# 机器学习

## —— 第7章 支持向量机 ——

授课：倪张凯

*zkni@tongji.edu.cn*

*https://eezkni.github.io/*

Question：Find a hyperplane in the sample space to separate samples of different categories.

# 二分类问题 Binary Classification



Question：Find a hyperplane in the sample space to separate samples of different categories.
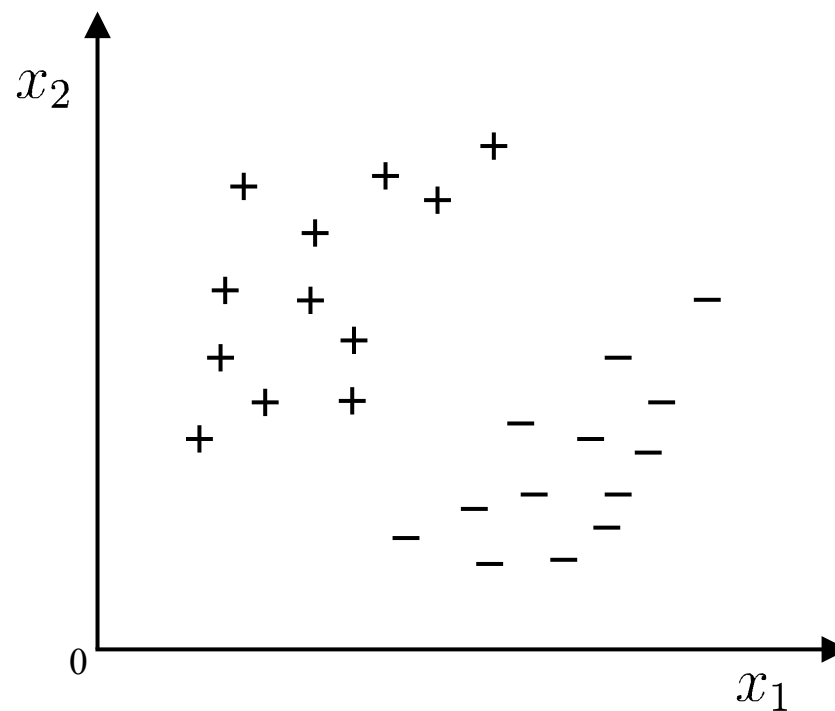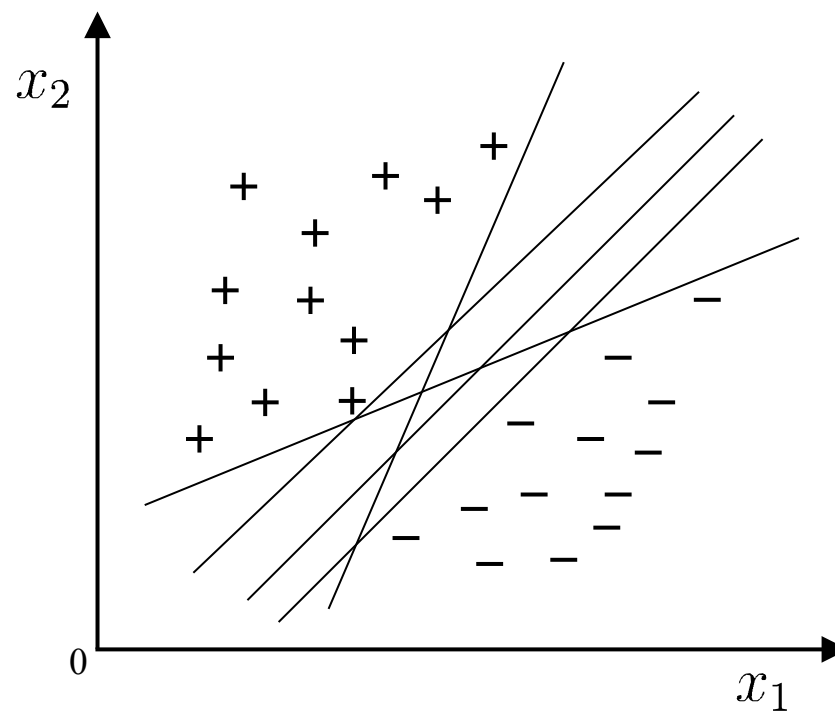
# 二分类问题 Binary Classification
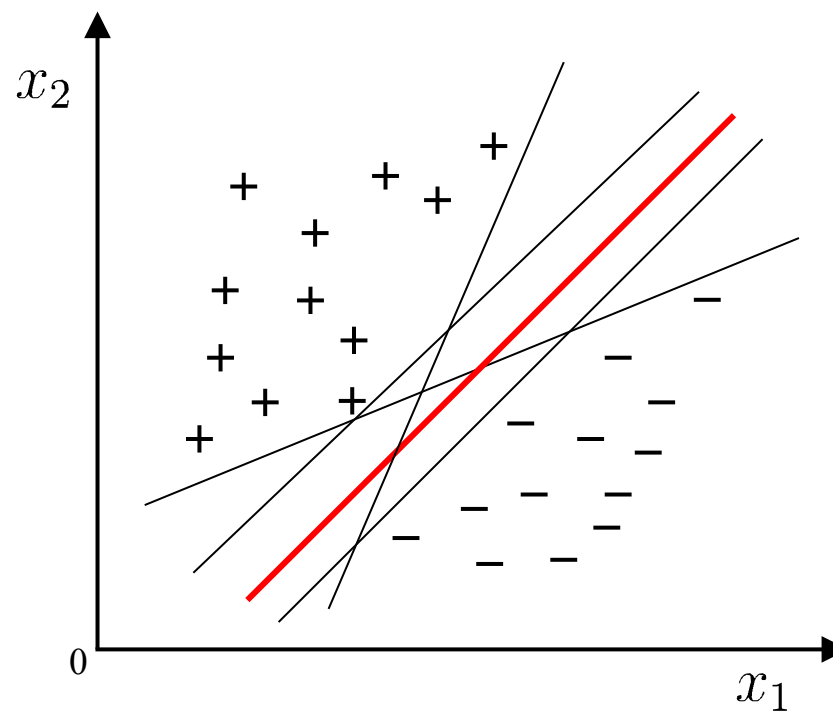


Question: Find a hyperplane in the sample space to separate samples of different categories.
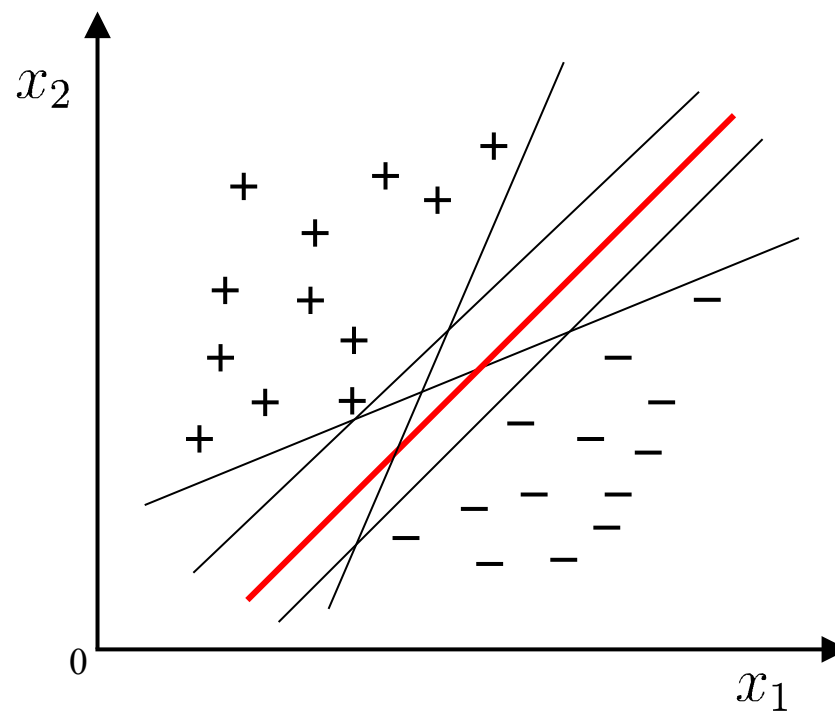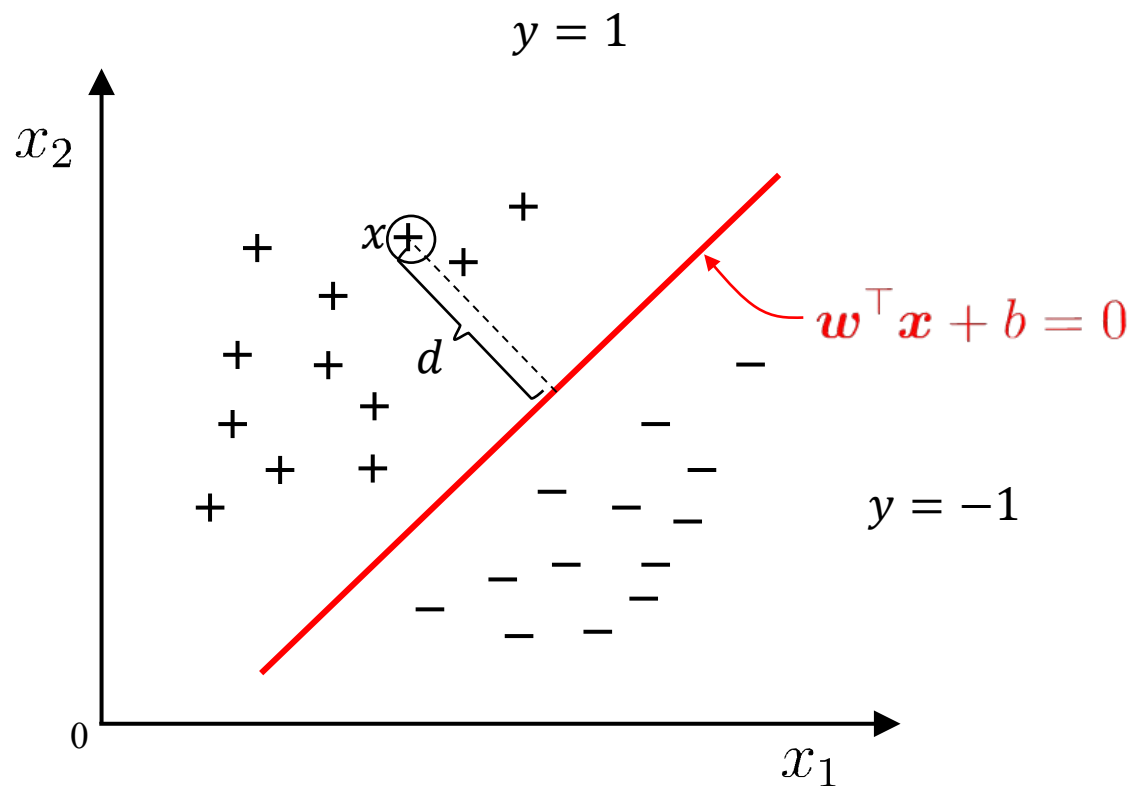
# 二分类问题 Binary Classification



It should choose "right in the middle", with good tolerance, high robustness and the strongest generalization ability

# 支持向量机 Support Vector Machine

$$y = 1$$

$x_2$

$x \bigoplus$

$d$

$\boldsymbol{w}^\top \boldsymbol{x} + b = 0$

$$y = -1$$

$x_1$

$0$

$$d = \frac{|w^{\mathrm{T}}x + b|}{||w||}$$

# 支持向量机 Support Vector Machine



margin

$y = 1$

$\boldsymbol{w}^\top \boldsymbol{x} + b = 1$

$\gamma = \dfrac{2}{\|\boldsymbol{w}\|}$

$\boldsymbol{w}^\top \boldsymbol{x} + b = 0$

$\boldsymbol{w}^\top \boldsymbol{x} + b = -1$

$y = -1$

$y_i(w^T x_i + b) \geq 1$

**maximum margin**

Support vectors

$x_2$

$x_1$

0

$$y = 1$$

$$\boldsymbol{w}^\top \boldsymbol{x} + b = 1$$

margin

$$\gamma = \frac{2}{\|\boldsymbol{w}\|}$$

$x_2$

$$\boldsymbol{w}^\top \boldsymbol{x} + b = 0$$

$$\boldsymbol{w}^\top \boldsymbol{x} + b = -1$$

$$y = -1$$

$$y_i(w^T x_i + b) \geq 1$$

**maximum margin**

Support vectors

凸二次规划

Convex Quadratic Programming

$0$

$x_1$

$$\arg\max_{\boldsymbol{w},b} \frac{2}{\|\boldsymbol{w}\|} \implies \arg\min_{\boldsymbol{w},b} \frac{1}{2}\|\boldsymbol{w}\|^2$$

$$\text{s.t.} \quad y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) \geq 1, \ i = 1, 2, \ldots, m.$$

# 对偶问题 Dual problem

在数学和优化领域中，对偶问题是与原始优化问题（原始问题）相对应的一个问题。通常，对偶问题是通过原始问题的一种变换来构建的，它可以帮助我们更好地理解原始问题的性质，提供额外的信息，或者用于求解原始问题。

- 对于最小化问题（Minimization Problem）： 对偶问题：max $g(\lambda,v)$其中，$g(\lambda,v)=\inf_x L(x,\lambda,v)$，表示对于给定$L(x,\lambda,v)$ 的最小值

- 对于最大化问题（Maximization Problem）： 对偶问题：min $g(\lambda,v)$ 其中，$g(\lambda,v)=\sup_x L(x,\lambda,v)$，表示对于给定$L(x,\lambda,v)$ 的最大值

在这里：

- $\lambda$ 和 $v$ 是对偶变量（Lagrange Multipliers），它们对应于原始问题的约束条件。

- $L(x,\lambda,v)$ 是称为拉格朗日函数（Lagrangian Function）的函数，它由原始问题的目标函数和约束条件组成，通常形式为：

   **$L(x,\lambda,v)$=目标函数$-\lambda T\cdot$(不等式约束)$-vT\cdot$(等式约束)**

通过求解对偶问题，我们可以获得原始问题的最优值的一个下界或上界。

# 拉格朗日乘数法的例子

- **问题**：假设你需要在平面上找到距离原点最近的点，但这个点必须位于直线 $y=2x+3$ 上。
- **解决方法**：

# 拉格朗日乘数法的例子

- **问题**：假设你需要在平面上找到距离原点最近的点，但这个点必须位于直线 $y=2x+3$ 上。

- **解决方法**：

  1. **目标函数**：需要最小化的目标函数是到原点的距离的平方，$f(x,y) = x^2 + y^2$。

  2. **约束**：点必须位于$y=2x+3$ 上，因此约束条件为$g(x,y) = y - 2x - 3 = 0$。

  3. **构建拉格朗日函数**：构建拉格朗日函数 $L(x, y, \lambda) = x^2 + y^2 + \lambda（y - 2x - 3）$。

  4. **求解**：通过 $L(x, y, \lambda)$ 关于 $x, y, \lambda$ 的偏导数求解并置为零，可以得到最优解。

# KKT （Karush-Kuhn-Tucker）条件的例子

- **问题**：假设你需要在一个盒子中放置最大体积的长方体，但其长、宽、高的和不能超过某个值，例如10
- **解决方法：**

# KKT（Karush-Kuhn-Tucker）条件的例子

- **问题**：假设你需要在一个盒子中放置最大体积的长方体，但其长、宽、高的和不能超过某个值，例如10

- **解决方法：**

    1.**目标函数**：最大化长方体的体积，即 $f(x,y,z)=xyz$。

    2.**约束**：长、宽、高的和不超过10，即 $g(x,y,z)=x+y+z-10\leq0$。这是一个不等式约束。

    3.**构建拉格朗日函数：拉格朗日函数 $L(x,y,z,\lambda)$ 结合了目标函数和约束条件，通过引入拉格朗日乘子 $\lambda$：**

$L(x,y,\lambda)=xyz-\lambda(x+y+z-10)$

# KKT（Karush-Kuhn-Tucker）条件的例子

$L(x,y,λ)=xyz−λ(x+y+z−10)$

## 4.应用KKT条件

- **梯度为零**：对 $L(x,y,z,λ)$ 对 $x,y,z,λ$ 的偏导数应为零。

  - $\frac{\partial L}{\partial x}=yz−λ=0$　　$\frac{\partial L}{\partial y}=xz−λ=0$　　　$\frac{\partial L}{\partial z}=xy−λ=0$

  - $\frac{\partial L}{\partial λ}=−x−y−z+10=0$

- **约束条件**：$x+y+z≤10$。

- **拉格朗日乘子非负**：$λ≥0$。

- **互补松弛性**：$λ(x+y+z−10)=0$。

## 5.求解：通过解上述方程组，我们可以找到满足约束的最优解。

- Using Lagrangian Multiplier Method and KKT Condition to solve the Optimal Value

$$\underset{\boldsymbol{w},b}{\arg\min} \quad \frac{1}{2}\|\boldsymbol{w}\|^2$$
$$\text{s.t.} \quad y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) \geq 1, \ i = 1, 2, \ldots, m.$$

- Integrated into: (Where $\alpha_i \geq 0$ is a Lagrangian multiplier)

$$L(w, b, a) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{m} \alpha_i(1 - y_i(w^T x_i + b))$$

- Let the partial derivative=0

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^{m} \alpha_i y_i x_i = 0 \qquad \frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^{m} \alpha_i y_i = 0$$

- then

$$\boldsymbol{w} = \sum_{i=1}^{m} \alpha_i y_i \boldsymbol{x}_i, \quad \sum_{i=1}^{m} \alpha_i y_i = 0.$$

$$L(w, b, a) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{m} \alpha_i(1 - y_i(w^T x_i + b))$$

$$= \frac{1}{2}w^T w - \sum_{i=1}^{m}\alpha_i y_i w^T x_i - b\sum_{i=1}^{m}\alpha_i y_i + \sum_{i=1}^{m}\alpha_i$$

$$= \frac{1}{2}w^T \sum_{i=1}^{m}\alpha_i y_i x_i - \sum_{i=1}^{m}\alpha_i y_i w^T x_i + + \sum_{i=1}^{m}\alpha_i$$

$$= \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\alpha_i y_i w^T x_i$$

Dual problem:

$$\max_{a} \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^T x_j = \min_{a} \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^{m}\alpha_i$$

$$s.t. \sum_{i=1}^{m}\alpha_i y_i = 0, \alpha_i \geq 0, i = 1,2,\dots,m$$

KKT Condition

Obtain the optima $\alpha$ (SMO, Sequential Minimal Optimization) algorithm
The optimal solution can be obtained

$$f(x) = w^T x + b = \sum_{i=1}^{m}\alpha_i y_i x_i^T x + b$$

$$\begin{cases} \alpha_i \geq 0, \\ y_i f(\boldsymbol{x}_i) \geq 1, \\ \alpha_i(y_i f(\boldsymbol{x}_i) - 1) = 0. \end{cases}$$

$y_i f(\boldsymbol{x}_i) > 1 \implies \alpha_i = 0$

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

$$L(w, b, a) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{m} \alpha_i(1 - y_i(w^T x_i + b))$$

$$= \frac{1}{2}w^T w - \sum_{i=1}^{m} \alpha_i y_i w^T x_i - b\sum_{i=1}^{m} \alpha_i y_i + \sum_{i=1}^{m} \alpha_i$$

$$= \frac{1}{2}w^T \sum_{i=1}^{m} \alpha_i y_i x_i - \sum_{i=1}^{m} \alpha_i y_i w^T x_i + + \sum_{i=1}^{m} \alpha_i$$

$$= \sum_{i=1}^{m} \alpha_i - \frac{1}{2}\sum_{i=1}^{m} \alpha_i y_i w^T x_i$$

Dual problem:

$$\max_{a} \sum_{i=1}^{m} \alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m} \alpha_i \alpha_j y_i y_j x_i^T x_j = \min_{a} \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m} \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^{m} \alpha_i$$

$$s.t. \sum_{i=1}^{m} \alpha_i y_i = 0, \alpha_i \geq 0, i = 1,2,\dots,m$$

KKT Condition

Obtain the optima $\alpha$ (SMO, Sequential Minimal Optimization) algorithm

The optimal solution can be obtained

$$f(x) = w^T x + b = \sum_{i=1}^{m} \alpha_i y_i x_i^T x + b$$

$$\begin{cases} \alpha_i \geq 0, \\ y_i f(\boldsymbol{x}_i) \geq 1, \\ \alpha_i(y_i f(\boldsymbol{x}_i) - 1) = 0. \end{cases}$$

$y_i f(\boldsymbol{x}_i) > 1 \implies \alpha_i = 0$

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

$$L(w, b, a) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{m} \alpha_i(1 - y_i(w^T x_i + b))$$

$$= \frac{1}{2}w^T w - \sum_{i=1}^{m} \alpha_i y_i w^T x_i - b\sum_{i=1}^{m} \alpha_i y_i + \sum_{i=1}^{m} \alpha_i$$

$$= \frac{1}{2}w^T \sum_{i=1}^{m} \alpha_i y_i x_i - \ldots$$

$$= \sum_{i=1}^{m} \alpha_i - \frac{1}{2}\sum_{i=1}^{m} \alpha_i y_i \ldots$$

Dual problem:

$$\max_{a} \sum_{i=1}^{m} \alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m} \alpha_i \alpha_j y_i y_j \ldots$$

$$s.t. \sum_{i=}^{m} \ldots$$

若 $\alpha_i$ = 0，则该样本将不会在左式 的求和中出现，也就不会对 f(x) 有任何影响;
若 $\alpha_i$ > 0,则必有执 f(Xi) = 1，所对应的样本点位于最大间隔边界上，是一个支持向量。

这显示出支持向量机的一个重要性质:训练完成后，大部分的训练样本都不需保留，最终模型仅与支持向量有关。

Obtain the optima $\alpha$ (SMO, Sequential Minimal Optimization) algorithm
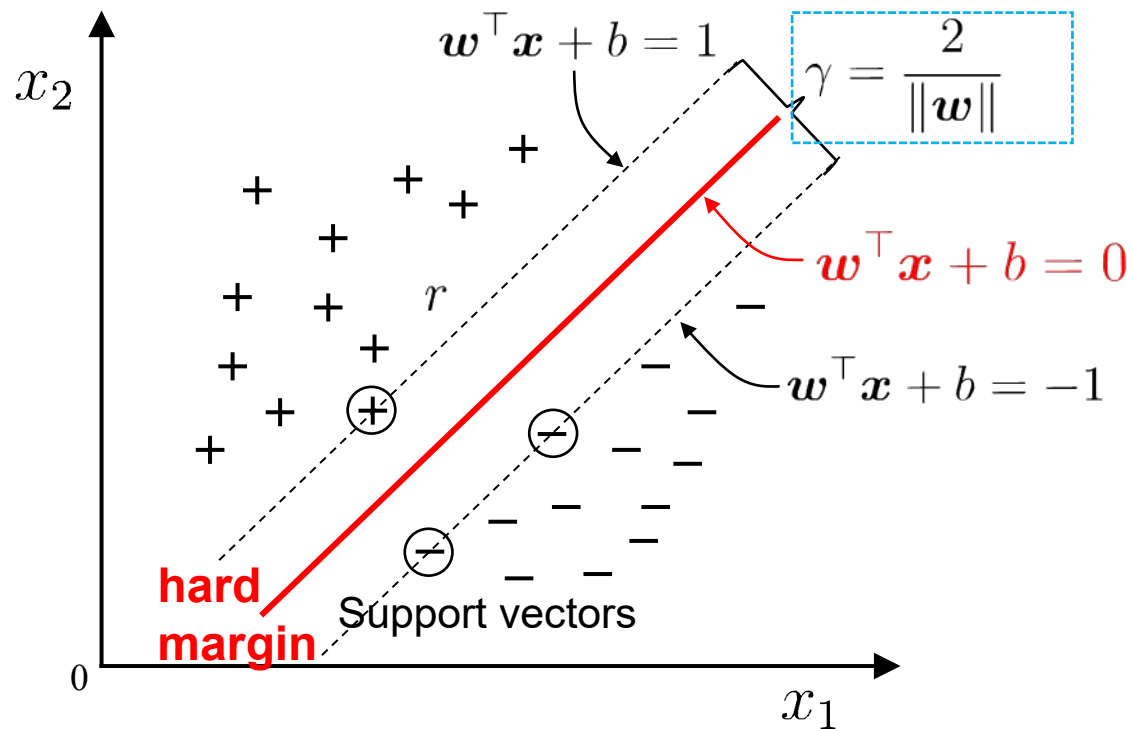
The optimal solution can be obtained

$$f(x) = w^T x + b = \sum_{i=1}^{m} \alpha_i y_i x_i^T x + b$$

$$\begin{cases} \alpha_i \geq 0, \\ y_i f(\boldsymbol{x}_i) \geq 1, \\ \alpha_i(y_i f(\boldsymbol{x}_i) - 1) = 0. \end{cases}$$

$$y_i f(\boldsymbol{x}_i) > 1 \quad \Longrightarrow \quad \alpha_i = 0$$
$$\alpha_i > 0 \quad \Longrightarrow \quad y_i f(x_i) = 1$$

# 支持向量机 Support Vector Machine



$$\boldsymbol{w}^{\top}\boldsymbol{x} + b = 1$$

$$\gamma = \frac{2}{\|\boldsymbol{w}\|}$$

$$\boldsymbol{w}^{\top}\boldsymbol{x} + b = 0$$

$$\boldsymbol{w}^{\top}\boldsymbol{x} + b = -1$$

**hard margin**

Support vectors

$$\underset{\boldsymbol{w},b}{\arg\max} \quad \frac{2}{\|\boldsymbol{w}\|} \quad \Longrightarrow \quad \underset{\boldsymbol{w},b}{\arg\min} \quad \frac{1}{2}\|\boldsymbol{w}\|^2$$

$$\text{s.t.} \quad y_i(\boldsymbol{w}^{\top}\boldsymbol{x}_i + b) \geq 1, \ i = 1, 2, \ldots, m.$$

-Q: It is difficult to determine a linearly separable hyperplane in the feature space; At the same time, it is difficult to determine whether a linearly separable result is caused by over fitting

-A: The concept of "soft margin" is introduced to allow the support vector machine to not meet the constraints on some samples

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2}\|w\|_2 + C\sum_i \varepsilon_i \qquad \text{s.t.} \quad y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0, \text{i=1,2,…m}$$

Integrated into: (Where $\alpha_i$ , $\varepsilon_i$ are Lagrangian multipliers)

$$L(w,b,a,\varepsilon,\mu) = \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{m}\varepsilon_i + \sum_{i=1}^{m}\alpha_i(1-\varepsilon_i - y_i(w^T x_i + b)) - \sum_{i=1}^{m}\mu_i\varepsilon_i$$

Let the partial derivative=0

$$\frac{\partial L(w,b,a,\varepsilon,\mu)}{\partial w} = w - \sum_{i=1}^{m}\alpha_i y_i x_i = 0 \qquad \frac{\partial L(w,b,a,\varepsilon,\mu)}{\partial b} = \sum_{i=1}^{m}\alpha_i y_i = 0$$

$$\frac{\partial L(w,b,a,\varepsilon,\mu)}{\partial \varepsilon_i} = C - \alpha_i - \mu_i = 0$$

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2}\|w\|_2 + C\sum_i \varepsilon_i \qquad \text{s.t.} \ \ y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0, \text{i=1,2,…m}$$

Integrated into: (Where $\alpha_i$ , $\varepsilon_i$ are Lagrangian multipliers)

$$L(w, b, a, \varepsilon, \mu) = \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{m}\varepsilon_i + \sum_{i=1}^{m}\alpha_i(1 - \varepsilon_i - y_i(w^T x_i + b)) - \sum_{i=1}^{m}\mu_i\varepsilon_i$$

Let the partial derivative=0

then

$$w = \sum_{i=1}^{m}\alpha_i y_i x_i \qquad\qquad \sum_{i=1}^{m}\alpha_i y_i = 0$$

$$C = \alpha_i + \mu_i$$

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2} \|w\|_2 + C \sum_i \varepsilon_i \qquad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0, \text{i=1,2,...m}$$

Integrated into: (Where $\alpha_i$ , $\varepsilon_i$ are Lagrangian multipliers)

$$L(w, b, a, \varepsilon, \mu) = \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{m}\varepsilon_i + \sum_{i=1}^{m}\alpha_i(1 - \varepsilon_i - y_i(w^T x_i + b)) - \sum_{i=1}^{m}\mu_i\varepsilon_i$$

Dual problem:

$$\max_{a} \sum_{i=1}^{m}\alpha_i - \sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^T x_j \quad s.t. \sum_{i=1}^{m}\alpha_i y_i = 0 , \qquad C \geq \alpha_i \geq 0, i = 1,2,...,m$$

Obtain the optima $\alpha$ (SMO, Sequential Minimal Optimization) algorithm

KKT Condition

The optimal solution can be obtained

$$f(x) = w^T x + b = \sum_{i=1}^{m}\alpha_i y_i x_i^T x + b$$

$$\begin{cases} \alpha_i \geq 0, \mu_i \geq 0 \\ y_i f(x_i) - 1 + \varepsilon_i \geq 0 \\ \alpha_i(y_i f(x_i) - 1 + \varepsilon_i) = 0 \\ \varepsilon_i \geq 0, \mu_i\varepsilon_i = 0 \end{cases}$$

# 软间隔支持向量机 Soft Margin Support Vector Machine

对任意训练样本，总有$\alpha_i$ =0或 $y_i f(x_i)$=1−$\varepsilon_i$

若$\alpha_i$= 0，则该样本将不会在左式 的求和中出现，也就不会对 f(x) 有任何影响;

若$\alpha_i$ > 0,则必有 $y_i f(x_i)$=1−$\varepsilon_i$ ，即该样本是一个支持向量。

若$\alpha_i$ < C，则$\mu_i$ > 0，进而有 $\varepsilon_i$ = 0，即该样本恰在最大间隔边界上;

若 $\alpha_i$ =C ， 则有$\mu_i$ = 0，此时若 $\varepsilon_i$ ≤1，则该样本落在最大间隔内部；若 $\varepsilon_i$ > 1 则该样本被错误分类。

由此可看出软间隔支持向量机的最终模型仅与支持向量有关，仍保持了稀疏性。

$+ b) \geq 1 - \varepsilon_i$

i=1,2,…m

$(w^T x_i + b)) - \sum_{i=1}^{m} \mu_i \varepsilon_i$

$C \geq \alpha_i \geq 0, i = 1, 2, \ldots, m$

KKT Condition

$$\begin{cases} \alpha_i \geq 0, \mu_i \geq 0 \\ y_i f(x_i) - 1 + \varepsilon_i \geq 0 \\ \alpha_i(y_i f(x_i) - 1 + \varepsilon_i) = 0 \\ \varepsilon_i \geq 0, \mu_i \varepsilon_i = 0 \end{cases}$$

Du

$f(x) = w^T x + b = \sum_{i=1}^{m} \alpha_i y_i x_i^T x + b$

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2}\|w\|_2 + C\sum_i \varepsilon_i \quad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

Dual Problem:

$$\max_a \sum_{i=1}^m \alpha_i - \frac{1}{2}\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j = \min_a \frac{1}{2}\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^m \alpha_i$$

$$s.t. \sum_{i=1}^m \alpha_i y_i = 0, C \geq \alpha_i \geq 0, i = 1,2,\dots,m$$

$$L(w, b, a) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{m} \alpha_i(1 - y_i(w^T x_i + b))$$

$$= \frac{1}{2}w^T w - \sum_{i=1}^{m}\alpha_i y_i w^T x_i - b\sum_{i=1}^{m}\alpha_i y_i + \sum_{i=1}^{m}\alpha_i$$

$$= \frac{1}{2}w^T \sum_{i=1}^{m}\alpha_i y_i x_i - \sum_{i=1}^{m}\alpha_i y_i w^T x_i + + \sum_{i=1}^{m}\alpha_i$$

$$= \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\alpha_i y_i w^T x_i$$

Dual problem:

$$\max_{a}\sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^T x_j = \min_{a}\frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^{m}\alpha_i$$

$$s.t. \sum_{i=1}^{m}\alpha_i y_i = 0 , \alpha_i \geq 0, i = 1,2,\dots,m$$

KKT Condition

Obtain the optima $\alpha$ (SMO, Sequential Minimal Optimization) algorithm

The optimal solution can be obtained

$$f(x) = w^T x + b = \sum_{i=1}^{m}\alpha_i y_i x_i^T x + b$$

$$\begin{cases} \alpha_i \geq 0, \\ y_i f(\boldsymbol{x}_i) \geq 1, \\ \alpha_i(y_i f(\boldsymbol{x}_i) - 1) = 0. \end{cases}$$

$y_i f(\boldsymbol{x}_i) > 1 \implies \alpha_i = 0$

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

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2}\|w\|_2 + C\sum_i \varepsilon_i \quad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

Dual Problem:

$$\max_a \sum_{i=1}^m \alpha_i - \frac{1}{2}\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j = \min_a \frac{1}{2}\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^m \alpha_i$$

$$s.t. \sum_{i=1}^m \alpha_i y_i = 0, C \geq \alpha_i \geq 0, i = 1,2,\ldots,m$$

# 软间隔支持向量机 Soft Margin Support Vector Machine

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2} \|w\|_2 + C \sum_i \varepsilon_i \quad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

在软间隔SVM中，引入正则化参数C控制了对分类错误的惩罚。
　　较小的C值会导致更大的间隔，容忍更多的分类错误，从而提高模型的容错性；
　　较大的C值会更强调正确分类，但可能导致对异常点更敏感。

$$\min_a \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \, \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^{m} \alpha_i$$

$$C \geq \alpha_i \geq 0, i = 1,2,\ldots,m$$

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2} \|w\|_2 + C \sum_i \varepsilon_i \quad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0, \text{i=1,2,…m}$$

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2} \|w\|_2 + C \sum_i \varepsilon_i \quad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0, i=1,2,\ldots m$$

$$\varepsilon^i = max(0, 1 - y^{(i)}f(x^{(i)}))$$

The slack variable $\varepsilon$ can be introduced :

$$\min_{w,b} \frac{1}{2}\|w\|_2 + C\sum_i \varepsilon_i \quad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0, \text{i=1,2,...m}$$

$$\boxed{\varepsilon^i = max(0, 1 - y^{(i)}f(x^{(i)}))}$$

The slack variable $\varepsilon$ indicating the extent to which the sample does not meet the constraint.
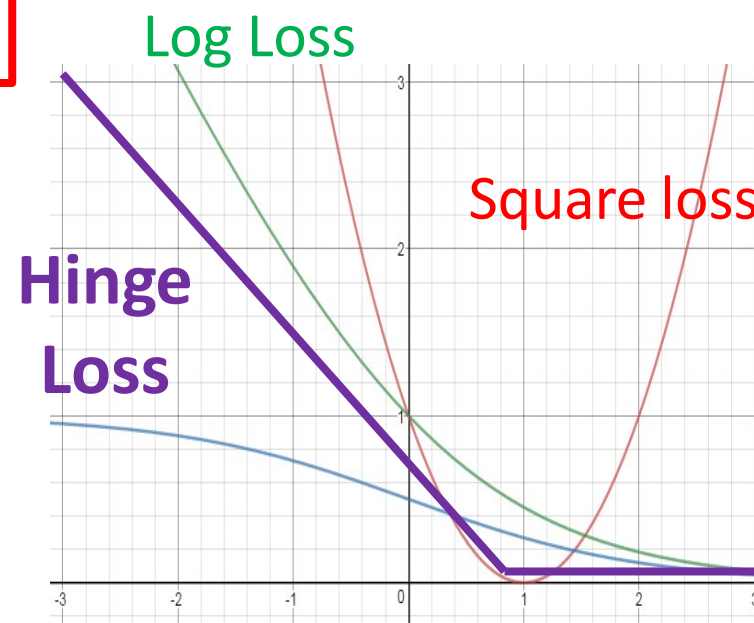
The slack variable $\varepsilon$ indicating the extent to which the sample does not meet the constraint.

$$\varepsilon^{i} = max(0, 1 - y^{(i)}f(x^{(i)}))$$

Hinge Loss function

Log Loss

Square loss

**Hinge Loss**

-Q: What if there is no hyperplane that can correctly divide two types of samples?

-A: The samples are mapped from the original space to a higher dimensional feature space, making the samples linearly separable in this feature space



$$\boldsymbol{x} \mapsto \phi(\boldsymbol{x})$$

$$L(w, b, a) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{m} \alpha_i(1 - y_i(w^T x_i + b))$$

$$= \frac{1}{2} w^T w - \sum_{i=1}^{m} \alpha_i y_i w^T x_i - b \sum_{i=1}^{m} \alpha_i y_i + \sum_{i=1}^{m} \alpha_i$$

$$= \frac{1}{2} w^T \sum_{i=1}^{m} \alpha_i y_i x_i - \sum_{i=1}^{m} \alpha_i y_i w^T x_i + + \sum_{i=1}^{m} \alpha_i$$

$$= \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i=1}^{m} \alpha_i y_i w^T x_i$$

Dual problem:

$$\max_a \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j y_i y_j x_i^T x_j = \min_a \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^{m} \alpha_i$$

$$s.t. \sum_{i=1}^{m} \alpha_i y_i = 0, \alpha_i \geq 0, i = 1,2,\dots,m$$

KKT Condition

Obtain the optima $\alpha$ (SMO, Sequential Minimal Optimization) algorithm
The optimal solution can be obtained

$$f(x) = w^T x + b = \sum_{i=1}^{m} \alpha_i y_i x_i^T x + b$$

$$\begin{cases} \alpha_i \geq 0, \\ y_i f(\boldsymbol{x}_i) \geq 1, \\ \alpha_i(y_i f(\boldsymbol{x}_i) - 1) = 0. \end{cases}$$

$$y_i f(\boldsymbol{x}_i) > 1 \implies \alpha_i = 0$$

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

sample $\boldsymbol{x} \mapsto \phi(\boldsymbol{x})$, then the hyperplane $f(\boldsymbol{x}) = \boldsymbol{w}^{\top}\phi(\boldsymbol{x}) + b$

Original question:

$$\min_{\boldsymbol{w},b} \quad \frac{1}{2}\|\boldsymbol{w}\|^2$$
$$\text{s.t.} \quad y_i(\boldsymbol{w}^{\top}\phi(\boldsymbol{x}_i) + b) \geq 1, \ i = 1, 2, \ldots, m.$$

Dual problem:

$$\min_{\boldsymbol{\alpha}} \quad \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j \phi(\boldsymbol{x}_i)^{\top}\phi(\boldsymbol{x}_j) - \sum_{i=1}^{m}\alpha_i$$
$$\text{s.t.} \quad \sum_{i=1}^{m}\alpha_i y_i = 0, \ \alpha_i \geq 0, \ i = 1, 2, \ldots, m.$$

Optimal solution:

$$f(\boldsymbol{x}) = \boldsymbol{w}^{\top}\phi(\boldsymbol{x}) + b = \sum_{i=1}^{m}\alpha_i y_i \phi(\boldsymbol{x}_i)^{\top}\phi(\boldsymbol{x}) + b$$

sample $\boldsymbol{x} \mapsto \phi(\boldsymbol{x})$, then the hyperplane $f(\boldsymbol{x}) = \boldsymbol{w}^\top \phi(\boldsymbol{x}) + b$

Original question:

Kernel Function

$$\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \phi(\boldsymbol{x}_i)^\top \phi(\boldsymbol{x}_j)$$

$$\min_{\boldsymbol{w}, b} \quad \frac{1}{2} \|\boldsymbol{w}\|^2$$
$$\text{s.t.} \quad y_i(\boldsymbol{w}^\top \phi(\boldsymbol{x}_i) + b) \geq 1, \ i = 1, 2, \ldots, m.$$

Dual problem:

$$\min_{\boldsymbol{\alpha}} \quad \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j y_i y_j \phi(\boldsymbol{x}_i)^\top \phi(\boldsymbol{x}_j) - \sum_{i=1}^{m} \alpha_i$$
$$\text{s.t.} \quad \sum_{i=1}^{m} \alpha_i y_i = 0, \ \alpha_i \geq 0, \ i = 1, 2, \ldots, m.$$

Optimal solution:

$$f(\boldsymbol{x}) = \boldsymbol{w}^\top \phi(\boldsymbol{x}) + b = \sum_{i=1}^{m} \alpha_i y_i \phi(\boldsymbol{x}_i)^\top \phi(\boldsymbol{x}) + b$$

- Given the training dataset of (data,label) pairs,

$$D = \left\{\left(x^{(i)}, y^{(i)}\right)\right\}_{i=1.2,\ldots,N}$$

let the machine learn a function from data to label

$$y^{(i)} \approx f_\theta\left(x^{(i)}\right)$$

- Function set $\{f_\theta(x^{(i)})\}$ is called hypothesis space
- Learning is referred to as updating the parameter $\theta$ to make the prediction closed to the corresponding label

# 监督学习 Supervised Learning

- Given the training dataset of (data,label) pairs,

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1.2,...,N}$$

let the machine learn a function from data to label

$$y^{(i)} \approx f_\theta(x^{(i)})$$

- How to learn?

  Update the parameter to make the prediction closed to the corresponding label

    1. What is the learning objective?

    2. How to update the parameters?

- Minimize the total loss

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^{N} L(y^{(i)}, f_\theta(x^{(i)}))$$

Loss function $L(y^{(i)}, f_\theta(x^{(i)}))$ measures the error between the label and prediction for single sample.

We have used squared loss：

$$\frac{1}{2}(y^{(i)} - f_\theta(x^{(i)}))^2$$

Log loss：

$$-y^{(i)} log \left( (f_\theta(x)) \right) - (1 - y^{(i)}) log \left( 1 - (f_\theta(x)) \right)$$

# 线性回归 Linear Regression

- Given the training dataset of (data,label) pairs,

$$D = \left\{\left(x^{(i)}, y^{(i)}\right)\right\}_{i=1.2,\ldots,N}$$

let the machine learn a function from data to label

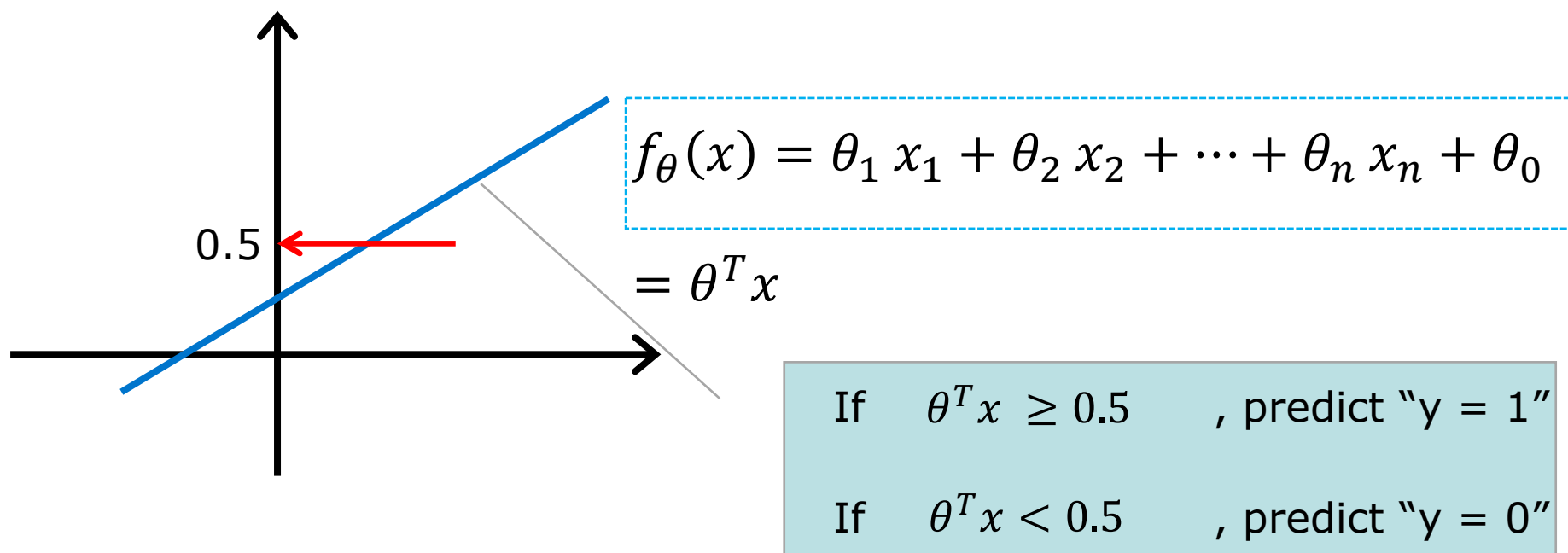$$y^{(i)} \approx \boxed{f_\theta\left(x^{(i)}\right)}$$

$$f_\theta(x) = \theta_1\, x_1 + \theta_2\, x_2 + \cdots + \theta_n\, x_n + \theta_0$$

- Function set $\{f_\theta(x^{(i)})\}$ is called hypothesis space

- Learning is referred to as updating the parameter $\theta$ to make the prediction closed to the corresponding label

$$y \in \{0, 1\}$$

0: "Negative Class" (如，坏瓜)
1: "Positive Class" (如, 好瓜)



$$f_\theta(x) = \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n + \theta_0$$

$$= \theta^T x$$

0.5

If $\quad \theta^T x \geq 0.5 \quad$ , predict "y = 1"

If $\quad \theta^T x < 0.5 \quad$ , predict "y = 0"

$$y \in \{0, 1\}$$

0: "Negative Class" (如，坏瓜)
1: "Positive Class" (如, 好瓜)



0.5

$$f_\theta(x) = \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n + \theta_0$$

$$= \theta^T x$$

If $\theta^T x \geq 0.5$, predict "y = 1"

If $\theta^T x < 0.5$, predict "y = 0"

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ 0 & f(x) < 0 \end{cases}$$

# 逻辑斯蒂回归 Logistic regression

$y \in \{0, 1\}$  0: "Negative Class" (如，坏瓜)
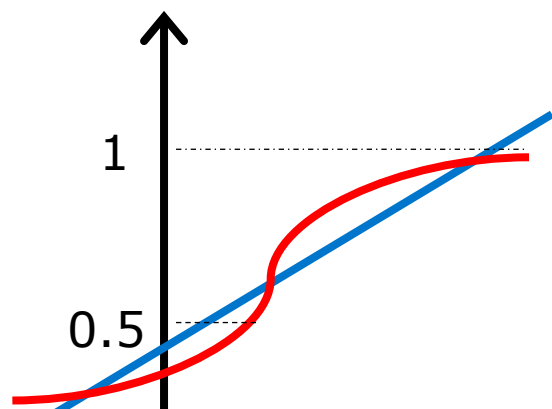1: "Positive Class" (如, 好瓜)



$$f_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$P(y = 1/x) = \frac{1}{1 + e^{-\theta^T x}}$$

= estimated probability
that y = 1,
given x, parameterized by $\theta$

$$P(y = 0/x) = \frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}}$$

= estimated probability
that y = 0,
given x, parameterized by $\theta$

$y \in \{0, 1\}$

0: "Negative Class" (如，坏瓜)
1: "Positive Class" (如, 好瓜)



$$f_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

= estimated probability
that y = 1,
given x, parameterized by $\theta$

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ 0 & f(x) < 0 \end{cases}$$

= estimated probability
that y = 0,
given x, parameterized by $\theta$

# 逻辑斯蒂回归 Logistic regression

- Given the training dataset of (data,label) pairs,

$$D = \left\{\left(x^{(i)}, y^{(i)}\right)\right\}_{i=1.2,\ldots,N}$$

let the machine learn a function from data to label

$$y^{(i)} \approx f_\theta\left(x^{(i)}\right) \quad \Rightarrow \quad f_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ 0 & f(x) < 0 \end{cases}$$

Cross-entropy loss：

$$-y^{(i)} log\left(\left(f_\theta(x)\right)\right) - (1 - y^{(i)}) log\left(1 - (f_\theta(x)\right)$$

- Step 1: Function (Model)

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ 0 & f(x) < 0 \end{cases}$$

- Step 1: Function (Model)

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ -1 & f(x) < 0 \end{cases}$$

- Step 1: Function (Model)

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ -1 & f(x) < 0 \end{cases}$$

- Step 2: Cost function

$$C(f) = \sum_n \delta\big(g\big(x^{(i)}\big) \neq y^{(i)}\big)\big)$$

- Step 1: Function (Model)

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ -1 & f(x) < 0 \end{cases}$$

- Step 2: Cost function

$$C(f) = \sum_n \delta\big(g(x^{(i)}) \neq y^{(i)})\big)$$

- Step 1: Function (Model)

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ -1 & f(x) < 0 \end{cases}$$

- Step 2: Cost function

$$C(f) = \sum_n \delta\big(g(x^{(i)}) \neq y^{(i)})\big)$$

The number of times get incorrect results on training data.

- Step 1: Function (Model)

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ -1 & f(x) < 0 \end{cases}$$

- Step 2: Cost function

$$C(f) = \sum_n \delta\big(g\big(x^{(i)}\big) \neq y^{(i)}\big))$$

- Step 3： Training by gradient descent is difficult

- Step 1: Function (Model)

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ -1 & f(x) < 0 \end{cases}$$

- Step 2: Cost function

$$C(f) = \sum_{n} L\big(f(x^{(i)}), y^{(i)}\big))$$

- Step 3：Training by gradient descent is difficult

# 二分类 Binary Classification

- Step 1: Function (Model)

$$g(x) = \begin{cases} 1 & f(x) \geq 0 \\ -1 & f(x) < 0 \end{cases}$$

$$\Rightarrow \quad \begin{array}{c} y^{(i)}f\big(x^{(i)}\big) \geq 0 \\ y^{(i)}f\big(x^{(i)}\big) < 0 \end{array}$$

- Step 2: Cost function

$$C(f) = \sum_n L\big(f\big(x^{(i)}\big), y^{(i)}\big))$$

- Step 3：Training by gradient descent is difficult

The slack variable $\varepsilon$ indicating the extent to which the sample does not meet the constraint.

$$\varepsilon^{\mathrm{i}} = max(0, 1 - y^{(i)}f(x^{(i)}))$$

Hinge Loss function

Log Loss

Square loss

**Hinge Loss**

# 损失函数 Loss function



Sigmoid +
cross entropy

Square loss

Ideal loss $\delta(g(x^{(i)}) \neq y^{(i)})$

Sigmoid +
Square loss

Larger value, smaller loss

$y^{(i)}f(x^{(i)})$

# 损失函数 Loss function

Square Loss:



Sigmoid +
cross entropy

Square loss

$$L(y^{(i)}, f(x^{(i)})) = (y^{(i)}f(x^{(i)}) - 1)^2$$

Sigmoid +
Square loss

Larger value, smaller loss

$$y^{(i)}f(x^{(i)})$$

Square Loss:

If $y^{(i)} = 1, f(x^{(i)})$ close to 1

If $y^{(i)} = -1, f(x^{(i)})$ close to -1

Sigmoid + cross entropy

Square loss

$$\mathbf{L(y^{(i)}, f(x^{(i)})) = (y^{(i)}f(x^{(i)}) - 1)^2}$$

Sigmoid + Square loss

If $y^{(i)} = 1$, L= $\mathbf{(f(x^{(i)}) - 1)^2}$

If $y^{(i)} = -1$, L= $\mathbf{(f(x^{(i)}) + 1)^2}$

Larger value, smaller loss

$y^{(i)}f(x^{(i)})$

Square Loss:

Sigmoid +
cross entropy

Square loss

$$\mathbf{L}(\mathbf{y^{(i)}}, \mathbf{f}(\mathbf{x^{(i)}})) = (\mathbf{y^{(i)}f(x^{(i)})} - \mathbf{1})^{\mathbf{2}}$$

Sigmoid +
Square loss

Larger value, smaller loss

$$y^{(i)}f(x^{(i)})$$

Square Loss:



Sigmoid +
cross entropy

Square loss

Sigmoid +
Square loss

$$\mathbf{L}(\mathbf{y^{(i)}}, \mathbf{f(x^{(i)})}) = (\mathbf{\sigma(y^{(i)}f(x^{(i)})) - 1})^2$$

Larger value, smaller loss

$$y^{(i)}f(x^{(i)})$$

Square Loss:

$$\text{If } y^{(i)} = 1, \sigma(f(x^{(i)})) \text{ close to } 1$$

$$\text{If } y^{(i)} = -1, \sigma(f(x^{(i)})) \text{ close to } 0$$

Sigmoid +

$$\text{If } y^{(i)} = 1, L = (\sigma(\mathbf{f}(\mathbf{x}^{(i)})) - \mathbf{1})^{\mathbf{2}}$$

$$\text{If } y^{(i)} = -1, L = (1 - \sigma(\mathbf{f}(\mathbf{x}^{(i)})) - \mathbf{1})^{\mathbf{2}}$$

loss

Sigmoid +
Square loss

$$\mathbf{L}(\mathbf{y}^{(\mathbf{i})}, \mathbf{f}(\mathbf{x}^{(\mathbf{i})})) = (\boldsymbol{\sigma}(\mathbf{y}^{(\mathbf{i})}\mathbf{f}(\mathbf{x}^{(\mathbf{i})})) - \mathbf{1})^{\mathbf{2}}$$

Larger value, smaller loss

$$y^{(i)} f(x^{(i)})$$

# 损失函数 Loss function

Cross entropy:

If $y^{(i)} = 1$,  cross entropy  ✗  $\boldsymbol{\sigma}(f(x^{(i)}))$

If $y^{(i)} = -1$,  ✗  $1 - \boldsymbol{\sigma}(f(x^{(i)}))$

Sigmoid + cross entropy

Square loss

$$\mathbf{L(y^{(i)}, f(x^{(i)})) = ln(1 + exp(- y^{(i)} f(x^{(i)})))}$$

Sigmoid + Square loss

Divided by ln2 here

Larger value, smaller loss

$$y^{(i)} \left( f(x^{(i)}) \right)$$

cross entropy

Cross entropy:

If $y^{(i)} = 1$, ❌ cross entropy $\boldsymbol{\sigma}(f(x^{(i)}))$

If $y^{(i)} = -1$, $1 - \boldsymbol{\sigma}(f(x^{(i)}))$

Sigmoid + cross entropy

Square loss

$$\mathbf{L}(\mathbf{y^{(i)}}, \mathbf{f(x^{(i)})}) = \mathbf{ln(1 + exp(-y^{(i)}f(x^{(i)})))}$$

Sigmoid + Square loss

Divided by ln2 here

Larger value, smaller loss

$$y^{(i)} \left( f(x^{(i)}) \right)$$

Sigmoid + cross entropy

Square loss

Sigmoid + Square loss

Divided by ln2 here

Larger value, smaller loss

$$y^{(i)}f(x^{(i)})$$

Sigmoid + cross entropy

Square loss

Sigmoid + Square loss

Divided by ln2 here

Larger value, smaller loss

$$y^{(i)} f(x^{(i)})$$

# 损失函数 Loss function



**Hinge Loss**

$$L(y^{(i)}, f(x^{(i)})) = \max(0, \ 1 - y^{(i)}f(x^{(i)}))$$

Sigmoid + cross entropy

Sigmoid + Square loss

Larger value, smaller loss

$$y^{(i)}f(x^{(i)})$$
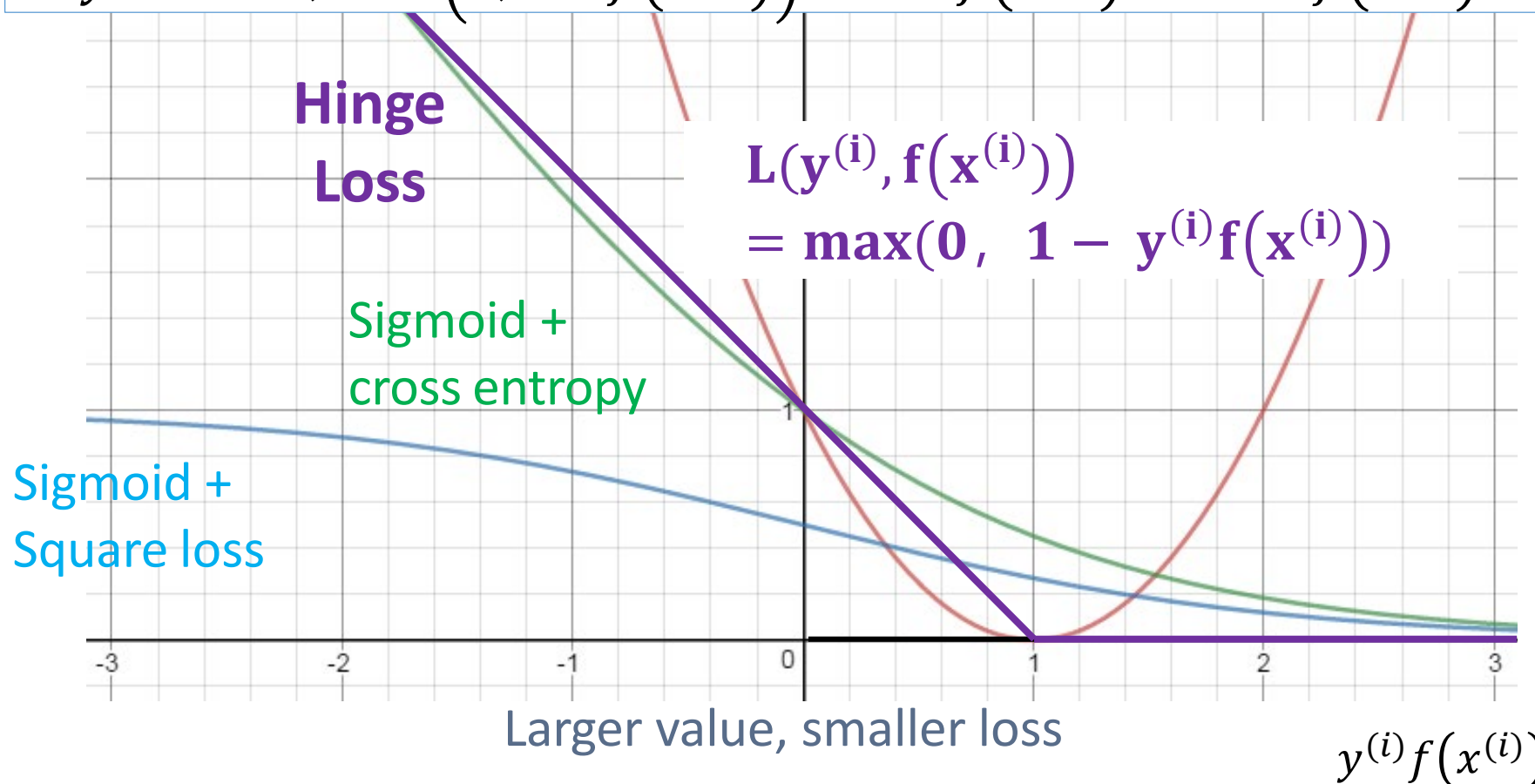
# 损失函数 Loss function

$$\text{If } y^{(i)} = 1, \quad \max\left(0, 1 - f(x^{(i)})\right) \quad 1 - f(x^{(i)}) < 0 \quad f(x^{(i)}) > 1$$

$$\text{If } y^{(i)} = -1, \max\left(0, 1 + f(x^{(i)})\right) \quad 1 + f(x^{(i)}) < 0 \quad f(x^{(i)}) < -1$$

**Hinge Loss**

**Sigmoid + cross entropy**

$$\mathbf{L(y^{(i)}, f(x^{(i)}))} = \mathbf{max(0, \ 1 - y^{(i)}f(x^{(i)}))}$$

Sigmoid + Square loss

Larger value, smaller loss

$$y^{(i)}f(x^{(i)})$$

# 损失函数 Loss function

$$\text{If } y^{(i)} = 1, \quad \max\left(0, 1 - f(x^{(i)})\right) \quad 1 - f(x^{(i)}) < 0 \quad f(x^{(i)}) > 1$$

$$\text{If } y^{(i)} = -1, \max\left(0, 1 + f(x^{(i)})\right) \quad 1 + f(x^{(i)}) < 0 \quad f(x^{(i)}) < -1$$



**Hinge Loss**

Square loss

Sigmoid + cross entropy

Sigmoid + Square loss

Good enough

# 损失函数 Loss function

$$\text{If } y^{(i)} = 1, \quad \max\left(0, 1 - f(x^{(i)})\right) \quad 1 - f(x^{(i)}) < 0 \qquad f(x^{(i)}) > 1$$

$$\text{If } y^{(i)} = -1, \max\left(0, 1 + f(x^{(i)})\right) \quad 1 + f(x^{(i)}) < 0 \qquad f(x^{(i)}) < -1$$

**Hinge Loss**

Square loss

Sigmoid + cross entropy

Sigmoid + Square loss

多好一点点

penalty   Good enough

# 损失函数 Loss function

$$\text{If } y^{(i)} = 1, \quad \max\left(0, 1 - f(x^{(i)})\right) \quad 1 - f(x^{(i)}) < 0 \quad f(x^{(i)}) > 1$$

$$\text{If } y^{(i)} = -1, \max\left(0, 1 + f(x^{(i)})\right) \quad 1 + f(x^{(i)}) < 0 \quad f(x^{(i)}) < -1$$

- Step 1: Function (Model)

$$f(x) = \sum_j w_j x_j + b$$

- Step 2: Cost function

$$C(f) = \sum_n L\big(f\big(x^{(i)}\big), y^{(i)}\big)\big)$$

- Step 1: Function (Model)

New w

$$f(x) = \sum_j w_j\, x_j + b \;=\; \begin{bmatrix} w \\ b \end{bmatrix} \cdot \begin{bmatrix} x \\ 1 \end{bmatrix} \;=\; w^T x$$

New x

- Step 2: Cost function

$$= \sum_i L\big(f\big(x^{(i)}\big), y^{(i)}\big) + \lambda\|w\|_2$$

$$L(y^{(i)}, f\big(x^{(i)}\big) = max(0,1 - \, y^{(i)} f\big(x^{(i)}\big))$$

Hinge loss

- Step 1: Function (Model)

$$f(x) = \sum_j w_j x_j + b = \begin{bmatrix} w \\ b \end{bmatrix} \cdot \begin{bmatrix} x \\ 1 \end{bmatrix}$$

- Step 2: Cost function

$$= \sum_i L\big(f(x^{(i)}), y^{(i)}\big) + \lambda\|w\|_2$$

convex

$$L(y^{(i)}, f(x^{(i)})) = max(0, 1 - y^{(i)}f(x^{(i)}))$$

- Step 3: gradient descent?

- Step 1: Function (Model)

$$f(x) = \sum_j w_j\, x_j + b \;=\; \begin{bmatrix} w \\ b \end{bmatrix} \cdot \begin{bmatrix} x \\ 1 \end{bmatrix}$$

- Step 2: Cost function  $$= \sum_i L\big(f\big(x^{(i)}\big), y^{(i)}\big) + \lambda\|w\|_2$$

convex

$$L\big(y^{(i)}, f\big(x^{(i)}\big)\big) = max\big(0, 1 - y^{(i)} f\big(x^{(i)}\big)\big)$$

Compared with logistic regression,
    linear SVM has different  (_____) ?

- Step 1: Function (Model)

$$f(x) = \sum_j w_j\, x_j + b = \begin{bmatrix} w \\ b \end{bmatrix} \cdot \begin{bmatrix} x \\ 1 \end{bmatrix}$$

- Step 2: Cost function

$$= \sum_i L\big(f(x^{(i)}), y^{(i)}\big) + \lambda\|w\|_2$$

convex

$$L(y^{(i)}, f(x^{(i)})) = max(0, 1 - y^{(i)} f(x^{(i)}))$$

Compared with logistic regression,
   linear SVM has different  (_____)  ?

Ignore regularization for simplicity

$$\sum_i L\big(f(x^{(i)}), y^{(i)}\big)$$

$$\frac{\partial L\big(f(x^{(i)}), y^{(i)}\big)}{\partial w_j} =$$

Ignore regularization for simplicity

$$\sum_i L\big(f(x^{(i)}), y^{(i)}\big)$$

$$\frac{\partial L\big(f(x^{(i)}), y^{(i)}\big)}{\partial w_j} = \frac{\partial L\big(f(x^{(i)}), y^{(i)}\big)}{\partial f(x^{(i)})} \frac{\partial f(x^{(i)})}{\partial w_j}$$

Ignore regularization for simplicity

$$\sum_i L\big(f(x^{(i)}), y^{(i)}\big) \qquad L(y^{(i)}, f(x^{(i)}) = max(0, 1 - y^{(i)}f(x^{(i)}))$$

$$\frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial w_j} = \frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial f(x^{(i)})} \underbrace{\frac{\partial f(x^{(i)})}{\partial w_j}}_{x_j^i} \qquad \boxed{\begin{array}{l} f(x^n) \\ = w^T \cdot x^n \end{array}}$$

Ignore regularization for simplicity

$$\sum_i L\big(f(x^{(i)}), y^{(i)}\big) \qquad L(y^{(i)}, f(x^{(i)}) = max(0, 1 - y^{(i)}f(x^{(i)}))$$

$$\frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial w_j} = \frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial f(x^{(i)})} \frac{\partial f(x^{(i)})}{\partial w_j} \qquad \boxed{\begin{array}{c} f(x^n) \\ = w^T \cdot x^n \end{array}}$$

$$x_j^i$$

$$\frac{\partial max(0, 1 - y^{(i)}f(x^{(i)}))}{\partial f(x^{(i)})} = \begin{cases} -y^{(i)} & \text{If } y^{(i)}f(x^{(i)}) < 1 \\ \\ 0 & \text{otherwise} \end{cases}$$

Ignore regularization for simplicity

$$\sum_i L\big(f(x^{(i)}), y^{(i)}\big)$$

$$L(y^{(i)}, f(x^{(i)}) = max(0, 1 - y^{(i)}f(x^{(i)}))$$

$$\frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial w_j} = \frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial f(x^{(i)})} \frac{\partial f(x^{(i)})}{\partial w_j}$$

$$x_j^i$$

$$\boxed{f(x^n) = w^T \cdot x^n}$$

$$\frac{\partial max(0, 1 - y^{(i)}f(x^{(i)}))}{\partial f(x^{(i)})} = \begin{cases} -y^{(i)} & \text{If } y^{(i)}f(x^{(i)}) < 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial \sum_i L(y^{(i)}, f(x^{(i)}))}{\partial w_j} = \sum_i -\delta\big(y^{(i)}f(x^{(i)}) < 1\big) y^{(i)} x_j^i$$

Ignore regularization for simplicity

$$\sum_i L\big(f(x^{(i)}), y^{(i)}\big) \qquad L(y^{(i)}, f(x^{(i)}) = max(0, 1 - y^{(i)}f(x^{(i)}))$$

$$\frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial w_j} = \frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial f(x^{(i)})} \underbrace{\frac{\partial f(x^{(i)})}{\partial w_j}}_{x_j^i} \qquad \boxed{\begin{array}{l} f(x^n) \\ = w^T \cdot x^n \end{array}}$$

$$\frac{\partial max(0, 1 - y^{(i)}f(x^{(i)}))}{\partial f(x^{(i)})} = \begin{cases} -y^{(i)} & \text{If } y^{(i)}f(x^{(i)}){<}1 \\ \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial \sum_i L(y^{(i)}, f(x^{(i)})}{\partial w_j} = \sum_i \underbrace{-\delta\big( y^{(i)}f(x^{(i)}) < 1\big)\, y^{(i)}}_{c^n(W)} x_j^i \qquad \boxed{w_j \leftarrow w_j - \eta \sum_i c^n(W)\, x_j^i}$$

Minimizing total loss function L:

$$min \sum_i \left( max(0, 1 - y^{(i)} f(x^{(i)})) \right) + \lambda \|w\|_2$$

Minimizing total loss function L:

$$min \sum_i (max(0, 1 - y^{(i)}f(x^{(i)}))) + \lambda\|w\|_2$$

$$\varepsilon^i = max(0, 1 - y^{(i)}f(x^{(i)}))$$

$\varepsilon^i$: slack variable

Minimizing total loss function L:

$$min \sum_i \varepsilon^i + \lambda\|w\|_2$$

$$\varepsilon^i = max(0, 1 - y^{(i)}f(x^{(i)}))$$

$\varepsilon^i$: slack variable

Minimizing total loss function L:

$$min \sum_i \varepsilon^i + \lambda\|w\|_2$$

$$\varepsilon^i = max(0, 1 - y^{(i)}f(x^{(i)}))$$

$\varepsilon^i$: slack variable

$$\varepsilon^i \geq 0$$
$$\varepsilon^i \geq 1 - y^{(i)}f(x^{(i)})$$

$$y^{(i)}f(x^{(i)}) \geq 1 - \varepsilon^i$$

Minimizing total loss function L:

$$min \sum_i \varepsilon^i + \lambda\|w\|_2$$

$\varepsilon^i = max(0, 1 - y^{(i)}f(x^{(i)}))$

$y^{(i)}(w^T x_j + b) \geq 1 - \varepsilon^i$

$\varepsilon^i$: slack variable

$$\varepsilon^i \geq 0$$
$$\varepsilon^i \geq 1 - y^{(i)}f(x^{(i)})$$

$y^{(i)}f(x^{(i)}) \geq 1 - \varepsilon^i$

Minimizing total loss function L:

$$min \sum_i \varepsilon^i + \lambda \|w\|_2$$

$$s.t. \ y^{(i)}(w^T x_j + b) \geq 1 - \varepsilon^i$$

$\varepsilon^i$: slack variable

Minimizing total loss function L:

$$min \sum_i \varepsilon^i + \lambda \|w\|_2$$

$$s.t. \ y^{(i)}(w^T x_j + b) \geq 1 - \varepsilon^i$$

$\varepsilon^i$: slack variable

$$min \frac{1}{2}\|w\|_2 + C \sum_i \varepsilon^i$$

$$s.t. \ y^{(i)}(w^T x_j + b) \geq 1 - \varepsilon^i$$

# 支持向量 Support vectors

$$w^{(*)} = \sum_n a_n^* x^n$$

$a^{(*)}$ may be sparse

Linear combination of data points

$x^{(i)}$ with non-zero $a^{(*)}$ are support vectors

Ignore regularization for simplicity

$$\sum_i L\big(f(x^{(i)}), y^{(i)}\big) \qquad L(y^{(i)}, f(x^{(i)}) = max(0, 1 - y^{(i)} f(x^{(i)}))$$

$$\frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial w_j} = \frac{\partial L(f(x^{(i)}), y^{(i)})}{\partial f(x^{(i)})} \underset{x_j^i}{\underline{\frac{\partial f(x^{(i)})}{\partial w_j}}} \qquad \boxed{\begin{array}{c} f(x^n) \\ = w^T \cdot x^n \end{array}}$$

$$\frac{\partial max(0, 1 - y^{(i)} f(x^{(i)}))}{\partial f(x^{(i)})} = \begin{cases} -y^{(i)} & \text{If } y^{(i)} f(x^{(i)}) < 1 \\ \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial \sum_i L(y^{(i)}, f(x^{(i)}))}{\partial w_j} = \sum_i -\underline{\delta\big( y^{(i)} f(x^{(i)}) < 1\big) y^{(i)}} x_j^i$$

$$c^n(W) \qquad\qquad \boxed{w_i \leftarrow w_i - \eta \sum_i c^n(W) x_j^i}$$

$$w^{(*)} = \sum_n a_n^* x^n$$

$a^{(*)}$ may be sparse

➡️ Linear combination of data points

$x^{(i)}$ with non-zero $a^{(*)}$ are support vectors

$$w_1 = w_1 - \sum_i c^n(w) \, x_1^n$$

$$w_i = w_i - \sum_i c^n(w) \, x_i^n$$

$$w_k = w_k - \sum_i c^n(w) \, x_k^n$$

$$w^{(*)} = \sum_n a_n^* x^n$$

$a^{(*)}$ may be sparse

➡️ Linear combination of data points

$x^{(i)}$ with non-zero $a^{(*)}$ are support vectors

$$w_1 = w_1 - \sum_i c^n(w)\, x_1^n$$

$$w_i = w_i - \sum_i c^n(w)\, x_i^n$$

$$w_k = w_k - \sum_i c^n(w)\, x_k^n$$

$$c^n(w) = \frac{\partial\, L\big(y^{(i)} f(x^{(i)})\big)}{\partial f(x^{(i)})}$$

Hinge loss: usually zero

$$w^{(*)} = \sum_n a_n^* x^n$$

$a^{(*)}$ may be sparse

Linear combination of data points

$x^{(i)}$ with non-zero $a^{(*)}$ are support vectors

$$w_1 = w_1 - \sum_i c^n(w) \, x_1^n$$

$$w_i = w_i - \sum_i c^n(w) \, x_i^n$$

$$w_k = w_k - \sum_i c^n(w) \, x_k^n$$

If w initialized as 0

$$c^n(w) = \frac{\partial \, L(y^{(i)} f(x^{(i)}))}{\partial f(x^{(i)})}$$

Hinge loss: usually zero

$$w^{(*)} = \sum_n a_n^* x^n$$

$a^{(*)}$ may be sparse

Linear combination of data points

$x^{(i)}$ with non-zero $a^{(*)}$ are support vectors

$$w_1 = w_1 - \sum_i c^n(w) x_1^n$$

$$w_i = w_i - \sum_i c^n(w) x_i^n$$

$$w_k = w_k - \sum_i c^n(w) x_k^n$$

If w initialized as 0

c.f. for logistic regression, it  is always non-zero

$$c^n(w) = \frac{\partial L(y^{(i)} f(x^{(i)}))}{\partial f(x^{(i)})}$$

Hinge loss: usually zero

$$J(\theta) = \frac{1}{N}\sum_{i=1}^{N}\left[\begin{array}{c} -y^{(i)}log\left((f_\theta(x^{(i)}))\right) \\ -(1-y^{(i)})\,log\left(1-(f_\theta(x^{(i)}))\right) \end{array}\right]$$

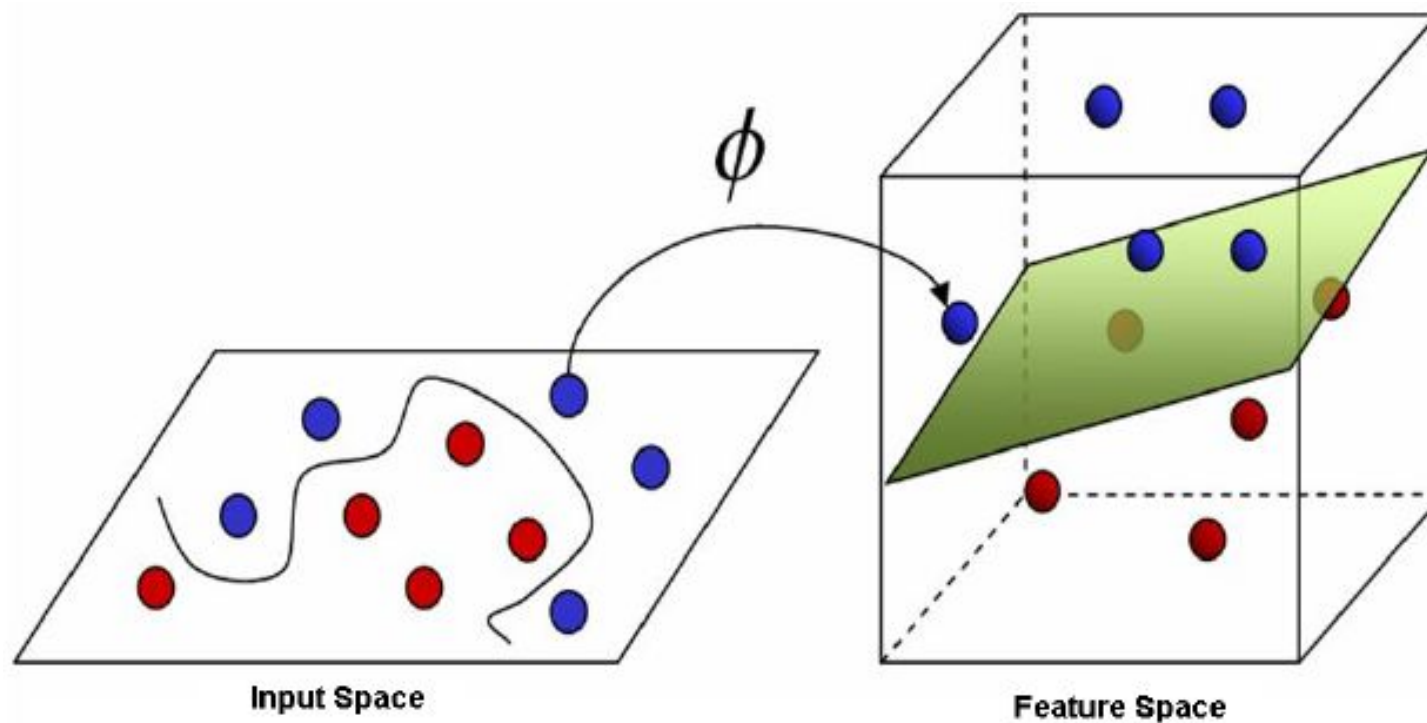Want $\left\{\begin{array}{l}\min_\theta J(\theta)\end{array}\right.$ :

Repeat

$$\theta_j := \theta_j - a\frac{1}{N}\sum_{i=1}^{N}(f_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$

(simultaneously update all $\theta_j$ )

}

Algorithm looks identical to linear regression!

# 核技巧 Kernel Trick



Input Space

$\phi$

Feature Space

$$w = \sum_n a_n x^n = Xa$$

$$X = \begin{bmatrix} x^1 & x^2 & \cdots\cdots & x^N \end{bmatrix} \qquad \boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix}$$

$$w = \sum_n a_n x^n = Xa$$

$$X = \begin{bmatrix} x^1 & x^2 & \cdots\cdots & x^N \end{bmatrix} \qquad \boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix}$$

Step 1: $\quad f(x) = w^T x$

$$w = \sum_n a_n x^n = Xa$$

$$X = \boxed{\begin{array}{cccc} x^1 & x^2 & \cdots\cdots & x^N \end{array}} \qquad \boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix}$$

$$w = X\boldsymbol{\alpha}$$

Step 1: $\quad f(x) = w^T x \quad \Longrightarrow \quad f(x) = \boldsymbol{\alpha}^T X^T x$

# 核技巧 Kernel Trick

$$w = \sum_n a_n x^n = Xa$$

$$X = \begin{bmatrix} x^1 & x^2 & \cdots\cdots & x^N \end{bmatrix} \quad \boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix}$$

$$w = X\boldsymbol{\alpha}$$

Step 1: $f(x) = w^T x \quad\longrightarrow\quad f(x) = \boldsymbol{\alpha}^T X^T x$

$$\begin{bmatrix} \alpha_1 & \cdots & \alpha_N \end{bmatrix}$$

$$\begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^N \end{bmatrix} \quad x$$

$$f(x) = \sum_n a_n (x^n . x)$$

$$\begin{bmatrix} x^1 \cdot x \\ x^2 \cdot x \\ \vdots \\ x^N \cdot x \end{bmatrix}$$

$$= \sum_n a_n K(x^n . x)$$

Step 1: $\quad f(x) \quad = \sum_n a_n K(x^n . x)$ Find $\quad a_1^*, a_2^*, \ldots a_n^*,$

Step 2, 3: Find $a_1^*, \ldots, a_n^*, \ldots, a_N^*$, minimizing loss function L

$$L(f) = \sum_i L\big(f(x^{(i)}), y^{(i)}\big)$$

$$= \sum_i L\left(\sum_n a_n K(x^n . x), y^{(i)}\right)$$

We only need to know the inner project between a pair of vectors x and z

Kernel Trick

- Liner kernel
$$K(x,z) = x^T z$$

- Polynomial kernel
$$K(x,z) = (x^T z)^d$$

- Gaussian kernel / Radial Basis Function Kernel
$$K(x,z) = \exp(-\frac{\|x-z\|^2}{2\sigma^2})$$

- Sigmoid kernel
$$K(x,z) = \tanh(\beta x^T z + \theta)$$

# RBF核 Radial Basis Function Kernel

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

$\boxed{\phi(*) \text{ has inf dim!!!}}$

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

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

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

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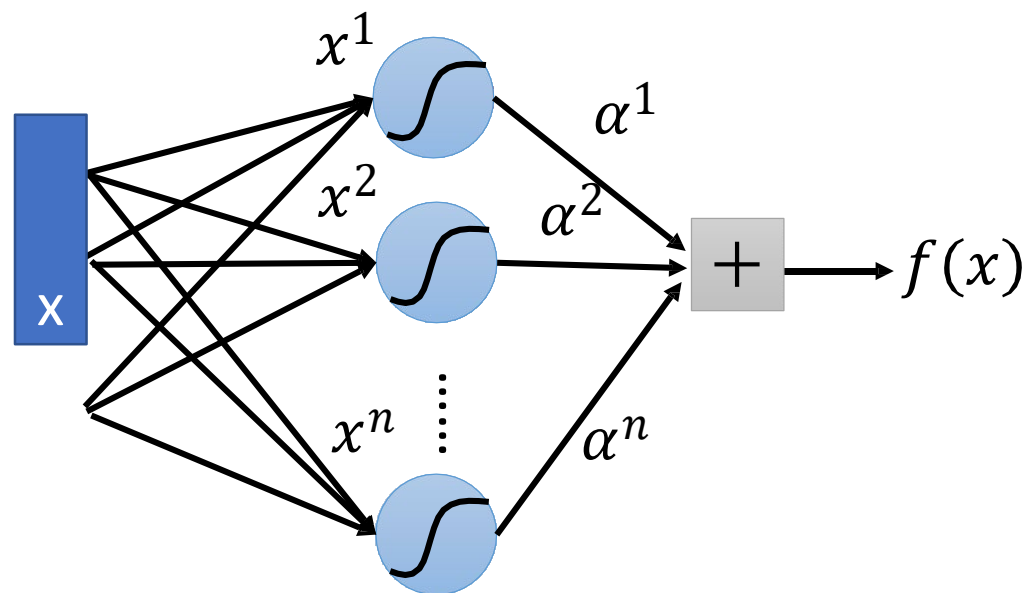
# Sigmoid核 Sigmoid Kernel

$$K(x,z) = tanh(x \cdot z)$$

- When using sigmoid kernel, we have a 1 hidden layer network.
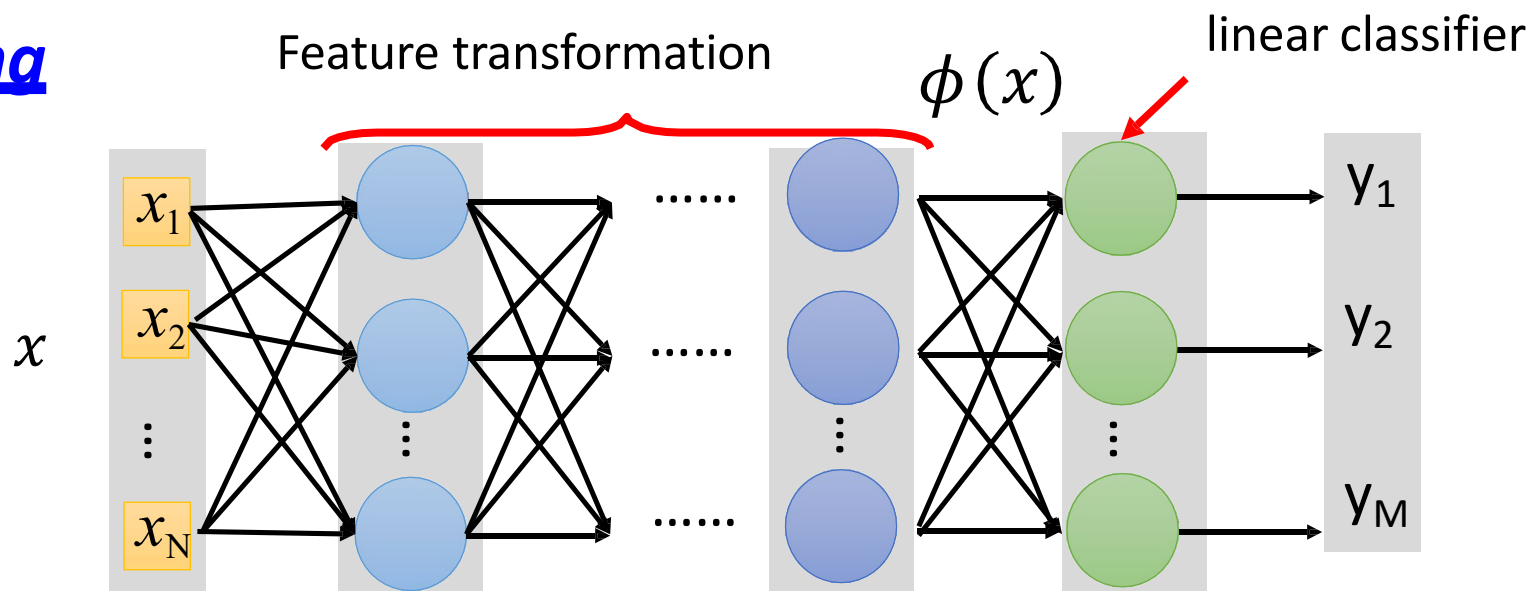
$$f(x) = \sum_n a_n K(x^n.x) = \sum_n a_n tanh(x^n.x)$$

The weight of each neuron is a data point
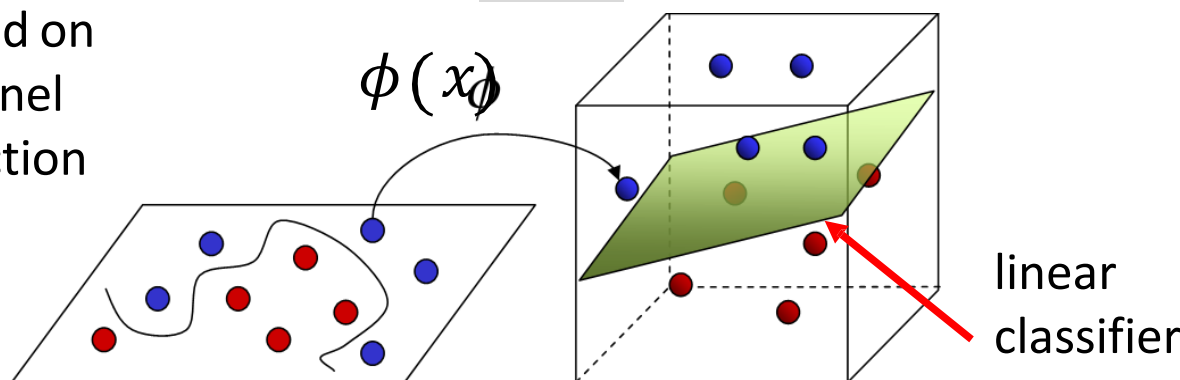
The number of support vectors is the number of neurons.

# 深度学习和支持向量机 Deep learning VS SVM

**Deep Learning**

Feature transformation

$\phi(x)$

linear classifier

$x_1$

$x$

$x_2$

$x_N$

$y_1$

$y_2$

$y_M$

**SVM**

Based on kernel function

$\phi(x)$

linear classifier

**Input Space**

**Feature Space**

Multiple Kernel learning [Alpaydin, Chapter 13.8]

# SVM 软件包

- LIBSVM
http://www.csie.ntu.edu.tw/~cjlin/libsvm/

- LIBLINEAR
http://www.csie.ntu.edu.tw/~cjlin/liblinear/

- SVM$^{light}$ 、 SVM$^{perf}$ 、 SVM$^{struct}$
http://svmlight.joachims.org/svm_struct.html

- Pegasos
http://www.cs.huji.ac.il/~shais/code/index.html


- Demo

- Support Vector Machine (dash.gallery)

- Support Vector Regression (SVR)
  - [Bishop chapter 7.1.4]

- Ranking SVM
  - [Alpaydin, Chapter 13.11]

- One-class SVM
  - [Alpaydin, Chapter 13.11]

- [Support Vector Machine (dash.gallery)](#)