

We will maximize this lower bound.

For best results, we have to choose  $\{z_{iy}\}$  that make this lower bound as close to original as possible.

Choose  $\{z_{iy}\}$  such that

$$\max_{\{z_{iy}\}} \sum_i \sum_y z_{iy} \log \left( \frac{\pi_y f(x^i, \theta_y)}{z_{iy}} \right) \quad \dots \text{concave in } \{z_{iy}\}$$

$\sum z_{iy} = 1$

$z_{iy} \geq 0$  use Lagrangian to enforce this constraint

For fixed  $\{\theta_y\}$

$\exists \gamma$  - For  $i = 5, z_5 = 1,$

$$\tilde{L} = \sum_i \sum_y z_{iy} \log \left( \frac{\pi_y f(x^i, \theta_y)}{z_{iy}} \right) + \sum_p \lambda_p \left( \sum_y z_{iy} - 1 \right)$$

$$\frac{\partial \tilde{L}}{\partial z_5} = \log \pi_1 f(x^5, \theta_1) - \log(z_5) - 1 + \lambda_5 = 0$$

$$\therefore z_5 = e^{\lambda_5 - 1} \pi_1 f(x^5, \theta_1)$$

$$\text{To eliminate } \lambda_5, \quad \sum_y z_{5y} = 1$$

$$\sum_y e^{\lambda_5 - 1} \pi_y f(x^5, \theta_y) = 1$$

$$\therefore z_5 = \frac{\pi_1 f(x^5, \theta_1)}{\sum_y \pi_y f(x^5, \theta_y)}$$

$$= \frac{P(y=1) P(x^5, y=1)}{\sum ( )}$$

$$\therefore z_{iy} = P(y | x^i)$$

$\therefore$  New objective :-

$$\max_{\theta_1, \dots, \theta_k, \pi_1, \dots, \pi_k} \sum_{i=1}^N \left[ z_{iy} \log \frac{\pi_y f(x^i, \theta_y)}{z_{iy}} + \sum_j \lambda_j (z_{iy} - 1) \right]$$

$\begin{matrix} z_{iy} \geq 0 \\ \sum z_{iy} = 1 \end{matrix}$

... concave in  $\{z_{iy}\}$  for fixed  $\{\theta\}, \{\pi\}, \{z_i\}$

$\{\theta\}, \{\pi\}, \{z_i\}$

### 'Expectation - Maximization (EM) Algorithm'

- Initially chosen arbitrary values for  $\{\theta\}, \{\pi\}, \{z_i\}$

$$\text{eg} - \pi_i = \frac{1}{k} \quad \forall i.$$

- while (change) {

E-step :-

for ( $i = 1$  to  $N$ ) {

for ( $y = 1$  to  $k$ ) {

$$z_{iy}^t = P(y | x^i, \theta^t, \pi^t) \quad \dots \text{soft assignment}$$

of points to clusters

$$= \frac{\pi_y^t f(x^i, \theta_y^t)}{\sum_{y=1}^k \pi_y^t f(x^i, \theta_y^t)}$$

M-step :- (maximization)

$$\{\pi^{t+1}\}, \{\theta^{t+1}\} = \arg \max \underbrace{\sum_i \sum_y z_{iy}^t \log \pi_y f(x^i, \theta_y)}_{\sum z_{iy}^t \text{ because max wrt } \theta \in \Pi \text{ only}} = Q^t(\theta, \pi)$$

P.T.O.

- First consider the  $\pi$  variables,

$$\frac{\partial \ell^t}{\partial \pi_y} = \frac{1}{\pi_y} \left[ \sum_i \sum_j z_{ij}^t \log(\pi_y f(x_i; \theta_j)) + \lambda (\sum_i \pi_i - 1) \right]$$

$$= \sum_i \frac{z_{iy}^t}{\pi_y} + \lambda = 0$$

$$\therefore \pi_y = - \frac{\sum z_{iy}^t}{\lambda}$$

$$= \frac{\sum z_{iy}^t}{\sum_k \sum_i z_{ik}^t}$$

$$= \frac{\sum z_{iy}^t}{\sum_i \underbrace{(\sum_k z_{ik}^t)}_{=1}}$$

$$= \boxed{\frac{\sum z_{iy}^t}{N}}$$

- For the  $\theta$  variables,

For instance,  $f$  is Gaussian

$$f(x^i; \theta_y) = \frac{1}{(2\pi)^{d/2} |\Sigma_y|} e^{-\frac{1}{2} (x^i - \mu_y)^T \Sigma_y^{-1} (x^i - \mu_y)}$$

$$\rightarrow \frac{\partial \ell^t}{\partial \mu_y} = \frac{\partial}{\partial \mu_y} \left( \sum_i \sum_j z_{ij}^t \left( \log \pi_y - \frac{1}{2} (x^i - \mu_y)^T \Sigma_y^{-1} (x^i - \mu_y) - \log |\Sigma_y|^{\frac{1}{2}} \right) \right)$$

$$= - \sum_{i=1}^N z_{iy}^t (x^i - \mu_y) = 0$$

$$\mu_y^{t+1} = \frac{\sum_{i=1}^N z_{iy}^t x^i}{\sum_i z_{iy}^t}$$

## BAGGING

- It is a type of ensemble learning
- Training multiple classifiers to predict final output.

• Consider 'm' classifiers  $C_i(x)$ , each giving a prediction  $\hat{y}_i$

- Predicted label = Majority ( $\hat{y}_1, \dots, \hat{y}_n$ ) =  $\hat{y}_B$

- Let  $p \triangleq$  probability that a classifier gives correct prediction  
 Assume classifiers to be independent.

$$\text{Probability that } k \text{ classifiers are correct} = {}^m C_k p^k (1-p)^{m-k}$$

$\exists$  We predict  $\hat{y}_B$  for majority (i.e. - when more than 50% of classifiers are correct)

$$\text{Probability that } \hat{y}_B \text{ is correct} = \sum_{k=m/2}^m {}^m C_k p^k (1-p)^{m-k}$$

Intuitively, for this probability ( $\hat{y}_B$  is correct) to be  $\geq p$ ,

$$p \geq \frac{1}{2}$$

→ Two methods of creating classifier committees :-

- 1) Bootstrap sampling
- 2) Random Forests (Attribute subsetting)

## I BOOTSTRAP SAMPLING

$$D = \{(x^1, y^1), \dots, (x^N, y^N)\}$$

Convert  $D$  into a point distribution (PD)

$$P_0(x, y) = \begin{cases} \frac{1}{N} & \text{for } (x, y) \in D \\ 0 & \text{otherwise} \end{cases}$$

- Obtain a "bootstrap sample" containing  $N$  elements  $(x, y)$  taken as  $N$  samples of above uniform distribution.
  - Every bootstrap sample is like an artificial dataset that we can run our classifier on
- Use a classifier on each bootstrap sample.

→ Expected number of distinct elements  $(x, y)$  in a bootstrap sample  
 $= N \left( \text{probability that an example is selected in any trial} \right)$   
 $= N \left( 1 - \text{probability that the example is not selected at all} \right)$

$$= N \underbrace{\left( 1 - \left( 1 - \frac{1}{N} \right)^N \right)}_{X}$$

$$\lim_{N \rightarrow \infty} X = \frac{1}{e} \approx 0.37$$

∴ On average, a bootstrap sample contains  $0.37N$  distinct instances.

→ Bagging is often useful in settings where each classifier is over-fitted on their bootstrap sample.

e.g. Create a decision tree without pruning.

Bagging smooths the predictions from multiple classifiers and suppresses ill-effects of over-fitting.

## II RANDOM FORESTS

Create artificial datasets by sampling attributes instead of instances.

No. of attributes for every  $X = d$

Randomly

For each decision tree you make :-

For each node of this tree :-

→ Randomly select  $\sqrt{d}$  no. of attributes from the total  $d$  attributes.

→ Select the best attributes from these  $\sqrt{d}$  attributes based on splitting criterion.

## BOOSTING

- Opposite of bagging

\* Bagging :- Use complicated, overfitted classifiers and smooth out their predictions

• Boosting :- Start with simple base classifiers and train them sequentially to get a more complicated decision surfaces

- Consider each instance in  $D$  to have a weight.

- Initially (for first classifier), all weights are equal.

- After each classifier, increase the weights of wrongly classified instances. Then train the next classifier.

\* Weight = Coefficient of error of that instance in total loss function.

- Final classifier =  $\underset{\text{Weighted}}{\lambda}$  Average of all classifiers

•  $D \rightarrow C_1(x) \rightarrow D_{x,1} \rightarrow C_2(x) \dots \dots \rightarrow D_{m-1} \rightarrow C_m(x)$

$\alpha_i$  = weight of classifier  $C_i(x)$

- Consider no. of classes =  $k = 2$ .

## I ADABOOST ALGORITHM

Algorithm - Initially assign equal weights to all examples (instances)

$$w_i = \frac{1}{N}$$

- for  $j$  in range (maxbooststages):

$C_j(x) \leftarrow$  Train classifier on weighted dataset  $D_{j-1} = \{(w_i, x_i, y_i)\}$

e.g. For loss regularizer classifier,

$$\min_{\theta} \sum_{i=1}^N w_i \text{loss}(y_i, C_j(x_i|\theta))$$

'Weighted loss'

$\epsilon_j \leftarrow$  Weighted error of  $C_j(x)$  on  $D_{j-1} = (w_i, D)$

$$= \sum_{i=1}^N w_i \delta(y_i \neq C_j(x_i))$$

$\hat{\gamma}_j \leftarrow$  weight of classifier

if  $\epsilon_j > \frac{1}{2}$ , stop

$\alpha_j \leftarrow$  weight of classifier

$$= \log \left( \frac{1 - \epsilon_j}{\epsilon_j} \right)$$

$w_i \leftarrow$  weight updation  $\forall i$

$$= w_i \left( \frac{1 - \epsilon_j}{\epsilon_j} \right) \quad \text{if } y_i \neq C_j(x_i)$$

Rescale  $w_i$  to make their sum 1.

Theorem Let  $f(x) = \sum_{j=1}^m \alpha_j c_j(x)$  and  $y \in \{-1, 1\}$

1 Choose a loss function which is exponential in the margin

$$\text{loss}(y f(x)) = e^{-y f(x)}$$

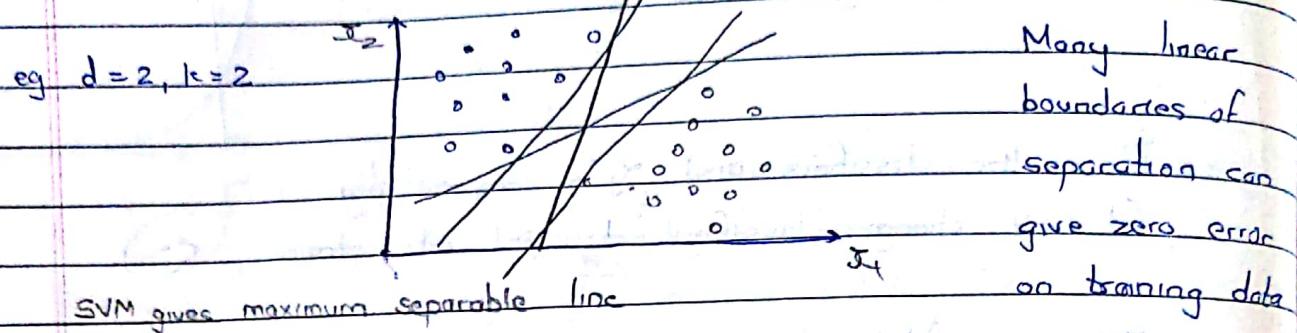
2 Train the classifiers and  $\alpha_j$ 's sequentially

(It don't change classifiers  $1 \dots j-1$  when training  $c_j$ )

Then :- AdaBoost is an optimal algorithm.

# SUPPORT VECTOR MACHINES

- Separate instances of two classes by hyperplanes in the original linear or a special high-dimensional space while maximizing the margin of separation between the classes.



- Hyperplanes in linear space :-  $w_1x_1 + w_2x_2 + \dots + w_dx_d + w_0 = 0$

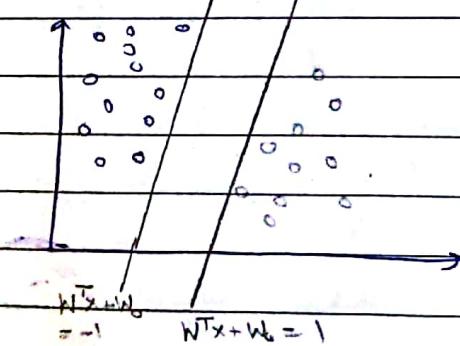
i) Separable case

$w$  s.t. for all positive points ( $y_i = 1$ ),

$$w^T x^i + w_0 \geq 1$$

and for all negative points ( $y_i = -1$ ),

$$w^T x^i + w_0 \leq -1$$



$$\text{Geometric distance between hyperplanes} = \frac{1 - (-1)}{\|w\|} = \frac{2}{\|w\|}$$

$\approx \frac{1}{\|w\|}$

Objective :-  $\underset{w, w_0}{\text{Max}} \frac{2}{\|w\|} \text{ s.t. } y_i(w^T x^i + w_0) \geq 1 \quad \forall i$ .

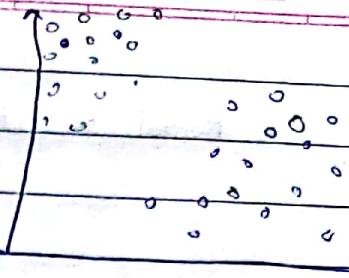
OR  $\underset{w, w_0}{\text{Min}} \frac{\|w\|}{2} \text{ s.t. } y_i(w^T x^i + w_0) \geq 1 \quad \forall i$

or  $\frac{\|w\|^2}{2}$  because  $\|w\|^2$  is convex

## 2) Non-separable case.

Allow some slack  $\epsilon_i \geq 0$  s.t.

$$y_i (w^T x_i + w_0) \geq 1 - \epsilon_i$$



Objective :-  $\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \epsilon_i$

Hyper-parameter (by trial and error)

$$\rightarrow \epsilon_i = \underbrace{\text{Hinge loss with } f(x) = w^T x_i + w_0}_{\max(0, 1 - y_i f(x))}$$

Proof:-  $\epsilon_i \geq 0$  (by definition)

$$\epsilon_i \geq 1 - y_i (w^T x_i + w_0) = 1 - y_i f(x)$$

$\therefore \epsilon_i \geq \text{Hinge Loss}$

Because we want to minimize the objective,

$$\epsilon_i = \text{Hinge Loss.}$$

- For non-separable case,

Objective =  $\min (C \cdot \text{HingeLoss} + \text{Regularizer})$

--- like loss regularizer.

T

## Non-linear SVMs.

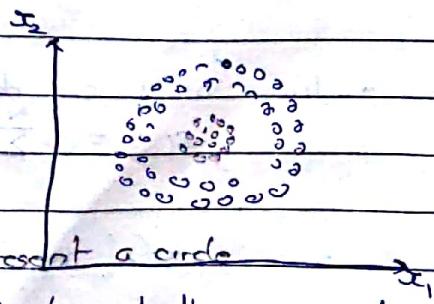
These points can be separated linearly

If they are embedded into a 5D space consisting of  $(x_1, x_2, x_1^2, x_2^2, x_1 x_2)$

$w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2$  to represent a circle

- SVMs can help do this using kernels, without explicitly representing the points in 5D space.

Work with kernels in original space itself.



- Kernel :-  $k(x, x') \equiv \underbrace{q(x) \cdot q(x')}_\text{Embedding vector of } x$

Consider  $x = (x_1, x_2)$  and  $x' = (x'_1, x'_2)$  (two points in 2D space)

$$\text{eg Kernel function} \equiv (x \cdot x' + 1)^2 = K(x, x')$$

$$= ((x_1, x_2) \cdot (x'_1, x'_2) + 1)^2$$

$$= (x_1 x'_1 + x_2 x'_2 + 1)^2$$

$$\star = x_1^2 x'^2 + x_2^2 x'^2 + 1 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x'_2 + 2x_2 x'_1$$

- Types of kernels Consider

$\Rightarrow$  Linear kernel

?

- Types of kernels

$$1) \text{Linear} \quad \Rightarrow \quad K(x, x') = x \cdot x'$$

$$2) \text{Quadratic} \quad \Rightarrow \quad (x \cdot x' + 1)^d$$

$$3) \text{Gaussian or RBF} \quad e^{-\sigma(x-x')^2} \quad \dots \text{popular}$$

- The SVM is defined in terms of the kernel function as

$$f(x) \equiv \sum_i \alpha_i k(x_i, x) y_i + w_0$$

$$\text{eg} \Rightarrow K(x_i, x) = x_i \cdot x$$

$x^i, y_i$  = training dataset point  
 $x$  = new point for testing

$$f(x) = \sum_{i=1}^N \alpha_i (x^i, x) y_i + w_0$$

$$= \underbrace{(\sum \alpha_i x^i y_i)}_{\star} + w_0$$

$\alpha_i$  = weights attached to points  $x^i$

- 31/10 • For Gaussian kernel, low  $\sigma$  :- bad decision boundary  
 high  $\sigma$  :- good decision boundary, so possibility of overfitting

- Think of kernel function  $f(x)$  to denote similarity between training point  $x$  and dataset points  $(x^i, y_i)$ .  
 $x$  is assigned class based on what points are in its vicinity.  
 → Gaussian :- High  $\sigma$  :- Every  $x^i$  has very less influence over  $x$ .

## II TRAINING SVMs

Objective

$$\underset{w, w_0, \xi_i}{\text{Minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum \xi_i \quad \text{s.t. } y_i (w^T x^i + w_0) \geq 1 - \xi_i$$

and  $\xi_i \geq 0$

- For hinge loss formulation:-

$$\underset{w, w_0}{\text{Minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum \max(0, 1 - y_i (w^T x^i + w_0))$$

→ Two methods

1 Hinge loss formulation:-

Use stochastic gradient descent algorithms

- But hinge loss is not differentiable in  $w$ . So, we will work with sub-gradients instead.

$$\star \min_{w, w_0} G(w, w_0) = \max_{w, w_0} (0, g'(w, w_0))$$

non-differentiable parts

$$\text{Sub-gradient at } (w, w_0) \equiv \begin{cases} 0 & \text{if } g'(w, w_0) \leq 0 \\ \nabla_{w, w_0} g'(w, w_0) & \text{otherwise} \end{cases}$$

2 Primal Dual Method

Starting with constrained QP formulation

$$\text{Primal (P)} \\ \min_w F(w)$$

$$\text{s.t. } g_k(w) \leq 0 \quad \forall k$$

where  $F(w), g_k(w)$  are convex

Lagrangian dual of P

$$(\alpha_1, \dots, \alpha_n) \\ \Theta = \min_w F(w) + \sum \alpha_k g_k(w)$$

$$\text{s.t. } \alpha_k \geq 0 \quad \forall k$$

- We will use Lagrangian dual of given primal question as it does not have constraints on  $g_k$ .

In our case,  $F(w)$

Theorem 'Weak Duality'

Let  $\hat{w}$  be a feasible solution of the primal  $\hat{P}$

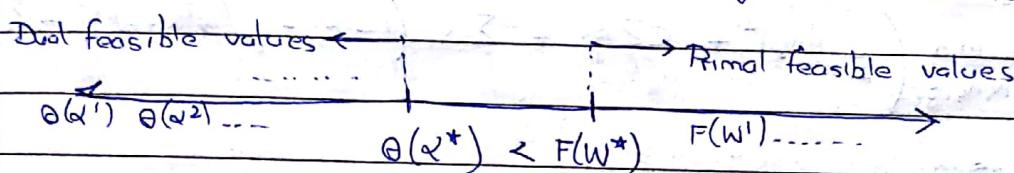
Then,  $\forall \alpha_k \geq 0$ , the dual  $\Theta(\vec{\alpha}) \leq F(\hat{w})$

$$\Theta(\vec{\alpha}) = \min_w F(w) + \sum_k \alpha_k g_k(w)$$

$$\leq F(\hat{w}) + \sum_k (\alpha_k g_k(\hat{w}))$$

Since  $\hat{w}$  is a solution of primal,  $g_k(\hat{w}) \leq 0$

$\leq F(\hat{w}) = \text{Primal objective value at } \hat{w}$



- $F(w^*) - \Theta(\vec{\alpha}^*)$  is known as 'duality gap', which needs to be minimized.

- $\Theta(\vec{\alpha})$  for all  $\alpha_j \geq 0 \leq F(\hat{w})$  for all  $w$  feasible

- New objective :- Maximize  $\Theta(\vec{\alpha})$

$$= \max_{\alpha_1, \dots, \alpha_k} \left[ \min_w \left( F(w) + \sum_k \alpha_k g_k(w) \right) \right]$$

for  $\forall \alpha_k \geq 0$  and for  $k \in \{1, \dots, k\}$

Theorem

'Strong Duality'

If  $F(w)$ ,  $g_k(w)$  are convex, then duality gap is zero.

$$\Theta(\vec{\alpha}^*) = F(w^*)$$

→ Applying Primal - Dual to SVMs :-

• Primal :-

$$\min_{w, w_0, \xi_i} \frac{1}{2} \|w\|^2 + C \sum \xi_i$$

$F(w)$  is convex  $\Rightarrow$

$$\text{such that} : \begin{cases} 1 - \xi_i - y_i (w^T x_i + w_0) \leq 0 \\ -\xi_i \leq 0 \end{cases} \quad \forall i \in \{1, 2, \dots, N\}$$

• New dual objective :-

$$\max_{\alpha_i, \mu_i} \min_{w, w_0, \xi_i} \left( \frac{1}{2} \|w\|^2 + C \sum \xi_i + \sum_i \alpha_i (1 - \xi_i - y_i (w^T x_i + w_0)) + \sum_i \mu_i (-\xi_i) \right)$$

$$\text{such that} : \begin{cases} \alpha_i \geq 0 \\ \mu_i \geq 0 \end{cases}$$

- Write as  $\max_{\alpha_i, \mu_i} \min_{w, w_0, \xi_i} L(w, w_0, \xi_i, \alpha_i, \mu_i)$  Nb. of parameters =  $d + 1 + 3N$

-  $\nabla_w L = w + \sum \alpha_i (-y_i x_i) = 0$

$$\Rightarrow -w^* = \sum_{i=1}^N \alpha_i y_i x_i$$

represents different hyperplanes for different values of  $\alpha_i^*$   
best hyperplane for optimum  $\alpha = \alpha^*$

-  $\nabla_{w_0} L = -\sum \alpha_i^* y_i = 0$

$$\therefore \sum_{y_i=1} \alpha_i^* = \sum_{y_i=-1} \alpha_i^* \quad \text{---- equal weights for both classes}$$

$$-\nabla_{\alpha_i} L = C - \alpha_i^* - \mu_i^* = 0$$

$$\alpha_i^* + \mu_i^* = C$$

Step 2 Substitute  $w = \sum \alpha_i y_i x^i$  in  $L(w, w_0, \xi_i, \alpha_i, \mu_i)$  and simplify to get

$$\begin{aligned} \min_{w, w_0, \xi_i} L(w, w_0, \xi_i, \alpha_i, \mu_i) &= \frac{-1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^N \alpha_i \\ \text{s.t. } \alpha_i \geq 0, \mu_i \geq 0 &\quad \text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0 \\ &\quad \alpha_i + \mu_i = C \\ &\quad \alpha_i \geq 0 \\ &\quad \mu_i \geq 0. \end{aligned}$$

Step 3 Eliminate  $\mu_i$

$$\max_{\alpha_1, \dots, \alpha_N} \frac{-1}{2} \left( \sum_{i=1}^N \alpha_i y_i \right) \quad \text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0 \\ \alpha_i \geq 0 \\ \alpha_i \leq C \text{ because } \alpha_i + \mu_i = C$$

•  $x^i - x^j$  in above objective is linear kernel

→ Extending the dual to arbitrary kernels using 'Kernel trick'

1 In the dual,  $x^i$  and  $x^j$  only appear as  $x^i \cdot x^j$

2 For a non-linear embedding  $Q(x)$  corresponding to some other kernel, we can use  $Q_k(x^i) \cdot Q_k(x^j) = K(x^i, x^j)$

• New kernelized dual objective :-

$$\text{Maximize}_{\alpha_1, \dots, \alpha_N} \frac{-1}{2} \left( \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(x^i, x^j) \right) + \sum \alpha_i$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad \sum \alpha_i y_i = 0$$

- Now,  $w = \sum_{i=1}^N x_i y_i x_i$

$$\begin{aligned}
 \text{Classifier} : w^T x + w_0 &= (\sum x_i y_i x_i) x + w_0 \\
 &= \sum x_i y_i K(x_i, x) + w_0
 \end{aligned}$$

→ Solving the objective.

- The objective written previously is concave in  $x_1, \dots, x_N$
- Gradient ascent, subject to constraints
- Time-consuming  $= O(N^2)$
- Cannot use single-coordinate gradient ascent because  $\sum \alpha_i$  won't be followed.

→ "SMO Algorithm"

- Two coordinate gradient ascent - Change two  $\alpha_i$ 's in each step

Initially,  $\alpha_i^0 = 0 \quad \forall i \in \{1, \dots, N\}$

while (change) {

- Pick two  $\alpha$ 's  $(\alpha_i^t, \alpha_j^t)$  to change

- Choose a change amount  $\lambda^t$ .

$$\alpha_i^{t+1} \leftarrow \alpha_i^t + \lambda^t$$

$$\alpha_j^{t+1} \leftarrow \alpha_j^t - \lambda^t$$

{ sign depends on whether  $y_i \neq y_j$

- Choosing a  $\lambda$  :- + if  $y_i \neq y_j$  and - if  $y_i = y_j$  to maintain  $\sum \alpha_i y_i = 0$

$$\arg \max_{\lambda} \Theta(\alpha_1^t, \dots, \alpha_i^t + \lambda, \alpha_j^t - \lambda, \dots, \alpha_N^t)$$

This can be solved in closed form by equating gradient of  $\Theta$  to 0.

$\lambda^t$  = value closest to  $\lambda^*$  while satisfying  $0 \leq \alpha_i^t + \lambda^t \leq C$   
 and  $0 \leq \alpha_j^t + \lambda^t \leq C$

e.g.  $t=0, \alpha_1, \dots, \alpha_N = 0$

Iteration 1 :- Choose  $(\alpha_1, \alpha_2)$  :- Say  $y_1 \neq y_2$

$$\Theta(\lambda) = \frac{-1}{2} (\lambda)(+\lambda) K(x_1, x_2) y_1 y_2 + 2\lambda$$

$$\frac{\partial \Theta}{\partial \lambda} = -\lambda K(x_1, x_2) (-1) + 2$$

$$= 0$$

$$\therefore \lambda^* = \frac{-2}{K(x_1, x_2)}$$

$$\text{Say } K(x_1, x_2) = -1 \quad \text{and } C = 1$$

$$\therefore \lambda^* = 2$$

To ensure that  $\alpha_i \& \alpha_j \in (0, C)$

$$\lambda^t = 1$$

$$\therefore \alpha_1 = 1, \alpha_2 = 1$$

### Theorem KKT conditions

Primal

$$\min_w F(w)$$

$$\text{s.t. } g_k(w) \leq 0 \quad \forall k=1 \dots k \quad \alpha_k \geq 0 \quad \forall k=1 \dots k$$

Dual

$$\max_{\alpha_1, \dots, \alpha_k} \left[ \min_w F(w) + \sum \alpha_k g_k(w) \right]$$

Primal &  
dual objectives  
are equivalent  
iff

When  $F(w)$ ,  $g_k(w)$  are convex, then  $(\alpha^*, w^*)$  are optimal solutions for the dual and primal respectively if and only if

- 1)  $\nabla_w F(w^*) + \sum \alpha_k^* \nabla g_k(w^*) = 0$
- 2)  $\alpha_k^* g_k(w) = 0 \quad \forall k=1 \dots k$
- 3)  $\alpha_k^* \geq 0$
- 4)  $g_k(w^*) \leq 0$

→ Apply KKT to SVMs

- i)  $w^* = \sum_i \alpha_i^* y_i x_i$
  - ii)  $\sum \alpha_i^* y_i = 0$
  - iii)  $\alpha_i^* + \mu_i^* = C$
  - iv)  $\mu_i^* \geq 0$
  - v)  $\alpha_i^* \geq 0$
  - vi)  $\xi_i^* \geq 0$
  - vii)  $\alpha_i^* (1 - y_i (w^* x_i + w_0^*) - \xi_i^*) = 0 \quad \forall i$
  - viii)  $\mu_i^* \xi_i^* = 0$
- {
1st KKT condition
{
Primal dual feasibility
{
2nd KKT condition

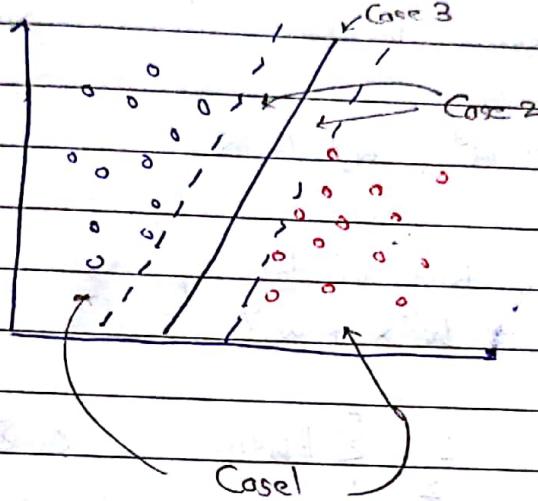
Case 1 Consider i s.t.  $y_i (w^* x_i + w_0^*) \geq 1$  and  $\xi_i^* = 0$   
 $\therefore$  From vii,  $\alpha_i^* = 0$

Case 2 Consider i s.t.  $\xi_i^* > 0$   
 $\therefore$  From viii,  $\mu_i^* = 0$   
 $\therefore$  From iii,  $\alpha_i^* = C$

Case 3 Consider i s.t  $0 < \alpha_i^* < C$

$$\Rightarrow \alpha_i^* > 0, \xi_i^* = 0$$

$$\Rightarrow y_i (w^T x^i + w_0) = 1$$



- In practice, most examples are outside the band ( $\alpha_i^* = 0$ )

$$f(x) = w^T x + w_0$$

$$\equiv \sum \alpha_i^* y_i K(x^i, x) + w_0$$

$\therefore$  We don't need to do this kernel computation for most points

→ How to compute  $w_0$ ?

Find an example  $i$  for which  $\alpha_i^*$  is  $\in (0, C)$

For this example,  $y_i (w^T x^i + w_0) = 1$

→ SVMs for multi-class classification

- k classes

- Like in logistic regression classifiers, allocate parameters  $w^y, w_0^y$  for each class  $y$ .

$$w^y \in \mathbb{R}^d, w_0^y \in \mathbb{R}$$

- Score for an example being in class  $y$  is

$$s(x, y) = w^y x + w_0^y$$

- Objective :- Minimize  $\frac{-1}{2} \sum_i^k \|w^y\|^2 + C \sum \xi_i$   
 $\{w^y, w_0^y\}, \{\xi_i\}$

such that  $w^{y_i} x^i + w_0^{y_i} \geq w^y x^i + w_0^y + 1 - \xi_i$   
 $\forall i = 1, N$

and  $\forall y \neq y_i \quad s(x^i, y) < 0$

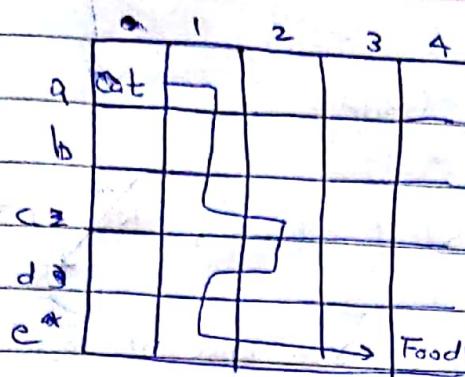
and  $\xi_i \geq 0$

- Hinge loss :-

$$\max_{y \neq y_i} (0, 1 - [s(x^i, y_i) - s(x^i, y)])$$

# REINFORCEMENT LEARNING

PAGE NO.	
DATE	/ /



Cat to reach Food

An 'agent' looks at current position and decides an action

L/R/U/D

→ Reward ( $R_t$ )

- Scalar feedback signal

- Indicates how well agent is doing at time 't'

- Agent tries to maximize cumulative reward.

→ Policy

= How an agent decides upon an action to take.

\* Longer paths could be made to have lower cumulative reward than shorter ones.

Use a 'discount factor'  $\gamma \in (0,1)$ . to account for this

Time	0	1	2	3	...	N	Cumulative
Reward	$r_0$	$r_1$	$r_2$	$r_3$		$r_N$	
Long-term reward	$r_0$	$\gamma r_1$	$\gamma^2 r_2$	$\gamma^3 r_3$	...	$\gamma^{N-1} r_N$	$\sum \gamma^i r_i$

→ Components of RL Agent

1 Policy - Agent's behaviour

2 Value function - Reward

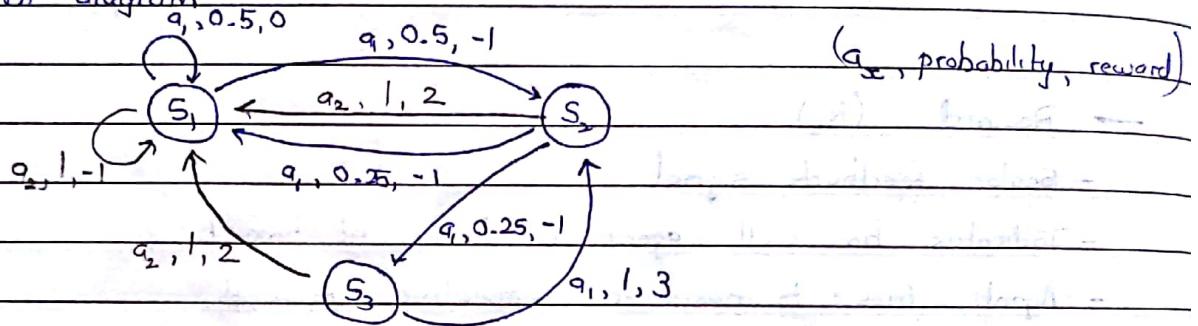
3 Model - Representation of environment

## A] Markov Decision Problem (MDP)

$S \equiv$  Set of states,  $A \equiv$  Set of actions,  $R(s, a) \equiv$  Reward function  
 $\downarrow$   
 $R(s, a, s')$

$T(s, a, s') \equiv$  Transition function,  $\gamma =$  Discount factor  
(Probability)

e.g. MDP diagram



$$S = \{S_1, S_2, S_3\} \quad A = \{a_1, a_2\}$$

$$g_1 = T(S_1, a_2, S_1) = 1$$

$$T(S_1, a_2, S_2) = 0$$

$$T(S_1, a_2, S_3) = 0$$

$$\text{e.g. } R(S_1, a_2, S_1) = -1$$

... and so on

- Trajectory:  $(s^0, a^0, r^0), (s^1, a^1, r^1) \dots (s^t, a^t, r^t)$

- Expected long-term reward:  $E[\sum \gamma^i r^i] \stackrel{\Delta}{=} V^\pi(s)$   
 $\Rightarrow$  'Value function'

$\pi =$  Policy  $\therefore$  e.g.:  $\pi = \{a_1, a_2, a_1, a_2\}$

$s_0 =$  Initial state

TFT - Most optimal policy  $\pi^*$

$$V^{\pi^*}(s) \geq V^\pi(s) \quad \forall s \in S, \forall \pi$$

Theorem  $\exists \pi^*$  for every MDP which is better than (at least just as good as) all other  $\pi$ s for all starting states  $s$ .

- $\pi^*$  may not be unique