



# Deep Learning for Computer Vision

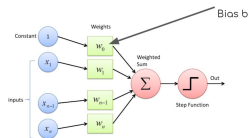
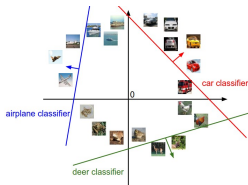
Dr. Konda Reddy Mopuri  
Mehta Family School of Data Science and Artificial Intelligence  
IIT Guwahati  
Aug-Dec 2022

# So far in the class..

- Image classification and Linear Classifier
- Perceptron



Dog  
Bird  
Car  
**Cat**  
Deer  
Truck



# Recap: Linear classifier



$$\textcircled{1} f(x) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

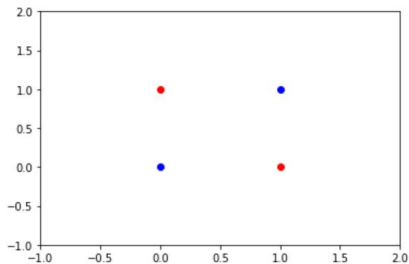
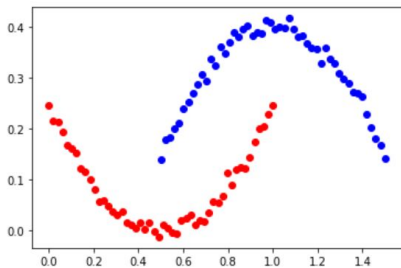
# Recap: Linear classifier



- ①  $f(x) = \sigma(\mathbf{w}^T \mathbf{x} + b)$
- ② Seen a couple of simple examples: MP neuron and Perceptron

# Linear Classifiers: Shortcomings

- Lower capacity: data has to be linearly separable
- Some times no hyper-plane can separate the data (e.g. XOR data)



# Pre-processing



- ① Sometimes, data specific pre-processing makes the data linearly separable

# Pre-processing

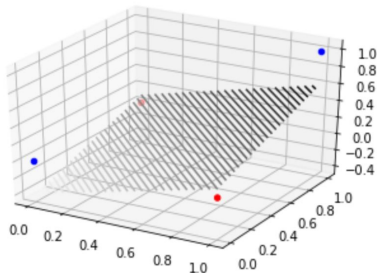


- ① Sometimes, data specific pre-processing makes the data linearly separable
- ② Consider the xor case

$$\phi(\mathbf{x}) = \phi(x_u, x_v) = (x_u, x_v, x_u x_v)$$

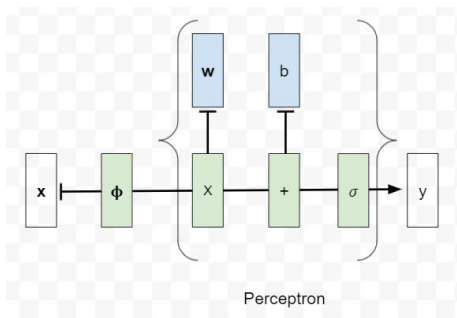
# Pre-processing

- ① Sometimes, data specific pre-processing makes the data linearly separable
- ② Consider the xor case  
 $\phi(\mathbf{x}) = \phi(x_u, x_v) = (x_u, x_v, x_u x_v)$
- ③ Consider the perceptron in the new space  $f(\mathbf{x}) = \sigma(\mathbf{w}^T \phi(\mathbf{x}) + b)$





# Pre-processing



# Pre-processing



- ① Feature design (or pre-processing) may also be another way to reduce the capacity without affecting (or improving) the bias

# Extending Linear Classifier



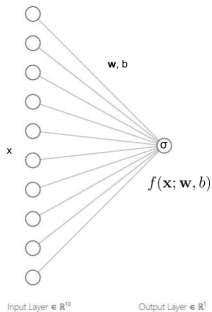
① Single class:  $f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$  from  $\mathcal{R}^D \rightarrow \mathcal{R}$  where  $\mathbf{w}$  and  $\mathbf{x} \in \mathcal{R}^D$



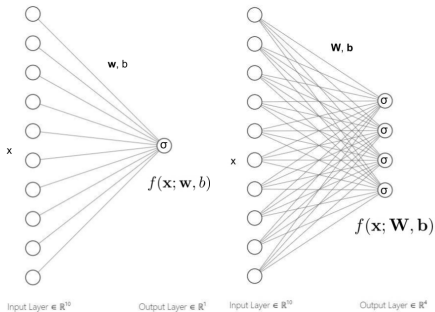
# Extending Linear Classifier

- ① Single class:  $f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$  from  $\mathcal{R}^D \rightarrow \mathcal{R}$  where  $\mathbf{w}$  and  $\mathbf{x} \in \mathcal{R}^D$
- ② Multi-class:  $f(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  from  $\mathcal{R}^D \rightarrow \mathcal{R}^C$  where  $\mathbf{W} \in \mathcal{R}^{C \times D}$  and  $\mathbf{b} \in \mathcal{R}^C$

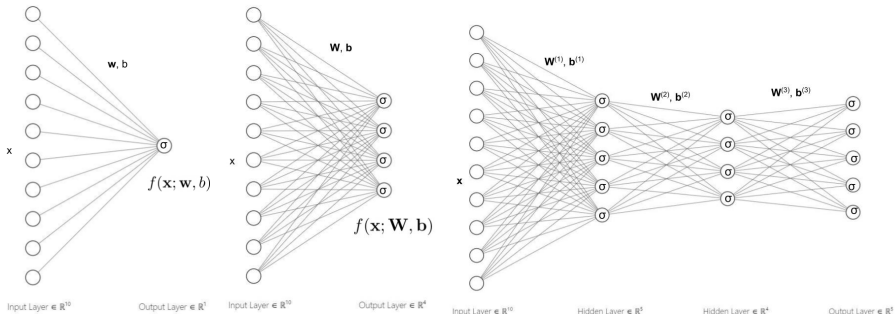
# Single unit to a layer of Perceptrons



# Single unit to a layer of Perceptrons



# Single unit to a layer of Perceptrons



# Formal Representation



- ① Latter is known as an MLP: Multi-Layered Perceptron (i.e, Multi-Layered network of Perceptrons)



# Formal Representation



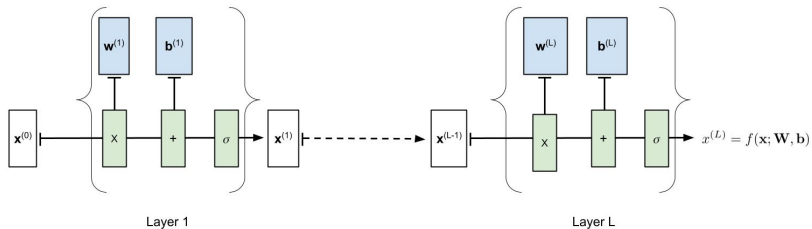
① Latter is known as an MLP: Multi-Layered Perceptron (i.e., Multi-Layered network of Perceptrons)

② can be represented as:

$$\mathbf{x}^{(0)} = \mathbf{x},$$

$$\forall l = 1, \dots, L, \quad \mathbf{x}^{(l)} = \sigma(\mathbf{W}^{(l)T} \mathbf{x}^{(l-1)} + \mathbf{b}^{(l)}), \text{ and}$$

# MLP



# Nonlinear Activation



- 1 Note that  $\sigma$  is nonlinear

# Nonlinear Activation

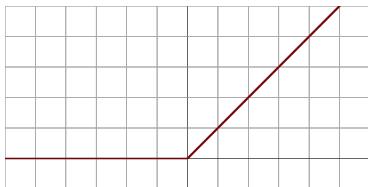
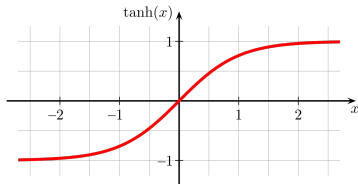


- ① Note that  $\sigma$  is nonlinear
- ② If it is an affine function, the full MLP becomes a complex affine transformation (composition of a series of affine mappings)

# Nonlinear Activation



## Familiar activation functions



Hyperbolic Tangent (Tanh)  $x \rightarrow \frac{2}{1+e^{-2x}} - 1$  and  
Rectified Linear Unit (ReLU)  $x \rightarrow \max(0, x)$  respectively



# Universal Approximation using ReLU functions

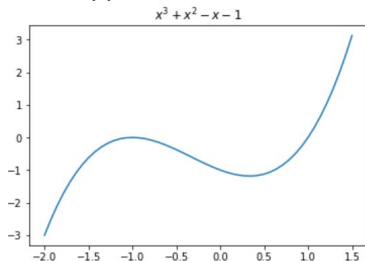
- ① We can approximate any function  $f$  from  $[a, b]$  to  $\mathcal{R}$  with a linear combination of ReLU functions

---

Example credits: Brendan Fortuner, and <https://towardsdatascience.com/>

# Universal Approximation using ReLU functions

- ① We can approximate any function  $f$  from  $[a, b]$  to  $\mathcal{R}$  with a linear combination of ReLU functions
- ② Let's approximate the following function using a bunch of ReLUs:

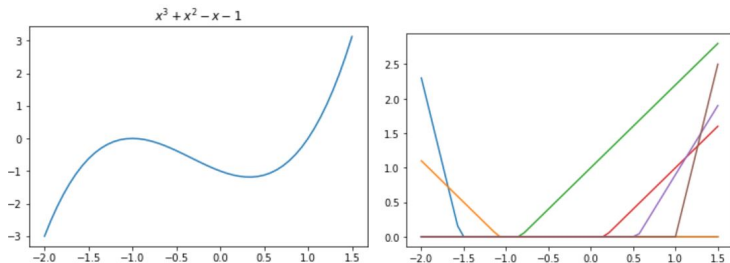


---

Example credits: Brendan Fortuner, and <https://towardsdatascience.com/>

# Universal Approximation using ReLU functions

- ①  $n_1 = \text{ReLU}(-5x - 7.7)$ ,  $n_2 = \text{ReLU}(-1.2x - 1.3)$ ,  $n_3 = \text{ReLU}(1.2x + 1)$ ,  $n_4 = \text{ReLU}(1.2x - 0.2)$ ,  $n_5 = \text{ReLU}(2x - 1.1)$ ,  $n_6 = \text{ReLU}(5x - 5)$



Example credits: Brendan Fortuner, and <https://towardsdatascience.com/>

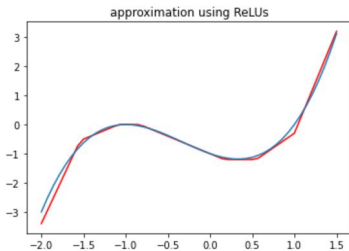


# Universal Approximation using ReLU functions



① Appropriate combination of these ReLUs:

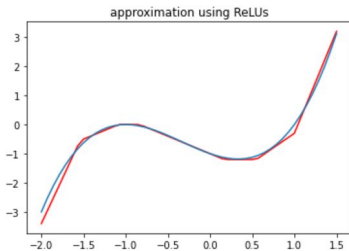
$$-n_1 - n_2 - n_3 + n_4 + n_5 + n_6$$



# Universal Approximation using ReLU functions



- ① Appropriate combination of these ReLUs:  
$$-n_1 - n_2 - n_3 + n_4 + n_5 + n_6$$
- ② Note that this also holds in case of other activation functions with mild assumptions.



# Universal Approximation Theorem



- ① We can approximate any continuous function  $\psi : \mathcal{R}^D \rightarrow \mathcal{R}$  with one hidden layer of perceptrons

# Universal Approximation Theorem



- ① We can approximate any continuous function  $\psi : \mathcal{R}^D \rightarrow \mathcal{R}$  with one hidden layer of perceptrons
- ②  $\mathbf{x} \rightarrow \mathbf{w}^T \sigma(W\mathbf{x} + \mathbf{b})$   
 $\mathbf{b} \in \mathcal{R}^C, W \in \mathcal{R}^{C \times D}, \mathbf{w} \in \mathcal{R}^C, \text{ and } \mathbf{x} \in \mathcal{R}^D$

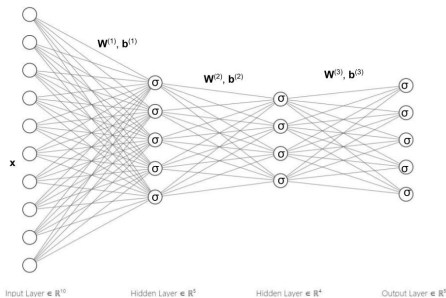
# Universal Approximation Theorem



- ① We can approximate any continuous function  $\psi : \mathcal{R}^D \rightarrow \mathcal{R}$  with one hidden layer of perceptrons
- ②  $\mathbf{x} \rightarrow \mathbf{w}^T \sigma(W\mathbf{x} + \mathbf{b})$   
 $\mathbf{b} \in \mathcal{R}^C, W \in \mathcal{R}^{C \times D}, \mathbf{w} \in \mathcal{R}^C$ , and  $\mathbf{x} \in \mathcal{R}^D$
- ③ However, the resulting NN
  - May require infeasible size for the hidden layer
  - May not generalize well

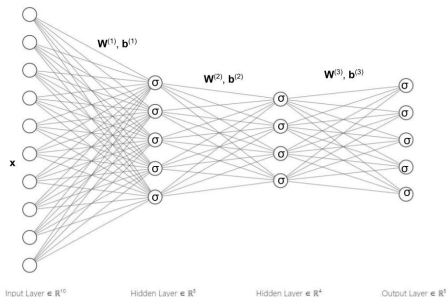
# MLP for regression

- ① Output is a continuous variable in  $\mathcal{R}^D$ 
  - Output layer has that many perceptrons (When  $D = 1$ , regresses a scalar value)
  - May employ a squared error loss



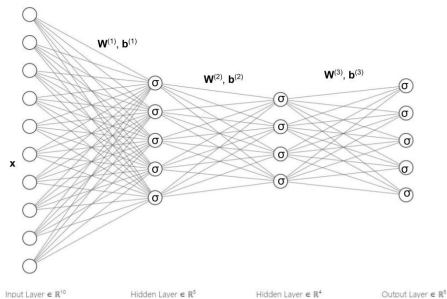
# MLP for regression

- ① Output is a continuous variable in  $\mathcal{R}^D$ 
  - Output layer has that many perceptrons (When  $D = 1$ , regresses a scalar value)
  - May employ a squared error loss
- ② Can have an arbitrary depth (number of layers)



# MLP for classification

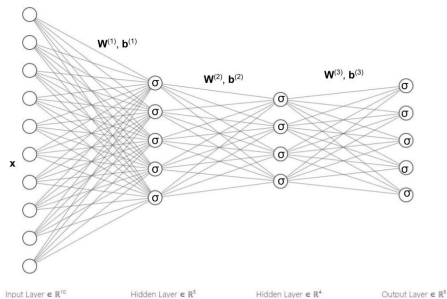
- ① Categorical output in  $\mathcal{R}^C$  where  $C$  is the number of categories





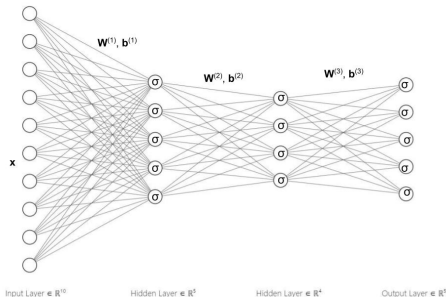
# MLP for classification

- ① Categorical output in  $\mathcal{R}^C$  where  $C$  is the number of categories
- ② Predicts the scores/confidences/probabilities towards each category
  - Then converts into a pmf
  - Employs loss that compares the probability distributions (e.g. cross-entropy)



# MLP for classification

- ① Categorical output in  $\mathcal{R}^C$  where  $C$  is the number of categories
- ② Predicts the scores/confidences/probabilities towards each category
  - Then converts into a pmf
  - Employs loss that compares the probability distributions (e.g. cross-entropy)
- ③ Can have an arbitrary depth



# Learning the NN parameters



- Gradient descent on the loss function?

# Learning the NN parameters



- Gradient descent on the loss function?
- Have to compute  $\frac{\partial L}{\partial W_{ij}^{(l)}}$  and  $\frac{\partial L}{\partial b_i^{(l)}}$ ,  $\forall W_{ij}^{(l)}, b_i^{(l)}$



# Learning the NN parameters

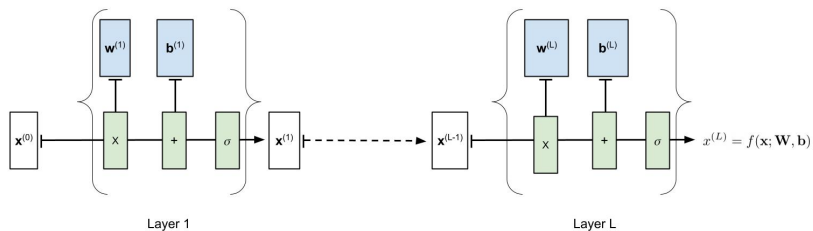
- Gradient descent on the loss function?
- Have to compute  $\frac{\partial L}{\partial W_{ij}^{(l)}}$  and  $\frac{\partial L}{\partial b_i^{(l)}}$ ,  $\forall W_{ij}^{(l)}, b_i^{(l)}$
- Almost impossible to derive these expressions analytically!

# Learning the NN parameters



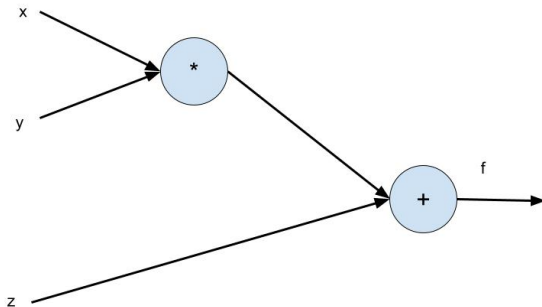
- Gradient descent on the loss function?
- Have to compute  $\frac{\partial L}{\partial W_{ij}^{(l)}}$  and  $\frac{\partial L}{\partial b_i^{(l)}}$ ,  $\forall W_{ij}^{(l)}, b_i^{(l)}$
- Almost impossible to derive these expressions analytically!
- specific to each loss function :- (

# Solution: Computational graphs



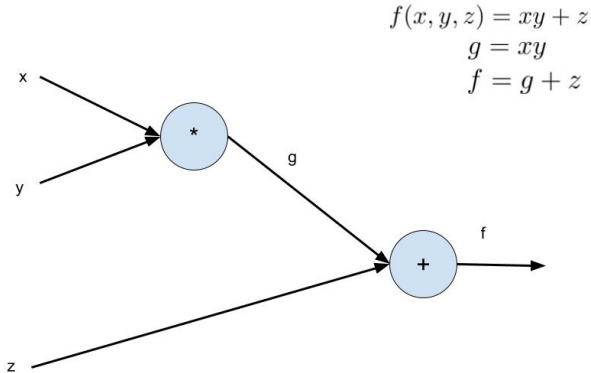
# E.g. Computational graph

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

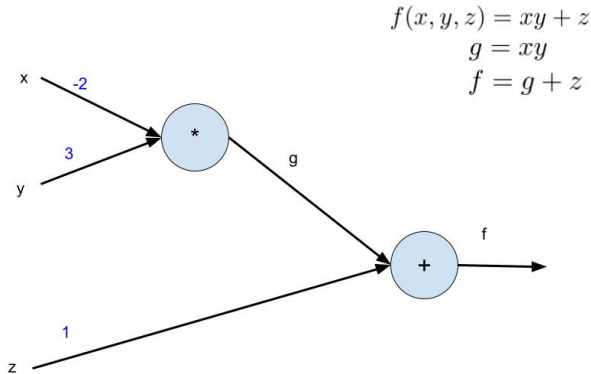




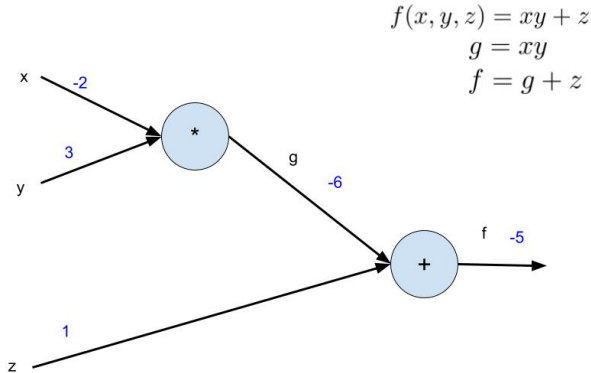
# E.g. Computational graph: Forward pass



# E.g. Computational graph: Forward pass



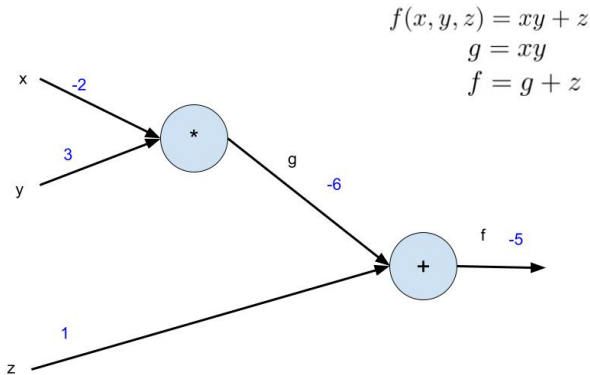
# E.g. Computational graph: Forward pass



# E.g. Computational graph: **Backward pass**



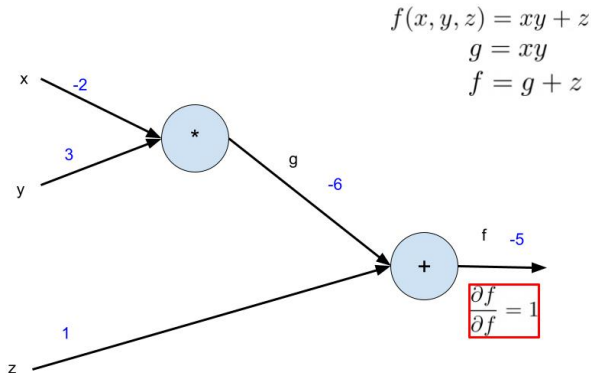
- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



# E.g. Computational graph: Backward pass



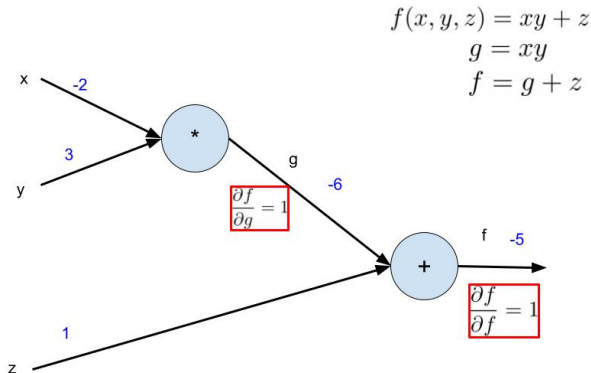
- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



# E.g. Computational graph: Backward pass



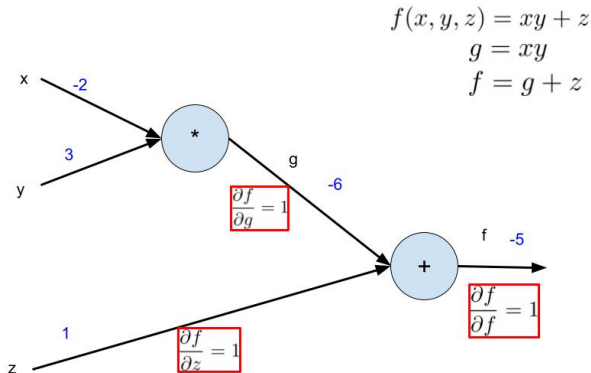
- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



# E.g. Computational graph: Backward pass

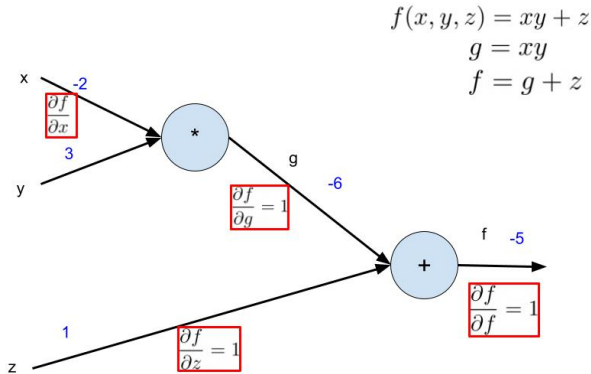


- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



# E.g. Computational graph: **Backward pass**

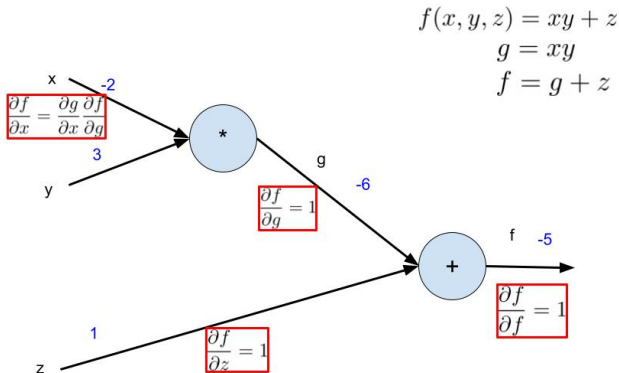
- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$





# E.g. Computational graph: **Backward pass**

- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



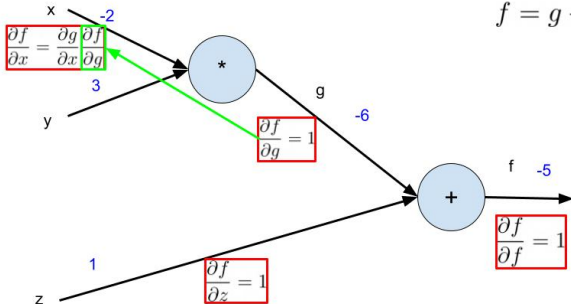
# E.g. Computational graph: **Backward pass**

- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$

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

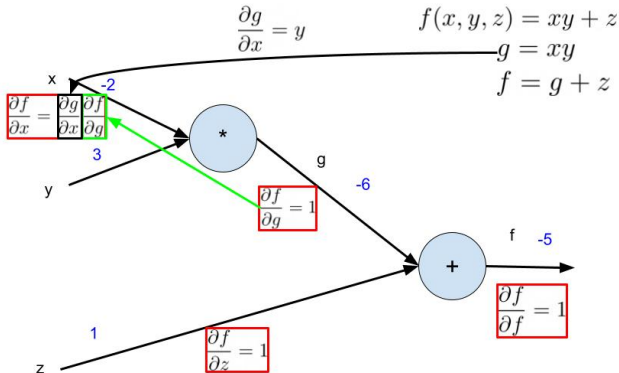
$$g = xy$$

$$f = g + z$$



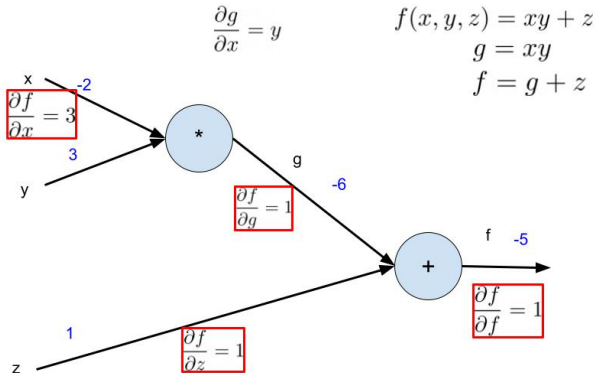
# E.g. Computational graph: **Backward pass**

- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



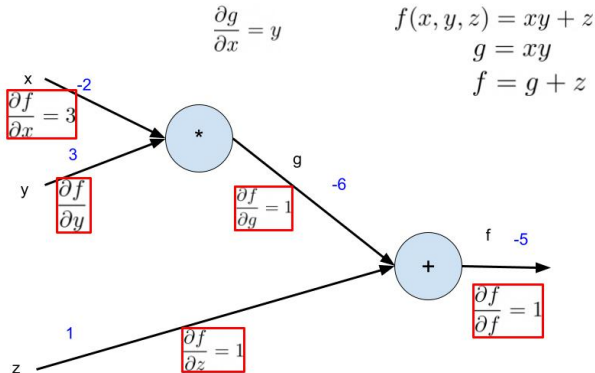
# E.g. Computational graph: **Backward pass**

- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



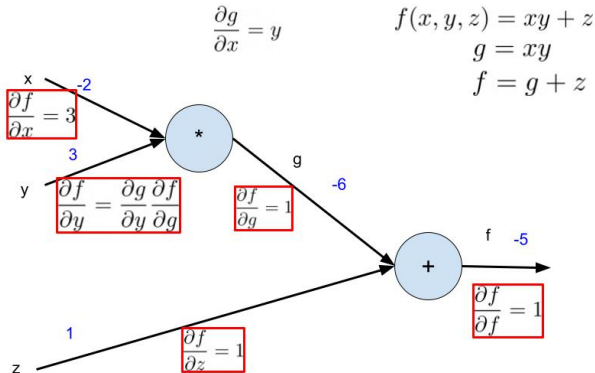
# E.g. Computational graph: **Backward pass**

- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



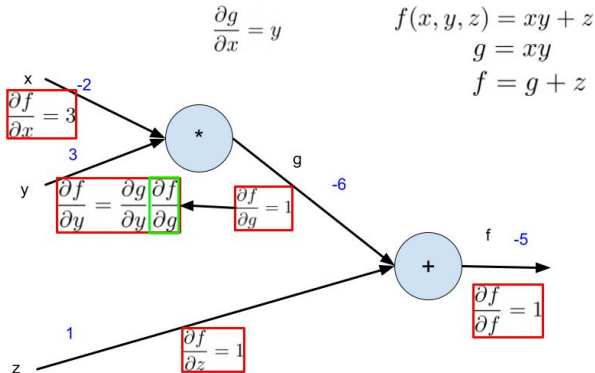
# E.g. Computational graph: **Backward pass**

- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



# E.g. Computational graph: **Backward pass**

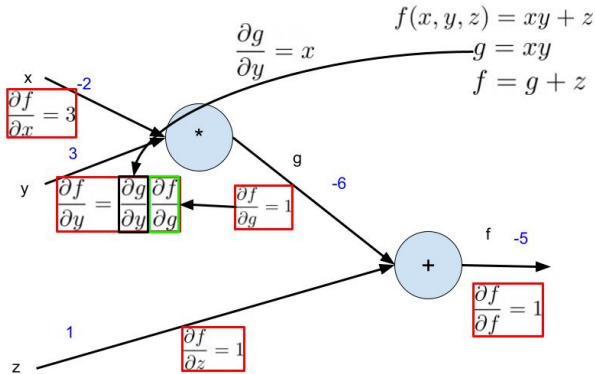
- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



# E.g. Computational graph: **Backward pass**



- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$

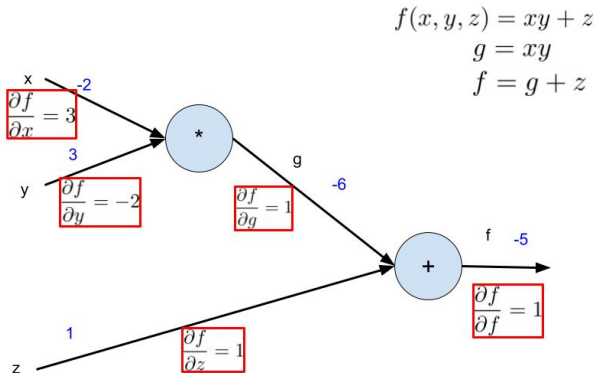




# E.g. Computational graph: Backward pass

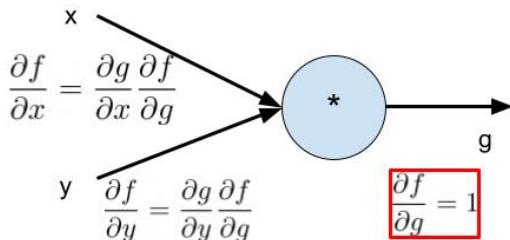


- Compute the derivatives  $\frac{\partial f}{\partial x}$ ,  $\frac{\partial f}{\partial y}$ ,  $\frac{\partial f}{\partial z}$



# Gradient flow

- down stream gradient = local gradient  $\times$  upstream gradient



$$g_d = g_l \times g_u$$

# References



- ① Cybenko G. 1989, [Approximation by superpositions of a sigmoidal function](#)