

Support vector machines

To separate 2 sets of points (cat pictures vs dog pictures, or spam vs not spam emails, ...)
 We aim to find a function

$$\begin{aligned} f(x_i) &> 0 \quad \text{for } x_i \in \text{Class 1} & i = 1, \dots, N \leftarrow \text{sample size} \\ f(x_i) &< 0 \quad \text{for } x_i \in \text{Class 2} \end{aligned}$$

Zero-level of f : $\{x \mid f(x)=0\}$ separates, or discriminates 2 classes

Linear separation function $f(x)$:

$$f(x) = a^T x - b$$

We seek a hyperplane to separate 2 classes

$$\textcircled{*} \quad \left[\begin{array}{l} a^T x_i - b > 0 \quad x_i \in \text{Class 1} \\ a^T x_i - b < 0 \quad x_i \in \text{Class 2} \end{array} \right] \quad \text{feasibility} \iff \left[\begin{array}{l} a^T x_i - b \geq 1 \quad x_i \in \text{Class 1} \\ a^T x_i - b \leq -1 \quad x_i \in \text{Class 2} \end{array} \right]$$

(since we can simultaneously rescale a and b)

Convex feasibility problem

$$\textcircled{\$} \quad \left[\begin{array}{l} \min_{a \in \mathbb{R}^n, b \in \mathbb{R}} 0 \\ a^T x_i - b \geq 0 \end{array} \right], \text{ where } X = \begin{bmatrix} x_1 & x_2 & \dots & x_m & x_1 & \dots & x_m \\ \underbrace{+1}_{N_1} & \underbrace{+1}_{N_1} & \dots & \underbrace{+1}_{N_1} & \underbrace{-1}_{N_2} & \dots & \underbrace{-1}_{N_2} \end{bmatrix}, \quad a_b = \begin{bmatrix} a_1 \\ \vdots \\ a_n \\ b \end{bmatrix} \quad (\text{restating})$$

When is it feasible?

Geometrically, when convex hulls $\text{conv} 1 \cap \text{conv} 2 = \emptyset$

$\text{conv} 1 = \text{conv} \{x_i \mid x_i \in \text{Class 1}\}$ do not intersect

$\text{conv} 2 = \text{conv} \{x_i \mid x_i \in \text{Class 2}\}$

Why? This follows from the strong alternatives to linear inequalities
 (recall Farkas lemma \rightarrow LP duality)

Thm Let $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$

$\exists x \in \mathbb{R}^n: \{Ax < b\}$ or strictly

$\exists \lambda \in \mathbb{R}^m: \{\lambda \neq 0, \lambda \geq 0, A^T \lambda = 0, \lambda^T b \leq 0\}$

(Boyd & Vanderberghe, §2.5; follows from a version of separability of convex sets theorem)
 p 50

Corollary \star is infeasible \Leftrightarrow

$$\left[\exists \lambda \in \mathbb{R}^N, \lambda \neq 0, \lambda \geq 0, \tilde{\lambda} \lambda = 0, 0 \leq 0 \right]$$

To restate, let $d = \begin{bmatrix} \bar{x} & | & \tilde{x} \\ -N_1 & & -N_2 \end{bmatrix}$

$$d_s = \begin{pmatrix} 1, \dots, 1 \\ -s \end{pmatrix} \in \mathbb{R}^S$$

$$\text{Then } \sum_{i=1}^N d_i \tilde{x}_i = 0 \Leftrightarrow \begin{cases} \sum_{i=1}^{N_1} \bar{d}_i \cdot \bar{x}_{ki} = \sum_{i=1}^{N_2} \tilde{d}_i \cdot \tilde{x}_{ki} \\ \sum_{i \in S} \bar{d}_i = \sum_{i \in S} \tilde{d}_i \end{cases}$$

We could rescale all d_i by the same number, e.g. making $\sum \bar{d}_i = \sum \tilde{d}_i = 1$
This implies the convex hull condition.

What if the problem is not feasible?

One could minimize the number of misclassified points instead

$$\begin{cases} \min_{a, b, \eta} \| \eta \|_0 \\ a, b, \eta \\ y_i (a^T x_i - b) \geq 1 - \eta_i \\ \eta_i \geq 0 \end{cases}$$

where $\eta \in \mathbb{R}^N$

$\eta_i = 0$ if a point is correctly classified	}
$\eta_i \in (0, 1)$ if a point is still correctly classified, but "uncertain"	
$\eta_i > 1$ if a point is misclassified	

and $y_i = +1$ if $x_i \in \text{Class 1}$
 $y_i = -1$ if $x_i \in \text{Class 2}$

Non-convex since $\| \cdot \|_0$ is non-convex

$\| \cdot \|_0 \rightarrow \| \cdot \|_1$ relaxation (convex envelope)

$$\begin{cases} \min_{a, b, \eta} \| \eta \|_1 \\ a, b, \eta \\ y_i (a^T x_i - b) \geq 1 - \eta_i \\ \eta_i \geq 0 \end{cases}$$

Feasible case: seeking the most robust solution from many

To pick one of them, one could specify the optimization task:

(1)

$$\max_{a, b, t} t$$

$$y_i(a^T x_i - b) \geq t$$

$$\|a\| \leq 1$$

"Max margin classifier"

we maximize the "gap" between the data and the separating hyperplane

normalization is required, otherwise $a \nearrow \infty$ increases $t \nearrow \infty$

1. Interpretation

This promotes robustness:

if a point's location is slightly (ϵ) off, it will be classified correctly

t maximizes the distance to $\{a^T x = b\}$

Why? Lemma

$$\text{dist}(v, a^T z = b) = \frac{|a^T v - b|}{\|a\|}$$

$v \in \mathbb{R}^n$ is a point

Proof of the lemma follows from

(Thm) Corollary of the optimality Thm: $\min_{x \in S \subseteq \mathbb{R}^n} f(x)$ convex $\Rightarrow x$ is optimal $\Leftrightarrow \nabla f(x)^T (y - x) \geq 0 \quad \forall y \in S$

$\begin{bmatrix} \min f(x) \\ Ax = b \end{bmatrix} \quad f \text{ is convex, } A \in \mathbb{R}^{m \times n}. \text{ For a feasible point } x,$

x is optimal $\Leftrightarrow \exists \mu: \nabla f(x) = A^T \mu.$

(Pf) Take $x: Ax = b$

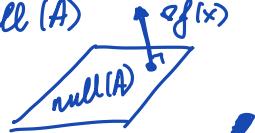
x is optimal $\Leftrightarrow \nabla f(x)^T (y - x) \geq 0 \quad \forall y: Ay = b$

$\Leftrightarrow \forall z \in \text{null}(A) \quad \nabla f^T(x) \cdot z \geq 0$

$\Leftrightarrow \forall z \in \text{null}(A) \quad \nabla^T f(x) \cdot z \leq 0 \quad (\text{consider } -z, \text{ it is also in nullspace})$

$\Leftrightarrow \nabla f(x)$ is orthogonal to $\text{null}(A)$

$\Leftrightarrow \nabla f(x) \in \text{row space of } A$



Pf) of the Lemma

Consider an optimization problem

$$\begin{bmatrix} \min \|v - z\|_2^2 \\ a^T z = b \end{bmatrix}$$

$$\left[\begin{array}{l} \min \sum_{i=1}^n v_i^2 - 2v_i z_i + z_i^2 \\ a^T z = b \end{array} \right]$$

$$\nabla f = 2z - 2v$$

$$\nabla^T f = a \mu \quad (\text{by Thm above})$$

$$2(z-v) = a \mu$$

$$\mu = \frac{2a^T(z-v)}{\|a\|^2}$$

$$\|a\|^2 = a^T a$$

Then,

$$\begin{aligned} \|v-z\| &= \frac{\|a\|}{2} |\mu| = \frac{\|a\|}{2} \frac{2|a^T(z-v)|}{\|a\|^2} \\ &= \frac{|b-a^T v|}{\|a\|^2} \end{aligned}$$

2. Properties

Thm (Properties of the max margin classifier)

① Optimal value of $\|a\|$ is achieved when $\|a\|=1$

② Solution to ① is unique

③ It is equivalent to the following problem

$$\left[\begin{array}{l} \min \|a\| \\ \text{a, b} \\ y_i(a^T x_i - b) \geq 1 \end{array} \right] \quad (2)$$

(Proof) ① If (t^*, a^*, b^*) gives optimal solution and $\|a^*\| \neq 1$, then

$$\text{consider } a := \frac{a^*}{\|a^*\|}, \quad b := \frac{b^*}{\|a^*\|}, \quad t = \frac{t^*}{\|a^*\|} > t^*$$

② Follows from ③: $\min \|a\| \Leftrightarrow \min \|a\|^2$ which is strictly convex, so min is unique in a.

if we have 2 ^{optimal}solutions (a, b^*) and (a, b^{**}) , consider $b = \frac{b^* + b^{**}}{2}$

let $b^{**} < b^*$ (without loss of generality), then $b^{**} < b < b^*$

$$y_i(a^T x_i - b^*) \geq 1 \quad \forall y_i = 1 \Rightarrow y_i(a^T x_i - b) > 1 \quad \forall y_i = 1$$

$$y_i(a^T x_i - b^*) \leq 1 \quad \forall y_i = -1 \Rightarrow y_i(a^T x_i - b) < 1 \quad \forall y_i = -1$$

Then we can shrink both a and b to improve on the solution (exercise: check this formally)

by defining $\epsilon := \min_i |1 - y_i(a^T x_i - b)|$

③ Equivalence: 2 problems are feasible at the same time

a)

If $a = 0$ for any feasible $(0, b, t)$ or (a, b) :

$-y_i b \geq t$ $\forall i$.

If there are indeed 2 classes ($y_i = +1$ and $y_i = -1$)
this leads to a contradiction
unless $t = 0$

So, both problems will have solution zero.

b) If not,

$$(a, b) \rightarrow \left(\frac{a}{\|a\|}, \frac{b}{\|a\|}, \frac{t}{\|a\|} \right)$$
$$\left(\frac{a}{t}, \frac{b}{t} \right) \leftarrow (a, b, t)$$

so, 2 problems are
feasible together,
 $t = \frac{1}{\|a\|}$,
so minimize $\|a\|$ is
equivalent to maximize t .

Note: problem (2) is easier for the analysis and has less variables
(e.g., uniqueness)

Infeasible case again: putting all together

① Support vector classifier

$$\begin{cases} \min_{a, b, \eta} \|a\| + \gamma \sum \eta_i \\ \text{such that } y_i (a^T x_i - b) \geq 1 - \eta_i \\ \eta_i \geq 0 \end{cases}$$

Minimizing (a) margin width and
(b) the number of
misclassifications
at the same time

γ is a parameter
(sets the relative importance
of the goals (a) and (b))

Margin width: $\frac{2}{\|a\|}$ (easy exercise by using the material above)

Other approaches to handle non-separability:

- ② Bayesian: define a model based on observations, via maximum likelihood
(logistic modeling)
- ③ Non-linear separating functions (e.g., quadratic)