

#1.(w)

- (a) [5pts] The goal of this question is to show that the AdaBoost algorithm changes the weights in order to force the weak learner to focus on difficult data points. Here we consider the case that the target labels are from the set $\{-1, +1\}$ and the weak learner also returns a classifier whose outputs belongs to $\{-1, +1\}$ (instead of $\{0, 1\}$). Consider the t -th iteration of AdaBoost, where the weak learner is

$$h_t \leftarrow \operatorname{argmin}_{h \in \mathcal{H}} \sum_{i=1}^N w_i \mathbb{I}\{h(\mathbf{x}^{(i)}) \neq t^{(i)}\},$$

the w -weighted classification error is

$$\text{err}_t = \frac{\sum_{i=1}^N w_i \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w_i},$$

and the classifier coefficient is $\alpha_t = \frac{1}{2} \log \frac{1-\text{err}_t}{\text{err}_t}$. (Here, \log denotes the natural logarithm.) AdaBoost changes the weights of each sample depending on whether the weak learner h_t classifies it correctly or incorrectly. The updated weights for sample i is denoted by w'_i and is

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

Show that the error w.r.t. (w'_1, \dots, w'_N) is exactly $\frac{1}{2}$. That is, show that

$$\text{err}'_t = \frac{\sum_{i=1}^N w'_i \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w'_i} = \frac{1}{2}.$$

Note that here we use the weak learner of iteration t and evaluate it according to the new weights, which will be used to learn the $t+1$ -st weak learner. What is the interpretation of this result?

Hints:

- i. Start from err_t and divide the summation to two sets of $E = \{i : h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}$ and its complement $E^c = \{i : h_t(\mathbf{x}^{(i)}) = t^{(i)}\}$.

- ii. Note that

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

$$\text{err}'_t = \frac{\sum_{i=1}^N w'_i \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w'_i}$$

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

By #3(b) the update data weights is equivalent to

$$w_i^{(t+1)} = w_i^{(t)} \cdot \exp(2\alpha_t \cdot \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}) \cdot c, \quad \text{where } c \text{ is a constant}$$

if misclassifies $\mathbf{x}^{(i)}$, i.e. $h_t(\mathbf{x}^{(i)}) \neq t^{(i)}$

then $\mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\} = 1$

$$\begin{aligned} w_i^{(t+1)} &= w_i^{(t)} \cdot \exp(2\alpha_t \cdot \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}) \cdot c \\ &= w_i^{(t)} \cdot c e^{2\alpha_t} \\ &= w_i^{(t)} \cdot c \cdot e^{\log \frac{1-\text{err}_t}{\text{err}_t}} \\ &= w_i^{(t)} \cdot c \cdot \frac{1-\text{err}_t}{\text{err}_t} \end{aligned}$$

if classifies $\mathbf{x}^{(i)}$ correctly, i.e. $h_t(\mathbf{x}^{(i)}) = t^{(i)}$

then $\mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\} = 0$

$$\begin{aligned} w_i^{(t+1)} &= w_i^{(t)} \cdot \exp(2\alpha_t \cdot \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}) \cdot c \\ &\xrightarrow{1} \\ &= w_i^{(t)} \cdot c \end{aligned}$$

$$w_i^{(t+1)} = \begin{cases} w_i^{(t)} \cdot c \cdot \frac{1-\text{err}_t}{\text{err}_t} & (\text{E}) \text{ if misclassifies } \mathbf{x}^{(i)}, \text{ i.e. } h_t(\mathbf{x}^{(i)}) \neq t^{(i)} \\ w_i^{(t)} \cdot c & (\text{E}^c) \text{ if classifies } \mathbf{x}^{(i)} \text{ correctly, i.e. } h_t(\mathbf{x}^{(i)}) = t^{(i)} \end{cases}$$

(E) if misclassifies $\mathbf{x}^{(i)}$, i.e. $h_t(\mathbf{x}^{(i)}) \neq t^{(i)}$

(E)^c if classifies $\mathbf{x}^{(i)}$ correctly, i.e. $h_t(\mathbf{x}^{(i)}) = t^{(i)}$

$$\text{err}_t' = \text{err}_{t+1} = \frac{\sum_{i \in E} w_i^{(t+1)}}{\sum_{i=1}^N w_i^{(t+1)}} = \frac{\sum_{i \in E} c w_i^{(t)} \cdot \frac{1 - \text{err}_t}{\text{err}_t}}{\sum_{i \in E} c w_i^{(t)} \cdot \frac{1 - \text{err}_t}{\text{err}_t} + \sum_{i \notin E} c w_i^{(t)}}$$

⋮

$$\text{err}_2 = \frac{\sum_{i \in E} w_i^{(2)}}{\sum_{i=1}^N w_i^{(2)}} = \frac{\sum_{i \in E} c w_i^{(1)} \cdot \frac{1 - \text{err}_1}{\text{err}_1}}{\sum_{i \in E} c w_i^{(1)} \cdot \frac{1 - \text{err}_1}{\text{err}_1} + \sum_{i \notin E} c w_i^{(1)}}$$

Since we initialize sample weights with $w_i^{(0)} = \frac{1}{N}$ for $i=1, \dots, N$

statistically, we have 50% probability to classify the data correctly, i.e. $\text{err}_1 = \frac{1}{2}$, then $\frac{1 - \text{err}_1}{\text{err}_1} = \frac{1/2}{1/2} = 1$

and we can then get $\text{err}_2 = \frac{1}{2}$. Along with this thread of thinking, $\text{err}_t = \frac{1}{2}$ and so $\text{err}_{t+1} = \frac{1}{2}$

Hence, we showed that the error w.r.t. (w_1, \dots, w_N) is exactly $\frac{1}{2}$

the $(t+1)^{\text{th}}$ w-weighted classification error is determined by the t^{th} w-weighted classification error

⋮

the 2nd w-weighted classification error is determined by the 1st w-weighted classification error

#1.(b)

(b) [2pts] Recall from lecture that we can rewrite the 0-1 loss as: $\mathbb{I}[h(x^{(n)}) \neq t^{(n)}] = \frac{1}{2}(1 - h(x^{(n)}) \cdot t^{(n)})$. Use this identity to prove that the weight update from part (a):

$$w'_i \leftarrow w_i \exp\left(-\alpha_t t^{(i)} h_t(\mathbf{x}^{(i)})\right).$$

is proportional, up to a constant factor, to weight update:

$$w'_i \leftarrow w_i \exp\left(2\alpha_t \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}\right).$$

What is the constant factor? Since we normalize the weights when computing err_t , these two updates are equivalent.

We want to show

$$\exp(-\alpha_t t^{(i)} h_t(\mathbf{x}^{(i)})) \propto \exp(2\alpha_t \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\})$$

By given identity, we can get

$$2 \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\} = 2 \cdot \frac{1}{2}(1 - h_t(\mathbf{x}^{(i)}) \cdot t^{(i)}) = 1 - h_t(\mathbf{x}^{(i)}) \cdot t^{(i)}$$

Then, we can get

$$\begin{aligned} e^{[2\alpha_t \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}]} &= (e^{[1 - h_t(\mathbf{x}^{(i)}) \cdot t^{(i)}]})^{\alpha_t} \\ &= (e \cdot e^{[-h_t(\mathbf{x}^{(i)}) \cdot t^{(i)}]})^{\alpha_t} \\ &= \underbrace{e^{\alpha_t}}_{\text{constant factor}} \cdot e^{[-\alpha_t h_t(\mathbf{x}^{(i)}) \cdot t^{(i)}]} \end{aligned}$$

Hence, we proved $\exp(-\alpha_t t^{(i)} h_t(\mathbf{x}^{(i)}))$ is proportional to $\exp(2\alpha_t \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\})$

The constant factor is e^{α_t}

[13pts] Fitting a Naïve Bayes Model

In this question, we'll fit a Naïve Bayes model to the MNIST digits using maximum likelihood. The starter code will download the dataset and parse it for you: Each training sample $(\mathbf{x}^{(i)}, \mathbf{t}^{(i)})$ is composed of a vectorized binary image $\mathbf{x}^{(i)} \in \{0, 1\}^{784}$, and 1-of-10 encoded class label $\mathbf{t}^{(i)}$, i.e., $t_j^{(i)} = 1$ means image i belongs to class c .

Given parameters π and θ , Naïve Bayes defines the joint probability of the each data point \mathbf{x} and its class label c as follows:

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

where $p(c | \pi) = \pi_c$ and $p(x_j = 1 | c, \theta_{jc}) = \theta_{jc}$. Here, θ is a matrix of probabilities for each pixel and each class, so its dimensions are 784×10 , and π is a vector with one entry for each class. (Note that in the lecture, we simplified notation and didn't write the probabilities conditioned on the parameters, i.e. $p(c | \pi)$ is written as $p(c)$ in lecture slides).

For binary data ($x_j \in \{0, 1\}$), we can write the Bernoulli likelihood as

$$p(x_j = 1 | c, \theta_{jc}) = \theta_{jc}^x (1 - \theta_{jc})^{(1-x_j)}, \quad (0.1)$$

which is just a way of expressing $p(x_j = 1 | c, \theta_{jc}) = \theta_{jc}$ and $p(x_j = 0 | c, \theta_{jc}) = 1 - \theta_{jc}$ in a compact form. For the prior $p(t | \pi)$, we use a categorical distribution (generalization of Bernoulli distribution to multi-class case).

$$p(t_c = 1 | \pi) = p(c | \pi) = \pi_c \text{ or equivalently } p(t | \pi) = \prod_{j=0}^9 \pi_i \quad \text{where } \sum_{i=0}^9 \pi_i = 1,$$

where $p(c | \pi)$ and $p(t | \pi)$ can be used interchangeably. You will fit the parameters θ and π using MLE and MAP techniques. In both cases, your fitting procedure can be written as a few simple matrix multiplication operations.

- (a) [3pts] First, derive the *maximum likelihood estimator* (MLE) for the class-conditional pixel probabilities θ and the prior π . Derivations should be rigorous.

Hint 1: We saw in lecture that MLE can be thought of as 'ratio of counts' for the data, so what should $\hat{\theta}_{jc}$ be counting?

Hint 2: Similar to the binary case, when calculating the MLE for π_j for $j = 0, 1, \dots, 8$,

write $p(t^{(i)} | \pi) = \prod_{j=0}^9 \pi_j^{t_j^{(i)}}$ and in the log-likelihood replace $\pi_0 = 1 - \sum_{j=0}^8 \pi_j$, and then take derivatives w.r.t. π_j . This will give you the ratio $\hat{\pi}_j / \hat{\pi}_0$ for $j = 0, 1, \dots, 8$. You know that $\hat{\pi}_j$'s sum up to 1.

let D = num of features (here, D=784)

let N = num of data (N=60000 for training data)

The likelihood func is

$$L(\pi, \theta) = \sum_{i=1}^N p(c, \bar{x})$$

The log-likelihood is

$$\begin{aligned} \ell(\pi, \theta) &= \log \left(\sum_{i=1}^N p(\bar{t}^{(i)}, \bar{x}^{(i)}) \right) \\ &= \sum_{i=1}^N \log p(\bar{t}^{(i)}, \bar{x}^{(i)}) \\ &= \sum_{i=1}^N \log \{ p(\bar{x}^{(i)} | \bar{t}^{(i)}) p(\bar{t}^{(i)}) \} \\ &= \sum_{i=1}^N \log \{ p(\bar{t}^{(i)}) \prod_{j=1}^D p(x_j^{(i)} | \bar{t}^{(i)}) \} \\ &= \sum_{i=1}^N \left[\log p(\bar{t}^{(i)}) + \sum_{j=1}^D \log p(x_j^{(i)} | \bar{t}^{(i)}) \right] \\ &= \underbrace{\sum_{i=1}^N \log p(\bar{t}^{(i)})}_{\text{Bernoulli log-likelihood}} + \underbrace{\sum_{i=1}^N \sum_{j=1}^D \log p(x_j^{(i)} | \bar{t}^{(i)})}_{\substack{\text{Bernoulli log-likelihood} \\ \text{of labels}}} \end{aligned}$$

For the prior, $p(\bar{t}^{(i)})$, we maximize $\sum_{i=1}^N \log p(\bar{t}^{(i)})$

$$\text{since } \pi_j = \frac{\sum_{i=1}^N I(t_j^{(i)} = 1)}{N} = \frac{\text{num of data of class } j}{\text{total num of data}}$$

then $p(t_j^{(i)} = 1) = \pi_j$

$$\begin{aligned} p(\bar{t}^{(i)}) &= \prod_{j=0}^9 p(t_j^{(i)} = 1) \\ &= \prod_{j=0}^9 \pi_j^{t_j^{(i)}} \\ &= \prod_{j=0}^9 \pi_j^{t_j^{(i)}} \cdot (\pi_1)^{t_1^{(i)}} \end{aligned}$$

$$\text{train-image.shape} = (60000, 784)$$

$$\text{train-label.shape} = (60000, 10)$$

$$\text{test-image.shape} = (1000, 784)$$

$$\text{test-label.shape} = (1000, 10)$$

$$\begin{aligned} p(t_c) &= p(c) = \pi_c \\ t_c &= \left[\begin{array}{c} 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right] \} 60000 \\ t &= \left[\begin{array}{c} t_0 \\ t_1 \\ \vdots \\ t_9 \end{array} \right] \end{aligned}$$

x data points

c = class label

θ (matrix, 784x10) prob for each pixel and each class

T (vector, 10x1) with one entry for each class

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

$$p(c | \pi) = \pi_c \text{ and } p(x_j = 1 | c, \theta_{jc}) = \theta_{jc} \quad \text{Pg. 11}$$

$$p(c) = \pi_c$$

$$p(x_j = 1 | c) = \theta_{jc}$$

For binary data ($x_j \in \{0, 1\}$), we can write the Bernoulli likelihood as

$$p(x_j = 1 | c, \theta_{jc}) = \theta_{jc}^x (1 - \theta_{jc})^{(1-x_j)}, \quad (0)$$

which is just a way of expressing $p(x_j = 1 | c, \theta_{jc}) = \theta_{jc}$ and $p(x_j = 0 | c, \theta_{jc}) = 1 - \theta_{jc}$

$$p(x_j = 1 | c) = \theta_{jc} \quad (1 - \theta_{jc})^{(1-x_j)}$$

$$p(x_j = 1 | c) = \theta_{jc}, \quad p(x_j = 0 | c) = 1 - \theta_{jc}$$

$t_j^{(i)}$: image i belongs to class j

$$p(t_j^{(i)} = 1 | \pi) = p(c | \pi) = \pi_c \text{ or equivalently } p(t | \pi) = \prod_{j=0}^9 \pi_j^{t_j^{(i)}} \text{ where } \sum_{i=0}^9 \pi_i = 1,$$

$$p(t_j^{(i)} = 1 | \pi) = \prod_{j=0}^9 \pi_j^{t_j^{(i)}} \cdot (\pi_1)^{t_1^{(i)}}$$

$$= \prod_{j=0}^8 \pi_j^{t_j^{(0)}} \cdot (1 - \sum_{j=0}^8 \pi_j)^{t_9^{(0)}}$$

$$\begin{aligned}
\sum_{i=1}^N \log P(t^{(0)}) &= \sum_{i=1}^N \log \left[\prod_{j=0}^8 \pi_j^{t_j^{(0)}} \cdot (\pi_9^{t_9^{(0)}}) \right] \\
&= \sum_{i=1}^N \left[\log \left(\prod_{j=0}^8 \pi_j^{t_j^{(0)}} \right) + \log (\pi_9^{t_9^{(0)}}) \right] \\
&= \sum_{i=1}^N \left[\sum_{j=0}^8 \log \pi_j^{t_j^{(0)}} + \log (\pi_9^{t_9^{(0)}}) \right] \\
&= \sum_{i=1}^N \sum_{j=0}^8 t_j^{(0)} \log \pi_j + \sum_{i=1}^N \log (\pi_9^{t_9^{(0)}}) \\
&= \sum_{i=1}^N \sum_{j=0}^8 t_j^{(0)} \log \pi_j + \sum_{i=1}^N t_9^{(0)} \log (\pi_9) \quad \text{since } \pi_9 = 1 - \sum_{j=0}^8 \pi_j \\
&= \sum_{i=1}^N \sum_{j=0}^8 t_j^{(0)} \log \pi_j + \sum_{i=1}^N t_9^{(0)} \log (1 - \sum_{j=0}^8 \pi_j) \quad t_9^{(0)} = 1 - \sum_{j=0}^8 t_j^{(0)} \\
&= \sum_{i=1}^N \sum_{j=0}^8 t_j^{(0)} \log \pi_j + \sum_{i=1}^N \left(1 - \sum_{j=0}^8 t_j^{(0)} \right) \log (1 - \sum_{j=0}^8 \pi_j)
\end{aligned}$$

$$\begin{aligned}
\frac{\partial}{\partial \pi_0} \left(\sum_{i=1}^N \log P(t^{(0)}) \right) &= \frac{\partial}{\partial \pi_0} \left(\sum_{j=0}^8 \log \pi_j \underbrace{\sum_{i=1}^N t_j^{(0)}} + \log (1 - \sum_{j=0}^8 \pi_j) \underbrace{\sum_{i=1}^N t_9^{(0)}} \right) \\
&= \frac{\partial}{\partial \pi_0} \sum_{j=0}^8 \log \pi_j \underbrace{\sum_{i=1}^N t_j^{(0)}} + \frac{\partial}{\partial \pi_0} \log (1 - \sum_{j=0}^8 \pi_j) \underbrace{\sum_{i=1}^N t_9^{(0)}} \\
&= \frac{\partial}{\partial \pi_0} \left(\log \pi_0 \underbrace{\sum_{i=1}^N t_0^{(0)}} + \sum_{j=1}^8 \log \pi_j \underbrace{\sum_{i=1}^N t_j^{(0)}} \right) + \underbrace{\sum_{i=1}^N t_9^{(0)}} \frac{\partial}{\partial \pi_0} \log (1 - \sum_{j=0}^8 \pi_j) \\
&= \frac{1}{\pi_0} \underbrace{\sum_{i=1}^N t_0^{(0)}} + \frac{1}{1 - \sum_{j=0}^8 \pi_j} \cdot (-1) \underbrace{\sum_{i=1}^N t_9^{(0)}} \\
&= \frac{1}{\pi_0} \underbrace{\sum_{i=1}^N t_0^{(0)}} - \frac{1}{1 - \sum_{j=0}^8 \pi_j} \underbrace{\sum_{i=1}^N t_9^{(0)}} \\
&= \frac{1}{\pi_0} \underbrace{\sum_{i=1}^N t_0^{(0)}} - \frac{1}{\pi_9} \underbrace{\sum_{i=1}^N t_9^{(0)}}
\end{aligned}$$

Similarly, the derivative for π_0 also applies to $\pi_1, \pi_2, \dots, \pi_8$

so, for $0 \leq j \leq 8$,

$$\frac{\partial}{\partial \pi_j} \left(\sum_{i=1}^N \log P(t^{(0)}) \right) = \frac{1}{\pi_j} \underbrace{\sum_{i=1}^N t_j^{(0)}} - \frac{1}{1 - \sum_{j=0}^8 \pi_j}$$

$$\text{In summary, } \frac{\partial}{\partial \pi_j} \left(\sum_{i=1}^N \log P(t^{(0)}) \right) = \begin{cases} \frac{1}{\pi_j} \underbrace{\sum_{i=1}^N t_j^{(0)}} - \frac{1}{\pi_9} \underbrace{\sum_{i=1}^N t_9^{(0)}} & \text{for } 0 \leq j \leq 8 \\ \frac{1}{\pi_j} \underbrace{\sum_{i=1}^N t_9^{(0)}} & \text{for } j = 9 \end{cases}$$

$$\text{let } \frac{d}{d\pi_j} \left(\sum_{i=1}^N \log P(t^{(0)}) \right) = 0$$

$$\text{that is, for } 0 \leq j \leq 8, \text{ let } \frac{1}{\pi_j} \underbrace{\sum_{i=1}^N t_j^{(0)}} = \frac{1}{\pi_9} \underbrace{\sum_{i=1}^N t_9^{(0)}} \Rightarrow \frac{\pi_j}{\pi_9} = \frac{\sum_{i=1}^N t_j^{(0)}}{\sum_{i=1}^N t_9^{(0)}}$$

$$\text{let } f(\pi_0) = 1 - \sum_{j=0}^8 \pi_j$$

$$\begin{aligned}
\frac{\partial \log(f(\pi_0))}{\partial \pi_0} &= \frac{1}{f(\pi_0)} \cdot \frac{\partial}{\partial \pi_0} (1 - \sum_{j=0}^8 \pi_j) \\
&= \frac{1}{1 - \sum_{j=0}^8 \pi_j} \cdot \frac{\partial}{\partial \pi_0} (1 - \pi_0 - \pi_1 - \dots - \pi_8) \\
&= \frac{1}{1 - \sum_{j=0}^8 \pi_j} \cdot (-1)
\end{aligned}$$

$$\left\{ \begin{array}{l} \frac{\hat{\pi}_j}{\hat{\pi}_q} = \frac{\frac{1}{N} \sum_{i=1}^N t_j^{(i)}}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} \Rightarrow \hat{\pi}_j = \hat{\pi}_q \frac{\frac{1}{N} \sum_{i=1}^N t_j^{(i)}}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} \quad \text{for } 0 \leq j \leq q \\ \sum_{j=0}^q \hat{\pi}_j = 1 \Rightarrow \hat{\pi}_q = 1 - \sum_{j=0}^q \hat{\pi}_j \end{array} \right.$$

$$\begin{aligned} \sum_{j=0}^q \hat{\pi}_j &= 1 \Rightarrow \sum_{j=0}^q \hat{\pi}_j + \hat{\pi}_q = 1 \\ \Rightarrow \hat{\pi}_q \frac{\frac{1}{N} \sum_{i=1}^N t_0^{(i)}}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} + \hat{\pi}_q \frac{\frac{1}{N} \sum_{i=1}^N t_1^{(i)}}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} + \dots + \hat{\pi}_q \frac{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} + \hat{\pi}_q &= 1 \\ \Rightarrow \hat{\pi}_q \left(\frac{\frac{1}{N} \sum_{i=1}^N t_0^{(i)}}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} + \frac{\frac{1}{N} \sum_{i=1}^N t_1^{(i)}}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} + \dots + \frac{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} \right) &= 1 \\ \Rightarrow \hat{\pi}_q &= \frac{1}{\frac{\sum_{j=0}^q \left(\frac{1}{N} \sum_{i=1}^N t_j^{(i)} \right)}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}}} = \frac{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}}{\sum_{j=0}^q \left(\frac{1}{N} \sum_{i=1}^N t_j^{(i)} \right)} \end{aligned}$$

$$\begin{aligned} \pi &= \begin{bmatrix} \pi_0 \\ \pi_1 \\ \vdots \\ \pi_g \\ \pi_q \end{bmatrix} \\ \Rightarrow \frac{\pi}{\pi_q} &= \begin{bmatrix} \hat{\pi}_0 \\ \hat{\pi}_1 \\ \vdots \\ \hat{\pi}_g \\ \hat{\pi}_q \end{bmatrix} = \frac{1}{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}} \begin{bmatrix} \hat{\pi}_0 \frac{\frac{1}{N} \sum_{i=1}^N t_0^{(i)}}{\hat{\pi}_q} / \hat{\pi}_q \\ \hat{\pi}_1 \frac{\frac{1}{N} \sum_{i=1}^N t_1^{(i)}}{\hat{\pi}_q} / \hat{\pi}_q \\ \vdots \\ \hat{\pi}_g \frac{\frac{1}{N} \sum_{i=1}^N t_g^{(i)}}{\hat{\pi}_q} / \hat{\pi}_q \\ \hat{\pi}_q \frac{\frac{1}{N} \sum_{i=1}^N t_q^{(i)}}{\hat{\pi}_q} / \hat{\pi}_q \end{bmatrix} = \frac{1}{\sum_{j=0}^q \left(\frac{1}{N} \sum_{i=1}^N t_j^{(i)} \right)} \begin{bmatrix} \frac{1}{N} \sum_{i=1}^N t_0^{(i)} \\ \frac{1}{N} \sum_{i=1}^N t_1^{(i)} \\ \vdots \\ \frac{1}{N} \sum_{i=1}^N t_g^{(i)} \\ \frac{1}{N} \sum_{i=1}^N t_q^{(i)} \end{bmatrix} = \frac{1}{\sum_{j=0}^q \left(\frac{1}{N} \sum_{i=1}^N t_j^{(i)} \right)} \begin{bmatrix} \mathbb{I}(t_0^{(i)} = 1) \\ \mathbb{I}(t_1^{(i)} = 1) \\ \vdots \\ \mathbb{I}(t_g^{(i)} = 1) \\ \mathbb{I}(t_q^{(i)} = 1) \end{bmatrix} \\ &= \boxed{\frac{1}{N} \begin{bmatrix} \mathbb{I}(t_0^{(i)} = 1) \\ \mathbb{I}(t_1^{(i)} = 1) \\ \vdots \\ \mathbb{I}(t_g^{(i)} = 1) \\ \mathbb{I}(t_q^{(i)} = 1) \end{bmatrix}} \end{aligned}$$

To maximize $\sum_{i=1}^N \log P(x_j^{(i)} | t^{(i)})$

$$\theta_{jc} = P(x_j^{(i)} = 1 | c)$$

$$P(x_j^{(i)} | c) = \theta_{jc}^{x_j^{(i)}} (1 - \theta_{jc})^{1-x_j^{(i)}}$$

θ_{jc} : num of class-c pixel in a pic
num of class-c pic in dataset

$$\begin{aligned} & \sum_{i=1}^N \log P(x_j^{(i)} | c^{(i)}) \\ &= \sum_{i=1}^N \sum_{k=0}^q c^{(i)} \left\{ \log \theta_{jk}^{x_j^{(i)}} + \log (1 - \theta_{jk})^{1-x_j^{(i)}} \right\} \\ &= \sum_{i=1}^N \sum_{k=0}^q c^{(i)} \left\{ x_j^{(i)} \log \theta_{jk} + (1 - x_j^{(i)}) \log (1 - \theta_{jk}) \right\} \end{aligned}$$

$$\begin{aligned} & \frac{\partial}{\partial \theta_{jc}} \left(\sum_{i=1}^N \log P(x_j^{(i)} | c^{(i)}) \right) \\ &= \frac{\partial}{\partial \theta_{jk}} \left(\sum_{i=1}^N \sum_{k=0}^q c^{(i)} \left\{ x_j^{(i)} \log \theta_{jk} + (1 - x_j^{(i)}) \log (1 - \theta_{jk}) \right\} \right) \\ &= \sum_{i=1}^N \sum_{k=0}^q c^{(i)} \left\{ x_j^{(i)} \frac{\partial}{\partial \theta_{jk}} \log \theta_{jk} + (1 - x_j^{(i)}) \frac{\partial}{\partial \theta_{jk}} \log (1 - \theta_{jk}) \right\} \\ \Rightarrow & \hat{\theta}_{jc} = \frac{\sum_{i=1}^N \mathbb{I}[x_j^{(i)} = 1 \text{ & } c^{(i)} = c]}{\sum_{i=1}^N \mathbb{I}[c^{(i)} = c]} \quad \text{where } c \in [0, q] \end{aligned}$$

(b) [1pt] Derive the log-likelihood $\log p(\mathbf{t}|\mathbf{x}, \theta, \pi)$ for a single training image.

$$p(\tilde{\mathbf{t}}|\tilde{\mathbf{x}}) = \frac{p(\tilde{\mathbf{t}}, \tilde{\mathbf{x}})}{p(\tilde{\mathbf{x}})} = \frac{p(\tilde{\mathbf{t}}) p(\tilde{\mathbf{x}}|\tilde{\mathbf{t}})}{\sum_{i=1}^q p(\tilde{\mathbf{t}}^{(i)}) p(\tilde{\mathbf{x}}|\tilde{\mathbf{t}}^{(i)})} = \frac{p(\tilde{\mathbf{t}}) \prod_{j=1}^D p(x_j|\tilde{\mathbf{t}})}{\sum_{i=1}^q p(\tilde{\mathbf{t}}^{(i)}) \prod_{j=1}^D p(x_j|\tilde{\mathbf{t}}^{(i)})}$$

$D=784$

$$p(c|\tilde{\mathbf{x}}) = \frac{p(c, \tilde{\mathbf{x}})}{p(\tilde{\mathbf{x}})} = \frac{p(c) p(\tilde{\mathbf{x}}|c)}{\sum_{c' \in \mathcal{C}} p(c') p(\tilde{\mathbf{x}}|c')} = \frac{p(c) \prod_{j=1}^D p(x_j|c)}{\sum_{c' \in \mathcal{C}} p(c') \prod_{j=1}^D p(x_j|c')} = \frac{\pi_c \prod_{j=1}^D \theta_{jc}^{x_j} (1-\theta_{jc})^{(1-x_j)}}{\sum_{c' \in \mathcal{C}} \pi_{c'} \prod_{j=1}^D p(x_j|c')}$$

$$\begin{aligned} \log p(c|\tilde{\mathbf{x}}) &= \log \frac{\pi_c \prod_{j=1}^D \theta_{jc}^{x_j} (1-\theta_{jc})^{(1-x_j)}}{\sum_{c' \in \mathcal{C}} \pi_{c'} \prod_{j=1}^D p(x_j|c')} \\ &= \log \pi_c + \log \prod_{j=1}^D \theta_{jc}^{x_j} (1-\theta_{jc})^{(1-x_j)} - \log \sum_{c' \in \mathcal{C}} \pi_{c'} \prod_{j=1}^D p(x_j|c') \\ &= \log \pi_c + \sum_{j=1}^D [\log \theta_{jc}^{x_j} + \log (1-\theta_{jc})^{(1-x_j)}] - \sum_{c' \in \mathcal{C}} [\log \pi_{c'} + \log \sum_{j=1}^D p(x_j|c')] \\ &= \log \pi_c + \sum_{j=1}^D [x_j \log \theta_{jc} + (1-x_j) \log (1-\theta_{jc})] - \sum_{c' \in \mathcal{C}} \sum_{j=1}^D \log p(x_j|c') \\ &= \log \pi_c + \sum_{j=1}^{784} [x_j \log \hat{\theta}_{jc} + (1-x_j) \log (1-\hat{\theta}_{jc})] - \sum_{c' \in \mathcal{C}} \sum_{j=1}^{784} \log p(x_j|c') \end{aligned}$$

$$\begin{aligned} \sum_{c' \in \mathcal{C}} \log \pi_{c'} &= \log \sum_{c' \in \mathcal{C}} \pi_{c'} \\ &= \log 1 \\ &= 0 \end{aligned}$$

since $\hat{\pi}_c = \frac{\sum_{i=1}^N \mathbb{I}(t_c^{(i)}=1)}{N}$

$$\hat{\theta}_{jc} = \frac{\sum_{i=1}^N \mathbb{I}[x_j^{(i)}=1 \& c^{(i)}=c]}{\sum_{i=1}^N \mathbb{I}[c^{(i)}=c]}$$

where $c \in [0, 9]$

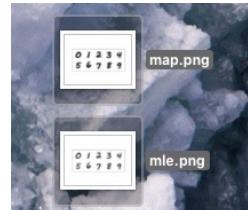
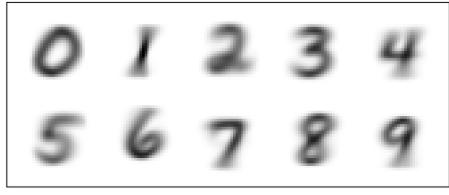
$$\begin{aligned} \log p(c|\tilde{\mathbf{x}}) &= \log \hat{\pi}_c + \sum_{j=1}^{784} [x_j \log \hat{\theta}_{jc} + (1-x_j) \log (1-\hat{\theta}_{jc})] - \sum_{c' \in \mathcal{C}} \sum_{j=1}^{784} \log p(x_j|c') \\ &= \log \frac{\sum_{i=1}^N \mathbb{I}(t_c^{(i)}=1)}{N} + \sum_{j=1}^{784} \left[x_j \log \frac{\sum_{i=1}^N \mathbb{I}[x_j^{(i)}=1 \& c^{(i)}=c]}{\sum_{i=1}^N \mathbb{I}[c^{(i)}=c]} + (1-x_j) \log \left(1 - \frac{\sum_{i=1}^N \mathbb{I}[x_j^{(i)}=1 \& c^{(i)}=c]}{\sum_{i=1}^N \mathbb{I}[c^{(i)}=c]}\right) \right] - \sum_{c' \in \mathcal{C}} \sum_{j=1}^{784} \log p(x_j|c') \end{aligned}$$

(c) [3pt] Fit the parameters θ and π using the training set with MLE, and try to report the average log-likelihood per data point $\frac{1}{N} \sum_{i=1}^N \log p(\mathbf{t}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\pi})$, using Equation (0.1). What goes wrong? (it's okay if you can't compute the average log-likelihood here).

```
In[2]: runfile('/Users/ChangyanXu/Desktop/naive_bayes.py', wdir='/Users/ChangyanXu/Desktop')
Backend MacOSX is interactive backend. Turning interactive mode on.
/Users/ChangyanXu/Desktop/naive_bayes.py:178: RuntimeWarning: divide by zero encountered in log
    x_dot_log_theta = image.dot(np.log(theta))
```

$\log(0)$ will result in an negative infinity.

(d) [1pt] Plot the MLE estimator $\hat{\theta}$ as 10 separate greyscale images, one for each class.



- (e) [2pt] Derive the *Maximum A posteriori Probability* (MAP) estimator for the class-conditional pixel probabilities θ , using a Beta(3, 3) prior on each θ_{jc} . Hint: it has a simple final form, and you can ignore the Beta normalizing constant.

$$\beta(a=3, b=3)$$

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \arg \max_{\theta} p(\theta | \mathcal{D}) \\ &= \arg \max_{\theta} p(\theta, \mathcal{D}) \\ &= \arg \max_{\theta} p(\theta) p(\mathcal{D} | \theta) \\ &= \arg \max_{\theta} \log p(\theta) + \log p(\mathcal{D} | \theta)\end{aligned}$$

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \arg \max_{\theta} p(\theta | \mathcal{D}) \\ &= \arg \max_{\theta} p(\theta, \mathcal{D}) \\ &= \arg \max_{\theta} p(\theta) p(\mathcal{D} | \theta) \\ &= \arg \max_{\theta} \log [p(\theta) p(\mathcal{D} | \theta)]\end{aligned}$$

$$\begin{aligned}\hat{\theta}_{jc, \text{MAP}} &= \arg \max_{\theta_{jc}} \prod_{i=1}^N \prod_{j=1}^{784} \theta_{jc}^{x_j^{(i)}} (1-\theta_{jc})^{1-x_j^{(i)}} \theta_{jc}^{a-1} (1-\theta_{jc})^{b-1} \\ &= \arg \max_{\theta_{jc}} \left(\frac{\sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} x_j^{(i)} + (a-1)}{(1-\theta_{jc})} \right) \quad \left(\frac{\sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} (1-x_j^{(i)}) + (b-1)}{(1-\theta_{jc})} \right)\end{aligned}$$

$$\text{let } A = \sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} x_j^{(i)} + (a-1), \quad B = \sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} (1-x_j^{(i)}) + (b-1)$$

$$\begin{aligned}\frac{d}{d\theta_{jc}} \log [\theta_{jc}^A (1-\theta_{jc})^B] &= \frac{d}{d\theta_{jc}} \left[A \log \theta_{jc} + B \log (1-\theta_{jc}) \right] \\ &= \frac{A}{\theta_{jc}} - \frac{B}{1-\theta_{jc}} = 0\end{aligned}$$

$$\hat{\theta}_{jc} = \frac{A}{A+B} = \frac{\sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} x_j^{(i)} + (a-1)}{\sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} + (a+b-2)}$$

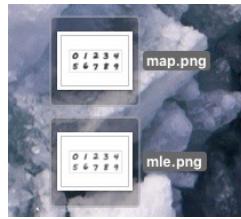
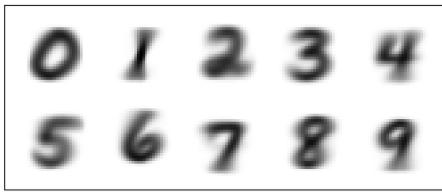
Since $\beta(a=3, b=3)$

$$\begin{aligned}\hat{\theta}_{jc, \text{MAP}} &= \frac{\sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} x_j^{(i)} + (a-1)}{\sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} + (a+b-2)} \\ &= \frac{\sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} x_j^{(i)} + 2}{\sum_{i=1}^N \mathbb{I}\{c^{(i)}=c\} + 4}\end{aligned}$$

(f) [2pt] Fit the parameters θ and π using the training set with MAP estimators from previous part, and report both the average log-likelihood per data point, $\frac{1}{N} \sum_{i=1}^N \log p(\mathbf{t}^{(i)} | \mathbf{x}^{(i)}, \hat{\theta}, \hat{\pi})$, and the accuracy on both the training and test set. The accuracy is defined as the fraction of examples where the true class is correctly predicted using $\hat{c} = \operatorname{argmax}_c \log p(t_c = 1 | \mathbf{x}, \hat{\theta}, \hat{\pi})$.

```
Average log-likelihood for MLE is nan
Average log-likelihood for MAP is  0.03364712014983577
Training accuracy for MAP is  0.0
Test accuracy for MAP is  0.0
```

(g) [1pt] Plot the MAP estimator $\hat{\theta}$ as 10 separate greyscale images, one for each class.



3. [4pts] Generating from a Naïve Bayes Model

Defining a joint probability distribution over the data lets us generate new data, and also lets us answer all sorts of queries about the data. This is why these models are called *Generative Models*. We will use the Naïve Bayes model trained in previous question to generate data.

if x_i and x_j are independent,

- (a) [1pt] True or false: Given this model's assumptions, any two pixels x_i and x_j where $i \neq j$ are independent given c .

$$p(x_i) p(x_j) = p(x_i, x_j)$$

False. Should be "conditional independent," rather than "independent"

- (b) [1pt] True or false: Given this model's assumptions, any two pixels x_i and x_j where $i \neq j$ are independent after marginalizing over c .

True

$$p(x_i, c) = p(c|x_i) \cdot p(x_i)$$

$$p(x_j, c) = p(c|x_j) \cdot p(x_j)$$

$$p(x_i|c) = \frac{p(x_i, c)}{p(c)} \quad p(x_j|c) = \frac{p(x_j, c)}{p(c)}$$

$$p(x_i|c) p(x_j|c) = \frac{p(x_i, c) p(x_j, c)}{p(c) p(c)} = \frac{p_{cc}(x_i) \cdot p(x_i) p_{cc}(x_j) \cdot p(x_j)}{p_{cc} p_{cc}}$$

- (c) [2pts] Using the parameters fit using MAP in Question 1, produce random image samples from the model. That is, randomly sample and plot 10 binary images from the marginal distribution $p(\mathbf{x}|\hat{\theta}, \hat{\pi})$. Hint: To sample from $p(\mathbf{x}|\hat{\theta}, \hat{\pi})$, first sample random variable c from $p(c|\hat{\pi})$ using `np.random.choice`, then depending on the value of c , sample x_j from $p(x_j|c, \hat{\theta}_{jc})$ for $j = 1, \dots, 784$ using `np.random.binomial(1, ...)`. These functions can take matrix probabilities as input, so your solution to this part should be a few lines of code.