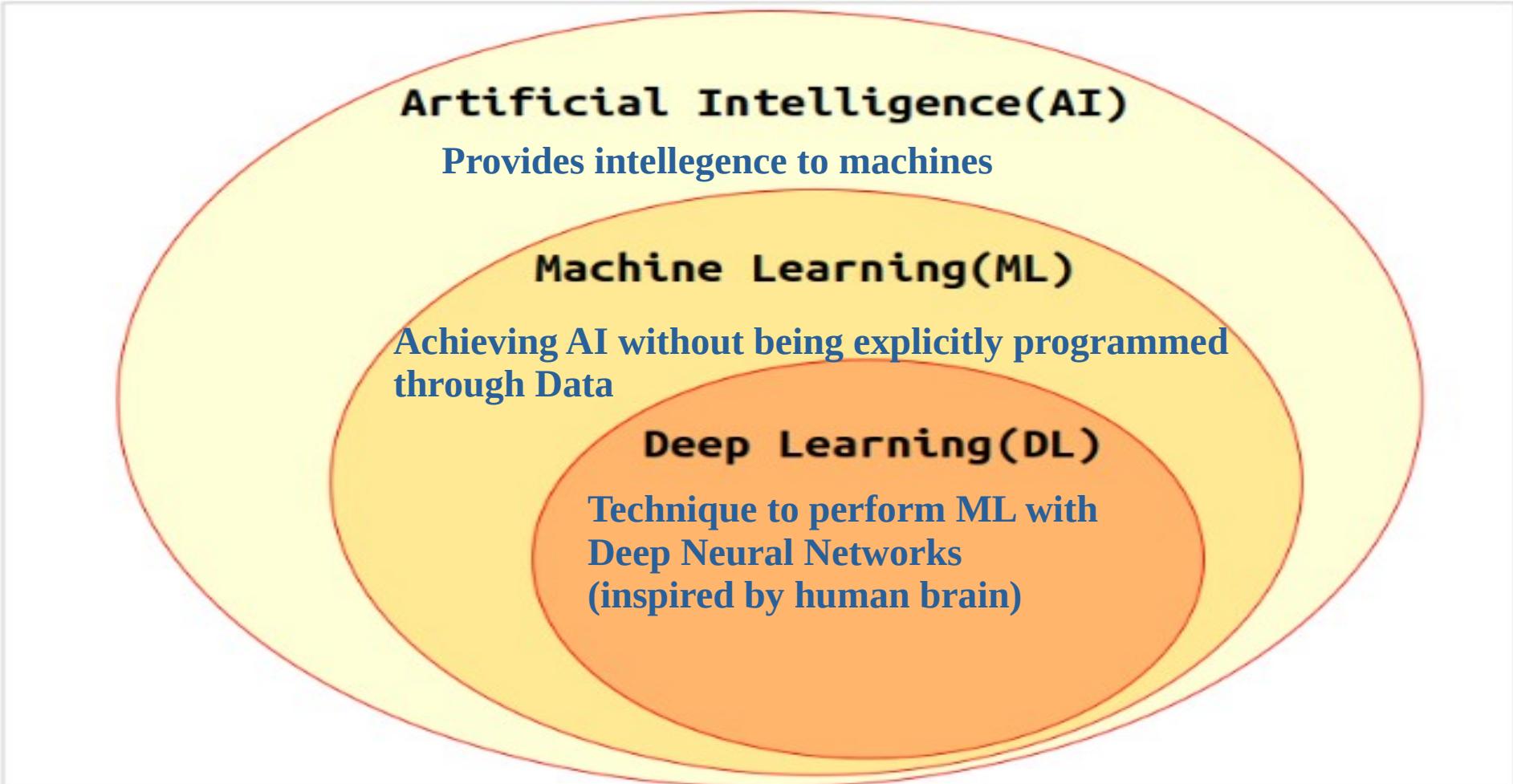


# Introduction to Deep Learning

# Deep Learning

- Deep learning is part of a broader family of machine learning methods based on artificial neural networks
- The word "deep" in "deep learning" refers to the number of layers through which the data is transformed
- Generally Deep Learning methods involve Neural Networks with more number of layers.

# What is the Difference between AI,ML and DL?



**Artificial Intelligence(AI)**

Provides intelligence to machines

**Machine Learning(ML)**

Achieving AI without being explicitly programmed through Data

**Deep Learning(DL)**

Technique to perform ML with Deep Neural Networks  
(inspired by human brain)

# What is the Difference between AI,ML and DL?

- AI, ML and DL does the same thing(providing intellegence to Machines).
- **Artificial Intelligence** involves machines that can perform tasks that are characteristic of human intelligence.
- like understanding language, recognizing objects and sounds, learning, and problem solving.
- **Machine Learning** is simply a way of achieving AI, the ability to learn without being explicitly programmed

# What is the Difference between AI,ML and DL?

- You can get AI without using machine learning
- But this would require building millions of lines of codes with complex rules and decision-trees.
- So instead of hard coding software routines with specific instructions to accomplish a particular task.
- Machine learning is a way of **training** an algorithm so that it can learn.
- **Training** involves feeding huge amounts of data to the algorithm and allowing the algorithm to adjust itself and improve

# What is the Difference between AI,ML and DL?

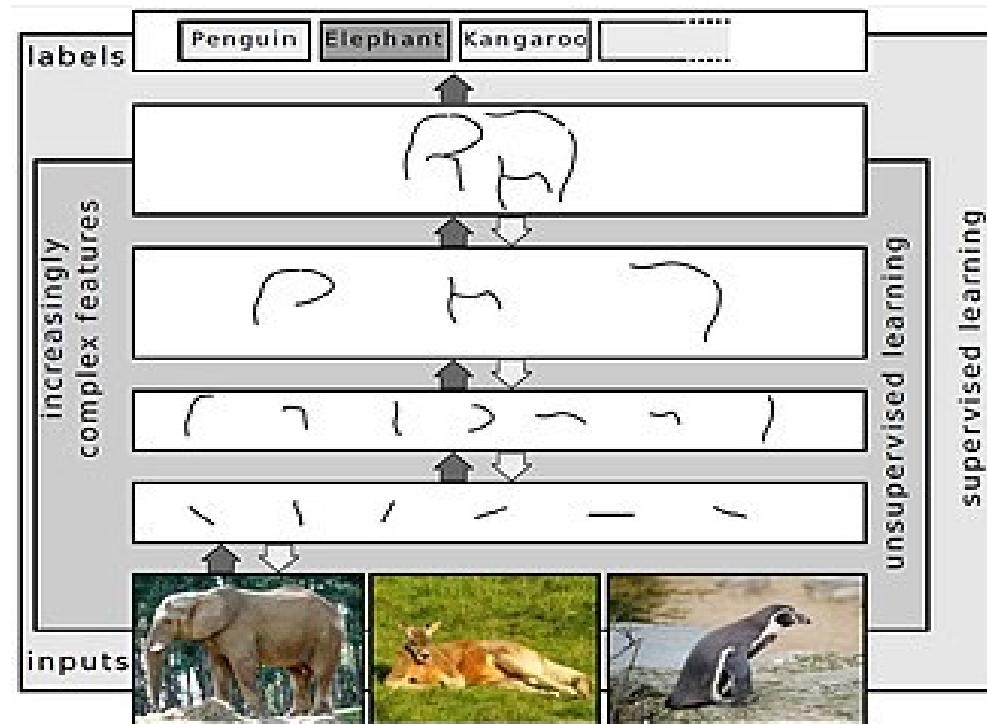
- Example: Computer Vision
- The ability of a machine to recognize an object in an image or video
- For example, the humans might tag pictures that have a cat in them versus those that do not
- algorithm tries to build a model that can accurately tag a picture as containing a cat or not as well as a human
- Once the accuracy level is high enough, the machine has now learned what a cat looks like

# What is the Difference between AI,ML and DL?

- **Deep learning** is one of many approaches to machine learning
- Other approaches include decision tree learning, inductive logic programming, clustering, reinforcement learning, and Bayesian networks.
- Deep learning was inspired by the structure and function of the brain, namely the interconnecting of many neurons.
- Artificial Neural Networks (ANNs) are algorithms that mimic the human brain
- ANNs are a collection of connected units or nodes called artificial neurons arranged as layers.

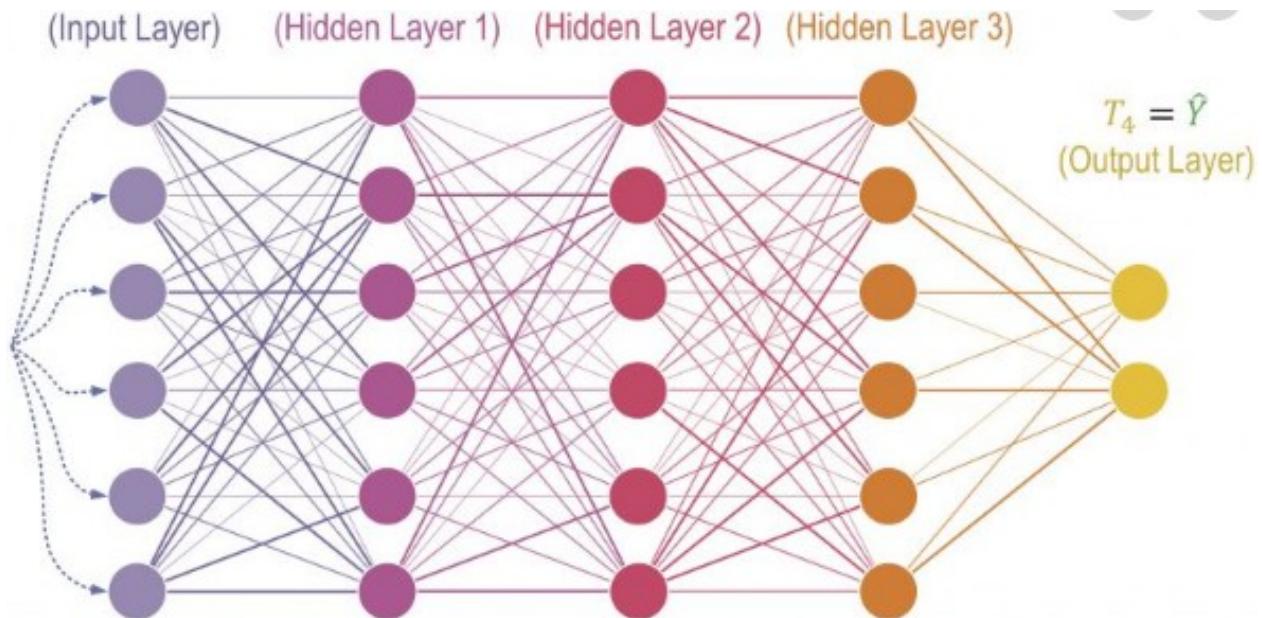
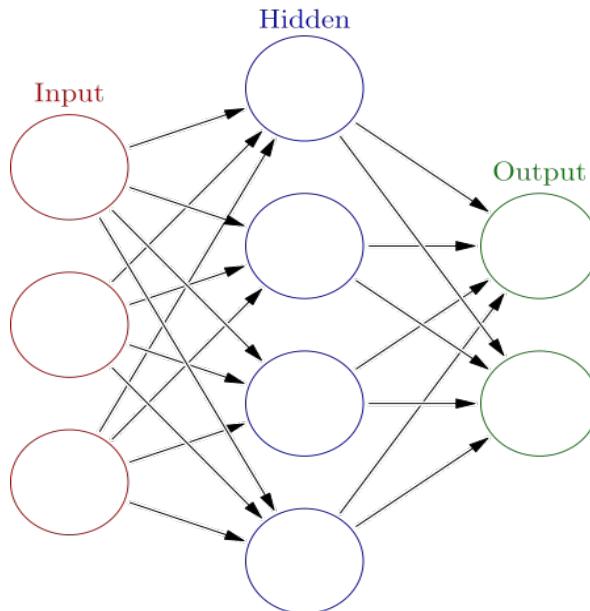
# What is the Difference between AI,ML and DL?

- Each layer picks out a specific feature to learn, such as curves/edges in image recognition



# What is the Difference between AI,ML and DL?

- The word "deep" in "deep learning" refers to the number of layers through which the data is transformed
- Depth is created by using multiple layers as opposed to a single layer

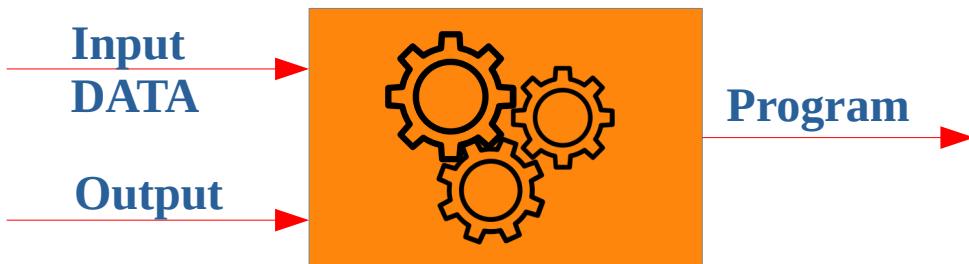


# Why we need ML based systems

- There are certain tasks which can't be programmed as there are no relevant logics describing the task.
- Even if we know the logic to write program the program might be very complicated
- **Example:**
- Compute the probability that a credit card transaction is fraudulent.
- Fraud is moving target.
- Program needs to keep changing

# Why we need ML based systems

- Instead of writing program by hand for each specific task
- We collect lot of examples that specify correct output for given input
- Machine learning Algorithm then takes these examples and produces program that does the job



# Why we need ML based systems

- ML systems learn (progressively improve their ability) to do tasks by considering examples, without task specific programming.
- If data changes program can change too by training on new data
- They have found most use in applications difficult to express with a traditional computer algorithm using rule-based programming.

# Problems best solved by learning

- Prediction of Future stock prices or currency exchange rates
- Prediction of Which movies will a person like?
- Recognizing Spoken words
- Recognizing facial expressions

# Machine Learning Paradigms

- **Supervised Learning**
  - Supervised learning is the task of inferring a function from labeled training data.
  - Each training data is a pair consisting of an input object (typically a vector) and a desired output value.
  - Regression, Decision Tree, Random Forest, KNN,SVM, Naive Bayes, Logistic Regression etc.
- **Unsupervised Learning**
  - It is that of trying to find hidden structure in unlabeled data. Since the examples given to the learner are unlabeled, there is no error or reward signal to evaluate a potential solution.
  - Apriori algorithm, K-means,GAN, Autoencoders.

# Machine Learning Paradigms

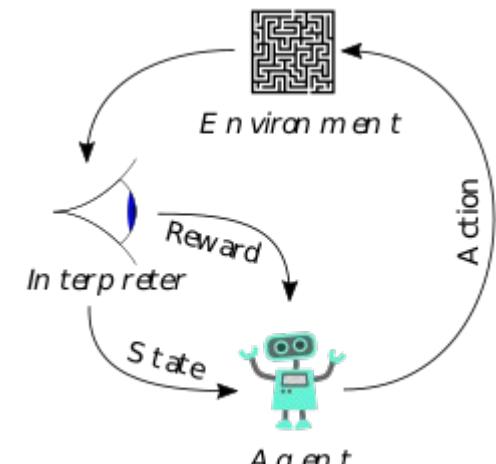
- **Reinforcement Learning**
  - an agent takes actions in an environment, which is interpreted into a reward and a representation of the state, which are fed back into the agent
  - The algorithms learn to react to an environment on their own



Supervised

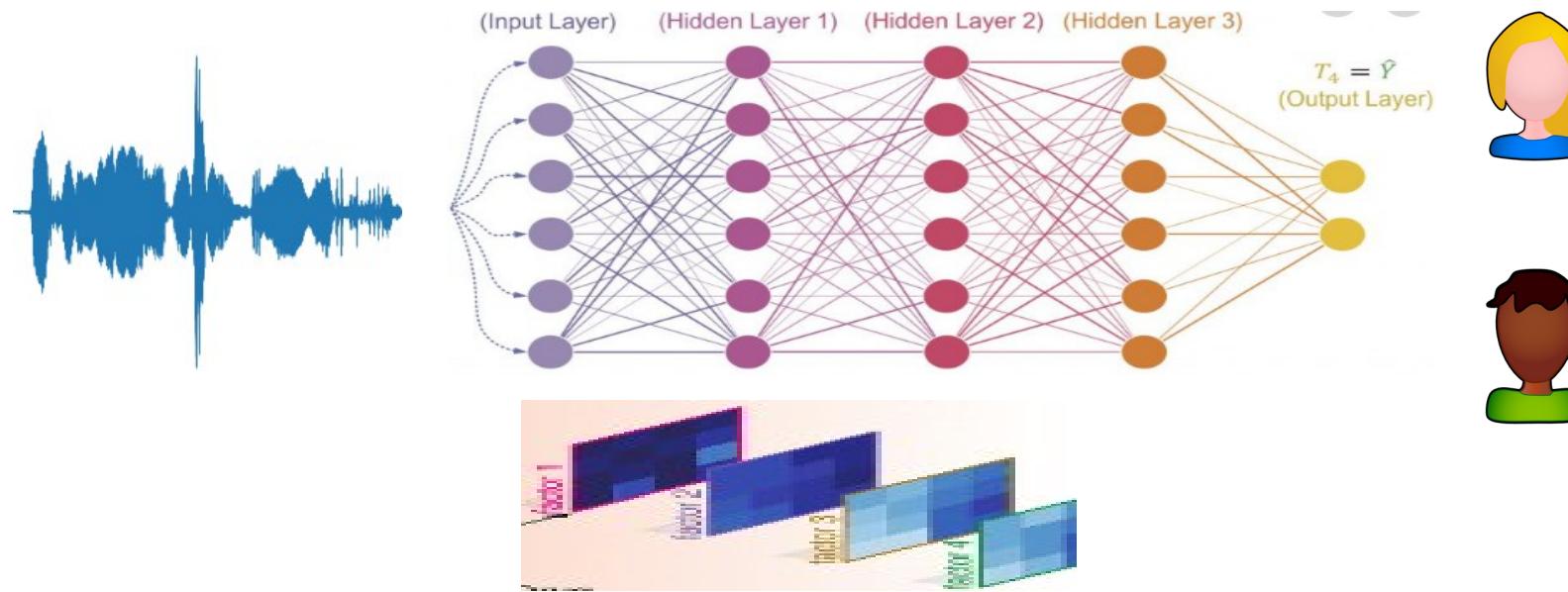
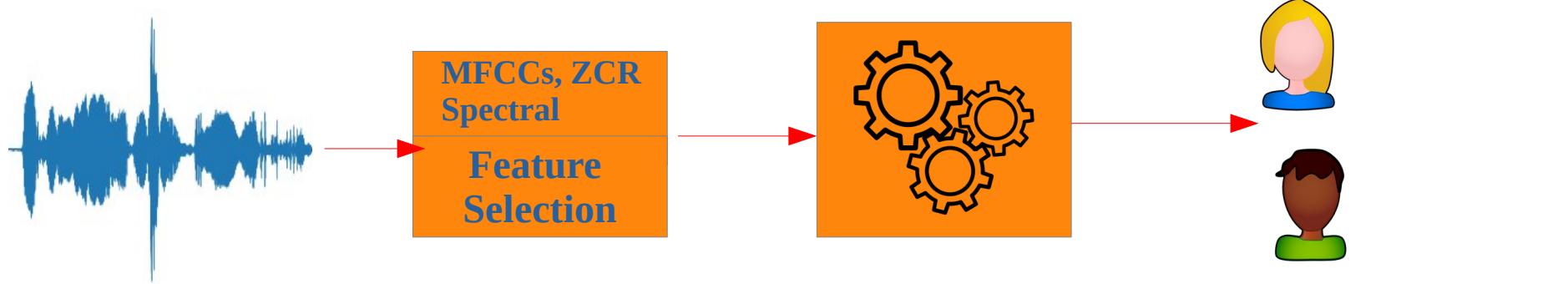


Unsupervised



Reinforcedment

# Traditional ML vs Deep Learning



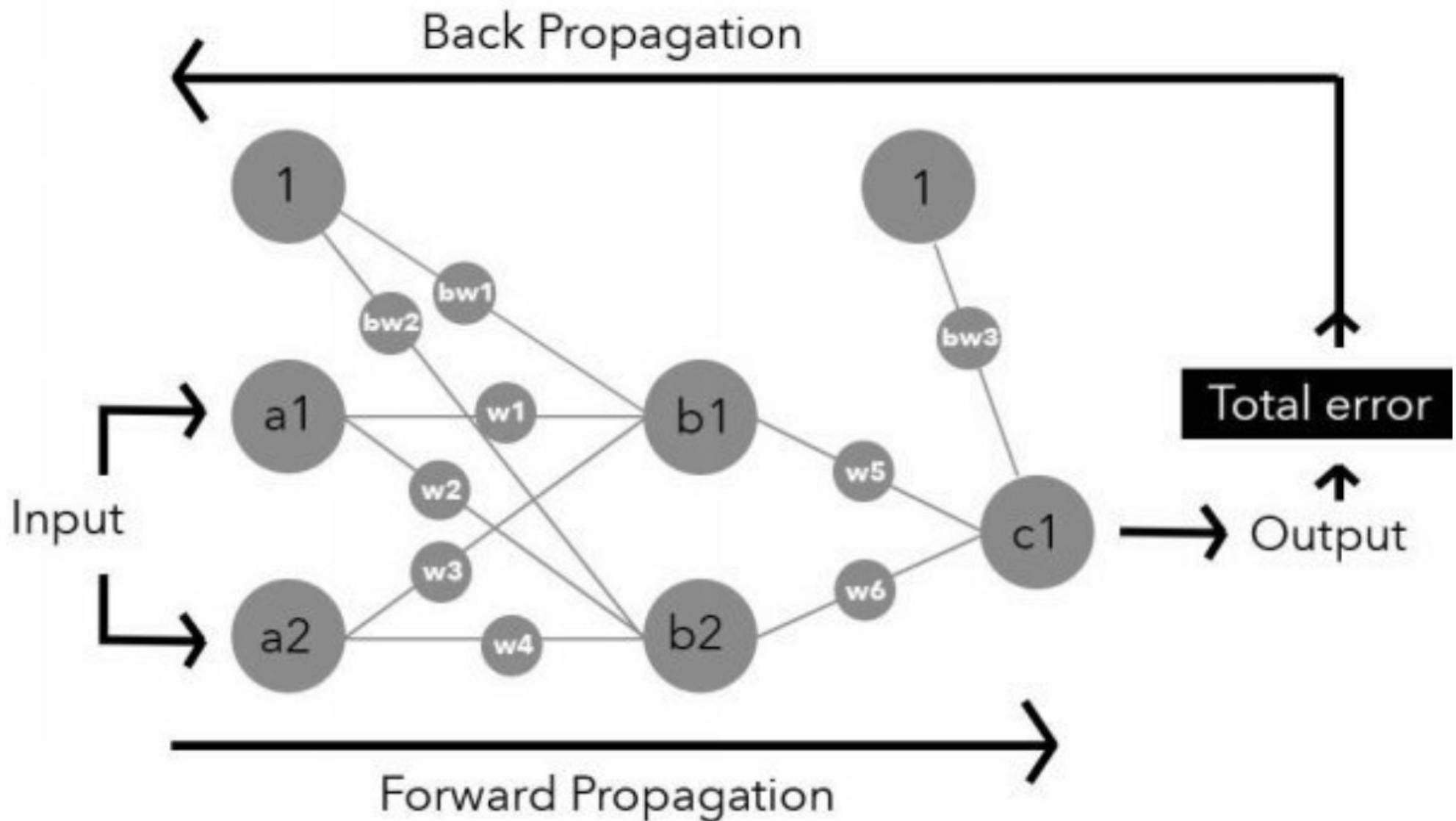
# Introduction to Neural Networks

# Applications of NN

- Make web searches better
- Organizing photos
- Speech translation
- Generate encryption
- Object detection
- Stock predictions

# What does inside NN?

- How do they accomplish these things
- They are mysterious and mind blowing
- On a high level a network learns just like we do through trial and error
- Regardless supervised, unsupervised



# Grandma's legendary soup recipe

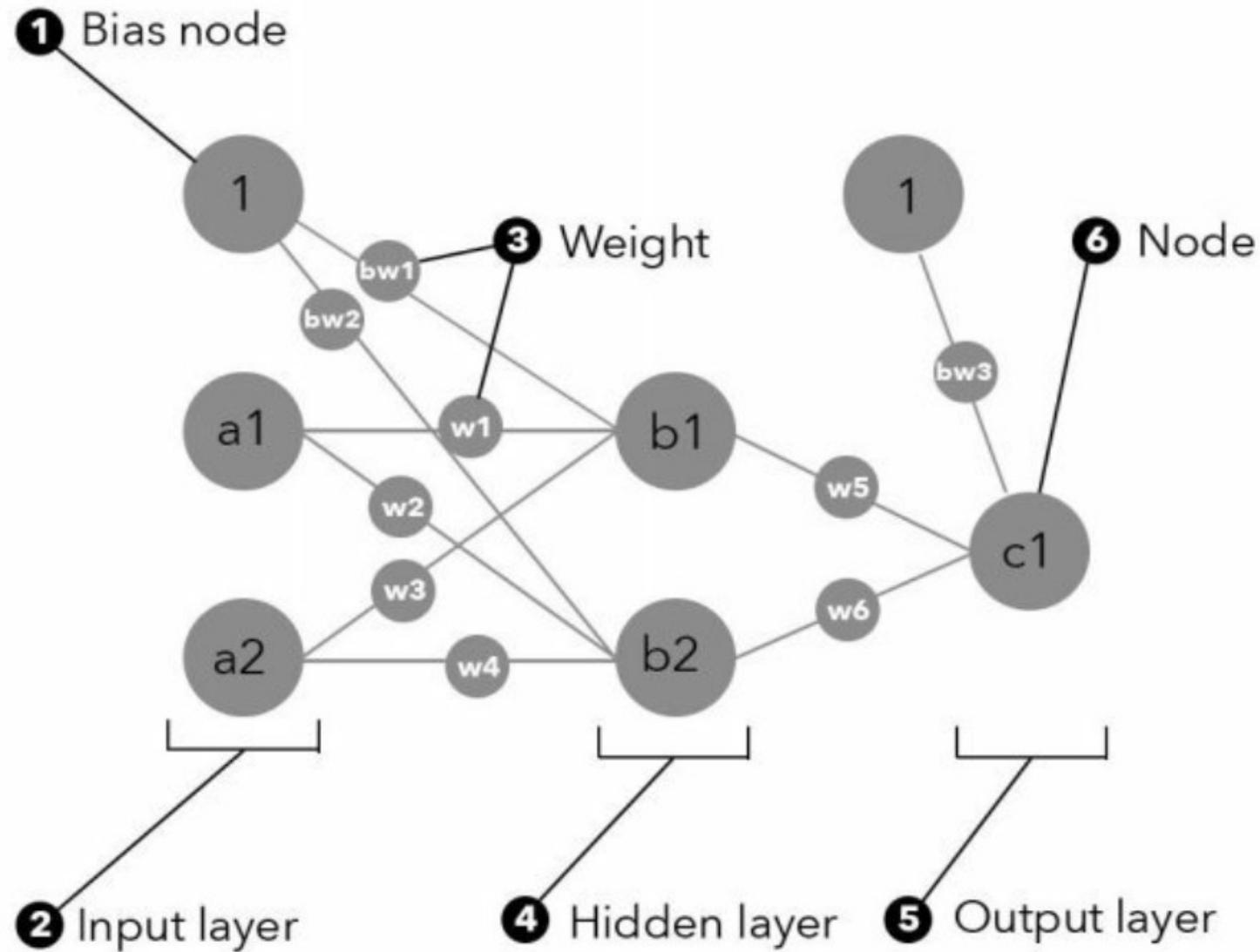
- You have 10 attempts to make the best attempt to grandma soup
- For each attempt your mom and your dad will taste the resulting soup
- And they will give feedback from your last soup
- You have to try change different quantities of ingredients to make soup quite close to grandma's
- Repeate same process 10 times
- At the last attempt you will have soup quite close to grandma's.

# Grandma's legendary soup recipe

- The inputs are the ingredients
- The weights are the quantities of each ingredient
- The neurons are your different soup attempts.
- The feedbacks are the losses.

# Creating Network Structure

- There are many ways to construct NN
- But all have a structure that is made up of similar parts
- These parts are called **hyper parameters**
- They can help a network successfully train
- hyper parameter Types
  - Required
  - Optional



# Creating Network Structure

- Require
  - Total no of input nodes
  - No of hidden layers
  - Total no of nodes in each hidden layer layer
  - Total no of output nodes
  - Weight(Synapse) values
  - Bias values
  - Learning rate

# Creating Network Structure

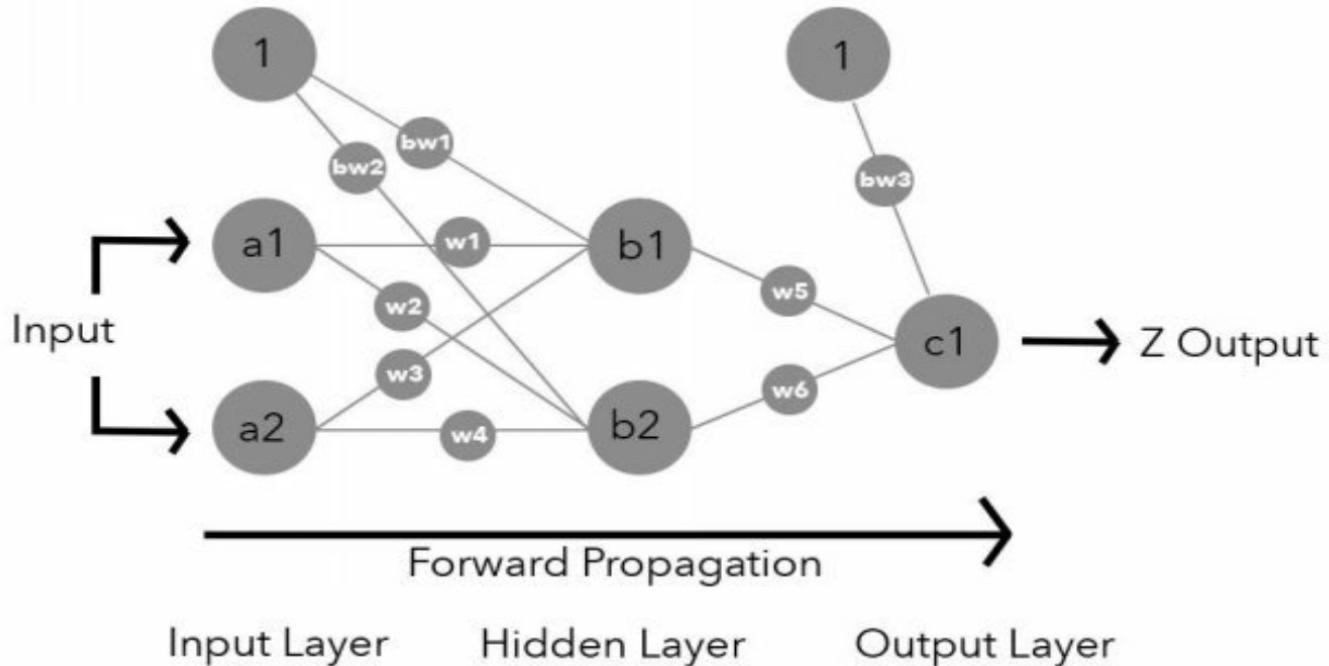
- Optimal
  - Learning rate decay
  - Momentum
  - Mini batch size
  - Weight decay
  - Dropout
- This list is not exhaustive

# Stages involved in NN implementation

- **Forward Propagation**
  - Summation operation
  - Activation Function
  - Calculate Total Error
    - Cost function
- **Back Propagation**
  - Calculate Gradients
    - Partial Derivatives
    - Chain Rule
  - Updating Weights
    - Weight update formula

# Forward Propagation

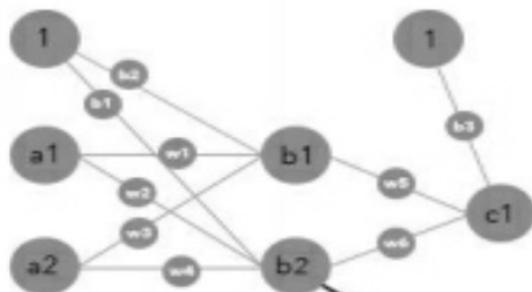
- The process of how input moves through the network to become output is called **Forward Propagation**



# Forward Propagation

- When input is passed into the network it moves from one layer to next until it passes through output layer
- Two mathematical functions are used to make all of this possible
- Mathematical functions
  - Summation operator
  - Activation function

# Single Node



1 A single node

A diagram showing the computation of a single node. It consists of two overlapping ovals. The left oval contains the summation operator  $\sum$ . The right oval contains the activation function  $\sigma(x)$ . A line from the text 'The summation operator' points to the  $\sum$  symbol. A line from the text 'The activation function...' points to the  $\sigma(x)$  symbol.

3 The activation function. In this example, it is the logistic function.

2 The summation operator

# Summation operator

- It sums all of a node's inputs to create a net input
- General Summation Operator

① The Greek letter sigma denotes the summation operator.

$$\sum_{i=1}^n x_i$$

This  $x_i$  represents the set of numbers that will be summed.

② This "i" is called the "index of summation." The numbers "1" and "n" are the lower and upper limits of the summation.

# Summation operator

- Calculating the net input of a node

$$\text{netinput} = b + \sum_{i=1}^n x_i w_i$$

**①** This “**b**” represents the input from a bias node.

**②** This “**i**” is called the “index of summation.” It begins with the first input node **(1)** and ends on input node “**n**”.

**③** The **X<sub>i</sub>** represents a unique node. The **W<sub>i</sub>** represents the unique weight situated on the nodes edge.

**④** This “**n**” represents the total number of input nodes.

# Summation Operator

- Bias values is always set to 1(one)
- So output of bias is always one
- Bias is always multiplied with edge weight to become input for a node
- The final output of a bias to a node is often its edge weight value

# Why this Summation Operator?

- Each hidden and output has multiple input values
- For these input values successfully move through the network and create an output they must be summed and turned into **single value**
- To do this summation operation make use of matrices
- Output is dot-product

# Activation function

- It receives output of the summation function and transform it into final output of a node
- It squashes input into certain range(i.e how much a node should fire)

# Types of Activation functions

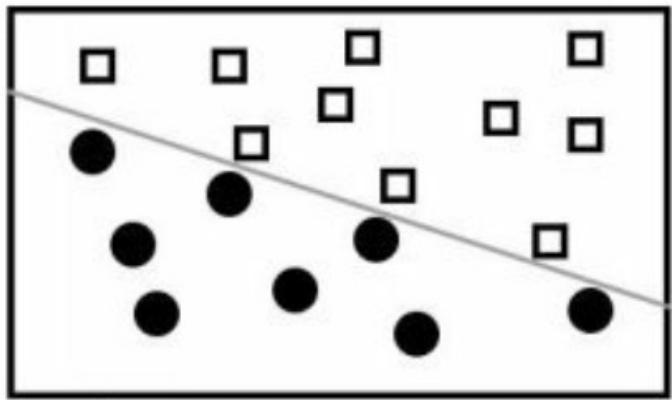
Activation Function	Graph	Equation
Linear		$f(x) = x$
Step (Heaviside)		$f(x) = \begin{cases} 0, & x < 0, \\ 1, & x \geq 0, \end{cases}$
Hyperbolic tangent		$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
Rectified Linear Unit (ReLU)		$f(x) = \max(0, x)$

① The “x” in every activation function equation represents the input of each function. The input is the *net input* calculated by the summation operator.

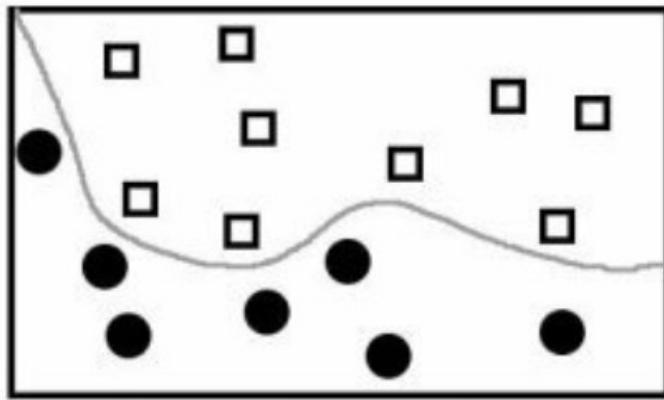
② Every “e” in this equation stands for a mathematical constant that is approximately 2.71828.

# Why Activation Functions?

- **Reason 1: Introducing Non-Linearity**
- Neural networks are often used to solve non-linear problems
- Non linear problems can't be expressed by linear functions like  $y = ax+b$ ,  $ax+by+c = 0$
- $ax^2 + bx + c = y$  or  $xy + 3x + 4=0$  are non-linear



1 There are two classes in this example. As you can see, they can be linearly divided.



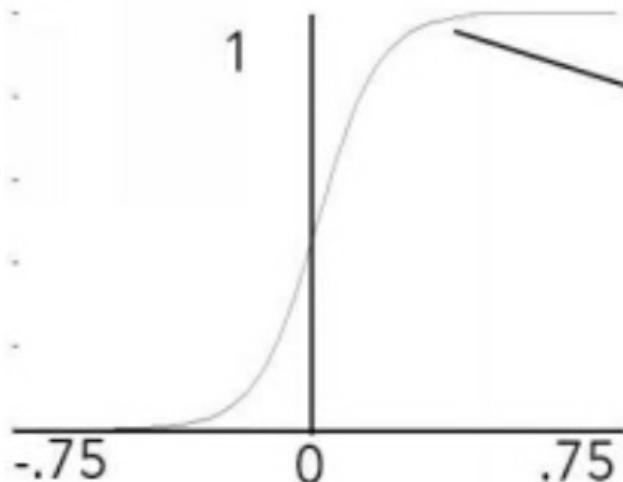
2 There are two classes in this example, but they cannot be divided linearly. Non-linearity must be introduced into the problem to successfully classify. The logistic (Sigmoid) and tanh functions both introduce this, and when used in conjunction with a multi-layer neural network can solve these problems.

# Why Activation Function

- Activation function such as tanh, logistic(sigmoid) break the linearity enable solve more complex problems.
- **Reason 2: Limiting Output**
- The tanh, sigmoid functions limit the output of a node to certain range
  - Tanh produces output between -1 and 1
  - The Logistic function produces output between 0 and 1
- Advantage of limiting output in certain range weight adjustment, faster training

# Advantages of Limiting Output

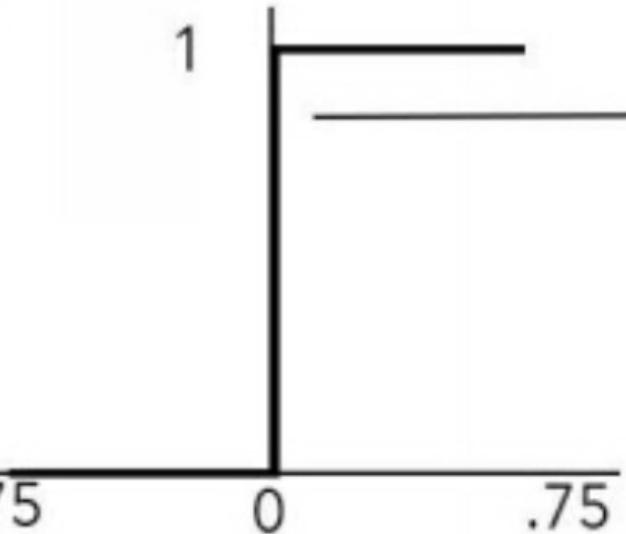
- Smoothness enables small changes in the weights and bias to produce a small change in the output.
- You can view this fine tuning that makes the task of “learning” much smoother and easier



1

1

- 1 The curved line represents the range of output of a logistic function. It's nice and smooth, and enables a small change in input to create a small change in output.



1

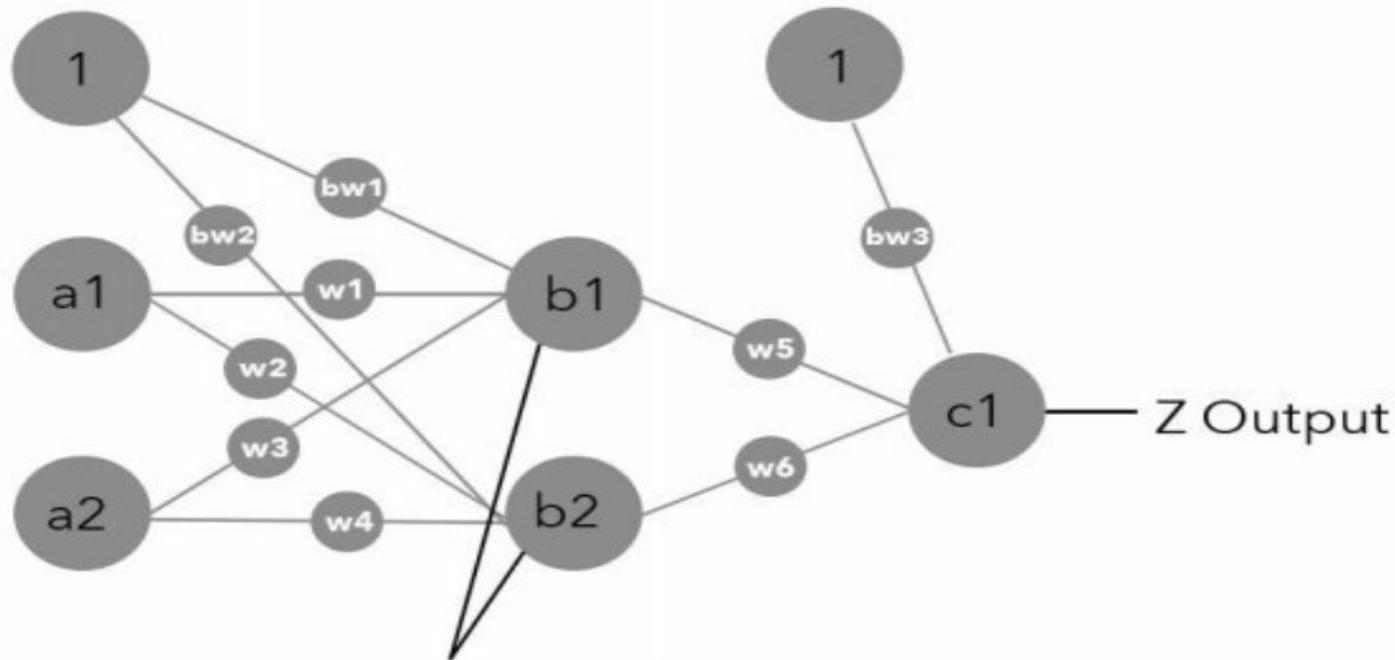
①

The **bolded** step represents the range of output for a Step (Heaviside) function. All output is either 0 (zero) or 1. A small change in input can create a large change in output.

# Advantages of Limiting Output

- The step function produces 1 or 0
- This means that small change in weight and bias is not reflected by a small change in output
- This can make learning difficult

# Forward Propagation



- 1 To calculate the net input of **b1** and **b2**, we need to multiply **a1** and **a2** by their respective weights and then sum the answers into **b1** and **b2**, respectively. The bias also needs to be added.

- 1 This  $x_{in}$  denotes the input into the hidden nodes **b1** and **b2**. It could be any notation.

$$x_{in} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} \cdot \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

2 The weight matrix.

3 The input vector.

1

The net input to node b1 (**b1net**), is calculated by multiplying the circled elements in the vector and matrix. The bias is also added.

$$x_{in} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} \cdot \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} + B = \begin{bmatrix} b_{1net} \end{bmatrix}$$

# Forward Propagation

- If we write in formula(Summation Operator)

$$b_{1\text{net}} = (a_1 \cdot w_1) + (a_2 \cdot w_3) + \text{bias}$$

$$b_{1\text{net}} = (a_1 \cdot w_1) + (a_2 \cdot w_3) + (\text{bias} \cdot b_{w1}) \text{ since bias is 1}$$

$$b_{1\text{net}} = (a_1 \cdot w_1) + (a_2 \cdot w_3) + b_{w1}$$

1

The net input to node b2 (**b2net**), is calculated by multiplying the circled elements in the vector and matrix. The bias is also added.

$$x_{in} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} \cdot \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} + B = \begin{bmatrix} b_{1net} \\ b_{2net} \end{bmatrix}$$

# Forward Propagation

- If we write in formula(Summation Operator)

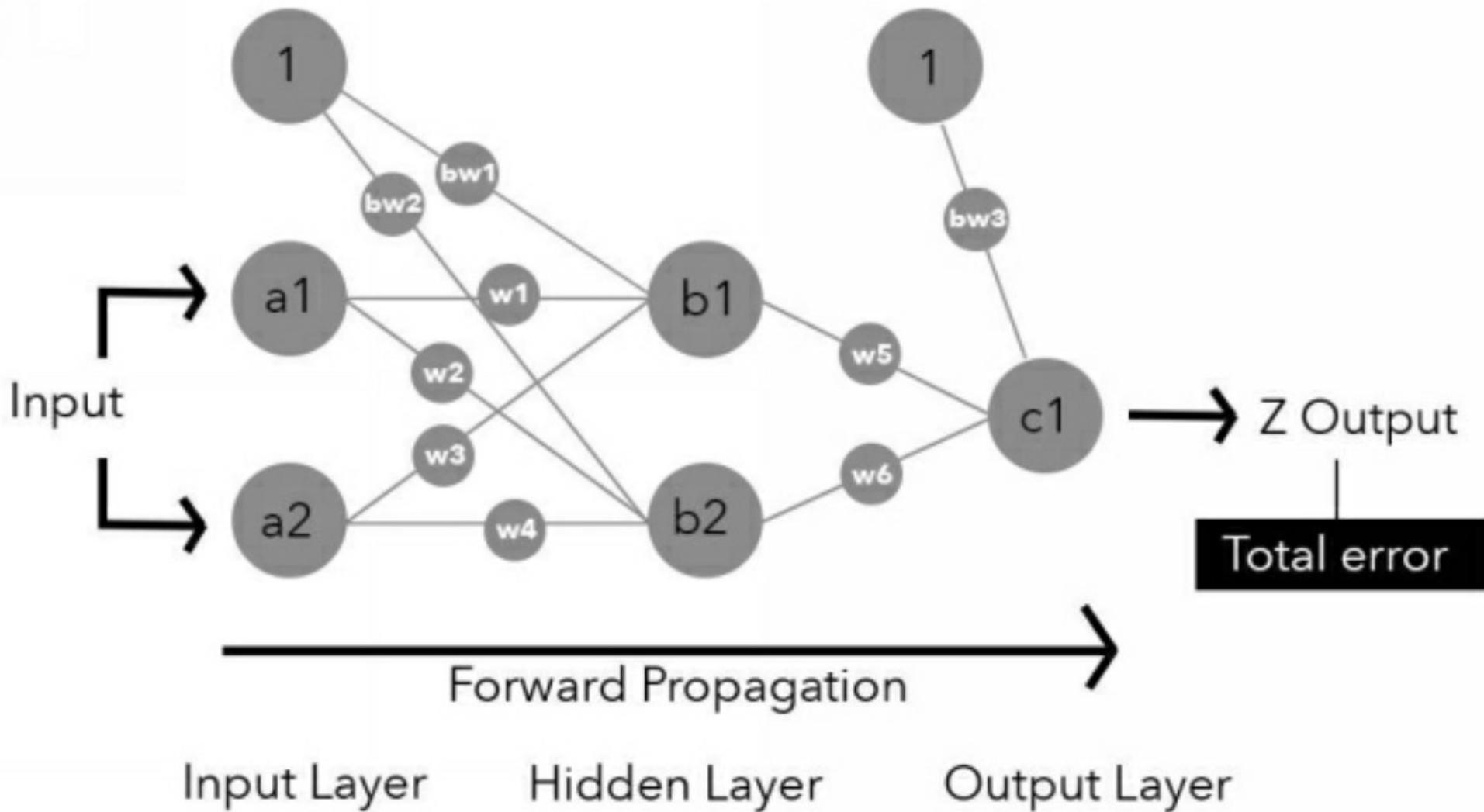
$$b_{2\text{net}} = (a_1 * w_2) + (a_2 * w_4) + \text{bias}$$

$$b_{2\text{net}} = (a_1 * w_2) + (a_2 * w_4) + (\text{bias} * b_{w2}) \text{ since bias is 1}$$

$$b_{2\text{net}} = (a_1 * w_2) + (a_2 * w_4) + b_{w2}$$

# Forward Propagation

- Total net input is then fed into an activation function
- Which transform the net input into new output
- This output sent over one or more edges and multiplies by a weight
- Cycle continues until the output layer calculations are completed



# Forward Propagation

- We successfully moved to output layer and computed Output
- Next step is to compute **total error** of the network
- Which will enable the network to adjust its weights and learn.
- Total Error is the difference between networks **actual output** and **target output**.

# Forward Propagation

- Once total error computed two things can happen
- Converge
  - The network has successfully converged(ie. Reached global minimum or local minum closed to global minimum)
  - At this point network ceases training
- Failed to Converge
  - The network has failed to converged(ie. global minimum not reached or local minimum not even close to global minimun)
  - Network continues to train

# Cost Functions

- Also referred as loss function or error function
- Cost function transforms everything that occurs in network into a number that represents **total error**
- It is measuring how wrong the network is
- Cost function is often referred as **objective function**
- This is because the objective of the network is to minimize the cost function

# Types of Cost Functions

- There are many types of cost functions some are:
  - Mean Squared Error
  - Squared Error
  - Root Mean Squared Error
  - Sum of Squares Error
  - Cross-Entropy
  - Exponential
  - Hellinger Distance
  - Kullback-Leibler Divergence

# Mean Squared Error(MSE)

- ❶ This sums and averages all of the squared output node errors.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (t_i - z_i)^2$$

- ❷ Total number of training examples.

- ❸ This squares the output node(s) error(s), which has multiple effects.

- ❹ The “ $t$ ” is the target output and the “ $z$ ” is the actual output of an output node. The “ $i$ ” refers to a unique value.

- ❺ This “ $i$ ” is called the “index of summation.” The numbers “1” and “ $n$ ” are the lower and upper limits of the summation. “ $n$ ” is equal to the number of training examples.

# Cost Functions

## Squared Error(SE)

$$SE = \frac{1}{2} \sum_{i=1}^n (t_i - z_i)^2$$

① Multiplied by 1/2  
instead of 1/n.

## Root Mean Square(RMS)

$$RMS = \sqrt{\frac{1}{n} \sum_{i=1}^n (t_i - z_i)^2}$$

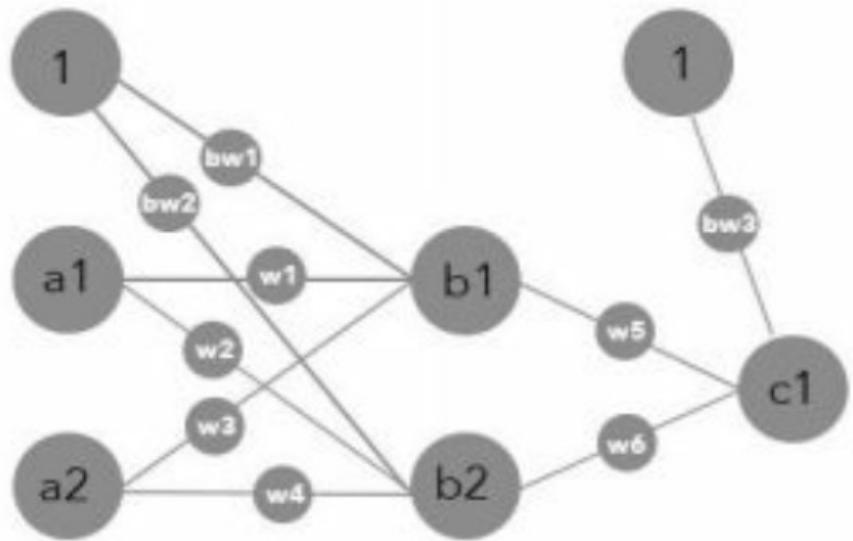
- Sum of Squares Error(SSE)

$$SSE = \sum_{i=1}^n (t_i - z_i)^2$$

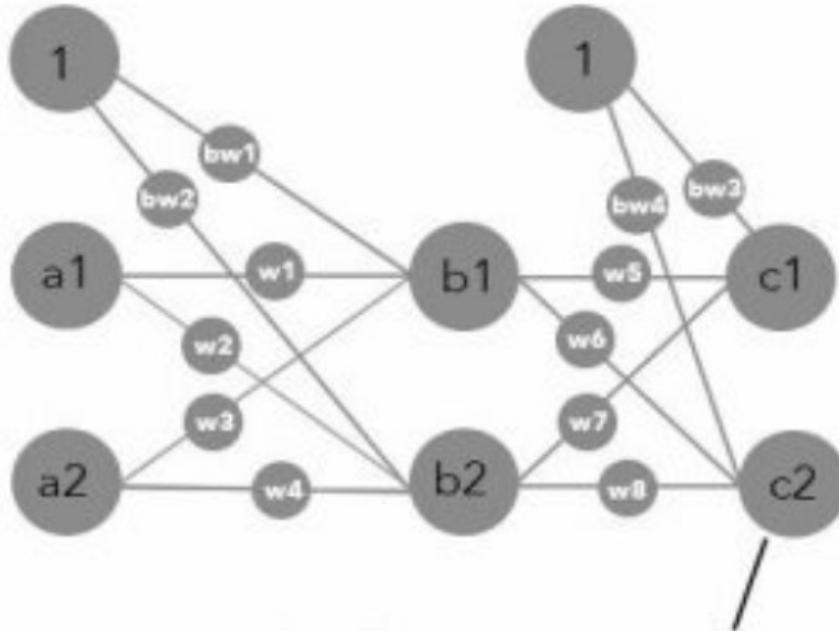
# Why Cost function?

- It provides the total error or difference between the **target output** and n/w **computed output**
- **1. Calculating the total error of each output node**
  - $(t_i - z_i)^2$
  - In order to find the total error the **local error** of each training example must be first calculated
  - **Local error** is the difference between a single training example's n/w **output** and **target output**
  - If there are multiple output layer nodes this means there are multiple outputs and target outputs
  - Hence local error calculated at each node and results are summed to create the **final local error for single training example**
  - $\sum(t_i - z_i)^2$

## Single Node Output



## Multiple Node Outputs



- 1 In this example there are multiple outputs. Therefore, the local error of each output is calculated, and the results are summed to create a final local error.

# Mean Squared Error(MSE)

- Single output(Single local Error)

$$e = (5-1)^2$$

- Actual output minus the target output

- Two output(Two local Error)

$$e1 = (5-1)^2$$

- The local error of node c1

$$e2 = (3-2)^2$$

- The local error of node c2

# Mean Squared Error(MSE)

- **2. Square each local error**
- Two uses with Squaring
  - **First**, Difference calculated is same whether it is positive or negetive
  - Avoid different signs from cancelling each other which can lead to a misreprasentation of wrong the network is
  - **Second**, squaring also helps the network converge faster

# Mean Squared Error(MSE)

- **3. Summing the local errors of all training examples**

$$\sum_{i=1}^n$$

- In this step all the final local errors of every training example are summed
- The **n** is training examples not individual nodes

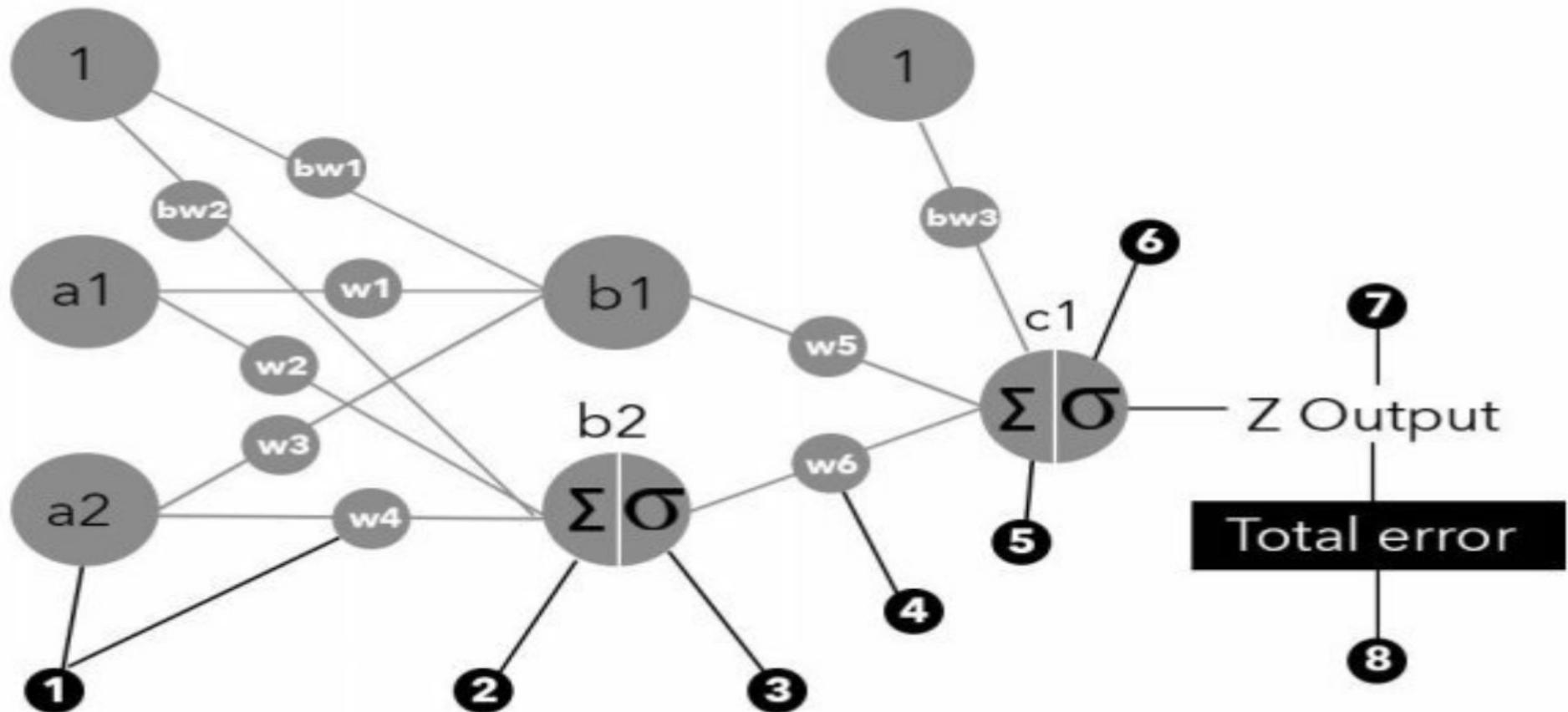
# Mean Squared Error(MSE)

- **4. Multiply and average**

$1/n$

- $1/n$  is divided by total number of training examples,  $n$
- The result is multiplied by sum of all local errors
- This normalizes the sum
- If there is a single training example no need to devide and multiply

# Forward Propagation : Big Picture



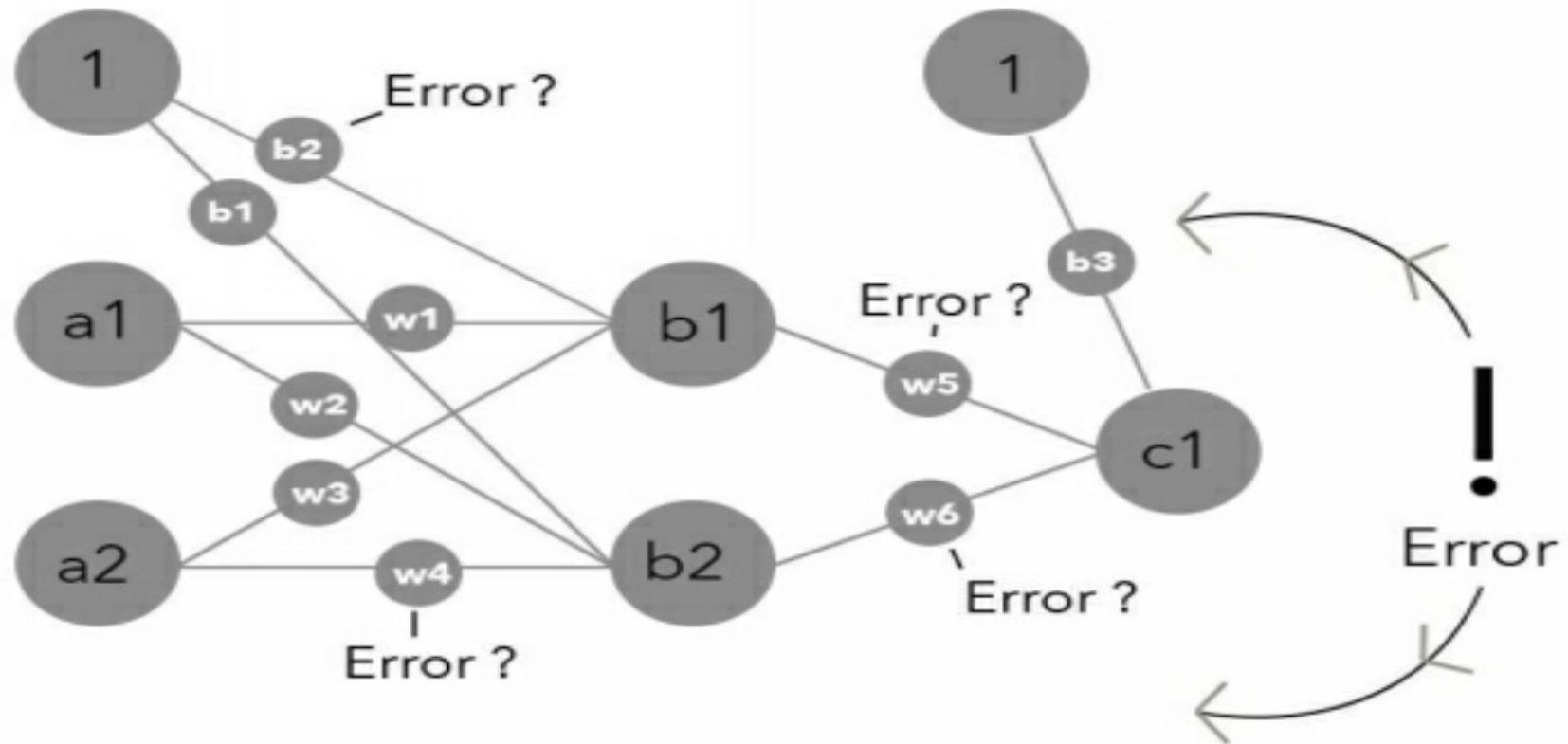
# Forward Propagation : Big Picture

- **Step1:** each input is multiplied by a weight as it travels over an edge to node in fallowing layer
- **Step2:** all inputs to node, including the bias are summed using summation operator
- **Step3-4:** output of **b2** multiplied by weight **w6**
- **Step5:** the result is summed with all other inputs to create the net input of node c, **netc**
- **Step6:** **Netc** becomes the input of activation function
- **Step7:** The activation fun output is the final outputbof the network
- **Step8:** the cost fun is applied to output of network, and a total errorbis caliculated.

# Back Propagation

- **Caliculating Gradients**
  - In previous step we caliculated total error of network
  - Find how this error is spred across every weight in the network
  - So that we can adjust the weights to minimize the error
  - To do this we are going to calculate the error of every weight in the network
  - Error of weight is technically gradient, and we will use the chain rule to caliculate this
  - The network try to minimize total error

# Back Propagation



# Mathematical Functions

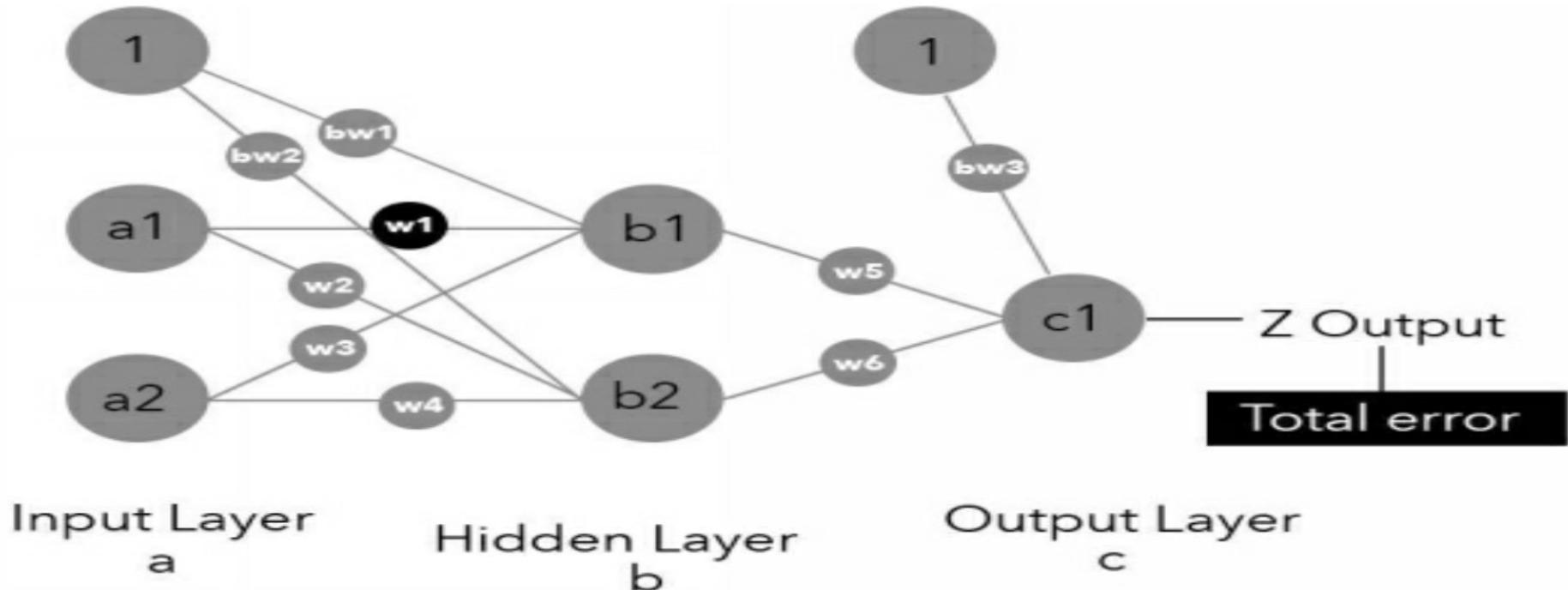
1) Partial Derivative

2) Chain Rule

- What is derivative?
  - Rate of change of a function at single point
- What is partial derivative?
  - It is the derivative of a function which has two or more variables **but with respect to only one variable**
  - All other variables are treated as constant.

# Partial Derivative

- Partial derivative enables you to measure how a single variable (out of many) impacts another single variable.



# Partial Derivative

- Take **w1** as example in above fig.
- The partial derivative allows us to discover how a change in weight **w1** affects the **total error**
- While all other weights remains constant
- The same can be calculated for **w2,w3,w4,bw1** etc.
- 

$$\frac{\partial f}{\partial x}$$

Used instead of "d" in usual  $\frac{df}{dx}$  notation to emphasize that this is a partial derivative.

# Partial Derivative

- Measuring a change in the **total error** by change in the weight **w1** would look like this

❶ The partial derivative of the **total error**. —————  $\frac{\partial E}{\partial w_1}$  ❷ The partial derivative of weight **w1**.

# Chain Rule

- The chain rule is used to find the partial derivative when an equation consists of a function inside another function.
- In Otherway It is used to differentiate the function of another function

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$$

# Gradient Importance

- Lets express our network in function

$$Z = \text{Sig}(w^2 \text{Sig}(w^1 a + b_1) + b_2)$$

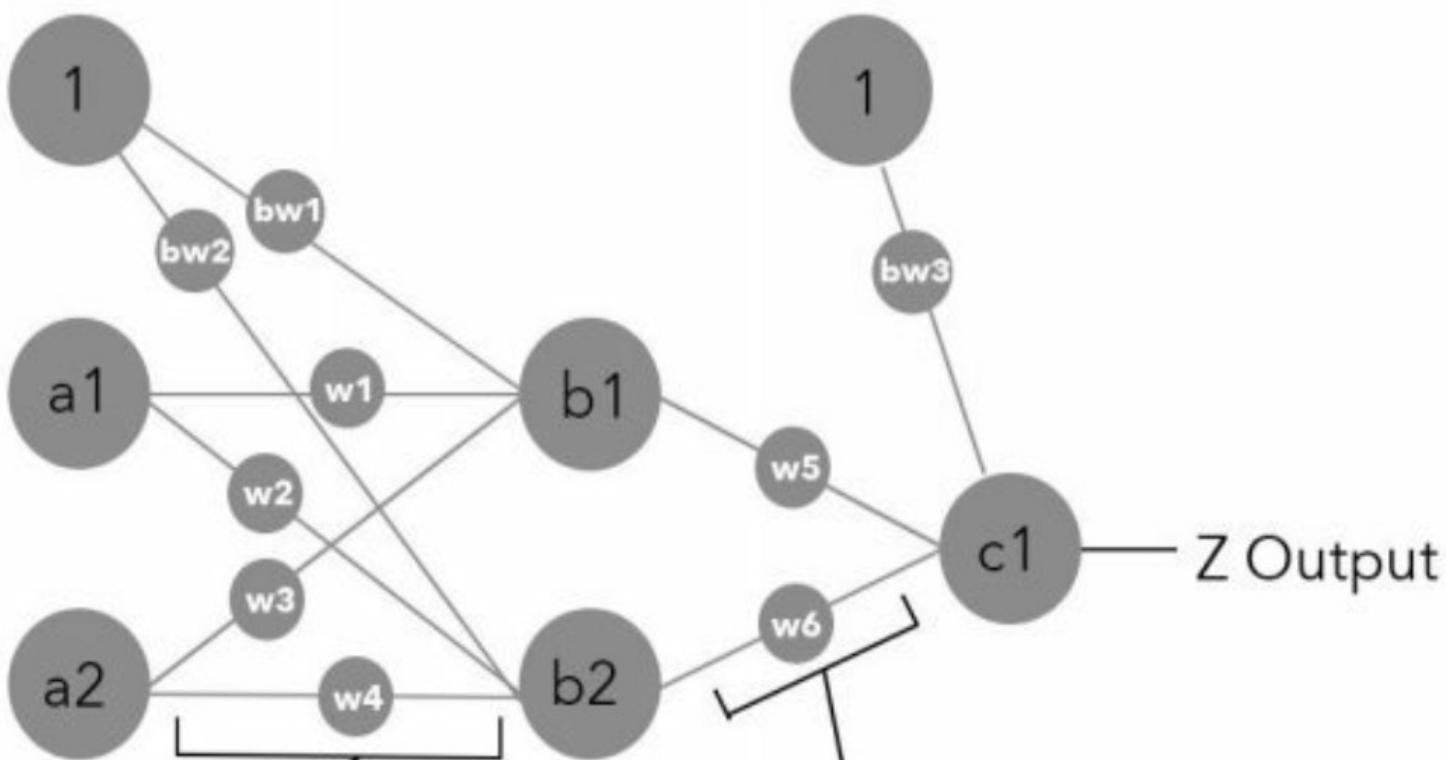
- Output of our network is product of weights and input
- The input is fixed so to minimize the total error our only option is to adjust the weights
- How much do we adjust weights?
- Do we increase or decrease weight?

# Gradient Importance

- What is the combination of weights that will produce a minimal **total error**?
- Answer is **Gradient Descent** which is an optimization method that helps us find the exact combination of weights and ultimately discover a **global minimum(total error)**.
- To perform **Gradient Descent** we need partial derivatives of each weight
- Then use the partial derivatives to update weights

# Gradient Descent

- There are two different type of weights
  - Hidden weights
  - Between input/hidden nodes
  - Output weights
  - Between hidden to output nodes
- How to perform Gradient Descent?



**1** The inside of the network. Weights here are not connected to an output node.

**2** The outside of the network. Weights here are connected to the output layer.

# Caliculating Gradients

- 1) Caliculating partial derivatives of output Layer Weights
- 2) Caliculating partial derivatives of output Layer Bias  
Weights
- 3) Caliculating partial derivatives of hidden Layer Weights
- 4) Caliculating partial derivatives of hidden Layer Bias  
Weights

# Partial derivatives of output weights

1 Partial derivative symbol.

$$\frac{\partial E}{\partial W_5}$$

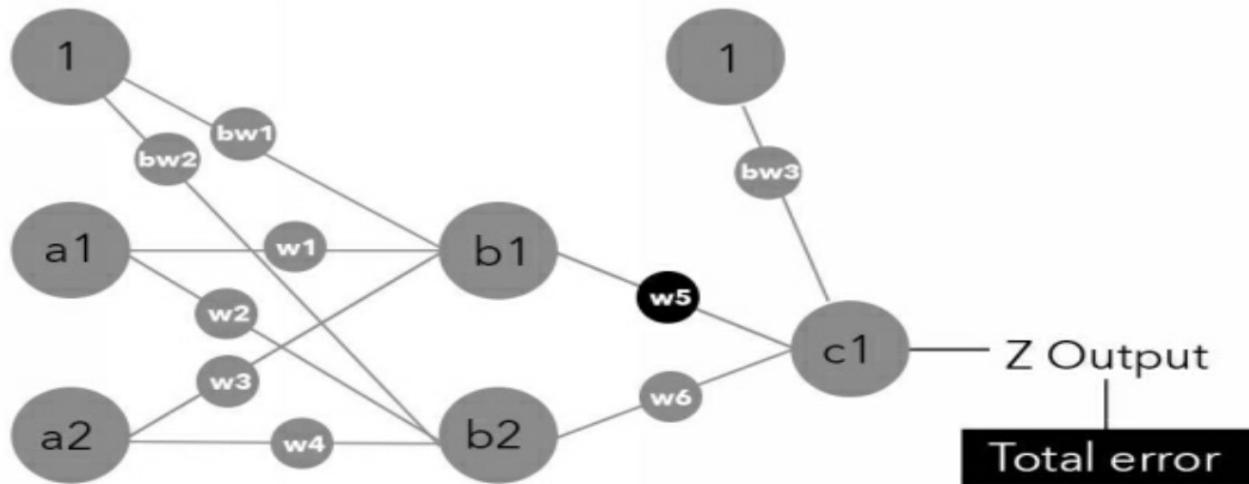
2 This "W" represents weight. It could be any other letter or number, but "W" is intuitive.

3 This "E" represents the total error of the network. It could be any other letter or number, but "E" is intuitive.

4 This "5" represents a specific weight that falls between the "b" and "c" layers. This is the weight we are finding the partial derivative for.

# Partial derivatives of output weights

- Partial derivative of the **total error** with respect to **w5**
- how much does change in in **w5** effect the **total error**



# Partial derivatives of output weights

$$\frac{\partial E}{\partial W_5} = \text{?} \rightarrow \text{Total Error}$$


# Partial derivatives of output weights

- But  $w_5$  not directly connected to **total error**
- $W_5$  is connected through a series of **partial derivatives** to the **total error**
- Here **total error** is computed by applying a **cost function** to the output  $z$
- So first, we need to know howmuch a change in the  $z$  affects the **total error**
- Second, we need to know howmuch a change in the  $w_5$  affects the  $z$

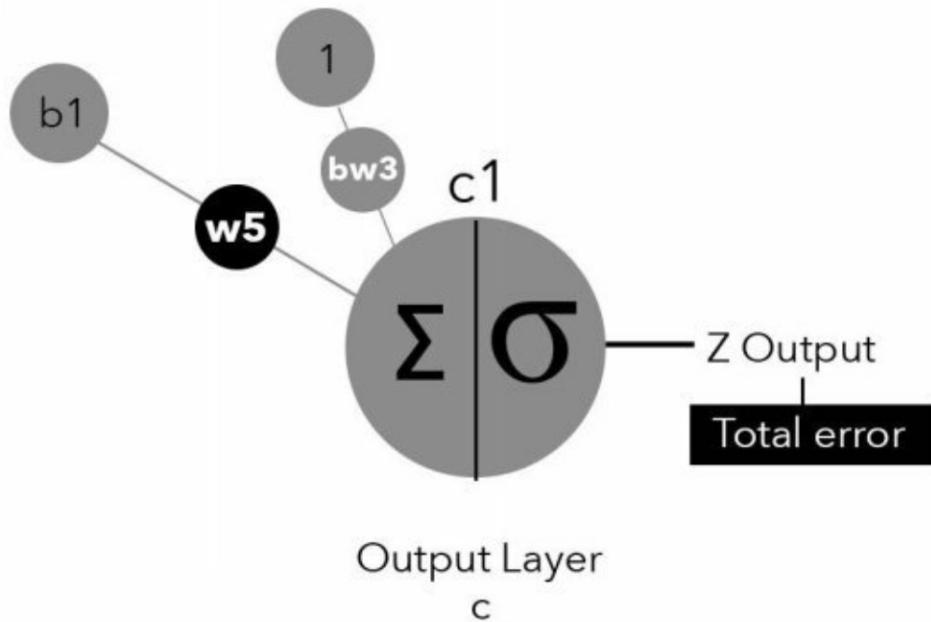
# Partial derivatives of output weights

- 1 The partial derivative of the **total error** with respect to the output **Z**.

$$\frac{\partial E}{\partial w_5} = \frac{\partial E}{\partial Z} \frac{\partial Z}{\partial w_5}$$

- 2 The partial derivative of the output **Z** with respect to weight **w5**.

# Partial derivatives of output weights



# Partial derivatives of output weights

- Output node **c1** is divided into two sections
- **Summation operator** which sums the net input into the node.
- **Logistic activation function** which takes output from summation operator and applies a function to it before passing it out as the final output  $z$
- There is still another intermediate variable inside of it
- Which tells us how the net input into node **c1** affect  $z$

# Partial derivatives of output weights

$$\frac{\partial Z}{\partial W_5} = \frac{\partial Z}{\partial Z_{net_c}} \frac{\partial Z_{net_c}}{\partial W_5}$$

- 1 The partial derivative of the total error with respect to the output  $\mathbf{z}$ .

$$\frac{\partial E}{\partial w_5} = \frac{\partial E}{\partial z} \frac{\partial z}{\partial net_c} \frac{\partial net_c}{\partial w_5}$$

- 3 The partial derivative of **netc** with respect to weight **w5..**

- 2 The partial derivative of  $\mathbf{z}$  with respect to **netc**.

# Calculating Partial Derivatives

$$\frac{\partial E}{\partial Z}$$

- 

$$MSE = \frac{1}{n} \sum_{i=1}^n (t_i - z_i)^2$$

$$\frac{\partial E}{\partial Z} = - (t - z)$$

1 This "t" represents the target output of a particular node.

2 This "z" represents the actual output of a particular node.

$$\frac{\partial E}{\partial Z} = (z - t)$$

# Calculating Partial Derivatives

$$\frac{\partial Z}{\partial \text{net}_c}$$

- 

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

# Calculating Partial Derivatives

$$\frac{\partial Z}{\partial \text{net}_c} = z(1 - z)$$

① Each "z" represents the actual output of a particular node.

# Calculating Partial Derivatives

$$\frac{\partial \text{net}_c}{\partial W_5}$$

- $\frac{\partial \text{net}_c}{\partial W_5} = \text{out}_{b1}$

# Calculating Partial Derivatives

$$\frac{\partial E}{\partial W_5} = (z-t) z(1 - z) \text{out}_{b1}$$

Annotations:

- ① Derivative 1: Points to the term  $(z-t)$ .
- ② Derivative 2: Points to the term  $z(1 - z)$ .
- ③ Derivative 3: Points to the term  $\text{out}_{b1}$ .

# Calculating Partial Derivatives

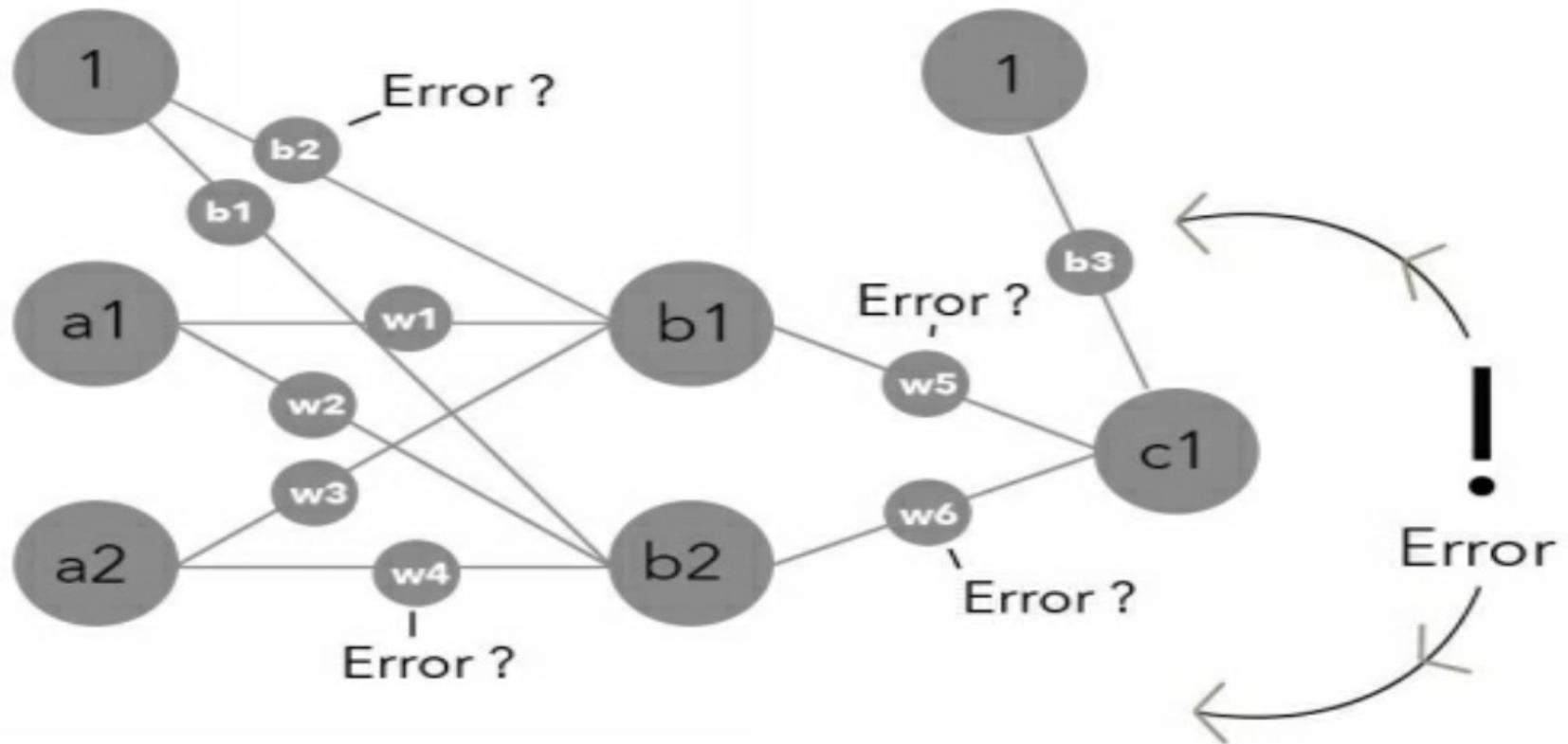
$$\delta_z = (z-t) z(1 - z)$$

$$\frac{\partial E}{\partial W_5} = \delta_z \text{out}_{b1}$$

① This  $\delta_z$  is the node delta for the output layer. See below for further explanation.

② This represents the output of the node **b1**.

# Back Propagation



# Caliculating partial derivatives of bias

- No outb1?
- Bias is no connected to previous layer

$$\frac{\partial E}{\partial B_{w3}} = (z-t) z(1 - z)$$

# Calculating partial derivatives of hidden weights

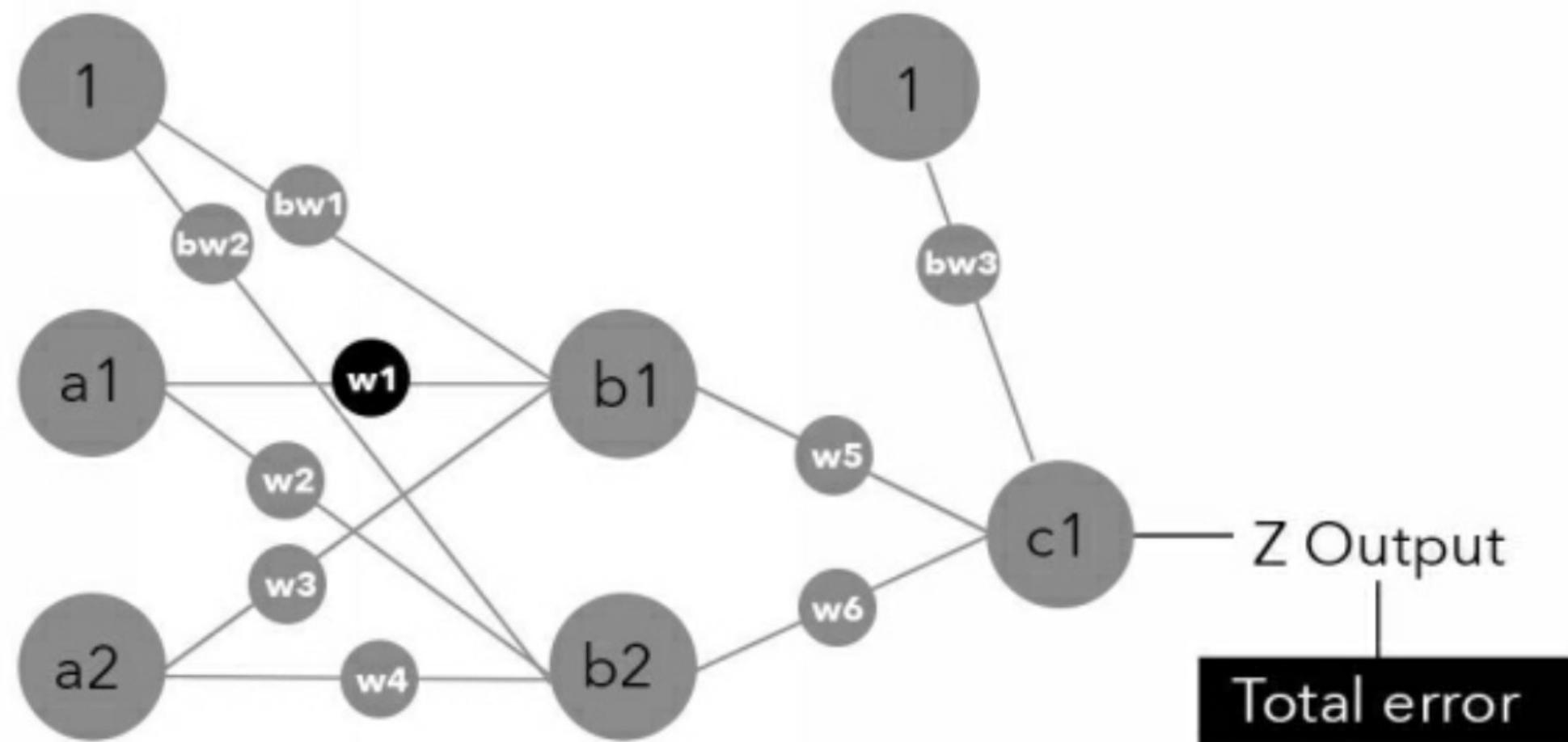
1 Partial derivative symbol.

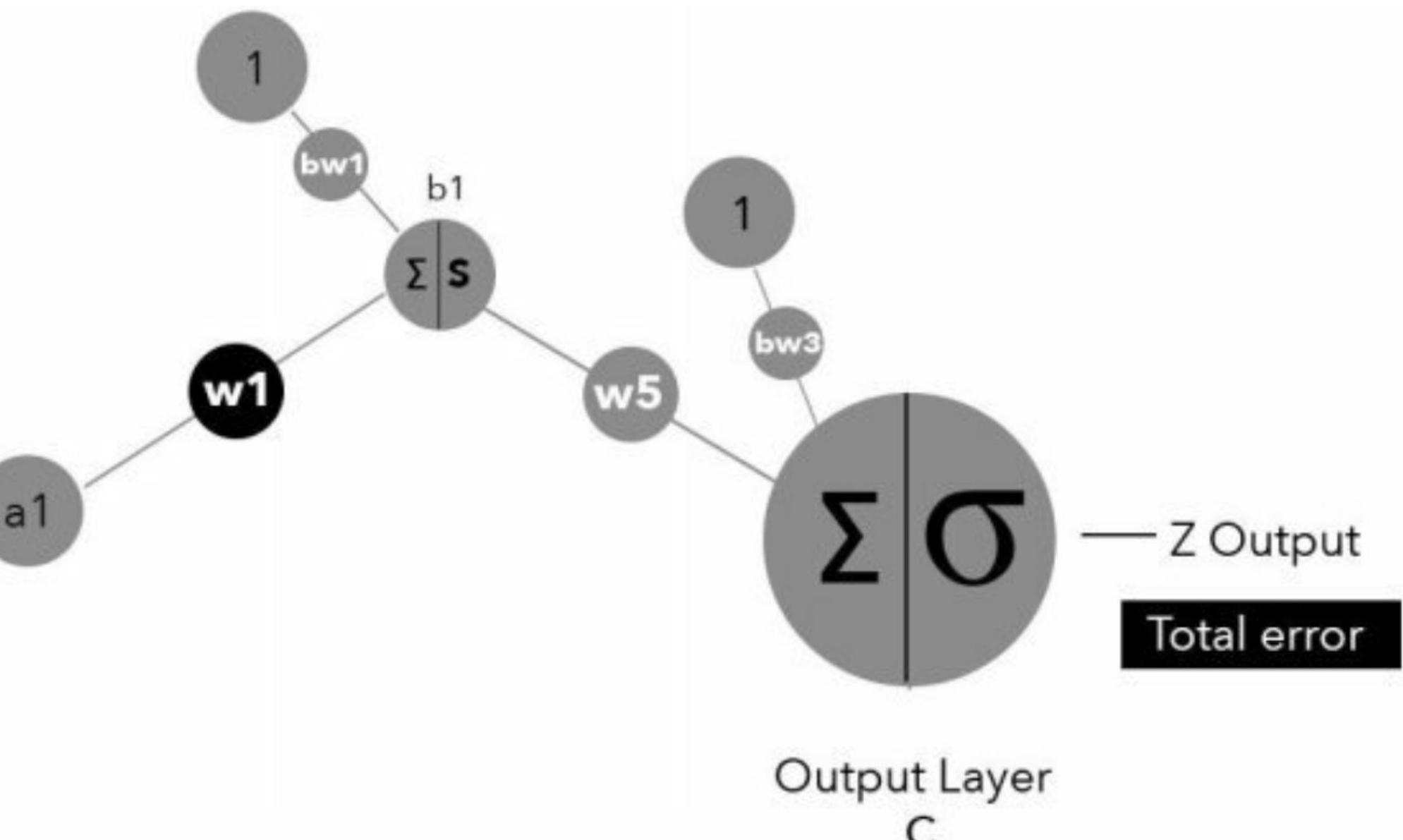
$$\frac{\partial E}{\partial W_1}$$

2 This "W" represents weight. It could be any other letter or number, but "W" is intuitive.

3 This "E" represents the total error of the network. It could be any other letter or number, but "E" is intuitive.

4 This "1" represents a specific weight that falls between the "a" and "b" layers. This is the weight we are finding the partial derivative for.





- ① The partial derivative of the **total error** with respect to **netb1**.

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial \text{net}_{b1}} \frac{\partial \text{net}_{b1}}{\partial w_1}$$

- ② The partial derivative of **netb1** with respect to weight **w1**.

# Caliculating partial derivatives of hidden weights

$$\frac{\partial E}{\partial \text{net}_{b1}} = \frac{\partial E}{\partial \text{net}_{c1}} \frac{\partial \text{net}_{c1}}{\partial \text{net}_{b1}}$$

$$\frac{\partial \text{net}_{c1}}{\partial \text{net}_{b1}} = \frac{\partial \text{net}_{c1}}{\partial \text{out}_{b1}} \frac{\partial \text{out}_{b1}}{\partial \text{net}_{b1}}$$

$$\frac{\partial E}{\partial W_1} = \frac{\overset{\textcircled{1}}{\partial E}}{\partial \text{net}_{c1}} \frac{\overset{\textcircled{2}}{\partial \text{net}_{c1}}}{\partial \text{out}_{b1}} \frac{\overset{\textcircled{3}}{\partial \text{out}_{b1}}}{\partial \text{net}_{b1}} \frac{\overset{\textcircled{4}}{\partial \text{net}_{b1}}}{\partial W_1}$$

# Caliculating partial derivatives of hidden weights

$$\frac{\partial E}{\partial \text{net}_{c1}}$$

# Caliculating partial derivatives of hidden weights

$$\frac{\partial E}{\partial \text{net}_{c1}} = \delta_z$$

$$\delta_z = (z-t) z(1 - z)$$

# Caliculating partial derivatives of hidden weights

$$\frac{\partial E}{\partial \text{net}_{c1}} = \delta_z$$

$$\delta_z = (z-t) z(1 - z)$$

# Caliculating partial derivatives of hidden weights

$$\frac{\partial \text{net}_{c1}}{\partial \text{out}_{b1}}$$

$$\frac{\partial \text{net}_{c1}}{\partial \text{out}_{b1}} = W_5$$

# Caliculating partial derivatives of hidden weights

$$\frac{\partial \text{out}_{b1}}{\partial \text{net}_{b1}}$$

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

# Caliculating partial derivatives of hidden weights

$$\frac{\partial \text{out}_{b1}}{\partial \text{net}_{b1}} = \text{out}_{b1}(1 - \text{out}_{b1})$$

# Caliculating partial derivatives of hidden weights

$$\frac{\partial \text{net}_{b1}}{\partial W_1}$$

$$\frac{\partial \text{net}_{b1}}{\partial W_1} = \text{out}_{a1}$$

# Caliculating partial derivatives of hidden weights

$$\frac{\partial E}{\partial W_1} = \left( \sum_c \delta_z W_5 \right) \overbrace{out_{b1}(1 - out_{b1})}^{Derivative\ 3} \overbrace{out_{a1}}^{Derivative\ 4}$$

Diagram illustrating the components of the derivative calculation:

- ① Derivative 1:  $\sum_c \delta_z W_5$
- ② Derivative 2:  $out_{b1}(1 - out_{b1})$
- ③ Derivative 3:  $out_{a1}$
- ④ Derivative 4:  $out_{b1}(1 - out_{b1})out_{a1}$

# Caliculating partial derivatives of hidden bias

$$\frac{\partial E}{\partial B_{w1}} = \left( \sum_c \delta_z W_5 \right) \text{out}_{b1} (1 - \text{out}_{b1})$$

# Weight updation

- ① The old weight **w5**, which is being updated.

$$W_5_{\text{new}} = W_5 - \eta * \frac{\partial E}{\partial W_5}$$

- ③ The partial derivative of **the total error** with respect to the weight **w5**.

- ② The greek letter eta represents the learning rate. Other symbols often used include alpha and epsilon.