

## ✓ COSE474-2024F: Deep Learning HW1

- 컴퓨터학과 2020320041 김석민

### ✓ 0.1 Installation

```
!pip install d2l==1.0.3
```



```
Requirement already satisfied: ipython-genutils in /usr/local/lib/python3.10/dist-packages (from ipykernel->
Requirement already satisfied: ipython>=5.0.0 in /usr/local/lib/python3.10/dist-packages (from ipykernel->
Requirement already satisfied: jupyter-client in /usr/local/lib/python3.10/dist-packages (from ipykernel->
Requirement already satisfied: tornado>=4.2 in /usr/local/lib/python3.10/dist-packages (from ipykernel->
Requirement already satisfied: widgetsnbextension~=3.6.0 in /usr/local/lib/python3.10/dist-packages (from
Requirement already satisfied: jupyterlab-widgets>=1.0.0 in /usr/local/lib/python3.10/dist-packages (from
Requirement already satisfied: prompt-toolkit!=3.0.0,!3.0.1,<3.1.0,>=2.0.0 in /usr/local/lib/python3.10,
Requirement already satisfied: pygments in /usr/local/lib/python3.10/dist-packages (from jupyter-console-
Requirement already satisfied: lxml in /usr/local/lib/python3.10/dist-packages (from nbconvert->jupyter=
Requirement already satisfied: beautifulsoup4 in /usr/local/lib/python3.10/dist-packages (from nbconvert-
Requirement already satisfied: bleach in /usr/local/lib/python3.10/dist-packages (from nbconvert->jupyter=
Requirement already satisfied: defusedxml in /usr/local/lib/python3.10/dist-packages (from nbconvert->ju
Requirement already satisfied: entrypoints>=0.2.2 in /usr/local/lib/python3.10/dist-packages (from nbcon
Requirement already satisfied: jinja2>=3.0 in /usr/local/lib/python3.10/dist-packages (from nbconvert->j
Requirement already satisfied: jupyter-core>=4.7 in /usr/local/lib/python3.10/dist-packages (from nbconve
Requirement already satisfied: jupyterlab-pygments in /usr/local/lib/python3.10/dist-packages (from nbco
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from nbconver
Requirement already satisfied: mistune<2,>=0.8.1 in /usr/local/lib/python3.10/dist-packages (from nbconve
Requirement already satisfied: nbclient>=0.5.0 in /usr/local/lib/python3.10/dist-packages (from nbconver
Requirement already satisfied: nbformat>=5.1 in /usr/local/lib/python3.10/dist-packages (from nbconvert->
```

```
Requirement already satisfied: jupyter-server<3,>=1.8 in /usr/local/lib/python3.10/dist-packages (from n
Requirement already satisfied: cffi>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from argon2-cffi-l
Requirement already satisfied: pycparser in /usr/local/lib/python3.10/dist-packages (from cffi>=1.0.1->ar
Requirement already satisfied: anyio<4,>=3.1.0 in /usr/local/lib/python3.10/dist-packages (from jupyter-
Requirement already satisfied: websocket-client in /usr/local/lib/python3.10/dist-packages (from jupyter-
Requirement already satisfied: sniffio>=1.1 in /usr/local/lib/python3.10/dist-packages (from anyio<4,>=3
Requirement already satisfied: exceptiongroup in /usr/local/lib/python3.10/dist-packages (from anyio<4,>=3
```

## ✓ 2.0. Preliminaries

### ✓ 2.1. Data Manipulation

```
import torch
```

```
x = torch.arange(12, dtype=torch.float32)
x
```

```
→ tensor([ 0.,  1.,  2.,  3.,  4.,  5.,  6.,  7.,  8.,  9., 10., 11.])
```

```
x.numel()
```

```
→ 12
```

```
x.shape
```

```
→ torch.Size([12])
```

```
X = x.reshape(3,4)
X
```

```
→ tensor([[ 0.,  1.,  2.,  3.],
          [ 4.,  5.,  6.,  7.],
          [ 8.,  9., 10., 11.]])
```

```
torch.zeros((2,3,4))
```

```
→ tensor([[[0., 0., 0., 0.],
           [0., 0., 0., 0.],
           [0., 0., 0., 0.]],
          [[0., 0., 0., 0.],
           [0., 0., 0., 0.],
           [0., 0., 0., 0.]])
```

```
torch.randn(3,4)
```

```
→ tensor([[ -0.3204, -1.4158,  0.6210,  0.4161],
          [-1.2358, -0.8195,  0.8359,  0.1204],
          [ 0.9744, -2.6877,  0.2139,  0.3217]])
```

```
torch.tensor([[2, 1, 4, 3], [1, 2, 3, 4], [4, 3, 2, 1]])
```

```
→ tensor([[2, 1, 4, 3],
          [1, 2, 3, 4],
          [4, 3, 2, 1]])
```

X[-1], X[1:3]

```
→ (tensor([ 8.,  9., 10., 11.]),
    tensor([[ 4.,  5.,  6.,  7.],
            [ 8.,  9., 10., 11.])))
```

X[1,2]=17

X

```
→ tensor([[ 0.,  1.,  2.,  3.],
          [ 4.,  5., 17.,  7.],
          [ 8.,  9., 10., 11.]])
```

X[:,2,:]=12

X

```
→ tensor([[12., 12., 12., 12.],
          [12., 12., 12., 12.],
          [ 8.,  9., 10., 11.]])
```

torch.exp(x)

```
→ tensor([162754.7969, 162754.7969, 162754.7969, 162754.7969, 162754.7969,
          162754.7969, 162754.7969, 162754.7969, 2980.9580, 8103.0840,
          22026.4648, 59874.1406])
```

x = torch.tensor([1.0, 2, 4, 8])

y = torch.tensor([2,2,2,2])

x+y, x-y, x\*y, x/y, x\*\*y

```
→ (tensor([ 3.,  4.,  6., 10.]),
    tensor([-1.,  0.,  2.,  6.]),
    tensor([ 2.,  4.,  8., 16.]),
    tensor([0.5000, 1.0000, 2.0000, 4.0000]),
    tensor([ 1.,  4., 16., 64.])))
```

X = torch.arange(12, dtype=torch.float32).reshape(3,4)

Y = torch.tensor([[2.0, 1, 4, 3], [1,2,3,4], [4,3,2,1]])

torch.cat((X,Y), dim=0), torch.cat((X,Y), dim=1)

```
→ (tensor([[ 0.,  1.,  2.,  3.],
          [ 4.,  5.,  6.,  7.],
          [ 8.,  9., 10., 11.],
          [ 2.,  1.,  4.,  3.],
          [ 1.,  2.,  3.,  4.],
          [ 4.,  3.,  2.,  1.])),
    tensor([[ 0.,  1.,  2.,  3.,  2.,  1.,  4.,  3.],
          [ 4.,  5.,  6.,  7.,  1.,  2.,  3.,  4.],
          [ 8.,  9., 10., 11.,  4.,  3.,  2.,  1.])))
```

X==Y

```
→ tensor([[False,  True, False,  True],
          [False, False, False, False],
          [False, False, False, False]])
```

```
a = torch.arange(3).reshape((3,1))
b = torch.arange(2).reshape((1,2))
a,b
```

```
↔ (tensor([[0],
          [1],
          [2]]),
    tensor([[0, 1]]))
```

```
a+b
```

```
↔ tensor([[0, 1],
          [1, 2],
          [2, 3]])
```

```
before = id(Y)
Y = Y+X
id(Y)==before
```

```
↔ False
```

```
Z = torch.zeros_like(Y)
print('id(Z): ',id(Z))
Z[:] = X+Y
print('id(Z): ',id(Z))
```

```
↔ id(Z): 134741912877664
   id(Z): 134741912877664
```

```
before = id(X)
X += Y
id(X) == before
```

```
↔ True
```

```
A = X.numpy()
B = torch.from_numpy(A)
type(A), type(B)
```

```
↔ (numpy.ndarray, torch.Tensor)
```

```
a = torch.tensor([3.5])
a, a.item(), float(a), int(a)
```

```
↔ (tensor([3.5000]), 3.5, 3.5, 3)
```

## ✓ 2.2. Data Preprocessing

```
import os

os.makedirs(os.path.join '..', 'data'), exist_ok=True)
data_file = os.path.join '..', 'data', 'house_tiny.csv')
with open(data_file, 'w') as f:
    f.write(''NumRooms,RoofType,Price
NA,NA,127500
2,NA,106000
4,Slate,178100
NA,NA,140000'')
```

```
import pandas as pd
```

```
data = pd.read_csv(data_file)
print(data)
```

```
↗
```

	NumRooms	RoofType	Price
0	NaN	NaN	127500
1	2.0	NaN	106000
2	4.0	Slate	178100
3	NaN	NaN	140000

```
inputs, targets = data.iloc[:,0:2], data.iloc[:,2]
inputs = pd.get_dummies(inputs, dummy_na=True)
print(inputs)
```

```
↗
```

	NumRooms	RoofType_Slate	RoofType_nan
0	NaN	False	True
1	2.0	False	True
2	4.0	True	False
3	NaN	False	True

```
inputs = inputs.fillna(inputs.mean())
print(inputs)
```

```
↗
```

	NumRooms	RoofType_Slate	RoofType_nan
0	3.0	False	True
1	2.0	False	True
2	4.0	True	False
3	3.0	False	True

```
import torch
```

```
X = torch.tensor(inputs.to_numpy(dtype=float))
y = torch.tensor(targets.to_numpy(dtype=float))
X,y
```

```
↗ (tensor([[3., 0., 1.],
          [2., 0., 1.],
          [4., 1., 0.],
          [3., 0., 1.]], dtype=torch.float64),
 tensor([127500., 106000., 178100., 140000.], dtype=torch.float64))
```

## ✓ 2.3. Linear Algebra

```
import torch
```

```
x = torch.tensor(3.0)
y = torch.tensor(2.0)
```

```
x+y, x*y, x/y, x**y
```

```
↔ (tensor(5.), tensor(6.), tensor(1.5000), tensor(9.))
```

```
x = torch.arange(3)
x
```

```
↔ tensor([0, 1, 2])
```

```
x[2]
```

```
↔ tensor(2)
```

```
len(x)
```

```
↔ 3
```

```
x.shape
```

```
↔ torch.Size([3])
```

```
A = torch.arange(6).reshape(3,2)
A
```

```
↔ tensor([[0, 1],
          [2, 3],
          [4, 5]])
```

```
A.T
```

```
↔ tensor([[0, 2, 4],
          [1, 3, 5]])
```

```
A = torch.tensor([[1, 2, 3], [2, 0, 4], [3, 4, 5]])
A == A.T
```

```
↔ tensor([[True, True, True],
          [True, True, True],
          [True, True, True]])
```

```
torch.arange(24).reshape(2, 3, 4)
```

```
↔ tensor([[[ 0,  1,  2,  3],
           [ 4,  5,  6,  7],
           [ 8,  9, 10, 11]],
          [[12, 13, 14, 15],
           [16, 17, 18, 19],
           [20, 21, 22, 23]]])
```

```
A = torch.arange(6, dtype=torch.float32).reshape(2, 3)
B = A.clone()
A, A + B
```

```
→ (tensor([[0., 1., 2.],
           [3., 4., 5.]]),
    tensor([[ 0., 2., 4.],
           [ 6., 8., 10.])))
```

A\*B

```
→ tensor([[ 0., 1., 4.],
          [ 9., 16., 25.]])
```

a = 2

X = torch.arange(24).reshape(2, 3, 4)

a + X, (a \* X).shape

```
→ (tensor([[[ 2, 3, 4, 5],
             [ 6, 7, 8, 9],
             [10, 11, 12, 13]],
           [[14, 15, 16, 17],
            [18, 19, 20, 21],
            [22, 23, 24, 25]]]),
    torch.Size([2, 3, 4]))
```

x = torch.arange(3, dtype=torch.float32)

x, x.sum()

```
→ (tensor([0., 1., 2.]), tensor(3.))
```

A.shape, A.sum()

```
→ (torch.Size([2, 3]), tensor(15.))
```

A.shape, A.sum(axis=0).shape

```
→ (torch.Size([2, 3]), torch.Size([3]))
```

A.sum(axis=0)

```
→ tensor([3., 5., 7.])
```

A.shape, A.sum(axis=1).shape

```
→ (torch.Size([2, 3]), torch.Size([2]))
```

A.sum(axis=[0,1])== A.sum()

```
→ tensor(True)
```

A.mean(), A.sum() / A.numel()

```
→ (tensor(2.5000), tensor(2.5000))
```

A.mean(axis=0), A.sum(axis=0) / A.shape[0]

```
→ (tensor([1.5000, 2.5000, 3.5000]), tensor([1.5000, 2.5000, 3.5000]))
```

```
sum_A = A.sum(axis=1, keepdims=True)
A, sum_A, sum_A.shape
```

```
↔ (tensor([[0., 1., 2.],
          [3., 4., 5.]]),
    tensor([[ 3.],
          [12.]]),
    torch.Size([2, 1]))
```

```
A/sum_A
```

```
↔ tensor([[0.0000, 0.3333, 0.6667],
          [0.2500, 0.3333, 0.4167]])
```

```
A.cumsum(axis=0)
```

```
↔ tensor([[0., 1., 2.],
          [3., 5., 7.]])
```

```
y = torch.ones(3, dtype=torch.float32)
x, y, torch.dot(x,y)
```

```
↔ (tensor([0., 1., 2.]), tensor([1., 1., 1.]), tensor(3.))
```

```
torch.sum(x*y)
```

```
↔ tensor(3.)
```

```
A.shape, x.shape, A, x, torch.mv(A,x), A@x
```

```
↔ (torch.Size([2, 3]),
    torch.Size([3]),
    tensor([[0., 1., 2.],
          [3., 4., 5.]]),
    tensor([0., 1., 2.]),
    tensor([ 5., 14.]),
    tensor([ 5., 14.]])
```

```
B = torch.ones(3,4)
torch.mm(A,B), A@B
```

```
↔ (tensor([[ 3.,  3.,  3.,  3.],
          [12., 12., 12., 12.]]),
    tensor([[ 3.,  3.,  3.,  3.],
          [12., 12., 12., 12.]])
```

```
u = torch.tensor([3.0, -4.0])
torch.norm(u)
```

```
↔ tensor(5.)
```

```
torch.abs(u).sum()
```

```
↔ tensor(7.)
```



```
torch.norm(torch.ones((4,9)))
```

```
↔ tensor(6.)
```

## ✓ 2.5. Automatic Differentiation

```
import torch
```

```
x = torch.arange(4.0)
```

```
x
```

```
↔ tensor([0., 1., 2., 3.])
```

```
x.requires_grad_(True)
```

```
x.grad
```

```
x
```

```
↔ tensor([0., 1., 2., 3.], requires_grad=True)
```

```
y = 2 * torch.dot(x,x)
```

```
y
```

```
↔ tensor(28., grad_fn=<MulBackward0>)
```

```
y.backward()
```

```
x.grad
```

```
↔ tensor([ 0.,  4.,  8., 12.])
```

```
x.grad == 4 * x
```

```
↔ tensor([True, True, True, True])
```

```
x.grad.zero_() #reset grad
```

```
y = x.sum()
```

```
print(y)
```

```
print(x)
```

```
y.backward()
```

```
x.grad
```

```
↔ tensor(6., grad_fn=<SumBackward0>)  
   tensor([0., 1., 2., 3.], requires_grad=True)  
   tensor([1., 1., 1., 1.])
```

```
x.grad.zero_()
```

```
y = x*x
```

```
y.backward(gradient = torch.ones(len(y)))
```

```
x.grad
```

```
↔ tensor([0., 2., 4., 6.])
```

```
x.grad.zero_()
y = x*x
u = y.detach()
z = u*x

z.sum().backward()
print(x.grad)
x.grad == u
```

```
→ tensor([0., 1., 4., 9.])
   tensor([True, True, True, True])
```

```
x.grad.zero_()
y.sum().backward()
x.grad==2*x
```

```
→ tensor([True, True, True, True])
```

```
def f(a):
    b = a*2
    while b.norm() < 1000:
        b = b*2
    if b.sum() > 0:
        c = b
    else:
        c = 100*b
    return c
```

```
a = torch.randn(size = (), requires_grad=True)
d = f(a)
d.backward()
```

```
a.grad == d/a
```

```
→ tensor(True)
```

## ✓ 3.0. Linear Neural Networks for Regression

### ✓ 3.1. Linear Regression

```
%matplotlib inline
import math
import time
import numpy as np
import torch
from d2l import torch as d2l
```

- Hypothesis

$$\hat{y} = \mathbf{w}^\top \mathbf{x} + b \text{ (vector - vector)}$$

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} + b \text{ (matrix - vector)}$$

- Loss function

$$l^{(i)}(\mathbf{w}, b) = \frac{1}{2} \left( \hat{y}^{(i)} - y^{(i)} \right)^2 \text{ (for example } i)$$

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \left( \mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)} \right)^2$$

(for entire training set)

- Training the model

$$\mathbf{w}^*, b^* = \underset{\mathbf{w}, b}{\operatorname{argmin}} L(\mathbf{w}, b)$$

- Minibatch Stochastic Gradient Descent

$$(\mathbf{w}, b) \leftarrow (\mathbf{w}, b) - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}_t} \partial_{(\mathbf{w}, b)} l^{(i)}(\mathbf{w}, b)$$

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}_t} \partial_{\mathbf{w}} l^{(i)}(\mathbf{w}, b) = \mathbf{w} - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}_t} \mathbf{x}^{(i)} \left( \mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)} \right)$$

$$b \leftarrow b - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}_t} \partial_b l^{(i)}(\mathbf{w}, b) = b - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}_t} \left( \mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)} \right).$$

- Vectorization for Speed

```
n = 10000
a = torch.ones(n)
b = torch.ones(n)

c = torch.zeros(n)
t = time.time()
for i in range(n):
    c[i] = a[i]+b[i]
f'{time.time() - t:.5f} sec'
```

↔ '0.28725 sec'

```
t = time.time()
d = a + b
f'{time.time() - t:.5f} sec'
```

↔ '0.00071 sec'

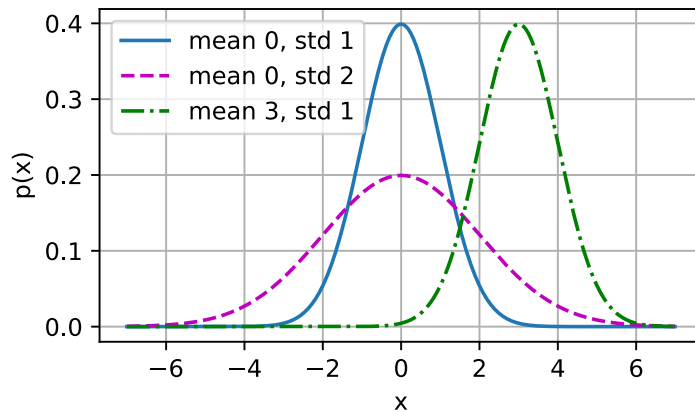
- The Normal Distribution and Squared Loss

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

```
def normal(x, mu, sigma):
    p = 1 / math.sqrt(2 * math.pi * sigma**2)
    return p * np.exp(-0.5 * (x - mu)**2 / sigma**2)
```

```
x = np.arange(-7, 7, 0.01)
```

```
params = [(0, 1), (0, 2), (3, 1)]
d2l.plot(x, [normal(x, mu, sigma) for mu, sigma in params], xlabel='x',
        ylabel='p(x)', figsize=(4.5, 2.5),
        legend=[f'mean {mu}, std {sigma}' for mu, sigma in params])
```



## ✓ 3.2. Object-Oriented Design for Implementation

```
import time
import numpy as np
import torch
from torch import nn
from d2l import torch as d2l
```

```
def add_to_class(Class):
    def wrapper(obj):
        setattr(Class, obj.__name__, obj)
    return wrapper
```

```
class A:
    def __init__(self):
        self.b = 1
```

```
a = A()
```

```
@add_to_class(A)
def do(self):
    print('Class attribute "b" is', self.b)
```

```
a.do()
```



```
Class attribute "b" is 1
```

```
class HyperParameters:
    def save_hyperparameters(self, ignore=[]):
        raise NotImplemented
```

```
class B(d2l.HyperParameters):
    def __init__(self, a, b, c):
        self.save_hyperparameters(ignore=['c'])
        print('self.a =', self.a, 'self.b =', self.b)
        print('There is no self.c =', not hasattr(self, 'c'))
```

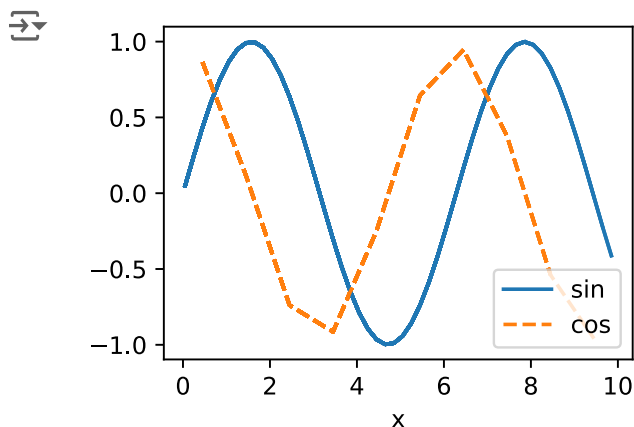
```
b = B(a=1, b=2, c=3)
```

```
self.a = 1 self.b = 2
There is no self.c = True
```

```
class ProgressBoard(d2l.HyperParameters):
    def __init__(self, xlabel=None, ylabel=None, xlim=None,
                  ylim=None, xscale='linear', yscale='linear',
                  ls=['-', '--', '-.', ':'], colors=['C0', 'C1', 'C2', 'C3'],
                  fig=None, axes=None, figsize=(3.5, 2.5), display=True):
        self.save_hyperparameters()

    def draw(self, x, y, label, every_n=1):
        raise NotImplemented
```

```
board = d2l.ProgressBoard('x')
for x in np.arange(0, 10, 0.1):
    board.draw(x, np.sin(x), 'sin', every_n=2)
    board.draw(x, np.cos(x), 'cos', every_n=10)
```



- models

```

class Module(nn.Module, d2l.HyperParameters):
    """The base class of models."""
    def __init__(self, plot_train_per_epoch=2, plot_valid_per_epoch=1):
        super().__init__()
        self.save_hyperparameters()
        self.board = ProgressBoard()

    def loss(self, y_hat, y):
        raise NotImplementedError

    def forward(self, X):
        assert hasattr(self, 'net'), 'Neural network is defined'
        return self.net(X)

    def plot(self, key, value, train):
        """Plot a point in animation."""
        assert hasattr(self, 'trainer'), 'Trainer is not initied'
        self.board.xlabel = 'epoch'
        if train:
            x = self.trainer.train_batch_idx / W
            self.trainer.num_train_batches
            n = self.trainer.num_train_batches / W
            self.plot_train_per_epoch
        else:
            x = self.trainer.epoch + 1
            n = self.trainer.num_val_batches / W
            self.plot_valid_per_epoch
        self.board.draw(x, value.to(d2l.cpu()).detach().numpy(),
                        ('train_' if train else 'val_') + key,
                        every_n=int(n))

    def training_step(self, batch):
        l = self.loss(self(*batch[:-1]), batch[-1])
        self.plot('loss', l, train=True)
        return l

    def validation_step(self, batch):
        l = self.loss(self(*batch[:-1]), batch[-1])
        self.plot('loss', l, train=False)

    def configure_optimizers(self):
        raise NotImplementedError

```

- Data

```

class DataModule(d2l.HyperParameters):
    """The base class of data."""
    def __init__(self, root='./data', num_workers=4):
        self.save_hyperparameters()

    def get_dataloader(self, train):
        raise NotImplementedError

    def train_dataloader(self):
        return self.get_dataloader(train=True)

    def val_dataloader(self):
        return self.get_dataloader(train=False)

```

- Training

```
class Trainer(d2l.HyperParameters):
    """The base class for training models with data."""
    def __init__(self, max_epochs, num_gpus=0, gradient_clip_val=0):
        self.save_hyperparameters()
        assert num_gpus == 0, 'No GPU support yet'

    def prepare_data(self, data):
        self.train_dataloader = data.train_dataloader()
        self.val_dataloader = data.val_dataloader()
        self.num_train_batches = len(self.train_dataloader)
        self.num_val_batches = (len(self.val_dataloader)
                                if self.val_dataloader is not None else 0)

    def prepare_model(self, model):
        model.trainer = self
        model.board.xlim = [0, self.max_epochs]
        self.model = model

    def fit(self, model, data):
        self.prepare_data(data)
        self.prepare_model(model)
        self.optim = model.configure_optimizers()
        self.epoch = 0
        self.train_batch_idx = 0
        self.val_batch_idx = 0
        for self.epoch in range(self.max_epochs):
            self.fit_epoch()

    def fit_epoch(self):
        raise NotImplementedError
```

### ✓ 3.4. Linear Regression Implementation from Scratch

```
%matplotlib inline
import torch
from d2l import torch as d2l
```

- Defining the model

```
class LinearRegressionScratch(d2l.Module):
    """The linear regression model implemented from scratch."""
    def __init__(self, num_inputs, lr, sigma=0.01):
        super().__init__()
        self.save_hyperparameters()
        self.w = torch.normal(0, sigma, (num_inputs, 1), requires_grad=True)
        self.b = torch.zeros(1, requires_grad=True)

@d2l.add_to_class(LinearRegressionScratch)
def forward(self, X):
    return torch.matmul(X, self.w) + self.b
```

- Defining the Loss Function

```
@d2l.add_to_class(LinearRegressionScratch)
def loss(self, y_hat, y):
    l = (y_hat - y) ** 2 / 2
    return l.mean()
```

- Defining the Optimization Algorithm

```
class SGD(d2l.HyperParameters):
    """Minibatch stochastic gradient descent."""
    def __init__(self, params, lr):
        self.save_hyperparameters()

    def step(self):
        for param in self.params:
            param -= self.lr * param.grad

    def zero_grad(self):
        for param in self.params:
            if param.grad is not None:
                param.grad.zero_()

@d2l.add_to_class(LinearRegressionScratch)
def configure_optimizers(self):
    return SGD([self.w, self.b], self.lr)
```

- Training

```
@d2l.add_to_class(d2l.Trainer)
def prepare_batch(self, batch):
    return batch

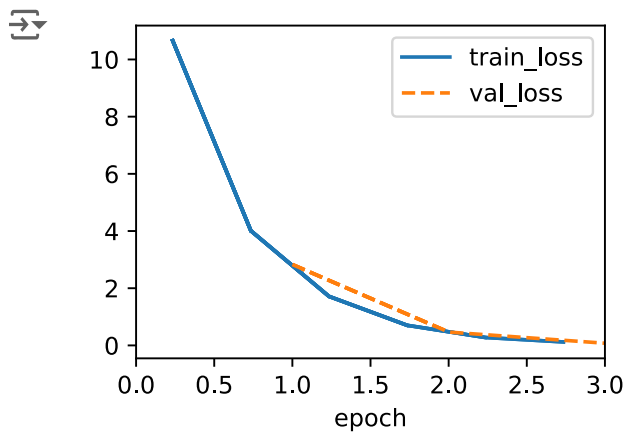
@d2l.add_to_class(d2l.Trainer)
def fit_epoch(self):
    self.model.train()
    for batch in self.train_dataloader:
        loss = self.model.training_step(self.prepare_batch(batch))
        self.optim.zero_grad()
        with torch.no_grad():
            loss.backward()
            if self.gradient_clip_val > 0: # To be discussed later
                self.clip_gradients(self.gradient_clip_val, self.model)
            self.optim.step()
        self.train_batch_idx += 1
    if self.val_dataloader is None:
        return
    self.model.eval()
    for batch in self.val_dataloader:
        with torch.no_grad():
            self.model.validation_step(self.prepare_batch(batch))
        self.val_batch_idx += 1
```



```

model = LinearRegressionScratch(2, lr=0.03)
data = d2l.SyntheticRegressionData(w=torch.tensor([2, -3.4]), b=4.2)
trainer = d2l.Trainer(max_epochs=3)
trainer.fit(model, data)

```



```

with torch.no_grad():
    print(f'error in estimating w: {data.w - model.w.reshape(data.w.shape)}')
    print(f'error in estimating b: {data.b - model.b}')

```

```

error in estimating w: tensor([ 0.1858, -0.2393])
error in estimating b: tensor([0.2372])

```

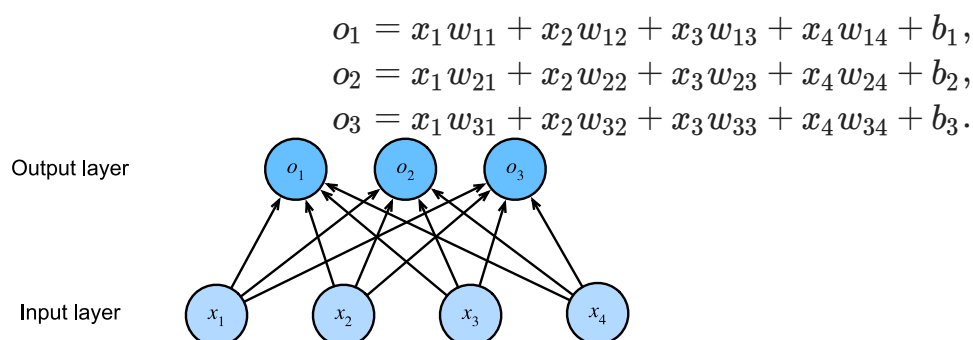
## ✓ 4.0. Linear Neural Networks for Classification

### ✓ 4.1. Softmax Regression

- Classification

1. Using integers: Label the categories as  $y \in \{1, 2, 3\}$ , which would be a natural choice for representing the categories. This approach works well if there's a natural order among the labels (e.g., age groups). Such problems can be treated as ordinal regression.
2. One-hot encoding: Since classification problems typically don't have natural orderings, one-hot encoding is used to represent categorical data. Here, each category is represented by a binary vector where one element is set to 1, and the rest are 0. For example, (1,0,0) for "cat," (0,1,0) for "chicken," and (0,0,1) for "dog."

- Linear Model



- Softmax

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{o}) \quad \text{where} \quad \hat{y}_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)}.$$

$$\underset{j}{\operatorname{argmax}} \hat{y}_j = \underset{j}{\operatorname{argmax}} o_j.$$

- Vectorization

$$\mathbf{O} = \mathbf{XW} + \mathbf{b},$$

$$\hat{\mathbf{Y}} = \text{softmax}(\mathbf{O}).$$

- Log-Likelihood

$$P(\mathbf{Y} | \mathbf{X}) = \prod_{i=1}^n P(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})$$

$$-\log P(\mathbf{Y} | \mathbf{X}) = \sum_{i=1}^n -\log P(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}) = \sum_{i=1}^n l(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)})$$

$$l(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{j=1}^q y_j \log \hat{y}_j$$

- Softmax and Cross-Entropy Loss

$$l(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{j=1}^q y_j \log \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)}$$

$$= \sum_{j=1}^q y_j \log \sum_{k=1}^q \exp(o_k) - \sum_{j=1}^q y_j o_j$$

$$= \log \sum_{k=1}^q \exp(o_k) - \sum_{j=1}^q y_j o_j.$$

$$\partial_{o_j} l(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)} - y_j = \text{softmax}(\mathbf{o})_j - y_j$$

- Entropy

$$H[P] = \sum_j -P(j) \log P(j)$$

## ✓ 4.2. The Image Classification Dataset

```
%matplotlib inline
import time
import torch
import torchvision
from torchvision import transforms
from d2l import torch as d2l

d2l.use_svg_display()
```

- Loading the Dataset

```
class FashionMNIST(d2l.DataModule):
    """The Fashion-MNIST dataset."""
    def __init__(self, batch_size=64, resize=(28, 28)):
        super().__init__()
        self.save_hyperparameters()
        trans = transforms.Compose([transforms.Resize(resize),
                                    transforms.ToTensor()])
        self.train = torchvision.datasets.FashionMNIST(
            root=self.root, train=True, transform=trans, download=True)
        self.val = torchvision.datasets.FashionMNIST(
            root=self.root, train=False, transform=trans, download=True)
```

```
data = FashionMNIST(resize=(32, 32))
len(data.train), len(data.val)
```

```
↗ (60000, 10000)
```

```
data.train[0][0].shape
```

```
↗ torch.Size([1, 32, 32])
```

```
@d2l.add_to_class(FashionMNIST)
def text_labels(self, indices):
    """Return text labels."""
    labels = ['t-shirt', 'trouser', 'pullover', 'dress', 'coat',
              'sandal', 'shirt', 'sneaker', 'bag', 'ankle boot']
    return [labels[int(i)] for i in indices]
```

- Reading a Minibatch

```
@d2l.add_to_class(FashionMNIST)
def get_dataloader(self, train):
    data = self.train if train else self.val
    return torch.utils.data.DataLoader(data, self.batch_size, shuffle=train,
                                       num_workers=self.num_workers)
```

```
X, y = next(iter(data.train_dataloader()))
print(X.shape, X.dtype, y.shape, y.dtype)
```

```
↗ /usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:557: UserWarning: This DataLoader
  warnings.warn(_create_warning_msg(
  torch.Size([64, 1, 32, 32]) torch.float32 torch.Size([64]) torch.int64
```

```
tic = time.time()
for X, y in data.train_dataloader():
    continue
f'{time.time() - tic:.2f} sec'
```

```
↗ '13.16 sec'
```

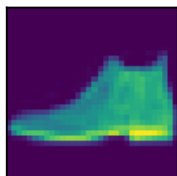
- Visualization

```
def show_images(imgs, num_rows, num_cols, titles=None, scale=1.5):
    """Plot a list of images."""
    raise NotImplementedError

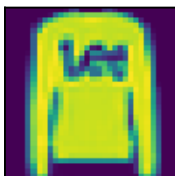
@d2l.add_to_class(FashionMNIST)
def visualize(self, batch, nrows=1, ncols=8, labels=[]):
    X, y = batch
    if not labels:
        labels = self.text_labels(y)
    d2l.show_images(X.squeeze(1), nrows, ncols, titles=labels)
batch = next(iter(data.val_dataloader()))
data.visualize(batch)
```



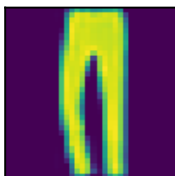
ankle boot



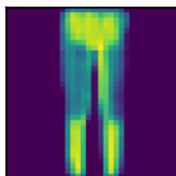
pullover



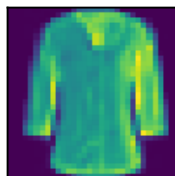
trouser



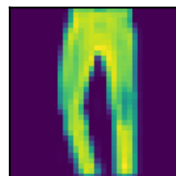
trouser



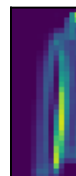
shirt



trouser



c



## ✓ 4.3. The Base Classification Model

```
import torch
from d2l import torch as d2l
```

- The Classifier Class

```
class Classifier(d2l.Module):
    """The base class of classification models."""
    def validation_step(self, batch):
        Y_hat = self(*batch[:-1])
        self.plot('loss', self.loss(Y_hat, batch[-1]), train=False)
        self.plot('acc', self.accuracy(Y_hat, batch[-1]), train=False)
```

```
@d2l.add_to_class(d2l.Module)
def configure_optimizers(self):
    return torch.optim.SGD(self.parameters(), lr=self.lr)
```

- Accuracy

```
@d2l.add_to_class(Classifier)
def accuracy(self, Y_hat, Y, averaged=True):
    """Compute the number of correct predictions."""
    Y_hat = Y_hat.reshape((-1, Y_hat.shape[-1]))
    preds = Y_hat.argmax(axis=1).type(Y.dtype)
    compare = (preds == Y.reshape(-1)).type(torch.float32)
    return compare.mean() if averaged else compare
```

## ✓ 4.4. Softmax Regression Implementation from Scratch

```
import torch
from d2l import torch as d2l
```

### • The Softmax

```
X = torch.tensor([[1.0, 2.0, 3.0], [4.0, 5.0, 6.0]])
X.sum(0, keepdims=True), X.sum(1, keepdims=True)
```

```
→ (tensor([[5., 7., 9.]]),
    tensor([[ 6.],
            [15.]])
```

```
def softmax(X):
    X_exp = torch.exp(X)
    partition = X_exp.sum(1, keepdims=True)
    return X_exp / partition
```

```
X = torch.rand((2, 5))
X_prob = softmax(X)
X_prob, X_prob.sum(1)
```

```
→ (tensor([[0.2695, 0.2052, 0.1874, 0.1914, 0.1465],
            [0.1540, 0.1453, 0.2855, 0.2425, 0.1727]]),
    tensor([1.0000, 1.0000]))
```

### • The Model

```
class SoftmaxRegressionScratch(d2l.Classifier):
    def __init__(self, num_inputs, num_outputs, lr, sigma=0.01):
        super().__init__()
        self.save_hyperparameters()
        self.W = torch.normal(0, sigma, size=(num_inputs, num_outputs),
                                requires_grad=True)
        self.b = torch.zeros(num_outputs, requires_grad=True)

    def parameters(self):
        return [self.W, self.b]
```

```
@d2l.add_to_class(SoftmaxRegressionScratch)
def forward(self, X):
    X = X.reshape((-1, self.W.shape[0]))
    return softmax(torch.matmul(X, self.W) + self.b)
```

### • The Cross-Entropy Loss

```
y = torch.tensor([0, 2])
y_hat = torch.tensor([[0.1, 0.3, 0.6], [0.3, 0.2, 0.5]])
y_hat[[0, 1], y]
```

```
→ tensor([0.1000, 0.5000])
```

```
def cross_entropy(y_hat, y):
    return -torch.log(y_hat[list(range(len(y_hat))), y]).mean()
```

```
cross_entropy(y_hat, y)
```

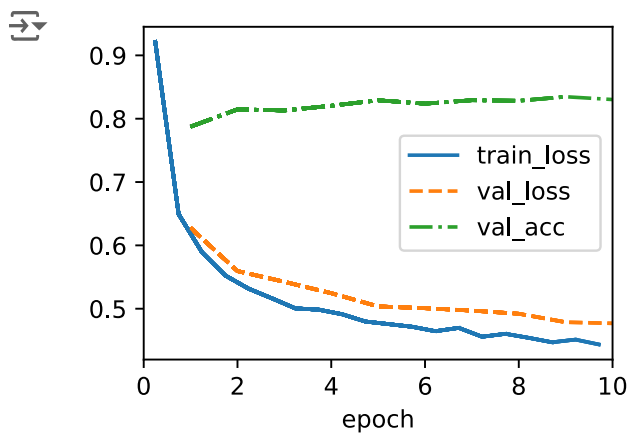
```
→ tensor(1.4979)
```

```
@d2l.add_to_class(SoftmaxRegressionScratch)
```

```
def loss(self, y_hat, y):
    return cross_entropy(y_hat, y)
```

### • Training

```
data = d2l.FashionMNIST(batch_size=256)
model = SoftmaxRegressionScratch(num_inputs=784, num_outputs=10, lr=0.1)
trainer = d2l.Trainer(max_epochs=10)
trainer.fit(model, data)
```



### • Prediction

```
X, y = next(iter(data.val_dataloader()))
preds = model(X).argmax(axis=1)
preds.shape
```

```
→ torch.Size([256])
```

```
wrong = preds.type(y.dtype) != y
X, y, preds = X[wrong], y[wrong], preds[wrong]
labels = [a+'Wn'+b for a, b in zip(
    data.text_labels(y), data.text_labels(preds))]
data.visualize([X, y], labels=labels)
```



## ✓ 5.0. Multilayer Perceptrons

### ✓ 5.1. Multilayer Perceptrons

```
%matplotlib inline
import torch
from d2l import torch as d2l
```

- Hidden Layer : From Linear to Nonlinear

$$\mathbf{H} = \sigma(\mathbf{XW}^{(1)} + \mathbf{b}^{(1)}) \text{ (Hidden layer)}$$

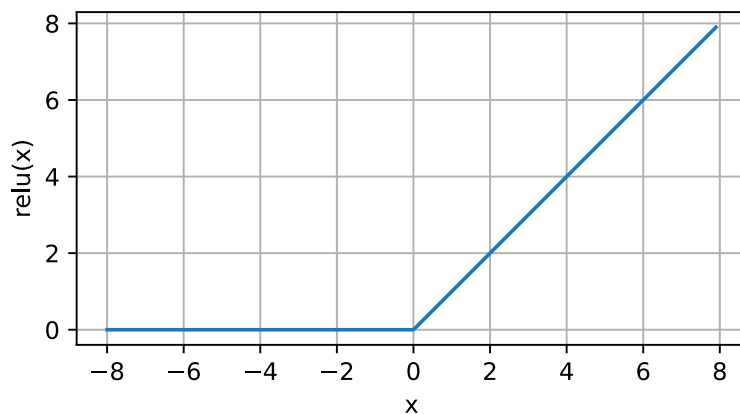
$$\mathbf{O} = \mathbf{HW}^{(2)} + \mathbf{b}^{(2)} \text{ (Output)}$$

$\sigma()$  : activation function

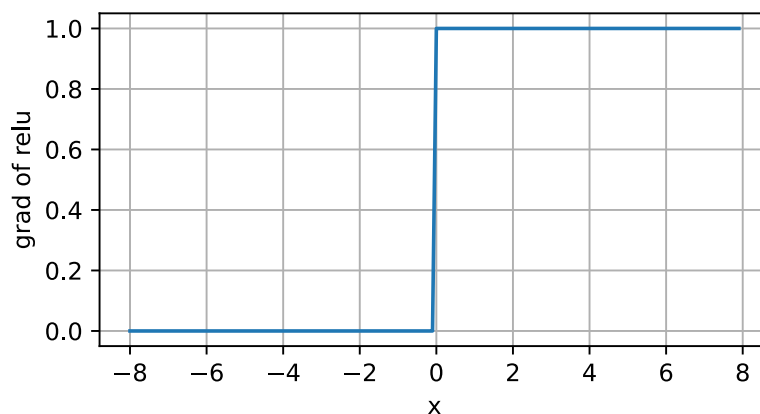
- Activation function : ReLU Function

$$\text{ReLU}(x) = \max(x, 0)$$

```
x = torch.arange(-8.0, 8.0, 0.1, requires_grad=True)
y = torch.relu(x)
d2l.plot(x.detach(), y.detach(), 'x', 'relu(x)', figsize=(5, 2.5))
```



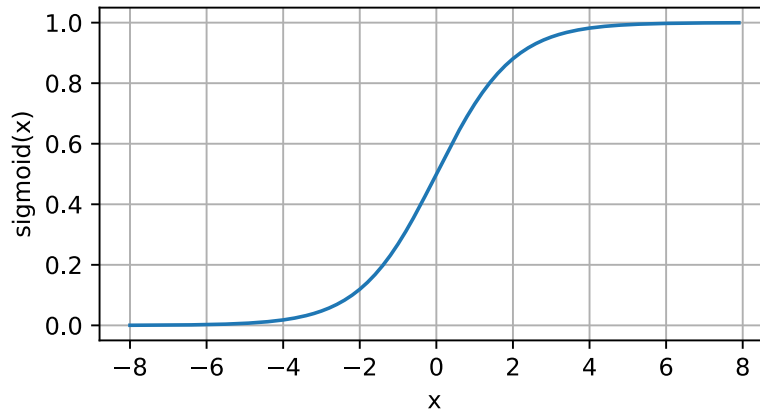
```
y.backward(torch.ones_like(x), retain_graph=True)
d2l.plot(x.detach(), x.grad, 'x', 'grad of relu', figsize=(5, 2.5))
```



- Activation function : Sigmoid Function

$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

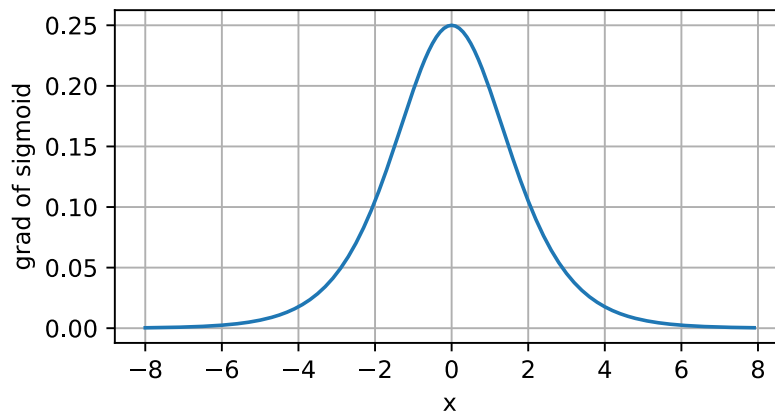
```
y = torch.sigmoid(x)
d2l.plot(x.detach(), y.detach(), 'x', 'sigmoid(x)', figsize=(5, 2.5))
```



- Derivative of sigmoid function

$$\frac{d}{dx} \text{sigmoid}(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2} = \text{sigmoid}(x) (1 - \text{sigmoid}(x))$$

```
x.grad.data.zero_()
y.backward(torch.ones_like(x), retain_graph=True)
d2l.plot(x.detach(), x.grad, 'x', 'grad of sigmoid', figsize=(5, 2.5))
```

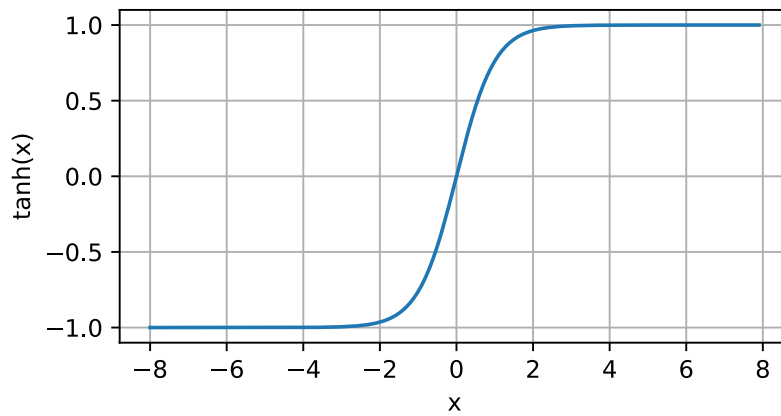


- Activation function : Tanh Function

$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)}$$

```
y = torch.tanh(x)
d2l.plot(x.detach(), y.detach(), 'x', 'tanh(x)', figsize=(5, 2.5))
```

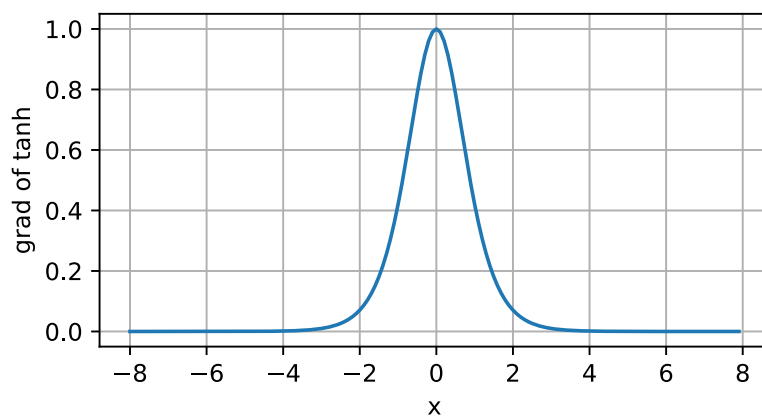




- Derivative of the tanh function

$$\frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

```
x.grad.data.zero_()
y.backward(torch.ones_like(x), retain_graph=True)
d2l.plot(x.detach(), x.grad, 'x', 'grad of tanh', figsize=(5, 2.5))
```



## ✓ 5.2. Implementation of Multilayer Perceptrons

```
import torch
from torch import nn
from d2l import torch as d2l
```

```
class MLPScratch(d2l.Classifier):
    def __init__(self, num_inputs, num_outputs, num_hiddens, lr, sigma=0.01):
        super().__init__()
        self.save_hyperparameters()
        self.W1 = nn.Parameter(torch.randn(num_inputs, num_hiddens) * sigma)
        self.b1 = nn.Parameter(torch.zeros(num_hiddens))
        self.W2 = nn.Parameter(torch.randn(num_hiddens, num_outputs) * sigma)
        self.b2 = nn.Parameter(torch.zeros(num_outputs))
```

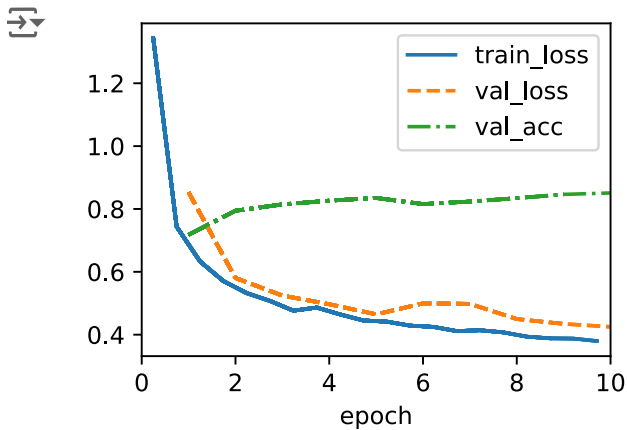
- Model

```
def relu(X):
    a = torch.zeros_like(X)
    return torch.max(X, a)

@d2l.add_to_class(MLPScratch)
def forward(self, X):
    X = X.reshape((-1, self.num_inputs))
    H = relu(torch.matmul(X, self.W1) + self.b1)
    return torch.matmul(H, self.W2) + self.b2
```

- Training

```
model = MLPScratch(num_inputs=784, num_outputs=10, num_hiddens=256, lr=0.1)
data = d2l.FashionMNIST(batch_size=256)
trainer = d2l.Trainer(max_epochs=10)
trainer.fit(model, data)
```



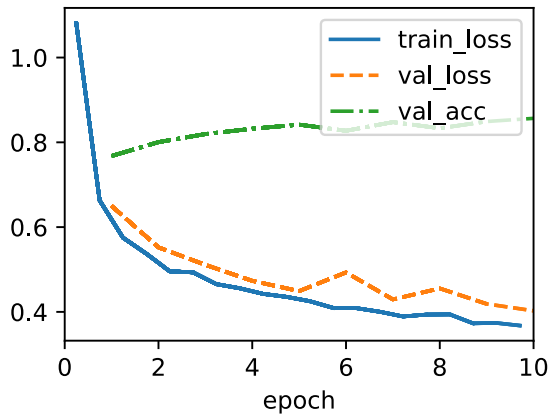
- Concise Implementation

- Model

```
class MLP(d2l.Classifier):
    def __init__(self, num_outputs, num_hiddens, lr):
        super().__init__()
        self.save_hyperparameters()
        self.net = nn.Sequential(nn.Flatten(), nn.LazyLinear(num_hiddens),
                                   nn.ReLU(), nn.LazyLinear(num_outputs))
```

- Training

```
model = MLP(num_outputs=10, num_hiddens=256, lr=0.1)
trainer.fit(model, data)
```



### ✓ 5.3. Forward Propagation, Backward Propagation, and Computational Graphs

- Forward Propagation : calculation and storage of intermediate variables (including outputs) for a neural network in order from the input layer to the output layer

1. Input Example: Suppose the input example is  $\mathbf{x} \in \mathbb{R}^d$ , and the hidden layer does not include a bias term. The intermediate variable is defined as:

$$\mathbf{z} = \mathbf{W}^{(1)}\mathbf{x},$$

2. Activation Function: The intermediate variable  $z$  is passed through the activation function  $\phi$  producing the hidden layer activation vector  $\mathbf{h} \in \mathbb{R}^h$ :

$$\mathbf{h} = \phi(\mathbf{z})$$

3. Hidden Layer Output: The hidden layer output  $h$  is also an intermediate variable. Assuming the output layer has a weight  $\mathbf{W}^{(2)} \in \mathbb{R}^{q \times h}$ , the output vector  $\mathbf{o}$  is obtained as:

$$\mathbf{o} = \mathbf{W}^{(2)}\mathbf{h}.$$

4. Loss Function: Given a loss function  $l$  and the label for the example  $y$ , the loss for a single data example is calculated as:

$$L = l(\mathbf{o}, y)$$

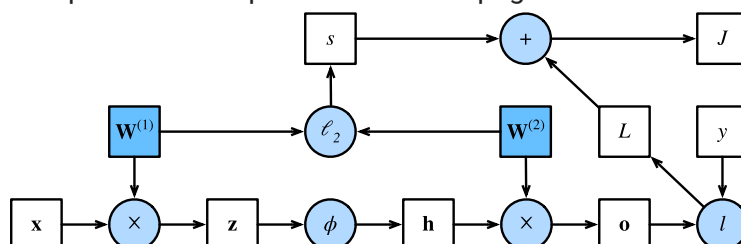
5. Regularization Term: When introducing  $l_2$  regularization later, and given the hyperparameter  $\lambda$ , the regularization term is defined as:

$$s = \frac{\lambda}{2} \left( \|\mathbf{W}^{(1)}\|_F^2 + \|\mathbf{W}^{(2)}\|_F^2 \right)$$

6. Regularized Loss: Finally, the model's regularized loss for a given data example is:

$$J = L + s$$

- Computational Graph of Forward Propagation



- Backpropagation

1. Objective function :  $J = L + s$

2. Derivative :  $\frac{\partial J}{\partial L} = 1$  and  $\frac{\partial J}{\partial s} = 1$

3. Gradient of objective

$$\frac{\partial J}{\partial \mathbf{o}} = \text{prod} \left( \frac{\partial J}{\partial L}, \frac{\partial L}{\partial \mathbf{o}} \right) = \frac{\partial L}{\partial \mathbf{o}} \in \mathbb{R}^q$$

4. Gradient of the regularization term

$$\frac{\partial s}{\partial \mathbf{W}^{(1)}} = \lambda \mathbf{W}^{(1)} \text{ and } \frac{\partial s}{\partial \mathbf{W}^{(2)}} = \lambda \mathbf{W}^{(2)}$$

5. Gradient of the model parameters closest to the output layer

$$\frac{\partial J}{\partial \mathbf{W}^{(2)}} = \text{prod} \left( \frac{\partial J}{\partial \mathbf{o}}, \frac{\partial \mathbf{o}}{\partial \mathbf{W}^{(2)}} \right) + \text{prod} \left( \frac{\partial J}{\partial s}, \frac{\partial s}{\partial \mathbf{W}^{(2)}} \right) = \frac{\partial J}{\partial \mathbf{o}} \mathbf{h}^\top + \lambda \mathbf{W}^{(2)}$$

6. Gradient with respect to the hidden layer output

$$\frac{\partial J}{\partial \mathbf{h}} = \text{prod} \left( \frac{\partial J}{\partial \mathbf{o}}, \frac{\partial \mathbf{o}}{\partial \mathbf{h}} \right) = \mathbf{W}^{(2)\top} \frac{\partial J}{\partial \mathbf{o}}$$

7. gradient of the intermediate variable  $z$

$$\frac{\partial J}{\partial \mathbf{z}} = \text{prod} \left( \frac{\partial J}{\partial \mathbf{h}}, \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \right) = \frac{\partial J}{\partial \mathbf{h}} \odot \phi'(\mathbf{z})$$

8. Gradient of the model parameters closest to the input layer

$$\frac{\partial J}{\partial \mathbf{W}^{(1)}} = \text{prod} \left( \frac{\partial J}{\partial \mathbf{z}}, \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} \right) + \text{prod} \left( \frac{\partial J}{\partial s}, \frac{\partial s}{\partial \mathbf{W}^{(1)}} \right) = \frac{\partial J}{\partial \mathbf{z}} \mathbf{x}^\top + \lambda \mathbf{W}^{(1)}.$$

- Training Neural Networks : When training deep learning models, forward propagation and backpropagation are interdependent, and training requires significantly more memory than prediction.

## ✓ Discussions and Exercises

### ✓ 2.2.4 : Complexities of data processing(Discussion)

- **Complexities of data processing** : Real-world data processing can be much more complicated. Data often arrives from multiple sources, such as relational databases, and can include various types like text, images, and audio. Efficient tools and algorithms are needed to manage this complexity and avoid bottlenecks in machine learning pipelines, especially in areas like computer vision and natural language processing. Additionally, attention to data quality is crucial since real-world datasets may contain errors, outliers, or faulty measurements. Data visualization tools like seaborn, Bokeh, and matplotlib can help in inspecting and understanding data before feeding it into a model.

### ✓ 2.5.5. : Power of automatic differentiation(Discussion)

- **Power of automatic differentiation** : The development of libraries that can calculate derivatives automatically and efficiently has greatly boosted productivity for deep learning practitioners. These tools allow the design of massive models, for which manual gradient computation would be too time-consuming. Interestingly, while autograd is used to optimize models statistically, optimizing autograd libraries themselves for computational efficiency is an important topic for framework designers. Compilers and graph manipulation tools are employed to compute results quickly and in a memory-efficient way. Basic principles to remember are: (i) attaching gradients to variables for which derivatives are needed, (ii) recording the computation of the target value, (iii) performing backpropagation, and (iv) accessing the resulting gradient

#### ✓ 4.1.4 : Computational aspects(Discussion)

- **Computational aspects** : it does not deeply explore computational aspects. Specifically, for a fully connected layer with  $d$  inputs and  $q$  outputs, the computational cost is  $O(dq)$ , which can be very high in practice. Fortunately, approximation and compression methods, such as those used in Deep Fried Convnets (Yang et al., 2015), can reduce this cost from quadratic to linear using techniques like