

Page 92 - WiFi A5 Notebook 2

Support Vector Machines

- * Works really well on some moderately sized datasets
- * Does not do ~~perfectly~~ all that well on large dataset (standard dataset)

Consider figure 8.1.

Is one of these classifiers better?

How do we quantify what we see?

What if we extend out from a decision bound?

The distance it takes to get to a data point is referred to a margin.

This region is a hyper-cylinder in n dimensions.

If we maximize this we have a maximum margin (linear) classifiers.

The data points that are closest to this line have a name as well:
Support Vectors!

Page 93 - WiFi A5 Notebook 2

We can now make two arguments:

1. margin should be as large as possible
2. Picking support vectors are the most useful datapoints because they are the one we might get wrong

After training we can throw away all but support vectors!

We've used this before and we will again

$$y = \vec{w} \cdot \vec{x} + b$$

if this is > 0 then we say this is the positive class (+)
otherwise it is the negative class (-)

Now we need to define the margin.

If \vec{x} is a support vector then call it \vec{x}^+ or \vec{x}^-



The book is pretty wrong on page 172. Here is the correct derivation:

For a training set, the margin of the i th sample is:

$$Y_i = \vec{w} \cdot \vec{x}_i + b$$

where \vec{w} is the weight vector.

Our problem is then:

$$\max_{\vec{w}, b} \min_i Y_i \quad \text{with } \|\vec{w}\| = 1$$

pick
support
vectors
as closest
to boundary

Let m be the maximum margin

$$\forall_i |\vec{w} \cdot \vec{x}_i + b| \geq m$$

Let t_i be $+1$ if this is a positive class. Let it be -1 if it is negative.

Then

$$\forall_i \quad t_i (\vec{w} \cdot \vec{x}_i + b) \geq m$$

$$\frac{t_i}{m} (\vec{w} \cdot \vec{x}_i + b) \geq 1$$

$$\parallel \vec{V} \parallel$$

By formulation we have

$$\parallel \vec{V} \parallel^2 = \left\| \frac{1}{m} \vec{w} \right\|^2$$

$$= \frac{1}{m^2} \parallel \vec{w} \parallel^2 = 1$$

margin

\therefore maximizing ~~margin~~ m is the

same as minimizing ~~margin~~

$$\parallel \vec{V} \parallel^2$$

So our problem is now:

$$\min \frac{1}{2} \|v\|^2$$

subject to $t_i (\vec{v} \cdot \vec{x}_i + b) \geq 1$
for all i

The book just uses \vec{w} to represent \vec{v} , so to stay consistent we will go back to their formulation which is equivalent.

$$\min \frac{1}{2} \vec{w}^T \cdot \vec{w} \quad \text{subject to} \quad t_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 \quad \forall i$$

This problem is actually one that can be solved with quadratic programming. We will not cover that here. — But we will formulate the problem.

When we find the optimal solution the following will be satisfied

↓ designates optimal

$$\lambda_i^* (1 - z_i (\vec{w}^{*T} x_i + b^*)) = 0$$

$$1 - z_i (\vec{w}^{*T} x_i + b) \leq 0$$

$$\lambda_i^* \geq 0$$

↑
Lagrangian multipliers which are positive values.

These are known as the Karush-Kuhn-Tucker (KKT) conditions.

The first condition says that for $\lambda_i \neq 0$

$$1 - z_i (\vec{w}^{*T} x_i + b^*) = 0$$

This is only true for support vectors so we only have to consider them.

In jargon these form the active set. For the SV the inequalities (\leq) are equalities so we can formulate the Lagrangian function \Rightarrow

$$L(\vec{w}, b, \lambda) = \frac{1}{2} \vec{w}^T \vec{w} + \sum_{i=1}^n \lambda_i (1 - t_i (\vec{w}^T \vec{x}_i + b))$$

We differentiate with respect to \vec{w} & b :

$$\nabla_{\vec{w}} L = \vec{w} - \sum_{i=1}^n \lambda_i t_i \vec{x}_i$$

$n \leftarrow \# \text{ of SVs}$

and

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^n \lambda_i t_i$$

If we set derivatives equal to 0, we get:

$$\vec{w}^* = \sum_{i=1}^n \lambda_i t_i \vec{x}_i$$

$$\sum_{i=1}^n \lambda_i t_i = 0$$

We can ~~substitute~~ substitute back into the equation and rearrange and get:

$$L(\vec{w}^*, b^*, \lambda) = \sum_{i=1}^n \lambda_i - \sum_{i=1}^n \lambda_i t_i - \frac{1}{2} (A)$$

$$A = \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j t_i t_j \vec{x}_i^T \vec{x}_j$$

This equation is known as the dual problem & we

aim to maximize it with

respect to λ_i .

constraints

$$\lambda_i \geq 0 \text{ for all } i$$

$$\text{and } \sum_{i=1}^n \lambda_i t_i = 0$$

So how do you find \vec{w}^* ?

It is right there in an equation.

How about b^* ?

consider all ~~the~~ Support Vectors:

$$b^* = \frac{1}{N_0} \sum_{\text{support vectors } j} \left(t_j - \sum_{i=1}^n \lambda_i t_i \bar{x}_i \bar{x}_j \right)$$

Once we know these we are all set!

But what about problems not linearly separable?

Add in slack variables.

see section 8.1.3.

Derivations are the same.
Coming up next: kernels!