

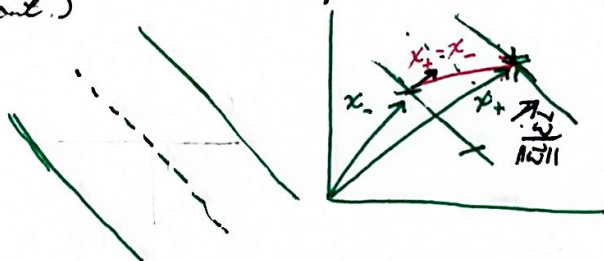
SUPPORT VECTOR MACHINES

SVM is an algorithm which fits the max-margin hyperplane in a certain feature space.

1.1 HARD MARGIN LINEAR SVM. GIVEN A LINEARLY SEPARABLE DATA SET WHERE EACH EXAMPLE IS EITHER FROM CLASS $y_i = +1$ OR $y_i = -1$, THEN A LINEAR SVM CAN ALWAYS FIND THE MAX-MARGIN CLASSIFIER THAT CORRECTLY CLASSIFIES ALL THE TRAINING DATA AS FOLLOWS.

$$\text{margin} = d^+ - d^- = \frac{w^T x_j}{\|w\|_2} - \frac{w^T x_k}{\|w\|_2}$$

SUPPOSE THAT WE KNOW THAT THERE ARE ONLY 2 TRAINING EXAMPLES ON THE MARGINS, I.E., THE SUPPORT VECTORS, $x_j, y_j = +1$ AND $x_k, y_k = -1$, AND THE PARAMETERS OF THE LINEAR SVM ARE w^* AND b^* . WRITE DOWN THE CONSTRAINTS THIS LINEAR SVM HAS TO SATISFY WITH x_j, x_k, w^*, b^* , AND FORM THE PROBLEM OF LINEAR SVM AS A CONSTRAINED OPTIMIZATION PROBLEM. (HINT: The optimization problem should use both d from eqn. 1, as well as the constraints we just worked out.)



DOT PRODUCT WITH UNIT VECTOR GIVES DISTANCE IN THAT DIRECTION LENGTH AND w IS NORMAL TO BOUNDARY SO $(x_k - x_j) \cdot \frac{w}{\|w\|} = \frac{x_k \cdot w}{\|w\|} - \frac{x_j \cdot w}{\|w\|}$ GIVES THE WIDTH OF THE MARGIN.

$w \cdot \vec{x} \geq c$
 $w \cdot \vec{x} + b \geq 0$ THEN +

① ADD CONSTRAINTS

$$w \cdot x_j + b \geq 1, w \cdot x_k + b \leq -1. y_i = \begin{cases} +1 & \text{for } + \text{ samples} \\ -1 & \text{for } - \text{ samples} \end{cases}$$

$$\Rightarrow y_i (w \cdot x_i + b) \geq 1 \Rightarrow y_i (w \cdot x_i + b) - 1 \geq 0$$

SO FOR x_j IN THE OUTER SET CONSTRAINT

$$y_j (w \cdot x_j + b) - 1 = 0$$

$$\Rightarrow w \cdot x_j + b - 1 = 0$$

$$\text{AND } w \cdot x_k + b + 1 = 0$$

(i)

② FORMULATE EQN TO BE SOLVED.

WE KNOW WE WANT TO MAXIMIZE THE WIDTH OF THE MARGIN $\frac{w^T x_j}{\|w\|} - \frac{w^T x_k}{\|w\|}$, and using (i) we CAN SOLVE FOR $w \cdot x_j$ AND $w \cdot x_k$, THEN PLUG THEM BACK INTO THE MARGIN EQN. AND SIMPLIFY TO GET

$$\text{WIDTH} = \frac{2}{\|w\|}, \text{ which we want to maximize.}$$

$$\Rightarrow \text{MINIMIZE } \|w\| \equiv \text{minimize } \frac{1}{2} \|w\|^2$$

③ FORM LINEAR SVM PROBLEM AS A CONSTRAINT OPTIMIZATION PROBLEM.

WE CAN COMBINE THE EQN TO BE MINIMIZED WITH OUR CONSTRAINTS USING LAGRANGIAN MULTIPLIERS

$$L = \frac{1}{2} \|\bar{w}\|^2 - \alpha_j [\bar{w} \cdot x_j + b - 1] - \alpha_k [\bar{w} \cdot x_k + b + 1]$$

1.2 However, in real world cases it is difficult to confirm if a dataset is linearly separable or not, even in kernel space. Therefore, soft margin SVM is introduced to handle those samples which can't correctly be classified.

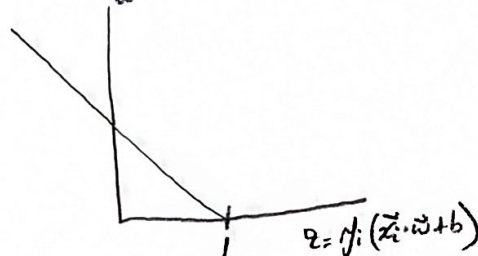
(a). PLEASE DESCRIBE HOW INCORRECTLY CLASSIFIED SAMPLES ARE HANDLED BY A SOFT-MARGIN SVM.

Instead of just minimizing the width of the margin, we now want to minimize the width of the margin + the number of mistakes. But not all mistakes are equally bad, so we'll use their position to the margin to weight the penalty.

$$\min \frac{1}{2} \|\bar{w}\|^2 + \sum_{i=1}^n \xi_i$$

↑
PARAM

The loss ξ_i for each x_i is calculated with



(b) Three types of examples will exist after training SVM based on the value of the slack variables. Describe them and state whether removing them would result in a change of decision boundary.

- ξ_i corresponding to $\alpha_i = 0$ IN A SUPPORT VECTOR ON THE MARGIN BOUNDARY
- ξ_i " " α_i s.t. $0 < \alpha_i < 1$, A SUPPORT VECTOR ON THE WRONG SIDE OF THE MARGIN BOUNDARY
- ξ_i " " $\alpha_i = 0$, x_i IS ON THE CORRECT SIDE OF THE DECISION BOUNDARY, AND IS NOT ON THE BOUNDARY.

THESE ARE SUPPORT VECTORS, AND BY DEF. WILL CHANGE THE DECISION BOUNDARY IF REMOVED

2 ADABOOST

2.1 ADABOOST IS AN ENSEMBLE OF WEAK CLASSIFIERS WHOSE ACCURACY IS SLIGHTLY OVER 50%. WHAT WILL HAPPEN TO THE ADABOOST ALGORITHM IF YOUR WEAK CLASSIFIER HAS EXACTLY 50% ACCURACY?

It will stop learning (or not learn anything, if this is the initial classifier h_0). This is because the voting pattern α is calculated by $\frac{1}{2} \ln(\frac{1-\epsilon}{\epsilon})$, which equals 0 when $\epsilon = 0.5$. No future classifiers will receive any voting power.

2.2 WHAT WILL HAPPEN IF YOU USE A WEAK BINARY CLASSIFIER WHOSE CLASSIFICATION ACCURACY IS LESS THAN 50%, SAY 45%?

It will learn to classify the opposite class, and you simply have to flip the sign.

0.5	1	2	3	4	5	6	7	8	9
1	1	1	-1	-1	-1	1	1	1	-1

TABLE 1: TRAINING SET

FOR THE WEAK CLASSIFIER C , WE WILL USE A SINGLE THRESHOLD θ , SO THAT

$$C(x) = \begin{cases} +1 & x < \theta \\ -1 & x \geq \theta \end{cases}$$

2.3 FOR THE FIRST ITERATION, WE LET THE THRESHOLD $\theta = 2.5$ WHICH MINIMIZES THE CLASSIFICATION ERROR AT THE CURRENT ITERATION.

SHOW US HOW YOU DERIVE THE WEIGHTS D_2 AFTER THIS ITERATION.

w_0	$\frac{1}{10}$	$\frac{1}{14}$
w_1	$\frac{1}{10}$	$\frac{1}{14}$
w_2	$\frac{1}{10}$	$\frac{1}{14}$
w_3	$\frac{1}{10}$	$\frac{1}{14}$
w_4	$\frac{1}{10}$	$\frac{1}{14}$
w_5	$\frac{1}{10}$	$\frac{1}{14}$
w_6	$\frac{1}{10}$	$\frac{1}{6}$
w_7	$\frac{1}{10}$	$\frac{1}{6}$
w_8	$\frac{1}{10}$	$\frac{1}{6}$
w_9	$\frac{1}{10}$	$\frac{1}{14}$

C misclassifies x_6, x_7, x_8

$$\epsilon = \frac{3}{10}$$

$$\text{So } w_6^{t+1} + w_7^{t+1} + w_8^{t+1} = \frac{1}{2}$$

w_6^t, w_7^t, w_8^t have ratio 1:1:1, so update weights to $\frac{1}{6}$ for each

Since for the 7 remaining weights (classified correctly, w_i^t equal ratio, w_i^{t+1} should sum to $\frac{1}{2}$), so update $w_i^{t+1} \forall i \in \text{CORRECT}$ to $\frac{1}{14}$

2.4 SHOW US HOW YOU DERIVE D_3 & D_4 , RESPECTIVELY.

SEE NEXT PAGE...

2.4 continued: show us how you derive D_3 & D_4 , EXHAUSTIVELY

	D_1	D_2	D_3	D_4
ω_0	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{8}$
ω_1	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{8}$
ω_2	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{8}$
ω_3	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{108}$
ω_4	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{108}$
ω_5	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{108}$
ω_6	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{108}$
ω_7	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{108}$
ω_8	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{108}$
ω_9	$\frac{1}{10}$	$\frac{1}{11}$	$\frac{1}{66}$	$\frac{1}{8}$

CONDITION FOR MINUS	MISCLASSIFIES	E_{D3}	E_{D4}	E_{D5}
$x \geq -0.5$	0, 1, 2, 6, 7, 8	$\frac{10}{11}$	$\frac{2}{6}$	
$x \leq 0.5$	6, 7, 8	$\frac{7}{11}$	$\frac{2}{6}$	
$x \leq 5.5$	3, 4, 5, 6, 7, 8	$\frac{10}{11}$	$\frac{5}{6}$	
$x \geq 8.5$	3, 4, 5	$\frac{3}{11}$	$\frac{1}{2}$	
$x \leq -0.5$	3, 4, 5, 9	$\frac{4}{11}$	$\frac{3}{6}$	
$x \geq 2.5$	0, 1, 2, 3, 4, 5, 9	$\frac{7}{11}$	$\frac{4}{6}$	
$x \leq 5.5$	0, 1, 2, 9	$\frac{4}{11}$	$\frac{2}{6}$	
$x \leq 8.5$	0, 1, 2, 6, 7, 8, 9	$\frac{10}{11}$	$\frac{1}{2}$	

For D_3 only:

$$\omega_0 = \omega_1 = \omega_2 = \frac{1}{6}$$

$$\omega_3 = \omega_4 = \omega_5 = \frac{1}{6}$$

For D_4 only:

$$\omega_6 = \frac{1}{2} \cdot \frac{1}{1-2^6} = \frac{1}{14 \cdot 2} \cdot \frac{1}{1-2^6} = \frac{1}{28} \cdot \frac{1}{1-2^6} = \frac{1}{11 \cdot 14 \cdot 2} = \frac{1}{308}$$

$$= \frac{1}{6 \cdot 2} \cdot \frac{1}{11} = \frac{1}{66} = \omega_6, \omega_7, \omega_8$$

D_4 :

$$h^3: x \leq 5.5$$

For D_4 only:

$$\omega_0 = \omega_1 = \omega_2 = \omega_3$$

$$\omega_6 + \omega_7 + \omega_8 + \omega_9 = \frac{1}{2}$$

$$\omega_0 + \omega_1 + \omega_2 + \omega_3 = \frac{1}{2}$$

For D_4 only:

$$\omega_3 = \frac{\omega_0}{2} \cdot \frac{1}{1-2^6}$$

$$\text{when } \omega_0 = \frac{1}{66}:$$

$$\frac{1}{66 \cdot 2} \cdot \frac{1}{54} = \frac{1}{108}$$

$$\text{when } \omega_0 = \frac{1}{66}$$

$$\frac{1}{66 \cdot 2} \cdot \frac{1}{54} = \frac{1}{108}$$

2.5 SHOW US WHAT HAPPENS ON THE FOURTH ITERATION

Based on the phrasing of the hint, all ϵ should be $\frac{1}{2}$, but that's not what I'm getting.

3. KNN CLASSIFIER

3.1 A LAZY CLASSIFIER. Consider an online learning setting where besides the observed training data, we have new training data coming in as time goes by, some classifiers have to be retrained from scratch.

NOTE: THE "HINT" FOR THIS PROBLEM ISN'T A HINT AT ALL. IT'S A RE-DEFINING OF WHAT IT MEANS TO "LEARN FROM SCRATCH", WHICH IS SILLY.

(c) BETWEEN SVM, KNN, + NB, WHICH HAVE TO BE RETRAINED "FROM SCRATCH" WHEN NEW DATA IS RECEIVED?

SVM has to be retrained when the new data falls on or on the wrong side of the margin. (IF THE SVM WAS MODELLING DATA GUARANTEED TO BE LINEARLY SEPARABLE, IT WOULD BE OK IF THE NEW DATA WAS ON THE MARGIN [i.e. not using slack vars]).

USING THIS DEFINITION, NB HAS TO BE RETRAINED "FROM SCRATCH".

KNN DO NOT HAVE A TRAINING PHASE, SO OF COURSE THEY DON'T NEED TO BE RETRAINED. (OF COURSE, YOU COULD ALSO CONSIDER ADDING NEW POINTS TO THE TRAIN SET AS "TRAINING", IN WHICH CASE, IT WOULD INCLUDE KNN $O(1)$)

(e)

FOR SVM, TIME COMPLEXITY FOR TESTS, WILL DEPEND ON IF IT IS A LINEAR SVM OR A KERNEL SVM.

$$O(\# \text{ OF SUPPORT VECTORS})$$

FOR NB,

$$O(1)$$

FOR KNN, testing is linear in terms of the # of training data.

$$O(n \cdot [\# \text{ of dimension}])$$

(NOTE THAT THIS IS USING MEDIAN OF MEDIAN TO FIND KTH SMALLEST DISTANCE).

3.2

(a)

ALGORITHM KNN

procedure predict

for test_sample in test_set:

diff = element-wise distance between test_sample + each train_sample in train_set

distance = L2-norm of diff matrix.

min_idx = indices of first k smallest values in distance vector

neighbors = train_targets[min_idx]

classification = most frequently occurring value in neighbors vector.

return all classifications

and procedure

(b) SEE ATTACHED IMAGE + CODE!