

# CSC311 Assignment3

Benson Li

April 25, 2023

## 1 Back-propagation

### 1.1 (a)

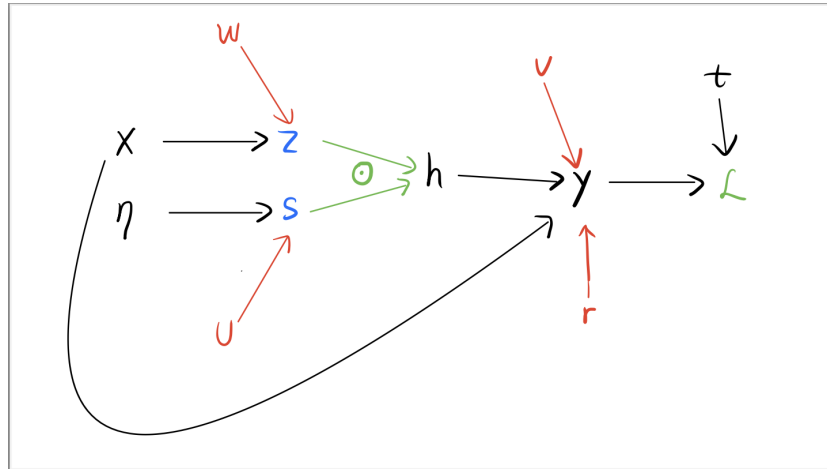


Figure 1: Computation graph

### 1.2 (b)

$$\frac{d\sigma(x)}{dx} = \frac{d}{dx} \left( \frac{1}{1 + e^{-x}} \right) \quad (1)$$

$$= \frac{-1}{(1 + e^{-x})^2} \cdot \frac{d}{dx} (1 + e^{-x}) \quad (2)$$

$$= \frac{-1}{(1 + e^{-x})^2} \cdot (-e^{-x}) \quad (3)$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2} \quad (4)$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \quad (5)$$

$$= \sigma(x)(1 - \sigma(x)) \quad (6)$$

### 1.3 (c)

First, we compute the error signal for the output  $y$ :

$$\frac{\partial L}{\partial y} = \frac{t}{y} - \frac{1-t}{1-y}$$

Then, we compute the error signal for the hidden layer activations  $h$ :

$$\frac{\partial L}{\partial h} = \left( \frac{\partial L}{\partial y} \right) \left( \frac{\partial y}{\partial h} \right) = \left( \frac{t}{y} - \frac{1-t}{1-y} \right) \sigma(v^\top h + r^\top x) (1 - \sigma(v^\top h + r^\top x)) v$$

Next, we compute the error signal for  $z$  and  $s$ , since  $a \odot b$  can be viewed as  $\text{diag}(a)\text{diag}(b)$ , we then have:

$$\begin{aligned} \frac{\partial L}{\partial z} &= \left( \frac{\partial L}{\partial h} \right) \left( \frac{\partial h}{\partial z} \right) = \left( \frac{\partial L}{\partial h} \right) \odot s \\ \frac{\partial L}{\partial s} &= \left( \frac{\partial L}{\partial h} \right) \left( \frac{\partial h}{\partial s} \right) = \left( \frac{\partial L}{\partial h} \right) \odot z \end{aligned}$$

Then, we compute the error signals for the weight matrices  $W$  and  $U$ :

$$\begin{aligned} \frac{\partial L}{\partial W} &= \left( \frac{\partial L}{\partial z} \right) \left( \frac{\partial z}{\partial W} \right) = \left( \frac{\partial L}{\partial z} \right) x^\top \\ \frac{\partial L}{\partial U} &= \left( \frac{\partial L}{\partial s} \right) \left( \frac{\partial s}{\partial U} \right) = \left( \frac{\partial L}{\partial s} \right) \eta^\top \end{aligned}$$

Next, we compute the error signals for the bias vectors  $r$  and  $v$ :

$$\begin{aligned} \frac{\partial L}{\partial r} &= \left( \frac{\partial L}{\partial y} \right) \left( \frac{\partial y}{\partial r} \right) = \left( \frac{\partial L}{\partial y} \right) \sigma(v^\top h + r^\top x) (1 - \sigma(v^\top h + r^\top x)) x \\ \frac{\partial L}{\partial v} &= \left( \frac{\partial L}{\partial y} \right) \left( \frac{\partial y}{\partial v} \right) = \left( \frac{\partial L}{\partial y} \right) \sigma(v^\top h + r^\top x) (1 - \sigma(v^\top h + r^\top x)) h \end{aligned}$$

Finally, we compute the error signals for the input vectors  $x$  and  $\eta$ :

$$\frac{\partial L}{\partial \eta} = \left( \frac{\partial L}{\partial s} \right) \left( \frac{\partial s}{\partial \eta} \right) = U^\top \left( \frac{\partial L}{\partial s} \right)$$

$$\begin{aligned} \frac{\partial L}{\partial x} &= \left( \frac{\partial L}{\partial y} \right) \left( \frac{\partial y}{\partial x} \right) + \left( \frac{\partial L}{\partial z} \right) \left( \frac{\partial z}{\partial x} \right) \\ &= \left( \frac{\partial L}{\partial y} \right) \sigma(v^\top h + r^\top x) (1 - \sigma(v^\top h + r^\top x)) r + W^\top \left( \frac{\partial L}{\partial z} \right) \end{aligned}$$

## 2 Fitting a Naive Bayes Model

### 2.1 (a)

$$L(\theta, \pi | X) = P(C | \pi) \prod_{j=1}^{784} P(x_j | C, \theta_{jC})$$

$$\ell(\theta, \pi | X) = \log(P(C | \pi)) + \sum_{j=1}^{784} \log(P(x_j | C, \theta_{jC}))$$

For  $\pi_{MLE}$ :

$$\log(P(C | \pi)) = \sum_{i=1}^N \left( \log\left(1 - \sum_{j=0}^8 \pi_j\right) t_{i9}^{(C)} + \sum_{j=0}^8 \log(\pi_j) \cdot t_{ij}^{(C)} \right)$$

For  $j \in [0, 8]$

$$\frac{\partial \log(P(C | \pi))}{\partial \pi_j} = \sum_{i=1}^N \left( \frac{-t_{i9}^{(C)}}{1 - \sum_{j=0}^8 \pi_j} + \frac{t_{ij}^{(C)}}{\pi_j} \right) \stackrel{\text{set}}{=} 0$$

$$\pi_9 \sum_{i=1}^N t_{ij}^{(C)} = \pi_j \sum_{i=1}^N t_{i9}^{(C)}$$

$$\frac{\hat{\pi}_j}{\hat{\pi}_9} = \frac{\sum_{i=1}^N t_{ij}^{(C)}}{\sum_{i=1}^N t_{i9}^{(C)}} = \frac{\# \text{ label } j}{\# \text{ label } 9}$$

Since  $\sum_{i=0}^9 \pi_i = 1$ , we then have:

$$\frac{1}{\hat{\pi}_9} = \frac{\sum_{i=0}^9 \hat{\pi}_i}{\hat{\pi}_9} = \frac{N}{\# \text{ label } 9} \Rightarrow \hat{\pi}_j = \frac{\# \text{ of label } j}{N}$$

For  $\theta_{MLE}$

$$\begin{aligned}\sum_{j=1}^{784} \log(p(x_j|c, \theta_{jc})) &= \sum_{i=1}^N \sum_{j=1}^{784} x_j^{(i)} \log(\theta_{jc}) + (1-x_j^{(i)}) \log(1-\theta_{jc}) \\ &= \sum_{i=1}^N \sum_{j=1}^{784} \sum_{c=0}^9 t_c(x_j^{(i)} \log(\theta_{jc}) + (1-x_j^{(i)}) \log(1-\theta_{jc}))\end{aligned}$$

$$\begin{aligned}\frac{\partial \sum_{j=1}^{784} \log(p(x_j|c, \theta_{jc}))}{\partial \theta_{jc}} &= \sum_{i=1}^N t_c^{(i)} \left( \frac{x_j^{(i)}}{\theta_{jc}} - \frac{1-x_j^{(i)}}{1-\theta_{jc}} \right) \\ &= \sum_{j=1}^N t_c^{(i)} \frac{x_j^{(i)} - \theta_{jc}}{\theta_{jc}(1-\theta_{jc})} \\ &\stackrel{\text{set}}{=} 0\end{aligned}$$

$$\begin{aligned}\sum_{i=1}^N t_c^{(i)} x_j^{(i)} &= \theta_{jc} \sum_{i=1}^N t_c^{(i)} \Rightarrow \\ \hat{\theta}_{jc} &= \frac{\sum_{i=1}^N t_c^{(i)} x_j^{(i)}}{\sum_{i=1}^N t_c^{(i)}} = \frac{\text{\# of feature } j \text{ with label } c}{\text{\# of label } c}\end{aligned}$$

Figure 2: Derive the maximum likelihood estimator

2.2 (b)

$$\begin{aligned}
 \log(t|x, \theta, \pi) &= \log\left(\frac{P(t|\theta, \pi) P(x|t, \theta, \pi)}{P(x|\theta, \pi)}\right) \\
 &= \log\left(\frac{P(t|\pi) P(x|t_c, \theta_c)}{P(x|\theta_c)}\right) \\
 &= \underbrace{\log(P(t|\theta, \pi))}_{(1)} + \underbrace{\log(P(t_c, \theta_c))}_{(2)} - \underbrace{\log(P(x|\theta_c))}_{(3)}
 \end{aligned}$$

$$(1) \log(P(t|\theta, \pi)) = \log(\pi_c)$$

$$\begin{aligned}
 (2) \log(P(x|t_c, \theta_c)) &= \log\left(\prod_{j=1}^{784} \theta_{jc}^{x_j} (1-\theta_{jc})^{1-x_j}\right) \\
 &= \sum_{j=1}^{784} x_j \log(\theta_{jc}) + (1-x_j) \log(1-\theta_{jc})
 \end{aligned}$$

$$(3) \log(P(x|\theta_c)) = \log\left(\sum_{c=0}^9 \pi_c \underbrace{\prod_{j=1}^{784} \theta_{jc}^{x_j} (1-\theta_{jc})^{(1-x_j)}}_{(4)}\right)$$

$$\begin{aligned}
 (4) \prod_{j=1}^{784} \theta_{jc}^{x_j} (1-\theta_{jc})^{(1-x_j)} &= e^{\log\left(\prod_{j=1}^{784} \theta_{jc}^{x_j} (1-\theta_{jc})^{(1-x_j)}\right)} \\
 &= e^{\sum_{j=1}^{784} x_j \log(\theta_{jc}) + (1-x_j) \log(1-\theta_{jc})}
 \end{aligned}$$

Figure 3: Derive the log-likelihood

### 2.3 (c)

The output is:

Average log-likelihood for MLE is nan

Average log-likelihood for MLE is not defined.

Possible Reason:  $\hat{\theta}_{jc}$  is zero, causes  $\log(\hat{\theta}_{jc})$  to be undefined.

### 2.4 (d)

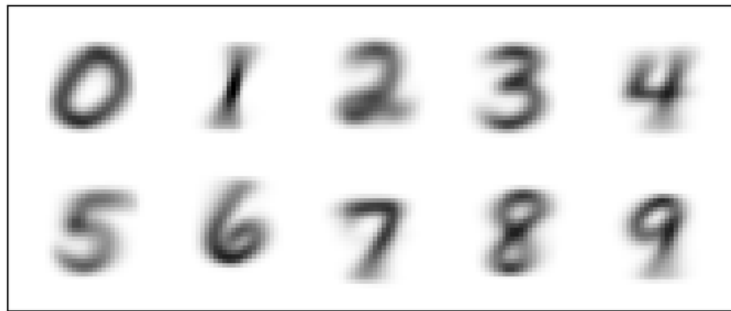


Figure 4: MLE images

2.5 (e)

According to the given prior

$$P(D|\theta) = \frac{\theta^2(1-\theta)^2}{B(3,3)} \propto \theta^2(1-\theta)^2$$

const

$$\begin{aligned} P(\theta) &= \prod_{i=1}^N t_c^{(i)} \theta_{jc}^{x_j^{(i)}} (1-\theta_{jc})^{1-x_j^{(i)}} \\ &= \theta_{jc}^{\sum_{i=1}^N t_c^{(i)} x_j^{(i)}} (1-\theta_{jc})^{\sum_{i=1}^N t_c^{(i)} (1-x_j^{(i)})} \\ &= \theta_{jc}^{\text{\# of label } c} (1-\theta_{jc})^{(\text{\# of label } c) - (\text{\# of feature } j \text{ with label } c)} \end{aligned}$$

Denote # of label  $c$  as  $N_c$

# of label  $c$  with feature  $j$  as  $N_{cj}$

By Bayes Rule:

$$\begin{aligned} P(\theta_{jc}|D) &\propto P(\theta_{jc}) P(D|\theta_{jc}) \\ &= \theta_{jc}^{N_c+2} (1-\theta_{jc})^{N_c-N_{cj}+2} \end{aligned}$$

$$\frac{\partial \log P(\theta_{jc}|D)}{\partial \theta_{jc}} = \frac{N_c+2}{\theta_{jc}} - \frac{N_c-N_{cj}+2}{1-\theta_{jc}} \stackrel{\text{set}}{=} 0$$

$$\hat{\theta}_{jc} = \frac{N_{cj}+2}{N_c+4}$$

Figure 5: Derive the MAP estimator

## 2.6 (f)

The output is:

```
Average log-likelihood for MAP is -3.3570631378602873
Training accuracy for MAP is 0.9670433333333334
Test accuracy for MAP is 0.9632
```

## 2.7 (g)

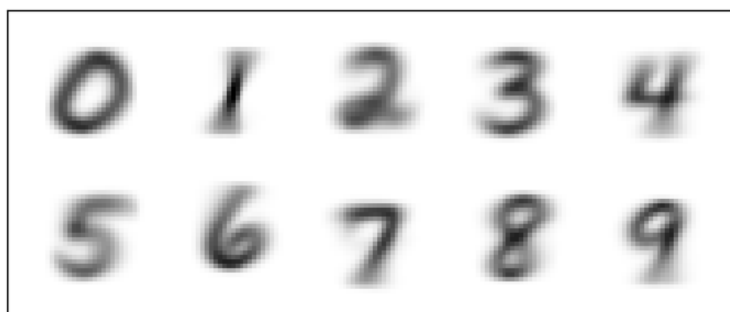


Figure 6: MAP images

## 2.8 (h)

Advantage of naive bayes for this problem:

Naive Bayes is a simple and easy to implement and understand. It's a fast algorithm that is computationally efficient for this task.

Disadvantage of naive bayes for this problem:

Naive Bayes assumes that all features are independent of each other, which is not always true for image data. Pixels in an image are often correlated with each other. If we move the main part of a image to its left by just one pixel, the result would be completely different. It would better to use convolutional neural network that can extract features from images.



### 3 Logistic Regression with Gaussian Prior.

#### 3.1 (a)

The log-likelihood of the parameter vector  $\theta$  for logistic regression is given by:

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^N \log p(y^{(i)}|x^{(i)}, \theta) \\ &= \sum_{i=1}^N y^{(i)} \log \frac{1}{1 + \exp(-x^{(i)T}\theta)} + (1 - y^{(i)}) \log \frac{\exp(-x^{(i)T}\theta)}{1 + \exp(-x^{(i)T}\theta)} \\ &= \sum_{i=1}^N y^{(i)} x^{(i)T} \theta - \log(1 + \exp(x^{(i)T}\theta)).\end{aligned}$$

To optimize the resulting log-likelihood, we can use gradient ascent algorithm. We can iteratively update the parameter vector until convergence. The gradient of the log-likelihood with respect to  $\theta$  is given by:

$$\begin{aligned}\nabla_{\theta} \ell(\theta) &= \sum_{i=1}^N (y^{(i)} - p(y^{(i)} = 1|x^{(i)}, \theta)) x^{(i)} \\ &= \sum_{i=1}^N (y^{(i)} - \frac{1}{1 + \exp(-x^{(i)T}\theta)}) x^{(i)}\end{aligned}$$

#### 3.2 (b)

The prior distribution over the parameters  $\theta$  is assumed to be a Gaussian distribution with zero mean and variance  $\sigma^2$ . Therefore, the prior density function  $p(\theta)$  is given by:

$$p(\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} \|\theta\|^2\right)$$

The logarithm of this function is:

$$\begin{aligned}\log p(\theta) &= -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\theta\|^2 \\ &= -\frac{1}{2\sigma^2} \|\theta\|^2 + \text{const}\end{aligned}$$

where const is a constant term that does not depend on  $\theta$ .

Therefore, The MAP log-likelihood of the parameter vector  $\theta$  for logistic regression with the Gaussian prior is given by:

$$\begin{aligned}\ell_{\text{MAP}}(\theta) &= \sum_{i=1}^N \log p(y^{(i)}|x^{(i)}, \theta) + \log p(\theta) \\ &= \sum_{i=1}^N y^{(i)} \log \frac{1}{1 + \exp(-x^{(i)T}\theta)} + (1 - y^{(i)}) \log \frac{\exp(-x^{(i)T}\theta)}{1 + \exp(-x^{(i)T}\theta)} - \frac{1}{2\sigma^2} \|\theta\|^2 + \text{const} \\ &= \sum_{i=1}^N y^{(i)} x^{(i)T} \theta - \log(1 + \exp(x^{(i)T}\theta)) - \frac{1}{2\sigma^2} \|\theta\|^2 + \text{const}.\end{aligned}$$

## 4 Gaussian Discriminant Analysis.

All of implementation can be referred to file `q4.py`.

### 4.1 (a)

Let  $X$  be a 64-dimensional feature vector, and  $y$  be the corresponding label. The probability density function (PDF) of the feature vector given the class  $k$  can be written as:

$$p(X|y = k; \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp \left( -\frac{1}{2} (X - \mu_k)^T \Sigma_k^{-1} (X - \mu_k) \right)$$

where  $\mu_k$  is the mean vector of class  $k$ ,  $\Sigma_k$  is the covariance matrix of class  $k$ , and  $d=64$  is the dimensionality of the feature vector.

The maximum likelihood estimation of the mean vector and covariance matrix for class  $k$  can be obtained by maximizing the log-likelihood function:

$$\ell(\mu_k, \Sigma_k) = \sum_{i=1}^n \log p(X^{(i)} | y^{(i)} = k; \mu_k, \Sigma_k)$$

where  $X^{(i)}$  is the  $i$ th feature vector and  $y^{(i)}$  is the corresponding label. Taking the derivative of  $\ell(\mu_k, \Sigma_k)$  with respect to  $\mu_k$  and  $\Sigma_k$  and setting them to zero, we get:

$$\mu_k = \frac{1}{n_k} \sum_{i=1}^n 1_{y^{(i)}=k} X^{(i)}$$

$$\Sigma_k = \frac{1}{n_k} \sum_{i=1}^n 1_{y^{(i)}=k} (X^{(i)} - \mu_k)(X^{(i)} - \mu_k)^T$$

where  $n_k$  is the number of training examples with label  $k$ .

We have added  $0.01\mathbf{I}$  to each covariance matrix to ensure numerical stability.

After computing the MLE parameters  $\mu_k, \Sigma_k$  for  $k=0, \dots, 9$ , we can compute the log-likelihood for an example  $X$  as:

$$\begin{aligned} \ell(\mu_k, \Sigma_k) &= \sum_{i=1}^n \log p(X^{(i)} | y^{(i)} = k; \mu_k, \Sigma_k) \\ &= \sum_{i=1}^n -\frac{1}{2} d \cdot \log(2\pi) - \frac{1}{2} |\Sigma_k| - \frac{1}{2} (X - \mu_k)^T \Sigma_k^{-1} (X - \mu_k) \end{aligned}$$

Using Bayes rule, we can compute the posterior probabilities for each class:

$$p(y = k|X) = \frac{p(X|y = k)p(y = k)}{\sum_{j=0}^9 p(X|y = j)p(y = j)}$$

Finally, we can make a prediction by choosing the class with the highest posterior probability.

$$\hat{y} = \arg \max_{k \in \{0, 1, \dots, 9\}} p(y = k|X)$$

To compute the average conditional log-likelihood, we can sum over the log-likelihoods of all examples and divide by the total number of examples.

I got the following output result as my average conditional log-likelihood:

Average Train LL: -2.4272095296626754

Average Test LL: -2.4992582962493013

## 4.2 (b)

I got the following output result for the prediction accuracy:

```
Train Accuracy:  0.9814285714285714
Test Accuracy:   0.97275
```

## 4.3 (c)

I got the following output under the assumption that the covariance matrix is diagonal:

```
Average Train LL: -3.5333505152668367
Average Test LL:  -3.589845459749884
Train Accuracy:   0.85
Test Accuracy:    0.84
```

### **Analyze:**

According to the output, the model without the assumption of diagonal covariance matrix performs better on both the train and test sets. Specifically, the average conditional log-likelihood is higher for the non-diagonal covariance matrix model, indicating that it is able to better capture the distribution of the data. The accuracy is also higher for the non-diagonal covariance matrix model on both train and test sets.

The reason is that, diagonal covariance matrix cannot model dependence between pixels. If we move the main part of a image to its left by just one pixel, the result would be completely different for the diagonal covariance matrix.