Neural Networks with Fully Connected Layers in Practice

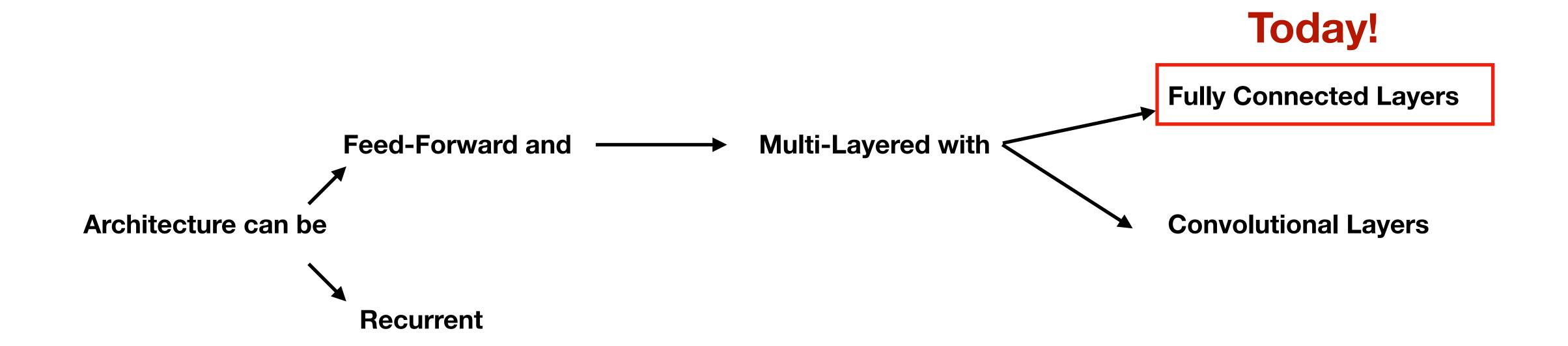
Fundamentals of Artificial Intelligence Fabien Cromieres Kyoto University

Program for today:

- 1- Discuss the mathematical representation of Fully-connected layers in Neural Networks
- 2- Define and train a real Neural Network in a Jupyter Notebook

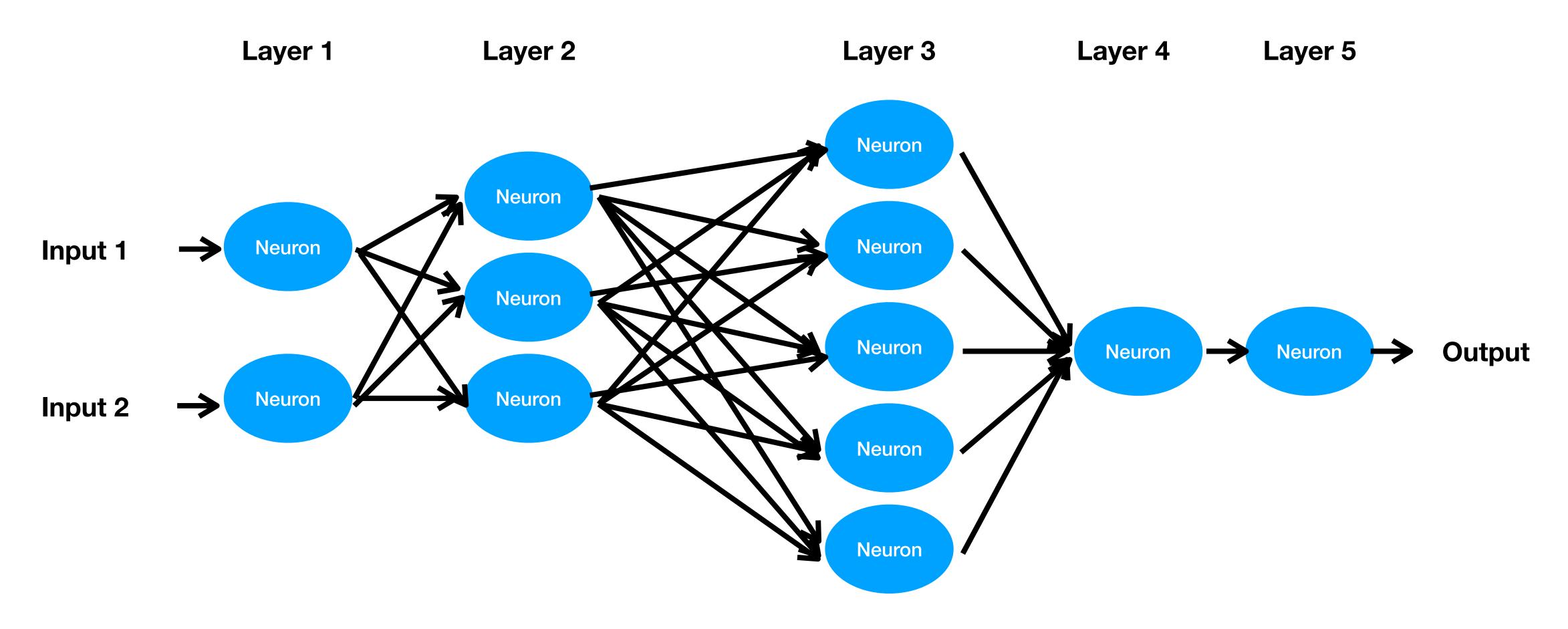
Neural Network Architectures

We are still here:



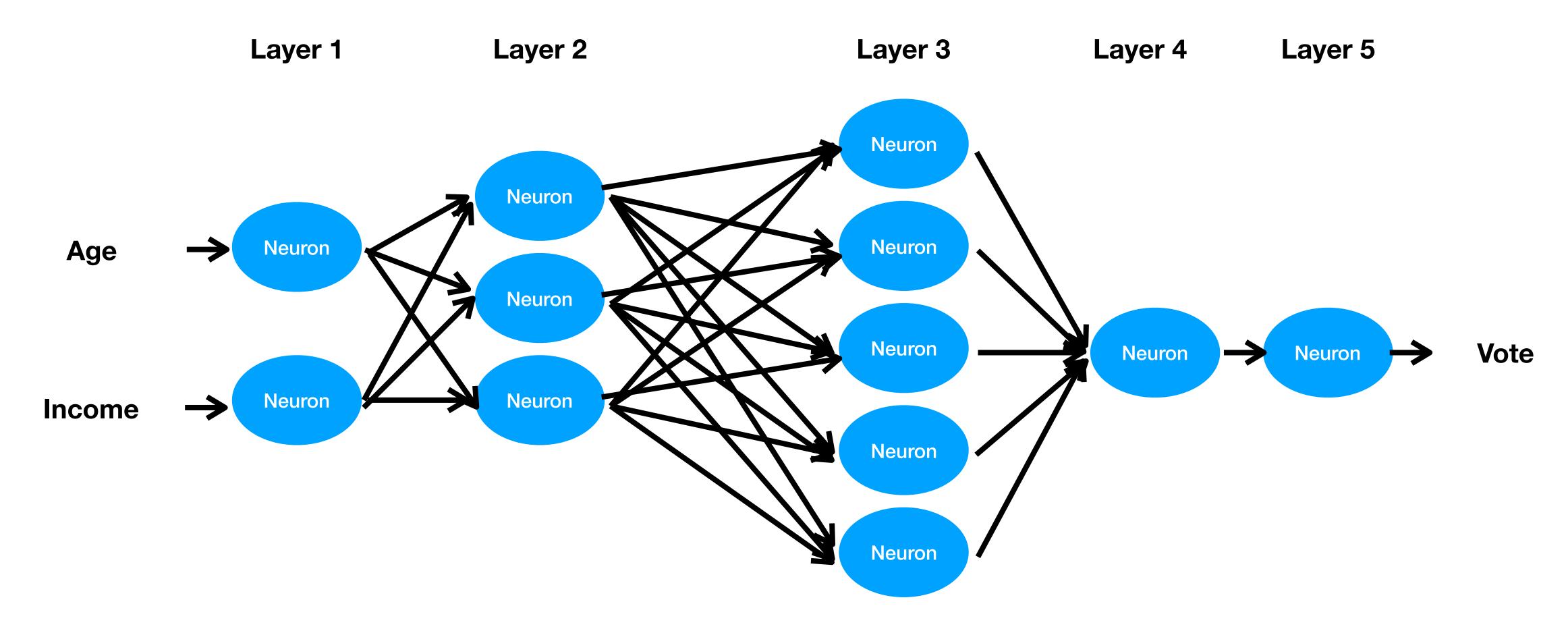
Feed-Forward networks with fully connected layers

• Therefore, we are going to consider this type of Neural Network:



Feed-Forward networks with fully connected layers

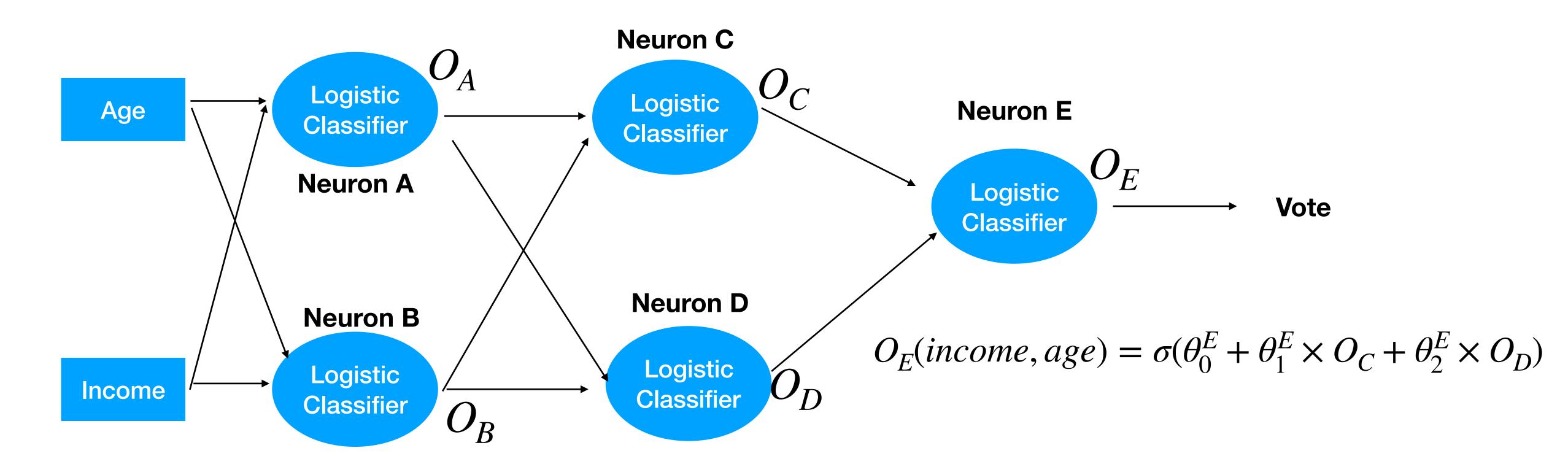
• Therefore, we are going to consider this type of Neural Network:



Keeping in mind what this type of graph mean

$$O_A(income, age) = \sigma(\theta_0^A + \theta_1^A \times income + \theta_2^A \times age)$$

$$O_C(income, age) = \sigma(\theta_0^C + \theta_1^C \times O_A + \theta_2^C \times O_B)$$



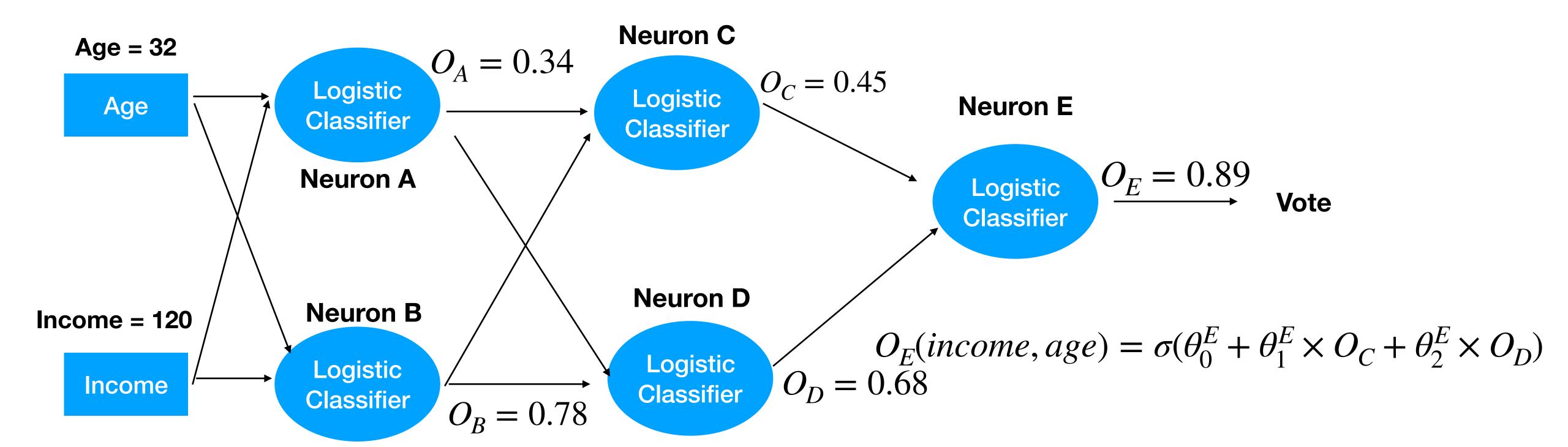
$$O_B(income, age) = \sigma(\theta_0^B + \theta_1^B \times income + \theta_2^B \times age)$$

$$O_D(income, age) = \sigma(\theta_0^D + \theta_1^D \times O_A + \theta_2^D \times O_B)$$

Keeping in mind what this type of graph mean

$$O_A(income, age) = \sigma(\theta_0^A + \theta_1^A \times income + \theta_2^A \times age)$$

$$O_C(income, age) = \sigma(\theta_0^C + \theta_1^C \times O_A + \theta_2^C \times O_B)$$



$$O_B(income, age) = \sigma(\theta_0^B + \theta_1^B \times income + \theta_2^B \times age)$$

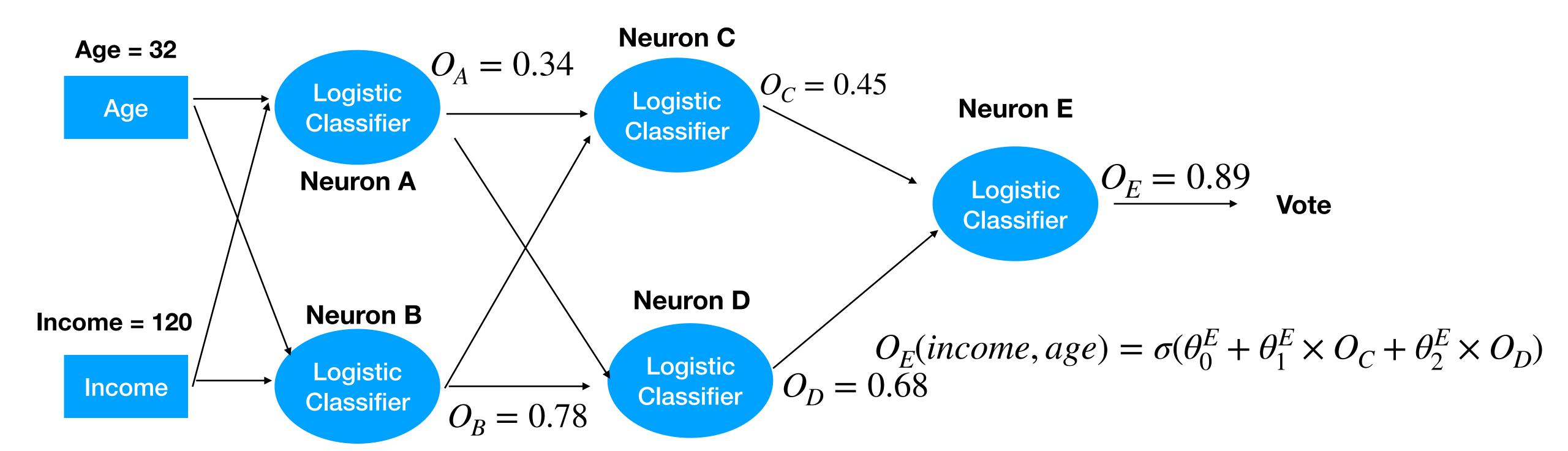
$$O_D(income, age) = \sigma(\theta_0^D + \theta_1^D \times O_A + \theta_2^D \times O_B)$$

Keeping in mind what this type of graph mean

-> Each Neural Network architecture defines a function of the input with parameters θ

$$O_A(income, age) = \sigma(\theta_0^A + \theta_1^A \times income + \theta_2^A \times age) \qquad O_C(income, age) = \sigma(\theta_0^C + \theta_1^C \times O_A + \theta_2^C \times O_B)$$

$$O_C(income, age) = \sigma(\theta_0^C + \theta_1^C \times O_A + \theta_2^C \times O_B)$$

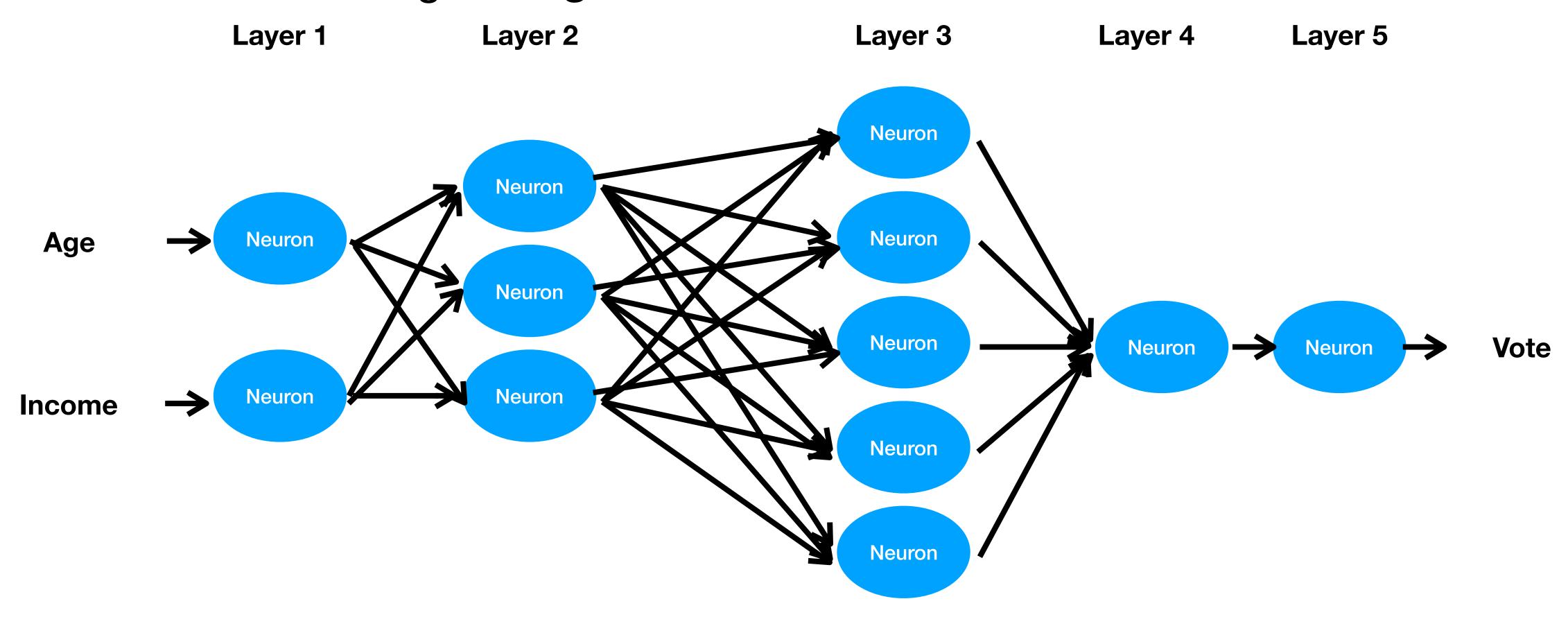


$$O_B(income, age) = \sigma(\theta_0^B + \theta_1^B \times income + \theta_2^B \times age)$$

$$O_D(income, age) = \sigma(\theta_0^D + \theta_1^D \times O_A + \theta_2^D \times O_B)$$

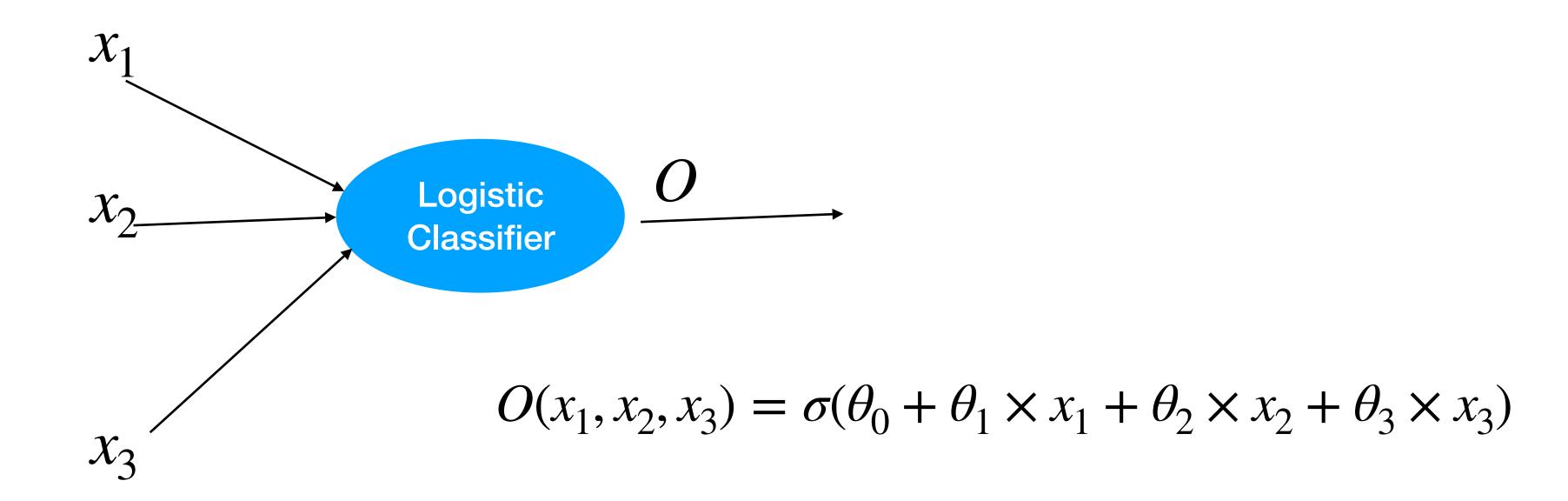
Feed-Forward networks with fully connected layers

- -> Each Neural Network architecture defines a function of the input with parameters θ
 - Therefore, this is just a <u>visual way</u> of defining a <u>complicated parameterized</u> function of *Vote* given *Age* and *Income*:



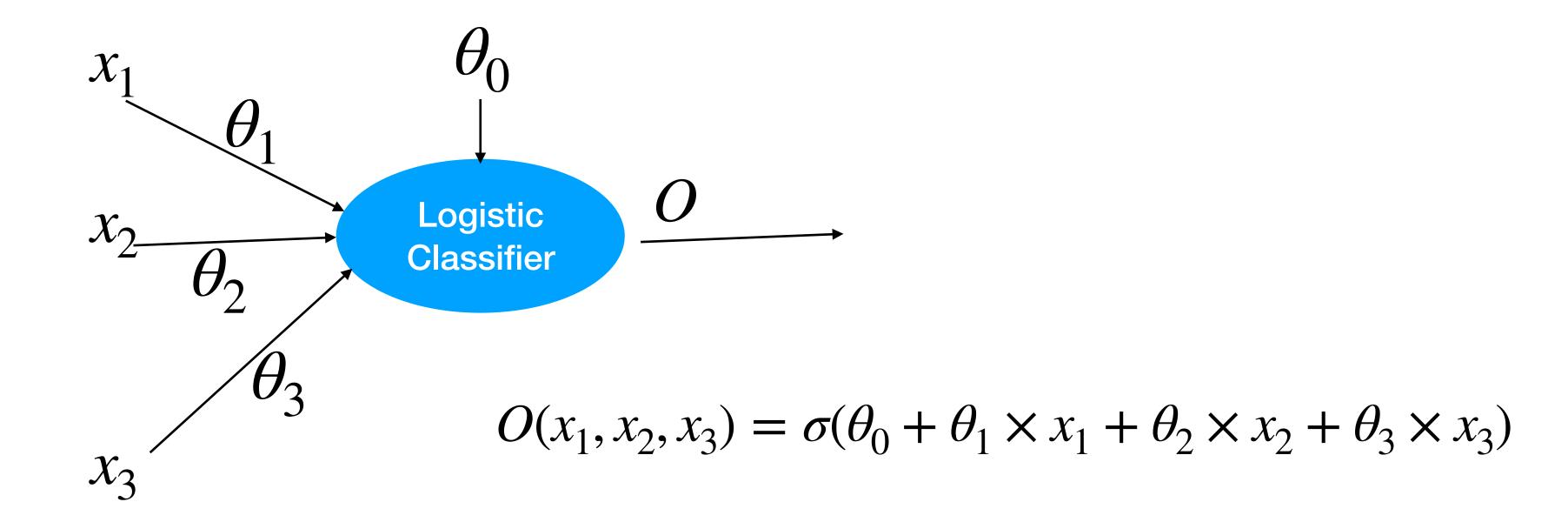
Parameters

• If a neuron has N inputs, it has N+1 parameters



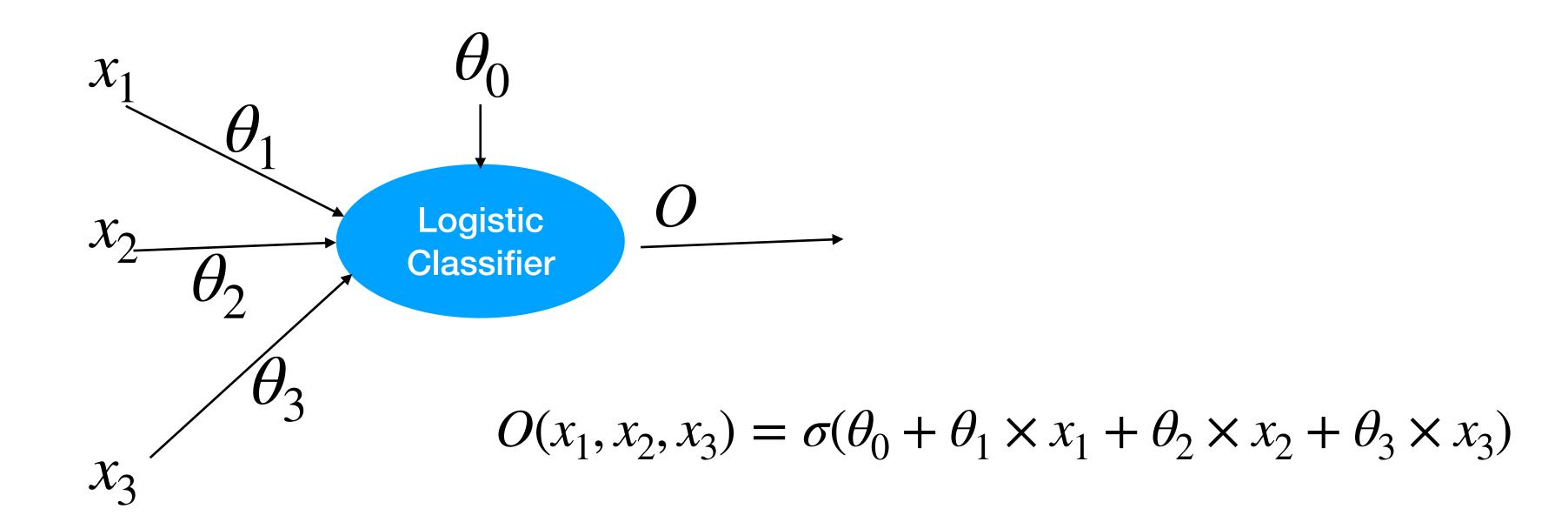
Parameters

- If a neuron has N inputs, it has N+1 parameters
 - Visually, we can associate a parameter to each input, and show θ_0 separately



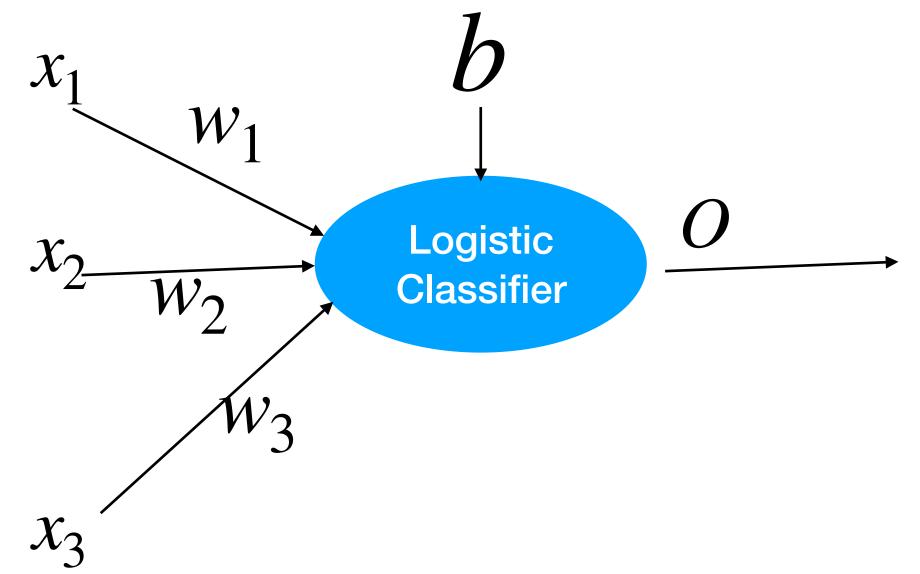
Parameters: Terminology

- θ_1 , θ_2 , θ_3 are often called the *weights* of the neuron
- θ_0 is often called the *bias* of the neuron



Parameters: Terminology

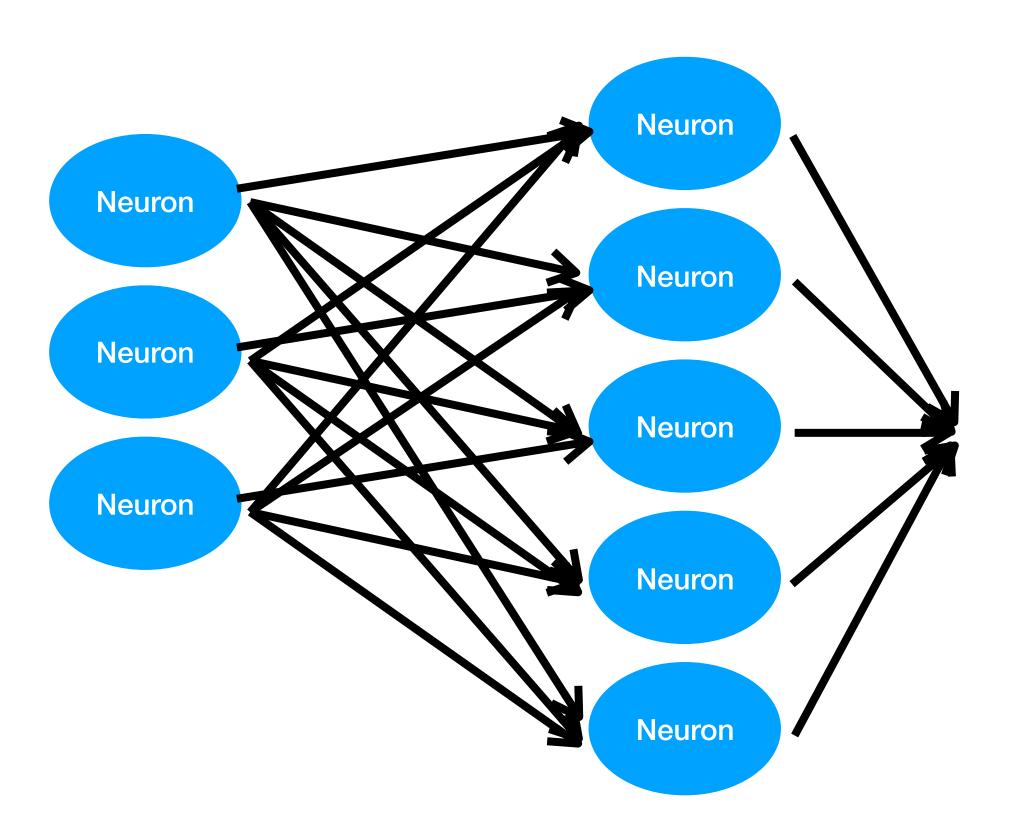
- θ_1 , θ_2 , θ_3 are often called the *weights* of the neuron
 - They are therefore often also noted w₁, w₂, w₃
- θ_0 is often called the **bias** of the neuron
 - It is often noted **b**



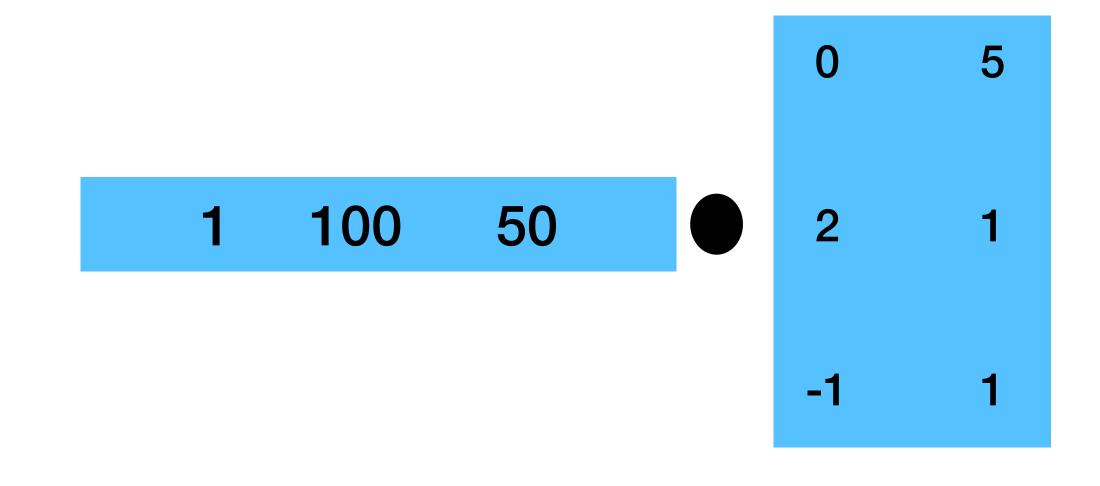
$$O(x_1, x_2, x_3) = \sigma(b + w_1 \times x_1 + w_2 \times x_2 + w_3 \times x_3)$$

"Connectionist" view vs "Mathematical" view

- "Connections" viewpoint
 - Fully Connected Layers

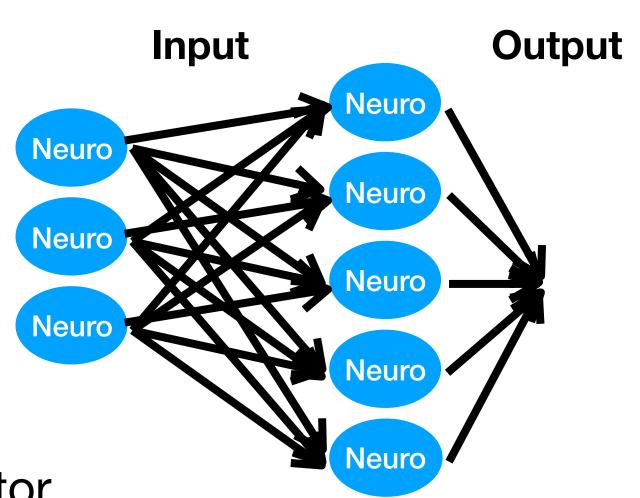


- Math operations viewpoint
 - Matrix multiplications



"Connectionist" view vs "Mathematical" view

- In practice, computing the output of a fully connected layer is equivalent to computing:
 - a Matrix multiplication
 - followed by an activation
- Equivalences:
 - Input of the layer (= output of the previous layer) <-> row vector
 - Output of the layer (= set of output of each neuron in a layer) <-> row vector
 - Weights of a neuron <-> column vector
 - Set weights of all neurons in a fully connected layer <-> matrix
 - Batched input of a layer (= several inputs at once) <-> matrix



Linear Algebra

- In order to discuss the way Neural Networks are actually implemented, we need to discuss some mathematical concepts
- We discuss:
 - Vector scalar products (a.k.a inner product)
 - Matrices
 - Matrix-vector multiplication
 - Matrix-Matrix multiplication
- Warning: it is a bit ambitious to explain that in 30 minutes. Hopefully you understand most of it.

Row Vectors and Column Vectors

- We have seen that *vectors* are just a "list" of numbers
- We are actually going to distinguish 2 "types" of vectors: row vectors and column vectors



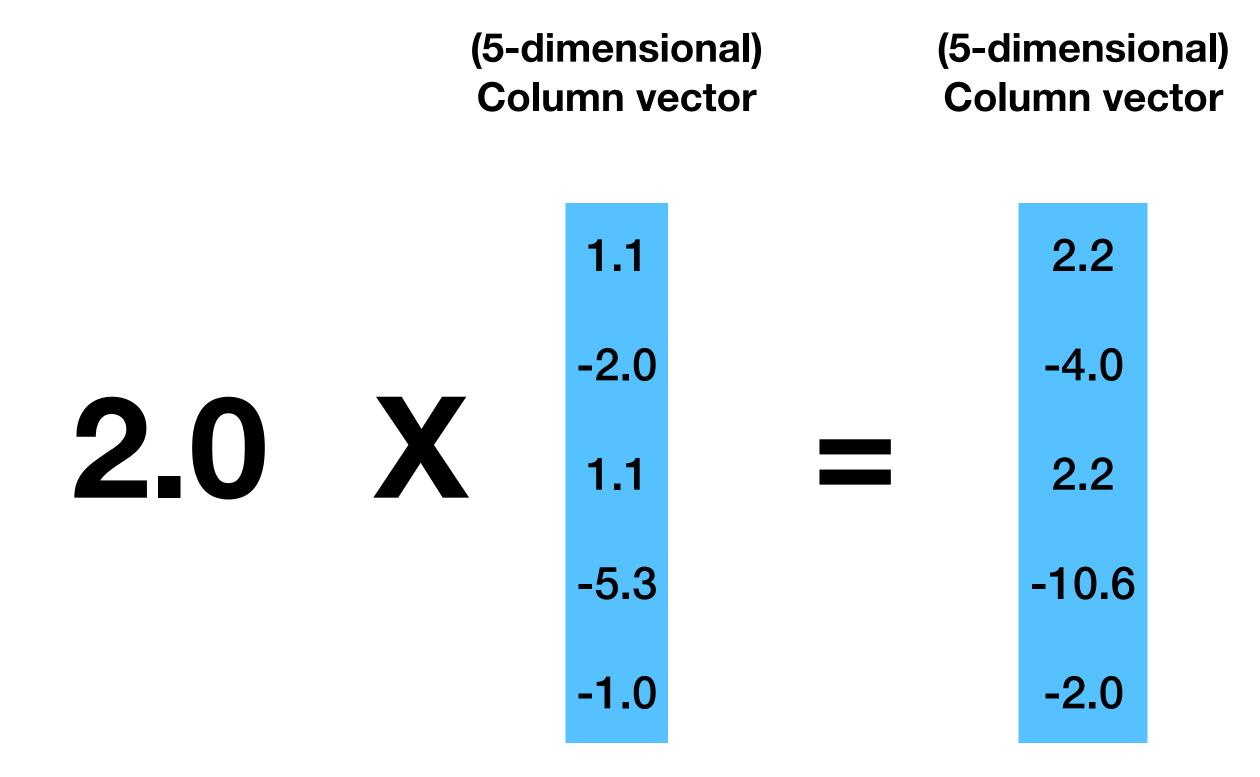
Column vectors

Column vectors of the same dimension can be added

(5-dime Columr		-	_	dimensional) olumn vector		(5-dimensional Column vectors		-
	0.2			1.1			1.3	
	3.6			-2.0			1.6	
	2.1	+		1.1			3.2	
	5.3			-5.3			0.0	
	-2.2			-1.0			-3.2	

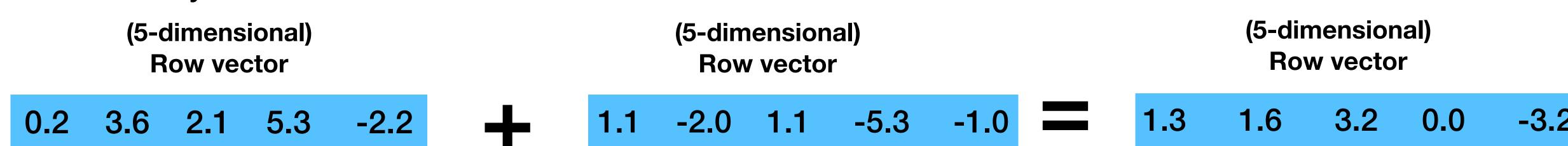
Column vectors

Column vectors can be multiplied by a number



Row vectors

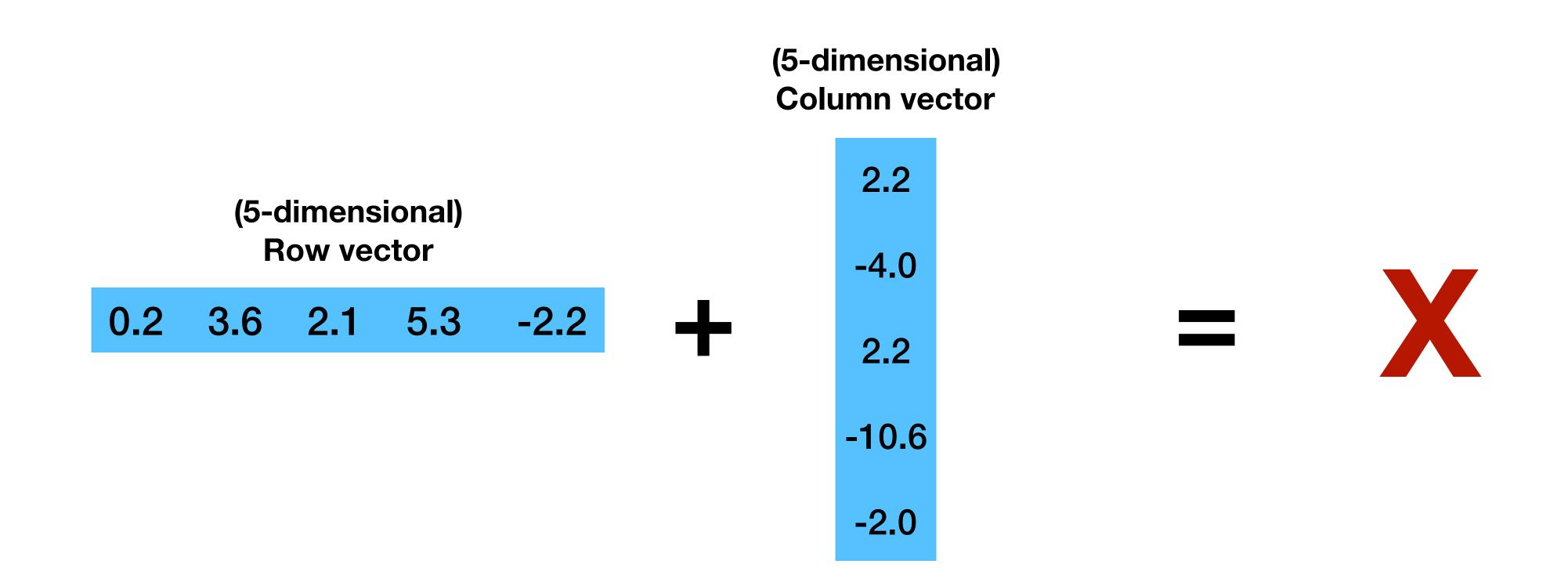
- Row vectors have the same operations as Column vectors
- They can be added:



They can be multiplied by a number:

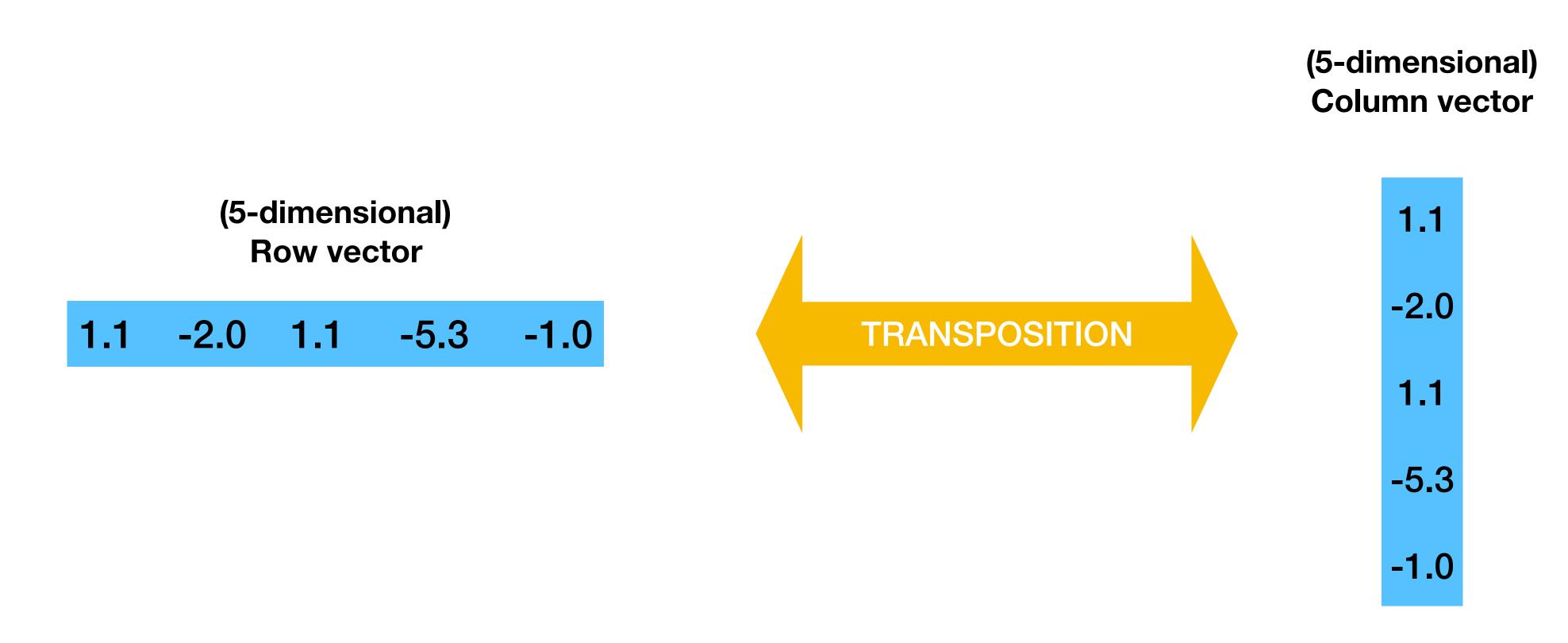
Row vectors and Column Vectors

Row vectors and Column vectors cannot be added together



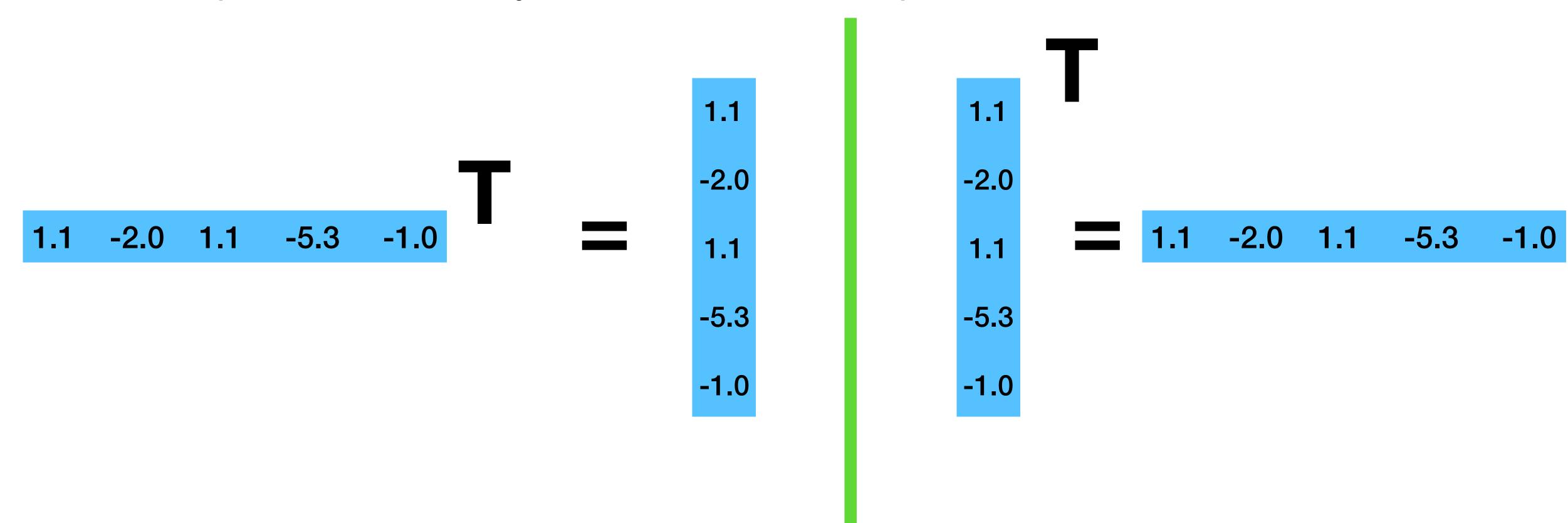
Transposition

 However, row vectors can be transformed in column vectors by an operation called transposition (and vice-versa)

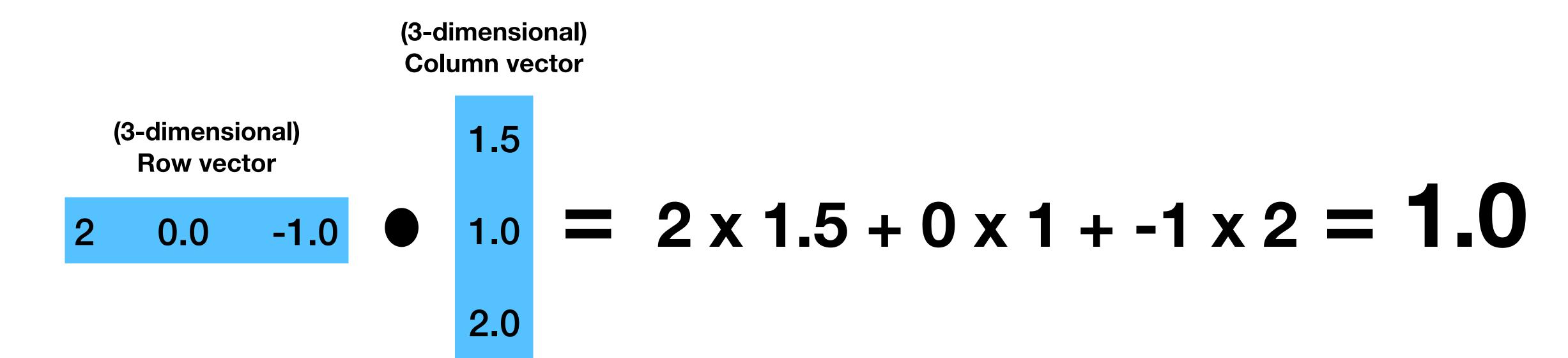


Transposition

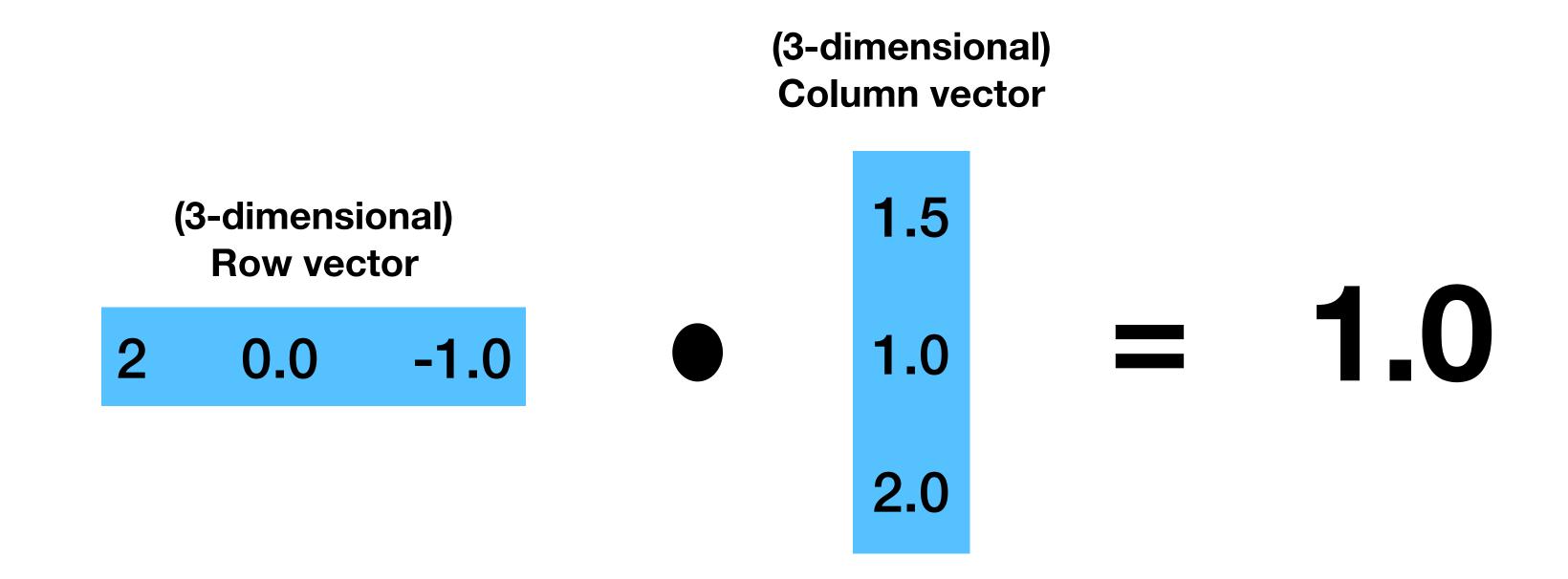
Transposition is usually noted with a *T* in exponent:

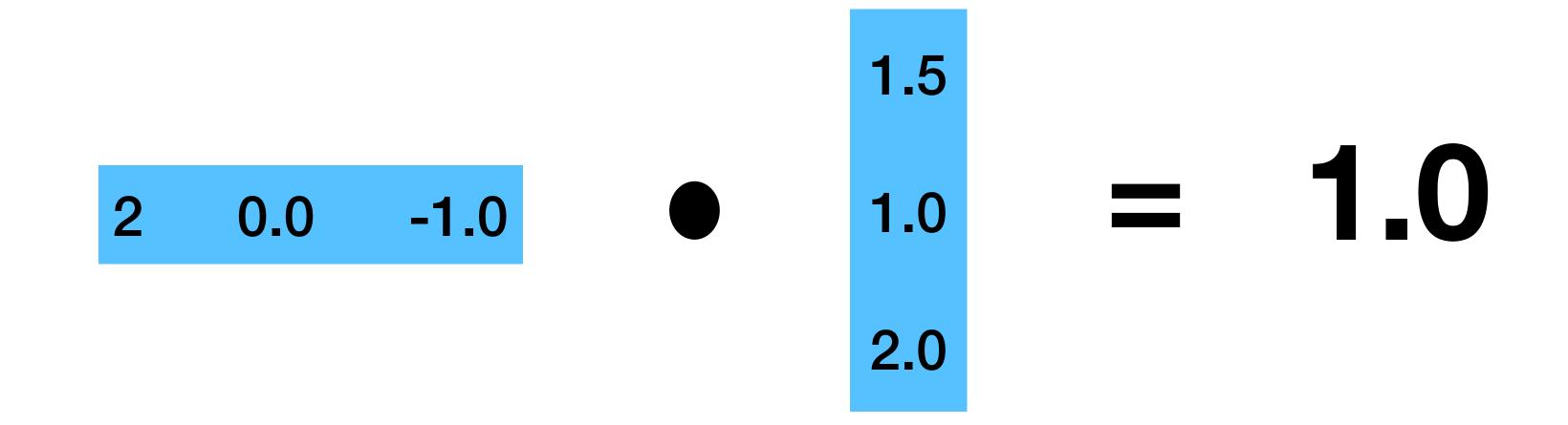


- We can compute the *inner product* of a *row vector* and a *column vector* to obtain a single number
- It consists in taking the product of the number dimension by dimension and then taking the sum

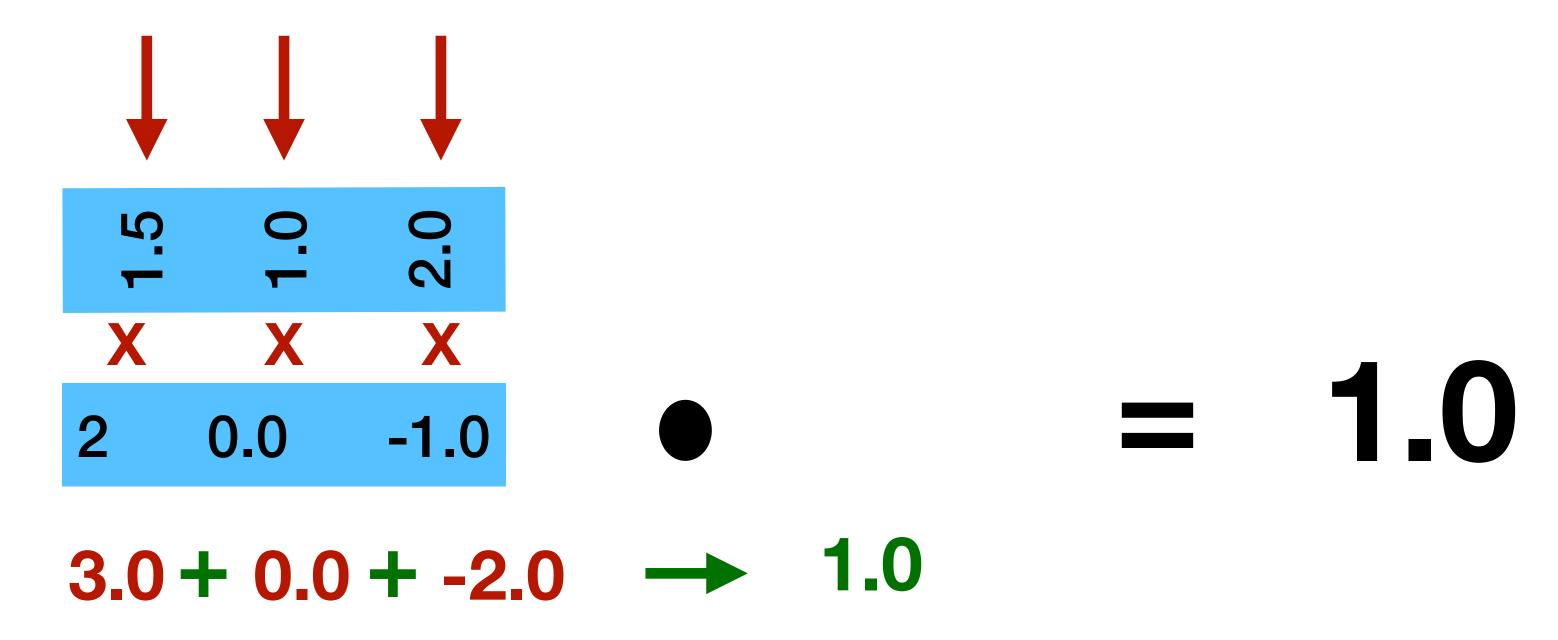


- The vectors should have the same dimension
- You should have the <u>row vector</u> on the <u>left</u>, and the <u>column vector</u> on the <u>right</u>
 - (if it is the opposite, the operation is called outer product and gives a different result)





You can try to visualize this as the column vector lying down on the row vector to produce the number



You can try to visualize this as the column vector lying down on the row vector to produce the number

Why inner product is important?

- It allows us to express linear functions efficiently
- And remember that linear functions are one of the fundamental components of Machine Learning

Linear regression

$$f(x, y) = \theta_0 + \theta_1 \times x + \theta_2 \times y$$

Logistic Classifier

$$score(x, y) = \theta_0 + \theta_1 \times x + \theta_2 \times y$$

 $prediction = \sigma(score)$

Neuron (same formula as Logistic Classifier)

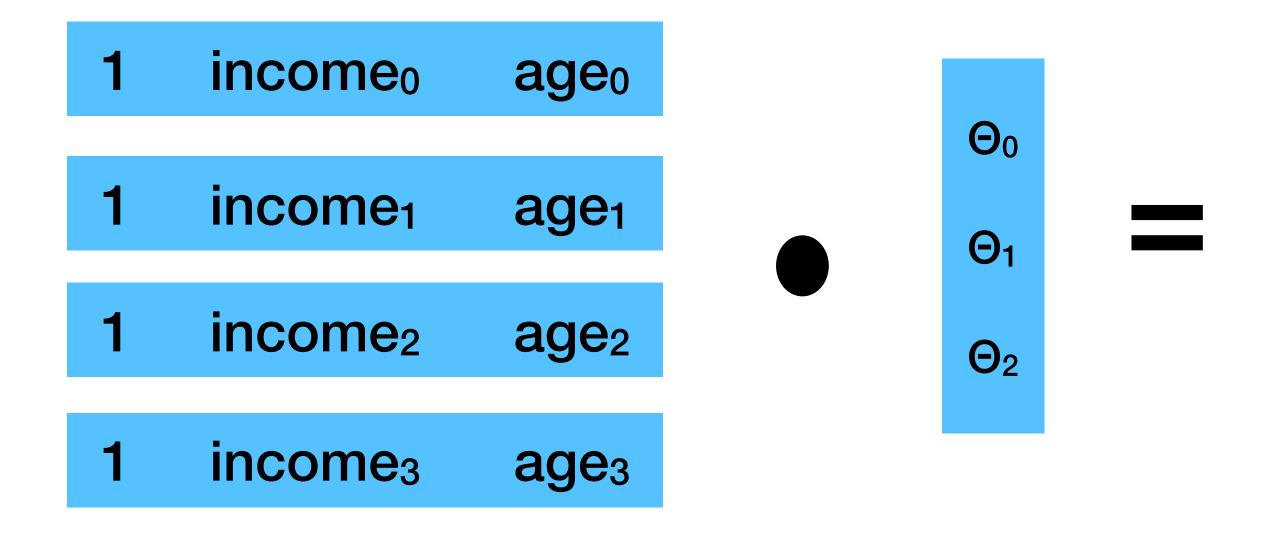
$$output(x, y) = \sigma(\theta_0 + \theta_1 \times x + \theta_2 \times y)$$

Why inner product is important?

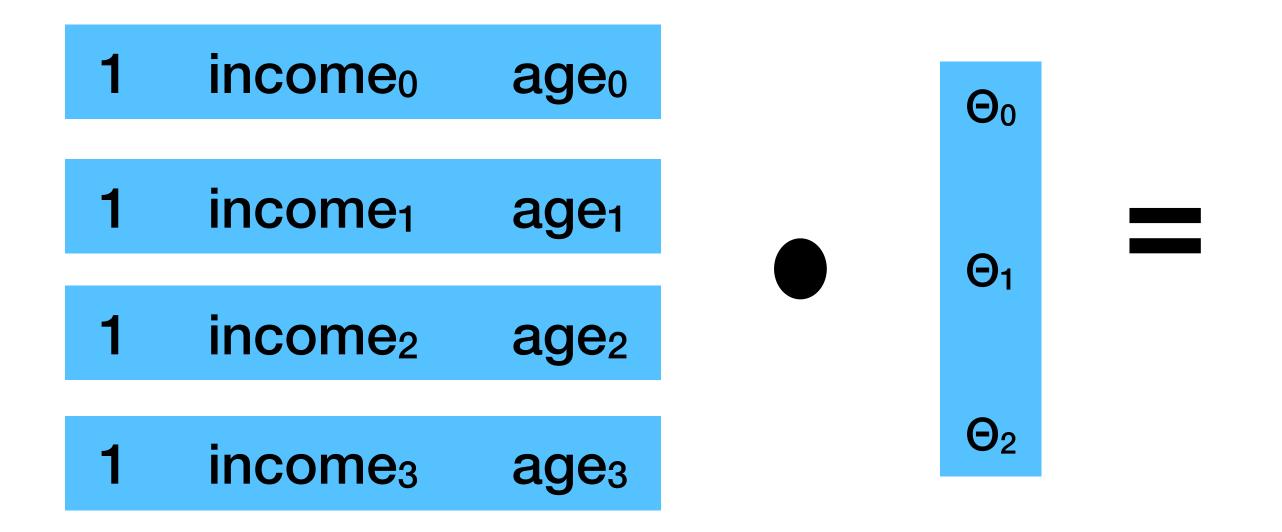
$$vote(age, income) = \sigma(\theta_0 + \theta_1 \times age + \theta_2 \times income)$$

1 income age Θ_0 Θ_1 Θ_2 Θ_0 Θ_0 Θ_0 Θ_0 Θ_0

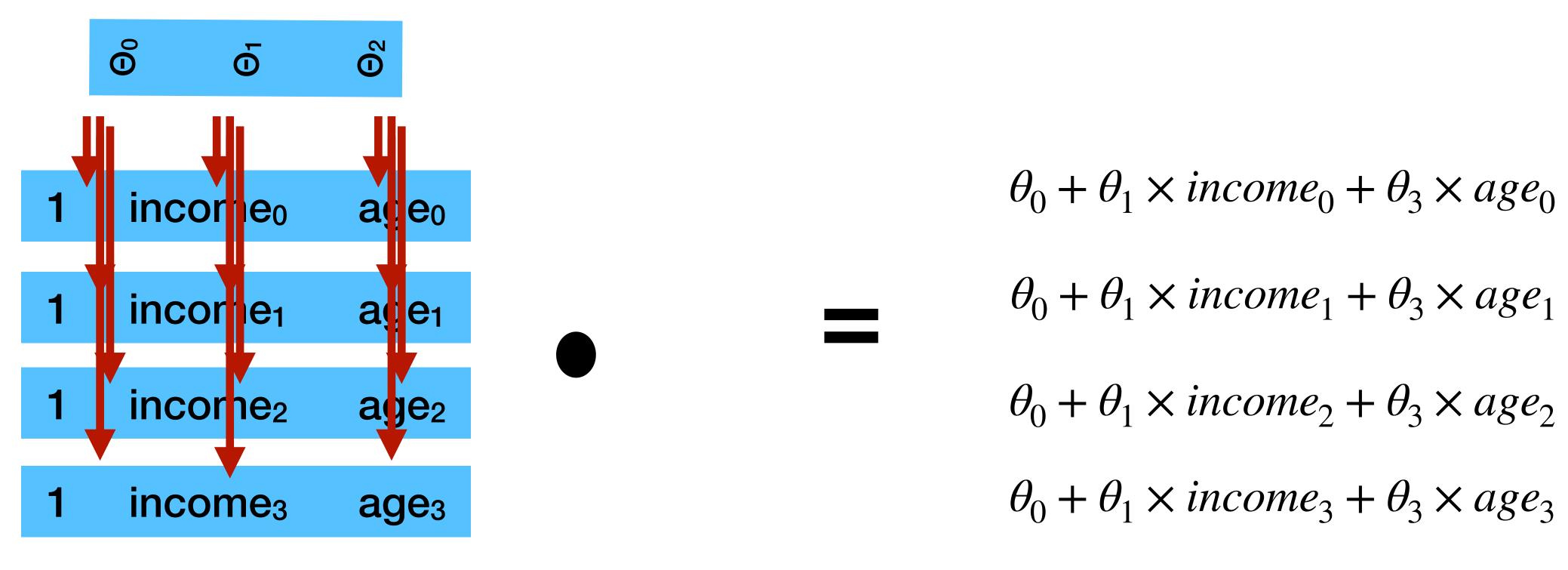
 $vote(age, income) = \sigma(\theta_0 + \theta_1 \times age + \theta_2 \times income)$



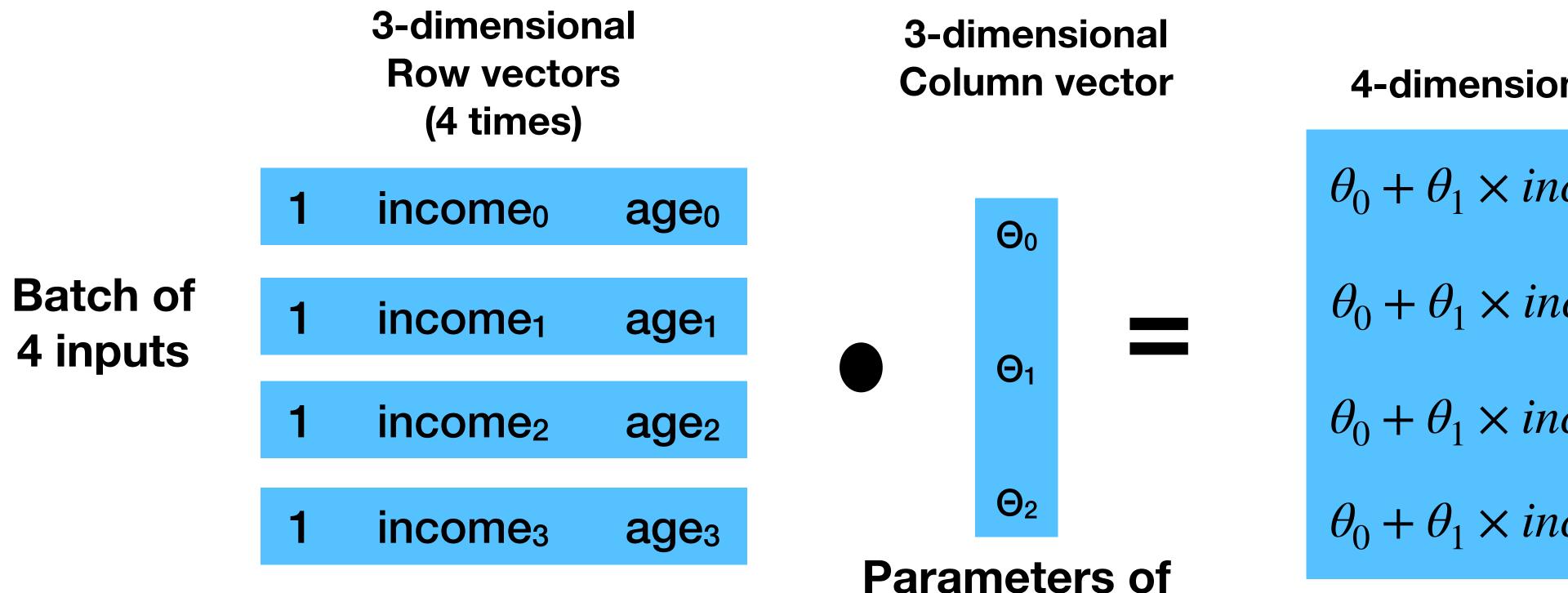
 $vote(age, income) = \sigma(\theta_0 + \theta_1 \times age + \theta_2 \times income)$



 $vote(age, income) = \sigma(\theta_0 + \theta_1 \times age + \theta_2 \times income)$



 $vote(age, income) = \sigma(\theta_0 + \theta_1 \times age + \theta_2 \times income)$



4-dimensional Column vector

$$\theta_{0} + \theta_{1} \times income_{0} + \theta_{3} \times age_{0}$$

$$\theta_{0} + \theta_{1} \times income_{1} + \theta_{3} \times age_{1}$$

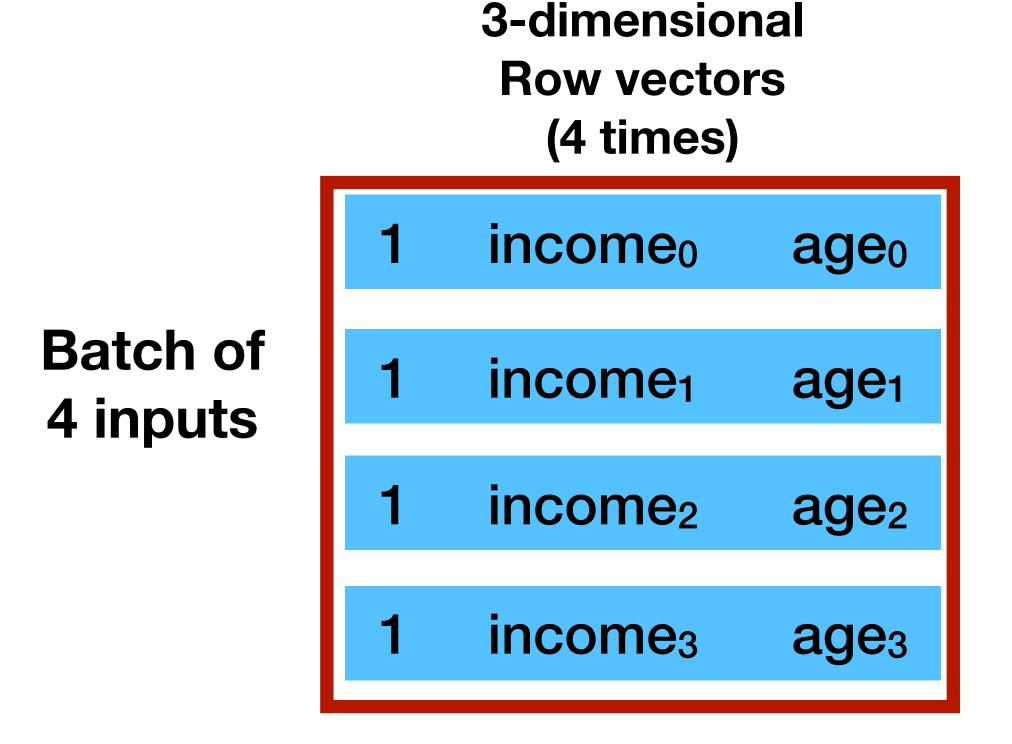
$$\theta_{0} + \theta_{1} \times income_{2} + \theta_{3} \times age_{2}$$

$$\theta_{0} + \theta_{1} \times income_{3} + \theta_{3} \times age_{3}$$

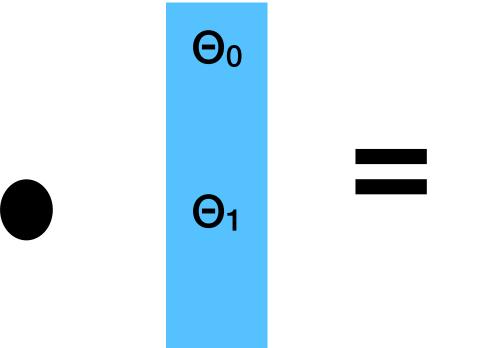
This represents the computation of the score output of a single neuron for 4 inputs at the same time!

one neuron

 $vote(age, income) = \sigma(\theta_0 + \theta_1 \times age + \theta_2 \times income)$



3-dimensional Column vector



 Θ_2

4-dimensional Column vector

$$\theta_0 + \theta_1 \times income_0 + \theta_3 \times age_0$$

 $\theta_0 + \theta_1 \times income_1 + \theta_3 \times age_1$
 $\theta_0 + \theta_1 \times income_2 + \theta_3 \times age_2$
 $\theta_0 + \theta_1 \times income_3 + \theta_3 \times age_3$

This is called a 4x3 matrix

 $vote(age, income) = \sigma(\theta_0 + \theta_1 \times age + \theta_2 \times income)$

	3-dimension Row vecto (4 times)	rs	3-dimensional Column vector				
1	incomeo	age₀		Θ_0			
1	income ₁	age ₁		Θ1			
1	income ₂	age ₂					
1	income ₃	age ₃		Θ_2			

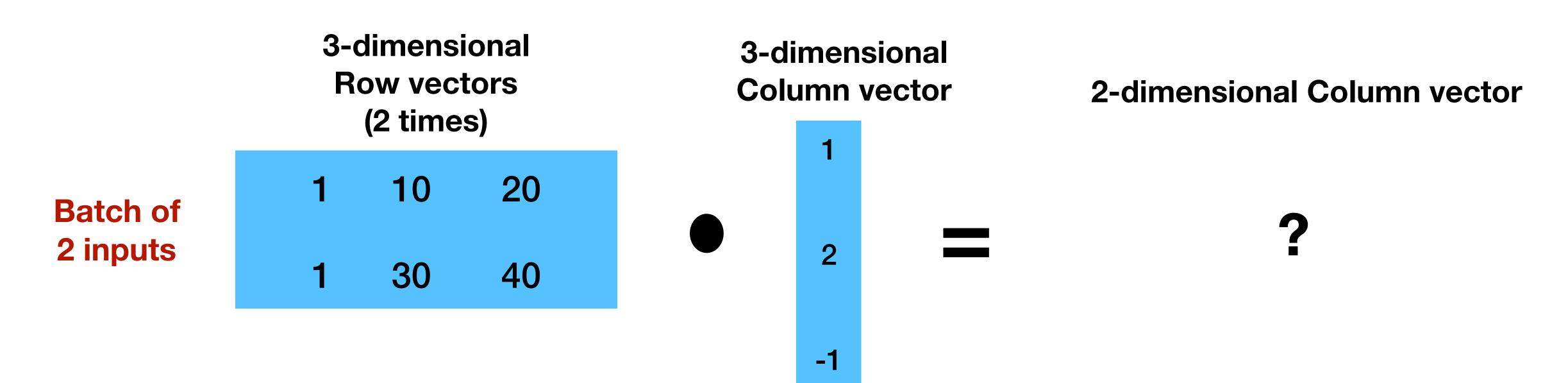
4-dimensional Column vector

$$\theta_0 + \theta_1 \times income_0 + \theta_3 \times age_0$$

 $\theta_0 + \theta_1 \times income_1 + \theta_3 \times age_1$
 $\theta_0 + \theta_1 \times income_2 + \theta_3 \times age_2$
 $\theta_0 + \theta_1 \times income_3 + \theta_3 \times age_3$

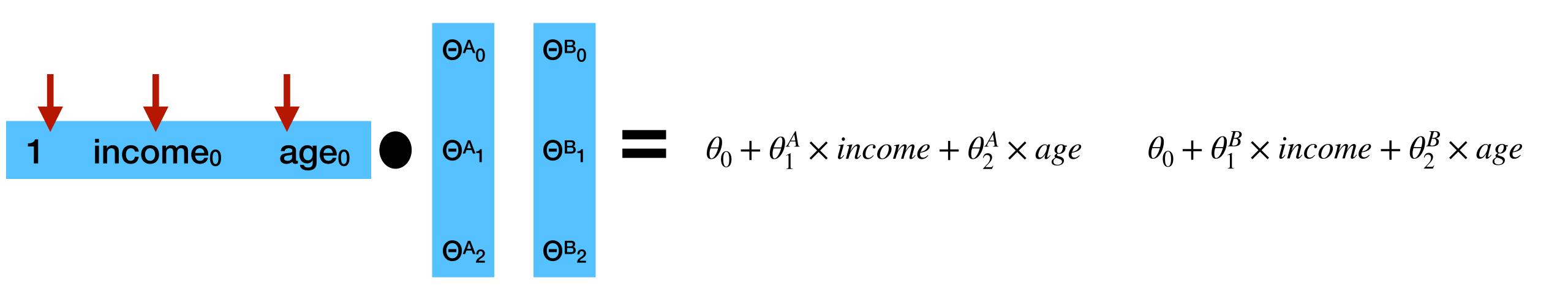
This is called a 4x3 matrix

 $vote(age, income) = \sigma(\theta_0 + \theta_1 \times age + \theta_2 \times income)$

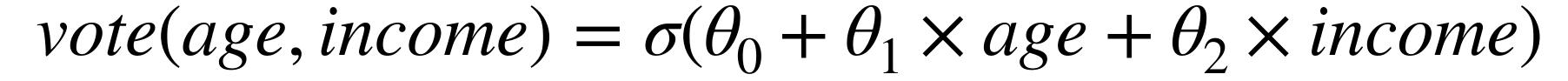


Parameters of one Neuron

$$vote(age, income) = \sigma(\theta_0 + \theta_1 \times income + \theta_2 \times age)$$



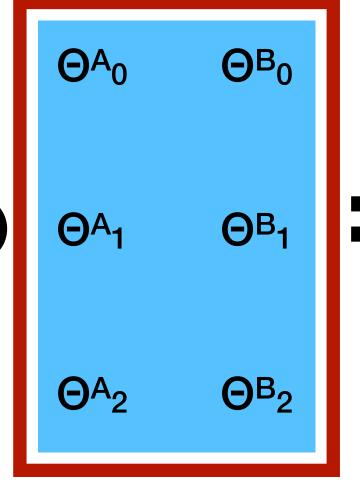
Parameters for two different neurons A and B



3-dimensional Column vector (2 times)

3-dimensional Row vectors

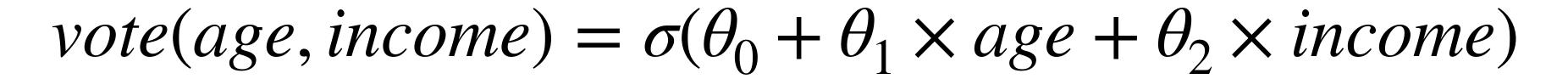
1 income₀ age₀

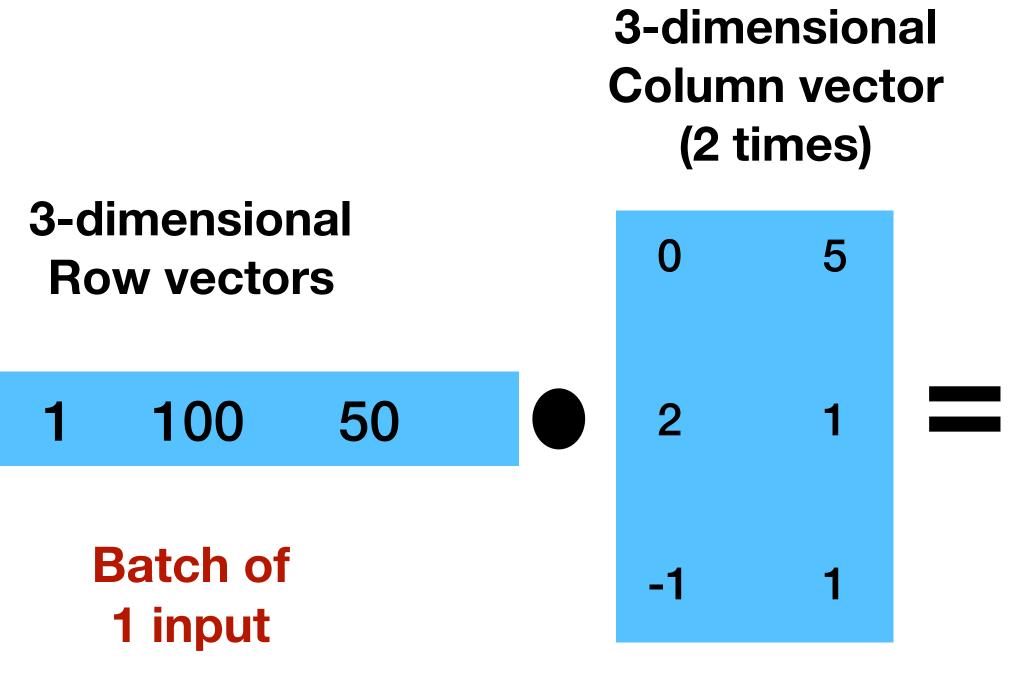


2-dimensional Row vector

$$\theta_0 + \theta_1^A \times income + \theta_2^A \times age$$
 $\theta_0 + \theta_1^B \times income + \theta_2^B \times age$

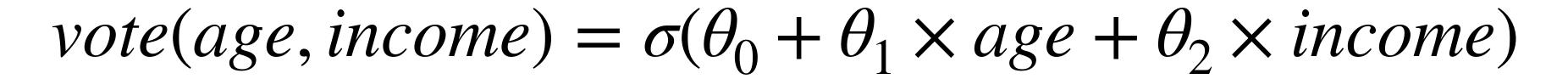
This is called a 3x2 matrix

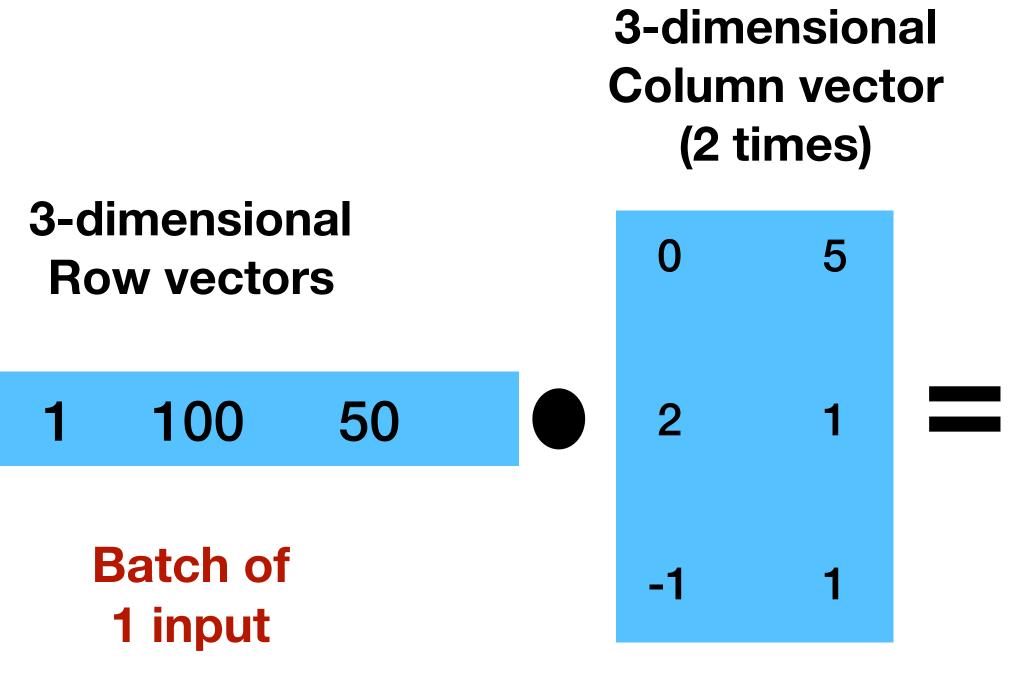




Parameters of two Neurons

2-dimensional Row vector

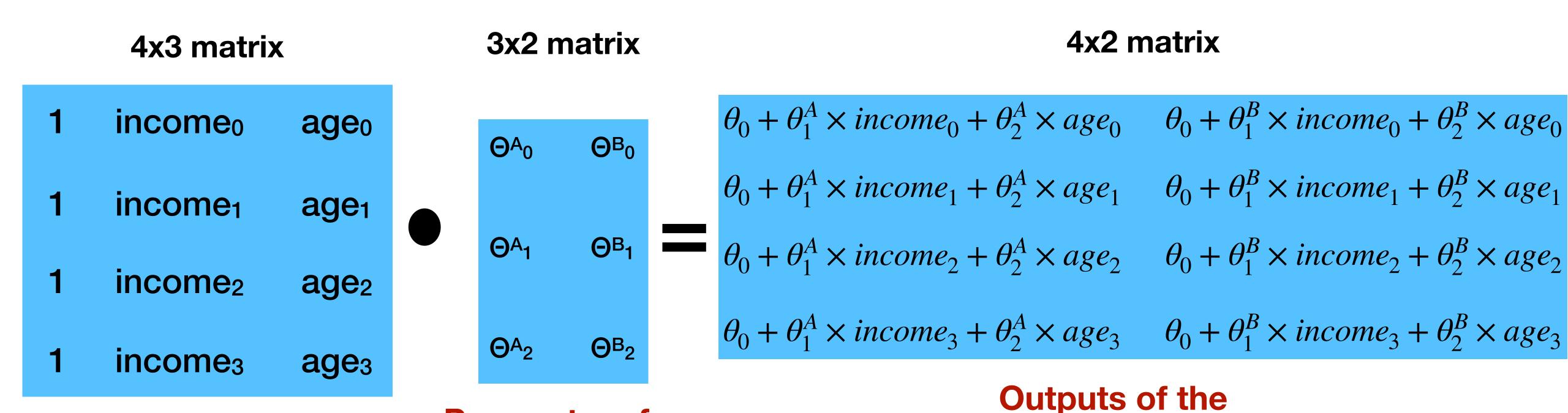




Parameters of two Neurons

2-dimensional Row vector

 If we combine the 2 aspects of having several inputs and several neurons, we have what is called a matrix multiplication



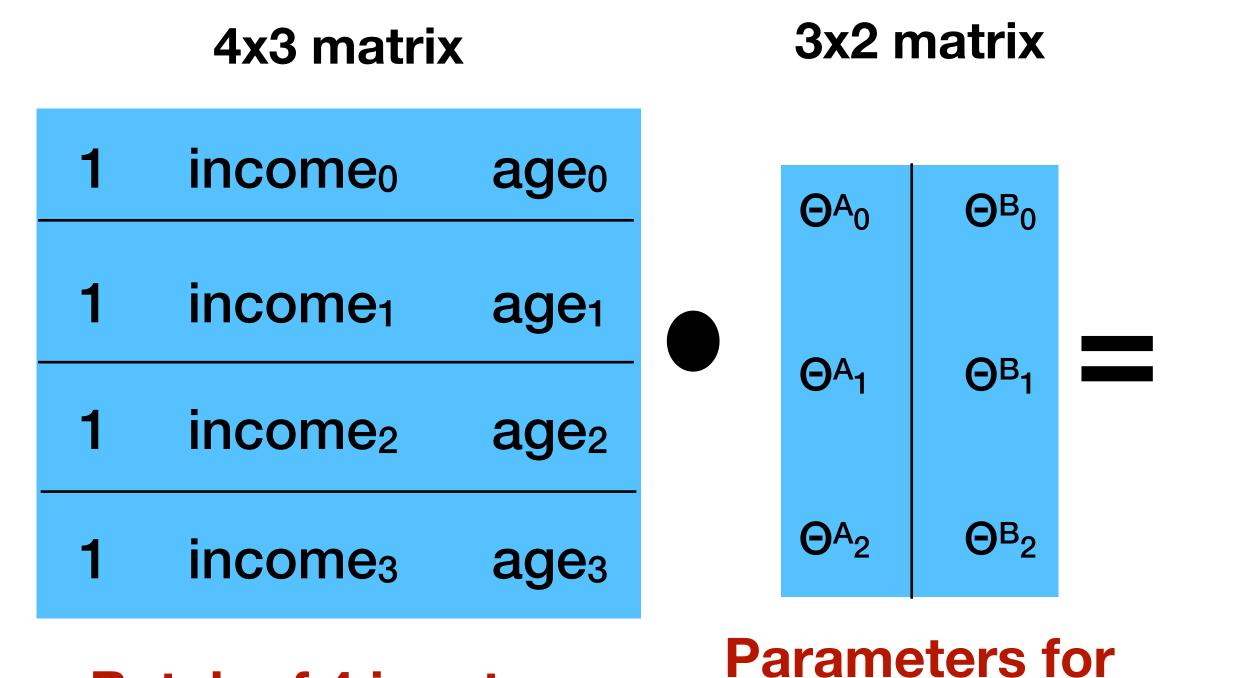
2 neurons for each of the 4 inputs

Parameters for

2 neurons

Batch of 4 inputs

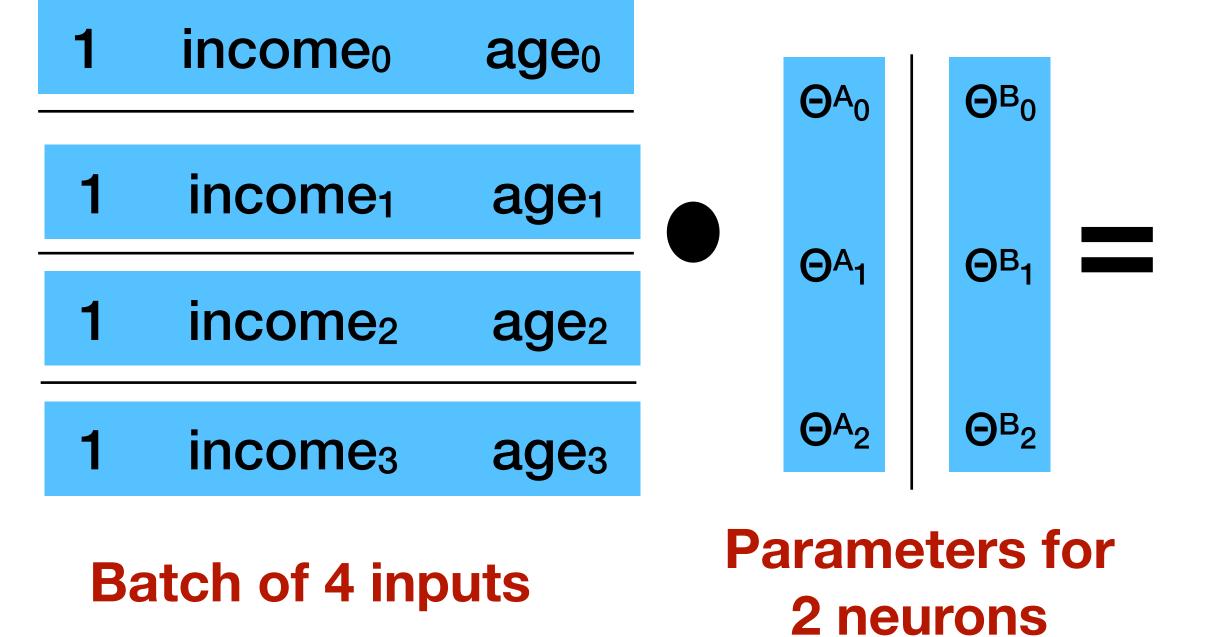
4x2 matrix

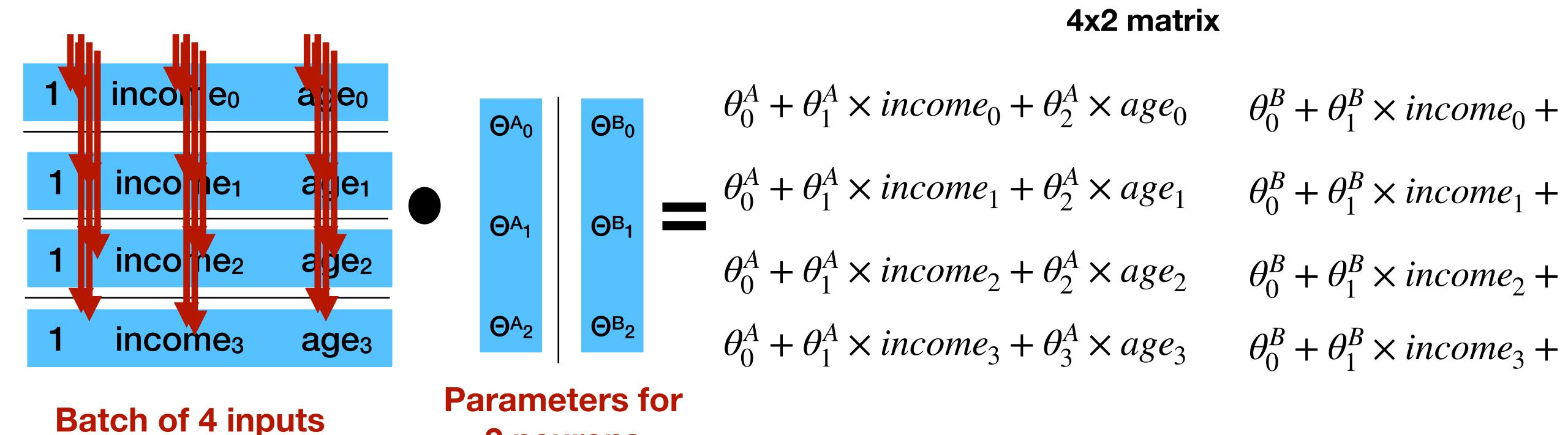


2 neurons

Batch of 4 inputs

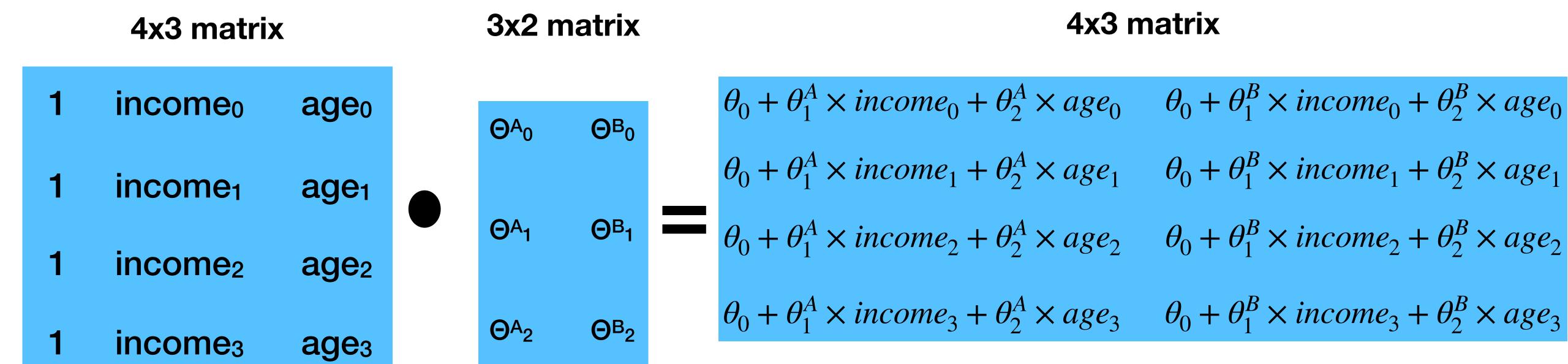
4x2 matrix





2 neurons

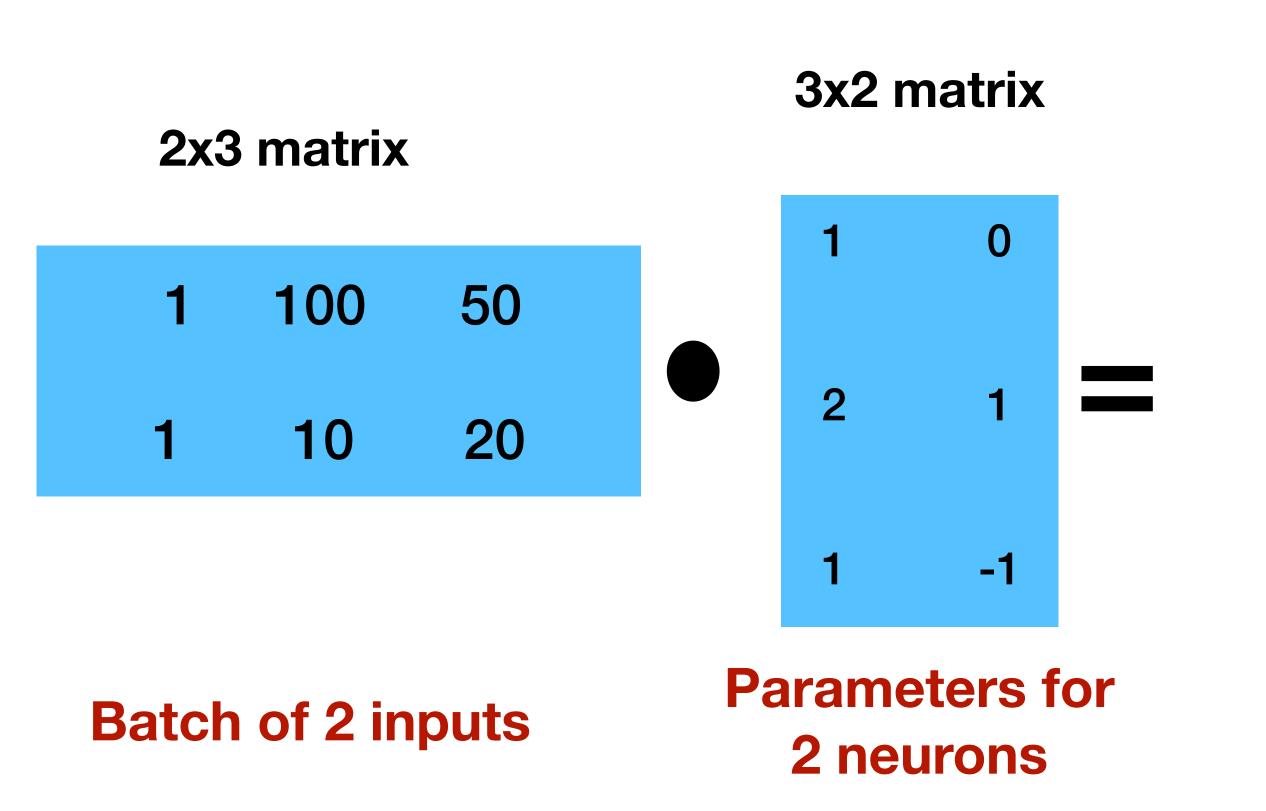
- We represent many inputs as a matrix
- We represent many neurons as a matrix
- Matrix multiplication compute the output score of all the neurons for all the inputs



Batch of 4 inputs

Parameters for 2 neurons

Outputs of the 2 neurons for each of the 4 inputs



2x2 matrix

Outputs of the 2 neurons for each of the 2 inputs

Computing the output of many neurons for many inputs

• In practice, we separate the weights from the bias

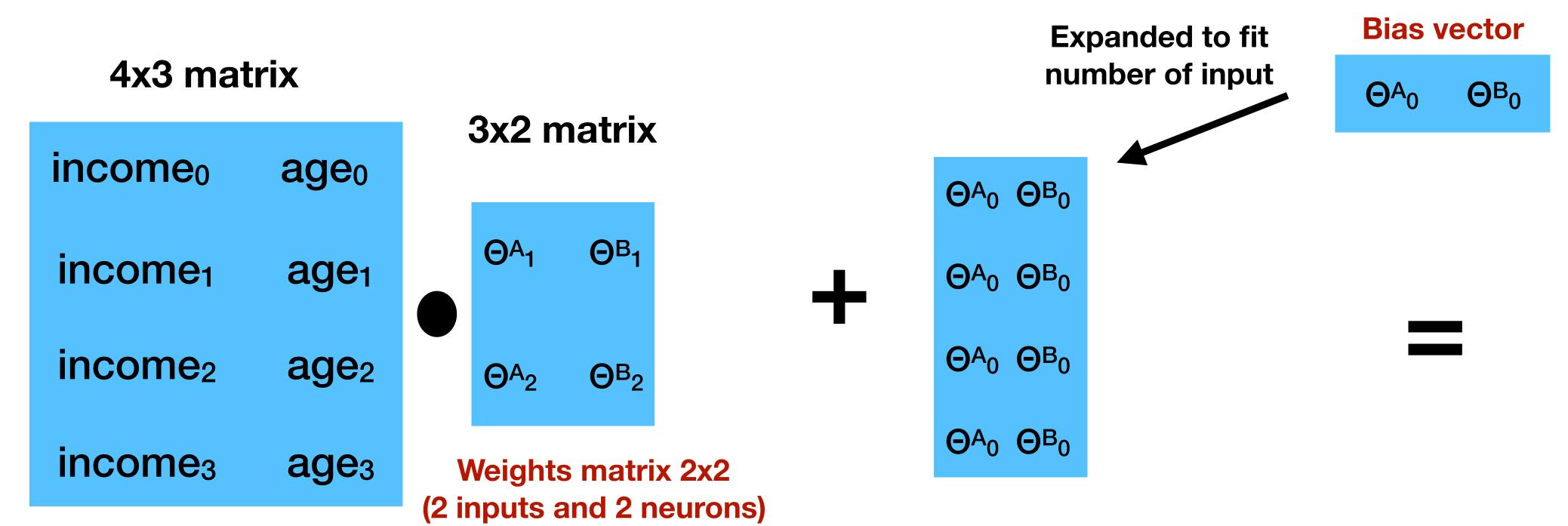
4x3 matrix 4x3 matrix Bias of 2 neurons $\theta_0 + \theta_1^A \times income_0 + \theta_2^A \times age_0$ $\theta_0 + \theta_1^B \times income_0 + \theta_2^B \times age_0$ income₀ age₀ Θ^{A_0} Θ^{B_0} $\theta_0 + \theta_1^A \times income_1 + \theta_2^A \times age_1$ $\theta_0 + \theta_1^B \times income_1 + \theta_2^B \times age_1$ income₁ age₁ $\theta_0 + \theta_1^A \times income_2 + \theta_2^A \times age_2$ $\theta_0 + \theta_1^B \times income_2 + \theta_2^B \times age_2$ Θ^{A_1} income₂ age₂ $\theta_0 + \theta_1^A \times income_3 + \theta_2^A \times age_3$ $\theta_0 + \theta_1^B \times income_3 + \theta_2^B \times age_3$

Batch of 4 inputs

Weights of 2 neurons

Outputs of the 2 neurons for each of the 4 inputs

Computing the output of many neurons for many inputs

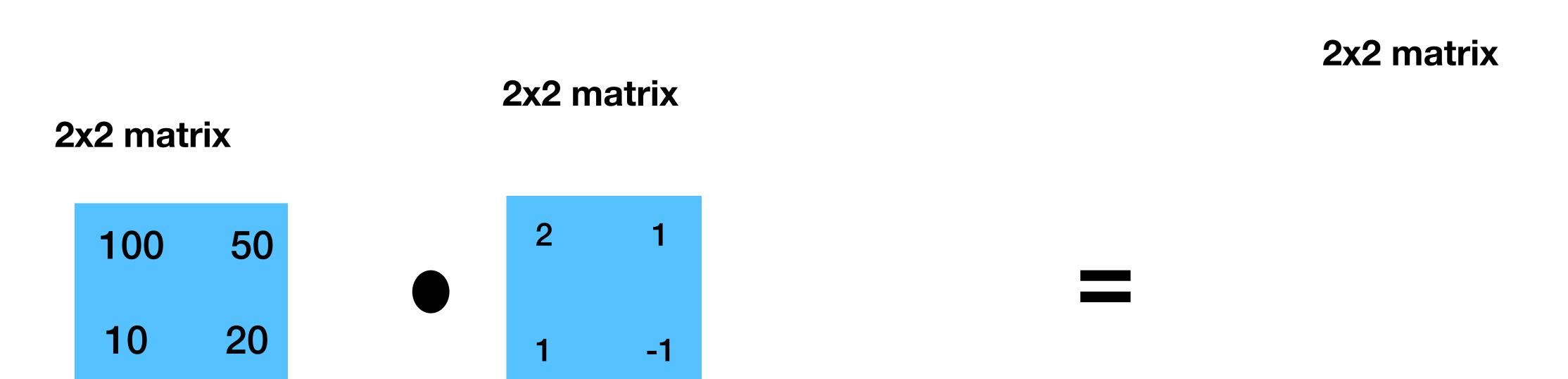


Batch input (4 inputs)

4x3 matrix

$$\begin{aligned} \theta_0 + \theta_1^A \times income_0 + \theta_2^A \times age_0 & \theta_0 + \theta_1^B \times income_0 + \theta_2^B \times age_0 \\ \theta_0 + \theta_1^A \times income_1 + \theta_2^A \times age_1 & \theta_0 + \theta_1^B \times income_1 + \theta_2^B \times age_1 \\ \theta_0 + \theta_1^A \times income_2 + \theta_2^A \times age_2 & \theta_0 + \theta_1^B \times income_2 + \theta_2^B \times age_2 \\ \theta_0 + \theta_1^A \times income_3 + \theta_2^A \times age_3 & \theta_0 + \theta_1^B \times income_3 + \theta_2^B \times age_3 \end{aligned}$$

Bias for 2 neurons 1 0



Batch of 2 inputs

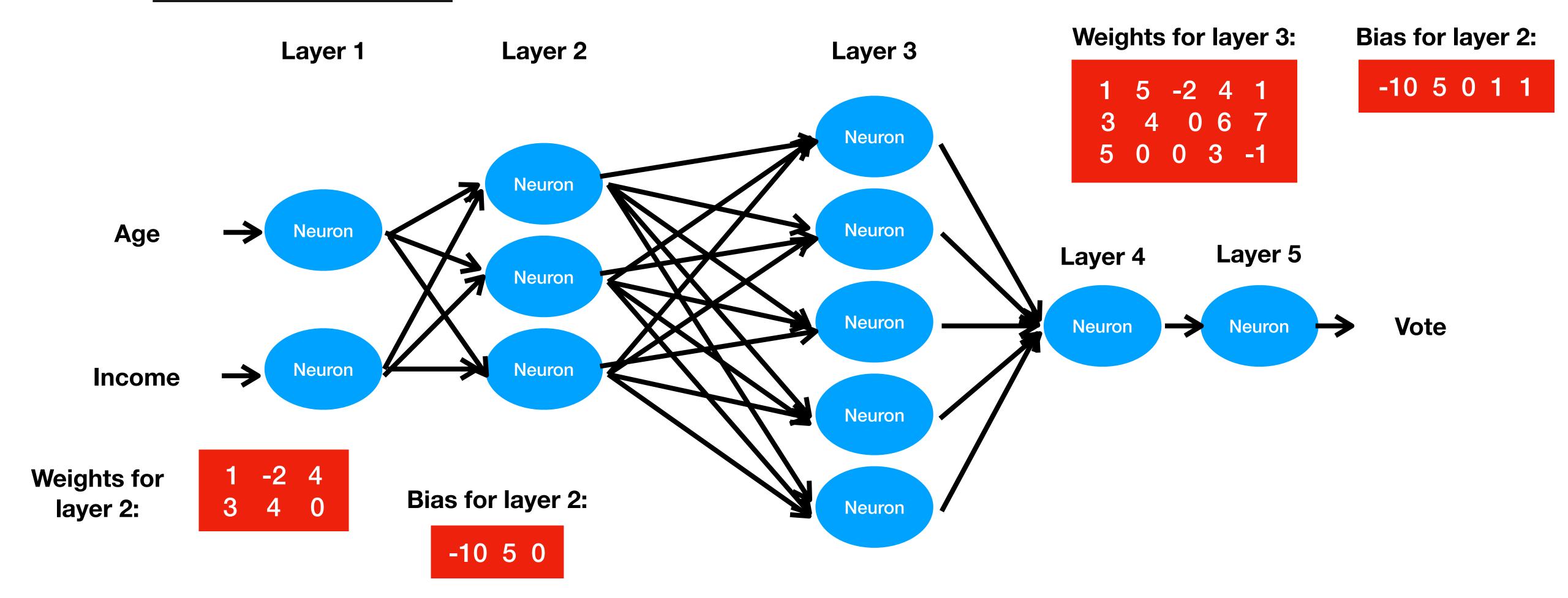
Weights for 2 neurons

Score outputs of the 2 neurons for each of the 2 inputs

- Important: using these matrix representations only work if all the neurons have the same input and are fully connected
- Then, in practice, a **Fully Connected layer of N neurons with K input** can be represented by a **matrix of weights W** of shape KxN, and a **bias vector b** of dimension N.

Feed-Forward networks with fully connected layers

 Then, in practice, a <u>Fully Connected layer of N neurons with K input</u> can be represented by <u>a matrix of weights W of shape KxN</u>, and a <u>bias vector</u>
 b of dimension N.



Yet another visualization

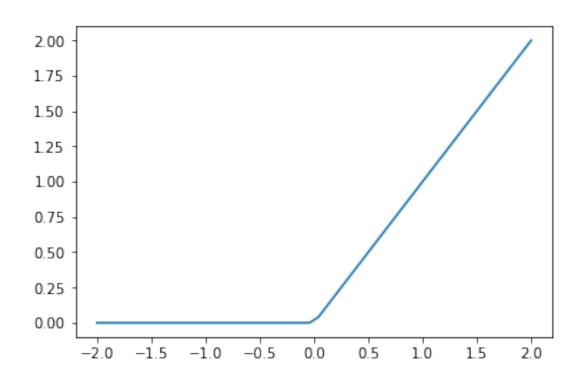
- An interactive Neural Network with Fully Connected Layers
 - From some nice people at Google
 - https://playground.tensorflow.org

Activations

Some possible activation functions:

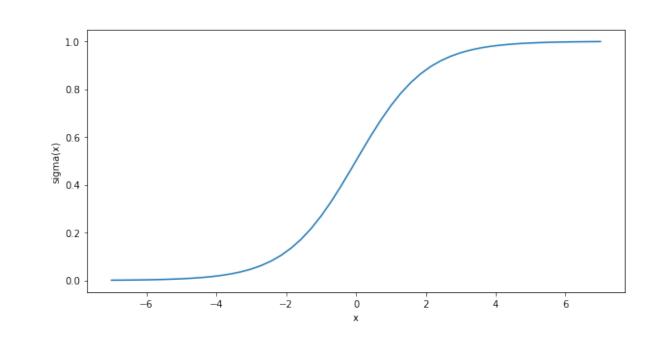
Rectified Linear Unit

$$ReLU(x) = max(x,0)$$



Sigmoid (a.k.a logistic function)

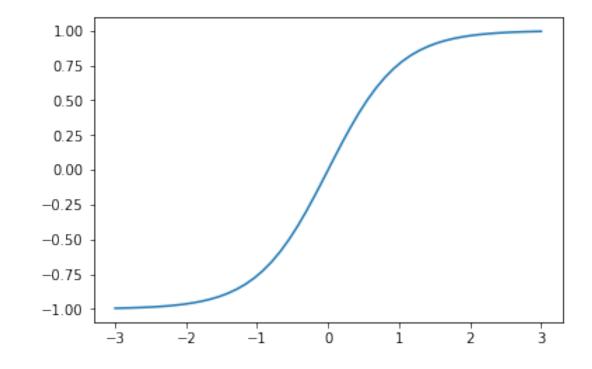
$$\sigma(x) = \frac{1}{1 + exp(-x)}$$

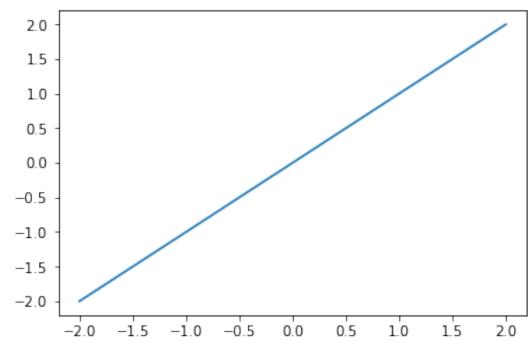


Hyperbolic Tangent Linear function (= a linear transformation (useless as activation function) of the sigmoid)

$$tanh(x) = 2\sigma(2x) - 1$$

$$f(x) = x$$





L1 and L2 Regularization

- We have discusses that L2 Regularization can be used for reducing overfitting
 - -> Change Regularization and Regularization Rate to observe the effect of Regularization

L2 Regularization

- We note $|\Theta|^2$ the sum of the square of all parameters Θ_k (this is called the "L2 Norm")
- Then we add this quantity to the loss we want to minimize:

$$Loss = MeanSquaredError = \frac{1}{N} \cdot \sum_{i} (model(cig_{i}, bmi_{i}, ismale_{i}) - age_{i})^{2}$$

$$Loss = \frac{1}{N} \cdot \sum_{i} (model(cig_{i}, bmi_{i}, ismale_{i}) - age_{i})^{2} + \lambda |\overrightarrow{\theta}|^{2}$$

Then we apply Gradient Descent to this new loss

L1 Regularization

- We note $|\Theta|^1$ the sum of the absolute value of all parameters Θ_k (this is called the "L1 Norm")
- Then we add this quantity to the loss we want to minimize:

$$Loss = MeanSquaredError = \frac{1}{N} \cdot \sum_{i} (model(cig_{i}, bmi_{i}, ismale_{i}) - age_{i})^{2}$$

$$Loss = \frac{1}{N} \cdot \sum_{i} (model(cig_{i}, bmi_{i}, ismale_{i}) - age_{i})^{2} + \lambda \mid \overrightarrow{\theta} \mid$$

Then we apply Gradient Descent to this new loss

Random restarts

- Because the weights of a neural network are initialized randomly, the result of a training will change every time we reinitialize the network
- Therefore, to obtain best performances, it is common to train a networks several times with different random initializations, and keep the best result
 - -> Press the reload button to restart a training with new random weight initialization

Regularization methods for Neural Network

- When we train a network with many neurons, the danger of overfitting is large
- There are a few technics that are very efficient at preventing this:
 - Early Stopping
 - Dropout
 - Weight Decay
 - Stochastic Gradient Descent

Early stopping

- We keep a validation set separate from the training data
- We fix a *patience* number (typically patience = 10 or 20)
- During training, if we see no improvement on validation set after *patience* measures, we stop the training
 - Check the evolution of test loss to detect when training should be stopped

Dropout

- During training, we add random noise to disturb the network
- In practice, we randomly "cut" a certain proportion of connections (typically 10% to 50%)
 - Check the effect of adding noise

Stochastic Gradient Descent

- Instead of computing the gradient of the loss for all the training data, we compute it for a subsampled part of the training data
- This is actually done anyway to get faster training
- But it is also beneficial to prevent overfitting even if you could afford to compute the gradient for all the examples at once.
 - -> change batch size to observe what happen when we compute gradient on more or less examples

Input features

- We sae that even for a simple model (like linear regression), adding new functions of the input can increase the capacity of the model
 - -> Add input features and see the effect when there are few neurons