



CAMBRIDGE

Lecture 3: Neural Networks

Stefan Bucher

MACHINE LEARNING IN ECONOMICS
UNIVERSITY OF CAMBRIDGE

Stefan Bucher

 Open in Colab

Prince (2023, chaps. 3, 4, 7).¹

- . Figures taken or adapted from Prince (2023). All rights belong to the original author and publisher. These materials are intended solely for educational purposes.

Shallow Neural Networks

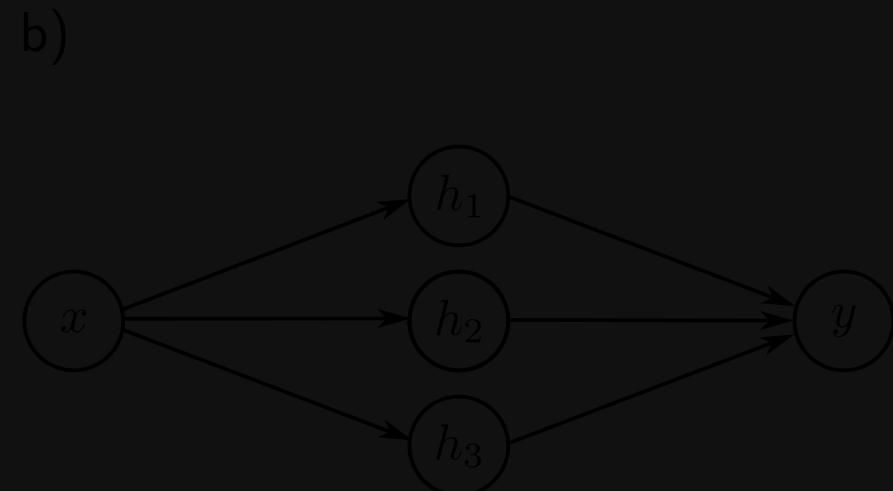
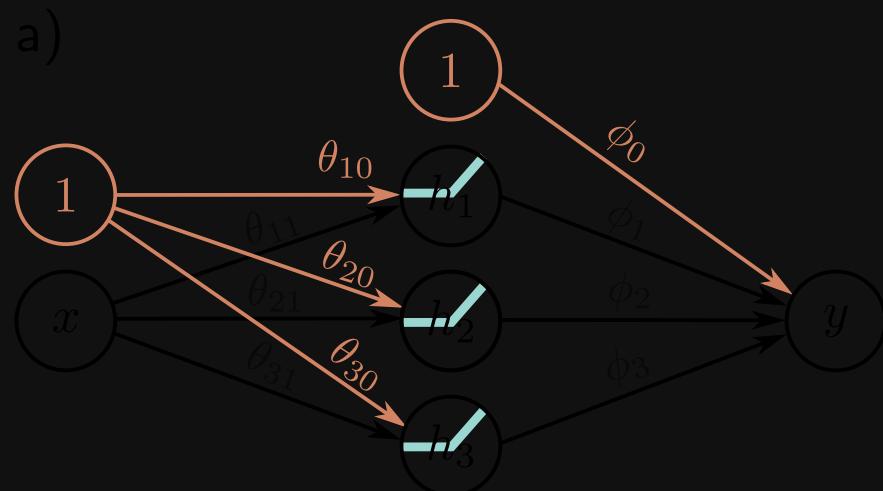
Prince (2023, chap. 3)

Example

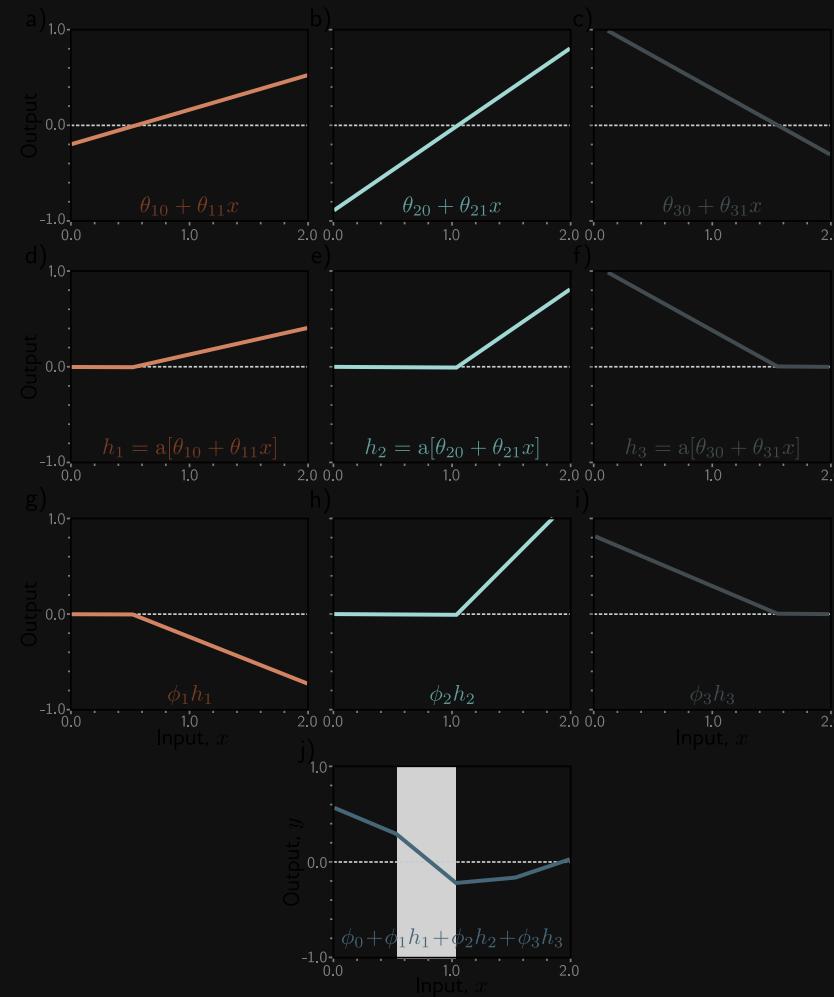
$$f[x, \phi] = \phi_0 + \phi_1 a(\underbrace{\theta_{10} + \theta_{11}x}_{h_1}) + \phi_2 a(\underbrace{\theta_{20} + \theta_{21}x}_{h_2}) + \phi_3 a(\underbrace{\theta_{30} + \theta_{31}x}_{h_3})$$

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

Neural networks (shallow & deep!) with ReLU represent continuous piecewise linear functions.



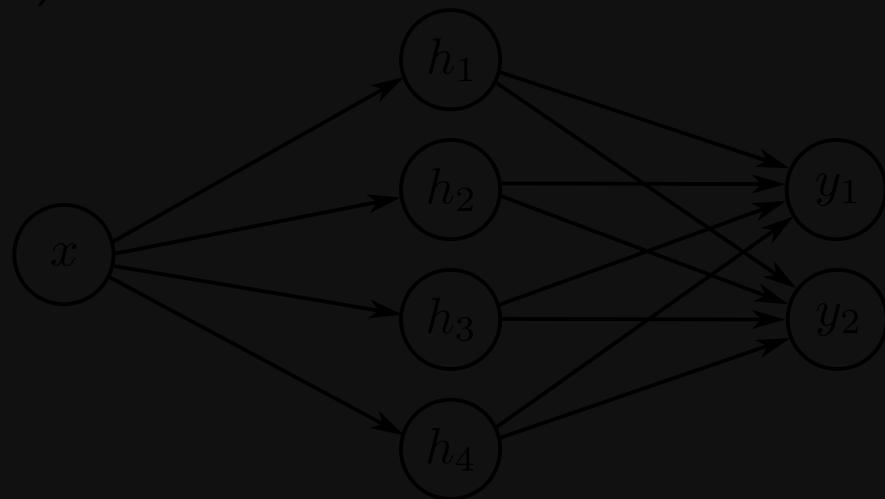
Activation Patterns



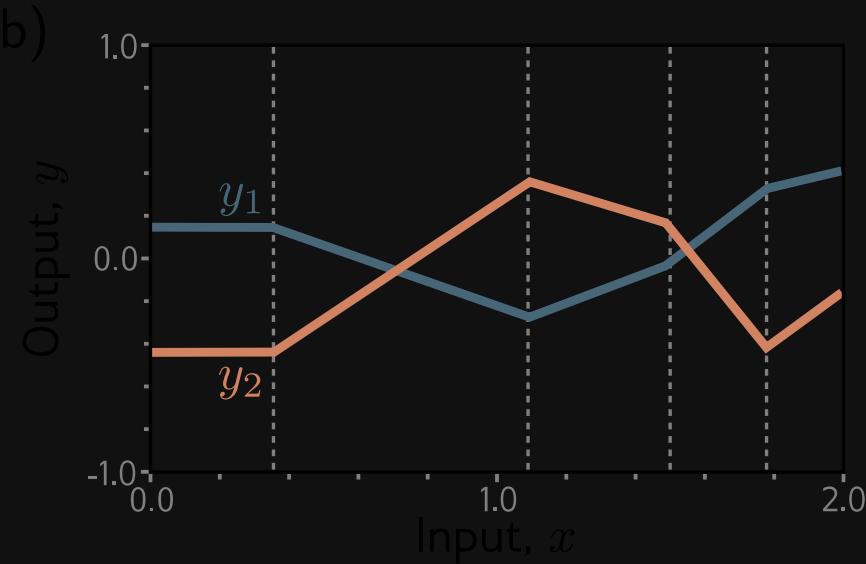
Joints where one hidden unit becomes (in)active.

Multivariate Output

a)



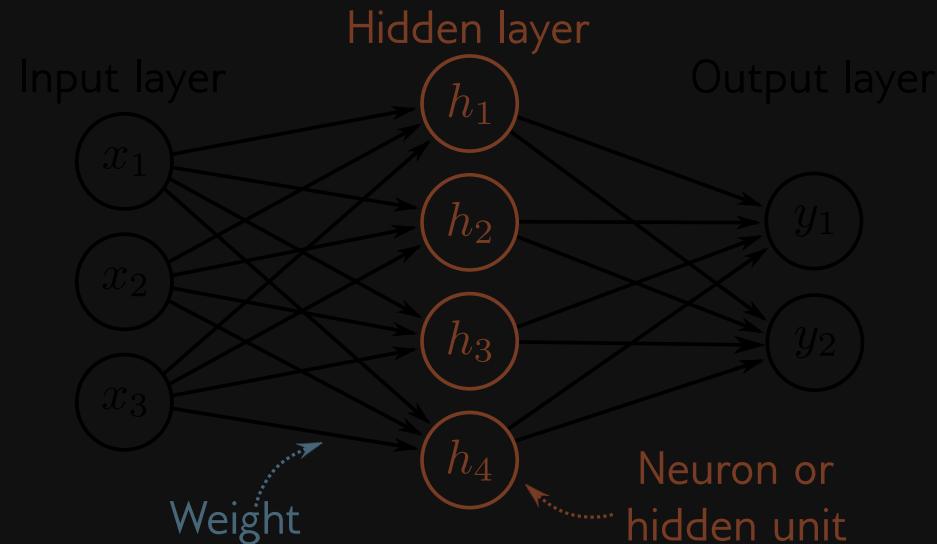
b)



Shallow Neural Networks

Map n_{in} -dimensional input to n_{out} -dimensional output.

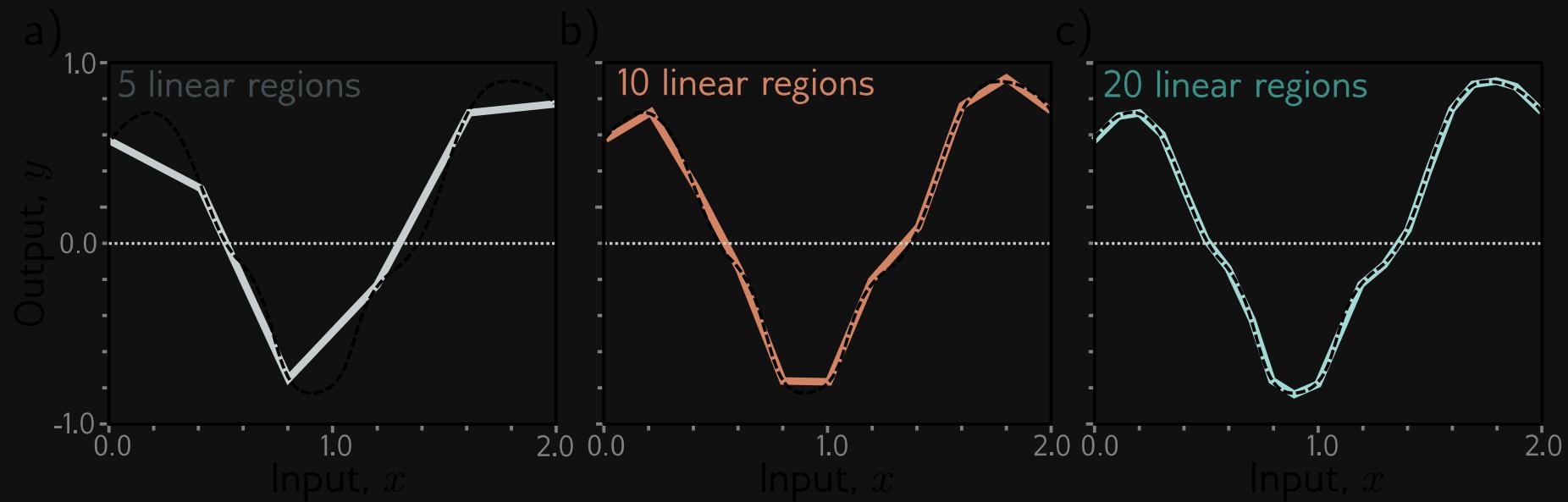
$$y_j = \phi_{j0} + \sum_{d=1}^D \phi_{jd} a(\theta_{d0} + \underbrace{\sum_{i=1}^{n_{in}} \theta_{di} x_i}_{h_d}), \quad j = 1, \dots, n_{out}$$



Shallow Networks (3.3 & 3.8)

Universal Approximation Thm

\exists shallow neural network with a single layer of sufficiently many hidden units that can approximate, with arbitrary precision, any continuous function from a compact subset of $R^{n_{in}}$ to $R^{n_{out}}$.



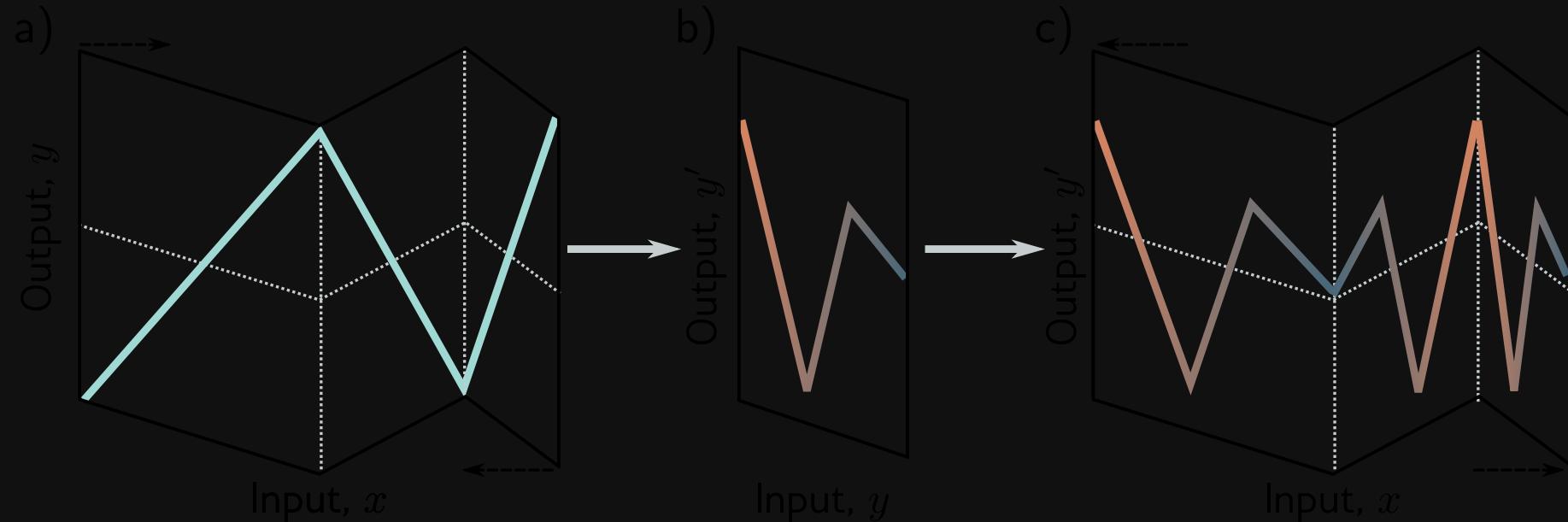
Deep Neural Networks

Prince (2023, chap. 4)

Why Depth? (4.1)

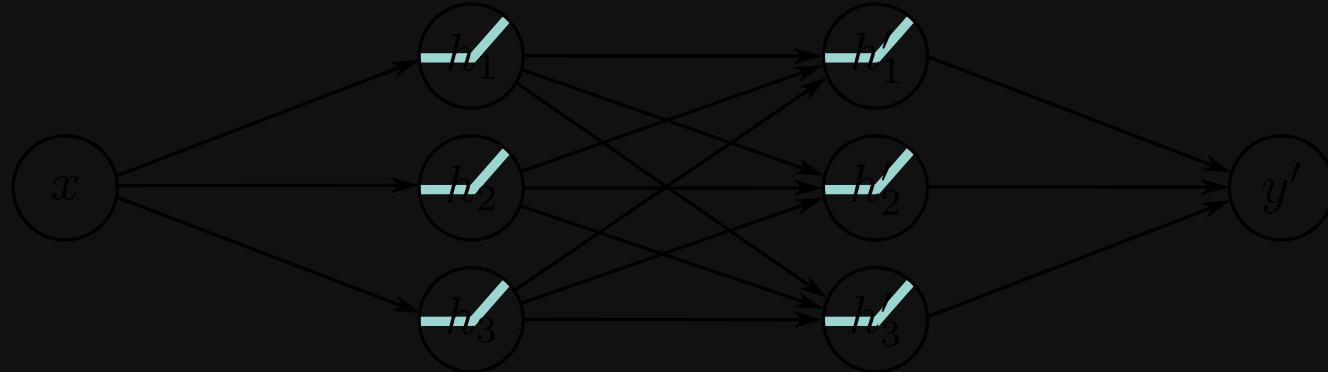
- Universal approximation theorem...
- Deep: more linear regions for same # parameters

Extra Layer is Folding



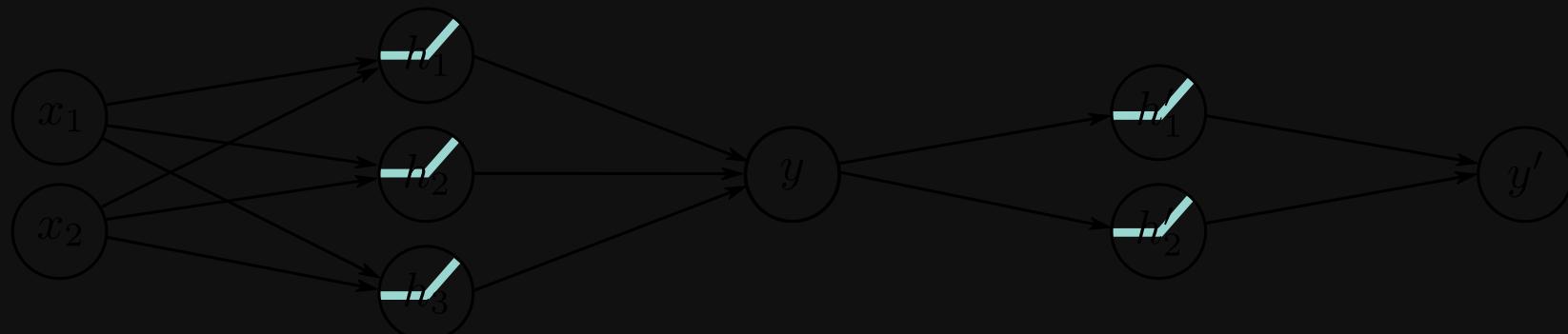
Extra Layer is Clipping (4.5)

Composing as Special Case

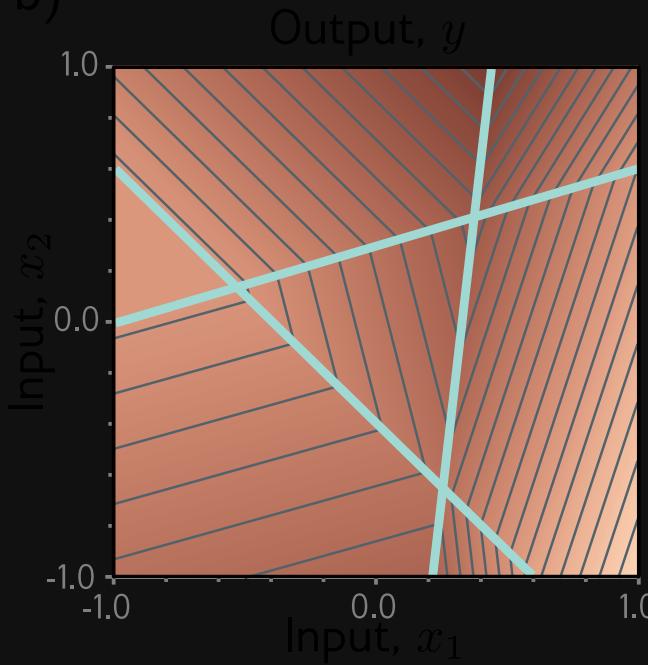


Composing in 2D

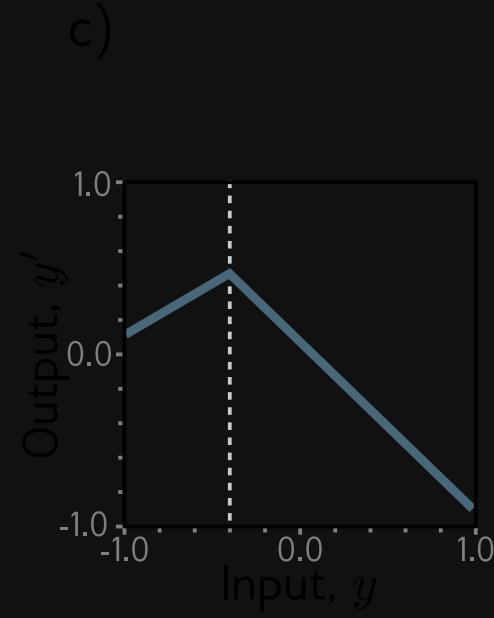
a)



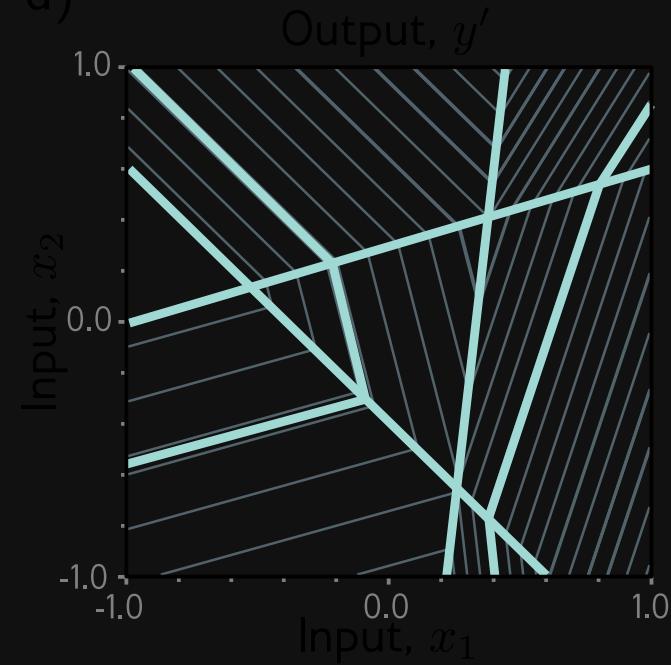
b)



c)



d)

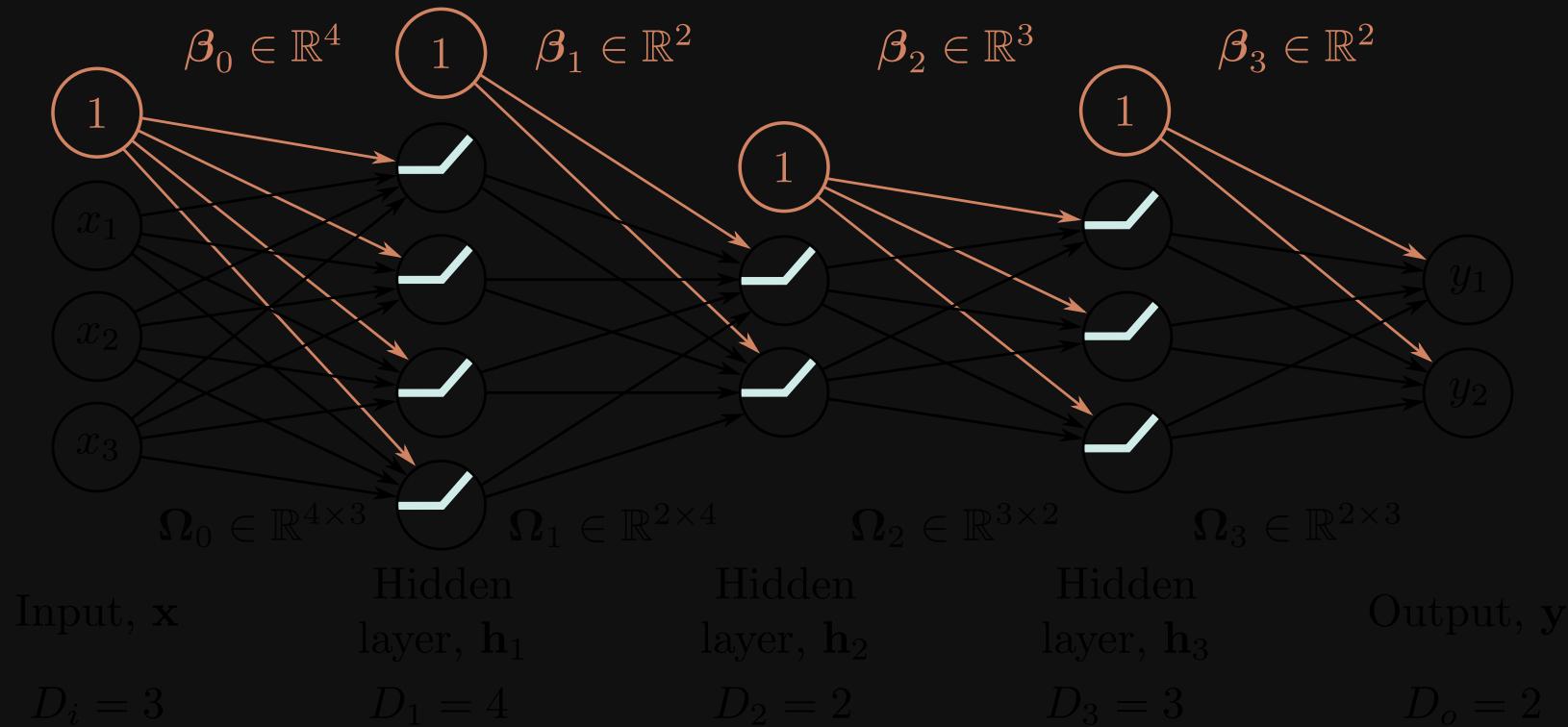


Deep Neural Network

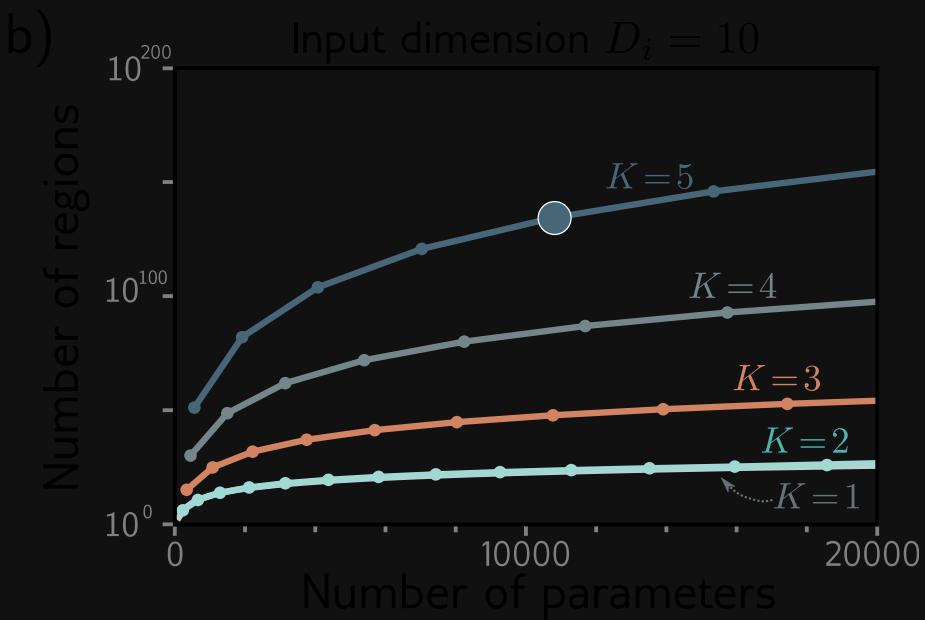
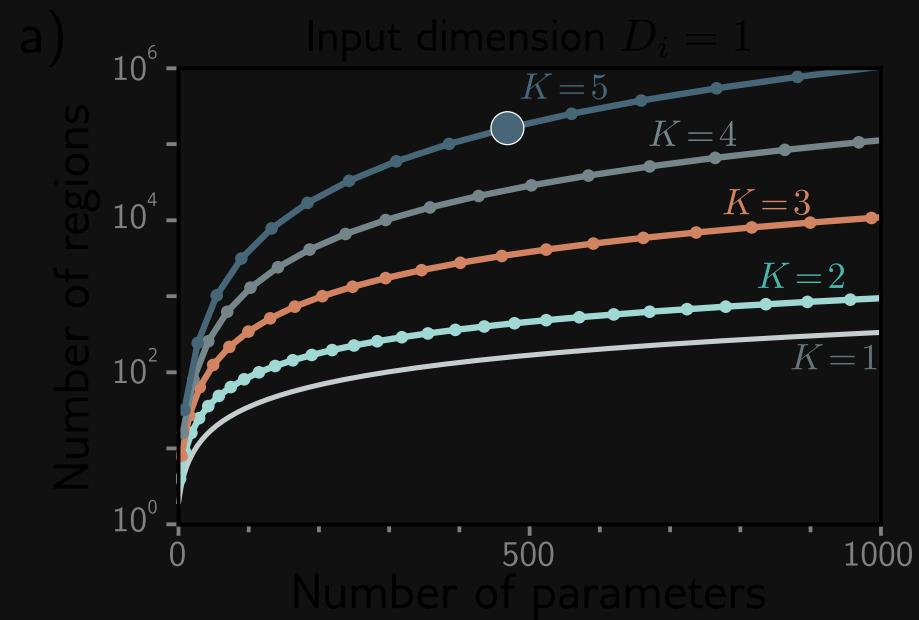
$$\mathbf{h}_1 = \mathbf{a}(\beta_0 + \boldsymbol{\Omega}_0 \mathbf{x})$$

$$\mathbf{h}_k = \mathbf{a}(\beta_{k-1} + \boldsymbol{\Omega}_{k-1} \mathbf{h}_{k-1}), \quad k = 2, \dots, K$$

$$\mathbf{y} = \beta_K + \boldsymbol{\Omega}_K \mathbf{h}_K$$



Depth Efficiency



Backpropagation

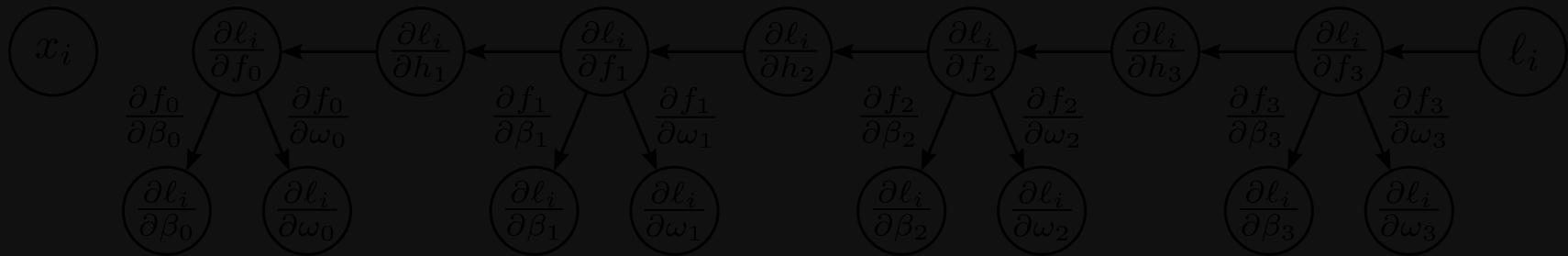
Prince (2023, chaps. 7.1–7.4)

Computing Gradients with Chain Rule

Stochastic Gradient Descent

$$\phi_{t+1} \leftarrow \phi_t - \alpha \sum_{i \in B_t} \nabla_\phi l_i(\phi_t)$$

requires computing $\frac{\partial l_i}{\partial \beta_k}$ and $\frac{\partial l_i}{\partial \Omega_k}$ at each layer $k \in \{0, \dots, K\}$.

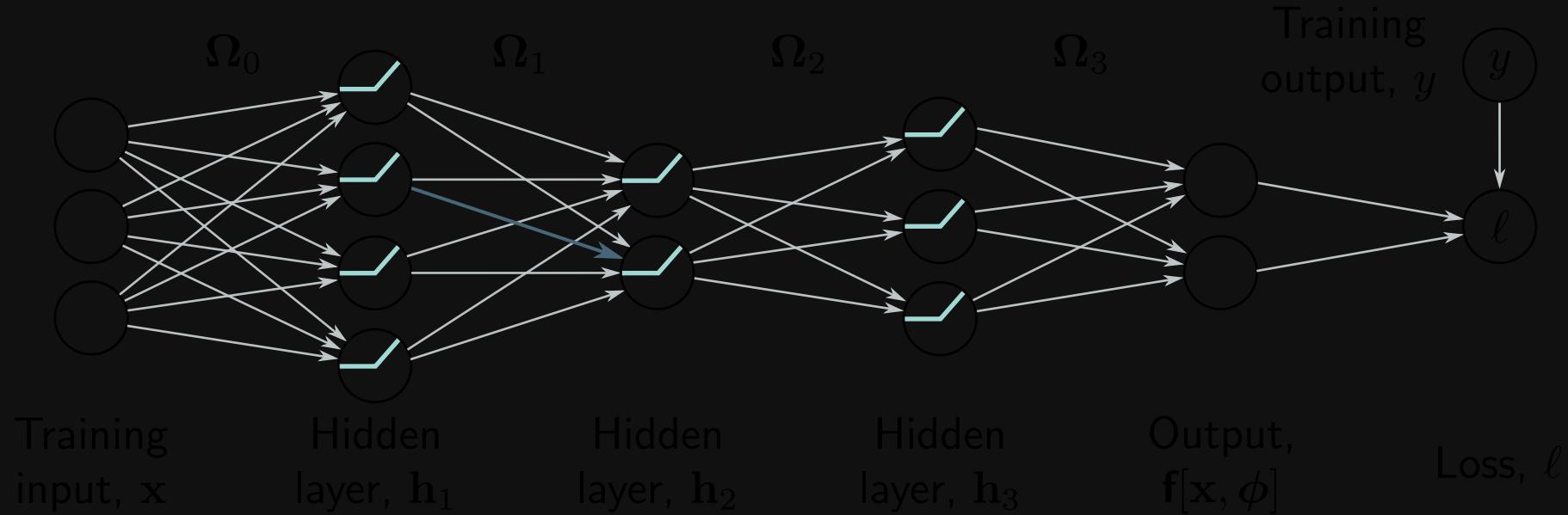


$$\frac{\partial l_i}{\partial \Omega_k} = \frac{\partial l_i}{\partial \mathbf{f}_k} \frac{\partial \mathbf{f}_k}{\partial \Omega_k} = \frac{\partial \mathbf{h}_{k+1}}{\partial \mathbf{f}_k} \frac{\partial \mathbf{f}_{k+1}}{\partial \mathbf{h}_{k+1}} \frac{\partial l_i}{\partial \mathbf{f}_{k+1}} \mathbf{h}_k^T$$

where activation \mathbf{h}_{k+1} results from pre-activation \mathbf{f}_k .

Forward Pass

Gradients of $l_i = l(\mathbf{f}[\mathbf{x}_i, \phi], \mathbf{y}_i)$ w.r.t. weights depend on activations \mathbf{h}_k , which we compute first.

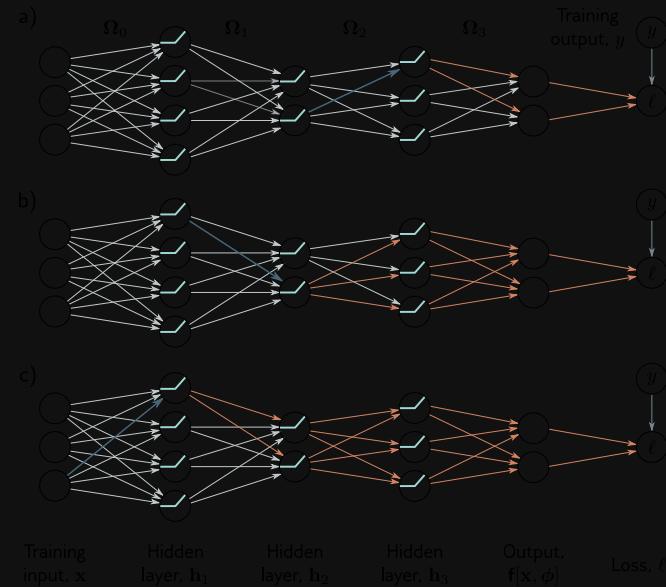


$$\mathbf{h}_k = \mathbf{a}(\mathbf{f}_{k-1})$$

$$\mathbf{f}_k = \beta_k + \boldsymbol{\Omega}_k \mathbf{h}_k$$

Backward Pass

Then compute gradients backwards using chain rule.



$$\frac{\partial l_i}{\partial \mathbf{f}_k} = \frac{\partial \mathbf{h}_{k+1}}{\partial \mathbf{f}_k} \frac{\partial \mathbf{f}_{k+1}}{\partial \mathbf{h}_{k+1}} \frac{\partial l_i}{\partial \mathbf{f}_{k+1}} = \mathbb{I}[\mathbf{f}_k > 0] \odot \boldsymbol{\Omega}_{k+1}^T \frac{\partial l_i}{\partial \mathbf{f}_{k+1}}$$

$$\frac{\partial l_i}{\partial \boldsymbol{\Omega}_k} = \frac{\partial \mathbf{f}_k}{\partial \boldsymbol{\Omega}_k} \frac{\partial l_i}{\partial \mathbf{f}_k} = \frac{\partial l_i}{\partial \mathbf{f}_k} \mathbf{h}_k^T$$

Backpropagation Algorithm

Forward pass computes activations

$$\mathbf{f}_0 = \beta_0 + \boldsymbol{\Omega}_0 \mathbf{x}_i$$

$$\mathbf{h}_k = \mathbf{a}(\mathbf{f}_{k-1}) \quad k \in \{1, \dots, K\}$$

$$\mathbf{f}_k = \beta_k + \boldsymbol{\Omega}_k \mathbf{h}_k \quad k \in \{1, \dots, K\}$$

Backward pass computes gradients

$$\frac{\partial l_i}{\partial \beta_k} = \frac{\partial l_i}{\partial \mathbf{f}_k} \quad k \in \{K, K-1, \dots, 0\}$$

$$\frac{\partial l_i}{\partial \boldsymbol{\Omega}_k} = \frac{\partial l_i}{\partial \mathbf{f}_k} \mathbf{h}_k^T \quad k \in \{K, K-1, \dots, 0\}$$

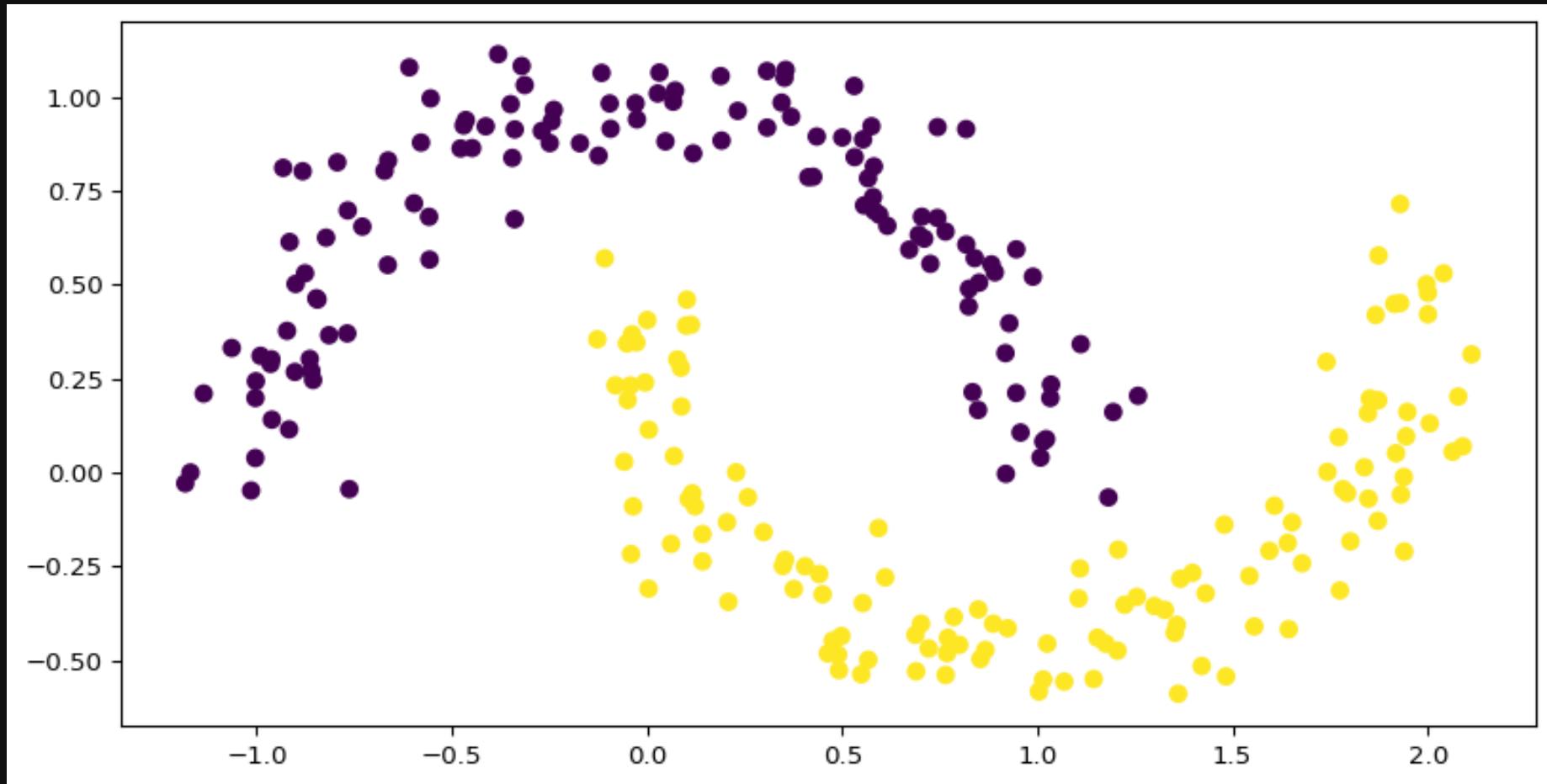
$$\frac{\partial l_i}{\partial \mathbf{f}_{k-1}} = \mathbb{I}[\mathbf{f}_{k-1} > 0] \odot \left(\boldsymbol{\Omega}_k^T \frac{\partial l_i}{\partial \mathbf{f}_k} \right) \quad k \in \{K, K-1, \dots, 1\}$$

where $l_i = l(\mathbf{f}[\mathbf{x}_i, \phi], \mathbf{y}_i)$ and $\mathbf{h}_0 = \mathbf{x}_i$ (i.e. do for each training example in batch).¹

- . This is done in parallel across a batch, so we get tensors instead of matrices.

PyTorch Implementation

Moons Dataset



Neural Network Implementation

```
1 import torch
2 import torch.nn as nn
3
4 class MyFirstNet(nn.Module):
5     def __init__(self):
6         super(MyFirstNet, self).__init__()
7         self.layers = nn.Sequential(
8             nn.Linear(2, 16),
9             nn.ReLU(),
10            nn.Linear(16, 2),
11        )
12
13    def forward(self, x):
14        return self.layers(x)
15
16    def predict(self, x):
17        output = self.forward(x)
18        return torch.argmax(output, 1)
19
20    def train(self, X, y):
21        loss_function = nn.CrossEntropyLoss()
22        optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
23
24        epochs = 15000
```

```
MyFirstNet(
(layers): Sequential(
(0): Linear(in_features=2, out_features=16, bias=True)
(1): ReLU()
(2): Linear(in_features=16, out_features=2, bias=True)
)
```

)

Sample input:

```
tensor([[ 0.8508,  0.5047],  
       [ 0.6876, -0.4327],  
       [ 1.8481, -0.0702],  
       [ 0.3535, -0.2334],  
       [ 0.9281,  0.3966]], device='mps:0')
```

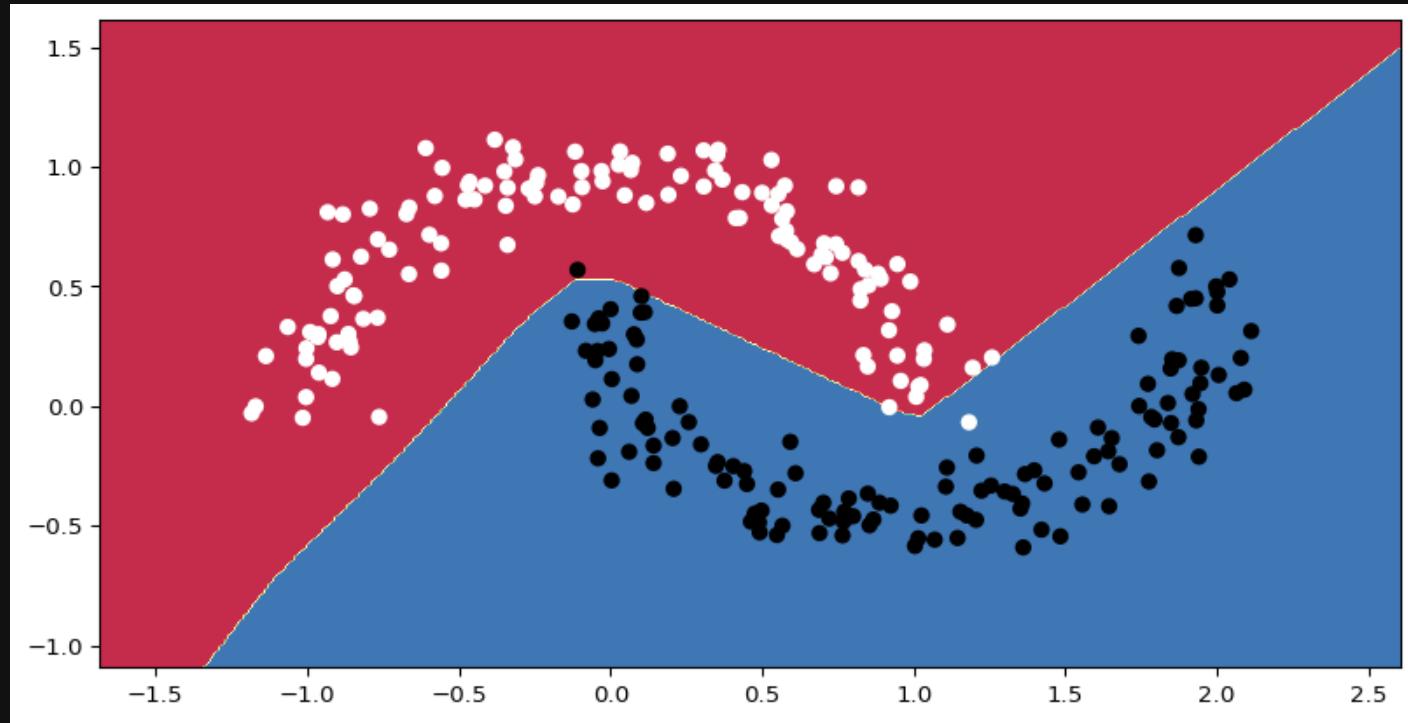
Network output:

```
tensor([[ 0.2646, -0.0653],  
       [ 0.0444,  0.1359],  
       [ 0.2497, -0.0940],  
       [ 0.0000,  0.0000],  
       [ 0.0000,  0.0000]])
```

Training

```
1 losses = model.train(X, y)
```

```
Epoch 0 loss is 0.7242650985717773
Epoch 1000 loss is 0.27691155672073364
Epoch 2000 loss is 0.24738825857639313
Epoch 3000 loss is 0.23405785858631134
Epoch 4000 loss is 0.22105270624160767
Epoch 5000 loss is 0.20637723803520203
Epoch 6000 loss is 0.18947520852088928
Epoch 7000 loss is 0.17046067118644714
Epoch 8000 loss is 0.15018796920776367
Epoch 9000 loss is 0.13044017553329468
Epoch 10000 loss is 0.11250948905944824
Epoch 11000 loss is 0.09699149429798126
Epoch 12000 loss is 0.0840226486325264
Epoch 13000 loss is 0.0734294205904007
Epoch 14000 loss is 0.06483542174100876
```



Another Example

```
1 import torch, torch.nn as nn
2 from torch.utils.data import TensorDataset, DataLoader
3 from torch.optim.lr_scheduler import StepLR
4
5 # input size, hidden layer size, output size
6 D_i, D_k, D_o = 10, 40, 5
7 # create model with two hidden layers
8 model = nn.Sequential(
9     nn.Linear(D_i, D_k),
10    nn.ReLU(),
11    nn.Linear(D_k, D_k),
12    nn.ReLU(),
13    nn.Linear(D_k, D_o))
14
15 # He initialization of weights
16 def weights_init(layer_in):
17     if isinstance(layer_in, nn.Linear):
18         nn.init.kaiming_normal_(layer_in.weight)
19         layer_in.bias.data.fill_(0.0)
20 model.apply(weights_init)
21 # choose least squares loss function
22 criterion = nn.MSELoss()
23 # construct SGD optimizer and initialize learning rate and momentum
24 optimizer = torch.optim.SGD(model.parameters(), lr = 0.1, momentum=0.9)
```

Epoch 99, loss 1.964

References

Prince, Simon J. D. 2023. *Understanding Deep Learning*. Cambridge, Massachusetts: The MIT Press.