# CS 471/571(Fall 2023): Introduction to Artificial Intelligence

# Lecture 25: Neural Nets

Thanh H. Nguyen
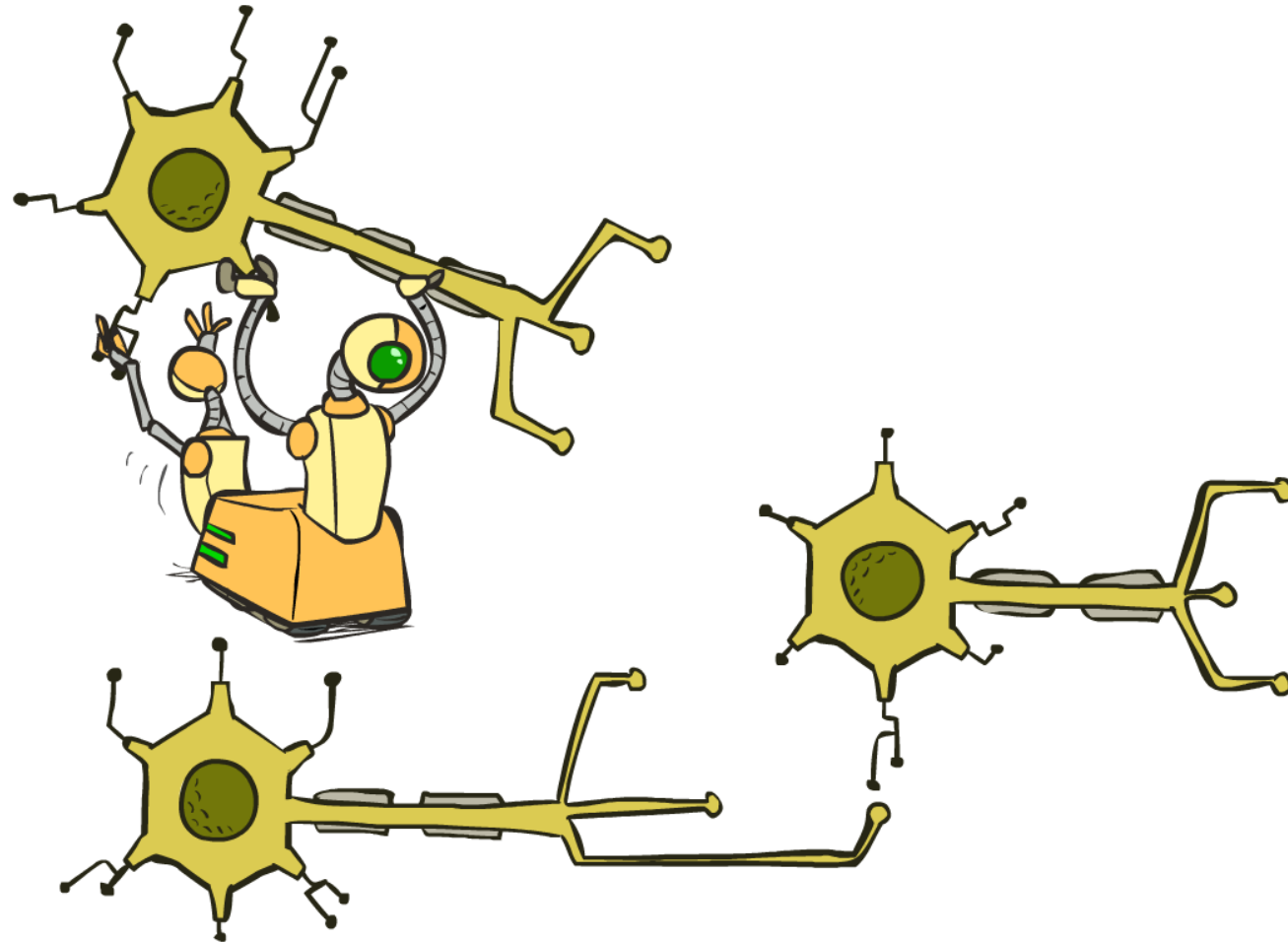
# Announcement & Reminder

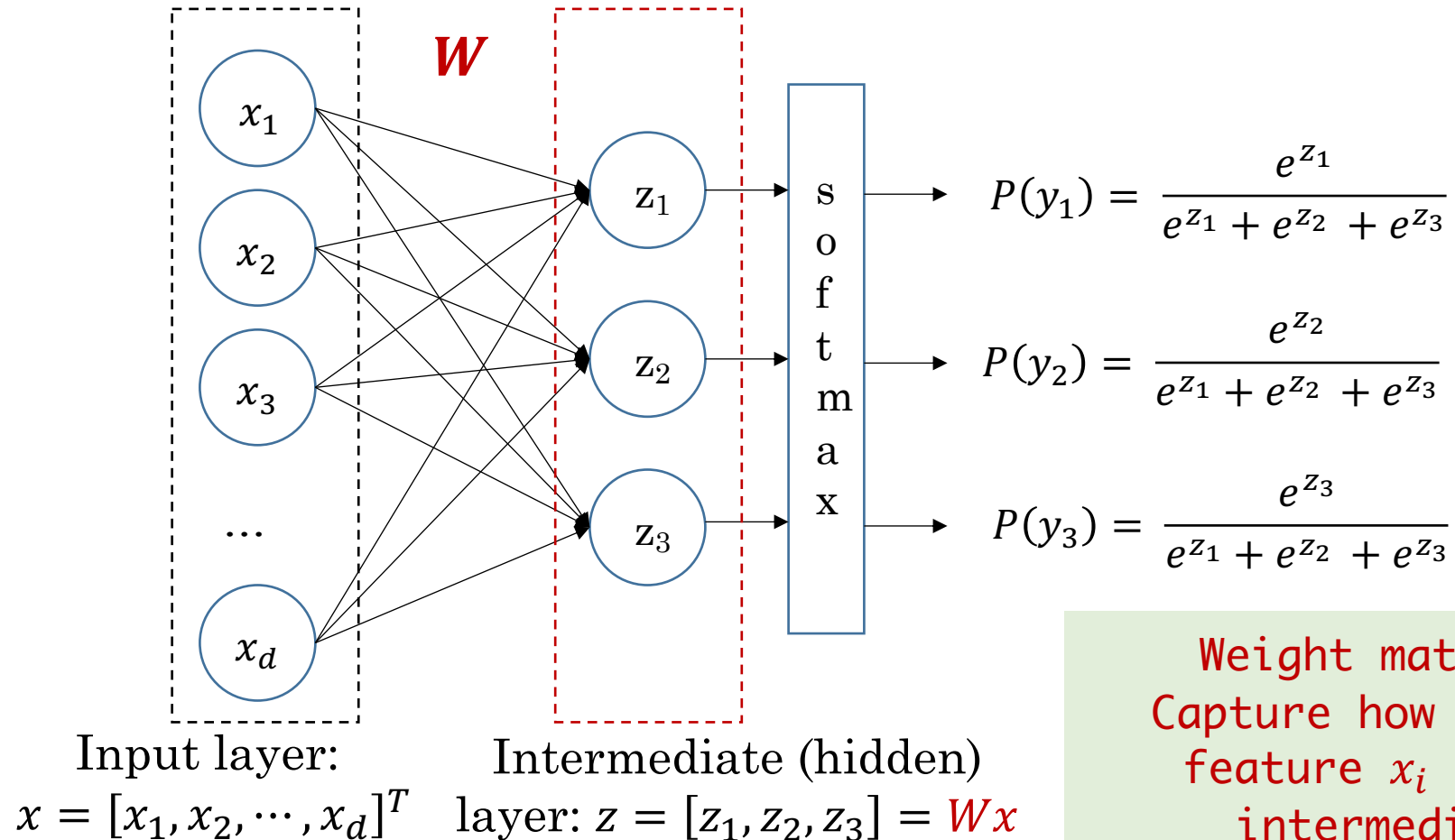- Exam Review
  - Friday, Dec 01$^{st}$, 2023

- Written assignment 4
  - Deadline: Wednesday, November 29$^{th}$, 2023.

- Student experience survey
  - Deadline: 06:00 AM on Monday, Dec 4$^{th}$, 2023
  - If >= 80% of students complete the survey, everyone will get an extra <u>2% credit</u> for your final grade
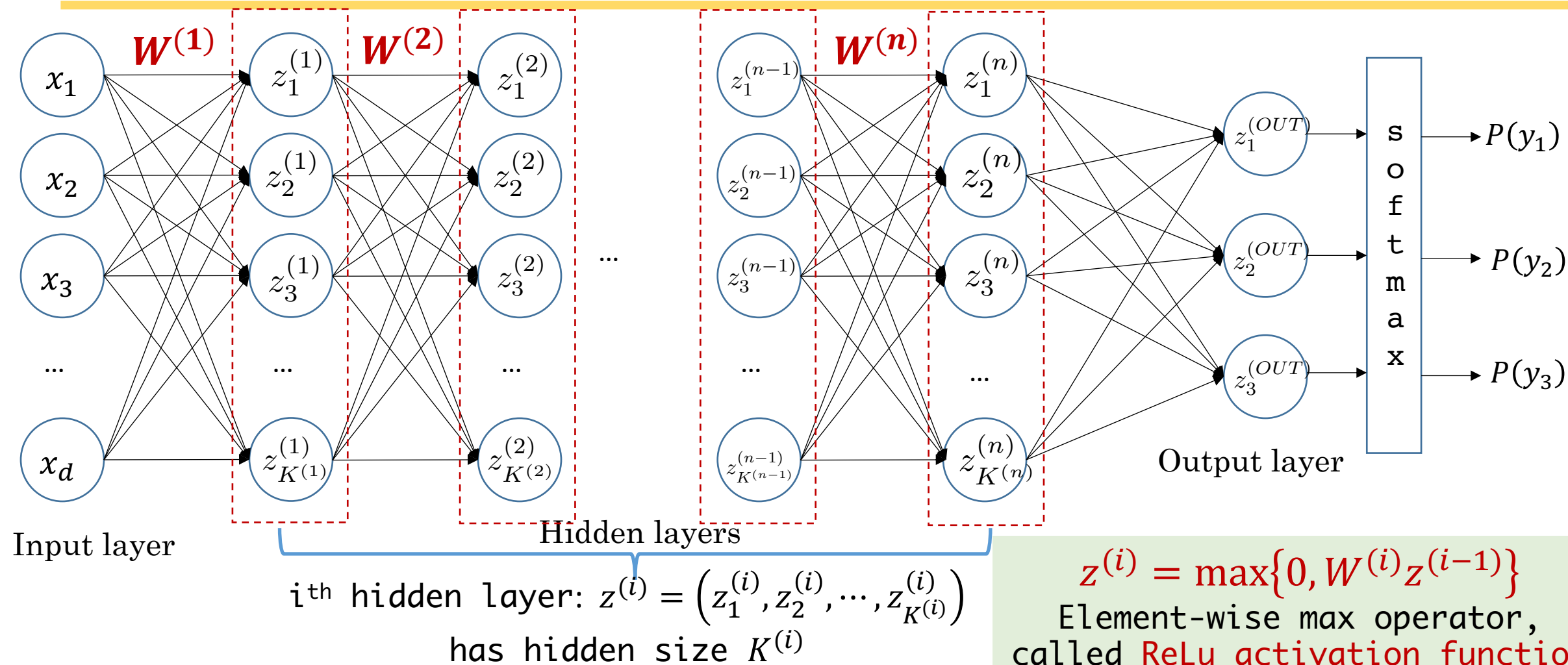
# Neural Networks

# Multi-class Logistic Regression

- = special case of neural network



$W$

$$P(y_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

Input layer:
$x = [x_1, x_2, \cdots, x_d]^T$

Intermediate (hidden)
layer: $z = [z_1, z_2, z_3] = Wx$

Weight matrix: $W = \{w_{ij}\}$
Capture how much each input
feature $x_i$ influence each
intermediate value $z_j$

# Deep Neural Network = Learn the Features!



$W^{(1)}$  $W^{(2)}$  $W^{(n)}$

Input layer

Hidden layers

Output layer

$i^{th}$ hidden layer: $z^{(i)} = \left( z_1^{(i)}, z_2^{(i)}, \cdots, z_{K^{(i)}}^{(i)} \right)$
has hidden size $K^{(i)}$

$$z^{(i)} = \max\{0, W^{(i)} z^{(i-1)}\}$$
Element-wise max operator,
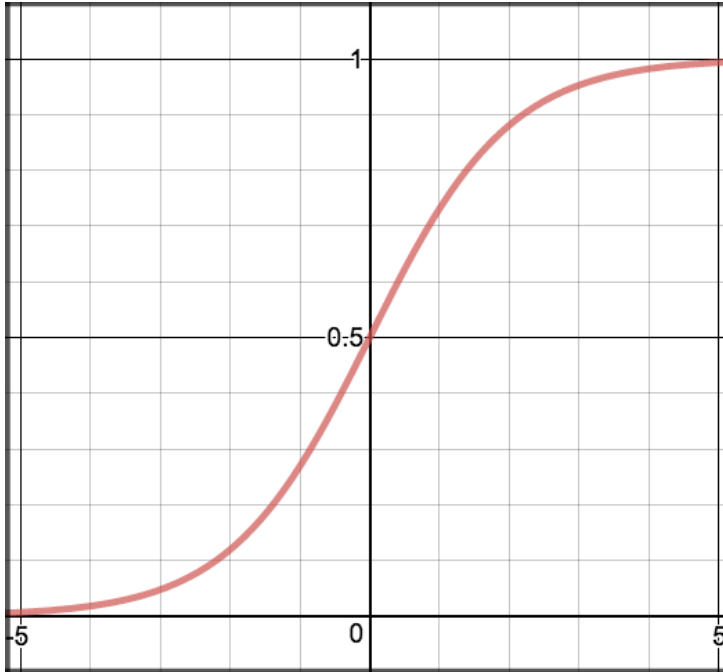called ReLu activation function

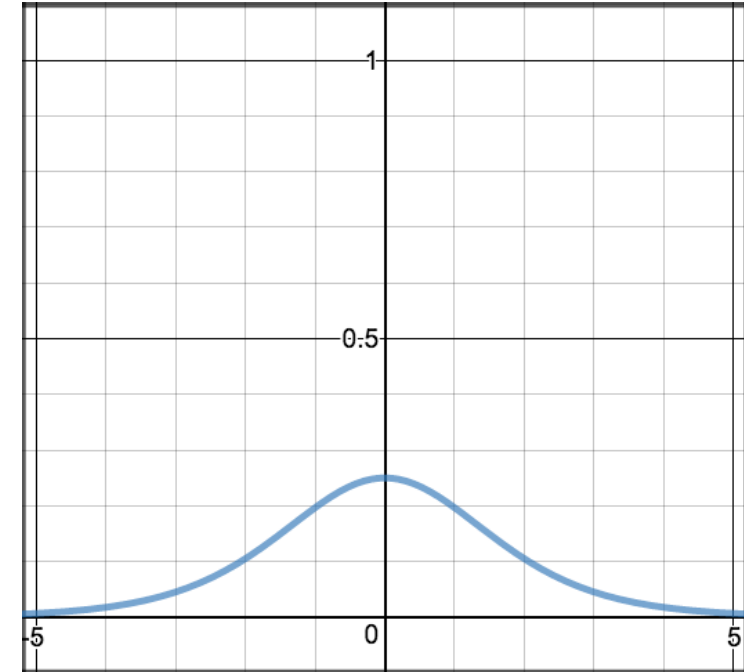# Common Activation Functions

Source: https://ml-cheatsheet.readthedocs.io/en/latest/activation_functions.html

# Sigmoid



Function: $z = \dfrac{1}{1 + e^{-x}}$



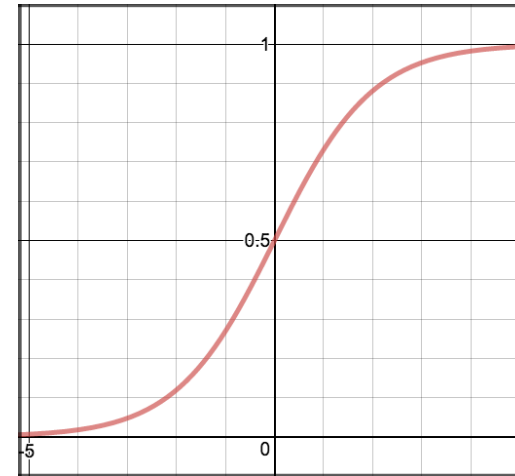Derivative: $\dfrac{dz}{dx} = z \cdot (1 - z)$

# Sigmoid

- **Pros**
  - Is nonlinear.
  - Has a smooth gradient.
  - Good for a classifier.
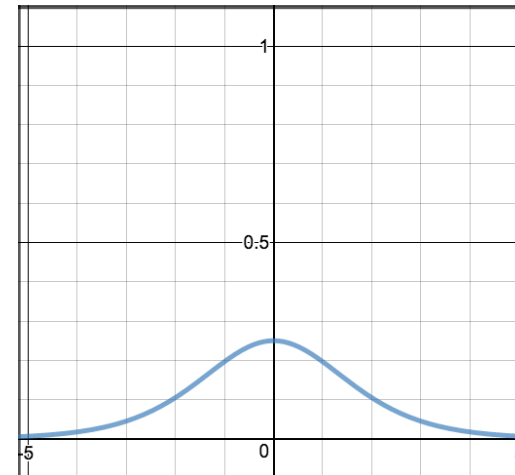  - The output is bounded within (0,1)
- **Cons**
  - Towards either end of the sigmoid function, the z values tend to respond very less to changes in x.
  - Gives rise to a problem of "vanishing gradients".
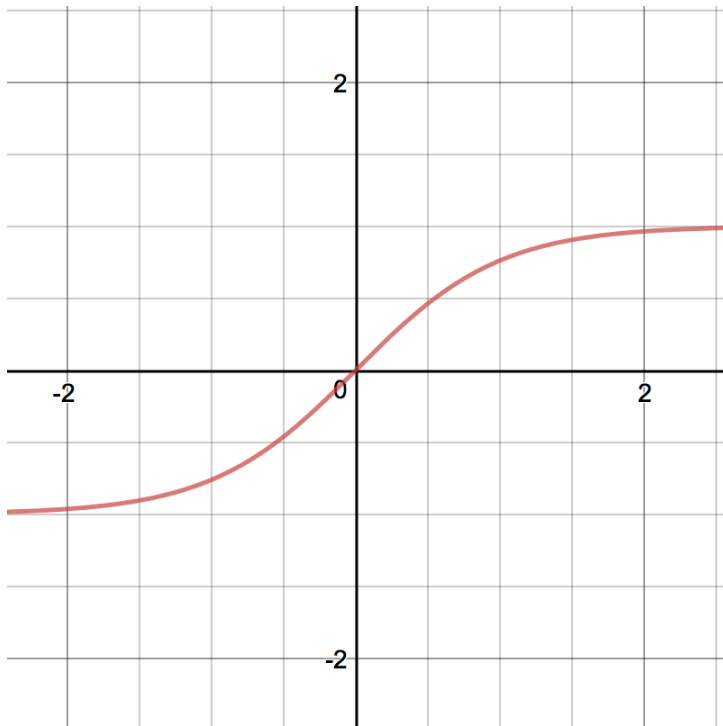
Function:
$$z = \frac{1}{1 + e^{-x}}$$

Derivative:
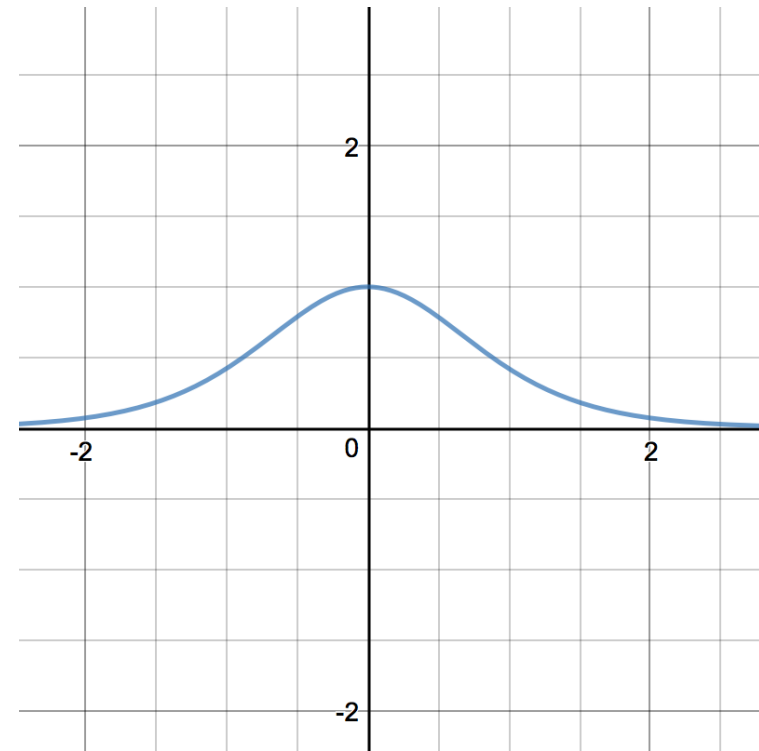$$\frac{dz}{dx} = z \cdot (1 - z)$$

# Tanh



Function: $z = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$



Derivative: $\dfrac{dz}{dx} = 1 - z^2$
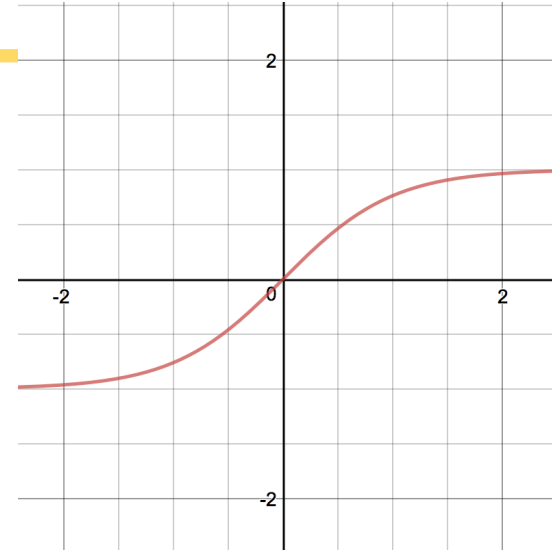
# Tanh

- **Pros**
  - The gradient is stronger for tanh than sigmoid ( derivatives are steeper).
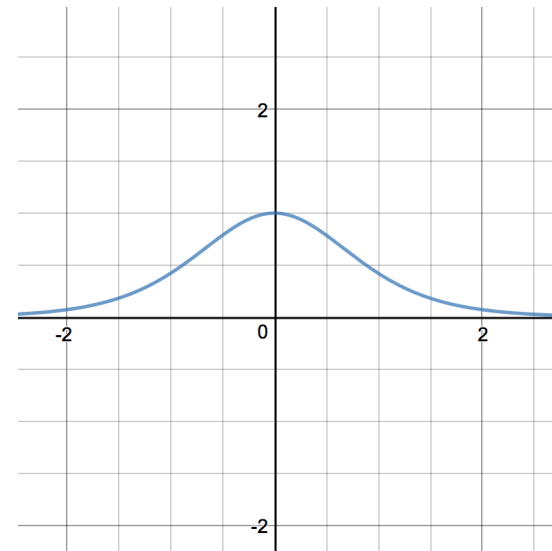
- **Cons**
  - Tanh also has the vanishing gradient problem.

Function:
$$z = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Derivative:
$$\frac{dz}{dx} = 1 - z^2$$

# ReLU



Function: $z = \max\{0, x\}$



Derivative: $\dfrac{dz}{dx} = \begin{cases} 1, \text{if } x > 0 \\ 0, \text{if } x < 0 \end{cases}$

# ReLU

- Pros
  - ReLu avoids and rectifies vanishing gradient problem.
  - ReLu is less computationally expensive

- Cons
  - ReLu should only be used within hidden layers.
  - Dying ReLu problem: When $(x < 0)$, gradient of ReLu is 0, meaning the weights will not get adjusted during gradient descent.
  - The range of ReLu is $[0, \infty)$. This means it can blow up the activation.

Function:
$$z = \max\{0, x\}$$

Derivative:
$$\frac{dz}{dx} = \begin{cases} 1, \text{if } x > 0 \\ 0, \text{if } x < 0 \end{cases}$$

# Train a Neural Net

- Train a neural networks: gradient descent

  - Forward pass: input is passed forward through the network to produce a prediction output. A loss metric is computed based on the difference between prediction and target (true output).

  - Backward pass: derivatives of this loss metric are calculated and propagated back through the network using a technique called backpropagation.

- We will talk more about this!!!

called neurons

input layer

hidden layer 1   hidden layer 2

output layer

# Common Loss Functions

# Cross Entropy

- Measures the performance of a classification model whose output is a probability value between 0 and 1

- Cross-entropy loss increases as the predicted probability diverges from the actual label.

- A perfect model would have a log loss of 0.



Log Loss when true label = 1

- Binary classification: $Loss = -(y \log p + (1-y) \log(1-p))$
- Multi-classification: $Loss = -\sum_l y_l \log p_l$
  - $y_l = 1$ if the true label is $l$, and $y_l = 0$, otherwise
  - $p_l$: predicted probability of having label $l$

# Mean Square Error (MSE)

- Commonly used for regression loss

- Measure the average of squared difference between predictions and actual observations

$$MSE = \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{n}$$

# Training a Neural Network

# Minimizing Error: Gradient Descent

- Main technique to minimize prediction error in neural network



- Challenge: computing gradient is non-trivial!!!

# (Bad Choice) Direct Computation

- Input:

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_d \end{bmatrix}$$



$x$

$\hat{y} = w^{(2)} \cdot h^{(1)}$

input layer

hidden layer

output layer

$h^{(1)} = \max(0, w^{(1)}x)$

Hing loss: $\max(0, 1 - \hat{y} \cdot y)$

- One hidden layer $h^{(1)} = \max(0, w^{(1)}x)$ with $K^{(1)}$ neurons
- Output layer $\hat{y} = w^{(2)} \cdot h^{(1)}$

$$w^{(1)} = \begin{bmatrix} w_{1,1}^{(1)} & \cdots & w_{1,K^{(1)}}^{(1)} \\ \vdots & \ddots & \vdots \\ w_{d,1}^{(1)} & \cdots & w_{d,K^{(1)}}^{(1)} \end{bmatrix}^T$$

$$w^{(2)} = \begin{bmatrix} w_1^{(2)} \\ w_2^{(2)} \\ \dots \\ w_{K^{(1)}}^{(2)} \end{bmatrix}$$

# (Bad Choice) Direct Computation

- Input data: $\left(x^{(1)}, y^{(1)}\right), \left(x^{(2)}, y^{(2)}\right), \cdots, \left(x^{(n)}, y^{(n)}\right)$

- Prediction (two-layer network): $\hat{y}^{(i)} = w^{(2)} \max\left(0, w^{(1)} x^{(i)}\right)$

- Total loss = prediction loss + regularization

$$L = \sum_i \max\left(0, 1 - \hat{y}^{(i)} \cdot y^{(i)}\right) + \lambda R\left(w^{(1)}\right) + \lambda R\left(w^{(2)}\right)$$

$$= \sum_i \max\left(0, 1 - \hat{y}^{(i)} w^{(2)} \max\left(0, w^{(1)} x^{(i)}\right)\right) + \lambda R\left(w^{(1)}\right) + \lambda R\left(w^{(2)}\right)$$

- Computing gradients: $\dfrac{dL}{dw^{(1)}}$ and $\dfrac{dL}{dw^{(2)}}$
  - So that we can learn $w^{(1)}$ and $w^{(2)}$

Direct computation of gradients is extremely inefficient!!!

# Computing Gradients

Computational Graphs
+
Backpropagation

# Computational Graphs

- The descriptive language of deep learning models

- Functional description of the required computation

- Can be instantiated to do two types of computation:

  - Forward computation

  - Backward computation

# Computational Graphs

- An acyclic directed graph

- Nodes: A node with an incoming edge is a function of that edge's tail node. This function can be basic arithmetic operations or any functions of which derivatives can be easily computed.

- Edges: An edge represents a function argument.

$x$

$z = x + y$

$+$

$y$

# Computing Gradients: Computational Graphs + Backpropagation

- Key ideas: applying gradient chain rule to unroll gradients through hidden layers in a neural nets

- Chain rule: $\frac{df(g(x))}{dx} = \frac{df(g(x))}{dg(x)} \cdot \frac{dg(x)}{dx}$

- Example: $f(u) = u^2 + 2u, \ g(x) = 3x + 1$

  - Direct computation: $f(g(x)) = g(x)^2 + 2g(x) = (3x+1)^2 + 2(3x+1) = 9x^2 + 12x + 3 \rightarrow \frac{df(g(x))}{dx} = 18x + 12$ <span style="color:red">(Inefficient)</span>

  - Chain rule: 

    <span style="color:red">(Much simpler)</span>

    $\frac{df(g(x))}{dx} = \frac{df(g(x))}{dg(x)} \cdot \frac{dg(x)}{dx} = (2g(x)+2)\times 3 = (6x+4)\times 3 = 18x + 12$

# Computational Graph + Backpropagation

$x$

Forward pass

$z$

$y$

Loss (L)

# Computational Graph + Backpropagation



Forward pass

$x$

$$\frac{dL}{dx} = \frac{dL}{dz} \cdot \frac{dz}{dx}$$

Downstream gradients

Local gradients

$$\frac{dz}{dx}$$

$$\frac{dz}{dy}$$

$z$

Loss (L)

$$\frac{dL}{dz}$$

Upstream gradient

$y$

$$\frac{dL}{dy} = \frac{dL}{dz} \cdot \frac{dz}{dy}$$

Backward pass

# Example

- $q = x + y$
- $f = q * z$

# Example

- $q = x + y$
- $f = q * z$

- Input: $x = -2, y = 5, z = -4$

- Goal: $\dfrac{df}{dx}, \dfrac{df}{dy}, \dfrac{df}{dz}$

# Example

- $q = x + y$
- $f = q * z$

- Input: $x = -2, y = 5, z = -4$

- Goal: $\frac{df}{dx}, \frac{df}{dy}, \frac{df}{dz}$

# Example

- $q = x + y$
- $f = q * z$

- Input: $x = -2, y = 5, z = -4$

- Goal: $\frac{df}{dx}, \frac{df}{dy}, \frac{df}{dz}$



$$\boxed{\frac{df}{df}}$$

# Example

- $q = x + y$
- $f = q * z$

- Input: $x = -2, y = 5, z = -4$

- Goal: $\dfrac{df}{dx}, \dfrac{df}{dy}, \dfrac{df}{dz}$



$x$ **−2**

$y$ **5**

$+$    $q$  **3**

$f$  **−12**

$*$

**1**

$z$ **−4**

**3**

$$\dfrac{df}{dz}$$

# Example

- $q = x + y$
- $f = q * z$

- Input: $x = -2, y = 5, z = -4$

- Goal: $\dfrac{df}{dx}, \dfrac{df}{dy}, \dfrac{df}{dz}$

# Example

- $q = x + y$
- $f = q * z$

- Input: $x = -2, y = 5, z = -4$

- Goal: $\frac{df}{dx}, \frac{df}{dy}, \frac{df}{dz}$



Chain rule

$$\frac{df}{dy} = \frac{df}{dq} \cdot \frac{dq}{dy}$$

Upstream gradient          Local gradient

# Example

- $q = x + y$
- $f = q * z$

- Input: $x = -2, y = 5, z = -4$

- Goal: $\dfrac{df}{dx}, \dfrac{df}{dy}, \dfrac{df}{dz}$



Chain rule

$$\frac{df}{dy} = \frac{df}{dq} \cdot \frac{dq}{dy}$$

Upstream gradient       Local gradient

$\dfrac{df}{dy}$
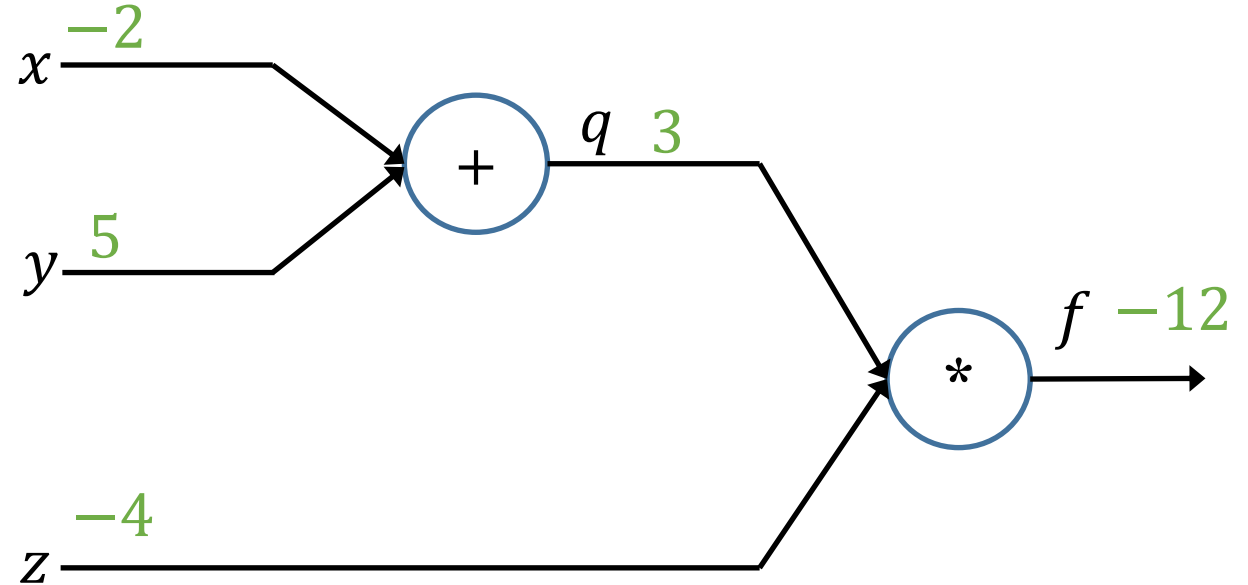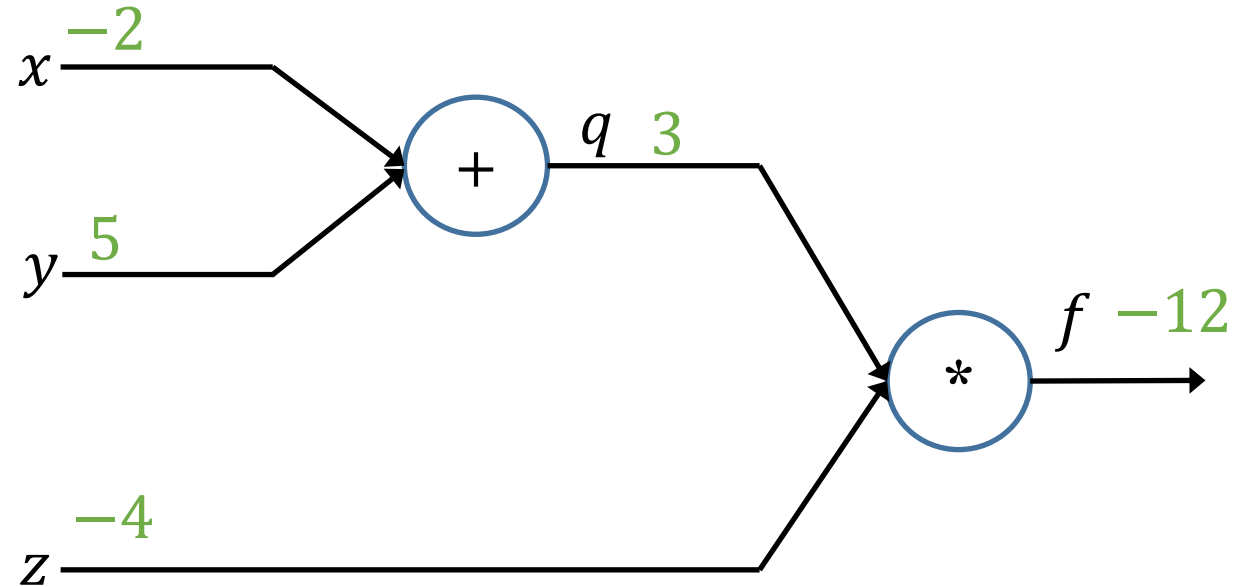
# Example

- $q = x + y$
- $f = q * z$

- Input: $x = -2, y = 5, z = -4$

- Goal: $\frac{df}{dx}, \frac{df}{dy}, \frac{df}{dz}$

$x$ —— $-2$
$-4$
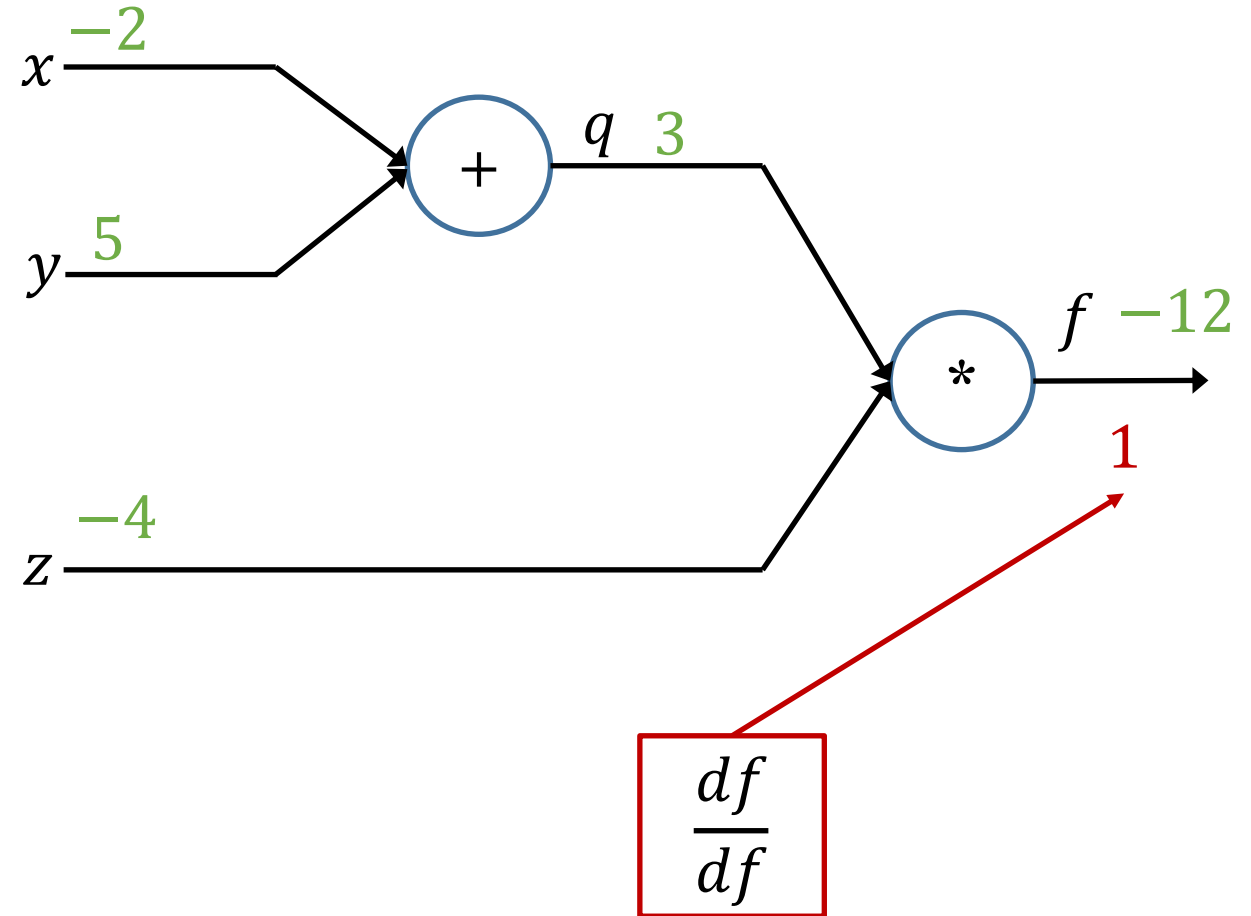
$y$ $5$
$-4$

$+$ → $q$ $3$
$-4$

$z$ $-4$
$3$

$*$ → $f$ $-12$
$1$

$\boxed{\frac{df}{dx}}$

Chain rule

$$\frac{df}{dx} = \frac{df}{dq} \cdot \frac{dq}{dx}$$

Upstream gradient       Local gradient

# Pytorch

- A Python-based scientific computing package of which goals are:
  - A replacement for NumPy to use the power of GPUs and other accelerators.
  - An automatic differentiation library that is useful to implement neural networks.

- GPU
  - A processor that has many smaller and more specialized cores
  - Has massive performance when a processing task can be divided up and processed across many cores.

Source: https://pytorch.org/

# Pytorch: Fundamental Concepts

- torch.Tensor: Tensors are the central data abstraction in PyTorch. Tensors are specialized data structure that are similar to arrays and matrices.

- torch.autograd: a built-in differentiation engine that supports automatic computation of gradient for any computational graph.

- torch.nn.Module: base class for all neural network modules.

- Pytorch tutorial:
  - https://pytorch.org/tutorials/beginner/pytorch_with_examples.html

# Pytorch: Example

Create data $\longrightarrow$

```python
1    import torch
2    import math
3
4    dtype = torch.float
5    device = torch.device("cpu")
6
7    x = torch.linspace(-math.pi, math.pi, 2000, device = device, dtype = dtype)
8    y = torch.sin(x)
9
10   a = torch.randn((), device = device, dtype = dtype, requires_grad=True)
11   b = torch.randn((), device = device, dtype = dtype, requires_grad=True)
12   c = torch.randn((), device = device, dtype = dtype, requires_grad=True)
13   d = torch.randn((), device = device, dtype = dtype, requires_grad=True)
14
15   learning_rate = 1e-6
16
17   for t in range(2000):
18       y_pred = a + b * x + c * x ** 2 + d * x ** 3
19       loss = (y_pred - y).pow(2).sum()
20
21       if t % 100 == 99:
22           print(t, loss.item())
23
24       loss.backward()
25       with torch.no_grad():
26           a -= learning_rate * a.grad
27           b -= learning_rate * b.grad
28           c -= learning_rate * c.grad
29           d -= learning_rate * d.grad
30
31           a.grad = None
32           b.grad = None
33           c.grad = None
34           d.grad = None
35
36   print(f'Result: y = {a.item()} + {b.item()} x + {c.item()} x^2 + {d.item()} x^3')
```

# Pytorch: Example

Create random weights →

```python
1    import torch
2    import math
3
4    dtype = torch.float
5    device = torch.device("cpu")
6
7    x = torch.linspace(-math.pi, math.pi, 2000, device = device, dtype = dtype)
8    y = torch.sin(x)
9
10   a = torch.randn((), device = device, dtype = dtype, requires_grad=True)
11   b = torch.randn((), device = device, dtype = dtype, requires_grad=True)
12   c = torch.randn((), device = device, dtype = dtype, requires_grad=True)
13   d = torch.randn((), device = device, dtype = dtype, requires_grad=True)
14
15   learning_rate = 1e-6
16
17   for t in range(2000):
18       y_pred = a + b * x + c * x ** 2 + d * x ** 3
19       loss = (y_pred - y).pow(2).sum()
20
21       if t % 100 == 99:
22           print(t, loss.item())
23
24       loss.backward()
25       with torch.no_grad():
26           a -= learning_rate * a.grad
27           b -= learning_rate * b.grad
28           c -= learning_rate * c.grad
29           d -= learning_rate * d.grad
30
31           a.grad = None
32           b.grad = None
33           c.grad = None
34           d.grad = None
35
36   print(f'Result: y = {a.item()} + {b.item()} x + {c.item()} x^2 + {d.item()} x^3')
```

# Pytorch: Example

Forward pass →

```python
1   import torch
2   import math
3
4   dtype = torch.float
5   device = torch.device("cpu")
6
7   x = torch.linspace(-math.pi, math.pi, 2000, device = device, dtype = dtype)
8   y = torch.sin(x)
9
10  a = torch.randn((), device = device, dtype = dtype, requires_grad=True)
11  b = torch.randn((), device = device, dtype = dtype, requires_grad=True)
12  c = torch.randn((), device = device, dtype = dtype, requires_grad=True)
13  d = torch.randn((), device = device, dtype = dtype, requires_grad=True)
14
15  learning_rate = 1e-6
16
17  for t in range(2000):
18      y_pred = a + b * x + c * x ** 2 + d * x ** 3
19      loss = (y_pred - y).pow(2).sum()
20
21      if t % 100 == 99:
22          print(t, loss.item())
23
24      loss.backward()
25      with torch.no_grad():
26          a -= learning_rate * a.grad
27          b -= learning_rate * b.grad
28          c -= learning_rate * c.grad
29          d -= learning_rate * d.grad
30
31          a.grad = None
32          b.grad = None
33          c.grad = None
34          d.grad = None
35
36  print(f'Result: y = {a.item()} + {b.item()} x + {c.item()} x^2 + {d.item()} x^3')
```

# Pytorch: Example

Backward pass →

```python
1   import torch
2   import math
3
4   dtype = torch.float
5   device = torch.device("cpu")
6
7   x = torch.linspace(-math.pi, math.pi, 2000, device = device, dtype = dtype)
8   y = torch.sin(x)
9
10  a = torch.randn((), device = device, dtype = dtype, requires_grad=True)
11  b = torch.randn((), device = device, dtype = dtype, requires_grad=True)
12  c = torch.randn((), device = device, dtype = dtype, requires_grad=True)
13  d = torch.randn((), device = device, dtype = dtype, requires_grad=True)
14
15  learning_rate = 1e-6
16
17  for t in range(2000):
18      y_pred = a + b * x + c * x ** 2 + d * x ** 3
19      loss = (y_pred - y).pow(2).sum()
20
21      if t % 100 == 99:
22          print(t, loss.item())
23
24      loss.backward()
25      with torch.no_grad():
26          a -= learning_rate * a.grad
27          b -= learning_rate * b.grad
28          c -= learning_rate * c.grad
29          d -= learning_rate * d.grad
30
31          a.grad = None
32          b.grad = None
33          c.grad = None
34          d.grad = None
35
36  print(f'Result: y = {a.item()} + {b.item()} x + {c.item()} x^2 + {d.item()} x^3')
```

Thanh H. Nguyen

# Pytorch: Example

```
1   import torch
2   import math
3
4   dtype = torch.float
5   device = torch.device("cpu")
6
7   x = torch.linspace(-math.pi, math.pi, 2000, device = device, dtype = dtype)
8   y = torch.sin(x)
9
10  a = torch.randn((), device = device, dtype = dtype, requires_grad=True)
11  b = torch.randn((), device = device, dtype = dtype, requires_grad=True)
12  c = torch.randn((), device = device, dtype = dtype, requires_grad=True)
13  d = torch.randn((), device = device, dtype = dtype, requires_grad=True)
14
15  learning_rate = 1e-6
16
17  for t in range(2000):
18      y_pred = a + b * x + c * x ** 2 + d * x ** 3
19      loss = (y_pred - y).pow(2).sum()
20
21      if t % 100 == 99:
22          print(t, loss.item())
23
24      loss.backward()
25      with torch.no_grad():
26          a -= learning_rate * a.grad
27          b -= learning_rate * b.grad
28          c -= learning_rate * c.grad
29          d -= learning_rate * d.grad
30
31          a.grad = None
32          b.grad = None
33          c.grad = None
34          d.grad = None
35
36      print(f'Result: y = {a.item()} + {b.item()} x + {c.item()} x^2 + {d.item()} x^3')
```

Gradient descent

# Pytorch: Example

```
1    import torch
2    import math
3
4    dtype = torch.float
5    device = torch.device("cpu")
6
7    x = torch.linspace(-math.pi, math.pi, 2000, device = device, dtype = dtype)
8    y = torch.sin(x)
9
10   a = torch.randn((), device = device, dtype = dtype, requires_grad=True)
11   b = torch.randn((), device = device, dtype = dtype, requires_grad=True)
12   c = torch.randn((), device = device, dtype = dtype, requires_grad=True)
13   d = torch.randn((), device = device, dtype = dtype, requires_grad=True)
14
15   learning_rate = 1e-6
16
17   for t in range(2000):
18       y_pred = a + b * x + c * x ** 2 + d * x ** 3
19       loss = (y_pred - y).pow(2).sum()
20
21       if t % 100 == 99:
22           print(t, loss.item())
23
24       loss.backward()
25       with torch.no_grad():
26           a -= learning_rate * a.grad
27           b -= learning_rate * b.grad
28           c -= learning_rate * c.grad
29           d -= learning_rate * d.grad
30
31           a.grad = None
32           b.grad = None
33           c.grad = None
34           d.grad = None
35
36   print(f'Result: y = {a.item()} + {b.item()} x + {c.item()} x^2 + {d.item()} x^3')
```

Reset gradient before next round of forward-backward pass

# Pytorch: Neural Network

Create data $\longrightarrow$

```
1   import torch
2   import math
3
4   x = torch.linspace(-math.pi, math.pi, 2000)
5   y = torch.sin(x)
6
7   p = torch.tensor([1, 2, 3])
8   xx = x.unsqueeze(-1).pow(p)
9
10  model = torch.nn.Sequential(
11      torch.nn.Linear(3, 1),
12      torch.nn.Flatten(0, 1)
13  )
14
15  loss_fn = torch.nn.MSELoss(reduction='sum')
16
17  learning_rate = 1e-6
18  for t in range(2000):
19      y_pred = model(xx)
20
21      loss = loss_fn(y_pred, y)
22      if t % 100 == 99:
23          print(t, loss.item())
24      model.zero_grad()
25
26      loss.backward()
27
28      with torch.no_grad():
29          for param in model.parameters():
30              param -= learning_rate * param.grad
31
32  linear_layer = model[0]
```

# Pytorch: Neural Network

Create predictive
model

```python
import torch
import math

x = torch.linspace(-math.pi, math.pi, 2000)
y = torch.sin(x)

p = torch.tensor([1, 2, 3])
xx = x.unsqueeze(-1).pow(p)

model = torch.nn.Sequential(
    torch.nn.Linear(3, 1),
    torch.nn.Flatten(0, 1)
)

loss_fn = torch.nn.MSELoss(reduction='sum')

learning_rate = 1e-6
for t in range(2000):
    y_pred = model(xx)

    loss = loss_fn(y_pred, y)
    if t % 100 == 99:
        print(t, loss.item())
    model.zero_grad()

    loss.backward()

    with torch.no_grad():
        for param in model.parameters():
            param -= learning_rate * param.grad

linear_layer = model[0]
```

# Pytorch: Neural Network

Define loss
function →

```python
1   import torch
2   import math
3
4   x = torch.linspace(-math.pi, math.pi, 2000)
5   y = torch.sin(x)
6
7   p = torch.tensor([1, 2, 3])
8   xx = x.unsqueeze(-1).pow(p)
9
10  model = torch.nn.Sequential(
11      torch.nn.Linear(3, 1),
12      torch.nn.Flatten(0, 1)
13  )
14
15  loss_fn = torch.nn.MSELoss(reduction='sum')
16
17  learning_rate = 1e-6
18  for t in range(2000):
19      y_pred = model(xx)
20
21      loss = loss_fn(y_pred, y)
22      if t % 100 == 99:
23          print(t, loss.item())
24      model.zero_grad()
25
26      loss.backward()
27
28      with torch.no_grad():
29          for param in model.parameters():
30              param -= learning_rate * param.grad
31
32  linear_layer = model[0]
```

# Pytorch: Neural Network

```
1    import torch
2    import math
3
4    x = torch.linspace(-math.pi, math.pi, 2000)
5    y = torch.sin(x)
6
7    p = torch.tensor([1, 2, 3])
8    xx = x.unsqueeze(-1).pow(p)
9
10   model = torch.nn.Sequential(
11       torch.nn.Linear(3, 1),
12       torch.nn.Flatten(0, 1)
13   )
14
15   loss_fn = torch.nn.MSELoss(reduction='sum')
16
17   learning_rate = 1e-6
18   for t in range(2000):
19       y_pred = model(xx)
20
21       loss = loss_fn(y_pred, y)
22       if t % 100 == 99:
23           print(t, loss.item())
24       model.zero_grad()
25
26       loss.backward()
27
28       with torch.no_grad():
29           for param in model.parameters():
30               param -= learning_rate * param.grad
31
32   linear_layer = model[0]
```

Forward pass

# Pytorch: Neural Network

```
1    import torch
2    import math
3
4    x = torch.linspace(-math.pi, math.pi, 2000)
5    y = torch.sin(x)
6
7    p = torch.tensor([1, 2, 3])
8    xx = x.unsqueeze(-1).pow(p)
9
10   model = torch.nn.Sequential(
11       torch.nn.Linear(3, 1),
12       torch.nn.Flatten(0, 1)
13   )
14
15   loss_fn = torch.nn.MSELoss(reduction='sum')
16
17   learning_rate = 1e-6
18   for t in range(2000):
19       y_pred = model(xx)
20
21       loss = loss_fn(y_pred, y)
22       if t % 100 == 99:
23           print(t, loss.item())
24       model.zero_grad()
25
26       loss.backward()
27
28       with torch.no_grad():
29           for param in model.parameters():
30               param -= learning_rate * param.grad
31
32   linear_layer = model[0]
```

Reset gradient and
run backward pass

# Pytorch: Neural Network

```python
1   import torch
2   import math
3
4   x = torch.linspace(-math.pi, math.pi, 2000)
5   y = torch.sin(x)
6
7   p = torch.tensor([1, 2, 3])
8   xx = x.unsqueeze(-1).pow(p)
9
10  model = torch.nn.Sequential(
11      torch.nn.Linear(3, 1),
12      torch.nn.Flatten(0, 1)
13  )
14
15  loss_fn = torch.nn.MSELoss(reduction='sum')
16
17  learning_rate = 1e-6
18  for t in range(2000):
19      y_pred = model(xx)
20
21      loss = loss_fn(y_pred, y)
22      if t % 100 == 99:
23          print(t, loss.item())
24      model.zero_grad()
25
26      loss.backward()
27
28      with torch.no_grad():
29          for param in model.parameters():
30              param -= learning_rate * param.grad
31
32  linear_layer = model[0]
```

Run gradient descent

# Pytorch: Optimizer

Choose an optimizer →

```
1   import torch
2   import math
3
4   x = torch.linspace(-math.pi, math.pi, 2000)
5   y = torch.sin(x)
6
7   p = torch.tensor([1, 2, 3])
8   xx = x.unsqueeze(-1).pow(p)
9
10  model = torch.nn.Sequential(
11      torch.nn.Linear(3, 1),
12      torch.nn.Flatten(0, 1)
13  )
14
15  loss_fn = torch.nn.MSELoss(reduction='sum')
16
17  learning_rate = 1e-3
18  optimizer = torch.optim.RMSprop(model.parameters(), lr=learning_rate)
19  for t in range(2000):
20
21      y_pred = model(xx)
22      loss = loss_fn(y_pred, y)
23
24      if t % 100 == 99:
25          print(t, loss.item())
26      model.zero_grad()
27
28      loss.backward()
29
30      optimizer.step()
31
```

# Pytorch: Optimizer

```python
1   import torch
2   import math
3
4   x = torch.linspace(-math.pi, math.pi, 2000)
5   y = torch.sin(x)
6
7   p = torch.tensor([1, 2, 3])
8   xx = x.unsqueeze(-1).pow(p)
9
10  model = torch.nn.Sequential(
11      torch.nn.Linear(3, 1),
12      torch.nn.Flatten(0, 1)
13  )
14
15  loss_fn = torch.nn.MSELoss(reduction='sum')
16
17  learning_rate = 1e-3
18  optimizer = torch.optim.RMSprop(model.parameters(), lr=learning_rate)
19  for t in range(2000):
20
21      y_pred = model(xx)
22      loss = loss_fn(y_pred, y)
23
24      if t % 100 == 99:
25          print(t, loss.item())
26      model.zero_grad()
27
28      loss.backward()
29
30      optimizer.step()
31
```

Update weights →

# Pytorch: Define a New Module

Define a neural network model as a module →

```python
import torch


class Net(torch.nn.Module):
    def __init__(self, d_in, d_hidden, d_out):
        super(Net, self).__init__()
        self.ln1 = torch.nn.Linear(d_in, d_hidden)
        self.ln2 = torch.nn.Linear(d_hidden, d_out)

    def forward(self, x):
        h_relu = self.ln1(x).clamp(min=0)
        y_pred = self.ln2(h_relu)
        return y_pred


train_size, d_in, d_hidden, d_out = 64, 1000, 100, 10
x = torch.randn(train_size, d_in)
y = torch.randn(train_size, d_out)

model = Net(d_in, d_hidden, d_out)
optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
loss_fn = torch.nn.MSELoss(reduction='sum')

for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)

    if t % 50 == 49:
        print(t, loss.item())

    model.zero_grad()
    loss.backward()

    optimizer.step()
```

# Pytorch: Define a New Module

Construct and train a model →

```python
1   import torch
2
3
4   class Net(torch.nn.Module):
5       def __init__(self, d_in, d_hidden, d_out):
6           super(Net, self).__init__()
7           self.ln1 = torch.nn.Linear(d_in, d_hidden)
8           self.ln2 = torch.nn.Linear(d_hidden, d_out)
9
10      def forward(self, x):
11          h_relu = self.ln1(x).clamp(min=0)
12          y_pred = self.ln2(h_relu)
13          return y_pred
14
15
16  train_size, d_in, d_hidden, d_out = 64, 1000, 100, 10
17  x = torch.randn(train_size, d_in)
18  y = torch.randn(train_size, d_out)
19
20  model = Net(d_in, d_hidden, d_out)
21  optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
22  loss_fn = torch.nn.MSELoss(reduction='sum')
23
24  for t in range(500):
25      y_pred = model(x)
26      loss = loss_fn(y_pred, y)
27
28      if t % 50 == 49:
29          print(t, loss.item())
30
31      model.zero_grad()
32      loss.backward()
33
34      optimizer.step()
35
```