

a)  $d=2, n=i$   
 $i=2$

$$x^{(1)} = [\cos(\pi), 0] = [-1, 0]$$

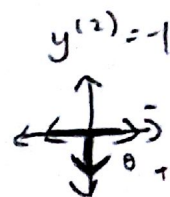
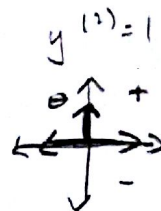
$$x^{(2)} = [0, \cos(2\pi)] = [0, 1]$$

First will always be mistake

$$x^{(1)} \Rightarrow \theta = y^{(1)} [-1, 0]$$



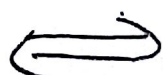
$$x^{(2)} \Rightarrow \theta = y^{(2)} [0, 1]$$



$$x^{(2)} \rightarrow \text{mistake}$$

$$\Rightarrow \theta = [y^{(1)}(-1), y^{(2)}(1)]$$

Same  $\theta$



$$x^{(1)} \Rightarrow \theta = [y^{(1)}(-1), y^{(2)}(1)]$$

We know  $\Rightarrow$  no matter what we start with, we end up with same  $\theta \Rightarrow 2$  corrections.

b) for every  $x^{(i)}$ ,  $x_i^{(i)} = \cos(i\pi) = \{-1, 1\}$   $x_j^{(i)} = 0$   
 $\Rightarrow$  only have one non-zero term per  $x^{(i)}$ .  $\Rightarrow$  first term is always mistake  $\theta = y^{(1)} x^{(i)}$ . Because, ~~only~~ only one term  $\theta_i$  switches, when we choose any  $x^{(i)}$  after, it will always lie on the linear separator  $\Rightarrow$  will always be zero.  $\Rightarrow$  Because other than  $x^{(i)}$ , the  $x_i$  term of all other  $x^{(i)}$ 's are zero.  
 $\Rightarrow$  ~~not~~  $\theta x^{(i)} = 0$  all the time  $\Rightarrow \theta = y^{(1)} x^{(1)} + y^{(2)} x^{(2)} + y^{(3)} x^{(3)} + \dots$   
 $y^{(1)} x^{(1)} \Rightarrow d$  corrections till convergence.

$$c) \theta = [y^{(1)} x_1^{(1)}, y^{(2)} x_2^{(2)}, \dots, y^{(d)} x_d^{(d)}]$$

$$= [-y^{(1)}, y^{(2)}, -y^{(3)}, \dots, (-1)^d y^{(d)}]$$

$\theta$  doesn't depend on ordering but does depend on labeling

$$d) \|x^{(i)}\| \leq R$$

$$\frac{y^{(i)}(\theta^* \cdot x^{(i)})}{\|\theta^*\|} \geq \gamma$$

$$\|x^{(i)}\| = 1$$

$$1 \leq R$$

$$\theta^* = [-y^{(1)}, y^{(2)}, -y^{(3)}, \dots, -1^d y^{(d)}] \cdot x^{(i)} = (-1^i y^i)(-1^i)$$

$$\|\theta^*\| = \sqrt{d}$$

$$= y^i$$

$$1 \leq R^2$$

$$\Rightarrow \frac{[y^{(i)}]^2}{\sqrt{d}} \geq \gamma$$

$$y^{(i)} = \{-1, +1\}$$

$$(y^{(i)})^2 = 1$$

$$\frac{1}{\sqrt{d}} \geq \gamma$$

$$\gamma^2 \leq \frac{1}{d}$$

$$d \geq \frac{1}{\gamma^2} \quad R^2 \geq 1$$

$$\frac{1}{\gamma^2} = d \quad R^2 = 1$$

$$\frac{R^2}{\gamma^2} = d \checkmark$$

$$2a) \frac{\lambda}{2} \|\theta - \theta^{(k)}\|^2 + \text{Loss}_h(y\theta \cdot x)$$

$$\theta^{(k+1)} = \theta^{(k)} + \eta yx$$

$\|\theta - \theta^{(k)}\|^2$  is how close  $\theta$  is to the current  $\theta^{(k)}$

$\Rightarrow$  Minimize  $\frac{\lambda}{2} \|\theta - \theta^{(k)}\|^2$ , if  $\frac{\lambda}{2} \uparrow$ ,  $\Rightarrow$  we want less  $\|\theta - \theta^{(k)}\|^2$ . This means that a higher learning rate leads our  $\theta^{(k)}$  to become more closer/accurate to  $\theta$ .

if  $\lambda$  is large, then learning rate is high.  $\Rightarrow \eta = \min \left\{ \frac{\text{Loss}_h(y\theta^{(k)}, x)}{\|x\|^2}, \frac{1}{\lambda} \right\}$

$\Rightarrow$  if  $\lambda \uparrow$   
 $\frac{1}{\lambda}$ ,  $\eta$  is small because it cannot be greater than  $\frac{1}{\lambda}$ .



- 2b)
- 1) d)  $\Rightarrow$  larger  $\lambda$  means  $\theta$  changes less
  - 2) b)  $\Rightarrow$  0-1 loss changes  $\theta$  more to point misclassified is now classified.
  - 3) a)
  - 4) c)

3b)

- 1) hinge loss & large  $\lambda$  means there is almost no or no change in  $\theta \Rightarrow$  d
- 2) hinge loss means very little to no change, but small  $\lambda$  means  $\theta$  will change a bit more than large  $\lambda \Rightarrow$  b
- 3) 0-1 loss causes misclassified to be classified  $\Rightarrow$  either a or c but larger  $\lambda$  is less change to  $\theta \Rightarrow$  a)
- 4) last on cases most change  $\theta \Rightarrow$  c

3a)  $\text{loss}_h = 1 - y \theta^{(k)} \cdot x$

~~$\frac{\lambda}{2} \|\theta - \theta^{(k)}\|^2 + (1 - y \theta \cdot x)$~~

$\Rightarrow \frac{\lambda}{2} \|\theta - \theta^{(k)}\|^2 + (1 - y \theta \cdot x) > 0$

$\frac{\lambda}{2} \|\theta - \theta^{(k)}\|^2 > y \theta \cdot x - 1$

$\|\theta - \theta^{(k)}\|^2 > \frac{2y \theta \cdot x - 2}{\lambda}$

$\|\theta\|^2 - 2\theta \cdot \theta^{(k)} + \|\theta^{(k)}\|^2 > \frac{2}{\lambda} (y \theta \cdot x - 1)$

~~$\|\theta\|^2 + \|\theta^{(k)}\|^2 + \frac{2}{\lambda} > (2\theta \cdot \theta^{(k)} + y \theta \cdot x)$~~

$\|\theta\|^2 - 2\|\theta\| \|\theta^{(k)}\| \cos \theta + \|\theta^{(k)}\|^2 > \frac{2}{\lambda} (y \theta \cdot x - 1)$

$\|\theta\|^k (\|\theta\|^k - 2\|\theta\|) > \frac{2}{\lambda} (y \theta \cdot x - 1)$

$$a) \quad \|\theta\|^2 - 2\|\theta\|\|\theta^k\| - \frac{2}{\lambda} (\gamma\theta \cdot x - 1) > 0$$

$$\|\theta^k\| = \frac{2\|\theta\| \pm \sqrt{4\|\theta\|^2 + 4(\gamma\|\theta\|\|x\| - 1)/\lambda}}{2}$$

$$= \|\theta\| \pm \sqrt{\|\theta\|(\|\theta\| + \gamma\|x\| - 1)/\lambda}$$

$$\|\theta^k\| = \|\theta\| \pm \sqrt{\|\theta\|^2 + (\gamma\|x\| - 1)}$$

Not sure how to continue & what not.

$$b) \quad \frac{\lambda}{2} \|\theta - \theta^{(k)}\|^2 + \text{Loss}_n(\gamma\theta \cdot x)$$

$$\theta^{(k+1)} = \theta^{(k)} + \eta \gamma x$$

$$\eta = \min \left\{ \frac{\text{Loss}}{\|\gamma x\|^2}, \frac{1}{\lambda} \right\} = \min \left\{ \text{Loss}, \frac{1}{\lambda} \right\}$$

Yes, because, even w/ updates different,  
only one  $x^{(i)}$  has their specific terms.