

## CS680 – Assignment 2

Ali ElSaid (20745892)

Q1.1:

Q1.1

- Since product of kernels is a kernel then,  $K_P$  is a kernel

For any path  $P$ .

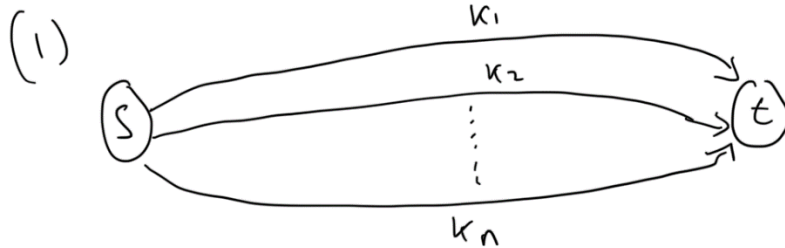
- & since the sum of two kernels is a kernel so

$\sum K_P$  is also a kernel

∴  $K_G$  is a kernel

Q1.2:

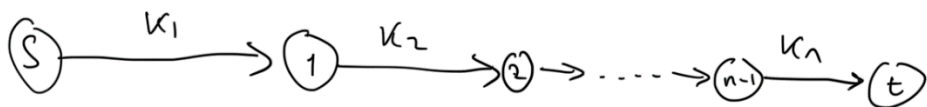
Q1.2



$$\therefore K_G = \sum_{P \in \text{Path}(s \rightarrow t)} K_P \quad \& \text{ every } P \text{ is one edge}$$

$$\text{So } \sum_{i=1}^n k_i = K_G$$

(2)



Only one path from  $s \rightarrow t \therefore K_G = \prod_{i=1}^n k_i$

Q1.3:

Q1.3

(1) Since only one path between  $s$  &  $c$

$$K_{G_1} = xz \cdot (xz-1)^3$$

(2) 2 paths between  $s$  &  $t$ , So

$$K_{G_2} = xz \cdot (xz-1)^3 \cdot e^{-(xz)^2/2} + xz \cdot e^{-|xz|} \cdot \tanh(xz+1)$$

Q1.4:

### Q1.4

Assuming there is one source ( $s$ ) & one sink ( $t$ ), & there's a path between them.

↑  
at least one

- then we can just start backwards from  $t$ , & store the result of each node & multiply it with the nodes that have edges going into it. This way we compute every nodes result once & pass by every edge once. (This is a form of dynamic programming, since we store the results of the nodes & reuse them later)
  - we can do this because the graph is a DAG & we can go backwards in the reverse topological sort ordering
- Since for a node  $u$ , all edges coming into  $u$  will multiply with its "result", i.e. the result of the back propagation.

#### • Algo:

We can use DFS and memo-ize the result of every node  $u$  after computing it & then reusing it when it is later called potentially from another path. With the base case being  $\text{result}(t) = 1$  the sink node. (This way we traverse each node & edge once, i.e.  $O(|V| + |E|)$ )

There's also the bottom up approach which starts from  $t$  & goes backwards, but the top down one with memo-ization is more intuitive.

Q1.5:

Q1.5:

if  $K$  is a kernel it is also positive semi-definite

function.  $\therefore K(s+\delta s, t+\delta t) > K(s, t+\delta t)$

$$\& K(s+\delta s, t+\delta t) > K(s+\delta s, t)$$

$\therefore$  the numerator is positive semi-definite too

• Since addition of kernels is a kernel,

the numerator is a kernel

• & since  $cK$  for a constant  $c$  is a kernel if  $K$  is a kernel

then  $\left(\frac{1}{\delta s \delta t}\right) \cdot (\text{numerator})$  is also a kernel  
 $\uparrow$   
 $c$

• Finally the limit of a kernel is also a kernel

$\therefore \lim_{\delta t \rightarrow 0, \delta s \rightarrow 0} \left(\frac{1}{\delta s \delta t}\right) \cdot (K(s+\delta s, t+\delta t) - K(s+\delta s, t) - K(s, t+\delta t) + K(s, t))$   
is a kernel as a whole

&  $\therefore$  the mixed derivative is also a kernel

Q2.1:

Q2.1

$$\bullet a = \cos(s) \quad , \quad \bullet b = \exp(a) = \exp(\cos(s))$$

$$\bullet c = \frac{1}{a^2 + 1} = \frac{1}{\cos(s)^2 + 1} \quad , \quad \bullet d = \sin(b) \cdot \ln c$$

$$= \sin(\exp(\cos(s))) \cdot \ln\left(\frac{1}{\cos(s)^2 + 1}\right)$$

$$\bullet e = d^2 = \left( \sin(\exp(\cos(s))) \cdot \ln\left(\frac{1}{\cos(s)^2 + 1}\right) \right)^2$$

## Q2.2

Q2.2:

Using Product rule & Chain rule we get

$$\begin{aligned}\frac{\partial t}{\partial s} &= 2(d) \cdot \frac{\partial d}{\partial s} = 2d \cdot \left( \sin b \cdot \frac{\partial \ln c}{\partial s} + \frac{\partial \sin b}{\partial s} \cdot \ln c \right) \\ &= 2d \left( \sin b \left( \frac{1}{c} \right) \cdot \frac{\partial c}{\partial s} \right) + \cos b \cdot \frac{\partial b}{\partial s} \cdot \ln c\end{aligned}$$

$$\cdot \frac{\partial c}{\partial s} = \frac{\partial \frac{1}{a^2+1}}{\partial a} \cdot \frac{\partial a}{\partial c} = -(a^2+1)^{-2} \cdot 2a \cdot \frac{\partial a}{\partial c}$$

$$\cdot \frac{\partial b}{\partial s} = \frac{\partial \exp(a)}{\partial a} \cdot \frac{\partial a}{\partial c} = e^a \cdot \frac{\partial a}{\partial c}$$

$$\cdot \frac{\partial a}{\partial c} = -\sin(c)$$

Plugging everything back, we get

$$\cdot \frac{\partial b}{\partial s} = \exp(\cos(s)) \cdot -\sin(s)$$

$$\cdot \frac{\partial c}{\partial s} = -(\cos(s)^2 + 1)^{-2} \cdot 2 \cos(s) \cdot -\sin(s)$$

$$\begin{aligned}\cdot \frac{\partial d}{\partial s} &= \sin(\exp(\cos(s))) \cdot (\cos(s)^2 + 1) \cdot \left( \frac{1}{\cos(s)^2 + 1} \right) \\ &\quad + \cos(\exp(\cos(s))) \cdot \ln\left(\frac{1}{\cos(s)^2 + 1}\right) \cdot \left( \frac{1}{\cos(s)^2 + 1} \right)\end{aligned}$$

$$\text{So } \rightarrow \frac{\partial t}{\partial s} = \left[ 2 \left( \sin(\exp(\cos(s))) \cdot \ln\left(\frac{1}{\cos(s)^2 + 1}\right) \right) \right]$$

$$\begin{aligned}&+ \left[ \sin(\exp(\cos(s))) \cdot (\cos(s)^2 + 1) \cdot (\cos(s)^2 + 1)^{-2} \cdot 2 \cos(s) \cdot \sin(s) \right. \\ &\quad \left. + \cos(\exp(\cos(s))) \cdot \ln\left(\frac{1}{\cos(s)^2 + 1}\right) \cdot (\exp(\cos(s)) \cdot -\sin(s)) \right]\end{aligned}$$

Q2.3:

Q 2.3:

Using the results from 2.2

$$\frac{\partial v_0}{\partial s} = 1, \quad \frac{\partial v_1}{\partial s} = -\sin(v_0) \cdot \frac{\partial v_0}{\partial s}$$

$$\frac{\partial v_2}{\partial s} = \exp(v_1) \cdot \frac{\partial v_1}{\partial s}, \quad \frac{\partial v_3}{\partial s} = -2v_1 (v_1 + 1)^{-2} \cdot \frac{\partial v_1}{\partial s}$$

$$\frac{\partial v_4}{\partial s} = \sin(v_2) \cdot \frac{1}{v_3} \cdot \frac{\partial v_3}{\partial s} + \cos(v_2) \cdot \frac{\partial v_2}{\partial s} \cdot \ln(v_3)$$

$$\frac{\partial v_5}{\partial s} = \frac{\partial t}{\partial s} = 2v_4 \frac{\partial v_4}{\partial s}, \quad \text{We can traverse once.}$$

Since we go from source to sink and compute along the way.



Q2.4:

Q2.4

$$\frac{\partial t}{\partial v_5} = 1, \quad \frac{\partial t}{\partial v_4} = \frac{\partial t}{\partial v_5} \cdot \underbrace{\frac{\partial v_5}{\partial v_4}}_{2v_4} \quad (\text{from Q1.3})$$

$$\frac{\partial t}{\partial v_3} = \frac{\partial t}{\partial v_4} \cdot \underbrace{\frac{\partial v_4}{\partial v_3}}_{\sin(v_2) \cdot \frac{1}{v_3}}, \quad \frac{\partial t}{\partial v_2} = \frac{\partial t}{\partial v_4} \cdot \underbrace{\frac{\partial v_4}{\partial v_2}}_{\cos(v_2) \cdot \ln(v_3)}$$

$$\frac{\partial t}{\partial v_1} = \frac{\partial t}{\partial v_2} \cdot \underbrace{\frac{\partial v_2}{\partial v_1}}_{-2v_1(v_1^2+1)^{-2}} + \frac{\partial t}{\partial v_3} \cdot \underbrace{\frac{\partial v_3}{\partial v_1}}_{\exp(v_1)}$$

$$\frac{\partial t}{\partial v_0} = \frac{\partial t}{\partial s} = \frac{\partial t}{\partial v_1} \cdot \underbrace{\frac{\partial v_1}{\partial v_0}}_{-\sin(v_1)}$$

$$\rightarrow \frac{\partial t}{\partial s} = \frac{\partial t}{\partial v_1} \cdot -\sin(s), \quad \text{we need to traverse two times, one to}$$

represent the final result in terms of intermediate ones & one to compute final result from them by substituting the calculations in the first iteration.


Q3.1:

Q3.1.a

Prove:  $\epsilon_t = 1 - \gamma_t \iff \epsilon_t + \gamma_t = 1$   
 $\quad \quad \quad \downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow$   
 $\quad \quad \quad dx1 \quad \quad \quad dx1 \quad \quad \quad dx1$

$$\epsilon_t + \gamma_t = M_-^T p_t + M_+^T p_t = (M_-^T + M_+^T) p_t$$

$$= |M|^T \cdot p_t = \underset{dxn}{1} \cdot \underset{n \times 1}{p_t} = \underset{dx1}{1}$$

 Since all the elements of  $M$  are either  $-1$  or  $1$

### Q3.1.b

Lets create a mapping from the Adaboost algorithm in class (ADA) to the one presented in the assignment (ASN). For the sake of notation, put  $\sim$  above a symbol to indicate that it is from the assignment algorithm.

First: Mapping between classifiers. Let each weak classifier in ADA be one of the  $d$  weak classifiers in ASN.

& we map the output as follows: (And we let  $T = \max\_Pass$ )

$$\underbrace{h_t^*(x_i)}_{\substack{\uparrow \\ \text{from ASN}}} = 2 \cdot \underbrace{h_t(x_i)}_{\substack{\uparrow \\ \text{from ADA}}} - 1, \text{ this maps:} \\ \text{the output} \quad \begin{array}{ccc} 0 & \longrightarrow & -1 \\ 1 & \longrightarrow & 1 \end{array} \\ \begin{array}{cc} \text{(ADA)} & \text{(ASN)} \end{array}$$

Second: Show that the weight assigned to each classifier in both algorithms is the same.

• ADA:

$$\text{weighted prediction per classifier: } \ln\left(\frac{1}{\beta_t}\right) \left(h_t(x_i) - \frac{1}{2}\right)$$

(transforming the prediction to be similar to ASN we get

$$\rightarrow \frac{1}{2} \ln\left(\frac{1}{\beta_t}\right) \cdot \underbrace{(2h_t(x_i) - 1)}_{\substack{\downarrow \\ h_t^*(x_i)}}$$

$\therefore$  the weight assigned to each classifiers prediction in ADA if we transform their prediction to  $\tilde{h}_t(x_i)$  is

$$\frac{1}{2} \ln\left(\frac{1}{\beta_t}\right)$$

## • ASN:

Since  $\tilde{w}_0$  is initialized to be zero, & we choose one classifier

~~let us say  $j_t$~~  each iteration. Then the final weight for a

specific classifier  $\tilde{h}_j$  is  $\rightarrow$

$$\tilde{w}_{mp,j} = \text{sum of } \tilde{\beta}_{t,j}$$

(Since only one entry of  $\alpha_t$  is 1 at a time,  $\therefore$  only one entry of  $w_t$  is updated at a time)  $\uparrow$   $mp = \text{max-pass}$  for all the times  $\tilde{h}_j$  was chosen

to map to ADA, let's consider each time  $\tilde{h}_j$  is chosen as a

"different" classifier with a different weight, so each "version"

at iteration  $t$  of  $\tilde{h}_j$  which we call  $\tilde{h}_{t,j}$  has weight

$$\begin{aligned}\tilde{\beta}_{t,j} &= \frac{1}{2} (\ln \tilde{\gamma}_{t,j} - \ln \tilde{\epsilon}_{t,j}) \\ &= \frac{1}{2} (\ln (1 - \tilde{\epsilon}_{t,j}) - \ln (\tilde{\epsilon}_{t,j})) \quad \left. \begin{array}{l} \text{since} \\ \epsilon_t + \gamma_t = 1 \end{array} \right\} \\ &= \frac{1}{2} \ln \left( \frac{\tilde{\epsilon}_t}{1 - \tilde{\epsilon}_t} \right)^{-1}\end{aligned}$$

• where  $\tilde{\epsilon}_{t,j}$  is the expected (weighted) loss for classifier  $j$  at

iteration  $t$  & likewise for  $\gamma_{t,j}$

• Note that  $\tilde{\epsilon}_{t,j} = \epsilon_t$ , i.e. the expected loss for the

classifier at time  $t$  is the same in both ADA and ASN

• if  $\tilde{\epsilon}_{t,j} = \epsilon_t$ , then  $\frac{\tilde{\epsilon}_{t,j}}{1 - \tilde{\epsilon}_{t,j}} = \frac{\epsilon_t}{1 - \epsilon_t} = \beta_t$ , making

ASN's weights  $\tilde{\beta}_{t,j}$  be also  $\rightarrow \frac{1}{2} \ln \left( \frac{1}{\beta_t} \right)$

So we have, weights of ADA = ASN  
for each classifier

• We prove this in the third part by proving that the updates to the weights of observations in both algorithms are essentially the same. ( $\tilde{p}_t$  &  $p_t$ ).

## Third: updates

ADA:  $w_{t+1} \leftarrow w_t \odot \beta_t^{\ell_t}$

if we multiply by a constant  $\beta_t^{-1/2}$  for all  $w_t$ , we change nothing

because of the normalization in the next iteration. ( $p_t = \frac{w_t}{\sum w_t}$ )

But the update becomes:  $\left\{ \begin{array}{l} \beta_t^1 \cdot \beta_t^{-1/2} = \beta_t^{1/2}, \text{ if } \ell_t = 1 \text{ (i.e. correct pred.)} \\ \beta_t^0 \cdot \beta_t^{-1/2} = \beta_t^{-1/2}, \text{ if } \ell_t = 0 \text{ (i.e. incorrect pred.)} \end{array} \right.$

where  $\beta_t = \frac{\epsilon_t}{1 - \epsilon_t}$

$$1 - \epsilon_t$$

ASN:

$$\tilde{p}_{t+1} \leftarrow \tilde{p}_t \odot \exp(-M(\tilde{\alpha}_t \odot \tilde{\beta}_t))$$

We focus only on the  $j$ th classifier (the one chosen at time  $t$ , since  $\tilde{\alpha}_t$  is zero at all entries corresponding to other classifiers & so no update happens)

$$\text{So, } \tilde{\alpha}_t \odot \tilde{\beta}_t = [\dots, \tilde{\beta}_{t,j}, \dots]$$

$$\& -M(\tilde{\alpha}_t \odot \tilde{\beta}_t) = \begin{bmatrix} 0 + 0 + \dots + y_i \tilde{h}_j(x_i) \cdot \tilde{\beta}_{t,j} + \dots + 0 + 0 \\ \vdots \\ 0 + 0 + \dots + y_n \tilde{h}_j(x_n) \tilde{\beta}_{t,j} + \dots + 0 + 0 \end{bmatrix}, \text{ i.e. } -M(\tilde{\alpha}_t \odot \tilde{\beta}_t) \text{ will only have entries corresponding to}$$

inspecting one element in the update corresponding to point  $(i)$

$$p_{t+1,i} = p_{t,i} \odot e^{-y_i \tilde{h}_j(x_i) \tilde{\beta}_{t,j}}$$

$$= e^{-\frac{1}{2} \ln\left(\frac{\tilde{\epsilon}_t}{1-\tilde{\epsilon}_t}\right) - y_i \tilde{h}_j(x_i)} = \left(\frac{\tilde{\epsilon}_t}{1-\tilde{\epsilon}_t}\right)^{\frac{1}{2} (y_i \tilde{h}_j(x_i))}$$

$$\text{which is } \begin{cases} \left(\frac{\tilde{\epsilon}_t}{1-\tilde{\epsilon}_t}\right)^{\frac{1}{2}} & , \text{ if correct prediction} \\ \left(\frac{\tilde{\epsilon}_t}{1-\tilde{\epsilon}_t}\right)^{-\frac{1}{2}} & , \text{ if incorrect prediction} \end{cases}$$

Inspecting both algorithms updates, we see that they are exactly the same due to the normalization & since both  $\tilde{\epsilon}_t$  &  $\epsilon_t$  are the respective weighted losses.

$$\therefore \tilde{\epsilon}_t = \epsilon_t \quad \& \quad \tilde{\beta}_{t,j} = \frac{1}{2} \ln \frac{1}{\beta_t}$$

Finally: That proves that both algorithms are equivalent.

Q3.2:

Q3.2

Algorithm: Since  $b_j$  is a threshold for  $x_j$ , we only need to try the  $n+1$  points between the values of  $x_j$ . At each of those points we also flip  $s_j$  &  $b_j$  to try flipping the current labeling  
(the sign of)  
& calculate the new loss as  $1 - \text{the current loss}$ , since we flipped all the predictions and  $\sum_{i=1}^n p_i = 1$ .

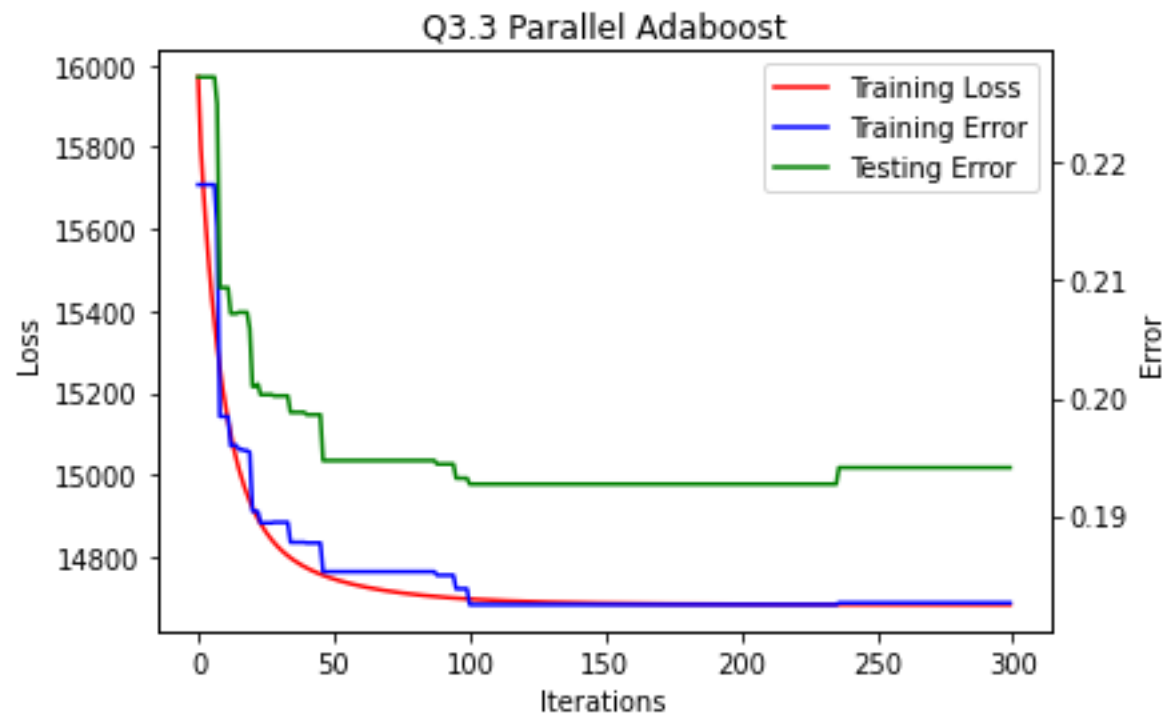
To efficiently compute the loss we do it as we go from the bottom most threshold to the top most, adding or subtracting the new  $p_i$  from the loss if the new point is misclassified or classified correctly respectively according to moving the threshold up once.

Finally, we report  $s_j$  &  $b_j$  corresponding to the minimum loss.

(we also take into consideration that the points can have duplicates & aren't sorted)

```
Best Feature: [6]
Test Error: [0.189]
Training Error: [0.1761]
sj: [1.]
bj: [-1.5]
Worst Feature: [14]
Test Error: [0.2276]
Training Error: [0.218]
sj: [-1.]
bj: [-19844.75]
```

Q3.3:



Q3.4:

