

CISC 484/684 — Intro to Machine Learning: Final Practice

True/False

For each question, check either True or False. If False, explain why.

6 points for True, 3 points for False, 3 points for an explanation for False.

1) The Kullback-Leibler divergence can be used as a "similarity" measure between two distributions. A 0 KL value indicates that the two distributions are identical.

True

False

Explanation if False:

2) A perceptron can only classify data that is linearly separable.

True

False

Explanation if False:

3) Adding more layers to a neural network always reduces training error.

True

False

Explanation if False:

4) In a CNN, the convolution operation reduces the spatial dimensions of the input.

True

False

Explanation if False:

5) Dropout works by randomly zeroing out some neurons only during inference to prevent overfitting.

True

False

Explanation if False:

6) In GAN training, the generator's objective is to minimize the discriminator's ability to correctly distinguish real from fake samples.

True

False

Explanation if False:

7) In a standard VAE, the encoder outputs both a mean vector and a variance vector for the latent distribution.

True

False

Explanation if False:

8) Autoencoders are generative models because they learn a latent representation of the input.

True

False

Explanation if False:

9) A perceptron uses the sigmoid activation function to compute its output.

True

False

Explanation if False:

10) In a CNN, increasing the number of filters in each layer always improves test accuracy.

True

False

Explanation if False:

5. At this point, what is the value of the k -means objective function? You do not need to worry about performing arithmetic here. In other words, an answer with a form resembling $(7.3 + 2.1)^3 + (8.3 + 1.2)^4$ is fine.

12) Neural Networks

Consider the neural network architecture shown above for a binary classification problem. The values for the weights are shown in the figure. We define:

$$a_1 = w_{11}x_1 + b_{11}$$

$$a_2 = w_{12}x_1 + b_{12}$$

$$a_3 = w_{21}z_1 + w_{22}z_2 + b_{21}$$

$$z_1 = \text{ReLU}(a_1)$$

$$z_2 = \text{ReLU}(a_2)$$

$$z_3 = \sigma(a_3), \sigma(x) = \frac{1}{1 + e^{-x}}$$

where $\text{ReLU}(x) = \max(0, x)$.

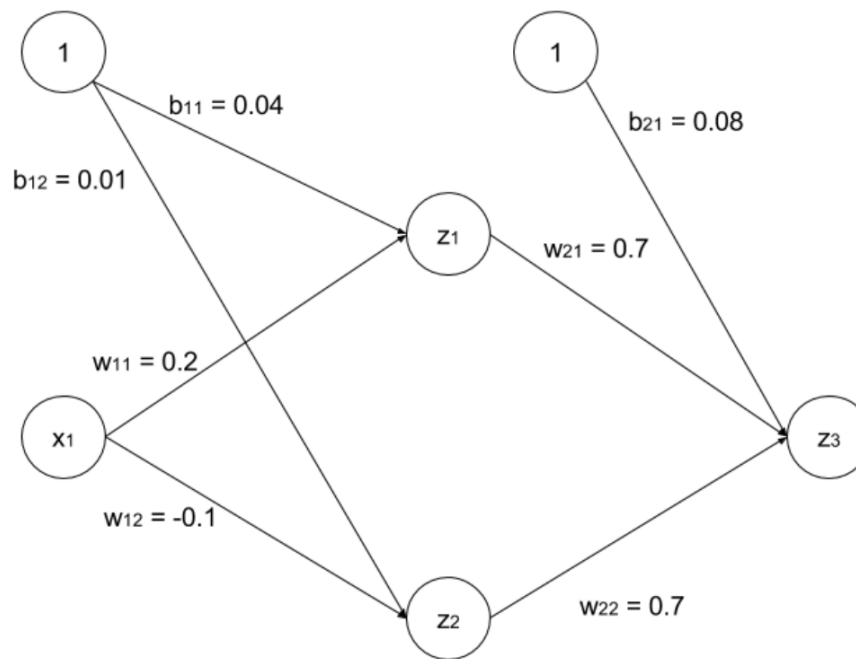


Figure 1: Neural Networks

1. For $x_1 = 0.3$, compute z_3 in terms of e .

2. **Select from 0 or 1:** Which class does the network predict for the data point $x_1 = 0.3$, assuming that $\hat{y} = 1$ if $z_3 > \frac{1}{2}$ and otherwise, $\hat{y} = 0$.

3. Perform backpropagation on the bias term b_{21} by deriving the expression for the gradient of the loss function $L(y, z_3)$ with respect to the bias term b_{21} , $\frac{\partial L}{\partial b_{21}}$.

Express your answer in terms of partial derivatives of the form $\frac{\partial \alpha}{\partial \beta}$, where α and β can be any of $L, z_i, a_i, b_{ij}, w_{ij}, x_1$ for all valid values of i, j . Your backpropagation algorithm should be as explicit as possible - that is, make sure each partial derivative $\frac{\partial \alpha}{\partial \beta}$ cannot be decomposed further into simpler partial derivatives. Do not evaluate the partial derivatives.

4. Perform backpropagation on the bias term b_{12} by deriving the expression for the gradient of the loss function $L(y, z_3)$ with respect to the bias term b_{12} , $\frac{\partial L}{\partial b_{12}}$. Express your answer in terms of partial derivatives of the form $\frac{\partial \alpha}{\partial \beta}$, where α and β can be any of $L, z_i, a_i, b_{ij}, w_{ij}, x_1$ for all valid values of i, j . Your backpropagation algorithm should be as explicit as possible - that is, make sure each partial derivative $\frac{\partial \alpha}{\partial \beta}$ cannot be decomposed further into simpler partial derivatives. Do not evaluate the partial derivatives.

13) Convolutional Neural Networks

Let's begin by considering some of the high-level components of a CNN kernel along with the basic motivation.

1. What is a kernel?
2. Why do we need stride, and what benefits/tradeoffs might different values of stride have on the output?
3. What functionality does padding add to the kernel? Why might we want to use it?
4. Consider the following image, filter, and output shape

$$X = \begin{array}{|c|c|c|c|c|c|} \hline 1 & 0 & -2 & 3 & 4 & 1 \\ \hline 2 & 9 & 5 & 6 & 0 & -1 \\ \hline 0 & -3 & 1 & 3 & 4 & 4 \\ \hline 6 & 5 & 2 & 0 & 6 & 8 \\ \hline -5 & 4 & -3 & 1 & 3 & -2 \\ \hline 4 & 1 & 2 & 8 & 9 & 7 \\ \hline \end{array}
 \quad
 F = \begin{array}{|c|c|c|} \hline -1 & -1 & -1 \\ \hline -1 & 8 & -1 \\ \hline -1 & -1 & -1 \\ \hline \end{array}
 \quad
 Y = \begin{array}{|c|c|c|c|} \hline a & b & c & d \\ \hline e & f & g & h \\ \hline i & j & k & l \\ \hline m & n & o & p \\ \hline \end{array}$$

The shape of this particular Y is that of a kernel using no padding and a stride of 1.

- (a) Suppose we decide that, instead of having our output shape be $(4, 4)$, we want a slightly smaller, $(3, 3)$ image as output for the kernel. In order for this to happen, what is the smallest combination of stride and padding that would work?
- (b) Let's make this a bit more general. Suppose our original image of shape (a, a) , and we want the shape of our final image to be of shape (b, b) , where $b \leq a$. Furthermore, the shape of the filter is (k, k) , the stride length is s , and the padding is p . Express b in terms of all defined variables.