

Convolutional nets

Lecture 08 — CS 577 Deep Learning

Instructor: Yutong Wang

Computer Science
Illinois Institute of Technology

October 9, 2024

Autograd-enabled tensors

Previously

- `ag.Scalar`

Today:

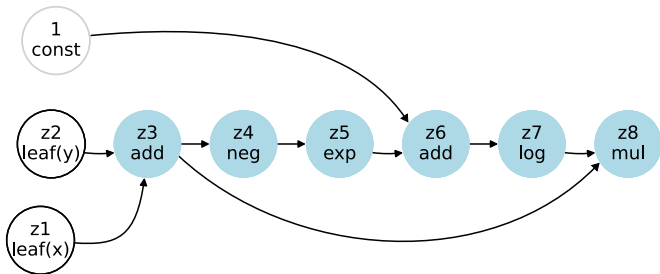
- Use `ag.Scalar` for training small relu networks
- `ag.Tensor`
- Overflow in cross entropy and how to avoid it
- Use `ag.Tensor` to train convolutional neural networks (CNNs)

Autograd-enabled scalars: `ag.Scalar`

```
1  class Scalar: # Scalars with grads
2      def __init__(self,
3                  value,
4                  op="",
5                  _backward= lambda : None,
6                  inputs=[],
7                  label=""):
8
9          self.value = float(value)
10         self.grad = 0.0
11
12         self._backward = _backward
13         self.inputs = inputs
14
15         self.op = op
16         self.label = label
```

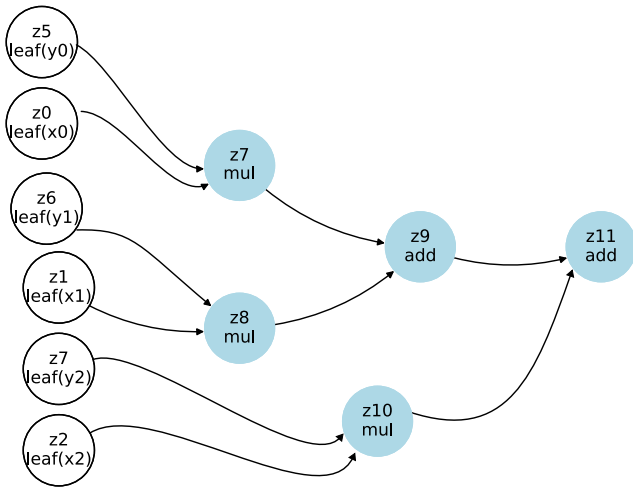
Computational graph

$$f(x, y) = \log(1 + e^{-(x+y)}) \cdot (x + y).$$



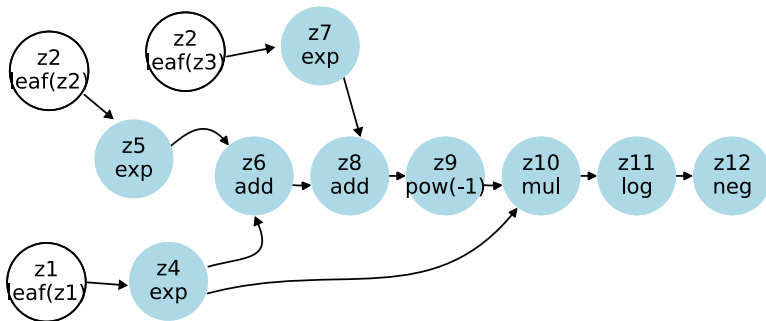
Computational graph

$$f(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y} = x_1 y_1 + x_2 y_2 + x_3 y_3$$



Computational graph

$$-\log(e^{z_1}/(e^{z_1} + e^{z_2} + e^{z_3}))$$



```
1 expz1 = ag.exp(z1)
2 expz2 = ag.exp(z2)
3 expz3 = ag.exp(z3)
4 final_output = -ag.log( expz1/(expz1+expz2+expz3))
```

What can we do with `ag.Scalar`?

PyTorch-“lite” consisting of:

- A `Model` class (params and forward)
- A `Loss` class (just Mean Squared Error (MSE) for now)
- An `Optimizer` class (just gradient descent for now)
- A training loop

We will fit a 1-hidden layer neural network with 1-dimensional input

$$f(x; \mathbf{w}_1, \mathbf{b}_1, \mathbf{w}_2, b_2) = \mathbf{w}_2^\top \text{relu}(\mathbf{w}_1 x + \mathbf{b}_1) + b_2$$

where \mathbf{w}_1 and \mathbf{w}_2 are both vectors of `n_hidden` dimensions.

Training Loop

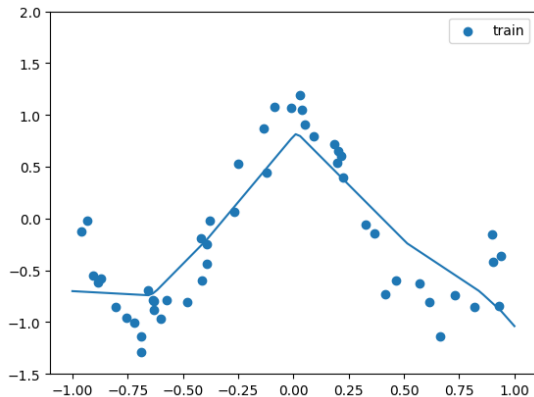
```
1 model = Model(n_hidden=20)
2 loss_fn = Loss()
3 optimizer = Optimizer(model.parameters, lr=0.1)
4
5 for epoch in range(100):
6     optimizer.zero_grad()
7     output = model.forward(xnp)
8     loss = loss_fn.mse(output, ynp)
9     loss.backward()
10    optimizer.step()
11    if epoch % 10 == 0:
12        print(f"Iteration {epoch}, Loss: {loss.value}")
```


Exercise 1

[lec08-in-class-ex1-framework.ipynb](#)

$$f(x; \mathbf{w}_1, \mathbf{b}_1, \mathbf{w}_2, b_2) = \mathbf{w}_2^\top \text{relu}(\mathbf{w}_1 x + \mathbf{b}_1) + b_2$$

where \mathbf{w}_1 and \mathbf{w}_2 are both vectors of `n_hidden` dimensions.



The Model class

```
1 class Model:
2     def __init__(self, n_hidden, rng_seed=42):
3         np.random.seed(rng_seed)
4
5         w1np = np.random.randn(n_hidden)
6         # ...
7         b2np = np.random.randn(1)
8
9         self.w1 = [ag.Scalar(val) for val in w1np]
10        # ...
11        self.b2 = [ag.Scalar(val) for val in b2np]
12
13        self.parameters = self.w1 + self.b1 + self.w2 + self.b2
14
15    def forward(self, x): # ...
16        # x is a 1-dimensional numpy array
```

The forward function

```
1  def forward(self, x):
2      # x is a 1-dimensional numpy array
3      # "upgrade" x into ag.Scalars
4      x_scalar = [ag.Scalar(val) for val in x]
5      n_samples = len(x_scalar)
6
7      # calculate the forward
8
9      ## YOUR CODE HERE
10     return [ag.Scalar(0.0) for i in range(n_samples)]
```

The Loss class

```
1 class Loss:
2     def mse(self, predictions, targets):
3         # mean squared error
4         assert len(predictions) == len(targets)
5         n_samples = len(predictions)
6         loss = ag.Scalar(0.0)
7
8         # YOUR CODE HERE
9
10    return loss
```

The Optimizer class

```
1 class Optimizer:
2     def __init__(self, parameters, lr=0.01):
3         self.parameters = parameters
4         self.lr = lr
5
6     def zero_grad(self):
7         for param in self.parameters:
8             param.grad = 0.0
9
10    def step(self):
11        for param in self.parameters:
12            param.value -= self.lr * param.grad
```

Check gradients!

```
1 Gradients w.r.t...
2 {'w1': array([-0.77812733, -0.          , ...
3
4  'b1': array([ 2.35879461, -0.          , ...
5
6  'w2': array([ 4.15815372e+00, -0.00000000e+00, ...
7
8  'b2': 3.1941792242464686}
```

Upshot

- It seems to work just fine.

Autograd-enabled tensors

Previously

- `ag.Scalar`

Today:

- ☒ Use `ag.Scalar` for training small relu networks
- `ag.Tensor`
- Use `ag.Tensor` to train convolutional neural networks (CNNs)
- Overflow in cross entropy and how to avoid it

ag.Scalar

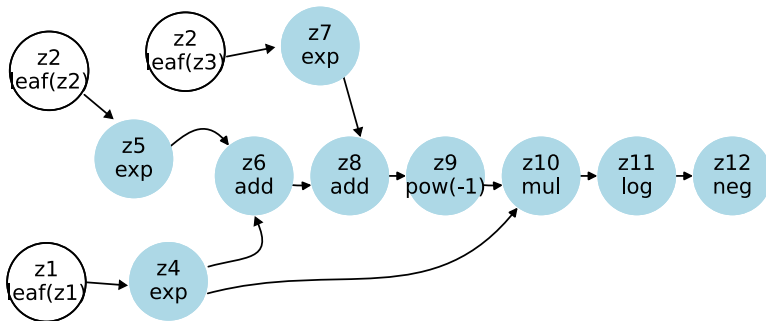
```
1  class Scalar: # Scalars with grads
2      def __init__(self,
3                      value,
4                      op="",
5                      _backward= lambda : None,
6                      inputs=[],
7                      label=""):
8
9      self.value = float(value)
10     self.grad = 0.0
```

ag.Tensor

```
1  class Tensor: # Tensor with grads
2      def __init__(self,
3                    value,
4                    op="",
5                    _backward= lambda : None,
6                    inputs=[],
7                    label=""):
8
9      if type(value) in [float ,int]:
10         value = np.array(value)
11         self.value = value # <-----\
12 #                                     |- same shape!
13         self.grad = np.zeros_like(value) # <-----/
```

Computational graph

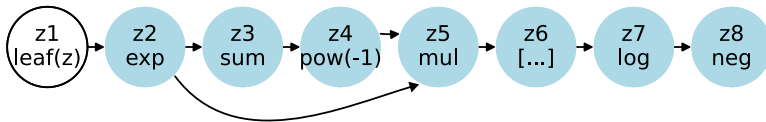
$$-\log(e^{z_1}/(e^{z_1} + e^{z_2} + e^{z_3}))$$



```
1 expz1 = ag.exp(z1)
2 expz2 = ag.exp(z2)
3 expz3 = ag.exp(z3)
4 final_output = -ag.log( expz1/(expz1+expz2+expz3))
```

Computational graph

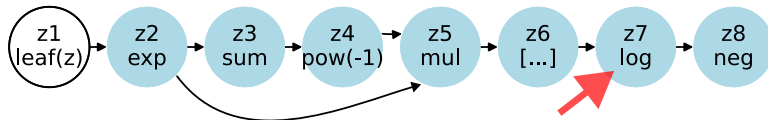
$$-\log(e^{z_1}/(e^{z_1} + e^{z_2} + e^{z_3})) = L(\mathbf{z}, 1) \quad \leftarrow \text{cross entropy}$$



```
1 z = ag.Tensor(np.random.randn(3), label="z1\nleaf(z)")
2 expX = ag.exp(z)
3 p = expX / ag.sum(expX)
4 nll = -ag.log(p[0]) # negative log-likelihood
```

Computational graph

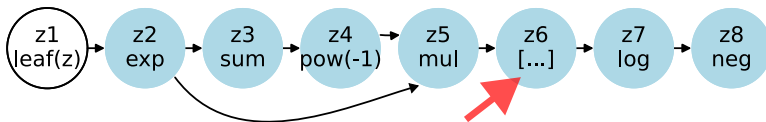
Entrywise ops



```
1 def log(input):  
2     output = ag.Tensor(np.log(input.value), inputs=[input], op="log")  
3     def _backward():  
4         input.grad += output.grad / input.value  
5         return None  
6     output._backward = _backward  
7     return output
```

Computational graph

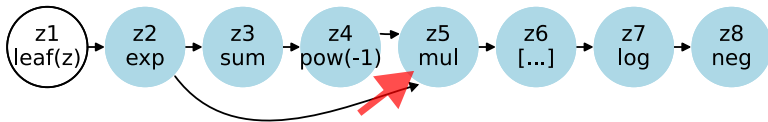
Slicing



```
1 def __getitem__(self, idx):
2     output = ag.Tensor(np.array(self.value[idx]),
3                           inputs = [self],
4                           op=f"[...]")
5
6     def _backward():
7         self.grad[idx] += output.grad # idx must not have repeats!
8         return None
9     output._backward = _backward
10    return output
```

Computational graph

Pairwise ops — multiplication (NOT matmul)



```
1 def __mul__(self, other):
2     # ... lines skipped
3     output = ag.Tensor(self.value * other.value,
4                          inputs=[self, other], op="mul")
5     def _backward():
6         self.grad += unbroadcast(output.grad*other.value, ax1, pad1)
7         other.grad += unbroadcast(output.grad*self.value, ax2, pad2)
```

What is “unbroadcast” and why define backward this way? ... next

Computational graph

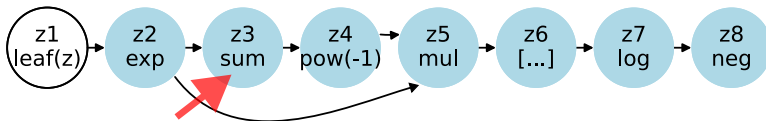
Pairwise ops — multiplication (NOT matmul)

What is “unbroadcast” and why define backward this way?

- `X1.shape = (2,3,1,5,6)`
- `X2.shape = (4,1,6)`
- What is `(X1 * X2).shape`?

Computational graph

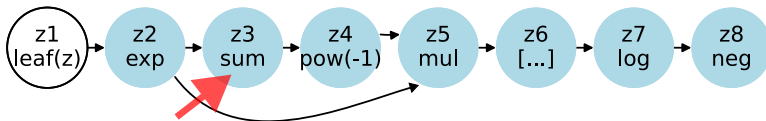
Reductive ops — sum



```
1 def sum(input,axis = None, keepdims = False):
2     output = ag.Tensor(np.sum(input.value, axis = axis, keepdims =
3     # ...
4     # ...
5     # ...
6     # ...
7     # ...
8     # ...
9     # ...
10    # ...
```

Computational graph

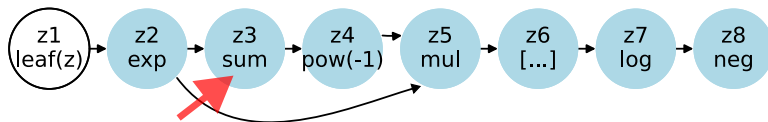
Reductive ops — sum



```
1 def sum(input,axis = None, keepdims = False):
2     output = ag.Tensor(np.sum(input.value, axis = axis, keepdims =
3         keepdims), inputs = [input], op='sum')
4     def _backward():
5         if axis == None or keepdims:
6             input.grad += output.grad
7         else:
8             # ...
9             # ...
10            # ...
```

Computational graph

Reductive ops — sum



```
1 def sum(input,axis = None, keepdims = False):
2     output = ag.Tensor(np.sum(input.value, axis = axis, keepdims =
3     keepdims), inputs = [input], op='sum')
4     def _backward():
5         if axis == None or keepdims:
6             input.grad += output.grad
7         else:
8             input.grad += np.expand_dims(output.grad, axis = axis)
9     return None
10 # ...
11 # ...
```

Exercises

- Reshape (Exercise 2.1)
 - Fix the backward function for reshape

Exercises

- Reshape (Exercise 2.1)
 - Fix the backward function for reshape
- Indexing (Exercise 2.2)
 - Fix the loss derivative autograd calculation

Exercises

- Reshape (Exercise 2.1)
 - Fix the backward function for reshape
- Indexing (Exercise 2.2)
 - Fix the loss derivative autograd calculation
- Implementing BCE (Exercise 2.3)
 - Implement the binary cross entropy.

Convolution

Class 2: Bird

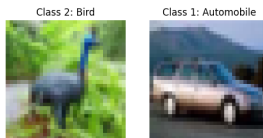


Class 1: Automobile

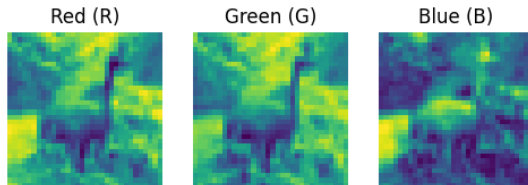


`X.shape = (N, C, H, W) = (2, 3, 32, 32)` (2 images from CIFAR10)

Convolution

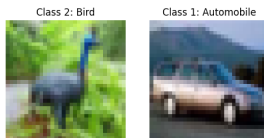


$$\mathbf{X}.\text{shape} = (N, C, H, W) = (2, 3, 32, 32)$$

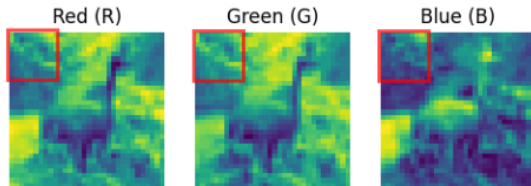


$$\mathbf{X}[0] . \text{shape} = (C, H, W) = (3, 32, 32)$$

Convolution

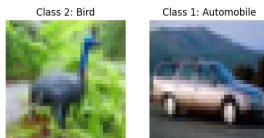


$$X.shape = (N, C, H, W) = (2, 3, 32, 32)$$

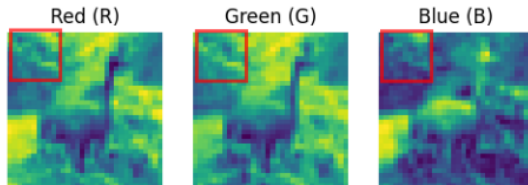


$$X[0].shape = (C, H, W) = (3, 32, 32)$$

Convolution

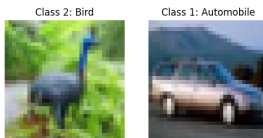


$$\mathbf{X}.\text{shape} = (N, C, H, W) = (2, 3, 32, 32)$$

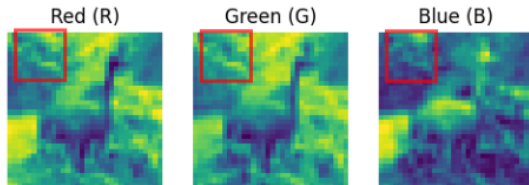


$$\mathbf{X}[0].\text{shape} = (C, H, W) = (3, 32, 32)$$

Convolution

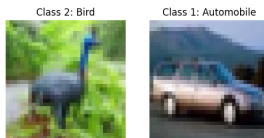


$$\mathbf{X}.\text{shape} = (N, C, H, W) = (2, 3, 32, 32)$$

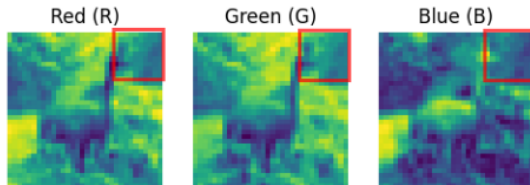


$$\mathbf{X}[0].\text{shape} = (C, H, W) = (3, 32, 32)$$

Convolution

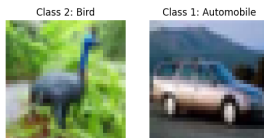


$$X.shape = (N, C, H, W) = (2, 3, 32, 32)$$

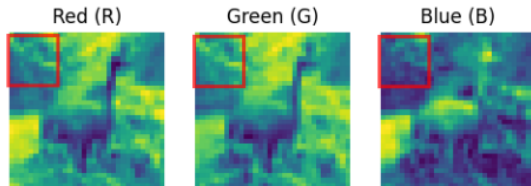


$$X[0].shape = (C, H, W) = (3, 32, 32)$$

Convolution



$$\mathbf{X}.\text{shape} = (N, C, H, W) = (2, 3, 32, 32)$$



$$\mathbf{X}[0].\text{shape} = (C, H, W) = (3, 32, 32)$$

LeNet5

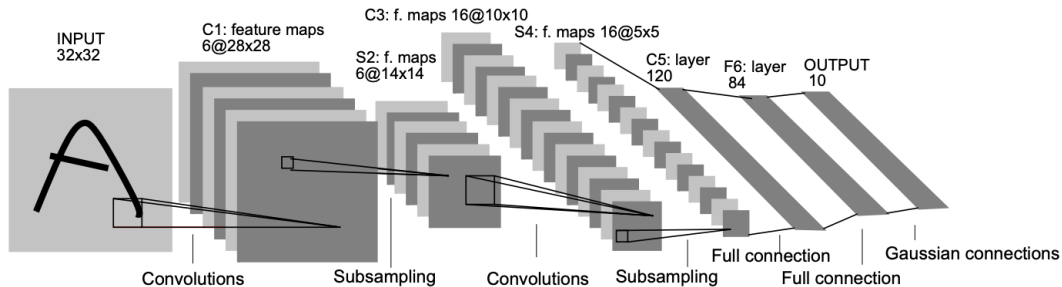
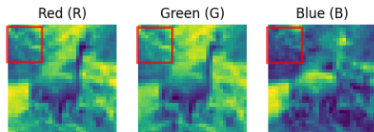
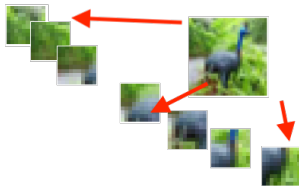


Image from LeCun et al. 1998

Convolution



$X[0].\text{shape} = (C, H, W) = (3, 32, 32)$



$X_{\text{patches}}[0].\text{shape} = (P, C, K, K) = (3, 32, 32)$

K = size of window, P = number of patches

A simple example

```
1 C, H, W = 2, 4, 5
2 K = 2 # Kernel size
3 N = 2 # Batch size
4 X = np.arange(N * C * H * W).reshape(N, C, H, W)
5 X[0]
6 array([[ 0,  1,  2,  3,  4],
7        [ 5,  6,  7,  8,  9],
8        [10, 11, 12, 13, 14],
9        [15, 16, 17, 18, 19]],
10
11       [[20, 21, 22, 23, 24],
12        [25, 26, 27, 28, 29],
13        [30, 31, 32, 33, 34],
14        [35, 36, 37, 38, 39]])
```

Goal: compute `X_patches`. Okay, but what is `X_patches`?

A simple example

```
1 X[0]
2 array([[ 0,  1,  2,  3,  4],
3        [ 5,  6,  7,  8,  9],
4        [10, 11, 12, 13, 14],
5        [15, 16, 17, 18, 19]],
6
7        [[20, 21, 22, 23, 24],
8         [25, 26, 27, 28, 29],
9         [30, 31, 32, 33, 34],
10        [35, 36, 37, 38, 39]])
```

Goal: compute `X_patches`. Okay, but what is `X_patches`?

```
1 X_patches[0,0]
2 array([[ 0.,  1.],
3        [ 5.,  6.]],
4
5        [[20., 21.],
6         [25., 26.]])
```

A simple example

```
1 X[0]
2 array([[ 0,  1,  2,  3,  4],
3        [ 5,  6,  7,  8,  9],
4        [10, 11, 12, 13, 14],
5        [15, 16, 17, 18, 19]],
6
7        [[20, 21, 22, 23, 24],
8        [25, 26, 27, 28, 29],
9        [30, 31, 32, 33, 34],
10       [35, 36, 37, 38, 39]])
```

Goal: compute `X_patches`. Okay, but what is `X_patches`?

```
1 X_patches[0,1]
2 array([[ 1.,  2.],
3        [ 6.,  7.]],
4
5        [[21., 22.],
6        [26., 27.]])
```

A simple example

```
1 X[0]
2 array([[[ 0,  1,  2,  3,  4],
3         [ 5,  6,  7,  8,  9],
4         [10, 11, 12, 13, 14],
5         [15, 16, 17, 18, 19]],
6
7        [[20, 21, 22, 23, 24],
8         [25, 26, 27, 28, 29],
9         [30, 31, 32, 33, 34],
10        [35, 36, 37, 38, 39]]])
```

Goal: compute `X_patches`. Okay, but what is `X_patches`?

```
1 X_patches[0,3]
2 array([[[ 3.,  4.],
3         [ 8.,  9.]],
4
5        [[23., 24.],
6         [28., 29.]])
```

A simple example

```
1 X[0]
2 array([[ 0,  1,  2,  3,  4],
3        [ 5,  6,  7,  8,  9],
4        [10, 11, 12, 13, 14],
5        [15, 16, 17, 18, 19]],
6
7        [[20, 21, 22, 23, 24],
8         [25, 26, 27, 28, 29],
9         [30, 31, 32, 33, 34],
10        [35, 36, 37, 38, 39]])
```

Goal: compute `X_patches`. Okay, but what is `X_patches`?

```
1 X_patches[0,4]
2 array([[ 5.,  6.],
3        [10., 11.]],
4
5        [[25., 26.],
6         [30., 31.]])
```

A simple example

```
1 X[1]
2 array([[40, 41, 42, 43, 44],
3        [45, 46, 47, 48, 49],
4        [50, 51, 52, 53, 54],
5        [55, 56, 57, 58, 59]],
6
7        [[60, 61, 62, 63, 64],
8        [65, 66, 67, 68, 69],
9        [70, 71, 72, 73, 74],
10       [75, 76, 77, 78, 79]])
```

Goal: compute `X_patches`. Okay, but what is `X_patches`?

```
1 X_patches[1,0]
2 array([[40., 41.],
3        [45., 46.]],
4
5        [[60., 61.],
6        [65., 66.]])
```

X_patches

$$\text{X_patches}[0].\text{shape} = (P, C, K, K) = (3, 32, 32)$$

K = size of window, P = number of patches

```
1 CHW = C * H * W
2 out_H = H - K + 1 # Output height
3 out_W = W - K + 1 # Output width
4 P = out_H * out_W # Total number of patches per image
```

Strategy:

- First flatten X to be a matrix, call it X_{flat}

X_patches

$$\text{X_patches}[0].\text{shape} = (P, C, K, K) = (3, 32, 32)$$

K = size of window, P = number of patches

```
1 CHW = C * H * W
2 out_H = H - K + 1 # Output height
3 out_W = W - K + 1 # Output width
4 P = out_H * out_W # Total number of patches per image
```

Strategy:

- First flatten X to be a matrix, call it X_{flat}
- Pick out the patches using this “im2col_mat” matrix

X_patches

$$\text{X_patches}[0].\text{shape} = (P, C, K, K) = (3, 32, 32)$$

K = size of window, P = number of patches

```
1 CHW = C * H * W
2 out_H = H - K + 1 # Output height
3 out_W = W - K + 1 # Output width
4 P = out_H * out_W # Total number of patches per image
```

Strategy:

- First flatten X to be a matrix, call it X_{flat}
- Pick out the patches using this “`im2col_mat`” matrix
- (there are two flavors: dense and sparse. More on that later)

X_patches

$$\text{X_patches}[0].\text{shape} = (P, C, K, K) = (3, 32, 32)$$

K = size of window, P = number of patches

```
1 CHW = C * H * W
2 out_H = H - K + 1 # Output height
3 out_W = W - K + 1 # Output width
4 P = out_H * out_W # Total number of patches per image
```

Strategy:

- First flatten X to be a matrix, call it X_{flat}
- Pick out the patches using this “`im2col_mat`” matrix
- (there are two flavors: dense and sparse. More on that later)
- Get $X_{\text{patches_flat}}$. Reshape them to get X_{patches} .

A simple example

```
1 C, H, W = 2, 4, 5
2 K,N = 2,2 # Kernel size, Batch size
3 X = np.arange(N * C * H * W).reshape(N, C, H, W)
4 X[0]
5 array([[[ 0,  1,  2,  3,  4],
6         [ 5,  6,  7,  8,  9],
7         [10, 11, 12, 13, 14],
8         [15, 16, 17, 18, 19]],
9
10        [[20, 21, 22, 23, 24],
11
12 # ...
13 X_flat = X.reshape(N, -1) # Shape (N, C*H*W)
```

A simple example

```
1 C, H, W = 2, 4, 5
2 K,N = 2,2 # Kernel size, Batch size
3 X = np.arange(N * C * H * W).reshape(N, C, H, W)
4 X[0]
5 array([[[ 0,  1,  2,  3,  4],
6         [ 5,  6,  7,  8,  9],
7         [10, 11, 12, 13, 14],
8         [15, 16, 17, 18, 19]],
9
10        [[20, 21, 22, 23, 24],
11
12 # ...
13 X_flat = X.reshape(N, -1) # Shape (N, C*H*W)
14 X_flat[0]
15 array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15,
16        16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32,
17        33, 34, 35, 36, 37, 38, 39])
```

A simple example

```
1 X_flat = X.reshape(N, -1)  # Shape (N, C*H*W)
2 X_flat[0]
3 array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15,
        16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32,
        33, 34, 35, 36, 37, 38, 39])
4
5 im2col_mat_dense = im2col_matrix_dense(X, K)  # assume this given
6 X_out_dense = np.matmul(X_flat, im2col_mat_dense)
7 X_out_dense
8 array([[ 0.,  1.,  5.,  6., 20., 21., 25., 26.,  1.,  2.,  6.,  7., 21.,
9         22., 26., 27.,  2.,  3.,  7.,  8., 22., 23., 27., 28.,  3.,  4.,
10        8.,  9., 23., 24., 28., 29.,  5.,  6., 10., 11., 25., 26., 30.,
11        31.,  6.,  7., 11., 12., 26., 27., 31., 32.,  7.,  8., 12., 13.,
12        ...])
```

A simple example

```
1 X_out_dense = np.matmul(X_flat, im2col_mat_dense)
2 X_out_dense
3 array([[ 0.,  1.,  5.,  6., 20., 21., 25., 26.,  1.,  2.,  6.,  7., 21.,
4         22., 26., 27.,  2.,  3.,  7.,  8., 22., 23., 27., 28.,  3.,  4.,
5         8.,  9., 23., 24., 28., 29.,  5.,  6., 10., 11., 25., 26., 30.,
6         31.,  6.,  7., 11., 12., 26., 27., 31., 32.,  7.,  8., 12., 13.,
7         ...
8 X_patches = X_patches_flat.reshape(N,P,C,K,K)
```

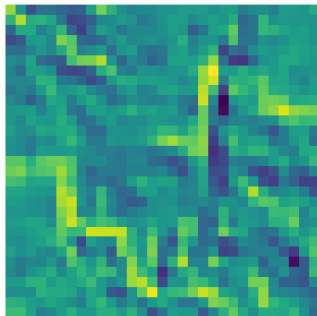
A simple example

```
1 X_out_dense = np.matmul(X_flat, im2col_mat_dense)
2 X_out_dense
3 array([[ 0.,  1.,  5.,  6., 20., 21., 25., 26.,  1.,  2.,  6.,  7., 21.,
4         22., 26., 27.,  2.,  3.,  7.,  8., 22., 23., 27., 28.,  3.,  4.,
5         8.,  9., 23., 24., 28., 29.,  5.,  6., 10., 11., 25., 26., 30.,
6         31.,  6.,  7., 11., 12., 26., 27., 31., 32.,  7.,  8., 12., 13.,
7         ...
8 X_patches = X_patches_flat.reshape(N,P,C,K,K)
9 X_patches[0,0]
10 array([[ 0.,  1.],
11         [ 5.,  6.]],
12
13         [[20., 21.],
14         [25., 26.]])
```

Exercise 3

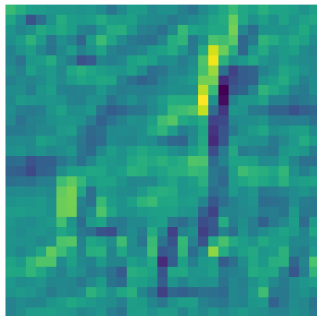
```
1      # input
2  X = np.arange(N * C * H * W).reshape(N, C, H, W)
3  X[0]
4  array([[[ 0,  1,  2,  3,  4],
5          [ 5,  6,  7,  8,  9],
6          [10, 11, 12, 13, 14],
7          [15, 16, 17, 18, 19]],
8         # ...
9
10     # desired output
11     # (2 times lower left - 1 times top right, ignore second channel)
12  X_convolved.shape = (2, 3, 4)
13  X_convolved[0]
14  array([[ 9., 10., 11., 12.],
15         [14., 15., 16., 17.],
16         [19., 20., 21., 22.]])
```

Smoothness



```
1 Conv_kernel = np.array([[[ 0, -1.], [ 1., 0.]], [[0., 0.], [0., 0.]],  
    [[0., 0.], [0., 0.]])
```


Smoothness



```
1 Conv_kernel = np.array([[[ 1, 0.], [ 0., -1.]], [[0., 0.], [0., 0.]],  
    [[0., 0.], [0., 0.]])
```

References I

- [LeC+98] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.