

# Lecture 1:

# Introduction to Neural Networks and Deep Learning



15.S04: Hands-on Deep Learning  
Spring 2024  
**Farias, Ramakrishnan**



# Welcome to Hands-on Deep Learning (HODL)!

# Prerequisites



- Familiarity with Python at this level
- Familiarity with fundamental machine learning concepts (such as training/validation/testing, overfitting/underfitting, and regularization).
- If you have taken 15.071/15.072 (or will be taking it concurrently) OR if you have other relevant coursework or work experience, you should be fine.
- If you aren't sure, email us.

The screenshot shows a web browser window for the Kaggle Learn Python course. The URL in the address bar is [kaggle.com/learn/python](https://kaggle.com/learn/python). The main title is "Python". Below it, a sub-headline says "Learn the most important language for data science.". A large "Begin Course" button is prominent, with a progress bar indicating "5 hours to go". Below the button, there are two tabs: "Course" (which is selected) and "Discussion". The "Lessons" section lists three items: 1. Hello, Python (A quick introduction to Python syntax, variable assignment, and numbers), 2. Functions and Getting Help (Calling functions and defining our own, and using Python's builtin documentation), and 3. Booleans and Conditionals. To the right of the lessons, there are "Tutorial" and "Exercise" tabs. On the far left, there is a vertical sidebar with various icons: a plus sign, a magnifying glass, a person, a document, a gear, and a downward arrow.

# Grading

## Grading

Your course grade will be based on 2 homework assignments, a final project, and class participation:

Class participation	10%
Homework assignments (25% x 2)	50%
Final project	40%

# Waitlist (Section B)

- Registered students
  - If you decide to drop HODL at any time, please let the TA know ASAP so we can offer a slot to the waitlist students
- Waitlist students:
  - Will send out update by EOD Wednesday based on drops.
  - If your plans change, let the TA know so we can take you off the waitlist
  - If we are unable to offer you a slot AND you are graduating in 2024, we will offer you the opportunity to register as a Listener

# “Housekeeping” (Section B)



- Class starts at 10:05 and ends at 11:25
- No mobile phones/laptops in class. OK to take notes on a tablet
- Lecture slides & links to Colabs will be posted to Canvas after class
- An absence is considered excused only if it is due to medical or religious reasons. Otherwise, we will mark it as unexcused

# TA – Gabriel Afriat (Section B)

- PhD student at the Operations Research Center working with Prof. Rahul Mazumder
- Master of Business Analytics (2022)
- Bachelor and Master's degree in Engineering from CentraleSupélec, France (2021)
- Research Interests: machine learning and optimization (sparsity in GAMs, model compression)
- Come see me for the best baguettes/croissants in Boston!



# Vivek Farias

Patrick J. McGovern (1959) Professor

## ➤ How to reach me outside of class:

- Please email me ([vivekf@mit.edu](mailto:vivekf@mit.edu)) to schedule an appointment.

- Ph.D. in EE at Stanford University. At MIT since 2007. Head of Operations Group
- **Research**
  - Reinforcement Learning and High-Dimensional Statistics with applications in Commerce and Biology.
  - Lab alums now professors @ MIT (x2), Columbia (x2), NYU (x2), CMU, INSEAD, UBC, LBS, HBS
- **Teaching**
  - (Co-) Developed: Intro to Ops (15.761, 15.778), Digital Product Management (and the PM Certificate) (15.785), Hands-On Deep Learning (15.773), and several PhD classes
  - MIT-wide Jamieson Prize for teaching, Sloan Teacher of the Year, etc.
- **Industry and Entrepreneurship**
  - Before MIT: Hedge Fund Industry @GMO. Co-architected first high-freq fund
  - Three ventures founded on my research that have collectively raised ~\$400M:
    - **Celect** (2014-19). Co-Founder/ CTO. *Supply Chain AI*. **Acquired by Nike** in August 2019.
    - **Seer Bio** (2017-). *High Throughput Proteomics*. **IPO** in 2020 (SEER) + Spin-out (Prognomiq)
    - **Cimulate** (2023-). Co-Founder/CTO. *Customer-GPT*. Seed Stage.
  - Misc: Active Angel and LP (eg. AIX Ventures). SAB member in multiple ventures.

# Why HODL?

# Why did we create HODL?



- DL is one of the most exciting and profound technology developments of our lifetimes
- It is important for Sloanies to understand how to use DL to transform businesses and create exciting new products/services
- While MIT has other (excellent) DL courses, we wanted one that was a better fit for Sloan

# HODL's “philosophy”

- Focus on the key concepts that underlie DL
- Skip the math\* (but we are happy to geek out in office hours and/or suggest readings for those who are interested)
- Focus on coding DL models ...
  - It is the only way to develop a visceral (not just intellectual) understanding of how this stuff works
  - Successful new products and services are often inspired by hands-on tinkering
- ... but only up to a point
  - We aren't trying to teach you how to be ML engineers
  - **But we want you to be able to build a V1.0 DL model by yourself without looking for a Data Scientist or ML Engineer to help you out**

---

\*If you are looking for a ‘mathy’ DL course, this course won’t be a good fit and you may want to consider dropping it.

## SCHEDULE (subject to change)

Mon	05 Feb	Introduction to Neural Networks and Deep Learning (motivation, smart representations, layers, activations, architectures); Training Deep NNs (loss functions, gradient descent, regularization)	
Wed	07 Feb	Training Deep NNs (continued); Introduction to Keras/Tensorflow (tensors, functional API, training, evaluation); Application to tabular data	
Mon	12 Feb	Deep Learning for Computer Vision – Building Convolutional Neural Networks from scratch	HW Assignment #1 posted
Wed	14 Feb	Deep Learning for Computer Vision – Transfer Learning and Fine-tuning; Introduction to HuggingFace	
Mon	19 Feb	<i>HOLIDAY (President's Day)</i>	
Tue	20 Feb	Deep Learning for Natural Language – The Basics (Tokenization, N-grams, Bag-of-Words)	HW Assignment #1 due
Wed	21 Feb	Deep Learning for Natural Language – Embeddings	HW Assignment #2 posted
Mon	26 Feb	Deep Learning for Natural Language – Transformers	Project Proposal Due
Wed	28 Feb	Deep Learning for Natural Language – Transformers, Self-Supervised Learning	
Mon	04 Mar	Generative AI - Large Language Models (LLMs) and Retrieval Augmented Generation (RAG)	HW Assignment #2 due
Wed	06 Mar	Generative AI – Adapting LLMs with Parameter-Efficient Fine-Tuning	
Mon	11 Mar	Generative AI – Text-to-Image Models	Projects due
Wed	13 Mar	Project Presentations	

# The field of Artificial Intelligence originated in 1956



In the back row from left to right are Oliver Selfridge, Nathaniel Rochester, Marvin Minsky, and John McCarthy. In front on the left is Ray Solomonoff; on the right, Claude Shannon. The identity of the person between Solomonoff and Shannon remained a mystery for some time. THE MINSKY FAMILY

<https://spectrum.ieee.org/dartmouth-ai-workshop>

# MIT was well-represented 🤘



Marvin  
Minsky

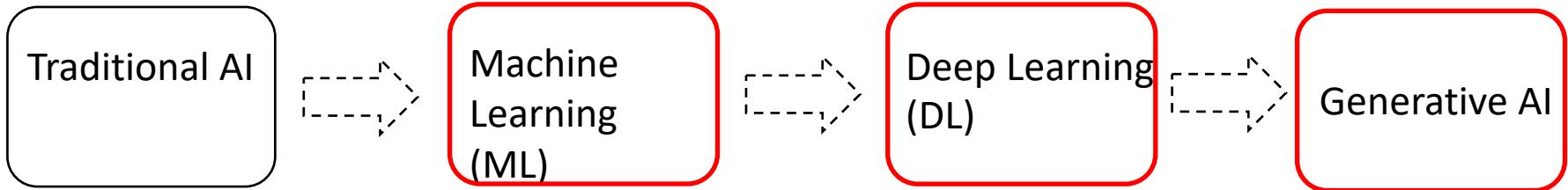
John McCarthy

Claude Shannon

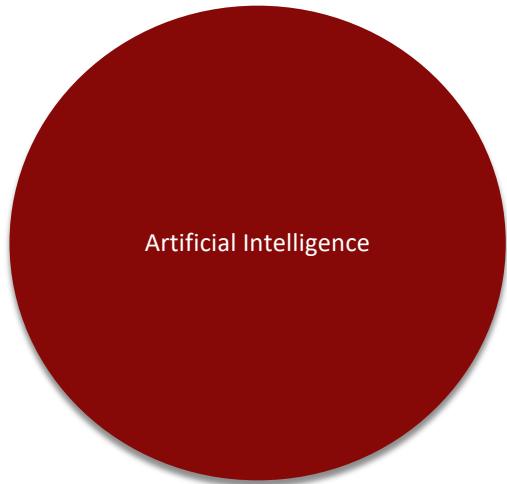
In the back row from left to right are Oliver Selfridge, Nathaniel Rochester, Marvin Minsky, and John McCarthy. In front on the left is Ray Solomonoff; on the right, Claude Shannon. The identity of the person between Solomonoff and Shannon remained a mystery for some time. THE MINSKY FAMILY

<https://spectrum.ieee.org/dartmouth-ai-workshop>

# In the decades since its founding, it has gone through several “breakthroughs”



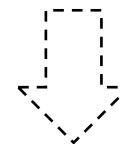
# The traditional approach to AI



*The Goal: Give computers the ability to do tasks that traditionally only humans have been able to do*

Traditional approach:

Ask human experts how they do it, write it down as IF-THEN rules, explicitly program these rules into the computer



Success in only a few areas

# Why is this so difficult?



- “We know more than we can tell” (Polanyi’s Paradox)
  - We can do lots of things easily but find it very hard to describe how exactly we do them
- We can’t write down if-then rules to cover all situations, edge cases etc. (i.e., we can’t generalize to new situations)

# To address this problem, a different approach was developed

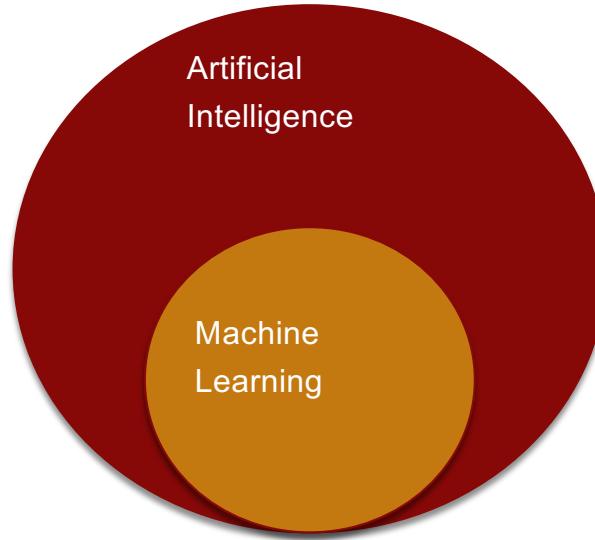
Instead of explicitly telling the computer what to do ...

# To address this problem, a different approach was developed

Instead of explicitly telling the computer what to do ...

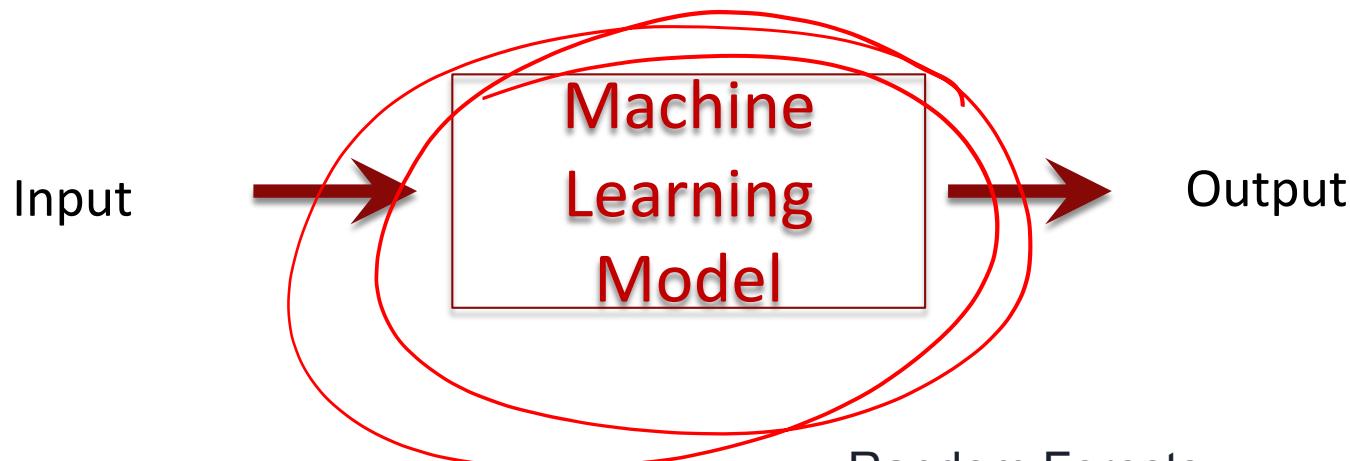
*Provide the computer with lots of examples of inputs-and-outputs and use statistical techniques to learn the relationship between inputs and outputs*

# This is Machine Learning



“learn from input-output examples using statistical techniques\*”

# There are numerous ways\* to create Machine Learning models



- Linear Regression
- Logistic Regression
- Classification and Regression Trees
- Support Vector Machines
- Random Forests
- Gradient Boosted Machines
- Neural Networks
- .....

---

\*Covered in detail in courses like The Analytics Edge

**Springer Series in Statistics**

Trevor Hastie  
Robert Tibshirani  
Jerome Friedman

# The Elements of Statistical Learning

Data Mining, Inference, and Prediction

Second Edition



Gupta, Michael Jordan, Shanti Gopalam, Radford Neal, Jorge Picazo, Bogdan Popescu, Olivier Renaud, Saharon Rosset, John Storey, Ji Zhu, Mu Zhu, two reviewers and many students read parts of the manuscript and offered helpful suggestions. John Kimmel was supportive, patient and helpful at every phase; MaryAnn Brickner and Frank Ganz headed a superb production team at Springer. Trevor Hastie would like to thank the statistics department at the University of Cape Town for their hospitality during the final stages of this book. We gratefully acknowledge NSF and NIH for their support of this work. Finally, we would like to thank our families and our parents for their love and support.

Trevor Hastie  
Robert Tibshirani  
Jerome Friedman  
Stanford, California  
May 2001

Chapter	What's new
1. Introduction	LAR algorithm and generalizations of the lasso
2. Overview of Supervised Learning	Lasso path for logistic regression
3. Linear Methods for Regression	Additional illustrations of RKHS
4. Linear Methods for Classification	
5. Basis Expansions and Regularization	
6. Kernel Smoothing Methods	
7. Model Assessment and Selection	Strengths and pitfalls of cross-validation
8. Model Inference and Averaging	
9. Additive Models, Trees, and Related Methods	
10. Boosting and Additive Trees	New example from ecology; some material split off to Chapter 16.
11. Neural Networks	Bayesian neural nets and the NIPS 2003 challenge
12. Support Vector Machines and Flexible Discriminants	Path algorithm for SVM classifier
13. Prototype Methods and Nearest-Neighbors	
14. Unsupervised Learning	Spectral clustering, kernel PCA, sparse PCA, non-negative matrix factorization, archetypal analysis, nonlinear dimension reduction, Google page rank algorithm, a direct approach to ICA
15. Random Forests	New
16. Ensemble Learning	New
17. Undirected Graphical Models	New
18. High-Dimensional Problems	New

Some further notes:

- Our first edition was unfriendly to colorblind readers; in particular, we tended to favor red/green contrasts which are particularly troublesome. We have changed the color palette in this edition to a large extent, replacing the above with an orange/blue contrast.
- We have changed the name of Chapter 6 from “Kernel Methods” to “Kernel Smoothing Methods”, to avoid confusion with the machine-learning kernel method that is discussed in the context of support vector machines (Chapter 12) and more generally in Chapters 5 and 14.
- In the first edition, the discussion of error-rate estimation in Chapter 7 was sloppy, as we did not clearly differentiate the notions of conditional error rates (conditional on the training set) and unconditional rates. We have fixed this in the new edition.

Machine Learning has had tremendous impact and is used worldwide across numerous applications (e.g., credit scoring, loan granting, disease prediction, demand forecasting, ....) but only if the input data is structured

---

# Structured input data = data that can be “numericalized” into a spreadsheet\*

INPUT						OUTPUT
Age	Smoker	Exercise	Cholesterol	Family History	Blood Pressure	Cardiac Arrest
30	No	120	190	Yes	120/80	No
45	Yes	30	220	No	130/90	Yes
50	No	60	210	Yes	125/85	No
35	Yes	45	230	No	135/88	Yes
40	No	150	180	Yes	118/78	No
55	Yes	10	240	Yes	140/92	Yes
28	No	180	170	No	115/75	No
60	Yes	20	250	Yes	145/95	Yes
48	No	90	200	No	128/82	No
53	Yes	35	235	Yes	133/89	Yes

\*informal definition

# But the situation is different for unstructured input data (images, videos, text, audio, ...)

**Images**



**Text**

*Four score and seven years  
ago our fathers brought forth,  
upon this continent, ...*

**Audio**

...

The reason: The “raw form” of unstructured data has no intrinsic meaning



	Red	Green	Blue
[,1]	[,2]	[,3]	[,4]
[1,]	147	131	138
[2,]	140	131	141
[3,]		[,1]	[,2]
[4,]		[,3]	[,4]
[5,]		[,5]	[,6]
[6,]		[,7]	[,8]
[7,]		[,9]	[,10]
[8,]			
[9,]			
[10,]			

To use ML on unstructured data, we had to manually create a better representation of the data first\*

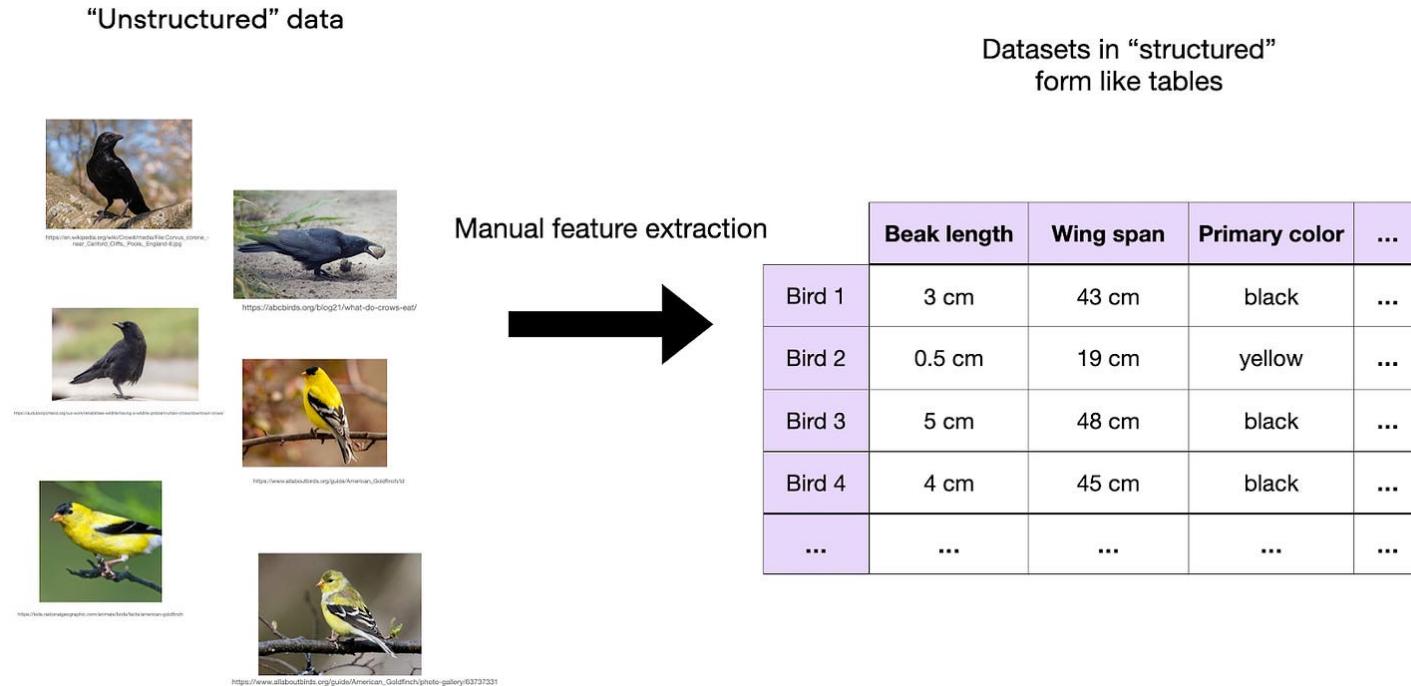


Image source: <https://magazine.sebastianraschka.com/p/ai-and-open-source-in-2023>

---

\*sometimes referred to as feature extraction

# Learning effective representations is vital

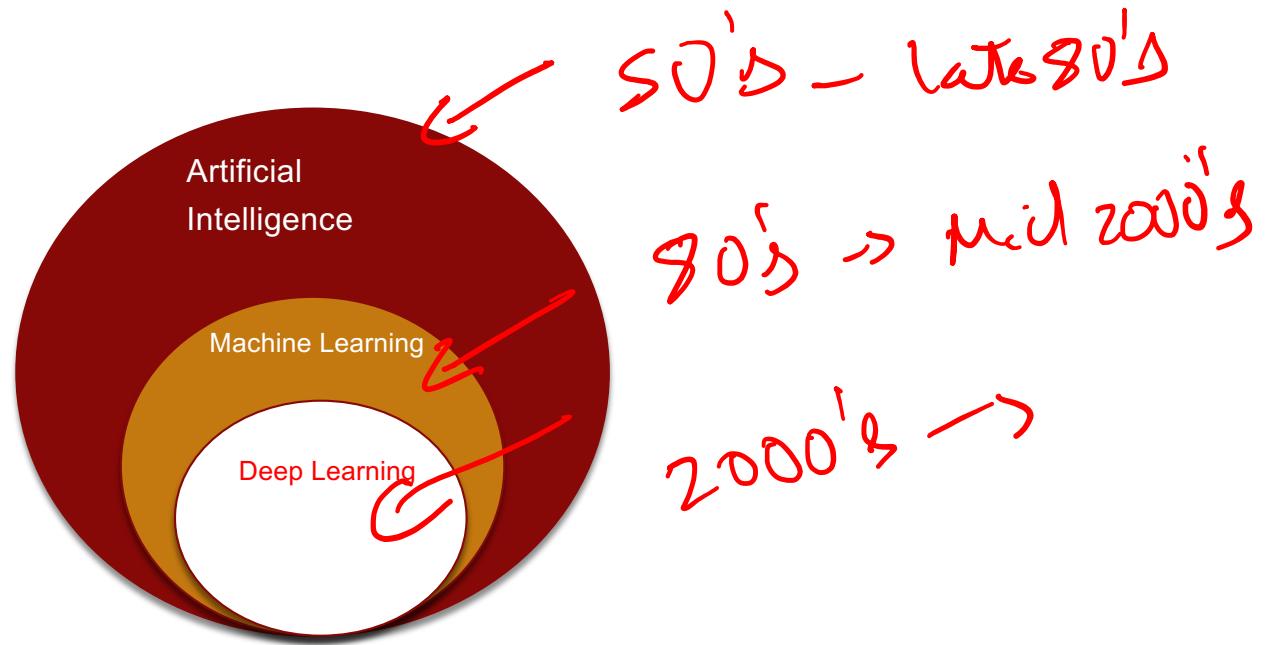


- The raw data has to somehow be transformed into a different representation
- Historically, researchers **manually** developed these representations and then fed them to traditional machine learning algorithms (often just linear/logistic regression!)
- But this required massive human effort and thus sharply limited the reach of Machine Learning



But developing good representations (before ML could be used) required massive human effort and this “human bottleneck” sharply limited the reach of Machine Learning

To address this problem, a different approach was developed – Deep Learning

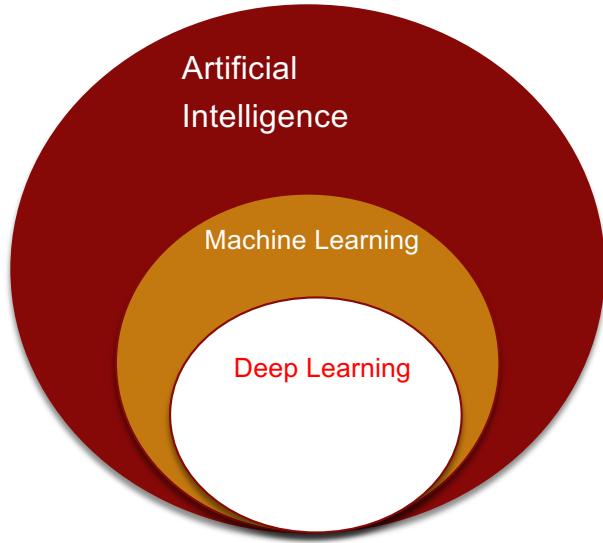


# What can Deep Learning do that traditional Machine Learning can't?

---

- It can automatically extract smart representations from raw, unstructured data.
- We can simply feed these smart representations to traditional models like linear regression or logistic regression and achieve amazing performance

# Deep Learning can handle unstructured input data without upfront manual preprocessing!!



Structured *and* **Unstructured** data → Deep Learning → Prediction

v	4.00	4.01	4.02	4.03	v	v
0	367	4.51	682	3.51	1	0
0	162	4.93	712	33.67	1	0
0	103	4.91	667	4.74	0	0
0	125	5.17	722	50.81	0	0
1	134	4.65	667	3.84	0	1
1	131	4.78	722	24.22	0	0
0	87	4.95	682	69.91	1	0
0	84	4.43	707	5.63	1	0
0	360	4.53	677	13.85	2	1
0	254	5.14	662	5.12	2	0
0	316	4.75	767	6.07	0	0
0	93	5	747	3.02	0	0



*Four score and seven years ago our fathers brought forth, upon this continent, ...*



This demolishes the “human bottleneck” for using  
Machine Learning with unstructured data

---

# Deep Learning

The breakthrough came from the confluence of three forces ...

- New algorithmic ideas
- Unprecedented amounts of data (due to the digitization of everything)
- Compute power (from the use of powerful Graphics Processing Units (GPUs))

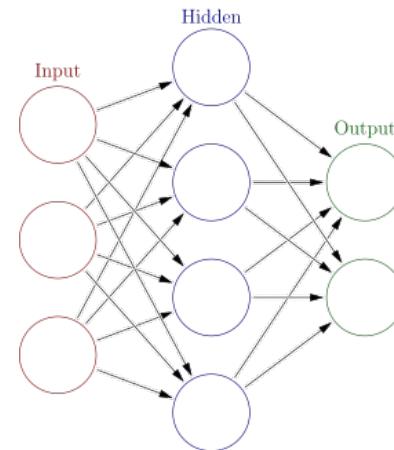
# Deep Learning

The breakthrough came from the confluence of three forces ...

- New algorithmic ideas
- Unprecedented amounts of data (due to the digitization of everything)
- Compute power (from the use of powerful Graphics Processing Units (GPUs))



... applied to an old idea:  
**Neural Networks**



[https://en.wikipedia.org/wiki/Artificial\\_neural\\_network](https://en.wikipedia.org/wiki/Artificial_neural_network)



What is the *immediate* application of Deep Learning?



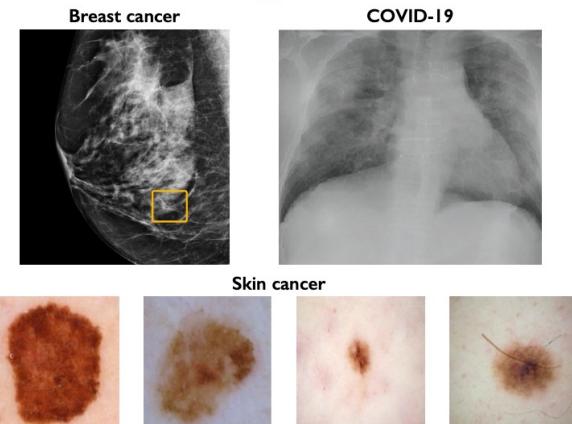
*Every “sensor” can be given the ability to detect, recognize and classify what it is sensing*

---

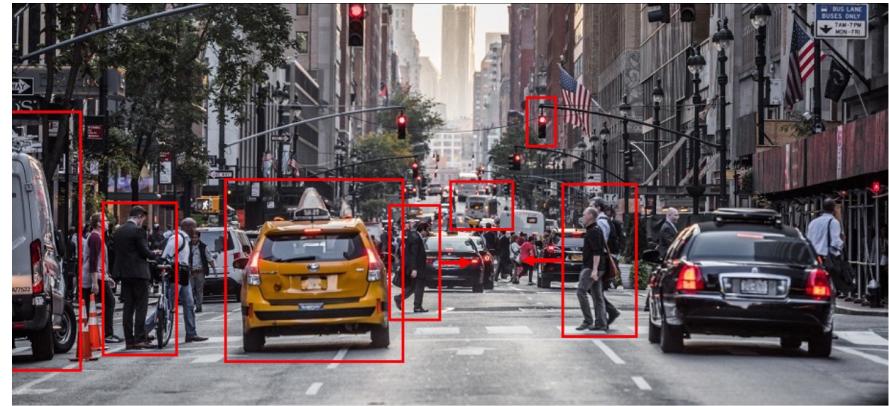
# Examples

## Use Face ID on your iPhone or iPad Pro

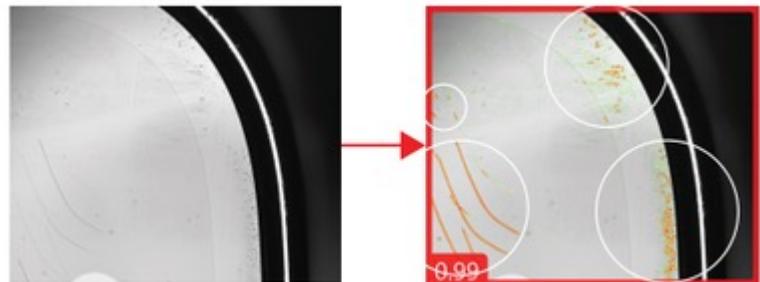
Face ID lets you securely unlock your iPhone or iPad, authenticate purchases, sign in to apps, and more — all with just a glance.



6S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com) [@MITDeepLearning](https://twitter.com/MITDeepLearning)



6S191 Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com) [@MITDeepLearning](https://twitter.com/MITDeepLearning)

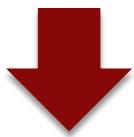


Smartphone glass defect

<https://dwfritz.com/smart-cosmetic-defect-detection-increases-productivity/>

Detected by system

*Every “sensor” can be given the ability to detect, recognize and classify what it is sensing*



*You can create dramatically better products and services by “attaching” DL to sensors*

---

# (Spotted last week!) Binoculars + DL = Smart Binoculars





+ DL =

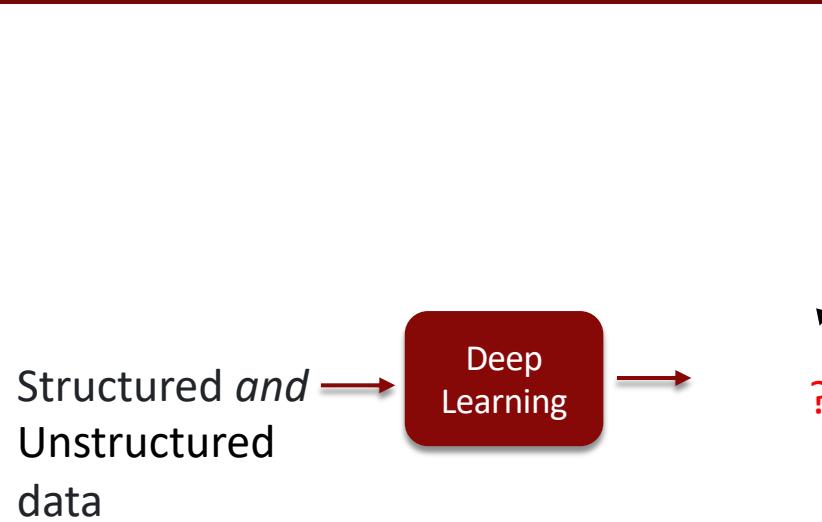


# Self-Driving!



<https://www.cnn.com/2021/11/03/cars/tesla-full-self-driving-fsd/index.html>

# Now, let's turn our attention to the output



# With Deep Learning, we could predict/generate structured outputs easily

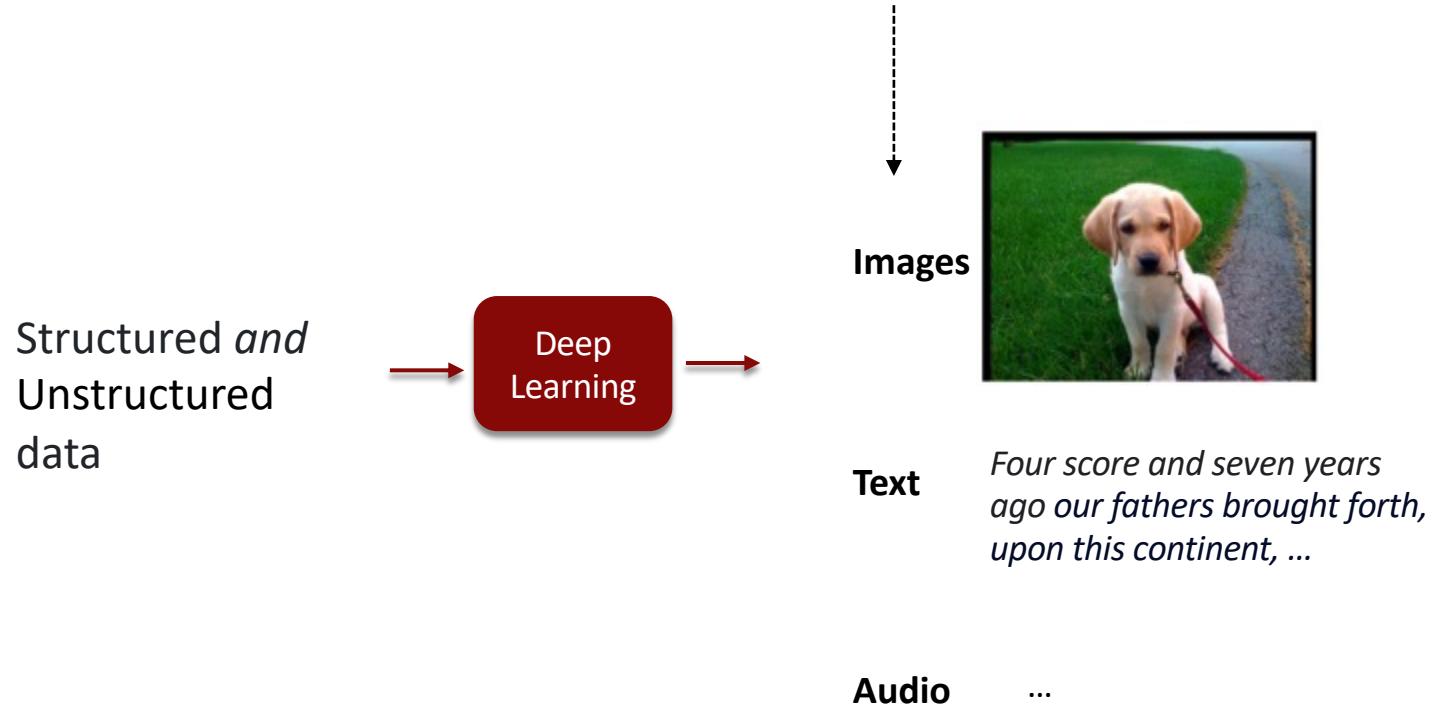
Structured and Unstructured data



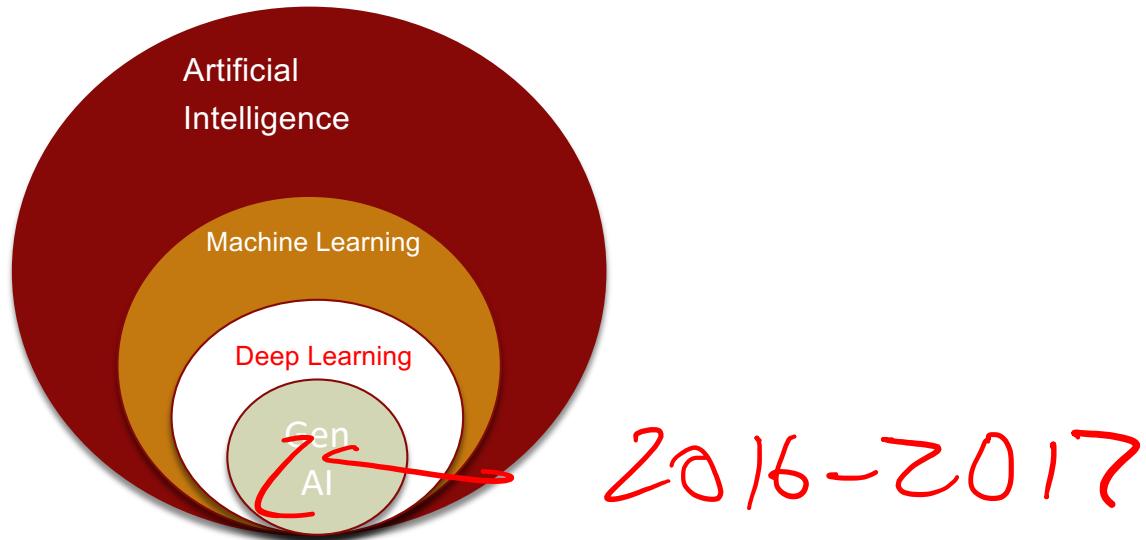
Examples:

- A single number
  - The probability that a borrower will repay a loan
  - The demand for a product next week
- A few numbers
  - The 4 probabilities that an image contains a chair, stool, table or sofa
  - The two GPS coordinates of a taxi

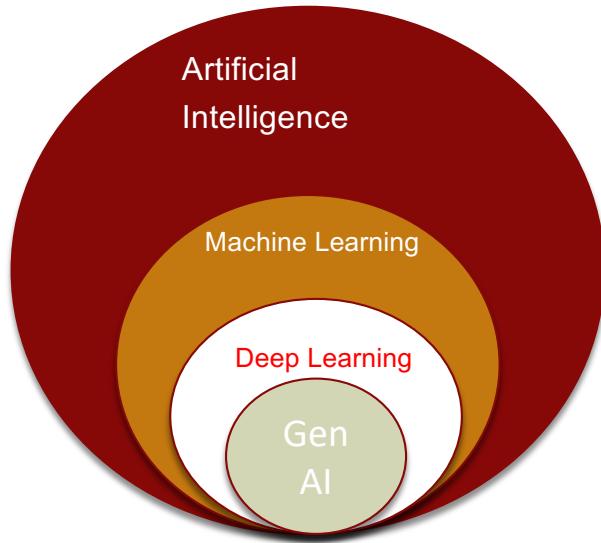
# But we couldn't generate unstructured output very well!



# And then Generative AI happened



# Gen AI can produce unstructured outputs

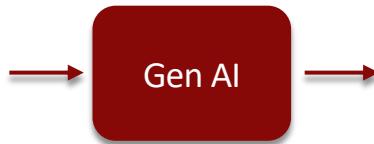


Unstructured input → Gen AI → Unstructured output

"generate a picture of a cute labrador retriever puppy"



# Image-to-text

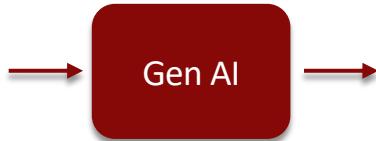


*two cats laying on  
a pink blanket  
with remotes*

<https://huggingface.co/spaces/nielsr/comparing-captioning-models>

# Text-to-image

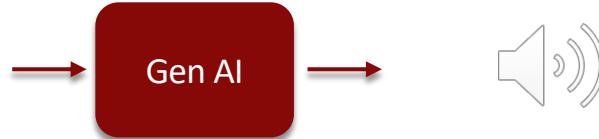
temple in ruins, forest, stairs, columns, cinematic, detailed, atmospheric, epic, concept art, Matte painting, background, mist, photo-realistic, concept art, volumetric light, cinematic epic + rule of thirds octane render, 8k, corona render, movie concept art, octane render, cinematic, trending on artstation, movie concept art, cinematic composition , ultra-detailed, realistic , hyper-realistic , volumetric lighting, 8k –ar 2:3 –test –uplight



<https://mpost.io/best-100-stable-diffusion-prompts-the-most-beautiful-ai-text-to-image-prompts/>

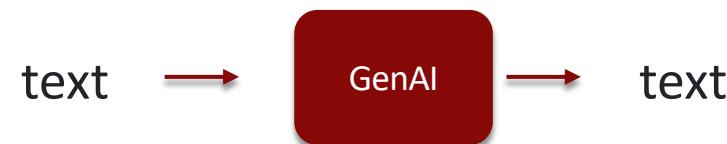
# Text-to-music

The main soundtrack of an arcade game. It is fast-paced and upbeat, with a catchy electric guitar riff. The music is repetitive and easy to remember, but with unexpected sounds, like cymbal crashes or drum rolls.



<https://google-research.github.io/seanet/musiclm/examples/>

# Text-to-text



# Multi-modal



# A fun multi-modal example

[text,

image]

→ ChatGPT

→



It's Wednesday at 4 pm. Can I park at this spot right now? Tell me in 1 line.

# A fun multi-modal example

[text,

image]



text

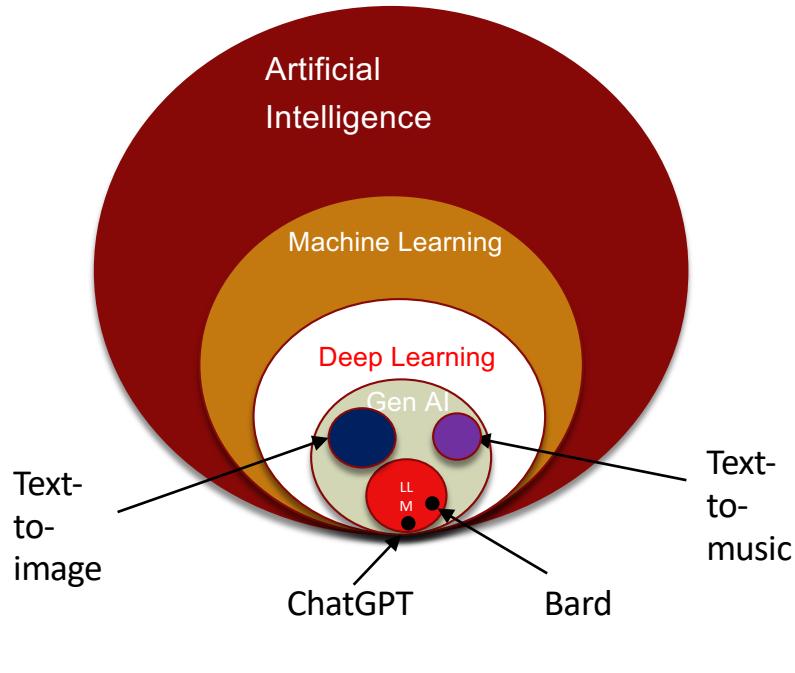


Yes, you can park for up to 1 hour starting at 4 pm.

It's Wednesday at 4 pm. Can I park at this spot right now? Tell me in 1 line.

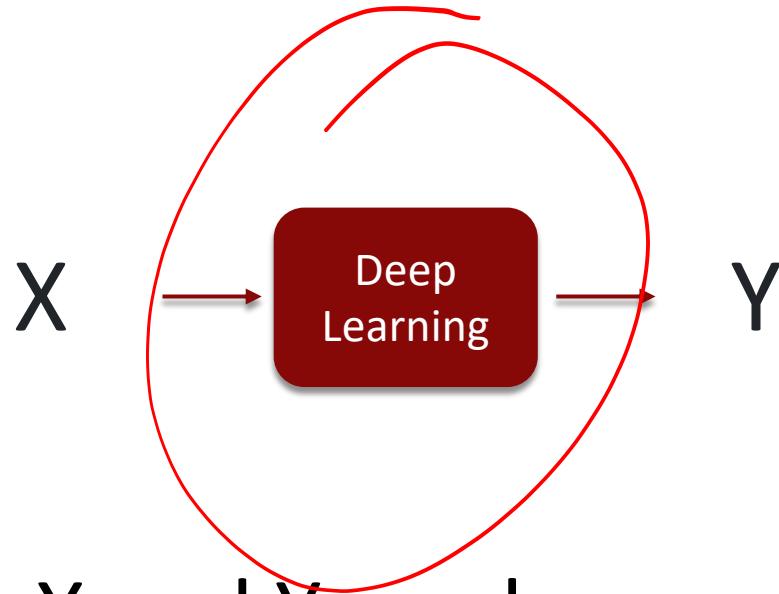
<https://twitter.com/petergyang/status/1707169696049668472>

# The landscape



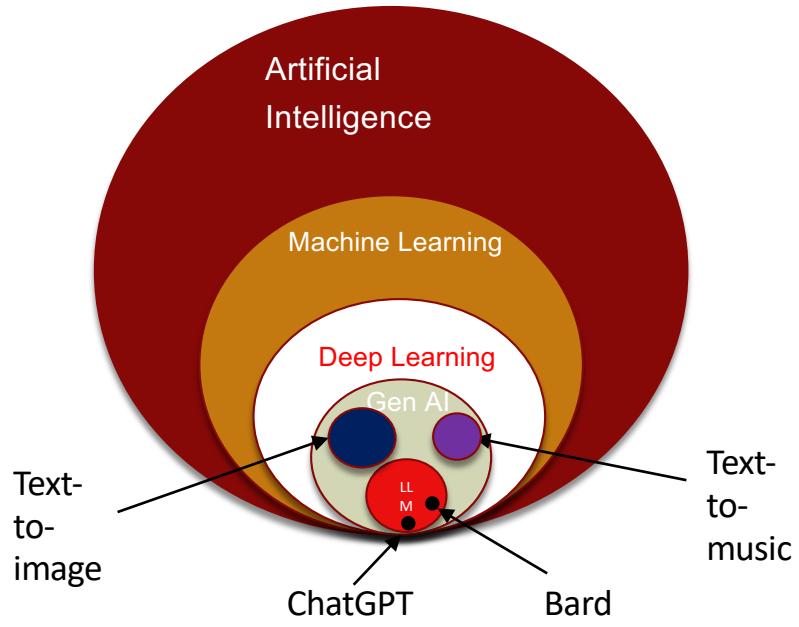
... and the  
circles inside  
GenAI are  
merging!

# Summary: X-to-Y



X and Y can be  
*anything and it can  
be multi-modal!*

# Note that ALL the AI excitement is due to the success of Deep Learning



If you understand Deep Learning,  
everything becomes possible! 😊

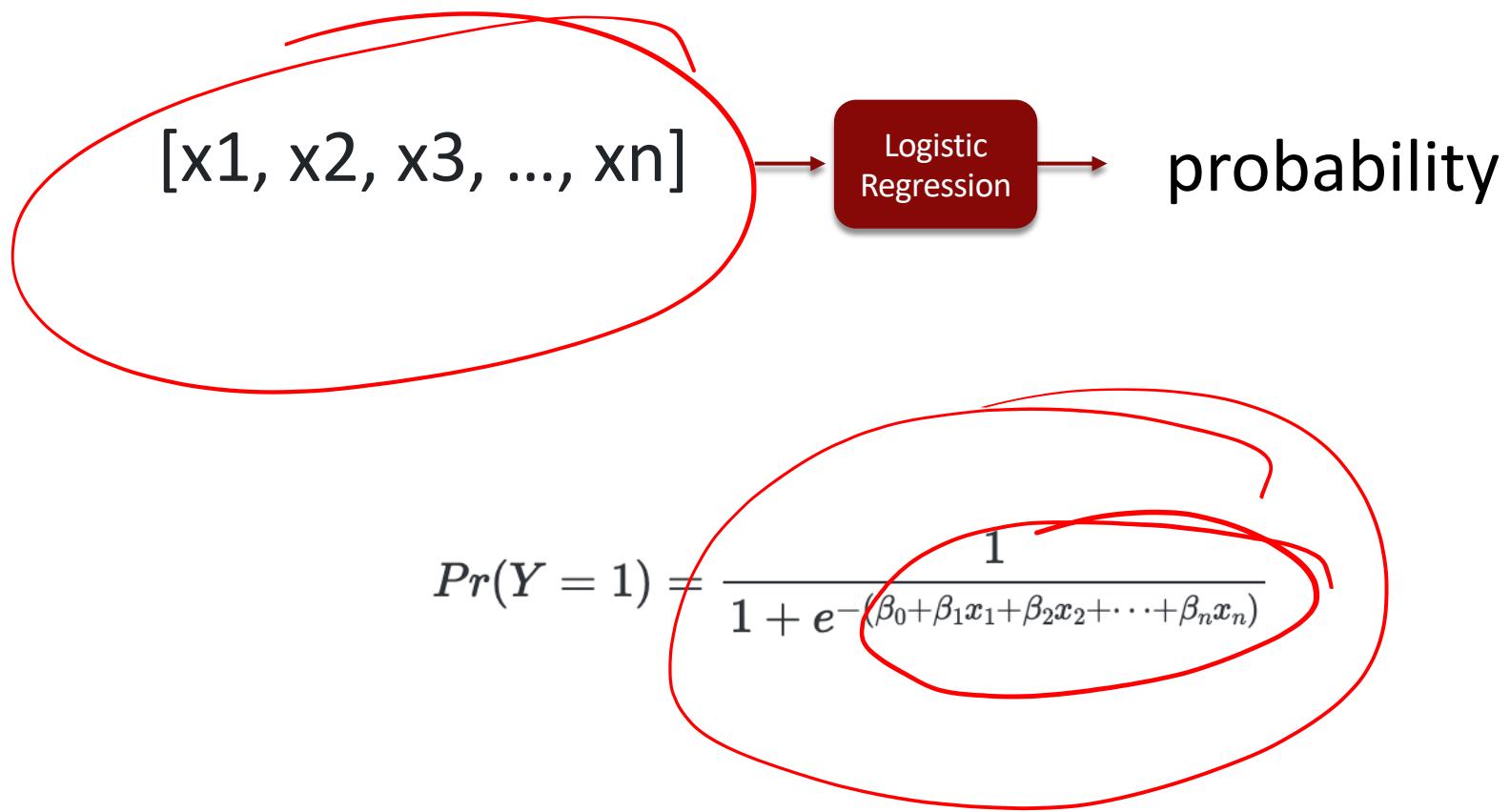
OK, let's start at the beginning.  
What's a Neural Network?

# We will work with this simple classification example

- Given independent variables ...
  - GPA
  - Experience
- ... predict who will be called for a job interview

	Interview	GPA	Experience
1	0	3.27	1.93
2	0	3.37	0.07
3	0	3.57	1.91
4	0	3.91	4.35
5	0	3.20	1.70
6	1	3.90	2.41
7	1	3.94	3.00
8	0	3.66	2.47
9	0	3.63	0.93
10	0	3.06	4.14
11	0	3.21	3.34
12	0	3.18	3.97
13	0	3.69	0.54
14	0	3.38	3.62
15	1	3.77	2.06
16	0	3.50	4.10
...	...	...	...
25	1	3.31	3.46
26	0	3.78	0.29
27	0	3.87	1.21
28	0	4.00	0.49
29	1	3.87	2.11

# Recall Logistic Regression

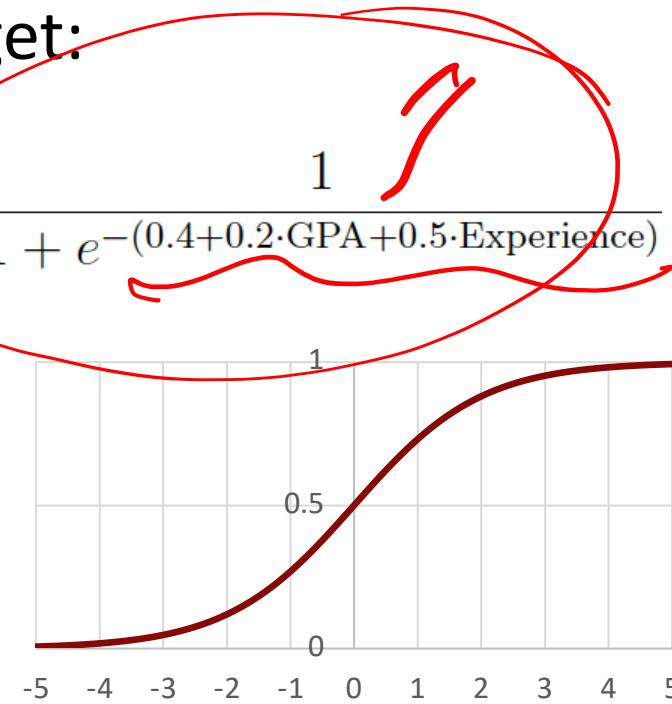


# We will work with this simple classification example

	Interview	GPA	Experience
1		0 3.27	1.93
2		0 3.37	0.07
3		0 3.57	1.91
4		0 3.91	4.35
5		0 3.20	1.70
6		1 3.90	2.41
7		1 3.94	3.00
8		0 3.66	2.47
9		0 3.63	0.93
10		0 3.06	4.14
11		0 3.21	3.34
12		0 3.18	3.97
13		0 3.69	0.54
14		0 3.38	3.62
15		1 3.77	2.06
16		0 3.50	4.10
...	...	...	...
25		1 3.31	3.46
26		0 3.78	0.29
27		0 3.87	1.21
28		0 4.00	0.49
29		1 3.87	2.11

We run this dataset through  
'glm' and get:

$$P(Y = 1) = \frac{1}{1 + e^{-(0.4 + 0.2 \cdot \text{GPA} + 0.5 \cdot \text{Experience})}}$$



We can re-write this formula as a **network** with input data flowing through two functions that have been connected

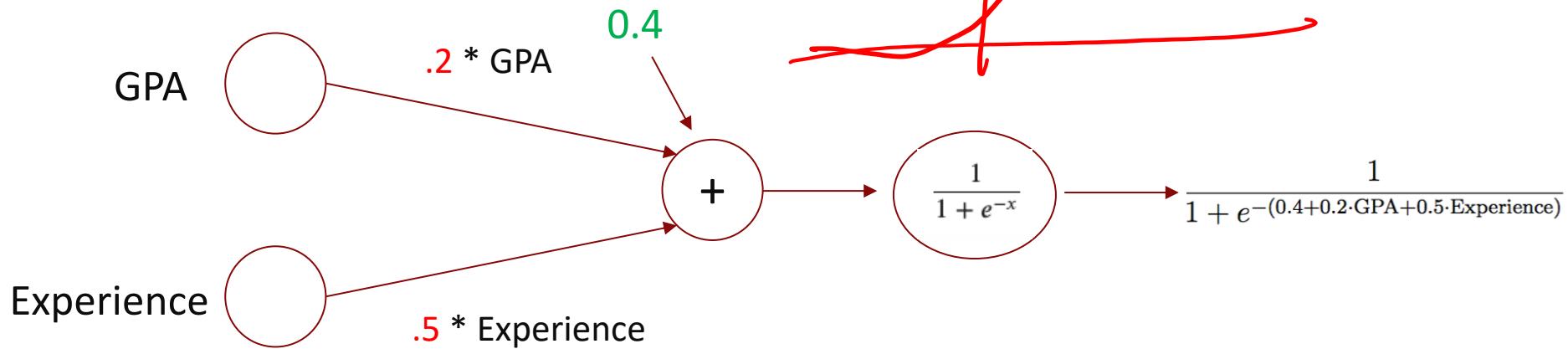
Model equation:

$$P(Y = 1) = \frac{1}{1 + e^{-(0.4 + 0.2 \cdot \text{GPA} + 0.5 \cdot \text{Experience})}}$$

We can re-write this formula as a **network** with input data flowing through two functions that have been connected

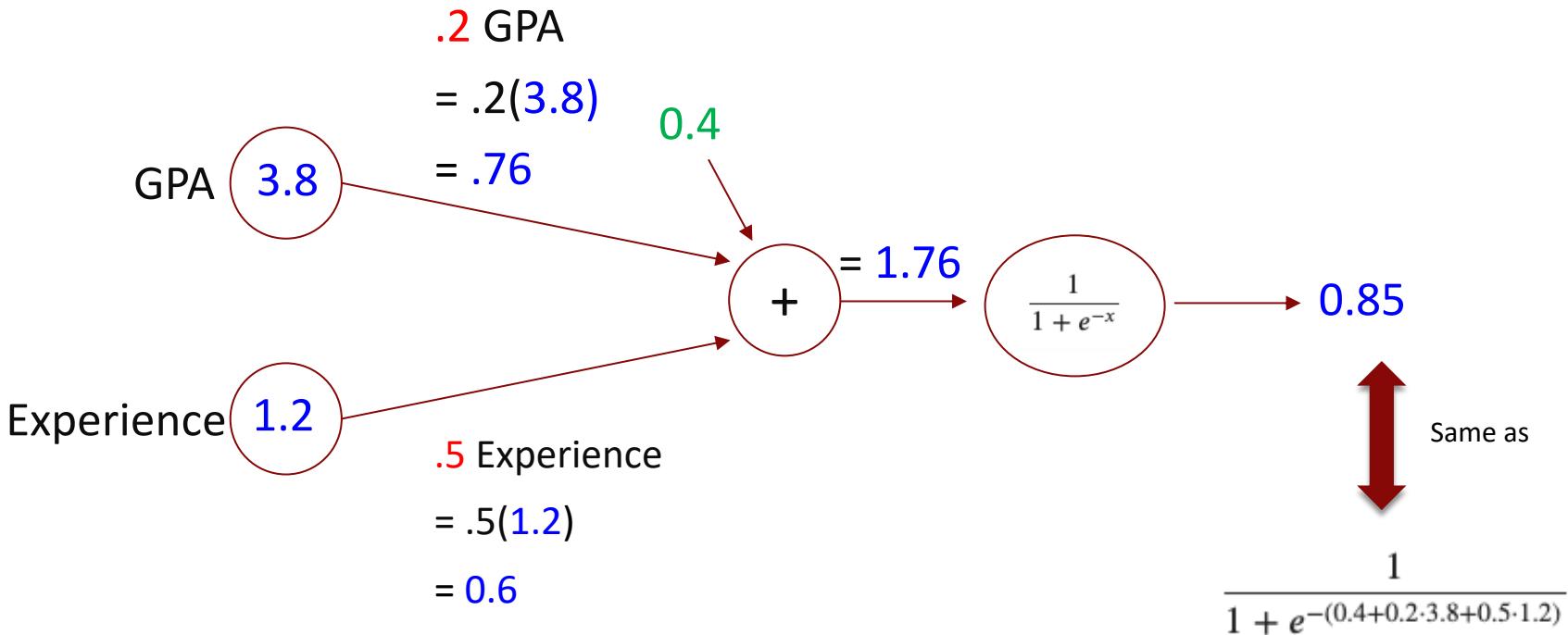
Model equation:

$$P(Y = 1) = \frac{1}{1 + e^{-(0.4 + 0.2 \cdot \text{GPA} + 0.5 \cdot \text{Experience})}}$$



# Let's make a prediction with this “network”

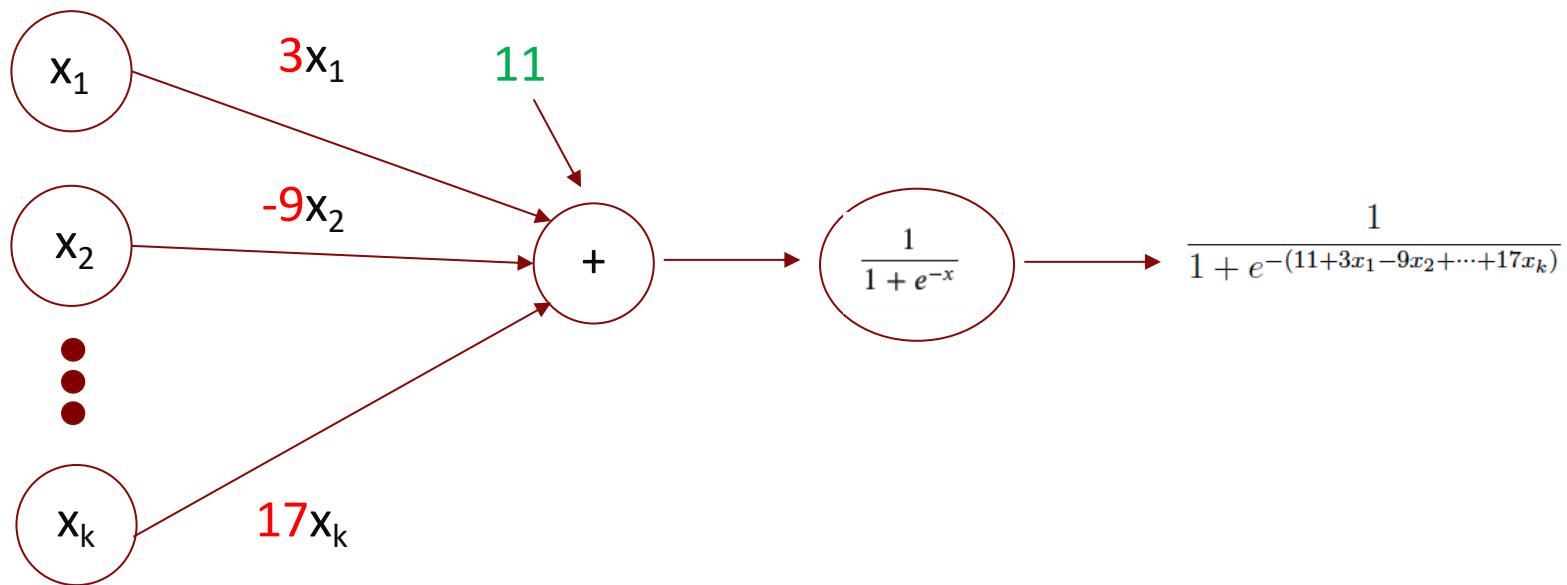
Consider a job applicant with a GPA of 3.8 and 1.2 years of experience.



# The general logistic regression model viewed through a network lens

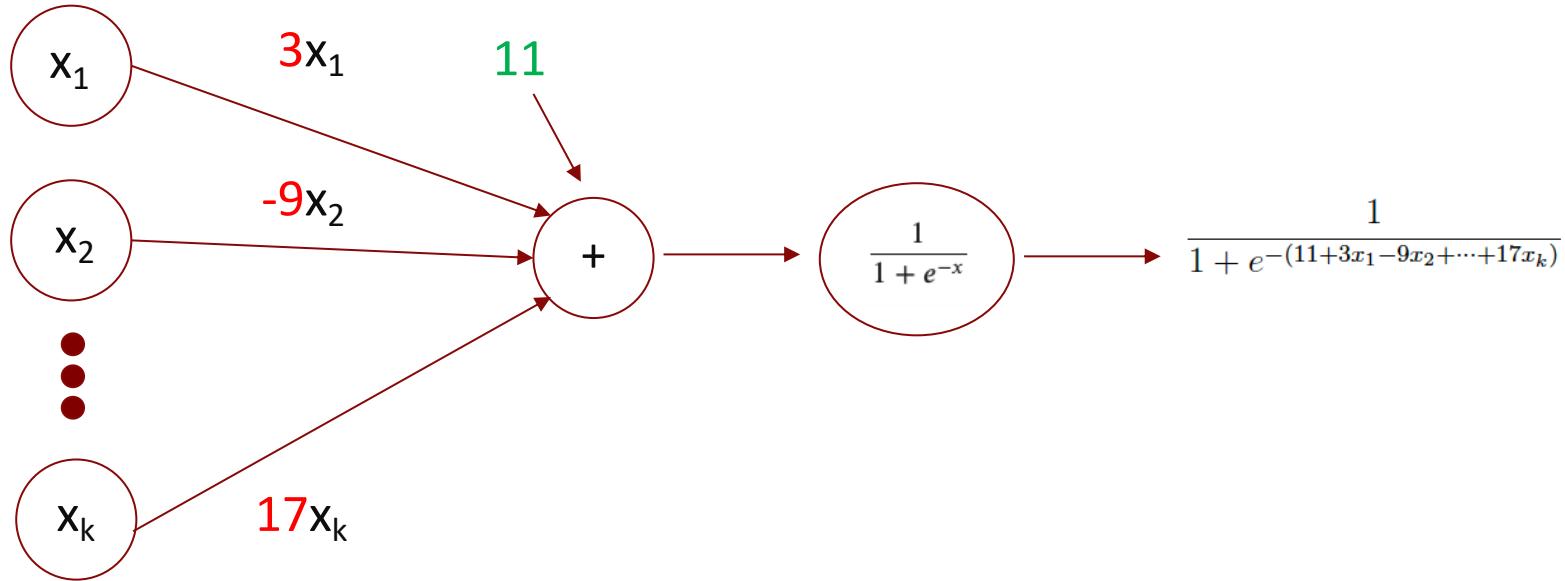
Model equation:

$$P(Y = 1) = \frac{1}{1 + e^{-(11 + 3x_1 - 9x_2 + \dots + 17x_k)}}$$



Notice how the data flows through the network from left to right

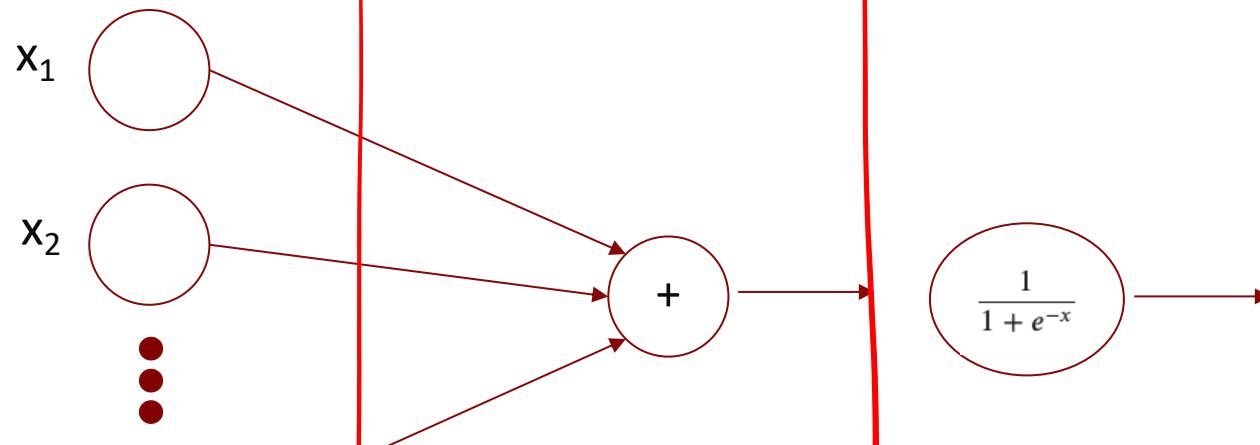
# Terminology



- Multipliers on values from each node = coefficients = **weights**
- Intercept = **bias**

What's the advantage of viewing through a network “lens”?

Recall the notion of learning smart representations of the input data

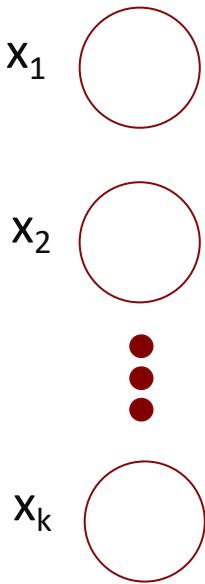


\*I am not showing the weights and biases to avoid clutter

To learn smart representations, we would like to transform the inputs one or more times before we do the prediction



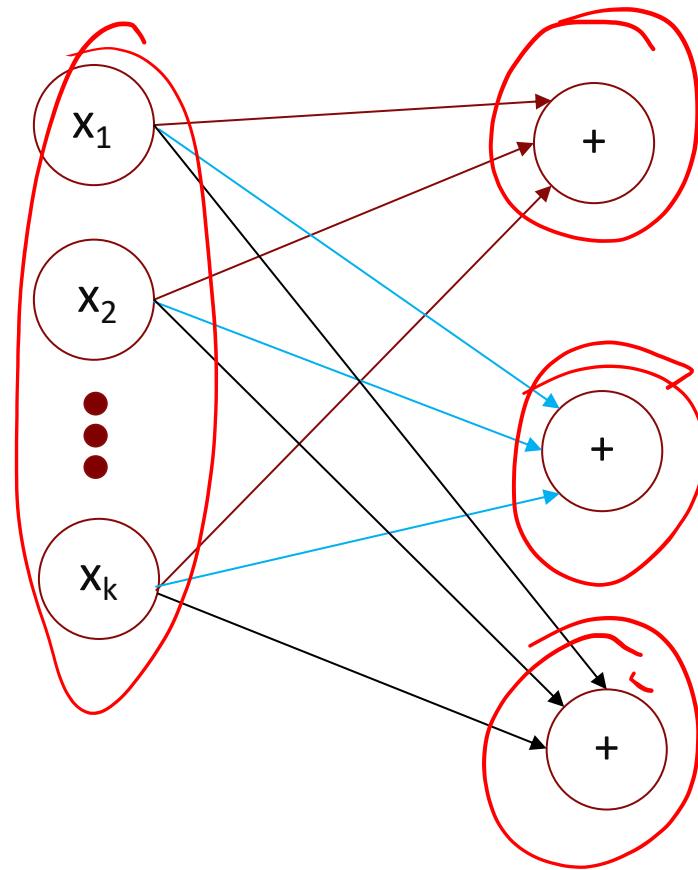
To learn smart representations, we would like to transform the inputs one or more times before we do the prediction



We can insert a stack of linear functions here



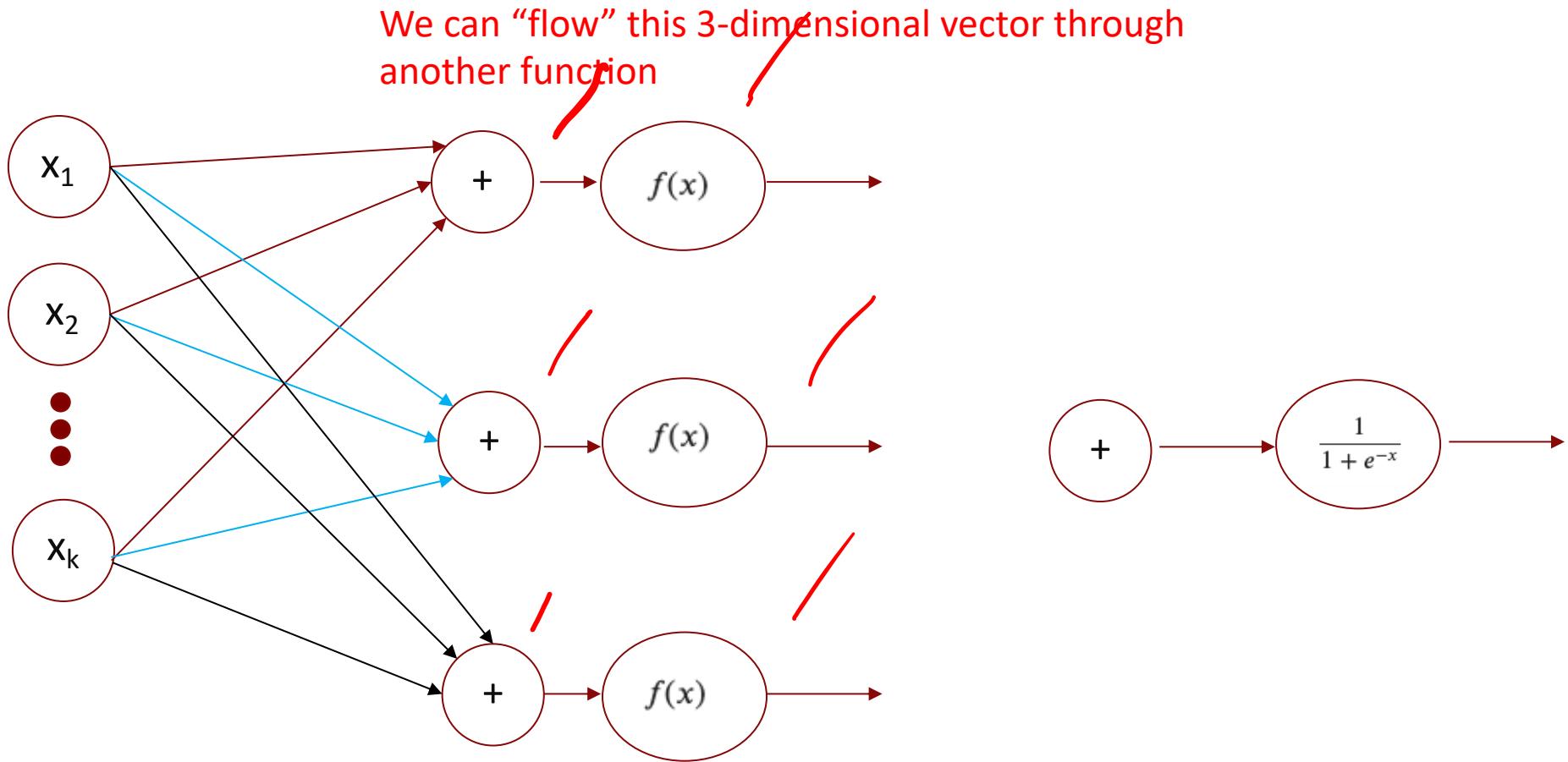
To learn smart representations, we would like to transform the inputs one or more times before we do the prediction



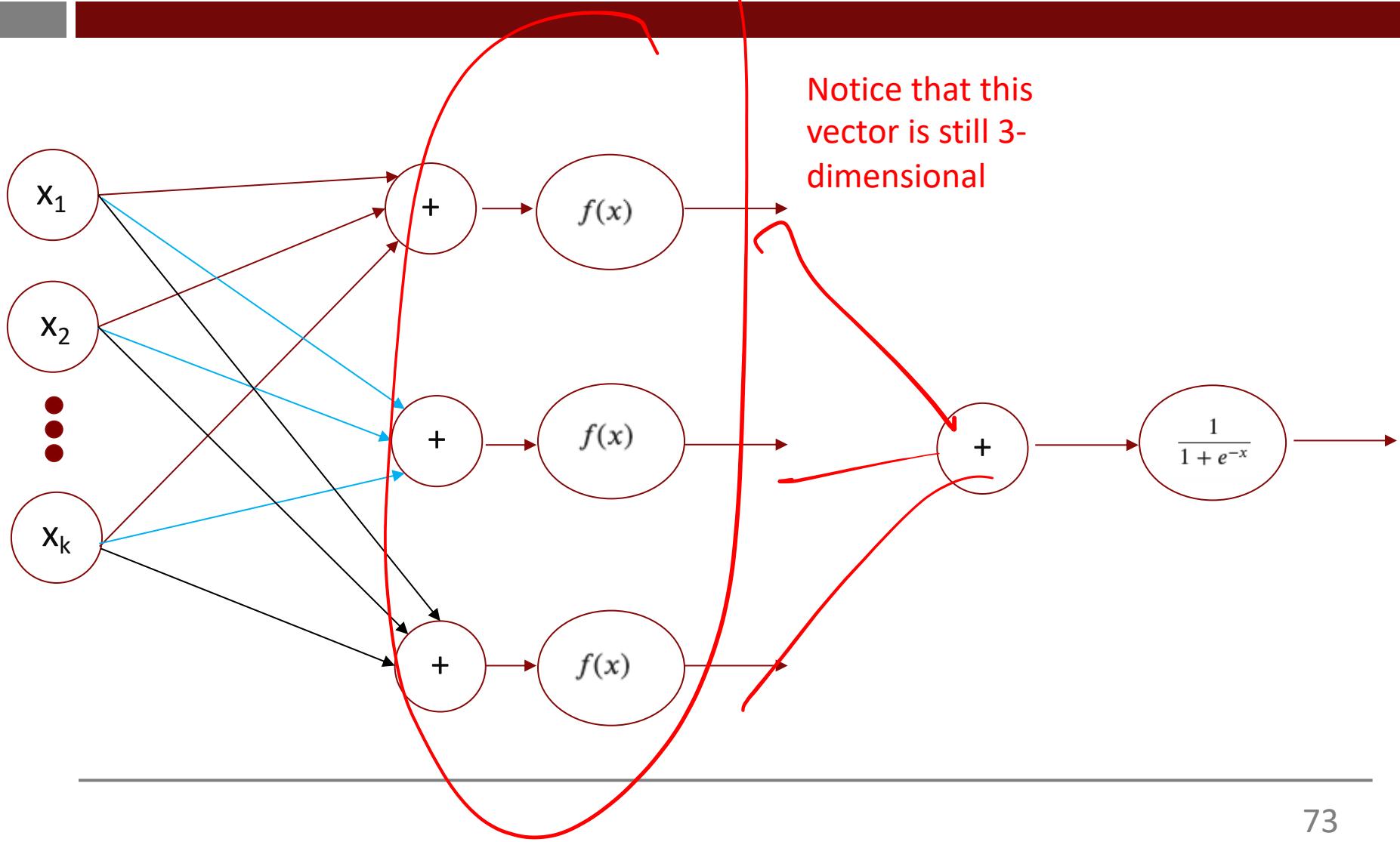
Notice: we have  
linearly  
transformed a  $k$ -  
dimensional  
input to a 3-  
dimensional  
vector

The diagram shows the final steps of the transformation. It consists of two sequential operations. The first operation is a circle containing a '+' sign, followed by a horizontal arrow pointing to the right. The second operation is a circle containing the mathematical expression  $\frac{1}{1 + e^{-x}}$ , followed by a horizontal arrow pointing to the right. This sequence represents the addition of bias terms and the application of a sigmoid activation function to produce the final output.

To learn smart representations, we would like to transform the inputs one or more times before we do the prediction

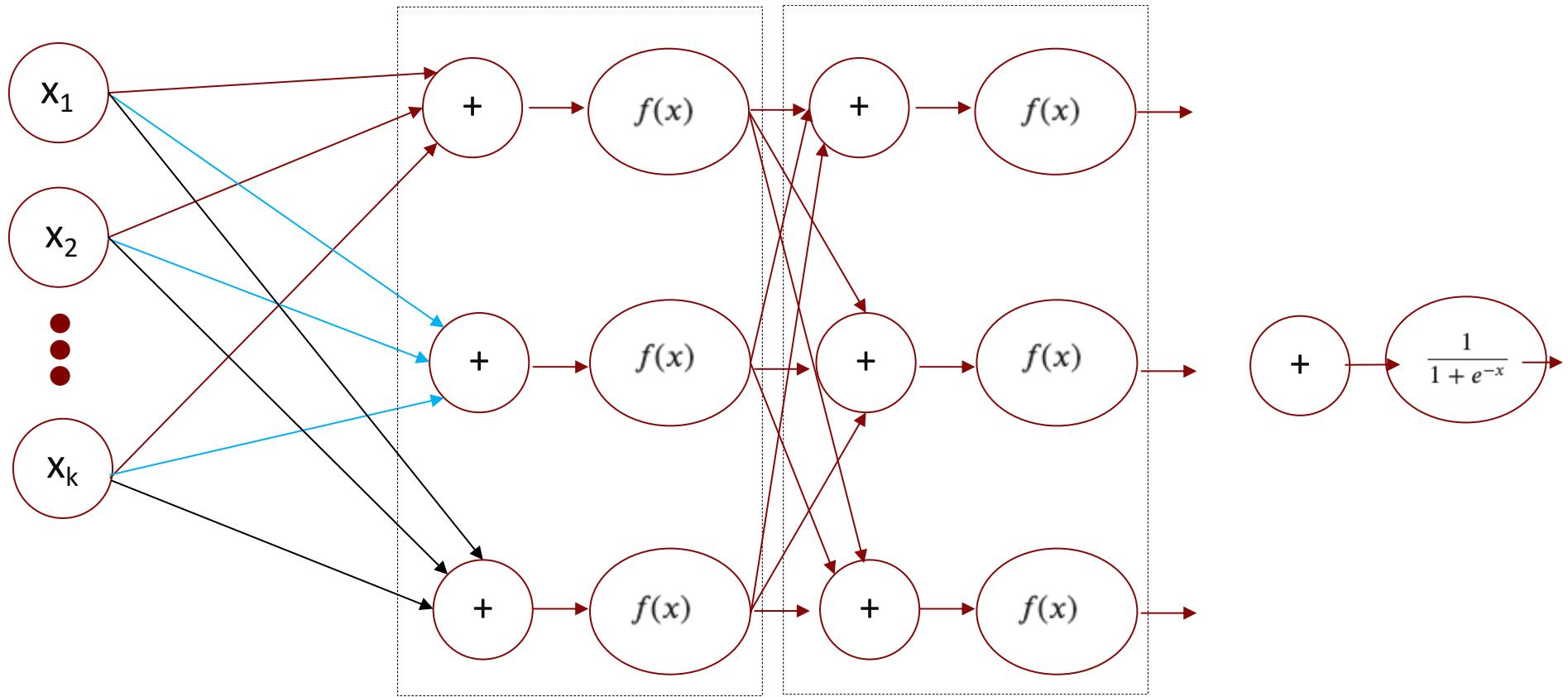


To learn smart representations, we would like to transform the inputs one or more times before we do the prediction

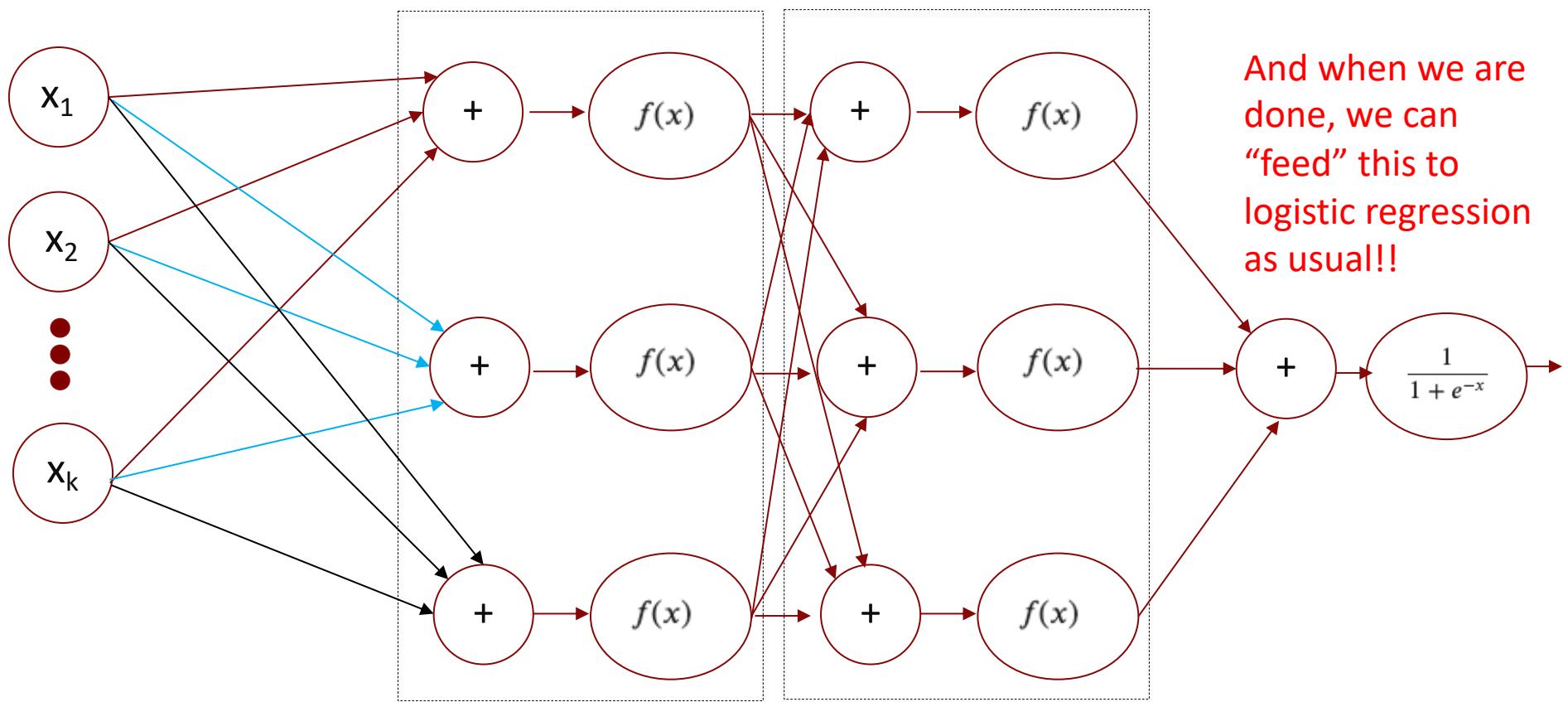


To learn smart representations, we would like to transform the inputs one or more times before we do the prediction

We can do this repeatedly

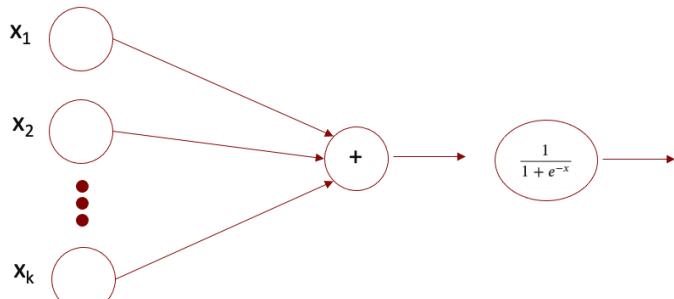


To learn smart representations, we would like to transform the inputs one or more times before we do the prediction

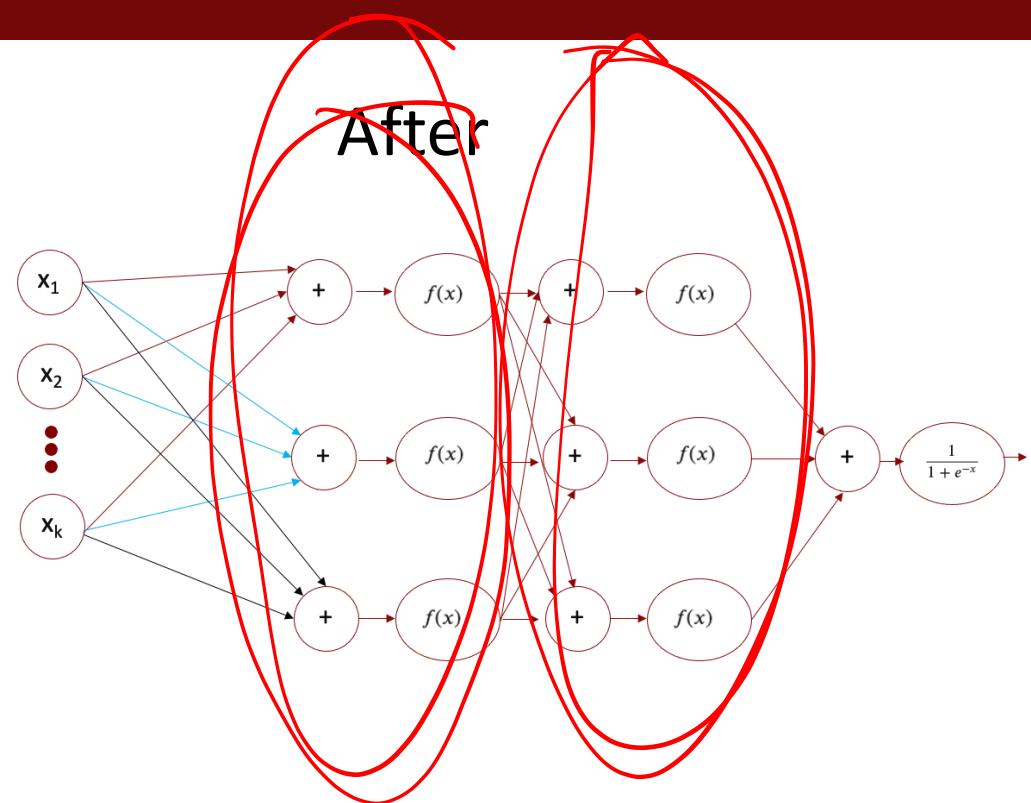


Key Takeaway: Instead of feeding the "raw" input to logistic regression, we feed a repeatedly transformed input

Before

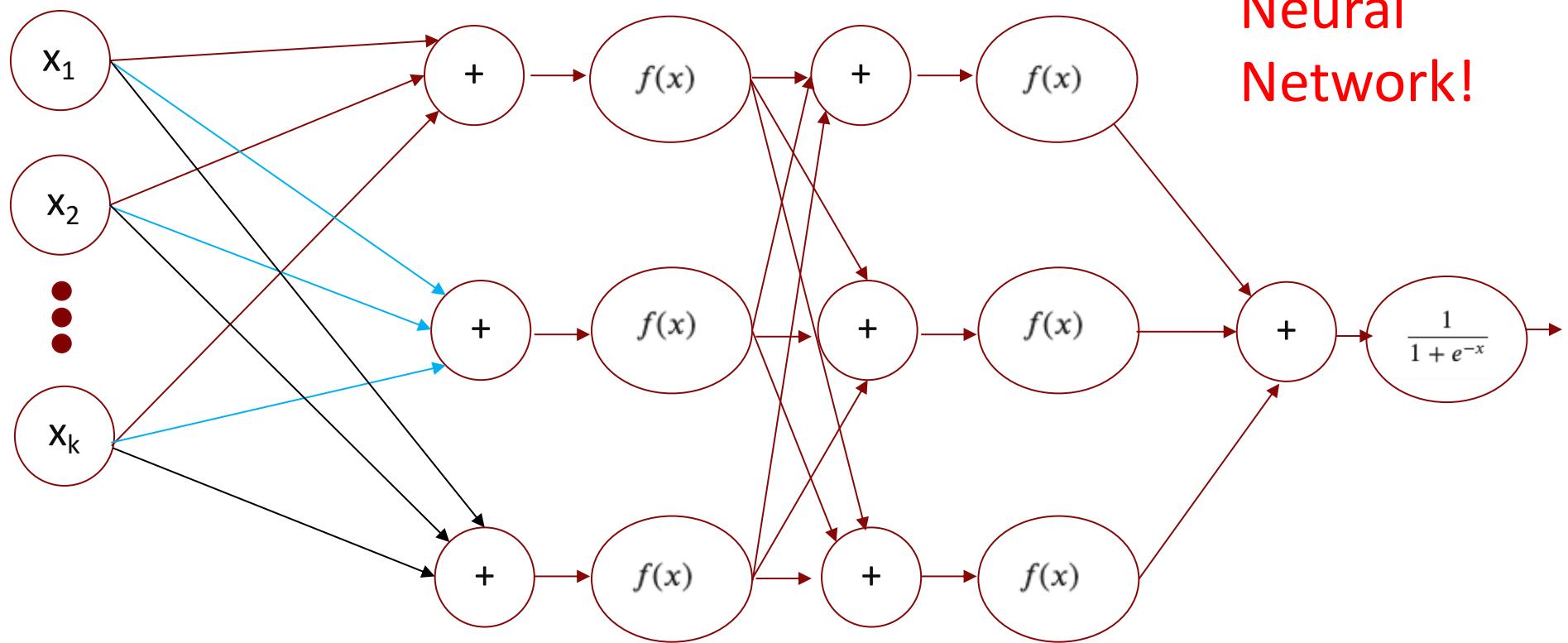


After



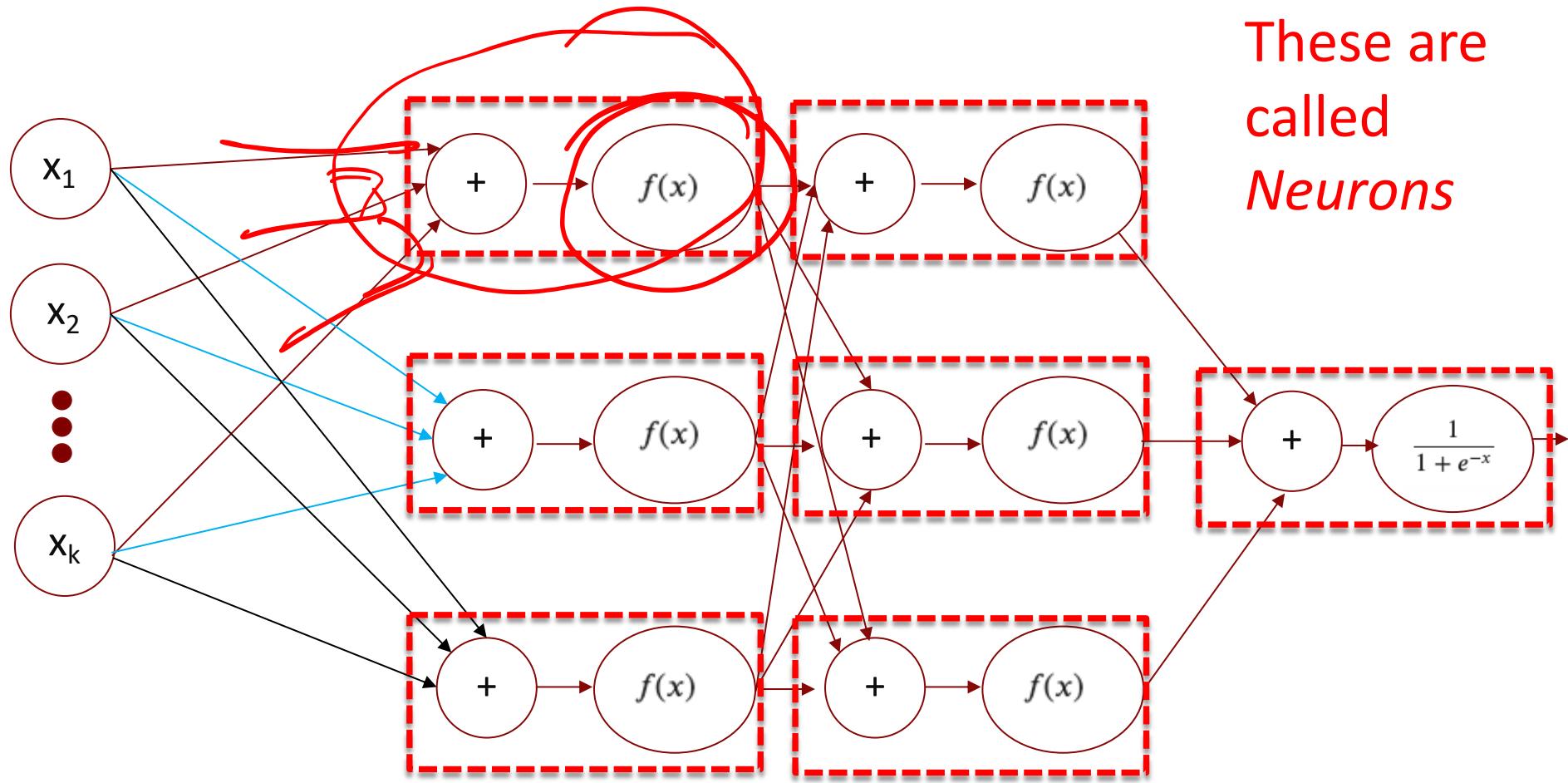
# Terminology

This is a  
Neural  
Network!

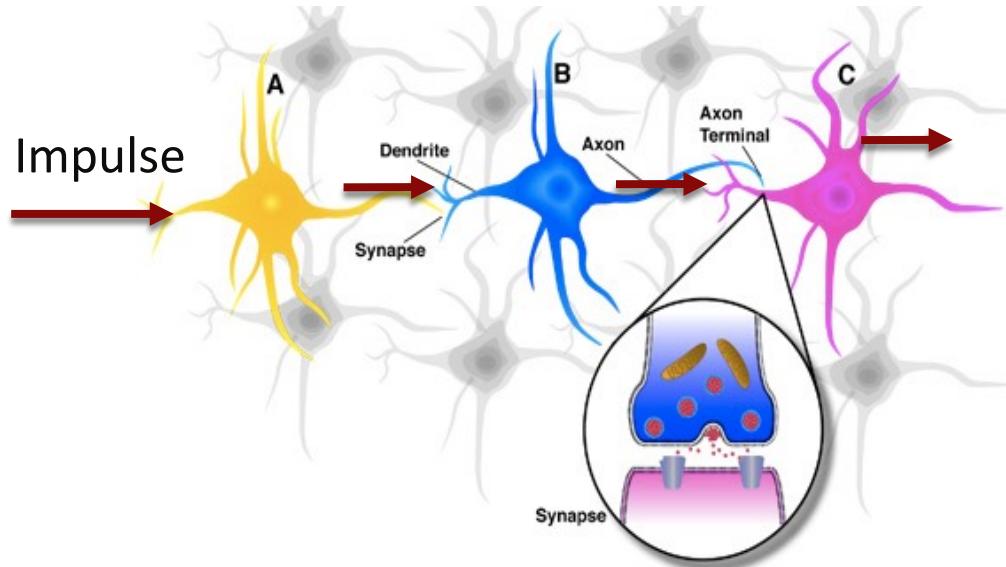


# Terminology

These are  
called  
*Neurons*

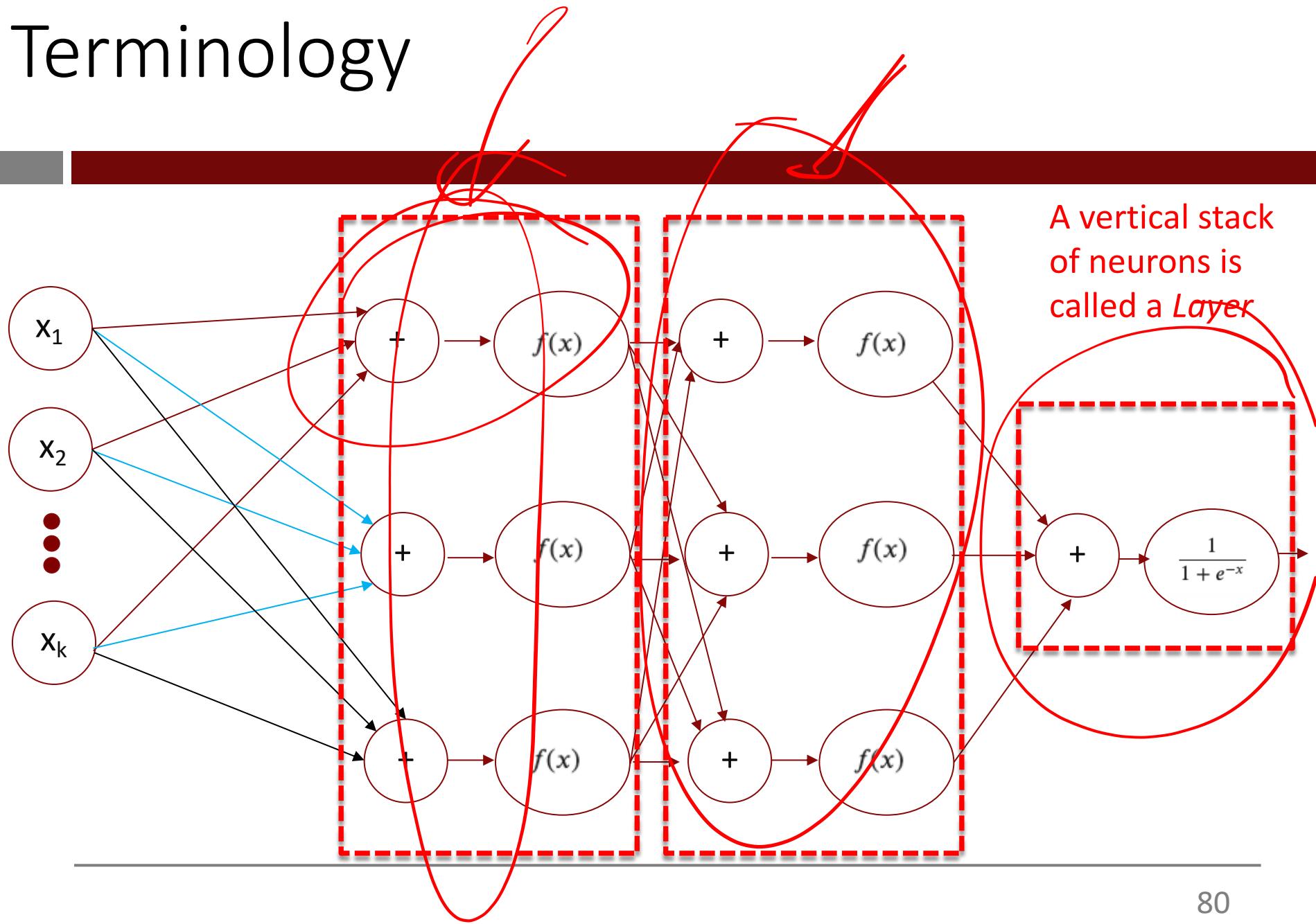


# Aside: The “neural” connection



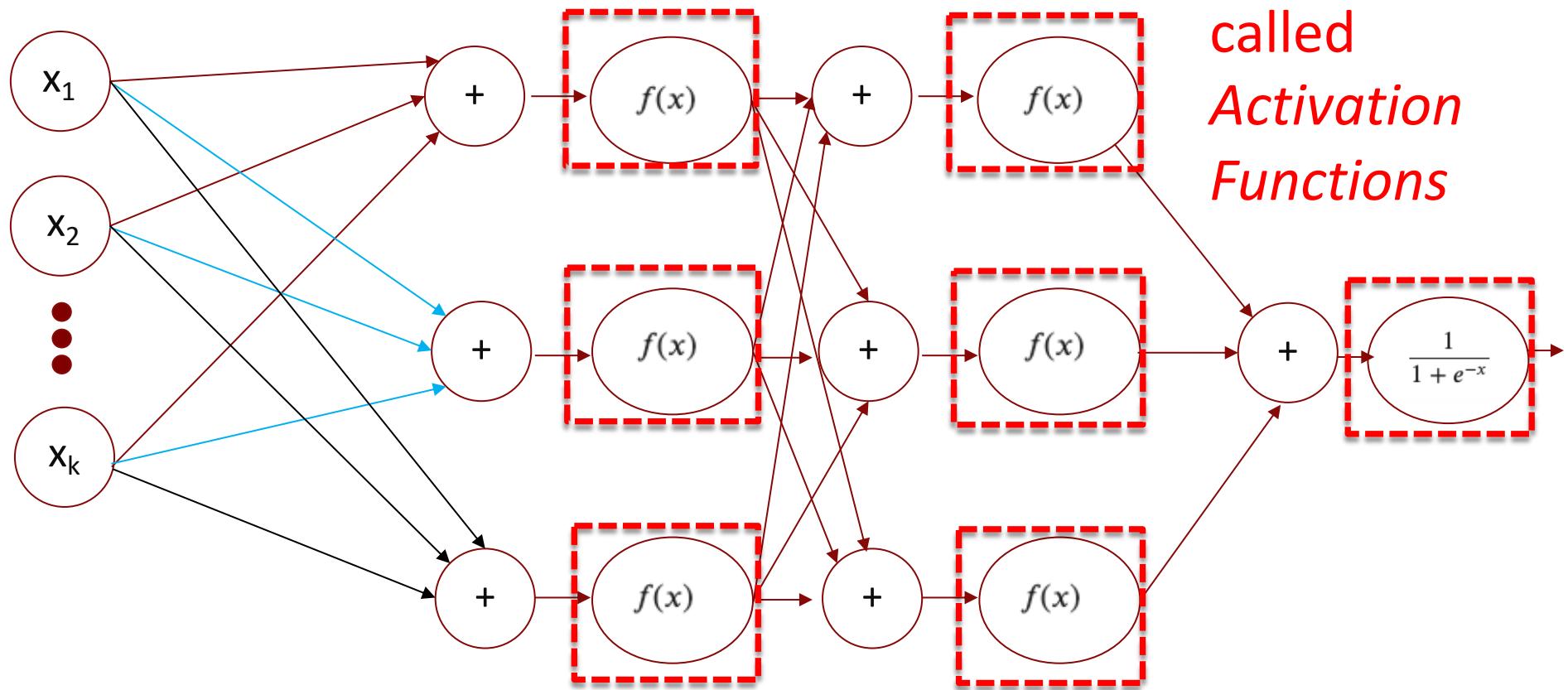
In the brain, the dendrites of one neuron connect to the axons of others at synapses, forming a network

# Terminology

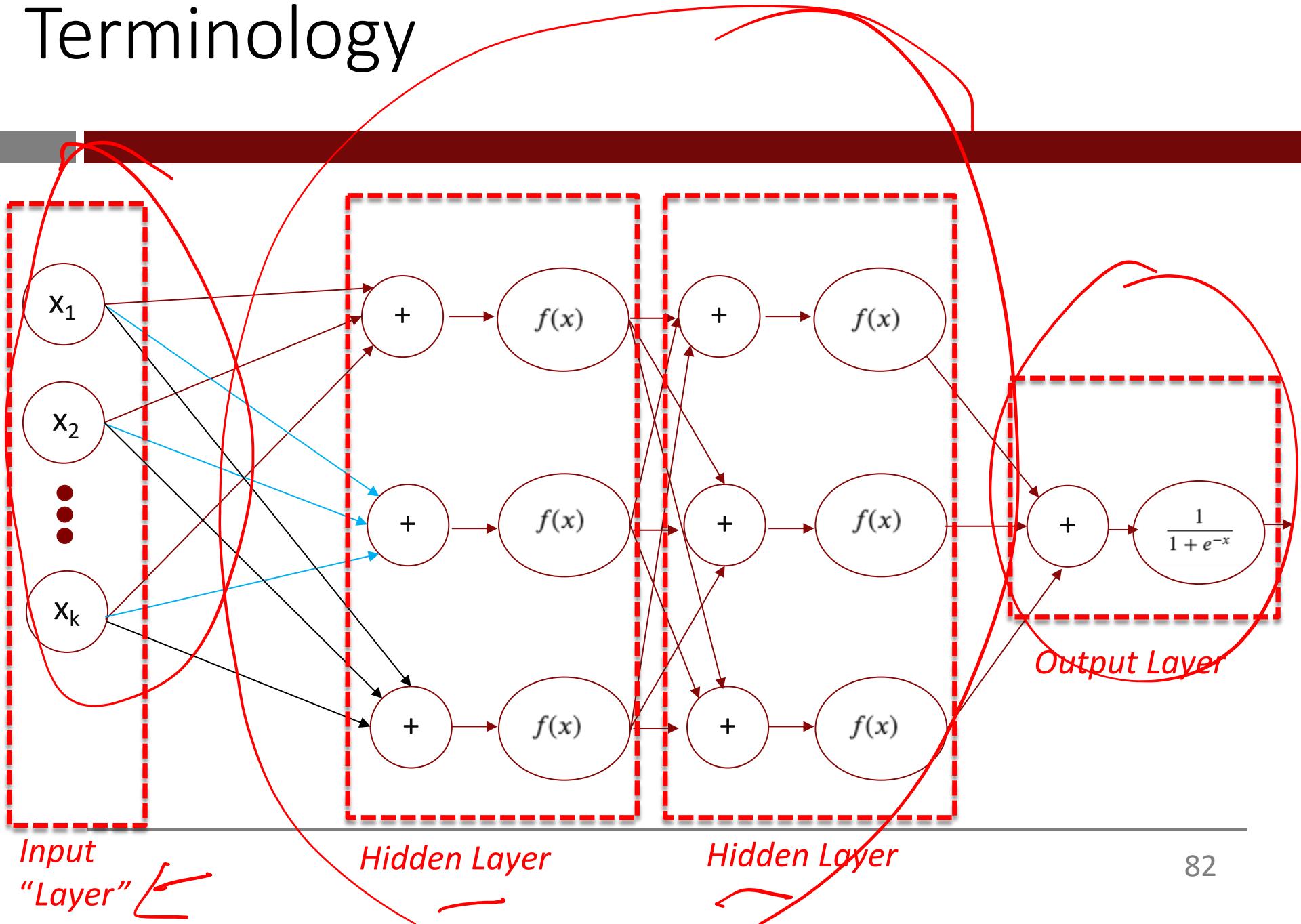


# Terminology

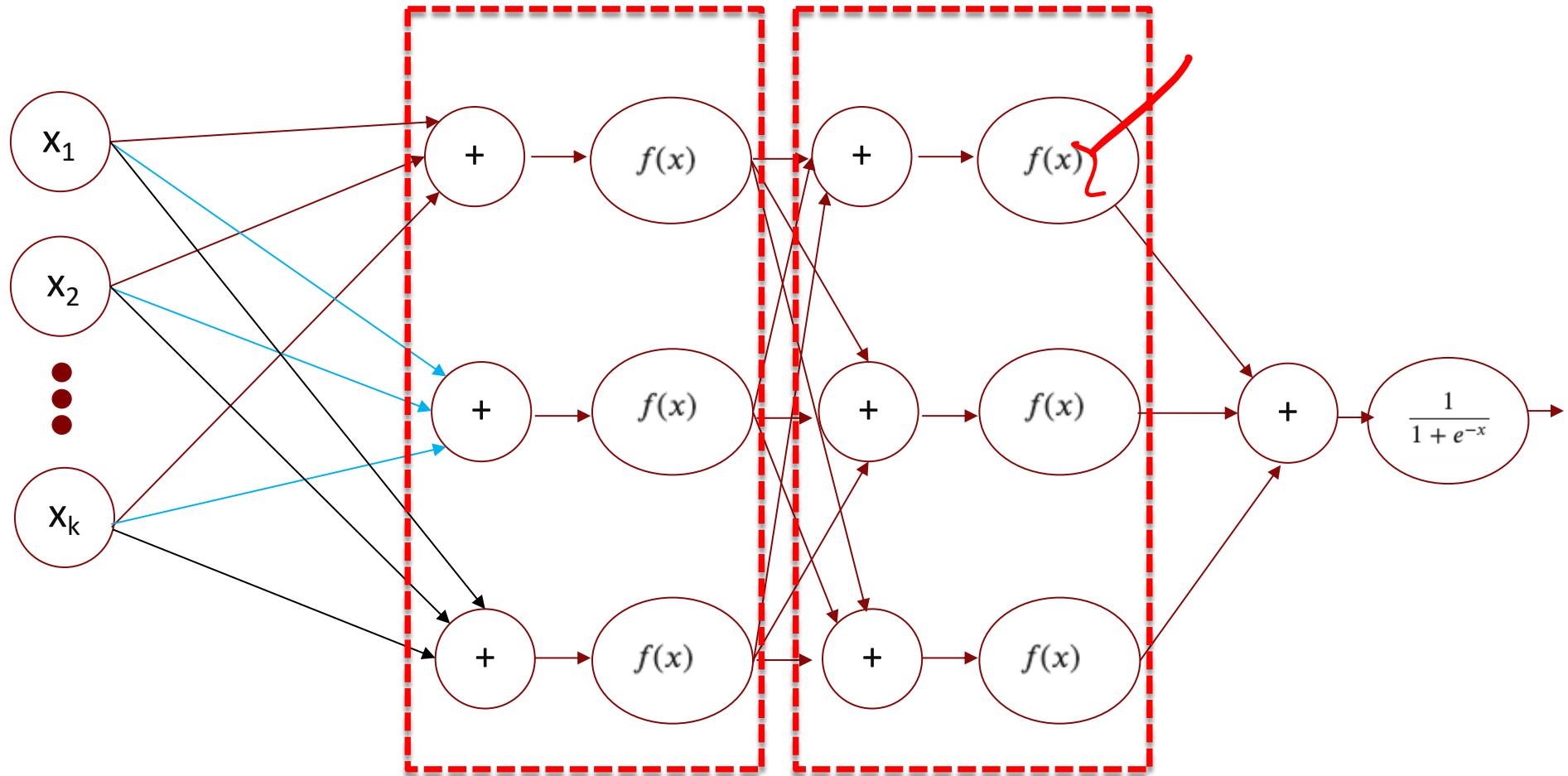
These are  
called  
*Activation  
Functions*



# Terminology

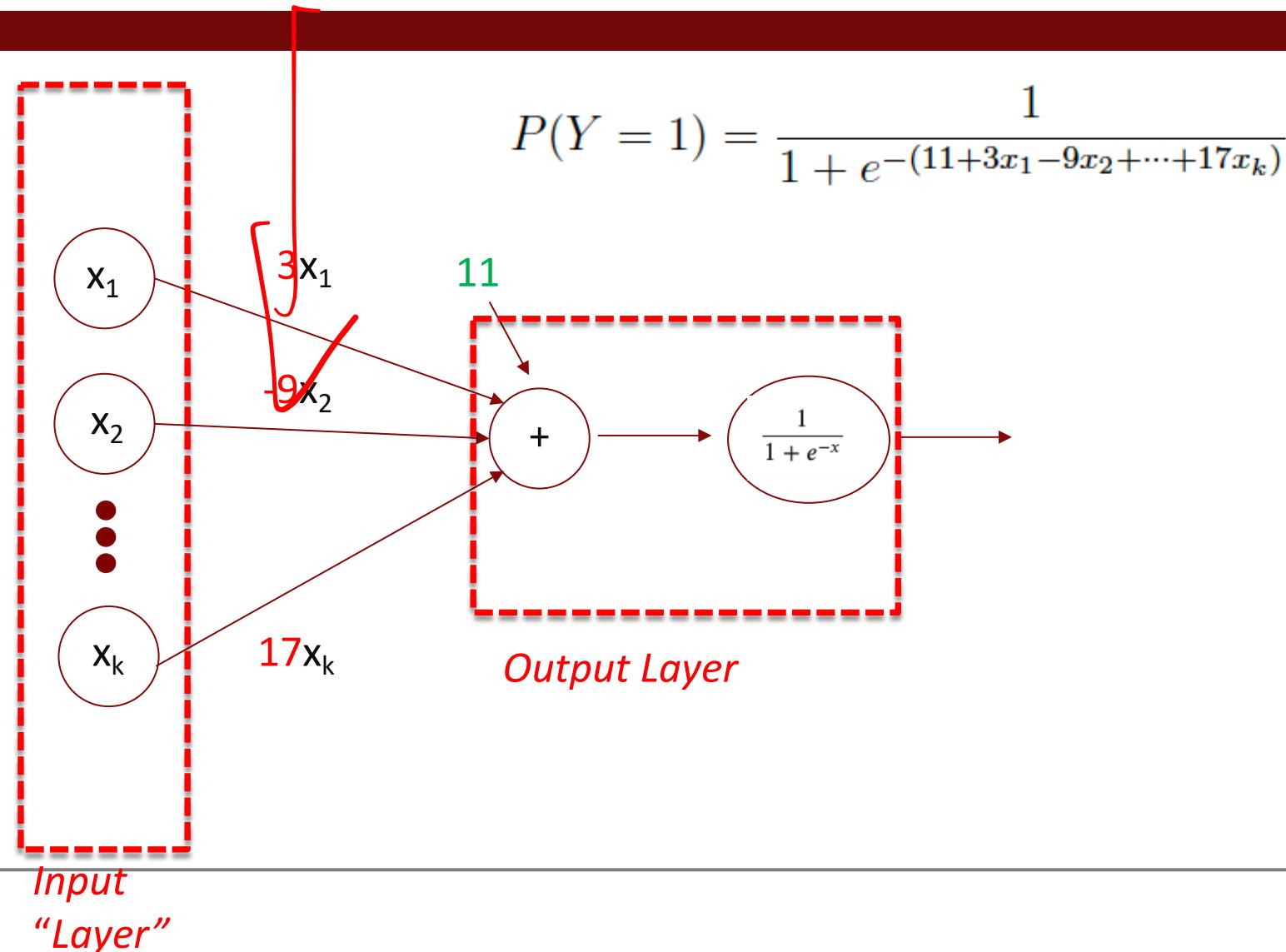


# Terminology



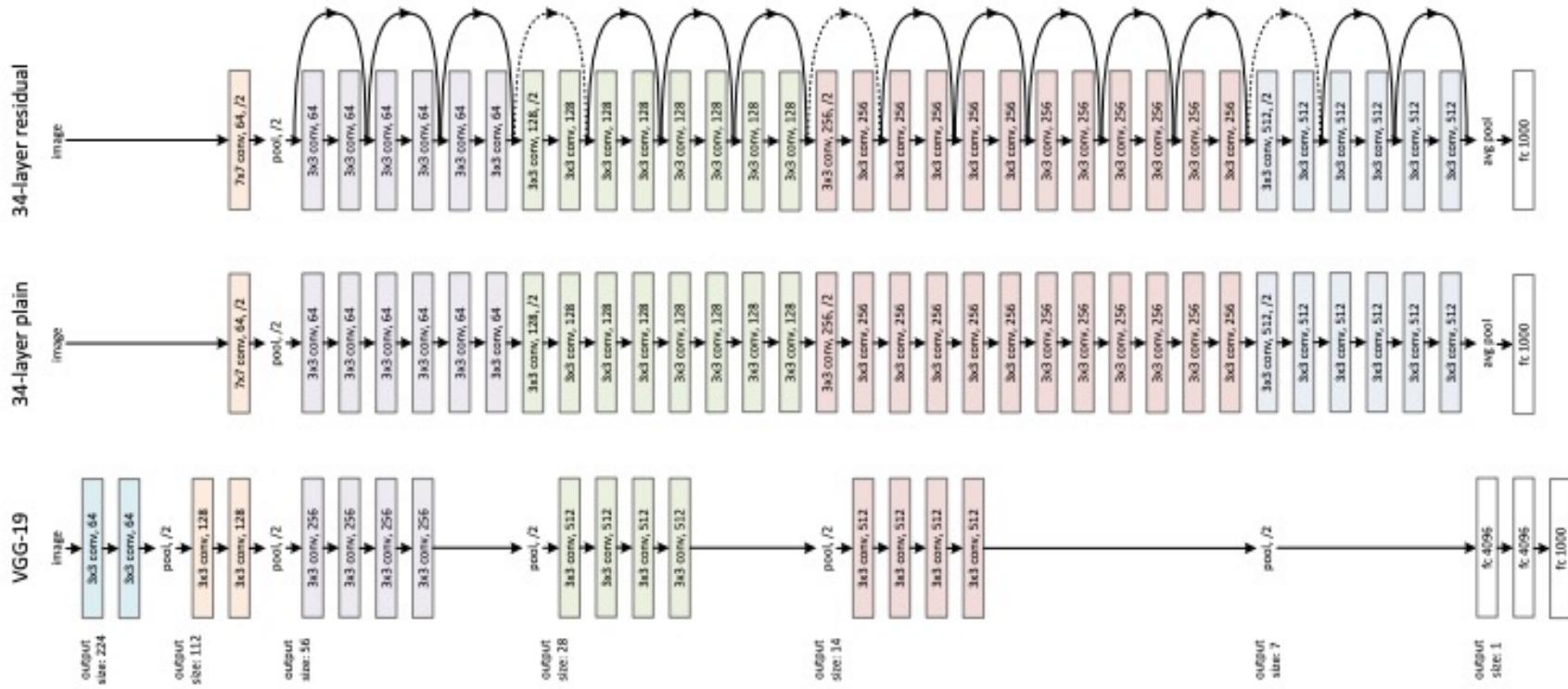
When every neuron in a layer is connected to every neuron in the next layer, it is called Dense or Fully Connected

The general logistic regression model is an NN  
but a simple one since it has no hidden layers



Deep Learning is *just* neural networks with lots and lots of ...

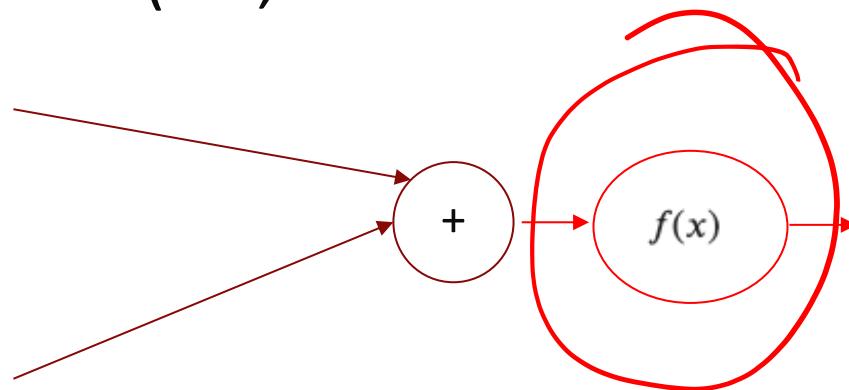
# Deep Learning is just neural networks with lots and lots of hidden layers



<https://arxiv.org/pdf/1512.03385.pdf>

# Activation functions

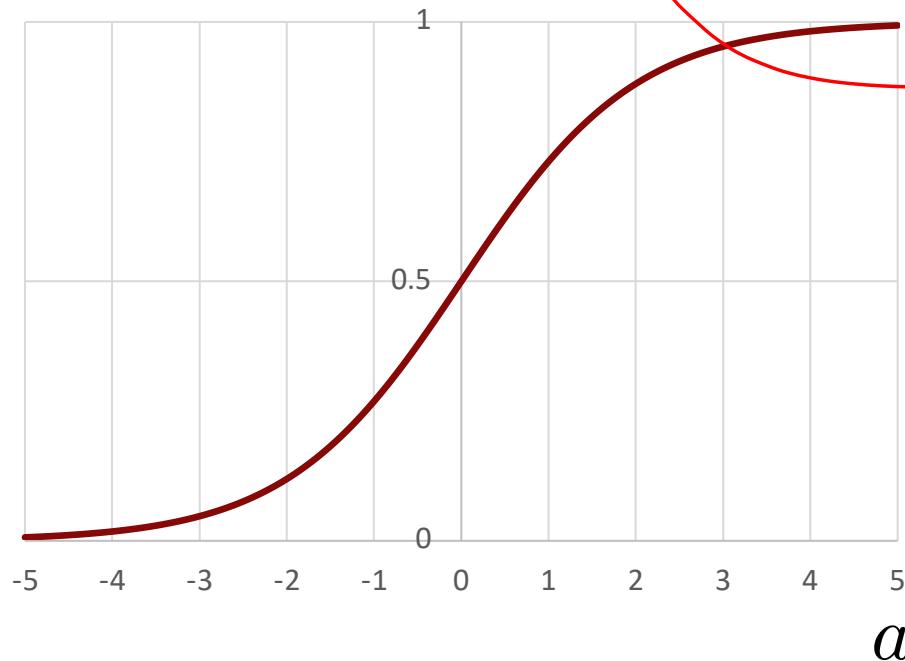
The activation function of a node is just a function that receives a single number and outputs a single number (i.e., scalar in  $\rightarrow$  scalar out)



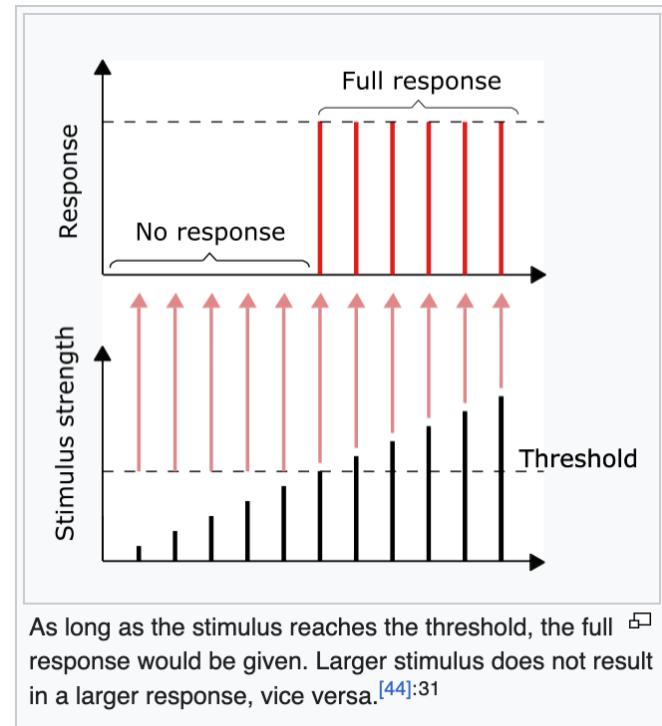
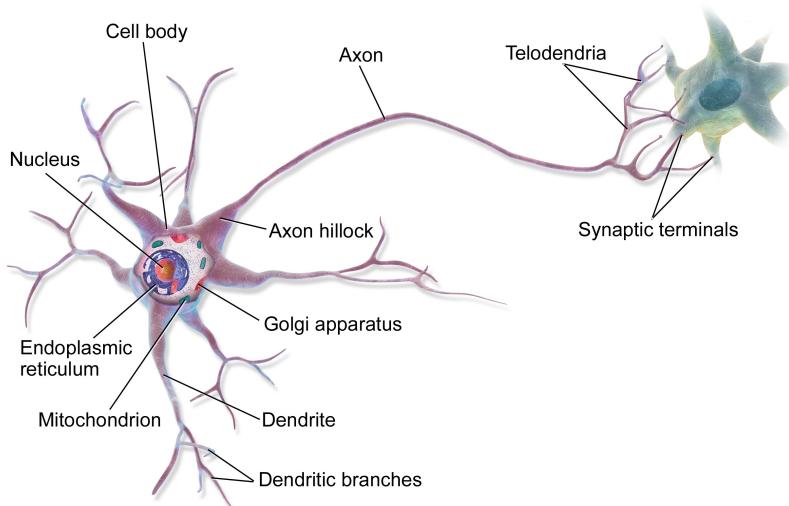
# Common Activation Functions

Sigmoid activation function:

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



# Aside: All-or-Nothing Principle

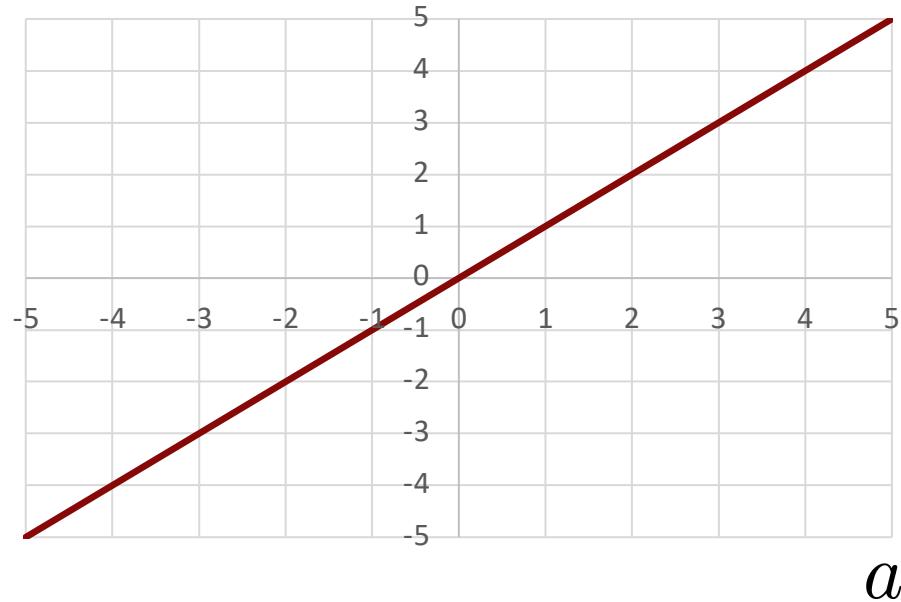


<https://en.wikipedia.org/wiki/Neuron>

# Common Activation Functions

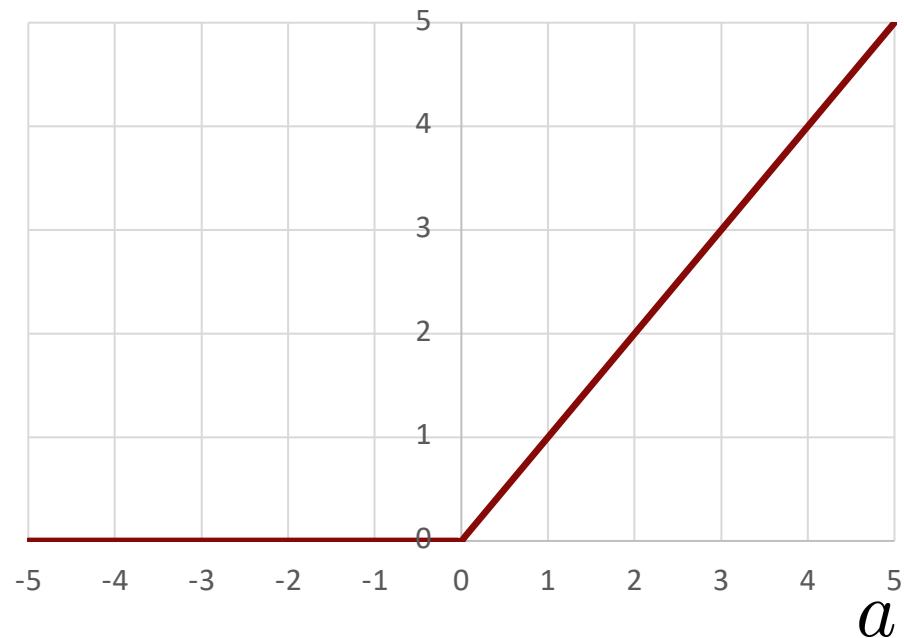
Linear activation function:

$$f(a) = a$$

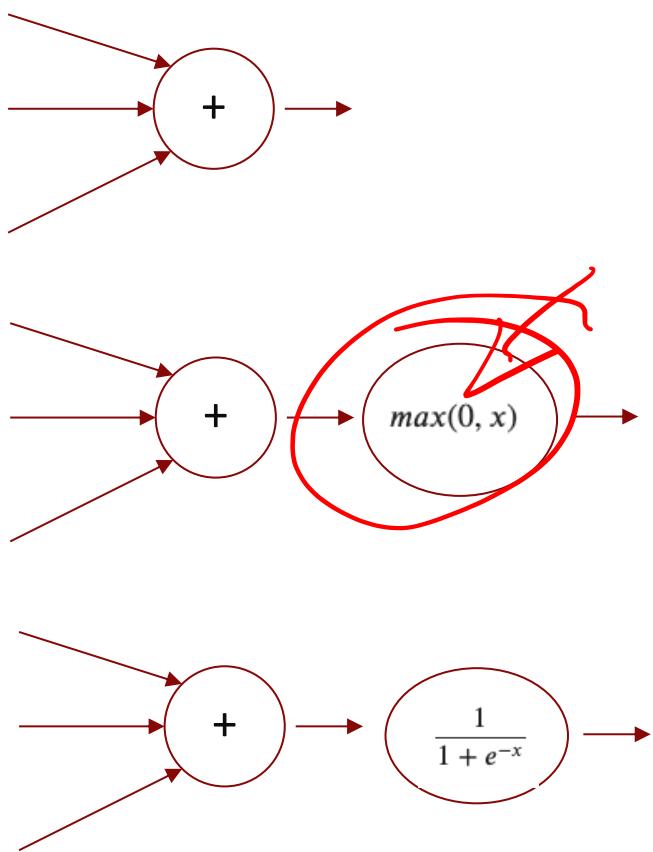


# Common Activation Functions

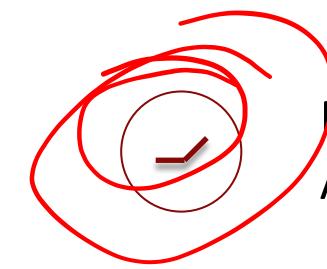
ReLU: 
$$g(a) = \max(0, a)$$
  
("Rectified Linear Unit")



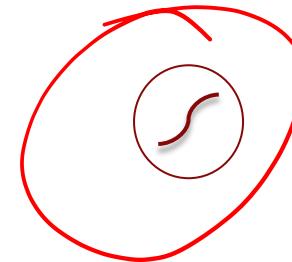
# We will use this “visual shorthand” from now on



Linear  
Activation

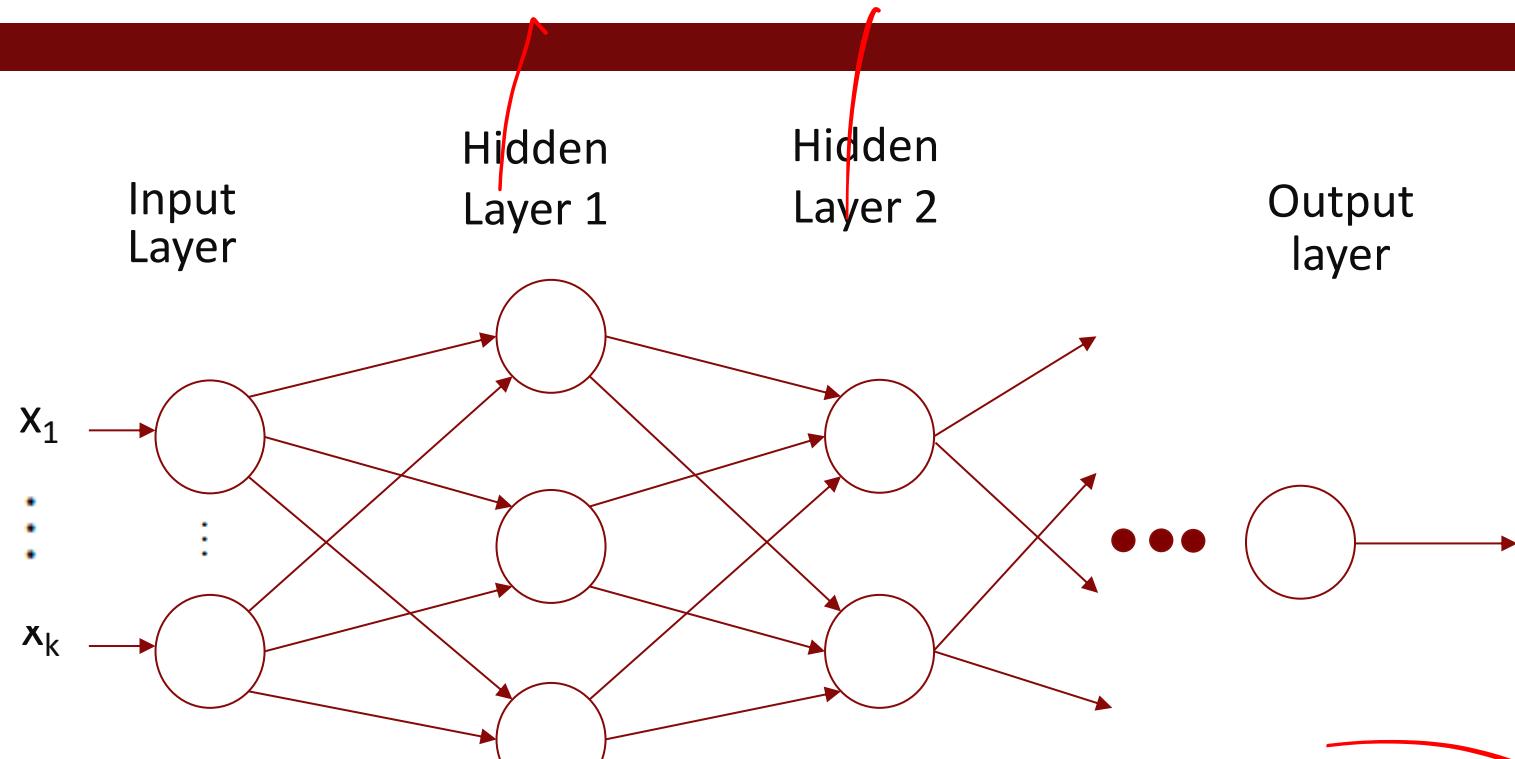


ReLU  
Activation



Sigmoid  
Activation

# Recap: Designing a DNN



User chooses the # of hidden layers, # neurons or units in each layer, the activation function(s) for the hidden layers and for the output layer

# Let's use a DNN for our interview classifier

- Recall the problem:
  - Two input variables (i.e., GPA and experience)
  - An output variable that should be between 0 and 1

User chooses the # of hidden layers, # units in each layer, the activation function(s) for the hidden layers and for the output layer

# Let's use a DNN for our interview classifier

- Recall the problem:
  - Two input variables (i.e., GPA and experience)
  - An output variable that should be between 0 and 1
- Design choices
  - We will use **one hidden layer** with **3 neurons (ReLU)**
  - Since the output is constrained to be in  $(0,1)$ , we will use the **sigmoid for the output layer**

User chooses the # of hidden layers, # units in each layer, the activation function(s) for the hidden layers and for the output layer

# Let's practice setting up a simple NN

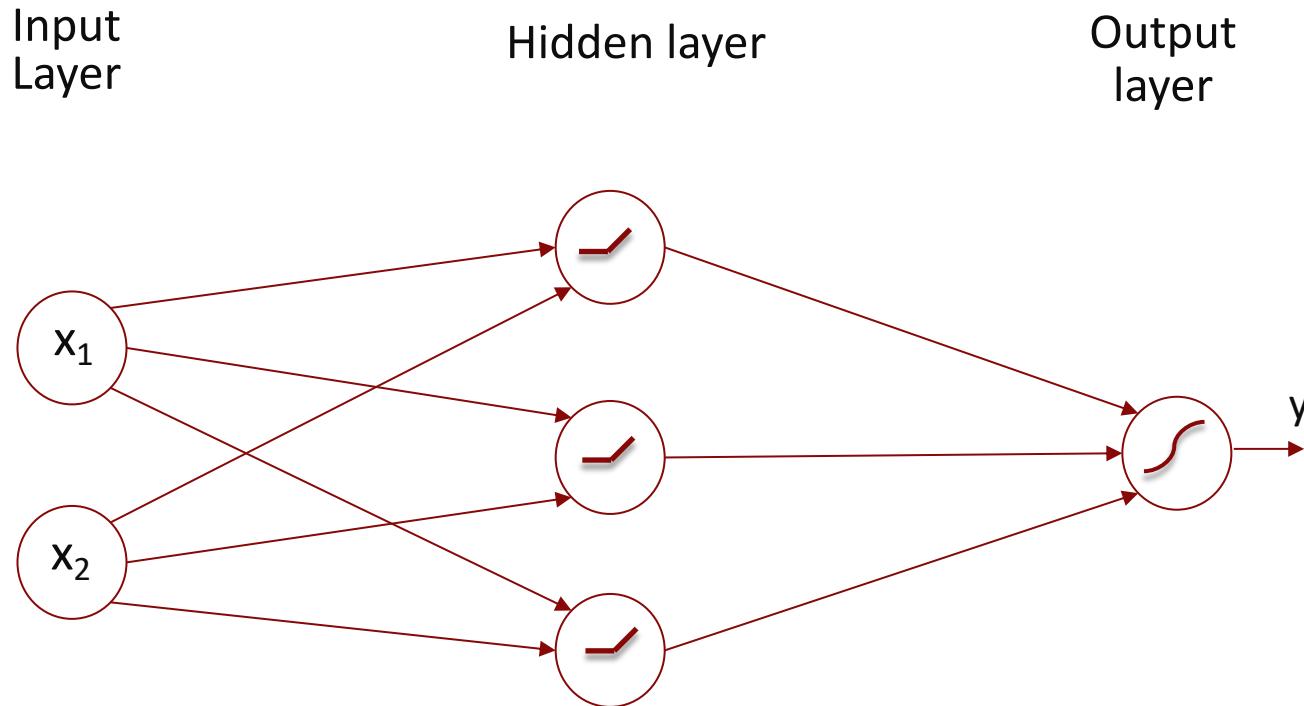
Input  
Layer

Hidden layer

Output  
layer

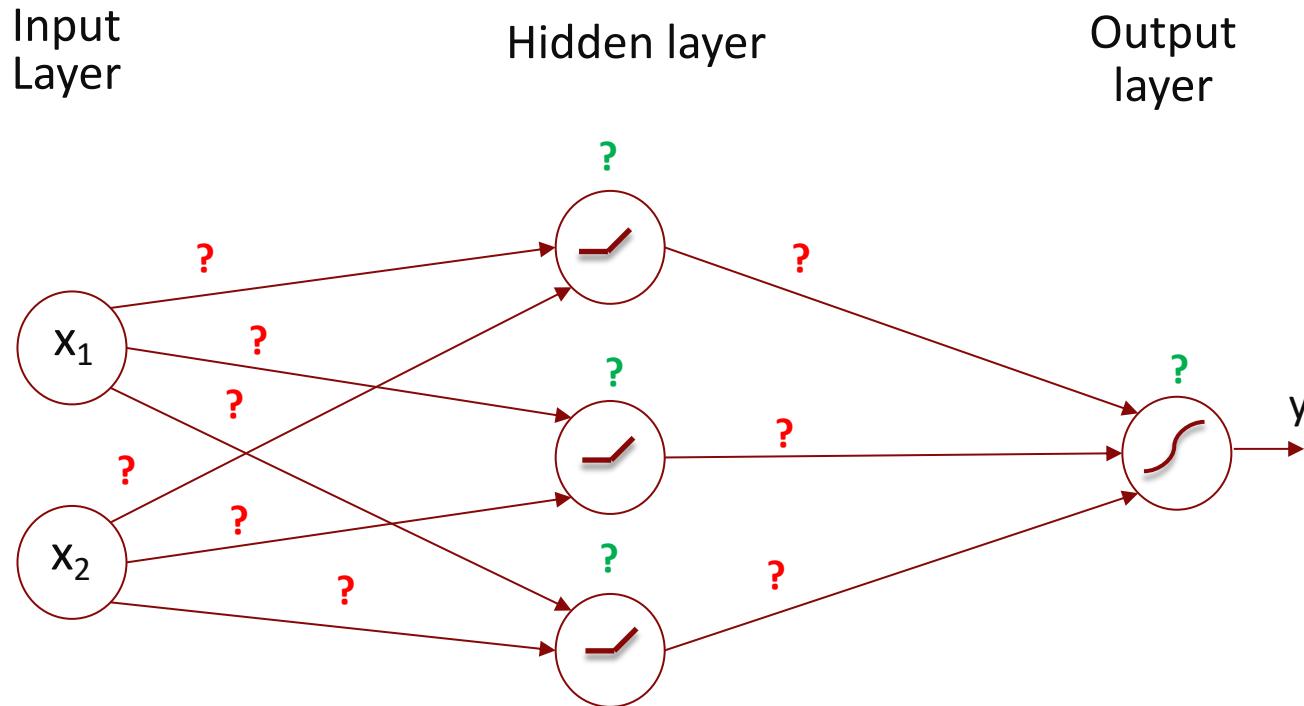
How many parameters (i.e., weights and biases) does this network have?

# Let's practice setting up a simple NN



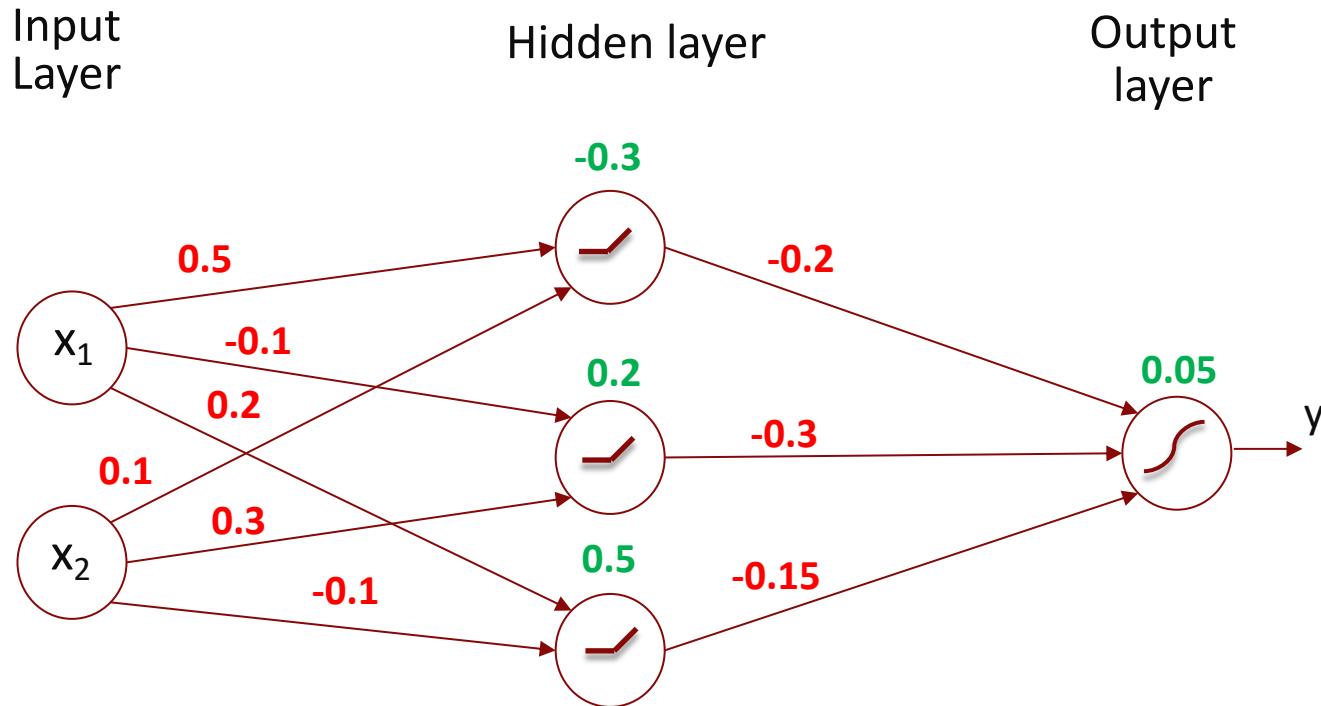
How many parameters (i.e., weights and biases) does this network have?

# Let's practice setting up a simple NN



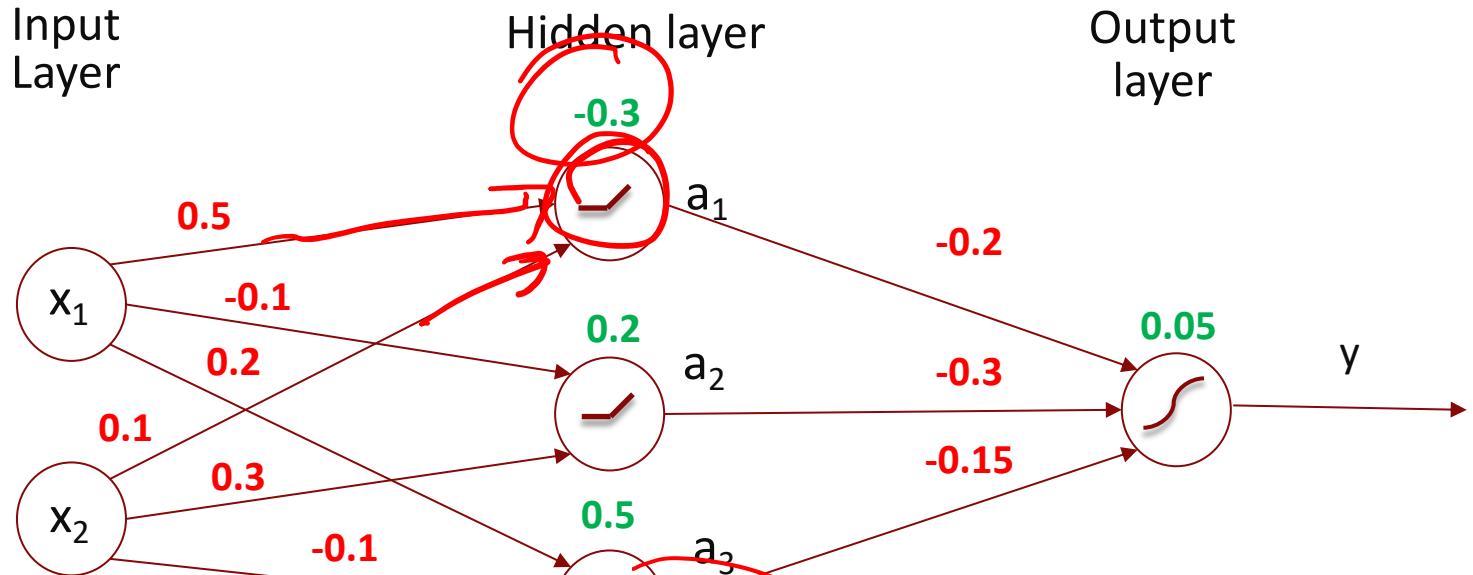
How many parameters (i.e., weights and biases) does this network have? **13**

Let's assume that we have trained\* this network on data and have found these values for the parameters



\*details coming up shortly

# Predicting with the NN

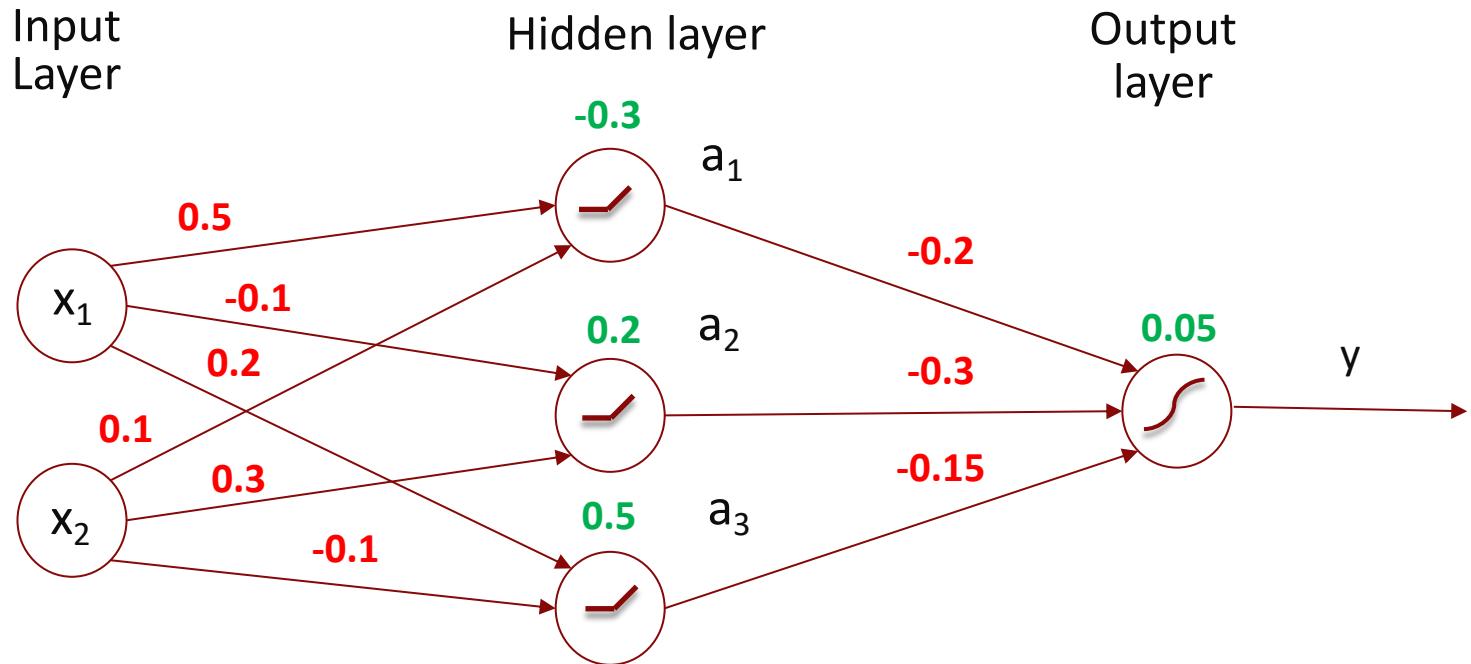


**Output of hidden layer :**

- Top node:
- Middle node:
- Bottom node:

$$\max(0, -0.3 + 0.5x_1 + 0.1x_2) = a_1$$
$$\max(0, 0.2 - 0.1x_1 + 0.3x_2) = a_2$$
$$\max(0, 0.5 + 0.2x_1 - 0.1x_2) = a_3$$

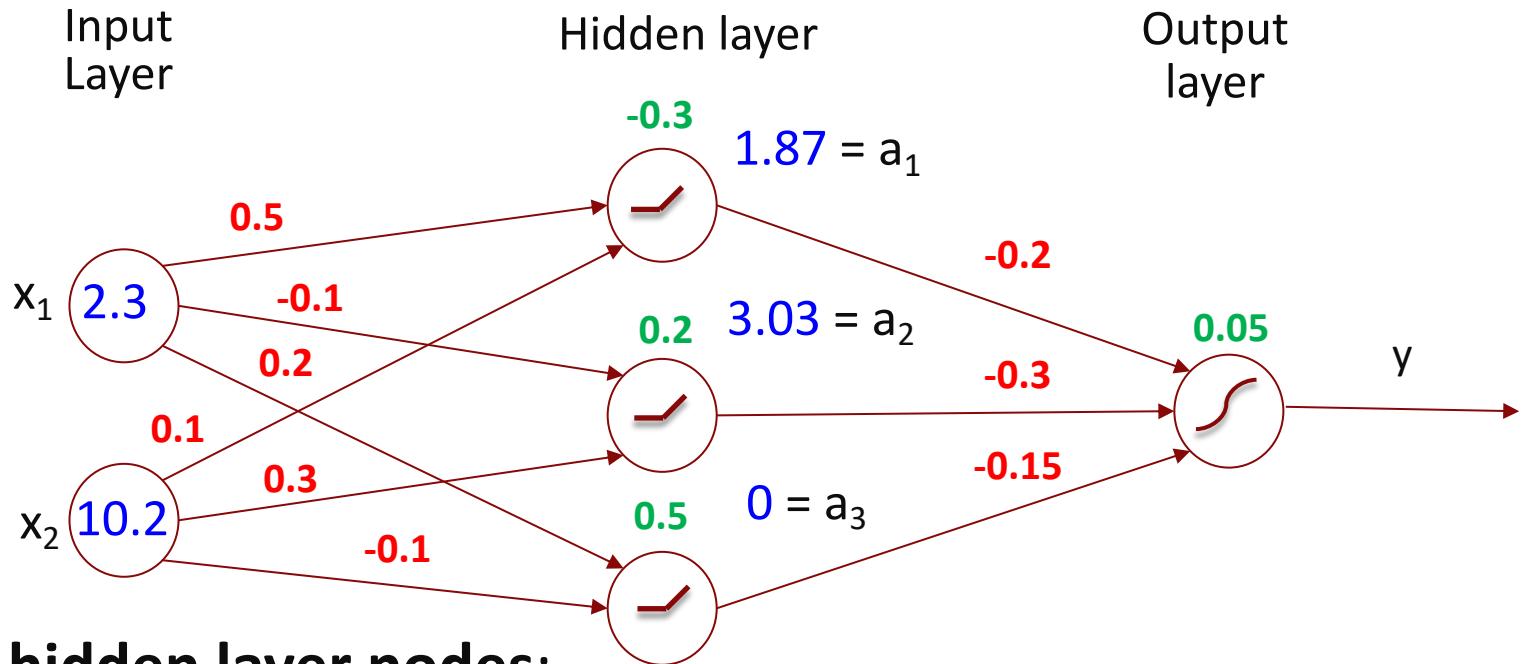
# Predicting with the NN



- Recall  $a_1$ ,  $a_2$ , and  $a_3$  are the output of the hidden layer nodes
- Output of **output layer node**:

$$\frac{1}{1 + e^{-(0.05 - 0.2a_1 - 0.3a_2 - 0.15a_3)}}$$

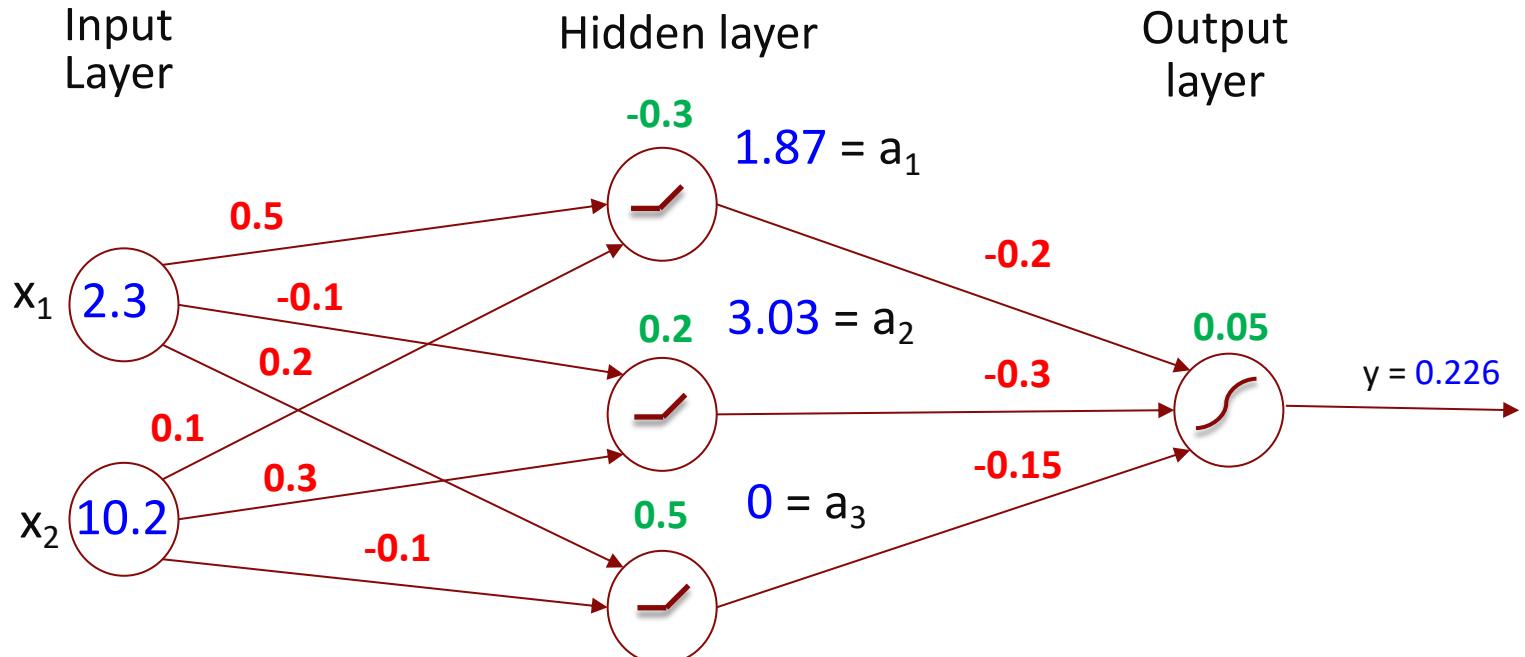
# Predicting with the NN



**Output of hidden layer nodes:**

- Top node:  $\max(0, -0.3 + 0.5*2.3 + 0.1*10.2) = 1.87 = a_1$
- Middle node:  $\max(0, 0.2 - 0.1*2.3 + 0.3*10.2) = 3.03 = a_2$
- Bottom node:  $\max(0, 0.5 + 0.2*2.3 - 0.1*10.2) = 0 = a_3$

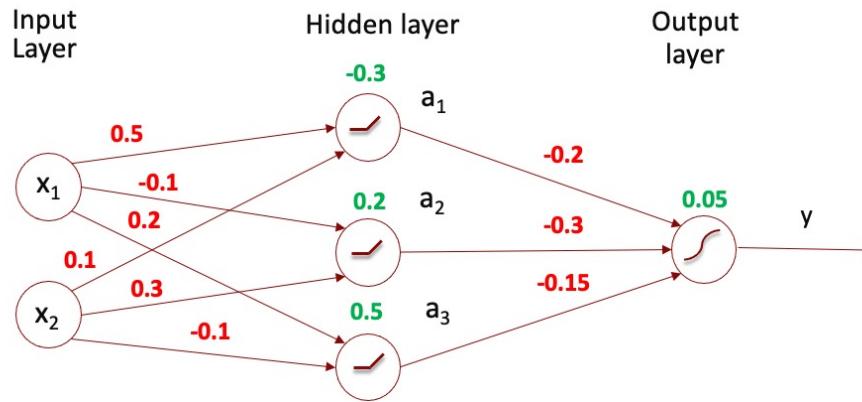
# Predicting with the NN



Output of **output layer node**:

$$\frac{1}{1 + e^{-(0.05 - 0.2 \cdot 1.87 - 0.3 \cdot 3.03 - 0.15 \cdot 0)}} = 0.226$$

# The Network can be written as this function

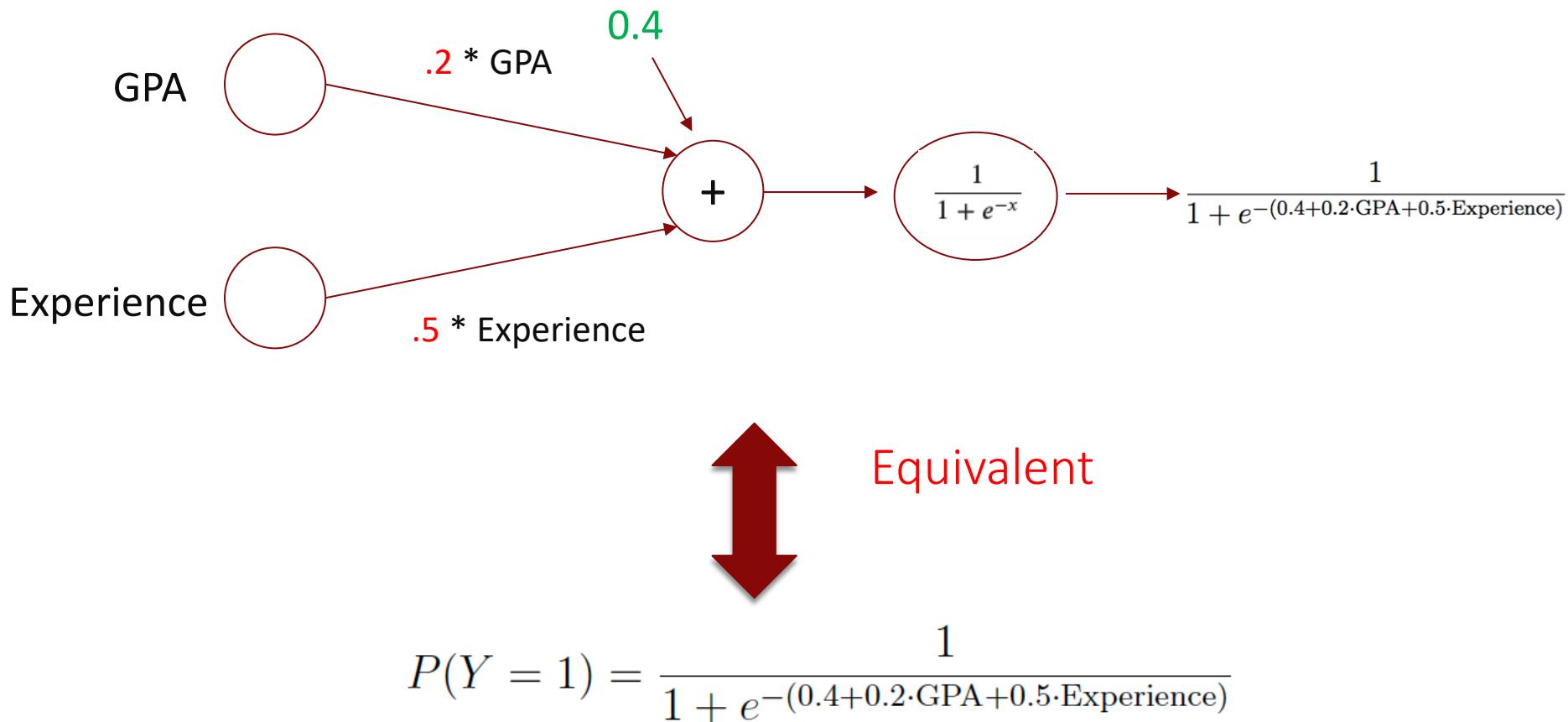


Equivalent

1

$$y = \frac{1}{1 + e^{-(.05 - 0.2(\max(0, -0.3 + 0.5x_1 + 0.1x_2)) - 0.3(\max(0, 0.2 - 0.1x_1 + 0.3x_2)) - 0.15(\max(0.5 + 0.2x_1 - 0.1x_2)))}}$$

## Contrast w/ the old model – no hidden layer



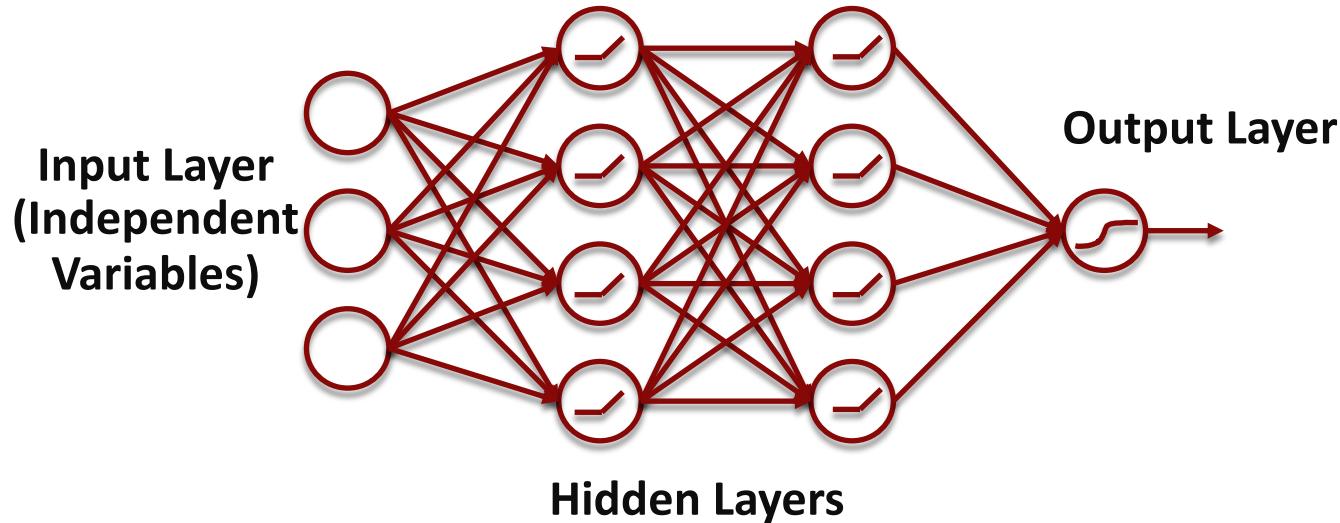
Note the complexity of even this simple network compared to a logistic regression model

$$y = \frac{1}{1 + e^{-(.05 - 0.2(\max(0, -0.3 + 0.5x_1 + 0.1x_2)) - 0.3(\max(0, 0.2 - 0.1x_1 + 0.3x_2)) - 0.15(\max(0.5 + 0.2x_1 - 0.1x_2)))}}$$

$$y = \frac{1}{1 + e^{-(0.4 + 0.2x_1 + 0.5x_2)}}$$

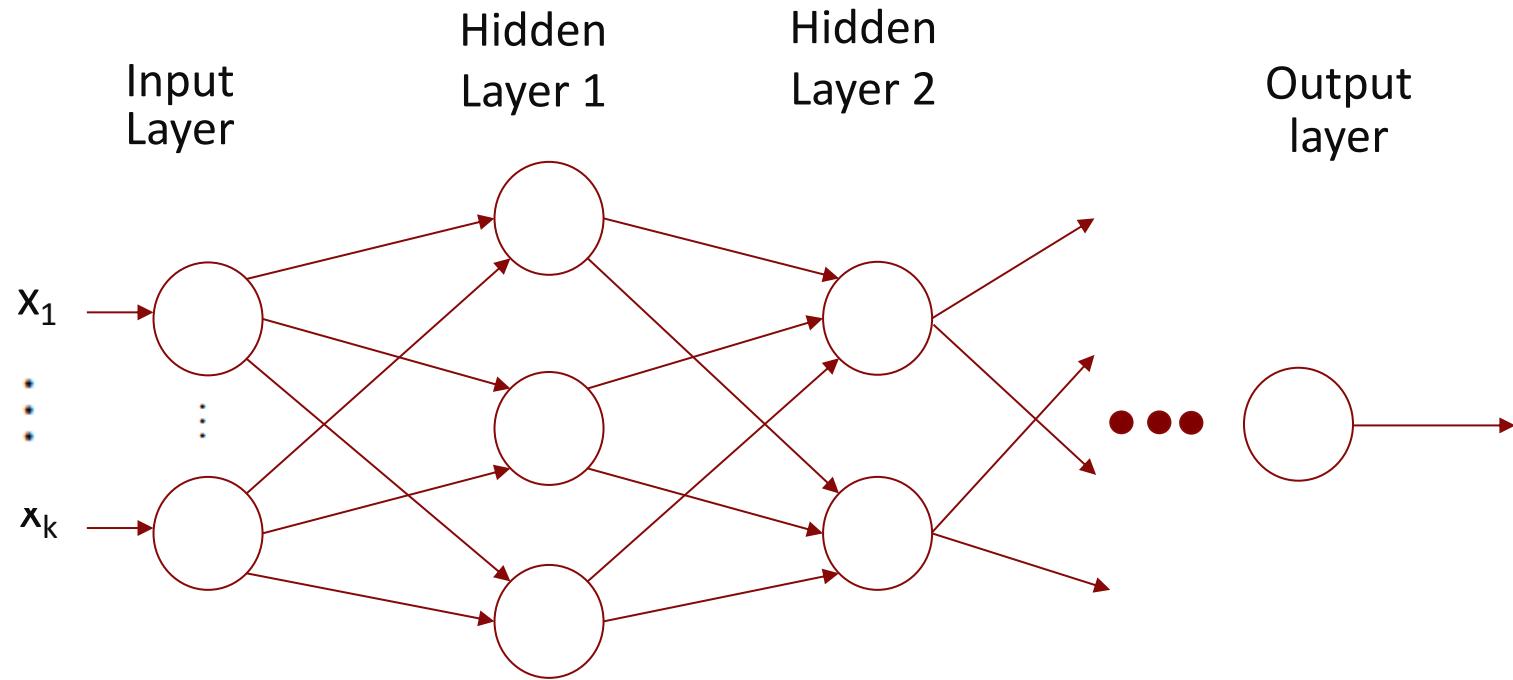
$y$  is a much more complex function of its inputs  $x_1$  and  $x_2$  compared to (say)  $\frac{1}{1 + e^{-(0.4 + 0.2 \cdot \text{GPA} + 0.5 \cdot \text{Experience})}}$  and it can capture more complex relationships between  $x$  and  $y$ .

# Summary: A Deep Neural Network



- This is a **feedforward** (or **vanilla**) neural network
- In general, the arrangement of neurons into layers, the activation functions, and the connections between layers are referred to as the network's **architecture**

# Summary: Designing a DNN



User chooses the # of hidden layers, # units in each layer, the activation function(s) for the hidden layers and for the output layer

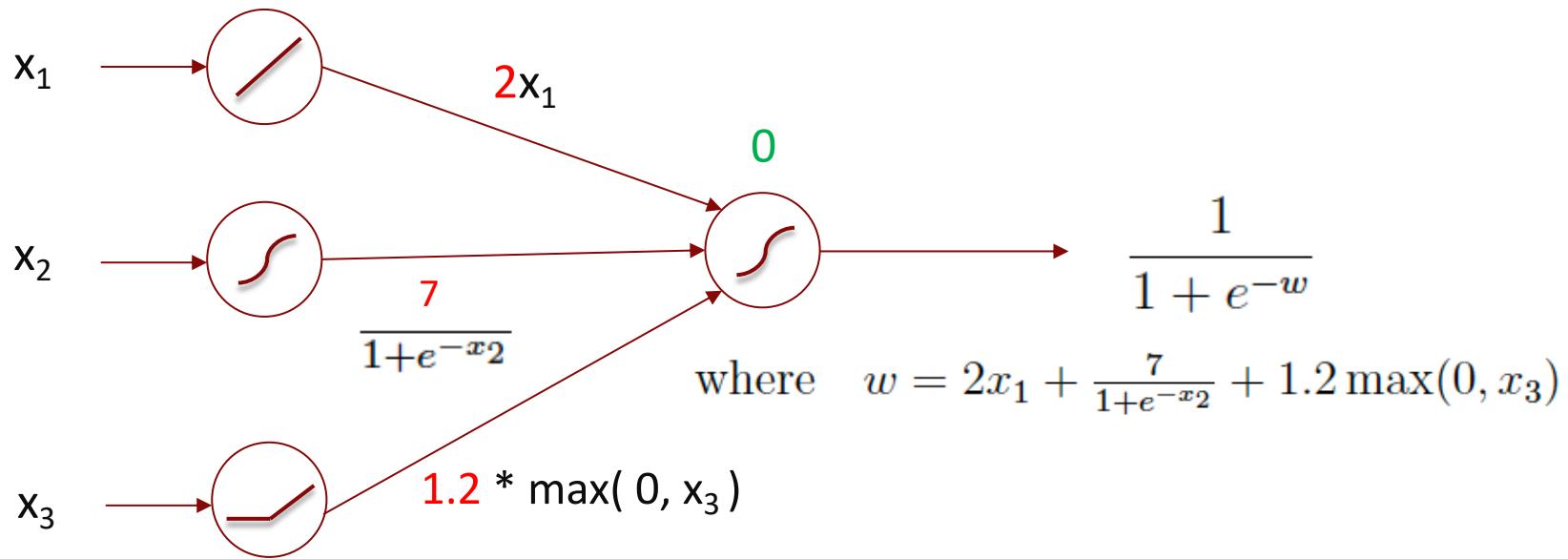
# Next Time

- How to ‘train’ the network?
- And we’ll code and train our first NN

# Appendix

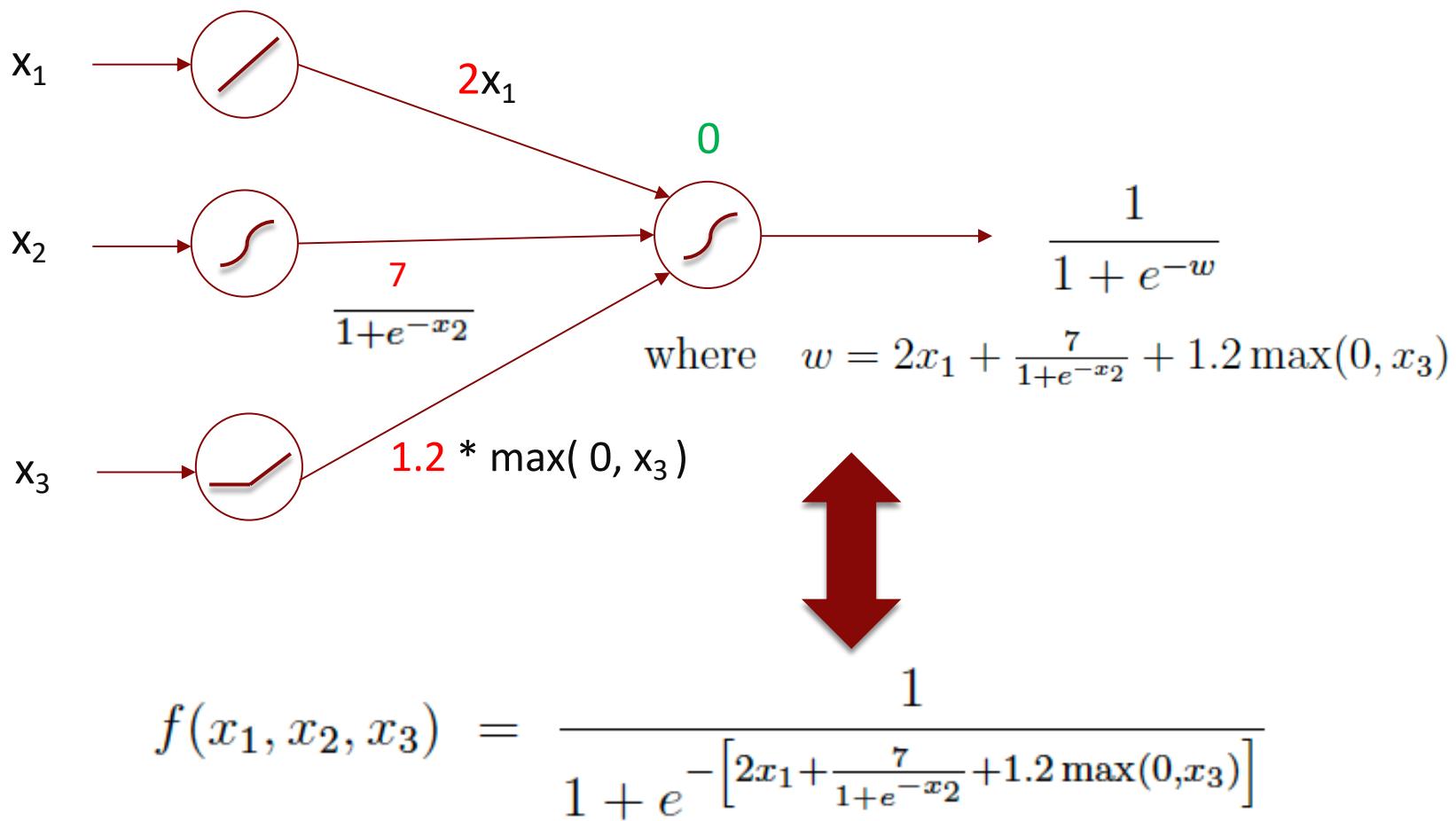


# Activation functions - Practice



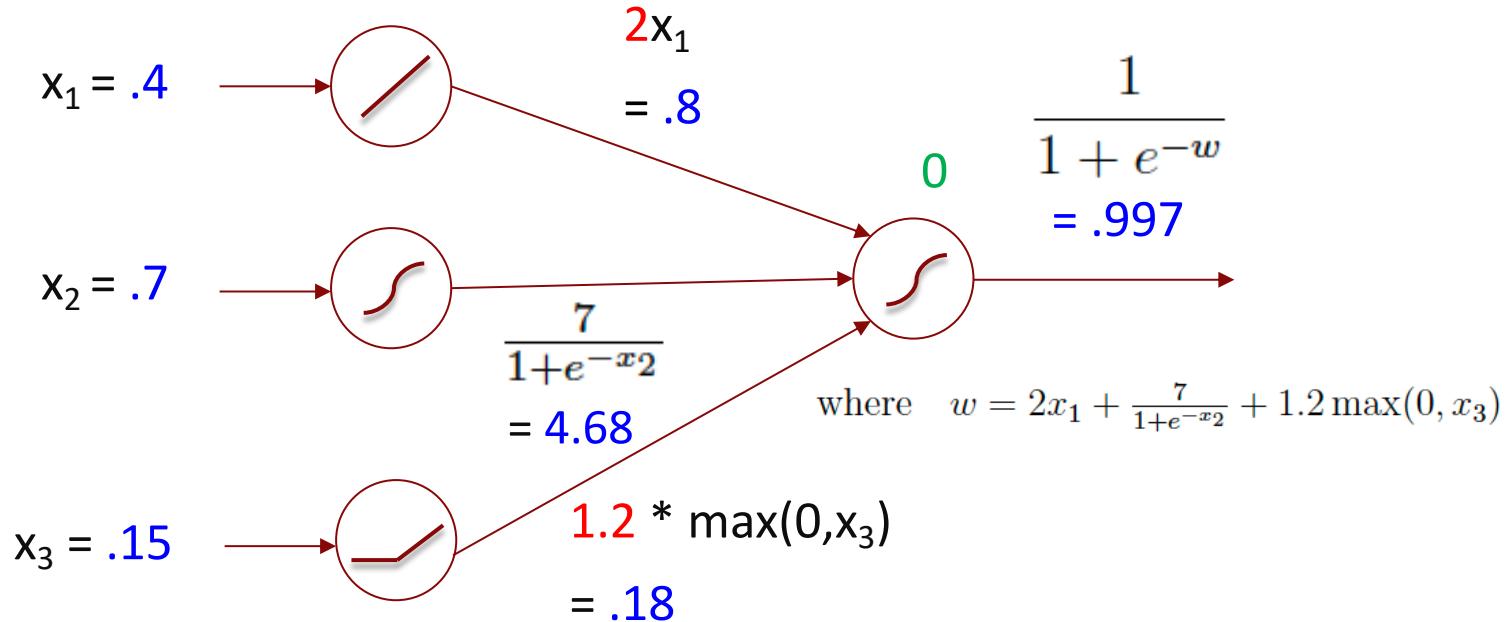
This is starting to look complicated. What does the final function look like?

# Activation functions - Practice



# Activation functions - Practice

Consider the observation  $(x_1, x_2, x_3) = (.4, .7, .15)$



# Recapping the three phases

Phase 1 (50s-90s): Artificial Intelligence

*“Vision: Computers to do tasks that only humans could do with if-then statements” but systems did not generalize*

Phase 2 (90s-10s): Statistical Machine Learning

*“With structure in your data, it is a statistical problem so linear regression, etc” but what about unstructured data?*

Phase 3 (10s-now): Deep Learning

*“Automatically extract useful representations from unstructured data”*