

Perceptrons (single layer)

Perceptron is an algorithm that generates a linear classifier for supervised learning, It, ideally, finds the optimal linear classifier for the data:

$$\text{data} = \left\{ \left(x^{(1)}, y^{(1)} \right), \dots, \left(x^{(n)}, y^{(n)} \right) \right\}$$

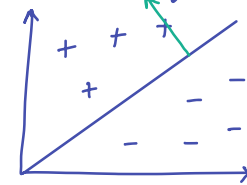
and $\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$
 $\underbrace{\quad}_{\text{weight}} \quad \underbrace{\quad}_{\text{offset/bias}}$

Linear classifier: $h(x; \theta, \theta_0)$

$$h(x; \theta, \theta_0) = \text{sign}(\theta^T x + \theta_0) = \begin{cases} +1 & \text{if } \theta^T x + \theta_0 > 0 \\ -1 & \text{otherwise} \end{cases}$$

return whether x is positive or negative

normal to the line at x



defines the line - separate positive + negative points



It may *not* find the optimal classifier if:

- the data is not linearly classifiable, eg ^^^
- there exists multiple linear classifiers (as soon as one is found, the weights stop adjusting)

The algorithm aims to find the hyperplane that separates the data by iteratively adjusting the weights of the hyperplane

The algorithm:

$\text{perceptron}(D, T)$

($D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ - training data)

$$\theta = \vec{0}$$

$$\theta_0 = 0$$

for $t = 1$ to T

for $i = 1$ to n

number of examples in data set

$$\text{if } y^{(i)} \cdot (\theta^T x^{(i)} + \theta_0) \leq 0$$

$$\theta = \theta + y^{(i)} x^{(i)}$$

$$\theta_0 = \theta_0 + y^{(i)}$$

(Loss function - only updates if values are wrongly classified)

if this is positive, the prediction is correct → ignore

if it is negative adjust!

return θ, θ_0

Using numpy in Python:

```
import numpy as np
import math

##### useful functions #####
def positive(x, th, th0):
    return np.sign(np.dot(np.transpose(th), x) + th0)

def score(data, labels, ths, th0s):
    pos = np.sign(np.dot(np.transpose(ths), data) + np.transpose(th0s))
    return np.sum(pos == labels, axis = 1, keepdims = True)

##### actual perceptron #####

def perceptron(data, labels, params = {}):
    # if T not in params, default to 100
    T = params.get('T', 100)
    (d, n) = data.shape

    theta = np.zeros((d, 1)); theta_0 = np.zeros((1, 1))
    for t in range(T):
        for i in range(n):
            x = data[:,i:i+1]
            y = labels[:,i:i+1]
            if y * positive(x, theta, theta_0) <= 0.0:
                theta = theta + y * x
                theta_0 = theta_0 + y

    return theta, theta_0

# Regular perceptron can be somewhat sensitive to the most recent examples
# that it sees. Instead, averaged perceptron produces a more stable output by
# outputting the average value of th and th0 across all iterations

def averaged_perceptron(data, labels, params = {}):
    # if T not in params, default to 100
    T = params.get('T', 100)
    (d, n) = data.shape

    theta = np.zeros((d, 1)); theta_0 = np.zeros((1, 1))
    thetas = np.zeros((d, 1)); theta_0s = np.zeros((1, 1))
    for t in range(T):
        for i in range(n):
            x = data[:,i:i+1]
            y = labels[:,i:i+1]
            if y * positive(x, theta, theta_0) <= 0.0:
                theta = theta + y * x
                theta_0 = theta_0 + y
            thetas = thetas + theta
            theta_0s = theta_0s + theta_0

    return thetas/(n*T), theta_0s/(n*T)
```

```
##### useful functions #####
def positive(x, th, th0):
    return np.sign(np.dot(np.transpose(th), x) + th0)

def score(data, labels, ths, th0s):
    pos = np.sign(np.dot(np.transpose(ths), data) + np.transpose(th0s))
    return np.sum(pos == labels, axis = 1, keepdims = True)
```

(nicer
colours)

```
##### actual perceptron #####
def perceptron(data, labels, params = {}):
    # if T not in params, default to 100
    T = params.get('T', 100)
    (d, n) = data.shape

    theta = np.zeros((d, 1)); theta_0 = np.zeros((1, 1))
    for t in range(T):
        for i in range(n):
            x = data[:,i:i+1]
            y = labels[:,i:i+1]
            if y * positive(x, theta, theta_0) <= 0.0:
                theta = theta + y * x
                theta_0 = theta_0 + y
    return theta, theta_0

# Regular perceptron can be somewhat sensitive to the most recent examples that it sees. Instead, averaged perceptron produces
# a more stable output by outputting the average value of th and th0 across all iterations
def averaged_perceptron(data, labels, params = {}):
    # if T not in params, default to 100
    T = params.get('T', 100)
    (d, n) = data.shape

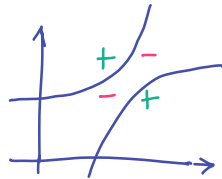
    theta = np.zeros((d, 1)); theta_0 = np.zeros((1, 1))
    thetas = np.zeros((d, 1)); theta_0s = np.zeros((1, 1))
    for t in range(T):
        for i in range(n):
            x = data[:,i:i+1]
            y = labels[:,i:i+1]
            if y * positive(x, theta, theta_0) <= 0.0:
                theta = theta + y * x
                theta_0 = theta_0 + y
                thetas = thetas + theta
                theta_0s = theta_0s + theta_0
    return thetas/(n*T), theta_0s/(n*T)
```

Idea: if the data is not linearly classifiable, we can try to change basis and check if it is now classifiable.

eg:



eg: $\gamma = [x, x^2]$



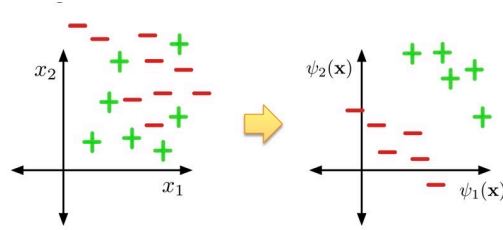
Polynomial basis:

Here is a table illustrating the k th order polynomial basis for different values of k .

Order	$d = 1$	in general
0	$[1]$	$[1]$
1	$[1, x]$	$[1, x_1, \dots, x_d]$
2	$[1, x, x^2]$	$[1, x_1, \dots, x_d, x_1^2, x_1 x_2, \dots]$
3	$[1, x, x^2, x^3]$	$[1, x_1, \dots, x_d^2, x_1 x_2, \dots, x_1 x_2 x_3, \dots]$

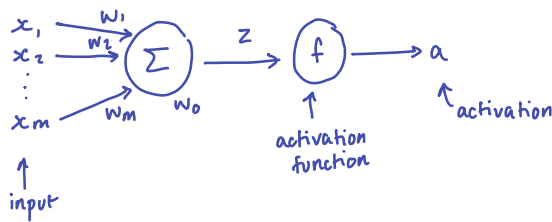
Multi-layer Perceptron

Idea: take the single-layer perceptron and add some 'hidden' layers that will make the data linearly classifiable



MLP from scratch:

Single neuron:

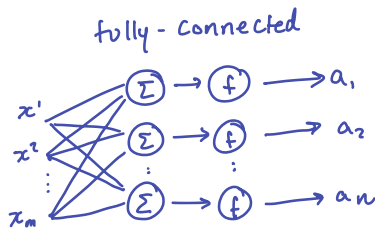


$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \in \mathbb{R}^m$$

$$a \in \mathbb{R}$$

$$a = f(z) = f\left(\left(\sum_{j=1}^m x_j w_j\right) + w_0\right)$$

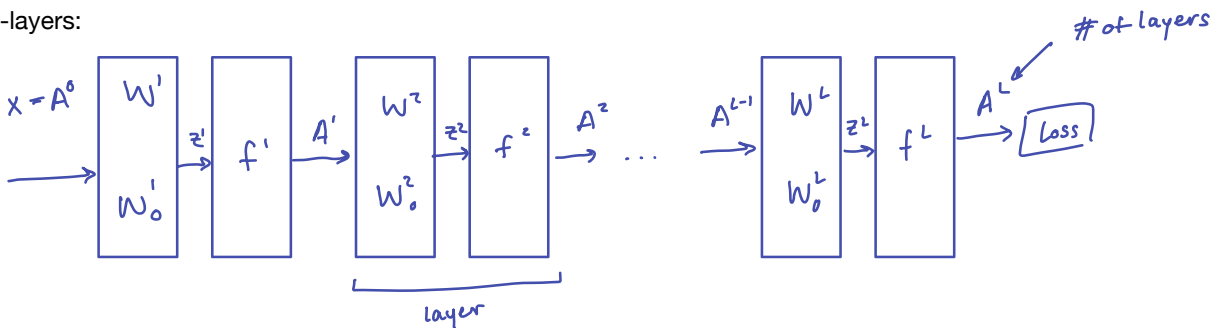
Single layer:



$W : m \times n$ matrix of vectors
 $W_0 : n \times 1$ vector
 $A = f(z) = f(W^T X + W_0)$

(f is applied element-wise)

L-layers:



$$A^L = W^L{}^T A^{L-1}$$

$$W^{L-1}{}^T A^{L-2}$$

$$\vdots$$

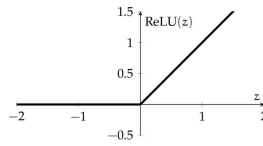
$$W^1{}^T X$$

$$A^L = W^{LT} W^{L-1T} \dots W^1 X$$

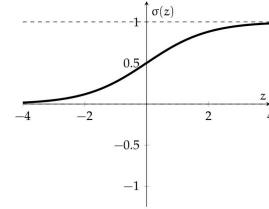
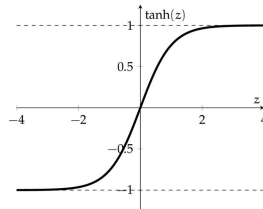
$$A^L = W^{total} X$$

Note: the layers are only useful if f is non-linear
(if f is linear, W_{total} is simply a linear transformation)

Possible f functions:



"Rule of thumb:
Start with ReLU activation. If
necessary, try tanh."
(University of Toronto)

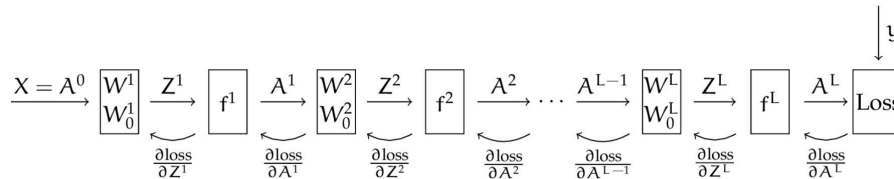


The idea now is to:

- forward pass, calculating the final output based on the original inputs and all the weights+activations
- compare to the actual output and calculate loss
- back-propagate through the layers, adjusting weights according to how much they are to "blame" for the loss

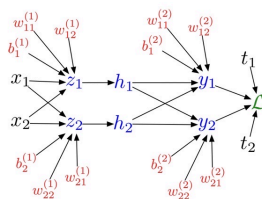
Error back-propagation:

Work backward and compute the gradient of the loss with respect to the weights in each layer



Another representation (notation is slightly different)

Backpropagation for a MLP (Vectorized)



Forward pass:

$$z = W^{(1)}x + b^{(1)}$$

$$h = \sigma(z)$$

$$y = W^{(2)}h + b^{(2)}$$

$$\mathcal{L} = \frac{1}{2} \|y - t\|^2$$

Backward pass:

$$\bar{\mathcal{L}} = 1$$

$$\bar{y} = \bar{\mathcal{L}}(y - t)$$

$$\bar{W}^{(2)} = \bar{y}h^T$$

$$\bar{b}^{(2)} = \bar{y}$$

$$\bar{h} = W^{(2)T} \bar{y}$$

$$\bar{z} = \bar{h} \circ \sigma'(z)$$

$$\bar{W}^{(1)} = \bar{z}x^T$$

$$\bar{b}^{(1)} = \bar{z}$$

Full algorithm:

SGD-NEURAL-NET($\mathcal{D}_n, T, L, (m^1, \dots, m^L), (f^1, \dots, f^L)$)

```

1  for l = 1 to L
2     $W_{ij}^l \sim \text{Gaussian}(0, 1/m^l)$ 
3     $W_{0j}^l \sim \text{Gaussian}(0, 1)$ 
4  for t = 1 to T
5    i = random sample from  $\{1, \dots, n\}$ 
6     $A^0 = x^{(i)}$ 
7    // forward pass to compute the output  $A^L$ 
8    for l = 1 to L
9       $Z^l = W^{lT} A^{l-1} + W_0^l$ 
10      $A^l = f^l(Z^l)$ 
11     loss = Loss( $A^L, y^{(i)}$ )
12   for l = L to 1:
13     // error back-propagation
14      $\partial \text{loss} / \partial A^l = \text{if } l < L \text{ then } \partial \text{loss} / \partial Z^{l+1} \cdot \partial Z^{l+1} / \partial A^l \text{ else } \partial \text{loss} / \partial A^L$ 
15      $\partial \text{loss} / \partial Z^l = \partial \text{loss} / \partial A^l \cdot \partial A^l / \partial Z^l$ 
16     // compute gradient with respect to weights
17      $\partial \text{loss} / \partial W^l = \partial \text{loss} / \partial Z^l \cdot \partial Z^l / \partial W^l$ 
18      $\partial \text{loss} / \partial W_0^l = \partial \text{loss} / \partial Z^l \cdot \partial Z^l / \partial W_0^l$ 
19     // stochastic gradient descent update
20      $W^l = W^l - \eta(t) \cdot \partial \text{loss} / \partial W^l$ 
21      $W_0^l = W_0^l - \eta(t) \cdot \partial \text{loss} / \partial W_0^l$ 

```

Note on initialisation:

- Weights should be chosen randomly -> if not, the symmetry between layers makes them less useful.
- Weights should not be too big -> the activation functions are only useful near zero, so gradient descent will not signal a useful direction to go if weights are too big.

Possible strategy

Choose each weight at random from a Gaussian (normal) distribution with mean 0 and standard deviation $(1/m)$ where m is the number of inputs to the unit.

Note on last-layer activation:

The activation function for the last layer of the network may be different from the one used in previous layers in order to make the data match the desired output type, eg:

last-layer function		other layers	purpose
$f(x) = x$	w/	squared loss	→ regression
sigmoid (interpret as probability)	w/	neg. log likelihood	→ multi-class classification

Using PyTorch in Python:

```
from collections import OrderedDict
import torch
import torch.nn as nn
import torch
import torchvision
import torchvision.transforms as transforms
import ssl
ssl._create_default_https_context = ssl._create_unverified_context

# Load the MNIST dataset
transform = transforms.Compose(
    [transforms.ToTensor(),
     transforms.Lambda(lambda x: x.view(-1))]) # Flatten images

train_dataset = torchvision.datasets.MNIST(root='./data', train=True,
transform=transform, download=True)
test_dataset = torchvision.datasets.MNIST(root='./data', train=False,
transform=transform, download=True)

train_loader = torch.utils.data.DataLoader(train_dataset, batch_size=64,
shuffle=True)
test_loader = torch.utils.data.DataLoader(test_dataset, batch_size=64,
shuffle=False)
```

```

model = nn.Sequential(OrderedDict([
    ('layer1', nn.Linear(784, 100)),
    ('activation1', nn.ReLU()),
    ('layer2', nn.Linear(100, 50)),
    ('activation2', nn.ReLU()),
    ('output', nn.Linear(50, 10)),
    ('outActivation', nn.Sigmoid()),
]))

loss_fn = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=0.001)

X, y = next(iter(train_loader))

# Training loop
num_epochs = 3
for epoch in range(num_epochs):
    model.train()
    total_loss = 0
    correct = 0
    total = 0

    for X, y in train_loader:
        y_pred = model(X)
        loss = loss_fn(y_pred, y)

        optimizer.zero_grad()
        loss.backward()      # apply backpropagation
        optimizer.step()

        total_loss += loss.item()

    # Compute accuracy
    _, predicted = torch.max(y_pred, 1) # Get the class with the highest score
    correct += (predicted == y).sum().item()
    total += y.size(0)

    train_accuracy = correct / total * 100
    avg_loss = total_loss / len(train_loader)
    print(f"Epoch [{epoch+1}/{num_epochs}], Loss: {avg_loss:.4f}, Accuracy: {train_accuracy:.2f}%")

# Evaluate on test set
model.eval()
correct = 0
total = 0
with torch.no_grad():
    for X, y in test_loader:
        y_pred = model(X)
        _, predicted = torch.max(y_pred, 1)
        correct += (predicted == y).sum().item()
        total += y.size(0)

test_accuracy = correct / total * 100
print(f"Final Test Accuracy: {test_accuracy:.2f}%")

```

```

10 # Load the MNIST dataset
11 transform = transforms.Compose([transforms.ToTensor(), transforms.Lambda(lambda x: x.view(-1))]) # Flatten images
12
13 train_dataset = torchvision.datasets.MNIST(root='./data', train=True, transform=transform, download=True)
14 test_dataset = torchvision.datasets.MNIST(root='./data', train=False, transform=transform, download=True)
15
16 train_loader = torch.utils.data.DataLoader(train_dataset, batch_size=64, shuffle=True)
17 test_loader = torch.utils.data.DataLoader(test_dataset, batch_size=64, shuffle=False)
18
19 model = nn.Sequential(OrderedDict([
20     ('layer1', nn.Linear(784, 100)),
21     ('activation1', nn.ReLU()),
22     ('layer2', nn.Linear(100, 50)),
23     ('activation2', nn.ReLU()),
24     ('output', nn.Linear(50, 10)),
25     ('outActivation', nn.Sigmoid()),
26 ]))
27
28 loss_fn = nn.CrossEntropyLoss()
29 optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
30
31 X, y = next(iter(train_loader))
32
33 # Training loop
34 num_epochs = 3
35 for epoch in range(num_epochs):
36     model.train()
37     total_loss = 0
38     correct = 0
39     total = 0
40
41     for X, y in train_loader:
42         y_pred = model(X)
43         loss = loss_fn(y_pred, y)
44
45         optimizer.zero_grad()
46         loss.backward() # apply backpropagation
47         optimizer.step()
48
49         total_loss += loss.item()
50
51         # Compute accuracy
52         _, predicted = torch.max(y_pred, 1) # Get the class with the highest score
53         correct += (predicted == y).sum().item()
54         total += y.size(0)
55
56     train_accuracy = correct / total * 100
57     avg_loss = total_loss / len(train_loader)
58     print(f"Epoch [{epoch+1}/{num_epochs}], Loss: {avg_loss:.4f}, Accuracy: {train_accuracy:.2f}%")
59
60 # Evaluate on test set
61 model.eval()
62 correct = 0
63 total = 0
64 with torch.no_grad():
65     for X, y in test_loader:
66         y_pred = model(X)
67         _, predicted = torch.max(y_pred, 1)
68         correct += (predicted == y).sum().item()
69         total += y.size(0)
70
71 test_accuracy = correct / total * 100
72 print(f"Final Test Accuracy: {test_accuracy:.2f}%")
73
74 torch.save(model.state_dict(), "my_model.pickle")

```