

Neural Networks

Krisztian Buza

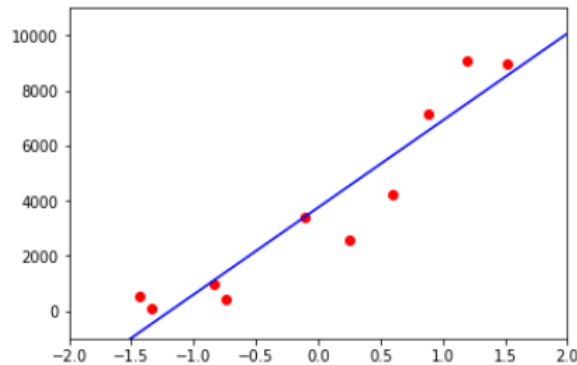
Department of Artificial Intelligence
Eötvös Loránd University
Budapest, Hungary

Announcements

Announcements

- Quick Test Replacement Test (QTRT)
 - In the same session with the „big test“
 - 11th Dec: „big test“ (60 minutes) + ~~QTRT (30 minutes)~~
 - Due to multiple requests, QTRT will be organized between 16th Dec and 20th Dec
 - You may collect maximal 14 points which replace the result of the 7 worst quick tests
- PIN code for querying the results of tests
- Form groups for presentation about ethical, social, economical, technical aspect of AI

Calculation of Root Mean Square Error

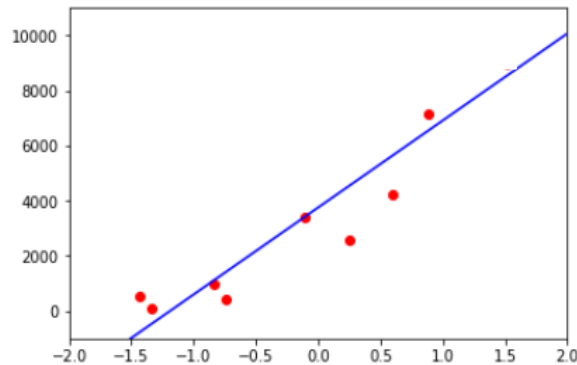


y	\hat{y}
549	-770
100	-470
976	1132
441	1433
3401	3435
2600	4537
4241	5638
7150	6539
9100	7541
9000	8542

$n = 10$

$$\frac{1}{10} \left((549 - (-770))^2 + (100 - (-470))^2 + \right. \\ \left. + (976 - 1132)^2 + (441 - 1433)^2 + \right. \\ \left. + (3401 - 3435)^2 + (2600 - 4537)^2 + \right. \\ \left. + (4241 - 5638)^2 + (7150 - 6539)^2 + \right. \\ \left. + (9100 - 7541)^2 + (9000 - 8542)^2 \right)$$

Calculation of Root Mean Square Error



y	\hat{y}
549	-770
100	-470
976	1132
441	1433
3401	3435
2600	4537
4241	5638
7150	6539
9100	7541
9000	8542

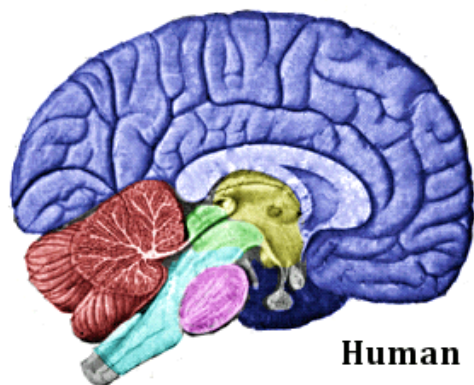
$n = 8$

$$\frac{1}{8} \left((549 - (-770))^2 + (100 - (-470))^2 + (976 - 1132)^2 + (441 - 1433)^2 + (3401 - 3435)^2 + (2600 - 4537)^2 + (4241 - 5638)^2 + (7150 - 6539)^2 + \right. \\ \left. + (9100 - 7541)^2 + (9000 - 8542)^2 \right)$$

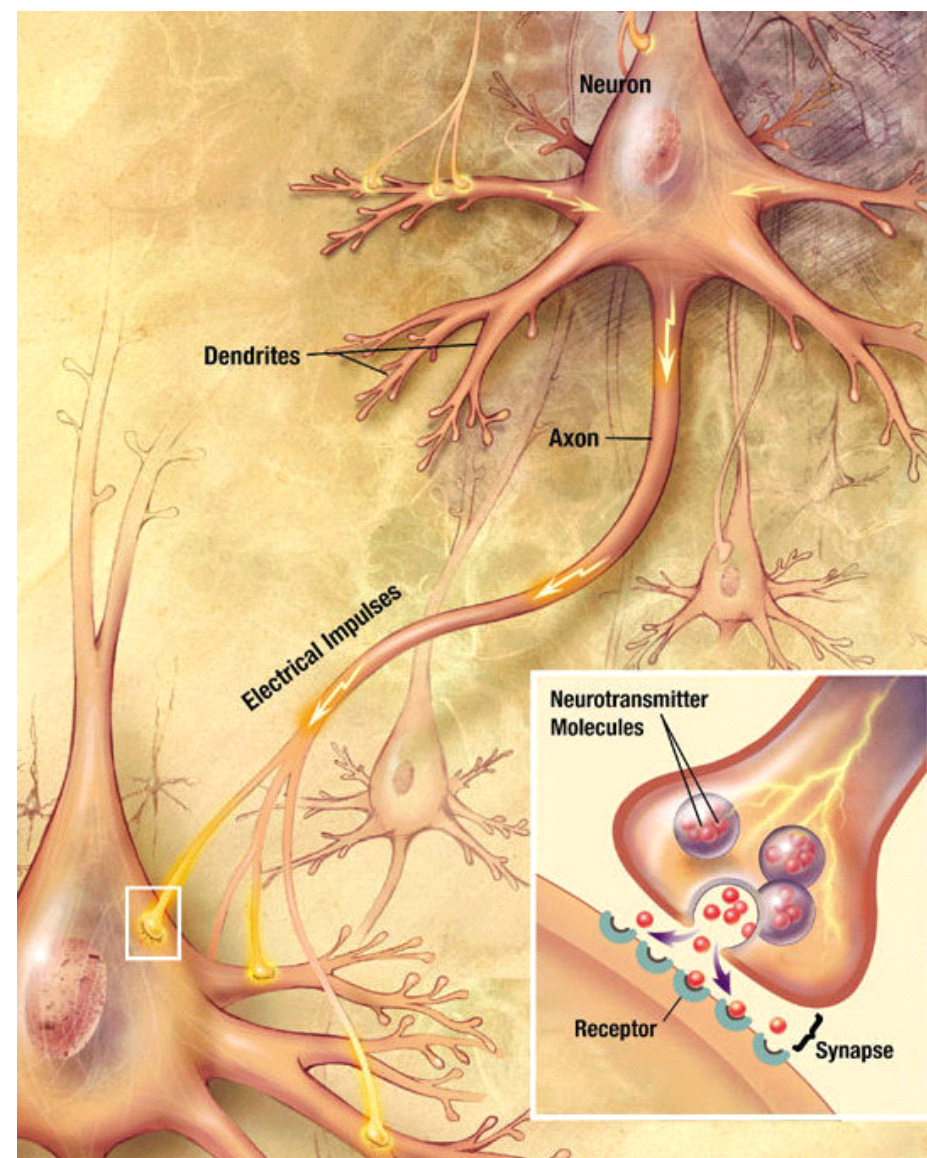
Artificial Neural Networks

Motivation

- Complex classification tasks with many variables (e.g. image classification)
- Biological motivation (single learning algorithm, structure, spikes, learned features)
- State-of-the-art learning algorithms

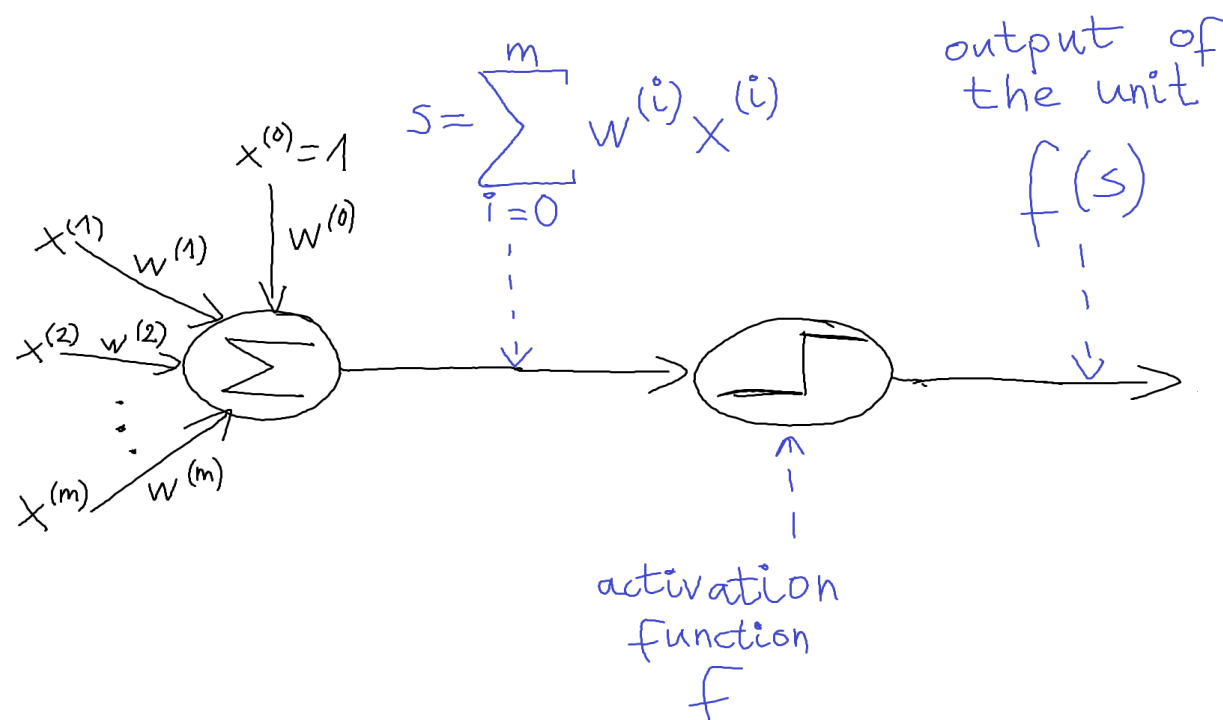


By Vertebrate-brain-regions.png: Looie496derivative work: Looie496 (Vertebrate-brain-regions.png) [Public domain], via Wikimedia Commons



By userLooie496 created file, US National Institutes of Health, National Institute on Aging created original [Public domain], via Wikimedia Commons

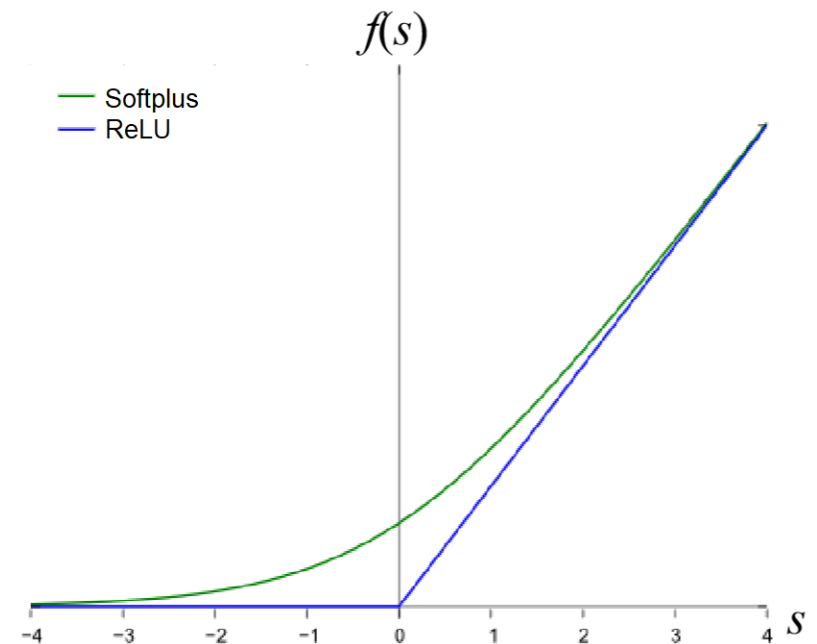
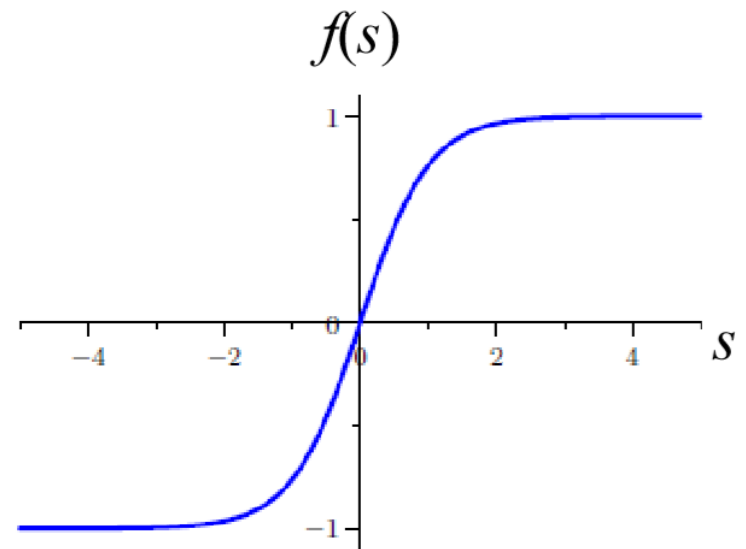
A Neural Unit



$x^{(1)}, x^{(2)}, \dots, x^{(m)}$ = inputs of the unit
 $w^{(1)}, w^{(2)}, \dots, w^{(m)}$ = weights of $x^{(1)}, x^{(2)}, \dots, x^{(m)}$
 $w^{(0)}$ = bias weight

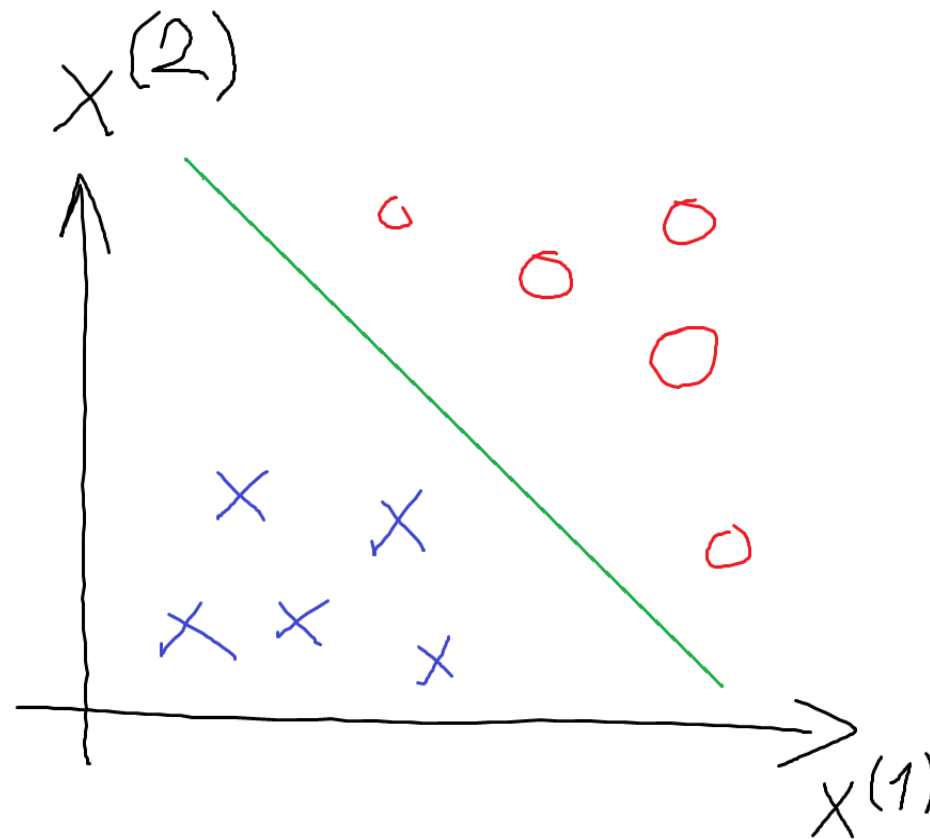
Activation Functions

- sigmoid \rightarrow logistic regression
- ReLU \rightarrow linear regression with nonnegative target



Activation Functions	
Linear	$f(s) = s$
Sigmoid	$f(s) = (1 + e^{-s})^{-1}$
Hyperbolic tangent	$f(s) = \tanh(s)$
Softsign	$f(s) = s((1 + s)^{-1})$
Rectifier Linear Unit (ReLU)	$f(s) = \max(0, s)$
Softplus	$f(s) = \ln(1 + e^s)$
...	...

Single Sigmoidal Unit – Decision Boundary

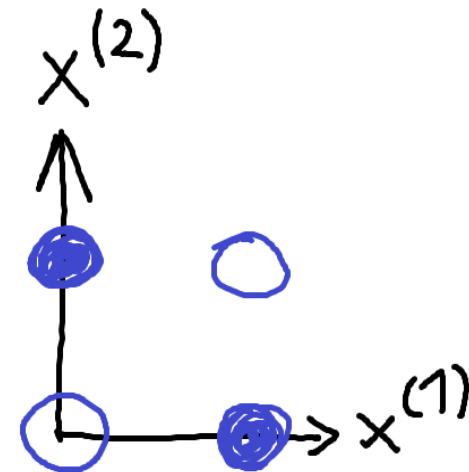


Can a Sigmoidal Unit Learn the XOR Function?

$X^{(1)}$	$X^{(2)}$	$X^{(1)} \text{ XOR } X^{(2)}$
0	0	0
0	1	1
1	0	1
1	1	0

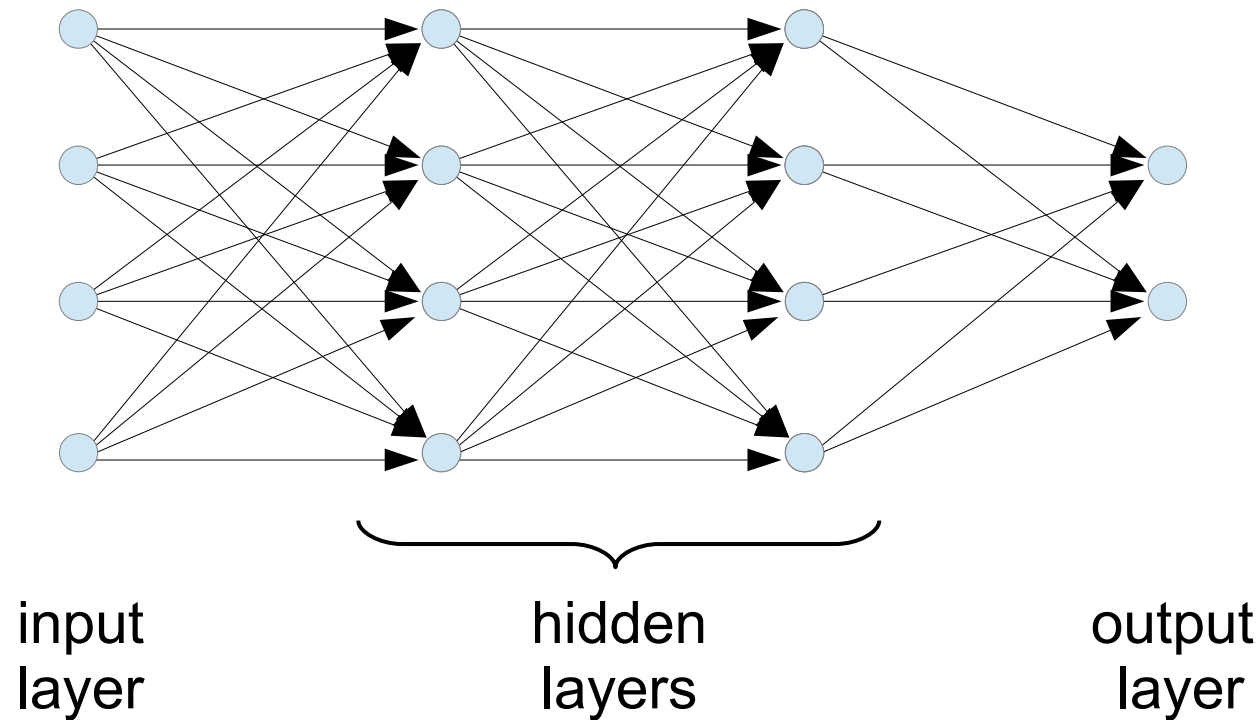
Can a Sigmoidal Unit Learn the XOR Function?

$x^{(1)}$	$x^{(2)}$	$x^{(1)} \text{ XOR } x^{(2)}$
0	0	0
0	1	1
1	0	1
1	1	0



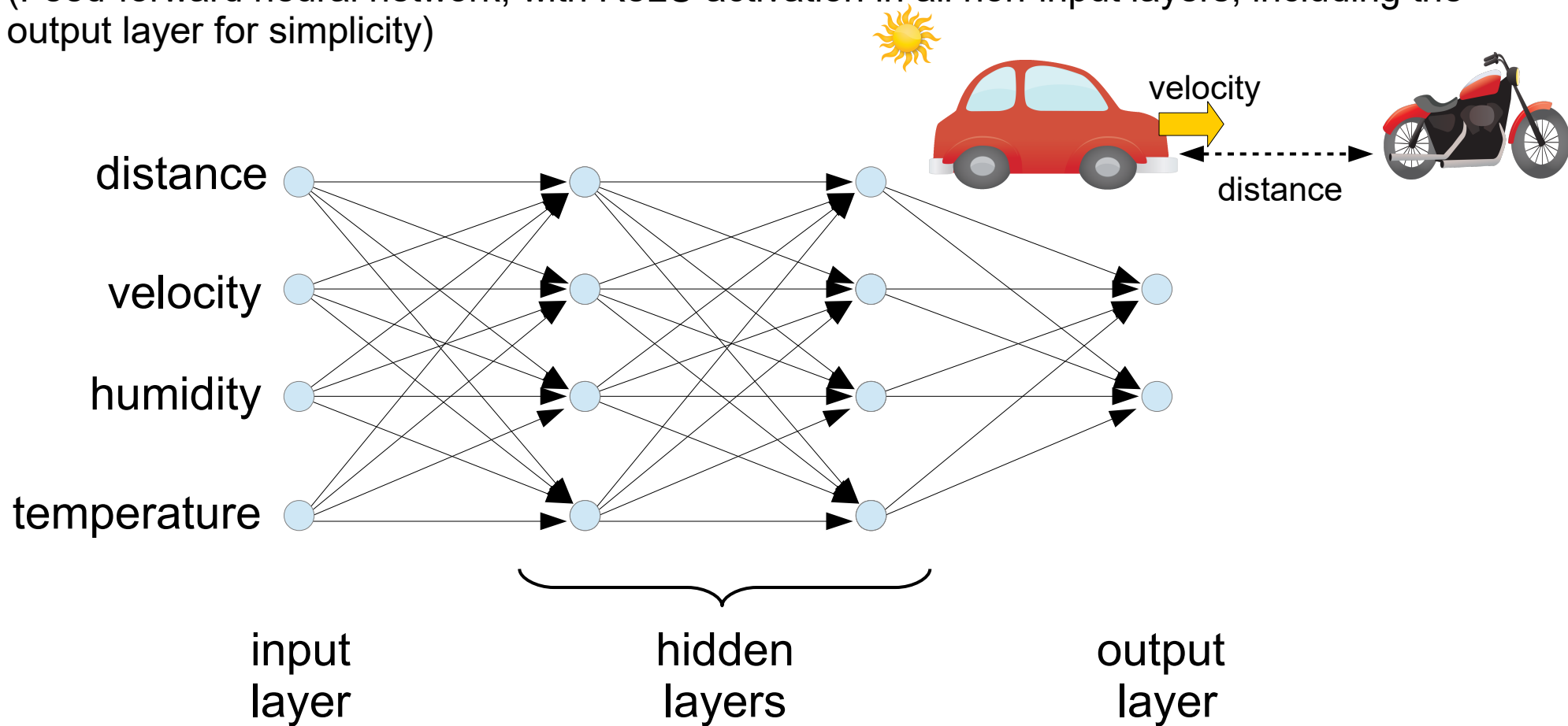
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



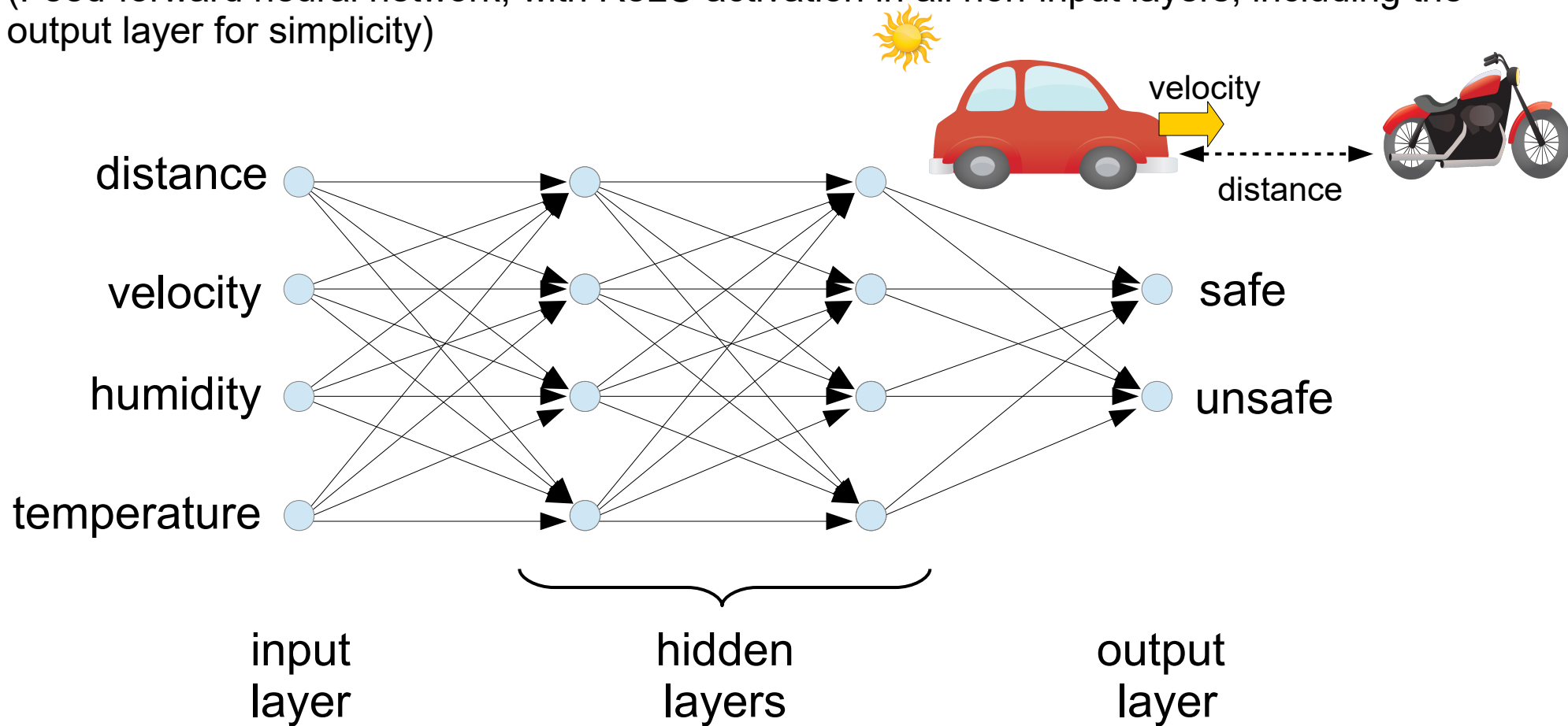
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



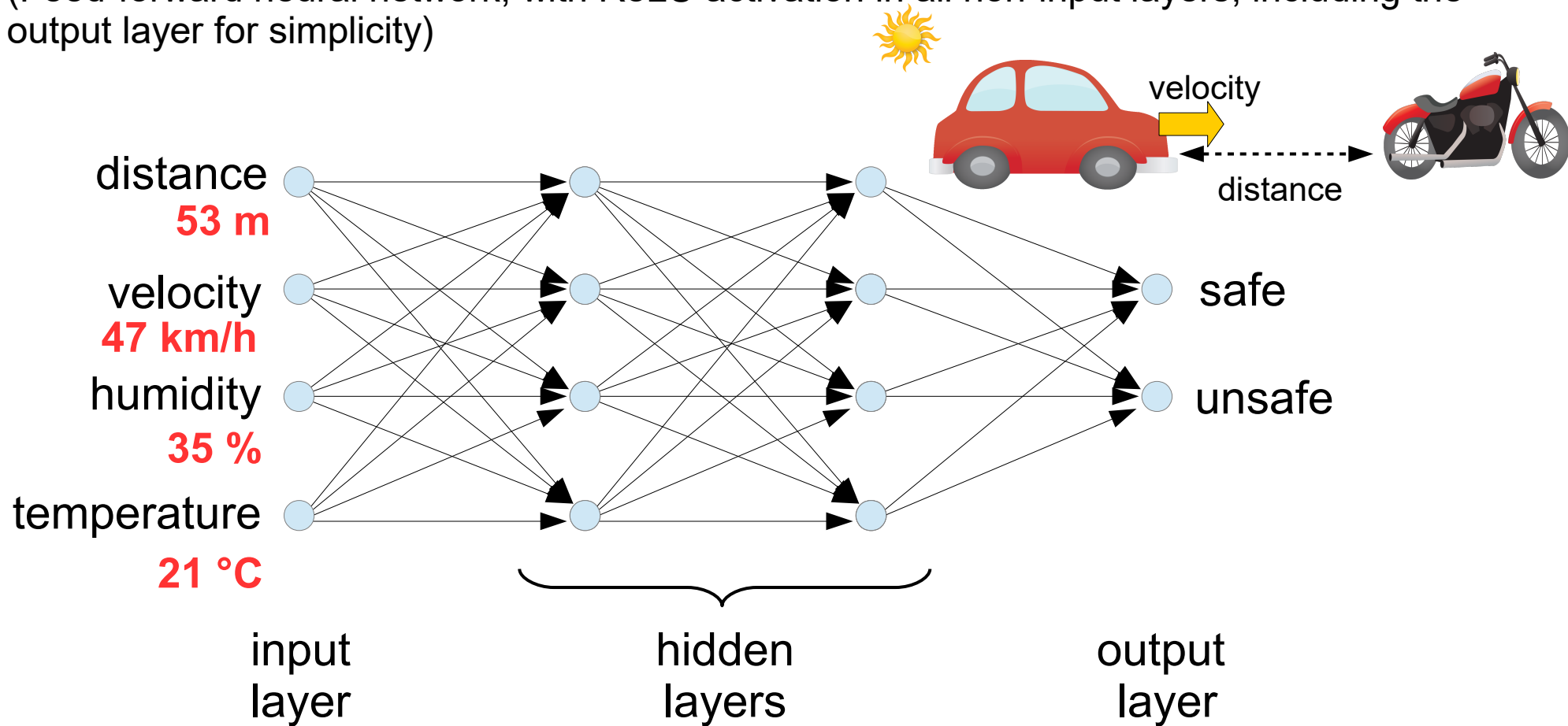
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



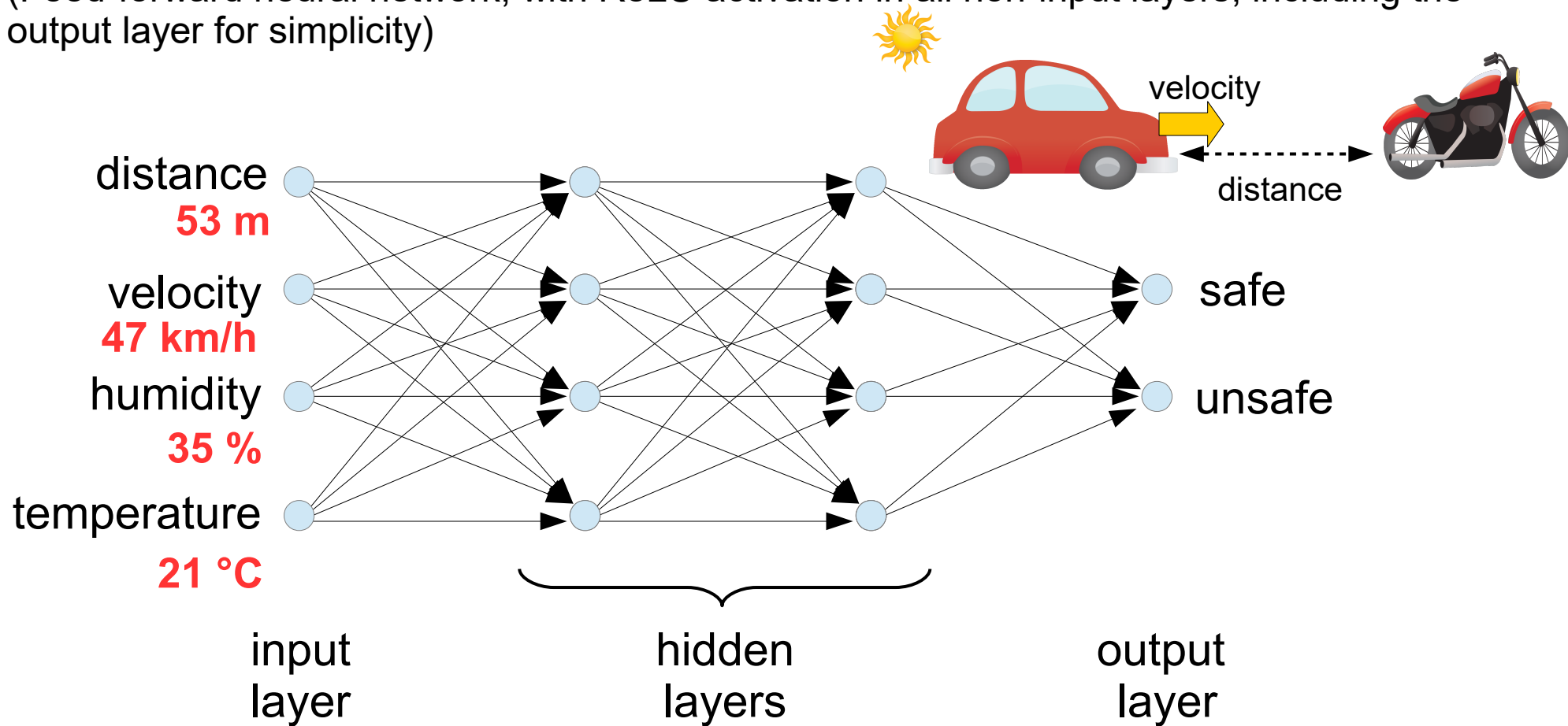
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



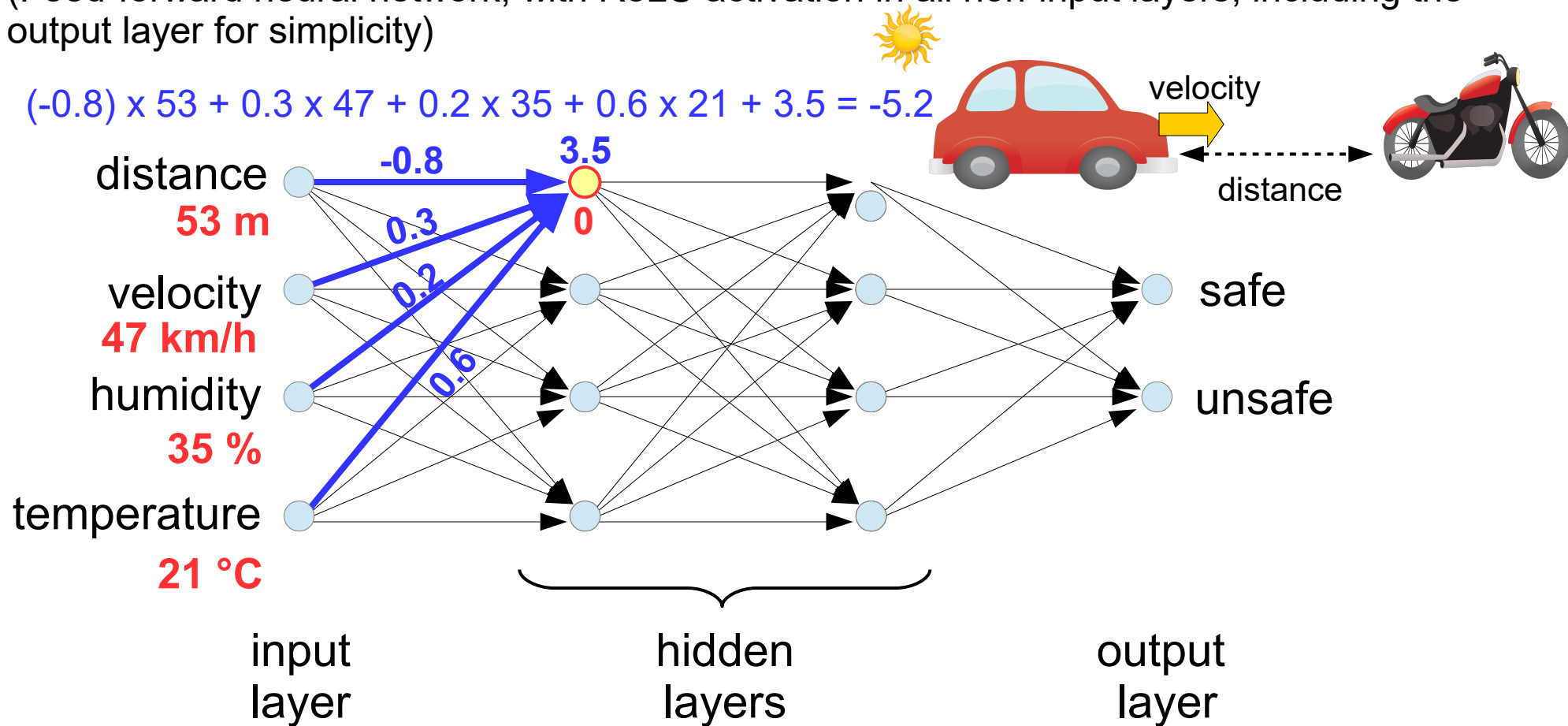
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



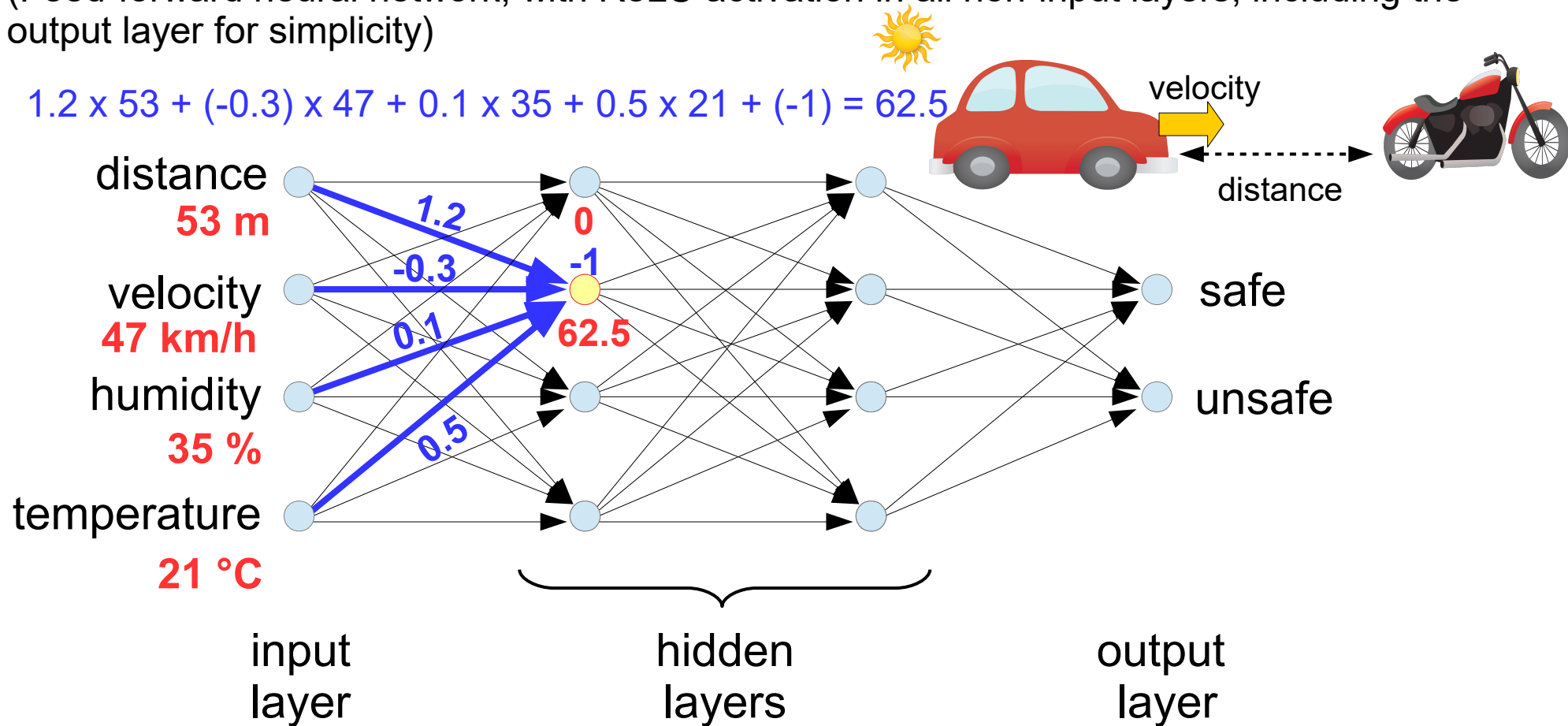
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



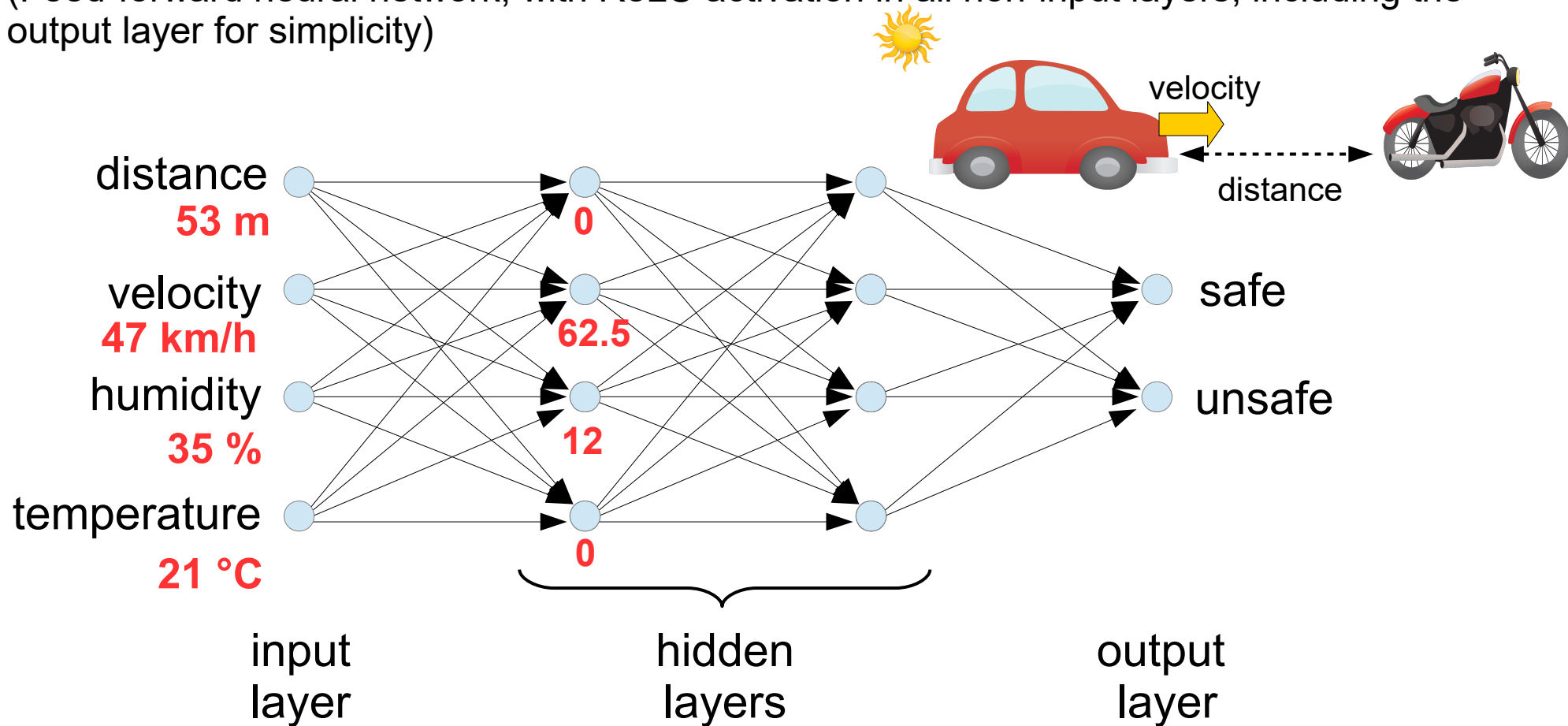
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



How does a Neural Network Work?

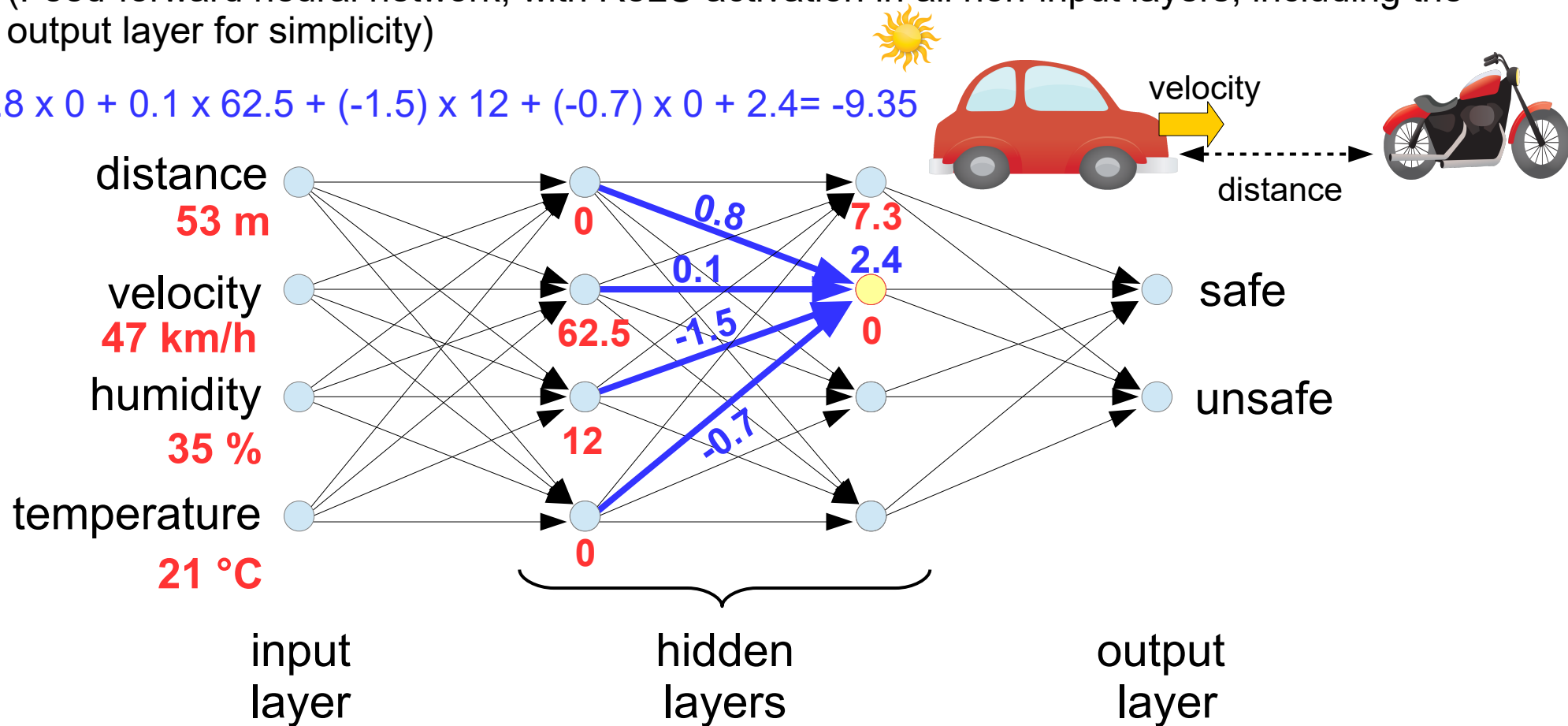
(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



How does a Neural Network Work?

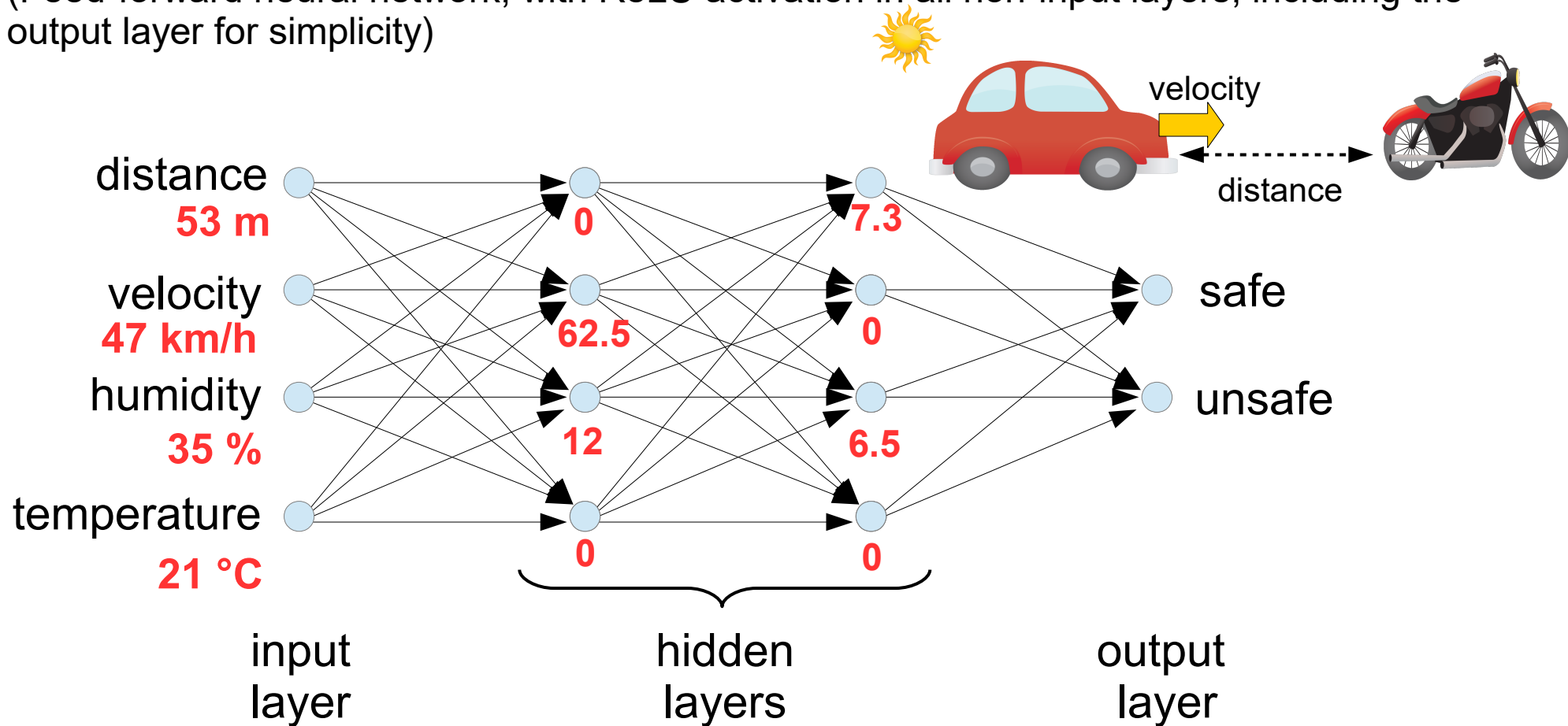
(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)

$$0.8 \times 0 + 0.1 \times 62.5 + (-1.5) \times 12 + (-0.7) \times 0 + 2.4 = -9.35$$



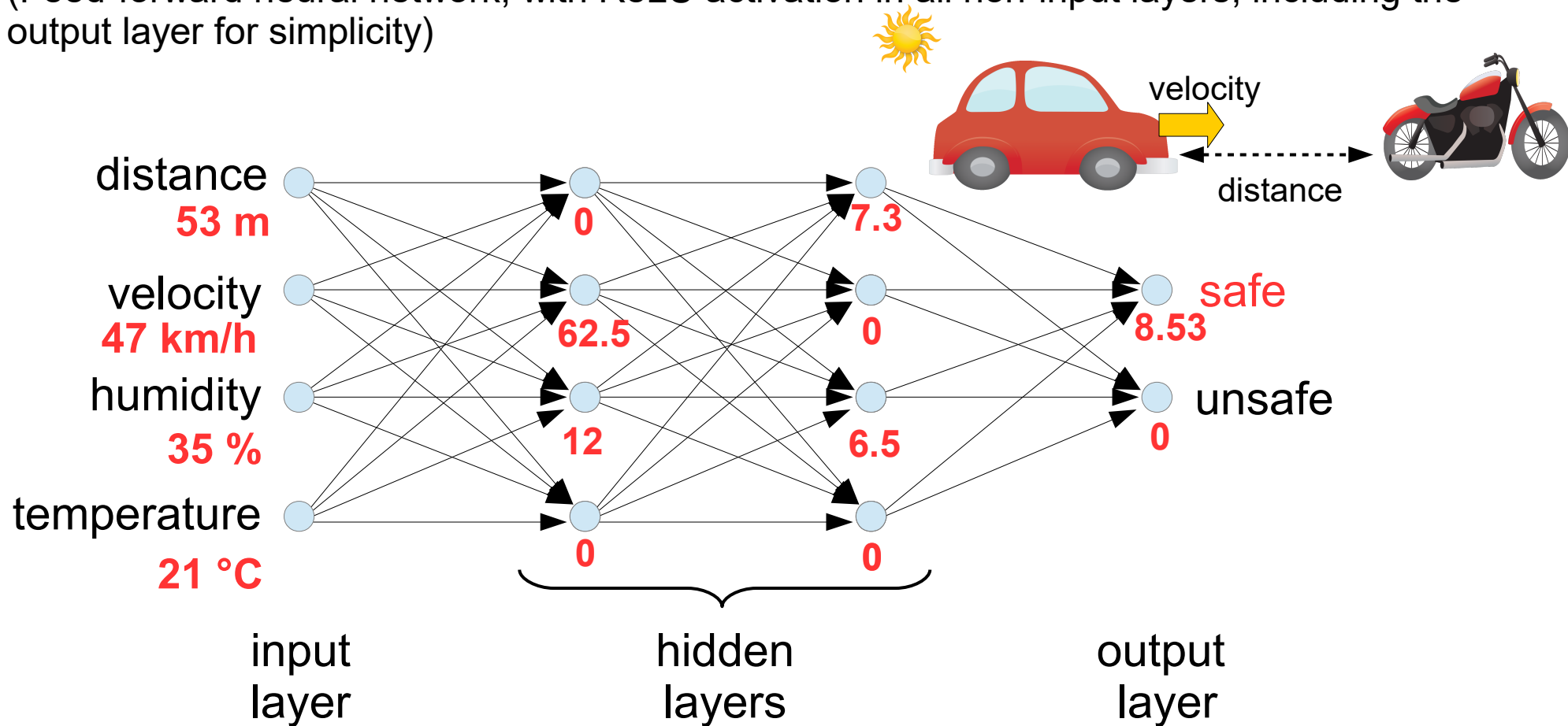
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



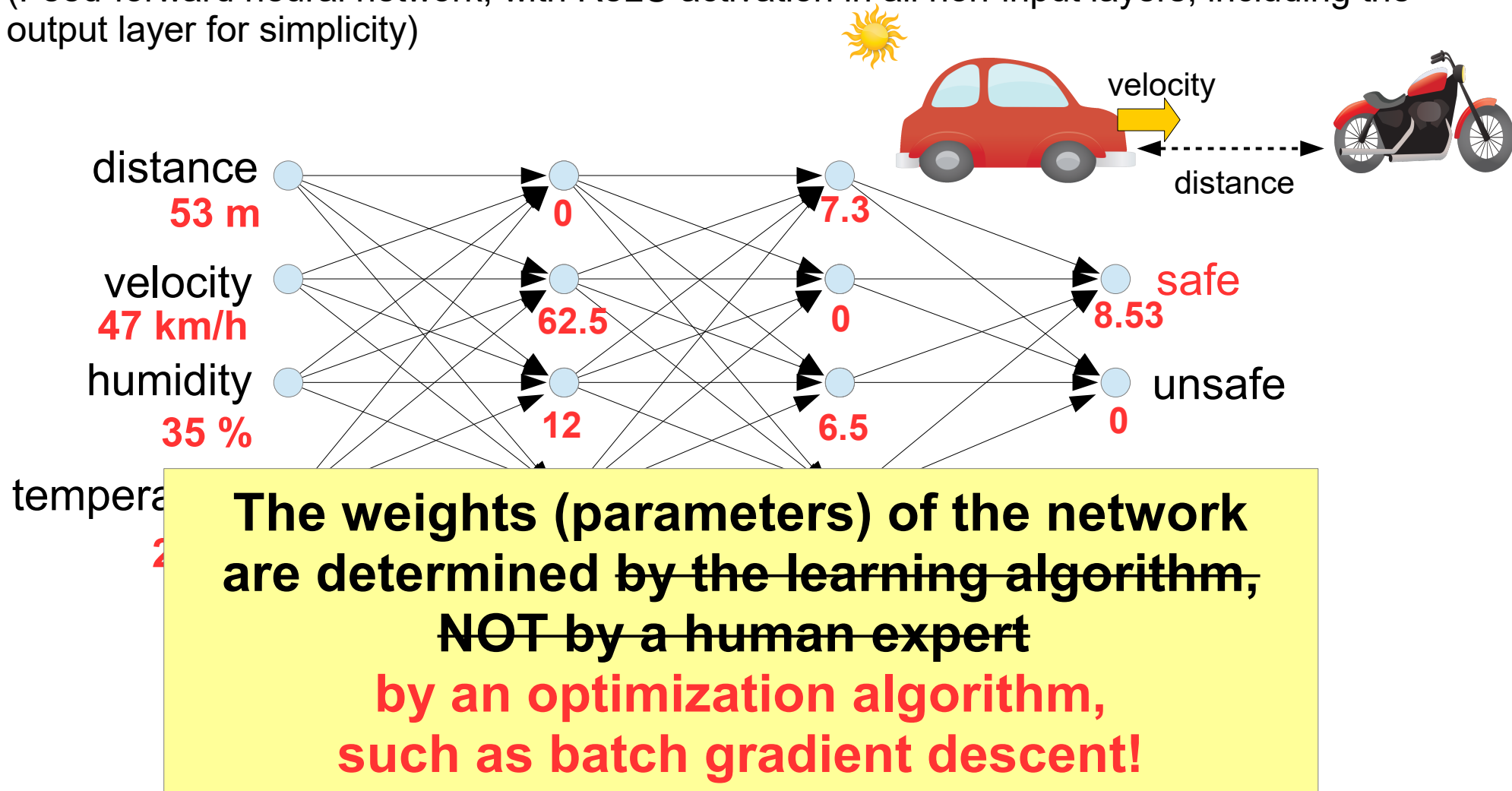
How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



How does a Neural Network Work?

(Feed-forward neural network, with ReLU activation in all non-input layers, including the output layer for simplicity)



How to Find Appropriate Weights?

- The error is a function of all the weights
- Calculate partial derivatives
- Use a gradient-based optimization technique (e.g. batch gradient descent)
- Backpropagation
 - Technique to calculate the gradients
 - Repeated application of the chain rule

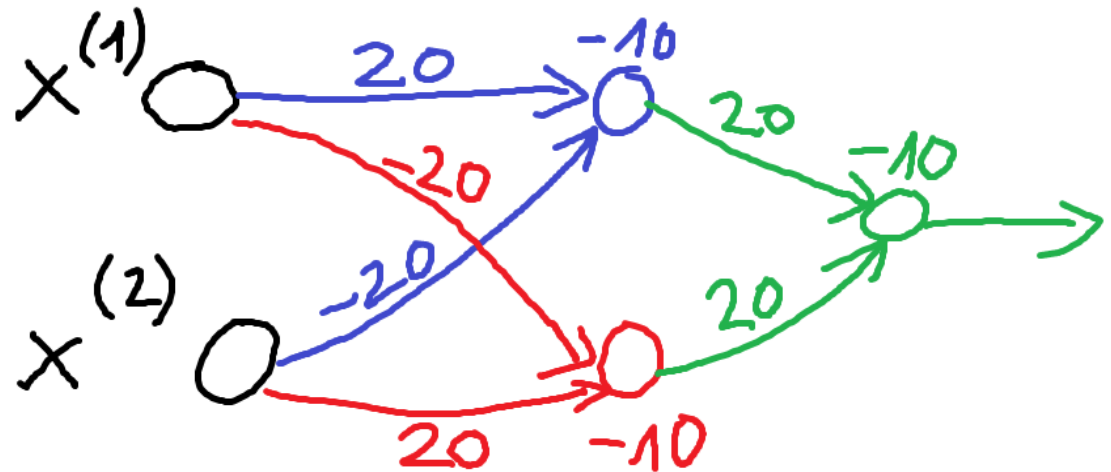
Neural Networks with One Hidden Layer

- Neural networks with one hidden layer with nonlinear activation function are **universal function approximators**.

Neural Networks with One Hidden Layer

- Neural networks with one hidden layer with nonlinear activation function are **universal function approximators**.
- Illustration for the case of binary data: XOR

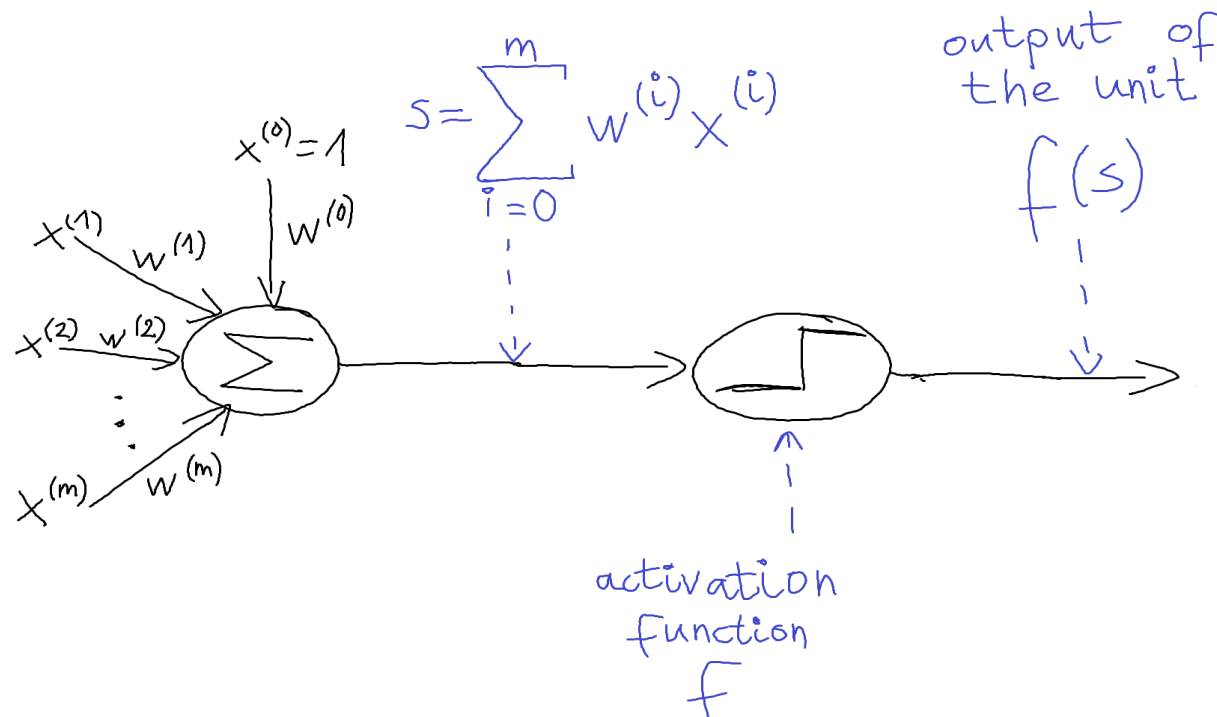
$X^{(1)}$	$X^{(2)}$	$X^{(1)} \text{ XOR } X^{(2)}$
0	0	0
0	1	1
1	0	1
1	1	0



Do We Need More Than One Hidden Layer?

- Neural networks with a single hidden layer are universal function approximators.
- But: in order to have a good approximation, we may need extremely large number of units if we only have one hidden layer.

How to Describe the Computations?



```

s = 0
for i = 0..m:
    s += x[i] * w[i]

```

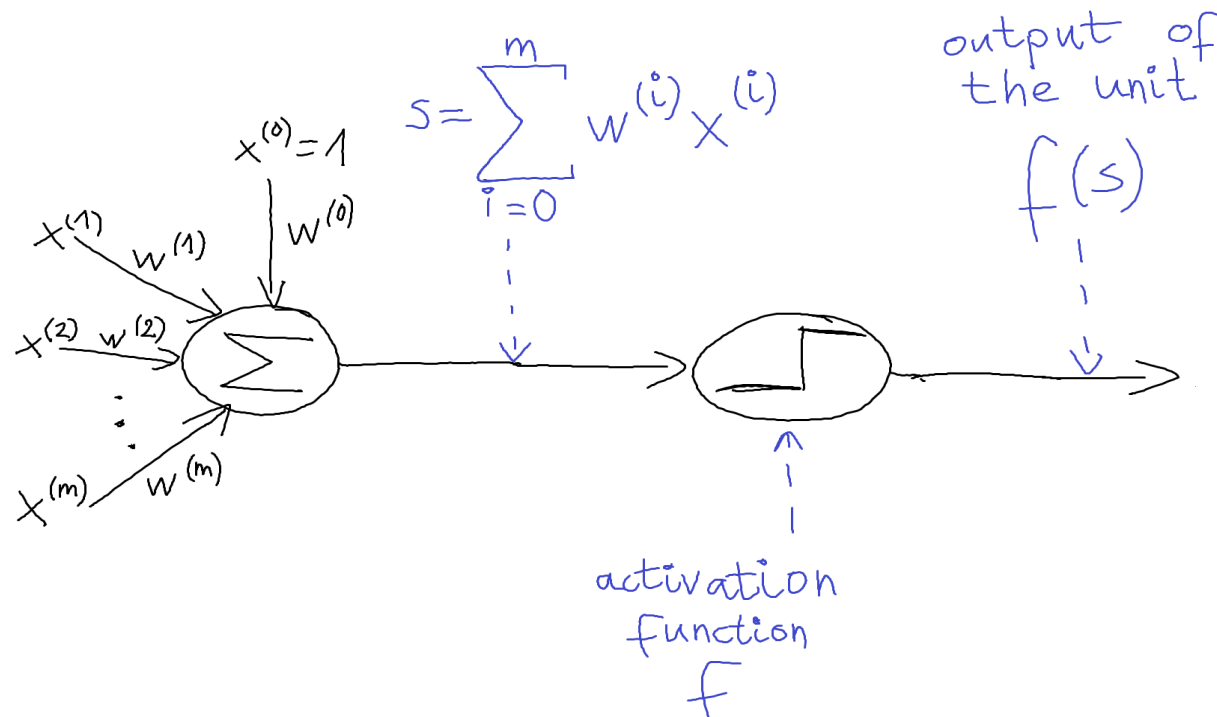
```

if s < 0:
    s = 0

```

$x^{(1)}, x^{(2)}, \dots, x^{(m)}$ = inputs of the unit
 $w^{(1)}, w^{(2)}, \dots, w^{(m)}$ = weights of $x^{(1)}, x^{(2)}, \dots, x^{(m)}$
 $w^{(0)}$ = bias weight

Compact Description of a Single Unit



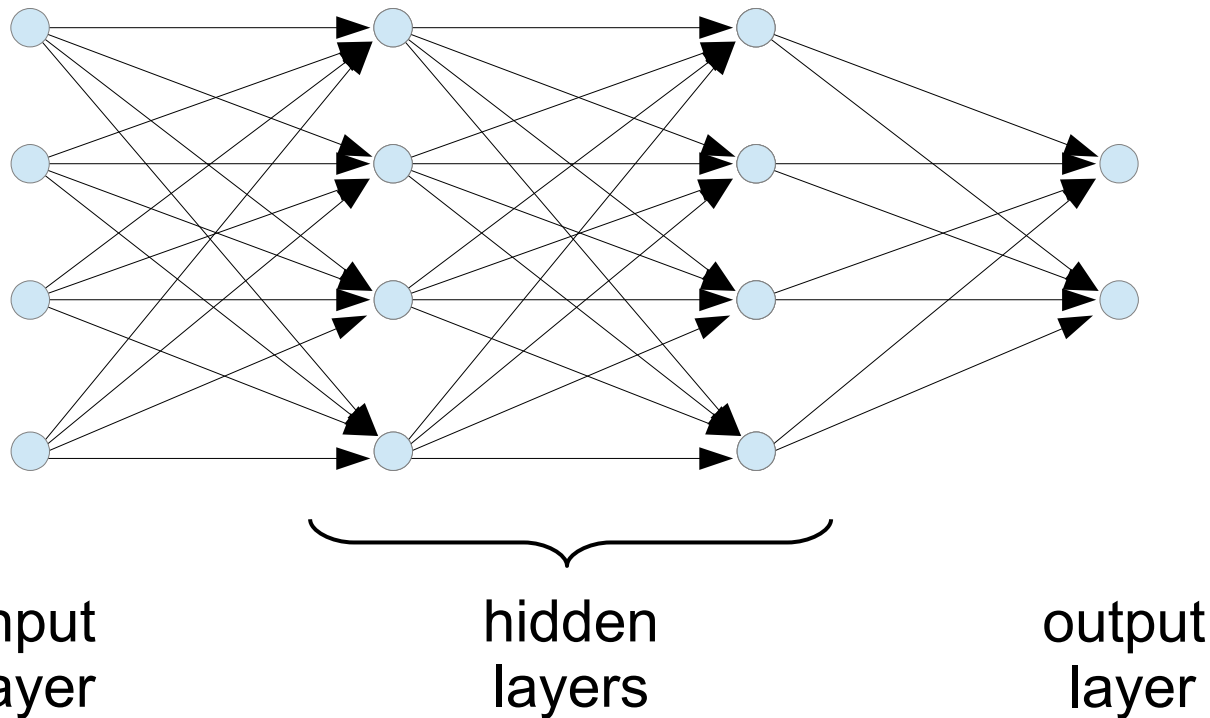
$$s = \text{np.dot}(x, w)$$

```
if s < 0:
    s = 0
```

$x^{(1)}, x^{(2)}, \dots, x^{(m)}$ = inputs of the unit
 $w^{(1)}, w^{(2)}, \dots, w^{(m)}$ = weights of $x^{(1)}, x^{(2)}, \dots, x^{(m)}$
 $w^{(0)}$ = bias weight

How to Describe the Computations?

dot product of
two vectors



```
a1 = np.dot(x, w1)
a2 = np.dot(x, w2)
...
```

```
if a1 < 0:
    a1 = 0
if a2 < 0:
    a2 = 0
...
```

```
b1 = np.dot(
    [a1, a2, a3, a4],
    w1_b )
...
```

Multiplication of Matrices

Example:

Multiplication of 2x2 matrices

$$\begin{bmatrix} 3 & 8 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} 0 & 2 \\ 5 & 7 \end{bmatrix}$$

Example:

Multiplication of 2x2 matrices

$$\begin{bmatrix} 3 & 8 \\ 2 & 0 \end{bmatrix} \times \begin{bmatrix} 0 & 2 \\ 5 & 7 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

Example:

Multiplication of 2x2 matrices

„row times column“

$$\begin{bmatrix} 3 & 8 \\ 2 & 0 \end{bmatrix} \times \begin{bmatrix} 0 & 2 \\ 5 & 7 \end{bmatrix} = \begin{bmatrix} \boxed{} & \\ & \end{bmatrix}$$

$$(3 \times 0) + (8 \times 5) = 0 + 40 = 40$$

Example:

Multiplication of 2x2 matrices

„row times column“

$$(3 \times 2) + (8 \times 7) = 6 + 56 = 62$$

$$\begin{bmatrix} 3 & 8 \\ 2 & 0 \end{bmatrix} \times \begin{bmatrix} 0 & 2 \\ 5 & 7 \end{bmatrix} = \begin{bmatrix} 40 & \boxed{} \\ & \end{bmatrix}$$

Example:

Multiplication of 2x2 matrices

„row times column“

$$\begin{bmatrix} 3 & 8 \\ 2 & 0 \end{bmatrix} \times \begin{bmatrix} 0 & 2 \\ 5 & 7 \end{bmatrix} = \begin{bmatrix} 40 & 62 \\ \boxed{} & \end{bmatrix}$$

$$(2 \times 0) + (0 \times 5) = 0 + 0 = 0$$

Example:

Multiplication of 2x2 matrices

„row times column“

$$\begin{bmatrix} 3 & 8 \\ 2 & 0 \end{bmatrix} \times \begin{bmatrix} 0 & 2 \\ 5 & 7 \end{bmatrix} = \begin{bmatrix} 40 & 62 \\ 0 & \boxed{} \end{bmatrix}$$

$$(2 \times 2) + (0 \times 7) = 4 + 0 = 4$$

Example:

Multiplication of 3x3 matrices

„row times column“

$$\begin{pmatrix} 8 & 6 & 12 \\ 3 & 9 & 1 \\ 0 & 7 & 21 \end{pmatrix} \times \begin{pmatrix} 1 & 9 & 2 \\ 5 & 10 & 1 \\ 6 & 3 & 7 \end{pmatrix} = \begin{pmatrix} & & \\ & & \\ & & \end{pmatrix}$$

$$(3 \times 1) + (9 \times 5) + (1 \times 6) = 3 + 45 + 6 = 54$$

Example:

Multiplication of 3x3 matrices

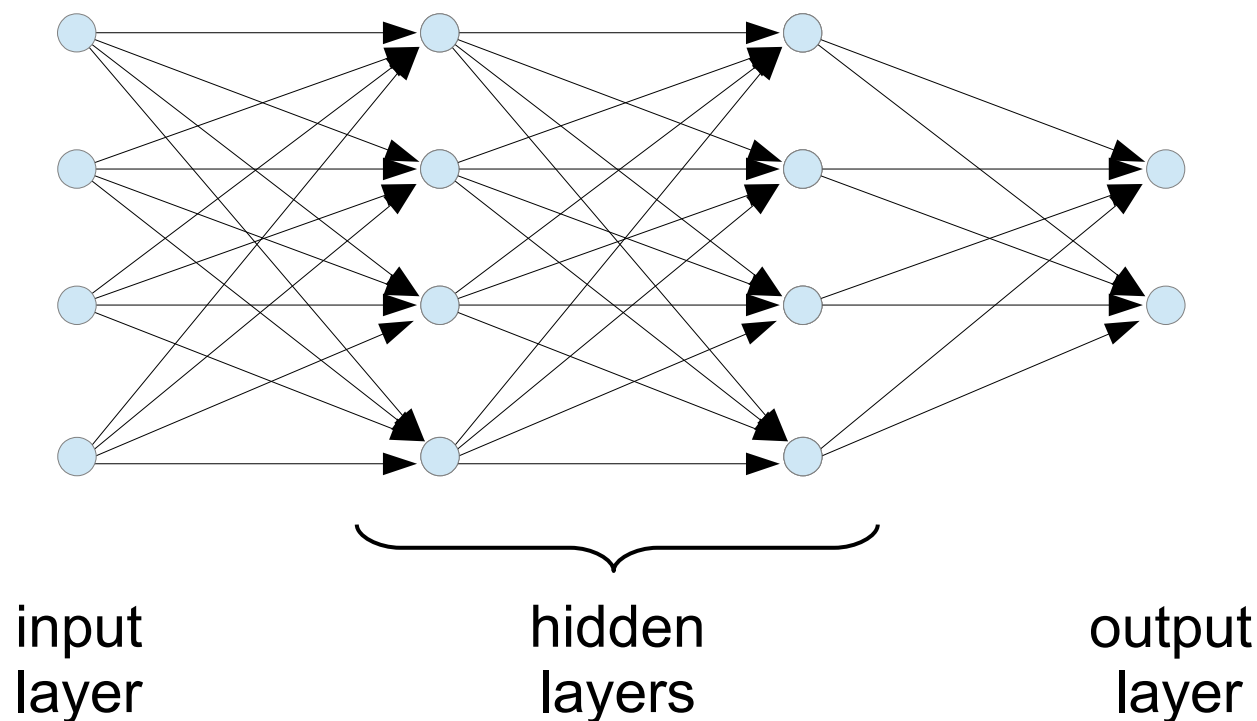
„row times column“

$$\begin{matrix} i\text{-th row} \rightarrow & \begin{bmatrix} 8 & 6 & 12 \\ 3 & 9 & 1 \\ 0 & 7 & 21 \end{bmatrix} & \times & \begin{matrix} j\text{-th column} \downarrow \\ \begin{bmatrix} 1 & 9 & 2 \\ 5 & 10 & 1 \\ 6 & 3 & 7 \end{bmatrix} \end{matrix} & = & \begin{bmatrix} & & \\ & \boxed{} & \\ & & \end{bmatrix} \end{matrix}$$

Element in the i -th row and j -th column:
product of the i -th row of the first matrix
and the j -th column of the second matrix

Compact Description of the Computations with Matrix Multiplication

How to Describe the Computations? ^{(row vector) x (matrix)}



$\mathbf{a} = \text{np.dot}(\mathbf{x}, \mathbf{W1})$

$\mathbf{a}[\mathbf{a} < 0] = 0$

$\mathbf{b} = \text{np.dot}(\mathbf{a}, \mathbf{W2})$

$\mathbf{b}[\mathbf{b} < 0] = 0$

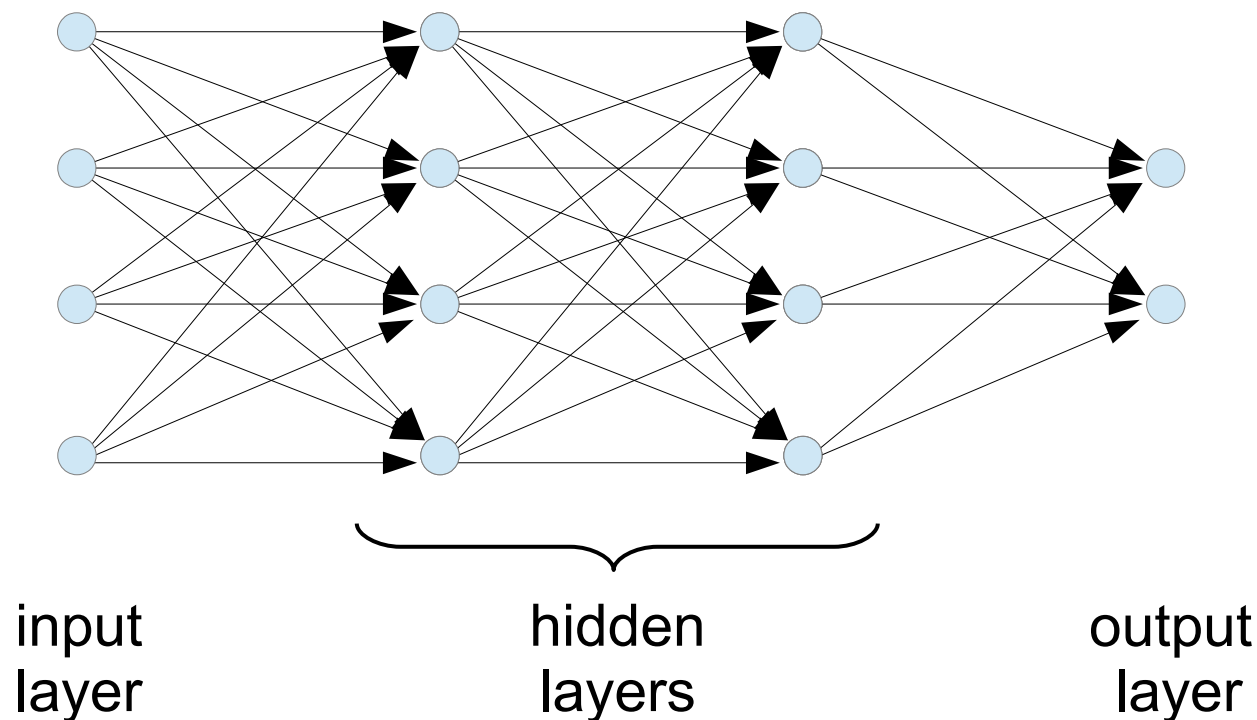
$\mathbf{c} = \text{np.dot}(\mathbf{b}, \mathbf{W3})$

$\mathbf{c}[\mathbf{c} < 0] = 0$

\mathbf{x} , \mathbf{a} , \mathbf{b} , \mathbf{c} are row vectors, not scalars!

$\mathbf{W1}$, $\mathbf{W2}$ and $\mathbf{W3}$ are matrices containing all the weights between (W1) input layer and first hidden layer, (W2) first hidden layer and second hidden layer, and (W3) second hidden layer and output layer, respectively. Each column of $\mathbf{W1}$, $\mathbf{W2}$ and $\mathbf{W3}$ contains all the weights of the incoming connections associated with the same hidden unit, e.g., first column of $\mathbf{W1}$ contains the weights of the first unit in the first hidden layer.

How to Describe the Computations?



(matrix) x (matrix)

```
A = np.dot(X, W1)
A[A<0]=0
```

```
B = np.dot(A, W2)
B[B<0]=0
```

```
C = np.dot(B, W3)
C[C<0]=0
```

X, **A**, **B**, **C** are matrices, neither vectors, nor scalars!

W1, **W2** and **W3** are matrices containing all the weights between (W1) input layer and first hidden layer, (W2) first hidden layer and second hidden layer, and (W3) second hidden layer and output layer, respectively. Each column of **W1**, **W2** and **W3** contains all the weights of the incoming connections associated with the same hidden unit, e.g., first column of W1 contains the weights of the first unit in the first hidden layer.

Deep Learning in a Nutshell

„What was wrong with backpropagation in 1986?“

(Geoff Hinton, „Deep Learning“, May 22, 2015)

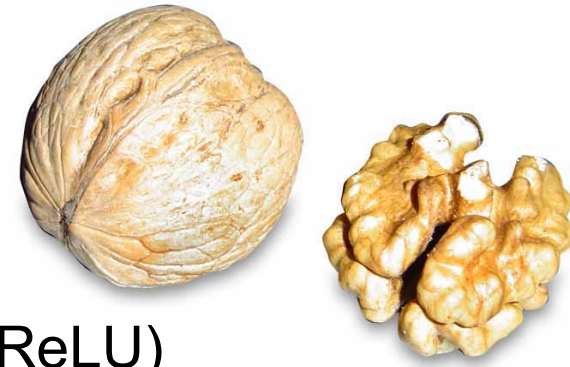
- Our labeled datasets were thousands of times too small.
- Our computers were millions of times too slow.
- We initialized the weights in a stupid way.
- We used the wrong type of non-linearity.



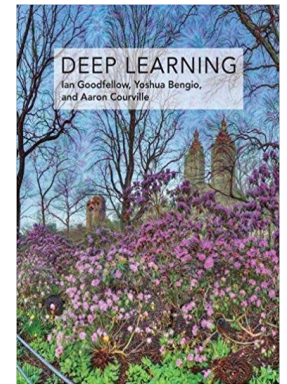
Image: https://upload.wikimedia.org/wikipedia/commons/3/34/Geoffrey_Hinton_at_UBC.jpg
Eviatar Bach [CC BY-SA 3.0 (<https://creativecommons.org/licenses/by-sa/3.0>)]

Deep Learning in a Nutshell

By Horst Frank - photo taken by Horst Frank, CC BY-SA 3.0,
<https://commons.wikimedia.org/w/index.php?curid=14209>



- Size and structure of the network:
few layers \rightarrow many layers
- Activation function: sigmoid \rightarrow rectified linear unit (ReLU)
- Loss function: quadratic loss \rightarrow cross-entropy
- Initialization of weights: random \rightarrow (unsupervised) pre-training
- Size of training data, much more memory,
distributed computation (HPC), GPUs, TPUs...
- New regularization techniques:
“sparsity-enforcing” regularisation terms,
drop-out, early stop,
parameter sharing (convolutional networks)
- New optimization techniques (e.g. ADAM)



<http://www.deeplearningbook.org>

Summary

Summary

- Mathematical operation: Matrix multiplication
- Neural Unit \rightarrow neural network
- Neural networks with one hidden layer with nonlinear activations are universal function approximators
- Overview of the techniques known under the umbrella of „deep learning“

Essential Concepts

- Decision boundary
- Neural Network, Feed-forward neural network
- Neural Unit
- Activation Function – Sigmoid, ReLU
- Matrix multiplication
- GPU, TPU („AI chip“)