

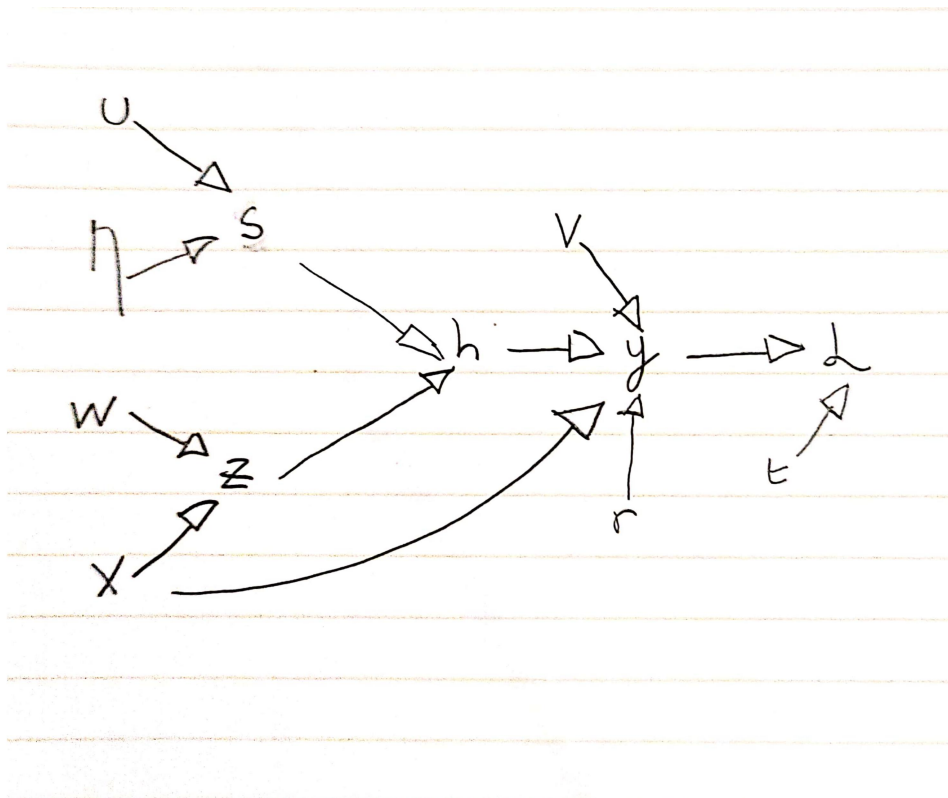
HW3 Writeup

Shiven Taneja 1005871013

November 2021

1 Backprop

1.a



1.b

Backprop Formulas:

$$\overline{\mathcal{L}} = 1$$

$$\begin{aligned}\overline{y} &= \overline{\mathcal{L}} \frac{\partial \mathcal{L}}{\partial y} \\ &= \overline{\mathcal{L}}(y - t)\end{aligned}$$

$$\begin{aligned}\overline{\mathbf{h}} &= \overline{y} \frac{\partial y}{\partial \mathbf{h}} \\ &= \overline{y} \mathbf{v}^\top\end{aligned}$$

$$\begin{aligned}\overline{\mathbf{s}} &= \overline{\mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{s}} \\ &= \overline{\mathbf{h}} \mathbf{z} \odot \sigma'(\mathbf{s})\end{aligned}$$

$$\begin{aligned}\overline{\mathbf{U}} &= \overline{\mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{U}} \\ &= \overline{\mathbf{s}} \boldsymbol{\eta}^\top\end{aligned}$$

$$\begin{aligned}\overline{\boldsymbol{\eta}} &= \overline{\mathbf{s}} \frac{\partial \mathbf{s}}{\partial \boldsymbol{\eta}} \\ &= \overline{\mathbf{s}} \mathbf{U}^\top\end{aligned}$$

$$\begin{aligned}\overline{\mathbf{z}} &= \overline{\mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \\ &= \overline{\mathbf{h}} \sigma(\mathbf{s})\end{aligned}$$

$$\begin{aligned}\overline{\mathbf{W}} &= \overline{\mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}} \\ &= \overline{\mathbf{z}} \mathbf{x}^\top\end{aligned}$$

$$\begin{aligned}\overline{\mathbf{x}} &= \overline{\mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}} + \overline{y} \frac{\partial y}{\partial \mathbf{x}} \\ &= \overline{\mathbf{z}} \mathbf{W}^\top + \overline{y} \mathbf{r}^\top\end{aligned}$$

2 Fitting a Naive Bayes Model

2.a

For $\hat{\theta}_{jc}$:

We start by finding the log likelihood:

$$\begin{aligned}\ell(\theta) &= \log(p(x^i, c^i | \theta, \pi)) = \sum_{i=1}^N \log(p(c^i | \pi) * \prod_{j=1}^{784} p(x_j^i | c^i, \theta_{jc})) \\ &= \sum_{i=1}^N (\log(\pi_{c^i}) + \log(\prod_{j=1}^{784} (\theta_{jc^i}^{x_j^i} (1 - \theta_{jc^i})^{1-x_j^i}))) \\ &= \sum_{i=1}^N (\log(\pi_{c^i}) + \sum_{j=1}^{784} \log(\theta_{jc^i}^{x_j^i} (1 - \theta_{jc^i})^{1-x_j^i})) \\ &= \sum_{i=1}^N (\log(\pi_{c^i}) + \sum_{j=1}^{784} (x_j^i \log(\theta_{jc^i}) + (1 - x_j^i) \log(1 - \theta_{jc^i})))\end{aligned}$$

Taking the derivative of θ_{jc} and setting to 0:

$$\begin{aligned}\frac{\partial \ell}{\partial \theta_{jc}} &= \sum_{i=1}^N \left(\frac{x_j^i}{\theta_{jc^i}} - \frac{(1-x_j^i)}{(1-\theta_{jc^i})} \right) = 0 \\ &= \sum_{i=1}^N (x_j^i - x_j^i \theta_{jc^i} - \theta_{jc^i} + x_j^i \theta_{jc^i}) = 0 \\ &= \sum_{i=1}^N (x_j^i - \theta_{jc^i}) = 0 \\ &= \sum_{i=1}^N (x_j^i - \theta_{jc^i}) * \mathbb{I}[t_c^i = c] = 0, \text{ where } c \text{ is a number from } 0-9\end{aligned}$$

$$\theta_{jc^i} = \frac{\sum_{i=1}^N x_j^i \mathbb{I}[t_c^i = c]}{\sum_{i=1}^N \mathbb{I}[t_c^i = c]}$$

This is the number of data points that are the digit 'c' and have the j^{th} pixel over the number of data points that are the digit 'c'

For $\hat{\pi}$:

We start by finding the log likelihood:

$$\begin{aligned}\ell(\pi) &= \log(p(t^i | \pi_j)) = \sum_{i=1}^N \log(\prod_{j=0}^9 \pi_j^{t_j^i}) \\ &= \sum_{i=1}^N \sum_{j=0}^9 t_j^i \log(\pi_j) \\ &= \sum_{i=1}^N (\sum_{j=0}^8 t_j^i \log(\pi_j) + t_9^i \log(1 - \sum_{j=0}^8 t_j^i))\end{aligned}$$

Taking the partial derivative of π and setting it to 0:

$$\begin{aligned}\frac{\partial \ell}{\partial \pi} &= \sum_{i=1}^N \left(\frac{t_j^i}{\pi_j} - \frac{t_9^i}{1 - \sum_{j=0}^8 t_j^i} \right) = \sum_{i=1}^N \left(\frac{t_j^i}{\pi_j} - \frac{t_9^i}{\pi_9} \right) = 0 \\ &= \sum_{i=1}^N \left(\frac{t_j^i}{\pi_j} \right) = \sum_{i=1}^N \left(\frac{t_9^i}{\pi_9} \right) \\ &= \sum_{i=1}^N \left(\frac{t_j^i}{t_9^i} \right) = \frac{\pi_j}{\pi_9}\end{aligned}$$

We know that $\hat{\pi}_j$ sums to 1 for $j = 1, \dots, 9$ so for $j = 1, \dots, 8$, $\hat{\pi}_j = 1 - \hat{\pi}_9$:

$$\sum_{i=1}^N \left(\frac{t_j^i}{t_9^i} \right) = \frac{1 - \pi_9}{\pi_9} = \frac{1}{\pi_9} - 1$$

$$\sum_{i=1}^N \left(\frac{t_j^i}{t_9^i} \right) + 1 = \frac{1}{\pi_9}$$

$$\sum_{i=1}^N \left(\frac{t_j^i + t_9^i}{t_9^i} \right) = \frac{1}{\pi_9}$$

$$\sum_{i=1}^N \left(\frac{t_9^i}{t_j^i} \right) = \pi_9$$

We know that $\hat{\pi}_j = 1 - \hat{\pi}_9$, so:

$$\hat{\pi}_j = 1 - \sum_{i=1}^N \left(\frac{t_9^i}{t_j^i} \right)$$

$$\hat{\pi}_j = \sum_{i=1}^N \left(\frac{t_j^i - t_9^i}{t_j^i} \right)$$

$$\hat{\pi}_j = \sum_{i=1}^N \left(\frac{\sum_{j=0}^8 t_j^i}{\sum_{j=0}^9 t_j^i} \right)$$

Since we chose $j = 8$ randomly, we could get the same result for any value of j and thus we get:

$$\hat{\pi}_j = \frac{1}{N} * \sum_{i=1}^N \sum_{j=0}^9 t_j^i$$

This is each class in the dataset divided by N

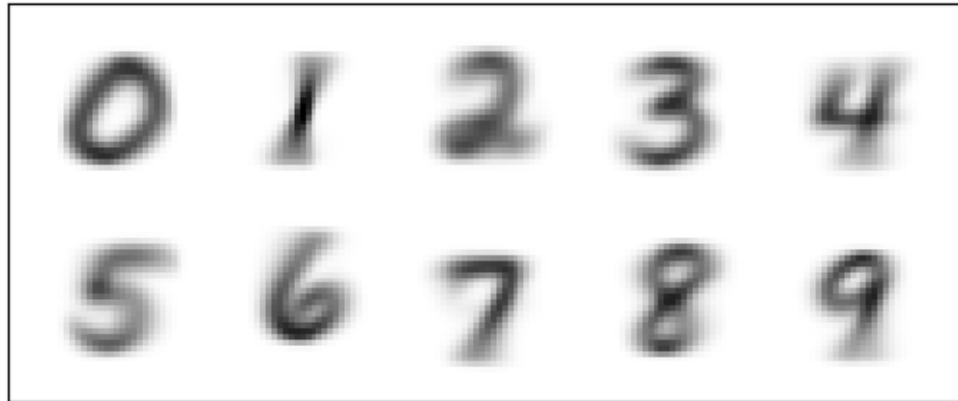
2.b

$$\begin{aligned}
 \log(P(c|x, \theta, \pi)) &= \log\left(\frac{P(c|x, \theta) P(x|c, \theta, \pi)}{\sum_{i=1}^q P(c_i|\pi) P(x|c_i, \theta, \pi)}\right) \quad \left[\begin{array}{l} \text{Where } c \in \{0, 1\} \\ \text{is the one-hot encoding} \\ \text{for } c \in \{0, \dots, q\} \end{array} \right] \\
 &= \log\left(\frac{P(c|\pi) \prod_{j=1}^{784} P(x_j|c, \theta, \pi)}{\sum_{i=1}^q P(c_i|\pi) \prod_{j=1}^{784} P(x_j|c_i, \theta, \pi)}\right) = \log\left(\frac{\pi_c \prod_{j=1}^{784} \theta_{jc}^{x_j} (1-\theta_{jc})^{(1-x_j)}}{\sum_{i=1}^q \pi_{c_i} \prod_{j=1}^{784} \theta_{jc_i}^{x_j} (1-\theta_{jc_i})^{(1-x_j)}}\right) \\
 &= \log(\pi_c) \cdot \sum_{j=1}^{784} (x_j \log(\theta_{jc}) + (1-x_j) \log(1-\theta_{jc})) - \sum_{i=1}^q \left(\log(\pi_{c_i}) + \right. \\
 &\quad \left. \sum_{j=1}^{784} (x_j \log(\theta_{jc_i}) + (1-x_j) \log(1-\theta_{jc_i})) \right) \\
 &= \log(\pi_c) \cdot \sum_{j=1}^{784} (x_j \log(\theta_{jc}) + (1-x_j) \log(1-\theta_{jc})) - \log \sum_{i=1}^q \exp\left(\log(\pi_{c_i}) + \right. \\
 &\quad \left. \sum_{j=1}^{784} (x_j \log(\theta_{jc_i}) + (1-x_j) \log(1-\theta_{jc_i})) \right)
 \end{aligned}$$

2.c

The average log-likelihood for MLE is nan. The reason for this may be due to the zero values in θ making the average close to 0.

2.d



2.e

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

$$\begin{aligned}
 \ell(\theta) &= \log(p(\theta | x, c, \pi)) = \log(p(\theta)) + \log(p(x, c | \theta, \pi)) \\
 &= \log(\theta^{3-1} \cdot (1-\theta)^{3-1}) + \sum_{c=0}^N \log(p(c | \pi) \prod_{j=1}^{784} p(x_j^i | c, \theta_{jc})) \\
 &= \log(\theta^2) + 2 \cdot \log(1-\theta) + \sum_{c=0}^N \left(\log(\pi_c) + \sum_{j=1}^{784} \log(\theta_{jc}^{x_j^i} (1-\theta_{jc})^{(1-x_j^i)}) \right) \\
 &= 2 \log(\theta) + 2 \log(1-\theta) + \sum_{c=0}^N \left(\log(\pi_c) + \sum_{j=1}^{784} (x_j^i \log(\theta_{jc}) + (1-x_j^i) \log(1-\theta_{jc})) \right)
 \end{aligned}$$

We can now take the derivative and set it to 0:

$$\frac{\partial \ell}{\partial \theta_{jc}} = \frac{2}{\theta_{jc}} - \frac{2}{1-\theta_{jc}} + \sum_{i=1}^N \left(\frac{x_j^i}{\theta_{jc}} - \frac{(1-x_j^i)}{(1-\theta_{jc})} \right)$$

$$0 = 2 - 2\theta_{jc} - 2\theta_{jc} + \sum_{i=1}^N (x_j^i - \theta_{jc})$$

$$0 = 2 - 2\theta_{jc} - 2\theta_{jc} + \sum_{i=1}^N \mathbb{1}[t_c^{(i)} = c] \cdot (x_j^i - \theta_{jc})$$

$$4\theta_{jc} + \sum_{i=1}^N \mathbb{1}[t_c^{(i)} = c] \theta_{jc} = 2 + \sum_{i=1}^N \mathbb{1}[t_c^{(i)} = c] x_j^i$$

$$\theta_{jc}^A = \frac{2 + \sum_{i=1}^N \mathbb{1}[t_c^{(i)} = c] x_j^i}{4 + \sum_{i=1}^N \mathbb{1}[t_c^{(i)} = c]}$$

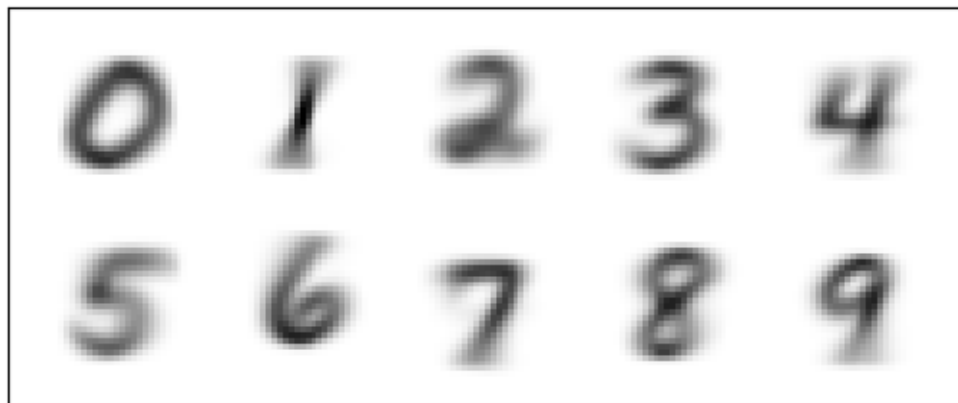
2.f

Average log-likelihood for MAP is -5.509798896626949

Training accuracy for MAP is 0.9383166666666667

Test accuracy for MAP is 0.9368

2.g



3 Categorical Distribution

3.a

Using Bayes Rule we know:

$$P(\theta | D) = \frac{P(\theta) P(D | \theta)}{\int P(\theta') P(D | \theta') d\theta'}$$

Thus:

$$P(\theta | D) \propto P(\theta) P(D | \theta)$$

Where:

$$P(\theta) \propto \theta_1^{a_1-1} \dots \theta_K^{a_K-1}$$

And since the data is independent and individually distributed:

$$P(x | \theta) = P(D | \theta) = \prod_{k=1}^K \theta_k^{x_k}$$

Therefore

$$P(\theta | D) \propto \prod_{k=1}^K \theta_k^{x_k}$$

The Dirichlet distribution is not a conjugate prior for the categorical distribution as when we combine the categorical and the Dirichlet prior, it does not have the same structure as a Dirichlet distribution

3.b

MAP Estimation

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | D)$$

$$= \arg \max_{\theta} p(\theta) p(D | \theta)$$

$$= \arg \max_{\theta} \log(p(\theta)) + \log(p(D | \theta))$$

$$\log(p(\theta | D)) = \log(p(\theta)) + \log(p(D | \theta))$$

$$= (a_K - 1) \log(\theta_K) + \sum_{k=1}^K (x_k^i \log(\theta_K))$$

$$= (a_K - 1) \log(\theta_K) + \sum_{k=1}^{K-1} (x_k^i \log(\theta_K)) + x_K^i \log(1 - \sum_{k=1}^{K-1} \theta_K)$$

Taking the partial derivative and setting to 0:

$$\frac{\partial \log}{\partial \theta_K} = \frac{a_K - 1}{\theta_K} + \frac{x_K^i}{\theta_K} - \frac{x_K^i}{\theta_K} = 0$$

$$= \frac{a_K - 1 + x_K^i}{\theta_K} - \frac{x_K^i}{\theta_K} = 0$$

$$= \theta_K (a_K - 1) + \theta_K x_K^i - \theta_K x_K^i = 0$$

$$= \theta_K = \frac{a_K - 1 + x_K^i}{x_K^i}$$

$$= \frac{1 - \theta_K}{\theta_K} = \frac{a_K - 1 + x_K^i}{x_K^i}$$

$$\frac{1}{\theta_K} - 1 = \frac{a_K - 1 + x_K^i}{x_K^i}$$

$$\frac{1}{\theta_K} = \frac{a_K - 1 + x_K^i}{x_K^i} + 1 = \frac{a_K - 1 + x_K^i + x_K^i}{x_K^i} = \frac{a_K - 1 + \sum_{k=1}^K x_K^i}{x_K^i}$$

$$\hat{\theta}_K = \frac{x_K^i}{a_{K-1} + x_K^i}$$

Since $K = k$ was chosen arbitrarily, we could do this for any value of k and thus $\hat{\theta}_K$ would then be:

$$\hat{\theta}_K = \frac{a_K - 1 + x_K^i}{N}$$

3.c

4 Gaussian Discriminant Analysis

4.a

Average conditional log likelihood of training set: -0.12462443666863293

Average conditional log likelihood of test set: -0.19667320325525828

4.b

Test accuracy for the training set: 0.9814285714285714

Test accuracy for the test set: 0.97275

4.c

