

Support Vector Machines

- 1 Support Vectors & Margins 517 - 522 5
- 2 Linear & Separable Data 523 - 527 4
- 3 Linear & non-separable data 527 - 532 5
- 3.2 Quadratic Loss 532 - 533 1
- 4 Non-linear Data 533 - 537 4
- 5 Stochastic Gradient Ascent 537 - 543 6
25

Support Vectors & Margins 5.1

D = Dataset w/ n points x_i
 in d-dimensional space,
 with class labels $y_i \in \{+1, -1\}$

Hyperplanes

A hyperplane in d-dimensions is given as
 a set of all points $x \in \mathbb{R}^{d \text{ (cols)}}$ that
 satisfy the equation $\pi h(x) = 0$

(2)

where $h(x)$ is a hyperplane function

$$h(x) = w^T x + b = w_1 x_1 + \dots + w_d x_d + b$$

w = weight vector, b = scalar bias

Points that create the hyperplane =

$$h(x) = w^T x + b = 0$$

Also defined as set of points that:

$$w^T x = -b$$

recall, the weight vector is orthogonal to the hyperplane $h(x) = 0$

- orthogonal means perpendicular

- $\frac{b}{w_i}$ = the offset where the hyperplane intersects the i th dimension ($w_i \neq 0$)

weight vector specifies the direction that is normal to the hyperplane

b fixes the offset to the hf in d-space

(3)

Separating Hyperplane

- A hp splits D-space into 2 half-spaces
- A dataset is said to be linearly separable if each half contains one class
- if $\text{lin sep} = h(x) = 0 : y_1 = 1 = h(x) > 0$
 $\& y_2 = 0 > 0$
- In this case the hp equation itself is a linear ~~stiff~~ classifier or linear discriminant
- That is to say it predicts class

$$y = \begin{cases} +1 & \text{if } h(x) > 0 \\ -1 & \text{if } h(x) \leq 0 \end{cases}$$

Distance of a point to the HP

Consider a data points $x \in \mathbb{R}^d$ which does not sit on the hyperplane

let x_p be the point x 's orthogonal projection onto the hyperplane
(Point on the hp)

Let this distance be r

$$r = x - x_p \quad (\text{the line})$$

This means:

$$x = x_p + r \quad (\text{transform } x_p \text{ to } x \text{ w/r})$$

$$x = x_p + r \frac{w}{\|w\|}$$

r (lowercase) is the directed distance of the point to the HP

r gives the offset of x from x_p in the unit weight vector
 $\frac{w}{\|w\|}$ = unit weight vector

r = Posit if r is in same dir as w

Neg if r is in opp dir to w

Explained

(5)

Example:

$$HP \ 2D = x - 2y + 3 = 0$$

$$\text{Point} = (4, 2)$$

goal = Express $x = (4, 2)$ as

$$x = xp + r \frac{\omega}{\|\omega\|}$$

where xp = projection of x onto the line

ω = weight vector (orth to HP)

r = signed distance (+/-)

Calculations

- $\omega = (1, -2)$ (coeff of HP)
- Projection from x onto the hyp = xp
to find = point on line that is closest
to x . Done using vector proj or
solving linear system of equations
- $xp = (2, 2, 25)$
- Signed Distance (r) calced as dot prod
$$r = (x - xp) \cdot \frac{\omega}{\|\omega\|}$$
 $\sqrt{1^2 + (-2)^2} \quad \|\omega\| = \sqrt{1^2 + (-2)^2}$

$r = x \Rightarrow x_p = \text{vector}$
 Back to textbook: one unit of orth

$$x = x_p + r$$

$$x = x_p + \left(r \frac{w}{\|w\|} \right)$$

Scalar of signed dist

- This is not a replace/plugin for r
- Instead it is meant to express x in terms of its projection & distances

How to calculate distance r ?

$$① x = x_p + r \cdot \frac{w}{\|w\|}$$

$$② \text{sub } x_p = x - x_p = r \cdot \frac{w}{\|w\|}$$

$$③ \text{Dot Prod both sides by } \frac{w}{\|w\|}$$

$$(x - x_p) \cdot \frac{w}{\|w\|} = r * \left(\frac{w}{\|w\|} \cdot \frac{w}{\|w\|} \right)$$

$$④ \text{Since } \frac{w}{\|w\|} = \text{unit vector w/mag=1 then } r = 1$$

$$⑤ \text{final} = r = (x - x_p) \cdot \left(\frac{w}{\|w\|} \right)$$

Note: w vector always points positive
 hence the need for a direction
 Distance to cover the neighborhood
 Space codes.

7

The reason we present x as a eq
 $d \quad x_p + r \frac{\omega}{\|\omega\|}$ is because we now
 plug this is the hyper-plane fun

$$HP = h(x) = \underline{w^T x + b} = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b$$

$$h(x) = h\left(x_p + r \frac{\omega}{\|\omega\|}\right)$$

$$= w^T \left(x_p + \frac{\omega}{\|\omega\|}\right) + b$$

$$= \underbrace{w^T x_p + b}_{h(x_p)} + r \frac{w^T \omega}{\|\omega\|}$$

$$= \underbrace{h(x_p)}_{r} + r \|\omega\| \quad h(x) = r \|\omega\|$$

Thus, the directed distance of a point to HP =

$$r = \frac{h(x)}{\|\omega\|}$$

To obtain Distance (no direction & thus non-neg) we multiply r by class label $\{h(x) \leq^{\circ} -1, h(x) \geq^{\circ} +1\}$

the distance of a point to the HP:

$$\delta = y \cdot r = \frac{y \cdot h(x)}{\|w\|}$$

for origin the directed Distances to the hyperplane =

$$r = h(0) = \frac{w^T 0 + b}{\|w\|} = \frac{b}{\|w\|}$$

(9)

Example:

2-dimensional, hyperplane = line, $X = (x_1, x_2)^T$

$$h(x) = w^T x + b = w_1 x_1 + w_2 x_2 + b = 0$$

Pearrange to isolate x_2 : Slope-intercept

$$x_2 = -\frac{w_1}{w_2} x_1 - \frac{b}{w_2}$$

- why isolate w_2 ? $w_2 = x$ -coordinate. Rearrange for $y = mx + b$ which is a typical representation of line, i.e. y ~~given~~ relationship to x

Here shows • $\frac{w_1}{w_2}$ to be the slope γ

- $\frac{b}{w_2}$ to be the intercept along 2nd Dimension (y)

Consider two points on the HP:

$$P = (P_1, P_2) = (4, 0) \quad Q = (q_1, q_2) = (2, 5)$$

$$-\frac{w_1}{w_2} = \frac{q_2 - P_2}{q_1 - P_1} = \frac{5 - 0}{2 - 4} = -\frac{5}{2}$$

 $\frac{\text{rise}}{\text{run}}$

(10)

So $w_1 = 5$ & $w_2 = 2$ on the HP

$$h(x) = 5x_1 + 2x_2 + b$$

Can derive b by using any point on the HP

$$x = 4, 0 \quad b = -5x_1 - 2x_2$$

$$b = -5 \cdot 4 - 2 \cdot 0$$

$$b = -20 - 0$$

$$b = -20$$

thus weight vector = $\begin{pmatrix} 5 \\ 2 \end{pmatrix}$ bias = -20

$$h(x) = w^T x + b = (5 \ 2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - 20 = 0$$

Distance of origin from HP =

$$\delta = y \cdot r = \underline{y \cdot h(x)}$$

$$= -1 \cdot r = \frac{y-b}{\|w\|} = \frac{-(-20)}{\sqrt{29}} = 3.71$$

(ii)

Margin & Support Vectors of a HP

given train set D of n points x_i with $y_i \in \{+1, -1\}$
 and hyperplane $h(x) = 0$

for each data point x_i , we can find its
 distance to the hyperplane using:

$$\delta_i = \frac{y_i \cdot h(x_i)}{\|w\|} = \frac{y_i (w^T x_i) + b}{\|w\|}$$

Over all of the points, the margin
 is the min distance point from HP

$$\delta^* = \min_{x_i} \left\{ \frac{y_i (w^T x_i) + b}{\|w\|} \right\}$$

$\delta \neq 0$ as is assumed to be a separating HP

- All points that achieve min DST = ^{Support} vectors
- Support vector = on margin

Note:

- $y_i (w^T x_i) + b$ = absolute Distance
- $\|w\|$ = makes it relative to unit vector

(12)

Canonical Hyperplane

Consider the hyperplane equation:

Multiplying on both sides by some scalar s yields an equivalent HP:

$$s \cdot h(x) = s \cdot w^T x + s \cdot b = (sw)^T x + (sb) = 0$$

$sw = 0 \quad \uparrow$

- Simply scaling will not change the HP
- Need to set the target $sw = 1$
- This creates a canonical hyperplane
- Not merely a new HP. It is a HP where the margin from separation = 1

$$(s \cdot y^*) \cdot (w^T x^* + b) = 1$$

- y^* = margin data point
 - need to find the scale that allows = 1
- implies:

$$s = \frac{1}{y^* (w^T x^* + b)} = \frac{1}{y^* h(x^*)}$$

(13)

for the canonical hyperplane,
for each support vector x_i^* (label y_i^*)
we now have $y_i^* h(x_i^*) = 1$

And hence, any point that is not
a support vector is ≤ 1

All points:

- $y_i (w^T x_i + b) \geq 1$ for all $x_i \in D$

14

Canonical HP: Example

Separating Hyperplane: $h(x) = \begin{pmatrix} 5 \\ 2 \end{pmatrix}^T x - 20 = 0$

Support Vector = $x^* = (2, 2)^T$ $y^* = -1$

To find the canonical hyperplane eq we need rescale the weight vector and bias by the scalar s

$$s = \frac{1}{y^* h(x^*)} = -1 \cdot \frac{\begin{pmatrix} 5 \\ 2 \end{pmatrix}^T \begin{pmatrix} 2 \\ 2 \end{pmatrix}}{-20} = \frac{1}{6}$$

thus rescaled w vector = $w = \frac{1}{6} \cdot \begin{pmatrix} 5 \\ 2 \end{pmatrix} = \begin{pmatrix} 5/6 \\ 2/6 \end{pmatrix}$ bias = $b = \frac{-20}{6}$

Canonical form of hyperplane:

$$h(x) = \begin{pmatrix} 5/6 \\ 2/6 \end{pmatrix}^T x - (-20/6) = \begin{pmatrix} 0.833 \\ 0.333 \end{pmatrix}^T x - 3.33$$

Margin:

$$\delta^* = \frac{y^* h(x)}{\|w\|} = \frac{-1}{\sqrt{\frac{5^2}{6} + \frac{2^2}{6}}} = \frac{1.14}{\sqrt{1.67}} = 1.14$$

relative margin
in terms of w

Absolute margin