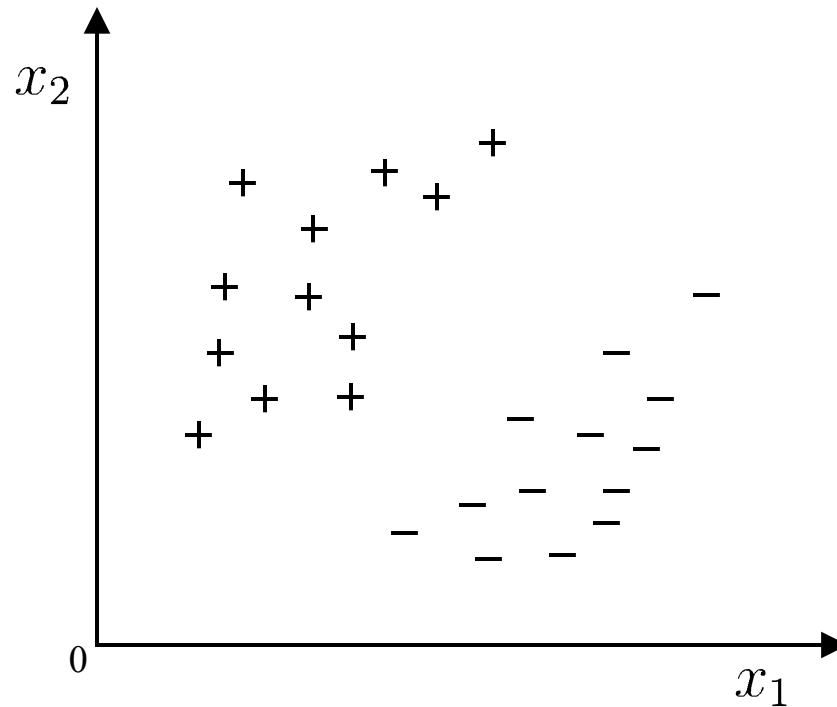# Machine Learning
# 机器学习

# Lecture8：支持向量机

李洁

nijanice@163.com

# 二分类问题
## Binary Classification



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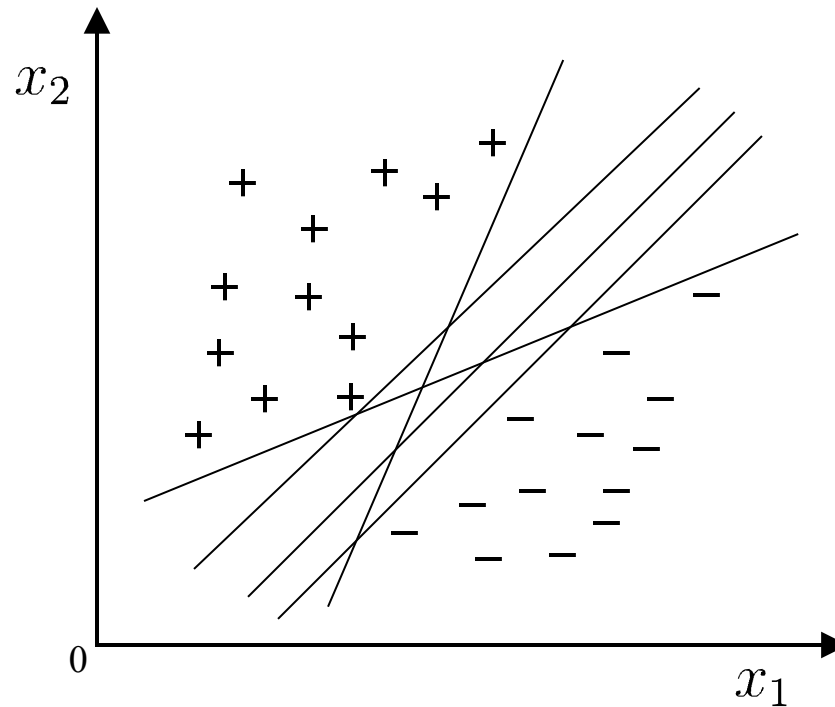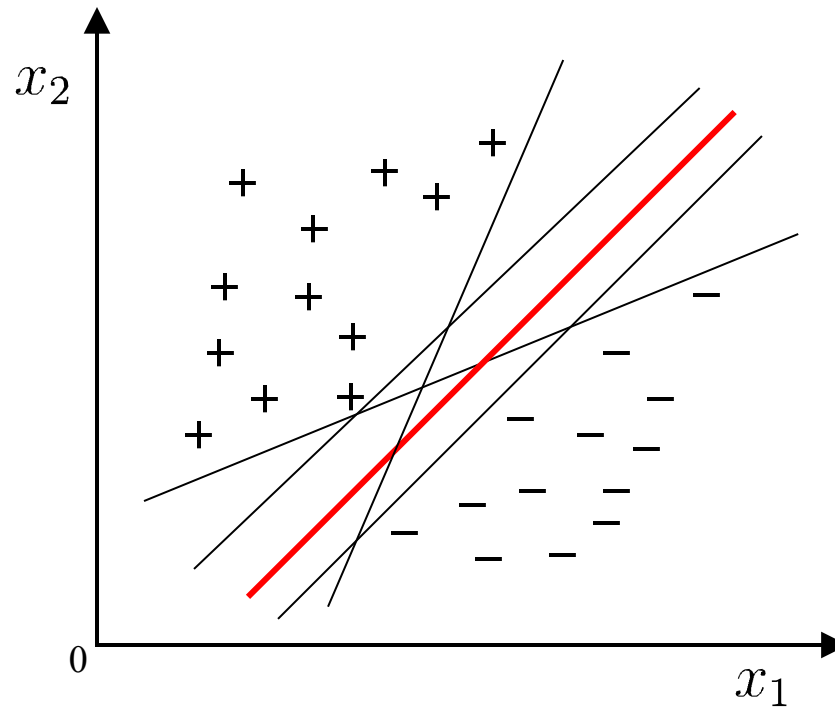
# 二分类问题
## Binary Classification



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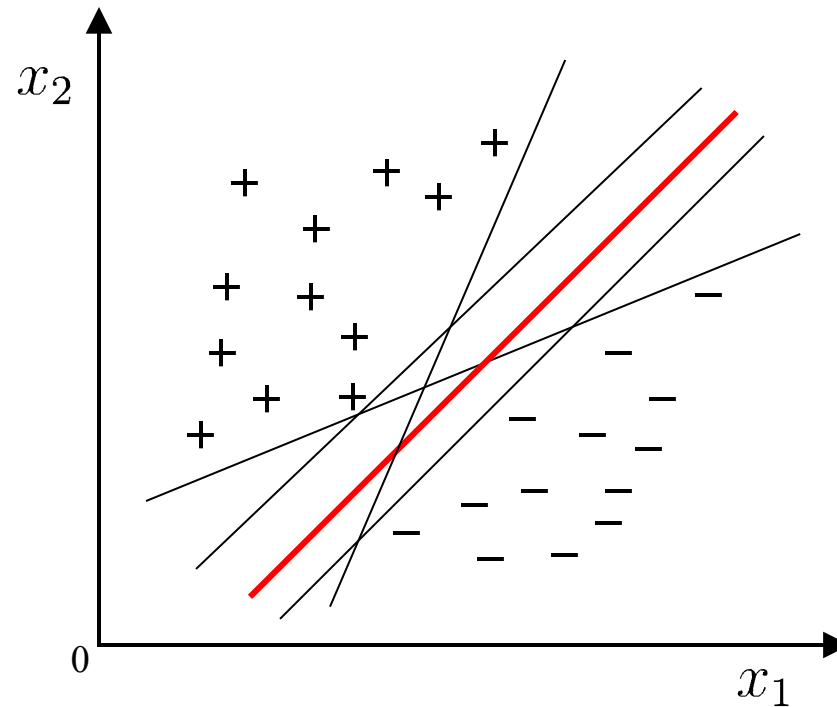
# 二分类问题
## Binary Classification



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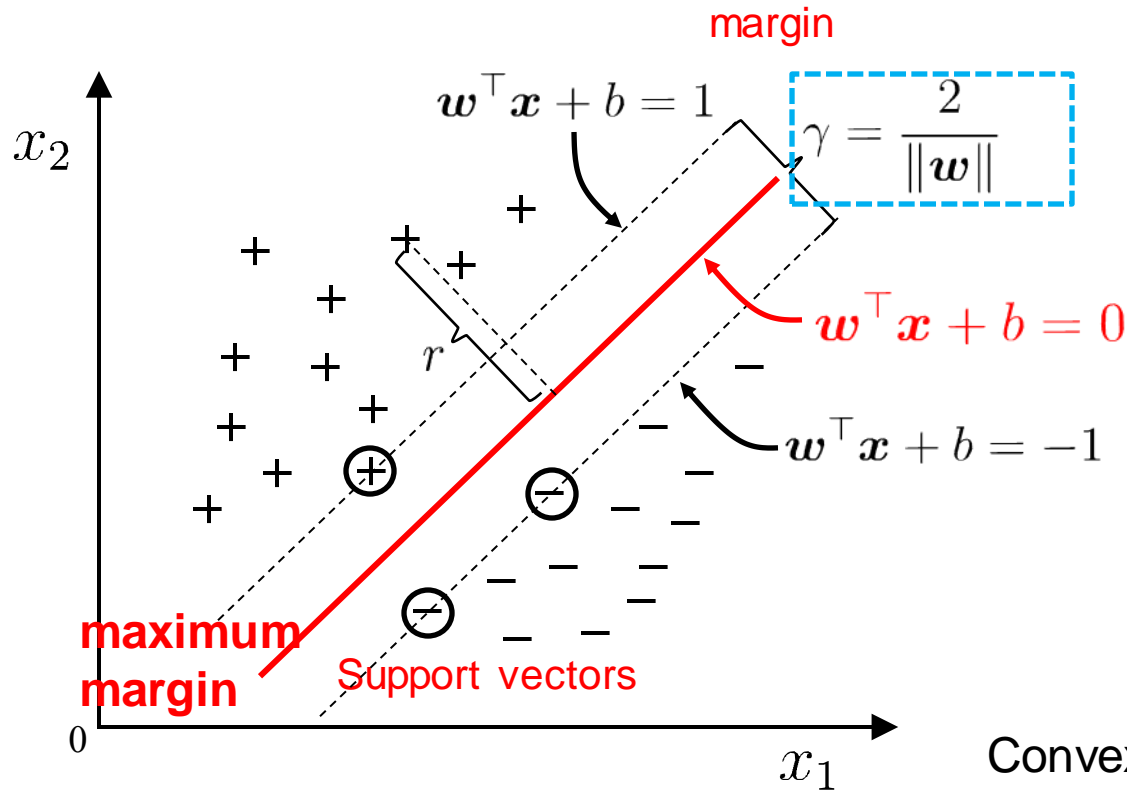
# 二分类问题
## Binary Classification



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

margin

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

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

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

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

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

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

$r$

$x_2$

$\oplus$

$\ominus$

$\ominus$

$-$

**maximum margin**

Support vectors

$0$

$x_1$

Convex Quadratic Programming

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

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

# 对偶问题
## Dual problem

A dual problem, in mathematics and optimization, is a counterpart to the primal optimization problem. It is constructed based on the constraints of the primal problem and aims to find either an upper bound (in case of minimization problems) or a lower bound (in case of maximization problems) on the optimal value of the primal problem.

For a minimization problem:

- Dual Problem: max $g(\lambda,v)$ where $g(\lambda,v)=\inf_x L(x,\lambda,v)$

For a maximization problem:

- Dual Problem: $\min g(\lambda,v)$ where $g(\lambda,v)=\sup_x L(x,\lambda,v)$

  Here:

- $\lambda$ and $v$ are dual variables (Lagrange Multipliers) corresponding to the constraints of the primal problem.

- $L(x,\lambda,v)$ is the Lagrangian function, which is composed of the primal objective function and the constraints, typically in the form

  $L(x,\lambda,v)=$objective function$-\lambda^T\cdot$(inequality constraints)$-v^T\cdot$(equality constraints)

By solving the dual problem, we obtain either an upper bound or a lower bound on the optimal value of the primal problem.

# 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
  $g(x,y,z)=x+y+z-10\leq 0$。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。

# 最优超平面求解
## Optimal hyperplane Solution

- 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

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

# 最优超平面求解
# Optimal hyperplane Solution

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

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

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

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

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

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

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

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

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

# 最优超平面求解
# Optimal hyperplane Solution

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

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

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

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

Dual problem:

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

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

KKT Condition

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

The optimal solution can be obtained

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

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

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

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

# 最优超平面求解
## Optimal hyperplane Solution

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

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

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

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

Dual problem:

$$\max_a \sum_{i=1}^{m} \alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m} \alpha_i\,\alpha_j y_i y$$

$$s.t. \sum$$

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

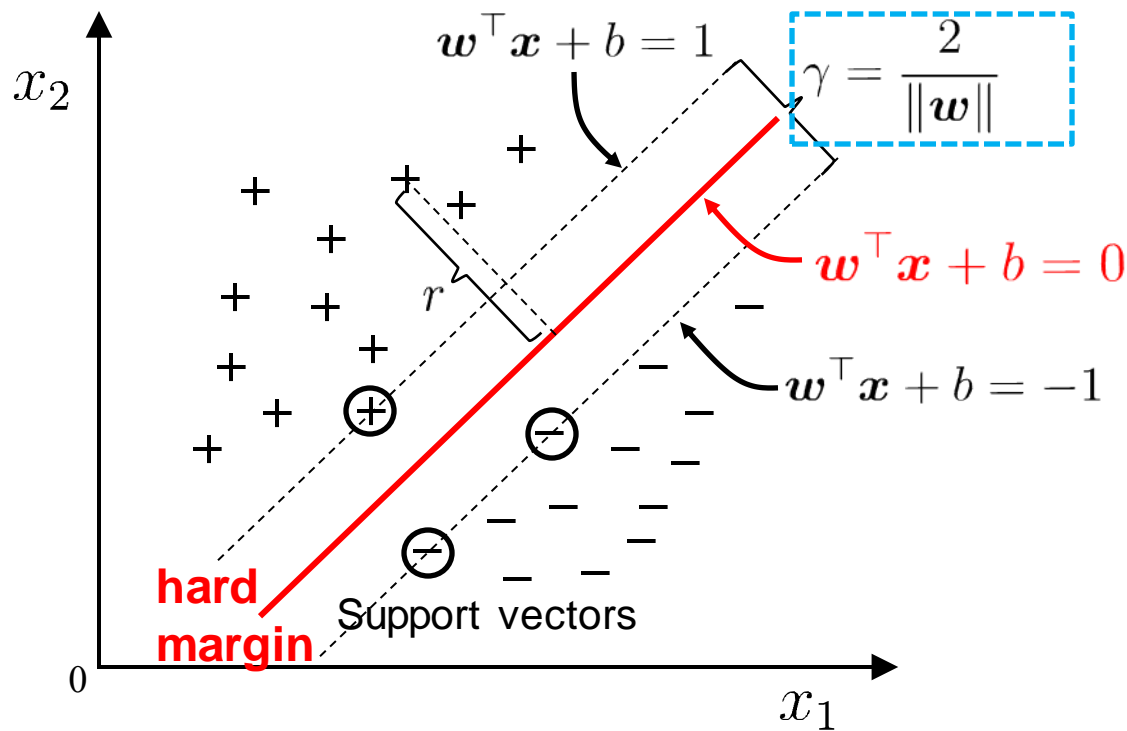The optimal solution can be obtained

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

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

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

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

# 支持向量机
# Support Vector Machine



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

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

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

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

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

$x_2$

$r$

**hard margin**

Support vectors

$0$

$x_1$

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

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

# 软间隔
# Soft Margin

-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

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

The slack variable $\varepsilon$ can be introduced :

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

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

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

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

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

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

The slack variable $\varepsilon$ can be introduced :

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

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

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

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

$$C = \alpha_i + \mu_i$$

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

The slack variable $\varepsilon$ can be introduced :

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

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

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

The optimal solution can be obtained

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

KKT Condition

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

# 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 $\varepsilon_i \leq 1$，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

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

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

KKT Condition

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

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

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

The slack variable $\varepsilon$ can be introduced :

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

Dual Problem:

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

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

# 最优超平面求解
# Optimal hyperplane Solution

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

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

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

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

Dual problem:

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

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

KKT Condition

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

The optimal solution can be obtained

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

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

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

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

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

The slack variable $\varepsilon$ can be introduced :

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

Dual Problem:

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

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

# 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。

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

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

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

The slack variable $\varepsilon$ can be introduced :

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

Dual Problem:

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

$$s.t. \sum_{i=1}^m \alpha_i y_i = 0 , C \geq \alpha_i \geq 0, i = 1,2,\dots,m$$
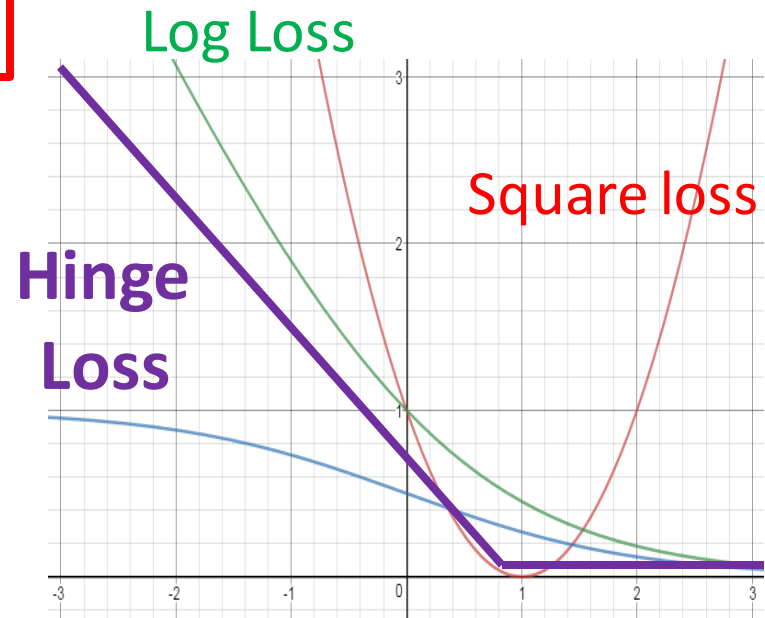
# 支持向量机
## Support Vector Machine

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

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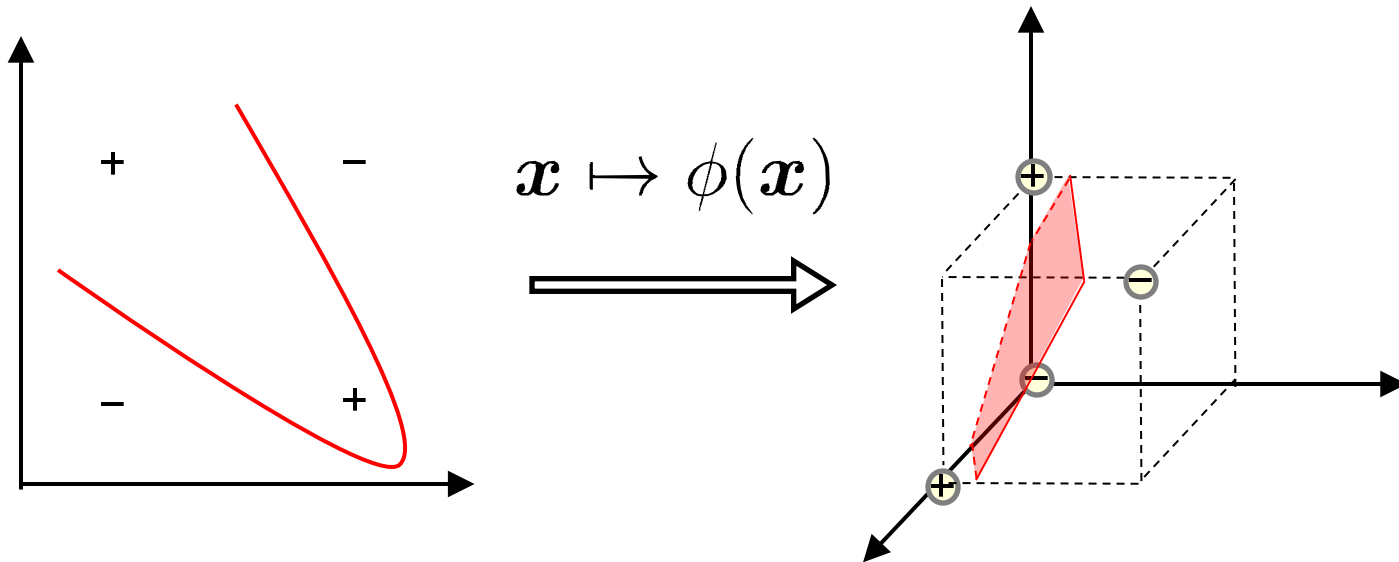
Hinge Loss function

Log Loss

Square loss

**Hinge Loss**

# 线性不可分
# linearly Inseparable Problem

-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



$$x \mapsto \phi(x)$$

# 最优超平面求解
# Optimal hyperplane Solution

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

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

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

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

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

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

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

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

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

# 核支持向量机
# Kernel SVM

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:

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

# 核函数
# Kernel Function

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

Original question:

<span style="color:red">Kernel Function</span>

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

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

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

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:

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

# 监督学习
# Supervised Learning

- 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,

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

let the machine learn a function from data to label

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

- 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?

# 学习目标
## Learning Objective

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

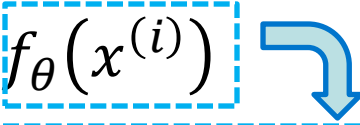Log loss: $-y^{(i)} log\left((f_\theta(x)\right) - (1 - y^{(i)}) log\left(1 - (f_\theta(x)\right)$

# 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

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

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

- 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

# 如何用于分类
# Classification task

$y \in \{0, 1\}$ 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    $\theta^T x \geq 0.5$    , predict "y = 1"

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

# 如何用于分类
# Classification task

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

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

# 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 $\theta$

# 逻辑斯蒂回归
# Logistic regression

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

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

= 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：

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

# 二分类
## Binary Classification

- Step 1: Function (Model)

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

# 二分类
## Binary Classification

- Step 1: Function (Model)

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

# 二分类
## Binary Classification

- Step 1: Function (Model)

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

- Step 2: Cost function

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

# 二分类
## Binary Classification

- Step 1: Function (Model)

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

- Step 2: Cost function

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

# 二分类
## Binary Classification

- Step 1: Function (Model)

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

- Step 2: Cost function

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

The number of times get incorrect results on training data.

# 二分类
# Binary Classification

- 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 3：Training by gradient descent is difficult

# 二分类
# Binary Classification

- Step 1: Function (Model)

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

- 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

# 二分类
# Binary Classification

- Step 1: Function (Model)

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

$\Longrightarrow$

$$y^{(i)} f\left(x^{(i)}\right) \geq 0$$
$$y^{(i)} f\left(x^{(i)}\right) < 0$$

- Step 2: Cost function

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

- Step 3： Training by gradient descent is difficult

# 支持向量机
# Support Vector Machine

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**

# 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

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

Sigmoid + Square loss

Larger value, smaller loss

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

# 损失函数
# Loss function

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

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

# 损失函数
# Loss function

Square Loss:



Sigmoid +
cross entropy

Square loss

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

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

Sigmoid + Square loss

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

Larger value, smaller loss

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

# 损失函数
# Loss function



Square Loss:

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

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

Sigmoid +

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

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

loss

Sigmoid +
Square loss

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

Larger value, smaller loss

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

# 损失函数
# Loss function

Cross entropy:

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

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

Sigmoid +
cross entropy

Square loss

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

Sigmoid +
Square loss

Divided by
ln2 here

Larger value, smaller loss

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

# 损失函数
# Loss function



Cross entropy:

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

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

Sigmoid +
cross entropy

Square loss

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

Sigmoid +
Square loss

Divided by
ln2 here

Larger value, smaller loss

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

# 损失函数
# Loss function



Sigmoid +
cross entropy

Square loss

Sigmoid +
Square loss

Divided by
ln2 here

3

2

1

-3  -2  -1  0  1  2  3

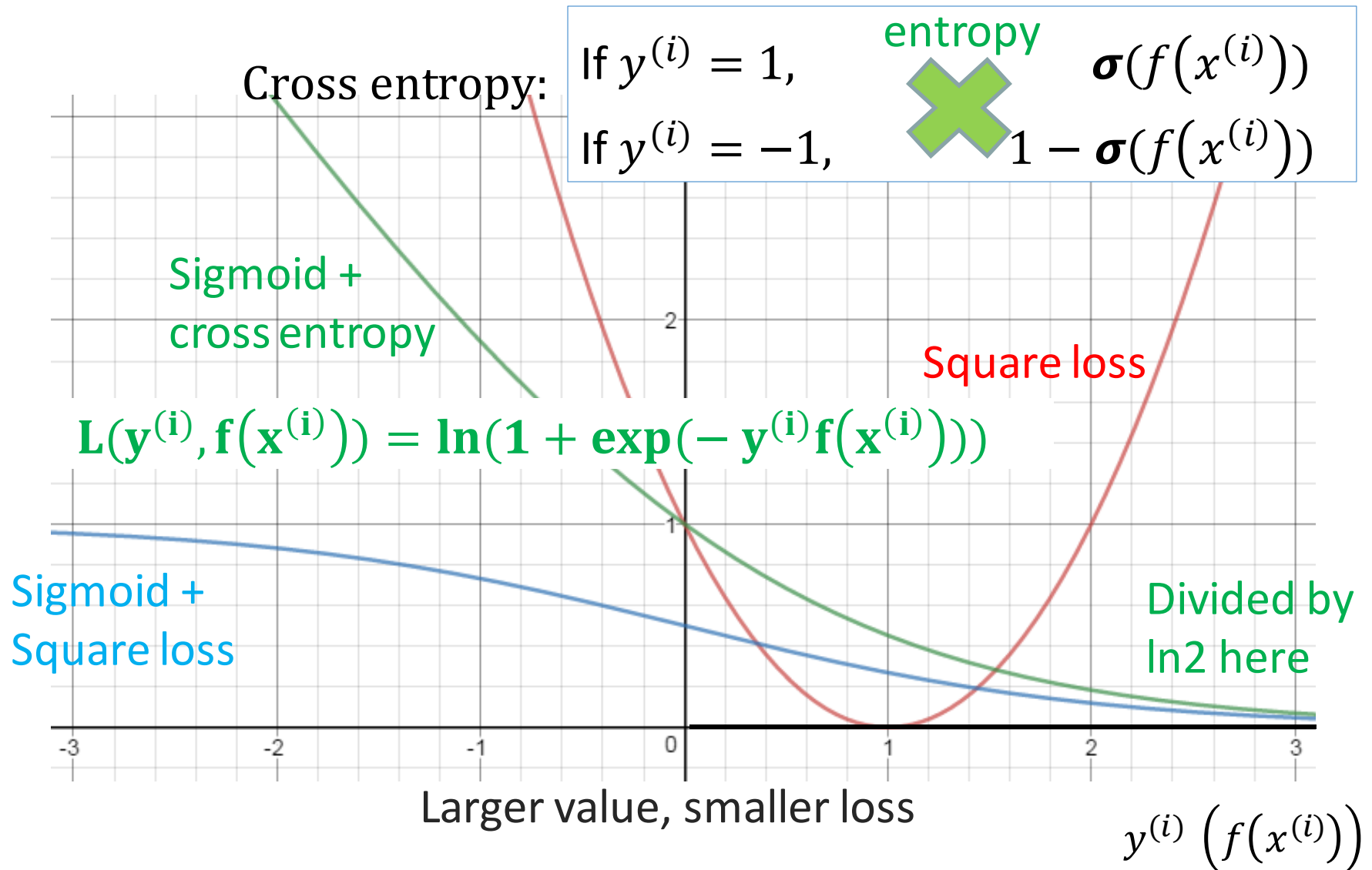Larger value, smaller loss

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

# 损失函数
# Loss function



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**

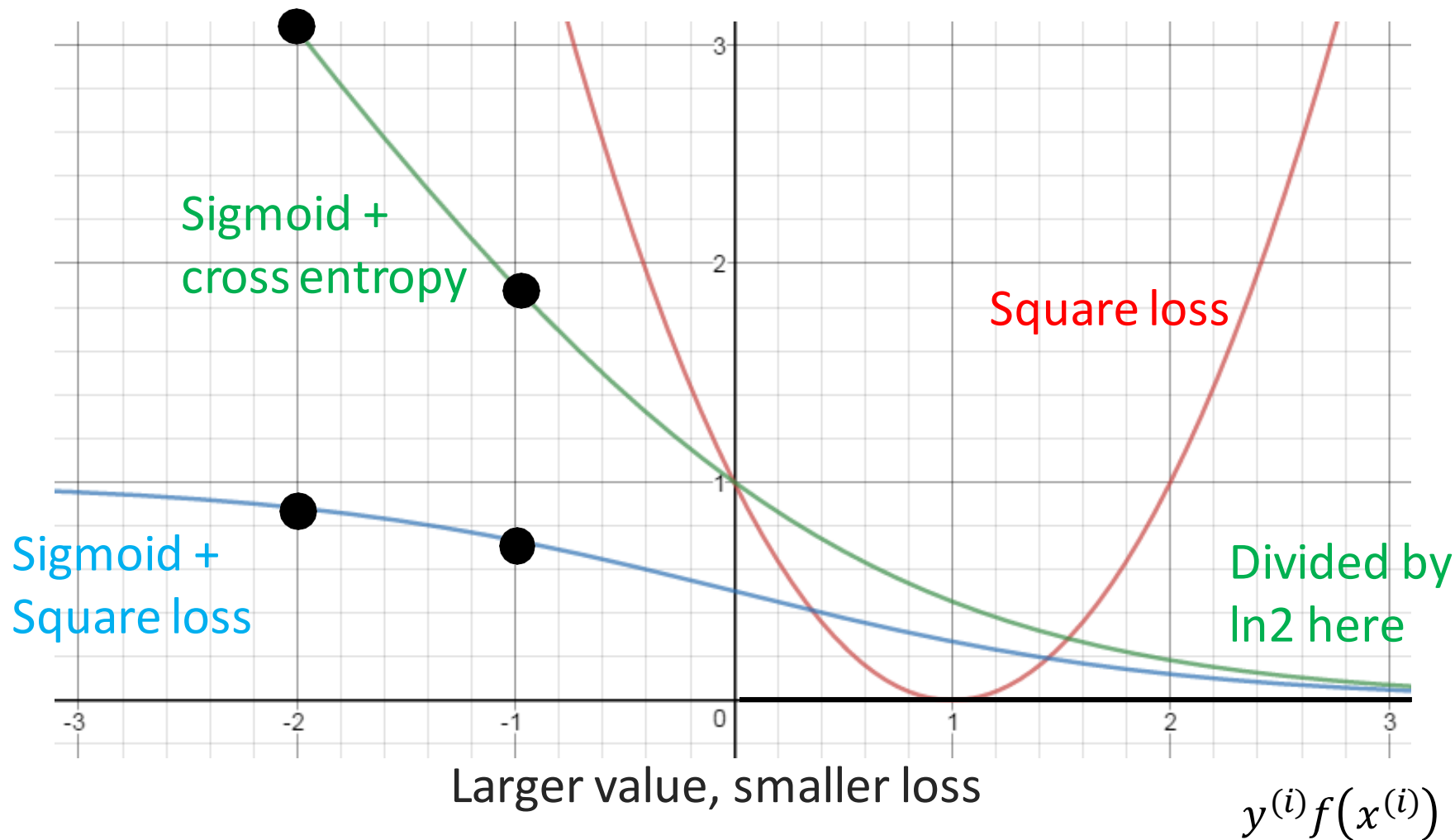$$\mathbf{L}(\mathbf{y}^{(i)}, \mathbf{f}(\mathbf{x}^{(i)}))$$
$$= \mathbf{max}(\mathbf{0},\ \mathbf{1} - \mathbf{y}^{(i)}\mathbf{f}(\mathbf{x}^{(i)}))$$

Sigmoid + cross entropy

Sigmoid + Square loss

Larger value, smaller loss

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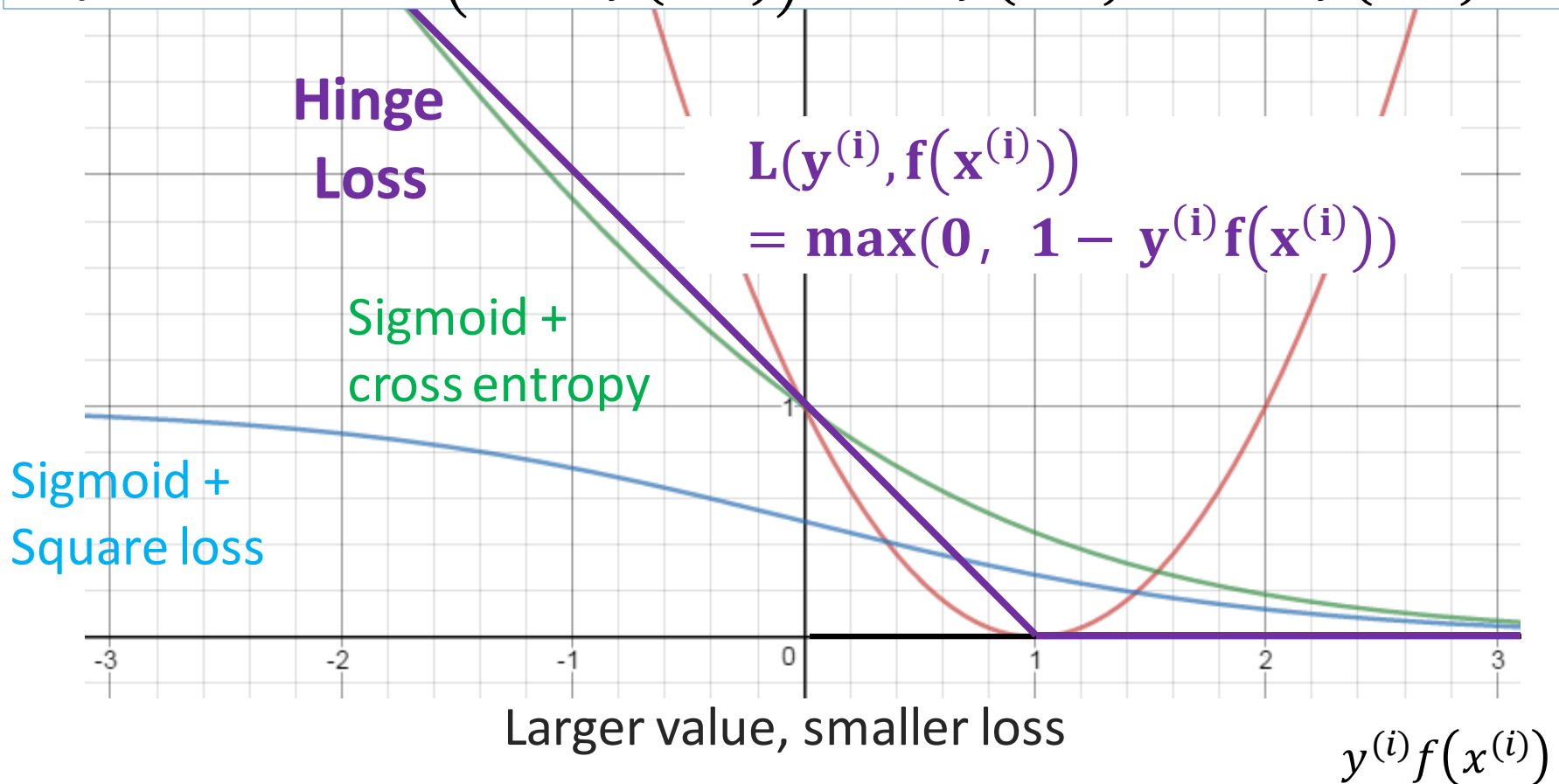
# 损失函数
# Loss function

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

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



**Hinge Loss**

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

Sigmoid + cross entropy

Sigmoid + Square loss

Larger value, smaller loss

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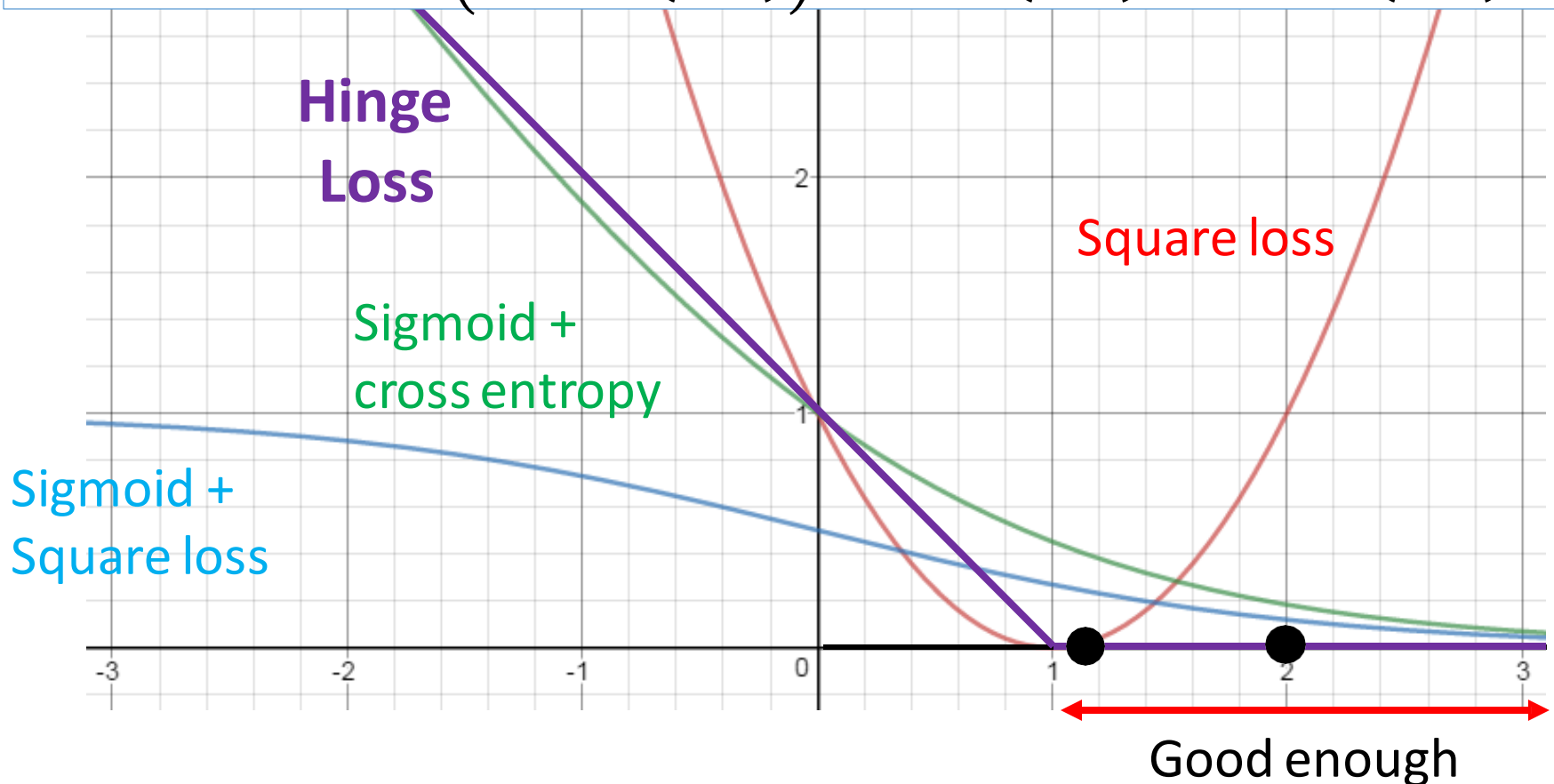
# 损失函数
# Loss function

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

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

**Hinge Loss**

Square loss

Sigmoid +
cross entropy

Sigmoid +
Square loss

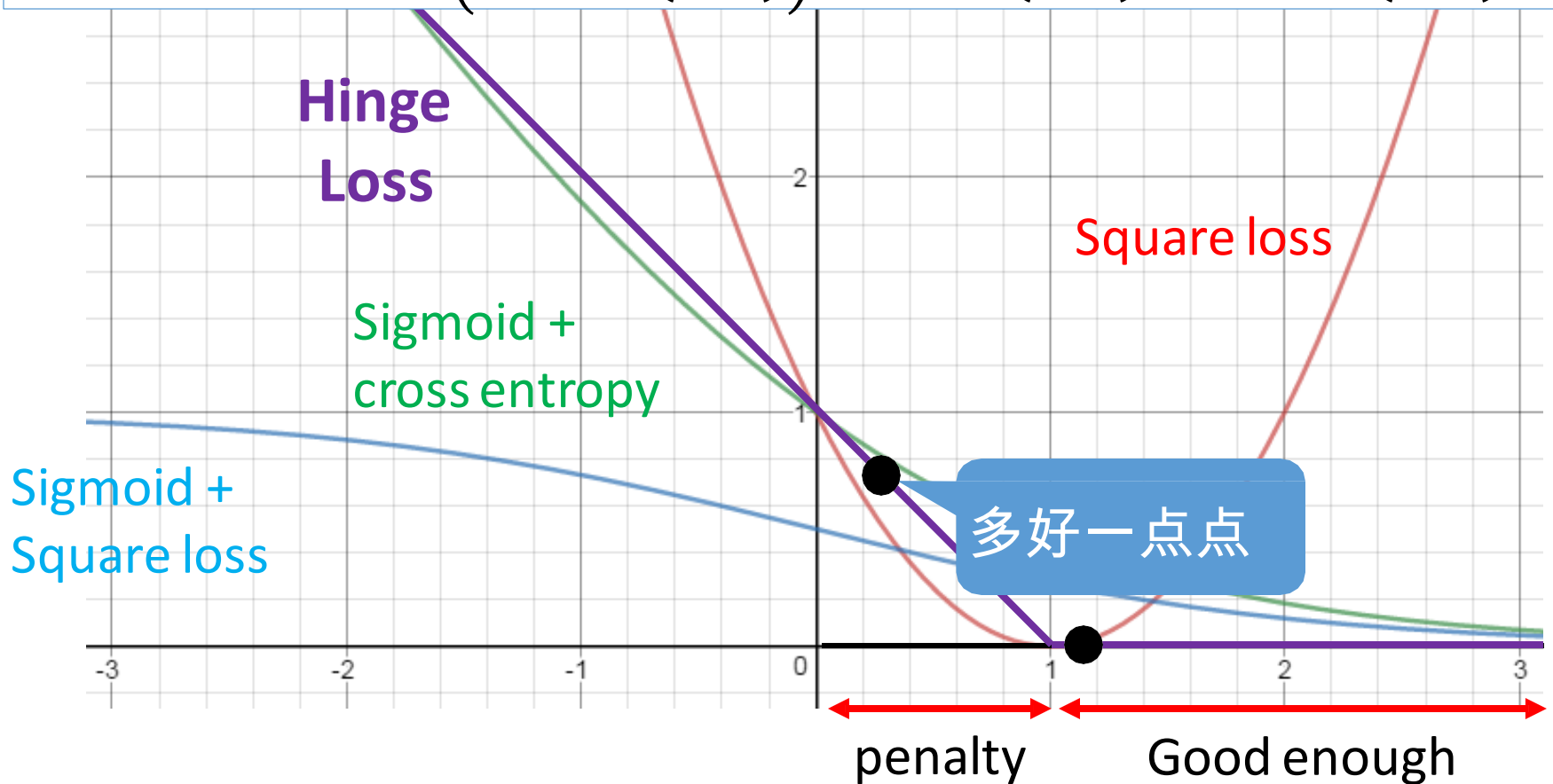Good enough

# 损失函数
# Loss function

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

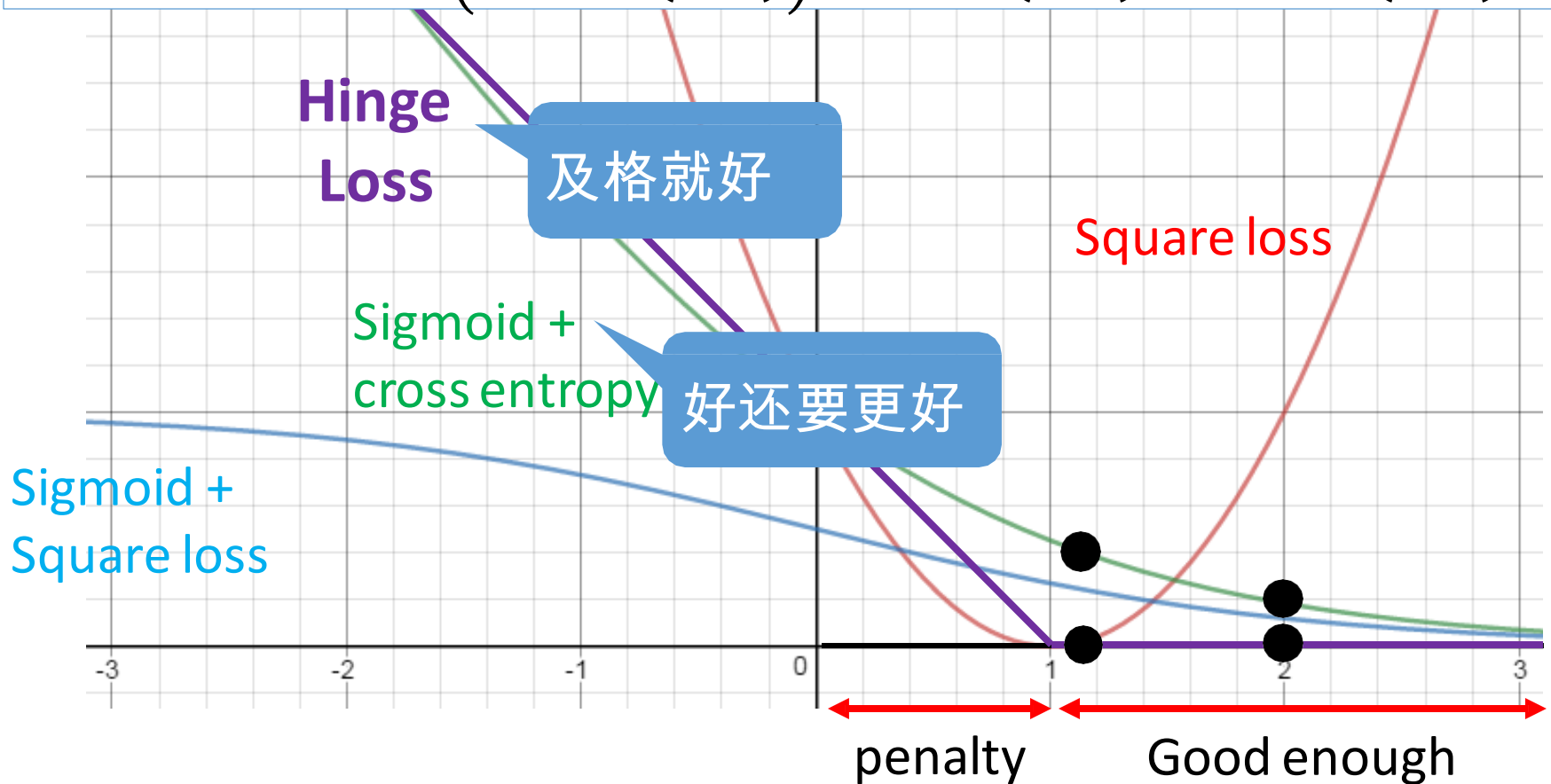If $y^{(i)} = -1$, $\max\left(0, 1 + f(x^{(i)})\right)$    $1 + f(x^{(i)}) < 0$    $f(x^{(i)}) < -1$



**Hinge Loss**

Square loss

Sigmoid + cross entropy

Sigmoid + Square loss

多好一点点

penalty      Good enough

# 损失函数
# Loss function

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

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**

及格就好

Square loss

Sigmoid + cross entropy

好还要更好

Sigmoid + Square loss

penalty

Good enough

# 二分类
## Binary Classification

- Step 1: Function (Model)

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

- Step 2: Cost function
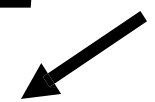
$$C(f) = \sum_n L\big(f\big(x^{(i)}\big), y^{(i)}\big)\big)$$

# 线性SVM
# Linear SVM

- 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

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

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

Hinge loss

# 线性SVM
# Linear SVM

- 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

$$L\big(y^{(i)}, f\big(x^{(i)}\big)\big) = max\big(0, 1 - y^{(i)} f\big(x^{(i)}\big)\big)$$
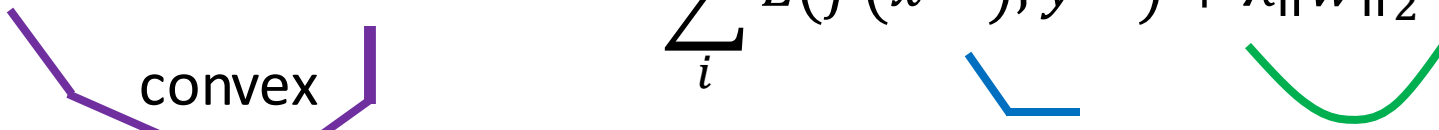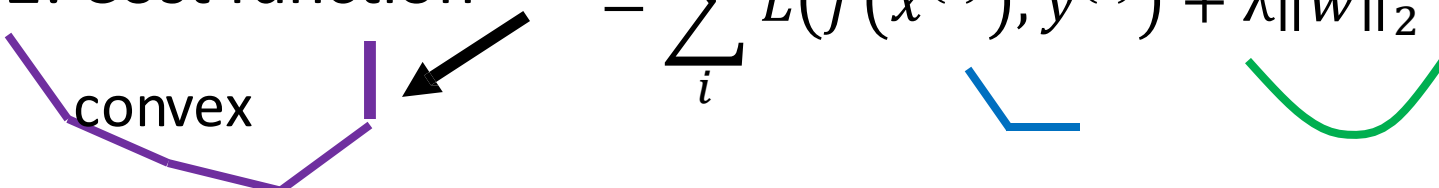
- Step 3: gradient descent?

# 线性SVM
# Linear SVM

- 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 ( _____ ) ?

# 线性SVM
# Linear SVM

- 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  (_____) ?

# 线性SVM
# Linear SVM

Ignore regularization for simplicity

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

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

# 线性SVM
# Linear SVM

Ignore regularization for simplicity

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

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

# 线性SVM
# Linear SVM

Ignore regularization for simplicity

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

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

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

$$x_j^i$$

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

# 线性SVM
# Linear SVM

Ignore regularization for simplicity

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

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

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

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

$$x_j^i$$

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

# 线性SVM
# Linear SVM

Ignore regularization for simplicity

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

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

$$x_j^i$$

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

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

# 线性SVM
# Linear SVM

Ignore regularization for simplicity

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

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

$$x_j^i$$

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

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

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

# 线性SVM
# Linear SVM

Minimizing total loss function L:

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

# 线性SVM
# Linear SVM

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

# 线性SVM
# Linear SVM

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

# 线性SVM
# Linear SVM

Minimizing total loss function L:

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

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

$\varepsilon^i$: slack variable

$$\boxed{\begin{aligned} \varepsilon^i &\geq 0 \\ \varepsilon^i &\geq 1 - y^{(i)}f(x^{(i)}) \end{aligned}}$$
➡ $y^{(i)}f(x^{(i)}) \geq 1 - \varepsilon^i$

# 线性SVM
# Linear SVM

Minimizing total loss function L:

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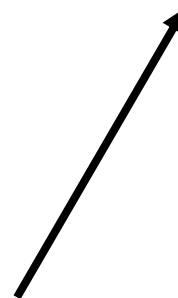
$$\varepsilon^i = max(0, 1 - y^{(i)}f(x^{(i)}))$$

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

$\varepsilon^i$: slack variable

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

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

# 线性SVM
# Linear SVM

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

# 线性SVM
# Linear SVM

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

# 线性SVM
# Linear SVM

Ignore regularization for simplicity

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

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

$$\underline{\phantom{xxxx}} \\ x_j^i$$

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

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

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

# 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

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

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

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

# 线性SVM
# Linear SVM

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

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

⮕ Linear combination of data points

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

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

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

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

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

Hinge loss: usually zero

# 线性SVM
# Linear SVM

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

$a^{(*)}$ 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

# 线性SVM
# Linear SVM

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

$a^{(*)}$ 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

# 逻辑斯蒂回归求解
# Logistic regression solution

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

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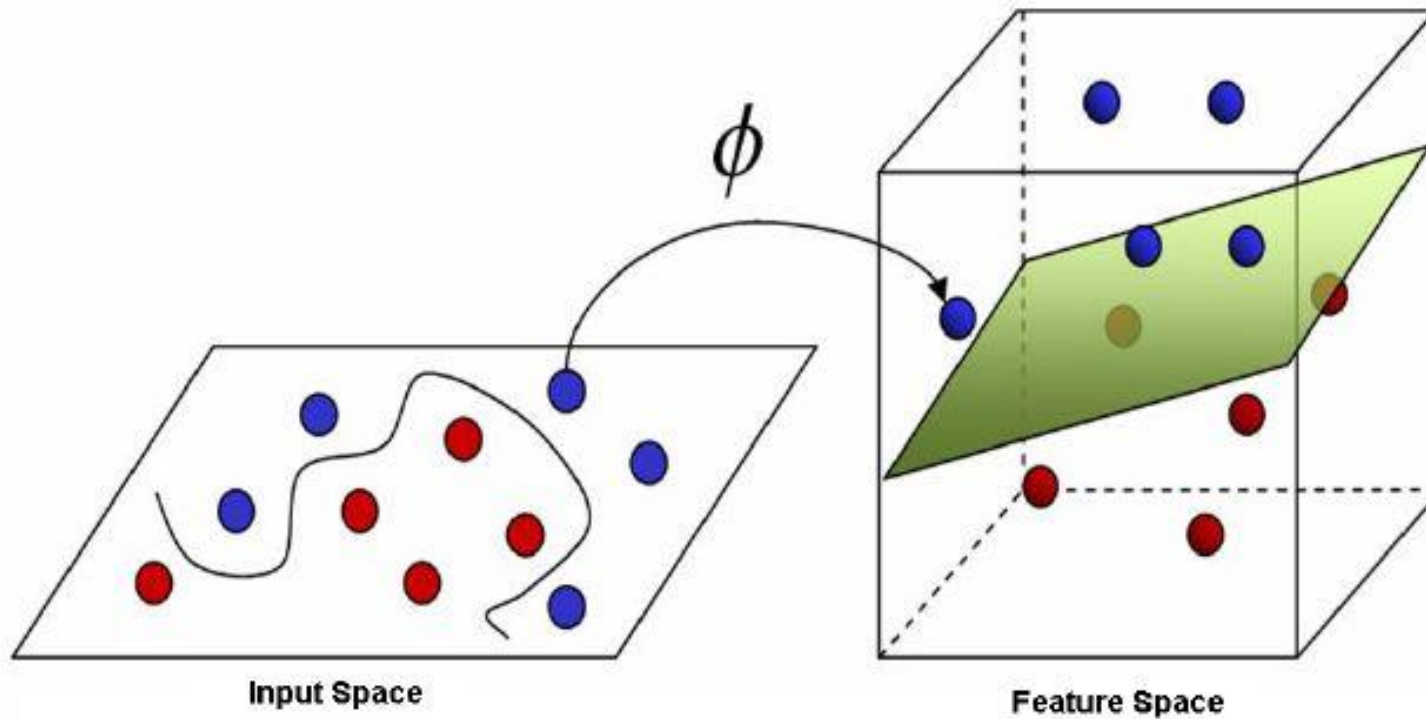
Repeat

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

(simultaneously update all $\theta_j$ )

$\}$

Algorithm looks identical to linear regression!

# 核技巧
# Kernel Trick

# 核技巧
## Kernel Trick

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

$$X = \begin{bmatrix} x^1 & x^2 & \cdots\cdots & x^N \end{bmatrix}$$

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix}$$

# 核技巧
# Kernel Trick

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

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

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

# 核技巧
# Kernel Trick

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

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

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

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

# 核技巧
# Kernel Trick

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

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

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

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

$$[\alpha_1 \quad \cdots \quad \alpha_N]$$

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

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

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

$$\begin{array}{|c|} \hline x^1 \\ \hline x^2 \\ \hline \vdots \\ \hline x^N \\ \hline \end{array} \qquad \boxed{x}$$

# 核技巧
# Kernel Trick

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

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

# 核技巧
# Kernel Trick

Directly computing $K(x, z)$ can be faster than "feature transformation + inner product" sometimes.
Kernel trick is useful when we transform all x to $\phi(x)$

$$K(x, z) = \phi(x) \cdot \phi(z) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$= x_1^2 z_1^2 + 2x_1x_2z_1z_2 + x_2^2 z_2^2$$

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

$$= (x_1z_1 + x_2z_2)^2 = \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \right)^2$$

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

# 核技巧
# Kernel Trick

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

$$= (x_1 z_1 + x_2 z_2 + \cdots + x_k z_k)^2$$

$$= x_1{}^2 z_1{}^2 + x_2{}^2 z_2{}^2 + \cdots + x_k{}^2 z_k{}^2$$

$$+ 2 x_1 x_2 z_1 z_2 + 2 x_1 x_3 z_1 z_3 + \cdots$$

$$+ 2 x_2 x_3 z_2 z_3 + 2 x_2 x_4 z_2 z_4 + \cdots$$

$$= \phi(x) \cdot \phi(z)$$

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

$$\phi(x) = \begin{bmatrix} x_1{}^2 \\ \vdots \\ x_k{}^2 \\ \sqrt{2} x_1 x_2 \\ \sqrt{2} x_1 x_3 \\ \vdots \\ \sqrt{2} x_2 x_3 \\ \vdots \end{bmatrix}$$

# 常用核函数
# Kernel Function

- 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

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

$\phi(*)$ has inf dim!!!

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

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

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

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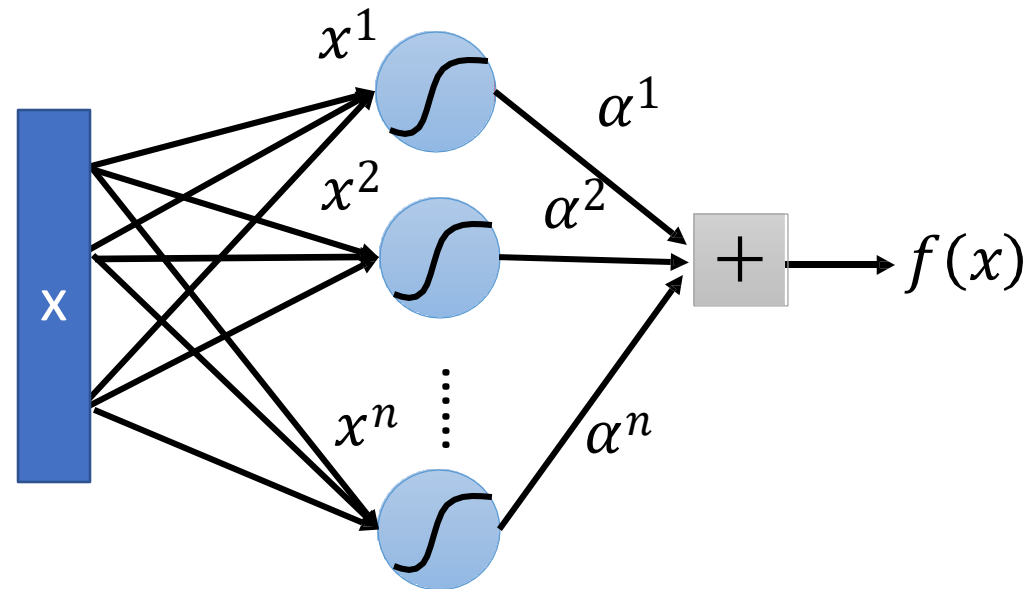
# Sigmoid核
# Sigmoid Kernel

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

- 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

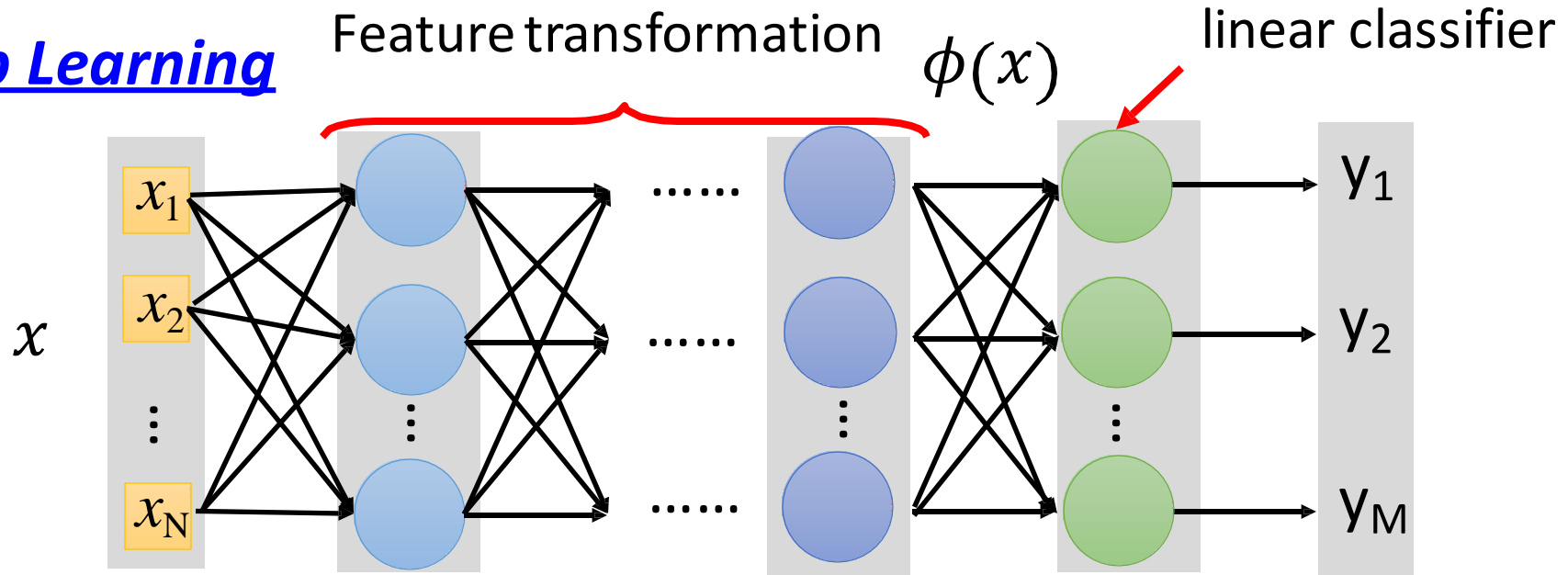The number of support vectors is the number of neurons.

# 深度学习和支持向量机
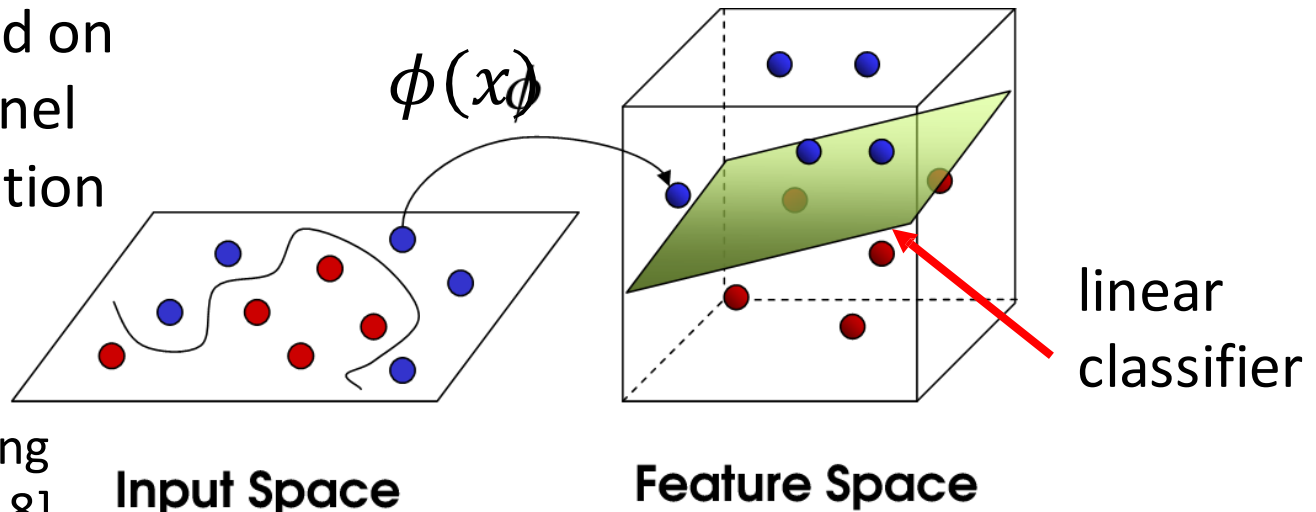# Deep learning VS SVM



**Deep Learning**

Feature transformation $\phi(x)$ linear classifier

$x$

$x_1$, $x_2$, $\vdots$, $x_N$

...... ......

$y_1$, $y_2$, $y_M$

**SVM**

Based on kernel function

$\phi(x)$

linear classifier

Multiple Kernel learning
[Alpaydin, Chapter 13.8]

**Input Space**　　**Feature Space**

# 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)

# SVM 方法
# SVM related methods

- 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)](dash.gallery)

简述一下本节课介绍的SVM的两个最重要的特点以及带来的好处？