

Binary classification

kernel : $\underset{\mathcal{X}}{\text{Data space}} \times \underset{\mathcal{X}}{\text{Data space}} \rightarrow \mathbb{C}$
(Hilbert space or \mathbb{R}^n)

- measures similarity in data

kernel methods: nonlinear generalizations
of algorithms that have dot product
forms

- Bayes decision rule / classifier
- Perceptron

Recap

Given

Data: $x_1, \dots, x_n \in \mathcal{X}$

y_1, \dots, y_n

$y_i = f(x_i)$ (true function)
target

Binary classification: $y_i = \pm 1$

Linear model:

$$h(x) = \text{sgn}(w \cdot x + b)$$

learned classifier

$$x = [x^{(1)}, \dots, x^{(d)}, 1]$$

$$w = [w, b]$$

$$h(x) = \text{sgn}(w \cdot x) \quad \begin{matrix} \text{(simplification:} \\ \text{affine} \rightarrow \\ \text{linear in} \\ \mathbb{R}^{d+1}) \end{matrix}$$

Bayes classifier

based on prob. assumptions on data

$y = \pm 1$ with eq. prob

$$P_Y(1) = P_Y(-1) = \frac{1}{2}$$

$P_{X|Y}(x|Y=1), P_{X|Y}(x|Y=-1)$

class conditional densities

$$P_{X|Y}(x_i|Y=1) = \begin{cases} 0 & y_i = -1 \\ \frac{1}{n_+} & f(x_i) = y_i = 1 \end{cases}$$

$$S^+ = \{i : i \in [n], f(x_i) = 1\}$$

$$S^- = \{i : i \in [n], f(x_i) = -1\}$$

$$n^+ = |S^+| \quad n^- = |S^-|$$

(# the pts)

$$P_{X|Y}(x_i|Y=-1) = \begin{cases} 0 & y_i = 1 \\ \frac{1}{n_-} & y_i = -1 \end{cases}$$

Bayes rule

$$h(x) = \begin{cases} 1 & P_{Y|X}(Y=1|X=x) > P_{Y|X}(Y=-1|X=x) \\ -1 & P_{Y|X}(Y=-1|X=x) > P_{Y|X}(Y=1|X=x) \end{cases}$$

$$P_{Y|X}(Y=1|X=x) = \frac{P_Y(Y=1) \cdot \bar{P}_{X|Y}(X=x|Y=1)}{P_X(x)}$$

$$P_{Y|X}(Y=-1|X=x) = \frac{P_Y(Y=-1) \cdot \bar{P}_{X|Y}(X=x|Y=-1)}{P_X(x)}$$

$$P_Y(Y=1) = P_Y(Y=-1) = \frac{1}{2}$$

$$h(x) = \begin{cases} 1 & P_{X|Y}(X=x|Y=1) > P_{X|Y}(X=x|Y=-1) \\ -1 & \text{o.w.} \end{cases}$$

Total probability

$$P_{X|Y}(X=x|Y=1) = \int \bar{P}_{X|Y, X'}(X=x|Y=1, X'=x') \frac{\bar{P}_{X'|Y}(X'=x'|Y=1)}{P_X(x)} dx'$$

\therefore

$$= \int x(x, x') P_{X|Y}(x'|Y=1) dx'$$

(expectation of $x(x, \cdot)$ w.r.t $P_{X|Y}(\cdot|Y=1)$)

$$\approx \frac{1}{n_+} \sum_{i \in S_+} x(x, x_i)$$

$$\approx \sum_{i=1}^n x(x, x_i) \underbrace{P_{X|Y}(x_i|Y=1)}_{P_{X|Y}(x|Y=1)}$$

$$h(x) = \begin{cases} 1 & \frac{1}{n_+} \sum_{i \in S_+} x(x, x_i) > \frac{1}{n_-} \sum_{i \in S_-} x(x, x_i) \\ -1 & \text{o.w.} \end{cases}$$

"Geometric" point of view

$$m_+ = \frac{1}{n_+} \sum_{i \in S_+} x_i \quad m_- = \frac{1}{n_-} \sum_{i \in S_-} x_i$$

(mean of the pts)

$$h_s(x) = \begin{cases} 1 & d(x, m_+) < d(x, m_-) \\ -1 & \text{o.w.} \end{cases}$$

For $x \in \mathbb{R}^n$,

$\therefore x \mapsto h_s(x)$ turns out to be a linear model

$$\begin{aligned} & \downarrow \\ & \text{---} \\ & \text{---} \\ & \text{---} \end{aligned}$$

$$h_s(x) = \text{sgn} \left((m_+ - m_-) \cdot x + \frac{1}{2} (\|m_-\|^2 - \|m_+\|^2) \right)$$

$$= \text{sgn} \left(\frac{1}{n_+} \sum_{i \in S_+} x_i \cdot x - \frac{1}{n_-} \sum_{i \in S_-} x_i \cdot x + b \right)$$

$h_s(x)$ is in

• Dot product form

• kernelizing: replacing dot product with kernel evaluations.

$$h_s^x(x) = \text{sgn} \left(\frac{1}{n_+} \sum_{i \in S_+} x(x_i, x) - \frac{1}{n_-} \sum_{i \in S_-} x(x_i, x) + \frac{1}{2} \left(\frac{1}{n_-^2} \sum_{i,j \in S_-} x(x_i, x_j) - \frac{1}{n_+^2} \sum_{i,j \in S_+} x(x_i, x_j) \right) \right)$$

$$\|m_-\|^2 = m_- \cdot m_- = \left(\frac{1}{n_-} \sum_{i \in S_-} x_i \right) \cdot \left(\frac{1}{n_-} \sum_{i \in S_-} x_i \right)$$

$$\rightarrow \frac{1}{(n_-)^2} \sum_{i,j \in S_-} x(x_i, x_j)$$

$$h_{\text{bayes}}(x) = \begin{cases} 1 & \frac{1}{n_+} \sum_{i \in S_+} x(x, x_i) > \frac{1}{n_-} \sum_{i \in S_-} x(x, x_i) \\ -1 & \text{o.w.} \end{cases}$$

$$h_{\text{bayes}}(x) = \text{sgn} \left(\frac{1}{n_+} \sum_{i \in S_+} x(x, x_i) - \frac{1}{n_-} \sum_{i \in S_-} x(x, x_i) \right)$$

• Mean-based geometric classifier h_s^x is nonlinear generalization of h_s (linear)

• h_s^x is also Bayes classifier when $b=0$

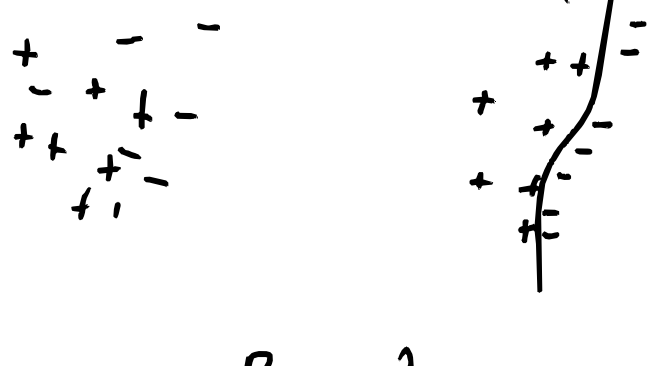
Linear separability

if true function $y_i = f(x_i)$

f has the form

$$f(x) = \text{sgn}(w \cdot x + b)$$

for some $w \in \mathcal{X}$



Perceptron algorithm

$$h_p(x) = \text{sgn}(w \cdot x) \quad (\text{basic neuron})$$

$$w_0 = 0 \in \mathcal{X}$$

If i^{th} data point, x_i , is misclassified, then,

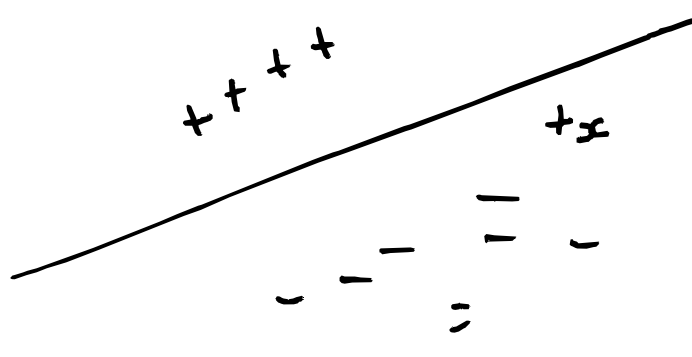
$$w_{i+1} = w_i + \eta y_i x_i$$

learning rate

O.W.

$$w_{i+1} = w_i$$

Intuition



$$\text{sgn}(w \cdot x) \neq y$$

$$\text{sgn}(y w \cdot x) = -1$$

$$y_i w_{i+1} \cdot x_i = y_i w_i \cdot x_i + \eta \|x_i\|^2$$

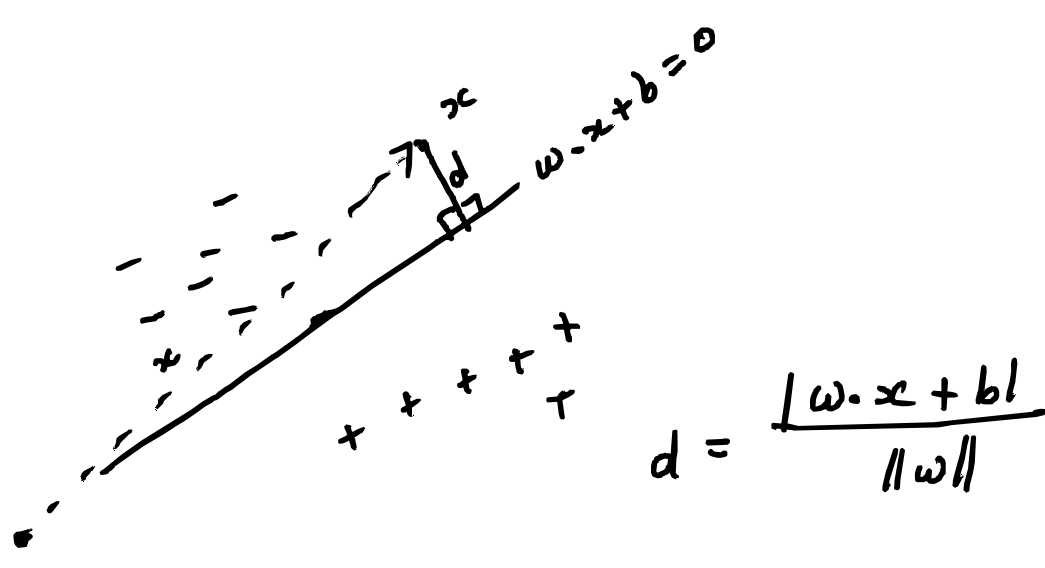
< 0 > 0

Convergence proof

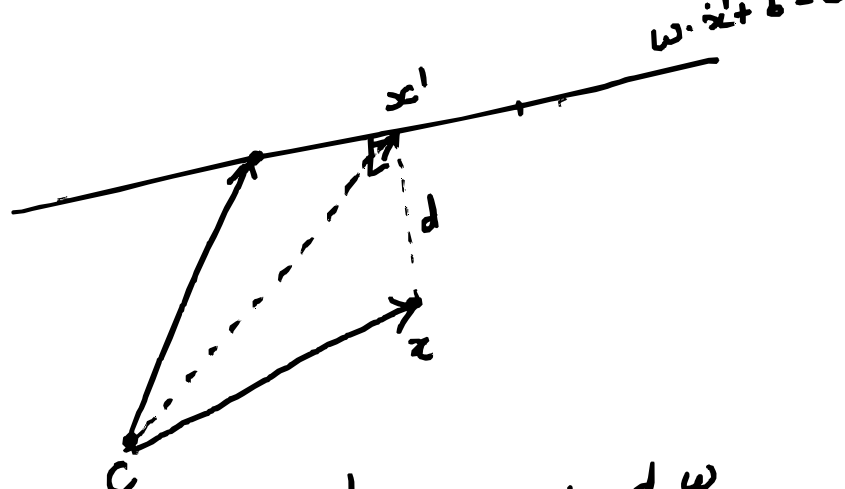
If $(x_i, y_i)_{i=1}^n$ are linearly separable, and $\|x_i\| < R \ \forall i$, then, the total number of updates, U , over n points satisfies (# of errors)

$$U \leq \frac{R^2}{\rho^2}$$

where $\rho = \min_i y_i \frac{w^* \cdot x_i}{\|w^*\|} > 0$ for some $w^* \in \mathcal{X}$.



$$d = \frac{|w \cdot x + b|}{\|w\|}$$



$$x' = x + d \frac{w}{\|w\|}$$

$$x' \cdot w + b = x \cdot w + d \|w\| + b$$

$$d = \left| \frac{(x \cdot w + b)}{\|w\|} \right|$$

Consequence of linear separability

$$\exists w \in \mathcal{X} \text{ s.t. } \min_i \frac{|w \cdot x_i + b|}{\|w\|} > 0$$

when $b = 0$,

$$\exists w \in \mathcal{X} \text{ s.t. } \min_i \frac{|w \cdot x_i|}{\|w\|} > 0$$

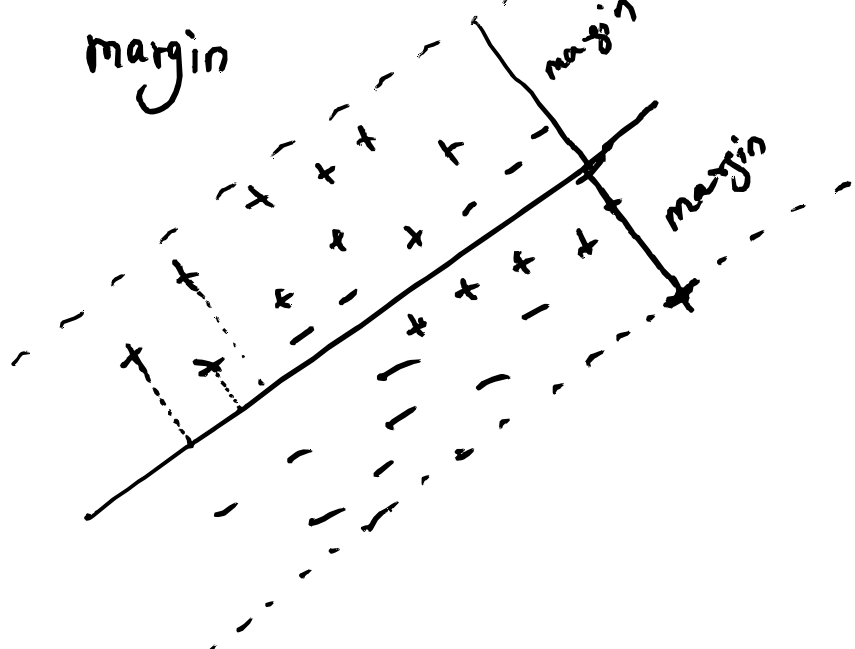
For a correctly classified pt,

$$\frac{y_i w \cdot x_i}{\|w\|} = \frac{|w \cdot x_i|}{\|w\|} = \text{distance of } x_i \text{ from plane.}$$

Under linear separability,

$$\frac{\rho}{\|w\|} = \min_i \frac{|w \cdot x_i + b|}{\|w\|} = \min_i \frac{y_i w \cdot x_i}{\|w\|}$$

margin



$$\rho = \min_i y_i \frac{w \cdot x_i}{\|w\|}$$

$$\begin{array}{c} \uparrow \\ U \rho \leq \left| \sum_{i \in I} y_i \frac{w_i \cdot x_i}{\|w_i\|} \right| \end{array}$$

of updates
or # of errors

set of
incorrect
of indices

(Cauchy
Schwarz)

$$\begin{aligned} &\leq \left\| \sum_{i \in I} y_i x_i \right\| \\ (\omega_{i+1} &= \omega_i + \eta y_i x_i) \\ &\leq \left\| \frac{1}{\eta} \sum_{i \in I} (\omega_{i+1} - \omega_i) \right\| \\ &= \frac{1}{\eta} \|\omega_{n+1}\| \\ &= \frac{1}{\eta} \sqrt{\|\omega_{n+1}\|^2} \end{aligned}$$

$$= \frac{1}{\eta} \sqrt{\sum_{i \in I} (\|\omega_{i+1}\|^2 - \|\omega_i\|^2)}$$

$$\leq \frac{1}{\eta} \sqrt{\sum_{i \in I} (\|\omega_i + \eta x_i y_i\|^2 - \|\omega_i\|^2)}$$

$$= \frac{1}{\eta} \sqrt{\sum_{i \in I} (\eta^2 \|x_i\|^2 + 2 \omega_i \cdot x_i y_i)}$$

$$\leq \sqrt{\sum_{i \in I} \|x_i\|^2}$$

$$\leq \sqrt{UR^2} = \sqrt{U} R$$

$$U \rho \leq \sqrt{U} R$$

$$\downarrow \quad U \leq \frac{R^2}{\rho^2}$$

iterations for convergence : independent