

# Hinge Loss for Binary Classifiers

专门处理二分类

核心思想：不仅要正确分类，还要确保分类结果有足够的“置信度”

↓  
(分类边界有足够的间隔)

# Objectives

- introduce disadvantage of squared error for classification
- introduce hinge loss cost function
- characteristics of hinge loss

Squared error "loss" can be problematic 2

Classifier design

$$\min_{\underline{w}} \ell(\underline{w}; \underline{A}, \underline{d}) + \lambda r(\underline{w}) \leftarrow \begin{array}{l} \text{regularizer} \\ \text{loss function} \end{array}$$

Squared error loss  $\ell(\underline{w}; \underline{A}, \underline{d}) = \|\underline{A}\underline{w} - \underline{d}\|_2^2$

Example: dwarf planet vs. planet

| object                          | Ceres | Eris | Pluto | Mercury | Earth | Jupiter |
|---------------------------------|-------|------|-------|---------|-------|---------|
| $x_i$ radius ( $\times 10^6$ m) | 1.0   | 2.3  | 2.4   | 4.9     | 12.8  | 143.0   |
| $d_i$ label                     | -1    | -1   | -1    | 1       | 1     | 1       |

$$\underline{A} = \begin{bmatrix} 1 & 1 \\ 2.3 & 1 \\ 2.4 & 1 \\ 4.9 & 1 \\ 12.8 & 1 \\ 143 & 1 \end{bmatrix}, \underline{d} = \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

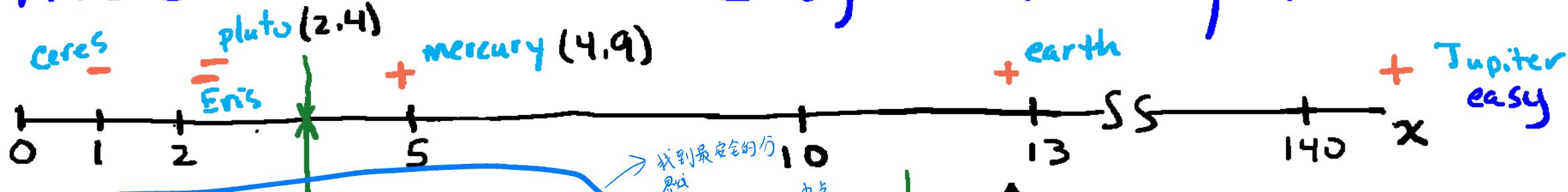
$$\underline{w}_{LS} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

$$\approx 0.01 \begin{bmatrix} 1 \\ -28 \end{bmatrix}$$

dwarf:  $x_i < 28$  (earth!)  
planet:  $x_i \geq 28$  (平方误差对异常值极敏感)

squared error  $\rightarrow$  poor classification

# Avoid loss due to "easy-to-classify" data

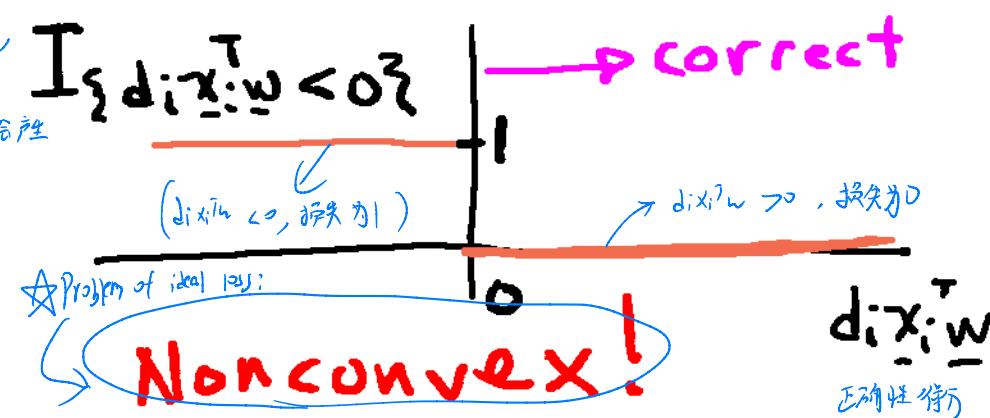
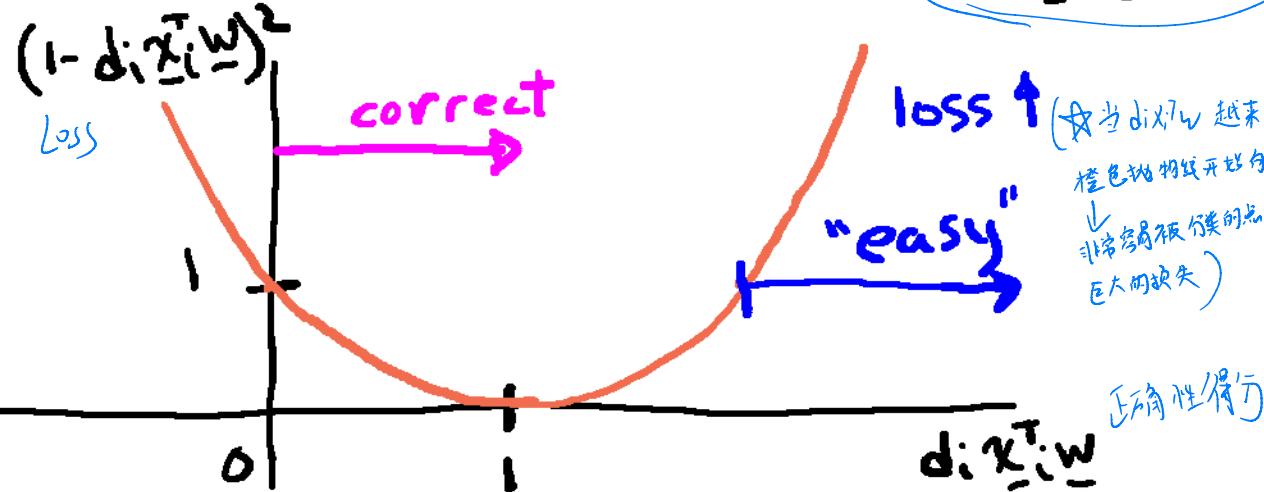


max margin classifier: midpoint  $\hat{d} = \text{sign}(x - 3.65)$

$$\text{Margin } 4.9 - 2.4 = 2.5 \text{ (class separation)}$$

Squared  $d$  error loss:  $\|\underline{A}\underline{w} - \underline{d}\|_2^2 = \sum_{i=1}^N (d_i - \underline{x}_i^\top \underline{w})^2 = \sum_i (1 - d_i \underline{x}_i^\top \underline{w})^2$

correct classification:  $d_i \underline{x}_i^\top \underline{w} > 0$



Since  $d$  is  $\pm 1$ , we can reform  
the formula like this

$$(d_i = \pm 1)$$

Ideal loss

$1 - d_i \underline{x}_i^\top \underline{w}$  is small

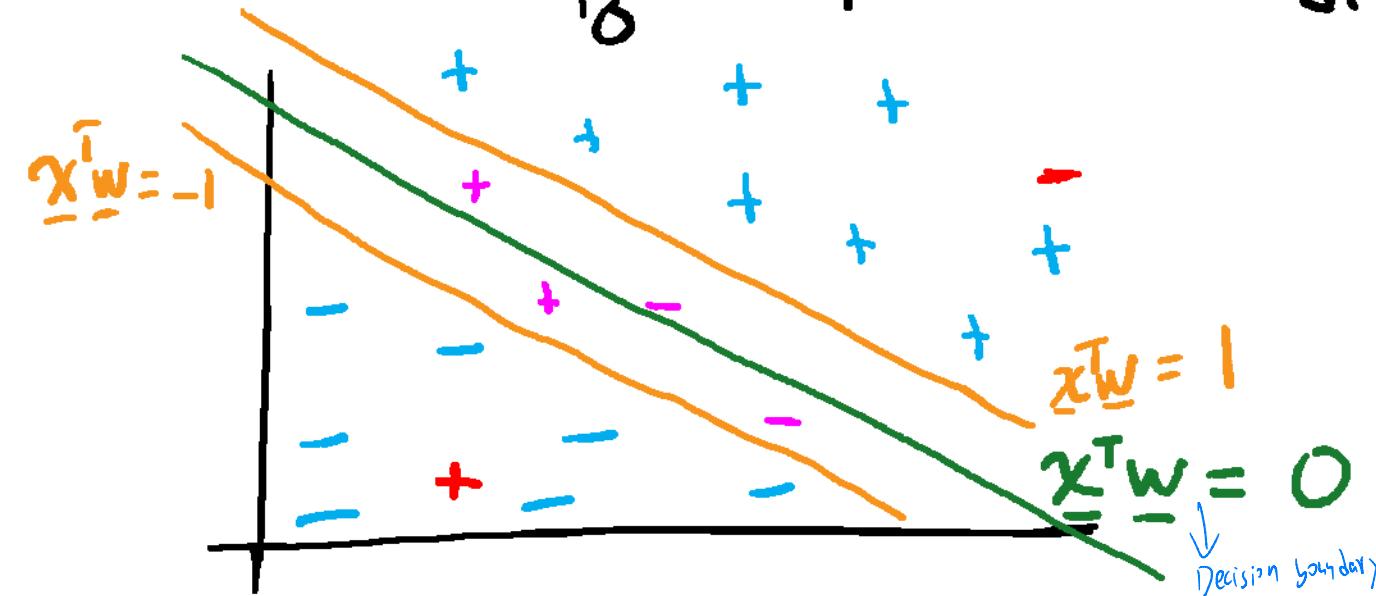
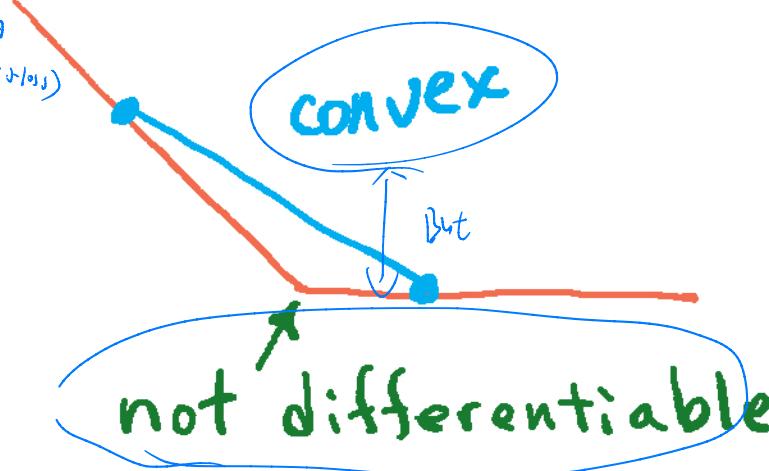
Hinge loss is convex and has no loss for easy 4

$$\ell(\underline{w}; \underline{A}, \underline{d}) = \sum_{i=1}^N (1 - d_i \underline{x}_i^\top \underline{w})_+$$



$$(\alpha)_+ = \begin{cases} \alpha & \alpha > 0 \\ 0 & \alpha \leq 0 \end{cases}$$

Take the positive part



- +,- no hinge loss → 太对了
- +,- small hinge loss → 高决策边界太近  
不够自信, 给小 loss (penalty)
- +,- large hinge loss ↓ 完全不对

Hinge loss better approximates ideal:  
number of misclassifications



Ideal, Hinge loss approximates this

Iterative algorithms required for  
finding minimum hinge loss classifier



Convex, but we can't solve in closed form

**Copyright 2019  
Barry Van Veen**