

› Installation

[] ↪ 숨겨진 셀 1개

✓ 2. 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.ones((2, 3, 4))
```

```
⇒ tensor([[[1., 1., 1., 1.],
           [1., 1., 1., 1.],
           [1., 1., 1., 1.]])
```

```
[1., 1., 1., 1.]],
[[1., 1., 1., 1.],
 [1., 1., 1., 1.],
 [1., 1., 1., 1.]])
```

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

```
⇒ tensor([[ 0.8052, -0.0158, -1.2208, -0.3906],
          [ 0.0487,  0.4455,  1.2307,  0.6482],
          [ 0.6113, -0.5732, -0.2939, -0.0564]])
```

```
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]])
```

```
X.sum()
```

```
⇒ tensor(66.)
```

```
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): 135739245029200
   id(Z): 135739245029200
```

```
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)
```

✓ Discussions & Exercises

2.1.1. Getting Started

To automatically infer one component of the shape, we can place a -1 for the shape component that should be inferred automatically. In our case, instead of calling `x.reshape(3, 4)`, we could have equivalently called `x.reshape(-1, 4)` or `x.reshape(3, -1)`.

2.1.4. Broadcasting

(i) Expand one or both arrays by copying elements along axes with length 1 so that after this transformation, the two tensors have the same shape.

(ii) Perform an elementwise operation on the resulting arrays.

✓ 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))
```

✓ Discussions & Exercises

2.2.2. Data Preparation

In supervised learning, we train models to predict a designated *target value*, given some set of *input* values. Our first step in processing the dataset is to separate out columns corresponding to input versus target values. We can select columns either by name or via integer-location based indexing (`iloc`).

Pandas replace all CSV entries with value NA with a special NaN (*not a number*) value. This can also happen whenever an entry is empty, e.g., “3,,,270000”. These are called *missing values* and they are the “bed bugs” of data science, a persistent menace that you will confront throughout your career. Depending upon the context, missing values might be handled either via *imputation* or *deletion*. Imputation replaces missing values with estimates of their values while deletion simply discards either those rows or those columns that contain missing values.

✓ 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.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)
```

```
sum_A, sum_A.shape
```

```
⇒ (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, torch.mv(A, x), A@x
```



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

✓ Discussions & Exercises

2.3.11. Norms

ℓ_2 norm is expressed as $\| \mathbf{x} \|_2 = \sqrt{\sum_{i=1}^n x_i^2}$.

ℓ_1 norm is also common and the associated measure is called the Manhattan distance:

$$\| \mathbf{x} \|_1 = \sum_{i=1}^n |x_i|.$$

Both the ℓ_2 and ℓ_1 norms are special cases of the more general ℓ_p norms:

$$\| \mathbf{x} \|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}.$$

The Frobenius norm behaves as if it were an ℓ_2 norm of a matrix-shaped vector:

$$\| \mathbf{X} \|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n x_{ij}^2}.$$

✓ 2.5. Automatic Differentiation

```
import torch
```

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

```
⇒ tensor([0., 1., 2., 3.])
```

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

```
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_()
y = x.sum()
y.backward()
x.grad
```

```
⇒ 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()
x.grad == u
```

```
⇒ 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)
```

✓ Discussions & Exercises

As we pass data through each successive function, the framework builds a *computational graph* that tracks how each value depends on others. To calculate derivatives, *automatic differentiation* (often shortened to *autograd*) works backwards through this graph applying the chain rule, called *backpropagation*.

✓ 3. 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

```

```

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'

```

```
⇒ '0.21952 sec'
```

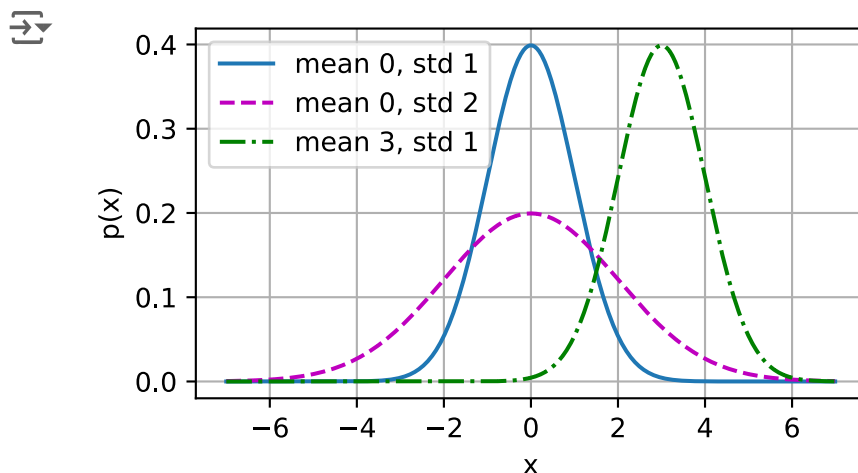
```
t = time.time()
d = a + b
f'{time.time() - t:.5f} sec'
```

```
⇒ '0.00052 sec'
```

```
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])
```



✓ Discussions & Exercises

Regression problems pop up whenever we want to predict a numerical value.

3.1.1.1. Model

Given a dataset, our goal is to choose the weights \mathbf{w} and the bias b that, on average, make our model's predictions fit the true values observed in the data as closely as possible.

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} + b$$

3.1.1.2. Loss Function

Loss functions quantify the distance between the *real* and *predicted* values of the target. The loss will usually be a nonnegative number where smaller values are better and perfect predictions incur a loss of 0. For regression problems, the most common loss function is the

squared error.

3.1.1.4. Minibatch Stochastic Gradient Descent

The key technique for optimizing nearly every deep learning model, is called *gradient descent*, which consists of iteratively reducing the error by updating the parameters in the direction that incrementally lowers the loss function. The most naive application of gradient descent consists of taking the derivative of the loss function, which is an average of the losses computed on every single example in the dataset.

The other extreme is to consider only a single example at a time and to take update steps based on one observation at a time, which is called *stochastic gradient descent* (SGD). Unfortunately, SGD has drawbacks, both computational and statistical.

The solution to both problems is to pick an intermediate strategy: rather than taking a full batch or only a single sample at a time, we take a *minibatch* of observations, which leads us to *minibatch stochastic gradient descent*.

✓ 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):
    """Register functions as methods in created 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:
    """The base class of 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'))
```

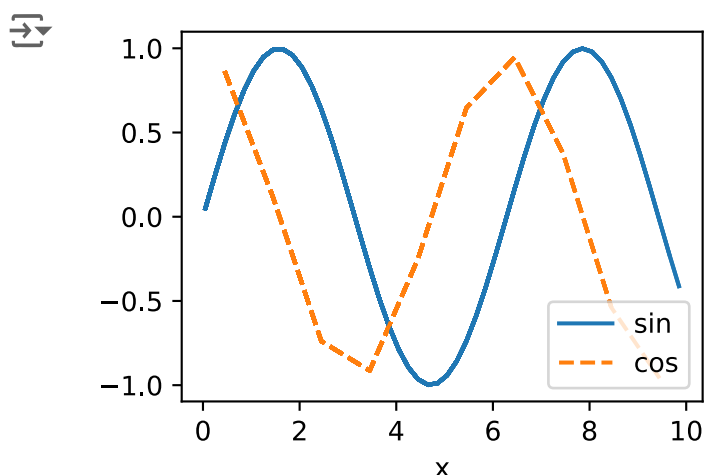
```
b = B(a=1, b=2, c=3)
```

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

```
class ProgressBoard(d2l.HyperParameters):
    """The board that plots data points in animation."""
    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)
```



```
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 / \
                self.trainer.num_train_batches
            n = self.trainer.num_train_batches / \
                self.plot_train_per_epoch
        else:
            x = self.trainer.epoch + 1
            n = self.trainer.num_val_batches / \
                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

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)

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

```

✓ Discussions & Exercises

3.2.2. Models

The `Module` class is the base class of all models we will implement. At the very least we need three methods: `__init__`, `training_step` and `configure_optimizers`.

3.2.3. Data

The `DataModule` class is the base class for data.

3.2.4. Training

The `Trainer` class trains the learnable parameters in the `Module` class with data specified in `DataModule`.

✓ 3.4. Linear Regression Implementation from Scratch

```

%matplotlib inline
import torch
from d2l import torch as d2l

```



```

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

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

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)

@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:
                self.clip_gradients(self.gradient_clip_val, self.model)
            self.optim.step()
        self.train_batch_idx += 1

```

```

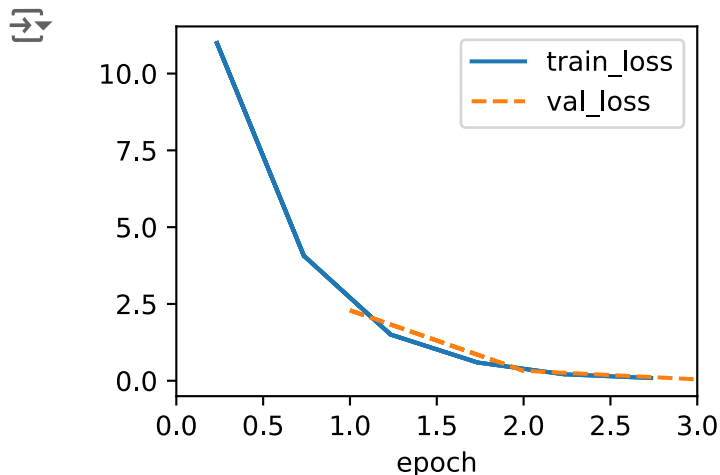
if self.val_dataloader is None:
    return
self.model.eval()
for batch in self.val_dataloader:
    with torch.no_grad():
        self.model.validation_step(self.prepare_batch(batch))
    self.val_batch_idx += 1

```

```

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([ 0.0969, -0.1792])
error in estimating b: tensor([0.2222])

```

✓ Discussions & Exercises

3.4.4. Training

In each *epoch*, we iterate through the entire training dataset, passing once through every example (assuming that the number of examples is divisible by the batch size). In each *iteration*, we grab a minibatch of training examples, and compute its loss through the model's `training_step` method. Then we compute the gradients with respect to each parameter. Finally, we will call the optimization algorithm to update the model parameters.

✓ 4. Linear Neural Networks for Classification

✓ 4.1. Softmax Regression

✓ Discussions & Exercises

Classification problems focus on *which category?* questions, describing two subtly different problems:

- (i) Those where we are interested only in hard assignments of examples to categories (classes).
- (ii) Those where we wish to make soft assignments.

Multi-label classification is a variant of the classification problem where multiple nonexclusive labels may be assigned to each instance.

4.1.1. Classification

Ordinal regression is a type of regression analysis used for predicting an ordinal variable, i.e. a variable whose value exists on an arbitrary scale where only the relative ordering between different values is significant.

But in general, classification problems do not come with natural orderings among the classes. A one-hot encoding is a vector with as many components as we have categories. The component corresponding to a particular instance's category is set to 1 and all other components are set to 0.

4.1.1.2. The Softmax

- There is no guarantee that the outputs o_i sum up to 1 in the way we expect probabilities to behave.
- There is no guarantee that the outputs o_i are even nonnegative, even if their outputs sum up to 1, or that they do not exceed 1.

Use an exponential function $P(y = i) \propto \exp o_i$ and then transform these values so that they add up to 1 by dividing each by their sum (*normalization*).

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

4.1.2.1. Log-Likelihood

The loss function l is $l(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{j=1}^q y_j \log \hat{y}_j$, called the *cross-entropy loss*.

4.1.3.1. Entropy

For a distribution P its *entropy*, $H[P]$, is defined as:

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

✓ 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()
```

```
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)
```

➡ Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/train-images-idx3-ubyte.gz>
 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/train-labels-idx1-ubyte.gz>
 100%|██████████| 26421880/26421880 [00:06<00:00, 4229522.60it/s]
 Extracting ../data/FashionMNIST/raw/train-images-idx3-ubyte.gz to ../data/FashionMNIST/raw/train-images-idx3-ubyte

Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/train-images-idx3-ubyte.gz>
 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/train-labels-idx1-ubyte.gz>
 100%|██████████| 29515/29515 [00:00<00:00, 195309.05it/s]
 Extracting ../data/FashionMNIST/raw/train-labels-idx1-ubyte.gz to ../data/FashionMNIST/raw/train-labels-idx1-ubyte

Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/t10k-images-idx3-ubyte.gz>
 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/t10k-labels-idx1-ubyte.gz>
 100%|██████████| 4422102/4422102 [00:01<00:00, 3693091.55it/s]
 Extracting ../data/FashionMNIST/raw/t10k-images-idx3-ubyte.gz to ../data/FashionMNIST/raw/t10k-images-idx3-ubyte

Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/t10k-labels-idx1-ubyte.gz>
 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/t10k-labels-idx1-ubyte.gz>
 100%|██████████| 5148/5148 [00:00<00:00, 5502619.01it/s]Extracting ../data/FashionMNIST/raw/t10k-labels-idx1-ubyte.gz to ../data/FashionMNIST/raw/t10k-labels-idx1-ubyte

```
(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]
```

```
@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: U:
  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'
```

```
↳ '12.12 sec'
```

```
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
					

▼ Discussions & Exercises

4.2.4. Summary

Data iterators are a key component for efficient performance. For instance, we might use GPUs for efficient image decompression, video transcoding, or other preprocessing. Whenever possible, rely on well-implemented data iterators that exploit high-performance computing to avoid slowing down the training loop.

✓ 4.3. The Base Classification Model

```
import torch
from d2l import torch as d2l

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)

@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
```

✓ Discussions & Exercises

4.3.2. Accuracy

When predictions are consistent with the label class y , they are correct. The classification accuracy is the fraction of all predictions that are correct. Although it can be difficult to optimize accuracy directly (it is not differentiable), it is often the performance measure that we care about the most.

✓ 4.4. Softmax Regression Implementation from Scratch

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

```
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([[0.2113, 0.1566, 0.3439, 0.1491, 0.1392],
         [0.2738, 0.1266, 0.1799, 0.2304, 0.1894]]),
    tensor([1., 1.]])
```

```
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)
```

```
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])
```

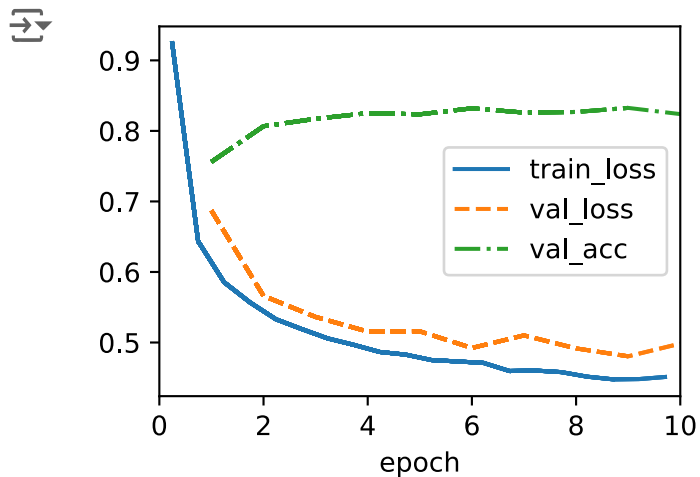
```
def cross_entropy(y_hat, y):
    return -torch.log(y_hat[list(range(len(y_hat))), y]).mean()

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)
```

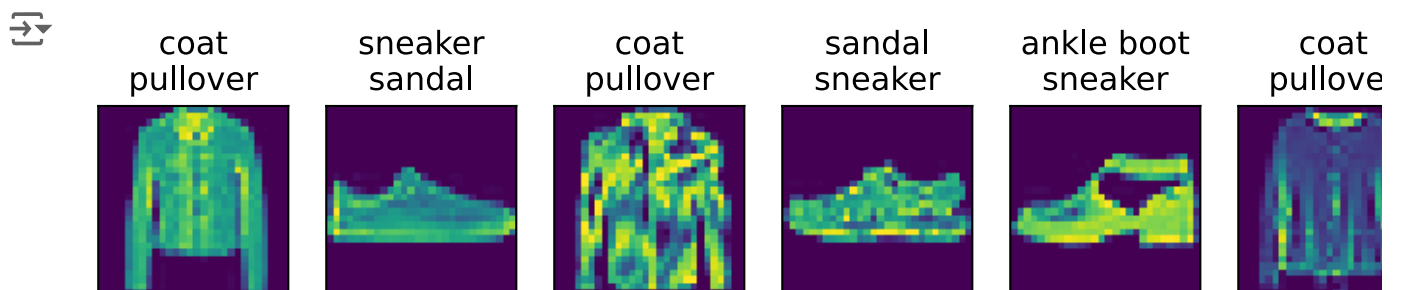
```
data = d2l.FashionMNIST(batch_size=256)
model = SoftmaxRegressionScratch(num_inputs=784, num_outputs=10, lr=0.1)
trainer = d2l.Trainer(max_epochs=10)
trainer.fit(model, data)
```



```
X, y = next(iter(data.val_data_loader()))
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+'\n'+b for a, b in zip(
    data.text_labels(y), data.text_labels(preds))]
data.visualize([X, y], labels=labels)
```



▼ Discussions & Exercises

4.4.1. The Softmax

Computing the softmax requires three steps:

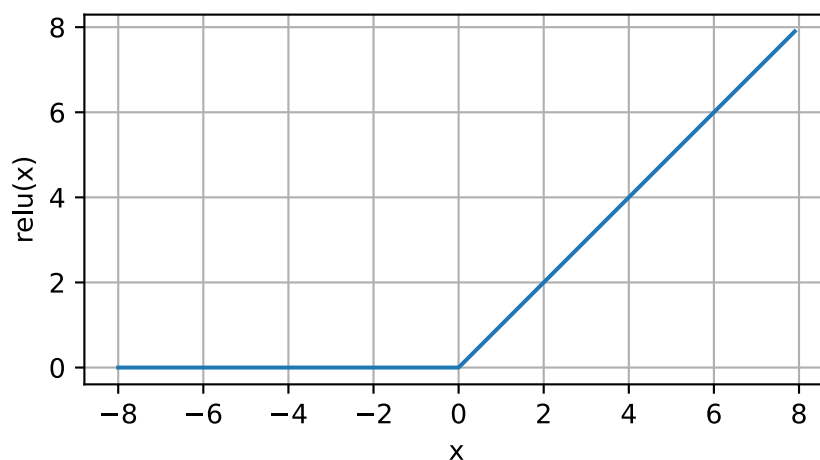
- (i) Exponentiation of each term.
- (ii) A sum over each row to compute the normalization constant for each example.
- (iii) Division of each row by its normalization constant, ensuring that the result sums to 1.

✓ 5. Multilayer Perceptrons

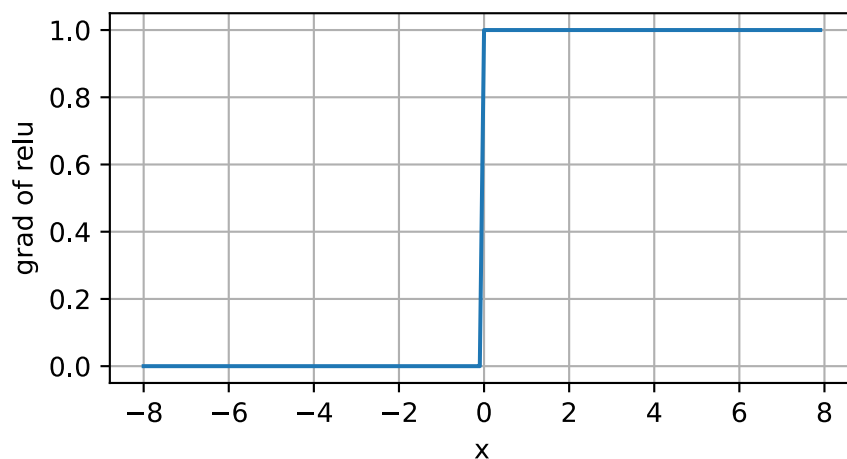
✓ 5.1. Multilayer Perceptrons

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

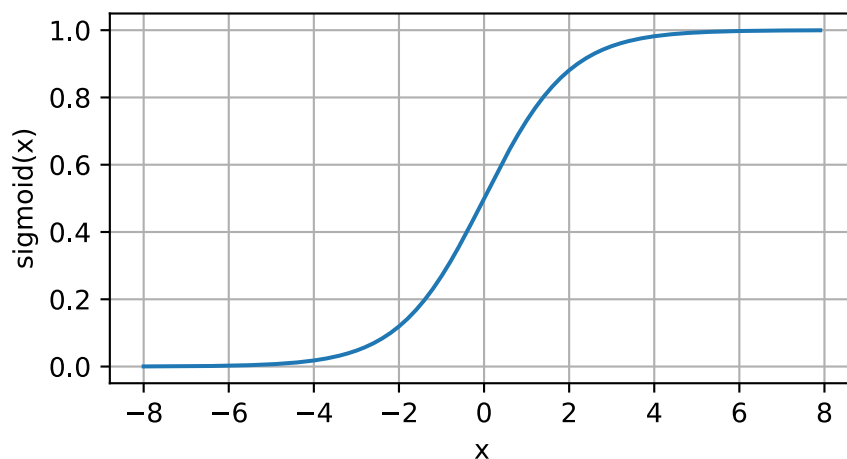
```
x = torch.arange(-8.0, 8.0, 0.1, requires_grad=True)
y = torch.relu(x)
d2l.plot(x.detach(), y.detach(), 'x', 'relu(x)', figsize=(5, 2.5))
```



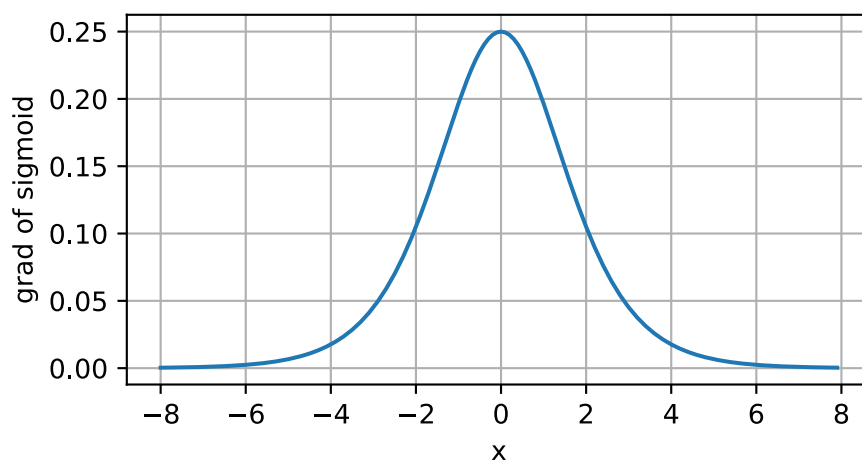
```
y.backward(torch.ones_like(x), retain_graph=True)
d2l.plot(x.detach(), x.grad, 'x', 'grad of relu', figsize=(5, 2.5))
```



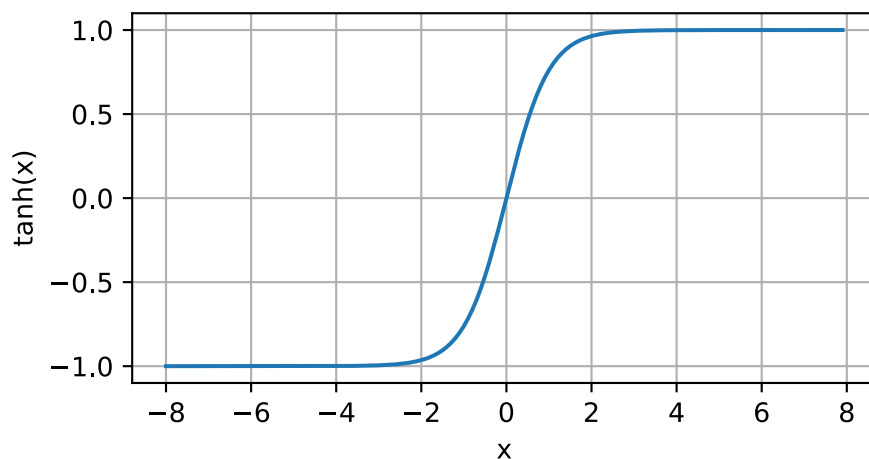
```
y = torch.sigmoid(x)
d2l.plot(x.detach(), y.detach(), 'x', 'sigmoid(x)', figsize=(5, 2.5))
```



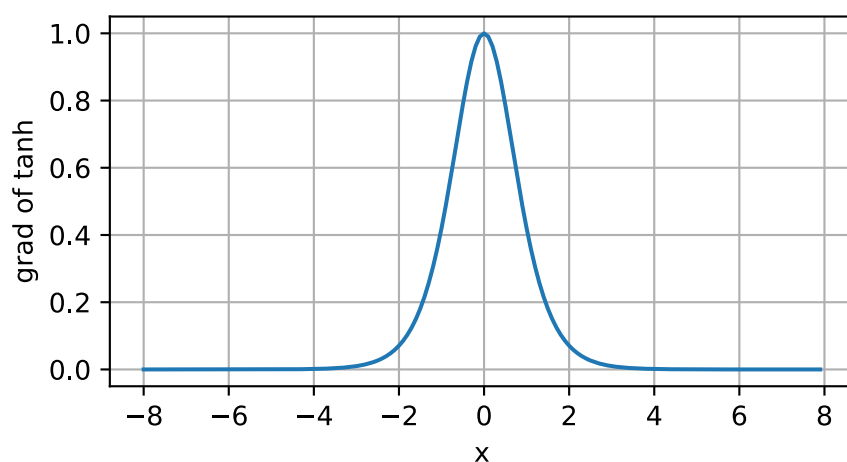
```
x.grad.data.zero_()
y.backward(torch.ones_like(x), retain_graph=True)
d2l.plot(x.detach(), x.grad, 'x', 'grad of sigmoid', figsize=(5, 2.5))
```



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



```
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))
```



✓ Discussions & Exercises

5.1.1.2. Incorporating Hidden Layers

We can overcome the limitations of linear models by incorporating one or more hidden layers. The easiest way to do this is to stack many fully connected layers on top of one another. Each layer feeds into the layer above it, until we generate outputs. We can think of the first $L - 1$ layers as our representation and the final layer as our linear predictor. This architecture is commonly called a *multilayer perceptron*, often abbreviated as *MLP*.

5.1.2. Activation Functions

Activation functions decide whether a neuron should be activated or not by calculating the weighted sum and further adding bias to it. They are differentiable operators for transforming input signals to outputs, while most of them add nonlinearity.

- ReLU Function

- Sigmoid Function

- Tanh Function

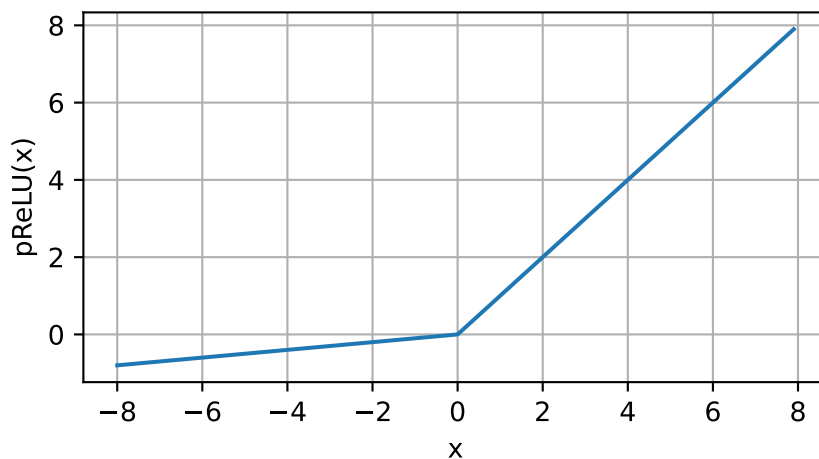
```
pReLU = lambda x, a: torch.max(torch.tensor(0), x) + a * torch.min(torch.tens
```

```
y = pReLU(x=x, a=0.1)
```

+ 코드

+ 텍스트

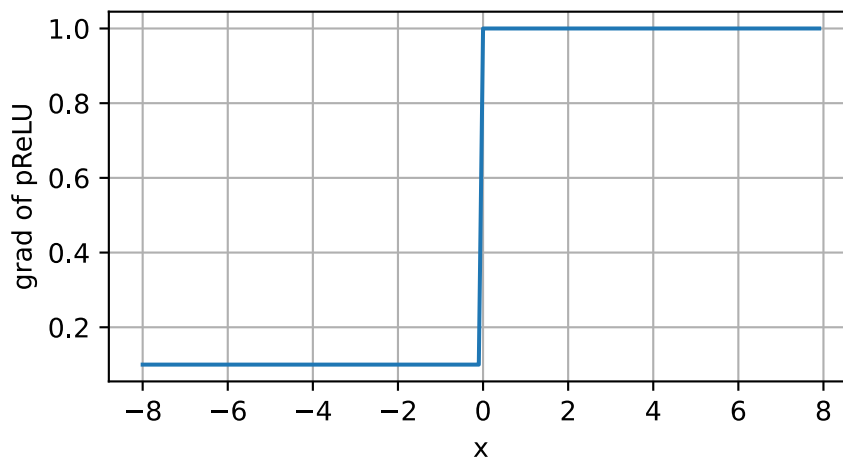
```
d2l.plot(x.detach(), y.detach(), 'x', 'pReLU(x)', figsize=(5,2.5))
```



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

```
y.backward(torch.ones_like(x), retain_graph=True)
```

```
d2l.plot(x.detach(), x.grad, 'x', 'grad of pReLU', 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))

```

```

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

```

```

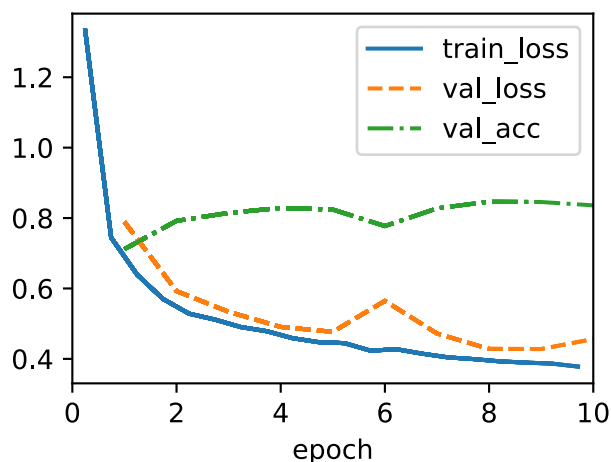
@d2l.add_to_class(MLPScratch)
def forward(self, X):
    X = X.reshape((-1, self.num_inputs))
    H = relu(torch.matmul(X, self.W1) + self.b1)
    return torch.matmul(H, self.W2) + self.b2

```

```

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)

```



```

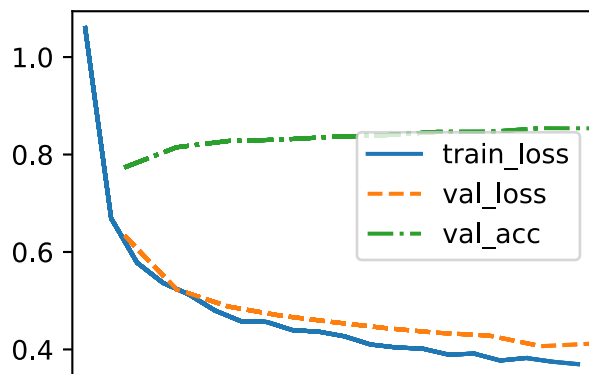
class MLP(d2l.Classifier):
    def __init__(self, num_outputs, num_hiddens, lr):
        super().__init__()
        self.save_hyperparameters()
        self.net = nn.Sequential(nn.Flatten(), nn.LazyLinear(num_hiddens),
                                  nn.ReLU(), nn.LazyLinear(num_outputs))

```

```

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

```



✓ Discussions & Exercises

5.2.3. Summary

Implementing MLPs from scratch is nonetheless messy: naming and keeping track of the model parameters makes it difficult to extend models. For instance, imagine wanting to insert another layer between layers 42 and 43. This might now be layer 42b, unless we are willing to perform sequential renaming. Moreover, if we implement the network from scratch, it is much more difficult for the framework to perform meaningful performance optimizations.

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

✓ Discussions & Exercises

5.3.1. Forward Propagation

Forward propagation (or forward pass) refers to the calculation and storage of intermediate variables (including outputs) for a neural network in order from the input layer to the output layer.