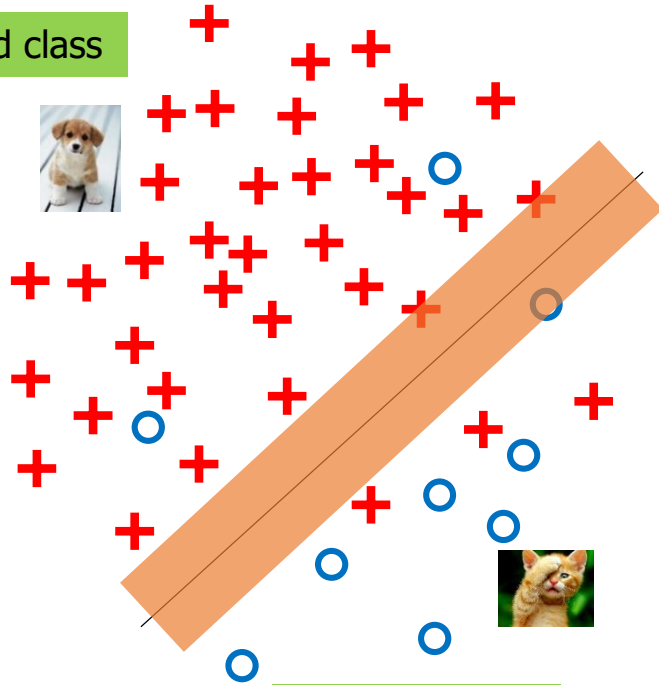
$$\forall i, \xi^{(i)} \geq 0$$

If the first class has much smaller size than the second class,

- apply different weights to the two classes: $C_1 > C_2$

SVM for Unbalanced Data

The second class



The first class

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

$$s.t. \quad \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 - \xi^{(i)}$$

$$\forall i, \xi^{(i)} \geq 0$$

If the first class has much smaller size than the second class,

- apply different weights to the two classes: $C_1 > C_2$

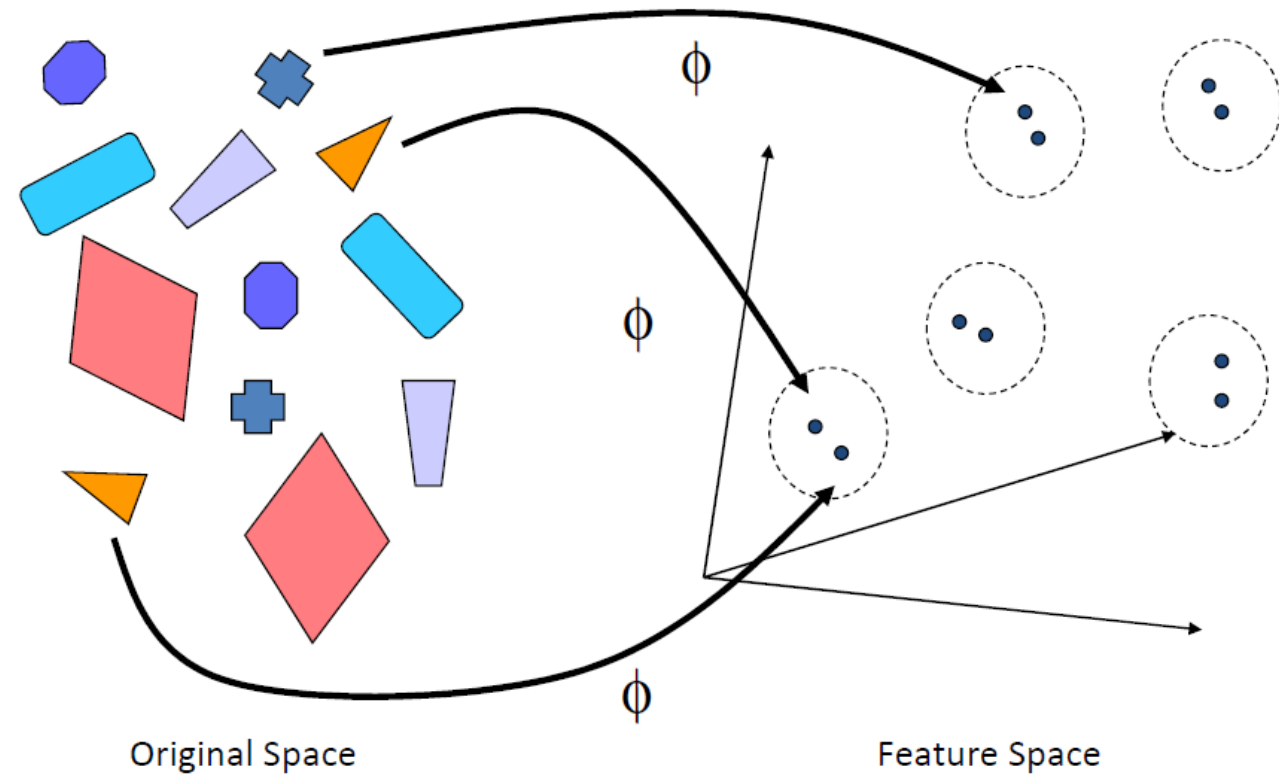
$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C_1 \sum_{x^{(i)} \in \text{Class 1}} \xi^{(i)} + C_2 \sum_{x^{(i)} \in \text{Class 2}} \xi^{(i)}$$

$$s.t. \quad \forall i, (w^T x^{(i)} + b) t^{(i)} \geq 1 - \xi^{(i)}$$

$$\forall i, \xi^{(i)} \geq 0$$

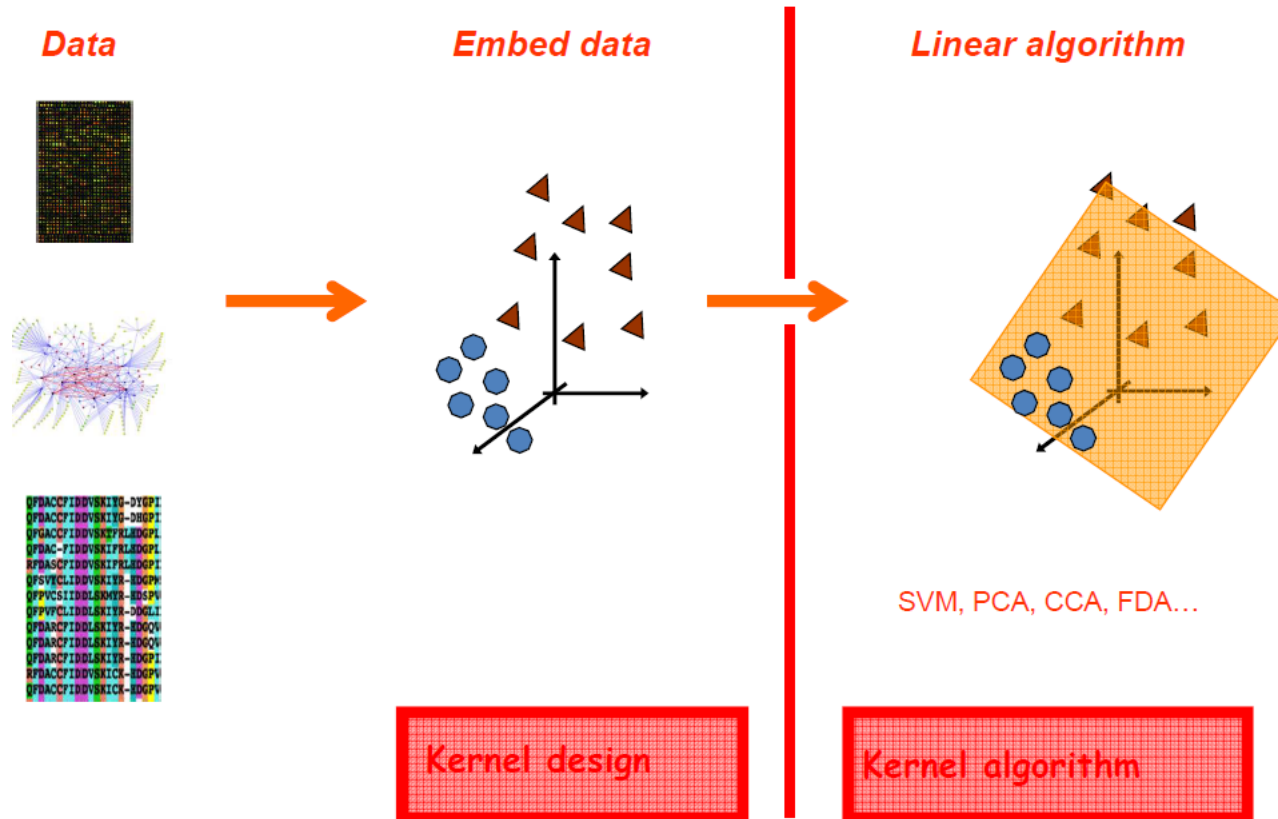
Kernel-based Learning

- Kernels are measures of similarity



Kernel-based Learning

- Kernels are measures of similarity



SVM

- Linear Regression

- Linear function + Square loss

$$\ell(w) = \frac{1}{2N} \sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2$$

- Logistic Regression

- Sigmoid function + Cross entropy loss

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

- SVM

- Linear function + Hinge loss + L2 norm

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$

SVM

We want to have a classifier with $f(x^{(i)})t^{(i)} \geq 1$ for all data points

- Linear Regression

- Linear function + Square loss

$$\ell(w) = \frac{1}{2N} \sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2$$

- Logistic Regression

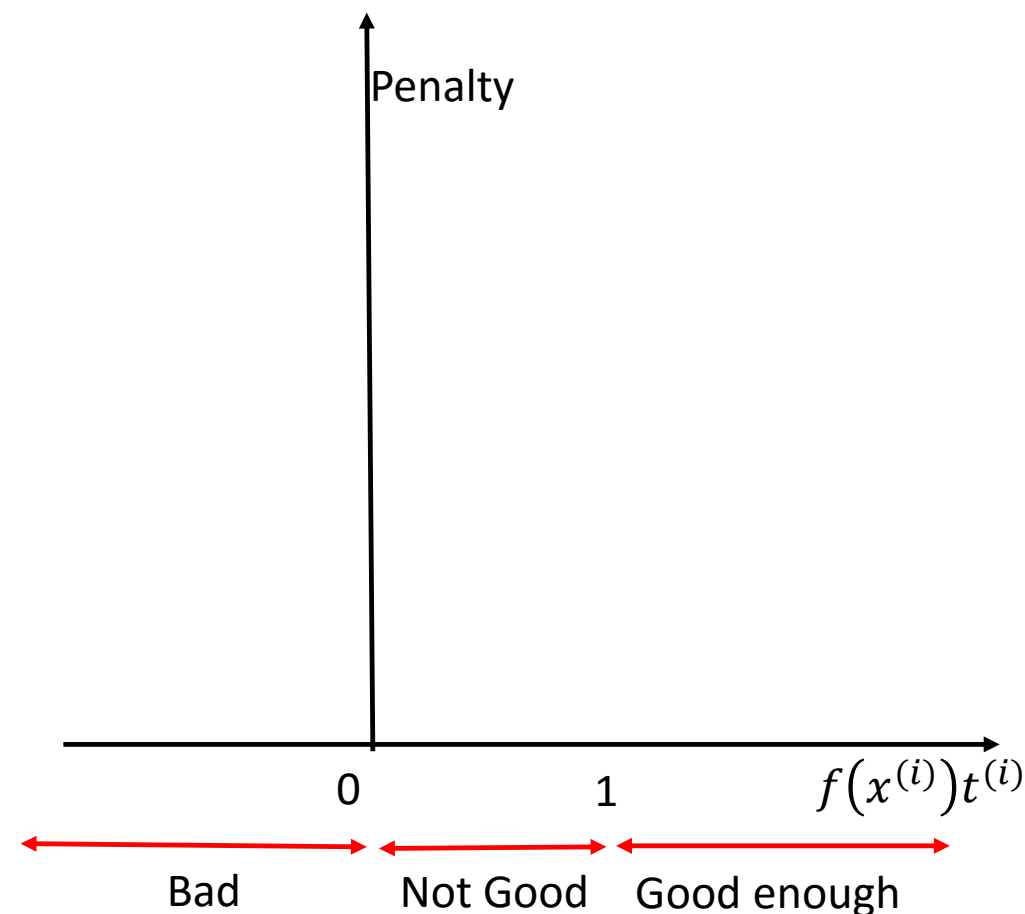
- Sigmoid function + Cross entropy loss

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

- SVM

- Linear function + Hinge loss + L2 norm

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$



SVM

We want to have a classifier with $f(x^{(i)})t^{(i)} \geq 1$ for all data points

- Linear Regression

- Linear function + Square loss

$$\ell(w) = \frac{1}{2N} \sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2$$

- Logistic Regression

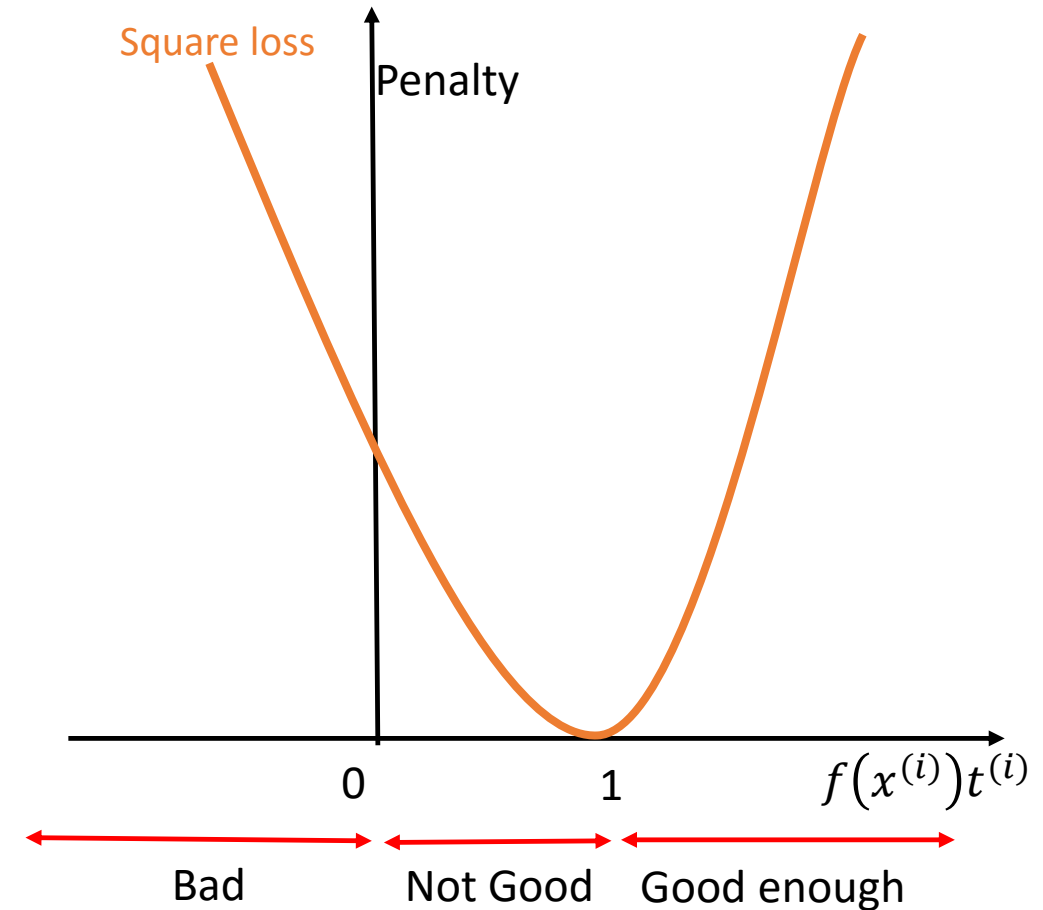
- Sigmoid function + Cross entropy loss

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

- SVM

- Linear function + Hinge loss + L2 norm

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$



SVM

We want to have a classifier with $f(x^{(i)})t^{(i)} \geq 1$ for all data points

- Linear Regression

- Linear function + Square loss

$$\ell(w) = \frac{1}{2N} \sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2$$

- Logistic Regression

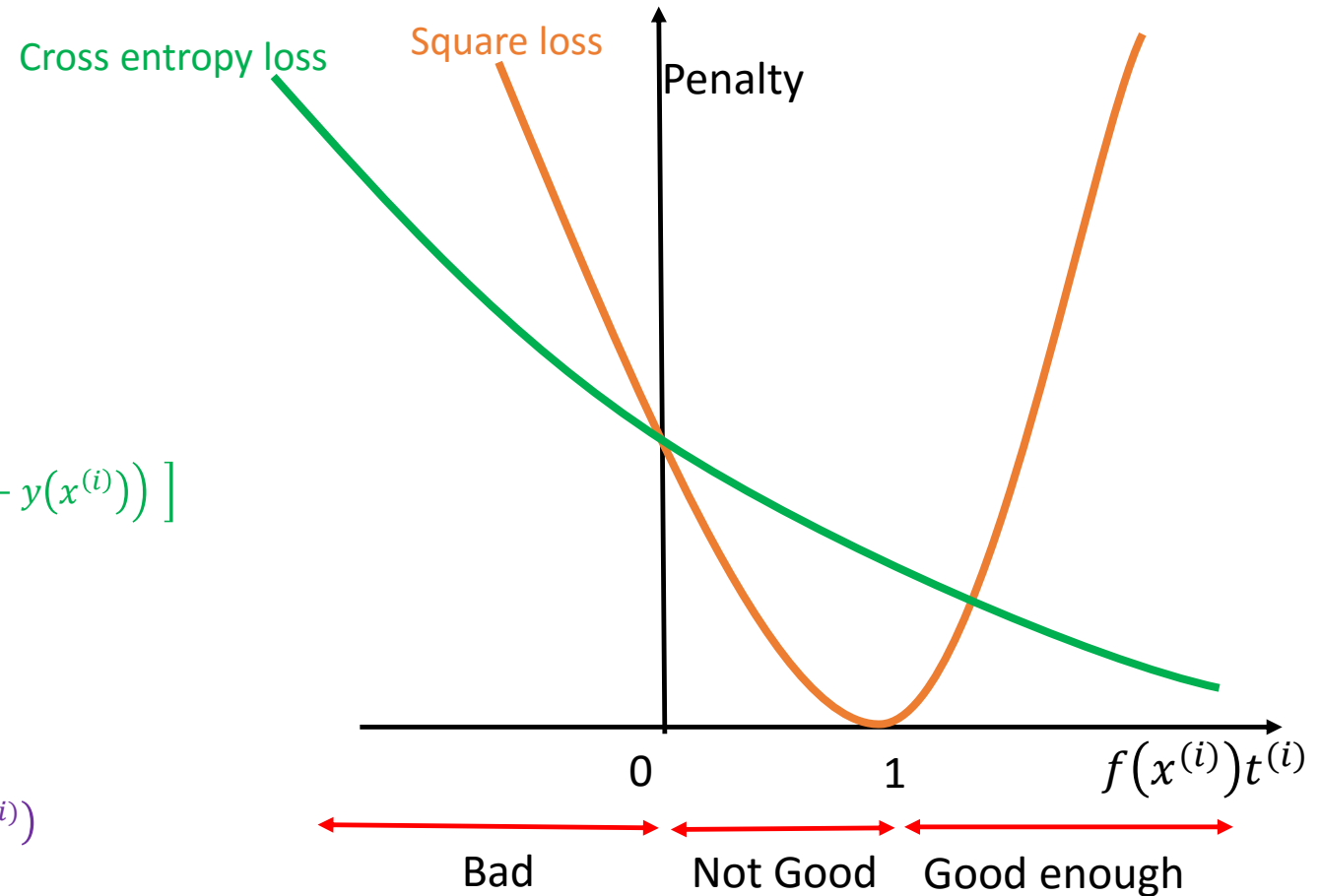
- Sigmoid function + Cross entropy loss

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

- SVM

- Linear function + Hinge loss + L2 norm

$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$



SVM

We want to have a classifier with $f(x^{(i)})t^{(i)} \geq 1$ for all data points

- Linear Regression

- Linear function + Square loss

$$\ell(w) = \frac{1}{2N} \sum_{i=1}^N [t^{(i)} - y(x^{(i)})]^2$$

- Logistic Regression

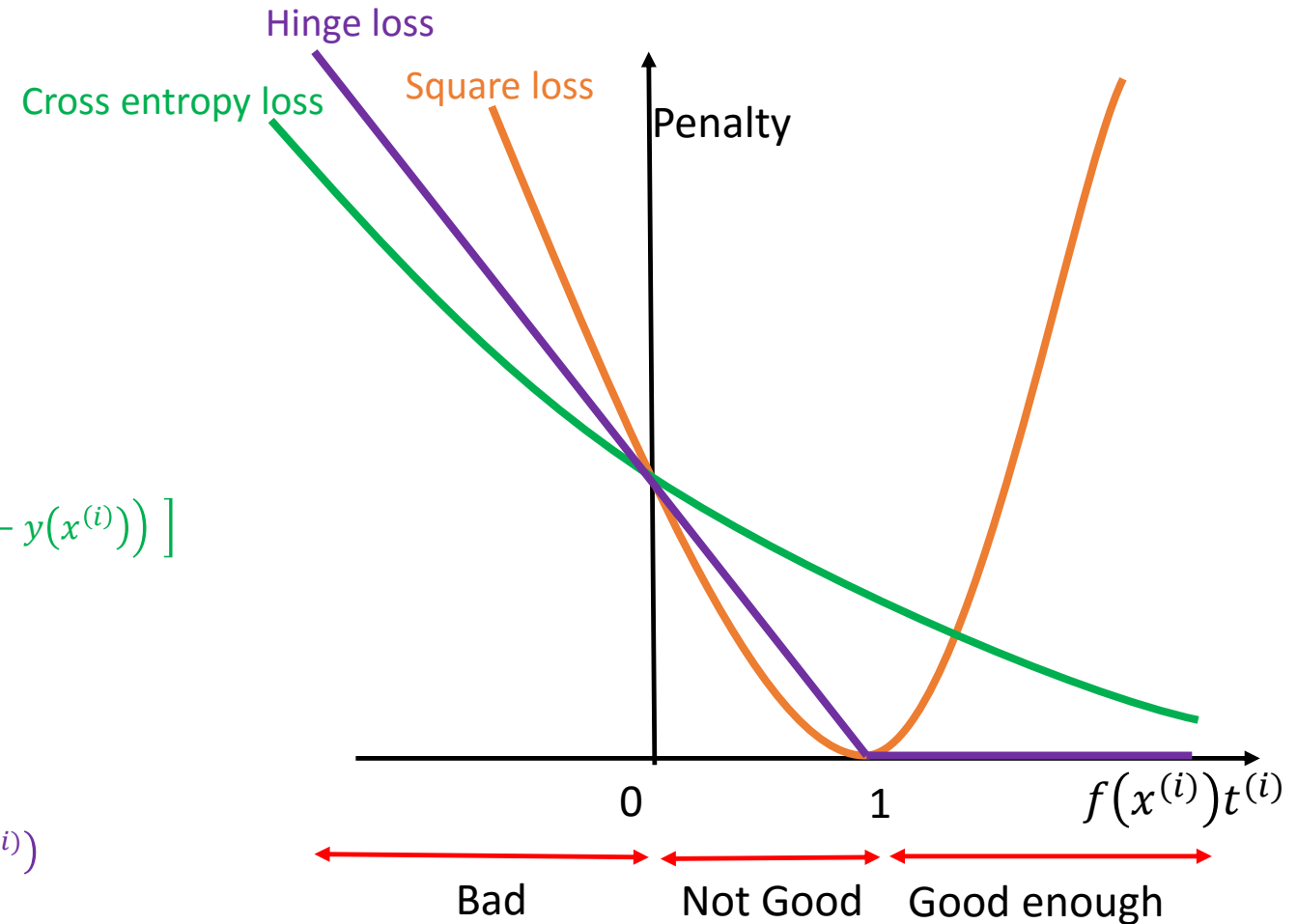
- Sigmoid function + Cross entropy loss

$$\ell(w) = -\frac{1}{N} \sum_{i=1}^N \left[t^{(i)} \log(y(x^{(i)})) + (1 - t^{(i)}) \log(1 - y(x^{(i)})) \right]$$

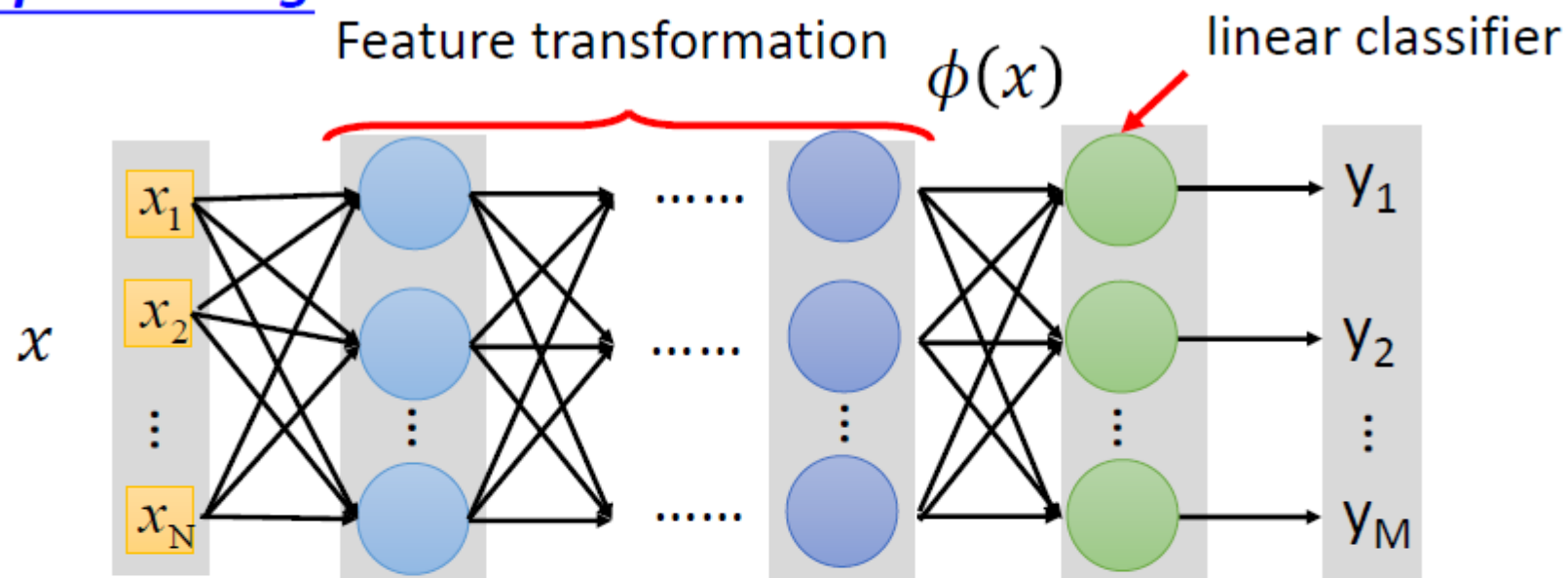
- SVM

- Linear function + Hinge loss + L2 norm

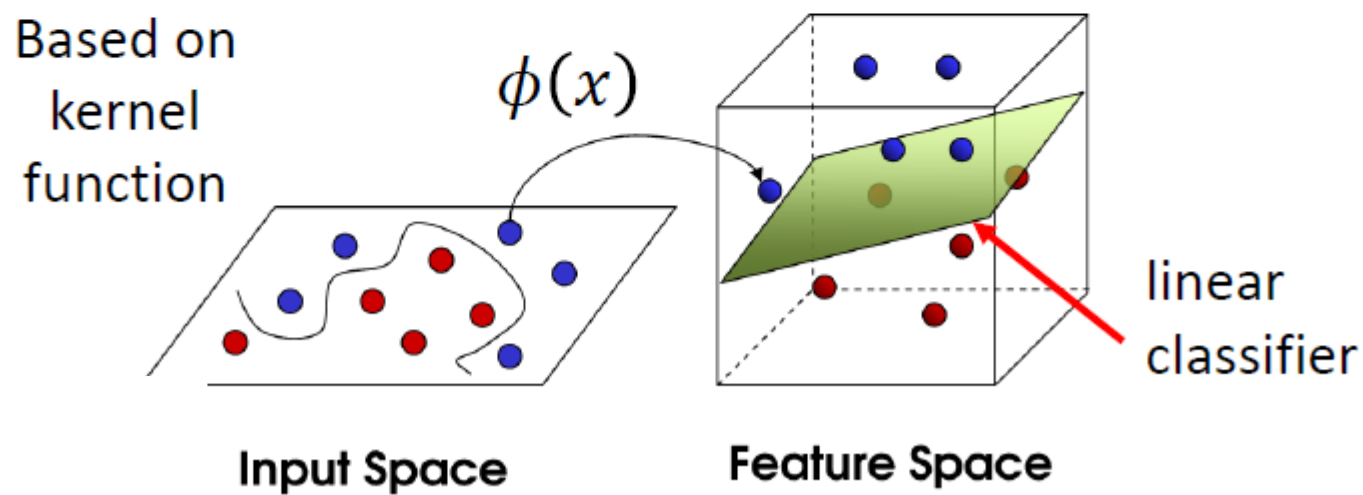
$$\ell(w, b) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \max(0, 1 - (w^T x^{(i)} + b)t^{(i)})$$



Deep Learning



SVM



Readings

- SVM implementations
 - <http://www.kernel-machines.org/software>
- Burges, Christopher JC. "A tutorial on support vector machines for pattern recognition." Data mining and knowledge discovery 2.2 (1998): 121-167.
- Sections 6.1-6.2 in the book "Pattern Recognition and Machine Learning", by Christopher M. Bishop, Springer, 2006.
- SVM - Understanding the math - Duality and Lagrange multipliers
 - <https://www.svm-tutorial.com/2016/09/duality-lagrange-multipliers/>
- MIT OpenCourseWare: Learning: Support Vector Machines
 - <https://www.youtube.com/watch?v=PwhiWxHK8o>
- Support Vector Regression (SVR): Section 7.1.4 in the book "Pattern Recognition and Machine Learning", by Christopher M. Bishop, Springer, 2006.