

SVM

感知机中, inductive bias 中认为二分任务中的边界是线性的

但若 if not linearly separable 呢?

这里介绍 kernel method (当然方法不限于此)

Def: A kernel K is a legal def of dot-product, i.e., 有一个隐式映射: Φ , s.t., $K(x, y) = \Phi(x) \cdot \Phi(y)$.

Eg. $K(x, y) = (x \cdot y + 1)^d$, $\Phi: n\text{维} \rightarrow n^d\text{维}$

这个定义看起来没啥用; 但 perceptron 中因 $x \cdot z$ 是跳不出 linear 的, 而 $x \cdot z$ 换为 $K(x, z)$ 就能拥抱高维

Formal Definition: $K(\cdot, \cdot)$ is a kernel if it can be viewed as a legal definition of inner product

• $\exists \Phi: X \rightarrow \mathbb{R}^N$ s.t. $K(x, z) = \Phi(x) \cdot \Phi(z)$, range of Φ : Φ -space
但 Φ 是隐式的! 这仅是为了方便 view kernel as inner product.

例: $K(x, z) = (x \cdot z)^d$ 对应: $\Phi(x) = (x_1^d, x_2^d, \sqrt{d} x_1 x_2)$

可见 $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$, $\Phi(x) \cdot \Phi(z) = (x_1^d, x_2^d, \sqrt{d} x_1 x_2) \cdot (z_1^d, z_2^d, \sqrt{d} z_1 z_2)$
 $= (x_1 z_1 + x_2 z_2)^2 = (x \cdot z)^2 = K(x, z)$

✱: 同一个 kernel, Φ 也许并不唯一! Eg. $(x, x) \rightarrow \Phi(x) = (x_1^2, x_2^2, x_1 x_2, x_2 x_1)$
应避免显式地扩到高维, i.e., Φ ! 因为 feature space 会长的飞快!!

因此可以尝试用 $K(x, z)$ 代替 $x \cdot z$ 。例如 Perceptron 中:

$w_t = \alpha_{i1} x_{i1} + \dots + \alpha_{ik} x_{ik}$ (可视为未正确分类的 sample 的线性组合).

对于第 t 的错误: mistake on positive: $\alpha_{it} \leftarrow 1$, store x_{it}

mistake on negative: $\alpha_{it} \leftarrow -1$, store x_{it}

$w_t \cdot x = \alpha_{i1} K(x_{i1}, x) + \dots + \alpha_{ik} K(x_{ik}, x)$



这个方法在 margin 恰是 Φ 维空间中 linear boundary 时表现很好!

可知: 若 margin 在 Φ -space 中, 则 Perceptron makes $(\frac{R}{\gamma})^2$ mistakes

Kernel Example:

Linear: $K(x, z) = x \cdot z$

Polynomial: $K(x, z) = (x \cdot z)^d$ or $K(x, z) = (1 + x \cdot z)^d$

Gaussian: $K(x, z) = \exp \left[-\frac{\|x - z\|^2}{2\sigma} \right]$

Laplace: $K(x, z) = \exp \left[-\frac{\|x - z\|}{2\sigma} \right]$

Properties of Kernel: K is kernel iff: ① K is symmetric

② For any training points x_1, \dots, x_m , and $\forall a_1, \dots, a_m \in \mathbb{R}$:

$$\sum_{i,j} a_i a_j K(x_i, x_j) \geq 0$$

i.e., $K = (K(x_i, x_j))_{i,j=1,\dots,n}$ 是半正定的: $a^T K a \geq 0$

若 $K_1(\cdot, \cdot)$, $K_2(\cdot, \cdot)$ 均为 kernels, 则 $\forall c_1, c_2 \geq 0$, $K(x, z) = c_1 K_1(x, z) + c_2 K_2(x, z)$ 也是 kernel! 且 $K_1(\cdot, \cdot)$, $K_2(\cdot, \cdot)$ 也是!

Proof: $\phi(x) = (\sqrt{c_1} \phi_1(x), \sqrt{c_2} \phi_2(x))$, $\phi(x) \phi(z) = c_1 \phi_1(x) \phi_1(z) + c_2 \phi_2(x) \phi_2(z)$

$\phi(x) = (\phi_{1,i}(x), \phi_{2,j}(x))_{i \in \{1, \dots, n\}, j \in \{1, \dots, m\}}$

$$\phi(x) \phi(z) = \sum_{i,j} \phi_{1,i}(x) \phi_{2,j}(x) \phi_{1,i}(z) \phi_{2,j}(z)$$

$$= \sum_i \phi_{1,i}(x) \phi_{1,i}(z) \left(\sum_j \phi_{2,j}(x) \phi_{2,j}(z) \right)$$

综合上述所有铺垫, SVM, Support Vector Machine 问世!

Def: Margin of example x w.r.t. linear separation w is the distance from x to plane $w \cdot x = 0$

Def: Margin γ_w of a set S is min margin over $x \in S$, w.r.t. a linear separator w .



Def: Margin γ of set S is the maximum γw over all linear sep. w

因此 SVM 目标为: Input: $S = \{(x_1, y_1) \dots (x_m, y_m)\}$

Find: some w and largest maximum γ where:

① $\|w\|^2 = 1$ ② For all i , $y_i w \cdot x_i \geq \gamma$

Output: Maximum marginal separator

等效为: $y_i \cdot \frac{w}{\gamma} \cdot x_i \geq 1$, 令 $w' = w/\gamma$. 可见是 $y_i w' \cdot x_i \geq 1$
满足前提下: minimize $\|w'\|^2$, i.e.,

$$\operatorname{argmin}_w \|w\|^2, \text{ s.t., } \forall i, y_i w \cdot x_i \geq 1$$

Dual Problem* (Equivalent to:)

$$\begin{aligned} \max_{\alpha} \sum \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) \\ = \max_{\alpha} \sum \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j K(x_i, x_j) \end{aligned}$$

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

最终: classifier 为: $w = \sum \alpha_i y_i x_i$

* 这一坨式子咋得到的?

在最大化间隔同时, 允许部分样本不满足约束 (尽可能少)

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \ell_{\alpha_1}(y_i(w^T x_i + b) - 1), \quad C: \text{constant}$$

$$\ell_{\alpha_1}(z) = \begin{cases} 1 & \text{if } z < 0 \\ 0 & \text{otherwise} \end{cases}, \text{ 取 hinge loss: } \ell_{\text{hinge}}(z) = \max(0, 1 - z)$$

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i(w^T x_i + b)), \text{ 引入 slack 变量 } \xi_i \geq 0:$$

$$= \min_{w, b, \xi_i} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \quad ①$$

三个约束条件

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

②

③

KOKUYO



则构造 Lagrange: $L(w, b, \alpha, \xi, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$
 $+ \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i (w^T x_i + b)) - \sum_{i=1}^m \mu_i \xi_i$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i \quad (1)$$

$$0 = \sum_{i=1}^m \alpha_i y_i \quad (2) \quad C = \alpha_i + \mu_i \quad (3) \quad \text{代回去:}^*$$

Dual Problem: $\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j$

s.t. $\sum_{i=1}^m \alpha_i y_i = 0$, $0 \leq \alpha_i \leq C$, $i=1, 2, \dots, m$ $= K(x_i, x_j)$

这是 $K(x, z) = x \cdot z$ 时; 若为 $\Phi(x)$, 则 $\Phi(x_i)^T \Phi(x_j) = \Phi(x_i) \cdot \Phi(x_j)$

Δ : 附: 使用拉格朗日乘子法得到的便是“对偶问题”

$$\begin{aligned} * : \inf_{w, b} L(w, b, \alpha) &= \frac{1}{2} w^T w + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i y_i w^T x_i - \sum_{i=1}^m \alpha_i y_i b \\ &= \frac{1}{2} w^T \sum_{i=1}^m \alpha_i y_i x_i - w^T \sum_{i=1}^m \alpha_i y_i x_i + \sum_{i=1}^m \alpha_i - b \sum_{i=1}^m \alpha_i y_i \\ &= -\frac{1}{2} w^T \sum_{i=1}^m \alpha_i y_i x_i + \sum_{i=1}^m \alpha_i \\ &= -\frac{1}{2} \sum_{i=1}^m \alpha_i y_i x_i^T \sum_{j=1}^m \alpha_j y_j x_j + \sum_{i=1}^m \alpha_i \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \end{aligned}$$

