

# CSC311 HW3 Write-up

Jaffa Romain (student #1005243407)

04/11/2020

# Question 1

## QUESTION 1

a)  $t^{(i)} \in \{1, -1\}$

$h_i \leftarrow \arg \min_{t \in \{1, -1\}} \sum_{i=1}^n w_i \mathbb{I}\{h(x^{(i)}) \neq t^{(i)}\}$

$\alpha_t = \frac{1}{2} \log \frac{1 - \text{error}_t}{\text{err}_t}$

$w_i' \leftarrow w_i \exp(-\alpha_t t^{(i)} h_t(x^{(i)}))$

w.t.s  $\text{err}_t' = \frac{1}{2}$

NOTE:

$E^c = \{i : h_t(x^{(i)}) = t^{(i)}\} = 0$

$\Rightarrow E^c : h_t(x_i) t_i = (1)(1) \text{ or } (-1)(-1) = 1$

$E = \{i : h_t(x^{(i)}) \neq t^{(i)}\} = 1$

$E : h_t(x_i) t_i = (-1)(1) = -1$

$\text{err}_t = \frac{\sum_{i \in E} w_i}{\sum_{i=1}^n w_i}$

$$\text{err}_t' = \frac{\sum_{i=1}^n w_i' \mathbb{I}\{h_t(x^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^n w_i'} = \frac{\sum_{i \in E} w_i \exp(-\alpha_t (-1))}{\sum_{i \in E} w_i \exp(-\alpha_t (-1)) + \sum_{i \in E^c} w_i \exp(-\alpha_t (1))} = \frac{\sum_{i \in E} w_i \exp(\alpha_t)}{\sum_{i \in E} w_i \exp(\alpha_t) + \sum_{i \in E^c} w_i \exp(-\alpha_t)}$$

$$= \frac{\sum_{i \in E} w_i \exp(\alpha_t)}{\sum_{i \in E} w_i \exp(\alpha_t) + \sum_{i \in E^c} w_i \exp(-\alpha_t)} = \frac{\sum_{i \in E} w_i \exp(\alpha_t)}{\sum_{i \in E} w_i \exp(\alpha_t) + \sum_{i \in E^c} w_i \exp(-\alpha_t)} = \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i \exp(\alpha_t) + \sum_{i \in E^c} w_i \exp(-\alpha_t)}$$

$$\frac{\sum_{i \in E} w_i}{\sum_{i=1}^n w_i} = \text{err}_t = \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i + \sum_{i \in E^c} w_i} = \text{err}_t$$

$\therefore \sum_{i \in E} w_i = \text{err}_t \cdot (\sum_{i \in E} w_i + \sum_{i \in E^c} w_i)$

$\therefore \frac{\sum_{i \in E} w_i}{\text{err}_t} = \sum_{i \in E} w_i + \sum_{i \in E^c} w_i$

$\therefore \frac{\sum_{i \in E} w_i}{\text{err}_t} - \sum_{i \in E} w_i = \sum_{i \in E^c} w_i$

using this:

$$\frac{\sum_{i \in E} w_i \exp(\alpha_t)}{\sum_{i \in E} w_i \exp(\alpha_t) + \sum_{i \in E^c} w_i \exp(-\alpha_t)} = \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i + \exp(-2\alpha_t) \cdot (\sum_{i \in E} w_i - \sum_{i \in E^c} w_i)} = \frac{\sum_{i \in E} w_i}{\sum_{i \in E} w_i + \frac{\text{err}_t}{1 - \text{err}_t} (\sum_{i \in E} w_i - \frac{\text{err}_t}{1 - \text{err}_t} \sum_{i \in E^c} w_i)}$$

$$= \frac{1}{1 + \frac{\text{err}_t}{1 - \text{err}_t} \cdot \frac{1}{\text{err}_t} - \frac{\text{err}_t}{1 - \text{err}_t}} = \frac{1}{1 + \frac{1}{1 - \text{err}_t} - \frac{\text{err}_t}{1 - \text{err}_t}} = \frac{1}{1 - \text{err}_t + 1 - \text{err}_t}$$

$$= \frac{1}{\left( \frac{2(1 - \text{err}_t)}{1 - \text{err}_t} \right)} = \boxed{\frac{1}{2}}$$

$\text{err}_t'$  being  $\frac{1}{2}$  means that for the  $t+1$ -th iteration, this will be the highest possible weighted error. This is the best error rate at  $t$ .

$$b) 0-1 \text{ loss: } \mathbb{I}(h(x^{(n)}) \neq t^{(n)}) = \frac{1}{2} (1 - h(x^{(n)}) \cdot t^{(n)})$$

$$\text{w.t.s } w_i' \leftarrow w_i \exp(-\alpha_t t^{(i)} h_t(x^{(i)})) \propto w_i' \leftarrow w_i \exp(2\alpha_t \mathbb{I}\{h_t(x^{(i)}) \neq t^{(i)}\})$$

$$w_i' \leftarrow w_i \exp(-\alpha_t t^{(i)} h_t(x^{(i)}))$$

$$= w_i \exp(-\alpha_t \mathbb{I}\{h_t(x^{(i)}) = t^{(i)}\}) + w_i \exp(-\alpha_t \mathbb{I}\{h_t(x^{(i)}) \neq t^{(i)}\})$$

$$= w_i \exp(-\alpha_t \cdot 1) + w_i \exp(-\alpha_t \cdot (-1))$$

$$= w_i \exp(-\alpha_t) + w_i \exp(\alpha_t)$$

$$= w_i \cdot \exp\left(-\frac{1}{2} \log\left(\frac{1 - \text{err}_t}{\text{err}_t}\right)\right) + w_i \cdot \exp\left(\frac{1}{2} \log\left(\frac{1 - \text{err}_t}{\text{err}_t}\right)\right)$$

$$= w_i \left( \frac{\sqrt{\text{err}_t}}{\sqrt{1 - \text{err}_t}} + \frac{\sqrt{1 - \text{err}_t}}{\sqrt{\text{err}_t}} \right) = w_i \left( \frac{\sqrt{\text{err}_t} \cdot (1 - \text{err}_t)}{\sqrt{1 - \text{err}_t}} + \frac{\sqrt{1 - \text{err}_t} \cdot (\text{err}_t)}{\sqrt{\text{err}_t}} \right)$$

$$= w_i \left( (\sqrt{\text{err}_t})(\sqrt{1 - \text{err}_t}) + (\sqrt{1 - \text{err}_t})(\sqrt{\text{err}_t}) \right)$$

$$= w_i \left( \sqrt{\text{err}_t - \text{err}_t^2} + \sqrt{\text{err}_t - \text{err}_t^2} \right)$$

$$= w_i (2 \sqrt{\text{err}_t - \text{err}_t^2}) = w_i (2 \sqrt{\text{err}_t (1 - \text{err}_t)})$$

The weight update is proportional to the constant factor  $2 \sqrt{\text{err}_t (1 - \text{err}_t)}$ .

## Question 2

### QUESTION 2

- $X^{(i)} \in \{0, 1\}^{784}$   
dimensions: 784x10  
NAIVE BAYES
- $p(x, c | \theta, \pi) = p(c | \theta, \pi) p(x | c, \theta, \pi) = p(c | \pi) \prod_{j=1}^{784} p(x_j | c, \theta_{jc})$
- $p(c | \pi) = \pi_c$
- $p(x_j = 1 | c, \theta, \pi) = \theta_{jc}$
- $(x_j \in \{0, 1\})$ :  $p(x_j | c, \theta_{jc}) = \theta_{jc}^{x_j} (1 - \theta_{jc})^{(1-x_j)}$   
 $p(x_j = 1 | c, \theta_{jc}) = \theta_{jc}$ ,  $p(x_j = 0 | c, \theta_{jc}) = 1 - \theta_{jc}$
- $p(t_c = 1 | \pi) = p(c | \pi) = \pi_c$      $p(t | \pi) = \prod_{j=0}^9 \pi_j^{t_j}$      $\sum_{j=0}^9 \pi_j = 1$

**a)** MLE FOR CLASS CONDITIONAL PROBABILITIES  $\theta$  AND PRIOR  $\pi$

①  $\hat{\theta}_{MLE}$   
 $p(x, c | \theta, \pi) = p(c | \theta, \pi) p(x | c, \theta, \pi)$

$$= p(c | \pi) \prod_{j=1}^{784} p(x_j | c, \theta_{jc})$$

$$L(\theta) = \prod_{i=1}^n \pi_c \prod_{j=1}^{784} \theta_{jc}^{x_j} (1 - \theta_{jc})^{(1-x_j)}$$

log-likelihood  $\ell(\theta)$ :

$$\ell(\theta) = \sum_{i=1}^n \left[ \log \pi_c + \sum_{j=1}^{784} (x_j \log \theta_{jc} + (1-x_j) \log (1-\theta_{jc})) \right]$$

$$\hat{\theta}_{MLE} = \max_{\theta \in [0,1]} \ell(\theta)$$

$$= \min_{\theta \in [0,1]} -\ell(\theta)$$

$$\begin{aligned} \therefore \text{to maximize } \ell(\theta): \quad \frac{\partial \ell}{\partial \theta} &= \frac{\partial}{\partial \theta} \sum_{i=1}^n \left[ \log \pi_c + \sum_{j=1}^{784} (x_j \log \theta_{jc} + (1-x_j) \log (1-\theta_{jc})) \right] \\ &= 0 + \sum_{i=1}^n \sum_{j=1}^{784} \frac{x_j}{\theta_{jc}} - \sum_{i=1}^n \sum_{j=1}^{784} \frac{(1-x_j)}{(1-\theta_{jc})} \\ 0 &= \sum_{i=1}^n \mathbb{I}(c^{(i)} = c) (x_j (1-\theta_{jc}) - (1-x_j) \theta_{jc}) \\ 0 &= \sum_{i=1}^n \mathbb{I}(c^{(i)} = c) (x_j - \theta_{jc} x_j - \theta_{jc} + \theta_{jc} x_j) \\ 0 &= \sum_{i=1}^n \mathbb{I}(c^{(i)} = c) (x_j - \theta_{jc}) \\ \sum_{i=1}^n \mathbb{I}(c^{(i)} = c) x_j &= \sum_{i=1}^n \mathbb{I}(c^{(i)} = c) \theta_{jc} \end{aligned}$$

$$\hat{\theta}_{MLE} = \frac{\sum_i \mathbb{I}(c^{(i)} = c) \cdot x_j}{\sum_i \mathbb{I}(c^{(i)} = c)}$$

for  $j = 0, 1, \dots, 9$   
 $i = 1, \dots, 784$

②  $\hat{\pi}_{MLE}$

$$p(t^{(i)} | \pi) = \prod_{j=0}^9 \pi_j^{t_j}, \quad \sum_{j=0}^9 \pi_j = 1$$

FINDING LOG LIKELIHOOD:

$$\begin{aligned} L(\theta) &= \sum_{i=1}^n \log \prod_{j=0}^9 \pi_j^{t_j} \\ &= \sum_{i=1}^n \sum_{j=0}^9 t_j \log(\pi_j) \\ &= \sum_{i=1}^n \left[ \sum_{j=0}^9 t_j \log \pi_j + t_9 \log \left( 1 - \sum_{j=0}^8 \pi_j \right) \right] \end{aligned}$$

TAKING THE DERIVATIVE OF  $L(\theta)$  AND SETTING IT TO 0:

$$0 = \sum_{i=1}^n \left[ \sum_{j=0}^8 \frac{t_j^{(i)}}{\pi_j} - \sum_{j=0}^8 \frac{t_9^{(i)}}{1 - \sum_{j=0}^8 \pi_j} \right] \quad 1 - \sum_{j=0}^8 \pi_j = \pi_9$$

$$0 = \sum_{i=1}^n \left[ \sum_{j=0}^8 \frac{t_j^{(i)}}{\pi_j} - \sum_{j=0}^8 \frac{t_9^{(i)}}{\pi_9} \right]$$

$$\sum_{i=1}^n \sum_{j=0}^8 \frac{t_j^{(i)}}{\pi_j} = \sum_{i=1}^n \sum_{j=0}^8 \frac{t_9^{(i)}}{\pi_9}$$

$$\sum_{i=1}^n \sum_{j=0}^8 \pi_j = \sum_{i=1}^n \sum_{j=0}^8 \frac{t_9^{(i)}}{t_9^{(i)}}$$

$$\hat{\pi}_j = \hat{\pi}_0 + \hat{\pi}_1 + \dots + \hat{\pi}_8 = 1$$

$$\therefore \frac{\hat{\pi}_1}{\hat{\pi}_9} + \dots + \frac{\hat{\pi}_8}{\hat{\pi}_9} = \frac{1}{\hat{\pi}_9}$$

$$\therefore \hat{\pi}_9 = \frac{\sum_{i=1}^n t_9^{(i)}}{\sum_{j=0}^8 t_j^{(i)}}$$

$$\therefore \hat{\pi}_j = \frac{\sum_{i=1}^n t_j^{(i)}}{\sum_{j=0}^8 t_j^{(i)}} = \frac{\sum_{i=1}^n t_j^{(i)}}{N}$$

$$\therefore \hat{\pi}_j = \frac{1}{N} \sum_{i=1}^n t_j^{(i)} \quad \text{for } j = 0, \dots, 8$$

b) log-likelihood  $p(t|x, \theta, \pi)$  for single training example:

$$p(t|x, \theta, \pi) = \frac{p(x, c | \theta, \pi)}{p(x | \theta, \pi)}$$

$$\Rightarrow = \log \left( \frac{p(c | \pi) \prod_{j=1}^{784} p(x_j | c, \theta_{jc})}{\sum_{c=0}^9 \pi_c \prod_{j=1}^{784} \theta_{jc}^{x_j} (1 - \theta_{jc})^{1-x_j}} \right)$$

$$= \log \left( \frac{\pi_c \prod_{j=1}^{784} \theta_{jc}^{x_j} (1 - \theta_{jc})^{1-x_j}}{\sum_{c=0}^9 \pi_c \prod_{j=1}^{784} \theta_{jc}^{x_j} (1 - \theta_{jc})^{1-x_j}} \right)$$

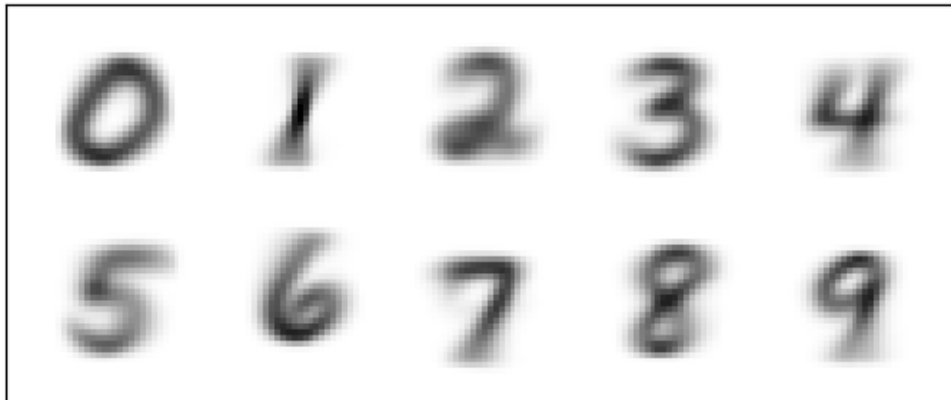
$$= \left( \log \pi_c + \sum_{j=1}^{784} (x_j \log \theta_{jc} + (1-x_j) \log (1 - \theta_{jc})) \right) - \sum_{c=0}^9 \left( \log \pi_c \sum_{j=1}^{784} (x_j \log \theta_{jc} + (1-x_j) \log (1 - \theta_{jc})) \right)$$

c)

When trying to report the log-likelihood for a single training example, there is a runtime error, where log is divided zero in the function and average log-likelihood is being returned as nan. This could be because of the number of zeros in theta. Since we are finding the log-likelihood of only one training example, there is only one entry that has 1 in the labels, and so  $\log(0)$  is being calculated in the process, which could be causing the error.

d)

MLE estimator  $\hat{\theta}$  as 10 separate greyscale images, one for each class.



e) MAP ESTIMATOR FOR  $\theta$  USING  $\text{Beta}(3,3)$  PRIOR ON EACH  $\theta_{jc}$ .

$$P(\theta | x, c, \pi) \propto p(\theta) p(x, c | \theta, \pi)$$

$$\theta \sim \text{Beta}(3, 3)$$

$$P(\theta | x, c, \pi) \propto p(\theta) p(x, c | \theta, \pi)$$

$$\begin{aligned} \ell(\theta) &= \log [\theta^{2-1} (1-\theta)^{3-1} \cdot p(c|\pi) \prod_{j=1}^{784} p(x_j | c, \theta_{jc})] \\ &= \log [\theta^2 (1-\theta)^2 \cdot \sum_{\hat{\lambda}=1}^n \left[ \pi_c \prod_{j=1}^{784} \theta_{jc}^{x_j} (1-\theta_{jc})^{1-x_j} \right]] \\ &= 2 \log \theta + 2 \log (1-\theta) + \sum_{\hat{\lambda}=1}^n \log \left( \pi_c \prod_{j=1}^{784} \theta_{jc}^{x_j} (1-\theta_{jc})^{1-x_j} \right) \\ &= 2 \log \theta + 2 \log (1-\theta) + \sum_{\hat{\lambda}=1}^n \left[ \log \pi_c + \sum_{j=1}^{784} x_j \log (\theta_{jc}) + (1-x_j) \log (1-\theta_{jc}) \right] \\ &= 2 \log \theta + 2 \log (1-\theta) + \sum_{\hat{\lambda}=1}^n \left[ \log \pi_c + \sum_{c \in \mathcal{O}} \mathbb{I}(c^{(\hat{\lambda})} = c) \left( \log (\theta_{jc}) + (1-\theta_{jc}) \log (1-\theta_{jc}) \right) \right] \end{aligned}$$

$$\frac{\partial \ell}{\partial \theta} = \frac{2}{\theta_{jc}} - \frac{2}{1-\theta_{jc}} + \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) \left( \frac{x_j}{\theta_{jc}} - \frac{1-x_j}{1-\theta_{jc}} \right)$$

$$0 = \frac{2}{\theta_{jc}} - \frac{2}{1-\theta_{jc}} + \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) \frac{x_j}{\theta_{jc}} - \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) \frac{1-x_j}{1-\theta_{jc}}$$

$$0 = 2 - 2\theta_{jc} - 2\theta_{jc} + \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) (x_j - x_j \theta_{jc} - \theta_{jc} + x_j \theta_{jc})$$

$$0 = 2 - 4\theta_{jc} + \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) (x_j) - \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) \theta_{jc}$$

$$4\theta_{jc} + \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) \theta_{jc} = 2 + \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) x_j$$

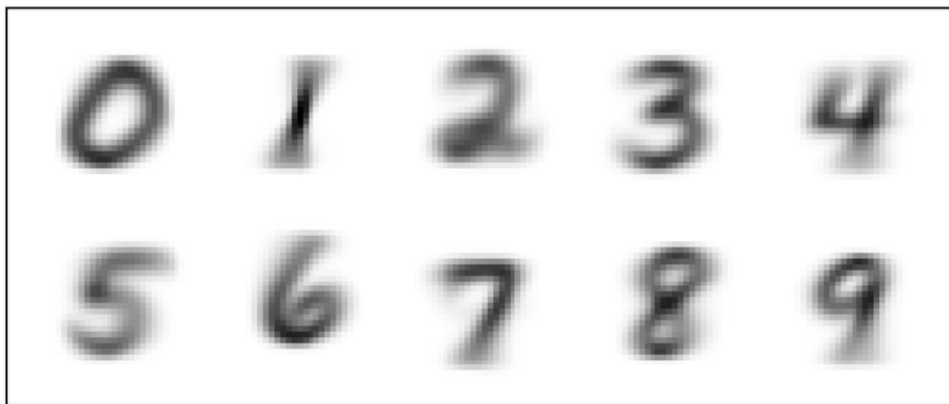
$$\hat{\theta}_{\text{MAP}} = \frac{2 + \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c) x_j}{4 + \sum_{\hat{\lambda}=1}^n \mathbb{I}(c^{(\hat{\lambda})} = c)}$$

$$j = 1, \dots, 784$$

$$\hat{\lambda} = 1, \dots, n$$

f)

Average log-likelihood for MAP = -3.3570631378602847 Training accuracy for MAP = 0.8352166666666667  
Test accuracy for MAP = 0.816



### Question 3

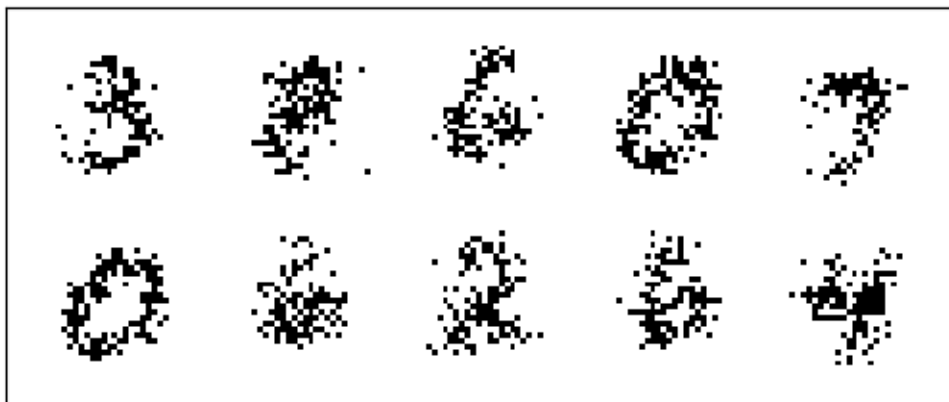
- a) True. We assume  $x_i$  and  $x_j$  are independent on the condition of  $c$ .
- b) False.

$$p(x_i, x_j) = \sum p(x_i, x_j | c) = \sum p(x_i, x_j | c) = \sum p(x_i | c) p(x_j | c).$$

$x_i$  and  $x_j$  are dependent when marginalized,  $p(x_i, x_j) \neq p(x_i)p(x_j)$ .

c)





c: [3 8 6 0 7 0 6 2 6 4]