



## Lecture 5 Quiz

[Back to Week 5](#)

**6/6** points earned  
(100%)

Quiz passed!



### Be Recognized for Your Achievements.

"Course Certificates give you the recognition you need to get the job, the material gives you the skills to do the job. It makes you look more valuable because you are more valuable." - Peter B., USA, Software Developer

**Showcase Your Accomplishment! Earn Your Course Certificate! CA\$64 >**



1 / 1  
points

1.

For which of the following tasks can we expect that the problem of "dimension hopping" will occur (given that the data is input correctly)? Check all that apply.



Estimating the risk that a patient will develop heart disease given their age, weight, blood pressure, and cholesterol level.



**Un-selected is correct**



Determining whether a wave has high frequency or low frequency. The input is a set of time values along with their corresponding vertical displacements.



**Correct**

Dimension hopping occurs when one can take the information contained in the dimensions of some input, and move this between dimensions while not changing the target. The canonical example is taking an image of a handwritten digit and translating it within the image. The dimensions that contain "ink" are now different (they have been moved to other dimensions), however the label we assign to the digit has not changed. Note that this is not something that happens consistently across the dataset, that is we may have a dataset containing two handwritten digits where one is a translated version of the other, however this still does not change the corresponding label of the digits.



Determining whether a given image shows a bike or a car. The bike or car might appear anywhere in the image. The input is the whole set of pixels for the image.



**Correct**

Dimension hopping occurs when one can take the information contained in the dimensions of some input, and move this between dimensions while not changing the target. The canonical example is taking an image of a handwritten digit and translating it within the image. The dimensions that contain "ink" are now different (they have been moved to other dimensions), however the label we assign to the digit has not changed. Note that this is not something that happens consistently across the dataset, that is we may have a dataset containing two handwritten digits where one is a translated version of the other, however this still does not change the corresponding label of the digits.



Predicting the next word in a sequence given the previous 3 words.



**Un-selected is correct**



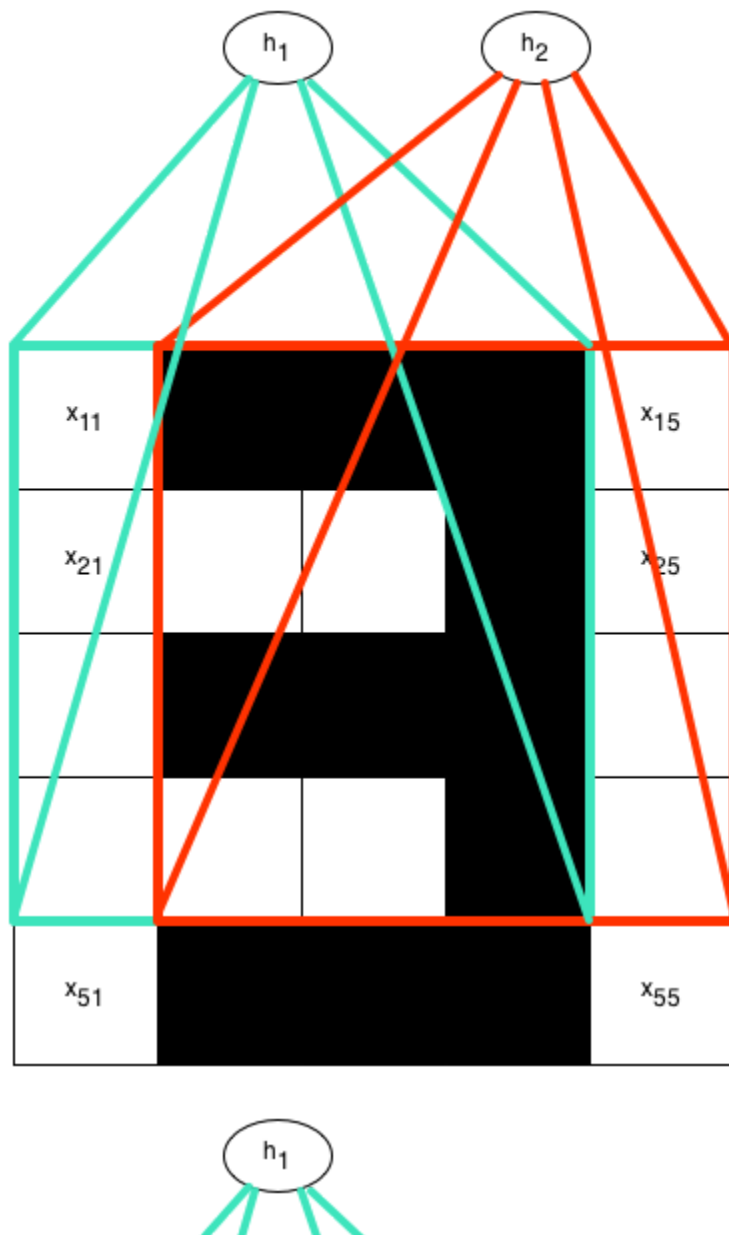
1 / 1  
points

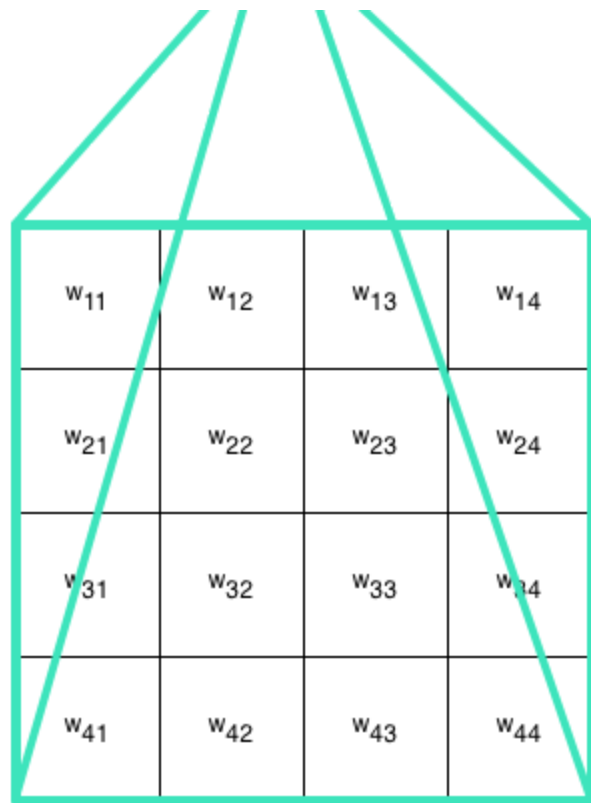
2.

We have a convolutional neural network for images of 5 by 5 pixels. In this network, each hidden unit is connected to a different 4 x 4 region of the input image:

- The first hidden unit,  $h_1$ , is connected to the upper left 4x4 portion of the input image (as shown below).
- The second hidden unit,  $h_2$ , is connected to the upper right 4x4 portion of the input image (as shown below).
- The third hidden unit,  $h_3$ , is connected to the lower left 4x4 portion of the input image (not shown in the diagram).
- The fourth hidden unit,  $h_4$ , is connected to the lower right 4x4 portion of the input image (not shown in the diagram).

Because it's a convolutional network, the weights (connection strengths) are the same for all hidden units: the only difference between the hidden units is that each of them connects to a different part of the input image. In the second diagram, we show the array of weights, which are the same for each of the four hidden units.





For  $h_1$ , weight  $w_{11}$  is connected to the top-left pixel, i.e.  $x_{11}$ , while for hidden unit  $h_2$ , weight  $w_{11}$  connects to the pixel that is one to the right of the top left pixel, i.e.  $x_{12}$ .

Imagine that for some training case, we have an input image where each of the black pixels in the top diagram has value 1, and each of the white ones has value 0. Notice that the image shows a "3" in pixels.

The network has no biases. The weights of the network are given as follows:

$$\begin{array}{llll}
 w_{11} = 1 & w_{12} = 1 & w_{13} = 1 & w_{14} = 0 \\
 w_{21} = 0 & w_{22} = 0 & w_{23} = 1 & w_{24} = 0 \\
 w_{31} = 1 & w_{32} = 1 & w_{33} = 1 & w_{34} = 0 \\
 w_{41} = 0 & w_{42} = 0 & w_{43} = 1 & w_{44} = 0
 \end{array}$$

The hidden units are *logistic*.

For the training case with that "3" input image, what is the output of each of the four hidden units?

- ☐  $h_1 = 0.982, h_2 = 0.881, h_3 = 1.000, h_4 = 0.982$
- ☐  $h_1 = 0.982, h_2 = 0.982, h_3 = 0.982, h_4 = 0.982$
- ☒  $h_1 = 0.982, h_2 = 1.000, h_3 = 0.881, h_4 = 0.982$

**Correct**

To get the answer we need to start by multiplying the weight matrix by the pixels that it is assigned to in order to get the *total input* to the hidden units. For hidden unit  $h_4$  this would mean:

$$\begin{aligned}
 z_4 &= w_{11}x_{22} + w_{12}x_{23} + w_{13}x_{24} + w_{14}x_{25} \\
 &\quad + w_{21}x_{32} + w_{22}x_{33} + w_{23}x_{34} + w_{24}x_{35} \\
 &\quad + w_{31}x_{42} + w_{32}x_{43} + w_{33}x_{44} + w_{34}x_{45} \\
 &\quad + w_{41}x_{52} + w_{42}x_{53} + w_{43}x_{54} + w_{44}x_{55} \\
 &= (1)(0) + (1)(0) + (1)(1) + (0)(0) \\
 &\quad + (0)(1) + (0)(1) + (1)(1) + (0)(0) \\
 &\quad + (1)(0) + (1)(0) + (1)(1) + (0)(0) \\
 &\quad + (0)(1) + (0)(1) + (1)(1) + (0)(0) \\
 &= 4
 \end{aligned}$$

Now we pass this through a sigmoid to get

$$y_4 = \frac{1}{1+\exp(-z_4)} = \sigma(z_4) = \sigma(4) \approx 0.982$$

☒  $h_1 = 0.881, h_2 = 0.982, h_3 = 0.982, h_4 = 1.000$



1 / 1  
points

3.

Recall that pooling is the process of combining the outputs of several hidden units to create a single hidden unit. This introduces some invariance to

local transformations in the input image.

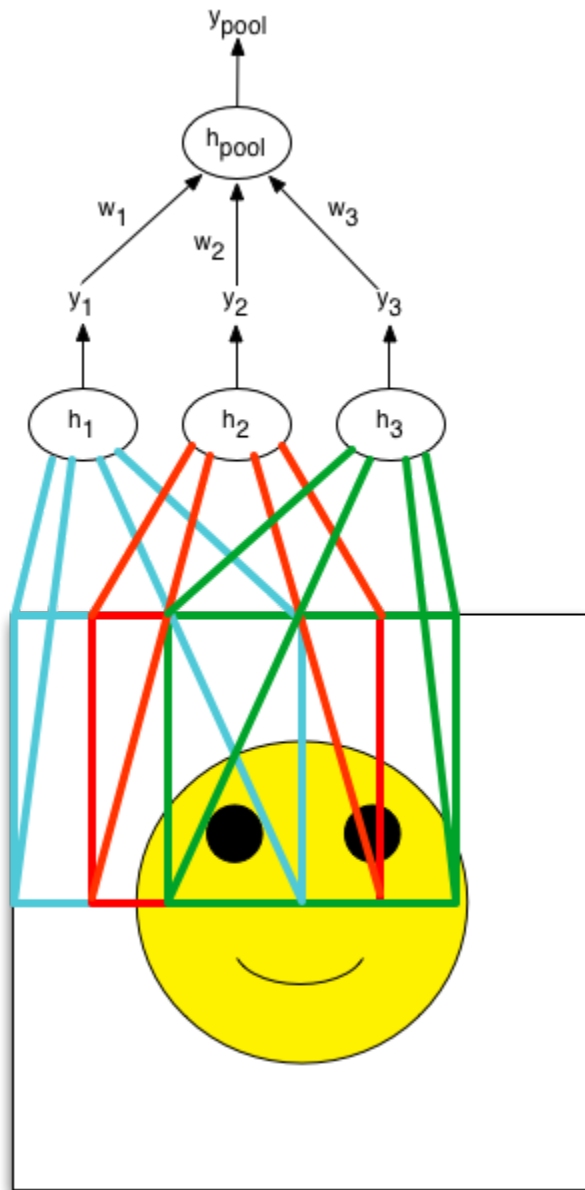
In this example we are pooling hidden units  $h_1$ ,  $h_2$ , and  $h_3$ . Let's denote the output of hidden unit  $h_i$  as  $y_i$ . The hidden unit

that we're creating by pooling is called  $h_{\text{pool}}$ , with output  $y_{\text{pool}}$ . **Sum pooling** is defined as  $y_{\text{pool}} = y_1 + y_2 + y_3$ .

This form of pooling can be equivalently represented by making  $h_{\text{pool}}$  a regular hidden unit that takes the output of the other three hidden

units, multiplies them by some weights  $w_1$ ,  $w_2$  and  $w_3$ , adds up the resulting numbers and then outputs some function of this sum. This is

illustrated in the figure below.



For this to be equivalent to **sum pooling**, what type of neuron should  $h_{\text{pool}}$  be, and what should the weights be?

- ☐  $w_1 = 1, w_2 = 1, w_3 = 1$ , and  $h_{\text{pool}}$  is a logistic neuron.
- ☒  $w_1 = 1, w_2 = 1, w_3 = 1$ , and  $h_{\text{pool}}$  is a linear neuron.

**Correct**

If  $w_1 = 1, w_2 = 1, w_3 = 1$ , then we can calculate:

$$z_{\text{pool}} = w_1 y_1 + w_2 y_2 + w_3 y_3 = y_1 + y_2 + y_3$$

and if  $h_{\text{pool}}$  is a linear neuron, then  $y_{\text{pool}} = z_{\text{pool}}$  so:

$$y_{\text{pool}} = y_1 + y_2 + y_3$$

Which is equivalent to the definition given above.

- ☐  $w_1 = \frac{1}{3}, w_2 = \frac{1}{3}, w_3 = \frac{1}{3}$ , and  $h_{\text{pool}}$  is a linear neuron.
- ☐  $w_1 = \frac{1}{3}, w_2 = \frac{1}{3}, w_3 = \frac{1}{3}$ , and  $h_{\text{pool}}$  is a logistic neuron.
- 



1 / 1  
points

4.

Suppose that we have a vocabulary of 3 words, "a", "b", and "c", and we want to predict the next word in a sentence given the previous two words. For

this network, we don't want to use feature vectors for words: we simply use the local encoding, i.e. a 3-component vector with one entry being 1

and the other two entries being 0.

In the language models that we have seen so far, each of the context words has its own dedicated section of the network, so we would encode this

problem with two 3-dimensional inputs. That makes for a total of 6 dimensions. For example, if the two preceding words (the "context" words) are "c"

and "b", then the input would be  $(0, 0, 1, 0, 1, 0)$ . Clearly, the more context words we want to include, the more input units our network must

have. More inputs means more parameters, and thus increases the risk of overfitting. Here is a proposal to reduce the number of parameters in the

model:

Consider a single neuron that is connected to this input, and call the weights that connect the input to this neuron  $w_1, w_2, w_3, w_4, w_5$ , and

$w_6$ .  $w_1$  connects the neuron to the first input unit,  $w_2$  connects it to the second input unit, etc. Notice how for every neuron, we need

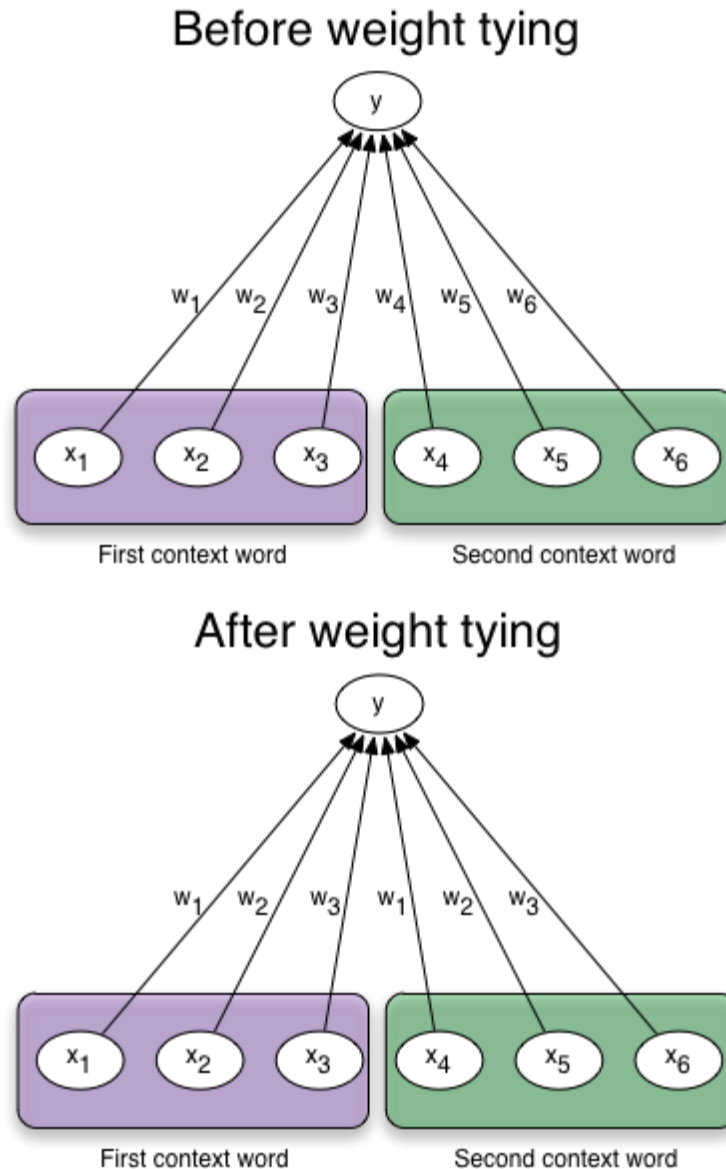
as many weights as there are input dimensions (6 in our case), which will be the number of words times the length of the context. A way to reduce the number of

parameters is to *tie* certain weights together, so that they share a

parameter. One possibility is to tie the weights coming from input units

that correspond to the same word but at different context positions. In our example that would mean that  $w_1 = w_4$ ,  $w_2 = w_5$ , and  $w_3 = w_6$

(see the "after" diagram below).



Are there any significant problems with this approach?

- ☐ No: the new model after weight tying is an example of a convolutional neural network, and these are more powerful than a non-convolutional network because they are invariant to small transformations in the data.
- ☐ No: this method is an appropriate solution in that it will reduce the number of parameters and therefore always improve generalization.



- ☐ Yes: weight tying only makes sense when we are working with images.
- ☒ Yes: the network loses the knowledge of the location at which a context word occurs, and that is valuable knowledge.

**Correct**

Hint: an equivalent network to the one with weight-tying is one with only 3 inputs and weights  $w_1, w_2, w_3$  (instead of 6 inputs), but instead of storing binary numbers, each dimension now stores the *counts* of each word. That is, each input stores the number of times the word appeared in the sequence. So the input [1 0 0 1 0 0] would be represented as [2 0 0]. What can we say about this network now?



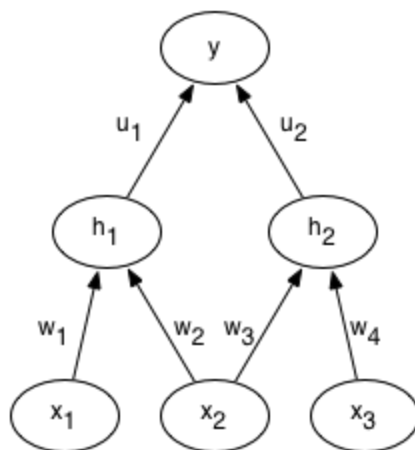
1 / 1  
points

5.

Let's look at what weight tying does to gradients, computed using the backpropagation algorithm. For this question, our network has three input units,

$x_1, x_2, x_3$ , two *logistic* hidden units  $h_1, h_2$ , four input to hidden weights  $w_1, w_2, w_3, w_4$ , and two hidden to output weights

$u_1, u_2$ . The output neuron  $y$  is a *linear neuron*, and we are using the *squared error* cost function.



Consider a single training case with target output  $t$ . The sequence of computations required to compute the error (this is called forward

propagation) are as follows:

propagation) are as follows:

$$z_1 = w_1x_1 + w_2x_2$$

$$z_2 = w_3x_2 + w_4x_3$$

$$h_1 = \sigma(z_1)$$

$$h_2 = \sigma(z_2)$$

$$y = u_1h_1 + u_2h_2$$

$$E = \frac{1}{2} (t - y)^2$$

$E$  is the error. Reminder:  $\sigma(x) = \frac{1}{1+\exp(-x)}$

$z_1$  and  $z_2$  are called the *total inputs* to hidden units  $h_1$  and  $h_2$  respectively.

Suppose we now decide to tie the weights so that  $w_2 = w_3 = w_{\text{tied}}$ . What is the derivative of the error  $E$  with respect to

$w_{\text{tied}}$ ?

Hint: forget about weight tying for a moment and simply compute the derivatives  $\frac{\partial E}{\partial w_2}$  and  $\frac{\partial E}{\partial w_3}$ .

Now set  $\frac{\partial E}{\partial w_{\text{tied}}} = \frac{\partial E}{\partial w_2} + \frac{\partial E}{\partial w_3}$ .

☐  $\frac{\partial E}{\partial w_{\text{tied}}} = -2(t - y)(u_2h_2(1 - h_2)x_2)$

☐  $\frac{\partial E}{\partial w_{\text{tied}}} = -2(t - y)(u_1h_1(1 - h_1)x_2)$

☐  $\frac{\partial E}{\partial w_{\text{tied}}} = -(t - y)(u_1h_1(1 - h_1)x_1 + u_2h_2(1 - h_2)x_3)$

☒  $\frac{\partial E}{\partial w_{\text{tied}}} = -(t - y)[u_1h_1(1 - h_1) + u_2h_2(1 - h_2)]x_2$

**Correct**

Starting from the error, backpropagation works by repeated application of the chain rule. Let's look at  $\frac{\partial E}{\partial w_2}$  and remember that we are ignoring weight tying for now:

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial w_2}$$

$$\frac{\partial E}{\partial y} = -(t - y)$$

$$\frac{\partial y}{\partial h_1} = u_1$$

$$\frac{\partial h_1}{\partial z_1} = h_1(1 - h_1)$$

$$\frac{\partial z_1}{\partial w_2} = x_2$$

$$\text{So } \frac{\partial E}{\partial w_2} = -(t - y)u_1h_1(1 - h_1)x_2$$

$$\text{We can similarly compute } \frac{\partial E}{\partial w_3} = -(t - y)u_2h_2(1 - h_2)x_2$$

To compute the derivative with tied weights, we now just need to add the derivatives:

$$\begin{aligned} \frac{\partial E}{\partial w_{\text{tied}}} &= \frac{\partial E}{\partial w_2} + \frac{\partial E}{\partial w_3} \\ &= -(t - y)u_1h_1(1 - h_1)x_2 - (t - y)u_2h_2(1 - h_2)x_2 \\ &= -(t - y)[u_1h_1(1 - h_1) + u_2h_2(1 - h_2)]x_2 \end{aligned}$$

How did we know which sequence of derivatives we would need for backpropagation? One answer is to look at the equations, the other is to look at the picture. Take  $w_2$  for example, and draw a path up the tree to the root. Notice how this path goes through  $h_1$  and not  $h_2$ ?

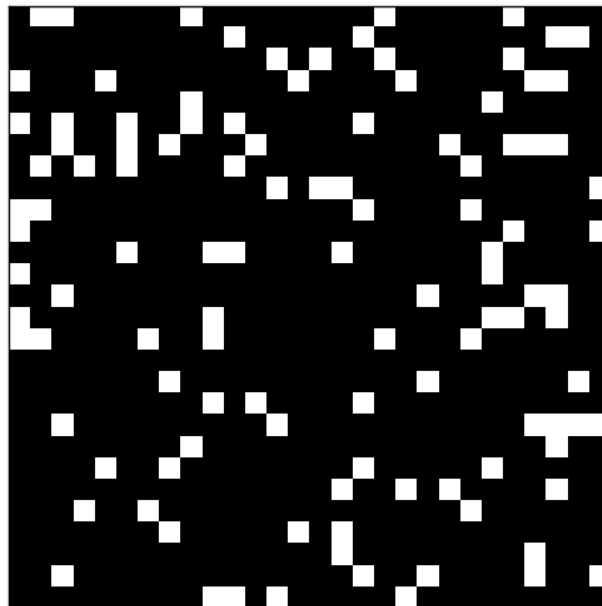
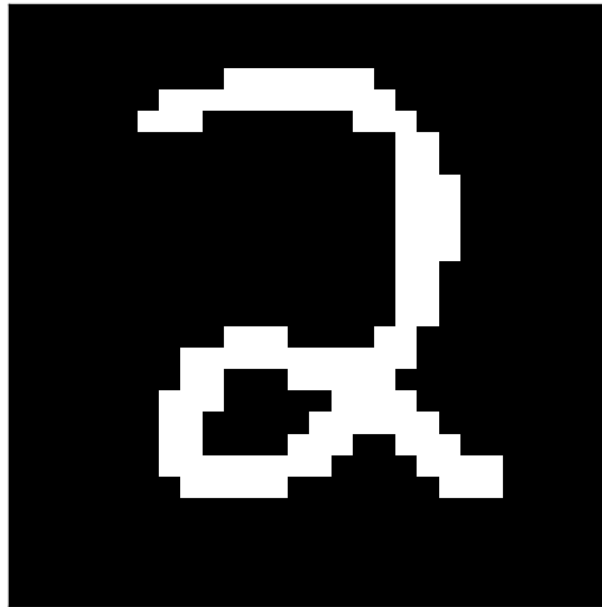


1 / 1  
points

6.

Oh no! Brian's done it again! Claire had a dataset of 28 x 28 pixel handwritten digits nicely prepared to train a neural network, but Brian has gone and accidentally scrambled the images by re-ordering the pixels in some totally meaningless way, and now they can't get the original dataset back! Below is an image of a handwritten '2' before and after being scrambled by Brian.

BeforeAfter



Luckily, all of the images (in both the training set and the test set) were changed in the same way. For example, if pixels number 1 and number 3 switched places in one image, then they switched places in every other image as well. Because of that, Claire thinks that perhaps she can train a network after all.

Whether Claire is right or not depends largely on the type of neural network that she has in mind. Which of the following neural networks will be at a **disadvantage** because of Brian's mistake? Check all that apply.

- ☐ A feed-forward neural network with one hidden layer of logistic units (and no convolution).



**Un-selected is correct**

- ☐ A convolutional neural network where the size of each weight filter is 10 x 10.



**Correct**

Imagine that the most important feature for being able to classify this '2' correctly is the little loopy part. A convolutional neural network will normally be able to identify this (provided the weight filter was big enough) because the important piece is located in a small spatial area. Once the dataset has been scrambled, the pixels for the loopy part will be potentially scattered all over the image, so a convolutional network that only has small weight filters will not be able to see enough of the loopy part to identify it.

- ☐ A convolutional neural network where the size of each weight filter is 8 x 8.



**Correct**

Imagine that the most important feature for being able to classify this '2' correctly is the little loopy part. A convolutional neural network will normally be able to identify this (provided the weight filter was big enough) because the important piece is located in a small spatial area. Once the dataset has been scrambled, the pixels for the loopy part will be potentially scattered all over the image, so a convolutional network that only has small weight filters will not be able to see enough of the loopy part to identify it.

- ☐ A feed-forward neural network with no hidden layer and linear units (and no convolution).



**Un-selected is correct**

