

SVM:

Motivation:

Assume a **linearly separable training set**

Find a plane  $f(x; w, b) = 0 \Leftrightarrow w^T x + b = 0$  that discriminates perfectly the training points and is "as far as possible" from all the training samples. For the math that follows, assume the labels  $y_i \in \{-1, +1\}$ .

A feasible plane should put the positive training samples in one side of the plane and the negative samples in the other side (semi-plane):  
 $f(x_1; w, b) > 0, \forall y_1 = +1$   
 $f(x_2; w, b) < 0, \forall y_2 = -1$

In a more compact way:  
 $y_i \times f(x_i; w, b) > 0, \forall x_i$

Note that since the algebraic distance of a training point  $x_i$  to the planes is  $\frac{f(x_i; w, b)}{\|w\|}$ , for feasible planes  $\frac{y_i \times f(x_i; w, b)}{\|w\|}$  gives the absolute value of the distance.

The SVM model is the feasible plane that gives the "best worst case".

The SVM model is the feasible plane that maximizes the minimum of the distances of the training samples to the plane:

$$\operatorname{argmax}_w \left( \min_i \frac{y_i \times f(x_i; w, b)}{\|w\|} \right)$$

S.t.:  
 $y_i \times f(x_i; w, b) > 0, \forall x_i$

Rewriting:

$$\operatorname{argmax}_w \left( \min_i \frac{y_i (w^T x_i + b)}{\|w\|} \right)$$

Prob(A)

Since the plane  $w^T x + b = 0$  is the same plane as  $w^T x/k + b/k = 0$ , for  $k \neq 0$ , one can work with the subset of equations (of planes) for which  $(\min_i y_i (w^T x_i + b) = 1)$

(see **Details A**)

The optimization problem simplifies to

$$\operatorname{argmax}_w \left( \frac{1}{\|w\|} \right)$$

$$y_i (w^T x_i + b) > 1, \forall x_i$$

Finally, it can be rewritten as:

$$\operatorname{argmin}_w (w^T w)$$

$$y_i (w^T x_i + b) > 1, \forall x_i$$

Nonlinearly separable datasets

Now, there is no plane to perfectly discriminate the points. Accepting that it's still better to stay with the linear model (maybe it's just due to the noise that we are not able to linearly discriminate the points).

However, now, prob(A) has no solution, the constraints cannot be satisfied for all points. This means that for some points  $y_i (w^T x_i + b)$  will be negative, which we can compensate by adding a nonnegative quantity  $\epsilon_i$ .

So, the constraints can be reformulated as

$$y_i (w^T x_i + b) + \epsilon_i > 0, \forall x_i$$

$$\epsilon_i \geq 0, \forall x_i$$

Without adapting the goal function these constraints lead to a trivial solution, very far away from the training points and where the points of one of the classes are all misclassified (with this trivial solution, the distance of the training samples to the plane are all very high).

The workaround is to penalize misclassifications in the goal function. One need to accept that some points have to be misclassified but we prefer solutions with just a few misclassifications.

$$\operatorname{argmin}_w \left( c \sum_{i=1}^N \epsilon_i + \frac{\|w\|}{\min_i y_i (w^T x_i + b)} \right)$$

C is a nonnegative constant weighing the cost of the misclassifications with the importance of achieving a large margin (the margin is the minimum of the distances).

Using the same simplification as before, one finally gets

$$\operatorname{argmax}_w \left( c \sum_{i=1}^N \epsilon_i + \|w\| \right)$$

$$y_i (w^T x_i + b) + \epsilon_i > 1, \forall x_i$$

$$\epsilon_i \geq 0, \forall x_i$$

Although valid, it's much more common in practice to relate the sum of the training errors with the square of the inverse of the margin, leading to the formulation widely adopted in SVMs:

$$\operatorname{argmax}_w \left( c \sum_{i=1}^N \epsilon_i + w^T w \right)$$

$$y_i (w^T x_i + b) + \epsilon_i > 1, \forall x_i$$

$$\epsilon_i \geq 0, \forall x_i$$

Note that the two constraints can be combined in a single one:

$$\epsilon_i \geq \max(0, 1 - y_i (w^T x_i + b)), \forall x_i$$

Leading to

$$\operatorname{argmax}_w \left( c \sum_{i=1}^N \epsilon_i + w^T w \right)$$

s.t.

$$\epsilon_i \geq \max(0, 1 - y_i (w^T x_i + b)), \forall x_i$$

Finally, there is always an optimal for which

$$\epsilon_i = \max(0, 1 - y_i (w^T x_i + b)), \forall x_i$$

(see **detailsEq**)

And therefore the constraints can be integrated in the goal function leading to

$$\text{Model} = \operatorname{argmax}_w \left( c \sum_{i=1}^N \max(0, 1 - y_i (w^T x_i + b)) + w^T w \right) \quad (\text{EqSVMmodel})$$

Hinge loss

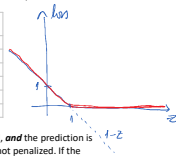
It is instructive to analyse the previous goal function as a sum of two terms. The first part, the data fitness,

$$\left( \sum_{i=1}^N \max(0, 1 - y_i (w^T x_i + b)) \right)$$

is penalizing errors in the training set, while the second is setting preference for 'simple' models.

Let  $z = y(w^T x + b)$ . If the prediction of the model

y	$w^T x + b$	$z = y(w^T x + b)$	$\max(0, 1 - y_i (w^T x_i + b)) = \max(0, 1 - z_i)$
+1	-5	-5	6
+1	5	5	0
+1	0.5	0.5	0.5
-1	5	-5	6
-1	-5	5	0
-1	-0.5	0.5	0.5



If the prediction agrees in sign with the sign of the true label, **and** the prediction is 'strong' (greater than 1 in absolute value), the prediction is not penalized. If the prediction agrees in sign with the sign of the true label **but** the prediction is weak (less than 1 in absolute value), in the sense that is very close to the boundary, the prediction is weakly penalized. If the prediction is with the wrong sign, the prediction is strongly penalized, proportional to absolute value of the prediction. This loss function - named the hinge loss - by penalizing predictions close to the boundary tries to push the values away from the boundary (and on the right side).

Kernel formulation: towards nonlinear models

So far, we have only been able to work with linear models (planes).

Towards the goal of achieving nonlinear models, we need to rewrite, again, our goal function.

Consider the goal function as given in (EqSVMmodel).

In optimal solution  $(w^*, b^*)$ , the  $w^*$  vector can always be written as a linear combination of the training observations.

As such, one can restrict the optimization process to the set of feasible solutions that can be written as a linear combination of the training samples.

$$\text{Model} = \operatorname{argmax}_{a_i, b} \left( c \sum_{i=1}^N \max(0, 1 - y_i (\sum_{j=1}^N a_j x_j^T x_i + b)) + \sum_{i=1}^N \sum_{j=1}^N a_i a_j x_i^T x_j \right)$$

The important part to note now is that input observations are only used as the inner product between pairs of observations:  $K(x_i, x_j) = x_i^T x_j$ .

This is true both when training (as seen in the previous equation) and for inference. The prediction for a test observation  $x_2$  is

$$w^T x_2 + b = \sum_{i=1}^N a_i x_i^T x_2 + b$$

(Note that b and  $a_i$  were learnt using only the inner products between pairs of training samples. When making prediction, one uses those values and the inner products between training samples and the test sample.)

Assume now that a linear model in the input space is not a good idea (see figure). One may then choose to apply the above described methodology, not directly in the input x-space but in the new (transformed) z-space.

The effort is not much. One just needs to replace the inner product between observations  $x_1$  and  $x_2$  by the corresponding inner product between  $z_1$  and  $z_2$ .

Assume as an example that  $z^{(1)} = x^{(1)}$ ,  $z^{(2)} = x^{(2)}$ ,  $z^{(3)} = (x^{(1)})^2 + (x^{(2)})^2$ . In here  $z^{(1)}$  is the 1-component of vector  $x$ .

In this case the inner product between vectors  $z_1$  and  $z_2$  is  $z_1^T z_2 = x_1^T x_2 + \|x_1\|^2 \|x_2\|^2$ . On can say that the inner product in the new, feature space, is given by a function computed over the input space:

$$z_1^T z_2 = K(x_1, x_2) = x_1^T x_2 + \|x_1\|^2 \|x_2\|^2.$$

This function is called the kernel function. The advantage is that often we can avoid the explicit computation of the feature z-space. By working with suitable kernel functions we are indeed exploring implicit feature transformations.

For sure, not every function between two vectors to a scalar  $K(x_1, x_2)$  will correspond to the inner product of some implicit feature space (see theory related with Mercer's theorem - [https://en.wikipedia.org/wiki/Mercer%27s\\_theorem](https://en.wikipedia.org/wiki/Mercer%27s_theorem)). But some well-known kernel functions are often used in practice with good results. RBF, polynomial, linear are among the most typical kernels.

It is a good idea to think of kernels as our measure of similarity between observations. This kernel can be application specific. For instance in computer vision, if an object is represented by the histogram of the grayscale values, one can use as kernel the correlation between histograms.

Support vectors

ZZZZ

ZZZZZZZ

Optimization

ZZZZZZZZZZZZ

Distance between a point and a plane

Let the function  $f(x; w) = w^T x + b, x \in \mathbb{R}^d$

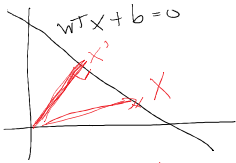
The set of points that make the equation  $f(x; w) = 0$  true is a plane.

The algebraic distance between a point  $x$  and the plane  $f(x; w) = 0$  is given by  $\frac{f(x; w)}{\|w\|} = \frac{w^T x + b}{\|w\|}$

Note that this is a signed distance: points in one of the semi-planes have positive distance; points in the other semi-plane have negative distance.

- If  $x$  is the origin, the distance is zero.

- If  $x$  is the origin, the proof A or the proof B for any  $x \in \mathbb{R}^d$ , see proof C.



$$\textcircled{A} \left\{ \begin{array}{l} \min \|x\| \\ \text{s.t. } w^T x + b = 0 \end{array} \right\} \Leftrightarrow \left\{ \begin{array}{l} \min x^T x \\ \text{s.t. } w^T x + b = 0 \end{array} \right\}$$

$$\textcircled{B} \frac{|<x, w>|}{\|w\|} = \frac{|w^T x|}{\|w\|} = \frac{|b|}{\|w\|} = \frac{|f(0; w)|}{\|w\|}$$

modify A accordingly

the gradient should be proportional

$$\left\{ \begin{array}{l} 2x = 2w \\ w^T x + b = 0 \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} 2w^T x = 2w^T w \\ w^T x + b = 0 \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} w^T x = b \\ 2 = \frac{2b}{w^T w} \end{array} \right\} \Rightarrow x = \frac{-b}{\|w\|^2} w$$

distance = projection of any X to plane over the w direction (w is orthogonal to the plane)

$$x_A \text{ and } x_B \in \text{plane} \\ w^T (x_A - x_B) = 0$$

$$x^T x =$$

$$= \left( x + \frac{b}{\|w\|} \hat{w} - \frac{b \hat{w}}{\|w\|} \right)^T \left( x + \frac{b}{\|w\|} \hat{w} - \frac{b \hat{w}}{\|w\|} \right)$$

$$= \left( x + \frac{b \hat{w}}{\|w\|} \right)^T \left( x + \frac{b \hat{w}}{\|w\|} \right) - \left( x + \frac{b \hat{w}}{\|w\|} \right)^T \cdot \frac{b \hat{w}}{\|w\|} - \frac{b \hat{w}^T}{\|w\|} \cdot \left( x + \frac{b \hat{w}}{\|w\|} \right) + \frac{b^2 \hat{w}^T \hat{w}}{\|w\|^2}$$

$$= \left( x + \frac{b \hat{w}}{\|w\|} \right)^T \left( x + \frac{b \hat{w}}{\|w\|} \right) + \frac{b^2 \hat{w}^T \hat{w}}{\|w\|^2}$$

$$\text{minimized for } x = - \frac{b \hat{w}}{\|w\|}$$