

0.1 Softmax [20 points]

1) [10 point] Prove that softmax is invariant to constant shifts in the input, i.e., for any input vector x and a constant scalar c , the following holds:

$$\text{softmax}(x) = \text{softmax}(x + c),$$

where $\text{softmax}(x)$ and $x + c$ means adding c to every dimension of x .

2) [10 point] Let $z = Wx + c$, where W and c are some matrix and vector, respectively. $\frac{\partial J}{\partial W}$.

Let

$$J = \log$$

① Given $\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_{i'} e^{x_{i'}}}$

$$\text{softmax}(x+c)_i = \frac{e^{x_i+c}}{\sum_{i'} e^{x_{i'}+c}} = \frac{e^{x_i} e^c}{\sum_{i'} e^{x_{i'}} e^c} = \frac{e^{x_i} \cancel{e^c}}{\cancel{e^c} \sum_{i'} e^{x_{i'}}}$$

Calculate the derivatives of J w.r.t. W and c , respectively, i.e., calculate $\frac{\partial J}{\partial W}$ and

$$= \frac{e^{x_i}}{\sum_{i'} e^{x_{i'}}} \equiv \text{softmax}(x)$$

② $J = \sum_i \log \text{softmax}(z)_i$

if $\text{softmax}(x)_i \triangleq \frac{e^{x_i}}{\sum_{i'} e^{x_{i'}}}$ then,

$$\text{softmax}(z)_i \triangleq \frac{e^{z_i}}{\sum_{i'} e^{z_{i'}}} \text{ and,}$$

$$J = \sum_i \log \frac{e^{z_i}}{\sum_{i'} e^{z_{i'}}} = \log \frac{e^{z_i}}{e^{z_{i'}}} = \log \frac{e^{z_i}}{e^{z_{i'}}} \quad \boxed{j=i'}$$

$$\frac{dJ}{dw} = \frac{d}{dw} \left(\log \frac{e^{z_i}}{j^2 e^z} \right) = 0 \text{ since all terms are constants}$$

then $j = i'$

$$\frac{dJ}{dc} = \frac{d}{dc} \left(\log \frac{e^{z_i}}{j^2 e^z} \right) = 0 \text{ since all terms are constants}$$

0.2 Logistic Regression with Regularization [20 points]

- 1) [10 point] Let the data be (X_i, Y_i) , where $X_i \in \mathbb{R}^d$ and $Y_i \in \{0, 1\}$. Logistic regression is a binary classification model, with the probability of Y_i being 1 as:

$$p(Y_i = 1; X_i) = \frac{1}{1 + e^{-\theta^T X_i}},$$

where θ is the model parameter. Assume we impose an L_2 regularization term on the parameter, defined as:

$$\mathcal{R}(\theta) = \frac{\lambda}{2} \theta^T \theta$$

$$\textcircled{1} p(y_i = 1, x_i, \theta) = \sigma(\theta^T x_i) \triangleq \frac{1}{1 + e^{-\theta^T x_i}}$$

$$\mathcal{R}(\theta) = \frac{\lambda}{2} \theta^T \theta$$

$$J(\theta) = \frac{1}{n} \left[Y^T \log(p(x, \theta)) + (1 - Y)^T \log(1 - p(x, \theta)) \right]$$

$$+ \frac{\lambda}{2} \theta^T \theta$$

L_2 Regularization

with a positive constant λ . Write out the final objective function for this logistic regression with regularization model.

$$\textcircled{2} \frac{dJ(\theta)}{d\theta} = \frac{d}{d\theta} \frac{1}{n} \sum_{i=1}^n \left[Y_i \log(\sigma(x_i)) + (1 - Y_i) \log(1 - \sigma(x_i)) \right]$$

- 2) [10 point] If we use gradient descent to solve the model parameter. Derive the

updating rule for θ . Your answer should contain the derivation, not just the final answer.

$$= \frac{1}{n} \sum_{i=1}^n \left[Y_i \frac{1}{\sigma(\theta^T x_i)} \sigma(\theta^T x_i) (1 - \sigma(\theta^T x_i)) \cdot \frac{d}{d\theta} (\theta^T x_i) + \right. \\ \left. - \frac{1 - Y_i}{1 - \sigma(\theta^T x_i)} (1 - \sigma(\theta^T x_i)) \sigma(\theta^T x_i) \frac{d}{d\theta} (\theta^T x_i) \right]$$

$$= \frac{1}{n} \sum x_i [y_i - y_i \sigma(\theta^T x_i) - \sigma(\theta^T x_i) + y_i \sigma(\theta^T x_i)]$$

$$= \frac{1}{n} \sum x_i [y_i - \sigma(\theta^T x_i)]$$

0.3 Derivative of the Softmax Function [30 points]

1) [10 point] Define the loss function as

$$J(\mathbf{z}) = - \sum_{k=1}^K y_k \log \tilde{y}_k ,$$

where $\mathbf{z} = (z_1, \dots, z_K)$ is a known probability vector. Derive the $\frac{dJ(\mathbf{z})}{d\mathbf{z}}$

Note $\mathbf{z} = (z_1, \dots, z_K)$ is a vector so $\frac{dJ(\mathbf{z})}{d\mathbf{z}}$ is in the form of a vector. Your answer should contain the derivation, not just the final answer.

2 [10 point] Assume the above softmax is the output layer of an FNN. Briefly explain how $\frac{\partial J}{\partial \mathbf{W}}$ the derivative is used in the backpropagation algorithm.

① $J(\mathbf{z}) = - \sum_{k=1}^K y_k \log \tilde{y}_k$ $\tilde{y}_k = \frac{e^{z_k}}{\sum_{k'=1}^K e^{z_{k'}}}$ $\frac{dJ(\mathbf{z})}{d\mathbf{z}} = \begin{bmatrix} dJ/dz_1 \\ \vdots \\ dJ/dz_K \end{bmatrix}$

so

$$\frac{dJ(\mathbf{z})}{d\mathbf{z}} = - \frac{d}{d\mathbf{z}} \sum_{k=1}^K y_k \log \tilde{y}_k$$

$$= - \frac{d}{d\mathbf{z}} \sum_{k=1}^K y_k \log \frac{e^{z_k}}{\sum_{k'=1}^K e^{z_{k'}}}$$

3) [10 points] Let $\mathbf{z} = \mathbf{W}\mathbf{h} + \mathbf{b}$, where $\sigma(\bullet)$ is the sigmoid function, \mathbf{W} is a matrix, \mathbf{b} and \mathbf{h} are

$$= - \frac{d}{d\mathbf{z}} \sum_{k=1}^K y_k z_k + \frac{d}{d\mathbf{z}} \sum_{k=1}^K y_k \log \sum_{k'=1}^K e^{z_{k'}}$$

vectors. Use the chain rule to calculate the gradient of \mathbf{W} and \mathbf{b} , i.e., $\frac{\partial J}{\partial \mathbf{W}}$ and $\frac{\partial J}{\partial \mathbf{b}}$, respectively (it is

$$= -y_k + \frac{d}{d\mathbf{z}} \sum_{k=1}^K y_k \log \sum_{k'=1}^K e^{z_{k'}}$$

enough to derive the gradients for one element of the matrix/vector parameter).

$$= -y_k + \sum_{k=1}^K y_k \cdot \frac{d}{d\mathbf{z}} \log \sum_{k'=1}^K e^{z_{k'}}$$

$$= -y_k + \sum_{k=1}^K y_k \cdot \frac{1}{\sum_{k'=1}^K e^{z_{k'}}} \cdot e^{z_k}$$

$$= -y_k + \frac{e^{z_k}}{\sum_{k'=1}^K e^{z_{k'}}} \left(\sum_{k=1}^K y_k \right)$$

$$= -y_k + \tilde{y}_k \quad \text{this means} \quad \frac{dJ(\tilde{z})}{dz} = \begin{bmatrix} -y_1 + \tilde{y}_1 \\ -y_2 + \tilde{y}_2 \\ \vdots \\ -y_k + \tilde{y}_k \end{bmatrix}$$

② The derivative is used in backpropagation to update the weights of a FNN during training to help minimize loss.

③ $z = \sigma(w^T h + b)$

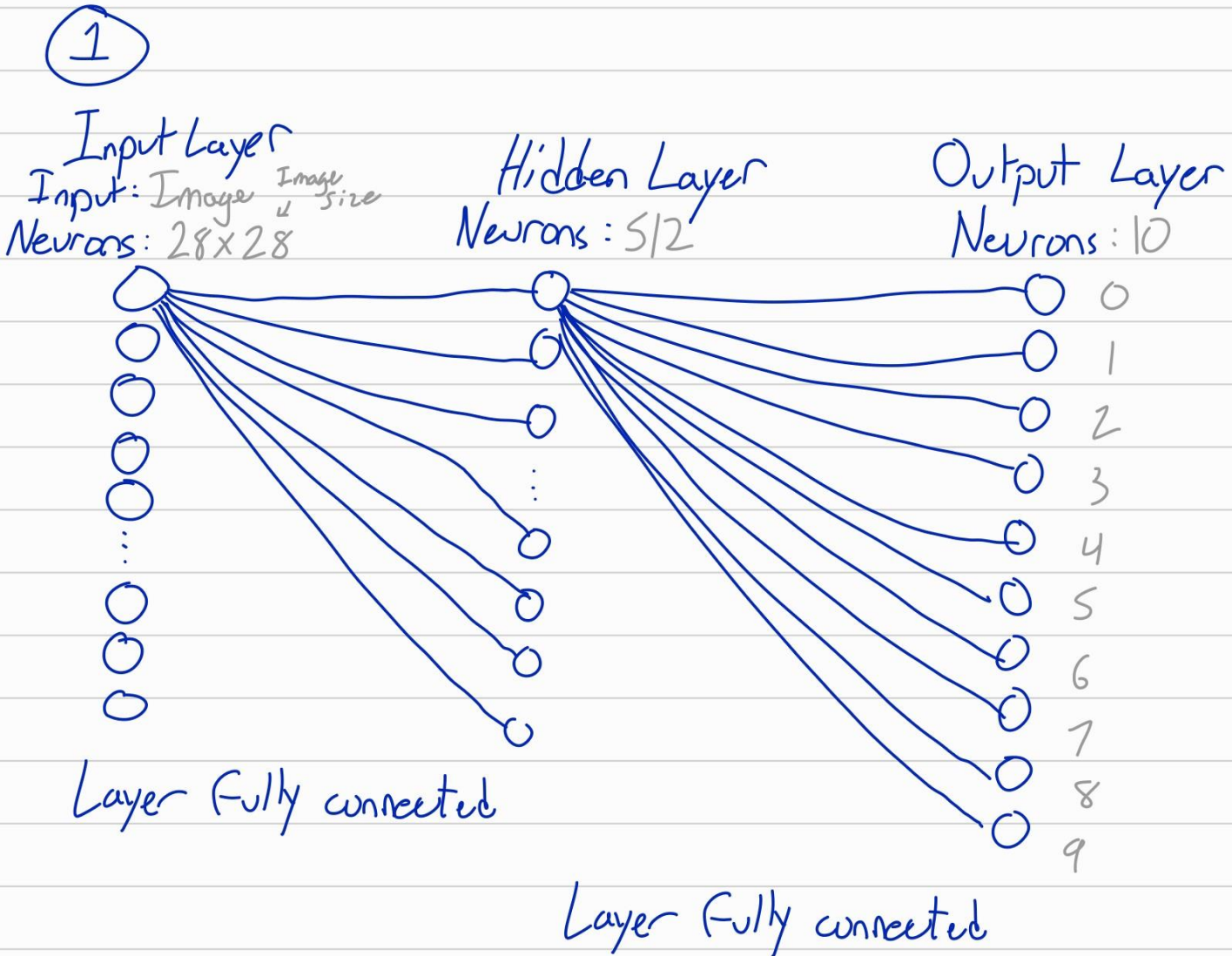
$$\frac{dJ}{dw} = ? \quad \frac{dJ}{db} = ?$$

$$\frac{dJ}{dw} = \frac{dJ}{dz} \cdot \frac{dz}{dw} \quad \frac{dJ}{dz} = \begin{bmatrix} -y_1 + \tilde{y}_1 \\ -y_2 + \tilde{y}_2 \\ \vdots \\ -y_k + \tilde{y}_k \end{bmatrix}$$

0.4 MNIST with F NN [30 points]

- 1) [10 points] Design an FNN for MNIST classification. Draw the computational graph of your model.
- 2) [20 points] Implement the model and plot two curves in one figure: i) training loss vs. training iterations; ii) test loss vs. training iterations.

- You can use any packages, but PyTorch is recommended.



References: Code adapted from [https://towardsdatascience.com/handwritten-digit-mnist-pytorch-](https://towardsdatascience.com/handwritten-digit-mnist-pytorch-977b5338e627)

[977b5338e627](https://towardsdatascience.com/handwritten-digit-mnist-pytorch-977b5338e627)

