

### THE UNIVERSITY OF BRITISH COLUMBIA

# **CPSC 425: Computer Vision**



Lecture 30: Neural Networks

# Menu for Today (November 23, 2020)

### **Topics:**

- Backpropagation
- Convolutional Layers

### **Redings:**

- Today's Lecture: N/A
- **Next** Lecture: N/A

### **Reminders:**

- Assignment 6: Deep Learning due Wednsday, November 2nd



### Pooling Layer



# Please fill out **Student Evaluations** (on Canvas)



— The basic unit of computation in a neural network is a neuron.

- A neuron accepts some number of input signals, computes their weighted sum, and applies an activation function (or non-linearity) to the sum.

- Common activation functions include sigmoid and rectified linear unit (ReLU)

### Lecture 29: Re-cap

A neural network comprises neurons connected in an acyclic graph The outputs of neurons can become inputs to other neurons Neural networks typically contain multiple layers of neurons



# **Neural** Network

### hidden layer

**Figure credit**: Fei-Fei and Karpathy

Example of a neural network with three inputs, a single hidden layer of four neurons, and an output layer of two neurons



### Lecture 29: Re-cap

Note: each neuron will have its own vector of weights and a bias, its easier to think of all neurons in a layer as a single entity with a matrix of weights (size = number of inputs x number of neurons) and a vector of biases (size = number of neurons)



# **Neural** Network

**Figure credit**: Fei-Fei and Karpathy

6





### Lecture 29: Re-cap

Note: each neuron will have its own vector of weights and a bias, its easier to think of all neurons in a layer as a single entity with a matrix of weights (size = number of inputs x number of neurons) and a vector of biases (size = number of neurons)



 $\hat{\mathbf{y}} = f(\mathbf{x}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2) =$ 

# **Neural** Network

**Figure credit**: Fei-Fei and Karpathy

$$= \sigma \left( \mathbf{W}_2^{(2 \times 4)} \sigma \left( \mathbf{W}_1^{(4 \times 3)} \mathbf{x} + \mathbf{b}_1^{(4)} \right) + \mathbf{b}_2^{(2)} \right)$$





When training a neural network, the final output will be some loss (error) function

- e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_j$  is the j-th element of the vector of class scores coming from neural net.

$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$

When training a neural network, the final output will be some loss (error) function

- e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_j$  is the j-th element of the vector of class scores coming from neural net.

Consider neural net which takes input vector  $\mathbf{x}_i$  and predicts scores for 3 classes, with true class being class 3:

$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$

When training a neural network, the final output will be some loss (error) function

- e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_j$  is the j-th element of the vector of class scores coming from neural net.

Consider neural net which takes input vector  $\mathbf{x}_i$  and predicts scores for 3 classes, with true class being class 3:

$$f$$
  
 $c_1 = -2.85$   
 $c_2 = 0.86$   
 $c_3 = 0.28$ 

$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$

When training a neural network, the final output will be some loss (error) function

- e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_j$  is the j-th element of the vector of class scores coming from neural net.

Consider neural net which takes input vector  $\mathbf{x}_i$  and predicts scores for 3 classes, with true class being class 3:



$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$

When training a neural network, the final output will be some loss (error) function

- e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_j$  is the j-th element of the vector of class scores coming from neural net.

Consider neural net which takes input vector  $\mathbf{x}_i$  and predicts scores for 3 classes, with true class being class 3:



$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$

 Normalize to sum to 1
 0.016

 0.631
 0.353

When training a neural network, the final output will be some loss (error) function

- e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_j$  is the j-th element of the vector of class scores coming from neural net.

Consider neural net which takes input vector  $\mathbf{x}_i$  and predicts scores for 3 classes, with true class being class 3:



$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$

probability of a class

Normalize to sum to 1

 $0.016 \\ 0.631 \\ 0.353$ 

When training a neural network, the final output will be some loss (error) function

– e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_i$ is the j-th element of the vector of class scores coming from neural net.

Consider neural net which takes input vector  $\mathbf{x}_i$  and predicts scores for 3 classes, with true class being class 3:



$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$

softmax function multi-class classifier

probability of a class

Normalize to sum to 1

0.0160.6310.353

When training a neural network, the final output will be some loss (error) function

- e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_i$ is the j-th element of the vector of class scores coming from neural net.

Consider neural net which takes input vector  $\mathbf{x}_i$  and predicts scores for 3 classes, with true class being class 3:



$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$

probability of a class

Normalize to sum to 1

0.016 $\longrightarrow 0.631$   $L_i = -\log(0.353) = 1.04$ 0.353



When training a neural network, the final output will be some loss (error) function

- e.g. cross-entropy loss:  $L_i = -$ 

which defines loss for i-th training example with true class index  $y_i$ ; and  $f_j$  is the j-th element of the vector of class scores coming from neural net.

We want to compute the **gradient** of the loss with respect to the network parameters so that we can incrementally adjust the network parameters

$$\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$$



\*slide adopted from V. Ordonex



1. Start from random value of  $\mathbf{W}_0, \mathbf{b}_0$ 

\*slide adopted from V. Ordonex



1. Start from random value of  $W_0, b_0$ 

\*slide adopted from V. Ordonex





### For k = 0 to max number of iterations

2. Compute gradient of the loss with respect to previous (initial) parameters:

 $\left. 
abla \, \mathcal{L}(\mathbf{W},\mathbf{b}) 
ight|_{\mathbf{W}=\mathbf{W}_k,\mathbf{b}=\mathbf{b}_k}$ 

\*slide adopted from V. Ordonex

. .

′∎



1. Start from random value of  $W_0, b_0$ 

For k = 0 to max number of iterations

2. Compute gradient of the loss with respect to previous (initial) parameters:

 $\nabla \mathcal{L}(\mathbf{W}, \mathbf{b})|_{\mathbf{W} = \mathbf{W}_k, \mathbf{b} = \mathbf{b}_k}$ 

\*slide adopted from V. Ordonex

. .

′∎



1. Start from random value of  $\mathbf{W}_0, \mathbf{b}_0$ 

For k = 0 to max number of iterations

2. Compute gradient of the loss with respect to previous (initial) parameters:

 $\left. 
abla \, \mathcal{L}(\mathbf{W}, \mathbf{b}) \right|_{\mathbf{W} = \mathbf{W}_k, \mathbf{b} = \mathbf{b}_k}$ 

3. Re-estimate the parameters

$$\mathbf{W}_{k+1} = \mathbf{W}_k - \lambda \left. \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}} \right|_{\mathbf{W} = \mathbf{W}_k}$$
$$\mathbf{b}_{k+1} = \mathbf{b}_k - \lambda \left. \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{b}} \right|_{\mathbf{W} = \mathbf{W}_k},$$

\*slide adopted from V. Ordonex

 $\mathbf{V}_k, \mathbf{b} = \mathbf{b}_k$ 

 $k, \mathbf{b} = \mathbf{b}_k$ 



1. Start from random value of  $\mathbf{W}_0, \mathbf{b}_0$ 

For k = 0 to max number of iterations

2. Compute gradient of the loss with respect to previous (initial) parameters:

 $\nabla \mathcal{L}(\mathbf{W}, \mathbf{b})|_{\mathbf{W} = \mathbf{W}_k, \mathbf{b} = \mathbf{b}_k}$ 

3. Re-estimate the parameters

$$\mathbf{W}_{k+1} = \mathbf{W}_k - \lambda \left. \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}} \right|_{\mathbf{W} = \mathbf{W}_k}$$
$$\mathbf{b}_{k+1} = \mathbf{b}_k - \lambda \left. \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{b}} \right|_{\mathbf{W} = \mathbf{W}_k, \mathbf{W}_k}$$

\*slide adopted from V. Ordonex

 $r_k, \mathbf{b} = \mathbf{b}_k$ 

 $\mathbf{b} = \mathbf{b}_k$ 



### 1. Start from random value of $\mathbf{W}_0, \mathbf{b}_0$

For k = 0 to max number of iterations

2. Compute gradient of the loss with respect to previous (initial) parameters:

$$\nabla \mathcal{L}(\mathbf{W}, \mathbf{b})|_{\mathbf{W} = \mathbf{W}_k, \mathbf{b} = \mathbf{b}_k}$$

3. Re-estimate the parameters

$$\mathbf{W}_{k+1} = \mathbf{W}_k - \lambda \left. \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}} \right|_{\mathbf{W} = \mathbf{W}_k}$$
$$\mathbf{b}_{k+1} = \mathbf{b}_k - \lambda \left. \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{b}} \right|_{\mathbf{W} = \mathbf{W}_k},$$

\*slide adopted from V. Ordonex

 $\mathbf{V}_k, \mathbf{b} = \mathbf{b}_k$ 

 $\mathbf{b} = \mathbf{b}_k$ 



 $\lambda$  - is the learning rate

### 1. Start from random value of $\mathbf{W}_0, \mathbf{b}_0$

For k = 0 to max number of iterations

2. Compute gradient of the loss with respect to previous (initial) parameters:

 $\left. 
abla \, \mathcal{L}(\mathbf{W},\mathbf{b}) 
ight|_{\mathbf{W}=\mathbf{W}_k,\mathbf{b}=\mathbf{b}_k}$ 

3. Re-estimate the parameters

$$\mathbf{W}_{k+1} = \mathbf{W}_k - \underline{\lambda} \left. \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}} \right|_{\mathbf{W} = \mathbf{W}_k}$$
$$\mathbf{b}_{k+1} = \mathbf{b}_k - \underline{\lambda} \left. \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial \mathbf{b}} \right|_{\mathbf{W} = \mathbf{W}_k},$$

\*slide adopted from V. Ordonex

 $V_k, \mathbf{b} = \mathbf{b}_k$ 

 $\mathbf{b} = \mathbf{b}_k$ 

Loss:



 $\mathbf{\hat{y}} = f(\mathbf{x}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2) =$ 

### $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = ||\mathbf{y} - \hat{\mathbf{y}}|| = ||\mathbf{y} - f(\mathbf{x}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2)||$

Figure credit: Fei-Fei and Karpathy

$$= \sigma \left( \mathbf{W}_2^{(2 \times 4)} \sigma \left( \mathbf{W}_1^{(4 \times 3)} \mathbf{x} + \mathbf{b}_1^{(4)} \right) + \mathbf{b}_2^{(2)} \right)$$

18



Loss:

**Gradient** Descent  

$$\mathbf{W}_{1,i,j} = \mathbf{W}_{1,i,j} - \lambda \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}_{1,i,j}}$$

$$\mathbf{b}_{1,i} = \mathbf{b}_{1,i} - \lambda \frac{\partial \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{b}_{1,i}}$$



### input layer

 $\mathbf{\hat{y}} = f(\mathbf{x}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2) =$ 

### $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = ||\mathbf{y} - \hat{\mathbf{y}}|| = ||\mathbf{y} - f(\mathbf{x}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2)||$

hidden layer

Figure credit: Fei-Fei and Karpathy

$$= \sigma \left( \mathbf{W}_2^{(2 \times 4)} \sigma \left( \mathbf{W}_1^{(4 \times 3)} \mathbf{x} + \mathbf{b}_1^{(4)} \right) + \mathbf{b}_2^{(2)} \right)$$

19







### **Problem:** For large datasets computing sum is expensive



### **Problem:** For large datasets computing sum is expensive

**Solution:** Compute approximate gradient with mini-batches of much smaller size (as little as 1-example sometimes)

$$\mathbf{gmoid}\left(\mathbf{W}^T\mathbf{x}^{(d)} + \mathbf{b}\right) - \mathbf{y}^{(d)}\right)^2$$



### **Problem:** For large datasets computing sum is expensive

**Solution:** Compute approximate gradient with mini-batches of much smaller size (as little as 1-example sometimes)

**Problem:** How do we compute the actual gradient?

$$\operatorname{gmoid}\left(\mathbf{W}^T\mathbf{x}^{(d)} + \mathbf{b}\right) - \mathbf{y}^{(d)}\right)^2$$

We can approximate the gradient numerically, using:

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \to 0}$$

 $\mathbf{1}_i$  - Vector of all zeros, except for one 1 in i-th location

 $\frac{f(\mathbf{x} + h\mathbf{1}_i) - f(\mathbf{x})}{h}$ 



We can approximate the gradient numerically, using:

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \to 0} \frac{f(\mathbf{x} + h\mathbf{1}_i) - f(\mathbf{x})}{h}$$

Even better, we can use central differencing:

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \to 0} \frac{f(\mathbf{x} + h\mathbf{1}_i) - f(\mathbf{x} - h\mathbf{1}_i)}{2h}$$

 $\mathbf{1}_i$  - Vector of all zeros, except for one 1 in i-th location



We can approximate the gradient numerically, using:

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \to 0} \frac{f(\mathbf{x} + h\mathbf{1}_i) - f(\mathbf{x})}{h}$$

Even better, we can use central differencing:

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \lim_{h \to 0} \frac{f(\mathbf{x} + h\mathbf{1}_i) - f(\mathbf{x} - h\mathbf{1}_i)}{2h}$$

However, both of theses suffer from rounding errors and are not good enough for learning (they are very good tools for checking the correctness of implementation though, e.g., use h = 0.000001).

 $\mathbf{1}_i$  - Vector of all zeros, except for one 1 in i-th location



We can approximate the gradient numerically, using:

$$\frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial w_{ij}} \approx \lim_{h \to 0} \frac{\mathcal{L}(\mathbf{W} + h\mathbf{1}_{ij}, \mathbf{b}) - \mathcal{L}(\mathbf{W}, \mathbf{b})}{h}$$

Even better, we can use central differencing:

$$\frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial w_{ij}} \approx \lim_{h \to 0} \frac{\mathcal{L}(\mathbf{W} + h\mathbf{1}_{ij}, \mathbf{b}) - \mathcal{L}(\mathbf{W} + h\mathbf{1}_{ij}, \mathbf{b})}{2h} \qquad \qquad \frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial b_j} \approx \lim_{h \to 0} \frac{\mathcal{L}(\mathbf{W}, \mathbf{b} + h\mathbf{1}_j) - \mathcal{L}(\mathbf{W}, \mathbf{b} + h\mathbf{1}_j)}{2h}$$

However, both of theses suffer from rounding errors and are not good enough for learning (they are very good tools for checking the correctness of implementation though, e.g., use h = 0.000001).

 $\mathbf{1}_i$  - Vector of all zeros, except for one 1 in i-th location  $\mathbf{1}_{ij}$  - Matrix of all zeros, except for one 1 in (i,j)-th location

$$\frac{\partial \mathcal{L}(\mathbf{W}, \mathbf{b})}{\partial b_j} \approx \lim_{h \to 0} \frac{\mathcal{L}(\mathbf{W}, \mathbf{b} + h\mathbf{1}_j) - \mathcal{L}(\mathbf{W}, \mathbf{b})}{h}$$





# **Symbolic** Differentiation

### Input function is represented as **computational graph** (a symbolic tree)



Implements differentiation rules for composite functions:



$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - sin(x_2)$$






### **Symbolic** Differentiation

### Input function is represented as **computational graph** (a symbolic tree)



Implements differentiation rules for composite functions:



**Problem:** For complex functions, expressions can be exponentially large; also difficult to deal with piece-wise functions (creates many symbolic cases)

\*slide adopted from T. Chen, H. Shen, A. Krishnamurthy CSE 599G1 lecture at UWashington

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$$







#### $y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - sin(x_2)$ **Automatic** Differentiation (AutoDiff)

**Intuition:** Interleave symbolic differentiation and simplification

**Key Idea:** apply symbolic differentiation at the elementary operation level, evaluate and keep intermediate results

\*slide adopted from T. Chen, H. Shen, A. Krishnamurthy CSE 599G1 lecture at UWashington

#### $y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin(x_2)$ **Automatic** Differentiation (AutoDiff)

**Intuition:** Interleave symbolic differentiation and simplification

**Key Idea:** apply symbolic differentiation at the elementary operation level, evaluate and keep intermediate results

Success of **deep learning** owes A LOT to success of AutoDiff algorithms (also to advances in parallel architectures, and large datasets, ...)

\*slide adopted from T. Chen, H. Shen, A. Krishnamurthy CSE 599G1 lecture at UWashington

The parameters of a neural network are learned using **backpropagation**, calculus

# which computes gradients via recursive application of the chain rule from

The parameters of a neural network are learned using **backpropagation**, which computes gradients via recursive application of the chain rule from calculus

Suppose f(x, y) = xy. What is the partial derivative of f with respect to x? What is the partial derivative of f with respect to y?

The parameters of a neural network are learned using **backpropagation**, which computes gradients via recursive application of the chain rule from calculus

is the partial derivative of f with respect to y?

$$\frac{\partial f}{\partial x} = y$$

Suppose f(x, y) = xy. What is the partial derivative of f with respect to x? What

$$\frac{\partial f}{\partial y} = x$$

What is the partial derivative of f with respect to y?

# Suppose f(x, y) = x + y. What is the partial derivative of f with respect to x?

# Suppose f(x, y) = x + y. What is the partial derivative of f with respect to x? What is the partial derivative of f with respect to y?

$$\frac{\partial f}{\partial x} = 1$$

$$\frac{\partial f}{\partial y} = 1$$

A trickier example:  $f(x, y) = \max(x, y)$ 

A trickier example:  $f(x, y) = \max(x, y)$ 

$$\frac{\partial f}{\partial x} = \mathbf{1}(x \ge y)$$

That is, the (sub)gradient is 1 on the input that is larger, and 0 on the other input

- For example, say x = 4, y = 2. Increasing y by a tiny amount does not change the value of f (f will still be 4), hence the gradient on y is zero.

$$\frac{\partial f}{\partial y} = \mathbf{1}(y \ge x)$$

applying the **chain rule** from calculus

# We can compose more complicated functions and compute their gradients by

We can compose more complicated functions and compute their gradients by applying the **chain rule** from calculus

to x? y? z?

Suppose f(x, y, z) = (x + y)z. What are the partial derivatives of f with respect

We can compose more complicated functions and compute their gradients by applying the **chain rule** from calculus

Suppose f(x, y, z) = (x + y)z. What to x? y? z?

For illustration we break this expression into q = x + y and f = qz. This is a sum and a product, and we have just seen how to compute partial derivatives for these.

Suppose f(x, y, z) = (x + y)z. What are the partial derivatives of f with respect

We can compose more complicated functions and compute their gradients by applying the **chain rule** from calculus

Suppose f(x, y, z) = (x + y)z. What to x? y? z?

For illustration we break this expression into q = x + y and f = qz. This is a sum and a product, and we have just seen how to compute partial derivatives for these.

By the chain rule

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = z \cdot 1 = z$$

Suppose f(x, y, z) = (x + y)z. What are the partial derivatives of f with respect

We can compose more complicated functions and compute their gradients by applying the **chain rule** from calculus

to x? y? z?

For illustration we break this expression into q = x + y and f = qz. This is a sum and a product, and we have just seen how to compute partial derivatives for these.

By the chain rule

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = z \cdot 1 = z$$

Suppose f(x, y, z) = (x + y)z. What are the partial derivatives of f with respect

 $\frac{\dot{z}}{\partial y} = \frac{\partial y}{\partial q} \frac{\partial y}{\partial y} = z \cdot 1 = z$  $\frac{J}{\partial z} = q$ 

f(x, y, z) = (x + y)z

**Computational graph** (a DAG) with variable ordering from topological sort, where each **node** is an input, intermediate, or output variable





**Computational graph** (a DAG) with variable ordering from topological sort, where each **node** is an input, intermediate, or output variable

Suppose the network input is: (x, y, y)

Then: 
$$q = x + y = 3$$
  $f = qz =$ 



$$z) = (-2, 5, -4)$$

-12(forward pass)



Suppose the network input is: (x, y, z) = (-2, 5, -4)

Then: q = x + y = 3 f = qz = -12

 $\frac{\partial f}{\partial a} = z = -4$  $\nabla q$ 



f(x, y, z) = (x + y)z

### (forward pass)

#### (**backward** pass)



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = \frac{\partial f}{\partial q} \cdot 1$$

Suppose the network input is: (x, y, z) = (-2, 5, -4)

Then: 
$$q = x + y = 3$$
  $f = qz =$ 

$$\frac{\partial f}{\partial q} = z = -4$$



f(x, y, z) = (x + y)z

#### -12(forward pass)

#### (**backward** pass)



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = \frac{\partial f}{\partial q} \cdot 1$$

Suppose the network input is: (x, y, z) = (-2, 5, -4)

Then: 
$$q = x + y = 3$$
  $f = qz =$ 

$$\frac{\partial f}{\partial q} = z = -4 \qquad \qquad \frac{\partial f}{\partial x} = -4$$



f(x, y, z) = (x + y)z

#### -12(forward pass)

#### (**backward** pass)



#### Back

propagation  

$$f(x, y, z) = (x + y)z$$

$$y$$

$$f(x, y, z) = (x + y)z$$

$$y$$

$$y$$

$$z$$

$$f(x, y, z) = (x + y)z$$

$$y$$

$$y$$

$$z$$

$$f(x, y, z) = (x + y)z$$

$$y$$

$$f(z, y, z) = (x + y)z$$

Suppose the network input is: (x, y, z) = (-2, 5, -4)

Then: 
$$q = x + y = 3$$
  $f = qz =$ 

$$\frac{\partial f}{\partial q} = z = -4 \qquad \qquad \frac{\partial f}{\partial x} = -4$$

-12(forward pass)

$$\frac{\partial f}{\partial y} = -4$$
  $\frac{\partial f}{\partial z} = 3$  (backward p













Lets create a neural network that will be able to differentiate (classify) these patterns of simple 3x3 pixel images



### We will need some labeled data







































Lets create a neural network that will be able to differentiate (classify) these patterns of simple 3x3 pixel images







First, lets re-formulate the problem



### What do we need to do?



Lets create a neural network that will be able to differentiate (classify) these patterns of simple 3x3 pixel images







First, lets re-formulate the problem



#### What do we need to do?



Lets create a neural network that will be able to differentiate (classify) these patterns of simple 3x3 pixel images







How many inputs should the network have? How neuron outputs?

### Now, lets build a **network**!



Lets create a neural network that will be able to differentiate (classify) these patterns of simple 3x3 pixel images













What else is missing for us to train it?



Lets create a neural network that will be able to differentiate (classify) these patterns of simple 3x3 pixel images









### **Output** Layer



#### Loss

 $L_i = -\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_{y_j}}}\right)$ 



Lets create a neural network that will be able to differentiate (classify) these patterns of simple 3x3 pixel images











**Output** Layer

Loss

 $L_{1} = -log\left(\frac{e^{\sum_{i=1}^{9}\sigma(w_{1,i}x_{i}+b_{1})}}{\sum_{i=1}^{3}e^{\sum_{i=1}^{9}\sigma(w_{1,i}x_{i}+b_{1})}}\right)$ 





## Fully Connected Layer



# **Example:** 200 x 200 image (small) x 40K hidden units

\* slide from Marc'Aurelio Renzato

## Fully Connected Layer



### **Example:** 200 x 200 image (small) x 40K hidden units

### = ~ 2 Billion parameters (for one layer!)

\* slide from Marc'Aurelio Renzato



## Fully Connected Layer



### **Example:** 200 x 200 image (small) x 40K hidden units

### = ~ 2 Billion parameters (for one layer!)

Spatial correlations are generally local

Waste of resources + we don't have enough data to train networks this large

\* slide from Marc'Aurelio Renzato




# Locally Connected Layer



# **Example:** 200 x 200 image (small) x 40K hidden units

## Filter size: 10 x 10

## = ~ 4 Million parameters

# Locally Connected Layer



# **Example:** 200 x 200 image (small) x 40K hidden units

## Filter size: 10 x 10

## = ~ 4 Million parameters

# **Stationarity** — statistics is similar at different locations



## **Example:** 200 x 200 image (small) x 40K hidden units

## **Filter size:** $10 \times 10$

## = ~ 4 Million parameters

## Share the same parameters across the locations (assuming input is stationary)

\* slide adopted from Marc'Aurelio Renzato





**Example:** 200 x 200 image (small) x 40K hidden units

**Filter size:**  $10 \times 10$ 

## = ~ 4 Million parameters

= 100 parameters

## Share the same parameters across the locations (assuming input is stationary)

\* slide adopted from Marc'Aurelio Renzato





































—  $\star$ 











# **Example:** 200 x 200 image (small) x 40K hidden units

## Filter size: 10 x 10

## **# of filters:** 20

## Learn multiple filters

\* slide from Marc'Aurelio Renzato

†*1* 



**Example:** 200 x 200 image (small) x 40K hidden units

## Filter size: 10 x 10

## **# of filters:** 20

= 2000 parameters

## Learn multiple filters

\* slide from Marc'Aurelio Renzato

†*1* 

32 x 32 x 3 image (note the image preserves spatial structure)



3 depth

## 32 x 32 x 3 **image**





## $5 \times 5 \times 3$ filter

**Convolve** the filter with the image (i.e., "slide over the image spatially, computing dot products")











## Filters always extend the full depth of the input volume

# 5 x 5 x 3 filter

**Convolve** the filter with the image (i.e., "slide over the image spatially, computing dot products"





## 32 x 32 x 3 **image**





**1 number:** the result of taking a dot product between the filter and a small 5 x 5 x 3 part of the image

$$\mathbf{W}^T \mathbf{x} + b$$
, where  $\mathbf{W}, \mathbf{x} \in \mathbb{R}^{75}$ 

## 32 x 32 x 3 **image**





**1 number:** the result of taking a dot product between the filter and a small 5 x 5 x 3 part of the image

$$\mathbf{W}^T \mathbf{x} + b$$
, where  $\mathbf{W}, \mathbf{x} \in \mathbb{R}^{75}$ 

## How many **parameters** does the layer have?

## 32 x 32 x 3 **image**





**1 number:** the result of taking a dot product between the filter and a small 5 x 5 x 3 part of the image

$$\mathbf{W}^T \mathbf{x} + b$$
, where  $\mathbf{W}, \mathbf{x} \in \mathbb{R}^{75}$ 

## How many **parameters** does the layer have? **76**

## 32 x 32 x 3 **image**





## activation map

\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

t

## 32 x 32 x 3 **image**





## activation map





this results in the "new image" of size 28 x 28 x 6!



- also affected by zero-padding
- input layer
- **Stride:** Controls spatial density. How far apart are depth columns?

The number of neurons in a layer is determined by depth and stride parameter

**Depth:** Controls number of neurons that connect to the same region of the

— a set of neurons connected to the same region is called a **depth column** 

# Convolutional Layer: Closer Look at Spatial Dimensions

## 32 x 32 x 3 **image**





## activation map



# Convolutional Neural Network (ConvNet)



3 depth


3 depth

### 28 height



3 depth



### **28** width



3 depth



3 depth



With padding we can achieve no shrinking (32 -> 28 -> 24); shrinking quickly (which happens with larger filters) doesn't work well in practice





As we go deeper in the network, filters learn and respond to increasingly specialized structures - The first layers may contain simple orientation filters, middle layers may respond to common substructures, and final layers may respond to entire objects

- **Convolutional neural networks** can be seen as learning a hierarchy of filters.

### What filters do networks learn?



Layer 1





[Zeiler and Fergus, 2013]



### What filters do networks learn?



[Zeiler and Fergus, 2013]





Let us assume the filter is an "eye" detector

How can we make detection spatially invariant (insensitive to position of the eye in the image)

\* slide from Marc'Aurelio Renzato



Let us assume the filter is an "eye" detector

How can we make detection spatially invariant (insensitive to position of the eye in the image)

> By "pooling" (e.g., taking a max) response over a spatial locations we gain robustness to position variations



\* slide from Marc'Aurelio Renzato

- Makes representation smaller, more manageable and spatially invariant
- Operates over each activation map independently



# e manageable and spatially invariant independently



- Makes representation smaller, more manageable and spatially invariant
- Operates over each activation map independently



# e manageable and spatially invariant independently



\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

How many **parameters**?

- Makes representation smaller, more manageable and spatially invariant
- Operates over each activation map independently



# e manageable and spatially invariant independently



### Max **Pooling**

### activation map





### max pool with 2 x 2 filter and stride of 2

6 8 3 4

### Average **Pooling**

### activation map





### avg pool with 2 x 2 filter and stride of 2

3.25 5.25 2 2