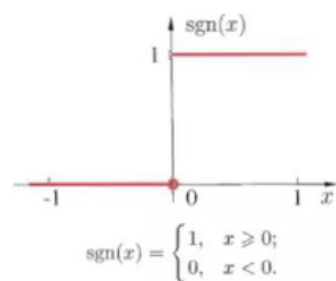
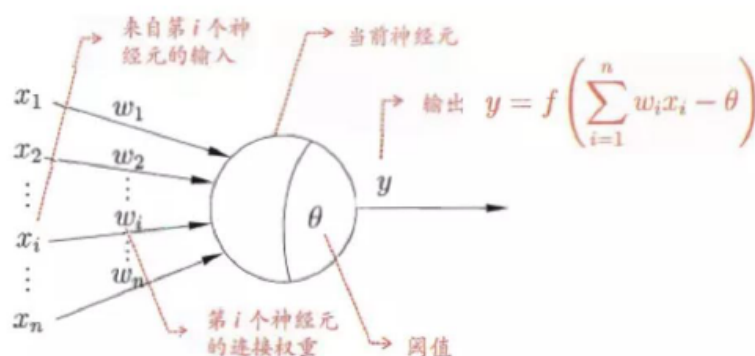


模式识别 U2 感知机

课堂内容

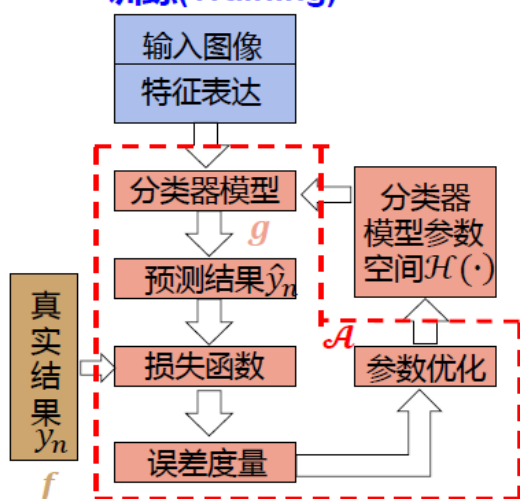
感知器模型参数空间

感知器 Perceptron

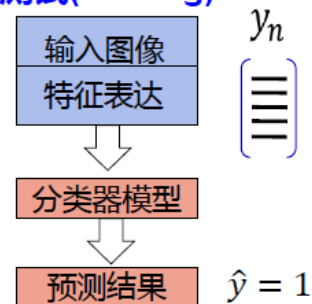


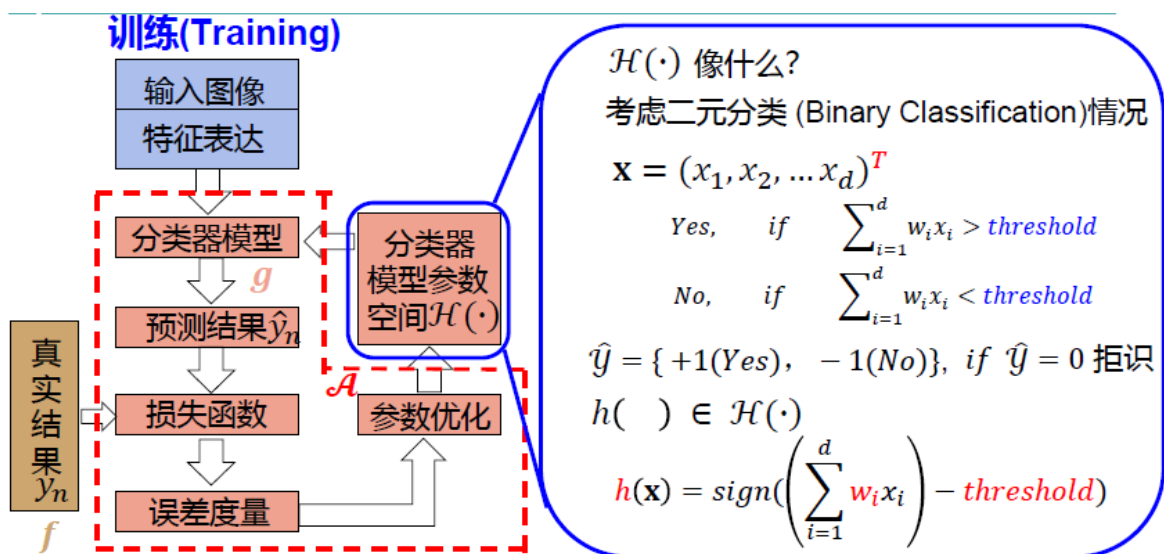
2.1 感知器模型参数空间

训练(Training)



测试(Testing)





用向量形式 (Vector Form) 来表示感知器模型

$$\begin{aligned}
 h(x) &= \text{sign}\left(\left(\sum_{i=1}^d w_i x_i\right) - threshold\right) \\
 &= \text{sign}\left(\left(\sum_{i=1}^d w_i x_i\right) + (-threshold) \cdot (+1)\right) \\
 &= \text{sign}\left(\sum_{i=0}^d w_i x_i\right) \\
 &= \text{sign}\left(\overrightarrow{w^T} \cdot \overrightarrow{x}\right)
 \end{aligned}$$

由上式我们可知，我们将阈值threshold扩展进了原来的w权重向量中，使其作为常数偏置存在；在进行这一操作时也在X中扩展出了一维全为1的增广X

我们称新的w, X为增广化后的 \mathbf{w}^T 、 \mathbf{X}

$$\mathbf{W} = [w_1, w_2, w_3, \dots, w_d, w_{d+1}]_{1 \times (d+1)}$$

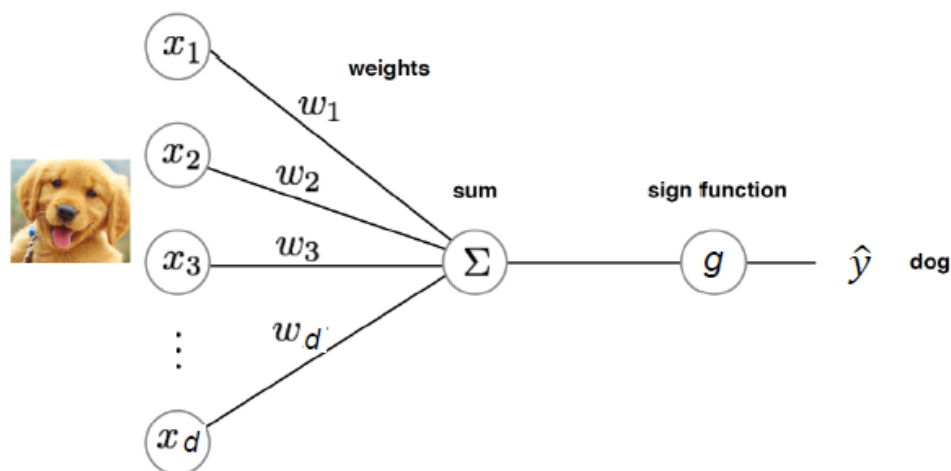
$$\overrightarrow{x_i} = [x_1, x_2, x_3, \dots, x_d, 1]_{1 \times (d+1)}$$

$$\mathbf{X} = \begin{bmatrix} \overrightarrow{x_1} \\ \overrightarrow{x_2} \\ \overrightarrow{x_3} \\ \dots \\ \overrightarrow{x_n} \end{bmatrix}_{n \times (d+1)}$$

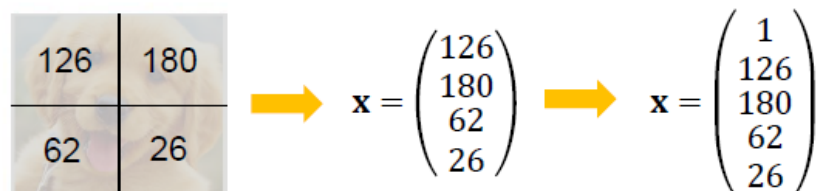
$$\begin{aligned} \text{则 } \mathbf{X} \cdot \mathbf{W}^T &= \begin{bmatrix} \overrightarrow{x_1} \cdot \mathbf{W}^T \\ \overrightarrow{x_2} \cdot \mathbf{W}^T \\ \overrightarrow{x_3} \cdot \mathbf{W}^T \\ \dots \\ \overrightarrow{x_n} \cdot \mathbf{W}^T \end{bmatrix}_{n \times 1} \\ &= \begin{bmatrix} x_1^{(1)} w_1 + x_2^{(1)} w_2 + \dots + x_d^{(1)} w_d + 1^{(1)} \cdot (w_{d+1}) \\ x_1^{(2)} w_1 + x_2^{(2)} w_2 + \dots + x_d^{(2)} w_d + 1^{(2)} \cdot (w_{d+1}) \\ x_1^{(3)} w_1 + x_2^{(3)} w_2 + \dots + x_d^{(3)} w_d + 1^{(3)} \cdot (w_{d+1}) \\ \dots \\ x_1^{(n)} w_1 + x_2^{(n)} w_2 + \dots + x_d^{(n)} w_d + 1^{(n)} \cdot (w_{d+1}) \end{bmatrix}_{n \times 1} \end{aligned}$$

\mathbf{W} 是 $1 \times (d + 1)$ 维, \mathbf{X} 是 $n \times (d + 1)$ 维

将感知器算法用于图像分类(image classification)示例:



将感知器算法用于图像分类(image classification)示例:

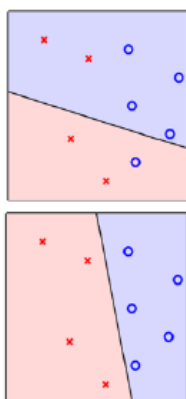


$$\text{if: } \mathbf{w}^T = (3.2 \quad 1.5 \quad 1.3 \quad 2.1 \quad 0.8)$$

$$\hat{y} = h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x}) = \text{sign}\left((3.2 \quad 0.5 \quad 1.3 \quad 2.1 \quad 0.8) \begin{pmatrix} 1 \\ 126 \\ 180 \\ 62 \\ 26 \end{pmatrix}\right) \\ = \text{sign}(451.2) = 1$$

在二维空间观察感知器的分类面(候选的 $\mathcal{H}(\cdot)$)

$$h(\mathbf{x}) = \text{sign}(w_0 + w_1 x_1 + w_2 x_2)$$



样本 \mathbf{x} → 用平面(或者 \mathbf{R}^d 超平面)上的点表示

标签 y → $\circ (+1), \times (-1)$

$h(\mathbf{x})$ → 平面上的线(或者 \mathbf{R}^d 上的超平面)

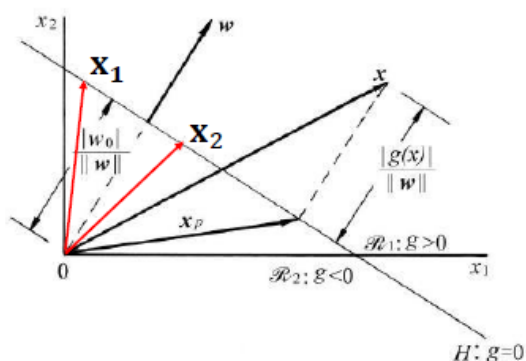
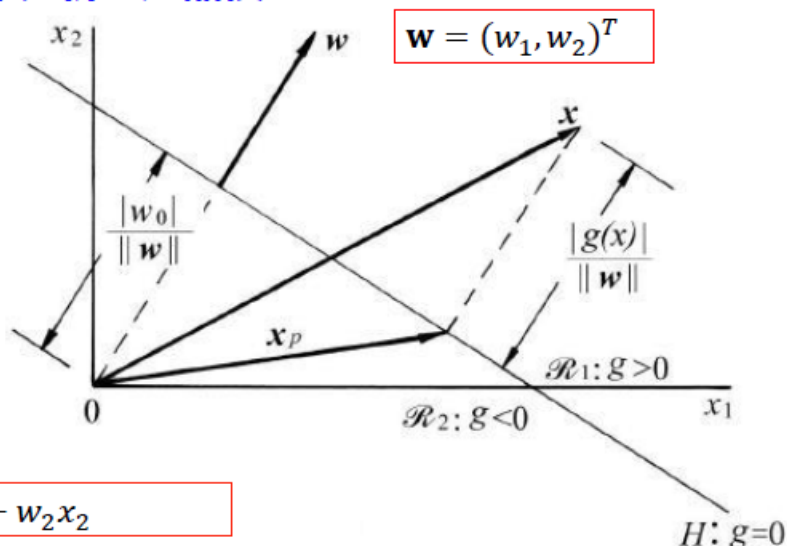
平面被分成两个区域, 分别表示+1类和-1类所在的区域
不同的线有可能将同一样本分到不同的类别中去

感知器也被称作二元线性分类器(binary linear classifiers)

在高维空间中感知器的分类面

$$\begin{aligned}
 h(x) &= \text{sign}(w_0 + w_1x_1 + w_2x_2 + \dots + w_dx_d) \\
 &= \text{sign}\left(\sum_{i=0}^d w_ix_i\right) \\
 &= \text{sign}(W \cdot \mathbf{x})
 \end{aligned}$$

用向量几何知识来分析感知器模型



$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad \mathbf{w} = (w_1, w_2)^T$$

假定 \mathbf{x}_1 和 \mathbf{x}_2 在分类面上:

$$g(\mathbf{x}_1) = \mathbf{w}^T \mathbf{x}_1 + w_0 = 0$$

$$g(\mathbf{x}_2) = \mathbf{w}^T \mathbf{x}_2 + w_0 = 0$$

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

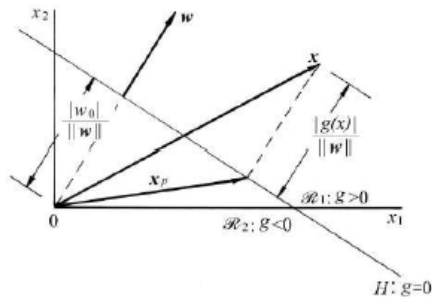
$$\mathbf{w} \perp (\mathbf{x}_1 - \mathbf{x}_2)$$

权向量 \mathbf{w} 垂直于分类面

用向量几何知识来分析感知器模型

如果 $\mathbf{x} \in D$ 在特征平面上, \mathbf{x}_p 在分类面上

r 表示 \mathbf{x} 与分类面 ($g(\mathbf{x}_p) = 0$) 之间的距离



$$g(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2$$

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

$$\begin{aligned} g(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + w_0 \\ &= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + w_0 \\ &= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\ &= r \|\mathbf{w}\| \end{aligned}$$



几何知识：二维中 点到直线的距离

$$l: ax_1 + bx_2 + c = 0$$

则距离

$$d = \frac{|ax_{p1} + bx_{p2} + c|}{\sqrt{a^2 + b^2}}$$

扩展到如今向量几何当中

$$d = r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

其中 $\frac{\mathbf{w}}{\|\mathbf{w}\|}$ 表示单位法向量, r 则可以用标量指示距离

$P. \|\mathbf{w}\|$ 指向量的模, 也可理解为向量 \mathbf{w} 的 L_2 范数

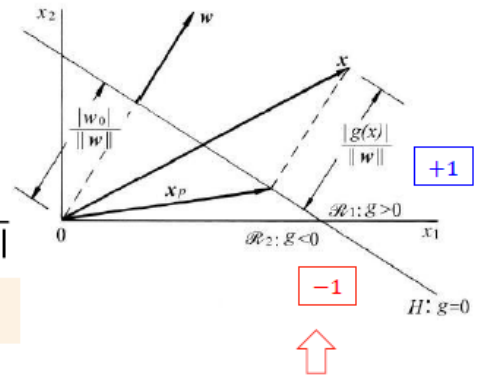
用向量几何知识来分析感知器模型

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$

当 $\mathbf{x} = \mathbf{0}$, $g(\mathbf{x}) = w_0 \longrightarrow r = \frac{w_0}{\|\mathbf{w}\|}$

即：原点到分类面的距离为 $r = \frac{w_0}{\|\mathbf{w}\|}$

- 当 $w_0 > 0$, 原点处于分类面的正区域
- 当 $w_0 < 0$, 原点处于分类面的负区域
- 当 $w_0 = 0$, 分类面穿过原点



$$h(\mathbf{x}) = \text{sign}(w_0 + w_1x_1 + w_2x_2) \\ = \text{sign}(g(\mathbf{x}))$$

有上述推导我们可得：

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$

其中

\mathbf{w} 是训练得到的感知器模型，其本质是可学习迭代的参数集合； $g(\mathbf{x})$ 则是将该数据点代入模型中取得的结果

$$\mathbf{W}^T \cdot \mathbf{X} = \|\mathbf{W}\| * \|\mathbf{X}\| \cos \theta$$

感知器算法 PLA

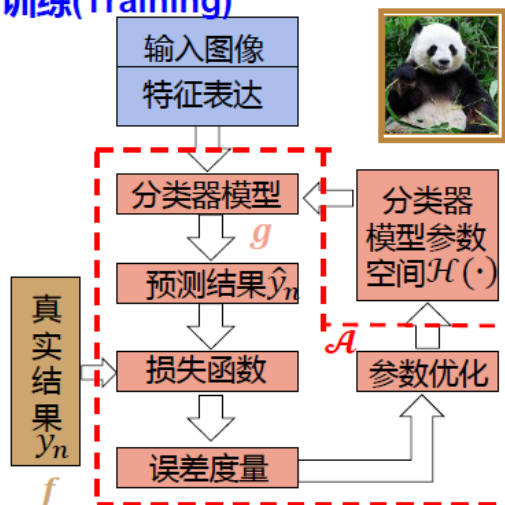
PLA (Perceptron Learning Algorithm)

算法思路

2.2 感知器算法 (PLA)



训练(Training)



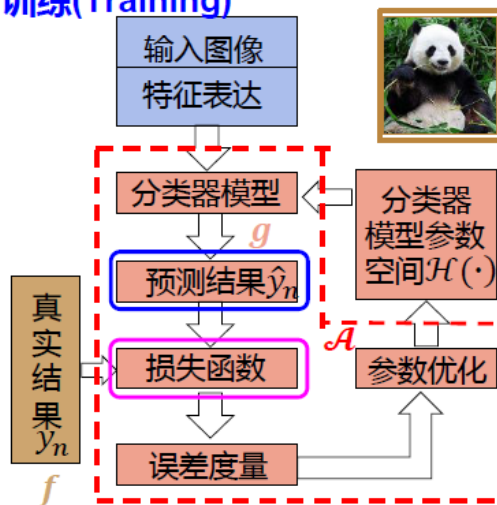
算法 \mathcal{A} 的目的是在 $\mathcal{H}(h(\cdot))$ 中找到最优结果作为分类器的模型 g

- 最优结果: $g \approx f$
- 挑战: f 未知
- 学习的资源: 在训练集 \mathcal{D} 上, 如果每一个样本都有: $g(\mathbf{x}_n) = y_n = f(\mathbf{x}_n)$, 则在训练集 \mathcal{D} 上做到了 $g \approx f$
- 困难: $\mathcal{H}(h(\cdot))$ 的候选模型无穷多

2.2 感知器算法 (PLA)



训练(Training)

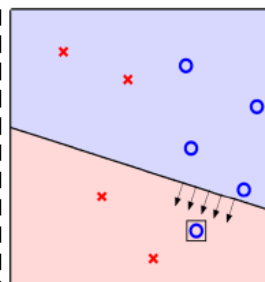


$$\hat{y}_{n(t)} = \text{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)})$$

$$L_{in} = \llbracket y_n \neq \hat{y}_{n(t)} \rrbracket$$

L_{in} 训练阶段

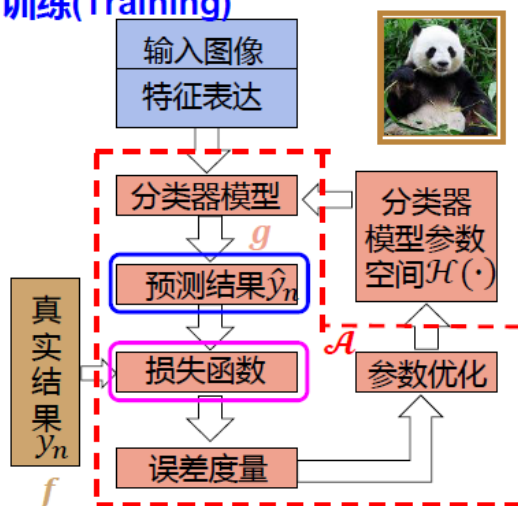
算法思路



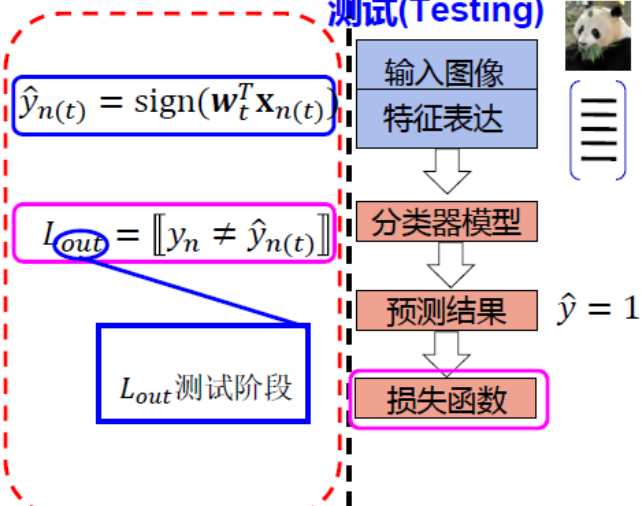
- 设置初始分类面 (权重) \mathbf{w}_0
- 如果有样本分错, 就修正权重

2.2 感知器算法 (PLA)

训练(Training)

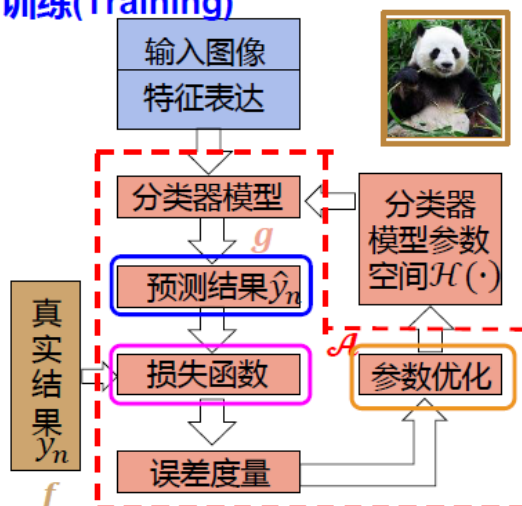


测试(Testing)

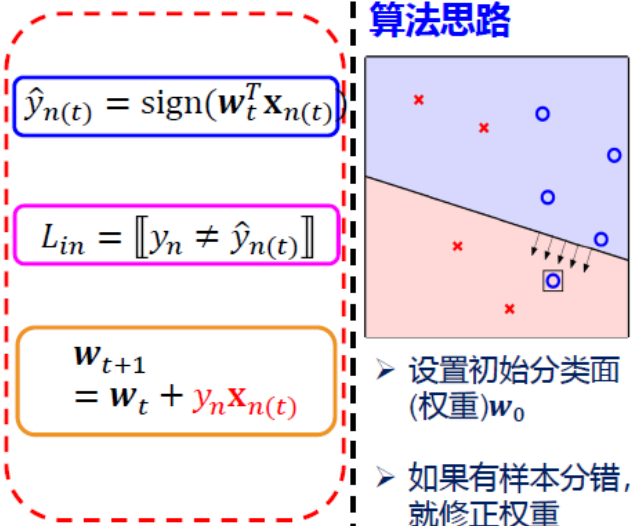


2.2 感知器算法 (PLA)

训练(Training)



算法思路



算法流程

- 对样本的特征向量 \mathbf{x} 和权向量 \mathbf{w} 增广化
- 初始化权向量 \mathbf{w}_0 (例如: $\mathbf{w}_0 = \mathbf{0}$)
- *for* $t = 0, 1, 2, \dots$ (t 代表迭代次数)
 - ① 进行到第 t 次迭代时权向量为 \mathbf{w}_t , 它对样本 $(\mathbf{x}_{n(t)}, y_{n(t)})$ 错分
 $\text{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)}) \neq y_n$
 - ② 通过下式对权向量 \mathbf{w}_t 进行更新: $\mathbf{w}_{t+1} = \mathbf{w}_t + y_n \mathbf{x}_{n(t)}$
 ...直到所有样本均能被正确分类, 此时的 \mathbf{w}_{t+1} 作为学到的 \mathbf{g}

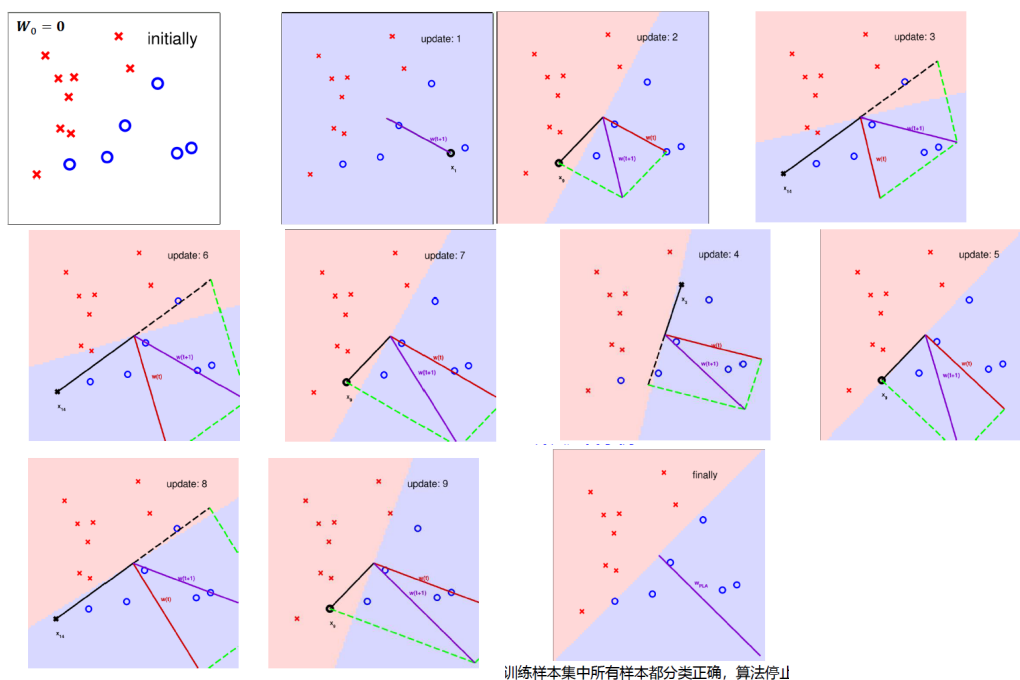


回顾感知器算法:

- 对样本的特征向量 \mathbf{x} 和权向量 \mathbf{w} 增广化
- 初始化权向量 \mathbf{w}_0 (例如: $\mathbf{w}_0 = \mathbf{0}$)
- *for* $t = 0, 1, 2, \dots$ (t 代表迭代次数)
 - ① 对某些样本 n , 通过下式对权向量 \mathbf{w}_t 进行更新:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + 1 \cdot (\llbracket \text{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)}) \neq y_n \rrbracket \mathbf{x}_{n(t)})$$
 ...直到满足停止条件, 此时的 \mathbf{w}_{t+1} 作为学到的 \mathbf{g}

算法迭代示例



算法问题

问题1：算法会收敛吗？

结果与输入样本顺序是否有关？

问题2：能学到 $g \approx f$ 吗？

- 在 \mathcal{D} 上，如果 $L_{in} = 0$ ， $g \approx f$ ？
- 不在 \mathcal{D} 上时， $L_{out} = 0$ ？ $g \approx f$ ？
- 算法不能收敛时， $g \approx f$ ？

感知器算法的收敛性

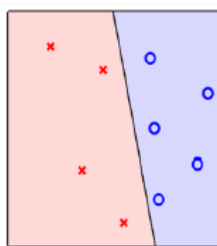
PLA收敛条件：数据集中所有样本线性可分

线性可分性 (linear separable)

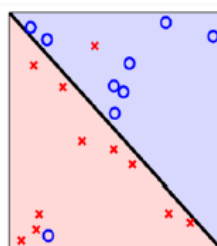
PLA算法收敛 \longleftrightarrow 算法停止 \longleftrightarrow w 对 \mathcal{D} 上所有样本正确分类



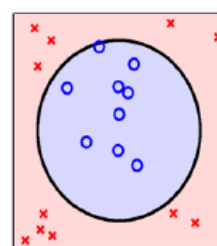
\mathcal{D} 上所有样本是线性可分的



线性可分



线性不可分



线性不可分

所有样本线性可分是否意味着PLA一定收敛?

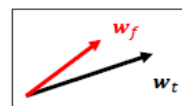
假设 w_f 是理想的分类面:

\mathcal{D} 是线性可分的 $\longleftrightarrow y_n = \text{sign}(w_f^T x_n)$

第 t 次迭代时, 任意一个样本 $x_{n(t)}$ 满足

$$y_{n(t)} w_f^T x_{n(t)} \geq \min_n y_n w_f^T x_n > 0$$

迭代次数增加



$$w_f^T w_{t+1} = w_f^T (w_t + y_{n(t)} x_{n(t)}) \geq w_f^T w_t + \min_n y_n w_f^T x_n > w_f^T w_t + 0$$

结论: 随着迭代次数增加, $w_f^T w_t$ 随之增加, 意味着 w_t 与 w_f 越来越接近

假设 w_t 是第 t 次迭代得到的分类面:

当 x_n 被分类错误时, 才更新 $w_t \longleftrightarrow \text{sign}(w_t^T x_n) \neq y_n \longleftrightarrow y_{n(t)} w_t^T x_{n(t)} \leq 0$

假设 x_n 在样本集中模值最大, 随着迭代次数增加, $\|w_t\|$?

$$\begin{aligned} \|w_{t+1}\|^2 &= \|w_t + y_{n(t)} x_{n(t)}\|^2 \\ &= \|w_t\|^2 + 2y_{n(t)} w_t^T x_{n(t)} + \|y_{n(t)} x_{n(t)}\|^2 \\ &\leq \|w_t\|^2 + 0 + \|y_{n(t)} x_{n(t)}\|^2 \leq \|w_t\|^2 + \max_n \|y_n x_n\|^2 \end{aligned}$$

结论: 随着迭代次数增加, w_t 模值增长不会太快, 意味着 w_t 与 w_f 的接近是方向上在靠近, 而非模值的贡献

课后证明 (作业) :

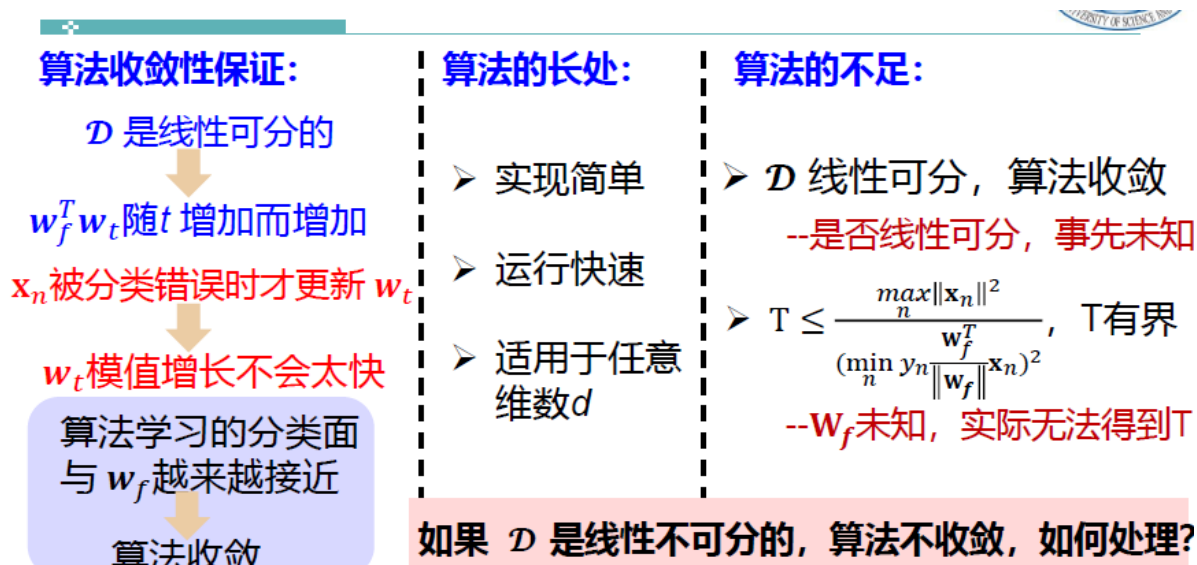
(1) 针对线性可分训练样本集, PLA 算法中, 当 $w_0 = 0$, 在对分错样本进行了 T 次修正后, 下式成立: $\frac{w_f^T}{\|w_f\|} \frac{w_T}{\|w_T\|} \geq \sqrt{T} \cdot constant$

(2) 针对线性可分训练样本集, PLA 算法中, 假设对分错样本进行了 T 次修正后, 得到的分类面不再出现错分状况, 定义:

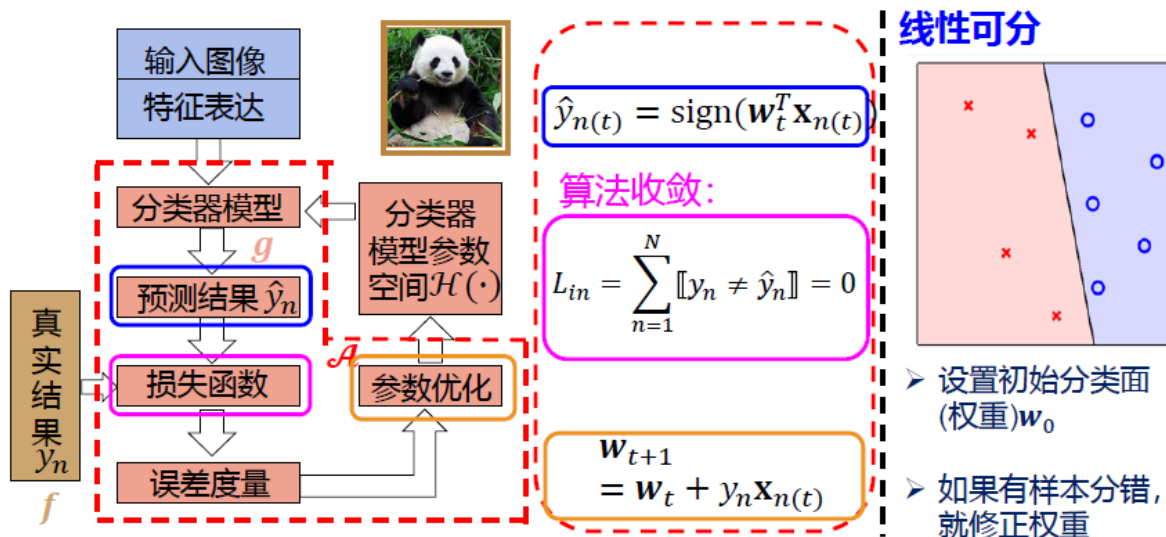
$$R^2 = \max_n \|x_n\|^2, \quad \rho = \min_n y_n \frac{w_f^T}{\|w_f\|} x_n, \quad \text{证明: } T \leq \frac{R^2}{\rho^2}$$

线性不可分情况

线性不可分分析



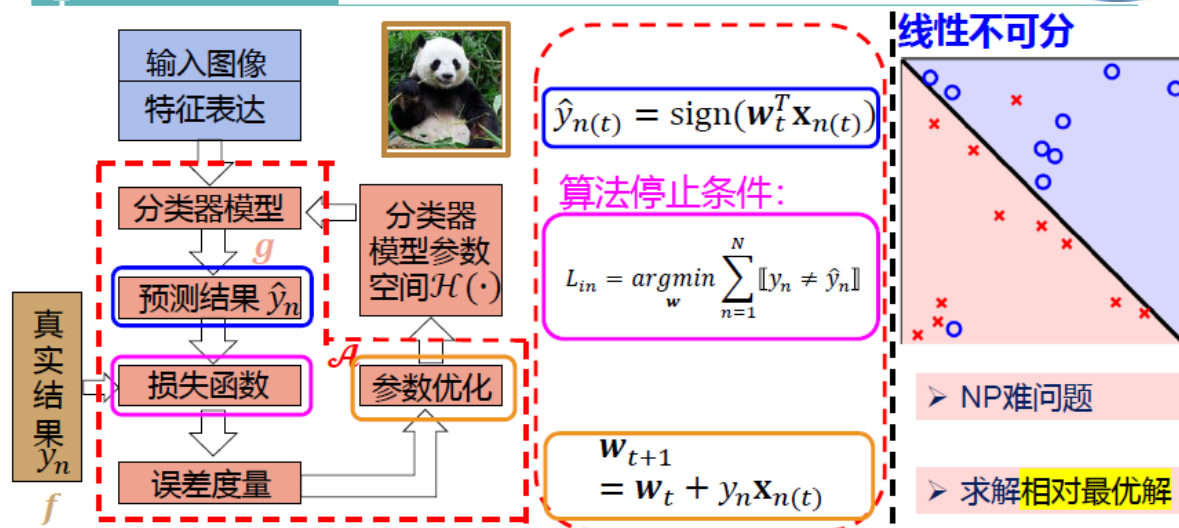
对于线性可分情况



模型的最终目的是实现收敛，即全部样本完全正确划分

对于线性不可分情况

2.4 线性不可分情况 (Non-separable Data)



调整模型算法停止条件为:损失函数最小

Pocket算法

为处理线性不可分情况而对PLA算法的修正

Pocket 算法 – 为处理线性不可分情况而对PLA算法的修正

- 对样本的特征向量 \mathbf{x} 和权向量 \mathbf{w} 增广化
- 初始化权向量 \mathbf{w}_0 (例如: $\mathbf{w}_0 = \mathbf{0}$), 并任意选一个“Pocket”向量 $\hat{\mathbf{w}}$
- **for** $t = 0, 1, 2, \dots$ (t 代表迭代次数)

① 进行到第 t 次迭代时权向量为 \mathbf{w}_t , 它对样本 $(\mathbf{x}_{n(t)}, y_{n(t)})$ 错分

$$\operatorname{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)}) \neq y_n$$

② 通过下式对权向量 \mathbf{w}_t 进行更新: $\mathbf{w}_{t+1} = \mathbf{w}_t + y_n \mathbf{x}_{n(t)}$

③ 如果 \mathbf{w}_{t+1} 在所有样本集上错分的样本少于 $\hat{\mathbf{w}}$, 则用 \mathbf{w}_{t+1} 代替 $\hat{\mathbf{w}}$, 并在错分样本中随机选一个对权向量进行更新

...达到指定的迭代次数

返回此时的“Pocket”向量 $\hat{\mathbf{w}}$ 作为算法学到的 g

Pocket算法

- ❖ 贪心思想对PLA的优化：每次更新权值前，先遍历所有点，如果错误点数量下降，则更新权值，否则不更新
 - 在选择修正点时存在随机性，迭代次数可能比PLA小也可能比PLA大，每次找到的直线也不相同
 - 执行时间比PLA长，这是因为在每次调整参数后需要去比对修正后与当前最优参数的结果好坏
 - Pocket算法可以保存最优参数，当得到相应迭代次数或者完全正确分类的时候停止
 - 在足够的迭代次数的情况下，便可以找到全局最优解，对于线性不可分数据集，也可以在迭代到相应次数时，停止算法，输出最优结果

小结

2.1 感知器模型参数空间

在 \mathbb{R}^d 空间的超平面的线性分类面

2.2 感知器算法 (PLA)

通过迭代的方式对错分样本的分类面进行修正

2.3 感知器算法的收敛性

如果训练样本集是线性可分的，算法能对所有的样本正确分类并停止

2.4 线性不可分情况

通过“Pocket”算法在设置的迭代次数下寻找相对最佳的分类面

作业部分

手写作业

1，假设训练样本集为 $D = \{(\mathbf{x}_1, y_1) = ((3,3)^T, 1), (\mathbf{x}_2, y_2) = ((4,3)^T, 1), (\mathbf{x}_3, y_3) = ((1,1)^T, -1)\}$ ，使用感知器算法设计分类面，并判断测试样本 $\mathbf{x} = (0,1)^T$ 属于哪个类别。↵

T1. \vec{x} y 初始化 $\vec{w}_0 = [0 \ 0 \ 0]_{1 \times 3}$
 解: $(3, 3) \ 1$ 增 $\vec{x} = \begin{bmatrix} 3 & 3 & 1 \\ 4 & 3 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vec{x}_3 \end{bmatrix}$
 $(4, 3) \ 1$
 $(1, 1) \ -1$

迭代 I1: $\text{sign}(\vec{w}_0 \cdot \vec{x}_1) = 0 \neq y_1 = 1$ 错分
 $\vec{w}_1 = \vec{w}_0 + y_1 \vec{x}_1 = [3 \ 3 \ 1]$
 I2: $\text{sign}(\vec{w}_1 \cdot \vec{x}_1) = 1 = y_1 = 1 \checkmark$
 $\text{sign}(\vec{w}_1 \cdot \vec{x}_2) = 1 = y_2 = 1 \checkmark$
 $\text{sign}(\vec{w}_1 \cdot \vec{x}_3) = 1 \neq y_3 = -1 \times$ 错分
 $\vec{w}_2 = \vec{w}_1 + y_3 \vec{x}_3 = [2 \ 2 \ 0]$
 I3: $\text{sign}(\vec{w}_2 \cdot \vec{x}_1) = 1 \checkmark$
 $\text{sign}(\vec{w}_2 \cdot \vec{x}_2) = 1 \checkmark$
 $\text{sign}(\vec{w}_2 \cdot \vec{x}_3) = 0 \neq y_3 = -1 \times$
 $\vec{w}_3 = \vec{w}_2 + y_3 \vec{x}_3 = [1 \ 1 \ -1]$

以此类推进行迭代
 $\vec{w}_4 = [0 \ 0 \ -2]$
 $\vec{w}_5 = [3 \ 3 \ -1]$
 $\vec{w}_6 = [2 \ 2 \ -2]$
 $\vec{w}_7 = [1 \ 1 \ -3]$

I8: $\vec{w}_7 = [1 \ 1 \ -3]$
 $\text{sign}(\vec{w}_7 \cdot \vec{x}_1) = 1 \checkmark$
 $\text{sign}(\vec{w}_7 \cdot \vec{x}_2) = 1 \checkmark$
 $\text{sign}(\vec{w}_7 \cdot \vec{x}_3) = -1 \checkmark$
 全部样本正确划分
 算法收敛

2, 对于感知器算法 (PLA), 假设第 t 次迭代时, 选择的是第 n 个样

本: $\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$, $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_n \mathbf{x}_n$, 下述那个式子正确? \leftarrow

(a) $\mathbf{w}_{t+1}^T \mathbf{x}_n = y_n \dots \dots \dots \leftarrow$

(b) $\text{sign}(\mathbf{w}_{t+1}^T \mathbf{x}_n) = y_n \leftarrow$

(c) $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n \geq y_n \mathbf{w}_t^T \mathbf{x}_n \leftarrow$

(d) $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n < y_n \mathbf{w}_t^T \mathbf{x}_n \leftarrow$

T2. 解. $\vec{w}_{t+1} = \vec{w}_t + y_n \vec{x}_n$

$$y_n \vec{w}_{t+1} \vec{x}_n = y_n (\vec{w}_t + y_n \vec{x}_n) \vec{x}_n = y_n \vec{w}_t \vec{x}_n + y_n^2 |\vec{x}_n|^2 \geq y_n \vec{w}_t \vec{x}_n$$

3, 证明: 针对线性可分训练样本集, PLA 算法中, 当 $\mathbf{w}_0 = \mathbf{0}$, 在对分错样本进行了 T 次纠正后, 下式成立: $\frac{\mathbf{w}_f^T \mathbf{w}_T}{\|\mathbf{w}_f\| \|\mathbf{w}_T\|} \geq \sqrt{T} \cdot \text{constant}$

T3 解

$$\mathbf{w}_f^T \mathbf{w}_{t+1} = \mathbf{w}_f^T (\mathbf{w}_t + y_i \vec{x}_i) = \mathbf{w}_f^T \mathbf{w}_t + y_i \mathbf{w}_f^T \vec{x}_i$$

$\because \mathbf{w}_f$ 是理想最佳权重模型 $\therefore y_i \mathbf{w}_f^T \vec{x}_i \geq 0 \quad \mathbf{w}_f^T \mathbf{w}_{t+1} \uparrow$

$$\mathbf{w}_f^T \mathbf{w}_{t+1} = \mathbf{w}_f^T \mathbf{w}_t + y_i \mathbf{w}_f^T \vec{x}_i \geq \mathbf{w}_f^T \mathbf{w}_t + \min y_n \mathbf{w}_f^T \vec{x}_n$$

迭代 $\Rightarrow \mathbf{w}_f^T \mathbf{w}_T \geq \mathbf{w}_f^T \mathbf{w}_0 + T \min y_n \mathbf{w}_f^T \vec{x}_n$

$$\|\mathbf{w}_{t+1}\|^2 = \|\mathbf{w}_t + y_i \vec{x}_i\|^2 \quad \text{错的} \Rightarrow \therefore < 0$$

$$= \|\mathbf{w}_t\|^2 + 2 y_i \mathbf{w}_t^T \vec{x}_i + \|y_i \vec{x}_i\|^2$$

$$\leq \|\mathbf{w}_t\|^2 + 0 + \|y_i \vec{x}_i\|^2 \leq \|\mathbf{w}_t\|^2 + \max \|\vec{x}_i\|^2$$

$$\therefore \|\mathbf{w}_T\| \leq \sqrt{\|\mathbf{w}_0\|^2 + T \max \|\vec{x}_i\|^2} = \sqrt{T \max \|\vec{x}_i\|^2}$$

$$\therefore \frac{\mathbf{w}_f^T \mathbf{w}_T}{\|\mathbf{w}_f\| \|\mathbf{w}_T\|} \geq \frac{T \min y_n \mathbf{w}_f^T \vec{x}_n}{\|\mathbf{w}_f\| \sqrt{T \max \|\vec{x}_i\|^2}} = \sqrt{T} \cdot \frac{\min y_n \mathbf{w}_f^T \vec{x}_n}{\|\mathbf{w}_f\| \max \|\vec{x}_i\|} = \text{constant}$$

常数

\therefore 记得 $\frac{\mathbf{w}_f^T \mathbf{w}_T}{\|\mathbf{w}_f\| \|\mathbf{w}_T\|} \geq \sqrt{T} \cdot \text{constant}$

4, 针对线性可分训练样本集, PLA 算法中, 假设对分错样本进行了 T 次纠正后得到的分类面不再出现错分状况, 定义: $R^2 = \max_n \|\mathbf{x}_n\|^2$,

$\rho = \min_n y_n \frac{\mathbf{w}_f^T}{\|\mathbf{w}_f\|} \mathbf{x}_n$, 试证明: $T \leq \frac{R^2}{\rho^2}$

T4. 解:

$$W_T = W_{T-1} + y_i \vec{x}_i \geq W_{T-1} + \min_n y_n \vec{x}_n$$

$$W_f^T W_T \geq W_f^T W_0 + T \min_n y_n W_f^T \vec{x}_n = T \min_n y_n W_f^T \vec{x}_n$$

$$\|W_T\|^2 \geq \|W_0\|^2 + 0 + T \max_n \|\vec{x}_n\|^2 \quad \|W_T\| \geq \sqrt{T} \cdot R^2$$

$$1 \geq \frac{W_f^T W_T}{\|W_f\| \|W_T\|} \geq \frac{1}{\|W_T\|} \cdot \frac{W_f^T W_T}{\|W_f\|} = \frac{1}{\sqrt{T} \cdot R} \cdot T \cdot \min_n y_n \frac{W_f^T \vec{x}_n}{\|W_f\|} = \sqrt{T} \cdot \frac{\rho}{R}$$

$$\therefore T \leq \frac{R^2}{\rho^2}$$

2/2

编程作业