



CONSTELLATE

# How does ChatGPT work?





CONSTELLATE

**ChatGPT**

**GPT: generative pre-trained transformer**



CONSTELLATE

# ChatGPT

GPT: generative pre-trained **transformer**



CONSTELLATE

ChatGPT

**Transformer: a multi-layer neural network that relies on the parallel multi-head attention mechanism.**



CONSTELLATE

ChatGPT

## Transformer

a multi-layer neural network that relies on the parallel multi-head attention mechanism.



CONSTELLATE

**ChatGPT**

**Part 1: multi-layer neural network**

**Part 2: multi-head attention mechanism**



CONSTELLATE

**ChatGPT**

## Part 1

# Multi-layer neural network



CONSTELLATE

# Part 1

## Learning objectives

- **Concepts**
  - Understand some basic concepts in neural networks
    - feature, weight, bias, vector, matrix.....
    - neuron, activation function, hidden layer.....
  - Understand a neuron as a computation unit
  - Understand the fundamental algorithms underlying NNs
- **Hands-on computation**
  - Know how to do matrix multiplication by hand





CONSTELLATE

# Roadmap

- One neuron, one layer
- Many neurons, one layer
- Many neurons, many layers



CONSTELLATE

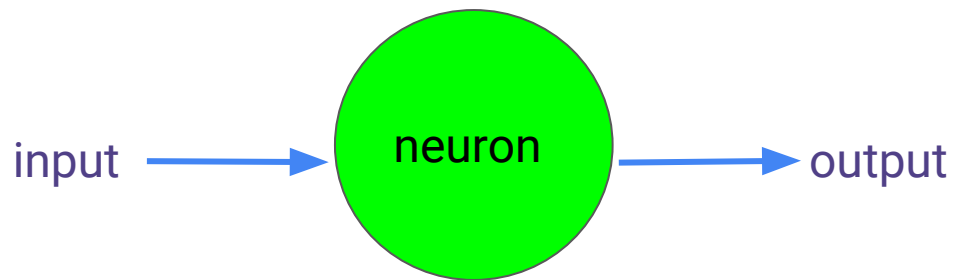
# Roadmap

- What will be covered
  - The machine learning algorithms
- What will not be covered
  - The optimization algorithms



CONSTELLATE

# Roadmap



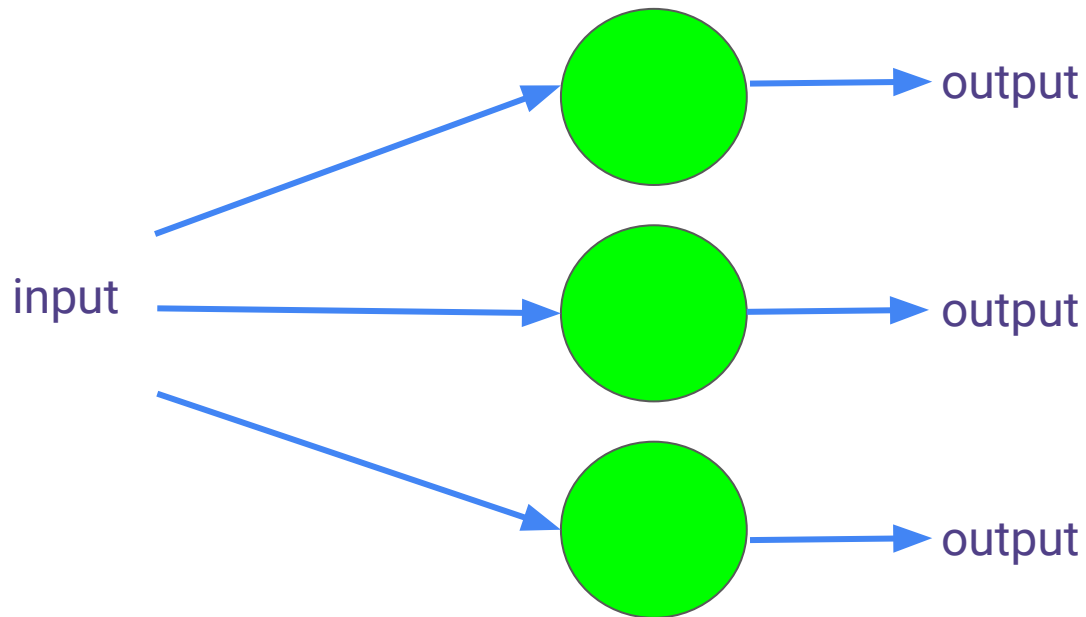
**One neuron**

**One layer**



CONSTELLATE

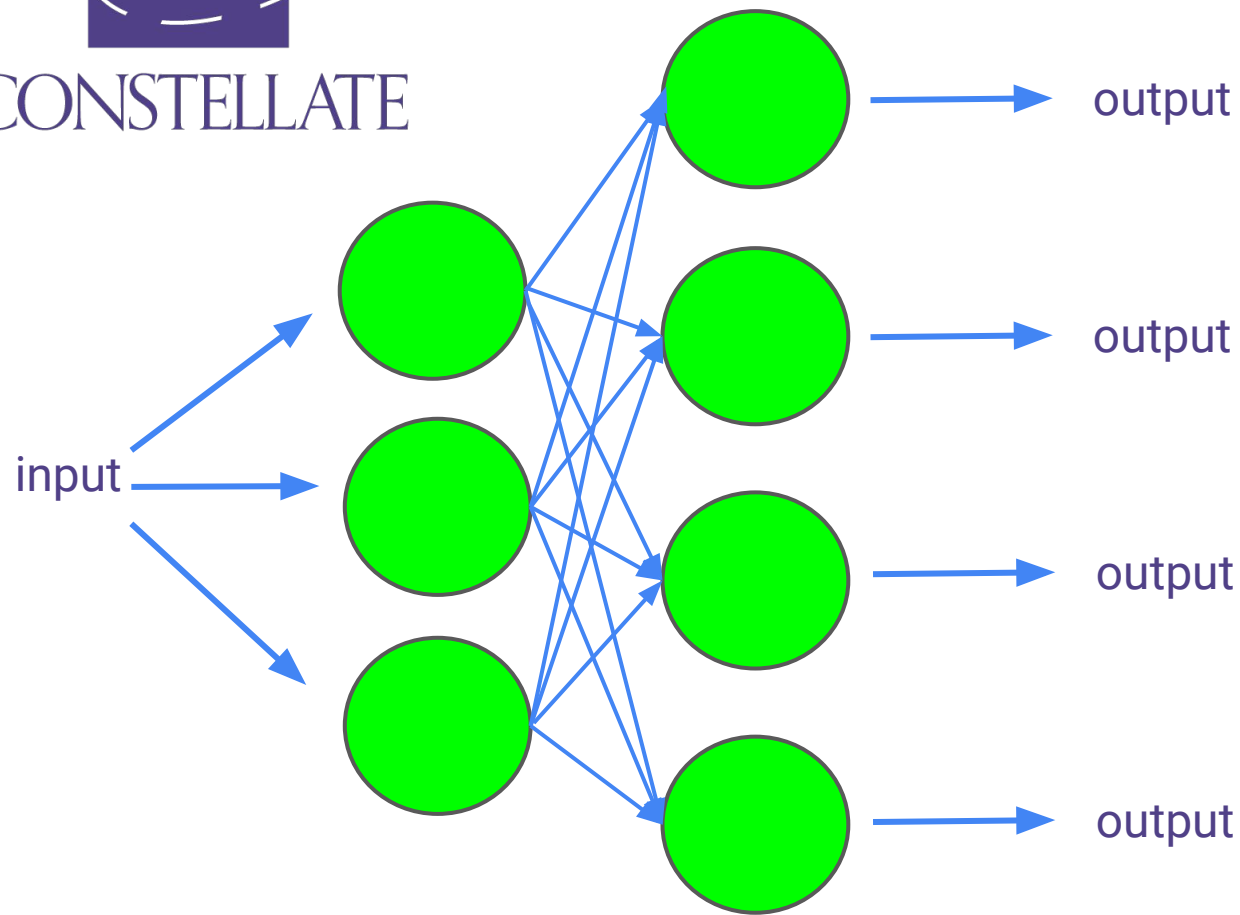
# Roadmap



**many neurons**  
**one layer**



CONSTELLATE



# Roadmap

**many neurons**  
**many layers**



CONSTELLATE

**One neuron, one layer**



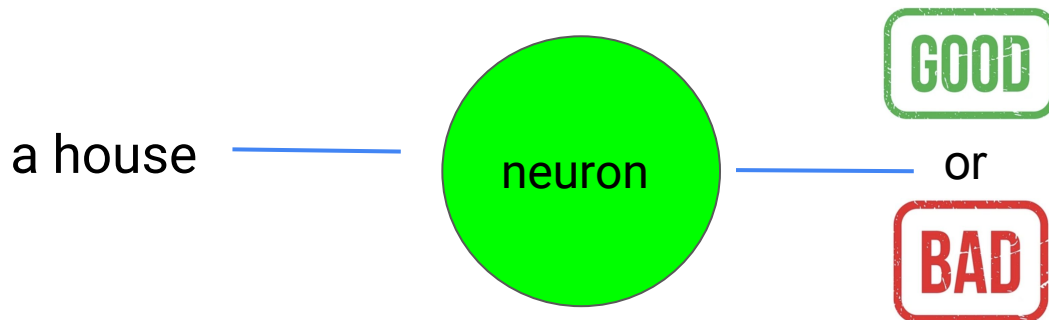
CONSTELLATE



Suppose you are looking to buy a home and would like to decide which houses are good and which houses are bad.



CONSTELLATE



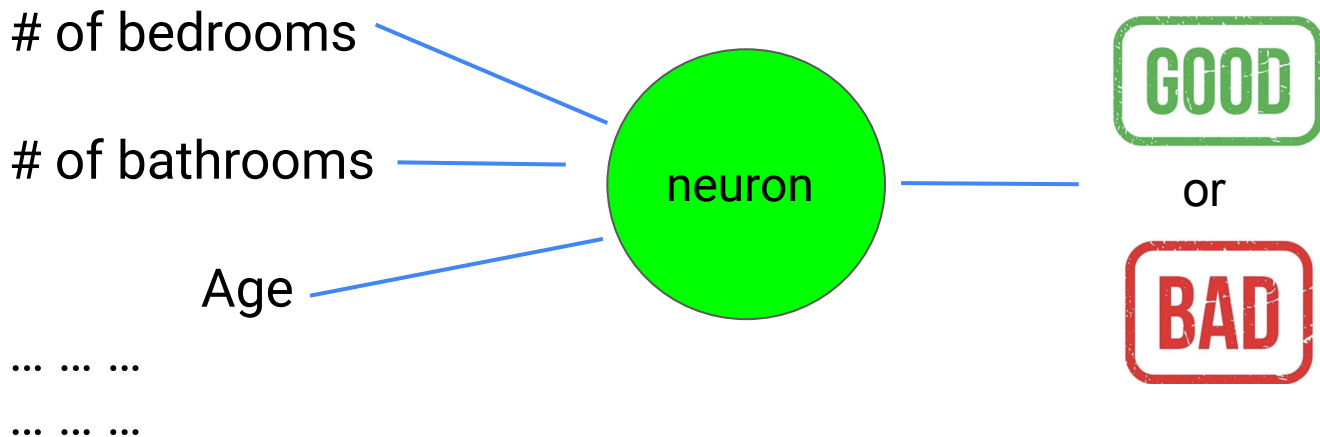
- You want a model that automatically takes in a house and outputs whether it is good or bad.
- You get some training data, i.e. houses, and hand label them as good or bad based on your domain knowledge of house quality and feed the training data to the model to let it learn the rules.
- After you find the best-performing model, you can apply it to new data, i.e. houses it has never seen before and let it decide whether the input house is good or bad!





CONSTELLATE

# Features



Features



CONSTELLATE

# Not all features are alike

# of bedrooms

# of bathrooms

Age

... ..

... ..

The features differ in how much they influence the house being good or bad

**Features**

# of bedrooms \*  $W_{\text{bed}}$

# of bathrooms \*  $W_{\text{bath}}$

Age \*  $W_{\text{age}}$

... ..

...

... ..

...

**Features**

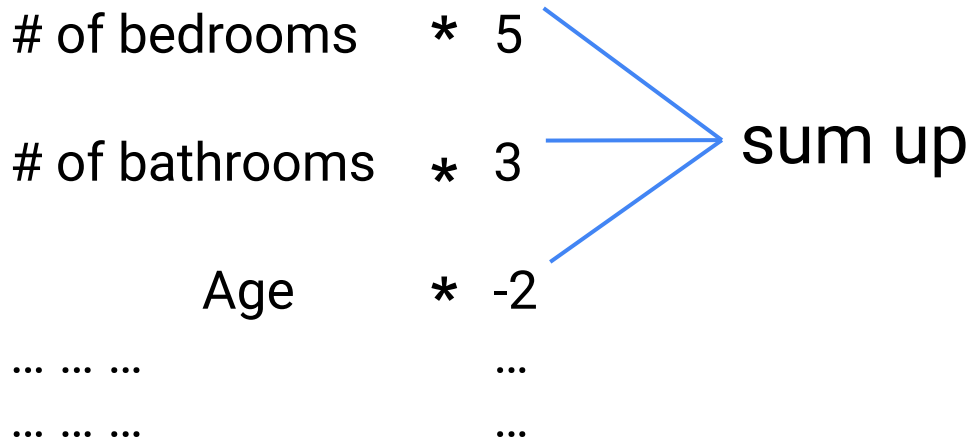
**weights**



CONSTELLATE

# Weight

For example



**Features**

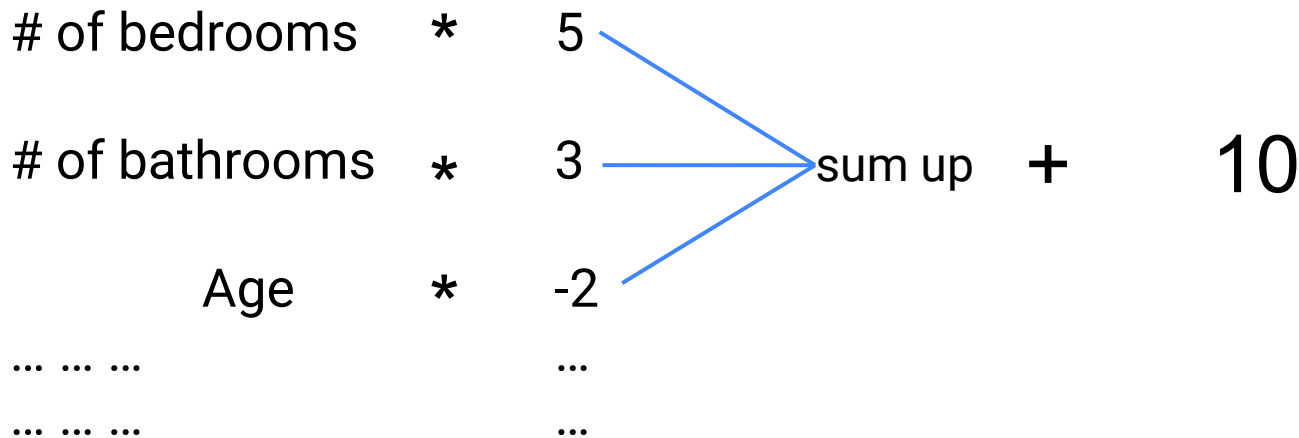
**weights**



CONSTELLATE

# Bias

For example



**Features**



**weights**

**+ Bias term**



CONSTELLATE

# Formalism

$$y = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$



# of  
bedrooms



# of  
bathrooms



age

$$\mathbf{x} = [x_1, x_2, \dots, x_n]$$

**Feature vector**

$$\mathbf{w} = [w_1, w_2, \dots, w_n]$$

**Weight vector**



CONSTELLATE

# Formalism

$$y = \boxed{w_1x_1 + w_2x_2 + \dots + w_nx_n} + b$$

↓ ↓ ↓  
# of # of age  
bedrooms bathrooms

$$\mathbf{x} = [x_1, x_2, \dots, x_n]$$

$$\mathbf{w} = [w_1, w_2, \dots, w_n]$$

$$\mathbf{x} \bullet \mathbf{w} = \boxed{w_1x_1 + w_2x_2 + \dots + w_nx_n}$$

$$y = \mathbf{x} \bullet \mathbf{w} + b$$



CONSTELLATE

# Vector multiplication

$$\mathbf{x} = [x_1, x_2, \dots, x_n]$$

$$\mathbf{w} = [w_1, w_2, \dots, w_n]$$

$$\mathbf{x} \bullet \mathbf{w} = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

$$[1 \quad 2 \quad 3] \cdot [3 \quad 2 \quad 1] = 1 \times 3 + 2 \times 2 + 3 \times 1 = 10$$





CONSTELLATE

# Formalism

$$y = \boxed{w_1x_1 + w_2x_2 + \dots + w_nx_n} + b$$

↓ ↓ ↓  
# of # of age  
bedrooms bathrooms

$$\mathbf{x} = [x_1, x_2, \dots, x_n]$$

$$\mathbf{w} = [w_1, w_2, \dots, w_n]$$

$$\mathbf{x} \bullet \mathbf{w} = \boxed{w_1x_1 + w_2x_2 + \dots + w_nx_n}$$

$$y = \mathbf{x} \bullet \mathbf{w} + b$$

# Sigmoid (activation function)

$$y = \mathbf{x} \cdot \mathbf{w} + b$$

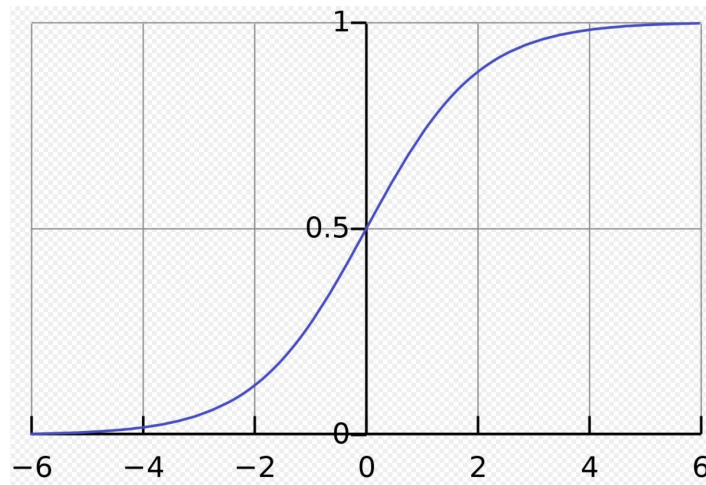
$$y \in (-\infty, +\infty)$$



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



$$\sigma(y) \in (0, 1)$$





CONSTELLATE

# Sigmoid (activation function)

$$y = \mathbf{x} \cdot \mathbf{w} + b$$

Raw score of a house

$$P(\text{Good}|x) = \sigma(\mathbf{x} \cdot \mathbf{w} + b)$$

How likely the house is good

$$P(\text{Bad}|x) = 1 - \sigma(\mathbf{x} \cdot \mathbf{w} + b)$$

How likely the house is bad



CONSTELLATE

# Classifier

$P(\text{Good}|x) = \sigma(\mathbf{x} \cdot \mathbf{w} + b)$  **How likely the house is good**

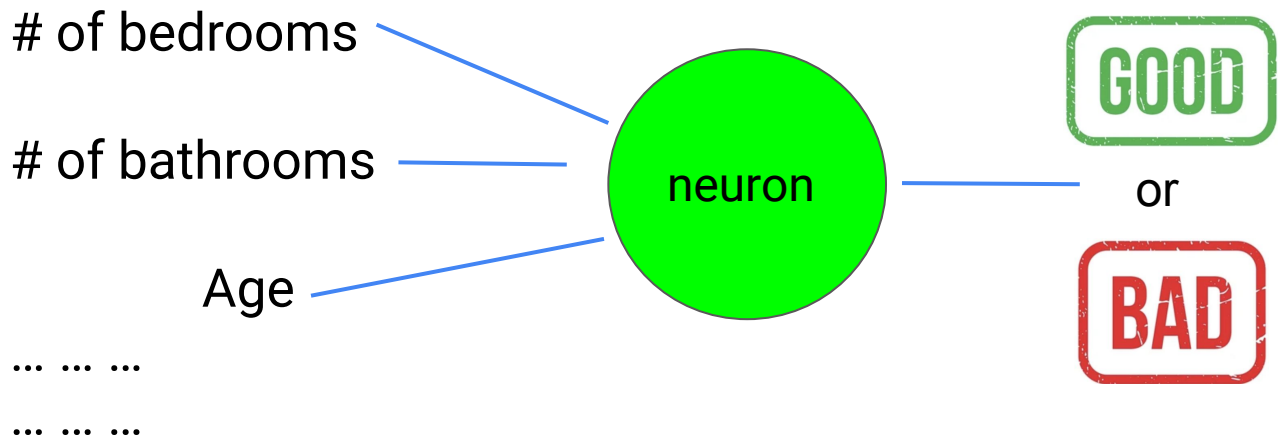
$P(\text{Bad}|x) = 1 - \sigma(\mathbf{x} \cdot \mathbf{w} + b)$  **How likely the house is bad**

$$\text{decision} = \begin{cases} \text{Good} & \text{if } P(\text{Good}|x) > 0.5 \\ \text{Bad} & \text{otherwise} \end{cases}$$



CONSTELLATE

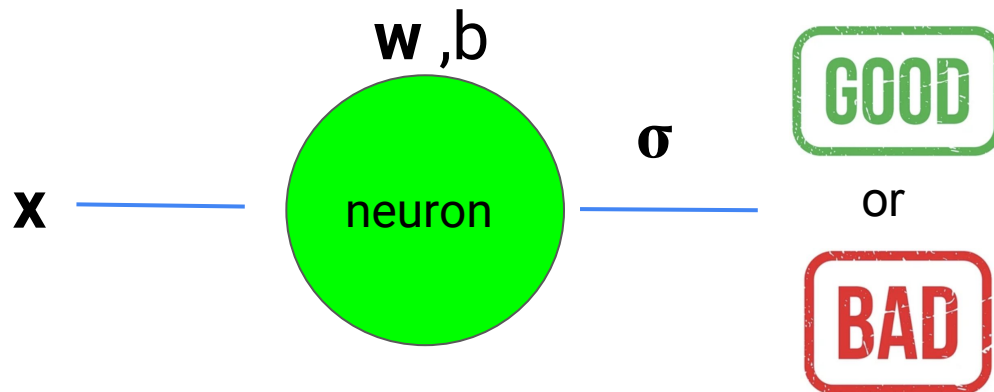
# Neuron as a computation unit





CONSTELLATE

# Neuron as a computation unit

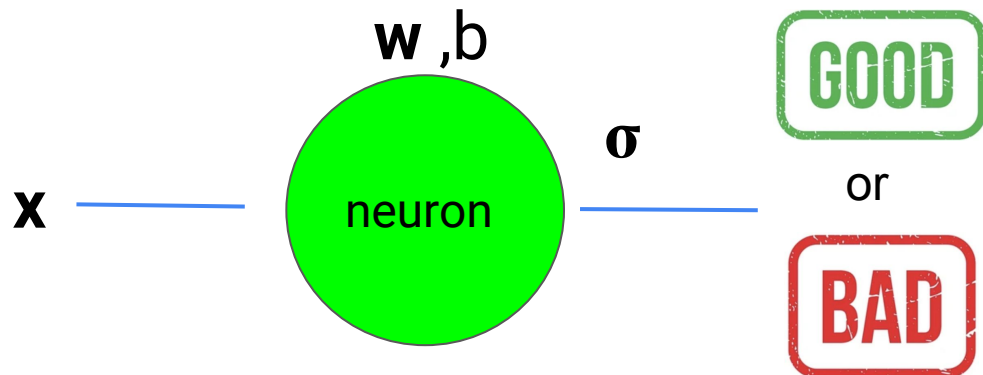


For each observation, the neuron outputs a prediction.  
In our example, for each house, the neuron outputs a prediction as to whether the house is good or bad.



CONSTELLATE

# error/loss/cost



For each observation, the neuron outputs a prediction.  
How does the prediction compare to the correct output?

**error/loss/cost:**

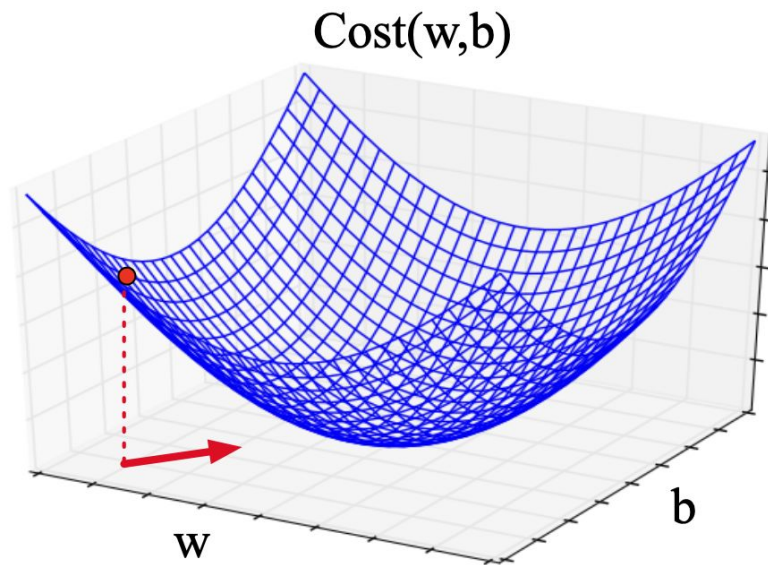
the distance between the predicted value and the correct value



CONSTELLATE

# Minimize error

Our task reduces to calculating the values of the weights  $w$  and the bias term  $b$  such that the error is minimized!

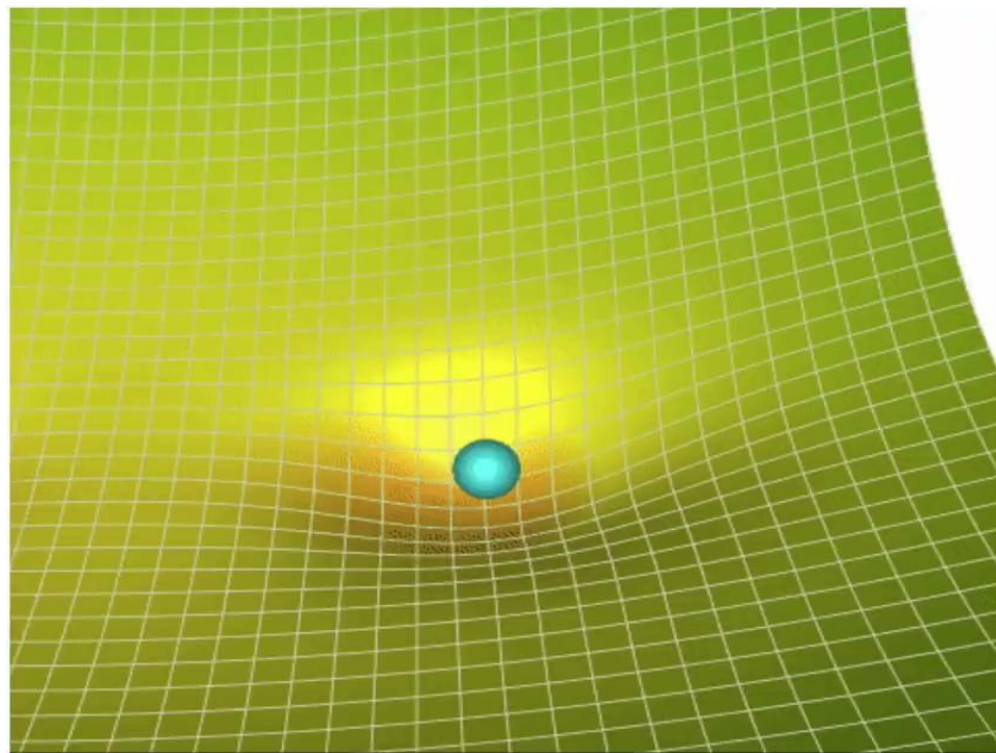






CONSTELLATE

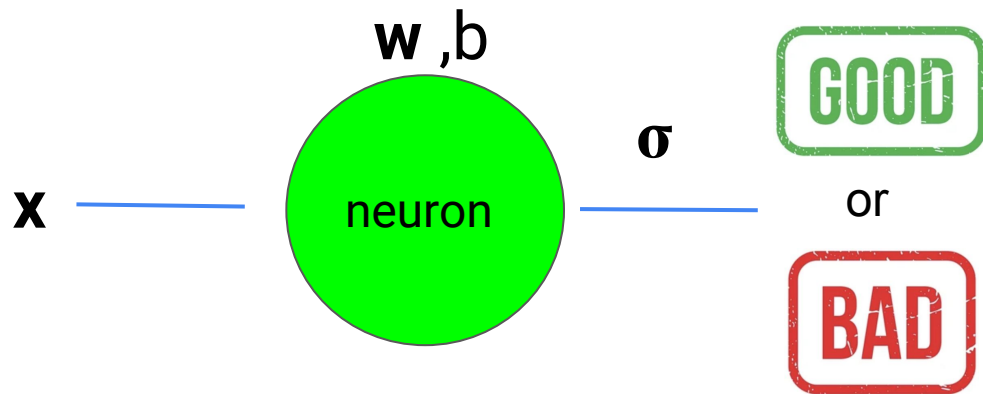
# Gradient descent





CONSTELLATE

# Use cases

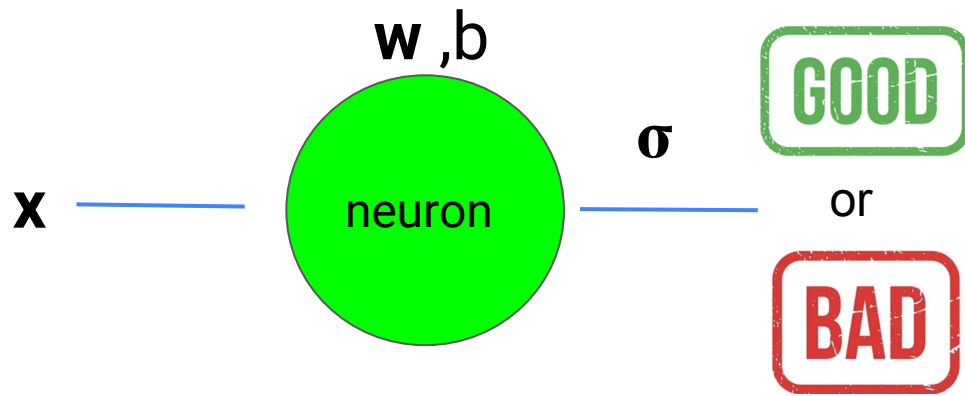


Sentiment analysis  
Medical diagnosis  
Period disambiguation  
Spam detection



CONSTELLATE

# Binomial logistic regression



- Binomial: outcome is dichotomous
- Logistic: sigmoid is a logistic function
- Regression: a statistical technique that relates a dependent variable (which, in our example, good or bad) to one or more independent variables (which, in our example, features of a house)



CONSTELLATE

# Exercise

Suppose we are doing binary sentiment classification of movie reviews. We'd like to assign a review to the positive class or the negative class.



# CONSTELLATE

## Exercise

features

Var	Definition
$x_1$	count(positive lexicon words $\in$ doc)
$x_2$	count(negative lexicon words $\in$ doc)
$x_3$	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
$x_4$	count(1st and 2nd pronouns $\in$ doc)
$x_5$	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$
$x_6$	$\ln(\text{word count of doc})$

It's **hokey**. There are virtually **no** surprises, and the writing is **second-rate**. So why was it so **enjoyable**? For one thing, the cast is **great**. Another **nice** touch is the music. **I** was overcome with the urge to get off the couch and start dancing. It sucked **me** in, and it'll do the same to **you**.

$x_1=3$     $x_2=2$     $x_3=1$     $x_4=3$     $x_5=0$     $x_6=4.19$

$$\begin{aligned} & [x_1, x_2, x_3, x_4, x_5, x_6] \\ &= [3, 2, 1, 3, 0, 4.19] \end{aligned}$$



CONSTELLATE

# Exercise

$$\begin{aligned} &[x_1, x_2, x_3, x_4, x_5, x_6] \\ &= [3, 2, 1, 3, 0, 4.19] \end{aligned}$$

**Suppose we have learned the values of the weights and the bias term**

$$\begin{aligned} &[w_1, w_2, w_3, w_4, w_5, w_6] \\ &= [2.5, -5.0, -1.2, 0.5, 2.0, 0.7] \quad b = 0.1 \end{aligned}$$



CONSTELLATE

# Exercise

$$\begin{aligned} [x_1, x_2, x_3, x_4, x_5, x_6] \\ = [3, 2, 1, 3, 0, 4.19] \end{aligned}$$

$$\begin{aligned} [w_1, w_2, w_3, w_4, w_5, w_6] \\ = [2.5, -5.0, -1.2, 0.5, 2.0, 0.7] \end{aligned}$$

Given the feature vector, weight vector and the bias term, do you know which class the logistic regression model assigns the movie review to, positive or negative?

$$b = 0.1$$

Steps to do this exercise:

1. Calculate a raw score for this review using the feature vector, weight vector and the bias term.
2. Use the sigmoid function(<https://www.vcalc.com/wiki/vcalc/sigmoid-function>) to calculate how likely the review is a positive review.
3. Use the decision boundary 0.5 to determine whether the review is a positive review or a negative review.



CONSTELLATE

Any questions?





CONSTELLATE

**Many neurons, one layer**



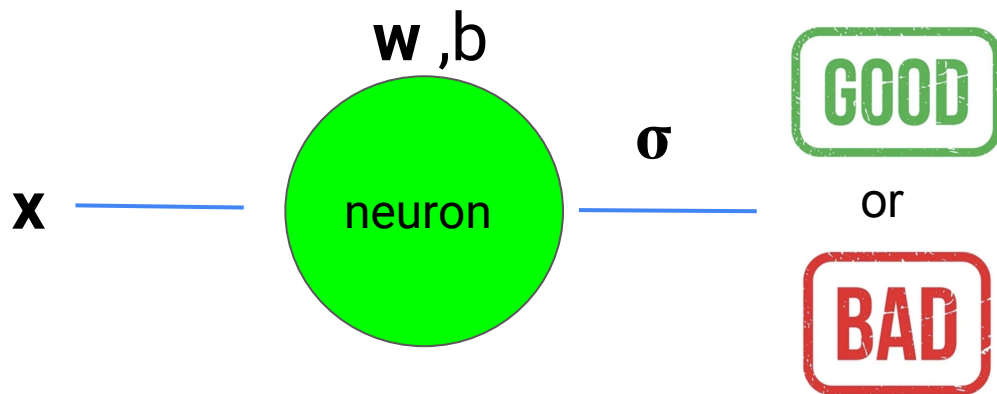
CONSTELLATE

**Recall that...**



CONSTELLATE

# Binomial logistic regression



one neuron  
one layer

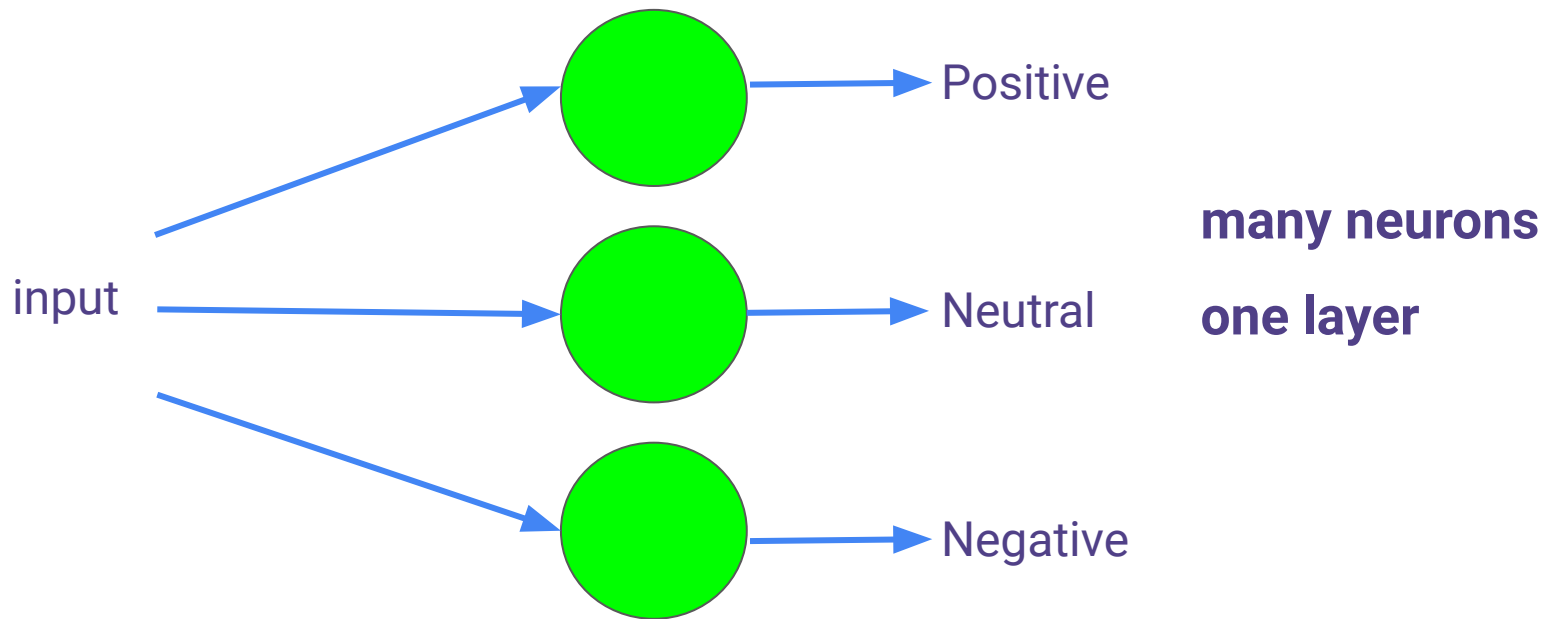
Output is dichotomous



CONSTELLATE

# More than two classes

Suppose we would like to classify some movie reviews into one of three categories: positive, negative or neutral



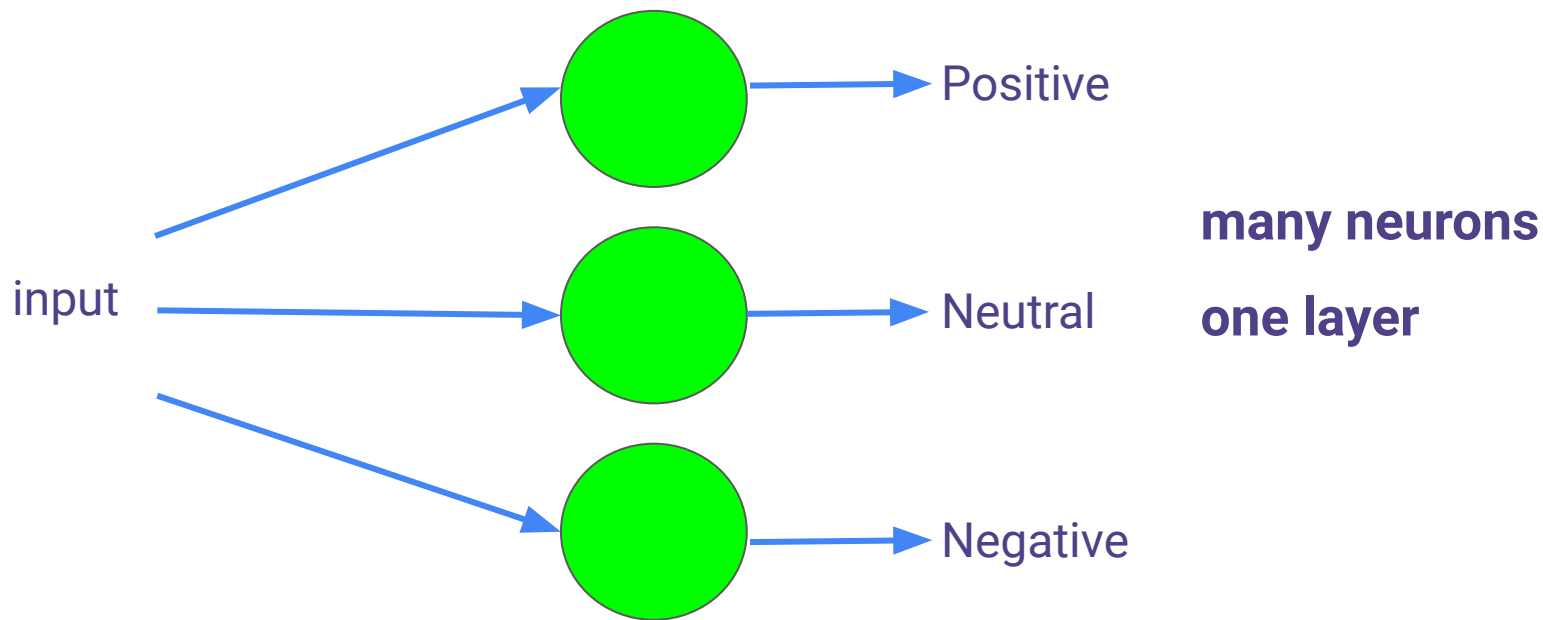
**Output has more than two classes**



CONSTELLATE

# Multinomial logistic regression

Suppose we would like to classify some movie reviews into one of three categories: positive, negative or neutral

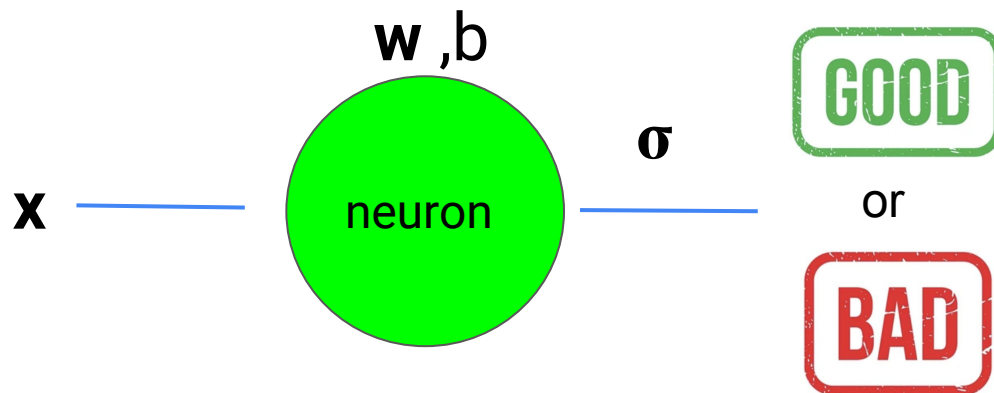


**Output has more than two classes**



CONSTELLATE

Before...

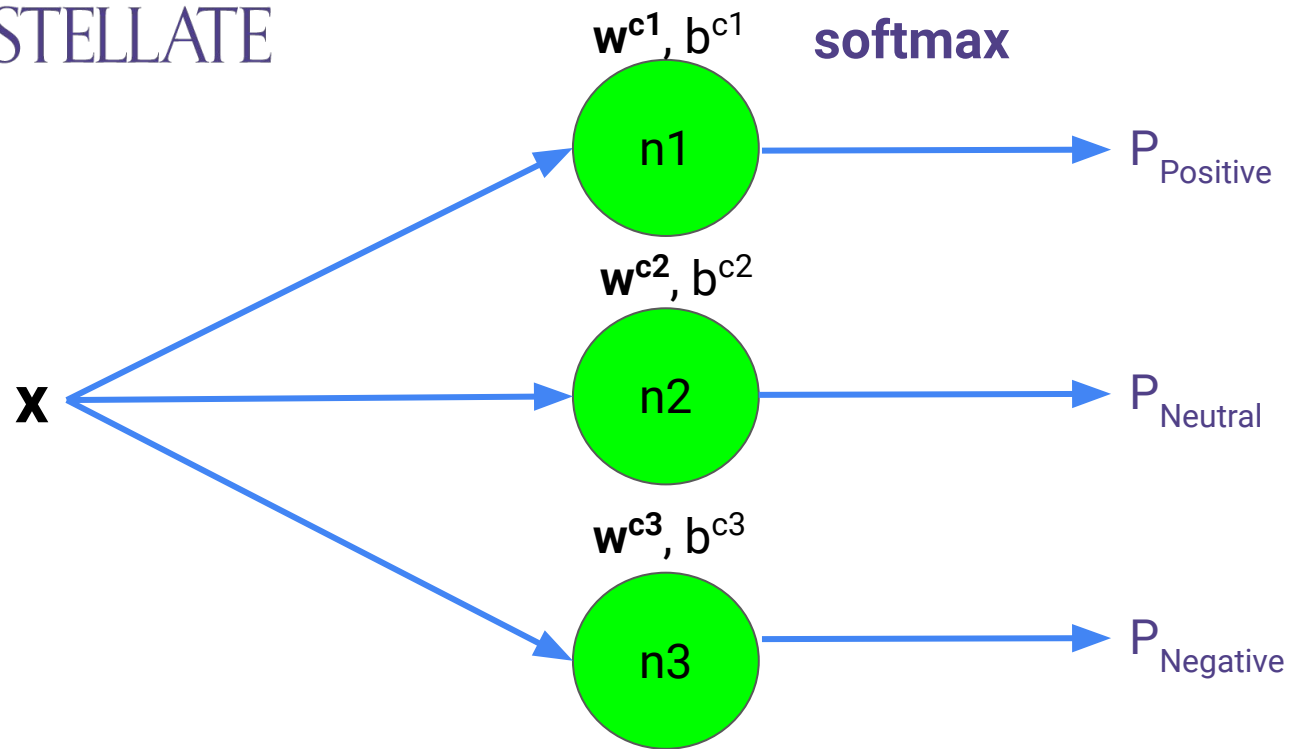


Neuron as a computation unit



Now

CONSTELLATE



many neurons  
one layer

Each neuron is a computation unit!



CONSTELLATE

# If we have $K$ classes...

$K$  different classes,  $K$  different weight vectors  
and bias terms, each for one of the classes!





CONSTELLATE

# Stacked vectors as a matrix

Suppose we have  $n$  features for a given movie review. Since each feature has a weight associated with it, each neuron has a weight vector of length  $n$ .

$$\begin{bmatrix} w_1^{c1} \\ w_2^{c1} \\ \dots \\ w_n^{c1} \end{bmatrix}$$

Weight vector from first neuron, which corresponds to the first class



CONSTELLATE

# Stacked vectors as a matrix

Suppose we have  $n$  features for a given movie review. Since each feature has a weight associated with it, each neuron has a weight vector of length  $n$ .

$$\begin{bmatrix} w_1^{c1} & w_1^{c2} \\ w_2^{c1} & w_2^{c2} \\ \dots & \dots \\ w_n^{c1} & w_n^{c2} \end{bmatrix}$$

Stacking the weight vectors of the first neuron and the second neuron



CONSTELLATE

# Stacked vectors as a matrix

Suppose we have  $n$  features for a given movie review. Since each feature has a weight associated with it, each neuron has a weight vector of length  $n$ .

$$\begin{bmatrix} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & \dots & \dots \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{bmatrix}$$

$$n \times K$$

Stacking the weight vectors from all  $K$  neurons



CONSTELLATE

# Given one observation...

**w**

$$\begin{array}{c} \mathbf{x} \\ [x_1 \ x_2 \ \dots \ x_n] \\ 1 \times n \end{array} \bullet \begin{array}{c} \left[ \begin{array}{cccc} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & \dots & \dots \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{array} \right] \\ n \times K \end{array} + \begin{array}{c} \mathbf{b} \\ [b^{c1} \ b^{c2} \ \dots \ b^{cK}] \\ 1 \times K \end{array}$$



## CONSTELLATE

$$\begin{matrix} & & & & & \mathbf{w} \\ & & & & & \\ & & & & & \\ \mathbf{x} & & & & & \\ [x_1 & x_2 & \dots & x_n] \bullet & \begin{bmatrix} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & & \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{bmatrix} & + & \begin{matrix} \mathbf{b} \\ [b^{c1} & b^{c2} & \dots & b^{cK}] \end{matrix} \end{matrix}$$



# CONSTELLATE

$$\begin{array}{c} \mathbf{x} \\ [x_1 \ x_2 \ \dots \ x_n] \\ 1 \times n \end{array} \bullet \begin{array}{c} \mathbf{w} \\ \left[ \begin{array}{cccc} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & \dots & \dots \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{array} \right] \\ n \times K \end{array}$$

$$= [x_1 \cdot w_1^{c1} + x_2 \cdot w_2^{c1} + \dots + x_n \cdot w_n^{c1}, \dots, x_1 \cdot w_1^{cK} + x_2 \cdot w_2^{cK} + \dots + x_n \cdot w_n^{cK}] \\ 1 \times K$$



# CONSTELLATE

$$\begin{array}{c} \mathbf{x} \\ [x_1 \ x_2 \ \dots \ x_n] \\ 1 \times n \end{array} \bullet \begin{array}{c} \mathbf{w} \\ \left[ \begin{array}{cccc} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & \dots & \dots \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{array} \right] \\ n \times K \end{array}$$
$$= [x_1 \cdot w_1^{c1} + x_2 \cdot w_2^{c1} + \dots + x_n \cdot w_n^{c1}, x_1 \cdot w_1^{c2} + x_2 \cdot w_2^{c2} + \dots + x_n \cdot w_n^{c2}, \dots]$$
$$1 \times K$$



# CONSTELLATE

$$\begin{array}{c} \mathbf{x} \\ [x_1 \ x_2 \ \dots \ x_n] \\ 1 \times n \end{array} \bullet \begin{array}{c} \mathbf{w} \\ \left[ \begin{array}{cccc} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & \dots & \dots \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{array} \right] \\ n \times K \end{array}$$

$$\begin{array}{c} = [x_1 \cdot w_1^{c1} + x_2 \cdot w_2^{c1} + \dots + x_n \cdot w_n^{c1}, \dots, x_1 \cdot w_1^{cK} + x_2 \cdot w_2^{cK} + \dots + x_n \cdot w_n^{cK}] \\ 1 \times K \end{array}$$





# CONSTELLATE

**W****X**

$$\begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}$$

 $\bullet$ **W<sup>c1</sup>****W<sup>c2</sup>****W<sup>cK</sup>**

$$\begin{bmatrix} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & \dots & \dots \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \cdot w_1^{c1} + x_2 \cdot w_2^{c1} + \dots + x_n \cdot w_n^{c1}, \dots, x_1 \cdot w_1^{cK} + x_2 \cdot w_2^{cK} + \dots + x_n \cdot w_n^{cK} \end{bmatrix}$$

**X·W<sup>c1</sup>****X·W<sup>cK</sup>**



# CONSTELLATE

$$\begin{aligned} & \mathbf{x} \cdot \mathbf{w} \\ & [x_1 \ x_2 \ \dots \ x_n] \bullet \begin{matrix} \mathbf{w}^{c1} & \mathbf{w}^{c2} & & \mathbf{w}^{cK} \\ \left[ \begin{array}{cccc} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & \dots & \dots \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{array} \right] \end{matrix} \\ & = [\mathbf{x} \cdot \mathbf{w}^{c1}, \mathbf{x} \cdot \mathbf{w}^{c2}, \dots, \mathbf{x} \cdot \mathbf{w}^{cK}] \end{aligned}$$



# Given one observation...

CONSTELLATE

$$[x_1 \ x_2 \ \dots \ x_n] \bullet \begin{bmatrix} w_1^{c1} & w_1^{c2} & \dots & w_1^{cK} \\ w_2^{c1} & w_2^{c2} & \dots & w_2^{cK} \\ \dots & \dots & \dots & \dots \\ w_n^{c1} & w_n^{c2} & \dots & w_n^{cK} \end{bmatrix} + [b^{c1} \ b^{c2} \ \dots \ b^{cK}]$$

$$[\mathbf{x} \cdot \mathbf{w}^{c1}, \mathbf{x} \cdot \mathbf{w}^{c2}, \dots, \mathbf{x} \cdot \mathbf{w}^{cK}]$$



CONSTELLATE

# Add the bias

$$\begin{aligned} & [\mathbf{x} \cdot \mathbf{w}^{c1}, \mathbf{x} \cdot \mathbf{w}^{c2}, \dots, \mathbf{x} \cdot \mathbf{w}^{cK}] + [b^{c1} \ b^{c2} \ \dots \ b^{cK}] \\ &= [\mathbf{x} \cdot \mathbf{w}^{c1} + b^{c1}, \mathbf{x} \cdot \mathbf{w}^{c2} + b^{c2}, \dots, \mathbf{x} \cdot \mathbf{w}^{cK} + b^{cK}] \end{aligned}$$



CONSTELLATE

# Add the bias

$$\begin{aligned} & [\mathbf{x} \cdot \mathbf{w}^{c1}, \mathbf{x} \cdot \mathbf{w}^{c2}, \dots, \mathbf{x} \cdot \mathbf{w}^{cK}] + [b^{c1}, b^{c2}, \dots, b^{cK}] \\ &= [\mathbf{x} \cdot \mathbf{w}^{c1} + b^{c1}, \mathbf{x} \cdot \mathbf{w}^{c2} + b^{c2}, \dots, \mathbf{x} \cdot \mathbf{w}^{cK} + b^{cK}] \end{aligned}$$



CONSTELLATE

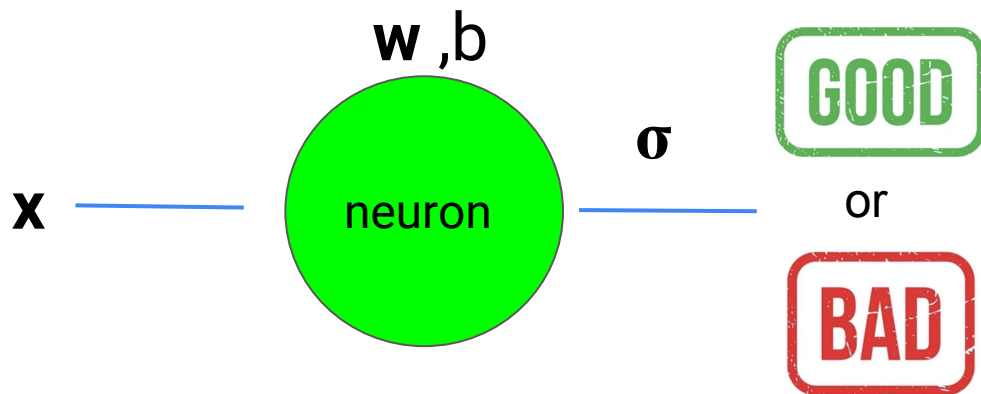
# Add the bias

$$\begin{aligned} & [\mathbf{x} \cdot \mathbf{w}^{c1}, \mathbf{x} \cdot \mathbf{w}^{c2}, \dots, \mathbf{x} \cdot \mathbf{w}^{cK}] + [b^{c1} \ b^{c2} \ \dots \ b^{cK}] \\ &= [\mathbf{x} \cdot \mathbf{w}^{c1} + b^{c1}, \mathbf{x} \cdot \mathbf{w}^{c2} + b^{c2}, \dots, \mathbf{x} \cdot \mathbf{w}^{cK} + b^{cK}] \end{aligned}$$



CONSTELLATE

# Recall that...



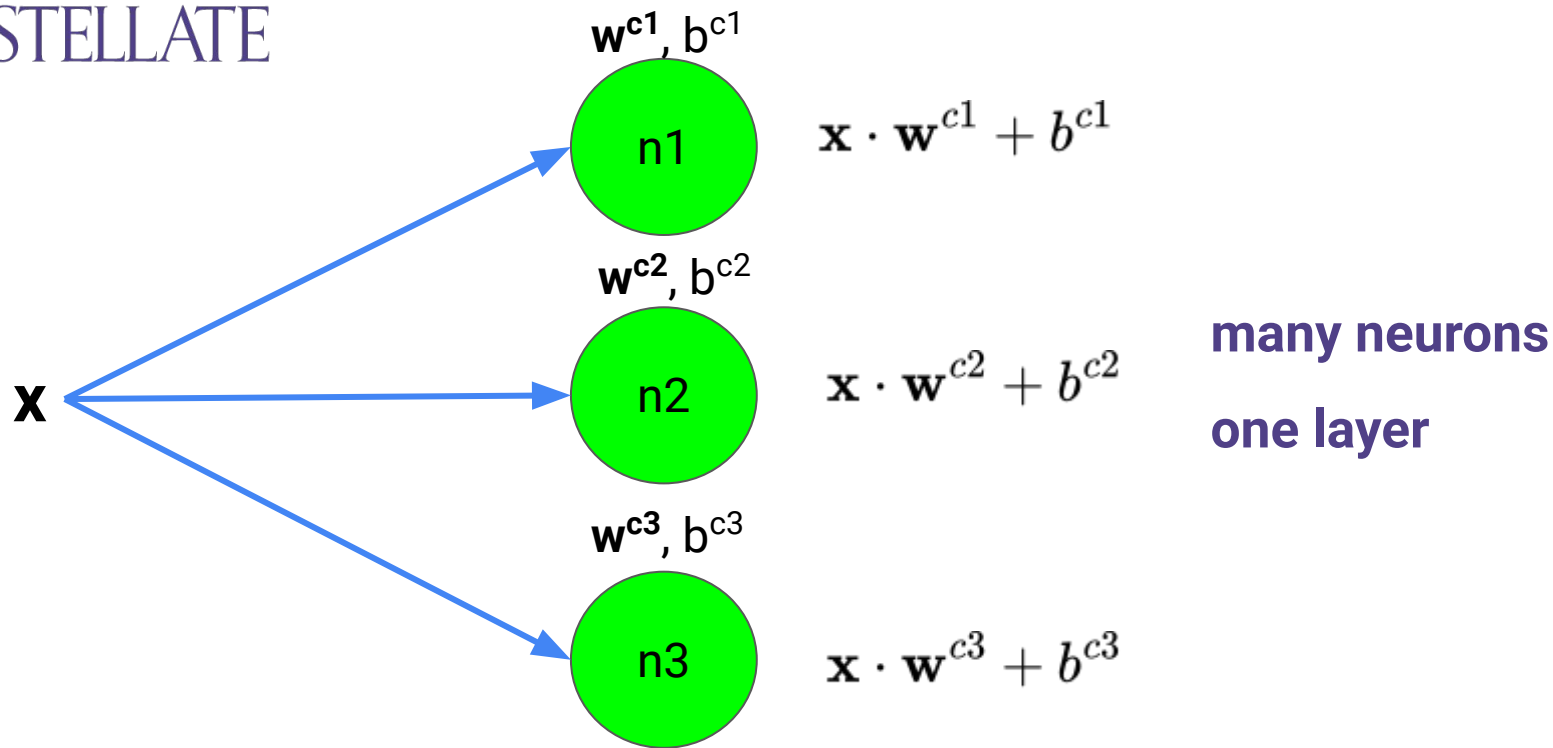
one neuron  
or  
one layer

We compute a raw score  $\mathbf{x} \cdot \mathbf{w} + b$



CONSTELLATE

# Multinomial logistic regression



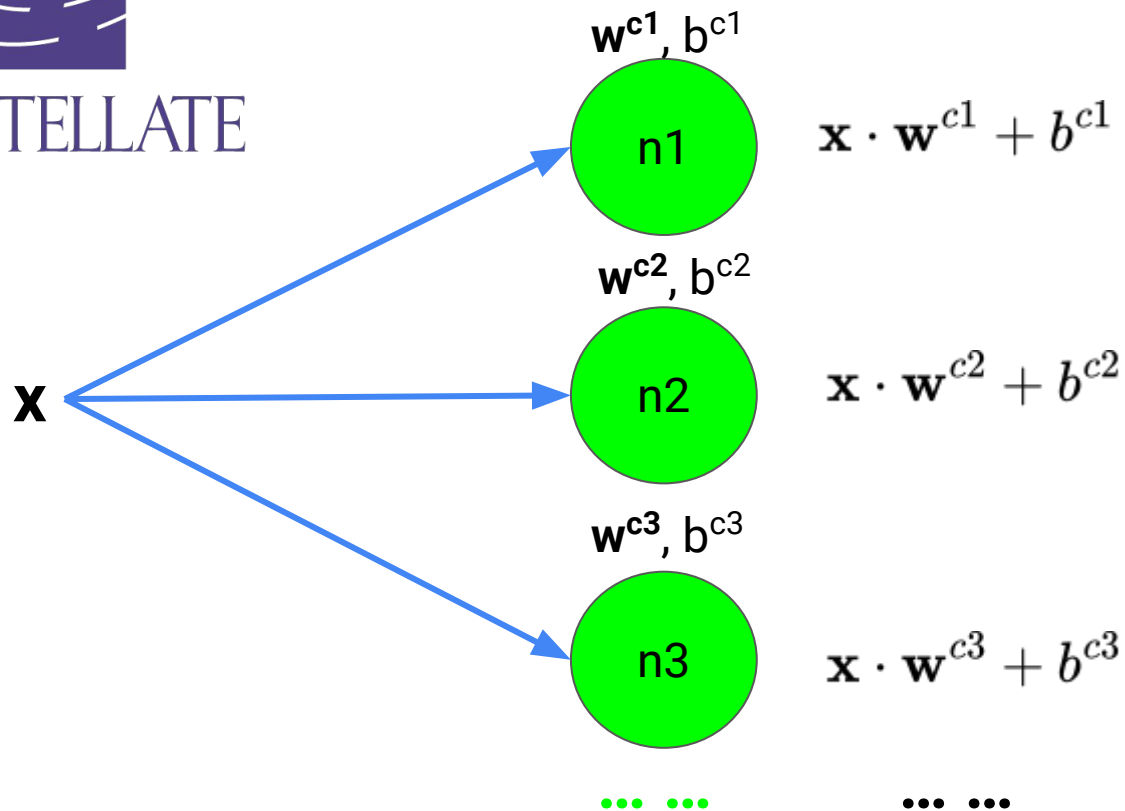
Each neuron is a computation unit!





CONSTELLATE

# Now what?



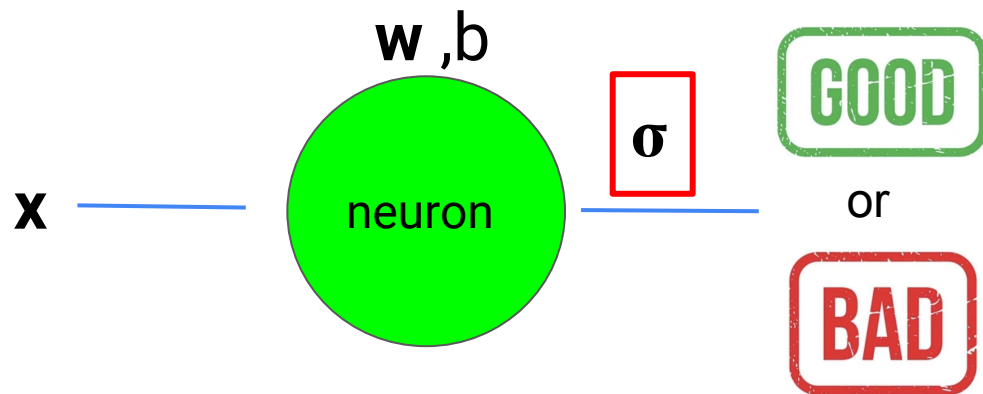
many neurons  
one layer

We have  $K$  raw scores, now what?



CONSTELLATE

# Before...



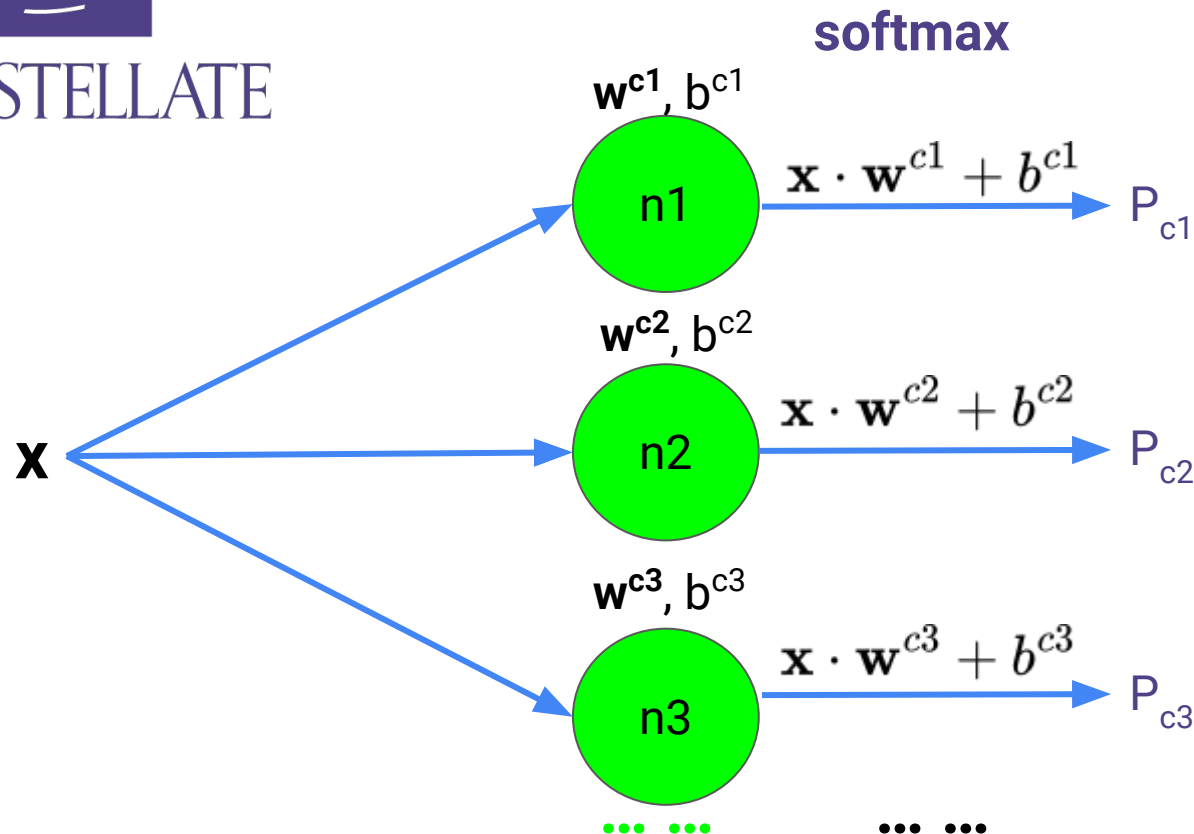
one neuron  
one layer

We use the sigmoid function  $\sigma$  to squash all the house values to probabilities of being good or bad



# Now

CONSTELLATE



many neurons  
one layer

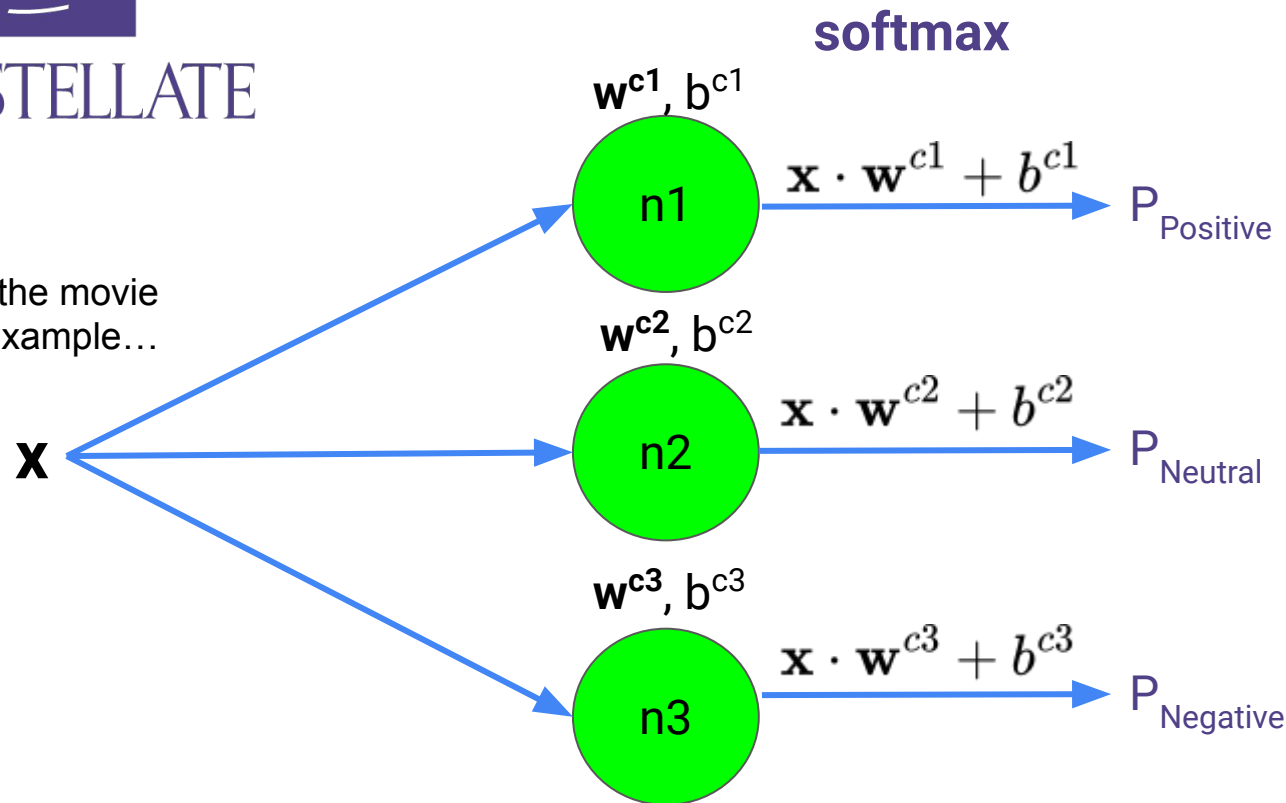
We use the softmax function to squash all K values to probabilities



CONSTELLATE

# For example

Back to the movie  
review example...



**many neurons  
one layer**

We use the softmax function to squash the 3 values to probabilities of being a positive, neutral or negative movie review



CONSTELLATE

# Softmax (activation function)

## Softmax function

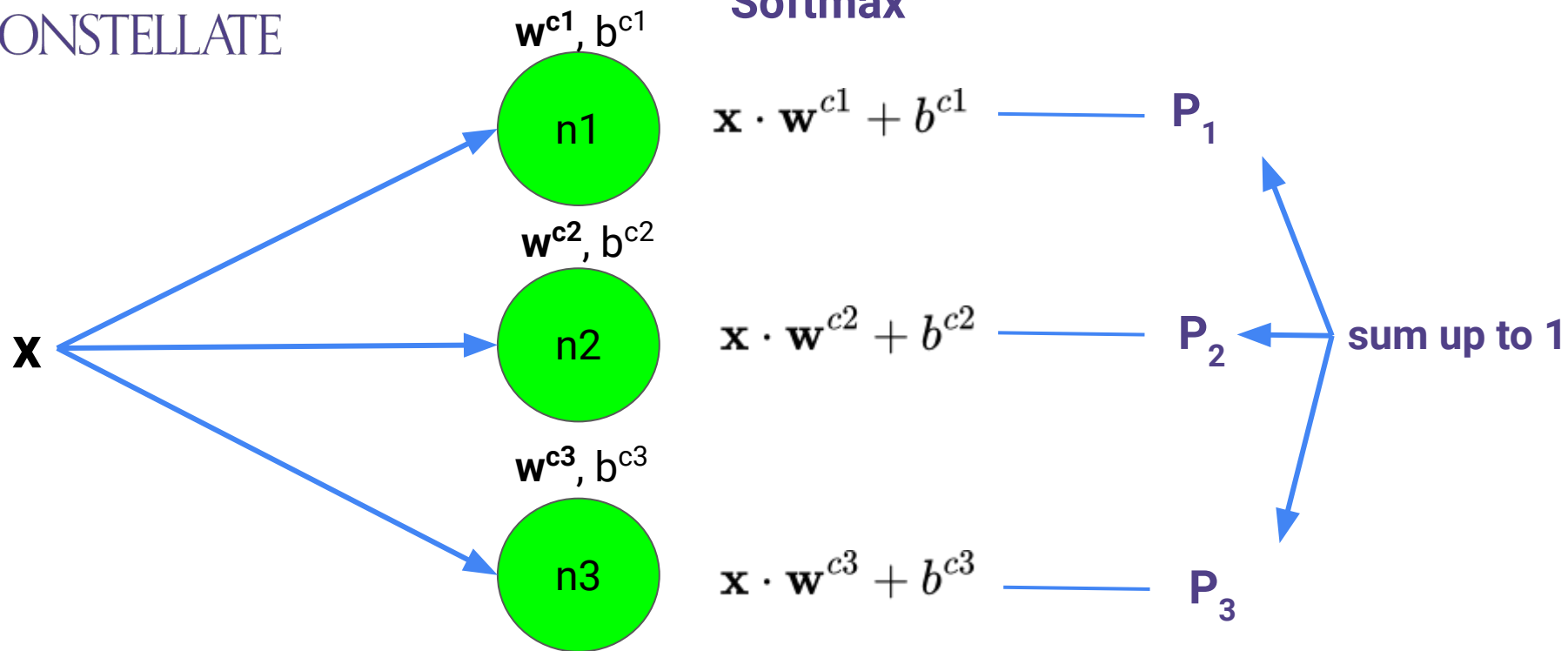
The softmax function takes a vector of  $K$  values  $[z_1, z_2, \dots, z_K]$ , and maps the values to a probability distribution where each value is in the range  $(0,1)$  and the probabilities sum up to one.

$$\text{softmax}(\mathbf{z}) = \left[ \frac{e^{z_1}}{\sum_{i=1}^K e^{z_i}}, \quad \frac{e^{z_2}}{\sum_{i=1}^K e^{z_i}}, \quad \dots, \quad \frac{e^{z_K}}{\sum_{i=1}^K e^{z_i}} \right]$$



CONSTELLATE

# Softmax





CONSTELLATE

**error**

Our task reduces to calculating the values of the weights **w** and the bias terms **b** such that the error is minimized!



CONSTELLATE

# Use cases

Multinomial logistic regression is appropriate for any situation where a limited number of outcome categories (more than two) are being modeled and where those outcome categories have no order.

Image classification

Voting choice in elections with multiple candidates

Career options by students

...





CONSTELLATE

# Exercise

$$[2 \quad 0 \quad -5] \bullet \begin{bmatrix} 2 & -1 \\ 2 & 3 \\ 1 & 1 \end{bmatrix}$$



CONSTELLATE

**Answer**

$$[2 \quad 0 \quad -5] \bullet \begin{bmatrix} 2 & -1 \\ 2 & 3 \\ 1 & 1 \end{bmatrix} = [-1 \quad -7]$$



CONSTELLATE

# Exercise

$$\begin{bmatrix} 2 & 0 & -5 \\ 4 & 1 & 3 \end{bmatrix} \bullet \begin{bmatrix} 2 & -1 \\ 2 & 3 \\ 1 & 1 \end{bmatrix}$$



CONSTELLATE

**Answer**

$$\begin{bmatrix} 2 & 0 & -5 \\ 4 & 1 & 3 \end{bmatrix} \bullet \begin{bmatrix} 2 & -1 \\ 2 & 3 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} -1 & -7 \\ 13 & 2 \end{bmatrix}$$



CONSTELLATE

# Interim Summary

Supervised Learning

**Data -> ML algorithm -> Quality metric**



CONSTELLATE

Any questions?

# References

- Jiang, L. (2020). A Visual Explanation of Gradient Descent Methods (Momentum, Ada-Grad, RMSProp, Adam). June-2020. [online]. Available: <https://towardsdatascience.com/a-visual-explanation-of-gradient-descent-methods-momentum-adagrad-rmsprop-adam-f898b102325c>
- Jurafsky, Daniel, and James H. Martin. (2023). [Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition.](#)
- McCormick, C. (2016). Word2vec tutorial-the skip-gram model. Apr-2016.[Online]. Available: <http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model>.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). [Distributed representations of words and phrases and their compositionality.](#) *Advances in neural information processing systems*, 26.