

## Assignment #2

In this assignment, we'll fit both generative and discriminative models to the MNIST dataset of handwritten numbers. Each datapoint in the MNIST [<http://yann.lecun.com/exdb/mnist/>] dataset is a 28x28 black-and-white image of a number in  $\{0 \dots 9\}$ , and a label indicating which number.

MNIST is the 'fruit fly' of machine learning - a simple standard problem useful for comparing the properties of different algorithms. Python code for loading and plotting MNIST is attached.

You can use whichever programming language you like, and libraries for loading and plotting data. You'll need to write your own initialization, fitting, and prediction code. You can use automatic differentiation in your code, but must still answer the gradient questions.

For this assignment, we'll *binarize* the dataset, converting the grey pixel values to either black or white (0 or 1) with  $> 0.5$  being the cutoff. When comparing models, we'll need a training and test set. Use the first 10000 samples for training, and another 10000 for testing. This is all done for you in the starter code. Hint: Also build a dataset of only 100 training samples to use when debugging, to make loading and training faster.

### Problem 1 (Basic Naïve Bayes, 10 points)

In this question, we'll fit a naïve Bayes model to the MNIST digits using maximum likelihood. Naïve Bayes defines the joint probability of the each datapoint  $\mathbf{x}$  and its class label  $c$  as follows:

$$p(\mathbf{x}, c | \boldsymbol{\theta}, \pi) = p(c | \pi) p(\mathbf{x} | c, \boldsymbol{\theta}_c) = p(c | \pi) \prod_{d=1}^{784} p(x_d | c, \theta_{cd}) \quad (1)$$

For binary data, we can use the Bernoulli likelihood:

$$p(x_d | c, \theta_{cd}) = \text{Ber}(x_d | \theta_{cd}) = \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)} \quad (2)$$

Which is just a way of expressing that  $p(x_d = 1 | c, \theta_{cd}) = \theta_{cd}$ .

For  $p(c | \pi)$ , we can just use a categorical distribution:

$$p(c | \pi) = \text{Cat}(c | \pi) = \pi_c \quad (3)$$

Note that we need  $\sum_{i=0}^9 \pi_i = 1$ .

- (a) Derive the *maximum likelihood estimate* (MLE) for the class-conditional pixel means  $\boldsymbol{\theta}$ . Hint: We saw in lecture that MLE can be thought of as 'counts' for the data, so what should  $\hat{\theta}_{cd}$  be counting?
- (b) Derive the *maximum a posteriori* (MAP) estimate for the class-conditional pixel means  $\boldsymbol{\theta}$ , using a Beta(2, 2) prior on each  $\theta$ . Hint: it has a simple final form, and you can ignore the Beta normalizing constant.
- (c) Fit  $\boldsymbol{\theta}$  to the training set using the MAP estimator. Plot  $\boldsymbol{\theta}$  as 10 separate greyscale images, one for each class.
- (d) Derive the predictive log-likelihood  $\log p(c | \mathbf{x}, \boldsymbol{\theta}, \pi)$  for a single training image.
- (e) Given parameters fit to the training set, and  $\pi_c = \frac{1}{10}$ , report both the average predictive log-likelihood per datapoint,  $\frac{1}{N} \sum_{i=1}^N \log p(c_i | x_i, \boldsymbol{\theta}, \pi)$  and the predictive accuracy on both the training and test set. The predictive accuracy is defined as the fraction of examples where the true class  $t = \arg\max_c p(c | \mathbf{x}, \boldsymbol{\theta}, \pi)$ .

The takeaway of this question is that we can automatically derive a learning algorithm just by first defining a joint probability!

- (a) In the following derivation,  $n$  refers only to samples where the ground truth label is of class  $c$ . Other samples are disregarded, since they are not involved in estimating  $\theta_{cd}$ .

$$\begin{aligned}
\mathbb{L}(\theta_{cd}; \mathbf{x}) &= \sum_{n=1}^N \log \left( p(x_d^{(n)} | c, \theta_{cd}) \right) \\
&= \sum_{n=1}^N \log \left( \theta_{cd}^{x_d^{(n)}} (1 - \theta_{cd})^{(1-x_d^{(n)})} \right) \\
&= \sum_{n=1}^N \log \theta_{cd}^{x_d^{(n)}} + \sum_{n=1}^N \log (1 - \theta_{cd})^{(1-x_d^{(n)})} \\
&= \log \theta_{cd} \sum_{n=1}^N x_d^{(n)} + \log (1 - \theta_{cd}) \sum_{n=1}^N (1 - x_d^{(n)}) \\
&= N_d^1 \log \theta_{cd} + N_d^0 \log (1 - \theta_{cd})
\end{aligned}$$

Where  $N_d^1$  and  $N_d^0$  are the number of times pixel  $x_d$  takes the value of 1 and 0, respectively.

$$\frac{\partial \mathbb{L}}{\partial \theta_{cd}} = \frac{N_d^1}{\theta_{cd}} - \frac{N_d^0}{1 - \theta_{cd}}$$

The MLE is found by equating  $\frac{\partial \mathbb{L}}{\partial \theta_{cd}}$  to 0.

$$\begin{aligned}
0 &= \frac{N_d^1}{\hat{\theta}_{cd}} - \frac{N_d^0}{1 - \hat{\theta}_{cd}} \\
\hat{\theta}_{cd} &= \frac{N_d^1}{N_d^1 + N_d^0} \\
\hat{\theta}_{cd} &= \frac{N_d^1}{N_d}
\end{aligned}$$

Where  $N_d$  is the total numbers of samples for class  $c$ . This result for the MLE corresponds to the class conditional mean of pixel  $d$ .

- (b)

$$\begin{aligned}
\mathbb{L}(\theta_{cd}; \mathbf{x}) &= \sum_{n=1}^N \log \left( p(x_d^{(n)} | c, \theta_{cd}) \right) + \log (p(\theta_{cd})) \\
&\propto \sum_{n=1}^N \log \left( \theta_{cd}^{x_d^{(n)}} (1 - \theta_{cd})^{(1-x_d^{(n)})} \right) + \log (\theta_{cd}(1 - \theta_{cd}))
\end{aligned}$$

Where the Beta normalization can be ignored since it will not affect the MAP

$$\begin{aligned}
&= N_d^1 \log \theta_{cd} + N_d^0 \log (1 - \theta_{cd}) + \log \theta_{cd} + \log (1 - \theta_{cd}) \\
&= (N_d^1 + 1) \log \theta_{cd} + (N_d^0 + 1) \log (1 - \theta_{cd})
\end{aligned}$$

$$\frac{\partial \mathbb{L}}{\partial \theta_{cd}} = \frac{N_d^1 + 1}{\theta_{cd}} - \frac{N_d^0 + 1}{1 - \theta_{cd}}$$

$$\begin{aligned}
0 &= \frac{N_d^1 + 1}{\hat{\theta}_{cd}} - \frac{N_d^0 + 1}{1 - \hat{\theta}_{cd}} \\
\hat{\theta}_{cd} &= \frac{N_d^1 + 1}{N_d^1 + N_d^0 + 2} \\
\hat{\theta}_{cd} &= \frac{N_d^1 + 1}{N_d + 2}
\end{aligned}$$

(c) Figure 1 shows the fit  $\theta_{cd}$  parameters for each class.

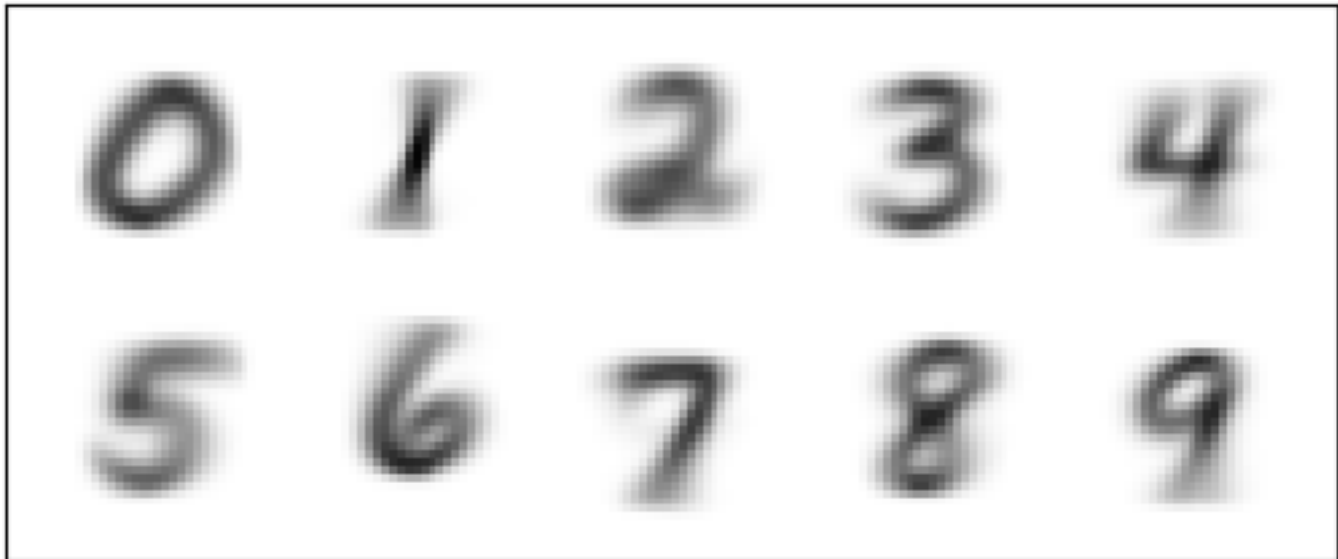


Figure 1:  $\theta_{cd}$  parameters for each class.

(d)

$$\begin{aligned}
 p(c | \mathbf{x}, \boldsymbol{\theta}, \pi) &= \frac{p(c, \mathbf{x} | \boldsymbol{\theta}, \pi)}{p(\mathbf{x} | \boldsymbol{\theta}, \pi)} \\
 &= \frac{p(c, \mathbf{x} | \boldsymbol{\theta}, \pi)}{\sum_{c'=0}^9 p(p(c', \mathbf{x} | \boldsymbol{\theta}, \pi))} \\
 &= \frac{p(c | \pi) \prod_{d=1}^{784} p(x_d | c, \theta_{cd})}{\sum_{c'=0}^9 p(c' | \pi) \prod_{d=1}^{784} p(x_d | c', \theta_{c'd})} \\
 &= \frac{\frac{1}{10} \prod_{d=1}^{784} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}{\frac{1}{10} \sum_{c'=0}^9 \prod_{d=1}^{784} \theta_{c'd}^{x_d} (1 - \theta_{c'd})^{(1-x_d)}} \\
 &= \frac{\prod_{d=1}^{784} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}{\sum_{c'=0}^9 \prod_{d=1}^{784} \theta_{c'd}^{x_d} (1 - \theta_{c'd})^{(1-x_d)}}
 \end{aligned}$$

$$\begin{aligned}
 \log(p(c | \mathbf{x}, \boldsymbol{\theta}, \pi)) &= \log \left( \prod_{d=1}^{784} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)} \right) - \log \left( \sum_{c'=0}^9 \prod_{d=1}^{784} \theta_{c'd}^{x_d} (1 - \theta_{c'd})^{(1-x_d)} \right) \\
 &= \sum_{d=1}^{784} \log \left( \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)} \right) - \log \left( \sum_{c'=0}^9 \prod_{d=1}^{784} \theta_{c'd}^{x_d} (1 - \theta_{c'd})^{(1-x_d)} \right)
 \end{aligned}$$

- (e) Training average predictive log-likelihood per datapoint: -3.35405  
 Test average predictive log-likelihood per datapoint: -3.02397  
 Training predictive accuracy: 83.5883%  
 Test predictive accuracy: 84.67%

**Problem 2** (Advanced Naïve Bayes, 10 points)

One of the advantages of generative models is that they can handle missing data, or be used to answer different sorts of questions about the model.

- (a) True or false: Given our model's assumptions, any two pixels  $x_i$  and  $x_j$  where  $i \neq j$  are independent given  $c$ .
- (b) True or false: Given our model's assumptions, any two pixels  $x_i$  and  $x_j$  where  $i \neq j$  are independent when marginalizing over  $c$ .
- (c) Using the parameters fit in question 1, produce random image samples from the model. That is, randomly sample and plot 10 binary images from the marginal distribution  $p(\mathbf{x}|\theta, \pi)$ . Hint: Use ancestral sampling.
- (d) Derive  $p(\mathbf{x}_{bottom}|\mathbf{x}_{top}, \theta, \pi)$ , the joint distribution over the bottom half of an image given the top half, conditioned on your fit parameters.
- (e) Derive  $p(\mathbf{x}_{i \in bottom}|\mathbf{x}_{top}, \theta, \pi)$ , the marginal distribution of a single pixel in the bottom half of an image given the top half, conditioned on your fit parameters.
- (f) For 20 images from the training set, plot the top half the image concatenated with the marginal distribution over each pixel in the bottom half.

(a) True

(b) False

(c) Figure 2 shows samples from the marginal distribution  $p(\mathbf{x}_{i \in bottom}|\mathbf{x}_{top}, \theta, \pi)$ .

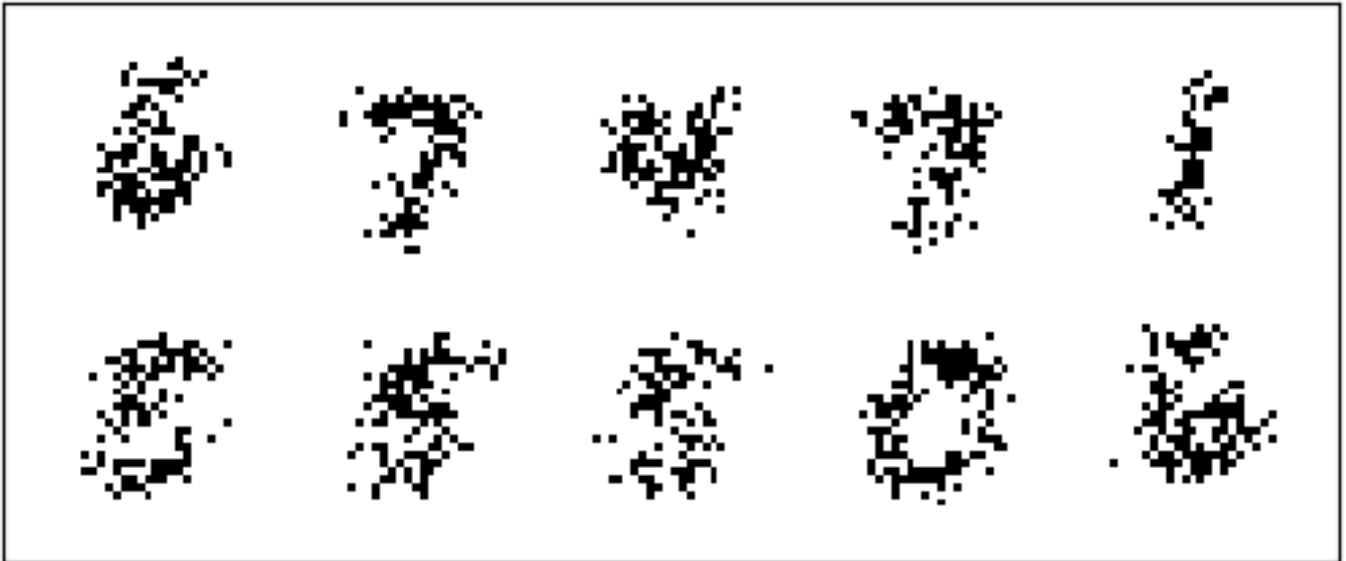


Figure 2: 10 samples from the marginal distribution  $p(\mathbf{x}_{i \in bottom}|\mathbf{x}_{top}, \theta, \pi)$ .

(d)

$$\begin{aligned}
p(\mathbf{x}_{bottom}|\mathbf{x}_{top}, \boldsymbol{\theta}, \pi) &= \frac{p(\mathbf{x}_{bottom}, \mathbf{x}_{top}|\boldsymbol{\theta}, \pi)}{p(\mathbf{x}_{top}|\boldsymbol{\theta}, \pi)} \\
&= \frac{p(\mathbf{x}|\boldsymbol{\theta}, \pi)}{p(\mathbf{x}_{top}|\boldsymbol{\theta}, \pi)} \\
&= \frac{\sum_{c=0}^9 p(c|\pi) p(\mathbf{x}|c, \theta_c)}{\sum_{c=0}^9 p(c|\pi) p(\mathbf{x}_{top}|c, \theta_c)} \\
&= \frac{\sum_{c=0}^9 \pi_c \prod_{d=1}^{784} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}{\sum_{c=0}^9 \pi_c \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}} \\
&= \frac{\frac{1}{10} \sum_{c=0}^9 \prod_{d=1}^{784} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}{\frac{1}{10} \sum_{c=0}^9 \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}} \\
&= \frac{\sum_{c=0}^9 \prod_{d=1}^{784} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}{\sum_{c=0}^9 \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}
\end{aligned}$$

(e)

$$\begin{aligned}
p(\mathbf{x}_{i \in bottom}|\mathbf{x}_{top}, \boldsymbol{\theta}, \pi) &= \frac{p(\mathbf{x}_{i \in bottom}, \mathbf{x}_{top}|\boldsymbol{\theta}, \pi)}{p(\mathbf{x}_{top}|\boldsymbol{\theta}, \pi)} \\
&= \frac{\sum_{c=0}^9 \pi_c \theta_{ci}^{x_i} (1 - \theta_{ci})^{(1-x_i)} \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}{\sum_{c=0}^9 \pi_c \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}} \\
&= \frac{\frac{1}{10} \sum_{c=0}^9 \theta_{ci}^{x_i} (1 - \theta_{ci})^{(1-x_i)} \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}{\frac{1}{10} \sum_{c=0}^9 \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}} \\
&= \frac{\sum_{c=0}^9 \theta_{ci}^{x_i} (1 - \theta_{ci})^{(1-x_i)} \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}{\sum_{c=0}^9 \prod_{d=1}^{392} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)}}
\end{aligned}$$

(f) Figure 3 shows the top half of 20 images in the training set concatenated with the marginal distribution over each pixel in the bottom half.

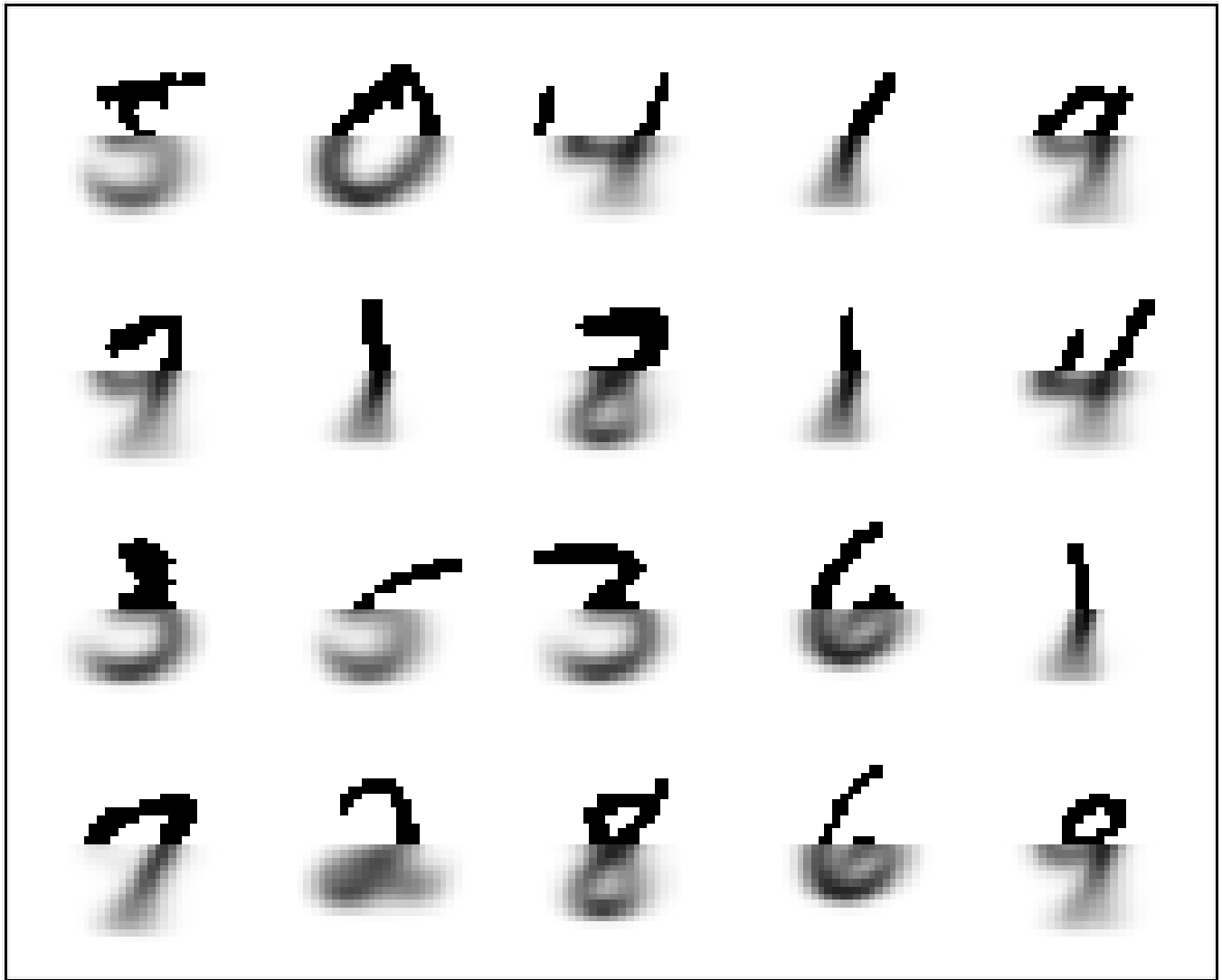


Figure 3: Top half of 20 images in the training set concatenated with the marginal distribution over each pixel in the bottom half.

**Problem 3** (Logistic Regression, 10 points)

Now, we'll fit a simple predictive model using gradient descent. Our model will be multiclass logistic regression:

$$p(c|\mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c'=0}^9 \exp(\mathbf{w}_{c'}^T \mathbf{x})} \quad (4)$$

You can ignore biases for this question.

- How many parameters does this model have?
- Derive the gradient of the predictive log-likelihood w.r.t.  $\mathbf{w}$ :  $\nabla_{\mathbf{w}} \log p(c|\mathbf{x}, \mathbf{w})$
- Code up a gradient-based optimizer of your choosing, it can be just vanilla gradient descent, and use it to optimize  $\mathbf{w}$  to maximize the log-likelihood of the training set, and plot the resulting parameters using one image per class. Since this objective is concave, you can initialize at all zeros. Using automatic differentiation is permitted, so you can use autograd to get gradients for use by your optimizer, and using minibatches is optional. However, you are not permitted to use optimizers which come built in to packages! Hint: Use `scipy.logsumexp` or its equivalent to make your code more numerically stable.
- Given parameters fit to the training set, report both the average predictive log-likelihood per datapoint, and the predictive accuracy on both the training and test set. How does it compare to Naïve Bayes?

(a)  $10(784) = 7840$

(b) Let  $s_c = \mathbf{w}_c^T \mathbf{x}$

$$p(c|\mathbf{x}, \mathbf{w}) = \frac{\exp(s_c)}{\sum_{c'=0}^9 \exp(s_{c'})}$$

$$\frac{\partial \log(p(c|\mathbf{x}, \mathbf{w}))}{\partial w_{ki}} = \frac{\partial \log(p(c|\mathbf{x}, \mathbf{w}))}{\partial p(c|\mathbf{x}, \mathbf{w})} \sum_{l=0}^9 \frac{\partial p(c|\mathbf{x}, \mathbf{w})}{\partial s_l} \frac{\partial s_l}{\partial w_{ki}}$$

- $\frac{\partial \log(p(c|\mathbf{x}, \mathbf{w}))}{\partial p(c|\mathbf{x}, \mathbf{w})} = \frac{1}{p(c|\mathbf{x}, \mathbf{w})}$
- $\frac{\partial s_l}{\partial w_{ki}} = \begin{cases} x_i & \text{if } k = 0 \\ 0 & \text{otherwise} \end{cases} \quad \therefore \text{only evaluate } \frac{\partial p(c|\mathbf{x}, \mathbf{w})}{\partial s_l} \text{ for } l = k$
- $\frac{\partial p(c|\mathbf{x}, \mathbf{w})}{\partial s_k} = \frac{\partial}{\partial s_k} \frac{\exp(s_c)}{\sum_{c'=0}^9 \exp(s_{c'})}$

if  $k = c$ ,

$$\begin{aligned} \frac{\partial}{\partial s_k} \frac{\exp(s_c)}{\sum_{c'=0}^9 \exp(s_{c'})} &= \frac{\exp(s_c) \sum_{c'=0}^9 \exp(s_{c'}) - \exp(s_c)^2}{\left( \sum_{c'=0}^9 \exp(s_{c'}) \right)^2} \\ &= p(c|\mathbf{x}, \mathbf{w}) - p(c|\mathbf{x}, \mathbf{w})^2 \\ &= p(c|\mathbf{x}, \mathbf{w}) (1 - p(c|\mathbf{x}, \mathbf{w})) \end{aligned}$$

if  $k \neq c$ ,

$$\begin{aligned} \frac{\partial}{\partial s_k} \frac{\exp(s_c)}{\sum_{c'=0}^9 \exp(s_{c'})} &= \frac{-\exp(s_k) \exp(s_c)}{\left( \sum_{c'=0}^9 \exp(s_{c'}) \right)^2} \\ &= -p(k|\mathbf{x}, \mathbf{w}) p(c|\mathbf{x}, \mathbf{w}) \end{aligned}$$

$$\begin{aligned}
&\therefore \text{Using (1), (2), and (3)} \\
&\text{if } k = c, \\
&\frac{\partial \log(p(c|\mathbf{x}, \mathbf{w}))}{\partial w_{ki}} = \frac{1}{p(c|\mathbf{x}, \mathbf{w})} p(c|\mathbf{x}, \mathbf{w}) (1 - p(c|\mathbf{x}, \mathbf{w})) x_i \\
&\quad = (1 - p(c|\mathbf{x}, \mathbf{w})) x_i \\
&\text{if } k \neq c, \\
&\frac{\partial \log(p(c|\mathbf{x}, \mathbf{w}))}{\partial w_{ki}} = \frac{-1}{p(c|\mathbf{x}, \mathbf{w})} p(k|\mathbf{x}, \mathbf{w}) p(c|\mathbf{x}, \mathbf{w}) x_i \\
&\quad = -p(k|\mathbf{x}, \mathbf{w}) x_i
\end{aligned}$$

In matrix form, where  $\mathbf{y}$  is a one-hot label vector,  $\mathbf{p}$  is the softmax output vector, and  $w_{ij}$  corresponds to the weight from the  $i$ 'th input to the  $j$ 'th softmax output,

$$\nabla_w \log(p(c|\mathbf{x}, \mathbf{w})) = \mathbf{x}(\mathbf{y} - \mathbf{p})$$

(c) Figure 4 shows the fit  $\theta_{cd}$  parameters for each class.

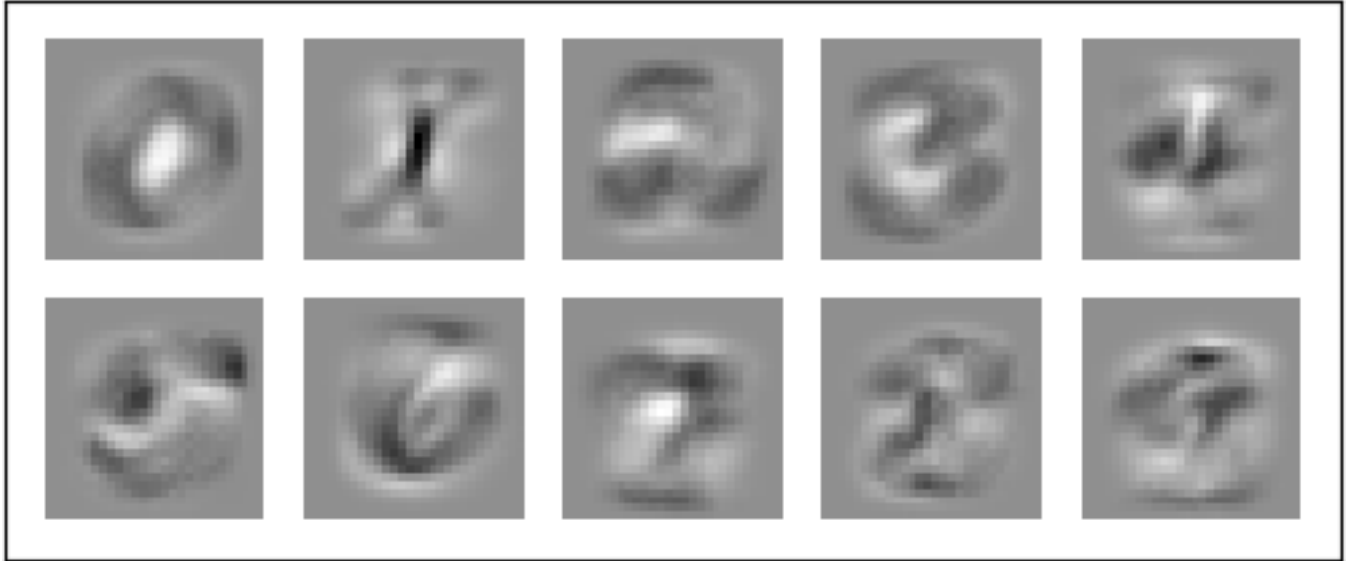


Figure 4:  $\theta_{cd}$  parameters for each class.

(d) Training average predictive log-likelihood per datapoint: -0.360213  
Test average predictive log-likelihood per datapoint: -0.341429  
Training predictive accuracy: 90.0167%  
Test predictive accuracy: 90.89%



**Problem 4** (Unsupervised Learning, 10 points)

Another advantage of generative models is that they can be trained in an unsupervised or semi-supervised manner. In this question, we'll fit the Naïve Bayes model without using labels. Since we don't observe labels, we now have a *latent variable model*. The probability of an image under this model is given by the marginal likelihood, integrating over  $c$ :

$$p(\mathbf{x}|\theta, \pi) = \sum_{c=1}^k p(\mathbf{x}, c|\theta, \pi) = \sum_{c=1}^k p(c|\pi) \prod_{d=1}^{784} p(x_d|c, \theta_{cd}) = \sum_{c=1}^k \text{Cat}(c|\pi) \prod_{d=1}^{784} \text{Ber}(x_d|\theta_{cd}) \quad (5)$$

It turns out that this gives us a mixture model! This model is sometimes called a “mixture of Bernoullis”, although it would be clearer to say “mixture of products of Bernoullis”. Again, this is just the same Naïve Bayes model as before, but where we haven't observed the class labels  $c$ . In fact, we are free to choose  $K$ , the number of mixture components.

- (a) Given  $K$ , how many parameters does this model have?
- (b) Derive the gradient of the log marginal likelihood with respect to  $\theta$ :  $\nabla_{\theta} \log p(\mathbf{x}|\theta, \pi)$
- (c) For a fixed  $\pi_c = \frac{1}{K}$  and  $K = 30$ , fit  $\theta$  on the training set using gradient based optimization. Note: you can't initialize at all zeros – you need to break symmetry somehow, which is done for you in starter code. Starter code reduces this problem to correctly coding the optimization objective. Plot the learned  $\theta$ . How do these cluster means compare to the supervised model?
- (d) For 20 images from the training set, plot the top half the image concatenated with the marginal distribution over each pixel in the bottom half. Hint: You can re-use the formula for  $p(\mathbf{x}_{i \in \text{bottom}}|\mathbf{x}_{\text{top}}, \theta, \pi)$  from before. How do these compare with the image completions from the supervised model?

(a) 784K

(b)

$$\begin{aligned}
\frac{\partial \log(p(\mathbf{x}|\boldsymbol{\theta}, \pi))}{\partial \theta_{ij}} &= \frac{\partial \log(p(\mathbf{x}|\boldsymbol{\theta}, \pi))}{\partial p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \frac{\partial p(\mathbf{x}|\boldsymbol{\theta}, \pi)}{\partial \theta_{ij}} \\
&= \frac{1}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \frac{\partial}{\partial \theta_{ij}} \sum_{c=1}^K p(\mathbf{x}, c|\boldsymbol{\theta}, \pi) \\
&= \frac{1}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \sum_{c=1}^K \frac{\partial}{\partial \theta_{ij}} p(\mathbf{x}, c|\boldsymbol{\theta}, \pi) \\
&= \frac{1}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \sum_{c=1}^K \frac{\partial}{\partial \theta_{ij}} \left( \text{Cat}(c|\pi) \prod_{d=1}^{784} \text{Ber}(x_d|\theta_{cd}) \right) \\
&= \frac{1}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \sum_{c=1}^K \text{Cat}(c|\pi) \prod_{d=1}^{784} \text{Ber}(x_d|\theta_{cd}) \frac{1}{\text{Ber}(x_j|\theta_{cj})} \frac{\partial}{\partial \theta_{ij}} (\text{Ber}(x_j|\theta_{cj})) \\
&\quad \text{Since } \frac{\partial}{\partial \theta_{ij}} (\text{Ber}(x_j|\theta_{cj})) = 0 \text{ when } c \neq i, \\
&= \frac{1}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \text{Cat}(i|\pi) \prod_{d=1}^{784} \text{Ber}(x_d|\theta_{id}) \frac{1}{\text{Ber}(x_j|\theta_{ij})} \frac{\partial}{\partial \theta_{ij}} (\text{Ber}(x_j|\theta_{ij})) \\
&= \frac{p(\mathbf{x}, i|\boldsymbol{\theta}, \pi)}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \frac{1}{\text{Ber}(x_j|\theta_{ij})} \frac{\partial}{\partial \theta_{ij}} (\text{Ber}(x_j|\theta_{ij})) \\
&= \frac{p(\mathbf{x}, i|\boldsymbol{\theta}, \pi)}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \frac{1}{\theta_{ij}^{x_j} (1 - \theta_{ij})^{(1-x_j)}} \left( x_j \theta_{ij}^{(x_j-1)} (1 - \theta_{ij})^{(1-x_j)} - (1 - x_j) (1 - \theta_{ij})^{(1-x_j-1)} \theta_{ij}^{x_j} \right) \\
&= \frac{p(\mathbf{x}, i|\boldsymbol{\theta}, \pi)}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \frac{\theta_{ij}^{(x_j-1)} (1 - \theta_{ij})^{(1-x_j-1)}}{\theta_{ij}^{x_j} (1 - \theta_{ij})^{(1-x_j)}} (x_j (1 - \theta_{ij}) - (1 - x_j) \theta_{ij}) \\
&= \frac{p(\mathbf{x}, i|\boldsymbol{\theta}, \pi)}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \frac{x_j - \theta_{ij}}{\theta_{ij} (1 - \theta_{ij})}
\end{aligned}$$

In matrix form, where  $p(\mathbf{x}, \mathbf{c}|\boldsymbol{\theta}, \pi)$  is a vector of length  $K$  composed of the values  $p(\mathbf{x}, c_i|\boldsymbol{\theta}, \pi)$  and  $\odot$  and  $\oslash$  denote elementwise multiplication and division, respectively,

$$\nabla_w \log(p(\mathbf{x}|\boldsymbol{\theta}, \pi)) = \frac{1}{p(\mathbf{x}|\boldsymbol{\theta}, \pi)} \left( p(\mathbf{x}, \mathbf{c}|\boldsymbol{\theta}, \pi) \mathbf{x}^\top - \text{diag}(p(\mathbf{x}, \mathbf{c}|\boldsymbol{\theta}, \pi)) \boldsymbol{\theta} \right) \oslash (\boldsymbol{\theta} \odot (1 - \boldsymbol{\theta}))$$

- (c) Figure 5 shows the fit  $\boldsymbol{\theta}$  parameters for each  $k$ . Compared to the results from the supervised model shown in Figure 1, these cluster means are less representative of digits in general and focus on more specific features. For instance, there are multiple clusters that resemble the digit 1, but they each look for specific types of 1's that people might write, including different stroke lengths and angles.
- (d) Figure 6 shows the top half of 20 images in the training set concatenated with the marginal distribution over each pixel in the bottom half. Compared to the results from the supervised model shown in Figure 3, these reconstructions are less convincing. Sometimes, the reconstruction is of an incorrect digit, as with the 4 in row 2 column 5 that was reconstructed in the bottom half as a 5. This did not occur with the reconstructions from the supervised model.

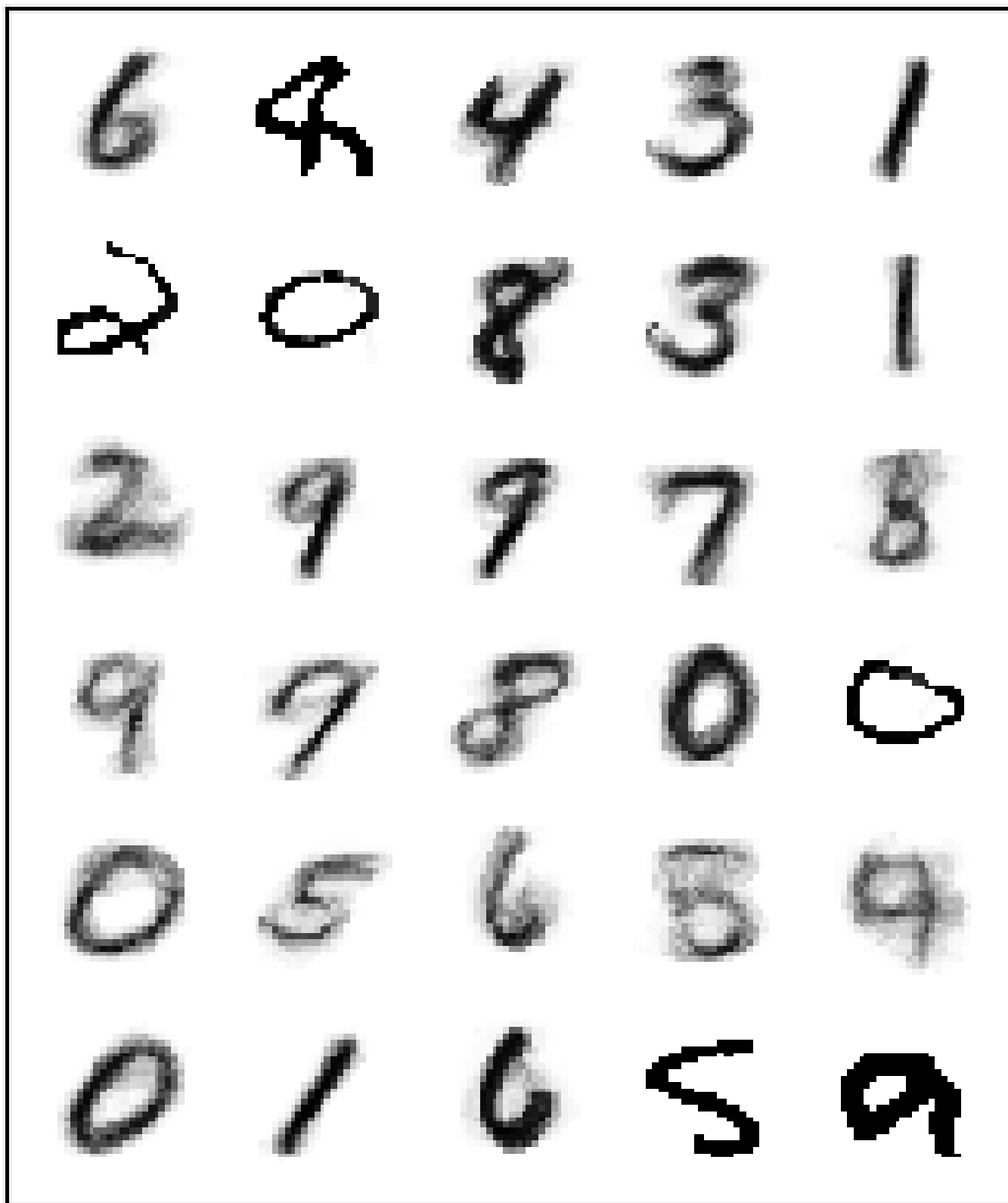


Figure 5:  $\theta$  parameters for each  $k$ .

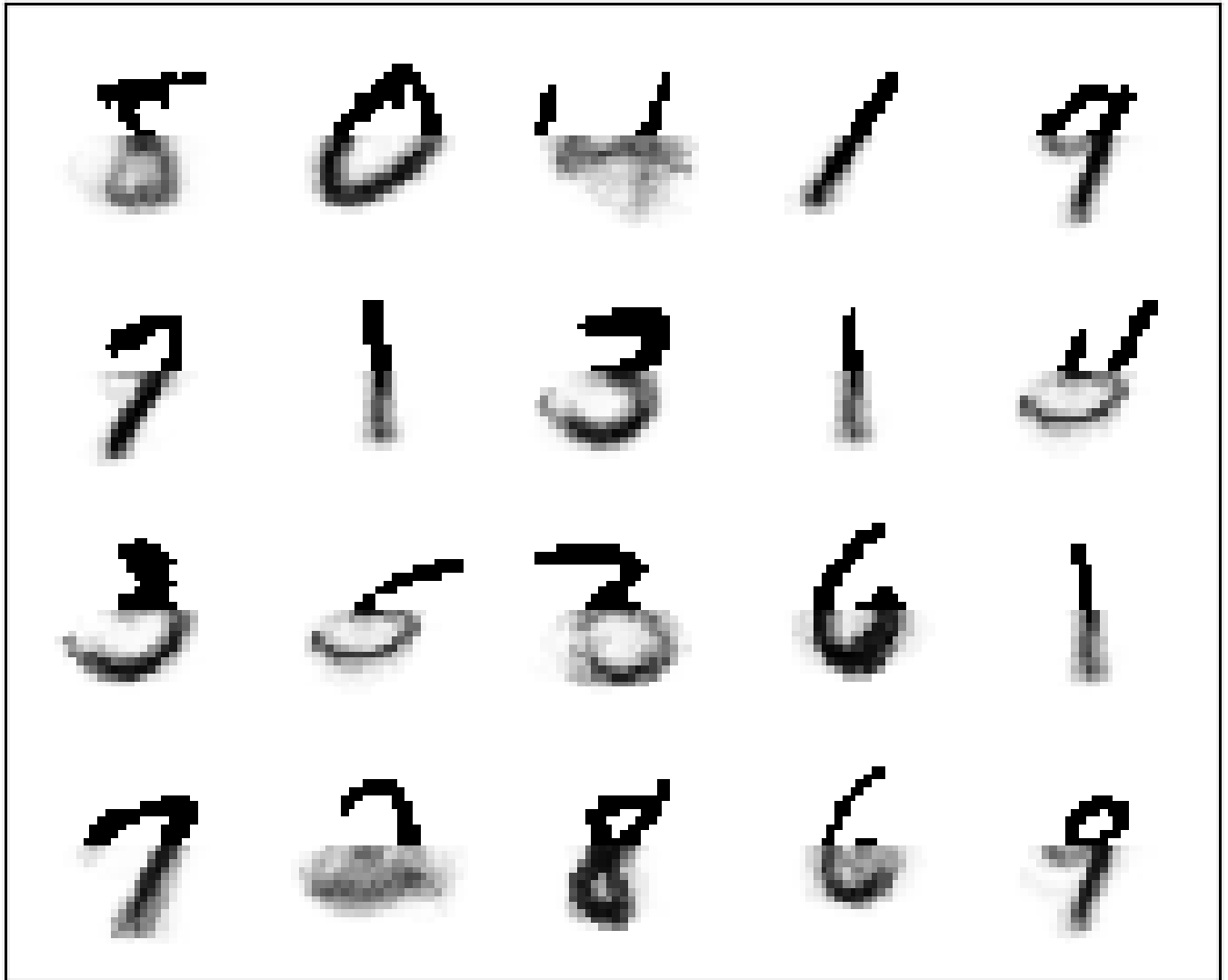


Figure 6: Top half of 20 images in the training set concatenated with the marginal distribution over each pixel in the bottom half.