

Chapter 1

Fundamentals of Neural Networks

Neural networks and deep learning

Artificial Neuron

An **artificial neuron** is the basic unit of neural networks.

Basic elements of an artificial neuron:

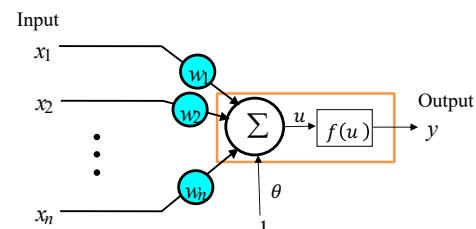
- A set of **input** signals: the input is a vector $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)^T$ where n is the number (or the dimension) of input signals. Inputs are also referred to as **features**.
- Inputs are connected to the neuron via synaptic connections whose strengths are represented by their **weights**.
- The weight vector $\mathbf{w} = (w_1 \ w_2 \ \cdots \ w_n)^T$ where w_i is the synaptic weight connecting i th input of the neuron.

© Nanyang Technological University

1

3

Artificial Neuron



Input vector $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)^T$
weight vector $\mathbf{w} = (w_1 \ w_2 \ \cdots \ w_n)^T$
 n is the number of inputs.

2

Notation

Vectors are denoted in **bold** and written as horizontally with a transpose (T).

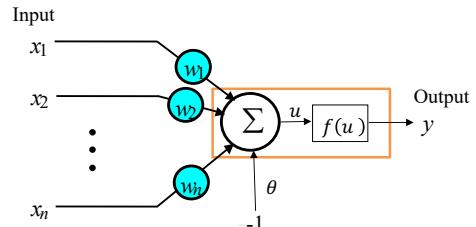
$$\mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} = (w_1 \ w_2 \ \cdots \ w_n)^T$$

Example

$$\mathbf{x} = \begin{pmatrix} 1.2 \\ -0.5 \\ 3.5 \\ 2.1 \end{pmatrix} = (1.2 \ -0.5 \ 3.5 \ 2.1)^T$$
$$\mathbf{x}^T = (1.2 \ -0.5 \ 3.5 \ 2.1)$$

4

Artificial Neuron



The total **synaptic input** u to the neuron is given by the sum of the products of the inputs and their corresponding connecting weights minus the **threshold** of the neuron.

The total synaptic input to a neuron, u is given by

$$u = w_1x_1 + w_2x_2 + \dots + w_nx_n - \theta = \sum_{i=1}^n w_i x_i - \theta$$

where θ is the threshold of the neuron.

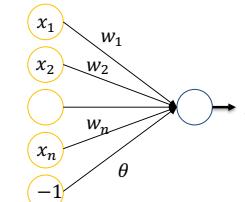
By using vector notations:

$$u = \mathbf{w}^T \mathbf{x} - \theta$$

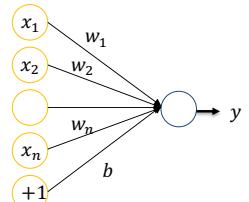
5

Bias vs Threshold

The threshold is often considered as a weight with a fixed input of -1 . Often the threshold is represented as a **bias** b that receives constant $+1$ input.



Threshold
 $u = \mathbf{w}^T \mathbf{x} - \theta$
 $y = f(u)$



Bias
 $u = \mathbf{w}^T \mathbf{x} + b$
 $y = f(u)$

Bias $b = -\theta$

7

Artificial Neuron

The **activation function** f relates synaptic input to the activation of the neuron.

$f(u)$ denotes the **activation** of the neuron.

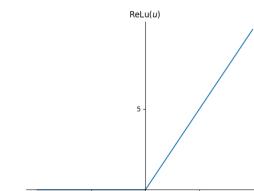
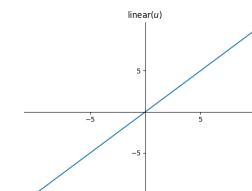
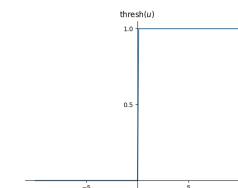
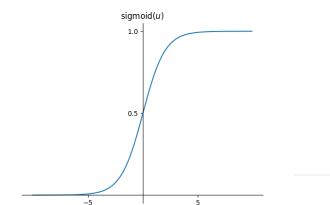
For some neurons, the **output** y is equal to the activation of the neuron.

$$y = f(u)$$

Note that activation is not generally equal to the output of the neuron.

6

Activation functions



8

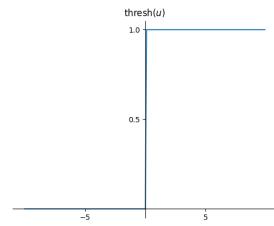
Activation functions

For **threshold (unit step) activation function**, the activation is given by

$$f(u) = \text{threshold}(u) = 1(u > 0)$$

where $1(\cdot)$ is the *Indictor function* or *Unit-step function*:

$$1(x) = \begin{cases} 1, & x \text{ is True} \\ 0, & x \text{ is False} \end{cases}$$

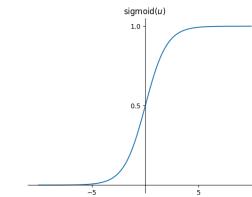


9

Sigmoid activation function

The sigmoidal is known as the **logistic function** or simply **sigmoid function**

$$f(u) = \text{sigmoid}(u) = \frac{1}{1 + e^{-u}}$$



In general, the **sigmoid** activation function can be written as

$$f(u) = \frac{a}{1 + e^{-bu}}$$

a is the gain (amplitude) and b is the slope.

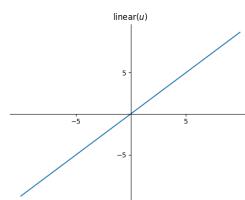
But often, $a = 1.0$ and $b = 1.0$.

11

Activation functions

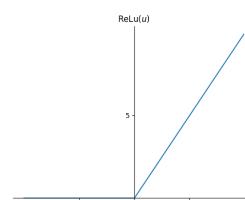
A neuron with **linear** activation function can be written as

$$f(u) = \text{linear}(u) = u$$



The **ReLU (rectified-linear unit)** activation function can be written as

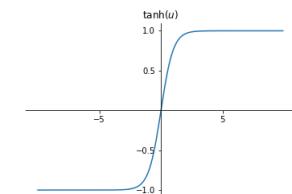
$$f(u) = \text{relu}(u) = \max\{0, u\}$$



10

Tanh activation function

$$f(u) = \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$



Tanh activation function has the same shape as sigmoidal and spans from -1 and +1. It is also known as **bipolar sigmoidal**.

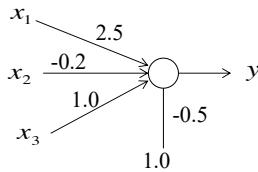
Sigmoidal is the most pervasive and biologically plausible activation function. Since sigmoid function is *differentiable*, it leads to mathematically attractive neuronal models.

12

Example 1

The artificial neuron in the figure receives 3-dimensional inputs $\mathbf{x} = (x_1 \ x_2 \ x_3)^T$ and has an activation function given by $f(u) = \frac{0.8}{1+e^{-1.2u}}$.

Find the synaptic input and the output of the neuron for inputs: $\begin{pmatrix} 0.8 \\ 2.0 \\ -0.5 \end{pmatrix}$ and $\begin{pmatrix} -0.4 \\ 1.5 \\ 1.0 \end{pmatrix}$.



PyTorch 2.0

PyTorch/Tensorflow is about processing of **tensors**. Tensor is a multidimensional array.

Rank refers to the number of dimensions and **shape** gives the sizes of each dimension of the tensor.

3. # a rank 0 tensor; a scalar with shape []:

[1., 2., 3.] # a rank 1 tensor; a vector with shape [3]

[[1., 2., 3.], [4., 5., 6.]] # a rank 2 tensor; a matrix with shape [2, 3]

[[[1., 2., 3.]], [[7., 8., 9.]]] # a rank 3 tensor with shape [2, 1, 3]

13

15

Example 1

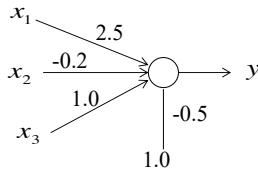
$$\mathbf{w} = \begin{pmatrix} 2.5 \\ -0.2 \\ 1.0 \end{pmatrix}, \quad b = -0.5$$

$$\text{Consider } \mathbf{x} = \begin{pmatrix} 0.8 \\ 2.0 \\ -0.5 \end{pmatrix}$$

$$\text{Synaptic input } u = \mathbf{w}^T \mathbf{x} + b = (2.5 \ -0.2 \ 1.0) \begin{pmatrix} 0.8 \\ 2.0 \\ -0.5 \end{pmatrix} - 0.5 = 0.6$$

$$\text{Output } y = f(u) = \frac{0.8}{1+e^{-1.2u}} = \frac{0.8}{1+e^{-1.2 \times 0.6}} = 0.538$$

Similarly, for $\mathbf{x} = \begin{pmatrix} -0.4 \\ 1.5 \\ 1.0 \end{pmatrix}$, $u = -0.8$ and output $y = f(u) = 0.222$



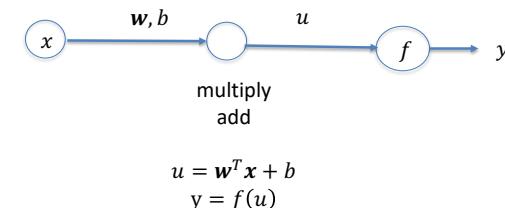
Computational Graph

PyTorch program involves building and evaluation of a **computational graph**.

A computational graph is a series of tensor **operations** arranged into a graph.

- Nodes of the graph represent tensor operations
- Edges represent values (tensors) that follow through the graph

Computational graph of a neuron



14

16

Torch Implementation of Example 1

```
import torch

# a class for neuron
class Neuron():
    # initiate a neuron class with weights and biases (initiate the object)
    def __init__(self):
        self.w = torch.tensor([2.5, -0.2, 1.0])
        self.b = torch.tensor(-0.5)

    # evaluate the neuron (implement a function)
    def __call__(self, x):
        u = torch.inner(self.w, x) + self.b
        y = 0.8/(1+torch.exp(-1.2*u))
        return u, y

# create a neuron
neuron = Neuron()

# evaluate
u, y = neuron(torch.tensor([0.8, 2.0, -0.5]))

# print: u = 0.600, y = 0.538
```

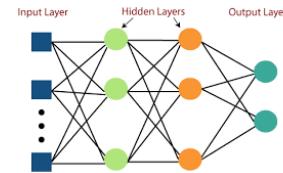
17

NN architectures: multilayer feedforward networks

- Comprised of more than one layer of neurons. Layers between input source nodes and output layer is referred to as *hidden layers*.

- Multilayer neural networks can handle *more complicated* and *larger scale problems* better than single-layer networks.

- However, training multilayer network may be *more difficult* and *time-consuming*.



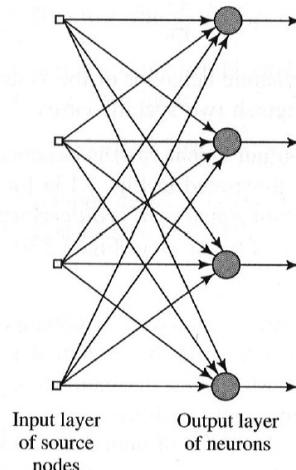
Three-layer network

19

NN Architectures: neuron layers

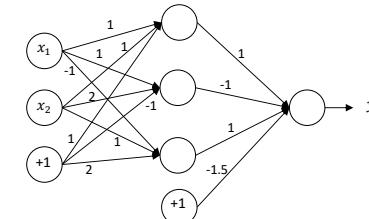
Single – layer of neurons

- Comprised of an input layer of source units that inject into an output layer of neurons.
- A fully-connected layer



18

Example 2

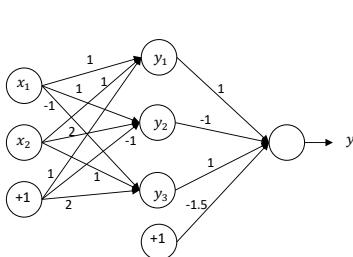


Two-layer neural network receives 2-dimensional inputs $(x_1, x_2) \in \mathbb{R}^2$ and has one output neuron and three hidden neurons. All the neurons have unit step activation functions. The weights of the connections are given in the figure. Find the space of inputs for which the output $y = 1.0$.

Find the output for inputs $(0.0, 0.0)$, $(2.0, 2.0)$, and $(-1.0, 1.0)$

20

Example 2



Synaptic inputs:

$$u_1 = x_1 + x_2 + 1$$

Output $y_1 = 1(u_1 > 0)$

$$u_2 = x_1 + 2x_2 - 1$$

$$y_2 = 1(u_2 > 0)$$

$$u_3 = -x_1 + x_2 + 2$$

$$y_3 = 1(u_3 > 0)$$

$$u = y_1 - y_2 + y_3 - 1.5$$

$$y = 1(u > 0)$$

Output layer neuron:

$$u = y_1 - y_2 + y_3 - 1.5$$

$$y = 1(u > 0)$$

Note that $y_1, y_2, y_3 \in \{0, 1\}$

y_1	y_2	y_3	u	y
0	0	0	-1.5	0
0	0	1	-0.5	0
0	1	0	-2.5	0
0	1	1	-1.5	0
1	0	0	-0.5	0
1	0	1	0.5	1
1	1	0	-1.5	0
1	1	1	-0.5	0

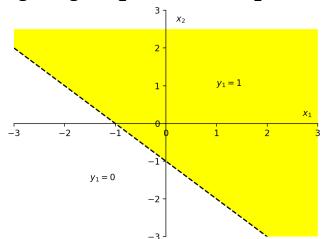
21

23

$$Y = Y_1 \bar{Y}_2 Y_3$$

$$y_1 = f(u_1) = 1(u_1 > 0)$$

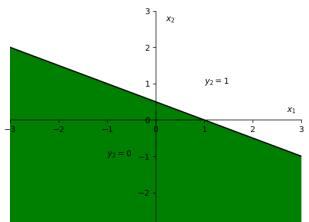
$$\text{Boundary: } u_1 = x_1 + x_2 + 1 = 0 \rightarrow x_2 = -x_1 - 1$$



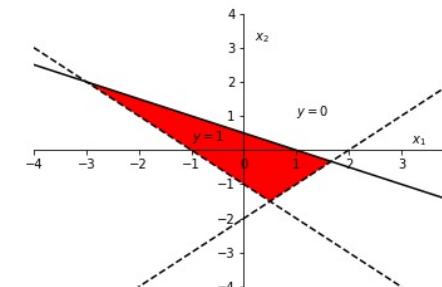
The boundary line is obtained by setting $u_1 = 0$; and for one side of the boundary, $y = 1$ and on other side $y_1 = 0$.

$$u_3 = x_2 - x_1 + 2 = 0 \rightarrow x_2 = x_1 - 2$$

$$u_2 = 2x_2 + x_1 - 1 = 0 \rightarrow x_2 = -0.5x_1 + 0.5$$



$y = 1$ region is given by the intersection of regions: $y_1=1$, $y_2=0$, and $y_3=1$.



$$x = (0,0,0) \rightarrow y = 1$$

$$x = (2,0,2) \rightarrow y = 0$$

$$x = (-1,0,1) \rightarrow y = 1$$

Note that networks of *discrete perceptrons* (neurons with threshold activation functions) can implement Boolean functions.

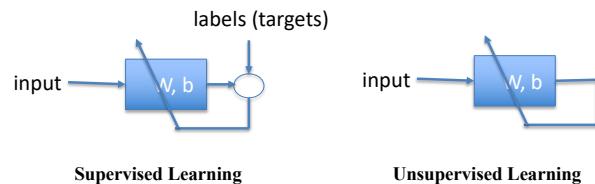
22

24

Training (or learning) of neural networks

Neural networks attain their operating characteristics through **learning** (or **training**) with training examples. During training, the weights or the strengths of connections are gradually adjusted iteratively to achieve their desirable labels (or **targets**).

Training may be either **supervised** or **unsupervised**.



Supervised learning

Learning of a neuron or neural network is usually performed in order to minimize a **cost function** (**loss function** or **error function**).

The cost function $J(\mathbf{W}, \mathbf{b})$ of an artificial neuron is typically a multi-dimensional function that depends on weights \mathbf{W} and the biases \mathbf{b} . The neuron learning attempts to find the optimal weights \mathbf{W}^* and biases \mathbf{b}^* that minimize the error function:

$$\mathbf{W}^*, \mathbf{b}^* = \arg \min_{\mathbf{W}, \mathbf{b}} J(\mathbf{W}, \mathbf{b})$$

Given a set of training patterns, the parameters (weights and biases) minimizing cost function are learned in an iterative procedure. In each iteration, small changes of weights $\Delta \mathbf{W}$ and biases $\Delta \mathbf{b}$ are made:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} + \Delta \mathbf{W} \\ \mathbf{b} &\leftarrow \mathbf{b} + \Delta \mathbf{b}\end{aligned}$$

Supervised and unsupervised learning

Supervised Learning:

For each training input pattern, the network is presented with the correct **target label** (the desired output).

Unsupervised Learning:

For each training input pattern, the network adjusts weights *without knowing* the correct target.

In unsupervised training, the network **self-organizes** to classify similar input patterns into clusters.

Gradient descent learning

The gradient descent procedure states that the weights \mathbf{W} (and biases \mathbf{b}) are updated during learning by searching in the direction of and proportional to the **negative gradient** of the cost function.

That is, the change of the weight vector:

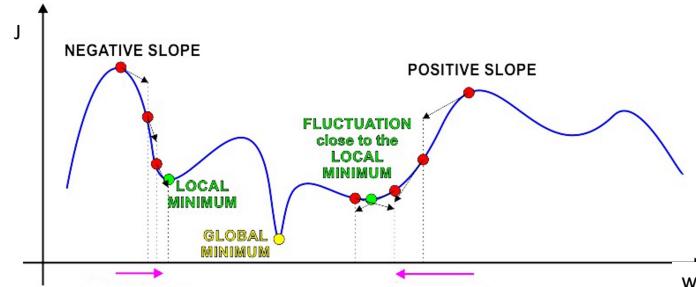
$$\begin{aligned}\Delta \mathbf{W} &\propto -\frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial \mathbf{w}} \\ \Delta \mathbf{W} &= -\alpha \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}}\end{aligned}$$

where $\frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}}$ is the gradient (partial derivative) of cost with respect to weight \mathbf{W} and α is **learning factor** or **learning rate**. $\alpha \in (0.0, 1.0]$.

The **gradient descent equations** for learning the weights is given by substituting above in

$$\mathbf{W} \leftarrow \mathbf{W} + \Delta \mathbf{W}$$

Gradient Descent Learning



29

Gradient descent learning

Given a set of training examples

Initialize weight \mathbf{W} and bias \mathbf{b}

Set the learning parameter α

Iterate until convergence:

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{b}} J$$

Convergence is achieved by observing one of the following:

1. No changes in weights and biases
2. No difference between the outputs and targets
3. No decrease in the cost function J

31

Gradient descent learning

The **gradient descent equations** for the weights is given by

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}}$$

Similarly, for the bias

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial \mathbf{b}}$$

Notation: $\nabla_{\mathbf{W}} J = \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial \mathbf{W}}$ and $\nabla_{\mathbf{b}} J = \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial \mathbf{b}}$.

Gradient descent learning is given by

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{b}} J\end{aligned}$$

Note: Will drop the arguments in J .

30

Training dataset

Training data are also referred to as **training examples** or **training patterns**. For supervised learning, a training pattern consists of a pair consisting of input pattern and the corresponding target.

A training dataset is a set of training examples: $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$ or $\{(\mathbf{x}_1, d_1), (\mathbf{x}_2, d_2), \dots, (\mathbf{x}_P, d_P)\}$

\mathbf{x}_p is the input (features) and d_p is the target (desired label) of p th training pattern. P is the number of examples in the dataset.

The input is usually n -dimensional and written as:

$$\mathbf{x}_p = (x_{p1}, x_{p2}, \dots, x_{pn})^T$$

32

Stochastic Gradient Descent (SGD) learning

Given training examples $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$

Set learning factor α

Initialize (\mathbf{W}, \mathbf{b})

Iterate until convergence:

for each pattern (\mathbf{x}_p, d_p) :

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J_p$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_b J_p$$

- In each **epoch** or cycle of iteration, the learning takes place individually over every pattern
- The cost J_p is individually computed from the output and the target of the p th training pattern.

(Batch) Gradient descent learning

Given a set of training patterns: (\mathbf{X}, \mathbf{d})

Set learning factor α

Initialize (\mathbf{W}, \mathbf{b})

Iterate until convergence:

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_b J$$

- The cost J is computed using all the training patterns. That is, using (\mathbf{X}, \mathbf{d})
- In each epoch, the weights are updated once considering all the input patterns.

Batches of Data

Inputs are often presented as a batch in a **data matrix X** and a **target vector d** . The input datapoints (or patterns) are written as rows in the data matrix and the targets are written into a single vector in the target vector.

Given a training dataset $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$.

Data matrix:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{P1} & x_{P2} & \cdots & x_{Pn} \end{pmatrix}$$

Target vector:

$$\mathbf{d} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_P \end{pmatrix}$$

Summary

- Analogy between biological and artificial neurons
- Transfer function of artificial neuron:

$$u = \mathbf{w}^T \mathbf{x} + b$$

$$y = f(u)$$
- Types of activation functions: sigmoid, threshold, linear, ReLU, and tanh.
- Given inputs, to find the outputs for simple feedforward networks
- Supervised and unsupervised learning
- Gradient descent learning:

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_b J$$
- Stochastic gradient descent (SGD) and batch gradient descent (GD)

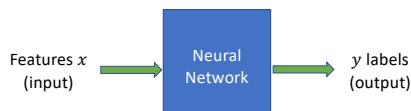
Chapter 2

Regression and classification

Neural networks and deep learning

© Nanyang Technological University

Regression and classification

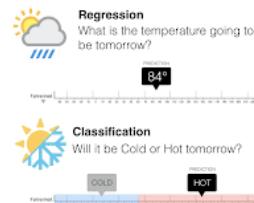


Primarily, neural network are used to *predict* output **labels** from input **features**.

Prediction tasks can be classified into two categories:

- Regression:** the labels are continuous (age, income, height, etc.)
- Classification:** the labels are discrete (sex, digits, type of flowers, etc.)

Training finds network weights and biases that are optimal for **prediction** of labels from features.



Linear neuron

Synaptic input u to a neuron is given by

$$u = \mathbf{w}^T \mathbf{x} + b$$

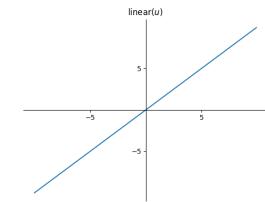
A linear neuron has a *linear activation function*. That is,

$$y = f(u) = u$$

A linear neuron with weights $\mathbf{w} = (w_1 \ w_2 \ \dots \ w_n)^T$ and bias b has an output:

$$y = \mathbf{w}^T \mathbf{x} + b$$

where input $\mathbf{x} = (x_1 \ x_2 \ \dots \ x_n)^T \in \mathbb{R}^n$ and output $y \in \mathbb{R}$.



Linear neuron performs linear regression

Representing a dependent (output) variable as a linear combination of independent (input) variables is known as **linear regression**.

The output of a linear neuron can be written as

$$y = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

where x_1, x_2, \dots, x_n are the inputs. That is, a linear neuron performs linear regression and the weights and biases (that is, b and w_1, \dots, w_n) act as regression coefficients. The above function forms a **hyperplane** in Euclidean space \mathbb{R}^n .

Given a training examples $\{(x_p, d_p)\}_{p=1}^P$ where input $x_p \in \mathbb{R}^n$ and target $d_p \in \mathbb{R}$, training a linear neuron finds a linear mapping $\phi: \mathbb{R}^n \rightarrow \mathbb{R}$ given by:

$$y = \mathbf{w}^T \mathbf{x} + b$$

Stochastic gradient descent (SGD) for linear neuron

The **cost function** J for regression is usually given as the *square error* (s.e.) between neuron outputs and targets.

Given a training pattern (\mathbf{x}, d) , $\frac{1}{2}$ square error cost J is defined as

$$J = \frac{1}{2}(d - y)^2$$

where y is neuron output for input pattern \mathbf{x} and d is the target label.

$$y = \mathbf{w}^T \mathbf{x} + b$$

The $\frac{1}{2}$ in the cost function is introduced to simplify learning equations and does not affect the optimal values of the parameters (weights and bias).

SGD for linear neuron

$$u = \mathbf{w}^T \mathbf{x} + b = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

$$\frac{\partial u}{\partial \mathbf{w}} = \begin{pmatrix} \frac{\partial u}{\partial w_1} \\ \frac{\partial u}{\partial w_2} \\ \vdots \\ \frac{\partial u}{\partial w_n} \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \mathbf{x}$$
(C)

$$\frac{\partial u}{\partial \mathbf{w}} = \mathbf{x}$$

Substituting (A) and (C) in (B),

$$\nabla_{\mathbf{w}} J = -(d - y) \mathbf{x} \quad (\text{D})$$

Similarly, since $\frac{\partial u}{\partial b} = 1$,

$$\nabla_b J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d - y) \quad (\text{E})$$

SGD for linear neuron



$$J = \frac{1}{2}(d - y)^2$$

$$y = u = \mathbf{w}^T \mathbf{x} + b$$

$$\frac{\partial J}{\partial u} = \frac{\partial J}{\partial y} = -(d - y) \quad (\text{A})$$

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial \mathbf{w}} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial \mathbf{w}} \quad (\text{B})$$

SGD for linear neuron

Gradient learning equations:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J \\ b &\leftarrow b - \alpha \nabla_b J \end{aligned}$$

By substituting from (D) and (E), SGD equations for a linear neuron are given by

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \alpha(d - y) \mathbf{x} \\ b &\leftarrow b + \alpha(d - y) \end{aligned}$$

SGD algorithm for linear neuron

Given a training dataset $\{(x_p, d_p)\}_{p=1}^P$

Set learning parameter α

Initialize w and b

Repeat until convergence:

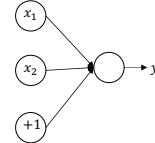
For every training pattern (x_p, d_p) :

$$y_p = w^T x_p + b$$

$$w \leftarrow w + \alpha(d_p - y_p)x_p$$

$$b \leftarrow b + \alpha(d_p - y_p)$$

Example 1



Let's initialize weights randomly and biases to zeros

$$w = \begin{pmatrix} 0.92 \\ 0.71 \end{pmatrix} \text{ and } b = 0.0$$

$$\alpha = 0.01$$

Need to shuffle the patterns before every epoch.

11

Example 1: SGD on a linear neuron

Train a linear neuron to perform the following mapping, using stochastic gradient descent (SGD) learning:

$x^T = (x_1, x_2)$	d
(0.54, -0.96)	1.33
(0.27, 0.50)	0.45
(0.00, -0.55)	0.56
(-0.60, 0.52)	-1.66
(-0.66, -0.82)	-1.07
(0.37, 0.91)	0.30

Use a learning factor $\alpha = 0.01$.

Example 1: epoch 1

Epoch 1 begins

Shuffle the patterns

First pattern $p = 1$ is applied:

$$x_p = \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix} \text{ and } d_p = 1.33$$

$$y_p = w^T x_p + b = (0.92 \quad 0.71) \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix} + 0.0 = -0.19$$

$$\text{s.e.} = (d_p - y_p)^2 = 2.292$$

$$w = w + \alpha(d_p - y_p)x_p = \begin{pmatrix} 0.92 \\ 0.71 \end{pmatrix} + 0.01 \times (1.33 + 0.19) \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix} = \begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix}$$

$$b = b + \alpha(d_p - y_p) = 0 + 0.01 \times (1.33 + 0.19) = 0.02$$

Example 1: epoch 1

Second pattern $p = 2$ is applied:

$$\mathbf{x}_p = \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix} \text{ and } d_p = -1.07$$

$$y_p = \mathbf{w}^T \mathbf{x}_p + b = (0.93 \quad 0.70) \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix} + 0.02 = -0.19$$

$$\text{s.e.} = (d_p - y_p)^2 = 0.01$$

$$\mathbf{w} = \mathbf{w} + \alpha(d_p - y_p)\mathbf{x}_p = \begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix} + 0.01 \times (-1.07 + 0.19) \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix} = \begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix}$$

$$b = b + \alpha(d_p - y_p) = 0.02 + 0.01 \times (1.33 + 0.19) = 0.02$$

Iterations continues for patterns $p = 3, \dots, 6$.

the second epoch starts

Shuffle the patterns

Apply patterns $p = 1, 2, \dots, 6$

Training epochs continue until convergence.

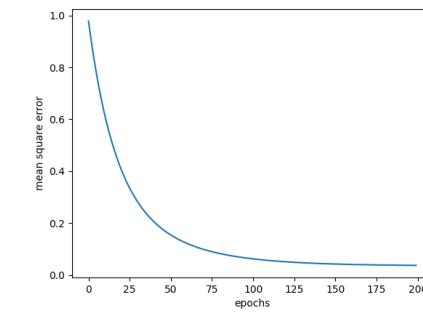
Example 1: epoch 200

\mathbf{x}_p	y_p	s.e.	\mathbf{w}	b
$\mathbf{x}_1 = \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix}$	1.49	0.03	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$\mathbf{x}_2 = \begin{pmatrix} 0.00 \\ -0.55 \end{pmatrix}$	0.22	0.12	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$\mathbf{x}_3 = \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix}$	-0.98	0.01	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$\mathbf{x}_4 = \begin{pmatrix} 0.37 \\ 0.91 \end{pmatrix}$	0.33	0.00	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$\mathbf{x}_5 = \begin{pmatrix} 0.27 \\ 0.50 \end{pmatrix}$	0.30	0.02	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$\mathbf{x}_6 = \begin{pmatrix} -0.60 \\ 0.52 \end{pmatrix}$	-1.45	0.04	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01

Example 1: epoch 1

\mathbf{x}_p	y_p	s.e.	\mathbf{w}	b
$\mathbf{x}_1 = \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix}$	-0.19	2.29	$\begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix}$	0.02
$\mathbf{x}_2 = \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix}$	-1.17	0.01	$\begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix}$	0.03
$\mathbf{x}_3 = \begin{pmatrix} 0.00 \\ -0.55 \end{pmatrix}$	-0.37	0.87	$\begin{pmatrix} 0.93 \\ 0.69 \end{pmatrix}$	0.03
$\mathbf{x}_4 = \begin{pmatrix} 0.27 \\ 0.50 \end{pmatrix}$	0.62	0.03	$\begin{pmatrix} 0.92 \\ 0.69 \end{pmatrix}$	0.02
$\mathbf{x}_5 = \begin{pmatrix} -0.60 \\ 0.52 \end{pmatrix}$	-0.17	2.21	$\begin{pmatrix} 0.93 \\ 0.69 \end{pmatrix}$	0.01
$\mathbf{x}_6 = \begin{pmatrix} 0.37 \\ 0.91 \end{pmatrix}$	0.98	0.45	$\begin{pmatrix} 0.93 \\ 0.68 \end{pmatrix}$	0.00

Example 1



$$\text{m.s.e.} = \frac{1}{6} \sum_{p=1}^6 (d_p - y_p)^2$$

Example 1

At convergence:

$$\mathbf{w} = \begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$$

$$b = -0.013$$

Mean square error = 0.037

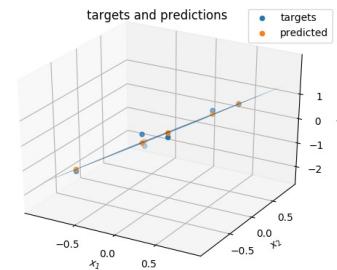
The regression equation:

$$y = \mathbf{x}^T \mathbf{w} + b = (x_1 \ x_2) \begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix} - 0.013$$

The mapping learnt by the linear neuron:

$$y = 2.00x_1 - 0.44x_2 - 0.013$$

Example 1



The mapping portrays a hyperplane in the 3-dimensional space:

$$y = 2.00x_1 - 0.44x_2 - 0.013$$

Example 1

inputs $x = (x_1, x_2)$	predictions $y = 2.00x_1 - 0.44x_2 - 0.01$	targets d
(0.54, -0.96)	1.49	1.33
(0.27, 0.50)	0.30	0.45
(0.00, -0.55)	0.22	0.56
(-0.60, 0.52)	-1.45	-1.66
(-0.66, -0.82)	-0.98	-1.07
(0.37, 0.91)	0.33	0.30

Gradient descent (GD) for linear neuron

Given a training dataset $\{(x_p, d_p)\}_{p=1}^P$, cost function J is given by the sum of square errors (s.s.e.):

$$J = \frac{1}{2} \sum_{p=1}^P (d_p - y_p)^2$$

where y_p is the neuron output for input pattern x_p .

$$J = \sum_{p=1}^P J_p \quad (\text{F})$$

where $J_p = \frac{1}{2}(d_p - y_p)^2$ is the square error for the p th pattern.

GD for linear neuron

GD for linear neuron

From (F):

$$\begin{aligned}
 \nabla_w J &= \sum_{p=1}^P \nabla_w J_p \\
 &= -\sum_{p=1}^P (d_p - y_p) \mathbf{x}_p \\
 &= -((d_1 - y_1) \mathbf{x}_1 + (d_2 - y_2) \mathbf{x}_2 + \dots + (d_P - y_P) \mathbf{x}_P) \\
 &= -(\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_P) \begin{pmatrix} (d_1 - y_1) \\ (d_2 - y_2) \\ \vdots \\ (d_P - y_P) \end{pmatrix} \\
 &= -\mathbf{X}^T (\mathbf{d} - \mathbf{y})
 \end{aligned} \tag{G}$$

from (D)

where $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix}$ is the data matrix, $\mathbf{d} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_P \end{pmatrix}$ is the target vector, and $\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}$ is the output vector.

Substituting (G) and (H) in gradient descent equations:

$$\begin{aligned}
 \mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_w J \\
 b &\leftarrow b - \alpha \nabla_b J
 \end{aligned}$$

We get GD learning equations for the linear neuron as

$$\begin{aligned}
 \mathbf{w} &\leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y}) \\
 b &\leftarrow b + \alpha \mathbf{1}_P^T (\mathbf{d} - \mathbf{y})
 \end{aligned}$$

And α is the learning factor.

Where:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + b \mathbf{1}_P$$

23

GD for linear neuron

Similarly, $\nabla_b J$ can be obtained by considering inputs of +1 and substituting a vector of +1 in (G):

$$\nabla_b J = -\mathbf{1}_P^T (\mathbf{d} - \mathbf{y}) \tag{H}$$

where $\mathbf{1}_P = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$ has P elements of 1.

The output vector \mathbf{y} for the batch of P patterns is given by

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_P \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \mathbf{w} + b \\ \mathbf{x}_2^T \mathbf{w} + b \\ \vdots \\ \mathbf{x}_P^T \mathbf{w} + b \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix} \mathbf{w} + b \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \mathbf{X}\mathbf{w} + b \mathbf{1}_P$$

GD for linear neuron

Given a training dataset (\mathbf{X}, \mathbf{d})

Set learning parameter α

Initialize \mathbf{w} and b

Repeat until convergence:

$$\begin{aligned}
 \mathbf{y} &= \mathbf{X}\mathbf{w} + b \mathbf{1}_P \\
 \mathbf{w} &\leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y}) \\
 b &\leftarrow b + \alpha \mathbf{1}_P^T (\mathbf{d} - \mathbf{y})
 \end{aligned}$$

22

24

GD and SGD for a linear neuron

GD	SGD
(\mathbf{X}, \mathbf{d})	(\mathbf{x}_p, d_p)
$J = \frac{1}{2} \sum_{p=1}^P (d_p - y_p)^2$	$J = \frac{1}{2} (d_p - y_p)^2$
$\mathbf{y} = \mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_P$	$y_p = u_p = \mathbf{x}_p^T \mathbf{w} + b$
$\mathbf{w} \leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y})$	$\mathbf{w} \leftarrow \mathbf{w} + \alpha (d_p - y_p) \mathbf{x}_p$
$b \leftarrow b + \alpha \mathbf{1}_P^T (\mathbf{d} - \mathbf{y})$	$b \leftarrow b + \alpha (d_p - y_p)$

25

SGD for perceptron



Cost function J is given by

$$J = \frac{1}{2} (d - y)^2$$

where $y = f(u)$ and $u = \mathbf{w}^T \mathbf{x} + b$

The gradient with respect to the synaptic input:

$$\frac{\partial J}{\partial u} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial u} = -(d - y) f'(u)$$

From (C), $\frac{\partial u}{\partial \mathbf{w}} = \mathbf{x}$ and $\frac{\partial u}{\partial b} = 1$.

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial \mathbf{w}} = -(d - y) f'(u) \mathbf{x} \quad (I)$$

$$\nabla_b J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d - y) f'(u) \quad (J)$$

Perceptron

Perceptron is a neuron having a **sigmoid** activation function and has an output

$$y = f(u)$$

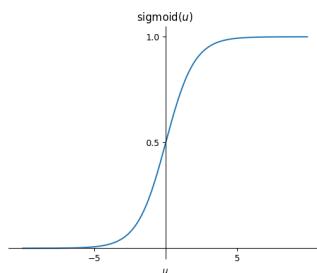
Where:

$$f(u) = \frac{1}{1 + e^{-u}} = \text{sigmoid}(u)$$

And $u = \mathbf{w}^T \mathbf{x} + b$

The square error is used as cost function for learning.

Perceptron performs a *non-linear regression* of inputs.



SGD for perceptron

Gradient learning equations:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J \\ b &\leftarrow b - \alpha \nabla_b J \end{aligned}$$

Substituting gradients from (I) and (J), SGD equations for a perceptron are given by

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \alpha (d - y) f'(u) \mathbf{x} \\ b &\leftarrow b + \alpha (d - y) f'(u) \end{aligned}$$

SGD algorithm for perceptron

Given a training dataset $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$

Set learning parameter α

Initialize \mathbf{w} and b

Repeat until convergence:

For every training pattern (\mathbf{x}_p, d_p) :

$$u_p = \mathbf{w}^T \mathbf{x}_p + b$$

$$y_p = f(u_p) = \frac{1}{1+e^{-u_p}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha(d_p - y_p)f'(u_p)\mathbf{x}_p$$

$$b \leftarrow b + \alpha(d_p - y_p)f'(u_p)$$

GD for perceptron

From (F):

$$\begin{aligned} \nabla_{\mathbf{w}} J &= \sum_{p=1}^P \nabla_{\mathbf{w}} J_p \\ &= - \sum_{p=1}^P (d_p - y_p)f'(u_p)\mathbf{x}_p && \text{From (J)} \\ &= -((d_1 - y_1)f'(u_1)\mathbf{x}_1 + (d_2 - y_2)f'(u_2)\mathbf{x}_2 + \dots + (d_P - y_P)f'(u_P)\mathbf{x}_P) \\ &= -(\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_P) \begin{pmatrix} (d_1 - y_1)f'(u_1) \\ (d_2 - y_2)f'(u_2) \\ \vdots \\ (d_P - y_P)f'(u_P) \end{pmatrix} \\ &= -\mathbf{X}^T(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \end{aligned} \tag{K}$$

$$\text{where } \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix}, \mathbf{d} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_P \end{pmatrix}, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}, \text{ and } f'(\mathbf{u}) = \begin{pmatrix} f'(u_1) \\ f'(u_2) \\ \vdots \\ f'(u_P) \end{pmatrix}$$

GD for perceptron

Given a training dataset $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$, cost function J is given by the sum of square errors (s.s.e) over all the patterns:

$$J = \frac{1}{2} \sum_{p=1}^P (d_p - y_p)^2 = \sum_{p=1}^P J_p \tag{F}$$

where $J_p = \frac{1}{2}(d_p - y_p)^2$ is the square error for the p th pattern.

GD for perceptron

Substituting \mathbf{X}^T by $\mathbf{1}_P^T$ in (K), we get

$$\nabla_b J = -\mathbf{1}_P^T(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \tag{L}$$

where $\mathbf{1}_P = (1 \quad 1 \quad \dots \quad 1)^T$.

The gradient descent learning is given by

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J \\ b &\leftarrow b - \alpha \nabla_b J \end{aligned}$$

Substituting (K) and (L), we get the learning equations:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \alpha \mathbf{X}^T(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \\ b &\leftarrow b + \alpha \mathbf{1}_P^T(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \end{aligned}$$

Note that \cdot is the element-wise product.

GD for perceptron

Given a training dataset (\mathbf{X}, \mathbf{d})

Set learning parameter α

Initialize \mathbf{w} and b

Repeat until convergence:

$$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_P$$

$$\mathbf{y} = f(\mathbf{u}) = \frac{1}{1+e^{-\mathbf{u}}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$$

$$b \leftarrow b + \alpha \mathbf{1}_P^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$$

Derivatives of Sigmoid

The activation function of the (continuous) perceptron is *sigmoid function* (i.e., unipolar sigmoidal function with $a = 1.0$ and $b = 1.0$):

$$y = f(u) = \frac{1}{1+e^{-u}}$$

The derivative is given by

$$f'(u) = \frac{-1}{(1+e^{-u})^2} = \frac{e^{-u}}{(1+e^{-u})^2} = \frac{1}{1+e^{-u}} - \frac{1}{(1+e^{-u})^2} = \mathbf{y}(1-\mathbf{y})$$

For *Tanh* function (bipolar sigmoid):

$$y = f(u) = \frac{e^{+u} - e^{-u}}{e^{+u} + e^{-u}}$$

$$f'(u) = \frac{(e^{+u} + e^{-u})(e^{+u} + e^{-u}) - (e^{+u} - e^{-u})(e^{+u} - e^{-u})}{(e^{+u} + e^{-u})^2} = \left(1 - \left(\frac{e^{+u} - e^{-u}}{e^{+u} + e^{-u}}\right)^2\right) = 1 - \mathbf{y}^2$$

33

Gradient descent for perceptron

GD	SGD
(\mathbf{X}, \mathbf{d})	(\mathbf{x}_p, d_p)
$J = \frac{1}{2} \sum_{p=1}^P (d_p - y_p)^2$	$J = \frac{1}{2} (d_p - y_p)^2$
$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_P$	$u_p = \mathbf{x}_p^T \mathbf{w} + b$
$\mathbf{y} = f(\mathbf{u})$	$y_p = f(u_p)$
$\mathbf{w} = \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$	$\mathbf{w} = \mathbf{w} + \alpha (d_p - y_p) f'(u_p) \mathbf{x}_p$
$b = b + \alpha \mathbf{1}_P^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$	$b = b + \alpha (d_p - y_p) f'(u_p)$

Example 2

Design a perceptron to learn the following mapping by using gradient descent (GD):

$\mathbf{x} = (x_1, x_2)$	d
(0.77, 0.02)	2.91
(0.63, 0.75)	0.55
(0.50, 0.22)	1.28
(0.20, 0.76)	-0.74
(0.17, 0.09)	0.88
(0.69, 0.95)	0.30
(0.00, 0.51)	-0.28

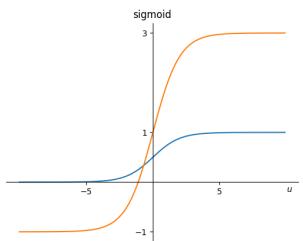
Use learning factor $\alpha = 0.01$.

34

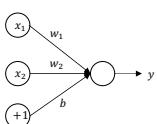
36

Example 2

$$X = \begin{pmatrix} 0.77 & 0.02 \\ 0.63 & 0.75 \\ 0.50 & 0.22 \\ 0.20 & 0.76 \\ 0.17 & 0.09 \\ 0.69 & 0.95 \\ 0.00 & 0.51 \end{pmatrix} \text{ and } d = \begin{pmatrix} 2.91 \\ 0.55 \\ 1.28 \\ -0.74 \\ 0.88 \\ 0.30 \\ -0.28 \end{pmatrix}$$



37



Initially, $w = \begin{pmatrix} 0.81 \\ 0.61 \end{pmatrix}$ and $b = 0.0$
 $\alpha = 0.01$

Output $y \in [-0.74, 2.91] \subset [-1.0, 3.0]$

Note that the sigmoidal should have an amplitude = 4 and shifted downwards by 1.0.

So, the activation function should be

$$y = f(u) = \frac{4.0}{1 + e^{-u}} - 1.0$$

$$f'(u) = \frac{4e^{-u}}{(1+e^{-u})^2} = (y+1) \frac{e^{-u}}{(1+e^{-u})} = (y+1) \left(1 - \frac{1}{1+e^{-u}}\right) = \frac{1}{4} (y+1)(3-y)$$

Example 2

$$f'(u) = \frac{1}{4} (y+1) \cdot (3-y) = \frac{1}{4} \left(\begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix} + \begin{pmatrix} 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \end{pmatrix} \right) \cdot \left(\begin{pmatrix} 3.0 \\ 3.0 \\ 3.0 \\ 3.0 \\ 3.0 \\ 3.0 \\ 3.0 \end{pmatrix} - \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix} \right) = \begin{pmatrix} 0.90 \\ 0.80 \\ 0.93 \\ 0.91 \\ 0.99 \\ 0.73 \\ 0.98 \end{pmatrix}$$

Epoch 1 begins ...

$$\mathbf{u} = \mathbf{Xw} + b\mathbf{1}_P = \begin{pmatrix} 0.77 & 0.02 \\ 0.63 & 0.75 \\ 0.50 & 0.22 \\ 0.20 & 0.76 \\ 0.17 & 0.09 \\ 0.69 & 0.95 \\ 0.00 & 0.51 \end{pmatrix} \begin{pmatrix} 0.81 \\ 0.61 \end{pmatrix} + 0.0 \begin{pmatrix} 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \end{pmatrix} = \begin{pmatrix} 0.64 \\ 0.97 \\ 0.54 \\ 0.63 \\ 0.19 \\ 1.14 \\ 0.32 \end{pmatrix}$$

$$y = f(\mathbf{u}) = \frac{4.0}{1 + e^{-u}} - 1.0 = \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix}$$

$$m.s.e. = \frac{1}{7} \sum_{p=1}^7 (d_p - y_p)^2 = \frac{1}{7} ((2.91 - 1.61)^2 + (0.55 - 1.90)^2 + \dots) = 2.11$$

38

Example 2

$$\mathbf{w} = \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \\ = \begin{pmatrix} 0.81 \\ 0.61 \end{pmatrix} + 0.01 \begin{pmatrix} 0.77 & 0.63 & 0.50 & 0.20 & 0.17 & 0.69 & 0.00 \\ 0.02 & 0.75 & 0.22 & 0.76 & 0.09 & 0.95 & 0.51 \end{pmatrix} \left(\begin{pmatrix} 2.91 \\ 0.55 \\ 1.28 \\ -0.74 \\ 0.88 \\ 0.30 \\ -0.28 \end{pmatrix} - \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix} \right) \cdot \begin{pmatrix} 0.90 \\ 0.80 \\ 0.93 \\ 0.91 \\ 0.99 \\ 0.73 \\ 0.98 \end{pmatrix} = \begin{pmatrix} 0.80 \\ 0.57 \end{pmatrix}$$

$$b = b + \alpha \mathbf{1}_P^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \\ = 0.0 + 0.01 \times (1.0 \ 1.0 \ 1.0 \ 1.0 \ 1.0 \ 1.0 \ 1.0) \left(\begin{pmatrix} 2.91 \\ 0.55 \\ 1.28 \\ -0.74 \\ 0.88 \\ 0.30 \\ -0.28 \end{pmatrix} - \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix} \right) \cdot \begin{pmatrix} 0.90 \\ 0.80 \\ 0.93 \\ 0.91 \\ 0.99 \\ 0.73 \\ 0.98 \end{pmatrix} = -0.05$$

40

Example 2

iter	u	y	$f'(u)$	mse	w	b
1	$\begin{pmatrix} 0.64 \\ 0.97 \\ 0.54 \\ 0.63 \\ 0.19 \\ 1.14 \\ 0.32 \end{pmatrix}$	$\begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix}$	$\begin{pmatrix} 0.90 \\ 0.80 \\ 0.93 \\ 0.91 \\ 0.99 \\ 0.73 \\ 0.98 \end{pmatrix}$	2.11	$\begin{pmatrix} 0.80 \\ 0.57 \end{pmatrix}$	-0.05
2	$\begin{pmatrix} 0.57 \\ 0.88 \\ 0.47 \\ 0.54 \\ 1.04 \\ 1.95 \\ 0.24 \end{pmatrix}$	$\begin{pmatrix} 1.56 \\ 1.83 \\ 1.46 \\ 1.52 \\ 1.13 \\ 1.95 \\ 1.24 \end{pmatrix}$	$\begin{pmatrix} 0.92 \\ 0.83 \\ 0.95 \\ 0.93 \\ 1.00 \\ 0.77 \\ 0.99 \end{pmatrix}$	1.96	$\begin{pmatrix} 0.79 \\ 0.52 \end{pmatrix}$	-0.11
1500	$\begin{pmatrix} 2.05 \\ -0.45 \\ 0.56 \\ -1.94 \\ -0.16 \\ -0.85 \\ -1.89 \end{pmatrix}$	$\begin{pmatrix} 2.54 \\ 0.56 \\ 1.55 \\ -0.50 \\ 0.84 \\ 0.20 \\ -0.48 \end{pmatrix}$	$\begin{pmatrix} 0.40 \\ 0.95 \\ 0.92 \\ 0.44 \\ 0.99 \\ 0.84 \\ 0.45 \end{pmatrix}$	0.046	$\begin{pmatrix} 3.35 \\ -2.80 \end{pmatrix}$	-0.47

41

Example 2

At convergence:

$$\mathbf{w} = \begin{pmatrix} 3.35 \\ -2.80 \end{pmatrix}$$

$$b = -0.47$$

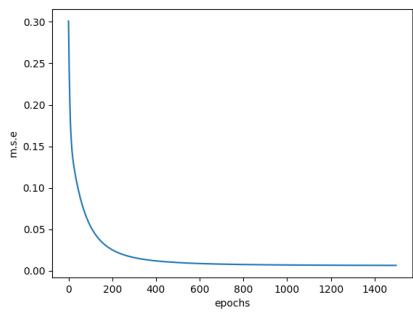
Mean square error = 0.01

$$u = \mathbf{x}^T \mathbf{w} + b = (x_1 \quad x_2) \begin{pmatrix} 3.35 \\ -2.80 \end{pmatrix} - 0.47 = 3.35x_1 - 2.8x_2 - 0.47$$

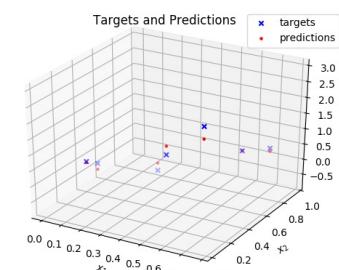
$$y = \frac{4.0}{1 + e^{-u}} - 1.0$$

$$y = \frac{4.0}{1 + e^{-3.35x_1 + 2.8x_2 + 0.47}} - 1.0$$

Example 2



Example 2



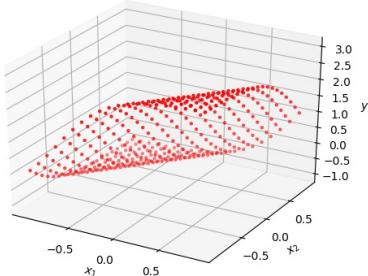
42

44

Example 2

Non-linear function learnt by the perceptron

$$y = \frac{4.0}{1 + e^{-3.35x_1 + 2.8x_2 + 0.47}} - 1.0$$



Classification Example

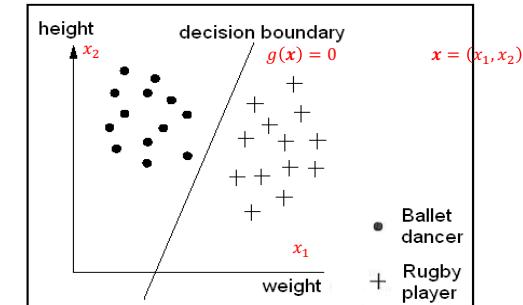


Figure: A linear classification decision boundary

47

Classification Example

Classification is to identify or distinguish classes or groups of objects.

Example: To identify *ballet dancers* from *rugby players*.

Two distinctive *features* that can aid in classification:

- *weight*
- *height*

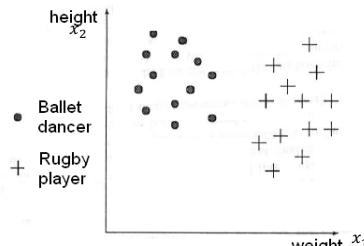


Figure: A 2-dimensional feature space

Let x_1 denote weight and x_2 denote height. Every individual is represented as a point $\mathbf{x} = (x_1, x_2)$ in the feature feature space.

Decision boundary

The **decision boundary** of the classifier, $g(\mathbf{x}) = 0$ where function $g(\mathbf{x})$ is referred to as the **discriminant function**.

A classifier finds a decision boundary separating the two classes in the feature space. On one side of the decision boundary, discriminant function is positive and on other side, discriminant function is negative.

Therefore, the following class definition may be employed:

If $g(\mathbf{x}) > 0 \Rightarrow$ Ballet dancer

If $g(\mathbf{x}) \leq 0 \Rightarrow$ Rugby player

Linear Classifier

If the two classes can be separated by a straight line, the classification is said to be **linearly separable**. For linear separable classes, one can design a **linear classifier**.

A linear classifier implements discriminant function or a decision boundary that is represented by a straight line (hyper plane) in the multidimensional **feature space**. Generally, the feature space is multidimensional. In the multidimensional space, a straight line or **hyperplane** is indicated by a linear sum of coordinates.

Given an input (features), $\mathbf{x} = (x_1 \ x_2 \ \dots \ x_n)^T$. A linear description function is given by

$$g(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

where $\mathbf{w} = (w_1 \ w_2 \ \dots \ w_n)^T$ are the coefficient/weights and w_0 is the constant term.

49

Classification error

In classification, the error is expressed as total mismatches between the target and output labels.

$$\text{Classification error} = \sum_{p=1}^P 1(d_p \neq y_p)$$

51

Discrete perceptron as a linear two-class classifier

The linear discriminant function can be implemented by the synaptic input to a neuron

$$g(\mathbf{x}) = u = \mathbf{w}^T \mathbf{x} + b$$

And with a threshold activation function $1(u)$:

$$\begin{aligned} u = g(\mathbf{x}) > 0 &\rightarrow y = 1 \rightarrow \text{class1} \\ u = g(\mathbf{x}) \leq 0 &\rightarrow y = 0 \rightarrow \text{class2} \end{aligned}$$

That is, two-class linear classifier (or a dichotomizer) can be implemented with an artificial neuron with a threshold (unit step) activation function (**discrete perceptron**).

The output, 0 or 1, of the binary neuron represents the **label** of the class.

50

Logistic regression neuron

A **logistic regression neuron** performs a binary classification of inputs. That is, it classifies inputs into two classes with labels '0' and '1'.

The activation function of the logistic regression neuron is given by the sigmoid function:

$$f(u) = \frac{1}{1 + e^{-u}}$$

where $u = \mathbf{w}^T \mathbf{x} + b$ is the synaptic input to the neuron.

The activation of a logistic regression neuron gives the probability of the neuron output belonging to class '1'.

$$P(y = 1|\mathbf{x}) = f(u)$$

Then

$$P(y = 0|\mathbf{x}) = 1 - P(y = 1|\mathbf{x}) = 1 - f(u)$$

52

Logistic Regression Neuron

A logistic regression neuron receives an input $\mathbf{x} \in \mathbb{R}^n$ and produces a class label $y \in \{0, 1\}$ as the output.

$$f(u) = \frac{1}{1 + e^{-u}}$$

When $u = 0$, $f(u) = P(y = 1|\mathbf{x}) = P(y = 0|\mathbf{x}) = 0.5$.
That is, $y = 1$ if $f(u) > 0.5$, else $y = 0$.

The output y of the neuron is given by:

$$y = 1(f(u) > 0.5) = 1(u > 0)$$

Note that for logistic neuron, the output and activation are different. It finds a linear boundary $u = 0$ separating the two classes.

SGD for logistic regression neuron

$$J = -d \log(f(u)) - (1-d)\log(1-f(u))$$

where $u = \mathbf{w}^T \mathbf{x} + b$ and $f(u) = \frac{1}{1+e^{-(u)}}$.

Gradient with respect to u :

$$\begin{aligned} \frac{\partial J}{\partial u} &= -\frac{\partial}{\partial f(u)} (d \log(f(u)) + (1-d)\log(1-f(u))) \frac{\partial f(u)}{\partial u} \\ &= -\left(\frac{d}{f(u)} - \frac{(1-d)}{1-f(u)}\right) f'(u) \end{aligned}$$

Substituting $f'(u) = f(u)(1-f(u))$ for sigmoid activation function,

$$\frac{\partial J}{\partial u} = -\frac{d-f(u)}{f(u)(1-f(u))} f(u)(1-f(u)) = -(d-f(u))$$

53

55

SGD for logistic regression neuron

Given a training pattern (\mathbf{x}, d) where $\mathbf{x} \in \mathbb{R}^n$ and $d \in \{0, 1\}$.

The cost function for classification is given by the **cross-entropy**:

$$J = -d \log(f(u)) - (1-d)\log(1-f(u))$$

where $u = \mathbf{w}^T \mathbf{x} + b$ and $f(u) = \frac{1}{1+e^{-u}}$.

The cost function J is minimized using the gradient descent procedure.

$$J = \begin{cases} -\log(f(u)) & \text{if } d = 1 \\ -\log(1-f(u)) & \text{if } d = 0 \end{cases}$$

SGD for logistic regression neuron

Substituting $\frac{\partial J}{\partial u}$, $\frac{\partial u}{\partial w} = \mathbf{x}$, and $\frac{\partial u}{\partial b} = 1$.

$$\nabla_w J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial w} = -(d-f(u)) \mathbf{x} \quad (\text{A})$$

$$\nabla_b J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d-f(u)) \quad (\text{B})$$

Gradient learning equations:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_w J \\ b &\leftarrow b - \alpha \nabla_b J \end{aligned}$$

Substituting $\nabla_w J$ and $\nabla_b J$ for logistic regression neuron

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \alpha(d-f(u)) \mathbf{x} \\ b &\leftarrow b + \alpha(d-f(u)) \end{aligned}$$

54

56

SGD for logistic regression neuron

Given a training dataset $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$

Set learning rate α

Initialize \mathbf{w} and b

Iterate until convergence:

For every pattern (\mathbf{x}_p, d_p) :

$$u_p = \mathbf{w}^T \mathbf{x}_p + b$$

$$f(u_p) = \frac{1}{1+e^{-u_p}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha (d_p - f(u_p)) \mathbf{x}_p$$

$$b \leftarrow b + \alpha (d_p - f(u_p))$$

GD for logistic regression neuron

$$\begin{aligned} \nabla_{\mathbf{w}} J &= \sum_{p=1}^P \nabla_{\mathbf{w}} J_p \\ &= - \sum_{p=1}^P (d_p - f(u_p)) \mathbf{x}_p \\ &= - ((d_1 - f(u_1)) \mathbf{x}_1 + (d_2 - f(u_2)) \mathbf{x}_2 + \dots + (d_P - f(u_P)) \mathbf{x}_P) \\ &= - (\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_P) \begin{pmatrix} (d_1 - f(u_1)) \\ (d_2 - f(u_2)) \\ \vdots \\ (d_P - f(u_P)) \end{pmatrix} \\ &= - \mathbf{X}^T (\mathbf{d} - f(\mathbf{u})) \end{aligned}$$

From (A)

$$\text{where } \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix}, \mathbf{d} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_P \end{pmatrix}, \text{ and } f(\mathbf{u}) = \begin{pmatrix} f(u_1) \\ f(u_2) \\ \vdots \\ f(u_P) \end{pmatrix}$$

57

59

GD for logistic regression neuron

Given a training dataset $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$ where $\mathbf{x}_p \in \mathbf{R}^n$ and $d_p \in \{0,1\}$.

The cost function for logistic regression is given by the *cross-entropy* (or *negative log-likelihood*) over all the training patterns:

$$J = - \sum_{p=1}^P d_p \log(f(u_p)) + (1 - d_p) \log(1 - f(u_p))$$

where $u_p = \mathbf{w}^T \mathbf{x}_p + b$ and $f(u_p) = \frac{1}{1+e^{-u_p}}$.

The cost function J can be written as

$$J = \sum_{p=1}^P J_p$$

where $J_p = -d_p \log(f(u_p)) - (1 - d_p) \log(1 - f(u_p))$ is cross-entropy due to p th pattern.

GD for logistic regression neuron

By substituting $\mathbf{1}_P$ for \mathbf{X} in above equation:

$$\nabla_b J = -\mathbf{1}_P^T (\mathbf{d} - f(\mathbf{u}))$$

Substituting the gradients in

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J \\ b &\leftarrow b - \alpha \nabla_b J \end{aligned}$$

the gradient descent learning for logistic regression neuron is given by

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - f(\mathbf{u})) \\ b &\leftarrow b + \alpha \mathbf{1}_P^T (\mathbf{d} - f(\mathbf{u})) \end{aligned}$$

Note that \mathbf{y} in the discrete perceptron is now replaced with $f(\mathbf{u})$ in logistic regression learning equations.

58

60

GD for logistic regression neuron

Given training data (\mathbf{X}, \mathbf{d})

Set learning rate α

Initialize \mathbf{w} and b

Iterate until convergence:

$$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_p$$

$$f(\mathbf{u}) = \frac{1}{1+e^{-\mathbf{u}}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - f(\mathbf{u}))$$

$$b \leftarrow b + \alpha \mathbf{1}_p^T (\mathbf{d} - f(\mathbf{u}))$$

Example 3: GD for logistic regression neuron

Train a logistic regression neuron to perform the following classification, using GD:

$$(1.33 \quad 0.72) \rightarrow \text{class } B$$

$$(-1.55 \quad -0.01) \rightarrow \text{class } A$$

$$(0.62 \quad -0.72) \rightarrow \text{class } A$$

$$(0.27 \quad 0.11) \rightarrow \text{class } A$$

$$(0.0 \quad -0.17) \rightarrow \text{class } A$$

$$(0.43 \quad 1.2) \rightarrow \text{class } B$$

$$(-0.97 \quad 1.03) \rightarrow \text{class } B$$

$$(0.23 \quad 0.45) \rightarrow \text{class } B$$

User a learning factor $\alpha = 0.04$.

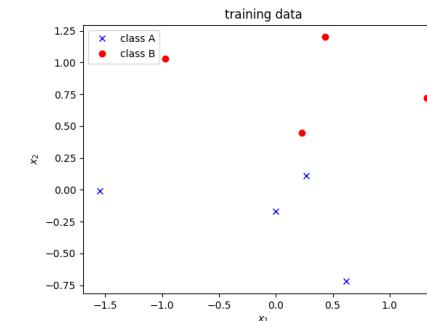
61

63

Learning for logistic regression neuron

GD	SGD
(\mathbf{X}, \mathbf{d})	(\mathbf{x}_p, d_p)
$J = -\sum_{p=1}^P d_p \log(f(u_p)) + (1-d_p) \log(1-f(u_p))$	$J_p = -d_p \log(f(u_p)) - (1-d_p) \log(1-f(u_p))$
$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_p$	$u_p = \mathbf{w}^T \mathbf{x}_p + b$
$f(\mathbf{u}) = \frac{1}{1+e^{-\mathbf{u}}}$	$f(u_p) = \frac{1}{1+e^{-u_p}}$
$\mathbf{y} = \mathbf{1}(f(\mathbf{u}) > 0.5)$	$y_p = \mathbf{1}(f(u_p) > 0.5)$
$\mathbf{w} \leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - f(\mathbf{u}))$	$\mathbf{w} \leftarrow \mathbf{w} + \alpha (d_p - f(u_p)) \mathbf{x}_p$
$b \leftarrow b + \alpha \mathbf{1}_p^T (\mathbf{d} - f(\mathbf{u}))$	$b \leftarrow b + \alpha (d_p - f(u_p))$

Example 3



62

64

Example 3

Let $y = 1$ for class A and $y = 0$ for class B.

$$X = \begin{pmatrix} 1.33 & 0.72 \\ -1.55 & 0.01 \\ 0.62 & -0.72 \\ 0.27 & 0.11 \\ 0.0 & -0.17 \\ 0.43 & 1.2 \\ -0.97 & 1.03 \\ 0.23 & 0.45 \end{pmatrix} \text{ and } \mathbf{d} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

Initially, $\mathbf{w} = \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix}$, $b = 0.0$ and $\alpha = 0.4$

65

Example 3

$$y = 1(f(\mathbf{u}) > 0.5) = 1 \left(\begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix} > 0.5 \right) = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ 1 \end{pmatrix}$$

$$\mathbf{d} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

Classification error = $\sum_{p=1}^8 1(d_p \neq y_p) = 5$

$$\begin{aligned} \text{Cross-entropy} &= -\sum_{p=1}^8 d_p \log(f(u_p)) + (1-d_p) \log(1-f(u_p)) \\ &= -\log(1-f(u_1)) - \log f(u_2) - \log f(u_3) - \dots - \log(1-f(u_8)) \\ &= -\log(1-0.74) - \log 0.23 - \log 0.61 - \dots - \log(1-0.55) \\ &= 6.653 \end{aligned}$$

67

Example 3

Epoch 1:

$$\mathbf{u} = X\mathbf{w} + b\mathbf{1} = \begin{pmatrix} 1.33 & 0.72 \\ -1.55 & 0.01 \\ 0.62 & -0.72 \\ 0.27 & 0.11 \\ 0.0 & -0.17 \\ 0.43 & 1.2 \\ -0.97 & 1.03 \\ 0.23 & 0.45 \end{pmatrix} \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.0 \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1.04 \\ -1.2 \\ 0.46 \\ 0.21 \\ 0.00 \\ 0.36 \\ -0.73 \\ 0.19 \end{pmatrix}$$

$$f(\mathbf{u}) = \frac{1}{1 + e^{(-\mathbf{u})}} = \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix}$$

66

Example 3

$$\mathbf{w} = \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - f(\mathbf{u}))$$

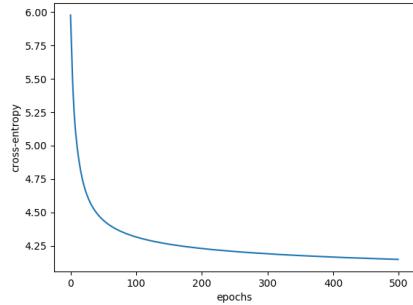
$$\begin{aligned} &= \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.04 \begin{pmatrix} 1.33 & -1.55 & 0.62 & 0.27 & 0 & 0.43 & -0.97 & 0.23 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} - \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix} \\ &= \begin{pmatrix} 0.69 \\ -0.2 \end{pmatrix} \end{aligned}$$

$$b = b + \alpha \mathbf{1}_p^T (\mathbf{d} - f(\mathbf{u}))$$

$$\begin{aligned} &= 0.0 + 0.04 (1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1) \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix} = -0.09 \end{aligned}$$

68

Example 3



Example 3

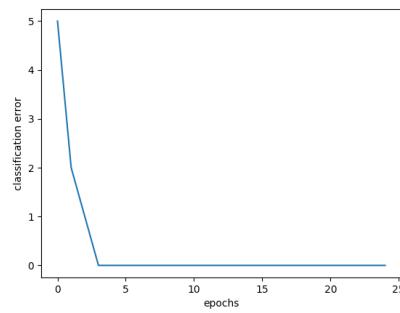
At convergence, $\mathbf{w} = \begin{pmatrix} -1.20 \\ -15.02 \end{pmatrix}, b = 4.47$

The decision boundary is given by: $u = \mathbf{x}^T \mathbf{w} + b = 0$
 $(x_1 \quad x_2)^T \begin{pmatrix} -1.20 \\ -15.02 \end{pmatrix} + 4.47 = 0$
 $-1.20x_1 - 15.02x_2 + 4.47 = 0$

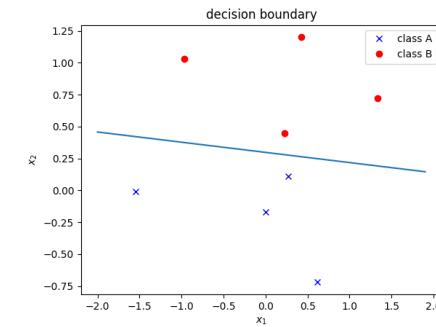
69

71

Example 3



Decision boundary:
 $-1.20x_1 - 15.02x_2 + 4.47 = 0$



70

72

Limitations of logistic regression neuron

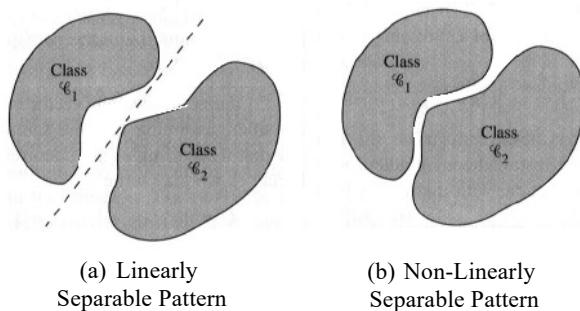
As long as the neuron is a *linear combiner* followed by a *non-linear activation function*, then regardless of the form of non-linearity used, the neuron can perform pattern classification *only on linearly separable patterns*.

Linear separability requires that the patterns to be classified must be sufficiently separated from each other to ensure that the decision boundaries are hyperplanes.

73

75

Limitations of logistic regression neuron



Discrete perceptron and logistic regression neuron can create only linear decision boundaries.

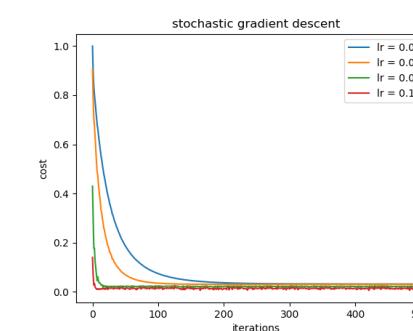
74

Example 4: effects of leaning rate

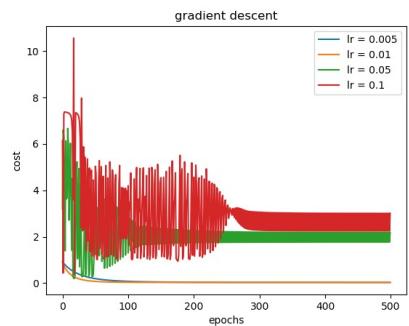
Design a perceptron to learn the following mapping by using gradient descent (GD):

$x = (x_1, x_2)$	d
(0.77, 0.02)	2.91
(0.63, 0.75)	0.55
(0.50, 0.22)	1.28
(0.20, 0.76)	-0.74
(0.17, 0.09)	0.88
(0.69, 0.95)	0.30
(0.00, 0.51)	-0.28

Use learning factor $\alpha = 0.01$.

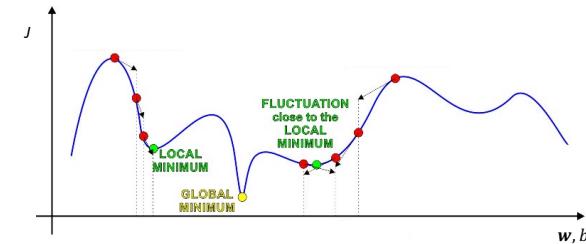


Example 4: Learning rates with GD



77

Local minima problem in gradient descent learning



Algorithm may get stuck in a local minimum of error function depending on the initial weights. Gradient descent gives a suboptimal solution and does not guarantee the optimal solution.

79

Learning rates

- At higher learning rates, convergence is faster but may not be stable.
- The *optimal learning rate* is the largest rate at which learning does not diverge.
- Generally, SGD converges to a better solution (lower error) as it capitalizes on randomness of data. SGD takes a longer time to converge
- Usually, GD can use a higher learning rate compared to SGD; The time for one add/multiply computation per a weight update is less when patterns are trained in a (small) batch.
- In practice, *mini-batch SGD* is used. Then, the time to train a network is dependent upon
 - the learning rate
 - the batch size

78

Summary: types of neurons

Role	Neuron
Regression (one dimensional)	Linear neuron Perceptron
Classification (two classes)	Logistic regression neuron

80

Summary: GD for neurons

$$\begin{aligned} & (\mathbf{X}, \mathbf{d}) \\ & \mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_P \\ & \mathbf{w} = \mathbf{w} - \alpha \mathbf{X}^T \nabla_u J \\ & b = b - \alpha \mathbf{1}_P^T \nabla_u J \end{aligned}$$

neuron	$f(\mathbf{u}), \mathbf{y}$	$\nabla_u J$
Logistic regression neuron	$f(\mathbf{u}) = \frac{1}{1 + e^{-\mathbf{u}}}$ $\mathbf{y} = \mathbf{1}(f(\mathbf{u}) > 0.5)$	$-(\mathbf{d} - f(\mathbf{u}))$
Linear neuron	$\mathbf{y} = \mathbf{u}$	$-(\mathbf{d} - \mathbf{y})$
Perceptron	$\mathbf{y} = f(\mathbf{u}) = \frac{1}{1 + e^{-\mathbf{u}}}$	$-(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$

Chapter 3

Neuron Layers

Neural networks and deep learning

81

© Nanyang Technological University

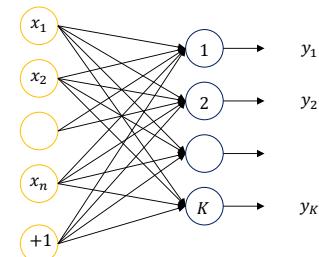
Summary: SGD for neurons

$$\begin{aligned} & (\mathbf{x}_p, d_p) \\ & u_p = \mathbf{w}^T \mathbf{x}_p + b \\ & \mathbf{w} = \mathbf{w} - \alpha \nabla_{u_p} J \mathbf{x}_p \\ & b = b - \alpha \nabla_{u_p} J \end{aligned}$$

neuron	$f(u_p), y_p$	$\nabla_{u_p} J$
Logistic regression neuron	$f(u_p) = \frac{1}{1 + e^{-u_p}}$ $y_p = \mathbf{1}(f(u_p) > 0.5)$	$-(d_p - f(u_p))$
Linear neuron	$y_p = u_p$	$-(d_p - y_p)$
Perceptron	$y_p = f(u_p) = \frac{1}{1 + e^{-u_p}}$	$-(d_p - y_p) \cdot f'(u_p)$

Weight matrix of a layer

Consider a layer of K neurons.



Let \mathbf{w}_k and b_k denote the weight vector and bias of k th neuron. Weights connected to a neuron layer is given by a weight matrix:

$$\mathbf{W} = (\mathbf{w}_1 \ \mathbf{w}_2 \ \cdots \ \mathbf{w}_K)$$

where columns are given by weight vectors of individual neurons.

And a bias vector \mathbf{b} where each element corresponds to a bias of a neuron:

$$\mathbf{b} = (b_1, b_2, \dots, b_K)^T$$

Synaptic input at a layer for single input

Given an input pattern $\mathbf{x} \in \mathbb{R}^n$ to a layer of K neurons.

Synaptic input u_k to k th neuron:

$$u_k = \mathbf{w}_k^T \mathbf{x} + b_k$$

\mathbf{w}_k and b_k denote the weight vector and bias of k th neuron.

Synaptic input vector \mathbf{u} to the layer :

$$\mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_K \end{pmatrix} = \begin{pmatrix} \mathbf{w}_1^T \mathbf{x} + b_1 \\ \mathbf{w}_2^T \mathbf{x} + b_2 \\ \vdots \\ \mathbf{w}_K^T \mathbf{x} + b_K \end{pmatrix} = \begin{pmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_K^T \end{pmatrix} \mathbf{x} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_K \end{pmatrix} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

$$\mathbf{u} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

where \mathbf{W} is the weight matrix and \mathbf{b} is the bias vector of the layer.

Activation at a layer for batch input

The synaptic input to the layer due to a batch of patterns:

$$\mathbf{U} = \mathbf{X} \mathbf{W} + \mathbf{B}$$

where rows of \mathbf{U} corresponds to synaptic inputs of the layer, corresponding to individual input patterns:

Activation of the layer:

$$f(\mathbf{U}) = \begin{pmatrix} f(\mathbf{u}_1^T) \\ f(\mathbf{u}_2^T) \\ \vdots \\ f(\mathbf{u}_P^T) \end{pmatrix} = \begin{pmatrix} f(\mathbf{u}_1)^T \\ f(\mathbf{u}_2)^T \\ \vdots \\ f(\mathbf{u}_P)^T \end{pmatrix}$$

where activations due to individual patterns are written as rows.

Synaptic input to a layer for batch input

Given a set $\{\mathbf{x}_p\}_{p=1}^P$ input patterns to a layer of K neurons where $\mathbf{x}_p \in \mathbb{R}^n$.

Synaptic input \mathbf{u}_p to the layer for an input pattern \mathbf{x}_p :

$$\mathbf{u}_p = \mathbf{W}^T \mathbf{x}_p + \mathbf{b}$$

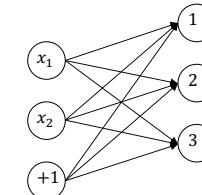
The synaptic input matrix \mathbf{U} to the layer for P patterns: $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$

$$\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_P^T \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \mathbf{W} + \mathbf{b}^T \\ \mathbf{x}_2^T \mathbf{W} + \mathbf{b}^T \\ \vdots \\ \mathbf{x}_P^T \mathbf{W} + \mathbf{b}^T \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix} \mathbf{W} + \begin{pmatrix} \mathbf{b}^T \\ \mathbf{b}^T \\ \vdots \\ \mathbf{b}^T \end{pmatrix} = \mathbf{X} \mathbf{W} + \mathbf{B}$$

where rows of \mathbf{U} are synaptic inputs corresponding to individual input patterns.

The matrix $\mathbf{B} = \begin{pmatrix} \mathbf{b}^T \\ \mathbf{b}^T \\ \vdots \\ \mathbf{b}^T \end{pmatrix}$ has bias vector propagated as rows.

Example 1: activations and outputs of a perceptron layer



A perceptron layer of 3 neurons shown in the figure receives 2-dimensional inputs $(x_1, x_2)^T$, and has a weight matrix \mathbf{W} and a bias vector \mathbf{b} given by

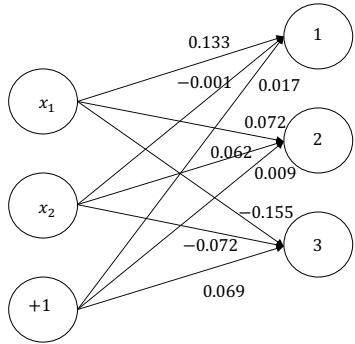
$$\mathbf{W} = \begin{pmatrix} 0.133 & 0.072 & -0.155 \\ -0.001 & 0.062 & -0.072 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} 0.017 \\ 0.009 \\ 0.069 \end{pmatrix}$$

Using batch processing, find the output for input patterns:

$$\begin{pmatrix} 0.5 \\ -1.66 \end{pmatrix}, \begin{pmatrix} -1.0 \\ -0.51 \end{pmatrix}, \begin{pmatrix} 0.78 \\ -0.65 \end{pmatrix}, \text{ and } \begin{pmatrix} 0.04 \\ -0.2 \end{pmatrix}.$$

Example 1

$$\mathbf{w} = \begin{pmatrix} 0.133 & 0.072 & -0.155 \\ -0.001 & 0.062 & -0.072 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} 0.017 \\ 0.009 \\ 0.069 \end{pmatrix}.$$



Example 1

$$\mathbf{U} = \begin{pmatrix} 0.085 & -0.059 & 0.111 \\ -0.115 & 0.094 & 0.26 \\ 0.121 & 0.024 & -0.005 \\ 0.022 & -0.001 & 0.077 \end{pmatrix}$$

For a perceptron layer

$$\mathbf{Y} = f(\mathbf{U}) = \frac{1}{1 + e^{-\mathbf{U}}} = \begin{pmatrix} 0.521 & 0.485 & 0.527 \\ 0.471 & 0.476 & 0.565 \\ 0.530 & 0.506 & 0.499 \\ 0.506 & 0.500 & 0.519 \end{pmatrix}$$

For example, third row corresponding to 3rd input:

$$\mathbf{x} = \begin{pmatrix} 0.78 \\ -0.65 \end{pmatrix}$$

And the corresponding output

$$\mathbf{y} = \begin{pmatrix} 0.530 \\ 0.506 \\ 0.499 \end{pmatrix}$$

2nd neuron output
for 3rd input pattern

Example 1

$$\mathbf{W} = \begin{pmatrix} 0.133 & 0.072 & -0.155 \\ -0.001 & 0.062 & -0.072 \end{pmatrix} \text{ and } \mathbf{B} = \begin{pmatrix} 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \end{pmatrix}.$$

Input as a batch of four patterns:

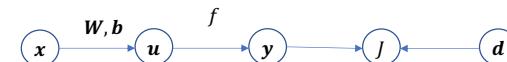
$$\mathbf{X} = \begin{pmatrix} 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix}$$

The synaptic input to the layer:

$$\begin{aligned} \mathbf{U} &= \mathbf{XW} + \mathbf{B} \\ &= \begin{pmatrix} 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix} \begin{pmatrix} 0.133 & 0.072 & -0.155 \\ -0.001 & 0.062 & -0.072 \end{pmatrix} + \begin{pmatrix} 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \end{pmatrix} \\ &= \begin{pmatrix} 0.085 & -0.059 & 0.111 \\ -0.115 & 0.094 & 0.26 \\ 0.121 & 0.024 & -0.005 \\ 0.022 & -0.001 & 0.077 \end{pmatrix} \end{aligned}$$

SGD for single layer

Computational graph processing input (\mathbf{x}, \mathbf{d}):



J denotes the cost function.

Need to compute gradients $\nabla_{\mathbf{W}} J$ and $\nabla_{\mathbf{b}} J$ to learn weight matrix \mathbf{W} and bias vector \mathbf{b} .

SGD for single layer

Consider k th neuron at the layer:

$$u_k = \mathbf{x}^T \mathbf{w}_k + b_k$$

And

$$\frac{\partial u_k}{\partial \mathbf{w}_k} = \mathbf{x}$$

The gradient of the cost with respect to the weight connected to k th neuron:

$$\nabla_{\mathbf{w}_k} J = \frac{\partial J}{\partial u_k} \frac{\partial u_k}{\partial \mathbf{w}_k} = \mathbf{x} \frac{\partial J}{\partial u_k} \quad (\text{A})$$

$$\nabla_{b_k} J = \frac{\partial J}{\partial u_k} \frac{\partial u_k}{\partial b_k} = \frac{\partial J}{\partial u_k} \quad (\text{B})$$

SGD for single layer

Similarly, by substituting $\frac{\partial J}{\partial b_k} = \frac{\partial J}{\partial u_k}$ from (B):

$$\nabla_{\mathbf{b}} J = \begin{pmatrix} \frac{\partial J}{\partial b_1} \\ \frac{\partial J}{\partial b_2} \\ \vdots \\ \frac{\partial J}{\partial b_K} \end{pmatrix} = \begin{pmatrix} \frac{\partial J}{\partial u_1} \\ \frac{\partial J}{\partial u_2} \\ \vdots \\ \frac{\partial J}{\partial u_K} \end{pmatrix} = \nabla_{\mathbf{u}} J \quad (\text{D})$$

$$\nabla_{\mathbf{b}} J = \nabla_{\mathbf{u}} J$$

SGD for single layer

Gradient of J with respect to $\mathbf{W} = (\mathbf{w}_1 \quad \mathbf{w}_2 \quad \cdots \quad \mathbf{w}_K)$:

$$\begin{aligned} \nabla_{\mathbf{W}} J &= (\nabla_{\mathbf{w}_1} J \quad \nabla_{\mathbf{w}_2} J \quad \cdots \quad \nabla_{\mathbf{w}_K} J) \\ &= \left(\mathbf{x} \frac{\partial J}{\partial u_1} \quad \mathbf{x} \frac{\partial J}{\partial u_2} \quad \cdots \quad \mathbf{x} \frac{\partial J}{\partial u_K} \right) \quad \text{From (A)} \\ &= \mathbf{x} \left(\frac{\partial J}{\partial u_1} \quad \frac{\partial J}{\partial u_2} \quad \cdots \quad \frac{\partial J}{\partial u_K} \right) \\ &= \mathbf{x} (\nabla_{\mathbf{u}} J)^T \end{aligned}$$

where

$$\nabla_{\mathbf{u}} J = \frac{\partial J}{\partial \mathbf{u}} = \begin{pmatrix} \frac{\partial J}{\partial u_1} \\ \frac{\partial J}{\partial u_2} \\ \vdots \\ \frac{\partial J}{\partial u_K} \end{pmatrix}$$

$$\text{That is, } \nabla_{\mathbf{W}} J = \mathbf{x} (\nabla_{\mathbf{u}} J)^T \quad (\text{C})$$

SGD for single layer

From (C) and (D),

$$\begin{aligned} \nabla_{\mathbf{W}} J &= \mathbf{x} (\nabla_{\mathbf{u}} J)^T \\ \nabla_{\mathbf{b}} J &= \nabla_{\mathbf{u}} J \end{aligned}$$

That is, by computing gradient $\nabla_{\mathbf{u}} J$ with respect to synaptic input \mathbf{u} , the gradient of cost J with respect to the weights and biases is obtained.

$$\begin{aligned} \mathbf{W} &= \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{u}} J)^T \\ \mathbf{b} &= \mathbf{b} - \alpha \nabla_{\mathbf{u}} J \end{aligned}$$

GD for single layer

Given a set of patterns $\{(\mathbf{x}_p, \mathbf{d}_p)\}_{p=1}^P$ where $\mathbf{x}_p \in \mathbf{R}^n$ and $\mathbf{d}_p \in \mathbf{R}^K$ for regression and $d_p \in \{1, 2, \dots, K\}$ for classification.

The cost J is given by the sum of cost due to individual patterns:

$$J = \sum_{p=1}^P J_p$$

Where Then,

$$\nabla_{\mathbf{W}} J = \sum_{p=1}^P \nabla_{\mathbf{W}} J_p$$

GD for single layer

$$J = \sum_{p=1}^P J_p$$

$$\begin{aligned} \nabla_b J &= \sum_{p=1}^P \nabla_b J_p \\ &= \sum_{p=1}^P \nabla_{\mathbf{u}_p} J_p && \text{Substituting from (D)} \\ &= \sum_{p=1}^P \nabla_{\mathbf{u}_p} J && \text{Since } \nabla_{\mathbf{u}_p} J = \nabla_{\mathbf{u}_p} J_p. \\ &= \nabla_{\mathbf{u}_1} J + \nabla_{\mathbf{u}_2} J + \dots + \nabla_{\mathbf{u}_P} J \\ &= (\nabla_{\mathbf{u}_1} J \quad \nabla_{\mathbf{u}_2} J \quad \dots \quad \nabla_{\mathbf{u}_P} J) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \\ &= (\nabla_{\mathbf{u}} J)^T \mathbf{1}_P \end{aligned} \tag{F}$$

where $\mathbf{1}_P = (1, 1, \dots, 1)^T$ is a vector of P ones.

GD for single layer

Substituting $\nabla_{\mathbf{W}} J_p = \mathbf{x}_p (\nabla_{\mathbf{u}_p} J_p)^T$ from (C) :

$$\begin{aligned} \nabla_{\mathbf{W}} J &= \sum_{p=1}^P \mathbf{x}_p (\nabla_{\mathbf{u}_p} J_p)^T \\ &= \sum_{p=1}^P \mathbf{x}_p (\nabla_{\mathbf{u}_p} J)^T && \text{since } \nabla_{\mathbf{u}_p} J = \nabla_{\mathbf{u}_p} J_p. \\ &= \mathbf{x}_1 (\nabla_{\mathbf{u}_1} J)^T + \mathbf{x}_2 (\nabla_{\mathbf{u}_2} J)^T + \dots + \mathbf{x}_P (\nabla_{\mathbf{u}_P} J)^T \\ &= (\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_P) \begin{pmatrix} (\nabla_{\mathbf{u}_1} J)^T \\ (\nabla_{\mathbf{u}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J)^T \end{pmatrix} \\ &= \mathbf{X}^T \nabla_{\mathbf{u}} J \end{aligned} \tag{E}$$

GD for single layer

From (E) and (F):

$$\begin{aligned} \nabla_{\mathbf{W}} J &= \mathbf{X}^T \nabla_{\mathbf{u}} J \\ \nabla_b J &= (\nabla_{\mathbf{u}} J)^T \mathbf{1}_P \end{aligned}$$

That is, by computing gradient $\nabla_{\mathbf{u}} J$ with respect to synaptic input, the weights and biases can be updated.

$$\begin{aligned} \mathbf{W} &= \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{u}} J \\ \mathbf{b} &= \mathbf{b} - \alpha (\nabla_{\mathbf{u}} J)^T \mathbf{1}_P \end{aligned}$$

Note that $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix}$ and $\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_P^T \end{pmatrix}$

Learning a single layer

SGD for perceptron layer

Learning a layer of neurons	
SGD	$\mathbf{W} = \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{u}} J)^T$ $\mathbf{b} = \mathbf{b} - \alpha \nabla_{\mathbf{u}} J$
GD	$\mathbf{W} = \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{u}} J$ $\mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{u}} J)^T \mathbf{1}_P$

To learn a given layer, we need to compute $\nabla_{\mathbf{u}} J$ for SGD and $\nabla_{\mathbf{u}} J$ for GD.

Those gradients with respect to synaptic inputs are dependent on the types of neurons in the layer.

Given a training pattern (\mathbf{x}, \mathbf{d})

Note $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbf{R}^n$ and $\mathbf{d} = (d_1, d_2, \dots, d_K)^T \in \mathbf{R}^K$.

The square-error cost function:

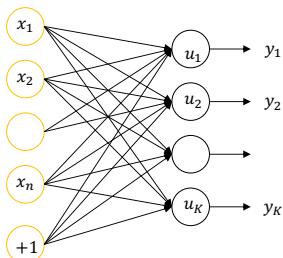
$$J = \frac{1}{2} \sum_{k=1}^K (d_k - y_k)^2$$

where $y_k = f(u_k) = \frac{1}{1+e^{-u_k}}$ and $u_k = \mathbf{x}^T \mathbf{w}_k + b_k$.

Gradient of J with respect to u_k :

$$\frac{\partial J}{\partial u_k} = \frac{\partial J}{\partial y_k} \frac{\partial y_k}{\partial u_k} = -(d_k - y_k) \frac{\partial y_k}{\partial u_k} = -(d_k - y_k) f'(u_k) \quad (G)$$

Perceptron layer



$$y_k = f(u_k) = \frac{1}{1 + e^{-u_k}}$$

SGD for perceptron layer

Substituting $\nabla_{u_k} J = \frac{\partial J}{\partial u_k} = -(d_k - y_k) f'(u_k)$ from (G):

$$\nabla_{\mathbf{u}} J = \begin{pmatrix} \nabla_{u_1} J \\ \nabla_{u_2} J \\ \vdots \\ \nabla_{u_K} J \end{pmatrix} = - \begin{pmatrix} (d_1 - y_1) f'(u_1) \\ (d_2 - y_2) f'(u_2) \\ \vdots \\ (d_K - y_K) f'(u_K) \end{pmatrix} = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \quad (H)$$

and ‘ \cdot ’ denotes element-wise multiplication.

$$\nabla_{\mathbf{u}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$$

A layer of perceptrons performs **multidimensional non-linear regression** and learns a multidimensional non-linear mapping:

$$\phi: \mathbf{R}^n \rightarrow \mathbf{R}^K$$

SGD for perceptron layer

Given a training dataset $\{(\mathbf{x}, \mathbf{d})\}$

Set learning parameter α

Initialize \mathbf{W} and \mathbf{b}

Repeat until convergence:

For every pattern (\mathbf{x}, \mathbf{d}) :

$$\mathbf{u} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

$$\mathbf{y} = f(\mathbf{u}) = \frac{1}{1+e^{-\mathbf{u}}}$$

$$\nabla_{\mathbf{u}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$$

$$\mathbf{W} = \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{u}} J)^T$$

$$\mathbf{b} = \mathbf{b} - \alpha \nabla_{\mathbf{u}} J$$

GD for perceptron layer

$$\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_P^T \end{pmatrix} \rightarrow \nabla_{\mathbf{U}} J = \begin{pmatrix} (\nabla_{\mathbf{u}_1} J)^T \\ (\nabla_{\mathbf{u}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J)^T \end{pmatrix} = \begin{pmatrix} (\nabla_{\mathbf{u}_1} J_1)^T \\ (\nabla_{\mathbf{u}_2} J_2)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J_P)^T \end{pmatrix}$$

From (H), substitute $\nabla_{\mathbf{U}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$:

$$\nabla_{\mathbf{U}} J = - \begin{pmatrix} ((\mathbf{d}_1 - \mathbf{y}_1) \cdot f'(\mathbf{u}_1))^T \\ ((\mathbf{d}_2 - \mathbf{y}_2) \cdot f'(\mathbf{u}_2))^T \\ \vdots \\ ((\mathbf{d}_P - \mathbf{y}_P) \cdot f'(\mathbf{u}_P))^T \end{pmatrix} = - \begin{pmatrix} (\mathbf{d}_1^T - \mathbf{y}_1^T) \cdot f'(\mathbf{u}_1^T) \\ (\mathbf{d}_2^T - \mathbf{y}_2^T) \cdot f'(\mathbf{u}_2^T) \\ \vdots \\ (\mathbf{d}_P^T - \mathbf{y}_P^T) \cdot f'(\mathbf{u}_P^T) \end{pmatrix}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$$

$$\text{where } \mathbf{D} = \begin{pmatrix} \mathbf{d}_1^T \\ \mathbf{d}_2^T \\ \vdots \\ \mathbf{d}_P^T \end{pmatrix}, \mathbf{Y} = \begin{pmatrix} \mathbf{y}_1^T \\ \mathbf{y}_2^T \\ \vdots \\ \mathbf{y}_P^T \end{pmatrix}, \text{ and } \mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_P^T \end{pmatrix}$$

23

GD for perceptron layer

Given a training dataset $\{(\mathbf{x}_p, \mathbf{d}_p)\}_{p=1}^P$

Note $\mathbf{x}_p = (x_{p1}, x_{p2}, \dots, x_{pn})^T \in \mathbf{R}^n$ and $\mathbf{d}_p = (d_{p1}, d_{p2}, \dots, d_{pK})^T \in \mathbf{R}^K$.

The cost function J is given by the sum of square errors (s.s.e.):

$$J = \frac{1}{2} \sum_{p=1}^P \sum_{k=1}^K (d_{pk} - y_{pk})^2$$

J can be written as the sum of cost due to individual patterns:

$$J = \sum_{p=1}^P J_p$$

where $J_p = \frac{1}{2} \sum_{k=1}^K (d_{pk} - y_{pk})^2$ is the square error for the p th pattern.

GD for perceptron layer

Given a training dataset (\mathbf{X}, \mathbf{D})

Set learning parameter α

Initialize \mathbf{W} and \mathbf{b}

Repeat until convergence:

$$\mathbf{U} = \mathbf{XW} + \mathbf{B}$$

$$\mathbf{Y} = f(\mathbf{U}) = \frac{1}{1+e^{-\mathbf{U}}}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$$

$$\mathbf{W} = \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

26

Learning a perceptron layer

Softmax layer

GD	SGD
(X, D)	(x, d)
$U = XW + B$	$u = W^T x + b$
$Y = f(U)$	$y = f(u)$
$\nabla_u J = -(\mathbf{d} - \mathbf{y}) \cdot f'(U)$	$\nabla_u J = -(\mathbf{d} - \mathbf{y}) \cdot f'(u)$
$W = W - \alpha X^T \nabla_u J$	$W = W - \alpha x (\nabla_u J)^T$
$b = b - \alpha (\nabla_u J)^T \mathbf{1}_P$	$b = b - \alpha \nabla_u J$

The K neurons in the softmax layer performs K class classification and represent K classes.

The activation of each neuron k estimates the probability $P(y = k|x)$ that the input x belongs to the class k :

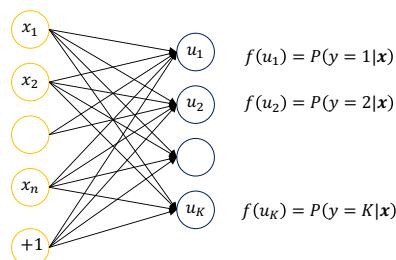
$$P(y = k|x) = f(u_k) = \frac{e^{u_k}}{\sum_{k'=1}^K e^{u_{k'}}}$$

where $u_k = \mathbf{w}_k^T x + b_k$, and \mathbf{w}_k is weight vector and b_k is bias of neuron k .

The above activation function f is known as **softmax activation function**.

Softmax layer

Softmax layer is the extension of logistic regression to **multiclass classification** problem, which is also known as *multinomial logistic regression*.



Each neuron in the softmax layer corresponds to one class label. The activation of a neuron gives the probability of the input belonging to that class label. Output is the class label having the maximum probability.

Softmax layer

The output y denotes the class label of the input pattern, which is given by

$$y = \operatorname{argmax}_k P(y = k|x) = \operatorname{argmax}_k f(u_k)$$

That is, the class label is assigned to the class with the maximum activation.

SGD for softmax layer

Given a training pattern (\mathbf{x}, d) where $\mathbf{x} \in \mathbb{R}^n$ and $d \in \{1, 2, \dots, K\}$.

The cost function for learning is by the *multiclass cross-entropy*:

$$J = - \sum_{k=1}^K 1(d=k) \log(f(u_k))$$

where u_k is the synaptic input to the k the neuron.

The cost function can also be written as

$$J = -\log(f(u_d))$$

where d is the target label of input \mathbf{x} .

Note that the logarithm here is natural: $\log = \log_e$

31

SGD for softmax layer

If $k = d$:

$$\begin{aligned} \frac{\partial f(u_d)}{\partial u_k} &= \frac{\partial}{\partial u_k} \left(\frac{e^{u_d}}{\sum_{k'=1}^K e^{u_{k'}}} \right) \\ &= \frac{(\sum_{k'=1}^K e^{u_{k'}}) e^{u_d} - e^{u_d} e^{u_d}}{\left(\sum_{k'=1}^K e^{u_{k'}} \right)^2} \\ &= \frac{e^{u_d}}{\sum_{k'=1}^K e^{u_{k'}}} \left(1 - \frac{e^{u_d}}{\sum_{k'=1}^K e^{u_{k'}}} \right) \\ &= f(u_d)(1 - f(u_d)) \\ &= f(u_d)(1(k=d) - f(u_d)) \end{aligned}$$

$$\frac{\partial (\sum_{k'=1}^K e^{u_{k'}})}{\partial u_k} = e^{u_k}$$

If $k \neq d$:

$$\begin{aligned} \frac{\partial f(u_d)}{\partial u_k} &= \frac{\partial}{\partial u_k} \left(\frac{e^{u_d}}{\sum_{k'=1}^K e^{u_{k'}}} \right) \\ &= -\frac{e^{u_d} e^{u_k}}{\left(\sum_{k'=1}^K e^{u_{k'}} \right)^2} \\ &= -f(u_d)f(u_k) \\ &= f(u_d)(1(k=d) - f(u_k)) \end{aligned}$$

$$1(k=d) = 0$$

SGD for softmax layer

$$J = -\log(f(u_d))$$

The gradient with respect to u_k is given by

$$\frac{\partial J}{\partial u_k} = -\frac{1}{f(u_d)} \frac{\partial f(u_d)}{\partial u_k} \quad (I)$$

where

$$\frac{\partial f(u_d)}{\partial u_k} = \frac{\partial}{\partial u_k} \left(\frac{e^{u_d}}{\sum_{k'=1}^K e^{u_{k'}}} \right)$$

The above differentiation need to be considered separately for $k = d$ and for $k \neq d$.

SGD for softmax layer

$$\frac{\partial f(u_d)}{\partial u_k} = f(u_d)(1(d=k) - f(u_k))$$

Substituting in (I):

$$\nabla_{u_k} J = \frac{\partial J}{\partial u_k} = -\frac{1}{f(u_d)} \frac{\partial f(u_d)}{\partial u_k} = -(1(d=k) - f(u_k))$$

Gradient J with respect to \mathbf{u} :

$$\nabla_{\mathbf{u}} J = \begin{pmatrix} \nabla_{u_1} J \\ \nabla_{u_2} J \\ \vdots \\ \nabla_{u_K} J \end{pmatrix} = -\begin{pmatrix} 1(d=1) - f(u_1) \\ 1(d=2) - f(u_2) \\ \vdots \\ 1(d=K) - f(u_K) \end{pmatrix} = -(1(\mathbf{k}=d) - f(\mathbf{u})) \quad (J)$$

where $\mathbf{k} = (1 \ 2 \ \dots \ K)^T$

SGD for softmax layer

For a softmax layer:

$$\nabla_{\mathbf{u}} J = -(1(\mathbf{k} = \mathbf{d}) - f(\mathbf{u}))$$

where:

$$1(\mathbf{k} = d) = \begin{pmatrix} 1(d=1) \\ 1(d=2) \\ \vdots \\ 1(d=K) \end{pmatrix} \text{ and } f(\mathbf{u}) = \begin{pmatrix} f(u_1) \\ f(u_2) \\ \vdots \\ f(u_K) \end{pmatrix}$$

Note that $1(\mathbf{k} = d)$ is a one-hot vector where the element corresponding to the target label d is ‘1’ and elsewhere is ‘0’.

GD for softmax layer

J can be written as the sum of cost due to individual patterns:

$$J = \sum_{p=1}^P J_p$$

where $J_p = -\log(f(u_{pd_p}))$ is the cross-entropy for the p th pattern.

GD for softmax layer

Given a set of patterns $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$ where $\mathbf{x}_p \in \mathbf{R}^n$ and $d_p \in \{1, 2, \dots, K\}$.

The cost function of the *softmax layer* is given by the *multiclass cross-entropy*:

$$J = -\sum_{p=1}^P \left(\sum_{k=1}^K 1(d_p = k) \log(f(u_{pk})) \right)$$

where u_{pk} is the synaptic input to the k the neuron for input \mathbf{x}_p .

The cost function J can also be written as

$$J = -\sum_{p=1}^P \log(f(u_{pd_p}))$$

GD for softmax layer

$$\nabla_{\mathbf{u}} J = \begin{pmatrix} (\nabla_{\mathbf{u}_1} J)^T \\ (\nabla_{\mathbf{u}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J)^T \end{pmatrix} = \begin{pmatrix} (\nabla_{\mathbf{u}_1} J_1)^T \\ (\nabla_{\mathbf{u}_2} J_2)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J_P)^T \end{pmatrix}$$

Substituting $\nabla_{\mathbf{u}} J = -(1(\mathbf{k} = \mathbf{d}) - f(\mathbf{u}))$ from (J):

$$\nabla_{\mathbf{u}} J = - \begin{pmatrix} (1(\mathbf{k} = d_1) - f(\mathbf{u}_1))^T \\ (1(\mathbf{k} = d_2) - f(\mathbf{u}_2))^T \\ \vdots \\ (1(\mathbf{k} = d_P) - f(\mathbf{u}_P))^T \end{pmatrix}$$

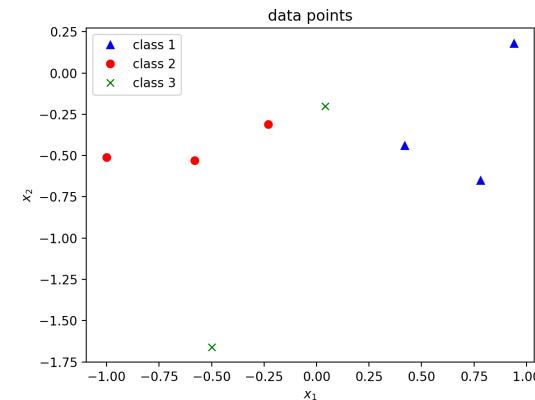
$$\nabla_{\mathbf{u}} J = -(\mathbf{K} - f(\mathbf{U}))$$

where $\mathbf{K} = \begin{pmatrix} 1(\mathbf{k} = d_1)^T \\ 1(\mathbf{k} = d_2)^T \\ \vdots \\ 1(\mathbf{k} = d_P)^T \end{pmatrix}$ is a matrix with every row is a one-hot vector.

Learning a softmax layer

GD	SGD
(X, D)	(x, d)
$U = XW + B$	$u = W^T x + b$
$f(U) = \frac{e^U}{\sum_{k'=1}^K e^{U_{k'}}}$	$f(u) = \frac{e^{u_k}}{\sum_{k'=1}^K e^{u_{k'}}}$
$y = \underset{k}{\operatorname{argmax}} f(U)$	$y = \underset{k}{\operatorname{argmax}} f(u)$
$\nabla_U J = -(K - f(U))$	$\nabla_u J = -(1(k=d) - f(u))$
$W = W - \alpha X^T \nabla_U J$	$W = W - \alpha x (\nabla_u J)^T$
$b = b - \alpha (\nabla_u J)^T \mathbf{1}_P$	$b = b - \alpha \nabla_u J$

Example 2



Example 2: GD of a softmax layer

Train a softmax regression layer of neurons to perform the following classification:

$$\begin{aligned}(0.94 & \quad 0.18) \rightarrow \text{class } A \\ (-0.58 & \quad -0.53) \rightarrow \text{class } B \\ (-0.23 & \quad -0.31) \rightarrow \text{class } B \\ (0.42 & \quad -0.44) \rightarrow \text{class } A \\ (0.5 & \quad -1.66) \rightarrow \text{class } C \\ (-1.0 & \quad -0.51) \rightarrow \text{class } B \\ (0.78 & \quad -0.65) \rightarrow \text{class } A \\ (0.04 & \quad -0.20) \rightarrow \text{class } C\end{aligned}$$

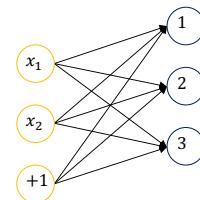
Use a learning factor $\alpha = 0.05$.

Example 2

$$\text{Let } y = \begin{cases} 1, \text{for class } A \\ 2, \text{for class } B \\ 3, \text{for class } C \end{cases}$$

$$X = \begin{pmatrix} 0.94 & 0.18 \\ -0.58 & -0.53 \\ -0.23 & -0.31 \\ 0.42 & -0.44 \\ 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix}, d = \begin{pmatrix} 1 \\ 2 \\ 2 \\ 1 \\ 3 \\ 2 \\ 1 \\ 3 \end{pmatrix}$$

$$K = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$



Example 2

$$\text{Initialize } \mathbf{W} = \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.23 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix},$$

$$\text{Then, } \mathbf{B} = \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{pmatrix}$$

1st epoch starts ...

$$\mathbf{U} = \mathbf{XW} + \mathbf{B} = \begin{pmatrix} 0.94 & 0.18 \\ -0.58 & -0.53 \\ -0.23 & -0.31 \\ 0.42 & -0.44 \\ 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix} \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.23 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{pmatrix}$$

Example 2

$$\mathbf{y} = \underset{k}{\operatorname{argmax}} \{f(\mathbf{U})\} = \underset{k}{\operatorname{argmax}} \left\{ \begin{pmatrix} 0.44 & 0.21 & 0.35 \\ 0.24 & 0.42 & 0.34 \\ 0.29 & 0.37 & 0.35 \\ 0.33 & 0.27 & 0.40 \\ 0.23 & 0.24 & 0.52 \\ 0.20 & 0.49 & 0.31 \\ 0.34 & 0.22 & 0.43 \\ 0.32 & 0.33 & 0.35 \end{pmatrix} \right\} = \begin{pmatrix} 1 \\ 2 \\ 2 \\ 3 \\ 3 \\ 2 \\ 3 \\ 3 \end{pmatrix}$$

$$\text{Errors} = \sum_{p=1}^8 1(d_p \neq y_p) = 2$$

$$\begin{aligned} \text{Entropy, } J &= -\sum_{p=1}^8 \log(f(u_{pd_p})) \\ &= -\log(0.44) - \log(0.42) - \dots - \log(0.35) \\ &= 7.26 \end{aligned}$$

$$\mathbf{d} = \begin{pmatrix} 1 \\ 2 \\ 2 \\ 1 \\ 3 \\ 2 \\ 1 \\ 3 \end{pmatrix}$$

Example 2

$$\mathbf{U} = \begin{pmatrix} 0.86 & 0.11 & 0.64 \\ -0.84 & -0.28 & -0.49 \\ -0.41 & -0.16 & -0.22 \\ -0.01 & -0.21 & 0.17 \\ -0.86 & -0.82 & -0.06 \\ -1.15 & -0.27 & -0.75 \\ 0.11 & -0.31 & 0.35 \\ -0.12 & -0.10 & -0.02 \end{pmatrix}$$

$$f(u_{12}) = \frac{e^{0.11}}{e^{0.86} + e^{0.11} + e^{0.64}}$$

$$f(\mathbf{U}) = \frac{e^{(\mathbf{U})}}{\sum_{k=1}^K e^{(\mathbf{U})}} = \begin{pmatrix} 0.44 & 0.21 & 0.35 \\ 0.24 & 0.42 & 0.34 \\ 0.29 & 0.37 & 0.35 \\ 0.33 & 0.27 & 0.40 \\ 0.23 & 0.24 & 0.52 \\ 0.20 & 0.49 & 0.31 \\ 0.34 & 0.22 & 0.43 \\ 0.32 & 0.33 & 0.35 \end{pmatrix}$$

Example 2

$$\begin{aligned} \nabla_{\mathbf{U}} J &= -(\mathbf{K} - f(\mathbf{U})) \\ &= -\left(\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 0.44 & 0.21 & 0.35 \\ 0.24 & 0.42 & 0.34 \\ 0.29 & 0.37 & 0.35 \\ 0.33 & 0.27 & 0.40 \\ 0.23 & 0.24 & 0.52 \\ 0.20 & 0.49 & 0.31 \\ 0.34 & 0.22 & 0.43 \\ 0.32 & 0.33 & 0.35 \end{pmatrix} \right) \\ &= \begin{pmatrix} -0.56 & 0.21 & 0.35 \\ 0.24 & -0.58 & 0.34 \\ 0.29 & -0.63 & 0.35 \\ -0.67 & 0.27 & 0.40 \\ 0.23 & 0.24 & -0.48 \\ 0.20 & -0.51 & 0.31 \\ -0.65 & 0.22 & 0.43 \\ 0.32 & 0.33 & -0.65 \end{pmatrix} \end{aligned}$$

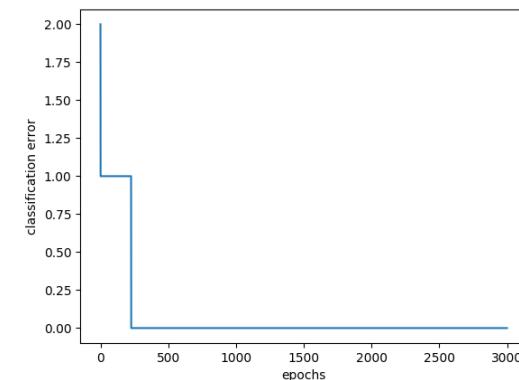
$$\mathbf{W} = \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J$$

$$= \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.23 \end{pmatrix} - 0.05 \begin{pmatrix} 0.94 & 0.18 \\ -0.58 & -0.53 \\ -0.23 & -0.31 \\ 0.42 & -0.44 \\ 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix}^T \begin{pmatrix} -0.56 & 0.21 & 0.35 \\ 0.24 & -0.58 & 0.34 \\ 0.29 & -0.63 & 0.35 \\ -0.67 & 0.27 & 0.40 \\ 0.23 & 0.24 & -0.48 \\ 0.20 & -0.51 & 0.31 \\ -0.65 & 0.22 & 0.43 \\ 0.32 & 0.33 & -0.65 \end{pmatrix}$$

$$= \begin{pmatrix} 0.85 & -0.06 & 0.63 \\ 0.76 & 0.50 & 0.22 \end{pmatrix}$$

$$\mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix} - 0.05 \begin{pmatrix} -0.56 & 0.21 & 0.35 \\ 0.24 & -0.58 & 0.34 \\ 0.29 & -0.63 & 0.35 \\ -0.67 & 0.27 & 0.40 \\ 0.23 & 0.24 & -0.48 \\ 0.20 & -0.51 & 0.31 \\ -0.65 & 0.22 & 0.43 \\ 0.32 & 0.33 & -0.65 \end{pmatrix}^T \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0.03 \\ 0.02 \\ -0.05 \end{pmatrix}$$

Example 2



49

Example 2

Example 2

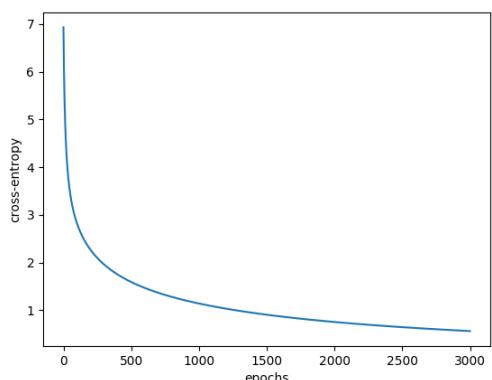
At convergence at 3000 iterations:

$$\mathbf{W} = \begin{pmatrix} 14.22 & -13.04 & 0.00 \\ 4.47 & -2.05 & -0.95 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} -0.53 \\ -0.47 \\ 1.00 \end{pmatrix}$$

Entropy = 0.562

Errors = 0



48

Initialization of weights

Random initialization is inefficient

At initialization, it is desirable that weights are small and near zero

- to operate in the linear region of the activation function
- to preserve the variance of activations and gradients.

Two methods:

- Using a uniform distribution within specified limits
- Using a truncated normal distribution

Initialization from a truncated normal distribution

$$w \sim \text{truncated_normal} \left[\text{mean} = 0, \text{std} = \frac{\text{gain}}{\sqrt{n_{in}}} \right]$$

In the truncated normal, the samples that are two s.d. away from the center are discarded and resampled again.

This is also known as **Kaiming normal initialization**.

Initialization from a uniform distribution (Xavier/Glorot uniform Initialization)

Uniformly draws weight samples $w \sim U(-a, +a)$:

where

$$a = \text{gain} \times \sqrt{\frac{6}{n_{in} + n_{out}}}$$

n_{in} is the number of input nodes and n_{out} is the number of neurons in the layer.

activation	linear	sigmoid	Tanh	ReLU	Leaky ReLU
gain	1	1	5/3	$\sqrt{2}$	$\sqrt{\frac{1}{1 + \text{slope}^2}}$

Iris dataset

Iris dataset:

<https://archive.ics.uci.edu/ml/datasets/Iris>

Three classes of iris flower:



Setosa

Versicolour

Virginica

Four features:

Sepal length, sepal width, petal length, petal width

Iris dataset

150 data points, 50 for each class

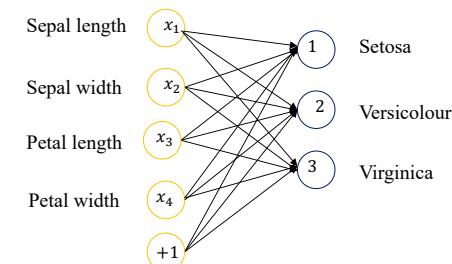
Features:

```
[ -7.4333333e-01  4.4600000e-01 -2.35866667e+00 -9.98666667e-01]
[ -9.4333333e-01 -5.4000000e-02 -2.35866667e+00 -9.98666667e-01]
[ -1.1433333e+00  1.4600000e-01 -2.45866667e+00 -9.98666667e-01]
```

Labels:

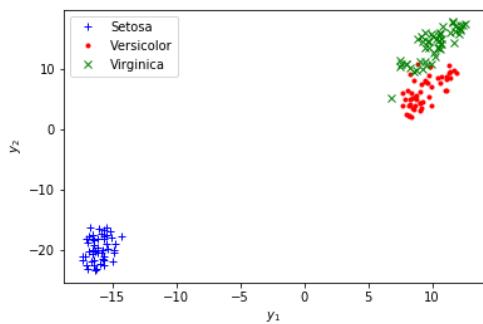
```
[0 0 0 0 0 0 0 0 0 ..1 1 1 1 1 1 1 1 1..... 2 2 2 2 2 2 2 2]
```

Example 3: Softmax classification of iris data

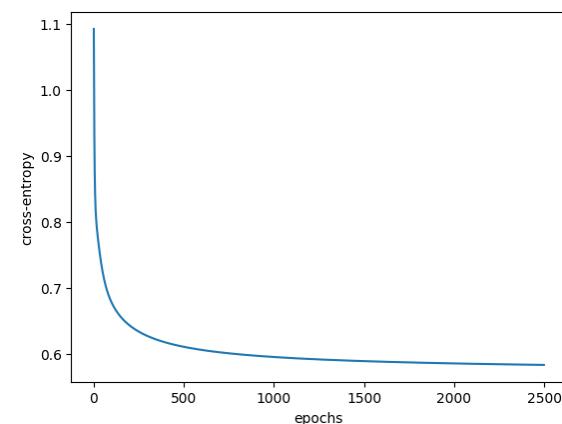


Iris data

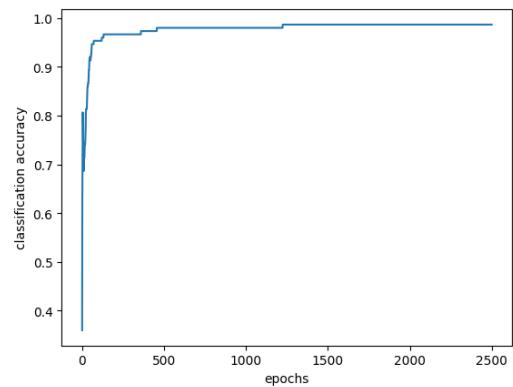
Display of data points after dimensionality reduction (from Four dimensions to Two) by t-SNE.



Example 3



Example 3



Final classification error = 5

Summary: GD for layers

$$\begin{aligned} & (\mathbf{X}, \mathbf{D}) \\ & \mathbf{U} = \mathbf{XW} + \mathbf{B} \\ & \mathbf{W} = \mathbf{W} - \alpha \mathbf{X}^T (\nabla_{\mathbf{U}} J) \\ & \mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P \end{aligned}$$

layer	$f(\mathbf{U}), \mathbf{Y}$	$\nabla_{\mathbf{U}} J$
Linear neuron layer	$\mathbf{Y} = f(\mathbf{U}) = \mathbf{U}$	$-(\mathbf{D} - \mathbf{Y})$
Perceptron layer	$\mathbf{Y} = f(\mathbf{U}) = \frac{1}{1 + e^{-\mathbf{U}}}$	$-(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$
Softmax layer	$f(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{U_k}}$ $\mathbf{y} = \underset{k}{\operatorname{argmax}} f(\mathbf{U})$	$-(\mathbf{K} - f(\mathbf{U}))$

Revision: neurons and layers

Classification	Logistic neurons	
Two-class	Logistic regression neuron	
Multiclass	Softmax layer	

Regression	Linear	Non-linear
One dimensional	Linear neuron	Perceptron
Multi-dimensional	Linear neuron layer	Perceptron layer

Summary: SGD for layers

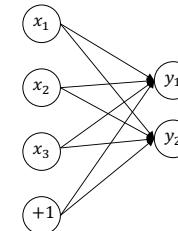
$$\begin{aligned} & (\mathbf{x}, \mathbf{d}) \\ & \mathbf{u} = \mathbf{W}^T \mathbf{x} + \mathbf{b} \\ & \mathbf{W} = \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{u}} J)^T \\ & \mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{u}} J) \end{aligned}$$

layer	$f(\mathbf{u}), \mathbf{y}$	$\nabla_{\mathbf{u}} J$
Linear neuron layer	$\mathbf{y} = f(\mathbf{u}) = \mathbf{u}$	$-(\mathbf{d} - \mathbf{y})$
Perceptron layer	$\mathbf{y} = f(\mathbf{u}) = \frac{1}{1 + e^{-\mathbf{u}}}$	$-(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$
Softmax layer	$f(\mathbf{u}) = \frac{e^{\mathbf{u}}}{\sum_{k'=1}^K e^{U_{k'}}}$ $\mathbf{y} = \underset{k}{\operatorname{argmax}} f(\mathbf{u})$	$-(1(k=d) - f(\mathbf{u}))$

Example 4

$$X = \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.22 \\ 0.20 & 0.76 & 0.17 \\ 0.09 & 0.69 & 0.95 \\ 0.00 & 0.51 & 0.81 \\ 0.61 & 0.72 & 0.29 \\ 0.92 & 0.71 & 0.54 \\ 0.14 & 0.37 & 0.67 \end{pmatrix} \text{ and } D = \begin{pmatrix} 0.37 & 0.47 \\ 0.36 & 0.38 \\ 0.35 & 0.25 \\ 0.48 & 0.42 \\ 0.36 & 0.29 \\ 0.44 & 0.52 \\ 0.60 & 0.52 \\ 0.28 & 0.37 \end{pmatrix}$$

Thank you.



Output $y_1, y_2 \in [0, 1]$

So, activation function for both neurons:

$$f(u) = \frac{1}{1 + e^{-u}}$$

$$f'(u) = y(1 - y)$$

Learning factor $\alpha = 0.1$.

Weights and biases are initialized:

$$W = \begin{pmatrix} 0.03 & 0.04 \\ 0.01 & 0.04 \\ 0.02 & 0.04 \end{pmatrix} \text{ and } b = \begin{pmatrix} 0.0 \\ 0.0 \end{pmatrix}$$

Example 4: GD of a perceptron layer

Design a perceptron layer to perform the following mapping using GD learning:

$x = (x_1, x_2, x_3)$	$d = (d_1, d_2)$
(0.77, 0.02, 0.63)	(0.37, 0.47)
(0.75, 0.50, 0.22)	(0.36, 0.38)
(0.20, 0.76, 0.17)	(0.35, 0.25)
(0.09, 0.69, 0.95)	(0.48, 0.42)
(0.00, 0.51, 0.81)	(0.36, 0.29)
(0.61, 0.72, 0.29)	(0.44, 0.52)
(0.92, 0.71, 0.54)	(0.60, 0.52)
(0.14, 0.37, 0.67)	(0.28, 0.37)

Use $\alpha = 0.1$.

Example 4

Epoch 1:

$$U = XW + B$$

$$U = \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.22 \\ 0.20 & 0.76 & 0.17 \\ 0.09 & 0.69 & 0.95 \\ 0.00 & 0.51 & 0.81 \\ 0.61 & 0.72 & 0.29 \\ 0.92 & 0.71 & 0.54 \\ 0.14 & 0.37 & 0.67 \end{pmatrix} \begin{pmatrix} 0.03 & 0.04 \\ 0.01 & 0.04 \\ 0.02 & 0.04 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} 0.03 & 0.06 \\ 0.03 & 0.06 \\ 0.02 & 0.05 \\ 0.03 & 0.07 \\ 0.02 & 0.05 \\ 0.03 & 0.07 \\ 0.04 & 0.09 \\ 0.02 & 0.05 \end{pmatrix}$$

$$Y = f(U) = \frac{1}{1+e^{-U}} = \begin{pmatrix} 0.51 & 0.51 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \\ 0.51 & 0.52 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \end{pmatrix}$$

$$\text{Mean square error} = \frac{1}{8} \sum_{p=1}^8 \sum_{k=1}^2 (d_{pk} - y_{pk})^2 = \frac{1}{8} \sum_{p=1}^8 (d_{p1} - y_{p1})^2 + (d_{p2} - y_{p2})^2 = 0.04$$

Example 4

$$f'(\mathbf{U}) = \mathbf{Y} \cdot (1 - \mathbf{Y}) = \begin{pmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \end{pmatrix}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U}) \\ = - \left(\begin{pmatrix} 0.37 & 0.47 \\ 0.36 & 0.38 \\ 0.35 & 0.25 \\ 0.48 & 0.42 \\ 0.36 & 0.29 \\ 0.44 & 0.52 \\ 0.60 & 0.52 \\ 0.28 & 0.37 \end{pmatrix} - \begin{pmatrix} 0.51 & 0.51 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \\ 0.51 & 0.52 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \end{pmatrix} \right) \cdot \begin{pmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \end{pmatrix} = \begin{pmatrix} 0.04 & 0.01 \\ 0.04 & 0.03 \\ 0.04 & 0.07 \\ 0.01 & 0.03 \\ 0.04 & 0.06 \\ 0.02 & 0.00 \\ -0.02 & 0.00 \\ 0.06 & 0.04 \end{pmatrix}$$

Example 4

Epoch	\mathbf{Y}	mse	\mathbf{W}	\mathbf{b}
2	$\begin{pmatrix} 0.50 & 0.51 \\ 0.50 & 0.51 \\ 0.50 & 0.50 \\ 0.50 & 0.51 \\ 0.50 & 0.50 \\ 0.50 & 0.51 \\ 0.50 & 0.51 \\ 0.50 & 0.50 \end{pmatrix}$	0.036	$\begin{pmatrix} 0.02 & 0.03 \\ -0.01 & 0.02 \\ 0.00 & 0.01 \end{pmatrix}$	(-0.04)
10000	$\begin{pmatrix} 0.34 & 0.46 \\ 0.39 & 0.42 \\ 0.32 & 0.29 \\ 0.48 & 0.40 \\ 0.36 & 0.34 \\ 0.45 & 0.42 \\ 0.59 & 0.55 \\ 0.31 & 0.33 \end{pmatrix}$	0.003	$\begin{pmatrix} 1.08 & 1.14 \\ 1.42 & 0.40 \\ 1.14 & 0.84 \end{pmatrix}$	(-2.24)

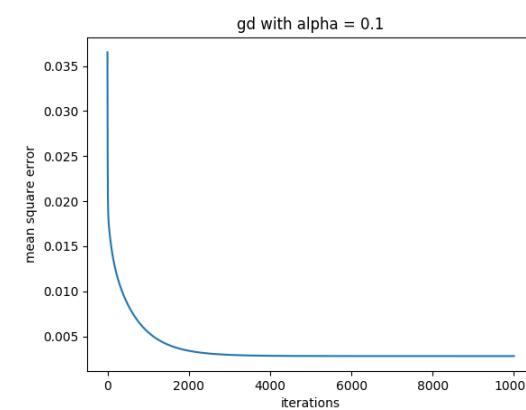
69

Example 4

$$\mathbf{W} = \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J \\ = \begin{pmatrix} 0.03 & 0.04 \\ 0.01 & 0.04 \\ 0.02 & 0.04 \end{pmatrix} - 0.1 \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.22 \\ 0.20 & 0.76 & 0.17 \\ 0.09 & 0.69 & 0.95 \\ 0.00 & 0.51 & 0.81 \\ 0.61 & 0.72 & 0.29 \\ 0.92 & 0.71 & 0.54 \\ 0.14 & 0.37 & 0.67 \end{pmatrix}^T \begin{pmatrix} 0.04 & 0.01 \\ 0.04 & 0.03 \\ 0.04 & 0.07 \\ 0.01 & 0.03 \\ 0.04 & 0.06 \\ 0.02 & 0.00 \\ -0.02 & 0.00 \\ 0.06 & 0.04 \end{pmatrix} = \begin{pmatrix} 0.02 & 0.04 \\ 0.00 & 0.03 \\ 0.01 & 0.03 \end{pmatrix}$$

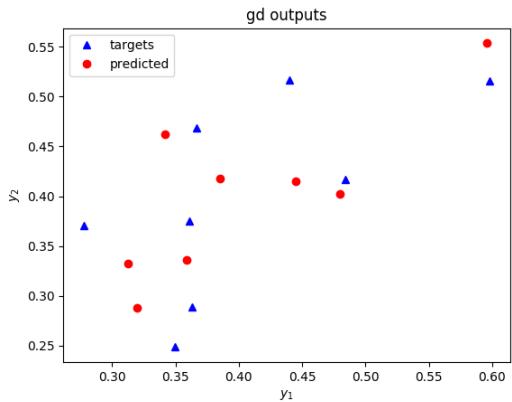
$$\mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P = \begin{pmatrix} 0.0 \\ 0.0 \end{pmatrix} - 0.1 \begin{pmatrix} 0.03 & 0.01 \\ 0.03 & 0.03 \\ 0.04 & 0.06 \\ 0.00 & 0.02 \\ 0.03 & 0.05 \\ 0.01 & 0.00 \\ -0.02 & 0.00 \\ 0.05 & 0.03 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} -0.02 \\ -0.02 \end{pmatrix}$$

Example 4



70

Example 4



71

Chain rule of differentiation

Let x, y , and $J \in \mathbf{R}$ be one-dimensional variables and

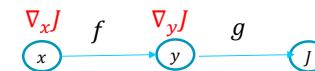
$$J = g(y)$$

$$y = f(x)$$

Chain rule of differentiation states that:

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial x}$$

$$\nabla_x J = \left(\frac{\partial y}{\partial x} \right) \nabla_y J$$



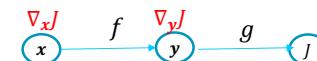
Note the transfer of gradient of J from y to x .

Chapter 4

Deep neural networks

Neural networks and deep learning

Chain rule in multidimensions



$\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbf{R}^n$, $\mathbf{y} = (y_1, y_2, \dots, y_K) \in \mathbf{R}^K$, $J \in \mathbf{R}$, and

$$\mathbf{y} = f(\mathbf{x})$$

$$J = g(\mathbf{y})$$

Then, the chain rule of differentiation states that:

$$\nabla_{\mathbf{x}} J = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} J$$

The matrix $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ is known as the **Jacobian** of the function f where $\mathbf{y} = f(\mathbf{x})$.

Chain rule in multidimensions

$$\nabla_x J = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_y J$$

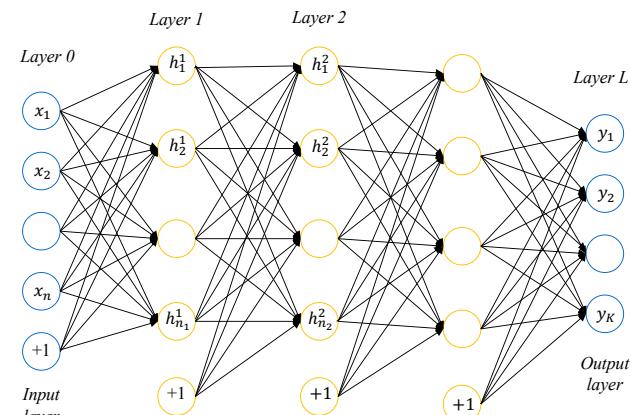
where

$$\nabla_x J = \begin{pmatrix} \frac{\partial J}{\partial x_1} \\ \frac{\partial J}{\partial x_2} \\ \vdots \\ \frac{\partial J}{\partial x_n} \end{pmatrix} \text{ and } \nabla_y J = \frac{\partial J}{\partial \mathbf{y}} = \begin{pmatrix} \frac{\partial J}{\partial y_1} \\ \frac{\partial J}{\partial y_2} \\ \vdots \\ \frac{\partial J}{\partial y_K} \end{pmatrix},$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_K}{\partial x_1} & \frac{\partial y_K}{\partial x_2} & \dots & \frac{\partial y_K}{\partial x_n} \end{pmatrix}$$

Note that differentiation of a scalar by a vector results in a vector and differentiation of a vector by a vector results in a matrix.

Deep neural networks (DNN): Also known as feedforward networks (FFN)



L layer DNN

Example 1: find Jacobian of a function

Let $\mathbf{x} = (x_1, x_2, x_3)$, $\mathbf{y} = (y_1, y_2)$, and $\mathbf{y} = f(\mathbf{x})$ where f is given by

$$\begin{aligned} y_1 &= 5 - 2x_1 + 3x_3 \\ y_2 &= x_1 + 5x_2^2 + x_3^3 - 1 \end{aligned}$$

Find the Jacobian of f .

$$\begin{aligned} \frac{\partial y_1}{\partial x_1} &= -2, & \frac{\partial y_1}{\partial x_2} &= 0, & \frac{\partial y_1}{\partial x_3} &= 3 \\ \frac{\partial y_2}{\partial x_1} &= 1, & \frac{\partial y_2}{\partial x_2} &= 10x_2, & \frac{\partial y_2}{\partial x_3} &= 3x_3^2 \end{aligned}$$

Jacobian:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \end{pmatrix} = \begin{pmatrix} -2 & 0 & 3 \\ 1 & 10x_2 & 3x_3^2 \end{pmatrix}$$

Feedforward Networks (FFN)

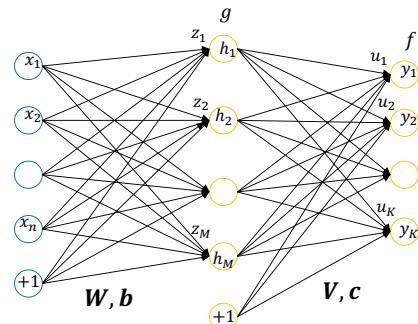
Feedforward networks (FFN) consists of several layers of neurons where activations propagate from input layer to output layer. The layers between the input and output layers are referred to as **hidden layers**.

The number of layers is referred to as the **depth** of the feedforward network. When a network has many hidden layers of neurons, feedforward networks are referred to as **deep neural networks (DNN)**. Learning in deep neural networks is referred to as **deep learning**. The number of neurons in a layer is referred to as the **width** of that layer.

The hidden layers are usually composed of perceptrons (sigmoidal units) or ReLU units and the output layer is usually

- A linear neuron layer for regression
- A softmax layer for classification

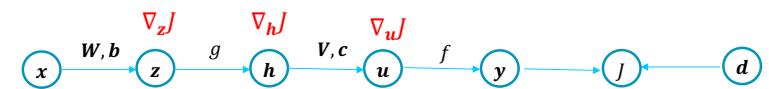
Two-layer FFN



Input $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)^T$
 Hidden-layer output $\mathbf{h} = (h_1 \ h_2 \ \cdots \ h_M)^T$
 Output $\mathbf{y} = (y_1 \ y_2 \ \cdots \ y_K)^T$

\mathbf{W}, \mathbf{b} – weight and bias of the hidden layer
 \mathbf{V}, \mathbf{c} – weight and bias of the output layer
 M is the number of hidden layer neurons

Backpropagation of gradients

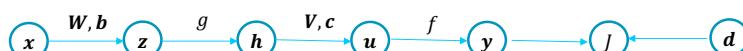


Since the targets appear at the output, the error gradient at the output layer is $\nabla_u J$ is known. Therefore, output weights and bias, \mathbf{V}, \mathbf{c} , can be learnt.

To learn hidden layer weights and biases, the gradients at the output layer are to be backpropagated to hidden layers.

Forward propagation of activations: single input pattern

Consider an input pattern (\mathbf{x}, \mathbf{d}) to 2-layer FFN:



Synaptic input \mathbf{z} to hidden layer:

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

Output \mathbf{h} of hidden layer:

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

g is hidden layer activation function.

Synaptic input \mathbf{u} to output layer:

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$$

Output \mathbf{y} of output layer:

$$\mathbf{y} = f(\mathbf{u})$$

Derivatives

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$$

Consider synaptic input to u_k to the k th neuron at the output layer. Let weight vector $\mathbf{v}_k = (v_{k1} \ v_{k2} \ \cdots \ v_{kM})^T$ and bias c_k .

The synaptic input u_k due to \mathbf{h} is given by

$$u_k = \mathbf{v}_k^T \mathbf{h} + c_k = v_{k1}h_1 + v_{k2}h_2 + \cdots + v_{kM}h_M + c_k$$

$$\frac{\partial u_k}{\partial h_j} = v_{kj} \quad \text{for all } j = 1, 2, \dots, K$$

Therefore

$$\frac{\partial \mathbf{u}}{\partial \mathbf{h}} = \begin{pmatrix} \frac{\partial u_1}{\partial h_1} & \frac{\partial u_1}{\partial h_2} & \cdots & \frac{\partial u_1}{\partial h_M} \\ \frac{\partial u_2}{\partial h_1} & \frac{\partial u_2}{\partial h_2} & \cdots & \frac{\partial u_2}{\partial h_M} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial u_K}{\partial h_1} & \frac{\partial u_K}{\partial h_2} & \cdots & \frac{\partial u_K}{\partial h_K} \end{pmatrix} = \begin{pmatrix} v_{11} & v_{12} & \cdots & v_{1M} \\ v_{21} & v_{22} & \cdots & v_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ v_{K1} & v_{K2} & \cdots & v_{KM} \end{pmatrix} = \mathbf{V}^T$$

That is,

$$\frac{\partial \mathbf{u}}{\partial \mathbf{h}} = \mathbf{V}^T \quad (A)$$

Derivatives

$$y = f(\mathbf{u})$$

Considering k th neuron:

$$y_k = f(u_k)$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{u}} = \begin{pmatrix} \frac{\partial y_1}{\partial u_1} & \frac{\partial y_1}{\partial u_2} & \dots & \frac{\partial y_1}{\partial u_K} \\ \frac{\partial y_2}{\partial u_1} & \frac{\partial y_2}{\partial u_2} & \dots & \frac{\partial y_2}{\partial u_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_K}{\partial u_1} & \frac{\partial y_K}{\partial u_2} & \dots & \frac{\partial y_K}{\partial u_K} \end{pmatrix} = \begin{pmatrix} f'(u_1) & 0 & \dots & 0 \\ 0 & f'(u_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f'(u_K) \end{pmatrix} = \text{diag}(f'(\mathbf{u}))$$

That is,

$$\frac{\partial \mathbf{y}}{\partial \mathbf{u}} = \text{diag}(f'(\mathbf{u})) \quad (\text{B})$$

where $\text{diag}(f'(\mathbf{u}))$ is a diagonal matrix composed of derivatives corresponding to individual components of \mathbf{u} in the diagonal.

Proof

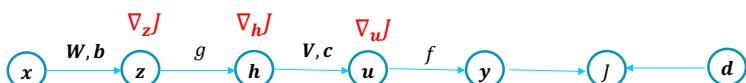
For a vector \mathbf{x} :

$$\text{diag}(f'(\mathbf{u}))\mathbf{x} = \begin{pmatrix} f'(u_1) & 0 & \dots & 0 \\ 0 & f'(u_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f'(u_K) \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{pmatrix} = \begin{pmatrix} f'(u_1)x_1 \\ f'(u_2)x_2 \\ \vdots \\ f'(u_K)x_K \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{pmatrix} \cdot \begin{pmatrix} f'(u_1) \\ f'(u_2) \\ \vdots \\ f'(u_K) \end{pmatrix}$$

That is:

$$\text{diag}(f'(\mathbf{u}))\mathbf{x} = \mathbf{x} \cdot f'(\mathbf{u}) = f'(\mathbf{u}) \cdot \mathbf{x}$$

Back-propagation of gradients: single pattern



Considering output layer,

$$\nabla_w J = \begin{cases} -(d - y) & \text{for a linear layer} \\ -(1(k=d) - f(u)) & \text{for a softmax layer} \end{cases}$$

From chain rule of differentiation,

$$\nabla_h J = \left(\frac{\partial \mathbf{u}}{\partial \mathbf{h}} \right)^T \nabla_w J = \mathbf{V} \nabla_w J \quad \text{From (A)}$$

$$\nabla_z J = \left(\frac{\partial \mathbf{h}}{\partial \mathbf{z}} \right)^T \nabla_h J = \text{diag}(g'(\mathbf{z})) \mathbf{V} \nabla_w J = \mathbf{V} \nabla_w J \cdot g'(\mathbf{z})$$

(C), from (B)

Back-propagation of gradients: single input pattern

From (C);

$$\nabla_z J = \mathbf{V} \nabla_w J \cdot g'(\mathbf{z})$$

That is, the gradients at output layer are multiplied by \mathbf{V} and back-propagated to hidden layer.

Note that hidden-layer activations are multiplied by \mathbf{V}^T (Note $\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$) in forward propagation and in back-propagation, the gradients are multiplied by \mathbf{V} .

Flow of gradients propagated in backward direction, hence named **back-propagation** (backprop) algorithm.

SGD of two-layer FFN

Output layer:

$$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) & \text{for linear layer} \\ -(1(\mathbf{k} = d) - f(\mathbf{u})) & \text{for softmax layer} \end{cases}$$

Hidden layer:

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z})$$

SGD for a two-layer FFN

Given a training dataset $\{(\mathbf{x}, \mathbf{d})\}$

Set learning parameter α

Initialize $\mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c}$

Repeat until convergence:

For every pattern (\mathbf{x}, \mathbf{d}) :

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

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

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$$

$$\mathbf{y} = f(\mathbf{u})$$

Forward propagation

$$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) \\ -(1(\mathbf{k} = d) - f(\mathbf{u})) \end{cases}$$

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z})$$

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \nabla_{\mathbf{u}} J$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{u}} J)^T$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{u}} J$$

Backward propagation

Forward propagation of activations: batch of inputs

Computational graph of 2-layer FFN for a batch of patterns (\mathbf{X}, \mathbf{D}) :



Synaptic input \mathbf{Z} to hidden layer:

$$\mathbf{Z} = \mathbf{X}\mathbf{W} + \mathbf{B}$$

Output \mathbf{H} of the hidden layer:

$$\mathbf{H} = g(\mathbf{Z})$$

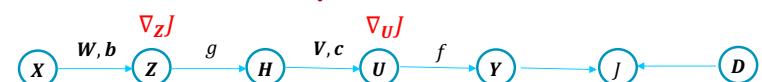
Synaptic input \mathbf{U} to output layer:

$$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C}$$

Output \mathbf{Y} of the output layer:

$$\mathbf{Y} = f(\mathbf{U})$$

Back-propagation of gradients: batch of patterns



$$\nabla_{\mathbf{U}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) & \text{for linear layer} \\ -(1(\mathbf{k} = d) - f(\mathbf{u})) & \text{for softmax layer} \end{cases}$$

$$\nabla_{\mathbf{Z}} J = \begin{pmatrix} (\nabla_{\mathbf{z}_1} J)^T \\ (\nabla_{\mathbf{z}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{z}_P} J)^T \end{pmatrix} = \begin{pmatrix} \left(\mathbf{V} \nabla_{\mathbf{u}_1} J \cdot g'(\mathbf{z}_1) \right)^T \\ \left(\mathbf{V} \nabla_{\mathbf{u}_2} J \cdot g'(\mathbf{z}_2) \right)^T \\ \vdots \\ \left(\mathbf{V} \nabla_{\mathbf{u}_P} J \cdot g'(\mathbf{z}_P) \right)^T \end{pmatrix}$$

Substituting from (C)

Back-propagation of gradients: batch of patterns

$$\nabla_{\mathbf{Z}} J = \begin{pmatrix} (\nabla_{u_1} J)^T \mathbf{V}^T \cdot (g'(\mathbf{z}_1))^T \\ (\nabla_{u_2} J)^T \mathbf{V}^T \cdot (g'(\mathbf{z}_2))^T \\ \vdots \\ (\nabla_{u_p} J)^T \mathbf{V}^T \cdot (g'(\mathbf{z}_p))^T \end{pmatrix}$$

$$(XY)^T = Y^T X^T$$

$$= \begin{pmatrix} (\nabla_{u_1} J)^T \\ (\nabla_{u_2} J)^T \\ \vdots \\ (\nabla_{u_p} J)^T \end{pmatrix} \mathbf{V}^T \cdot \begin{pmatrix} (g'(\mathbf{z}_1))^T \\ (g'(\mathbf{z}_2))^T \\ \vdots \\ (g'(\mathbf{z}_p))^T \end{pmatrix}$$

$$= (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$$

GD for a two-layer FFN

Given a training dataset (\mathbf{X}, \mathbf{D})

Set learning parameter α

Initialize $\mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c}$

Repeat until convergence:

$$\mathbf{Z} = \mathbf{X}\mathbf{W} + \mathbf{B}$$

$$\mathbf{H} = g(\mathbf{Z})$$

$$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C}$$

$$\mathbf{Y} = f(\mathbf{U})$$

Forward propagation

$$\nabla_{\mathbf{U}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(K - f(\mathbf{U})) \end{cases}$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$$

Backward propagation

$$\begin{aligned} \mathbf{V} &\leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J \\ \mathbf{C} &\leftarrow \mathbf{C} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P \\ \mathbf{W} &\leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J \\ \mathbf{B} &\leftarrow \mathbf{B} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P \end{aligned}$$

GD of two-layer FFN

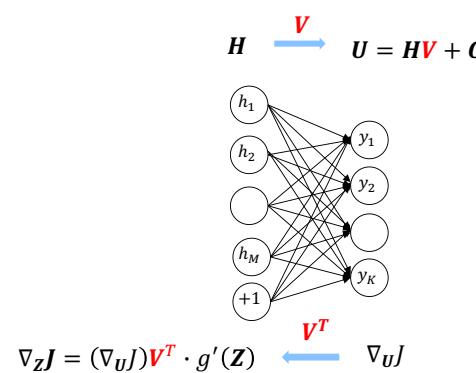
Output layer:

$$\nabla_{\mathbf{U}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(K - f(\mathbf{U})) \end{cases} \quad \begin{array}{l} \text{for linear layer} \\ \text{for softmax layer} \end{array}$$

Hidden layer:

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z}) \quad (D)$$

Back-propagation



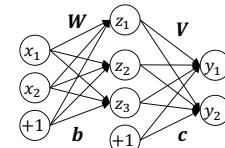
The error gradient can be seen as propagating from the output layer to the hidden layer and so learning in feedforward networks is known as the *back-propagation* algorithm

Learning in two-layer FFN

GD	SGD
(X, D)	(x, d)
$Z = XW + B$	$z = W^T x + b$
$H = g(Z)$	$h = g(z)$
$U = HV + C$	$u = V^T h + c$
$Y = f(U)$	$y = f(u)$
$\nabla_{\mathbf{U}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(\mathbf{K} - f(\mathbf{U})) \end{cases}$	$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) \\ -(1(\mathbf{k} = d) - f(\mathbf{u})) \end{cases}$
$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$	$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z})$
$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J$	$\mathbf{W} \leftarrow \mathbf{W} - \alpha x (\nabla_{\mathbf{z}} J)^T$
$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P$	$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{z}} J$
$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J$	$\mathbf{V} \leftarrow \mathbf{V} - \alpha h (\nabla_{\mathbf{u}} J)^T$
$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$	$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J$

Example 2

$$X = \begin{pmatrix} 0.77 & 0.02 \\ 0.63 & 0.75 \\ 0.50 & 0.22 \\ 0.20 & 0.76 \\ 0.17 & 0.09 \\ 0.69 & 0.95 \\ 0.00 & 0.51 \\ 0.81 & 0.61 \end{pmatrix} \text{ and } D = \begin{pmatrix} 0.44 & -0.42 \\ 0.84 & 0.43 \\ 0.09 & -0.72 \\ -0.25 & 0.35 \\ -0.12 & -0.13 \\ 0.24 & 0.03 \\ 0.30 & 0.20 \\ 0.61 & 0.04 \end{pmatrix}$$



Output layer is a linear neuron layer

Hidden layer is a sigmoidal layer

Initialized (weights using a uniform distribution):

$$\mathbf{W} = \begin{pmatrix} -3.97 & 1.10 & 0.42 \\ 2.79 & -2.64 & 3.13 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}, \mathbf{V} = \begin{pmatrix} 3.58 & -1.58 \\ -3.58 & -1.75 \\ -3.38 & 2.88 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

Example 2: Two-layer FFN

Design a two-layer FFN, using gradient descent to perform the following mapping. Use a learning factor = 0.05 and three perceptrons in the hidden-layer.

Inputs $x = (x_1, x_2)$	Targets $d = (d_1, d_2)$
(0.77, 0.02)	(0.44, -0.42)
(0.63, 0.75)	(0.84, 0.43)
(0.50, 0.22)	(0.09, -0.72)
(0.20, 0.76)	(-0.25, 0.35)
(0.17, 0.09)	(-0.12, -0.13)
(0.69, 0.95)	(0.24, 0.03)
(0.00, 0.51)	(0.30, 0.20)
(0.81, 0.61)	(0.61, 0.04)

Example 2

Epoch 1:

$$Z = XW + B = \begin{pmatrix} 0.77 & 0.02 \\ 0.63 & 0.75 \\ 0.50 & 0.22 \\ 0.20 & 0.76 \\ 0.17 & 0.09 \\ 0.69 & 0.95 \\ 0.00 & 0.51 \\ 0.81 & 0.61 \end{pmatrix} \begin{pmatrix} -3.97 & 1.10 & 0.42 \\ 2.79 & -2.64 & 3.13 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} -3.00 & 0.8 & 0.39 \\ -0.42 & -1.27 & 2.61 \\ -1.35 & -0.04 & 0.91 \\ 1.34 & -1.79 & 2.46 \\ -0.42 & -0.05 & 0.35 \\ -0.05 & -1.76 & 3.27 \\ 1.42 & -1.34 & 1.60 \\ -1.51 & -0.72 & 2.25 \end{pmatrix}$$

$$H = g(Z) = \frac{1}{1 + e^{-z}} = \begin{pmatrix} 0.05 & 0.69 & 0.60 \\ 0.40 & 0.22 & 0.93 \\ 0.21 & 0.49 & 0.71 \\ 0.79 & 0.14 & 0.92 \\ 0.40 & 0.49 & 0.59 \\ 0.49 & 0.15 & 0.96 \\ 0.80 & 0.21 & 0.83 \\ 0.18 & 0.33 & 0.91 \end{pmatrix}$$

Example 2

$$\mathbf{Y} = \mathbf{H}\mathbf{V} + \mathbf{C} = \begin{pmatrix} 0.05 & 0.69 & 0.60 \\ 0.40 & 0.22 & 0.93 \\ 0.21 & 0.49 & 0.71 \\ 0.79 & 0.14 & 0.92 \\ 0.40 & 0.49 & 0.59 \\ 0.49 & 0.15 & 0.96 \\ 0.80 & 0.21 & 0.83 \\ 0.18 & 0.33 & 0.91 \end{pmatrix} \begin{pmatrix} 3.58 & -1.18 \\ -3.58 & -1.78 \\ -3.38 & 2.88 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} -4.32 & 0.44 \\ -2.52 & 1.67 \\ -3.43 & 0.87 \\ -0.79 & 1.15 \\ -2.32 & 0.21 \\ -2.04 & 1.75 \\ -0.67 & 0.76 \\ -3.59 & 1.75 \end{pmatrix}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) = - \left(\begin{pmatrix} 0.44 & -0.42 \\ 0.84 & 0.43 \\ 0.09 & -0.72 \\ -0.25 & 0.35 \\ -0.12 & -0.13 \\ 0.24 & 0.03 \\ 0.30 & 0.20 \\ 0.61 & 0.04 \end{pmatrix} - \begin{pmatrix} -4.32 & 0.44 \\ -2.52 & 1.67 \\ -3.43 & 0.87 \\ -0.79 & 1.15 \\ -2.32 & 0.21 \\ -2.04 & 1.75 \\ -0.67 & 0.76 \\ -3.59 & 1.75 \end{pmatrix} \right) = \begin{pmatrix} -4.76 & 0.85 \\ -3.35 & 1.24 \\ -3.52 & 1.59 \\ -0.54 & 0.80 \\ -2.20 & 0.34 \\ -2.28 & 1.72 \\ -0.98 & 0.56 \\ -4.20 & 1.70 \end{pmatrix}$$

Example 2

Output layer:

$$\nabla_{\mathbf{V}} J = \mathbf{H}^T \nabla_{\mathbf{U}} J = \begin{pmatrix} -6.23 & 3.22 \\ -8.81 & 2.85 \\ -17.07 & 7.40 \end{pmatrix}$$

$$\nabla_{\mathbf{C}} J = (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P = \begin{pmatrix} -21.83 \\ 8.81 \end{pmatrix}$$

Hidden layer:

$$\nabla_{\mathbf{W}} J = \mathbf{X}^T \nabla_{\mathbf{Z}} J = \begin{pmatrix} -8.43 & 7.81 & 7.79 \\ -8.21 & 4.56 & 3.76 \end{pmatrix}$$

$$\nabla_{\mathbf{b}} J = (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P = \begin{pmatrix} -15.22 \\ 13.11 \\ 13.92 \end{pmatrix}$$

Example 2

$$g'(\mathbf{Z}) = \mathbf{H} \cdot (\mathbf{1} - \mathbf{H}) = \begin{pmatrix} 0.04 & 0.21 & 0.24 \\ 0.24 & 0.17 & 0.06 \\ 0.16 & 0.25 & 0.20 \\ 0.16 & 0.12 & 0.07 \\ 0.24 & 0.25 & 0.24 \\ 0.25 & 0.13 & 0.04 \\ 0.16 & 0.16 & 0.14 \\ 0.15 & 0.22 & 0.09 \end{pmatrix}$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$$

$$= \begin{pmatrix} -4.76 & 0.85 \\ -3.35 & 1.24 \\ -3.52 & 1.59 \\ -0.54 & 0.80 \\ -2.20 & 0.34 \\ -2.28 & 1.72 \\ -0.98 & 0.56 \\ -4.20 & 1.70 \end{pmatrix} \begin{pmatrix} 3.58 & -1.58 \\ -3.58 & -1.75 \\ -3.38 & 2.88 \end{pmatrix}^T \cdot \begin{pmatrix} 0.04 & 0.21 & 0.24 \\ 0.24 & 0.17 & 0.06 \\ 0.16 & 0.25 & 0.20 \\ 0.16 & 0.12 & 0.07 \\ 0.24 & 0.25 & 0.24 \\ 0.25 & 0.13 & 0.04 \\ 0.16 & 0.16 & 0.14 \\ 0.15 & 0.22 & 0.09 \end{pmatrix}$$

Example 2

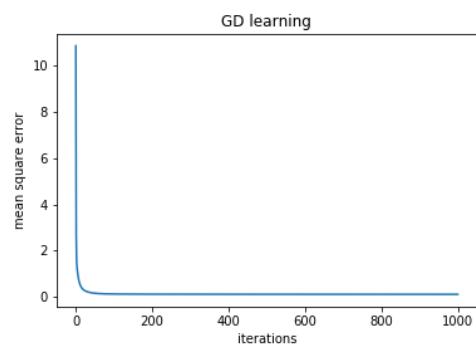
$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \nabla_{\mathbf{V}} J = \begin{pmatrix} 3.89 & -1.74 \\ -3.14 & -1.89 \\ -2.53 & 2.51 \end{pmatrix}$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{c}} J = \begin{pmatrix} 1.09 \\ -0.44 \end{pmatrix}$$

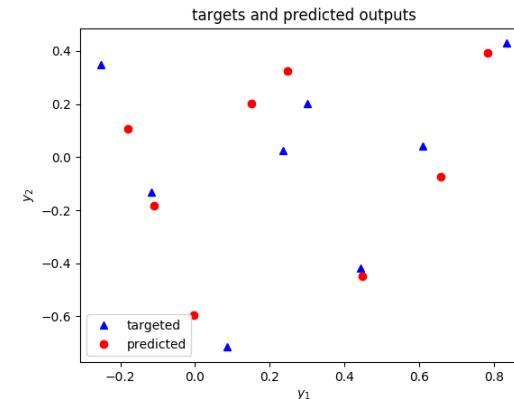
$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J = \begin{pmatrix} -3.55 & 0.72 & 0.03 \\ 3.21 & -2.87 & 2.94 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{b}} J = \begin{pmatrix} 0.76 \\ -0.66 \\ -0.70 \end{pmatrix}$$

Example 2



Example 2



After 20,000 iterations

Example 2

After 1,000 epochs,

m. s. e = 0.107

$$\mathbf{W} = \begin{pmatrix} -2.04 & -0.52 & -1.88 \\ 3.55 & -2.6 & 2.43 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 0.22 \\ -0.93 \\ 0.15 \end{pmatrix}$$

$$\mathbf{V} = \begin{pmatrix} 2.14 & -1.14 \\ -2.93 & -1.78 \\ 3.59 & 2.28 \end{pmatrix}$$

$$\mathbf{c} = \begin{pmatrix} 1.25 \\ -0.43 \end{pmatrix}$$

Preprocessing of inputs

If inputs have similar variations, better approximation of inputs or prediction of outputs is achieved. Mainly, there are two approaches to normalization of inputs.

Suppose i th input $x_i \in [x_{i,min}, x_{i,max}]$ and has a mean μ_i and a standard deviation σ_i .

If \tilde{x}_i denotes the normalized input.

1. **Scaling** the inputs such that $\tilde{x}_i \in [0, 1]$:

$$\tilde{x}_i = \frac{x_i - x_{i,min}}{x_{i,max} - x_{i,min}}$$

2. **Normalizing** the input to have standard normal distributions $\tilde{x}_i \sim N(0, 1)$:

$$\tilde{x}_i = \frac{x_i - \mu_i}{\sigma_i}$$

Preprocessing of Outputs

Linear activation function:

The convergence is usually improved if each output is **normalized** to have zero mean and unit standard deviation: $\tilde{y}_k \sim N(0,1)$

$$\tilde{y}_k = \frac{y_k - \mu_k}{\sigma_k}$$

Sigmoid activation function:

Since sigmoidal activation range from 0 to 1.0, you can **scale** $\tilde{y}_k \in [0,1]$:

$$\begin{aligned}\tilde{y}_k &= \frac{y_k - y_{k,min}}{y_{k,max} - y_{k,min}} \\ &= \frac{1}{y_{k,max} - y_{k,min}} y_k - \frac{y_{k,min}}{y_{k,max} - y_{k,min}}\end{aligned}$$

California housing dataset

<https://developers.google.com/machine-learning/crash-course/california-housing-data-description>

9 variables

The problem is to predict the housing prices using the other 8 variables.

20540 samples

Train: 14448 samples

Test: 6192 samples

longitude	A measure of how far west a house is; a more negative value is farther west
latitude	A measure of how far north a house is; a higher value is farther north
housingMedianAge	Median age of a house within a block; a lower number is a newer building
totalRooms	Total number of rooms within a block
totalBedrooms	Total number of bedrooms within a block
population	Total number of people residing within a block
households	Total number of households, a group of people residing within a home unit, for a block
medianIncome	Median income for households within a block of houses (measured in tens of thousands of US Dollars)
medianHouseValue	Median house value for households within a block (measured in US Dollars)

Example 3: Two-layer FFN predicting housing prices in California

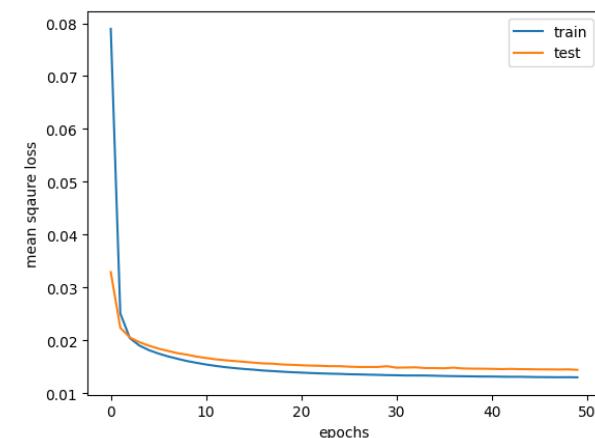
Thirteen input variables, One output variable

We use FFN with one hidden layer with 10 neurons
Network size: [8, 10, 1]

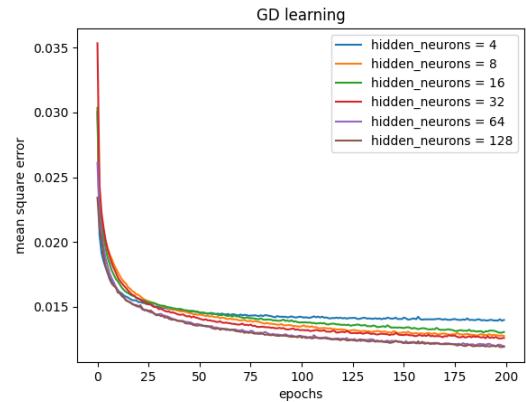
```
class FFN(nn.Module):
    def __init__(self, no_features, no_hidden, no_labels):
        super().__init__()
        self.relu_stack = nn.Sequential(
            nn.Linear(no_features, no_hidden),
            nn.ReLU(),
            nn.Linear(no_hidden, no_labels),
        )

    def forward(self, x):
        logits = self.relu_stack(x)
        return logits
```

Example 3a



Example 3b: no of hidden neurons



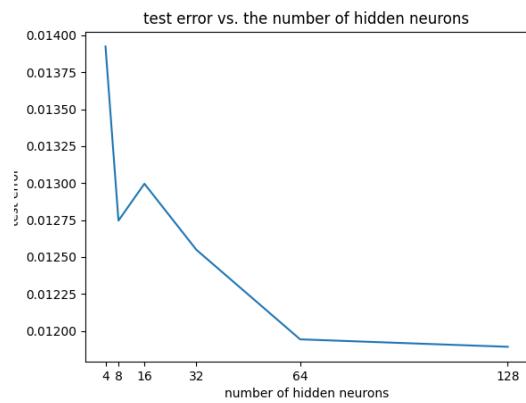
Width of hidden Layers

The number of parameters of the network increases with the **width** of layers. Therefore, the network attempts to remember the training patterns with increasing number of parameters. In other words, the network aims at minimizing the training error at the expense of its generalization ability on unseen data.

As the number of hidden units increases, the test error decreases initially but tends to increase at some point. The optimal number of hidden units is often determined empirically (that is, by trial and error).

40

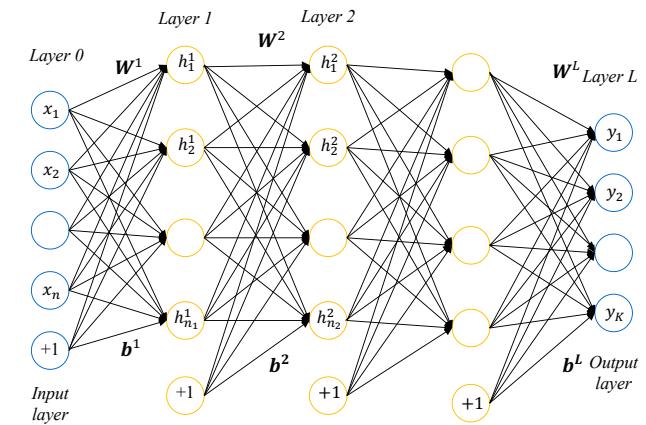
Example 3b: width of the hidden layer



Optimum number of hidden neurons 64

41

Deep neural networks (DNN)



Depth = L

DNN notations

Input layer $l = 0$:

Width = n

Input \mathbf{x}, \mathbf{X}

Hidden layers $l = 1, 2, \dots, L - 1$

Width: n_l

Weight matrix \mathbf{W}^l , bias vector \mathbf{b}^l

Synaptic input $\mathbf{u}^l, \mathbf{U}^l$

Activation function f^l

Output $\mathbf{h}^l, \mathbf{H}^l$

Output layer $l = L$

Width: K

Synaptic input $\mathbf{u}^L, \mathbf{U}^L$

Activation function f^L

Output \mathbf{y}, \mathbf{Y}

Desired output \mathbf{d}, \mathbf{D}

Back-propagation of gradients in DNN: single pattern

if $l = L$:

$$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) & \text{for linear layer} \\ -(1(\mathbf{k} = d) - f^L(\mathbf{u}^L)) & \text{for softmax layer} \end{cases}$$

else:

$$\nabla_{\mathbf{u}} J = \mathbf{W}^{l+1} (\nabla_{\mathbf{u}^{l+1}} J) \cdot f^l' (\mathbf{u}^l) \quad \text{from (C)}$$

$$\nabla_{\mathbf{W}} J = \mathbf{h}^{l-1} (\nabla_{\mathbf{u}^l} J)^T$$

$$\nabla_{\mathbf{b}} J = \nabla_{\mathbf{u}^l} J$$

Gradients are backpropagated from the output layer to the input layer

Forward propagation of activation in DNN: single pattern

Input (\mathbf{x}, \mathbf{d})

$$\mathbf{u}^1 = \mathbf{W}^{1^T} \mathbf{x} + \mathbf{b}^1$$

For layers $l = 1, 2, \dots, L - 1$:

$$\mathbf{h}^l = f^l(\mathbf{u}^l)$$

$$\mathbf{u}^{l+1} = \mathbf{W}^{l+1^T} \mathbf{h}^l + \mathbf{b}^{l+1}$$

$$\mathbf{y} = f^L(\mathbf{u}^L)$$

Forward propagation of activation in DNN: batch of patterns

Input (\mathbf{X}, \mathbf{D})

$$\mathbf{U}^1 = \mathbf{X} \mathbf{W}^1 + \mathbf{B}^1$$

For layers $l = 1, 2, \dots, L - 1$:

$$\mathbf{H}^l = f^l(\mathbf{U}^l)$$

$$\mathbf{U}^{l+1} = \mathbf{H}^l \mathbf{W}^{l+1} + \mathbf{B}^{l+1}$$

$$\mathbf{Y} = f^L(\mathbf{U}^L)$$

Back-propagation of gradients in DNN: batch of patterns

If $l = L$:

$$\nabla_{\mathbf{u}^L} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(\mathbf{K} - f^L(\mathbf{u}^L)) \end{cases}$$

Else:

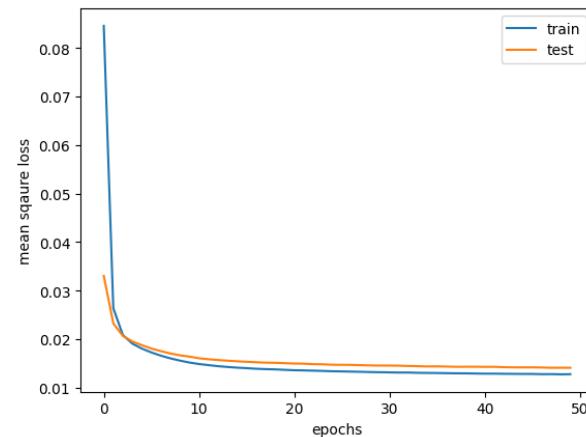
$$\nabla_{\mathbf{u}^l} J = (\nabla_{\mathbf{u}^{l+1}} J) \mathbf{W}^{l+1T} \cdot f^{l'}(\mathbf{u}^l) \quad \text{from (D)}$$

$$\nabla_{\mathbf{W}} J = \mathbf{H}^{l-1T} (\nabla_{\mathbf{u}^l} J)$$

$$\nabla_{\mathbf{b}} J = (\nabla_{\mathbf{u}^l} J)^T \mathbf{1}_P$$

Gradients are backpropagated from the output layer to the input layer

Example 4



50

Example 4: DNN on California Housing data

California housing data:

<https://developers.google.com/machine-learning/crash-course/california-housing-data-description>

Predicting housing price from other 8 variables

DNN with two hidden layers:

[8, 10, 5, 1]

```
relu_stack = nn.Sequential(
    nn.Linear(no_features, no_hidden1),
    nn.ReLU(),
    nn.Linear(no_hidden1, no_hidden2),
    nn.ReLU(),
    nn.Linear(no_hidden2, no_labels),
)
```

Example 4b: Varying the depth of DNN

Architectures:

One hidden layer: [8, 5, 1]

Two hidden layers: [8, 5, 5, 1]

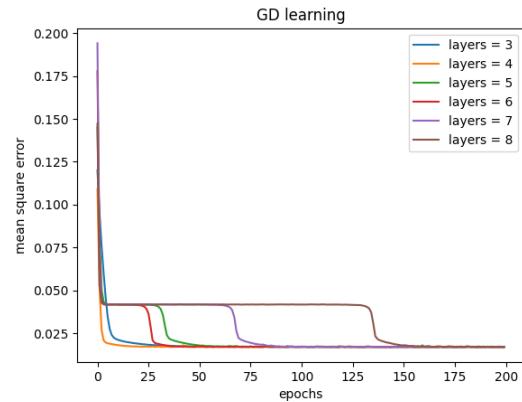
Three hidden layers: [8, 5, 5, 5, 1]

Four hidden layers: [8, 5, 5, 5, 5, 1]

Five hidden layers: [8, 5, 5, 5, 5, 5, 1]

50

Example 4

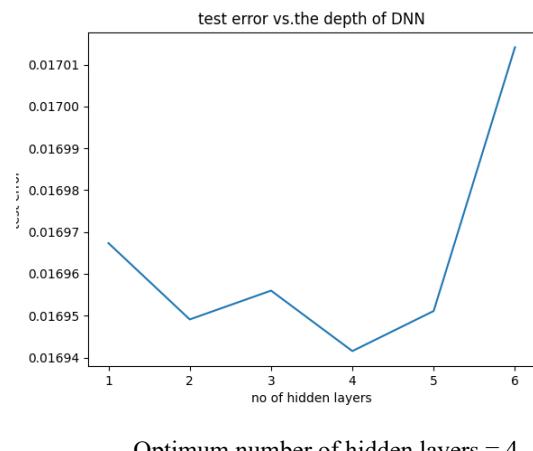


Depth of DNN

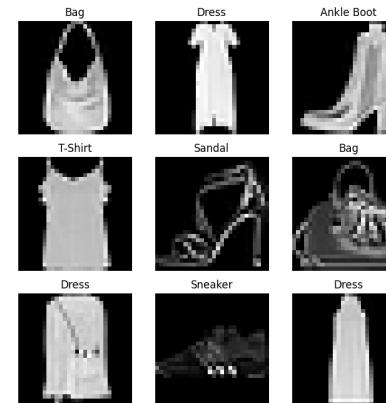
The deep networks extract features at different levels of complexity for regression or classification. However, the **depth** or the number of layers that you can have for the networks depend on the number of training patterns available. The deep networks have more parameters (weights and biases) to learn, so need more data to train. Deep networks can learn complex mapping accurately if sufficient training data is available.

The optimal number of layers is determined usually through experiments. The optimal architecture minimizes the error (training, test, and validation).

Example 4



Fashion MNIST dataset

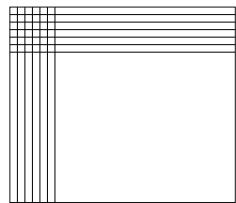


The Fasion-MNIST database of gray level images of fashion items from 10 classes:
<https://github.com/zalandoresearch/fashion-mnist>

Each image is 28x28 size.
Intensities are in the range [0, 255].

Training set: 60,000
Test set: 10,000

MNIST images



An image is divided into rows and columns and defined by its pixels.

Size of the image = rows x columns pixels

Pixels of **grey-level image** are assigned intensity values: For example, integer values between 0 and 255 assigned as intensities (grey-values) for pixels with 0 representing ‘black’ and 255 representing ‘white’.

Color images has three color channels: red, green, and blue. A pixel in a color image is a vector (r, g, b) denoting intensity in red, green, and blue channels.

```
class NeuralNetwork(nn.Module):
    def __init__(self):
        super(NeuralNetwork, self).__init__()
        self.flatten = nn.Flatten()
        self.softmax_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512),
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10),
            nn.Softmax(dim=1)
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.softmax_relu_stack(x)
        return logits

model = NeuralNetwork()
```

58

Example 5: Classification of Fashion-MNIST images

No of inputs $n = 28 \times 28 = 784$ (after flattening)

Inputs were normalized to $[0.0, 1.0]$

Use a 3-layer FFN

- Hidden-layer-1 is a perceptron layer
- Hidden-layer-2 is perceptron layer
- Output-layer is a softmax layer

Input-layer size $n = 784$

Hidden-layer-1 size $n_1 = 512$

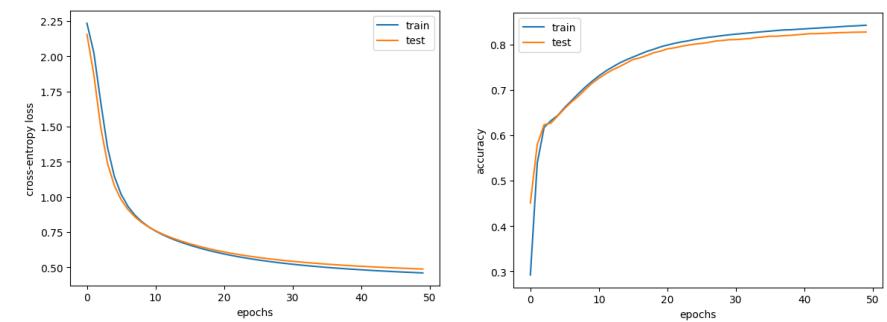
Hidden-layer-2 size $n_2 = 512$

Output-layer size $K = 10$

Training:

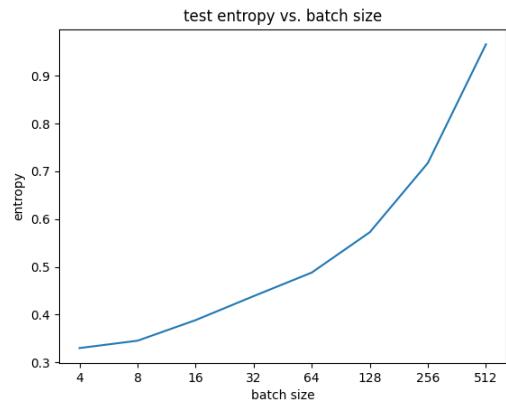
Batch size = 64

Learning rate $\alpha = 0.001$

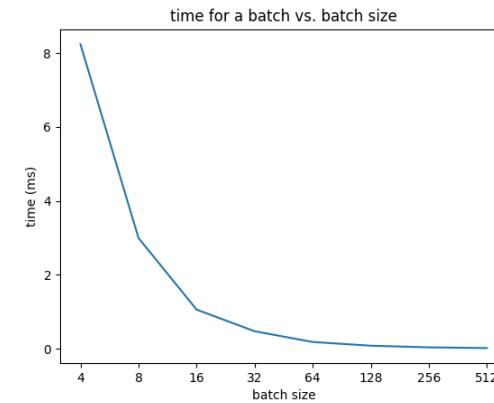


Example 5

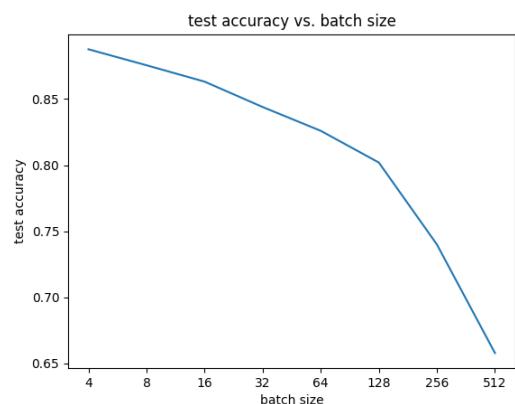
Example 5b: Effect of batch size



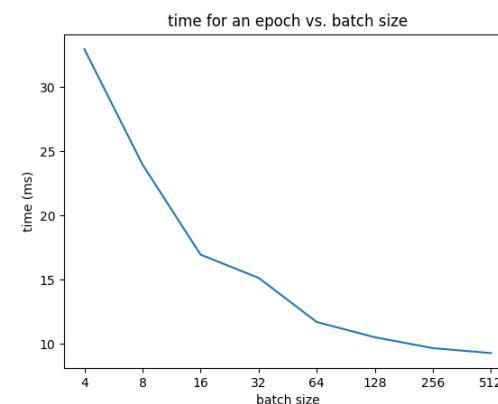
Example 5: Effect of batch size



Example 5: Effect of batch size



Example 5: Effect of batch size



Mini-batch SGD

In practice, gradient descent is performed on *mini-batch* updates of gradients within a *batch* or *block* of data of size B . In mini-batch SGD, the data is divided into blocks and the gradients are evaluated on blocks in an epoch in random order.

$B = 1$: stochastic (online) gradient descent

$B = P$ (size of training data): (batch) gradient descent

$1 < B < P \rightarrow$ mini-batch stochastic gradient descent

When B increases, more add-multiply operations per second, taking advantages of parallelism and matrix computations. On the other hand, as B increases, the number of computations per update (of weights, biases) increases.

Therefore, the curve of the time for weight update against batch size usually take a U-shape curve. There exists an optimal value of B – that depends on the sizes of the caches as well.

Summary

- Chain rule for backpropagation of gradients: $\nabla_x J = \left(\frac{\partial y}{\partial x}\right)^T \nabla_y J$
- FFN with one hidden layer (Shallow FFN)
- Backpropagation for FFN with one hidden layer:

$$\nabla_z J = (\nabla_u J) V^T \cdot g'(\mathbf{Z})$$

- Backpropagation learning for deep FFN (DNN)
- Training deep neural networks (GD and SGD):
 - Forward propagation of activation
 - Backpropagation of gradients
 - Updating weights
- Parameters to be decided: depth, width, and batch size

Selection of batch size

For SGD, it is desirable to randomly sample the patterns from training data in each epoch. In order to efficiently sample blocks, the training patterns are shuffled at the beginning of every training epoch and then blocks are sequentially fetched from memory.

Typical batch sizes: 16, 32, 64, 128, and 256.

The batch size is dependent on the size of caches of CPU and GPUs.

Chapter 5

Model selection and overfitting

Neural networks and deep learning

Model Selection



In neural networks, there exist several free parameters: learning rate, batch size, no of layers, number of neurons, etc.

Every set of parameters of the network leads to a specific model.

How do we determine the “optimum” parameter(s) or the model for a given regression or classification problem?

Selecting the best model with the best parameter values

2

Performance estimates



How do we measure the performance of the network?

Some metrics:

1. Mean-square error/Root-mean square error for **regression** — the mean-squared error or its square root. A measure of the deviation from actual.

$$MSE = \frac{1}{P} \sum_{p=1}^P \sum_{k=1}^K (d_{pk} - y_{pk})^2 \text{ and } RMSE = \sqrt{\frac{1}{P} \sum_{p=1}^P \sum_{k=1}^K (d_{pk} - y_{pk})^2}$$

2. Classification error of a **classifier**.

$$\text{classification error} = \sum_{p=1}^P 1(d_p \neq y_p)$$

where d_p is the target and y_p is the predicted output of pattern p . $1(\cdot)$ is the indicator function.

3

True error or apparent error?

Apparent error (training error): the error on the training data.
What the learning algorithm tries to optimize.

True error: the error that will be obtained in use (i.e., over the whole *sample space*). What we want to optimize *but unknown*.

However, the **apparent error** is not always a good estimate of the **true error**. It is just an optimistic measure of it.

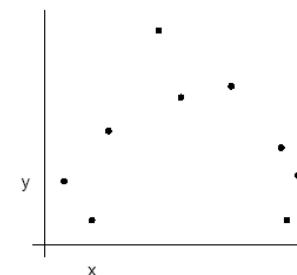
Test error: (out-of-sample error) an estimate of the true error obtained by testing the network on some independent data.

Generally, a larger test set helps provide a greater confidence on the accuracy of the estimate.

4

Example: Regression

Given the following sample data:

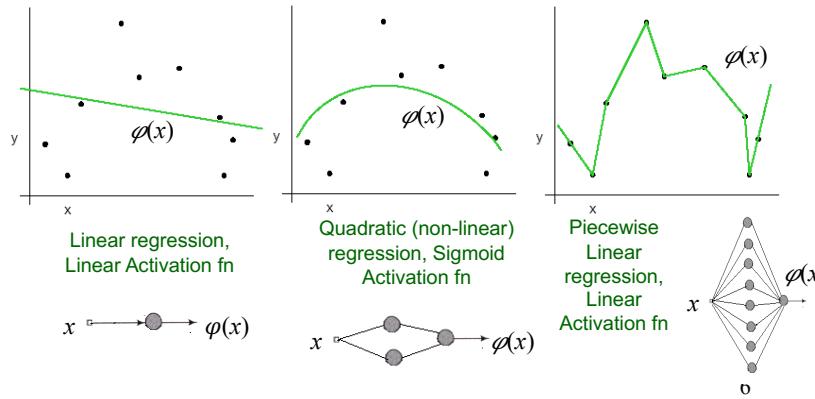


Approximate $y \approx \varphi(x)$ using an NN

5

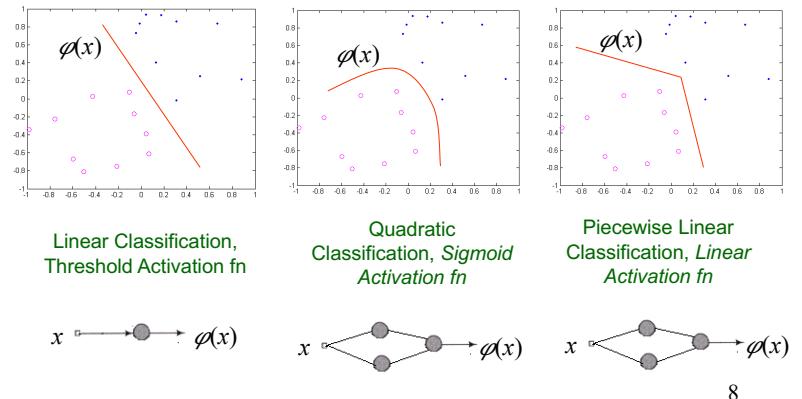
Example: Regression

Which one is the best approximation model?



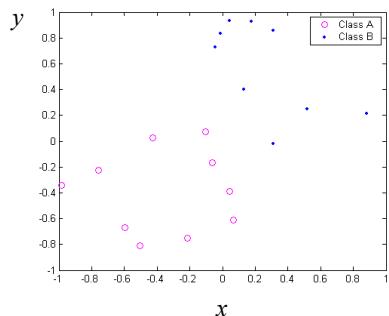
Example: Classification

Which one is the best classification model?



Example: Classification

Given the following sample data:



Classify the following data using an NN.

Estimation of true error

Intuition: Choose the model with the best fit to the data?

Meaning: Choose the model that provides the lowest error rate on the entire sample population. Of course, that error rate is the *true error rate*.

“However, to choose a model, we must first know how to **estimate the error** of a model.”

The entire **sample population** is often unavailable and only **example data** is available.

Validation

In real applications, we only have access to a finite set of examples, usually smaller than we wanted.

Validation is the approach to use the entire example data available to build the model and estimate the error rate. The validation uses a part of the data to select the model, which is known as the *validation set*.

Validation attempts to solve fundamental problems encountered:

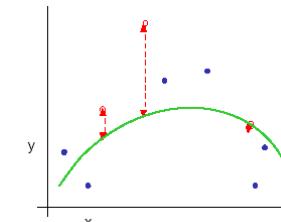
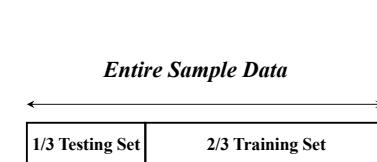
- The model tends to *overfit* the training data. It is not uncommon to have 100% correct classification on training data.
- There is no way of knowing how well the model performs on *unseen data*
- The *error rate estimate* will be overly optimistic (usually lower than the true error rate). Need to get an unbiased estimate.

10

Holdout Method

Split entire dataset into two sets:

- **Training set** (blue): used to train the classifier
- **Testing set** (red): used to estimate the error rate of the trained classifier on unseen data samples



12

Validation Methods

An effective approach is to split the entire data into subsets, i.e., Training/Validation/Testing datasets.

Some Validation Methods:

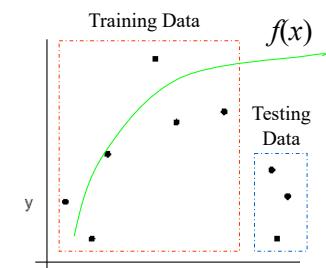
- The Holdout (1/3 - 2/3 rule for test and train partitions)
- Re-sampling techniques
 - Random Subsampling
 - K-fold Cross-Validation
 - Leave one out Cross-Validation
- Three-way data splits (train-validation-test partitions)

11

Holdout Method

The holdout method has two basic drawbacks:

- By setting some samples for testing, the training dataset becomes smaller
- Use of a single train-and-test experiment, could lead to misleading estimate if an “unfortunate” split happens



An “unfortunate” split:
training data may not cover
the space of testing data

13

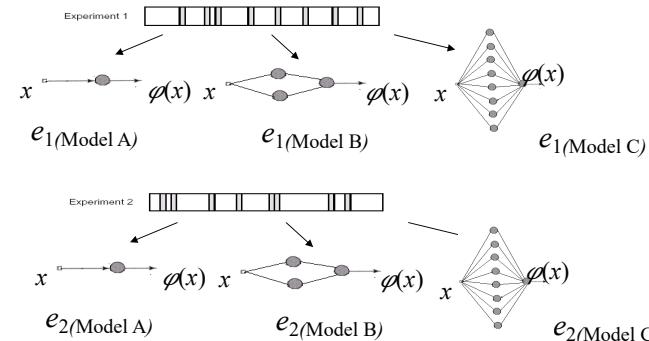
Random Sampling Methods

Limitations of the holdout can be overcome with a family of resampling methods at the expense of more computations:

- Random Subsampling
- K-Fold Cross-Validation
- Leave-one-out (LOO) Cross-Validation

14

Example: K=2 Data Splits Random SubSampling



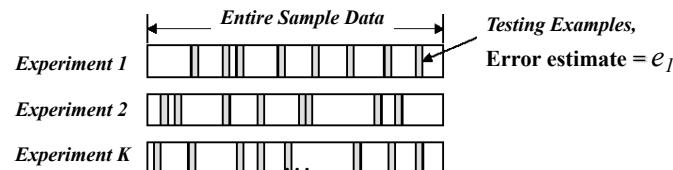
$$\text{Test error}_{(M)} = \frac{1}{2}(e_{1(M)} + e_{2(M)})$$

Choose the model with the best test error, i.e., lowest average test error!

16

K Data Splits Random SubSampling

Random Subsampling performs K data splits of the dataset for training and testing.



Each split randomly selects a (fixed) no. of examples. For each data split we retrain the classifier from scratch with the training data. Let the error estimate obtained for i th split (experiments) be e_i .

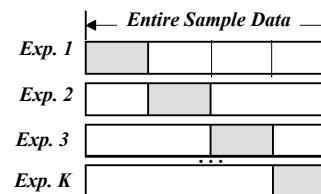
$$\text{Average test error} = \frac{1}{K} \sum_{i=1}^K e_i$$

15

K-fold Cross-Validation

Create a **K-fold** partition of the the dataset:

- For each of K experiments, use $K-1$ folds for training and the remaining one-fold for testing



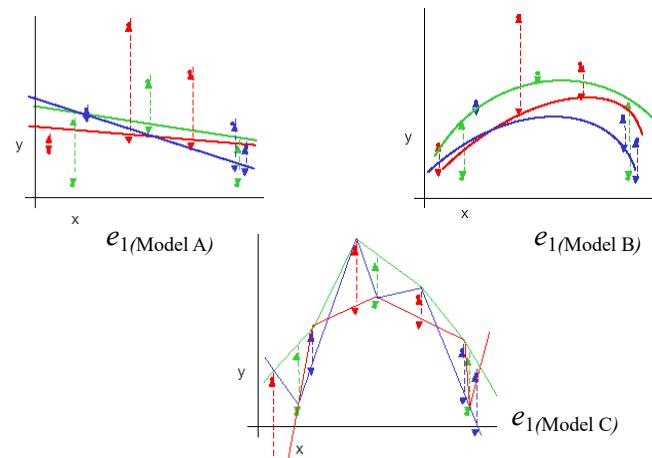
Let the error estimate for i th experiment on test partition be e_i .

$$\text{Cross-validation error} = \frac{1}{K} \sum_{i=1}^K e_i$$

K-fold cross validation is similar to Random Subsampling. The *advantage* of K-Fold Cross validation is that all examples in the dataset are eventually used for both training and testing

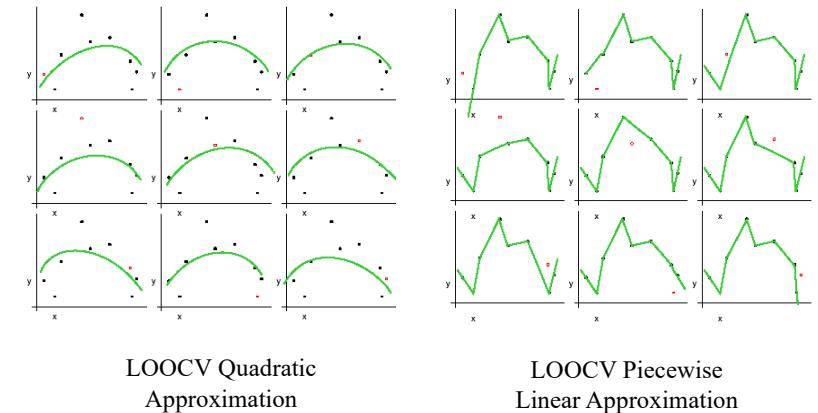
17

Example: 3-fold Cross-Validation



18

Example: Leave-One-Out Cross-Validation

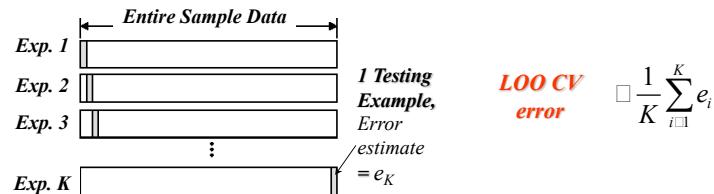


20

Leave-One-Out (LOO) Cross-Validation

Leave-One-Out is the degenerate case of **K**-Fold Cross Validation, where **K** is chosen as the total number of examples:

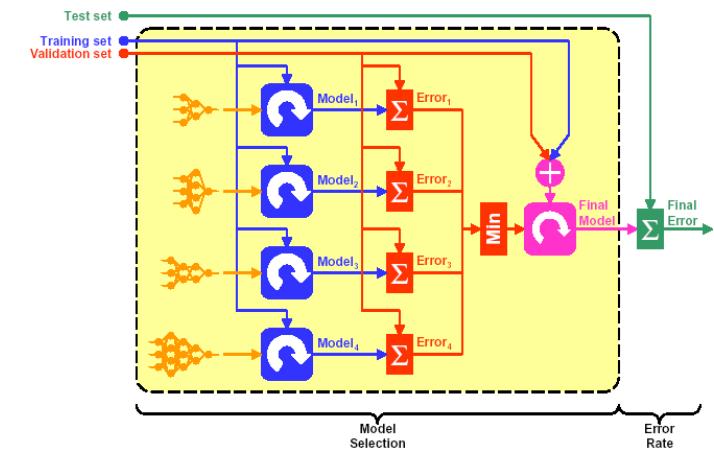
- For a dataset with N examples, perform N experiments, i.e., $N=K$.
- For each experiment use $N-1$ examples for training and the remaining one example for testing.



19

Three-Way Data Splits Method

Dataset is partitioned into training set, validation set, and testing set.



21

Three-Way Data Splits Method

If model selection and true error estimates are to be computed simultaneously, the data needs to be divided into three disjoint sets:

- **Training set:** examples for *learning* to fit the parameters of several possible classifiers. In the case of DNN, we would use the training set to find the “optimal” weights with the gradient descent rule.
- **Validation set:** examples to *determine* the error J_m of different models m , using the validation set. The optimal model m^* is given by
$$m^* = \operatorname{argmin}_m J_m$$
- **Training + Validation set:** combine examples used to re-train/redesign m^* , and find new “optimal” weights and biases.
- **Test set:** examples used only to *assess* the performance of a *trained model* m^* . We will use the test data to estimate the error rate after we have trained the final model with train + validation data.

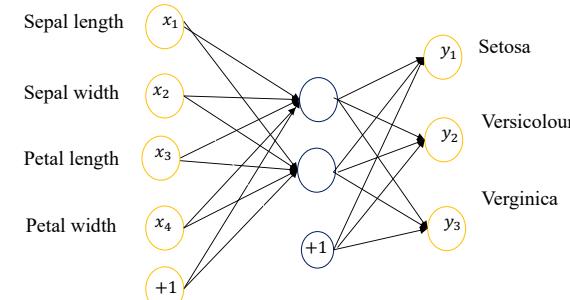
22

Examples 1-3: Iris dataset

<https://archive.ics.uci.edu/ml/datasets/Iris>

Three classes of iris flower: Setosa, versicolour, and virginica
Four features: Sepal length, sepal width, petal length, petal width
150 data points

DNN with one hidden layer. **Determine number of hidden neurons?**



24

Three-Way Data Splits Method

Why separate test and validation sets?

- The error rate estimate of the final model on validation data will be biased (smaller than the true error rate) since the validation set is also used to select in the process of final model selection.
- After assessing the final model, an independent test set is required to estimate the performance of the final model.

“NO FURTHERING TUNNING OF THE MODEL IS ALLOWED!”

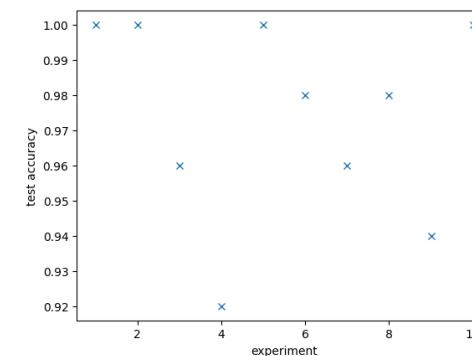
23

Example 1a: Random subsampling

150 data points

In each experiment, 50 points for testing and 100 for training

Example: 5 hidden neurons, 10 experiments

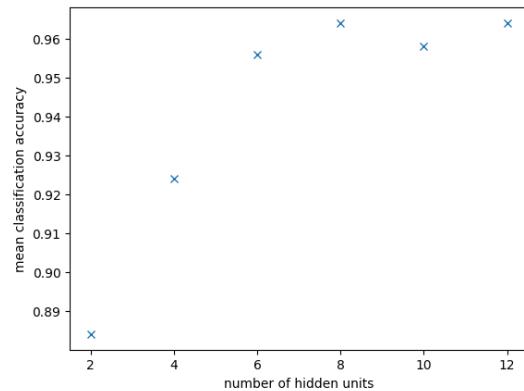


Mean accuracy = 97.4%

25

Example 1b

For different number of hidden units, misclassification errors in 10 experiments

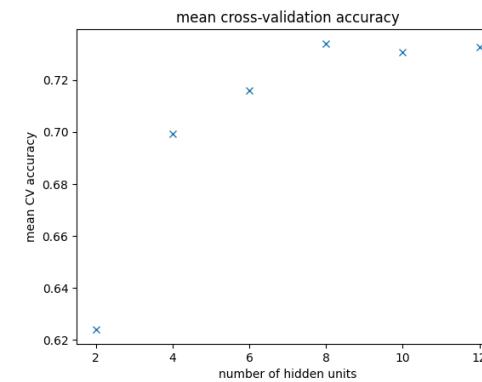


Optimum number of hidden units = 8
Accuracy = 96.4%

26

Example 2b

Mean CV error for 10 experiments for different number of hidden units:



Optimum number of hidden neurons = 8
Cross-validation accuracy = 73.4%

28

Example 2a: Cross validation

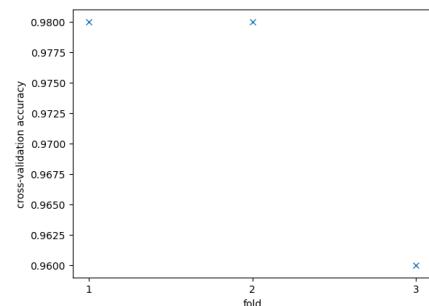
150 data points



3-fold cross validation: 50 data points in one fold.

Two folds are used for training and the remaining fold for testing

Example: hidden number of units = 5



3-fold cross-validation (CV) accuracy = 97.3%

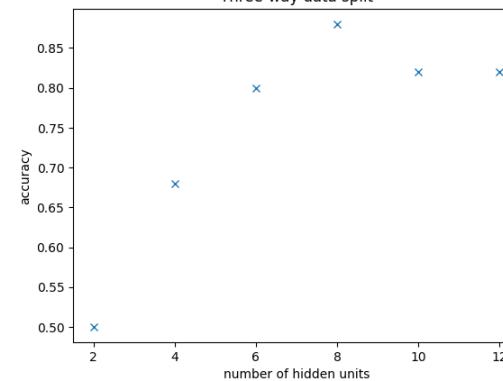
27

Example 3a: Three-way data splits

150 data points

50 data points each for training, for validation, and for testing.

Three-way data split

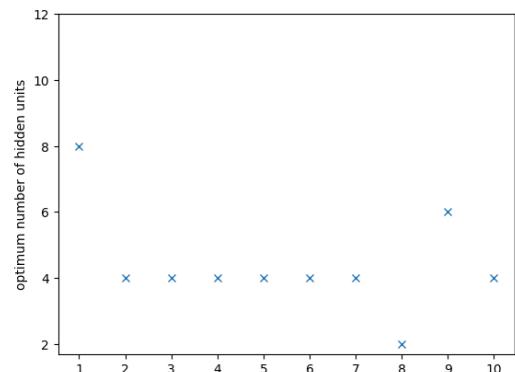


Optimum number of hidden neurons = 8
Accuracy = 88.0%

29

Example 3b

One way to further improve the results is to repeat it for many experiments, say for 10 experiments:



Optimum number of hidden neurons = 4

30

Model Complexity

Complex models: models with many adjustable weights and biases will

- more likely to be able to solve the required task,
- more likely to memorize the training data without solving the task.

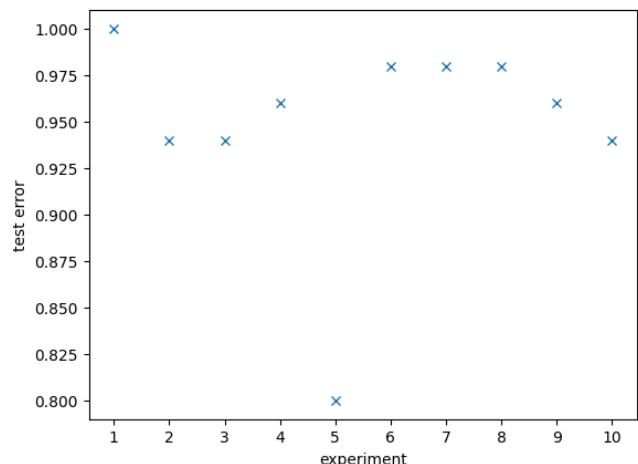
Simple models: The simpler the model that can learn from the training data is more likely to generalize over the entire sample space. May not learn the problem well.

This is the fundamental trade-off:

- Too simple — cannot do the task because not enough parameters to learn. e.g., 5 hidden neurons. (*underfitting*)
- Too complex — cannot generalize from small and noisy datasets well, e.g., 20 hidden neurons. (*overfitting*)

32

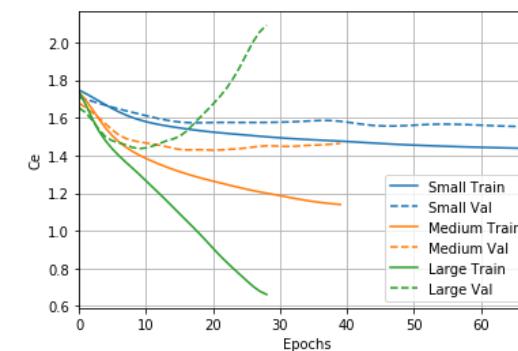
Example 3b



Test accuracy = 93.4% (average)

31

Overfitting and underfitting



Small model is **underfitting**

Medium model seems just right

Large model is **overfitting**

33

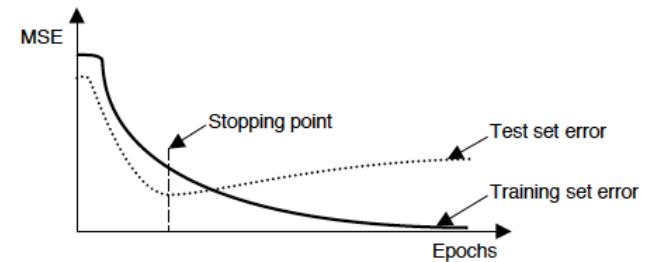
Overfitting

Overfitting is one of the problems that occur during training of neural networks, which drives the training error of the network to a very small value at the expense of the test error. The network learns to respond correctly to the training inputs by remembering them too much but is unable to generalize to produce correct outputs to novel inputs.

- Overfitting happens when the amount of training data is inadequate in comparison to the number of network parameters to learn.
- Overfitting occurs when the weights and biases become too large and are fine-tuned to remember the training patterns too much.
- Even your model is right, too much training can cause overfitting.

34

Early stopping



Training of the network is to be stopped when the validation error starts increasing. Early stopping can be used in test/validation by stopping when the validation error is minimum.

36

Methods to overcome overfitting

1. Early stopping
2. Regularization of weights
3. Dropouts

35

Regularization of weights

During overfitting, some weights attain large values to reduce training error, jeopardizing its ability to generalize. In order to avoid this, a *penalty term* (*regularization term*) is added to the cost function.

For a network with weights $\mathbf{W} = \{w_{ij}\}$ and bias \mathbf{b} , the penalized (or regularized) cost function $J_1(\mathbf{W}, \mathbf{b})$ is defined as

$$J_1 = J + \beta_1 \sum_{i,j} |w_{ij}| + \beta_2 \sum_{ij} (w_{ij})^2$$

where $J(\mathbf{W}, \mathbf{b})$ is the standard cost function (i.e., m.s.e. or cross-entropy),

$$L^1 - norm = \sum_{i,j} |w_{ij}|$$

$$L^2 - norm = \sum_{ij} (w_{ij})^2$$

And β_1 and β_2 are known as L^1 and L^2 regularization (penalty) constants, respectively. These penalties discourage weights from attaining large values.

37

L2 regularization of weights

Regularization is usually not applied on bias terms. L^2 regularization is most popular on weights.

$$J_1 = J + \beta_2 \sum_{ij} (w_{ij})^2$$

Gradient of the regularized cost wrt weights:

$$\nabla_{\mathbf{W}} J_1 = \nabla_{\mathbf{W}} J + \beta_2 \frac{\partial (\sum_{ij} w_{ij}^2)}{\partial \mathbf{W}} \quad (\text{A})$$

38

L2 regularization of weights

Substituting in (A):

$$\nabla_{\mathbf{W}} J_1 = \nabla_{\mathbf{W}} J + \beta_2 \frac{\partial (\sum_{ij} w_{ij}^2)}{\partial \mathbf{W}} = \nabla_{\mathbf{W}} J + 2\beta_2 \mathbf{W}$$

For gradient decent learning that uses regularized cost function:

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J_1$$

Substituting above:

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha (\nabla_{\mathbf{W}} J + \beta \mathbf{W})$$

where $\beta = 2\beta_2$

β is known as the *weight decay parameter*.

That is for L^2 regularization, the weight matrix is weighted by decay parameter and added to the gradient term.

40

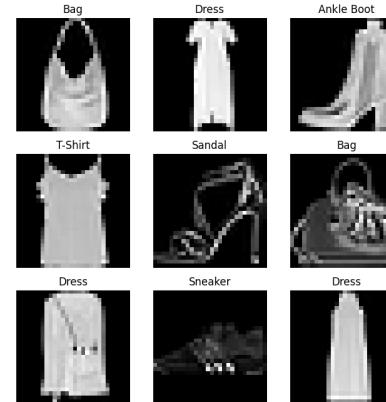
L2 regularization of weights

$$L^2 \text{ norm} = \sum_{ij} (w_{ij})^2 = w_{11}^2 + w_{12}^2 + \dots + w_{Kn}^2$$

$$\frac{\partial (\sum_{ij} (w_{ij})^2)}{\partial \mathbf{W}} = \begin{pmatrix} \frac{\partial \sum (w_{ij})^2}{\partial w_{11}} & \frac{\partial \sum (w_{ij})^2}{\partial w_{12}} & \dots & \frac{\partial \sum (w_{ij})^2}{\partial w_{1k}} \\ \frac{\partial \sum (w_{ij})^2}{\partial w_{21}} & \frac{\partial \sum (w_{ij})^2}{\partial w_{22}} & \dots & \frac{\partial \sum (w_{ij})^2}{\partial w_{2k}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \sum (w_{ij})^2}{\partial w_{n1}} & \frac{\partial \sum (w_{ij})^2}{\partial w_{n2}} & \dots & \frac{\partial \sum (w_{ij})^2}{\partial w_{nK}} \end{pmatrix} = 2\mathbf{W}$$

39

Fashion MNIST dataset



<https://github.com/zalandoresearch/fashion-mnist>

41

Example 4: Early stopping and weight decay

DNN with [784, 400, 400, 400, 10] architecture.

Let's implement L2 regularization with $\beta = 0.0001$

Use early stopping to terminate learning.

```
class NeuralNetwork(nn.Module):
    def __init__(self, hidden_size=500):
        super(NeuralNetwork, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(32*32*3, hidden_size),
            nn.ReLU(),
            nn.Linear(hidden_size, hidden_size),
            nn.ReLU(),
            nn.Linear(hidden_size, hidden_size),
            nn.ReLU(),
            nn.Linear(hidden_size, 10),
            nn.Softmax(dim=1)
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
```

42

Example 4

```
model = NeuralNetwork()
loss_fn = nn.CrossEntropyLoss()
optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate, weight_decay=0.001)
early_stopper = EarlyStopper(patience=patience, min_delta)
```

```
if early_stopper.early_stop(test_loss):
    print("Done!")
    break
```

44

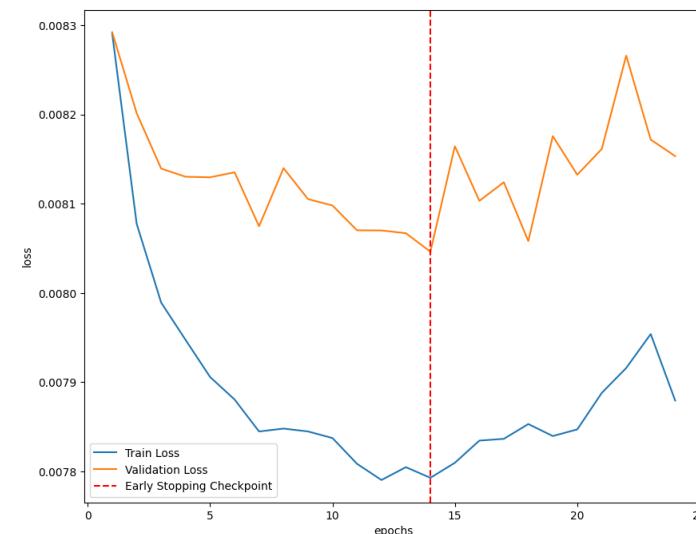
Example 4

```
class EarlyStopper:
    def __init__(self, patience=5, min_delta=0):
        self.patience = patience
        self.min_delta = min_delta
        self.counter = 0
        self.min_validation_loss = np.inf

    def early_stop(self, validation_loss):
        if validation_loss < self.min_validation_loss:
            self.min_validation_loss = validation_loss
            self.counter = 0
        elif validation_loss > (self.min_validation_loss + self.min_delta):
            self.counter += 1
        if self.counter >= self.patience:
            return True
        return False
```

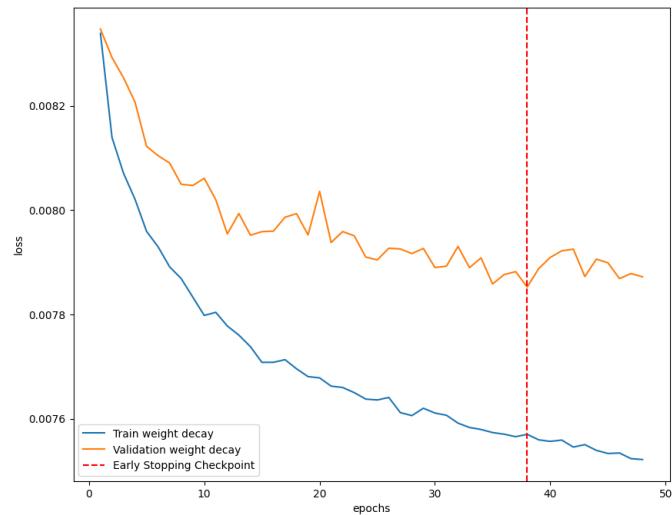
43

Example 4



45

Example 4

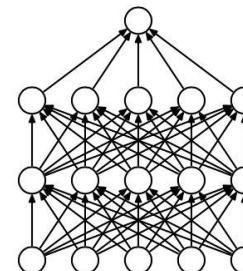


46

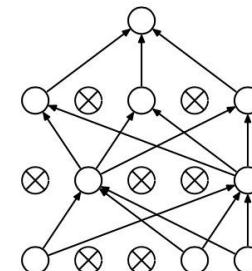
Dropouts

Overfitting can be avoided by training only a fraction of weights in each iteration. The key idea of ‘dropouts’ is to randomly drop neurons (along with their connections) from the networks during training.

This prevents neurons from co-adapting and thereby reduces overfitting.



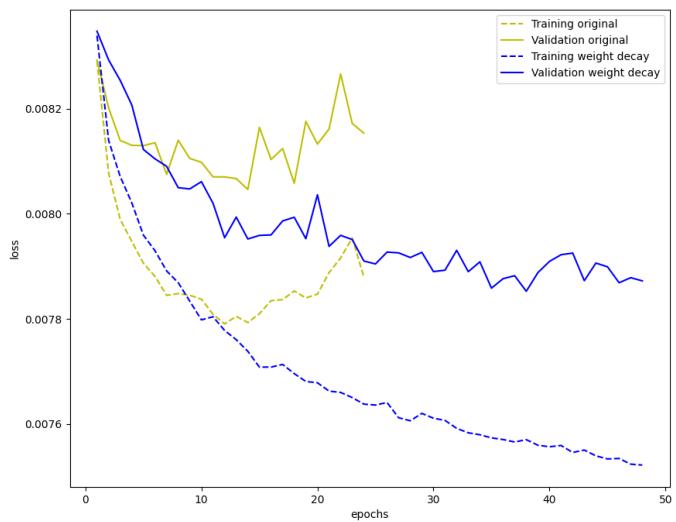
Fully connected network



Network with dropouts

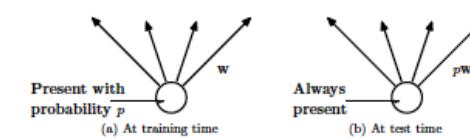
48

Example 4



47

Dropouts



At the training time, the units (neurons) are present with a probability p and presented to the next layer with weight W to the next layer.

This results in a scenario that at test time, the weights are always present, and presented to the network with weights multiplied by probability p . That is, The output at the test time is multiplied by $\frac{1}{p}$.

Applying dropouts result in a ‘thinned network’ that consists of only neurons that survived. This minimizes the redundancy in the network.

49

Dropout ratio

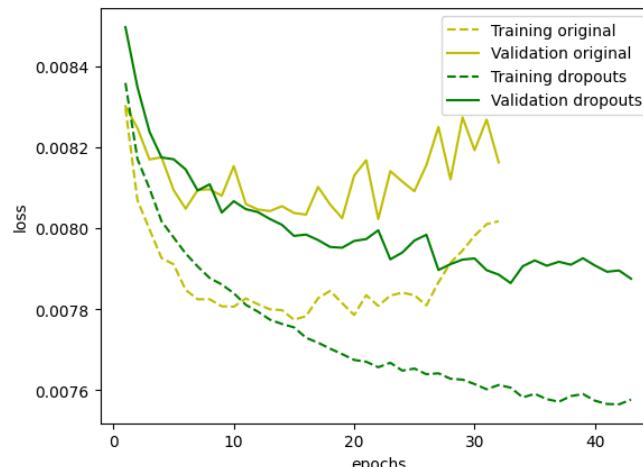
Dropout ratio is the fraction of neurons to be dropped out at one forward step.

Dropout ratio (p) has to be specified with `nn.Dropout()` in torch after activation function.

```
nn.Dropout(p=0.2)
```

50

Example 5



52

Example 5: dropouts

When training with dropouts, we train on a thinned network dropping out units in each mini-batch.

```
class NeuralNetwork_dropout(nn.Module):
    def __init__(self, hidden_size = 100, drop_out=0.5):
        super(NeuralNetwork_dropout, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(32*32*3, hidden_size),
            nn.ReLU(),
            nn.Linear(hidden_size, hidden_size),
            nn.ReLU(),
            nn.Linear(hidden_size, hidden_size),
            nn.ReLU(),
            nn.Dropout(p=drop_out),
            nn.Linear(hidden_size, 10),
            nn.Softmax(dim=1)
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
```

51

Summary

● Model selection

- Holdout method
- Resampling methods
 - Random subsampling
 - K-fold cross-validation
 - LOO cross-validation
- Three-way data split

● Methods to overcome overfitting

- Early stopping
- Weight regularization
- Dropouts

53