

Convolutional Neural Networks*(find and work with a partner)*

1. *Activation Functions.* So far we have seen three different activation functions:

- Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$
- Hyperbolic Tangent: $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- ReLU: $f(x) = \max(0, x)$

If the input has value exactly 0, what is the activation of each function above?

2. *Neural Network Architectures.* Say we have a fully connected neural network with 3 layers. The input data has shape $(n, p) = (100, 3)$, the first and second hidden layers each have 4 units (i.e. $p_1 = p_2 = 4$), and we have one output. If we use biases for all 3 layers, how many total parameters do we need to optimize?

3. *Cross-Entropy Loss.* We define cross-entropy between two discrete probability distributions p and q (with the same set of inputs \mathcal{X}) as

$$H(p, q) = - \sum_{x \in \mathcal{X}} p(x) \log_2 q(x).$$

When using cross-entropy as a classification loss function between the true one-hot class label vector y and the output probabilities for each class \hat{y} , we have

$$H(y, \hat{y}) = - \sum_{k=1}^K y_k \log_2 \hat{y}_k,$$

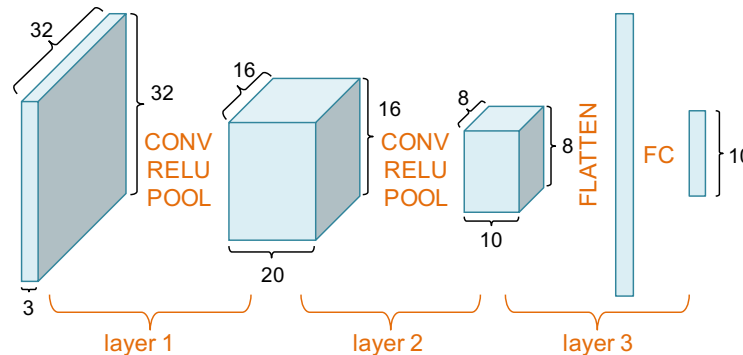
where K is the number of classes. Say we have $K = 3$. For a given data point \vec{x} , say the true label is class 2, and our neural network produced the probabilities $\hat{y} = [\frac{1}{4}, \frac{1}{2}, \frac{1}{4}]$. What is the cross-entropy loss for this data point?

If instead our method produced probabilities $\hat{y} = [\frac{3}{8}, \frac{1}{8}, \frac{1}{2}]$, what is the loss for this data point?

Note: for multiple examples, we usually use the loss function:

$$J(\text{all weights}) = \frac{1}{n} \sum_{i=1}^n H(y_i, \hat{y}_i)$$

4. Say we use a 3-layer CNN with architecture shown below. Our inputs have shape $(32 \times 32 \times 3)$.



- In the CONV step of the first layer, we use 20 filters, each (5×5) in width and height, but all the way through the depth. We use a stride of 1 in each dimension, and use padding = “SAME” to make sure the input and output width/height are the same. After CONV, we apply a RELU non-linearity, then apply POOL using a max-pooling strategy with (2×2) filters and stride 2 in width/height. This reduces the width and height by a factor of 2.
 - In the second layer we use 10 filters, each (3×3) in width and height, but all the way through the depth. The stride and padding follow the same procedure as the first layer. ReLU and pooling also follow the same strategy.
 - Finally, we flatten the volume in preparation for the full connected layer. The FC layer transforms the flattened volume into scores for 10 classes.
- (a) Which steps (i.e. CONV, RELU, POOL, FLATTEN, FC) require parameter learning through gradient descent? Which steps don't?
- (b) How many parameters do we need to learn for the first layer? What if we also included a bias for each filter?
- (c) How many parameters do we need to learn for the second layer? What if we also included a bias for each filter?
- (d) How many parameters do we need to learn for the third (FC) layer? What if we also included a bias for each class?
- (e) Assuming we keep the biases for each layer, how many parameters total do we need to learn?