

Tensor (aka multidimensional array) manipulation

Lecture 05 — CS 577 Deep Learning

Instructor: Yutong Wang

Computer Science
Illinois Institute of Technology

September 18, 2024

Multiclass classification

Let $i = 1, \dots, N$ (the sample index)

- Training samples $\mathbf{x}^{(i)} \in \mathcal{X} \subseteq \mathbb{R}^d$
- Labels $y^{(i)} \in \mathcal{Y} = \{1, \dots, K\}$
- $f(\cdot; \boldsymbol{\theta}) : \mathcal{X} \rightarrow \mathcal{Y}$

Multiclass classification

```
1 import numpy as np
2 from sklearn.datasets import load_iris
3
4 n_classes = 3 # Number of classes
5
6 X, y = load_iris(return_X_y = True)
7 n = X.shape[0]
8 X = X[:, :2] # Keep only the first two features
9 X = X - np.mean(X, axis=0) # Center the data
```

Note: Shape of data

1-layer neural network

To define $f(\cdot; \boldsymbol{\theta}) : \mathcal{X} \rightarrow \mathcal{Y}$

$$\mathbf{h}^{(i)} = \mathbf{W}^{(1)\top} \mathbf{x}^{(i)} + \mathbf{b}^{(1)} \quad \text{and} \quad \mathbf{z}^{(i)} = \mathbf{W}^{(2)\top} g(\mathbf{h}^{(i)}) + \mathbf{b}^{(2)}$$

```
1 def init_params(input_dim, hidden_dim, output_dim):  
2     np.random.seed(0) # Ensure reproducibility  
3     W1 = np.random.randn(input_dim, hidden_dim)  
4     b1 = np.random.randn(hidden_dim)  
5     W2 = np.random.randn(hidden_dim, output_dim)  
6     b2 = np.random.randn(output_dim)  
7     return W1, b1, W2, b2
```

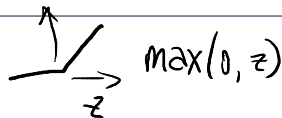
1-layer neural network

$$\mathbf{h}^{(i)} = \mathbf{W}^{(1)\top} \mathbf{x}^{(i)} + \mathbf{b}^{(1)} \quad \text{and} \quad \mathbf{z}^{(i)} = \mathbf{W}^{(2)\top} g(\mathbf{h}^{(i)}) + \mathbf{b}^{(2)}$$

```
1 input_dim = 2           # Number of input dimension/features
2 hidden_dim = 10         # Number of hidden neurons
3 output_dim = 3          # Number of classes
4
5 W1, b1, W2, b2 = init_params(input_dim, hidden_dim, output_dim)
6
7 theta = {
8     "W1": W1,
9     "b1": b1,
10    "W2": W2,
11    "b2": b2
12 }
```

pack & unpack

Relu



$$\mathbf{h}^{(i)} = \mathbf{W}^{(1)\top} \mathbf{x}^{(i)} + \mathbf{b}^{(1)} \quad \text{and} \quad \mathbf{z}^{(i)} = \mathbf{W}^{(2)\top} g(\mathbf{h}^{(i)}) + \mathbf{b}^{(2)}$$

1 `def relu(z):`

2 `return np.maximum(0, z)`

number



tensor

1 input



2 inputs (or more)



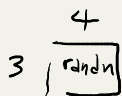
Note: What is the difference between `np.max` vs `np.maximum`?

↑ ↑
(also in pytorch)

Pairwise vs Reduction operations

- `np.maximum` is a type of pairwise operation like `+` and `*`

```
1 a = 12*np.random.randn(3,4)
2 b = np.arange(12).reshape((3,4))
3 # array([[ 0.,  1.,  2.,  3.],
4 #        [ 4.,  5.,  6.,  7.],
5 #        [ 8.,  9., 10., 11.]])
6 np.maximum(a,b)
7 # array([[ 0.,  3.88881425,  2.,  3.],
8 #        [ 4., 23.30701124,  7.05094358,  7.],
9 #        [ 8.,  9., 10., 12.53350356]])
```



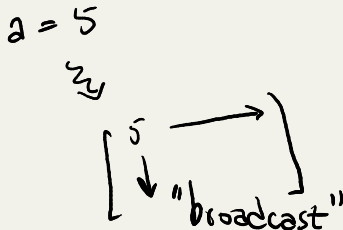
- `np.max` is a type of reduction operation like `np.sum` and `np.mean`

```
1 np.max(a,b) # will throw an error
```

Pairwise operations between tensors

- `np.maximum` is a type of pairwise operation like `+` and `*`

```
1 a = 5
2 b = np.arange(12).reshape((3,4))
3 # array([[ 0,  1,  2,  3],
4 #        [ 4,  5,  6,  7],
5 #        [ 8,  9, 10, 11]])
6 np.maximum(a,b)
7 # array([[ 5,  5,  5,  5],
8 #        [ 5,  5,  6,  7],
9 #        [ 8,  9, 10, 11]])
```

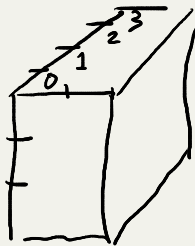


What is a **tensor** (aka multidimensional array) anyways?

24

↑ multi table of number

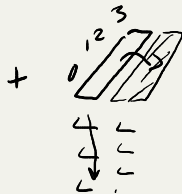
```
1 a = np.arange(3*2*4).reshape([3,2,4])
2
3
4 # array([[[ 0,  1,  2,  3],
5 #         [ 4,  5,  6,  7]],
6 #
7 #        [[ 8,  9, 10, 11],
8 #         [12, 13, 14, 15]],
9 #
10 #       [[16, 17, 18, 19],
11 #        [20, 21, 22, 23]])
```



What is broadcasting?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(4)
3 a + b
4 # array([[[ 0,  2,  4,  6],
5 #          [ 4,  6,  8, 10]],
6 #
7 #          [[ 8, 10, 12, 14],
8 #          [12, 14, 16, 18]],
9 #
10 #          [[16, 18, 20, 22],
11 #          [20, 22, 24, 26]]])
```

$[3, 2, 4]$
 $[4]$



Broadcast:

check if
one of the 2 dim
is = 1.
otherwise not compatible

Starts at the end
▷ if dim matches
then entrywise
▷ if no match

What is broadcasting?

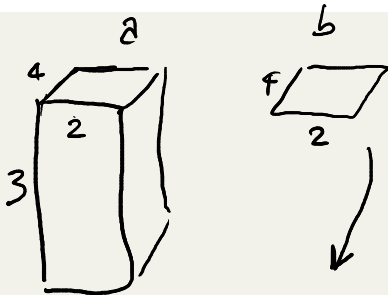
```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(4)
3 a + b
4 # array([[[ 0,  2,  4,  6],
5 #          [ 4,  6,  8, 10]],
6 #
7 #          [[ 8, 10, 12, 14],
8 #          [12, 14, 16, 18]],
9 #
10 #          [[16, 18, 20, 22],
11 #           [20, 22, 24, 26]]])
```

What is a tensor (aka multidimensional array) anyways?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2
3
4 # array([[[ 0,  1,  2,  3],
5 #         [ 4,  5,  6,  7]],
6 #
7 #        [[ 8,  9, 10, 11],
8 #         [12, 13, 14, 15]],
9 #
10 #        [[16, 17, 18, 19],
11 #         [20, 21, 22, 23]])
```

What is broadcasting?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(2*4).reshape([2,4])
3 a + b
4 # array([[[ 0,  2,  4,  6],
5 #         [ 8, 10, 12, 14]],
6 #
7 #        [[ 8, 10, 12, 14],
8 #         [16, 18, 20, 22]],
9 #
10 #       [[16, 18, 20, 22],
11 #        [24, 26, 28, 30]]])
```



What is broadcasting?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(3*4).reshape([3,4])
3 a + b
4 # error
5 #
6 #
7 #
8 #
9 #
10 #
11 #
```

X

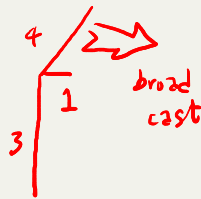
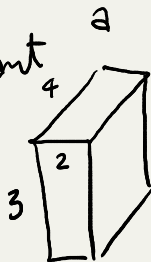
neither

2 or 3
is = 1.
Incompatible

What is broadcasting?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(3*4).reshape([3,1,4])
3 a + b
4 # array([[[ 0,  2,  4,  6],
5 #         [ 4,  6,  8, 10]],
6 #
7 #        [[12, 14, 16, 18],
8 #         [16, 18, 20, 22]],
9 #
10 #       [[24, 26, 28, 30],
11 #        [28, 30, 32, 34]]])
```

redundant



What is a tensor (aka multidimensional array) anyways?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2
3
4 # array([[[ 0,  1,  2,  3],
5 #         [ 4,  5,  6,  7]],
6 #
7 #        [[ 8,  9, 10, 11],
8 #         [12, 13, 14, 15]],
9 #
10 #       [[16, 17, 18, 19],
11 #        [20, 21, 22, 23]])
```


What is broadcasting?

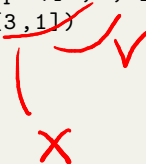
```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(3).reshape([3])
3 a + b
4 # ?
5 #
6 #
7 #
8 #
9 #
10 #
11 #
```

Handwritten annotations:

- A red "no" is written next to line 4.
- A red "X" is written next to line 3, with a red arc connecting it to the `[3]` in line 2.

What is broadcasting?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(3).reshape([3,1])
3 a + b
4 # ?
5 #
6 #
7 #
8 #
9 #
10 #
11 #
```



What is broadcasting?

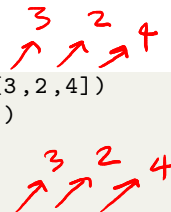
```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(3).reshape([3,1,1])
3 a + b
4 # ?
5 #
6 #
7 #
8 #
9 #
10 #
11 #
```

What is broadcasting?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(3).reshape([3,1,1])
3 a + b
4 # array([[[ 0,  1,  2,  3],
5 #          [ 4,  5,  6,  7]],
6 #
7 #          [[ 9, 10, 11, 12],
8 #          [13, 14, 15, 16]],
9 #
10 #          [[18, 19, 20, 21],
11 #          [22, 23, 24, 25]]])
```

Which of these two works?

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 b = np.arange(2).reshape([2,1])
3 a + b
4 # ?
5 #
6 a = np.arange(3*2*4).reshape([3,2,4])
7 b = np.arange(2).reshape([1,2,1])
8 a + b
9 # ?
10 #
```



Broadcasting rule applies

Note: Back to `np.maximum...`

- `np.maximum` is a type of pairwise operation like `+` and `*`

```
1 a = np.array([5,5,6,6])
2 b = np.arange(12).reshape((3,4))
3 # array([[ 0,  1,  2,  3],
4 #        [ 4,  5,  6,  7],
5 #        [ 8,  9, 10, 11]])
6 np.maximum(a,b)
7 # array([[ 5,  5,  6,  6],
8 #        [ 5,  5,  6,  7],
9 #        [ 8,  9, 10, 11]])
```

← np 1 number

Broadcasting rule applies

Note: Back to `np.maximum...`

- `np.maximum` is a type of pairwise operation like `+` and `*`

```
1 a = np.array([2,6,9])
2 b = np.arange(12).reshape((3,4))
3 # array([[ 0,  1,  2,  3],
4 #        [ 4,  5,  6,  7],
5 #        [ 8,  9, 10, 11]])
6 np.maximum(a,b)
7 # ?
8 #
9 #
```

(3, 1)
(3)

Broadcasting rule applies

Note: Back to `np.maximum...`

- `np.maximum` is a type of pairwise operation like `+` and `*`

```
1 a = np.array([2,6,9]).reshape((3,1))
2 b = np.arange(12).reshape((3,4))
3 # array([[ 0,  1,  2,  3],
4 #        [ 4,  5,  6,  7],
5 #        [ 8,  9, 10, 11]])
6 np.maximum(a,b)
7 # array([[ 2,  2,  2,  3],
8 #        [ 6,  6,  6,  7],
9 #        [ 9,  9, 10, 11]])
```

add fake dim

Broadcasting rule applies

Note: Back to `np.maximum...`

- `np.maximum` is a type of pairwise operation like `+` and `*`

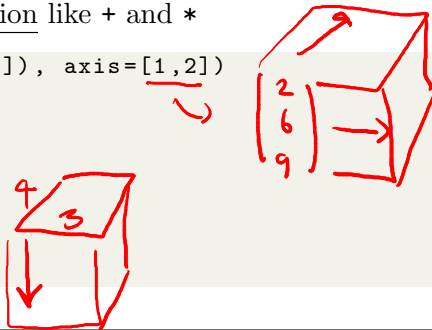
```
1 a = np.expand_dims(np.array([2,6,9]), axis=1)
2 b = np.arange(12).reshape((3,4))
3 # array([[ 0,  1,  2,  3],
4 #        [ 4,  5,  6,  7],
5 #        [ 8,  9, 10, 11]])
6 np.maximum(a,b)
7 # array([[ 2,  2,  2,  3],
8 #        [ 6,  6,  6,  7],
9 #        [ 9,  9, 10, 11]])
```

(3)
↑
0-th axis

Broadcasting rule applies

- `np.maximum` is a type of pairwise operation like `+` and `*`

```
1 a = np.expand_dims(np.array([2,6,9]), axis=[1,2])
2 # a.shape == (3,1,1)
3 b = np.arange(12).reshape((3,4))
4 # array([[ 0,  1,  2,  3],
5 #        [ 4,  5,  6,  7],
6 #        [ 8,  9, 10, 11]])
7 np.maximum(a,b)
8 # ???
```



Note: Breakout session: explain the output.

Broadcasting rule applies

- `np.maximum` is a type of pairwise operation like `+` and `*`

```
1 a = np.expand_dims(np.array([2,6,9]), axis=[1,2])
2 # a.shape == (3,1,1)
3 b = np.arange(12).reshape((3,4))
4 # array([[ 0,  1,  2,  3],
5 #        [ 4,  5,  6,  7],
6 #        [ 8,  9, 10, 11]])
7 np.maximum(a,b)
8 # ???
```

Note: Breakout session: explain the output.

Broadcasting rule applies

- `np.maximum` is a type of pairwise operation like `+` and `*`

```
1 a = np.expand_dims(np.array([2,6,9]), axis=[1,2])
2 # a.shape == (3,1,1)
3 b = np.arange(12).reshape((3,4))
4 # array([[ 0,  1,  2,  3],
5 #        [ 4,  5,  6,  7],
6 #        [ 8,  9, 10, 11]])
7 np.maximum(a,b)
8 # ???
```

Note: Breakout session: explain the output.

Reduction

- `np.max` is a type of reduction operation like `np.sum` and `np.mean`

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 # array([[[ 0,  1,  2,  3],
3 #         [ 4,  5,  6,  7]],
4 #
5 #        [[ 8,  9, 10, 11],
6 #         [12, 13, 14, 15]],
7 #
8 #        [[16, 17, 18, 19],
9 #         [20, 21, 22, 23]])
10 np.max(a)
11 # 23
```

Reduction

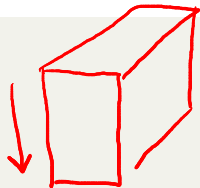
- `np.max` is a type of reduction operation like `np.sum` and `np.mean`

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 # array([[[ 0,  1,  2,  3],
3 #         [ 4,  5,  6,  7]],
4 #
5 #        [[ 8,  9, 10, 11],
6 #         [12, 13, 14, 15]],
7 #
8 #        [[16, 17, 18, 19],
9 #         [20, 21, 22, 23]])
10 np.max(a)
11 # 23
```

Reduction

- `np.max` is a type of reduction operation like `np.sum` and `np.mean`

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 # array([[[ 0,  1,  2,  3],
3 #        [ 4,  5,  6,  7]],
4 #
5 #        [[ 8,  9, 10, 11],
6 #        [12, 13, 14, 15]],
7 #
8 #        [[16, 17, 18, 19],
9 #        [20, 21, 22, 23]])
10 np.max(a, axis=0)
11 # a.shape == (2,4)
12 # array([[16, 17, 18, 19],
13 #        [20, 21, 22, 23]])
```

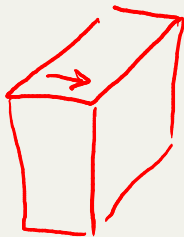


reduce axis = 0

Reduction

- `np.max` is a type of reduction operation like `np.sum` and `np.mean`

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 # array([[[ 0,  1,  2,  3],
3 #         [ 4,  5,  6,  7]],
4 #
5 #        [[ 8,  9, 10, 11],
6 #         [12, 13, 14, 15]],
7 #
8 #        [[16, 17, 18, 19],
9 #         [20, 21, 22, 23]])
10 np.max(a, axis=1)
11 # a.shape == (3,4)
12 # array([[ 4,  5,  6,  7],
13 #        [12, 13, 14, 15],
14 #        [20, 21, 22, 23]])
```



Reduction

- `np.max` is a type of reduction operation like `np.sum` and `np.mean`

```
1 a = np.arange(3*2*4).reshape([3,2,4])
2 # array([[[ 0,  1,  2,  3],
3 #         [ 4,  5,  6,  7]],
4 #
5 #        [[ 8,  9, 10, 11],
6 #         [12, 13, 14, 15]],
7 #
8 #        [[16, 17, 18, 19],
9 #         [20, 21, 22, 23]])
10 np.max(a, axis=-1)
11 # a.shape == (3,2)
12 # array([[ 3,  7],
13 #        [11, 15],
14 #        [19, 23]])
```

Forward

$n = \# \text{ of samples}$

$$\mathbf{h}^{(i)} = \mathbf{W}^{(1)\top} \mathbf{x}^{(i)} + \mathbf{b}^{(1)} \quad \text{and} \quad \mathbf{z}^{(i)} = \mathbf{W}^{(2)\top} g(\mathbf{h}^{(i)}) + \mathbf{b}^{(2)}$$

$\# \text{ of neurons}$

```
1 def relu(z):  
2     return np.maximum(0, z)  
3  
4 def forward(X, theta):  
5     W1, b1, W2, b2 = theta["W1"], theta["b1"], theta["W2"], theta["b2"]  
6     h = relu(np.dot(X, W1) + b1)  
7     z = np.dot(h, W2) + b2  
8     return h, z
```

$b_1 \text{ shape} = (m)$

$(X @ W_1) \text{ shape} = (n, m)$

broadcast

Note: Batched or vectorized operation

Forward (Non-vectorized)

$$\mathbf{h}^{(i)} = \mathbf{W}^{(1)\top} \mathbf{x}^{(i)} + \mathbf{b}^{(1)} \quad \text{and} \quad \mathbf{z}^{(i)} = \mathbf{W}^{(2)\top} g(\mathbf{h}^{(i)}) + \mathbf{b}^{(2)}$$

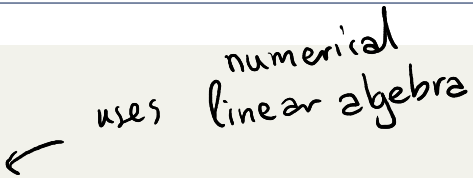
```
1 def forward_nonvectorized(X, theta):
2     W1, b1, W2, b2 = theta["W1"], theta["b1"], theta["W2"], theta["b2"]
3     h = np.zeros((X.shape[0], W1.shape[1])) # Initialize h
4     z = np.zeros((X.shape[0], W2.shape[1])) # Initialize z
5
6     for i in range(X.shape[0]):
7         h[i, :] = relu(np.dot(X[i, :], W1) + b1)
8         z[i, :] = np.dot(h[i, :], W2) + b2
9     return h, z
```

Note: Non-vectorized or row-wise operation

Timing the Forward Functions

```
1 import time
2
3 # Vectorized forward function
4 start_time = time.time()
5 h_vec, z_vec = forward(X, theta)
6 end_time = time.time()
7 print(f"Vectorized time: {end_time - start_time:.5f} seconds")
8
9 # Non-vectorized forward function
10 start_time = time.time()
11 h_nonvec, z_nonvec = forward_nonvectorized(X, theta)
12 end_time = time.time()
13 print(f"Non-vectorized time: {end_time - start_time:.5f} seconds")
```

uses numerical linear algebra



Note: Measure the time for vectorized vs. non-vectorized forward pass

Training loss / risk

1.

$$J(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N J_i(\boldsymbol{\theta})$$

and

$$J_i(\boldsymbol{\theta}) := L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$$

↑
cross entropy

Backward

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(2)}} = \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}} \quad \Longleftarrow$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(2)}} = g(\mathbf{h}^{(i)}) \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(1)}} = (\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})$$


$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)\top}$$

Cross entropy

Cross entropy

$$L(\mathbf{z}, y) = -\log(\text{softmax}(\mathbf{z})_y)$$

Derivative


$$\frac{\partial L}{\partial \mathbf{z}}(\mathbf{z}, y) = \underbrace{\mathbf{p}}_{\text{softmax}(\mathbf{z})} - \underbrace{\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}}_{\text{1-hot vector}} \leftarrow y\text{-th position}$$

Softmax

$$\text{softmax}(\mathbf{z}^{(i)}) = \frac{1}{\sum_{j=1}^K \exp(z_j^{(i)})} \begin{bmatrix} \exp(z_1^{(i)}) \\ \vdots \\ \exp(z_K^{(i)}) \end{bmatrix} = \begin{bmatrix} \text{softmax}(\mathbf{z}^{(i)})_1 \\ \vdots \\ \text{softmax}(\mathbf{z}^{(i)})_K \end{bmatrix} = \begin{bmatrix} p_1 \\ \vdots \\ p_K \end{bmatrix} = \mathbf{p}$$

```
1 p = np.zeros((n,n_classes))
2 h,z = forward(X,theta)
3 for i in range(n):
4     p[i,:] = np.exp(z[i,:])/np.sum(np.exp(z[i,:]))
```

Note: Breakout session: compute this without the for-loop

Softmax

$$\text{softmax}(\mathbf{z}^{(i)}) = \frac{1}{\sum_{j=1}^K \exp(z_j^{(i)})} \begin{bmatrix} \exp(z_1^{(i)}) \\ \vdots \\ \exp(z_K^{(i)}) \end{bmatrix} = \begin{bmatrix} \text{softmax}(\mathbf{z}^{(i)})_1 \\ \vdots \\ \text{softmax}(\mathbf{z}^{(i)})_K \end{bmatrix} = \begin{bmatrix} p_1 \\ \vdots \\ p_K \end{bmatrix} = \mathbf{p}$$

```
1 h,z = forward(X,theta)
2 expz = np.exp(z)
3
4 p = expz / np.sum(expz,axis=-1)[:,np.newaxis] # your code here
```

data \downarrow $\begin{bmatrix} z \\ e \end{bmatrix} \xrightarrow{\text{np.sum}(-,axis=-1)} \begin{bmatrix} \sum e^z \end{bmatrix}$

pairwise op also broadcast

Loss derivative

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}}_{\text{pick } y\text{-th column}} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow y\text{-th position}$$

Identity matrix of size `n_classes`

```
1 E = np.eye(n_classes)
2 E[y[i], :]
```

 1-hot vector

Cross entropy

Derivative

$$\frac{\partial L}{\partial \mathbf{z}}(\mathbf{z}, y) = \underbrace{\mathbf{p}}_{\text{softmax}(\mathbf{z})} - \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow y\text{-th position} = \mathbf{p} - \mathbf{e}^{(y)}$$

```
1 loss_der = np.zeros((n,n_classes))
2 E = np.eye(n_classes)
3 h,z = forward(X,theta)
4 for i in range(n):
5     loss_der[i,:] = p[i,:] - E[y[i], :]
```

loss_der = p -

no for-loops

Note: Vectorize this

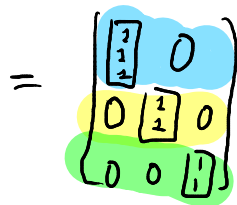
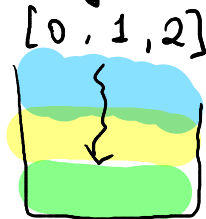
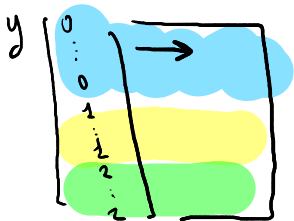
Cross entropy

cast to float
↑

$$\frac{\partial J_i(\theta)}{\partial \mathbf{b}^{(2)}} = \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}$$

↗ bool

```
1 y_one_hot = 1*(y[:,np.newaxis] == np.arange(n_classes))  
2  
3 loss_der = p - y_one_hot  
4  
5 dJdb2 = np.mean(loss_der, axis=0)
```



1-layer neural network (multicategory data)

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(2)}} = \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}} \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(2)}} = g(\mathbf{h}^{(i)}) \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top \quad \Longleftarrow$$

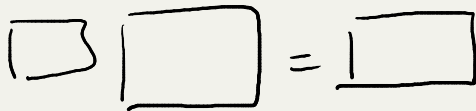
$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(1)}} = (\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)\top}$$

Need to talk about matrix mult

Matrix multiplication

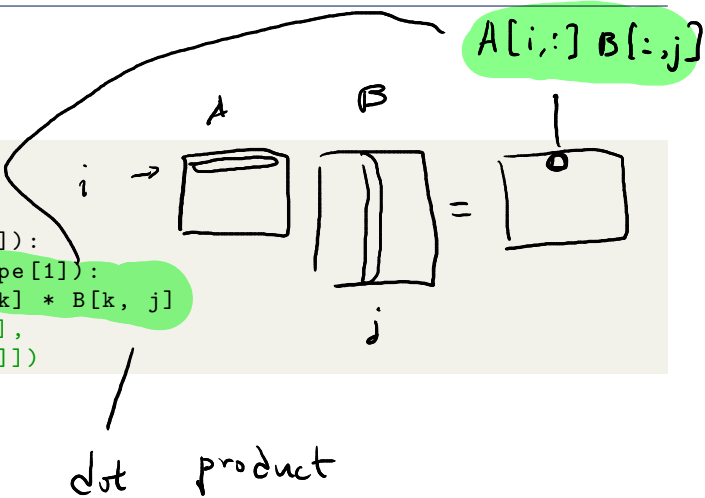
```
1 np.random.seed(42)
2 A = np.random.randint(0,9,(2,3))
3 B = np.random.randint(0,9,(3,4))
4 A@B
5 # array([[85, 70, 79, 68],
6 #        [74, 72, 42, 52]])
```



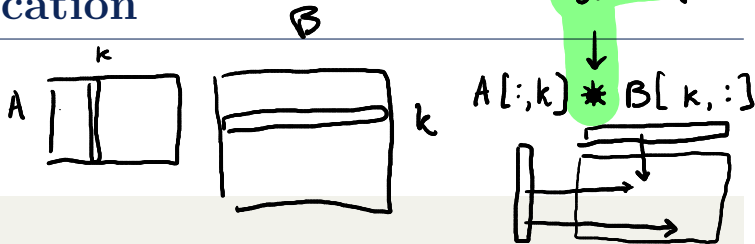
Matrix multiplication

inner product

```
1 C = np.zeros((2, 4))
2
3 for i in range(A.shape[0]):
4     for j in range(B.shape[1]):
5         for k in range(A.shape[1]):
6             C[i, j] += A[i, k] * B[k, j]
7 # array([[85., 70., 79., 68.],
8 #        [74., 72., 42., 52.]])
```



Matrix multiplication



```
1 D = np.zeros((2, 4))
2
3 for k in range(A.shape[1]): # iterate over the elements to multiply
4     D += np.outer(A[:, k], B[k, :])
5 # array([[85., 70., 79., 68.],
6 #        [74., 72., 42., 52.]])
```


Cross entropy

$$\sum_i \frac{\partial J_i(\theta)}{\partial W^{(2)}}$$

single
comp

$$\frac{\partial J_i(\theta)}{\partial W^{(2)}} = g(\mathbf{h}^{(i)}) \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top$$

$\text{reln}(\mathbf{h})$

(n, m)

$(n, m, 1)$

loss_der

$(n, 1, K)$

(n, m, K)

$\hookrightarrow (n, m, K)$

```
1 dJdW2 = np.zeros_like(W2)
2 # h,z = forward(X,theta) # already computed previously
3
4 for i in range(n):
5     dJdW2 += np.outer(reu(h[i,:]), loss_der[i,:])
6
7 dJdW2/n
```

Note: “vectorize” this

Cross entropy

$$\text{relu}(h) = \begin{bmatrix} \text{sample } i \\ \text{---} \end{bmatrix}$$

$$\text{loss_der} = \begin{bmatrix} \text{---} \end{bmatrix}$$

$$\frac{\partial J_i(\theta)}{\partial \mathbf{W}^{(2)}} = g(\mathbf{h}^{(i)}) \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top$$

$$\text{relu}(h)^\top = \begin{bmatrix} \text{sample } i \\ \text{---} \end{bmatrix}$$

```
1 dJdW2 = np.matmul(relu(h).T, loss_der)/n
```

exactly — is $g(h^{(i)})$

$\begin{matrix} \text{---} \\ \text{column} \end{matrix} \star \begin{matrix} \text{---} \\ \text{row} \end{matrix}$
 \uparrow
 outer prod

1-layer neural network (multicategory data)

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(2)}} = \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}} \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(2)}} = g(\mathbf{h}^{(i)}) \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(1)}} = (\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)}) \quad \Leftarrow$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)\top}$$

Cross Entropy

inner layer

$$\frac{\partial J_i(\theta)}{\partial \mathbf{b}^{(1)}} = (\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})$$

```
1 def relu_derivative(z):  
2     return 1.0*(z > 0)  
3  
4 dJdb1 = np.zeros_like(b1)  
5  
6 # h,z = forward(X,theta)  
7 for i in range(n):  
8     dJdb1 += (W2@(loss_der[i,:])) * relu_derivative(h[i,:])  
9 dJdb1 /=n
```

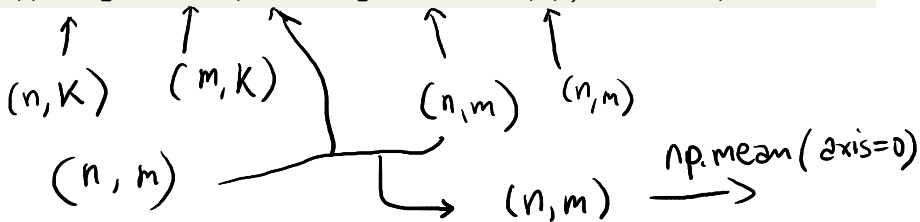
Note: Vectorize this

Cross Entropy

$$W^{(2)} \quad (m, K) \quad \swarrow \quad \mathcal{D} \quad (n, K)$$

$$\frac{\partial J_i(\theta)}{\partial \mathbf{b}^{(1)}} = (W^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})$$

```
1 dJdb1 = np.mean((loss_der@W2.T) * relu_derivative(h), axis = 0)
```



1-layer neural network (multicategory data)

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(2)}} = \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}} \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(2)}} = g(\mathbf{h}^{(i)}) \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(1)}} = (\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)}) \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)\top} \quad \Longleftarrow$$

Cross Entropy

X^T
(2, n)

$$\frac{\partial J_i(\theta)}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)T}$$

$\text{loss_der} @ W_2 * g'(h)$

\parallel
 G
 \parallel

(n, m)

```
1 dJdW1 = np.zeros_like(W1)
2
3 # h, z = forward(X, theta)
4 for i in range(n):
5     dJdW1 += np.outer((W2@loss_der[i,:]) *
6                        relu_derivative(h[i,:]), X[i,:]).T
7 dJdW1 /= n
```

Note: Vectorize this

Cross Entropy

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)\top}$$

```
1 dJdW1 = np.zeros_like(W1)
2
3 # h,z = forward(X,theta)
4 for i in range(n):
5     dJdW1 += np.outer((W2@loss_der[i,:]) *
6                        relu_derivative(h[i,:]), X[i,:]).T
7 dJdW1 /=n
```

Note: Vectorize this

Cross Entropy

$$=: G$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)\top}$$

```
1 G = (loss_der@W2.T) * relu_derivative(h)
2 dJdW1 = np.matmul(G.T, X).T/n
```

Backprop for 1-layer
= matrix multiply

1-layer neural network (multicategory data)

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(2)}} = \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}} \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(2)}} = g(\mathbf{h}^{(i)}) \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(1)}} = (\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)}) \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)\top} \quad (\checkmark)$$

1-layer neural network (multicategory data)

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(2)}} = \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}} \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(2)}} = g(\mathbf{h}^{(i)}) \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(1)}} = (\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)}) \quad (\checkmark)$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{W}^{(1)}} = ((\mathbf{W}^{(2)} \frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}) \odot g'(\mathbf{h}^{(i)})) \mathbf{x}^{(i)\top} \quad (\checkmark)$$

$$x^{(1)}, x, \dots, x^{(n)} \in \mathbb{R}^d$$

$$y^{(1)}, y, \dots, y^{(n)} \in \{\pm 1\}$$

binary.

Standard setup

Every "data point" is now

C different data points

$$x^{(i)}$$

$C = \text{context e.g.} = 64$

$$\{ x^{(i,1)}, x^{(i,2)}, \dots, x^{(i,C)} \}$$

$$\cap \mathbb{R}^d$$

$$\cap \mathbb{R}^d$$

$$\cap \mathbb{R}^d$$

$$y^{(i)} \in \{\pm 1\}$$

Every data point is
a prompt

$$\{x^{(i,1)}, \dots, x^{(i,C)}\}$$

$$\overset{\parallel}{x^{(i)}} \in \mathbb{R}^{d \times C}$$

Consider this model

$$\overbrace{(x^{(i)})^T}^{K^T} \overbrace{W^{(1)} W^{(2)} x^{(i)}}^Q$$

↗
?

$$\cap \\ \mathbb{R}^{d \times q}$$

$$\cap \\ \mathbb{R}^{d \times q}$$

$$C \times C$$

$$x^{(i)} \text{softmax}(\underbrace{K^T Q}_{C \times C})$$

\mathbb{R}^d dimensional vector

$$\equiv w^{(3)T} x^{(i)} \text{softmax}(K_{(i)}^T Q_{(i)})$$

$$\hat{y}^{(i)} \quad \theta = \{w^{(1)}, w^{(2)}, w^{(3)}\}$$

$$L(y^{(i)}, \hat{y}^{(i)})$$

How to implement the loss derivative

$$\frac{\partial L(\mathbf{z}^{(i)}, y)}{\partial \mathbf{z}}^\top$$

```
1  h, z = forward(X, theta)
2  p = np.exp(z[i, :]) / np.sum(np.exp(z[i, :]))
3  loss_der = p - E[y[i], :]
4  dJdW1 += np.outer((W2 @ loss_der) * relu_derivative(h[i, :]), X[i, :]).T
5  dJdb1 += (W2 @ (loss_der)) * relu_derivative(h[i, :])
6  dJdW2 += np.outer(relu(h[i, :]), loss_der)
7  dJdb2 += loss_der
```

Stochastic gradient descent (SGD)

Let $\eta_t > 0$ be learning rates, $t = 1, 2, \dots$

Let $m \geq 1$ be an integer

- Initialize $\boldsymbol{\theta}$
- While not converged ($t = \text{iteration counter}$):
 - Select m samples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ and matching labels $\{y^{(1)}, \dots, y^{(m)}\}$
 - Compute gradient $\mathbf{g} \leftarrow \nabla_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}, \boldsymbol{\theta}), y^{(i)})$
 - Compute update $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta_t \mathbf{g}$

Stochastic gradient descent (SGD)

Let $\eta_t > 0$ be learning rates, $t = 1, 2, \dots$

Let $m \geq 1$ be an integer

- Initialize $\boldsymbol{\theta}$
- While not converged ($t = \text{iteration counter}$):
 - Select 1 sample $\{\mathbf{x}^{(t)}\}$ and its label $\{y^{(t)}\}$
 - Compute gradient $\mathbf{g} \leftarrow \nabla_{\boldsymbol{\theta}} L(f(\mathbf{x}^{(t)}, \boldsymbol{\theta}), y^{(t)})$
 - Compute update $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta_t \mathbf{g}$

References I
