# CSE 555
# Assignment 3
# Name: Rajiv Ranjan
# Ub id:50249099

1.  Write code to train a multi-class support vector classifier with dot-product kernel and 1-norm soft margin using the MNIST training data set. Then reporting the performance using MNIST test data set. There is a hyper-parameter that sets the trade-off between the margin and the training error --- tune this hyper-parameter through cross-validation.

Soln:  The code is being sent as an attachment in the file svm1.py and svm2.py

Set up for the code: Cross validation Score has been taken as 5. For validating the training data set through Cross Validation often an odd number is taken.

I have written 2 set of codes: First one is SVM1.py which runs on the entire MNIST data set, as per the requirement a dot-product kernel has been used that means the kernel was selected to be linear. The cache size was defined to be 1000. Also it was required that the hyper parameter C and gama be varied and it was varied. Below is the output for various values of C and gama.

O/p
scaling
grid search
Fitting 5 folds for each of 20 candidates, totalling 100 fits
[CV] C=100000.0, gamma=0.03125 ......................................
[CV] C=100000.0, gamma=0.03125 ......................................
[CV] C=100000.0, gamma=0.03125 ......................................
[CV] C=100000.0, gamma=0.03125 ......................................
[CV] .............. C=100000.0, gamma=0.03125, score=0.781745 -25.8min
[CV] C=100000.0, gamma=0.03125 ......................................
[CV] .............. C=100000.0, gamma=0.03125, score=0.785336 -26.0min
[CV] C=100000.0, gamma=0.0625 ......................................
[CV] .............. C=100000.0, gamma=0.03125, score=0.782500 -26.0min
[CV] C=100000.0, gamma=0.0625 ......................................
[CV] .............. C=100000.0, gamma=0.03125, score=0.785090 -26.1min
[CV] C=100000.0, gamma=0.0625 ......................................
[CV] .............. C=100000.0, gamma=0.03125, score=0.784328 -24.9min
[CV] C=100000.0, gamma=0.0625 ......................................
[CV] .............. C=100000.0, gamma=0.0625, score=0.778250 -62.1min
[CV] C=100000.0, gamma=0.0625 ......................................
[CV] .............. C=100000.0, gamma=0.0625, score=0.779753 -62.2min
[CV] C=100000.0, gamma=0.125 ......................................
[CV] .............. C=100000.0, gamma=0.0625, score=0.780841 -62.8min
[CV] C=100000.0, gamma=0.125 ......................................
[CV] .............. C=100000.0, gamma=0.0625, score=0.776494 -63.4min

[CV] C=100000.0, gamma=0.125 .........................................
[CV] ............... C=100000.0, gamma=0.0625, score=0.680077 -64.5min
[CV] C=100000.0, gamma=0.125 .........................................
[CV] ............... C=100000.0, gamma=0.125, score=0.692795 -147.0min
[CV] C=100000.0, gamma=0.125 .........................................
[CV] ............... C=100000.0, gamma=0.125, score=0.679853 -146.6min
[CV] C=100000.0, gamma=0.25 .........................................
[CV] ............... C=100000.0, gamma=0.125, score=0.688000 -146.7min
[CV] C=100000.0, gamma=0.25 .........................................
[CV] ............... C=100000.0, gamma=0.125, score=0.683721 -142.5min
[CV] C=100000.0, gamma=0.25 .........................................
[CV] ............... C=100000.0, gamma=0.125, score=0.700717 -151.0min
[CV] C=100000.0, gamma=0.25 .........................................
[CV] ................ C=100000.0, gamma=0.25, score=0.441233 -168.4min
[CV] C=100000.0, gamma=0.25 .........................................
[CV] ................ C=100000.0, gamma=0.25, score=0.376487 -160.6min
[CV] C=1000000.0, gamma=0.03125 ....................................
[CV] ............. C=1000000.0, gamma=0.03125, score=0.785090 -16.0min
[CV] C=1000000.0, gamma=0.03125 ....................................
[CV] ................ C=100000.0, gamma=0.25, score=0.419333 -151.4min
[CV] C=1000000.0, gamma=0.03125 ....................................
[CV] ............. C=1000000.0, gamma=0.03125, score=0.785336 -17.1min
[CV] C=1000000.0, gamma=0.03125 ....................................
[CV] ............. C=1000000.0, gamma=0.03125, score=0.782500 -17.1min
[CV] C=1000000.0, gamma=0.03125 ....................................
[CV] ............. C=1000000.0, gamma=0.03125, score=0.781745 -16.0min
[CV] C=1000000.0, gamma=0.0625 ....................................
[CV] ............. C=1000000.0, gamma=0.03125, score=0.784328 -16.4min
[CV] C=1000000.0, gamma=0.0625 ....................................
[CV] ................ C=100000.0, gamma=0.25, score=0.398150 -122.4min
[CV] C=1000000.0, gamma=0.0625 ....................................
[Parallel(n_jobs=-1)]: Done  24 tasks      | elapsed: 508.6min

Inference: This shows the various C and gama values used: they were calculated using the below line of code which is self explanatory:
parameters = {'C':10. ** np.arange(1,5), 'gamma':2. ** np.arange(-5, -1)}

Also we see that the code ran for a lot of time approximately 10 hours. Also the best prediction score was found to be around 0.78 i.e around 78 percent accurate, which is good for a linear SVM classifier.

The various gamma values were: 0.25, 0.03125, 0.125, 0.0625.

Various C values was: 10000,100000 etc

The accuracy score for a total run of 10 hrs has been shared above.

2nd Part: <mark>I also wrote a second code to test the accuracy. The code for this is in SVM2.py</mark> This code was run a small data set to check how well the classifier can distinguish between 2 different digits, let's see if it can differentiate between 8 and 9 how well. Cross validation score was taken to be: 5

For this purpose, load_digits dataset from the sklearn list of datasets was used.

It is a subset of the MNIST data set.

Also for this optunity package was used.

1)Now firstly it was run with default parameters and the accuracy was found to be 0.7655589359455676.

2)The program was run for all the 3 models linear kernel, polynomial kernel and the rbf kernel. The various adjusted best hyperparametrs found were :

Optimal parameters {'kernel': 'rbf', 'C': 5.807812500000001, 'coef0': None, 'degree': None, 'logGamma': -3.4678653190880104}

Best Score of tuned SVM: 0.985

But the assignment specifically asks to run for linear i.e. dot product:

So the output for that was also run and could be seen below:

| C | coef0 | degree | kernel | logGamma | value |
|---|---|---|---|---|---|
| 1.742604 | NaN | NaN | linear | NaN | 0.962069 |
| 0.117604 | NaN | NaN | linear | NaN | 0.962069 |
| 0.210938 | NaN | NaN | linear | NaN | 0.962069 |
| 1.835938 | NaN | NaN | linear | NaN | 0.962069 |
| 0.024271 | NaN | NaN | linear | NaN | 0.962069 |
| 1.649271 | NaN | NaN | linear | NaN | 0.962069 |

The different C hyperparameter which tells how much to avoid misclassifying and the different score values of accuracy are also shown.

2. (8 points) Identify the Lagrange dual problem of the following primal problem:

Given features $(X_1, Y_1), \cdots, (X_N, Y_N)$, where $Y_1, \cdots, Y_N \in \{-1, 1\}$,

Minimize $w^T \cdot w + C \sum\limits_{i=1}^{N} \xi_i$, the weighted sum between the squared length of the separating vector and the errors, where

w is the separating vector, $w^T \cdot w$ is the dot product, and $\xi_i$ is the error made by separating vector w on feature $(X_i, Y_i)$.

Subject to $Y_i \cdot (w^T \cdot X_i) \geq 1 - \xi_i$ and $\xi_i \geq 0$ for $i = 1, \cdots, N$. In other words, if the "normalized feature" $Y_i X_i$ has a margin less than 1,

$w^T \cdot (Y_i X_i) \leq 1$, we add a slackness term to make it 1.

Point out what is the "margin" in both the primal formulation and the dual formulation, what are the benefits of maximizing the margin. Characterize the support vectors. Point out the benefit of solving the dual problem instead of the primal problem.


Soln:

(1)

Given is

$(x_1, y_1), \dots (x_N, y_N)$ where $y_1, \dots y_N \in \{-1, 1\}$

Soln:- is basically N points where $x_1 \dots x_n$ is a feature vector, $y_N$ is a binary class can have value either as $-1$ or $1$

Part 1. Minimize $w^T \cdot w + C \sum_{i=1}^{N} \xi_i$

where $w$ is the separating vector, $w^T w$ is the dot product & $\xi_i$ is the error made by separating vector $w$ on feature $(x_i, y_i)$

subject to. $y_i \cdot (w^T \cdot x_i) \geq 1 - \xi_i$

$\xi_i \geq 0$ for $i = 1 \dots N$

Sol$^n$. So, the thing is the actual SVM is a mixture of 2



(empirical error)

tendency to overfit = false minimum

flexibility increases as we go right

here $m$ is the margin assumed is width of the margin

-- so, for actual SVM we have already aggregated over all the empirical error, and we now add something that will quantify our tendency to overfit,

so coming back to this formula here, this was given the margin $m$.

$\underset{w, b, \xi}{\arg\min} \sum_{i=1}^{n} \xi_i$ such that $y_i (w^T x_i + b) \geq 1 - \xi_i$

$\xi_i \geq 0$

$\|w\| \leq \frac{1}{m}$

so we will consider $\sigma(w) = w^T w \propto \|w\|$  (2)

But suppose we don't know $m$ from advance, in that case, I define it as in a form such that

$$\underset{w, b, \xi}{\arg\min} \; \lambda\, \sigma(w) + \sum_i \xi_i \quad s.t.$$
$$y_i\,(w^T x_i + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0$$

such the length of my normal vector, in $\lambda$. so I no longer specify the margin $\|m\|$, instead I have introduced a new user designed parameter $\lambda$, given this parameter $\lambda$, I can solve this problem, corresponding to overfitting, term. so this would be a tendency to overfit and as you have seen it increases with decreasing "$m$" margin bcz $\|w\| \leq \frac{1}{m}$, is it in something that increases with $\|w\|$. And what people mostly use in the quadratic squared euclidean length of the normal vector

$$\sigma(w) = w^T w$$

→ So, for the best model, we were trying to find a model that minimizes $\underset{(w, b)}{\arg\min} \; \sum_{i=1}^{n} \xi_i$

that dependend on $w$ & $b$, but we know, bcz of overfitting thing, we have to add an extra term $\lambda \cdot \sigma(w)$ to get the most/best classifier

so, important choice in the one that you see here, so this would be the squared $l_2$ - Norm

or another choice would be to just take the $l_1$ norm of $w$

so, $\sigma(w) = w^T w$

so, we will consider, $\sigma(w) = w^T w$ or $\|w\|$     (2)

choosing L2 regularized term

we will derive the optimization & actions

→ so the tendency to overfit is the regularization term whereas the other thing here the slack variable ($\xi$) is the data term

given $\lambda$

$$\underset{w, b, \xi}{\arg\min} \quad \lambda \cdot \sigma(w) + \sum_{i=1}^{n} \xi_i \quad s.t. \quad y_i(w^T x_i + b) \geq 1 - \xi_i$$

regularization / tendency to overfit     data term    (A)     $\xi_i = 0$    $i = 1..N$

— so we have this formula, we have these formula constraints, we now need to make a specific choice as to the nature of this regularization term. we use L2 regularization $\sigma(w) = w^T w$ — (B)

so the formula (A) became using (B)

$$\underset{w, b, \xi}{\arg\min} \quad \frac{1}{2} w^T w + c \sum_{i=1}^{n} \xi_i \quad s.t. \quad y_i(w^T x_i + b) \geq 1 - \xi_i;$$
$$\xi_i \geq 0$$

→ so this formula is in primal domain. now with optimization problems, such the as these, we here have a quadratic program (quadratic objective) function & linear constraints. It would be perfectly fine to optimize it here, but sometimes we need to go go dual domain for solving certain questions.

→, So since there are inequalities the constraints, so instead of using the Lagrangian formalism. we will use KKT equation or

also called "Lagrangian Inequalities"
- So I am setting up a Lagrangian function here
which involves both primal variables $(w, b \& \xi)$
and new dual variables the lagrange multipliers
$\alpha$ k u.

so now the equation becomes ↗ penalization of the slack variable

$$L(w, b, \xi, \alpha, \mu) = \frac{1}{2} w^T w + c \sum_{i=1}^{\hat{n}} \xi_i - \sum_{i=1}^{N} \alpha_i (y_i (w^T x_i + b) - 1 + \xi_i) - \sum_{i=1}^{n} \mu_i \xi_i$$

regularization
term

◯①

↳ and these other terms came from from the constraints
of inequalities, so what I have done here is that
since I have 1 equation for each observation here and
I'm now bringing what is on the RHS to the LHS

so I have $y_i (w^T x_i + b) - 1 + \xi_i$ and I am
multiplying this with the Lagrange multiplier $\alpha_i$,
and now since I have equation ① summed over that
many Lagrange multipliers $\alpha_i's$, here

so

$$\sum_{i=1} \alpha_i (y_i (w^T x_i + b) - 1 + \xi_i) - \sum \mu_i \xi_i$$

↳ similarly we have $\mu_i$ for this
-- so, what it says is that there $\xi \geq 1$, we want to
look for an optimum, it has to fulfill the following
properties, the Lagrangian ~ this equation here
must be at a minimum w.r.t to the primal variables
$w, b, \& \xi$ and it must be at a maximum with
respect to also an extremum but at a
maximum w.r.t to the dual variables and at

this saddle point it is minimum w.r.t to a primal
maximum w.r.t to dual variables &#9315;

—optimum given by saddle point : minimum w.r.t.
primal variable & maximum w.r.t dual
variables



$\angle(w,b,\xi)$    & $\alpha_i, \mu_i \geq 0$

$(\alpha, \mu)$

$(w^\circ, b^\circ, \xi^\circ, \alpha^\circ, \mu^\circ)$

— this is where we have an optimum, I can indicate
the optimum parameters by the superscript 0,
these lagrange multipliers have to be non-ve. ie.
$\alpha_i, \mu_i \geq 0$, we can go & find the extremum by
differentiating this Lagrangian w.r.t its primal
variables, so we differentiate w.r.t to
the normal vector (w), bias (b) & slack variables

(&#958;) & we set the derivative to 0

$$\frac{\partial L}{\partial w} = w - \sum_i \alpha_i y_i x_i = 0 \qquad \boxed{w = \sum_{i=1}^{n} \alpha_i y_i x_i}$$

$\hookrightarrow$ (A)

$$\frac{\partial L}{\partial b} = -\sum_i \alpha_i y_i = 0 \qquad \sum_{i=1}^{N} \alpha_i y_i = 0 \;-(2)$$

$$\frac{\partial L}{\partial \xi_i} = c - \alpha_i - \mu_i = 0 \qquad \alpha_i + \mu_i = c \;-(3)$$

$$\boxed{\text{so, the margin is } w \quad w = \sum_{i=1}^{n} \alpha_i y_i x_i}$$

— Now I'm starting with the first term going
. @ we will see that that the sum here is
going to be sparse, so I don't need all training
examples to express my natu normal vector
& this is one of the great things about SVM

So, we have found the conditions that describe how we can find a minimum w.r.t the primal variables what remains is to maximize the Wolfe dual w.r.t the dual varibles.

so

$$\underset{\alpha, u}{\text{argmax}} \quad L(w^\circ, b^\circ, \mathcal{E}^\circ, \alpha, u) = \tfrac{1}{2} w^\circ{}^T w^\circ$$

optimal values have already been put using equations (A), (2) & (3)

so using (1), (A), (2) & (3)

we get

$$L(w^\circ; b^\circ; \mathcal{E}^\circ, \alpha^\circ, u^\circ) = \tfrac{1}{2} w^T w + c \sum_{i=1}^{N} \mathcal{E}_i - \sum_{i=1}^{N} \alpha_i (y_i (w^T x_i + b) - 1 + \mathcal{E}_i) - \sum_{i=1}^{N} u_i \mathcal{E}_i$$

solving & substituting the value of from (A), (2)
& (3), so

$$\underset{\alpha, u}{\text{argmax}} \quad L(w^\circ, b^\circ, \mathcal{E}^\circ, \alpha, u) = \tfrac{1}{2} w^\circ{}^T w^\circ + \sum \mathcal{E}_i (\alpha_i + y_i) - w^\circ{}^T w^\circ - \underbrace{\sum_{i=1}^{N} \alpha_i (y_i b - 1 - \mathcal{E}_i)}_{A} - \sum_{i=1}^{N} u_i \mathcal{E}_i$$

we still have the term
here that depends on b
i.e. b(0)

$$\underset{\alpha, u}{\text{argmax}} \quad L(w^\circ, b^\circ, \mathcal{E}^\circ, \alpha, u) = \tfrac{1}{2} w^\circ{}^T w^\circ + \sum_i \alpha_i \text{ such}$$

$$\oplus \alpha_i \geq 0, \quad u_i \geq 0$$

that this term here is quadratic in w & the sum over $\alpha_i$ from the summing one here & this is subject to the equations that we had here, and to the dual variables being - non - +e have eqng (A), (2) & (3) where we have replaced w by $\sum_{i=1}^{n} \alpha_i x_i y_i$.

$$= \operatorname*{argmax}_{\alpha, u} -\frac{1}{2} \sum_i \sum_j y_i y_j \alpha_i \alpha_j x_i^T x_j + \sum \alpha_i \qquad (4)$$

such that

$$\alpha_i \geq 0$$
$$u_i \geq 0$$
$$\alpha_i + u_i = c \quad \Big\{ \begin{array}{l} 0 \leq \alpha_i \leq c \\ \sum \alpha_i y_i = 0 \end{array}$$

-- These along with equation ① & ② say that the lagrange multiplier $\alpha_i$, they anyway have to be bounded below by 0, but they must also not becomes larger than C which was the user-defined constraint that we put into backoff the regularization versus the data term.

" So this is the optimization problem that we need to solve

i.e.

$$\mathcal{L}(\alpha, u, \omega^0, b^0, \xi^0) =$$
$$\operatorname*{argmax}_{\alpha, u} -\frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i y_j \alpha_i \alpha_j \underbrace{x_i^T x_j}_{k(x_i, x_j)} + \sum \alpha_i \quad \text{such}$$

that

$$\alpha_i \geq 0 \quad \alpha_i + u_i = c \Big\}$$
$$u_i \geq 0 \quad 0 \leq \alpha_i \leq c$$
$$\sum_{i=1}^{n} \alpha_i y_i = 0$$

→ so, just like the primal, here also we have a product of 2 dual variables $\alpha_i \& \alpha_j$ but we can write this whole thing in matrix form as follows:- so rather than maximizing, we may wish to minimize to but this to a standard form for a quadratic programming problem, then we have to invert the sign of and we had an inner product between two vectors $x_i \&$ $x_j$. I am inserting this inner products. I am

collecting all of these inner products in a kernel matrix, so the element i,j of this kernel matrix is going to be $x_i$ transpose times $x_j$ which is written here, then I have these vectors $\alpha$ summarizing all of my langarrian multipliers and one more term in there, we have a missing term here, the sum over the $\alpha$ can be written as 1 transpose alpha, where there 1 is vector of all 1's

$$= \underset{\alpha}{\arg\min} \quad +\frac{1}{2}\alpha^T \underbrace{Y K Y}_{q} \alpha + L^T\alpha \qquad \text{such that } 0 \leq \alpha_i \leq C$$
$$\alpha \cdot y = 0$$

standard QP problem

kernel matrix
$(K)_{ij} = x_i^T x_j$

diagonal matrix holding labels

→ one thing to note here than I have not mentioned in v so, I can collect here in a diagonal matrix, i.e. holding the $v$ labels, which must either be $+1$ or $-1$ along with the diagonals.

→ so this is something that you can put into any standard QP solver specifically. here I have underlined in green the inputs that you need to provide - so, this is what you want to solve for but $V$ the labels are given in the training set, the kernel matrix you can compute using this recipe and $c$ is a constant supplied by the user, to specify how strongly do we want to regularize = how closely we want to fit the data, ok - so this you can really put into it's over and once you have done that you get out a collection of $\alpha$'s and now how do you relate them to the solution in the primal domain, you can again remember that (so normal vector $w$ = the simply sum over ? $\alpha$, $w = \sum_i^N \alpha_i x_i y_i$.

after optimizing for the alpha here, here you can insert what you found & then get the normal vector to the primal domain. Now I mentioned earlier that perhaps we don't need to involve all of the observations indeed at this is important for SVM, & the reason is given by KTT called complementary slackness condition which states that the product of Lagrange multiplier and the equality ~ inequality tries to enforce the product of these two has to be zero, now for this product to be zero

$$\alpha_i \left( y_i (a^T x_i + b) + 1 - \xi_i \right) = 0$$

**2) part2) Point out what is the "margin" in both the primal formulation and the dual formulation,**

**Soln: As suggested in the above derivation**

So, we can say $w^Tx+b > 0$    $y=+1$ (y is +ve)    4

       If $w^Tx+b < 0$    $y=-1$ (y is -ve)

       $w^Tx+b = 0$ (can by anything)
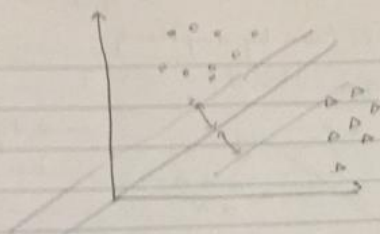
## Idea of max margin



→ 2 types of points are there ooo & ∆∆

How to define margin ? → First we define the slope, let's say we have the slope, what we do in the first we draw a line with that slope going through origin then we start moving it in either direction, we will move it up little bit llel & whenever we hit a point on either side, maybe we will hit a circle first or a ∆gle first, may be we will hit both at the same time., but which ever, we hit, we stop there, this is 1 margin, bez after



same distance

this point I will start making mistakes & then I draw a symmetric

line over the other side with the same distance, then the line i.e on the middle of it (need not go through the origin) will be my decision boundary, so, this will be called $w^Tx+b = 0$,



$w^Tx+b = -1$
$w^Tx+b = 0$
$w^Tx+b = 1$

But ...

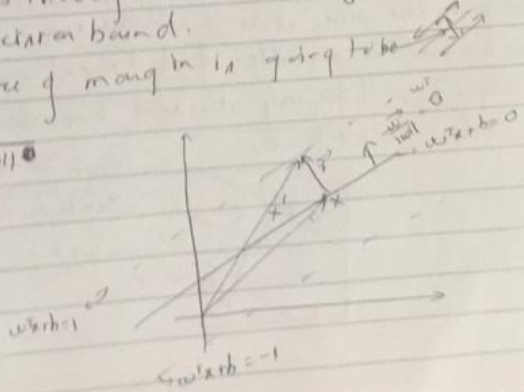So, I am defining these as my boundaries, so now I want to find w & b such that this whole distance is maximized



Let's say, that again, let's try to fit w, any w, then I start going up & down, so with same increment I keep going up I say hey at this point, I hit a +ve point so, that will be my one side margin, then I am going to rescale my w & b, so that the equation of the line like this $w^T x + b = cL$ ... something

- Remember, if $w^T x + b = L$ in a line, then you can always multiply it with same constant it will be the same line, so. I basically adjust my weights such that this $w^T x + b = L$ in the equation of my line, then I look at a line which is opposite to it. $w^T x + b = -1$, & then the line that goes through it middle of it $w^T x + b = 0$ will be my decision boundary.

So the length/size of margin is going to be

$$= \frac{2}{||w||}$$



$w^T x + b = 1$

$w^T x + b = -1$

$\hat{r} = r \frac{\vec{w}}{\|w\|}$ magnitude

$\hat{u} = \frac{\vec{w}}{\|w\|}$

so $\vec{r}' = \vec{x}' - \vec{x_0}$

$\vec{w}^T \vec{r}' = w^T x' - w^T x_0$

$= \hat{u}^T \vec{r} \frac{w}{\|w\|} = -b - e(1-eb) \cdot (1-b) - b(-b)$

$= \frac{r \, w^T w}{\|w\|} = +1 \Big]$    $r\|w\| = +1$

$\gamma = \frac{+1}{\|w\|}$

So, now, when you derive the same quantity for this side

this length $= \frac{1}{\gamma}$



size of the margin $= 2\gamma = \frac{2}{\|w\|}$

Now, the task has become like this, we want to find a line $w$ & $b$, which not only separates the 2 types of points

So, we want to learn $w,b$ such that for all points for which $y_i = +1$, $w^T x_i + b > 0$ &for all points $x_i$ for which $y_i = -1$ $\}$ zero training error $w^T x_i + b < 0$

and $\boxed{\frac{2}{\|w\|}}$ is as large as possible $\Big]$ → support vector machines

If for this we can write $y_i(w_i^T x_i + b) > 0$ bcz if $y_i$ is $-ve$ then we want $(w_i^T x_i + b)$ to be

+ve as well.

How do we learn this $w$ & $b$?
→ We are going to assume that the data is linearly separable, that is, we accommodate no loss.
in past

An how are we going to learn $w$&$b$? → either we try to minimize some loss or maximize some likelihood, here we are going to pose it in a way — we are going to maximize the margin i.e $\frac{2}{\|w\|}$ → minimize $\|w\|$

Optimization: maximize margin while incurring zero training error.

Max $\frac{2}{\|w\|}$ with 0 loss → minimize $\frac{\|w\|}{2}$ with 0 loss.

→ Minimize $\frac{\|w\|_2^2}{2}$ with 0 loss

$\frac{w^Tw}{2}$ $\frac{\sqrt{w^Tw}^2}{2}$

Minimize $\frac{w^Tw}{2}$ with 0 loss.

→ This is what the optimization criteria is.
i.e minimize $\frac{\|w\|^2}{2}$
w.r.t

subject to $y_n(w^Tx_n + b) \geq 1$, $n=1,...,N$ (for all the training data points, this should be always $\geq 1$

→ so we not only want the data points to be on right side, correct side of the margin, i.e. if it is +ve we want it be on +ve side and if it is -ve, on the -ve side. but we actually want it be to on the other side of margin. We don't want any points inside the margin at all, so that is determined by

$(y_n)(w^Tx_n+b) \geq 1$ — this

bcz, if this were 0, all I am saying is that, 6
try to find a line such that, maximizes the margin
& points are on the correct side of the decision boundary

1)

So, thisisn a constrainted optimization, subject to certain
contraint.

so, Optimization with N linear inequality constraint
↓ no. of data points

→ So, large margin = small $|w|$
→ small $|w|$ ⟹ regularized/simple solutions
→ simple solution ⟹ Better generalizability

so      $\underset{w,b}{\text{minimize}} \ \dfrac{\|w\|^2}{2}$

subject to $y_n(w^T x_n + b) \geq 1, \quad n=1 \cdots N$

→ so, this is an quadratic objective function to minimized
   N constraints

Basic Optimization

$\underset{x,y}{\text{minimize}} \ f(x,y) = 2 - x^2 - 2y^2$

subject to $h(x,y) = x + y - 1 = 0$  } our in equality constraint

How to do it?
An so, there is a notion of Lagrangian multiplier ($\beta$), that
Lets you combine the two equation in L.

so, what you do in that you define a new objective function
which   $\underset{x,y,\beta}{\text{minimize}} \ L(x,y,\beta) = f(x,y) - \beta g(x,y)$

→ It has been shown that, if you find a solution to this

$\dfrac{\partial y}{} \quad \cdots y \cdot v \cdot v \quad \dfrac{\partial L}{\partial \beta} = x + y - 1 = 0$

Let's go back to SVMs and re-write our optimization program as the following convex optimization:

$$\min \frac{1}{2}\|w\|^2 \quad s.t. :$$
$$y^i(w \cdot x^i + b) \geq 1, \quad i = 1, ..., m$$
$$g_i(w, b) = -y^i(w \cdot x^i + b) + 1 \leq 0$$

Following the KKT conditions, we get $\alpha_i > 0$ only for points in the training set which have a margin of exactly 1. These are the Support Vectors of the training set. Figure 10.4 shows a maximal margin classifier and its support vectors.

Let's construct the Lagrangian for this problem:
$L(w, b, \alpha) = \frac{1}{2}\|w\|^2 - \sum_{i=1}^{m} \alpha_i [y^i(w \cdot x^i + b) - 1]$.

Now we will find the dual form of the problem. To do so, we need to first minimize $L(w, b, \alpha)$ with respect to $w$ and $b$ (for a fixed $\alpha$), in order to get $\Theta_D$. We will do this by setting the derivatives of $L$ with respect to $w$ and $b$ to zero. We have:

$$\nabla_w L(w, b, \alpha) = w - \sum_{i=1}^{m} \alpha_i y^i x^i = 0, \tag{10.20}$$

which implies that:

$$w^* = \sum_{i=1}^{m} \alpha_i y^i x^i. \tag{10.21}$$

This gives a formula for the optimal $w$ given $\alpha$.

which can be found out by finding w* is the optimum w here because w*=summation of alpha * y* x (all three quantities going from 1 to N). So given alpha we can find the w
. and we know that margin is **2/|w|** .

When we take the derivative with respect to $b$ we get:

$$\frac{\partial}{\partial b}L(w,b,\alpha) = \sum_{i=1}^{m} \alpha_i y^i = 0 \tag{10.22}$$

This identity gives us a restriction on $\alpha$.

We then take the definition of $w^*$, which we derived in (10.21), plug it back into the Lagrangian, and we get:

$$L(w^*,b^*,\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2}\sum_{i,j=1}^{m} y^i y^j \alpha_i \alpha_j x^i x^j - b\sum_{1=1}^{m} \alpha_i y^i. \tag{10.23}$$

From (10.22) we get that the last term equals zero. Therefore:

$$L(w^*,b^*,\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2}\sum_{i,j=1}^{m} y^i y^j \alpha_i \alpha_j x^i x^j = W(\alpha). \tag{10.24}$$

We end up with the following dual optimization problem:

$$\max W(\alpha) \quad s.t. :$$
$$\alpha_i \geq 0, i = 1, ..., m$$
$$\sum_{i=1}^{m} \alpha_i y^i = 0$$

The KKT conditions hold, so we can solve the dual problem, instead of solving the primal problem, by finding the $\alpha^*$'s that maximize $W(\alpha)$ subject to the constraints. Assuming we found the optimal $\alpha^*$'s we define:

$$w^* = \sum_{i=1}^{m} \alpha_i^* y^i x^i \tag{10.25}$$

which is the solution to the primal problem. We still need to find $b^*$. To do that, let's assume $x^i$ is a support vector. We get:

$$1 = y^i(w^* \cdot x^i + b^*) \tag{10.26}$$
$$y^i = w^* \cdot x^i + b^* \tag{10.27}$$
$$b^* = y^i - w^* \cdot x^i \tag{10.28}$$

## 2)Part3 what are the benefits of maximizing the margin.
## Ans: Because of following reasons:

1)A large margin effectively corresponds to a regularization of SVM weights which prevents overfitting. Hence, we prefer a large margin (or the right margin chosen by cross-validation) because it helps us generalize our predictions and perform better on the test data by not overfitting the model to the training data.

2) "Maximizing the margin seems good because points near the decision surface represent very uncertain classification decisions: there is almost a 50% chance of the classifier deciding either way. A classifier with a large margin makes no low certainty classification decisions. This gives you a classification safety margin: a slight error in measurement or a slight document variation will not cause a mis-classification."

3) SVM maximizes margin, so the model is slightly more robust (compared to linear regression) but more importantly: SVM supports kernels, so you can model even non-linear relations

## 2)Part4 Characterize the support vectors.
## Soln:

The optimality criterion (13) characterizes which training examples become support vectors. Recalling (10) and (14)

$$g_k - y_k \rho = 1 - y_k \sum_{i=1}^{n} y_i \alpha_i K_{ik} - y_k b^* = 1 - y_k \, \hat{y}(\mathbf{x}_k). \tag{16}$$

Replacing in (13) gives

$$\begin{cases} \text{if } y_k \, \hat{y}(\mathbf{x}_k) < 1 & \text{then } \alpha_k = C, \\ \text{if } y_k \, \hat{y}(\mathbf{x}_k) > 1 & \text{then } \alpha_k = 0. \end{cases} \tag{17}$$

This result splits the training examples in three categories[2]:

- Examples $(\mathbf{x}_k, y_k)$ such that $y_k \, \hat{y}(\mathbf{x}_k) > 1$ are not support vectors. They do not appear in the discriminant function because $\alpha_k = 0$.

- Examples $(\mathbf{x}_k, y_k)$ such that $y_k \, \hat{y}(\mathbf{x}_k) < 1$ are called *bounded support vectors* because they activate the inequality constraint $\alpha_k \leq C$. They appear in the discriminant function with coefficient $\alpha_k = C$.

- Examples $(\mathbf{x}_k, y_k)$ such that $y_k \, \hat{y}(\mathbf{x}_k) = 1$ are called *free support vectors*. They appear in the discriminant function with a coefficient in range $[0, C]$.

Let $\mathcal{B}$ represent the best error achievable by a linear decision boundary in the chosen feature space for the problem at hand. When the training set size $n$ becomes large, one can expect about $\mathcal{B}n$ misclassified training examples, that is to say $y_k \hat{y} \leq 0$. All these misclassified examples[3] are bounded support vectors. Therefore the number of bounded support vectors scales at least linearly with the number of examples.

When the hyperparameter $C$ follows the right scaling laws, Steinwart (2004) has shown that the total number of support vectors is asymptotically equivalent to $2\mathcal{B}n$. Noisy problems do not lead to very sparse SVMs. Collobert et al. (2006) explore ways to improve sparsity in such cases.

This is also explained below as part of solution of the above mentioned question:

i.e.

$$\arg\max_{\alpha,u} -\frac{1}{2}\sum_i\sum_j y_i y_j \alpha_i \alpha_j \underbrace{x_i^T x_j}_{k(x_i,x_j)} + \sum_i \alpha_i \quad \text{such that}$$

$$\alpha_i \geq 0 \qquad \alpha_i + u_i = c \quad \Big\} \quad 0 \leq \alpha_i \leq c$$
$$u_i \geq 0 \qquad \qquad \sum_i \alpha_i y_i = 0$$

all this equation is a bit redundant here and in fact this can be written even more succinctly as follows to bring out even more clearly the quadratic structure we have. So this is also a quadratic objective just like we had in the primal here, we have a product of 2 dual variables $\alpha_i$ & $\alpha_j$ but we can write this whole thing in matrix form as follows:

So, rather than maximizing we may wish to minimize to put this to standard form for a quadratic programming problem then we have to invert the sign of course and now here we had an inner product between two vectors $x_i^T \& x_j$, I am inverting these inner products. I'm collecting all of these inner products in a kernel matrix so the element $i,j$ of this kernel matrix is going to be $x_i$ transpose time $x_j$ which is written here, then I have these vectors $\alpha$ summarizing all of my lagrangian multipliers and one more term is there we have a missing from here the sum over the $\alpha$ can be written as 1 transpose $\alpha$, where this is a vector of all 1's

$$\text{argmin}_{\alpha} + \tfrac{1}{2}\alpha^T \underline{\Gamma \bar{K} \bar{Y}} \alpha + \bar{1}^T \alpha \quad \text{s.t. } 0 \le \alpha_i \le C$$
$$\alpha, y \ge 0$$

standard QP problem — kernel matrix $(K)_{ij} = x_i^T x_j$ — diagonal matrix holding labels

- one thing that I have not mentioned in y for 1 can collect these in a diagonal matrix that is holding the m labels, which must be $+1$ or $-1$ along with its diagonals

- so this is something that you can put into any standard QP solver. specifically, here I have underlined in green the inputs that you need to provide so. this is what you want to solve for but $y$ the labels are given in the training set, the kernel matrix you can compute using this recipe. and $c$ is a constant supplied by the user to specify how strongly do we want to regularize ~ how closely do we want to fit the data. & so that you can really put into it's over and over and once you have done that you get out a collection of $\alpha$'s and now how do you relate them to the solution in the primal domain. if one can again remember that for normal vector $w$, the simply sum over $i$ alpha

$$w^* = \sum \alpha_i y_i x_i \quad \text{after optimizing for the}$$

alphas here, you can inset what you found.

⟶ & this gets normal vector in the primal domains. Now I mentioned earlier that

perhaps we don't need to involve all of the observations, indeed this is one of the interesting perhaps of SVM and the reason is given by KKT conditions called complementary slackness. conditions which states that the product of Lagrange multiplier and the equality or inequality tries to enforce the, product of these two has to be zero in our case this product to be zero
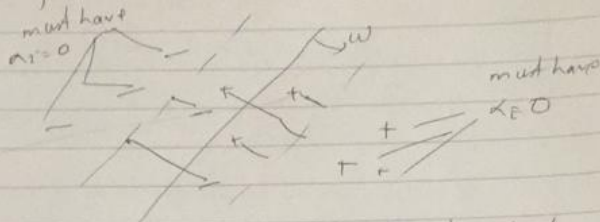
$$\alpha_i \left( y_i (w^T x_i + b) - 1 + \xi_i \right) = 0$$

larger than $\xi_i$ true for all training points on correct side of their margin

→ obviously one of the two factors has to be zero and let's check these in turn this is lagrange in multiplier ; $\alpha_i$ on the other hand. (lagrange in multiplier.

$y_i (w^T x_i + b)$ = was the score so how much I was lying on the correct or wrong side of the decision boundary, this is minus 1, the margin that I want to subtract and if a point lies on the right side of the margin then this thing here cannot be 0 ; $\xi_i$ (here, then there will be nonslack), so if slack is required and remember that slack appeared in the primal objective function the program was explicitly asked to minimize. If so thing will be $\xi_i = 0$ & the term — $(y_i (w^T x_i + b) - 1)$ will be some term

larger than 0

this is exactly what makes for sparsity what
this very interesting property of SVM, no, the slack
of a couple of variables, I have indicated here



this would be the amount of slack ~ the distance
that would has from the correct side of margin
~ that the positive point has from it's side of the
margin

Imp
 — My normal vector $w$, will now be specified as a
linear combination of only of the support vectors and
not of these other points out here.

$$K_{i,j} = x_i^T x_j$$

 — for matrix $k$ is a measure of similarity for observations
 $i \& j$
 → can use other similarity measures in particular
 for non-vectorial data you will want to use their
 own similarity measures the only thing that they have
 to satisfy is that this matrix $k$ should remain
 +ve semi-definite bcz as long as it is we have
 a convex optimization problem and a unique optimum
 if this matrix is not +ve semi-definite then
 optimization would much, much harder. $k$

# 2)Part5
# Point out the benefit of solving the dual problem instead of the primal problem.

## Soln: Below is the answer:

**1)** **Understanding the dual problem leads to specialized algorithms for some important classes of linear programming problems.** Examples include the transportation simplex method, the Hungarian algorithm for the assignment problem, and the network simplex method. Even column generation relies partly on duality.

2) **The dual can be helpful for sensitivity analysis.** Changing the primal's right-hand side constraint vector or adding a new constraint to it can make the original primal optimal solution infeasible. However, this only changes the objective function or adds a new variable to the dual, respectively, so the original dual optimal solution is still feasible (and is usually not far from the new dual optimal solution)

3. **Sometimes finding an initial feasible solution to the dual is much easier than finding one for the primal.** For example, if the primal is a minimization problem, the constraints are often of the form $Ax \geq b$, $x \geq 0$, for $b \geq 0$. The dual constraints would then likely be of the form $A^T y \leq c$, $y \geq 0$, for $c \geq 0$. The origin is feasible for the latter problem but not for the former.

4. **The dual variables give the shadow prices for the primal constraints.** Suppose you have a profit maximization problem with a resource constraint $i$. Then the value $y_i$ of the corresponding dual variable in the optimal solution tells you that you get an increase of $y_i$ in the maximum profit for each unit increase in the amount of resource $i$ (absent degeneracy and for small increases in resource $i$).

# 5)

- One of the important advantage of using the dual form in SVM is that it allow us to apply kernels. Kernel search an optimal separating hyperplane in a higher dimensional space without increasing the computational complexity much. Kernel can be applied, if the algorithm takes the features in terms of its inner product ($x_i^T x_k$). Please see the dual form of the objective function (eq.2) which inputs the features in terms of its inner product.

- See the equation 3 (i.e. Complementary slackness condition). (i) If $\alpha_i \geq 0$, then $y_i(x_i^T \beta + \beta_0) = 1$. In other words, $x_i$ is on the boundary of the hyperplane. (ii) If $y_i(x_i^T \beta + \beta_0) \geq 1$, then $x_i$ is not on the boundary of the hyperplane and $\alpha_i = 0$. From equation 4, we can say that $\beta$ is defined in terms of linear combination of $x_i$ (which $\alpha_i \geq 0$, i.e $x_i$ lies on the boundary of the hyperplane.) So, $x_i$ are called as support points. We knew that usually the number of support points for SVM are few. It means that most of the $\alpha$ (dual variable) are zero.

- There are some algorithms like SMO(Sequential Minimal Optimization) solves the dual problem efficiently.

**A proper explanation is below:**

...we now need to make a specific choice as
to the nature of this regularization term, we use $l_2$ regularized

so, the formula $\sim$ norm $\sigma(w) = w^T w$

so the formula remains same except that $\|\cdot\|$ is $\sigma(w)$
$= w^T w$ replaced.

so given. C

$$\underset{w, b, \xi_i}{\arg\min} \quad \frac{1}{2} w^T w + c \sum_i \xi_i \quad s.t. \quad y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

→ so the formula $\sim$ is no called primal domain,
now, with optimization problems, such as these, we
here have a quadratic program ('quadratic objective
function), and linear constrains : It would be
perfectly okay to optimize it and in this primal
domain, however it turns out to sometimes be
beneficial to go to the dual domain and I will
say in a moment what that dual is, the reasons
are that if it is it can be more efficient if the
no. of no. of features is much much larger than
& no. of observations $p >> n$

features        observations

So, if we have 10,000 features, then each of these inner products here and involving the normal vector

$$y_i(w^Tx_i+b)$$

will always involve, well 10,000 numbers and even after training, if you want to make predictions, I will always have to take an inner product between my observations will would be 10,000 dimensional and this normal vector and in the dual domain it turns out that we operate with variables with as many variables at most as we have observations, rather than features.

→ In addition, we will discuss, so called kernel trick which allows to use all of this formalism so that so far has involved hyperplanes only, to use this formalism also for non-linear classification problems.

→ But up here if $y_i(w^Tx_i+b) \geq 1-\varepsilon_i$

$$\varepsilon_i \geq 0$$

were equalities, we could have used the Lagrangian formalism but we have inequalities and the for inequalities the proper approach is to use the KKT equations, also called "Lagrangian to inequalities".

So, I am setting up a Lagrangian function here, which involves both primal variables and, $w$, $b$ and $\varepsilon$, and new dual variables: the Lagrange multipliers $\alpha$ & $\mu$.

So, I am copying the same terms here

2. (8 points) Identify the Lagrange dual problem of the following primal problem:

Given features $(x_1, y_1), \cdots, (x_N, y_N)$, where $y_1, \cdots, y_N \in \{-1, 1\}$,

Minimize $w^T \cdot w + C \sum_{i=1}^{N} \xi_i$, the weighted sum between the squared length of the separating vector and the errors, where

w is the separating vector, $w^T \cdot w$ is the dot product, and $\xi_i$ is the error made by separating vector w on feature $(x_i, y_i)$.

Subject to $y_i \cdot (w^T \cdot x_i) \geq 1 - \xi_i$ and $\xi_i \geq 0$ for $i = 1, \cdots, N$. In other words, if the "normalized feature" $y_i x_i$ has a margin less than 1,

$w^T \cdot (y_i x_i) \leq 1$, we add a slackness term to make it 1.

Point out what is the "margin" in both the primal formulation and the dual formulation, what are the benefits of maximizing the margin. Characterize the support vectors. Point out the benefit of solving the dual problem instead of the primal problem.

**2ⁿᵈ solution: My second at taking a different approach to the solution of problem 2 of this assignment**

Given features $(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$

where $y_1, \ldots y_n \in [-1, 1]$

Minimize the $w^T \cdot w + c \sum_{i=1}^{N} \xi_i$

$$f(x) = c \sum_{i=1}^{N} \xi_i \quad (say)$$

dual function $L^*(u) = \min_{\xi \in p} f(\xi) + w^T w(\xi)$

$$= \min_{\xi \in p} \lambda \left[ -f(\xi) + w_1^T w_g(\xi) \right] + [1-\lambda] \left[ f(\xi) + w_2^T g(\xi) \right]$$

$$\geq \lambda \left[ \min_{\xi \in p} f(\xi) + w_1^T w(\xi) \right] + (1-\lambda) \left[ \min_{\xi \in p} (f(\xi) + w_2^T g(\xi)) \right]$$
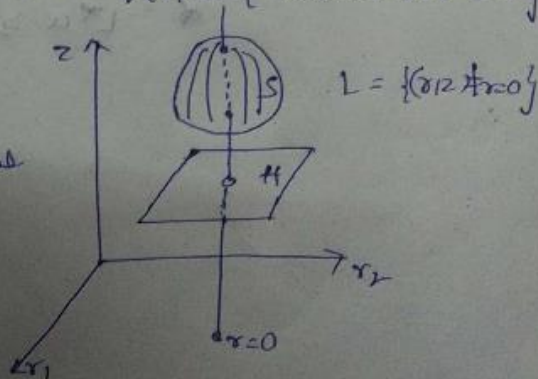
$$= \lambda L^*(u_1) + (1-\lambda) L^*(u_2)$$

let us consider primal problem from a "resource cost" point of view for each $\xi \in p$ we have array of resource and costs associated with $\xi$

$$\binom{r}{z} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_m \\ z \end{bmatrix} = \begin{bmatrix} w_1(\xi) \\ w_2(\xi) \\ \vdots \\ w_m(\xi) \\ f(\xi) \end{bmatrix} = \begin{bmatrix} g w(\xi) \\ f(\xi) \end{bmatrix}$$

we call the region $s$

$$S = \left\{ (r, z) \in R^{m+1} \mid (r, z) = (w(\xi), f(\xi)) \text{ for } \xi \in p \right\}$$



The column geometry of primary and dual problem

$$L = \{(r, z) \mid r = 0\}$$

The feasible solution with lowest cost corresponds to lowest point on intersection of S and L. The feasible lowest value is exactly the value of primal problem.

→ In the figes shown hyperplane H, lies below the region S

Suppose we try to find that hyperplane H lying below S whose intersection with L is large as possible

let formulate it

$$H = H_{w,\alpha} = \left\{ (\gamma, z) \in \mathbb{R}^{m+1} \mid z + w^T \gamma = \alpha \right\}$$

$$L = \left\{ (\gamma z) \in \mathbb{R}^{m+1} \mid \gamma = 0 \right\}$$

The intersection $H_{w,\alpha}$ with L is $(H \cap L)$ occur at point $(\gamma, z) = (0, \alpha)$

→ the $H_{w,\alpha}$ lying below S whose intersection with L is highest

maximum $w, \alpha$     $\alpha$
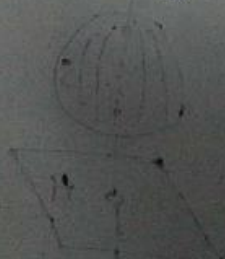
State that $H_{w,\alpha}$ lies below s.

Conditions
For max $w, \alpha$

$$z + w^T \gamma \geq \alpha \text{ for all } (\gamma, z) \in S$$

$$f(\xi) + w^T w \geq \alpha \text{ for all } \xi \in P$$

$$L^* u \geq \alpha$$

**3. (Optional) Formulate the primal problem and derive the dual problem if there are multiple classes.**
**Ans: The solution remains more and less the same other than a few changes mentioned below:**

# Bonus Question Solution

Soln :- For a multiclass classification problem, a new multi-class classifier called NHCMC - Nonparallel hyperplanes classifier to multiclass classification

Consider the multiple classification problem with the training set :

$$T = \{(x_1, y_1), \ldots (x_l, y_l)\}, \qquad -(14)$$

where $x_i \in R^n$, $i = 1 \ldots l$, $y_i \in \{1, \ldots k\}$ is the corresponding pattern of $x_i$.

For multiple classification, we seek K nonparallel hyperplanes :

$$(w_k \cdot x) + b_k = 0, k = 1 \ldots, k \qquad -(15)$$

For convenience, we denote the no. of each class of the training set (14) as $l_k$ and the points belonging to $k^{th}$ class as $A_k \in R^{l_k \times n}$, $k = 1, \ldots k$. Besides, we define the matrix

$$B_k = [A_1^T, \ldots, A_{k-1}^T, A_{k+1}^T, \ldots, A_k^T]^T$$

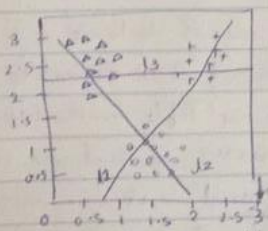as all the points except for the points belonging to $k^{th}$ class



Fig: A toy example learned by the linear NHCMC

## Linear case.

### The primal Problem
We seek to construct k nonparallel hyperplanes (15) by solving the following convex quadratic programming problems (QPB).

$$\min_{w_k, bb, \eta_k, \xi_k} \frac{1}{2} C_i \|w_k\|^2 + \frac{1}{2} \eta_k^\top \eta_k + C_2 e_n^\top \xi_k$$

s.t. $\quad B_k w_k + e_k bb + \eta_k$,

$$(A_k w_k + e_{k2} b_k) + \xi_k \geq e_{k2} \qquad - (17)$$

$$\xi_k \geq 0$$

where $\eta_k \in R^{(l_1 - l_k)}$ is a variable and $\xi_k$ is a slack variable, $e_{k_1} \in R^{(l_1 - l_k)}$ and $e_{k_2} \in R^{l_k}$ are the vectors of one. $C_1 \geq 0$ and $C_2 \geq 0$ are penalty parameters.

In order to illustrate the primal problem of NHCMC, we generated an artificial two dimensional three-class dataset. The geometric interpretation of above problem with $x \in R^2$ is shown in Figure above, where minimizes the sum of the squared distance from the hyperplane of $k-1$ classes, that is all classes except for those of the $k$-th class, and the points of the $k$-th class are far from the $i$th hyperplane. Take the "×" class in Figure 1 as an example. We hope the hyperplane of the the "×" class $i$ is far from the "×" points and closest the "o" and triangular points. In order to minimize classification, the points of the $k$-th class are at distance 1 from the hyperplane, and we minimize the sum of error variables with soft margin loss.

The differences between multiple birth support vector machine (MBSVM) and NHCMC are that we introduce a regularization term to implement structural risk minimization (SRM) principle and a variable is introduced to make a term of objective function to be constraints. These changes have many +ve effects on original NHCMC.

The Dual Problem

In order to get the solution of problem (17), we need to derive its dual problem. The Lagrangian of the problem (17) is given by:-

$$L(w_k, b_k, \eta_k, \varepsilon_k, \alpha, \beta, \lambda) = \frac{1}{2}(c_1 \|w_k\|^2 + \frac{1}{2}\eta_k^T \eta_k$$

$$+ c_2 \varepsilon_k^T \varepsilon_k + \lambda^T(B_k w_k + e_{b1} b_k - \eta_k) -$$

$$\alpha^T(A_k w_k + e_{b2} b_k + \varepsilon_k + -e_{b2}) - \beta^T \varepsilon_k,$$

where $\alpha = (\alpha_1, \dots, \alpha_{1k})^T$, $\beta = (\beta_1, \dots, \beta_{1k})^T$, $\lambda = (\lambda_1, \dots, \lambda_{1-1k})^T$ are the Lagrange multiplier vectors. The KTT conditions (13) for $w_k, b_k, \eta_k, \varepsilon_k$ and $\alpha, \beta, \lambda$ are given by:-

$$\nabla_{w_k} L = c_1 w_k + B_k^T \lambda - A_b^T \alpha = 0 \qquad -19$$

$$\nabla_{b_k} L = e_{b1}^T \lambda - e_{b2}^T \alpha = 0 \qquad -20$$

$$\nabla_{\eta_k} L = \eta_k - \lambda = 0 \qquad -21$$

$$\Delta_{\varepsilon_k} L = c_2 \varepsilon_k - \alpha - \beta = 0 \qquad -22$$

$$B_k w_k + e_{b1} b_k = \eta_k, \qquad -23$$

$$(A_k w_k + e_{b2} b_k) + \varepsilon_k \geq e_{b2}, \varepsilon_k \geq 0 \qquad (24)$$

$$\alpha^T((A_k w_k + e_{b2} b_k) + \varepsilon_k) = 0, \beta^T \varepsilon_k = 0 \qquad 25$$

$$\alpha \geq 0, \beta \geq 0 \qquad 26$$

Since $\beta \geq 0$, from (22) we have

$$0 \leq \alpha \leq c_2 e_{k2} \qquad -27$$

And from (19), we have

$$w_k = -\frac{1}{c_1}(B_n^T \lambda - A_k^T \alpha), \qquad -28$$

Then putting (28) and (21) into Lagrangian and using 19u 26 we obtain the dual problem of problem (17)

$$\min_{\hat{\alpha}} \tfrac{1}{2} \hat{\alpha}^T \hat{\Lambda} \hat{\alpha} + \hat{E}^T \hat{\alpha}, \quad - \quad (29)$$

$$\text{s.t. } e_{k_1}^T \lambda - e b_2^T \alpha = 0,$$

$$\hat{C}_1 \le \hat{\alpha} \le \hat{C}_2,$$

where

$$\hat{\alpha} = (\lambda^T, \alpha^T)^T \qquad \qquad 30$$

$$\hat{E} = (0, -(ceb_2^T)^T, \qquad \qquad 31$$

$$\hat{C}_1 = (-\infty eb_1^T, 0)^T, \qquad \qquad 32$$

$$\hat{C}_2 = (+\infty eb_1^T, (ceb_2^T)^T, \qquad \qquad 33$$

$$\hat{\Lambda} = \begin{pmatrix} \hat{Q}_1 & \hat{Q}_2 \\ \hat{Q}_2^T & \hat{Q}_3 \end{pmatrix} \qquad - \qquad 34$$

$$\hat{Q}_1 = B_b B_b^T + C_1 I \qquad - \qquad 35$$

$$\hat{Q}_2 = B_b A_b^T \qquad \qquad 36$$

$$\hat{Q}_3 = A_b A_b^T \qquad \qquad 37$$

where $I$ is the identity matrix.

After getting the solution of the problem (29), we can obtain $w_b^*$ and $b_b^*$ with (19) and (23). A new point $x \in R^n$ is assigned to class $k (k \in 1, \ldots, k)$, depending on which of the $k$ hyperplanes given by (15) it lies farthest to. The decision function is defined as

$$f(x) = \arg\max_{b = 1 \ldots k} \frac{|w_b^*\cdot x) + b_b^*|}{\|w_b^*\|},$$

where $|\cdot|$ is the perpendicular distance of point $x$ from the hyperplanes $(w_b \cdot x) + bb = 0, b = 1, \ldots, k$.