

Birla Institute of Technology & Science - Pilani, Hyderabad Campus
First Semester 2020-21

CS F441: Selected Topics from Computer Science (Deep Learning)
Comprehensive Examination

Type: Open (Online)

Time: 120 mins

Max Marks: 68

Date: 21.12.2020

All parts of the same question should be answered together.

1. Suppose you are asked to solve the binary classification task of classifying images as dog vs. non-dog. Suppose you design a CNN with a single output neuron. Let the output of this neuron be z . The final output of your network, y is given by:

$$y = \text{sigmoid}(\text{ReLU}(z))$$

You would like classify all inputs with a final value $y \geq 0.5$ as dog images. What major problem are you going to encounter?

[8 Marks]

Sol:

we show that the network always outputs the image as 'dog' irrespective of the value z . Equivalently we show that $y \geq 0.5$ irrespective of the input image if the final output of the network is given by $y = \text{sigmoid}(\text{ReLU}(z))$.

Case (i): $z \geq 0$

$$z \geq 0 \Rightarrow e^z \geq e^0 \Rightarrow e^z \geq 1 \Rightarrow 2e^z \geq 1 + e^z$$

$$\Rightarrow \frac{e^z}{1+e^z} \geq \frac{1}{2} \Rightarrow \frac{1}{1+e^{-z}} \geq \frac{1}{2}$$

$$\Rightarrow \text{sigmoid}(\text{ReLU}(z)) \geq 0.5 \Rightarrow y \geq 0.5$$

Case (ii): $z < 0$

$$z < 0 \Rightarrow \text{ReLU}(z) = 0$$

$$\text{sigmoid}(\text{ReLU}(z)) = \text{sigmoid}(0)$$

$$\Rightarrow y = \frac{1}{1+e^0} \Rightarrow y = \frac{1}{1+1}$$

$$\Rightarrow y = 0.5$$

Thus we proved that irrespective of value of z , $y \geq 0.5$. And hence the output of the network is always a 'dog'.

2. Suppose you are given 'N' training examples. Consider a univariate regression $y = wx$ where $w \in \mathbb{R}$, and $x \in \mathbb{R}^{1 \times m}$. Let the classical error function i.e., $L(w) = \text{sum of squares of errors} / m$. Find the first derivative of L with respect to w , i.e., dL/dw . The final expression should be as compact as possible.

[4 Marks]

Sol:

Let $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be N training examples.

$$x_n = (x_{n1}, x_{n2}, \dots, x_{nm})$$

The regression model to fit is

$$y = wx, \quad w \in \mathbb{R}^{1 \times m}$$

It is evident that $y \in \mathbb{R}_T$

$$y_n = (y_{n1}, y_{n2}, \dots, y_{nm})$$

The squared error due to n^{th} training example is

$$\sum_{i=1}^m (wx_{ni} - y_i)^2$$

$$\text{Hence } L(w) = \frac{1}{N} \sum_{n=1}^N \left(\sum_{i=1}^m (wx_{ni} - y_i)^2 \right)$$

$$\frac{\partial L}{\partial w} = \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^m 2 (wx_{ni} - y_i) (x_{ni})$$

$$= \frac{2}{N} \sum_{n=1}^N \sum_{i=1}^m (wx_{ni} - y_i) (x_{ni})$$

$$= \frac{2}{N} \sum_{n=1}^N x_n^T (wx_n - y_n)$$

$$= \frac{2}{N} \sum_{n=1}^N x_n^T (\hat{y}_n - y_n) \quad \left(\hat{y}_n \text{ is the predicted outcome of } x_n \right).$$

3. Explain the significance of tensors in deep learning and provide a case study in which its usage is maximally exploited?

[4 Marks]

Sol: Tensors are multidimensional arrays in Deep Learning that are used to represent data. They represent the data with higher dimensions. Due to the high-level nature of the programming languages, the syntax of tensors are easily understood and broadly used.

Images data – $(i, j, (R, G, B))$

4. Suppose $y = (y_1, y_2, \dots, y_n)$ be the outcome of a layer after tanh is applied on (z_1, z_2, \dots, z_n) . i.e., $(y_1, y_2, \dots, y_n) = (\tanh z_1, \tanh z_2, \dots, \tanh z_n)$. Write down the Jacobian matrix dy/dz . You may give your answer in terms of $\tanh'(z)$, the univariate derivative of the tanh function.

[6 Marks]

Sol:

$$\begin{aligned}
 (y_1, y_2, \dots, y_n) &= (\tanh(z_1), \tanh(z_2), \dots, \tanh(z_n)) \\
 \text{Jacobian } \frac{dy}{dz} &= \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \frac{\partial y_1}{\partial z_2} & \dots & \frac{\partial y_1}{\partial z_n} \\ \frac{\partial y_2}{\partial z_1} & \frac{\partial y_2}{\partial z_2} & \dots & \frac{\partial y_2}{\partial z_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial z_1} & \frac{\partial y_n}{\partial z_2} & \dots & \frac{\partial y_n}{\partial z_n} \end{bmatrix} \\
 &= \begin{bmatrix} \tanh'(z_1) & 0 & \dots & 0 \\ 0 & \tanh'(z_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \tanh'(z_n) \end{bmatrix}
 \end{aligned}$$

5. One of the difficulties with the logistic activation function is that of saturated units. Briefly explain the problem, and whether switching to tanh fixes the problem. (You may refer to your answer from the above question or sketch the activation functions.) [6 Marks]

Sol: No, switching to tanh does not fix the problem. The derivative of sigmoid(z) is small for large negative or positive z . The same problem persists in $\tanh(z)$. Both function has a sigmoidal shape. We can see tanh is effectively a scaled and translated sigmoid function: $\tanh(z) = 2 \text{ sigmoid}(2z) - 1$.

6. Give an example of a data augmentation technique that would be useful for classifying images of cats vs. dogs, but not for classifying handwritten digits. Briefly explain your answer. [4 Marks]

Sol: Flipping the image horizontally; doing this to a dog image would result in a plausible dog image, but not so for an image of a digit.

7. Suppose that you have a model that provides around 80% accuracy on the training as well as on the out-of-sample data test data. Would you recommend increasing the amount of data or adjusting the model to improve accuracy? Please provide suitable reasoning in support of your answer. [4 Marks]

Sol: As the error in the training data and testing data (unseen data) is more or less the same the chosen model is probably not overfitting. So to improve the accuracy, adjusting the model (for example, in the case of polynomial regression, changing the degree of the polynomial etc.) might improve the accuracy.

8. We have seen that multilayer perceptrons are universal for the set of functions mapping binary-valued input vectors to binary valued outputs. [4 Marks]

(a). What do we mean by universal?

Sol: For any function in the set, there is some network that computes it.

(b). If multilayer perceptrons are universal, why do we still consider other architectures?

Sol: The function might not be compactly representable, i.e. it might require exponentially many units. This means it may be prohibitively expensive to compute, and training such a large architecture wouldn't generalize to new data.

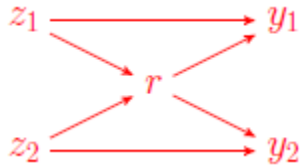
9. We know that the softmax function takes a vector (z_1, z_2, \dots, z_M) and returns a vector (y_1, y_2, \dots, y_M). We can express it in the following term:

$$r = (e^{z_1} + e^{z_2} + \dots + e^{z_M}) \text{ and } y_i = e^{z_i} / r$$

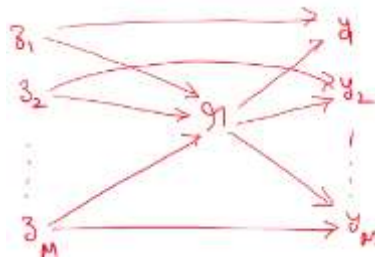
Suppose $M = 2$, i.e., just two inputs and outputs to the softmax. Draw the computation graph relating z_1, z_2, r, y_1, y_2 .

For any generic M , draw the computation graph relating z_1, z_2, r, y_1, y_2 . [6 Marks]

Sol: Consider $M = 2$, i.e. just two inputs and outputs to the softmax. Draw the computation graph relating z_1, z_2, r, y_1 , and y_2 .



For a generic M , the computational graph is as shown below:



10. Can you think of top two hurdles that computer vision algorithms face in dealing with images? i.e. two characteristics of image formation that make it difficult to recover the image content.

[4 Marks]

Sol: There are several examples, such as: (1) they can be from various poses and perspectives, (2) they may vary in color and brightness, (3) there might be occlusion or ambiguity of objects

11. Consider a 1-dimensional time-series with values 2, 1, 3, 4, 7. Perform a convolution with a 1-dimensional filter 1, 0, 1. Also, for a one-dimensional time series of length L and a filter of size F , what is the length of the output? [4 Marks]

Sol: If a 1-dimensional filter 1, 0, 1 is applied on the 1-dimensional time-series with values 2, 1, 3, 4, 7 the the output is 5, 5, 10.

If a filter of size F is applied on a one-dimensional time series of length L then the length of the output is $L - F + 1$.

12. If a data block in a convolutional network has dimension $H \times W \times D = 300 \times 300 \times 256$, and we apply a convolutional filter (or kernel) to it of dimensions $H_F \times W_F \times D = 11 \times 11 \times 256$, what is the dimension of the output data block? How many weight parameters are there in this simple operation with and without bias? [6 points]

Sol: The size of the output data block is $290 \times 290 \times 1$.

The weight parameters $(11 \times 11 \times 256)$ without bias parameter and $(11 \times 11 \times 256) + 1$ with bias parameter.

13. Can you mention couple of reasons to support the argument that convolutional layers more commonly used than fully-connected layers for image processing? [4 Marks]

Sol: Convolutional layers are better suited for processing images because: (a) they assume that nearby pixels share similar structure, (b) the parameters for a single filter are shared across the entire input image in the sense that a single filter “slides” across an image and produces a new

output with one image channel. This allows the parameters to be invariant to the location of objects in images.

14. Can you identify one option amongst the below mentioned four option to answer the question - When should multi-task learning be used? You should provide appropriate reasoning to support the selected option. [4 Marks]

- (i) When your problem involves more than one class label.
- (ii) When two tasks have the same dataset.
- (iii) When you have a small amount of data for a particular task that would benefit from the large dataset of another task.
- (iv) When the tasks have datasets of different formats (text and images).

Sol: Option (ii) is correct. The prior belief is “among the factors that explain the variations observed in the data associated with the different tasks, some are shared across two or more tasks”.

Option (iii) typically refers to refers to transfer learning problem. However with decent explanation, a partial credit can be given.