

LINEAR CLASSIFIER

$H =$ Space of hypotheses.

$=$ Half-spaces or Linear Separators.

$$\theta_1 x_{i1} + \theta_2 x_{i2} + \dots + \theta_d x_{id} \leq b$$

$$\text{sign}(\vec{\theta} \cdot \vec{x}_i - b) = l_i, \quad l_i \in \{+1, -1\} \quad i=1, \dots, n.$$

$$\text{Projective Space} \Rightarrow \vec{w} = (\vec{\theta}, -b) \begin{cases} \vec{a}_i = (\vec{x}_i, +1) & \text{if } l_i = +1 \\ \vec{a}_i = (-\vec{x}_i, -1) & \text{if } l_i = -1 \end{cases}$$

$$\text{sign}(\vec{w} \cdot \vec{a}_i) = \begin{cases} \text{sign}(\vec{\theta} \cdot \vec{x}_i - b) = 1 & \text{if } l_i = +1 \\ \text{or } \text{sign}(-\vec{\theta} \cdot \vec{x}_i + b) = -\text{sign}(\vec{\theta} \cdot \vec{x}_i - b) \\ = 1 & \text{if } l_i = -1. \end{cases}$$

Given a sample

$$S' = \{(\vec{x}_1, l_1), (\vec{x}_2, l_2), \dots, (\vec{x}_n, l_n)\} \subseteq \mathbb{R}^d \times \{+1, -1\}$$

Rewrite

$$S = \{\vec{a}_1, \vec{a}_2, \dots, \vec{a}_n\} \subseteq \mathbb{R}^{d+1}, \quad |\vec{a}_i| = 1$$

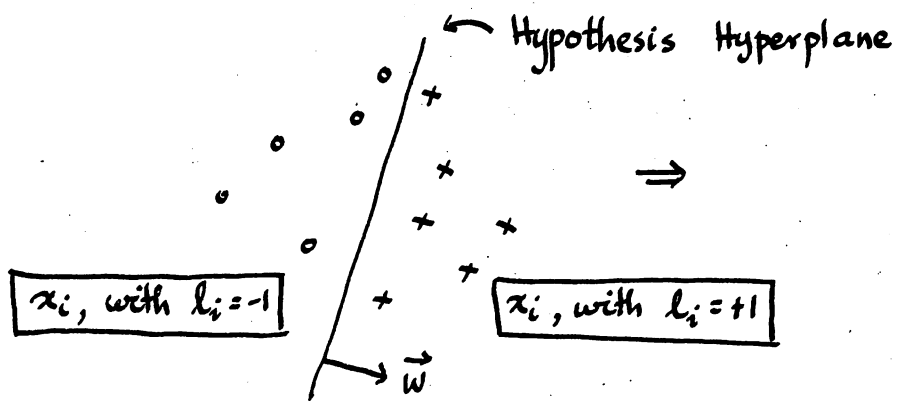
Find a solution vector $\vec{w} \in \mathbb{R}^{d+1}$, $|\vec{w}| = 1$ such that

$$\forall i \in \{1, n\} \quad \boxed{\vec{w} \cdot \vec{a}_i > 0} \quad \square.$$

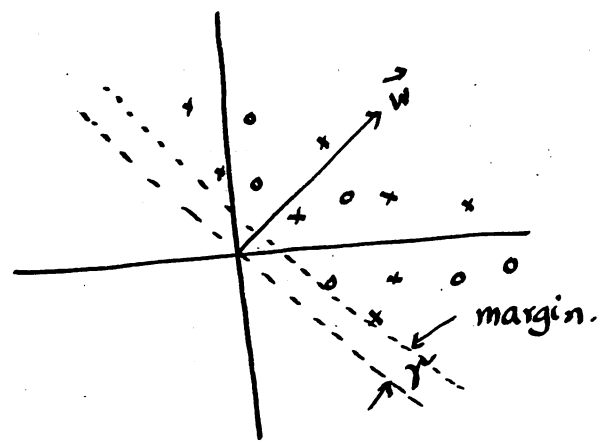
OR SOLVE THE FOLLOWING OPTIMIZATION PROBLEM.

$$\max_{\vec{w}^*} \min_j \frac{\vec{w}^* \cdot \vec{a}_j}{|\vec{w}^*| |\vec{a}_j|} = \text{margin}$$

CHECK IF THE MARGIN > 0 .



Transformed problem.



$$\gamma_P = \max_{\vec{w}^*} \min_i \frac{\vec{w}^* \cdot \vec{a}_i}{\|\vec{w}^*\|_2 \|\vec{a}_i\|_2}$$

$$\gamma_W = \max_{\vec{w}^*} \min_i \frac{\vec{w}^* \cdot \vec{a}_i}{\|\vec{w}^*\|_1 \|\vec{a}_i\|_\infty}$$

Check $\gamma_P > 0$ (iff) $\gamma_W > 0$.

Related by some function of d .
 $d = \text{const.}$

PROBLEMS MULTIPLICATIVE WEIGHTING.

Simulate a repeated game $w^* \approx$ Mixed Strategy (row player) } WINNOW ALGORITHM

$$M(i, j) = -a_{ij}$$

Problem 1. $w_i \geq 0$, weight must be positive

$$S = \{ (\vec{a}_1, -\vec{a}_1), \dots, (\vec{a}_n, -\vec{a}_n) \} \in \mathbb{R}^{2(d+1)}$$

$$\begin{aligned} \vec{w} \cdot \vec{a}_i &= w_1 a_{i1} + w_2 a_{i2} + \dots + w_{d+1} a_{i(d+1)} \\ &= (w_1^+ - w_1^-) a_{i1} + \dots + (w_{d+1}^+ - w_{d+1}^-) a_{i(d+1)} \end{aligned}$$

$$\vec{w} \cdot \vec{a}_i = \langle w_1^+, w_2^+, \dots, w_{d+1}^+, w_1^-, \dots, w_{d+1}^- \rangle \cdot \langle \vec{a}_i, -\vec{a}_i \rangle$$

The Algorithm

v_i = weight in the direction of the i^{th} feature

$$v_i^0 = 1 \quad \forall i = 1, \dots, d+1$$

Round t

$$\sigma_{i,t} = \frac{v_i^{t-1}}{\sum v_i^{t-1}}$$

$$w^t = \langle \sigma_{1,t}, \dots, \sigma_{d+1,t} \rangle$$

Check $\exists_{j_t} \vec{w}^t \cdot \vec{a}_{j_t} \leq 0$ [a_j = misclassified]

If ~~$\exists_{j_t} \vec{w}^t \cdot \vec{a}_{j_t} \leq 0$~~

If $\forall_j \vec{w}^t \cdot \vec{a}_j > 0 \Rightarrow$ SUCCESS.

Else

$$v_i^t = v_i^{t-1} (1 - \epsilon)^{-a_{ij_t}} \approx v_i^{t-1} (1 + \epsilon a_{ij_t})$$

(Neglect $O(\epsilon^2)$ terms)

$$(A) \sum_t M(\sigma_{r,t}, \sigma_{c,t}) = \sum_t \sum_i \sigma_{i,t} M(i, j_t) = -\sum_t \sum_i \sigma_{i,t} a_{ij_t}$$

$$= -\sum_t \vec{w}^t \cdot \vec{a}_{j_t} \geq 0$$

$$(B) \sum_t M(\sigma_{r,t}, \sigma_{c,t}) \leq \frac{\ln n}{\epsilon} + (1 + \epsilon) \sum_t M(\sigma_r^*, \sigma_{c,t})$$

$$= \frac{\ln n}{\epsilon} + (1 + \epsilon) \sum_t -w^* \cdot a_{j_t}$$

$$\leq \frac{\ln n}{\epsilon} + \epsilon T - r_w T$$

Take $\epsilon = \frac{3r_w}{4}$ (ϵ Needs to be searched by binary search) $\Rightarrow r_w \leq \epsilon \leq \frac{r_w}{2}$ (107)

$$\therefore \frac{4 \ln n}{3r_w} + \frac{3r_w}{4} T - r_w T \geq 0$$

$$\frac{r_w T}{4} \leq \frac{4 \ln n}{3r_w} \Rightarrow T \leq \frac{16}{3} \frac{\ln n}{r_w^2} \quad (\text{More correct estimate takes into account Binary Search})$$

$$T = O\left(\frac{\ln n}{r_w^2}\right)$$

A HEURISTIC APPROACH \rightarrow PERCEPTRON

Additive Update \rightarrow Instead of Multiplicative.

a) SIMPLE

b) NO need to worry about $w_i^t \geq 0$

The Algorithm.

$$w^0 \leftarrow (0, 0, \dots, 0)$$

Round t .

$$\text{Check } \exists_{j_t} \vec{w}^t \cdot \vec{a}_{j_t} \leq 0 \quad [\vec{a}_{j_t} = \text{misclassified}]$$

$$\text{If } \forall_j \vec{w}^t \cdot \vec{a}_j > 0 \Rightarrow \text{SUCCESS.}$$

Else

$$\vec{w}^{t+1} = \vec{w}^t + \vec{a}_{j_t}$$

$$\begin{aligned} \text{Note } \vec{w}^{t+1} \cdot \vec{a}_{j_t} &= \vec{w}^t \cdot \vec{a}_{j_t} + \vec{a}_{j_t} \cdot \vec{a}_{j_t} \\ &= \vec{w}^t \cdot \vec{a}_{j_t} + 1 > \vec{w}^t \cdot \vec{a}_{j_t} \end{aligned}$$

However, there is no guarantee that it will not fail for some sample point that was correctly classified so far!

Let's define a potential function

$$\Phi_t = \text{cosine}(\vec{\omega}^t, \vec{\omega}^*) = \frac{\vec{\omega}^t \cdot \vec{\omega}^*}{\|\vec{\omega}^t\|_2 \|\vec{\omega}^*\|_2} = \frac{N_t}{D_t}$$

← Numerator
← Denominator

$$N_{t+1} = \vec{\omega}^{t+1} \cdot \vec{\omega}^* = \vec{\omega}^t \cdot \vec{\omega}^* + \vec{a}_{j_t} \cdot \vec{\omega}^* \\ \geq N_t + r_p \|\vec{\omega}^*\|_2$$

$$\frac{D_{t+1}^2}{\|\vec{\omega}^*\|_2^2} = \|\vec{\omega}^{t+1}\|_2^2 = (\vec{\omega}^t + \vec{a}_{j_t}) \cdot (\vec{\omega}^t + \vec{a}_{j_t}) \\ = \|\vec{\omega}^t\|_2^2 + 2(\vec{\omega}^t \cdot \vec{a}_{j_t}) + \|\vec{a}_{j_t}\|_2^2 \\ \leq \frac{D_t^2}{\|\vec{\omega}^*\|_2^2} + 1.$$

$$\therefore N_T \geq T \cdot r_p \|\vec{\omega}^*\|_2$$

$$D_T \leq \sqrt{T} \|\vec{\omega}^*\|_2$$

$$\therefore \Phi_T = \text{cosine}(\vec{\omega}^T, \vec{\omega}^*) \geq \frac{T \cdot r_p}{\sqrt{T}} = \sqrt{T} r_p$$

$$1 \geq \Phi_T \geq \sqrt{T} r_p \quad T \leq \frac{1}{r_p^2}$$

$$T = O\left(\frac{1}{r_p^2}\right)$$

□

In this case Perceptron Algorithm does as well as the more sophisticated learning algorithm.

SUPPORT VECTOR MACHINE:

(109)

The two algorithms described earlier do not guarantee optimality of the margin computed.

Optimization Problem:

$$\gamma_P = \max_{\vec{w}^*} \min_i \frac{\vec{w}^* \cdot \vec{a}_i}{\|\vec{w}^*\|_2 \|\vec{a}_i\|_2} \quad \leftarrow \text{Nonconvex ...}$$

Since ~~the~~ we can normalize the sample points such that

$$\|\vec{a}_i\|_2 = 1$$

$$\gamma_P = \max_{\vec{w}^*} \min_i \frac{\vec{w}^* \cdot \vec{a}_i}{\|\vec{w}^*\|_2}$$

$$\text{Let } \vec{v} = \frac{\vec{w}}{\|\vec{w}\| \gamma_P}$$

An equivalent problem is

$$\text{minimize } |\vec{v}| \quad \text{subject to } \underbrace{\vec{v} \cdot \vec{a}_i \geq 1}_{\text{Linear constraints}} \quad \forall_i$$

↓
convex but not smooth

SVM - Large Margin Classifier Problem is

$$\text{minimize } |\vec{v}|^2$$
$$\text{subject to } \vec{v} \cdot \vec{a}_i \geq 1 \quad \forall_i$$

LAGRANGE DUALITY.

(110)

SVM \rightarrow Maximum Margin Classifiers

Constrained Optimization Problem.

$$\min_w f(w) \text{ subject to } h_i(w) = 0 \quad \forall i \ i=1, \dots, l.$$

Lagrangian

$$L(w, \beta) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

\hookrightarrow Lagrange Multipliers.

$$\frac{\partial L}{\partial w_i} = 0 \quad \& \quad \frac{\partial L}{\partial \beta_i} = 0 \quad \forall i \ i=1, \dots, l.$$

Generalized Lagrangian

$$\min_w f(w) \text{ s.t. } \begin{array}{ll} g_i(w) \leq 0 & \forall i=1, \dots, k \\ h_i(w) = 0 & \forall i=1, \dots, l. \end{array}$$

Generalized Lagrangian:

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

P: Primal \rightarrow

$$\Theta_P(w) = \max_{\alpha, \beta: \alpha_i \geq 0} L(w, \alpha, \beta)$$

$\forall w, w \neq \text{feasible}$

$$\begin{array}{l} g_i(w) > 0 \\ h_i(w) \neq 0 \end{array}$$

$$\Theta_P(w) = \max_{\alpha, \beta: \alpha_i \geq 0} L(w, \alpha, \beta)$$

$$= \max_{\alpha, \beta: \alpha_i \geq 0} f(w) + \sum_{i=1}^k \alpha_i g_i + \sum_{i=1}^l \beta_i h_i$$

$\rightarrow \infty$

$$\Theta_p(\omega) = \begin{cases} f(\omega) & \text{if } \omega \text{ satisfies primal constraints} \\ \infty & \text{(if } \omega = \text{infeasible)} \end{cases} \quad (111)$$

$$\begin{aligned} \min_{\omega} \Theta_p(\omega) &= \min_{\omega} \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(\omega, \alpha, \beta) \\ &= p^* \leftarrow \text{Value of the primal problem.} \end{aligned}$$

DUALITY

$$\Theta_D(\omega) = \min_{\omega} \mathcal{L}(\omega, \alpha, \beta)$$

$$\begin{aligned} \max_{\alpha, \beta: \alpha_i \geq 0} \Theta_D(\omega) &= \max_{\alpha, \beta: \alpha_i \geq 0} \min_{\omega} \mathcal{L}(\omega, \alpha, \beta) \\ &= d^* \leftarrow \text{Value of the dual problem.} \end{aligned}$$

$$\boxed{\max \min \leq \min \max}$$

$$d^* = \max_{\substack{\alpha, \beta \\ \alpha_i \geq 0}} \min_{\omega} \mathcal{L}(\omega, \alpha, \beta) \leq \min_{\omega} \max_{\substack{\alpha, \beta \\ \alpha_i \geq 0}} \mathcal{L}(\omega, \alpha, \beta) = p^*$$

THEOREM (Karush-Kuhn-Tucker: KKT)

$$\left. \begin{array}{l} 1) f \text{ and } g_i = \text{convex} \\ 2) h_i = \text{affine} \\ 3) g_i\text{'s are strictly feasible} \\ \exists \omega \ g_i(\omega) < 0 \text{ for all } i \end{array} \right\} \Rightarrow \begin{array}{l} d^* = p^* = \mathcal{L}(\omega^*, \alpha^*, \beta^*) \\ \omega^* = \text{sols to the primal problem} \\ \alpha^*, \beta^* = \text{sols to the dual problem} \end{array}$$

$$\Leftrightarrow \text{Satisfy KKT cond}^n\text{'s} \left\{ \begin{array}{l} \frac{\partial}{\partial \omega_i} \mathcal{L}(\omega, \alpha, \beta) \Big|_{\omega^*, \alpha^*, \beta^*} = 0 \quad \forall_i \\ \frac{\partial}{\partial \beta_i} \mathcal{L}(\omega, \alpha, \beta) \Big|_{\omega^*, \alpha^*, \beta^*} = 0 \quad \forall_i \end{array} \right. \Leftrightarrow \begin{array}{l} \alpha_i^* g_i(\omega^*) = 0 \quad \forall_i \\ g_i(\omega^*) \leq 0 \quad \forall_i \\ \alpha_i^* \geq 0 \quad \forall_i \end{array}$$

Optimal Margin Classifier:

$$\min \frac{1}{2} \|\vec{v}\|_2^2 \quad \text{s.t.} \quad \vec{v} \cdot \vec{a}_i \geq 1 \quad \forall_i$$

$$\mathcal{L}(\vec{v}, \alpha) = \frac{1}{2} \vec{v}^T \cdot \vec{v} - \sum_i \alpha_i [\vec{v} \cdot \vec{a}_i - 1]$$

$$\nabla_{\vec{v}} \mathcal{L} = \vec{v} - \sum_i \alpha_i \vec{a}_i = 0 \quad \Rightarrow \quad \vec{v} = \sum_i \alpha_i \vec{a}_i, \quad \alpha_i \geq 0$$

↳ Positive combination of \vec{a}_i 's

$$\begin{aligned} \mathcal{L}(\vec{v}, \alpha) &= \frac{1}{2} \left(\sum_i \alpha_i \vec{a}_i \right) \cdot \left(\sum_i \alpha_i \vec{a}_i \right) - \sum_i \alpha_i \left[\left(\sum_i \alpha_i \vec{a}_i \right) \cdot \vec{a}_i - 1 \right] \\ &= \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j (\vec{a}_i \cdot \vec{a}_j) - \sum_{i,j} \alpha_i \alpha_j (\vec{a}_i \cdot \vec{a}_j) + \sum_i \alpha_i \\ &= \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j (\vec{a}_i \cdot \vec{a}_j) \end{aligned}$$

DUAL

$$\Rightarrow \max_{\alpha} W(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j (\vec{a}_i \cdot \vec{a}_j)$$

subject to $\alpha_i \geq 0$

$$\text{Soln} = \alpha^*$$

$$\vec{v}^* = \sum_i \alpha_i^* \vec{a}_i = \frac{\vec{w}}{\|\vec{w}\|_2 \gamma_P}$$

$$\text{Recover } \vec{w} = (\vec{\theta}, -b) \rightarrow \vec{\theta} \cdot \vec{x}_i - b = l_i \quad \square$$