

Time: 45 minutes

Max marks: 10 8

1. (2+1=3 points) The sigmoid function is given as $\sigma(z) = \frac{1}{1+e^{-z}}$. Logistic regression can be implemented by a single neuron neural network that has a sigmoid activation function. (i) Compute the gradient of the output of this neural network with respect to the weight parameters $w_i, i = 0, 1, 2, \dots, d$, where w_0 is the bias term and the input is \mathbf{x} with a 1 appended. The *logit* is computed as $z = \mathbf{w}^\top \mathbf{x}$. (ii) Write this derivative expression in terms of $\sigma(z)$.

$$\mathbf{x} \in \mathbb{R}^d \quad P(y=1|\mathbf{x}) = a = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}$$

$$\mathbf{w}^\top \mathbf{x} = w_0 + w_1 x_1 + \dots + \dots + w_d x_d$$

$$\begin{aligned} \text{(i)} \quad \frac{\partial a}{\partial w_i} &= \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_i} \\ \frac{\partial a}{\partial z} &= \frac{\partial}{\partial z} \left(\frac{1}{1+e^{-z}} \right) = \frac{-1}{(1+e^{-z})^2} \cdot \frac{\partial}{\partial z} (1+e^{-z}) = \frac{-(-e^{-z})}{(1+e^{-z})^2} \\ &= \frac{e^{-z}}{1+e^{-z}} \cdot \frac{1}{1+e^{-z}} = \left(\frac{1+e^{-z}-1}{1+e^{-z}} \right) \cdot \left(\frac{1}{1+e^{-z}} \right) \\ &= \left(1 - \frac{1}{1+e^{-z}} \right) \left(\frac{1}{1+e^{-z}} \right) = (1 - \sigma(z)) \sigma(z) \end{aligned}$$

$$\frac{\partial z}{\partial w_i} = \frac{\partial}{\partial w_i} (w_0 + w_1 x_1 + w_2 x_2 + \dots + w_i x_i + \dots + w_d x_d)$$

\therefore The gradient vector

$$\nabla_w a = (1 - \sigma(\mathbf{w}^\top \mathbf{x})) \sigma(\mathbf{w}^\top \mathbf{x}) \cdot$$

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_d \end{bmatrix}$$

(ii) See derivation above,

$$\sigma(\mathbf{w}^\top \mathbf{x}) (1 - \sigma(\mathbf{w}^\top \mathbf{x})) \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_d \end{bmatrix}$$

2. (1+1=2 points) One of my new students has trained a binary classifier and is achieving a (i) training accuracy of around 95% and a validation accuracy of around 64% (ii) Exactly 99.9% accuracy on both training and validation datasets. The classifier was trained for a **weapon / no-weapon** classification at a building security screening facility, and the class distribution was 999000 samples of **no-weapon** category and 1000 samples of the **weapon** category in the training set and 49950 samples of **no-weapon** category and 50 samples of the **weapon** category in the validation set. I claim that the classification model isn't trained correctly in both cases. Can you suggest to my student what might be the problems in the two cases above?

- (i) The model seems to be overfitting the training set.
More diverse training samples or regularization may help improve the model performance
- (ii) It appears that the classifier might be trivially classifying all samples as the "no-weapon" category. It may help to look at the False Negatives (FN) or the confusion matrix instead of the accuracy metric. Reporting the Precision & Recall or F1 score may also be more useful.

3. (2+1=3 points) I have a CNN model that works with RGB images of size 512×512 . The first layer uses 5×5 kernels with a stride of 2 and generates 100 feature maps, and the second layer uses a 3×3 kernel with a stride of 4 and generates 200 feature maps. The activations in both layers are sigmoid. (i) What are the number of learnable parameters in layer 1 and layer 2? (ii) Would the number of parameters change if we change the stride? Would the feature map size change if we change the stride?

(i) Layer 1 :
 $\# \text{ input channels} = 3$, $\# \text{ params for each kernel} = 3 \times 5 \times 5 + 1$
 $\# \text{ bias}$

$$\# \text{ kernels} = \# \text{ of feature maps} = 100$$

$$\# \text{ learnable parameters} = (3 \times 25 + 1) \times 100 = 7600$$

Layer 2 :
 $\# \text{ input channels} = 100$, $\# \text{ params for each kernel} = 100 \times 3 \times 3 + 1$
 $\# \text{ bias}$ $= 901$

$$\# \text{ kernels} = \# \text{ of feature maps} = 200$$

$$\# \text{ learnable params} = 901 \times 200 = 180200$$

(ii) No, the # parameters won't change with change in stride

Yes, the size of feature maps will change with a change in the stride.

4. (Extra credit - 2 points) Let x be a discrete random variable that can take integer values 0 through 9. Let $P(x)$ and $Q(x)$ be two distributions over x . The self-entropy is given as $H(P) = -\mathbf{E}_{x \sim P} \log P(x)$ and the cross-entropy is given as $H(P, Q) = -\mathbf{E}_{x \sim P} \log Q(x)$, where $\mathbf{E}[\cdot]$ is the expectation operator. Given that the KL Divergence is defined as $D_{KL}(P||Q) = -\sum_x P(x) \log \frac{P(x)}{Q(x)}$ ¹, show that

$$H(P, Q) = H(P) + D_{KL}(P||Q)$$

$$D_{KL}(P||Q) = \sum_n P(n) [\log P(n) - \log Q(n)]$$

Adding $H(P)$ on both sides, we get

$$\begin{aligned} D_{KL}(P||Q) + H(P) &= \sum_n P(n) [\log P(n) - \log Q(n)] \\ &\quad - \sum_n P(n) \log P(n) \\ &= \cancel{\sum_n P(n) \log P(n)} - \sum_n P(n) \log Q(n) \\ &\quad - \cancel{\sum_n P(n) \log P(n)} \\ &= - \sum_n P(n) \log (Q(x)) \\ &= - \mathbf{E}_{x \sim P} \log (Q(x)) = H(P, Q) \end{aligned}$$

¹The question paper had a mistake with the $D_{KL}(P||Q)$ defined as $D_{KL}(P||Q) = -\sum_x P(x) \frac{\log P(x)}{\log Q(x)}$. This question will be treated as extra-credit