# Deep Learning

**Feed Forward Neural Networks**
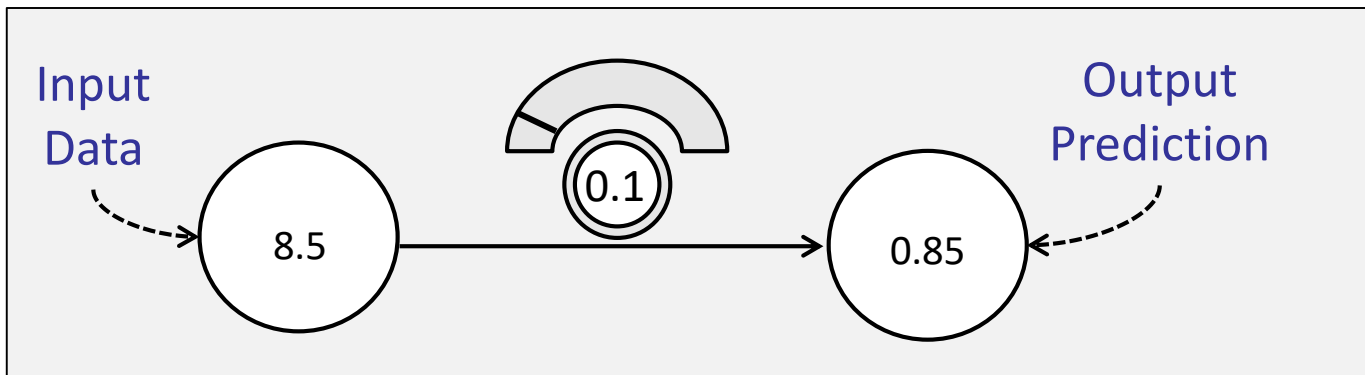
Tuesday 5th November

Dr. Nicholas Cummins

- **Neural Networks**
  - **Weights represent knowledge**
    - A measure of sensitivity between the input and the prediction
    - It uses *knowledge* captured in the weights to *interpret* the input data to *predict* a certain outcome



```python
weight = 0.1
def neural_network(input, weight):
    prediction = input * weight
    return prediction
number_of_offsides = [8.5, 9.5, 10, 9]
input = number_of_offsides[0]
pred = neural_network(input,weight)
print(pred)
```

- **Updating weights to make predictions more accurate**
  - **Compare using a loss function**
    - Evaluate how well the network performed

    ```
    error = ((input * weight)- goal_pred) ** 2
    ```

  - **Learn weights via Gradient Descent**
    - Adjusting each weight to reduce the error
    - Using the derivate of `weight` and `error` relationship defined through the loss function to adjust the weights

    ```
    weight = weight - (alpha*derivative)
    ```

- **Gradient Descent**
    - <span style="color:red">Learning is adjusting the weight to reduce the `error` to 0</span>
        - The function conforms to the patterns in the data

        ```
        error = (pred - goal_pred) ** 2
        ```
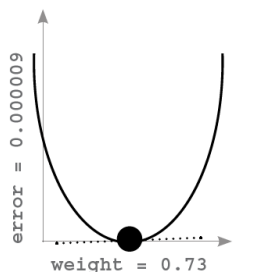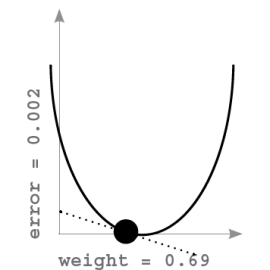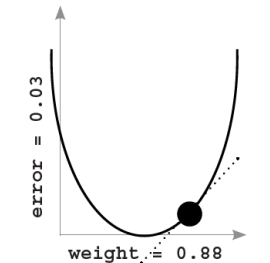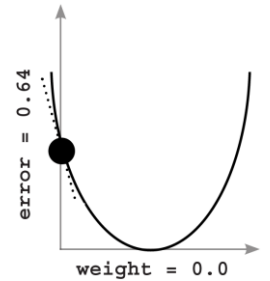
    - Derivative of the error functions defines the amount that error changes when you change weight

        ```
        derivative = input * (pred - goal_pred)
        ```

    - How can we use the derivative to minimise the error?
        - Adjust the weights in the opposite direction of the derivative
            1. Calculate the derivative of weight with respect to error
            2. Change weight in the opposite direction of that slope.

error = 0.64
weight = 0.0

error = 0.03
weight = 0.88

error = 0.002
weight = 0.69

error = 0.000009
weight = 0.73

- **Gradient Descent for a single network**

```
pred = input * weight

error = (pred - goal_pred) ** 2

derivative = input * (pred - goal_pred)

weight = weight - (alpha * derivative)
```

- `error` is a measure of how much the network missed by
  - We define `error` to be always positive
- `derivative` is the derivate of `weight` and `error`
  - Predicts both direction and amount to adjust the weights
- `Alpha` scales the weight update
  - Helps minimise divergence effects when the input is large

Image source: https://medium.com

- **Gradient Descent**
  - Networks with multiple inputs
  - Networks with multiple outputs
  - Networks with multiple inputs and outputs
- **Correlation**
  - Learning Correlation
  - Creating Correlation
- **Backpropagation**
- **Summary**

Image source: https://medium.com

- **Gradient Descent**
  - **Networks with multiple inputs**
  - Networks with multiple outputs
  - Networks with multiple inputs and outputs
- Correlation
  - Learning Correlation
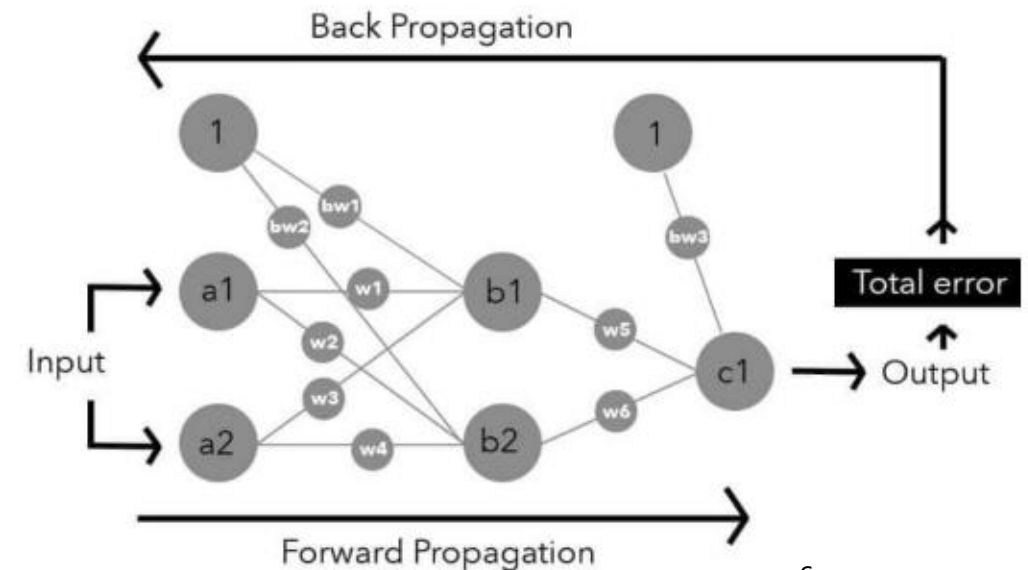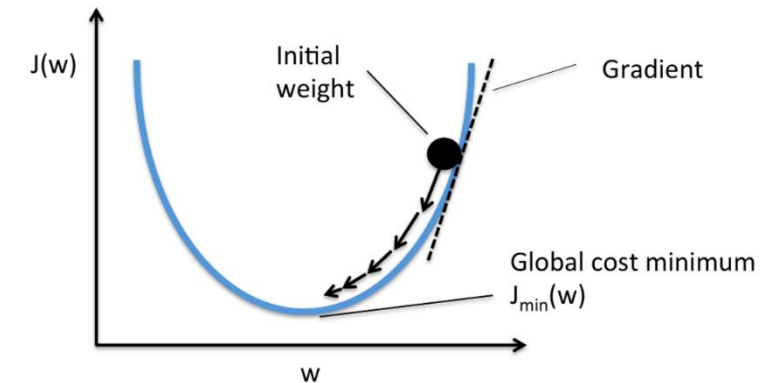  - Creating Correlation
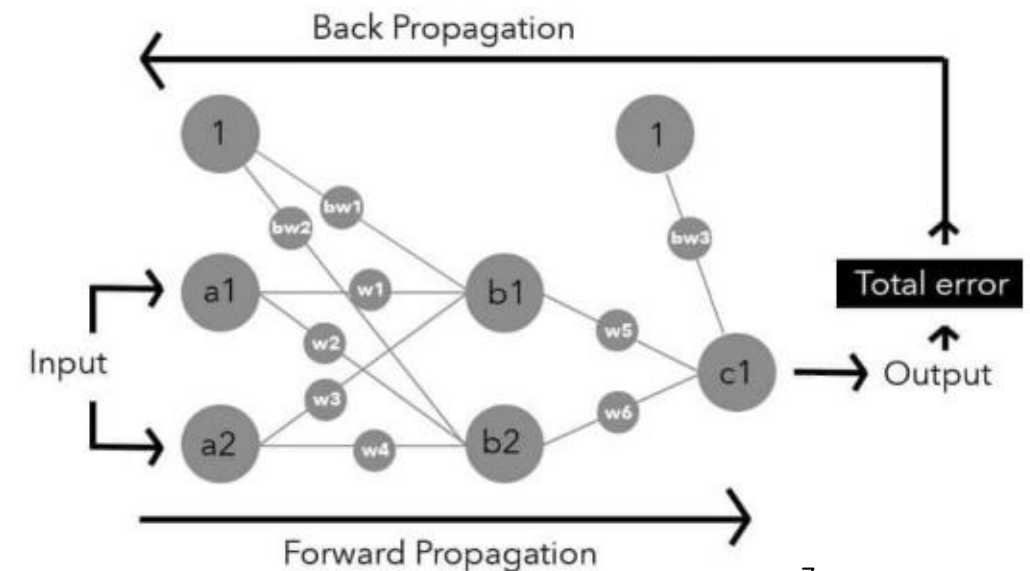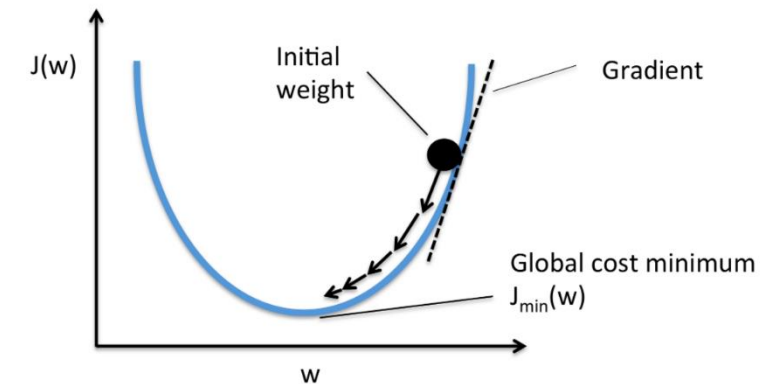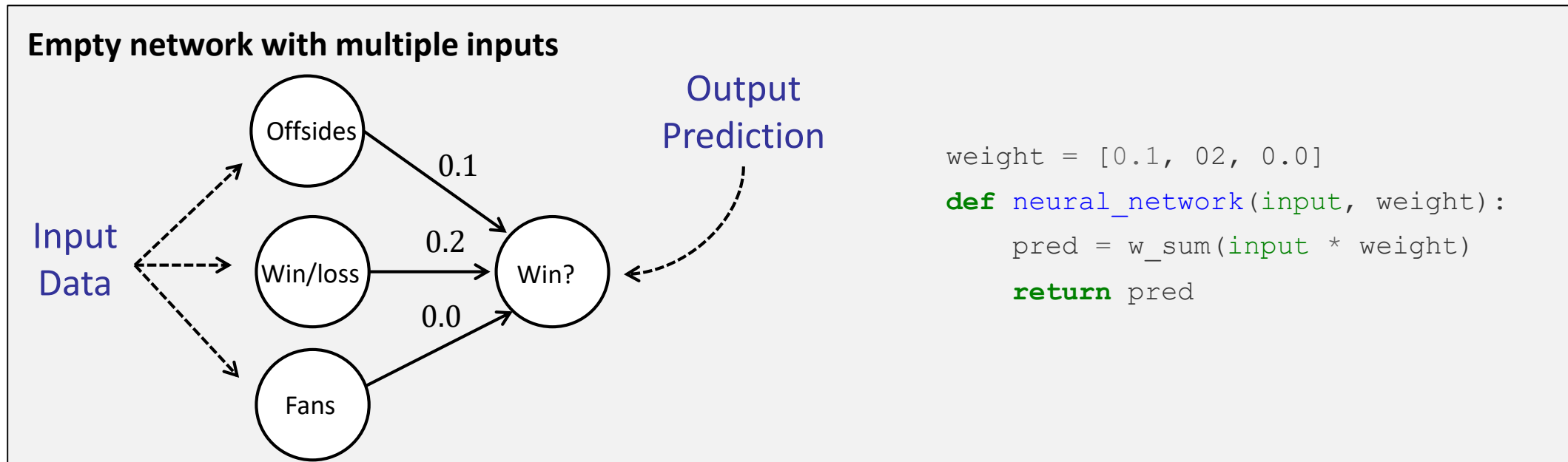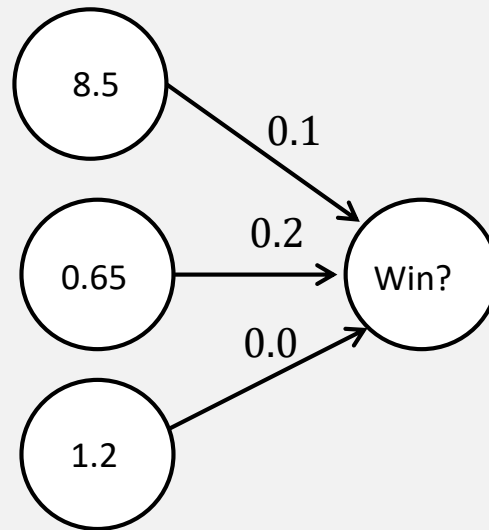- Backpropagation
- Summary

- **Neural networks can combine intelligence from multiple data points**
  - This allows the network to combine various forms of information to make better-informed decisions

**Empty network with multiple inputs**



```
weight = [0.1, 02, 0.0]
def neural_network(input, weight):
    pred = w_sum(input * weight)
    return pred
```

# Multiple each input by its *own* weight

- Fundamental weighting mechanism has not changed
- New property is that there are multiple weights
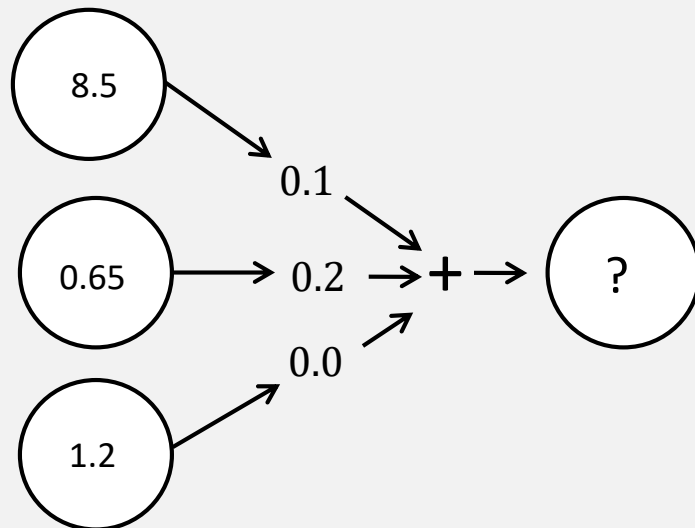
**Inserting a vector**



```python
def w_sum(a,b):
    assert(len(a) == len(b))
    output = 0
    for i in range(len(a)):
        output += (a[i] * b[i])
    return output
```

# • **Output is a weighted sum**

– Output gained by summing the respective predictions

• Weighted sum of the input, achieved via the dot product
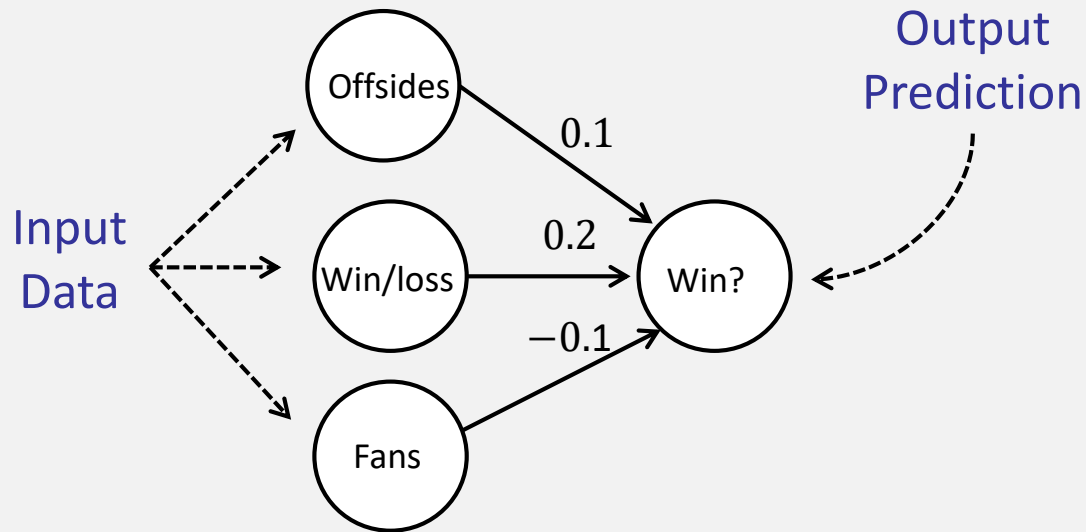
**Performing weighted sum**

| Inputs | | Weights | | Local Predictions | | |
|---|---|---|---|---|---|---|
| (8.5 | * | 0.1) | = | 0.85 | = | Offside prediction |
| (0.65 | * | 0.2) | = | 0.13 | = | Win/loss prediction |
| (1.20 | * | 0.0) | = | 0.00 | = | Number of fans prediction |

| Offside prediction | + | Win/loss prediction | + | Number of fans prediction | = | Final Prediction |
|---|---|---|---|---|---|---|
| 0.85 | + | 0.13 | + | 0.0 | = | 0.98 |

8.5
0.65
1.2
0.1
0.2
0.0
+
?

- **Gradient descent also works with multiple inputs**
  - The same techniques for updating single networks can be used to update a network containing multiple weights.

**Empty network with multiple inputs**
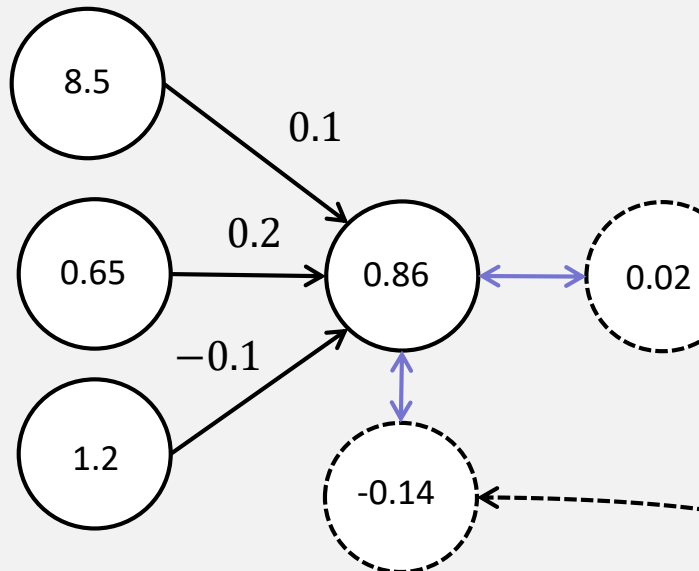


```
def w_sum(a,b):
    assert(len(a) == len(b))
    output = 0
    for i in range(len(a)):
        output += (a[i] * b[i])
    return output


weights = [0.1, 0.2, -0.1]
def neural_network(input, weight):
    pred = w_sum(input * weight)
    return pred
```

- **Basic compare steps have not changed**
  - Make a prediction, and calculate error and delta

**Output the prediction, calculate error & delta**



```
offsides = [8.5, 9.5, 9.9, 9.0]
wlrec = [0.65, 0.8, 0.8, 0.9]
nfans = [1.2, 1.3, 0.5, 1.0]
win_or_lose_binary = [1, 1, 0, 1]


true = win_or_lose_binary[0]

input = [offsides[0],wlrec[0],nfans[0]]

pred = neural_network(input,weight)

error = (pred - true) ** 2
delta = pred - true
```
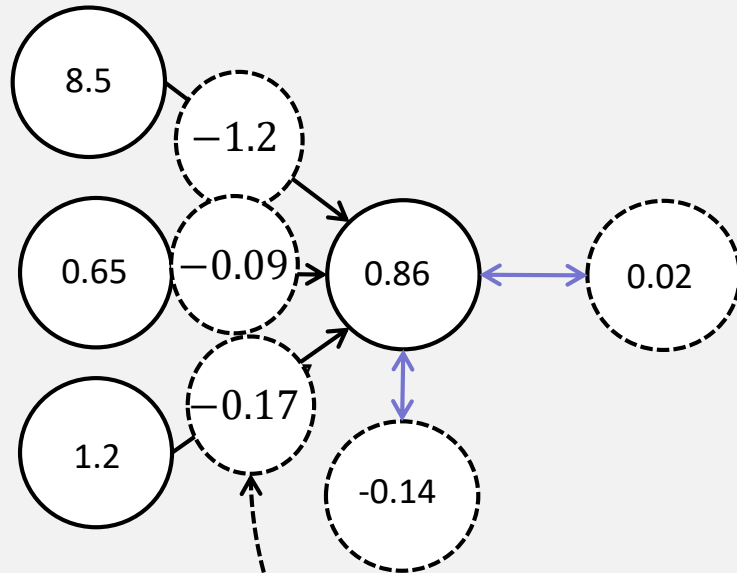
- ## **Basic learn steps have not changed**
  - Calculate a `weight_delta` and update each weight

**Output the prediction, calculate error & delta**



```
def ele_mul(number,vector):
    output = [0,0,0]
    assert(len(output) == len(vector))
    for i in range(len(vector)):
        output[i] = number * vector[i]
    return output
pred = neural_network(input,weight)
error = (pred - true) ** 2
delta = pred - true
weight_deltas = ele_mul(delta,input)
```
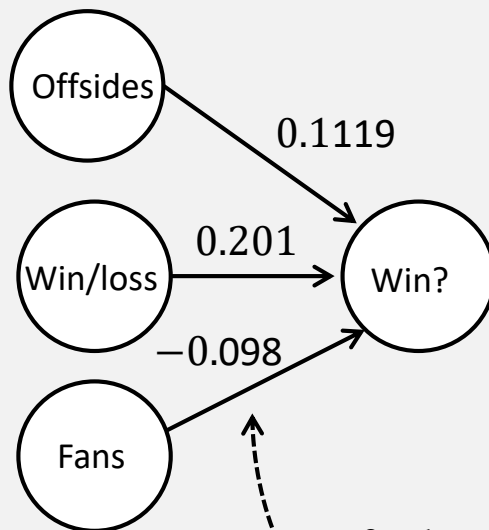
```
8.5 * -0.14 = -1.19 = weight_deltas[0]
0.65 * -0.14 = -0.091 = weight_deltas[1]
1.2 * -0.14 = -0.168 = weight_deltas[2]
```

- **Basic learn steps have not changed**
  - Last step, update the weights



**Update the weights**

```
input = [offsides[0],wlrec[0],nfans[0]]
pred = neural_network(input,weight)
error = (pred - true) ** 2
delta = pred - true
weight_deltas = ele_mul(delta,input)


alpha = 0.01

for i in range(len(weights)):
        weights[i] -= alpha * weight_deltas[i]
```

```
0.1 - (-1.19 * 0.01) = 0.1119 = weights[0]
0.2 - (-.091 * 0.01) = 0.2009 = weights[1]
-0.1 - (-.168 * 0.01) = -0.098 = weights[2]
```

- **Single `delta` to multiple `weight_delta` variables**
  - Reminder of their roles and purposes
  - **`delta`**: A measure of how much higher or lower a node's output value should be
    - Computed by a direct subtraction between the node's predicted value and the node's true value
  - **`weighted_delta`**: A derivative-based estimate of the direction and amount we need to move a `weight` by
    - Computed by multiplying `delta` by a weight's `input`
    - Accounts for effects of scaling, negative reversal, and stopping

- **Single `delta` to multiple `weight_delta` variables**
  - Reminder of their roles and purposes
  - Consider a single weight update

**delta:** Hey, inputs. Next time, predict a little higher.

**Single weight:**
- **Stopping** *if* my input was 0, then my weight wouldn't have mattered, and I wouldn't change a thing
- **Negative Reversal:** *if* my input was negative, then I'd want to decrease my weight instead of increase it
- **Scaling:** *my current input* is positive and quite large, so my personal prediction was important  I'm going to move my weight up a lot to compensate

$0.1 \rightarrow 1.29 (* alpha)$

8.5

0.65    0.2    0.86    0.02

1.2    $-0.1$    -0.14

- **Single `delta` to multiple `weight_delta` variables**
  - Each weight has a unique `input` but a common `delta`
  - Each respective `weight_delta` is calculated by multiplying the unique input multiplied by delta
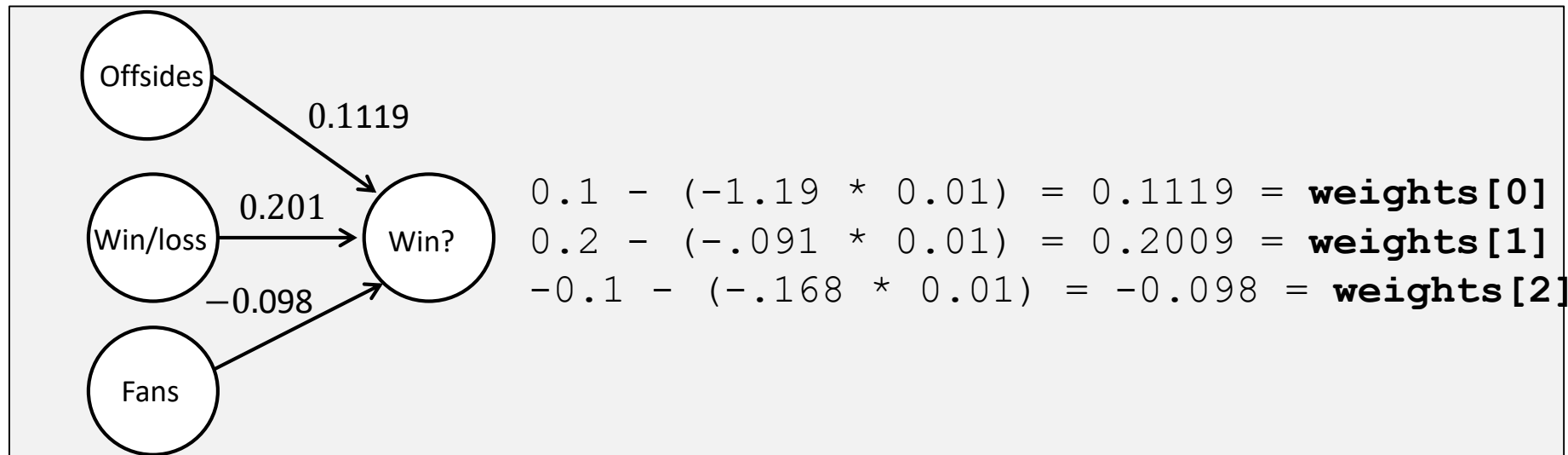


```
8.5 * -0.14 = -1.19 = weight_deltas[0]
0.65 * -0.14 = -0.091 = weight_deltas[1]
1.2 * -0.14 = -0.168 = weight_deltas[2]
```

- **Multiple weight updates**
  - Identical procedure to the single-input network,
    - Except performed over multiple weights
  - Once the individual `weight_delta` values have been calculated, multiply by `alpha` and subtract from the `weight`



```
0.1 - (-1.19 * 0.01) = 0.1119 = weights[0]
0.2 - (-.091 * 0.01) = 0.2009 = weights[1]
-0.1 - (-.168 * 0.01) = -0.098 = weights[2]
```

# Multiple iterations of learning

```python
def neural_network(input, weights):
    out = 0
    for i in range(len(input)):
        out += (input[i] * weights[i])
    return out

def ele_mul(scalar, vector):
    out = [0,0,0]
    for i in range(len(out)):
        out[i] = vector[i] * scalar
    return out

offsides = [8.5, 9.5, 9.9, 9.0]
wlrec = [0.65, 0.8, 0.8, 0.9]
nfans = [1.2, 1.3, 0.5, 1.0]
win_or_lose_binary = [1, 1, 0, 1]

true = win_or_lose_binary[0]
alpha = 0.01
weights = [0.1, 0.2, -.1]
input = [offsides[0],wlrec[0],nfans[0]]

for iter in range(3):
    pred = neural_network(input,weights)
    error = (pred - true) ** 2
    delta = pred - true
    weight_deltas=ele_mul(delta,input)
    for i in range(len(weights)):
        weights[i]-=alpha*weight_deltas[i]
```
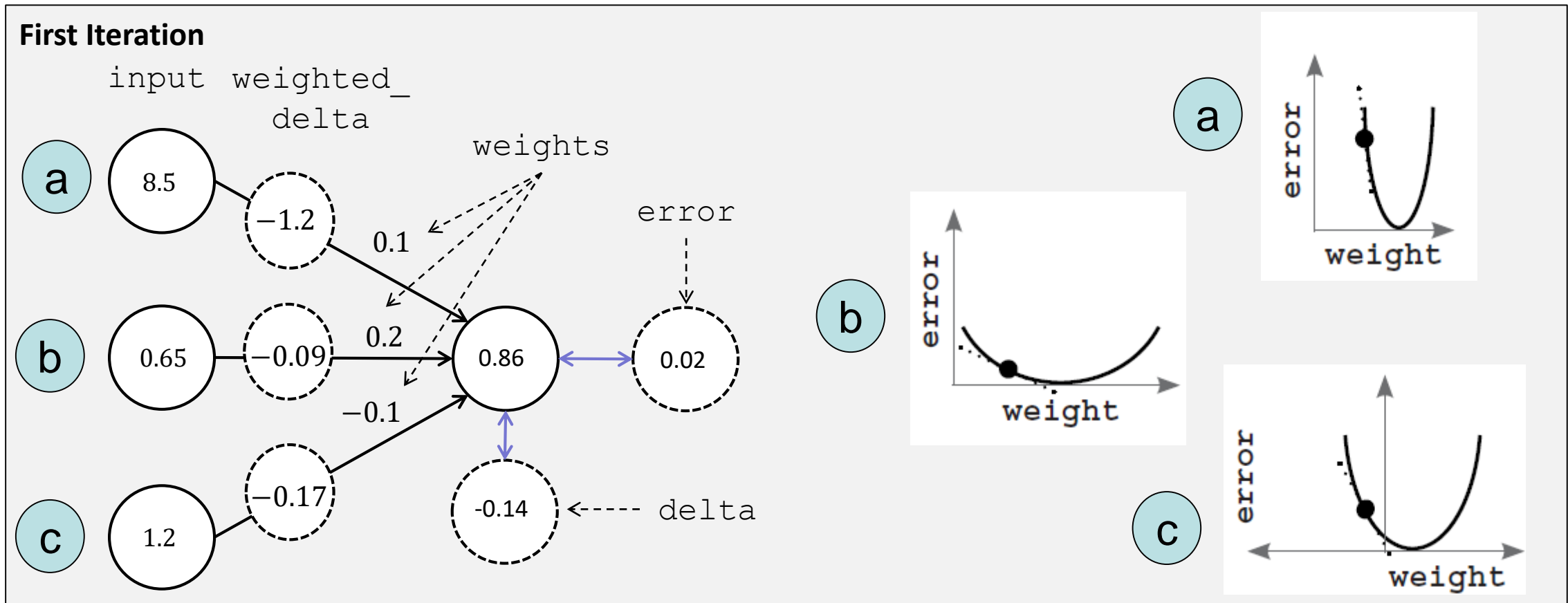
- **Multiple iterations of learning**

# Multiple iterations of learning



**First Iteration**

input weighted_delta

weights

a 8.5 −1.2 0.1

error

b 0.65 −0.09 0.2 0.86 ⟷ 0.02

−0.1

c 1.2 −0.17 -0.14 ⟵---- delta

a

b

c

# Multiple iterations of learning

- **Multiple iterations of learning**

Image source: https://medium.com
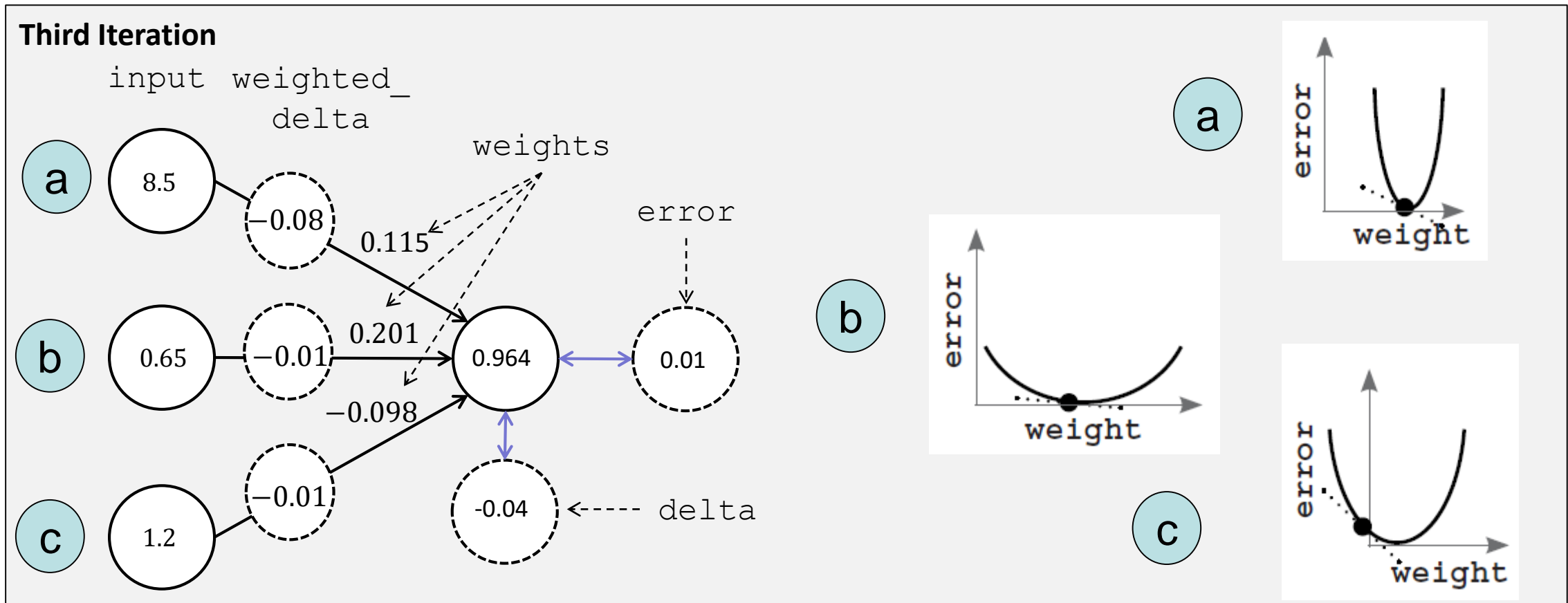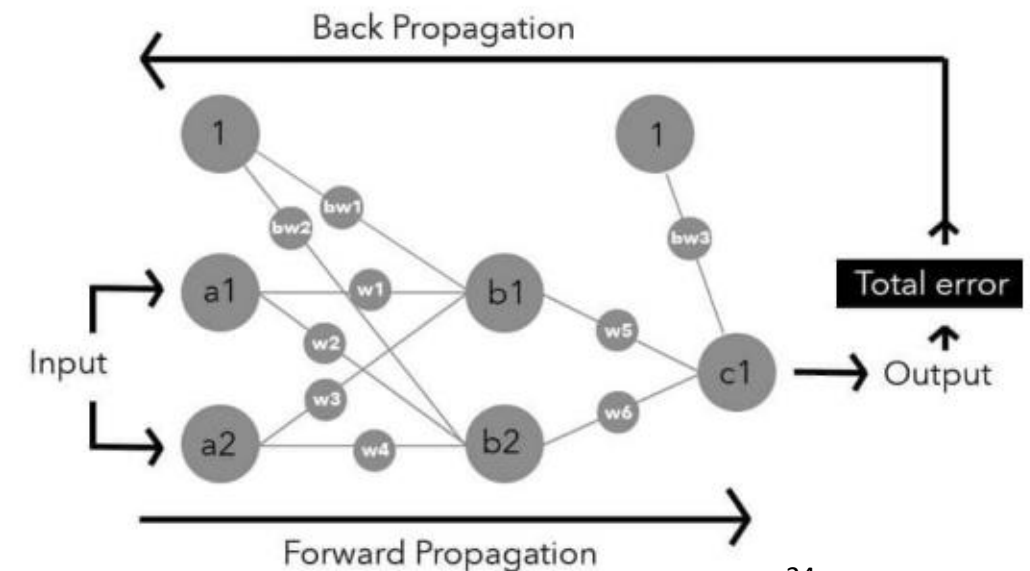
- # Gradient Descent

  - Networks with multiple inputs

  - ## Networks with multiple outputs

  - Networks with multiple inputs and outputs

- Correlation

  - Learning Correlation

  - Creating Correlation

- Backpropagation

- Summary

# • **Single input → multiple predictions**

– Prediction occurs the same as if there were multiple disconnected single-weight neural networks

**Empty network with multiple outputs**

Input Data ---→ Win/loss

0.3 → Injuries?

0.2 → Win?

0.9 → Happy?

Output Predictions

```
weight = [0.3, 0.2, 0.9]
def neural_network(input, weight):
    pred = ele_mul(input * weight)
    return pred
```

- **Independence**
  - Network behaves as three independent components, each receiving the same input data

**Elementwise multiplication**



```python
def ele_mul(number,vector):
    output = [0,0,0]
    assert(len(output) == len(vector))
    for i in range(len(vector)):
        output[i] = number * vector[i]
    return output
```

| Inputs | | Weights | | Final Predictions | | |
|--------|---|---------|---|-------------------|---|---|
| (0.65 | * | 0.3) | = | 0.195 | = | Injury prediction |
| (0.65 | * | 0.2) | = | 0.13 | = | Win prediction |
| (0.65 | * | 0.9) | = | 0.585 | = | Happy prediction |

# Single input → multiple predictions

– Making individual predictions and calculate individual `error` and `delta`



```
input = wlrec[0]
true = [hurt[0], win[0], sad[0]]
pred = neural_network(input,weights)

error = [0, 0, 0]
delta = [0, 0, 0]
for i in range(len(true)):
    error[i] = (pred[i] - true[i]) ** 2
    delta[i] = pred[i] - true[i]
```

- **Single input → multiple predictions**
  - Calculate each `weighted_delta` and update weights
    - As before the `weight_delta` values are computed by multiplying the input value with the node delta for each weight



```python
def scalar_ele_mul(number,vector):
    output = [0,0,0]
    assert(len(output) == len(vector))
    for i in range(len(vector)):
        output[i] = number * vector[i]
    return output

weight_deltas = scalar_ele_mul(input,weights)
```
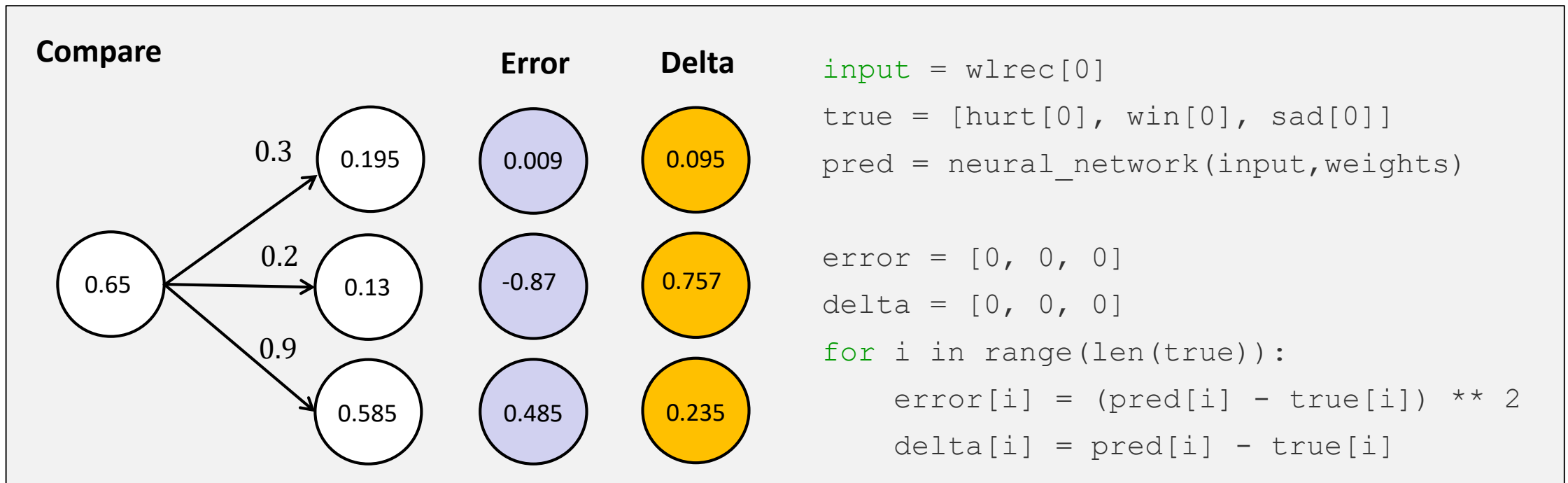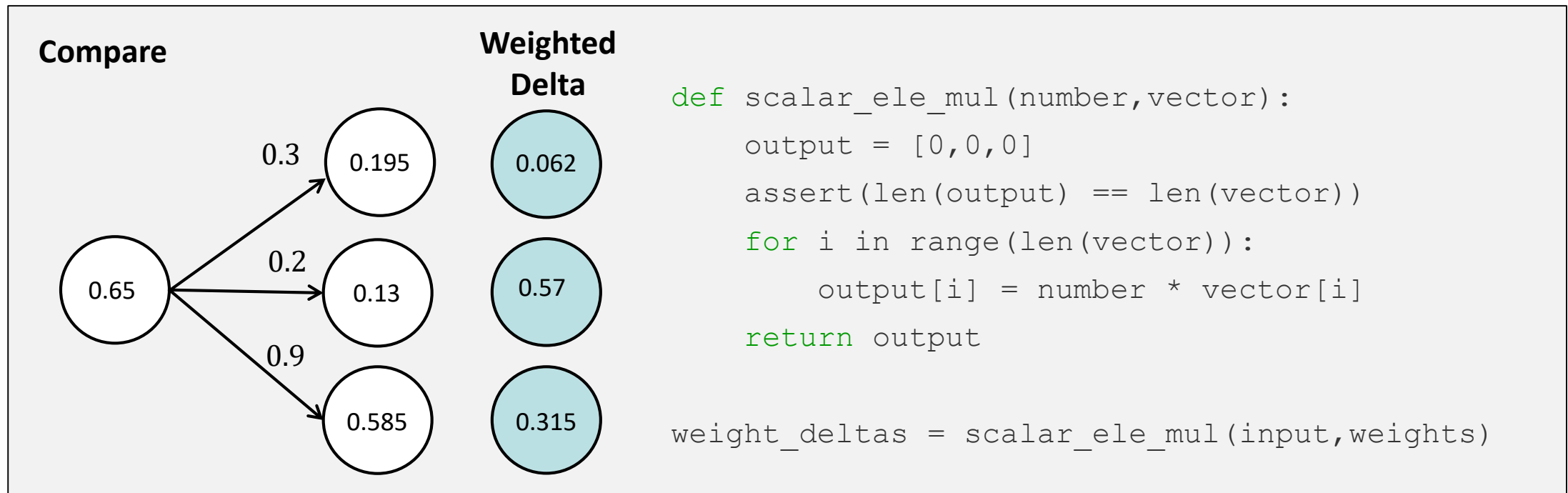
- **Single input → multiple predictions**
  - Perform each weight update



```
alpha = 0.1
for i in range(len(weights)):
    weights[i] -= (weight_deltas[i] * alpha)
```
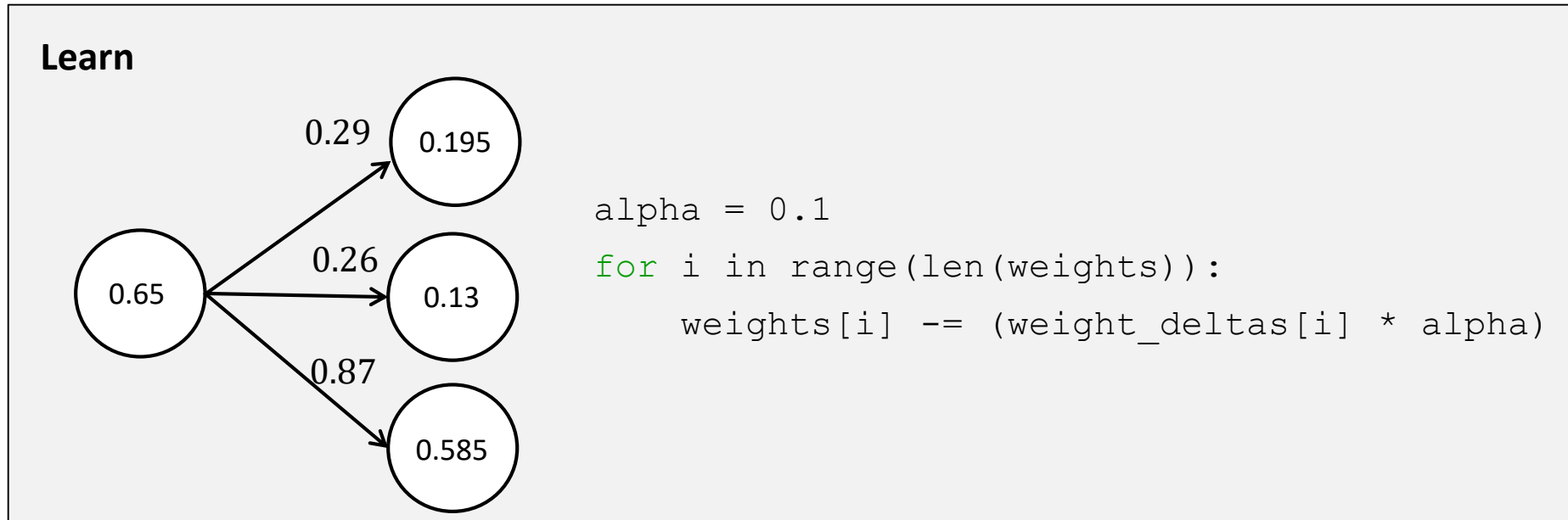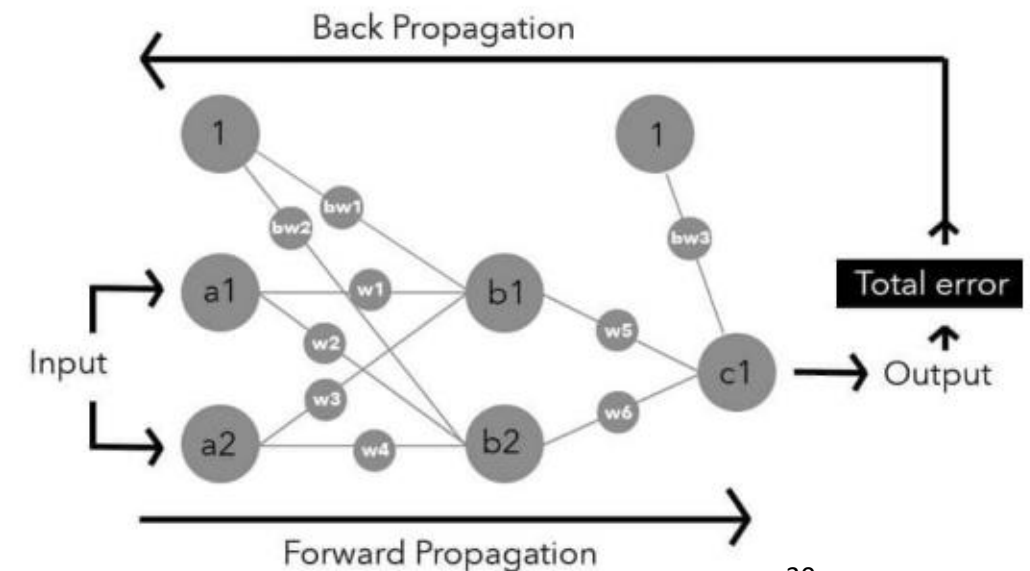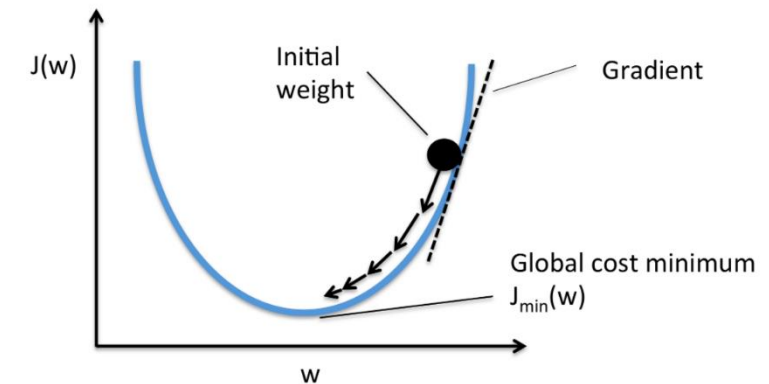
Image source: https://medium.com

- # Gradient Descent
  - Networks with multiple inputs
  - Networks with multiple outputs
  - **Networks with multiple inputs and outputs**
- Correlation
  - Learning Correlation
  - Creating Correlation
- Backpropagation
- Summary

- ## **Combining the lot!**
  - – Of course, it is straightforward to combine multiple input and multiple output networks to build a network that has both multiple inputs and multiple outputs

**Empty network with multiple inputs and outputs**



```
                    #offside #win #fans
weights = [ [0.1, 0.1, -0.3], #injuries?
            [0.1, 0.2, 0.0], #win?
            [0.0, 1.3, 0.1] ] #happy?


def neural_network(input, weight):
    pred = vec_mat_mul(input * weight)
    return pred
```

- Each output node takes its own weighted sum of the input and makes a prediction

**For each output, perform a weighted sum of inputs**



```
def w_sum(a,b):
    assert(len(a) == len(b))
    output = 0
    for i in range(len(a)):
        output += (a[i] * b[i])
    return output

def vect_mat_mul(vect,matrix):
    assert(len(vect) == len(matrix))
    output = [0,0,0]
    for i in range(len(vect)):
        output[i] = w_sum(vect,matrix[i])
    return output
```

- **Combining the lot!**
  - Gradient descent generalizes to arbitrarily large networks

**Empty network with multiple inputs and outputs**
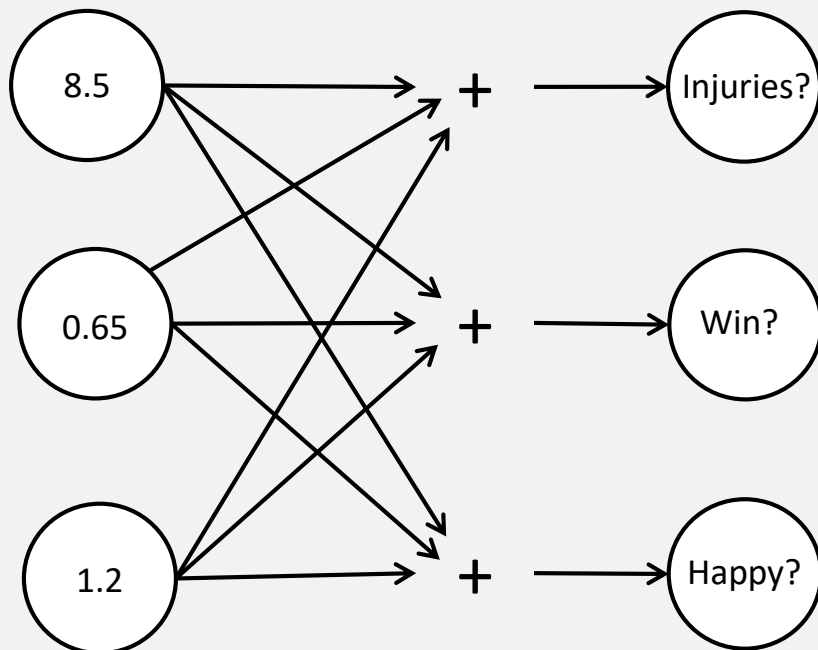


Input Data

Output Predictions

```
               #offside #win #fans
weights = [ [0.1, 0.1, -0.3],   #injuries?
            [0.1, 0.2, 0.0],    #win?
            [0.0, 1.3, 0.1] ]   #happy?
def vect_mat_mul(vect,matrix):
    assert(len(vect) == len(matrix))
    output = [0,0,0]
    for i in range(len(vect)):
        output[i] = w_sum(vect,matrix[i])
    return output
def neural_network(input, weights):
    pred = vect_mat_mul(input,weights)
    return pred
```
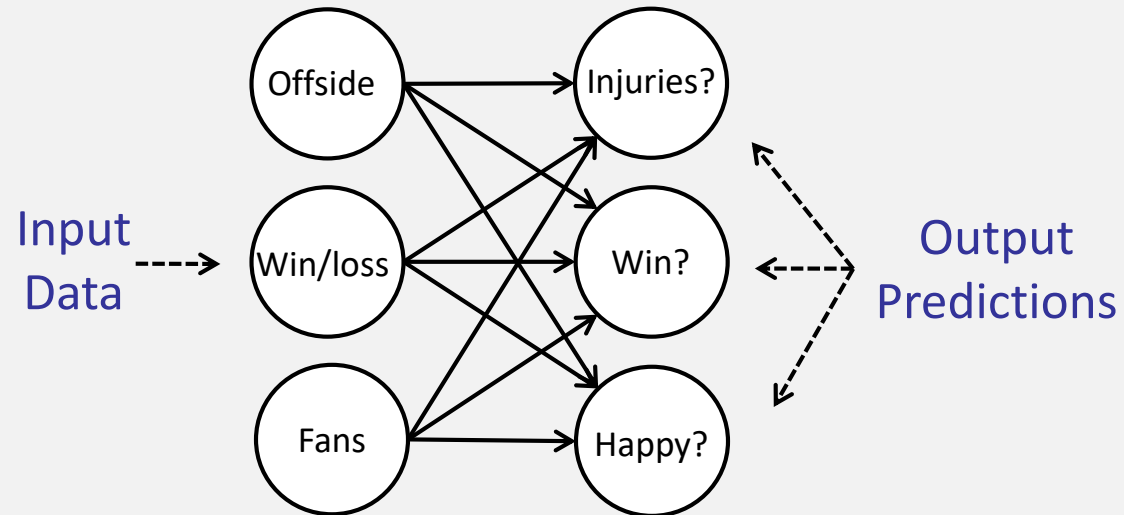
- **Combining the lot!**

**PREDICT: Making a prediction and calculating error and delta**



```
Offsides = [8.5, 9.5, 9.9, 9.0]
wlrec = [0.65,0.8, 0.8, 0.9]
nfans = [1.2, 1.3, 0.5, 1.0]
hurt = [0.1, 0.0, 0.0, 0.1]
win = [ 1, 1, 0, 1]
happy = [0.1, 0.0, 0.1, 0.2]
alpha = 0.01

input = [offsides[0],wlrec[0],nfans[0]]
true = [hurt[0], win[0], happy[0]]
pred = neural_network(input,weights)
error = [0, 0, 0]
delta = [0, 0, 0]
for i in range(len(true)):
    error[i] = (pred[i] - true[i]) ** 2
    delta = pred[i] - true[i]
```

- ## **Combining the lot!**



COMPARE: Calculating each weight_delta and putting it on each weight

Weighted Delta for
Win/loss node

```
def outer_prod(vec_a, vec_b):
    out = zeros_matrix(len(a),len(b))
    for i in range(len(a)):
        for j in range(len(b)):
            out[i][j] = vec_a[i]*vec_b[j]
    return out
weight_deltas = outer_prod(input,delta)
```

# Combining the lot!

**LEARN: Updating the weights**

**New weights for middle node**



```
for i in range(len(weights)):
        for j in range(len(weights[0])):
            weights[i][j] -= alpha * \ weight_deltas[i][j]
```

Image source: https://medium.com

- Gradient Descent
  - Networks with multiple inputs
  - Networks with multiple outputs
  - Networks with multiple inputs and outputs

- **Correlation**
  - **Learning Correlation**
  - Creating Correlation

- Backpropagation

- Summary

- **Key messages relating to weights so far**
  - **Weights represent knowledge**
    - A network uses this knowledge to interpret the input data
    - Weights are a measure of sensitivity between the input data of the network and its prediction
  - **Each weight tries to reduce the error**
    - We can compute the relationship between the error and any one of the weights so that we know how changing the weight changes the error. You can then use this to reduce the error to 0.
  - **What do they learn in aggregate?**

Deep Learning

- **What do these weights learn?**



784 Input Nodes — 10 Output Nodes

**Highest Prediction**
Network thinks this image is most likely a 2

**Small Error**
Network also thinks this image kind of looks like a 9

- **What do these weights learn?**
  - **Correlations between input and output**
    - If a weight is high, it means the model believes there's a high degree of correlation between that input and the prediction.
    - If the number is very low (negative), then the network believes there is a very low correlation between that input and the prediction
  - **Why is this?**
    - Weights are found via dot products
    - Dot product encodes similarity

**High Prediction**
High correlation/similarity
between inputs and weights

dot → **0.98**

dot → **0.1**

**Low Prediction**
Low correlation /similarity
between inputs and weights

- How does a network learns on entire datasets?
  - Introduce a new toy problem: The streetlight problem
    - Problem: learning to use unfamiliar streetlight in a foreign country



    - What we wish to learn: when is safe to cross
    - How is this achieved: interpreting the streetlight patterns
      - Observe and record the different streetlight combinations for when people stop or walk

  Observing the correlation between light combination and safe crossing

Deep Learning

- **Streetlight observations**
  - Observe correlation between the middle light and walking

Deep Learning

- **How does a network do this?**
  - First step convert observation into datasets
    - What we know, what we wish to know

| What we know | What we wish to know | | What we know | What we wish to know |
|---|---|---|---|---|
| ⊘ ○ STOP | | → | 101 | 0 |
| ● ○ WALK | | → | 011 | 1 |
| ○ STOP | | → | 001 | 0 |
| ⊘ ● ○ WALK | | → | 111 | 1 |
| ● ○ WALK | | → | 011 | 1 |
| ⊘ ○ STOP | | → | 101 | 0 |

```python
import numpy as np
weights = np.array([0.5,0.48,-0.7])
alpha = 0.1
streetlights = np.array( [ [ 1, 0, 1 ],
                           [ 0, 1, 1 ],
                           [ 0, 0, 1 ],
                           [ 1, 1, 1 ],
                           [ 0, 1, 1 ],
                           [ 1, 0, 1 ] ] )
walk_vs_stop = np.array( [ 0, 1, 0, 1, 1, 0 ] )

input = streetlights[0]
goal_prediction = walk_vs_stop[0]

for iteration in range(20):
    prediction = input.dot(weights)
    error = (goal_prediction - prediction) ** 2
    delta = prediction - goal_prediction
    weights = weights - (alpha * (input * delta))

print("Error:" + str(error) + " Prediction:" + str(prediction))
```

We are only learning from one sample

```python
import numpy as np
weights = np.array([0.5,0.48,-0.7])
alpha = 0.1
streetlights = np.array( [ [ 1, 0, 1 ],
                           [ 0, 1, 1 ],
                           [ 0, 0, 1 ],
                           [ 1, 1, 1 ],
                           [ 0, 1, 1 ],
                           [ 1, 0, 1 ] ] )
walk_vs_stop = np.array( [ 0, 1, 0, 1, 1, 0 ] )

for iteration in range(40):
    error_for_all_lights = 0
    for row_index in range(len(walk_vs_stop)):
        input = streetlights[row_index]
            goal_prediction = walk_vs_stop[row_index]
            prediction = input.dot(weights)
            error = (goal_prediction - prediction) ** 2
            error_for_all_lights += error
            delta = prediction - goal_prediction
            weights = weights - (alpha * (input * delta))
            print("Prediction:" + str(prediction))
        print("Error:" + str(error_for_all_lights) + "\n")
```

```
Prediction:-0.19999999999999996
Prediction:-0.1999999999999996
Prediction:-0.55999999999999999
Prediction:0.616000000000001
Prediction:0.1727999999999995
Prediction:0.17552
Error:2.6561231104

Prediction:0.14041599999999999
Prediction:0.3066464
Prediction:-0.34513824
Prediction:1.006637344
Prediction:0.478503475199999
Prediction:0.26700416768
Error:0.9628701776715985
.
.
.
Prediction:-0.0022410273814405524
Prediction:0.9978745386023716
Prediction:-0.01672126442988494947
Prediction:1.015112745989381212
Prediction:0.9969492081270097
Prediction:-0.00262561932978312525
Error:0.00053373677328488
```

- **What did our light crossing network learn?**
  - It learnt the middle light was most important
  - Note that:
    - The middle weight is near one
    - The left and right weights are near zero
  - The network identified a correlation between the middle node & the output
  - It also identified the randomness between the left and right nodes & the output
    - Reflected by the near zero weights

- **How did the network identify this correlation?**
  - The process of gradient descent can be thought of as applying *upwards* or *downwards* pressure on the weights
    - On average the middle weight received more upward pressure
    - On average the outer weights received more downward pressure

| Training Data | | Weight Pressure | |
|---|---|---|---|
| 1 0 1 | 0 | − 0 − | 0 |
| 0 1 1 | 1 | 0 + + | 1 |
| 0 0 1 | 0 | 0 0 − | 0 |
| 1 1 1 | 1 | + + + | 1 |
| 0 1 1 | 1 | 0 + + | 1 |
| 1 0 1 | 0 | − 0 − | 0 |

- **This pressure comes from the data**
  - Each node is individually trying to correctly predict the outcome
  - Only communication between nodes is the shared error
    - Weight update is multiplying shared error by respective inputs
  - Concept of *Error Attribution*
    - Given the share error, the network must decide weights to update & which to leave

| Training Data | | Weight Pressure | |
|---|---|---|---|
| 1 0 1 | 0 | − 0 − | 0 |
| 0 1 1 | 1 | 0 + + | 1 |
| 0 0 1 | 0 | 0 0 − | 0 |
| 1 1 1 | 1 | + + + | 1 |
| 0 1 1 | 1 | 0 + + | 1 |
| 1 0 1 | 0 | − 0 − | 0 |

- **This pressure comes from the data**
  - Weight pressure table shows this effect
    - In the first training sample the left and right nodes are uncorrelated with the output so their weights experience downwards pressure
    - In the second training sample the middle and right nodes are correlated with the output and their weights experience upwards pressure

| Training Data | | Weight Pressure | |
| --- | --- | --- | --- |
| 1 0 1 | 0 | $-$ 0 $-$ | 0 |
| 0 1 1 | 1 | 0 + + | 1 |
| 0 0 1 | 0 | 0 0 $-$ | 0 |
| 1 1 1 | 1 | + + + | 1 |
| 0 1 1 | 1 | 0 + + | 1 |
| 1 0 1 | 0 | $-$ 0 $-$ | 0 |

- **How did the network identify this correlation?**
  - Reminder: prediction is a weighted sum of the inputs
  - During training
    - The learning algorithm rewards inputs that correlate with the output with upward pressure on their weight
      - Drives the weights towards 1
    - At the same time, the learning algorithm penalises inputs with discorrelation with downward pressure
      - Drives the weights towards 0
  - Learning rewards correlation with larger weights
  - Learning finds correlation between the two datasets

Image source: https://medium.com

- Gradient Descent
  - Networks with multiple inputs
  - Networks with multiple outputs
  - Networks with multiple inputs and outputs

- **Correlation**
  - Learning Correlation
  - **Creating Correlation**

- Backpropagation

- Summary

- ## **What happens when there is no correlation?**
  - Considering the following dataset:
    - There is no correlation between *any* input column and the output column.
    - Every weight has an equal amount of upward pressure and downward pressure
  - How can the network learn?

|  Training Data  | | | |
| --- | --- | --- | --- |
| 1 | 0 | 1 | 1 |
| 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 0 |
| 1 | 1 | 1 | 0 |

$\longrightarrow$

| Weight Pressure | | | |
| --- | --- | --- | --- |
| + | 0 | + | 1 |
| 0 | + | + | 1 |
| 0 | 0 | − | 0 |
| − | − | − | 0 |

## • Key message

- Computing the relationship between the error and weights allows us to change the weights to reduce the error to 0.

- Adjusting the weights to reduce the error over a series of training examples ultimately searches for correlation between the input and the output layers

- **If no correlation exists, then the error will never reach 0**

- This is a major limiting factor, therefore we need to change our approach to cope with this

- **What happens when there is no correlation?**
  - Learnt today that neural network is an instrument that searches for correlation between input and output datasets
  - In reality the network searches for correlation between its input and output layers
  - So to learn when there is no correlation, just use more networks (i.e., add another layer)
    - First layer creates a layer with limited correlation with the output
    - Second layer uses this limited correlation to correctly predict the output
  - Hidden layer(s) can be thought of as creating intermediate dataset(s) that has correlation with the output

Walk/Stop

layer_2

weights_1_2

This will be the intermediate data representation

layer_1

weights_0_1

layer_0

- **Stacked neural networks**
  - Output of the first network is the input to the next
  - Prediction is identical to what we have already learnt
  - Top half of the network can be trained via gradient descent using methods we have learnt
  - But how to update the weights for the hidden layer?

- **Linear vs. nonlinear**
  - Without an addition of non-linear transformations in-between layers, the hidden layers adds nothing new to the network
  - Why?
    - **Any two multiplications can (in general) by accomplished using a single multiplication**



$$1 \times 10 \times 10 = 100$$
$$5 \times 20 = 100$$

$$1 \times 0.25 \times 0.9 = 0.225$$
$$1 \times 0.225 = 0.225$$

- **Linear vs. nonlinear**
  - Stacking linear neural networks does not give and more power
  - It just gives a more computationally expensive version of a single weighted sum
  - Anything that the three-layer linear network can do, the two-layer linear network can also do

- **Linear vs. nonlinear**
  - Moreover, in linear networks the hidden nodes do not have a correlation of their own
  - They are more or less correlated to the various input nodes
  - Ideally we want the middle layer to sometimes correlate with the input and sometimes not correlate
  - Such an effect means the hidden layers can give a correlation of their own
    - This is known as conditional correlation
  - This can be done by introducing non-linearities

## Activation Functions

- Uses the weighted input value to determine the level of output activation
  - Introduces nonlinearities into network
- Typical activation functions include
  - Identity → $f(x) = x$
  - Sigmoid → $f(x) = {}^1\!/_{(1+\,e^{-x})}$
  - Tanh → $f(x) = \tanh(x) = {}^{(e^x - e^{-x})}\!/_{(e^x + e^{-x})}$
  - Rectified Linear Unit → $f(x) = \begin{cases} 0 \ for \ x \leq 0 \\ x \ for \ x > 0 \end{cases}$

**Sigmoid**

**Tanh**

**ReLU**

- **What's the point of creating intermediate datasets that have correlation?**
  - Deep learning is all about creating intermediate layers
  - Each node in an intermediate layer represents the presence or absence of a different configuration of inputs
  - No individual input has to correlate directly with the target
  - Middle layer attempts to identify different configurations of the input that may or may not correlate with the output
  - These many different configurations will give the final layer the information (correlation) it needs to perform the prediction

Image source: https://medium.com

- Gradient Descent
  - Networks with multiple inputs
  - Networks with multiple outputs
  - Networks with multiple inputs and outputs

- Correlation
  - Learning Correlation
  - Creating Correlation

- **Backpropagation**

- Summary

# • **Neural Networks**

– Use layered combinations of perceptrons make complex predictions



$$\varphi\left((w_1 \quad w_2 \quad \dots \quad b)\begin{pmatrix} x_1 \\ x_2 \\ \dots \\ 1 \end{pmatrix}\right) = \varphi(w_1 x_1 + w_2 x_2 + \dots + b) = y_k$$

Rectified Linear unit
$$\varphi(x) = \max(0, x)$$

# Forward Propagation

– Information flows from input to output to make a predication

- **Forward Propagation**
  - Each neuron is a function of the previous one connected to it
    - Output is a composite function of the weights, inputs, and activations
      - Change any one of these and ultimately the output well change

| Input | Hidden 1 | Hidden 2 | Output |
|-------|----------|----------|--------|

$$output = act\left(W^3 \times act(W^2 \times act(W^1 \times input))\right)$$

- **Calculate Error/Loss**
  - Calculate difference between predication and target



$$J = \frac{1}{n} \sum_{i=1}^{n} \left( P_i - \hat{P}_i \right)^2$$

**Error**

- **Propagating error back through the network**
  - A neural network learns via a error function
    - This function defines the relationship between the error and the weights of the network
    - Update weights via derivative of this function

$$W^{(t+1)} = W^{(t)} - \alpha \nabla J(W^{(t)}; x, y)$$

  - Backpropagation is a method for computing the derivatives with respect to every weight in a network
    - With this derivative we perform our gradient descent update
  - It is known as backpropagation as we use it to travers the output error back through the network

- **Perform Gradient descent**
  - Output value effected by weights at all layers



Input      Hidden 1      Hidden 2      Output      Loss

UPDATE $\longleftarrow$ UPDATE $\longleftarrow$ UPDATE

$$J = \frac{1}{n}\sum_{i=1}^{n}(P_i - \hat{P}_i)^2$$

**Error**

# Backpropagation

– Tool to calculate the gradient of the loss function

- ## **Calculating gradient for arbitrary weight**
  - Iteratively apply the chain rule
  - Note: Error is now a function of the output and hence a function of the input, weights, and activation functions

| Input | Hidden 1 | Hidden 2 | Output | Loss |
|---|---|---|---|---|

$$\frac{\partial error}{\partial W^1} = \frac{\partial error}{\partial output} \times \frac{\partial output}{\partial hidden2} \times \frac{\partial hidden2}{\partial hidden1} \times \frac{\partial hidden1}{\partial W^1}$$

- **Calculating gradient via backpropagation**
  - Local gradients can be calculated before starting the backpropagation process
  - Iteratively apply the chain rule
    - The derivative of the output of the network with respect to a local variable is found by multiplying the local gradient with the upstream gradient

Upstream Gradient $\times$ Local Gradient

$$\frac{df}{dx} \times \frac{dx}{da}$$

$a$

Upstream Gradient

$$\frac{df}{dx}$$

$x$

$+$

$$\frac{df}{dx} \times \frac{dx}{db}$$

$b$

- **Calculating gradient via backpropagation**
  - Multivariate chain rule: Gradients add at branches

$$\frac{df}{db} = \frac{df}{dx} \times \frac{dx}{db} + \frac{df}{dy} \times \frac{dy}{db}$$

- **Example: Calculating** $\dfrac{\partial error}{\partial W_{ij}^1} = \dfrac{\partial j}{\partial W_{ij}^1}$



$W_{ij}^1 \rightarrow$ denotes an arbitrary weight between input and first hidden layer

- **Example: Calculating** $\dfrac{\partial j}{\partial W^1_{ij}}$

- **Sigmoid activation function**

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

$$\frac{d\sigma(x)}{d(x)} = \sigma(x)(1 - \sigma(x))$$

- **Example: Calculating** $\dfrac{\partial j}{\partial W_{ij}^1}$

- **Error function $J(W)$**

  – Output estimations denoted as vector $p$

  – Actual (target) outputs as vector $\mathbf{a}$

$$J = \frac{1}{2}(p - a)^2$$

$$\frac{dj}{dp} = (p - a)$$

- **Example: Calculating $\frac{\partial j}{\partial W_{ij}^1}$**

- **Notation (Generalised)**

  - $x_i, x_j$ - any arbitrary input/output pair for the input layer

    - Note as this is the input: $x_i = x_j$

  - $y_i, y_j$ - any arbitrary input/output pair for hidden layer 1

    - Note: $y_j = \sigma(y_i) = \sigma(\sum_j x_j)$

  - $z_i, z_j$ - any arbitrary input/output pair for hidden layer 2

  - $p_i, p_j$ - input/output pair for output layer

    - Note: sigmoid activation used on output

- **Example: Calculating** $\frac{\partial j}{\partial W_{ij}^1}$

- **Notation (Generalised)**

  Note: $i$'s and $j$'s in the weights and other variables are completely different
  - For $x/y/z/p$ they represent input/output
  - For $W$ they indicate neuron connectivity

  – The weights are organised into three separate variables: $\boldsymbol{W}^1, \boldsymbol{W}^2\ \boldsymbol{W}^3$

    - Each $\boldsymbol{W}$ is a matrix denote all weights at the respective layer

  – $W_{ij}^L$ is any single arbitrary weight at a given layer

    - $\boldsymbol{W}_i^L$ denotes all the weights that connect arbitrary neuron $i$ at the preceding layer

    - $W_{ij}^L$ is the weight that connects arbitrary neuron $i$ at the preceding layer to an arbitrary neuron $j$ at the next layer.

- **Example: Calculating** $\dfrac{\partial j}{\partial W_{ij}^1}$



$W_{ij}^1 \rightarrow$ denotes an arbitrary weight between input and first hidden layer

- **Calculating $\frac{\partial J}{\partial W_{ij}^1} \rightarrow$ starting at $p_j$**



$W_{ij}^1 \rightarrow$ denotes an arbitrary weight between input and first hidden layer

- Starting at $p_j$

$$\frac{\partial J}{\partial W_{ij}^1} = \frac{\partial J}{\partial p_j} \frac{\partial p_j}{\partial W_{ij}^1}$$

  - Noting $\frac{dj}{dp} = (p - a)$

$$\frac{\partial J}{\partial W_{ij}^1} = (p_j - a) \frac{\partial p_j}{\partial W_{ij}^1}$$

- **Calculating $\dfrac{\partial J}{\partial W_{ij}^1} \rightarrow$ propagating $p_j$ to $p_i$**



$W_{ij}^1 \rightarrow$ denotes an arbitrary weight between input and first hidden layer

Deep Learning

- **Propagating one layer back until to $p_i$**

$$\frac{\partial J}{\partial W_{ij}^1} = (p_j - a)\frac{\partial p_j}{\partial W_{ij}^1} = (p_j - a)\frac{\partial p_j}{\partial p_i}\frac{\partial p_i}{\partial W_{ij}^1}$$

- Noting $\dfrac{\partial p_j}{\partial p_i} = p_j(1 - p_j)$ $\longleftarrow$ Derivative of sigmoid with simplified notation

$$\frac{\partial J}{\partial W_{ij}^1} = (p_j - a)p_j(1 - p_j)\frac{\partial p_i}{\partial W_{ij}^1}$$

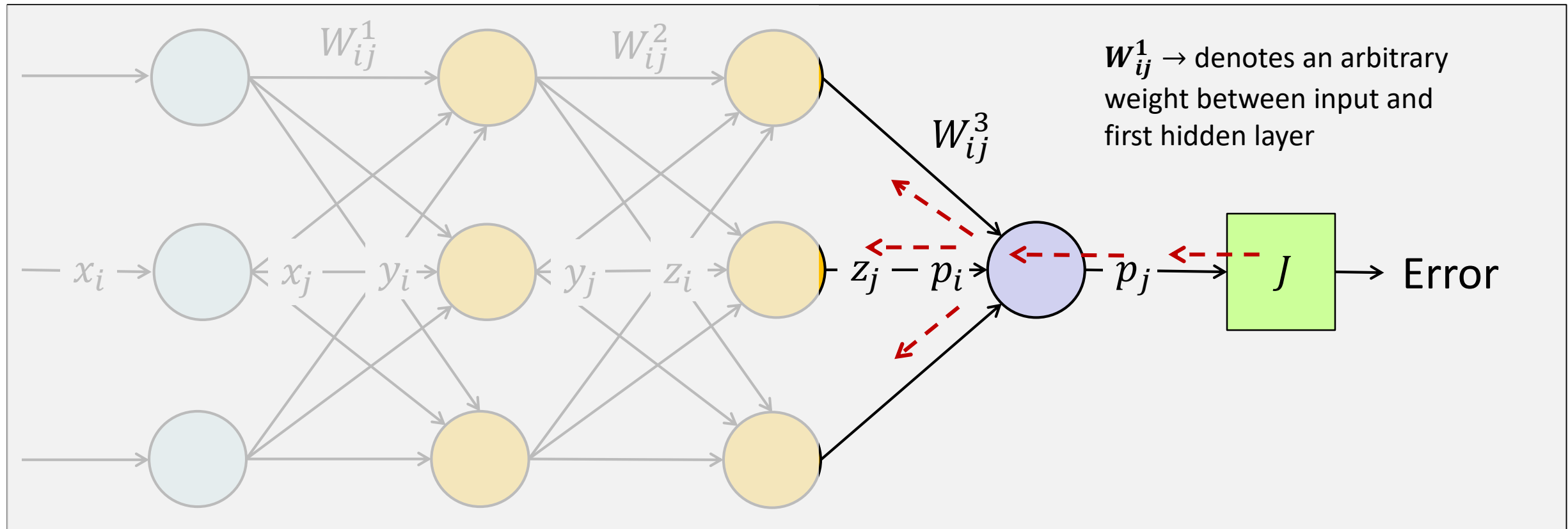- **Calculating $\dfrac{\partial J}{\partial W_{ij}^1} \rightarrow$ propagating $p_i$ to $z_j$**



$W_{ij}^1 \rightarrow$ denotes an arbitrary weight between input and first hidden layer

- **Propagating from $p_i$ to $z_j$**

  - Need to take into account action of the different weights

$$\frac{\partial p_i}{\partial W_{ij}^1} = \sum_j \frac{\partial p_i}{\partial z_j} \frac{\partial z_j}{\partial W_{ij}^1}$$

Summation reflects that we can choose any one of $j$ paths to propagate along

- Noting $\frac{\partial p_i}{\partial z_j} = W_{ij}^3$, where $\frac{\partial p_i}{\partial z_j}$ is derivate w.r.t. any *arbitrary* $z_j$

$$\frac{\partial p_i}{\partial W_{ij}^1} = \sum_j \left[ W_{ij}^3 \frac{\partial z_j}{\partial W_{ij}^1} \right]$$

$W_{ij}^3$ is the connecting weight on our (arbitrary) path back to $W_{ij}^1$

- **Calculating $\dfrac{\partial J}{\partial W_{ij}^1} \rightarrow$ propagating from $z_j$ to $z_i$**



$\boldsymbol{W_{ij}^1} \rightarrow$ denotes an arbitrary weight between input and first hidden layer

- **Propagating from $z_j$ to $z_i$**

$$\frac{\partial \boldsymbol{z}_j}{\partial \boldsymbol{W}_{ij}^1} = \frac{\partial \boldsymbol{z}_j}{\partial \boldsymbol{z}_i} \frac{\partial \boldsymbol{z}_i}{\partial \boldsymbol{W}_{ij}^1}$$

– Noting $\frac{\partial \boldsymbol{z}_j}{\partial \boldsymbol{z_i}} = \boldsymbol{z}_j(1 - \boldsymbol{z}_j)$ $\longleftarrow$ Derivative of sigmoid with simplified notation

$$\frac{\partial \boldsymbol{z}_j}{\partial \boldsymbol{W}_{ij}^1} = \boldsymbol{z}_j(1 - \boldsymbol{z}_j) \frac{\partial \boldsymbol{z}_i}{\partial \boldsymbol{W}_{ij}^1}$$

- **Calculating $\dfrac{\partial J}{\partial W_{ij}^1} \rightarrow$ propagating from $z_i$ to $y_j$**



$W_{ij}^1 \rightarrow$ denotes an arbitrary weight between input and first hidden layer

- **Propagating from $z_i$ to $y_j$**

  - Again, need to take into account action of the different weights

$$\frac{\partial z_i}{\partial W_{ij}^1} = \sum_j \frac{\partial z_i}{\partial y_j} \frac{\partial y_j}{\partial W_{ij}^1}$$
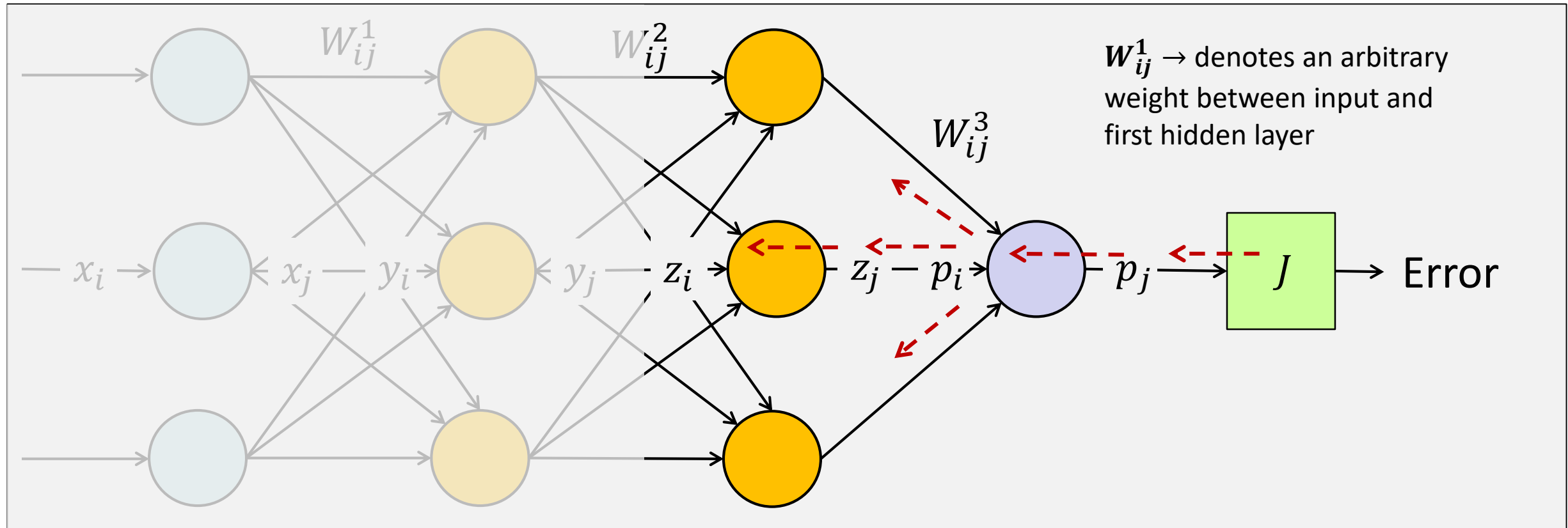
Summation reflects that we can choose any one of $j$ paths to propagate along

- Noting $\frac{\partial z_i}{\partial y_j} = W_{ij}^2$, note $\frac{\partial z_i}{\partial y_j}$ is derivate wrt any arbitrary $y_j$

$$\frac{\partial z_i}{\partial W_{ij}^1} = \sum_j \left[ W_{ij}^2 \frac{\partial y_j}{\partial W_{ij}^1} \right]$$

$W_{ij}^2$ is the connecting weight on our (arbitrary) path back to $W_{ij}^1$

- **Calculating** $\dfrac{\partial J}{\partial W_{ij}^1} \rightarrow$ **propagating from** $y_j$ **to** $y_i$
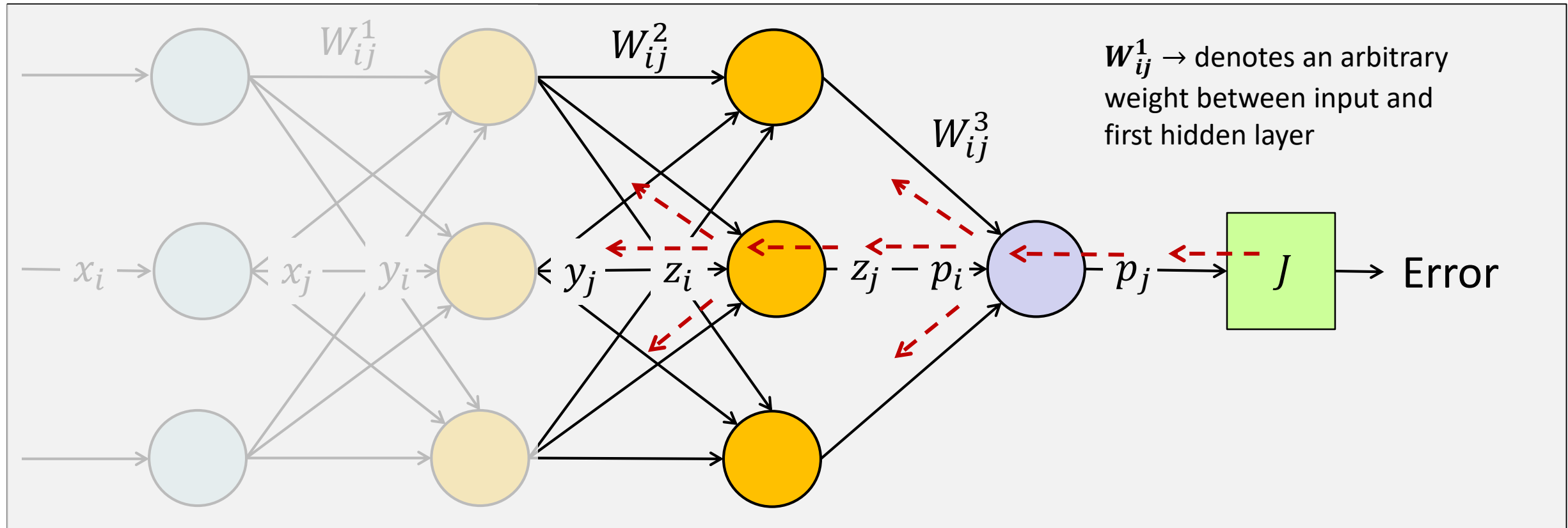


$W_{ij}^1 \rightarrow$ denotes an arbitrary weight between input and first hidden layer

- **Propagating from $y_j$ to $y_i$**

$$\frac{\partial y_j}{\partial W_{ij}^1} = \frac{\partial y_j}{\partial y_i} \frac{\partial y_i}{\partial W_{ij}^1}$$

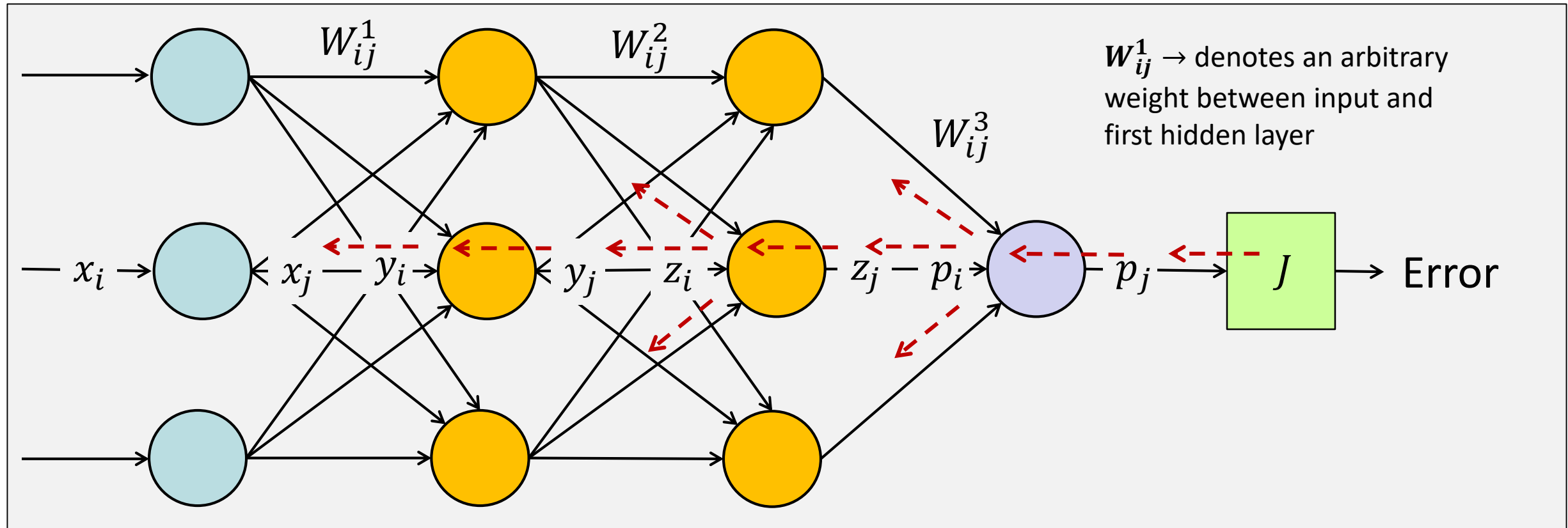– Noting $\dfrac{\partial y_j}{\partial y_i} = y_j(1 - y_j)$ ⟵ Derivative of sigmoid with simplified notation

$$\frac{\partial y_j}{\partial W_{ij}^1} = y_j(1 - y_j) \frac{\partial y_i}{\partial W_{ij}^1}$$

- **Calculating $\dfrac{\partial J}{\partial W_{ij}^1} \rightarrow$ propagating from $y_i$ to $x_j$**



$W_{ij}^1 \rightarrow$ denotes an arbitrary weight between input and first hidden layer

- **Propagating from $y_i$ to $x_j$**

  – We have to pass through $W^1_{ij}$ → what we are backpropagating to

We are now directly deriving $y_j$ w.r.t $W^1_{ij}$. We are **not** passing through $W^1_{ij}$ as in the previous layers

$$\frac{\partial y_i}{\partial W^1_{ij}} = \sum_j x_j$$

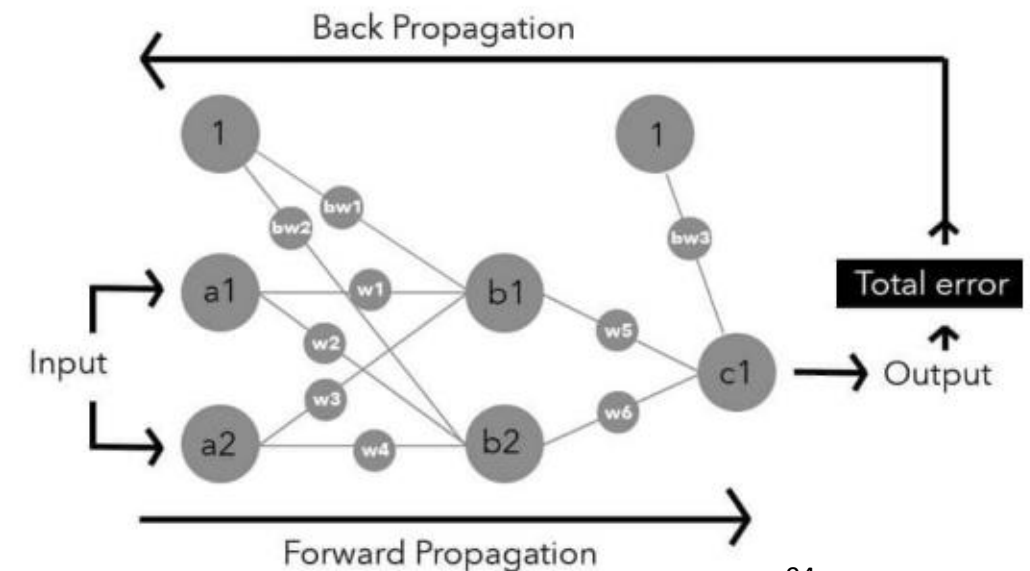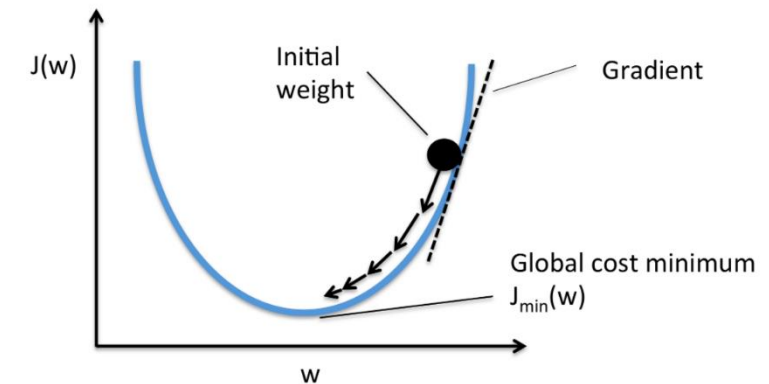Summation reflects that we can choose any one of $j$ paths to propagate along

  – $\sum_j x_j$ will be a numerical value, so cannot be derived further

  – So putting it all together:

$$\frac{\partial J}{\partial W^1_{ij}} = (p_j - a)p_j(1 - p_j)\sum_j\left[W^3_{ij}z_j(1 - z_j)\sum_j\left[W^2_{ij}y_j(1 - y_j)\sum_j x_j\right]\right]$$
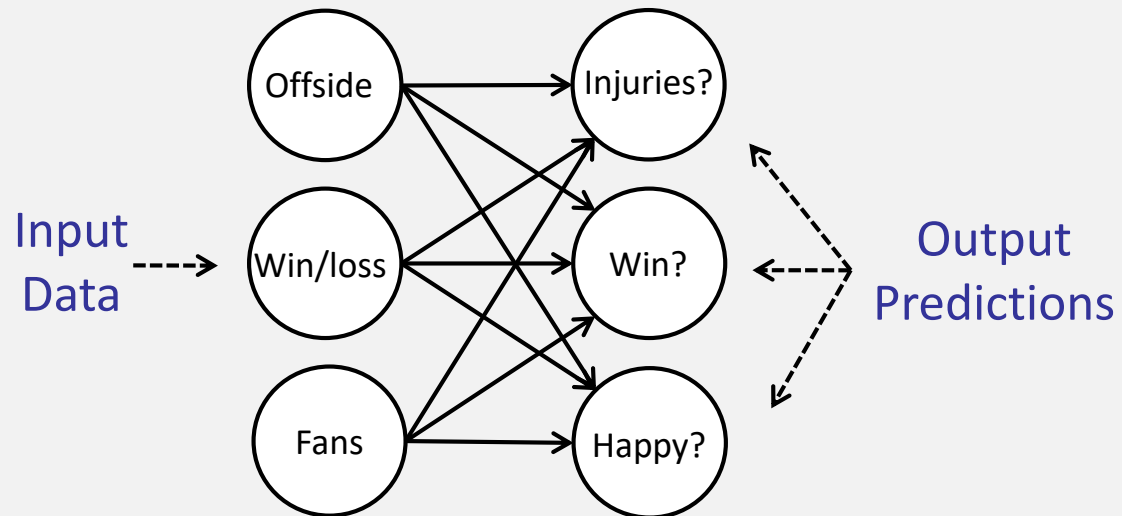
Image source: https://medium.com

- Gradient Descent
  - Networks with multiple inputs
  - Networks with multiple outputs
  - Networks with multiple inputs and outputs
- Correlation
  - Learning Correlation
  - Creating Correlation
- Backpropagation
- Summary

- **Gradient descent generalies to arbitrarily large networks**

**Empty network with multiple inputs and outputs**



```
pred = neural_network(input,weights)
error = [0, 0, 0]
delta = [0, 0, 0]
for i in range(len(true)):
    error[i] = (pred[i] - true[i]) ** 2
    delta = pred[i] - true[i]


weight_deltas = outer_prod(input,delta)

for i in range(len(weights)):
    for j in range(len(weights[0])):
        weights[i][j] -= alpha * \ weight_deltas[i][j]
```

- **Weights learn correlations between input and output**
  - Prediction is a weighted sum of the inputs – the dot product
  - During training
    - The learning algorithm rewards inputs that correlate with the output with upward pressure on their weight
      - Drives the weights towards 1
    - At the same time, the learning algorithm penalises inputs with discorrelation with downward pressure
      - Drives the weights towards 0
  - Learning rewards correlation with larger weights
  - Learning finds correlation between the two datasets
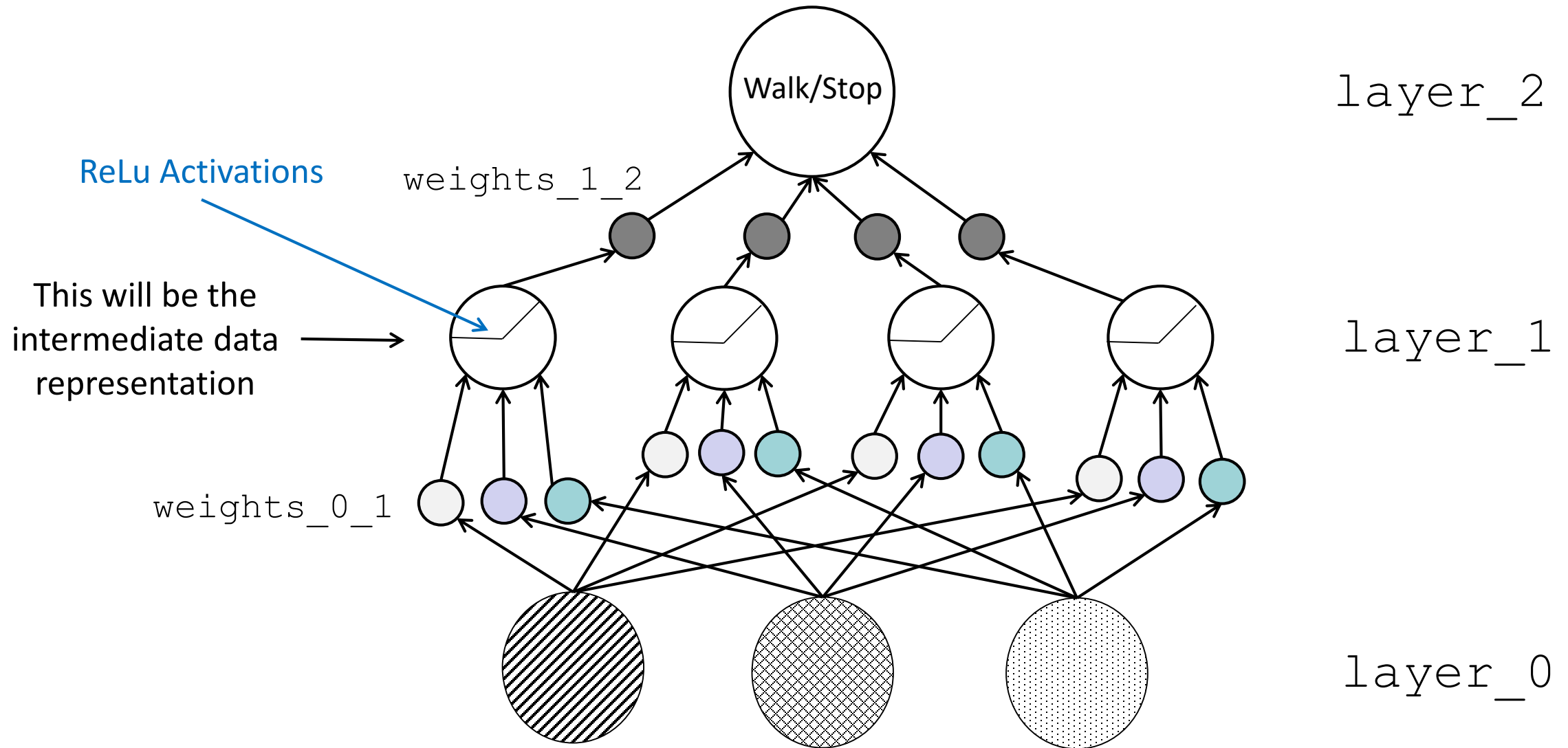
- **Creating Correlation**
  - To learn when there is no correlation, just use more layers
    - First layer creates a layer with limited correlation with the output
    - Second layer uses this limited correlation to correctly predict the output
  - Hidden layer(s) can be thought of as creating intermediate dataset(s) that has correlation with the output
  - Stacking linear neural networks does not give and more power
    - A more computationally expensive version of a single weighted sum
    - Use non-linear activation functions to induce conditional correlation between layer

Walk/Stop

layer_2

ReLu Activations

weights_1_2

This will be the intermediate data representation

layer_1

weights_0_1

layer_0

Deep Learning

- **Backpropagation**
  - Update weights to minimise loss function
    - This is achieved by taking the gradient of loss function with respect to the weights
  - Not a trivial process as neural networks are a series of layers
    - The output is a composite function of the weights, inputs, and activation functions
  - Backpropagation is a tool to calculate the gradient of the loss function with respect to any single weight in a network
    - Allows us to calculate the gradient at all weights
    - This is achieved through iterative applications of the chain rule

- **Next week**
  - Loss functions
    - The difference between what was predicted and what it should have been predicted
  - Activation functions
    - A function applied to neurons in a layer during prediction
  - Optimisers
    - Different algorithms for updating neural network weights
  - Recurrent Neural Networks
    - Model *sequences* of data