

Linear Classification

Training Set = $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$
 $\in \mathbb{R}^{d+1}$ ← $\left\{ \begin{matrix} x \\ \text{add } 1 \end{matrix} \right.$ $\rightarrow \in \{+1, -1\}$

Goal: Compute a hypothesis $h_\theta : \mathbb{R}^{d+1} \rightarrow \{+1, -1\}$
 $(\theta \in \mathbb{R}^{d+1})$ such that

$$h_\theta(x) = \begin{cases} +1 & \text{if } x^T \theta > 0 \\ -1 & \text{if } x^T \theta \leq 0. \end{cases}$$

$\rightarrow \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_d x_d$

and

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n l(\theta, x^{(i)}, y^{(i)})$$

is minimized

0-1 Loss function

$$l_{01}(\theta, x, y) = \begin{cases} 1 & \text{if } h_\theta(x) \neq y \\ 0 & \text{o.w.} \end{cases}$$

- Very difficult to compute (NP-hard)

Desire: A loss function that is "continuous" so that we can use SGD.

$$l_h(\theta, x, y) = 1[y \theta^T x \leq 0] \quad \rightarrow [1(b) = \begin{cases} 0 & \text{if } b = \text{false} \\ 1 & \text{if } b = \text{true} \end{cases}$$

$\forall b$ $l_{01}(\theta, x, y) = 1$ then $l_h(\theta, x, y) = 1$

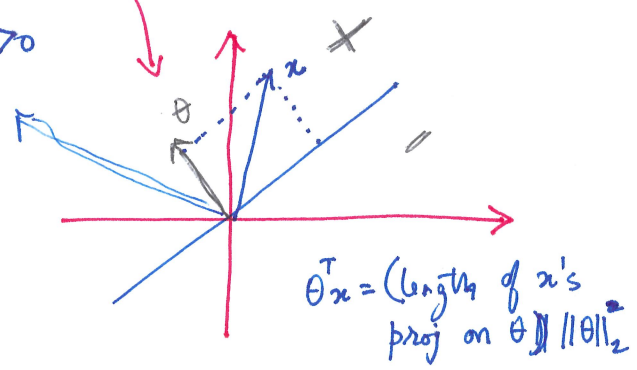
$\forall b$ $\theta^T x = 0$ then $l_h(\theta, x, y) = 1$ no matter what y is.

For x s.t. $\theta^T x \neq 0$, $l_{01}(\theta, x, y) = l_h(\theta, x, y)$.

$$l(\theta, x, y) = \begin{cases} -y \theta^T x \\ 0 \end{cases}$$

$$\begin{cases} y \theta^T x \leq 0 \\ y \theta^T x > 0 \end{cases}$$

$$\nabla_{\theta} l(\theta, x, y) = \begin{cases} -yx & y \theta^T x \leq 0 \\ 0 & y \theta^T x > 0 \end{cases}$$



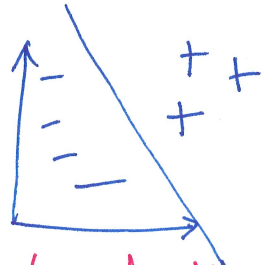
Perceptron: $\theta = 0$
Repeat until θ does not change

for $i \in \{1, 2, \dots, n\}$

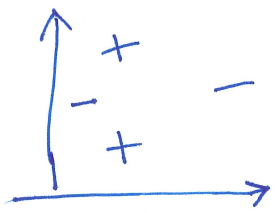
$$\theta \leftarrow \theta - \alpha \frac{1}{n} \nabla_{\theta} l(\theta, x^{(i)}, y^{(i)})$$

~~if~~ if $(y \theta^T x \leq 0)$
 $\theta \leftarrow \theta + y^{(i)} x^{(i)}$

Linear Separability: A training set is linearly separable if there is a hyperplane that separates the +1 and -1 examples.



Linearly separable



Not linearly separable.

$\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ is linearly separable if there is θ s.t.
 $y^{(i)} \theta^T x^{(i)} > 0 \quad \forall i \in \{1, 2, \dots, n\}$.

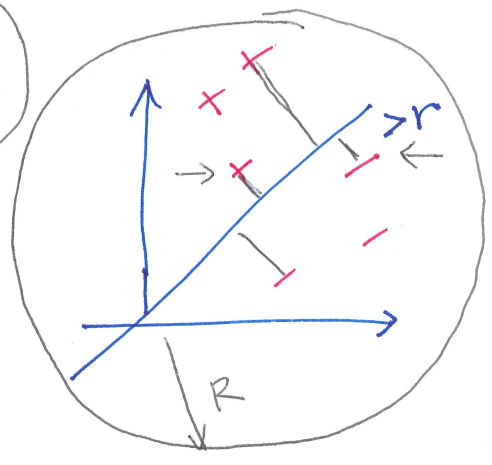
Signed distance of x w.r.t θ : $\frac{\theta^T x}{\|\theta\|_2}$

If θ is such that $y^{(i)} \theta^T x^{(i)} > 0 \quad \forall i \in \{1, 2, \dots, n\}$ then distance of $x^{(i)}$ from θ is $\frac{y^{(i)} \theta^T x^{(i)}}{\|\theta\|_2}$

Margin: Let θ be such that $y^{(i)} \theta^T x^{(i)} > 0$. The margin of θ w.r.t. $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$

$$= \min_{i \in \{1, 2, \dots, n\}} \frac{y^{(i)} \theta^T x^{(i)}}{\|\theta\|_2}$$

$$\frac{y^{(i)} \theta^T x^{(i)}}{\|\theta\|_2}$$



Perceptron Theorem: Suppose $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ that is linearly separable, i.e., $\exists \theta$ s.t. $y^{(i)} \theta^T x^{(i)} > 0 \quad \forall i$.

Let $\forall i \quad \frac{y^{(i)} \theta^T x^{(i)}}{\|\theta\|_2} > \gamma$.

Let us assume $R \geq \max_{i \in \{1, 2, \dots, n\}} \|x^{(i)}\|_2$

Then the number of times θ is changed in the perceptron algorithm is $\leq \left(\frac{R}{\gamma}\right)^2$.

Logistic Regression

$$l_{\log}(\theta, x, y) = \underbrace{\ln}_{\text{natural log } \log_e} (1 + \underbrace{e^{-y\theta^T x}}_{\text{euler's constant}})$$