

Perceptron

Key idea: Try to learn a hyperplane

设 \vec{a} 为列向量 (by default), $\vec{a}^T \vec{b} = \sum_{d=1}^D a_d b_d$

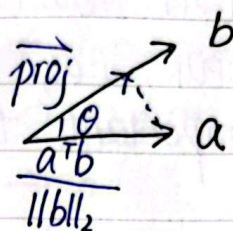
L_2 -norm: $\|\vec{a}\|_2 = \sqrt{\vec{a}^T \vec{a}}$

则 \vec{a} 在 \vec{b} 上的投影为?

$$\Delta: \vec{a} \cdot \vec{b} = \|\vec{b}\| \|\vec{a}\| \cos \theta$$

$$\text{而 } \vec{a}^T \vec{b} = \|\vec{b}\| \cdot \|\text{proj}\|, \quad \|\text{proj}\| = \frac{\vec{a}^T \vec{b}}{\|\vec{b}\|_2}$$

故: $\text{proj} = \frac{(\vec{a}^T \vec{b})}{\|\vec{b}\|_2^2} \vec{b}$, 因为方向为 \vec{b} 方向



在 2D 中, $w_1 x_1 + w_2 x_2 + b = 0$ 为 line; 3D 中为 plane

4+D 中表示为 hyperplane, $w^T x + b = 0$. 值得注意的是:

w 向量指向总是 $w^T x + b > 0$ 的部分!

因此, hyperplane creates 2 halfspaces.

Perceptron 算法: 如何更新 w 与 b ?

初始化: $w = \vec{0}$, $b = 0$

$$\text{input: } x^{(t)} y^{(t)} = \{\pm 1\}, \quad \hat{y} = \text{sign}(w^T x + b) = \begin{cases} +1 & \text{if } w^T x + b \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

若 misclassify - 个 +1 例子 ($y^{(t)} = +1, \hat{y} = -1$):

$$w \leftarrow w + x^{(t)} \quad b \leftarrow b + 1$$

若 misclassify - 个 -1 例子 ($y^{(t)} = -1, \hat{y} = +1$):

$$w \leftarrow w - x^{(t)} \quad b \leftarrow b - 1$$

为何这样更新? 若: $b + w^T x < 0$, $y^{(t)} = +1$, 欲 $w^T x \uparrow$.

$$w \leftarrow w + x^{(t)}, \quad w^T x = (w + x^{(t)})^T x = w^T x + \underbrace{x^{(t)T} x}_{>0} > 0, \uparrow$$

[反之亦然] 同时 b 也 \uparrow , $b + w^T x \uparrow$

而更新逻辑可简化为:

$$\begin{cases} w \leftarrow w + y^{(t)} x^{(t)} \\ b \leftarrow b + y^{(t)} \end{cases}$$

可见, w 可视为 $x^{(1)}, \dots, x^{(n)}$ 的线性组合



同时还有一个 trick: $\vec{x} = \begin{bmatrix} 1 \\ \vec{x} \end{bmatrix}$ $\vec{\theta} = \begin{bmatrix} b \\ \vec{w} \end{bmatrix}$. 则:

$$\hat{y} = \text{sign}(\vec{w}^T \vec{x} + b) = \text{sign}(\vec{\theta}^T \vec{x})$$

$$\theta \leftarrow \theta + y^{(t)} \vec{x}^{(t)}, \text{ where } \vec{x}^{(t)} = \begin{bmatrix} 1 \\ \vec{x}^{(t)} \end{bmatrix}$$

那么 Perceptron 的 Inductive Bias 是什么? Decision Boundary should be linear, i.e., hyperplane, 且 Recent mistakes are more important than older ones.

而当 θ, b 训练到不断产生正确预测时, 我们就认为感知机收敛 (converge) 了。同时, Perceptron 也可能受 overfitting 影响。之前介绍到, 学的其实是一个超平面; 在 binary 分类中, 如果两类点之间的分隔不是“面”而是“space”呢?



Def: 对二分类问题, 样本集 S is linearly separable if there exists a linear boundary that separate the points

“space”大小又引出下面定义: Def: The margin γ for dataset D is the greatest possible distance between a linear separator and the closest point in D to the linear separator

Theorem: 若 Data margin 为 γ . 所有点都在半径为 R 的球内, 则 online perceptron 算法 makes $\leq (R/\gamma)^2$ mistakes.*

那么 margin 如何计算? 相当于是 $(\vec{x}'' - \vec{x}')$

向量在 $\vec{w}/\|\vec{w}\|_2$ 上的投影长度 (i.e., 内积)

$$\text{则 } \left| \frac{\vec{w}^T (\vec{x}'' - \vec{x}')}{\|\vec{w}\|_2} \right| = \frac{|\vec{w}^T \vec{x}'' - \vec{w}^T \vec{x}'|}{\|\vec{w}\|_2}$$

$$= \frac{|\vec{w}^T \vec{x}'' + b|}{\|\vec{w}\|_2}$$

