

(1)

5.2 SVM: Linear & Separable

given D , n points $x_i \in \mathbb{R}^d$ w/ labels
 $y_i \in \{-1, +1\}$

Assume that a hyperplane exists that
 can perfectly separate into $\{-1, +1\}$
 $h(x) < 0 \quad h(x) \geq 0$

Should this be the case then there is an
 infinite number of HPs in the space

Marginal

Max Margin Hyperplane

fundamental to SVM is to choose the
 canonical hyperplane (w, b) that yields
 the max margin among all possible
 separating hyperplanes

If S_n^* represents the margin for the HP
 $h(x)=0$ the the goal is to find the
 optimal h^*

$$h^* = \underset{h}{\operatorname{argmax}} \{ \delta_h^* \} = \underset{w,b}{\operatorname{argmax}} \left\{ \frac{1}{\|w\|} \right\}$$

\curvearrowleft Altered using scaling factor

SVM task is to find hyperplane that maximizes the margin $\frac{1}{\|w\|}$

Subject to constraints:

$$\text{all } y_i (w^T x_i + b) \geq 1 \text{ for all points } x_i$$

Notice instead of maxing margin $\frac{1}{\|w\|}$

we can minimize just $\|w\|$

new object function: $\min_{w,b} \left\{ \frac{\|w\|^2}{2} \right\}$
 w/ same constraint

we can directly solve the minimization problem using standard optim algs

But it is more common to solve the Dual problem - obtained via the use of Lagrange multipliers

introduce a Lagrange multiplier a_i for each constraint (data point) which satisfies the Karush-Kuhn-Tucker (KKT) conditions

KKT = framework for solving constraint opt prob

- transform constrained prob into set of equations & inequalities that can be solved by opt techniques
- for SVMs KKT conditions are essential for deriving the dual formulation of the opt prob

$$a_i(y_i(w^T x_i + b) - 1) = 0$$

and $a_i \geq 0$

Incorp all n constraints into obj func

$$\min L = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n a_i (y_i (w^T x_i + b) - 1)$$

L should be min w/ respect to w & b

& maximized with respect to a_i

Take the derivative w/ res to w & b
and set to 0

$$\frac{\partial}{\partial w} L = w - \sum_{i=1}^n a_i y_i x_i = 0$$

$$w = \sum_{i=1}^n a_i y_i x_i$$

$$\frac{\partial}{\partial b} L = \sum_{i=1}^n a_i y_i = 0$$

How does the lagrange
change affect
change b

▲
This shows the optimal weight vector can be
shown as a linear combination of the x_i
data points w/ Lagrange multipliers as
coefficients, a_i, y_i

2nd eq implies the sum of lang
multi must be zero

Plug these constraints into the
L function to obtain the dual
Lagrangian Function

~~Dual~~

$$L_{\text{dual}} = \frac{1}{2} w^T w - \sum_{i=1}^n a_i (y_i (w^T x_i + b) - 1)$$

~~margin~~ ~~constraints~~

Allt under the hood Expands simple terms

$$L_{\text{dual}} = \frac{1}{2} w^T w - w^T (\underbrace{\sum a_i y_i x_i}_w) - b \sum a_i y_i + \sum a_i$$

$\frac{1}{2} w^T w - w^T w$

$$= -\frac{1}{2} w^T w + \sum a_i$$

- Expand, Transpose
- Distribute terms

$$= \sum a_i - \frac{1}{2} \sum \sum a_i a_j y_i y_j x_i^T x_j$$

This gives right to a new set of obj funcs

Started w/ Primal: $\min_{w,b} \frac{1}{2} \|w\|^2$

The new dual objective:

$$\text{func: Max}_{\alpha} L_{\text{dual}} = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\text{Constraints: } \alpha_i \geq 0, \forall i \in D, \sum \alpha_i y_i = 0$$

$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$ vector of lag multipliers

L_{dual} is convex quadratic programming prob (over α)
which can be solved using standard opt techs; i.e.
gradient based method for solve dual formulation

Dual = $\min \rightarrow \max$

Cheaper computationally esp for large data + ^{Standard} Techs
uses kernel tricks (???) good for non-linear data

Dual = Expressed in terms of Lagrange

Dual = Same solution to Dual as primal

we do not solve w or b

only the new α multipliers

Weight Vectors & bias

Once we have obtained the α_i values for $i=1, \dots, n$ we can solve for the weight and bias

According to KKT we have:

$$\alpha_i (y_i (w^T x_i + b)) = \underline{0}$$

this means two things must occur:

$$(1) \alpha_i = 0$$

$$(2) \underbrace{y_i (w^T x_i + b)}_{\text{must be } 1} - 1 = 0$$

Important because:

- if $\alpha_i > 0$ then $y_i (w^T x_i + b) = 1$
- if $= 1$ then x_i point is a support vector
- otherwise $> 1 \Rightarrow \alpha_i$ has to be 0
- if not support then $\alpha_i = 0$

Once we have solved for all α_i points we can compute the w vector

$$\text{weight} = \omega = \sum_{a_i > 0} a_i y_i w_i$$

note $a_i = 0$ works as a switch to turn off the non-support vectors

In other words ω is obtained as a linear combination of its support vectors

The other points do not determine ω

Computing b

$$\text{recall, } a_i(y_i(\omega^T x_i + b) - 1) = 0$$

- b is computed per support vectors, so $a_i \neq 0$
- Meaning inside has to $= 0$ for ≤ 0

$$\text{so, } y_i(\omega^T x_i + b_i) = 1$$

rearrange:

$$b_i = \frac{1}{y_i} - \omega^T x_i = y_i - \omega^T x_i$$

then compute b as the average over all b_i ^{sup}

$$\cdot b = \arg \left\{ \min_{a_i > 0} \{ b_i \} \right\}$$

SVM Classifier

Once you have found the optimal hyperplane function $h(x) = w^T x + b$

for any new data point & class can be predicted as

$$\hat{y} = \text{Sign}(h(z)) = \text{Sign}(w^T z + b)$$

$\text{Sign}(\cdot)$ function returns +1 if argument is positive, -1 if neg

i.e. depending on the side of the HP it falls on

Much the same as Separable
but a new Slack & error
for to compute those with &
over the margin

Hinge vs Quad loss to gone
Change obj function