



Episode 1: Core AI Concepts

Alastair Droop

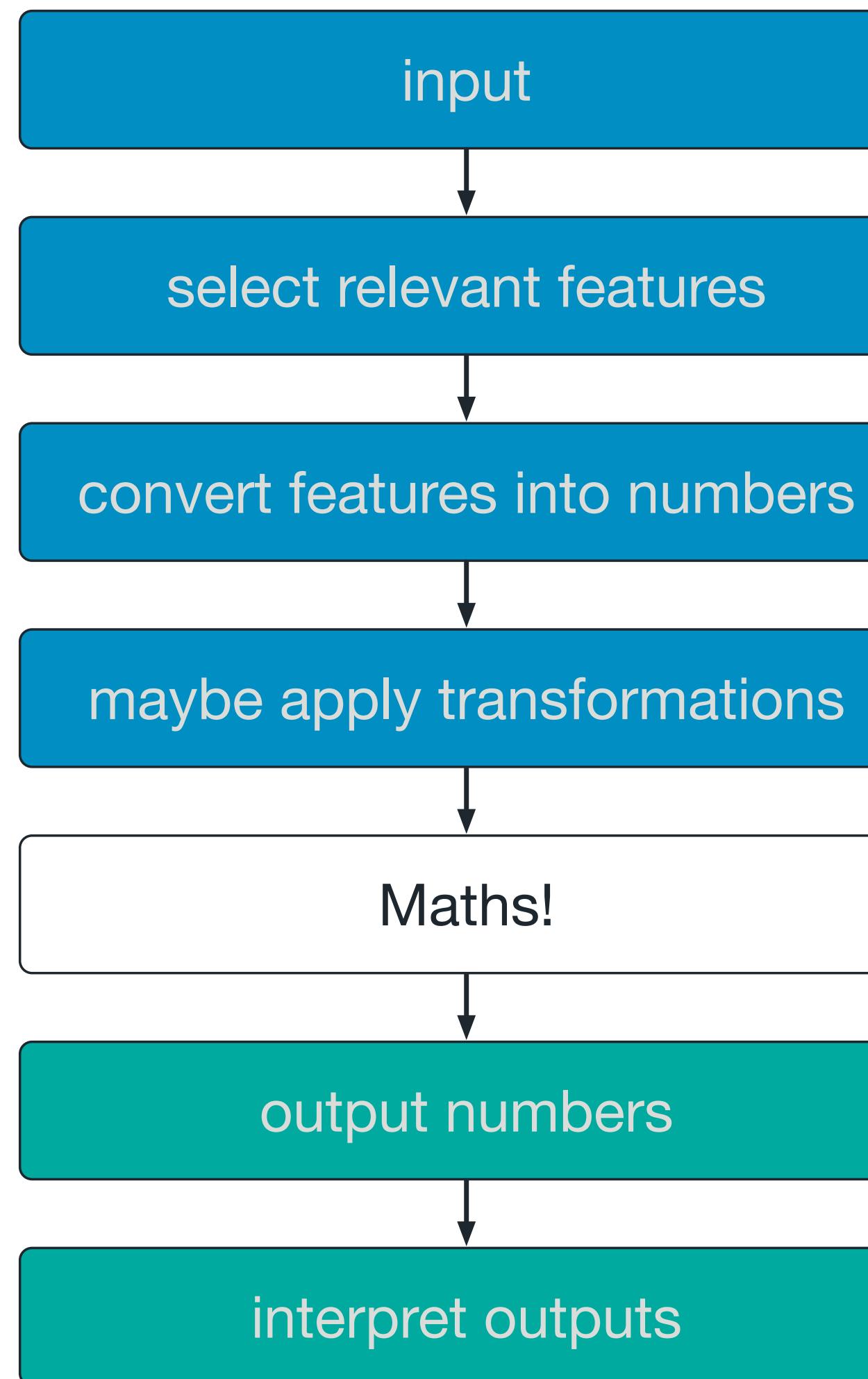
AI & Machine Learning: From Theory to Practice



“The capacity of computers or other machines to exhibit or simulate intelligent behaviour; the field of study concerned with this.”



Some Terminology...



The inputs to an AI process can be pretty much anything we can measure

Features are the things (called variables) that we want to measure (eg. Blood pressure, disease status, etc...)

We encode features into numbers so we can use maths to process them. Different kinds of features will use different encodings to do this

Transformations change the values of the features in defined ways to make them easier to use. Whatever transformations we do, we're changing all the values in the same way

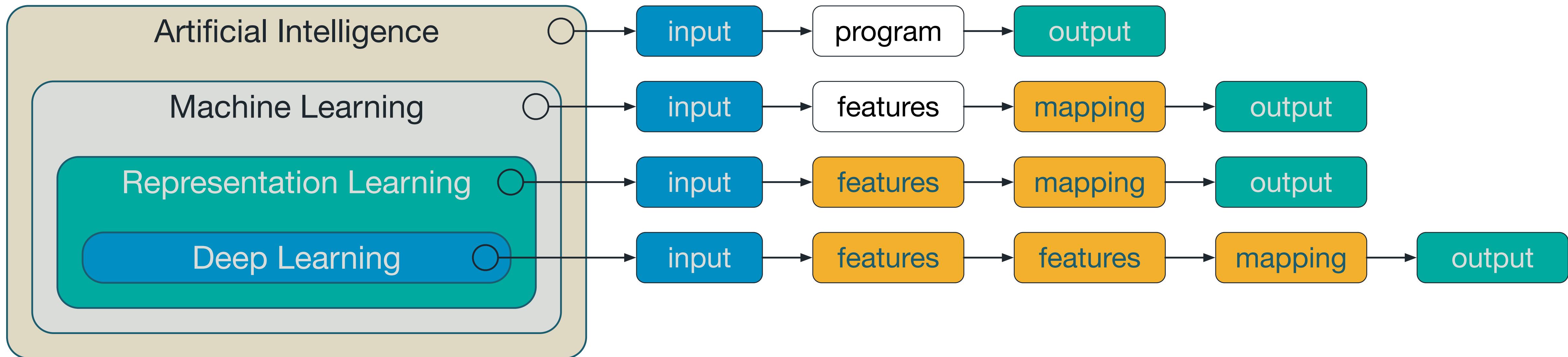
We input the transformed input features into a mathematical function (an algorithm) that generates the outputs for us

The AI function generates numerical values as outputs

We need to interpret the numerical values to work out what the AI has told us



Types of AI



The boxes in yellow are where we allow the computer to learn

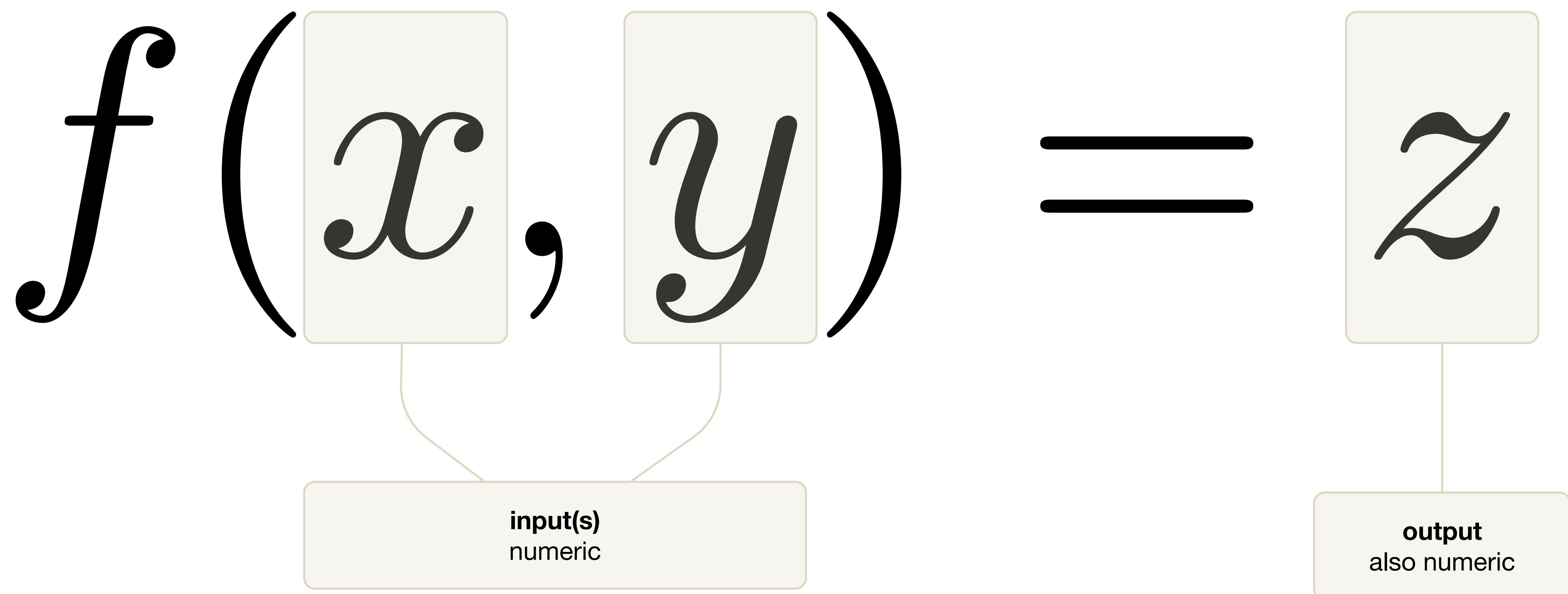
Machine learning: the computer learns how to modify & combine (*ie. map*) the inputs

Representation learning: the computer learns the features to extract as well as the mappings



Functions

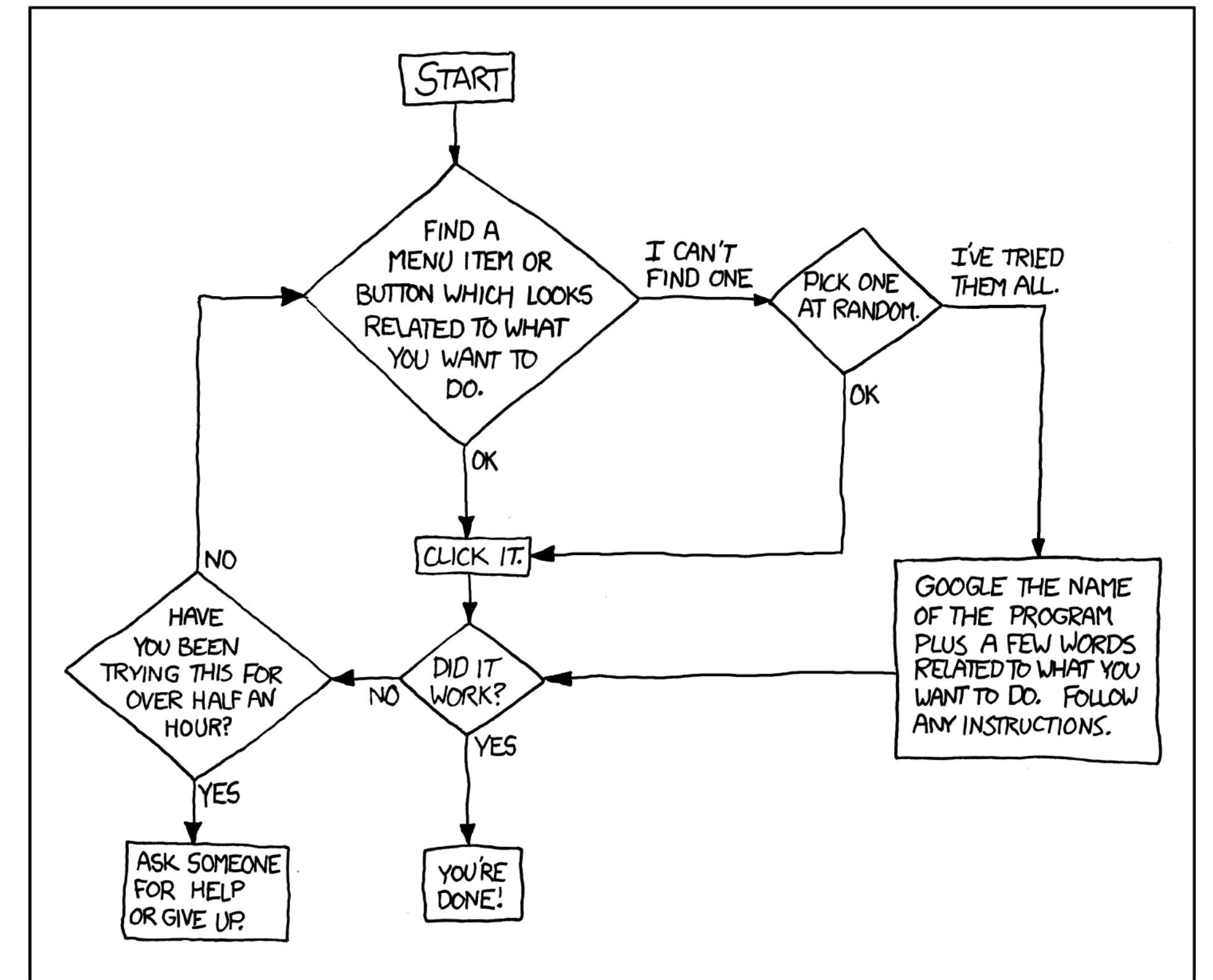
A function takes zero or more inputs and returns an output

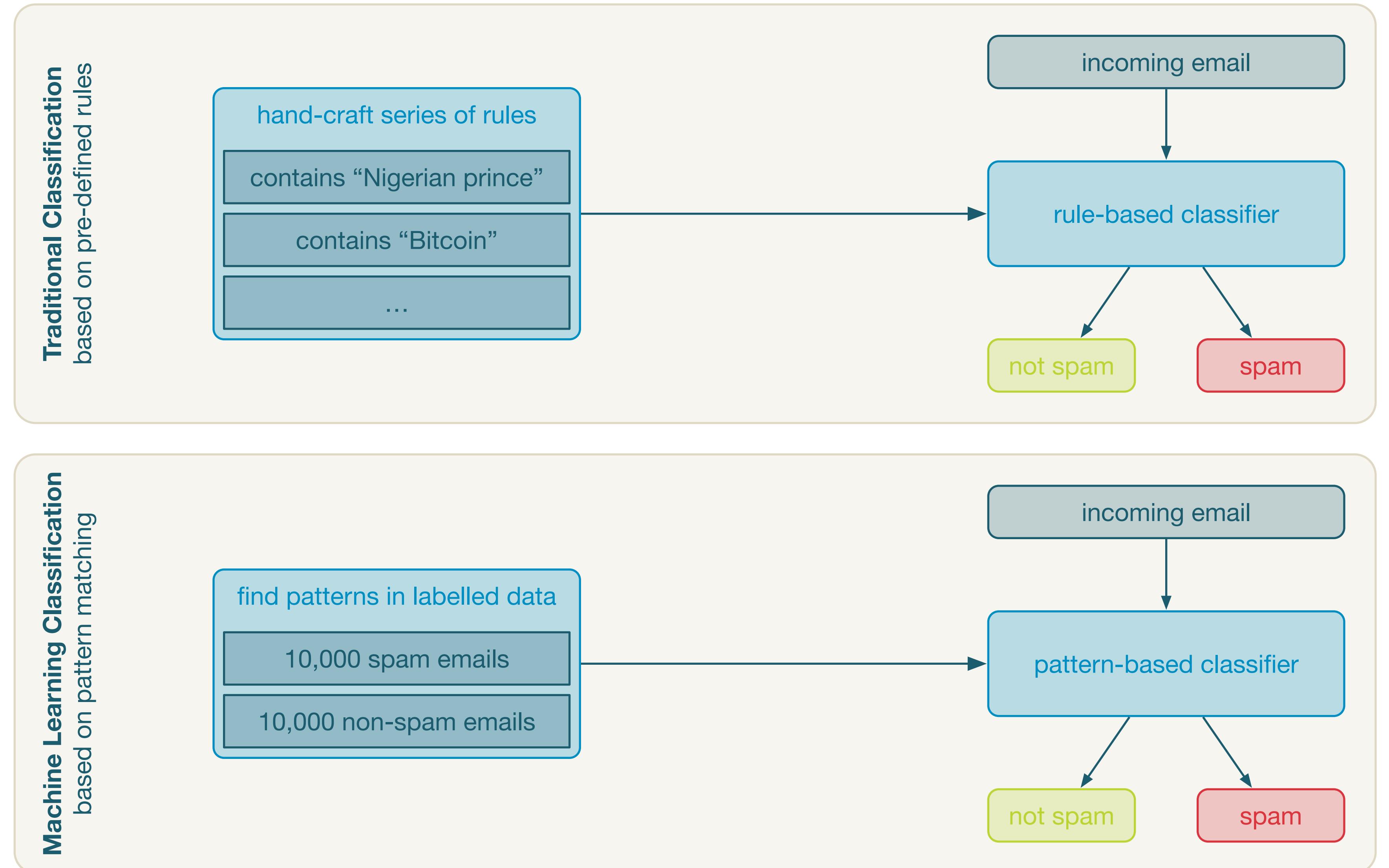




Algorithms

A series of specific & reproducible step-by-step instructions to perform a specific task







Classification

A very common problem in data science is classification

- Given a set of data, place each datapoint into one of several categories

Unsupervised Learning (eg. k-means & hierarchical clustering)

- No requirement to train on known data before use

Supervised Learning (eg. Linear regression, SVM, Decision Trees, Neural Networks)

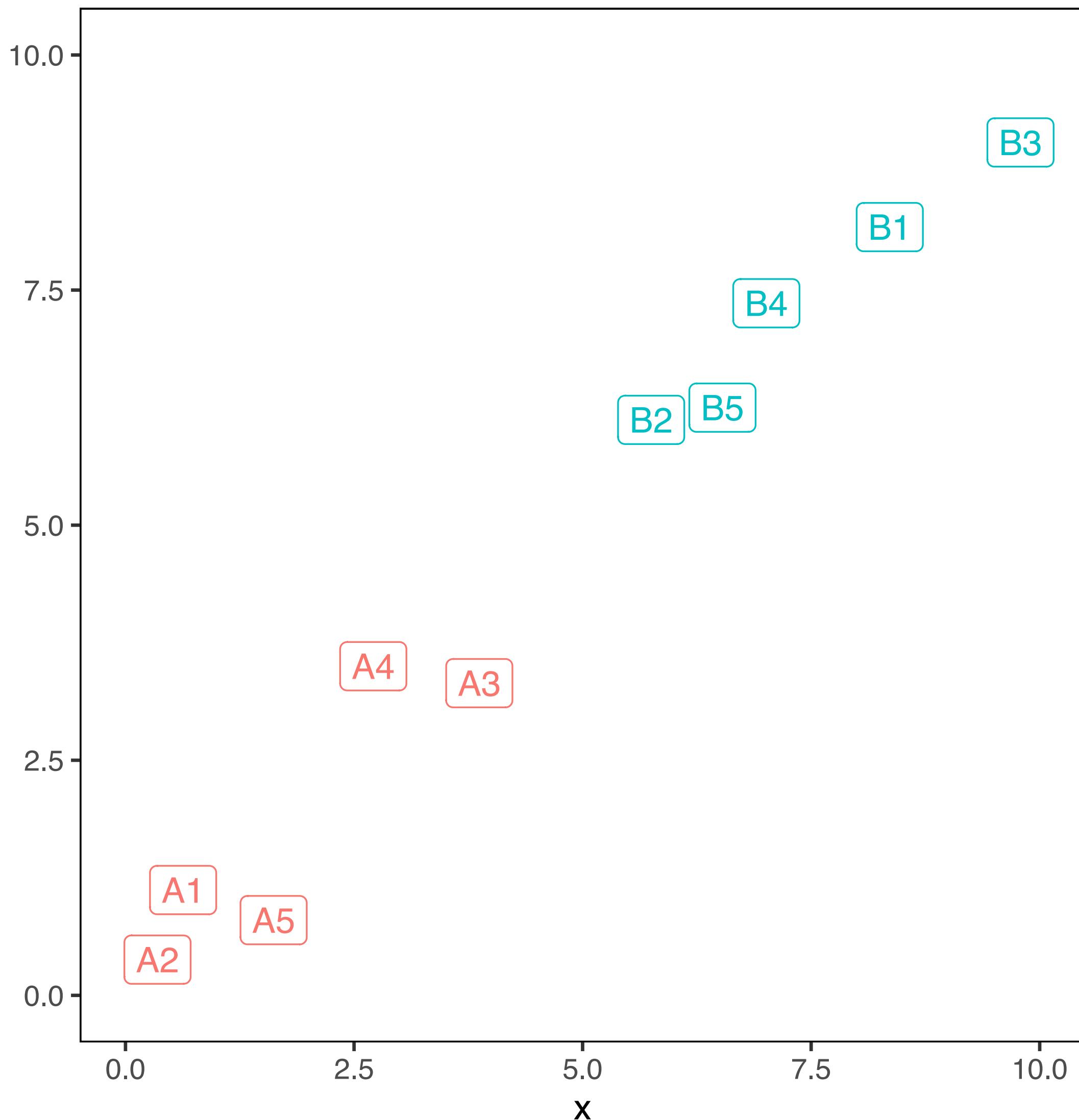
- Trained on known, labelled data before being used to classify unknown data



Imagine we have some 2D data

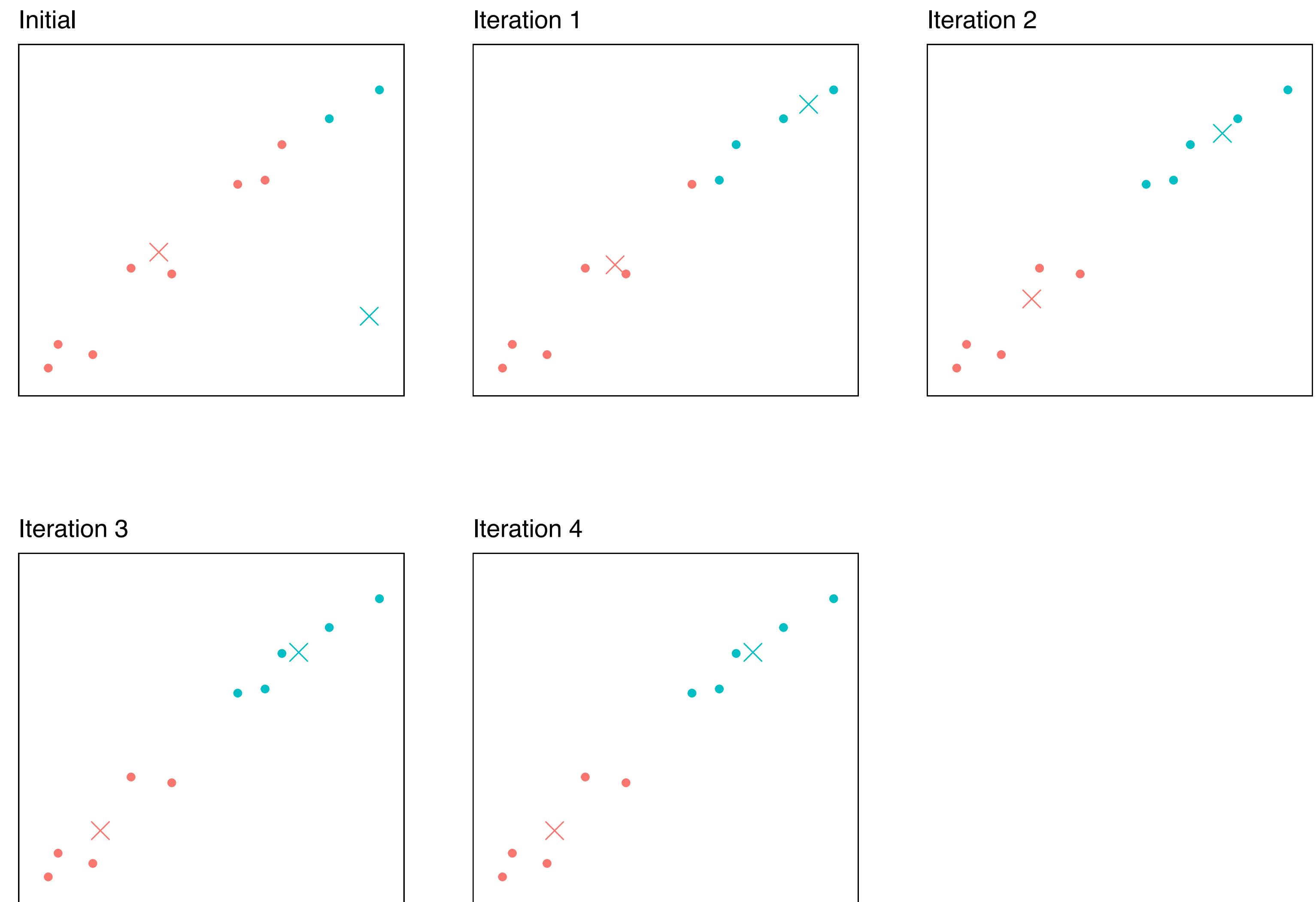
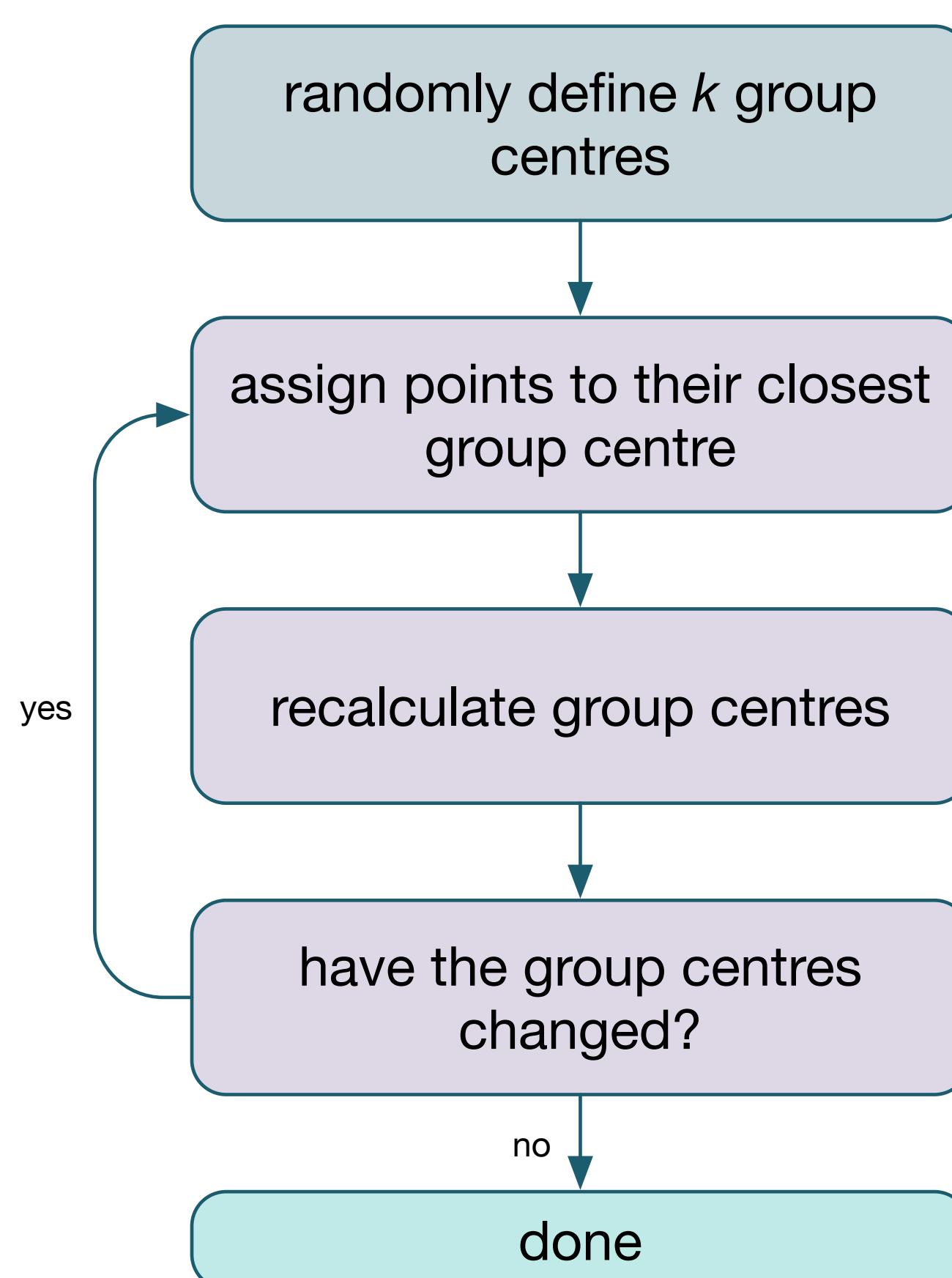
- Dimensions are simply variables that we measure
- Each data point has a value for each dimension

We wish to place each datapoint into a group



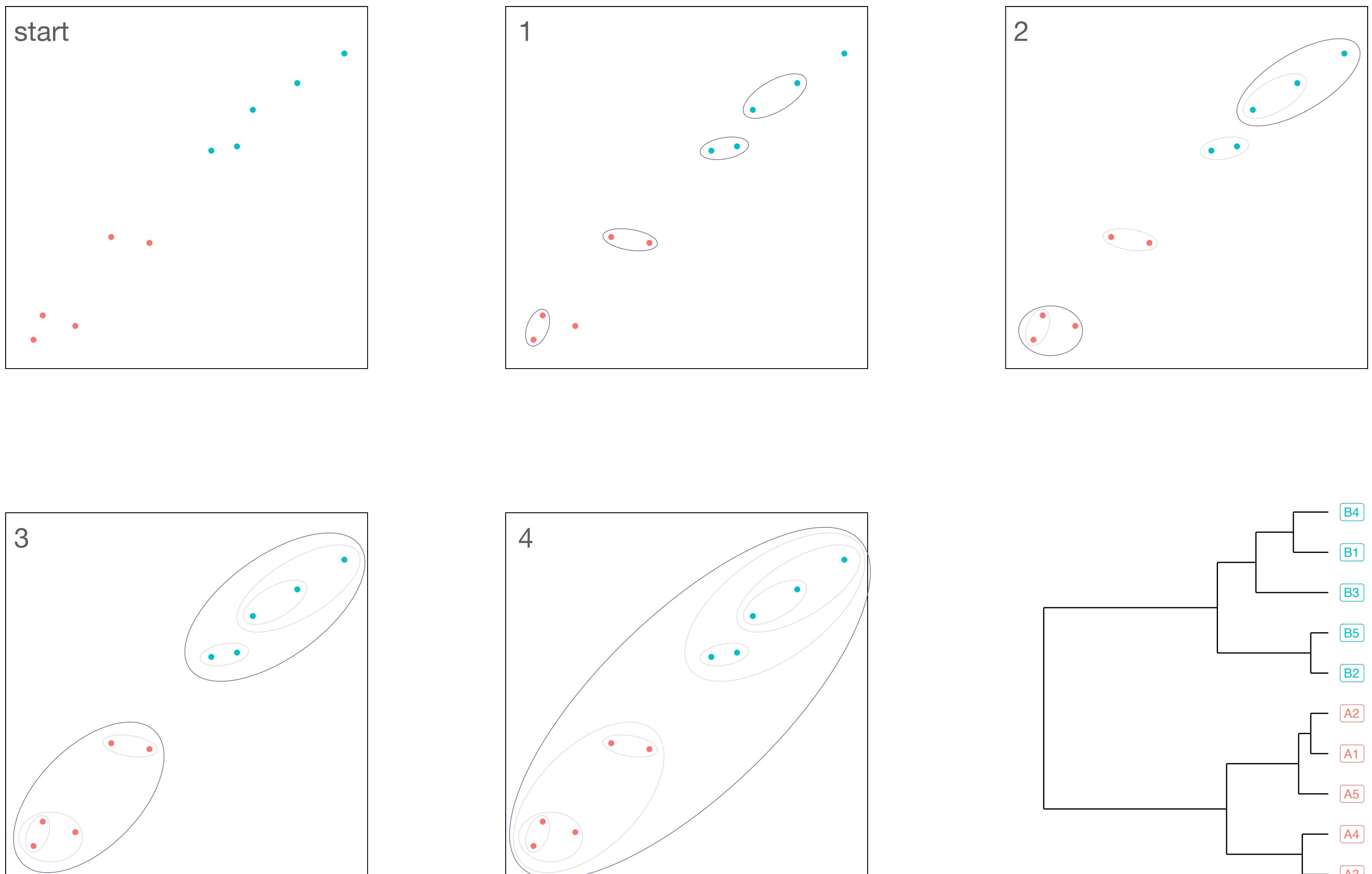
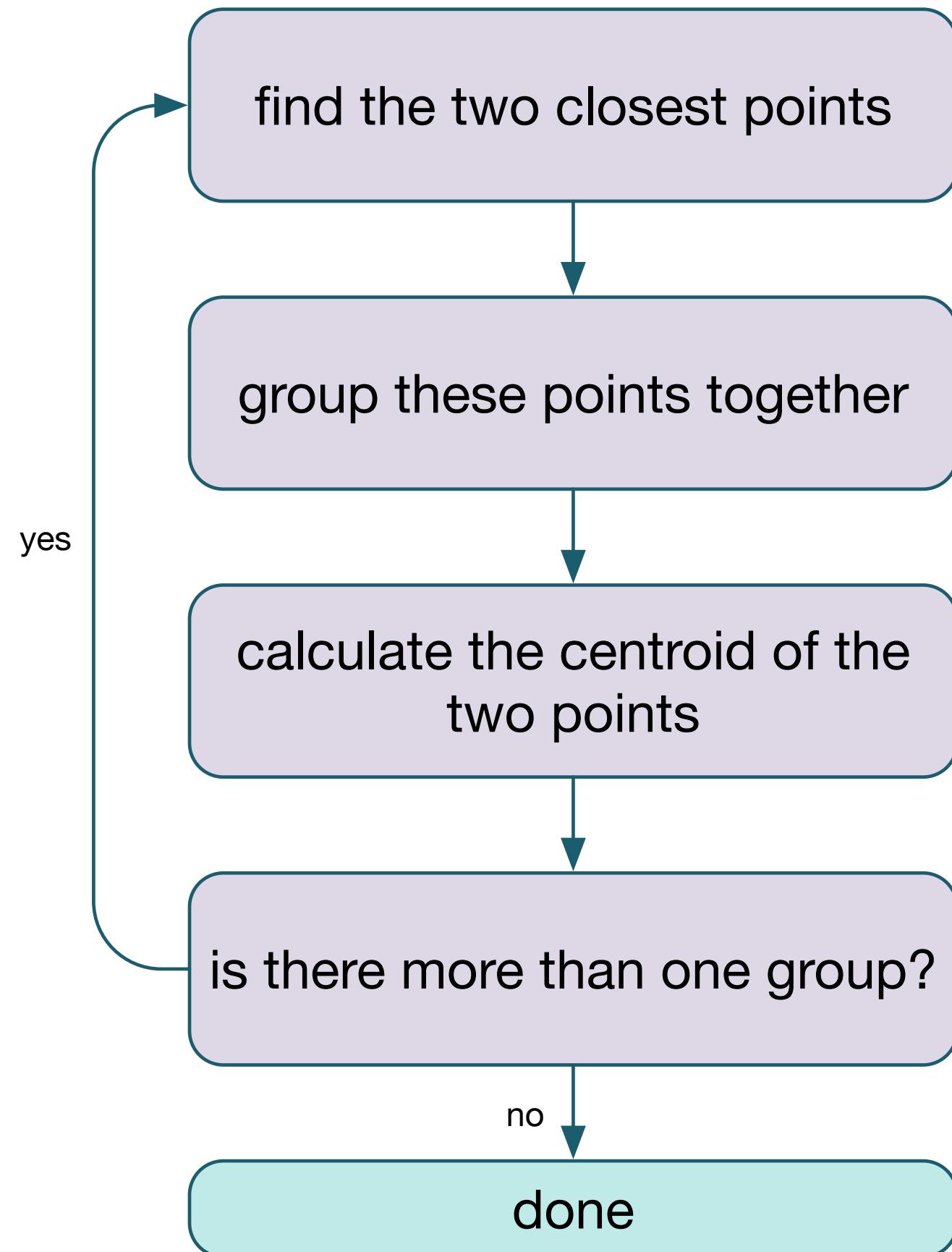


k-Means Clustering





Hierarchical Clustering





Encoding

Encoding data changes the labels to make it easier for computers to work with

Encodings need to be unambiguous, consistent and (as much as possible) complete

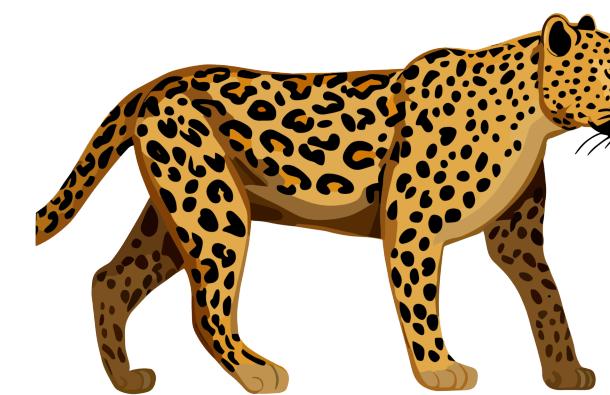
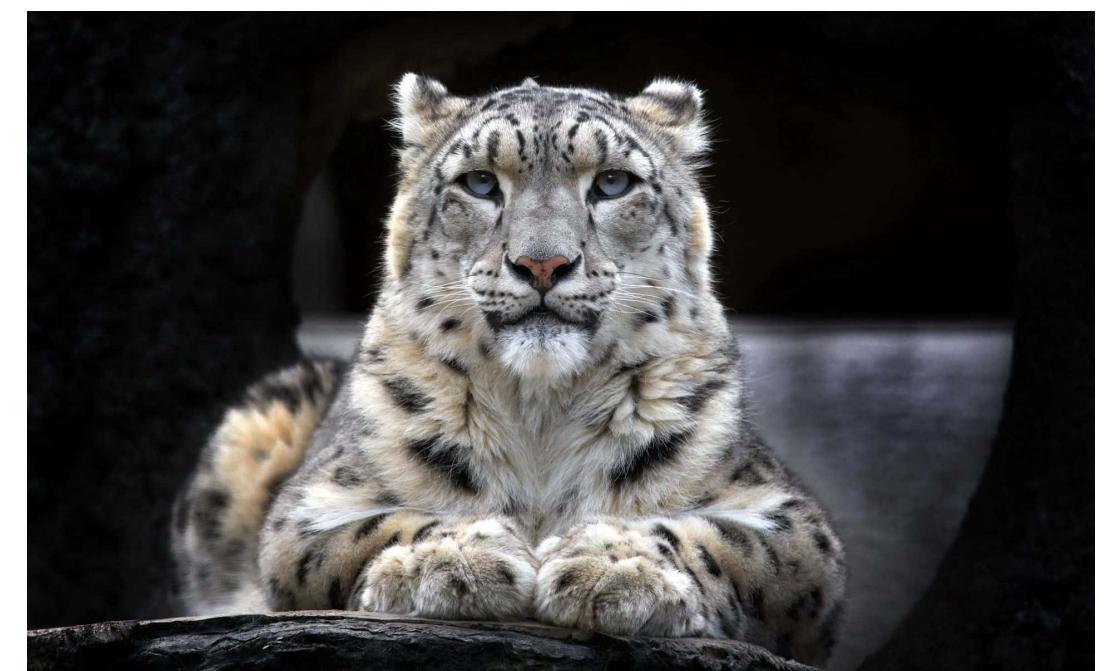
Karyotype	Name	Incidence (live births)	Encoding
XX	Female	1 in 2	1
XY	Male	1 in 2	2
X	Turner syndrome	1 in 4,000	3
XXX	Trisomy X	1 in 2,000	4
XXY	Klinefelter syndrome	1 in 1,000 – 1 in 2,000	5
XYY	Jacobs syndrome	1 in 1,000	6
XXXY	48,XXXY syndrome	1 in 100,000	7
XXYY	48,XXYY syndrome	1 in 80,000	8
XXXXY	49,XXXXY syndrome	1 in 200,000	9
other	–	–	10
unknown	–	–	0



Real Data Are Messy

Many problems are not simple to solve using the above techniques

Which of these images show leopards?





Dimensionality Reduction

When performing experiments, we measure multiple variables for each of multiple samples

- Expression of multiple genes
- Fluorescence of multiple markers

We wish to make statements about sample groups, taking into account all of the variables

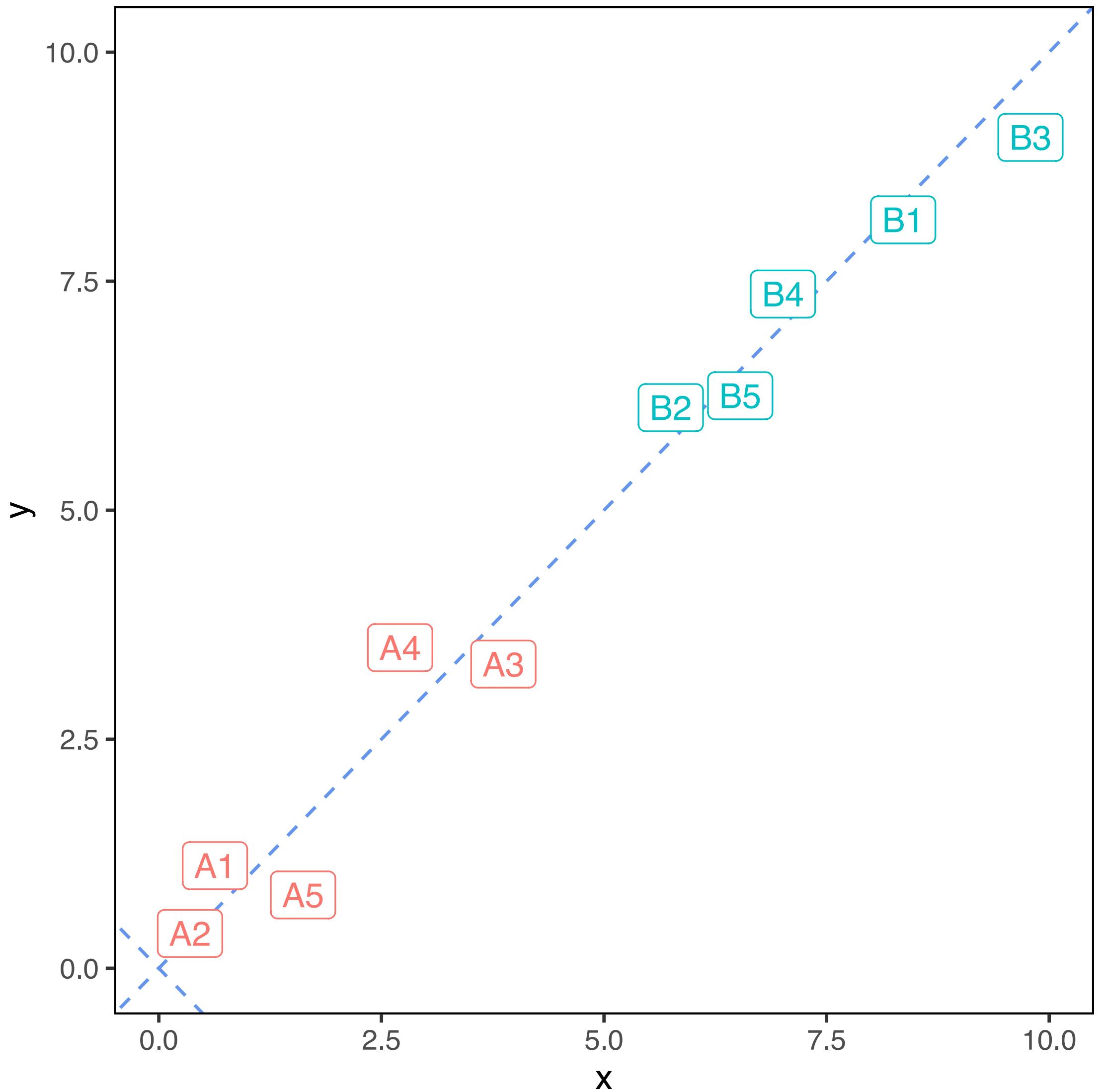
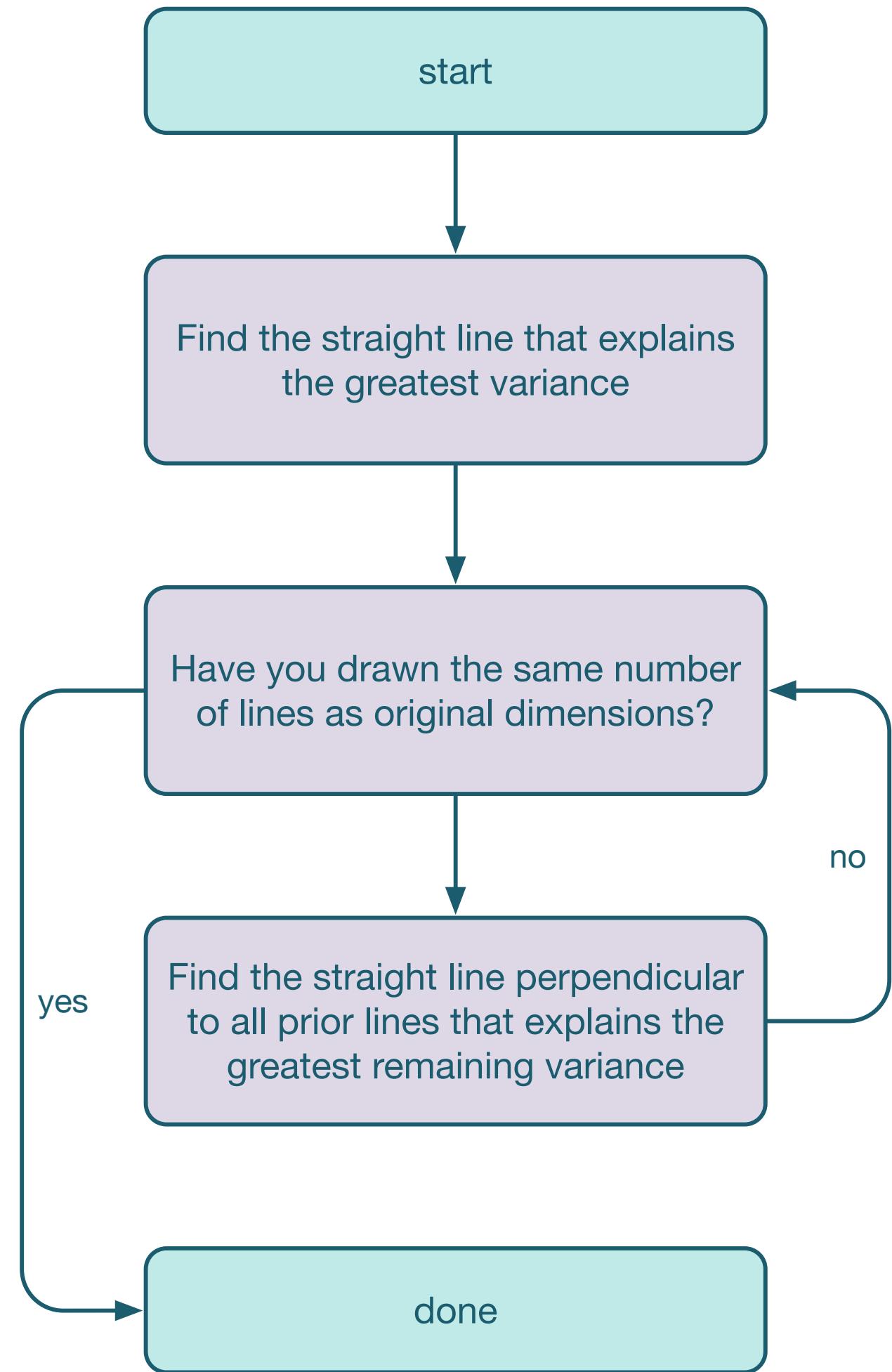
The variables we've measured often don't correspond to the features we care about

For example:

- If we wish to identify UK cities, we might record latitude and longitude.
- We've used two variables to describe a single hidden variable (which city)



PCA



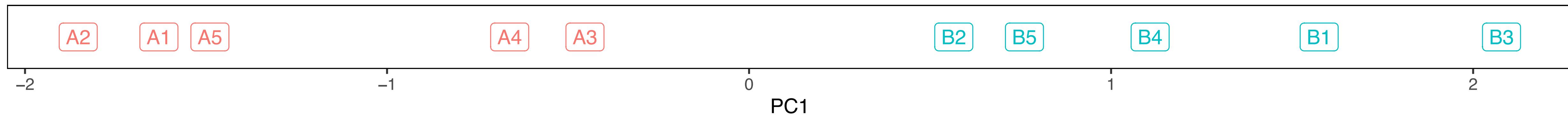


Each new dimension we get from PCA is a transformation of the original coordinate space

But, PCA ranks these new axes by the amount of variance they explain

This means that we can decide how many principle components to use

In our case, one dimension will do



Note that we've lost some information in doing this, but often this is OK (or even desired)



Another approach is to rescale our data so that the distances in a lower dimension are the same as those in the higher-dimension space

This is Multidimensional scaling (MDS)

A2 A1 A5

A4 A3

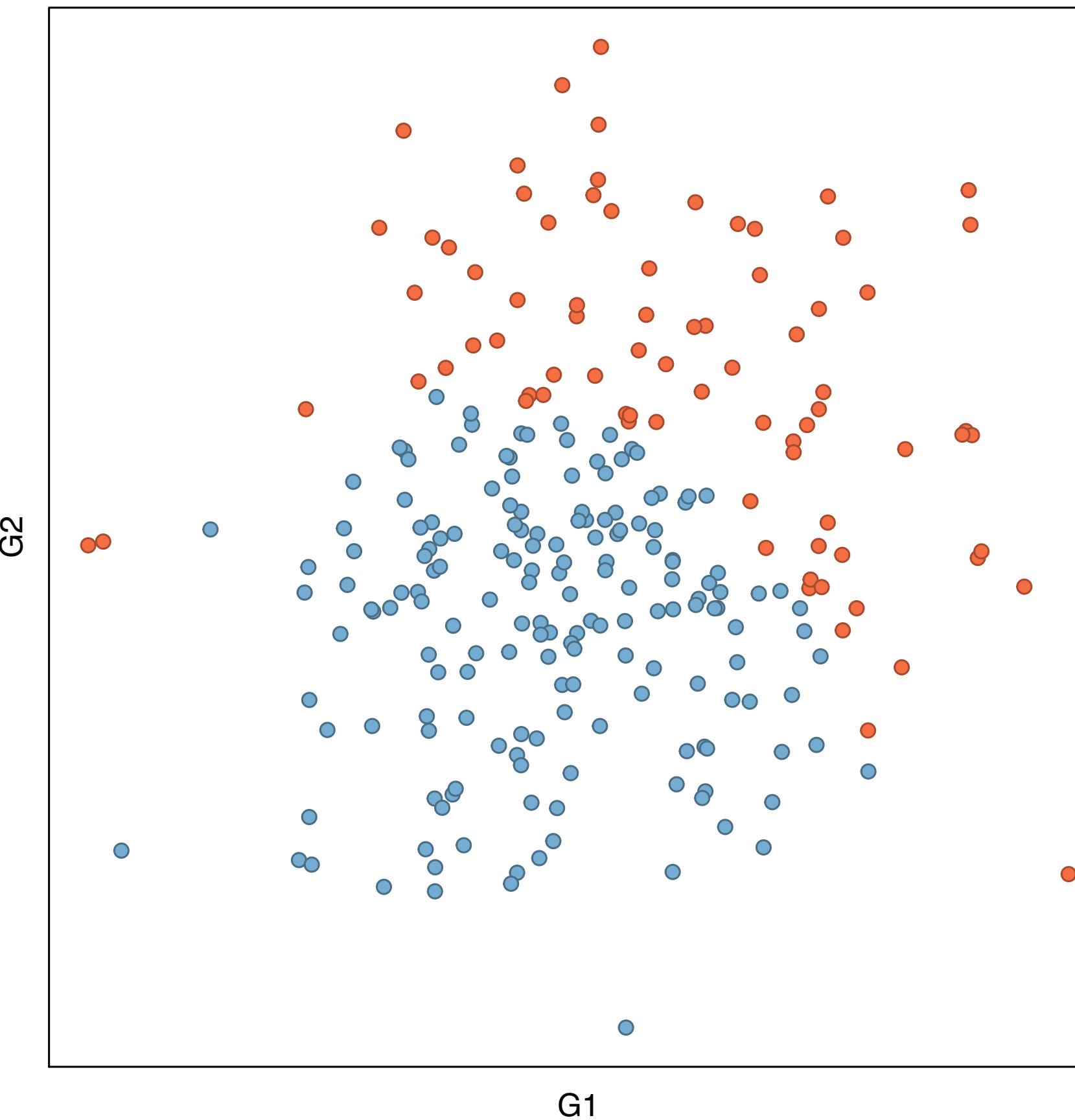
B2 B5 B4

B1 B3

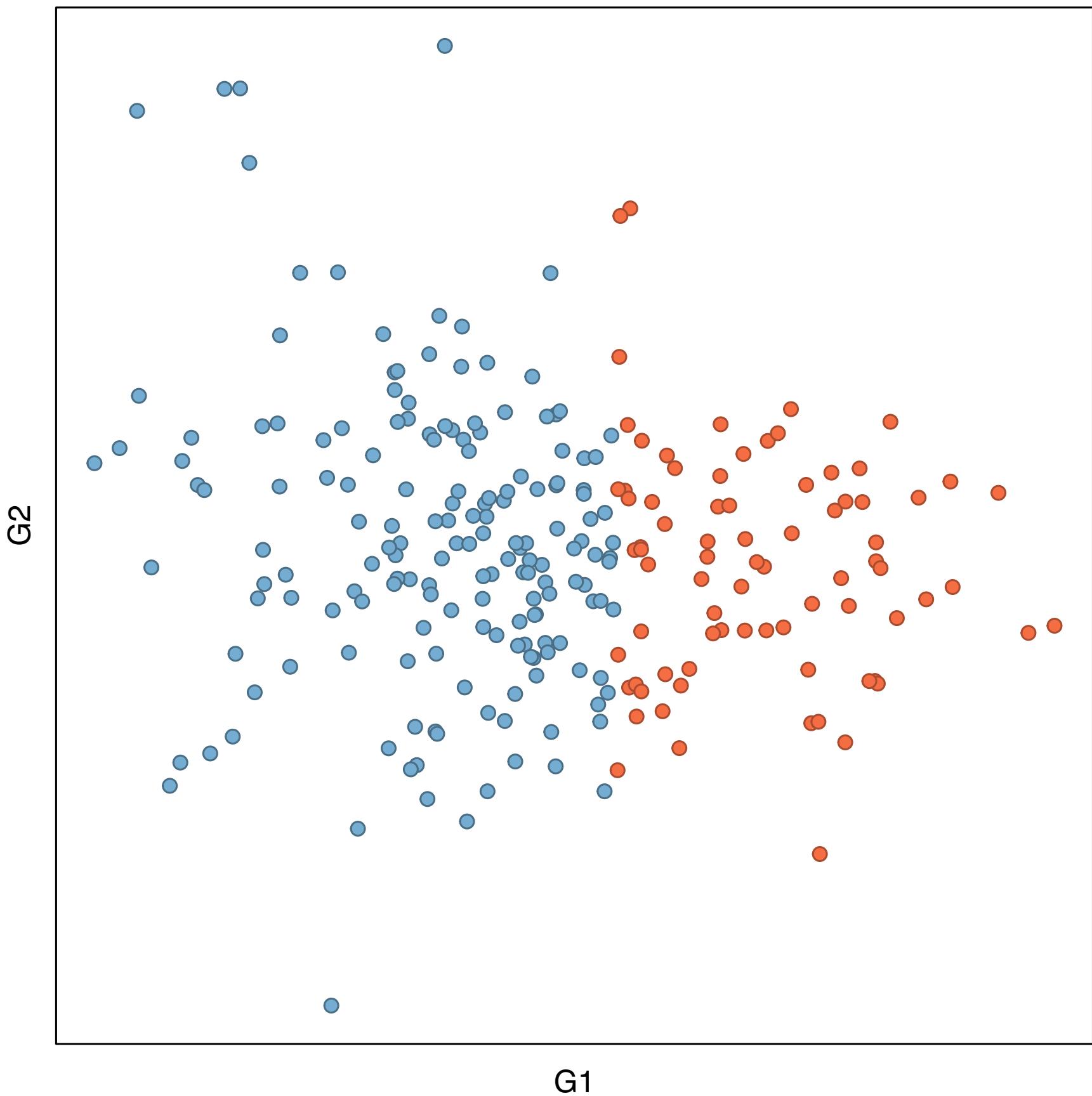


Transformation

We can transform our data before starting to analyse it to make the maths easier



Polar coordinate transform





Neural Networks

Neural Networks are a family of methods for classifying complex datasets

Allow us to make sense of complex and messy data

Originally inspired by biological neural networks

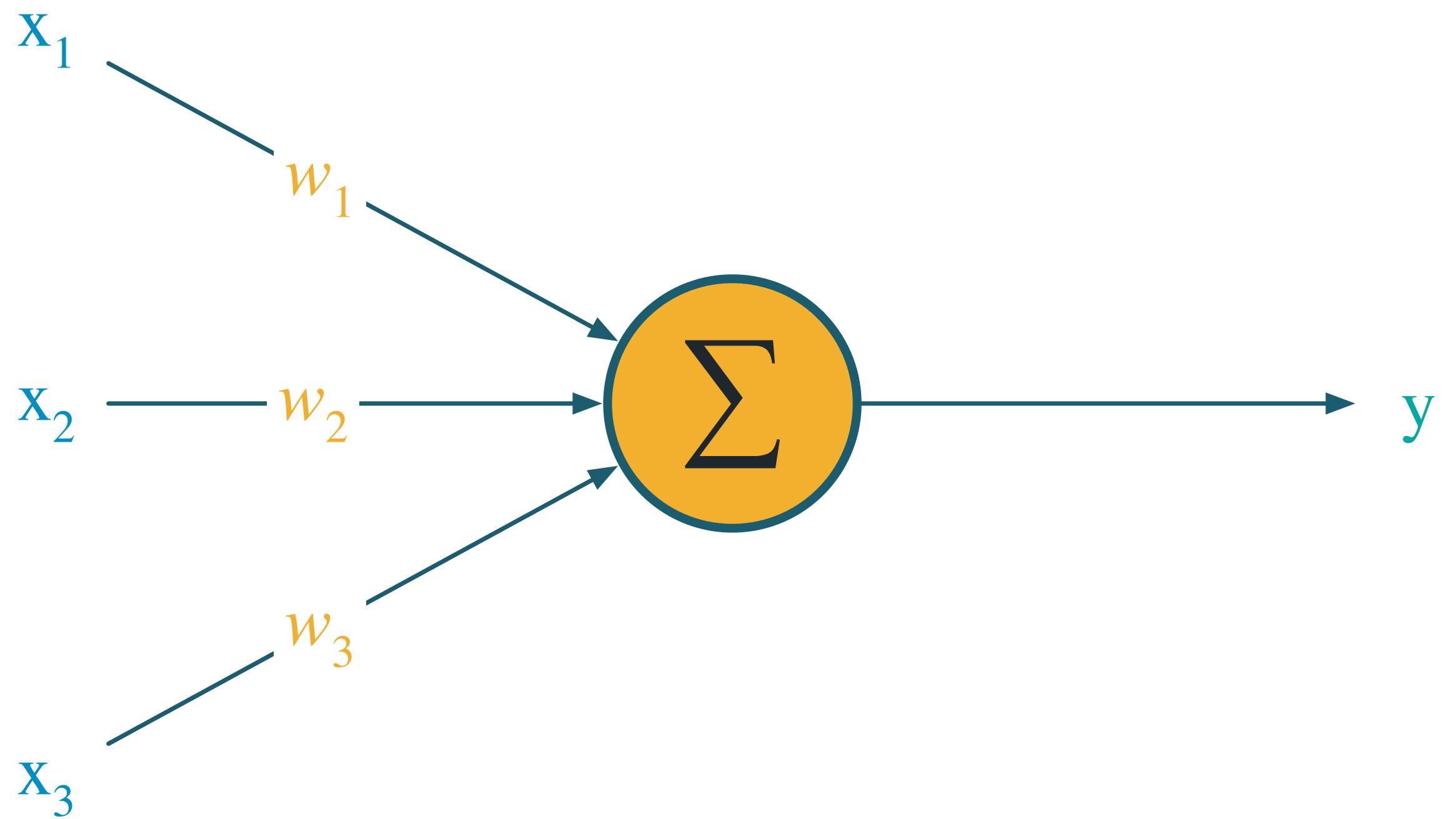
Individual neurones (*perceptrons*) process their individual inputs based on learned weights

Allows us to train classifiers without knowing how the classifiers work



Perceptrons

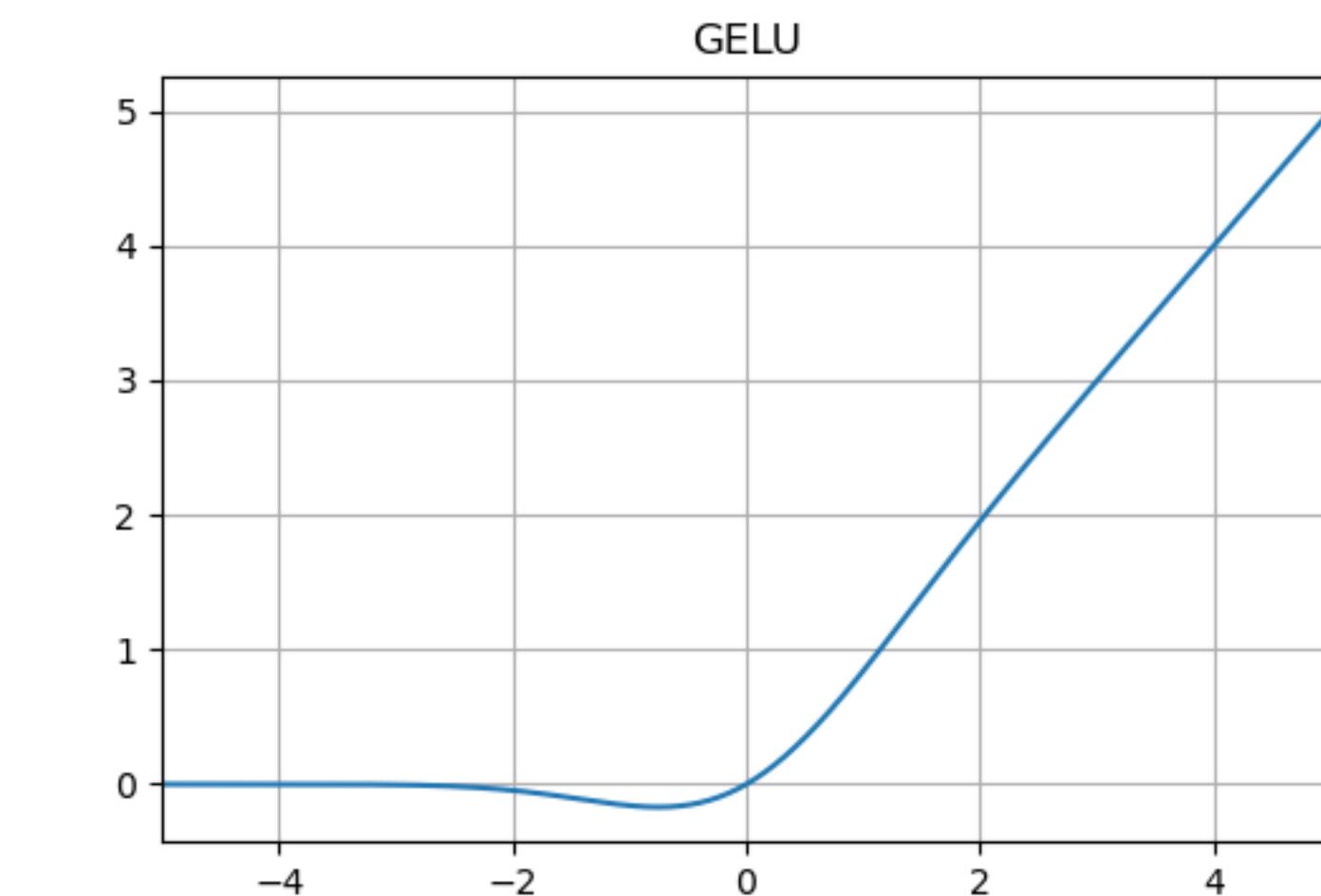
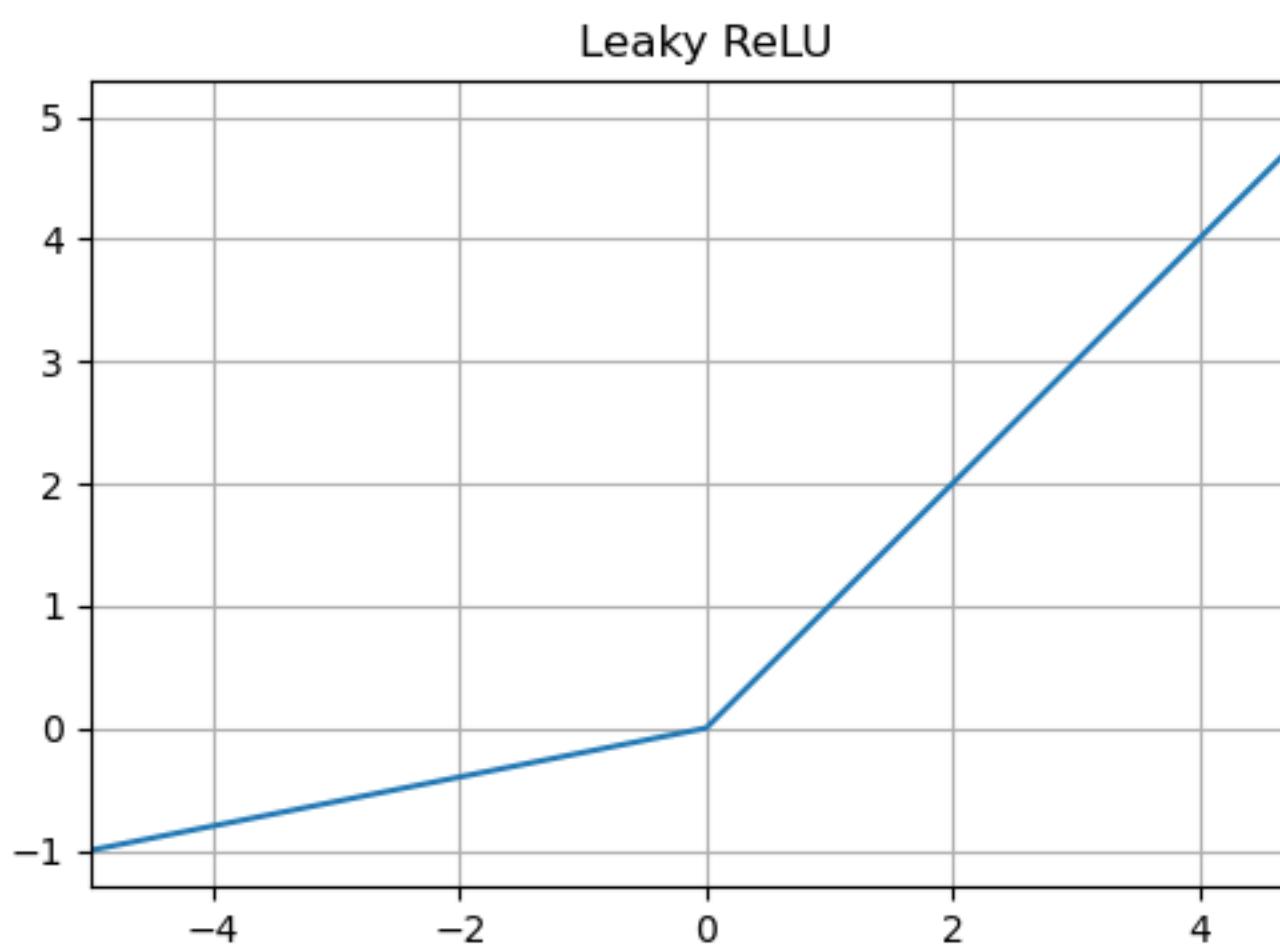
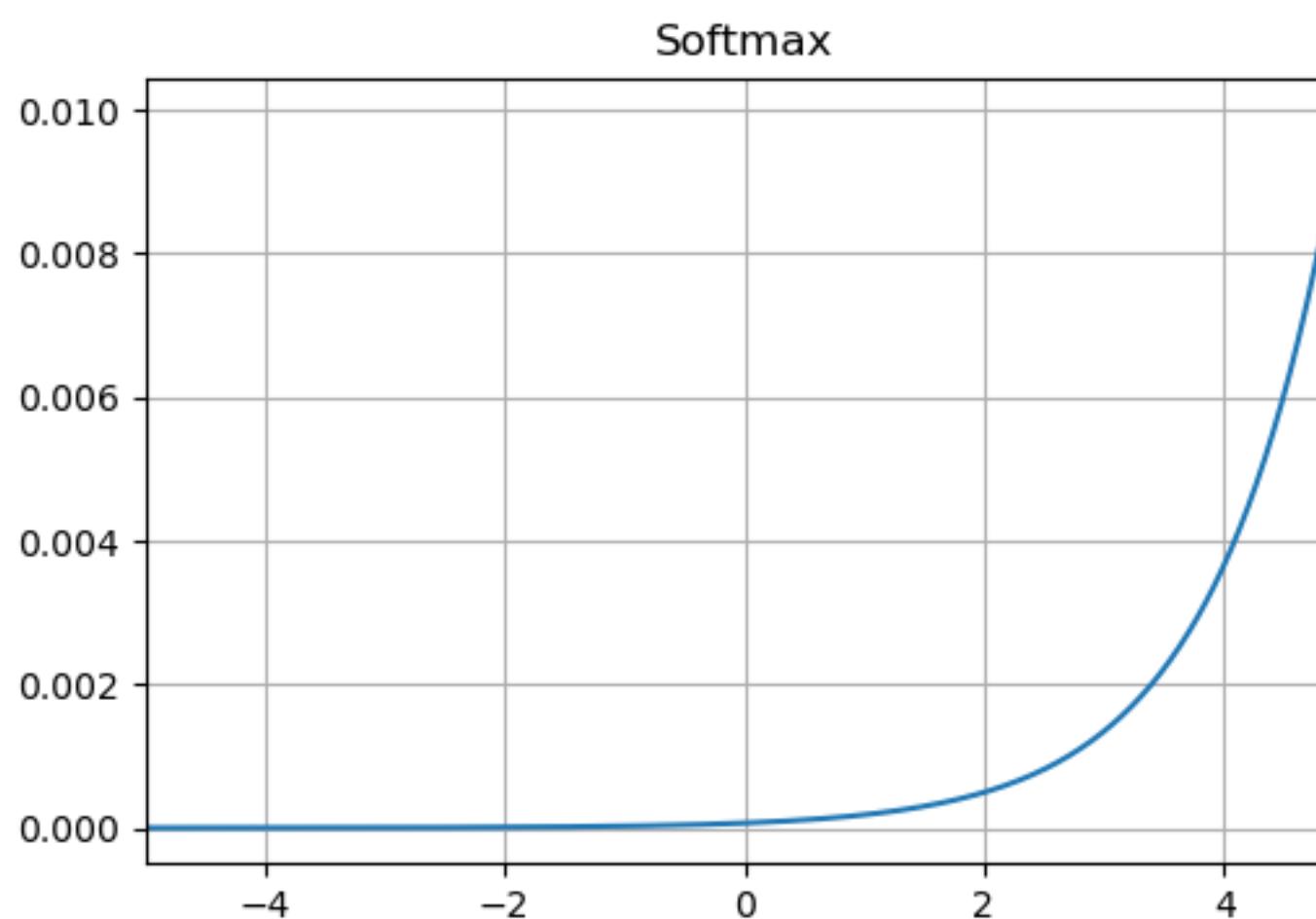
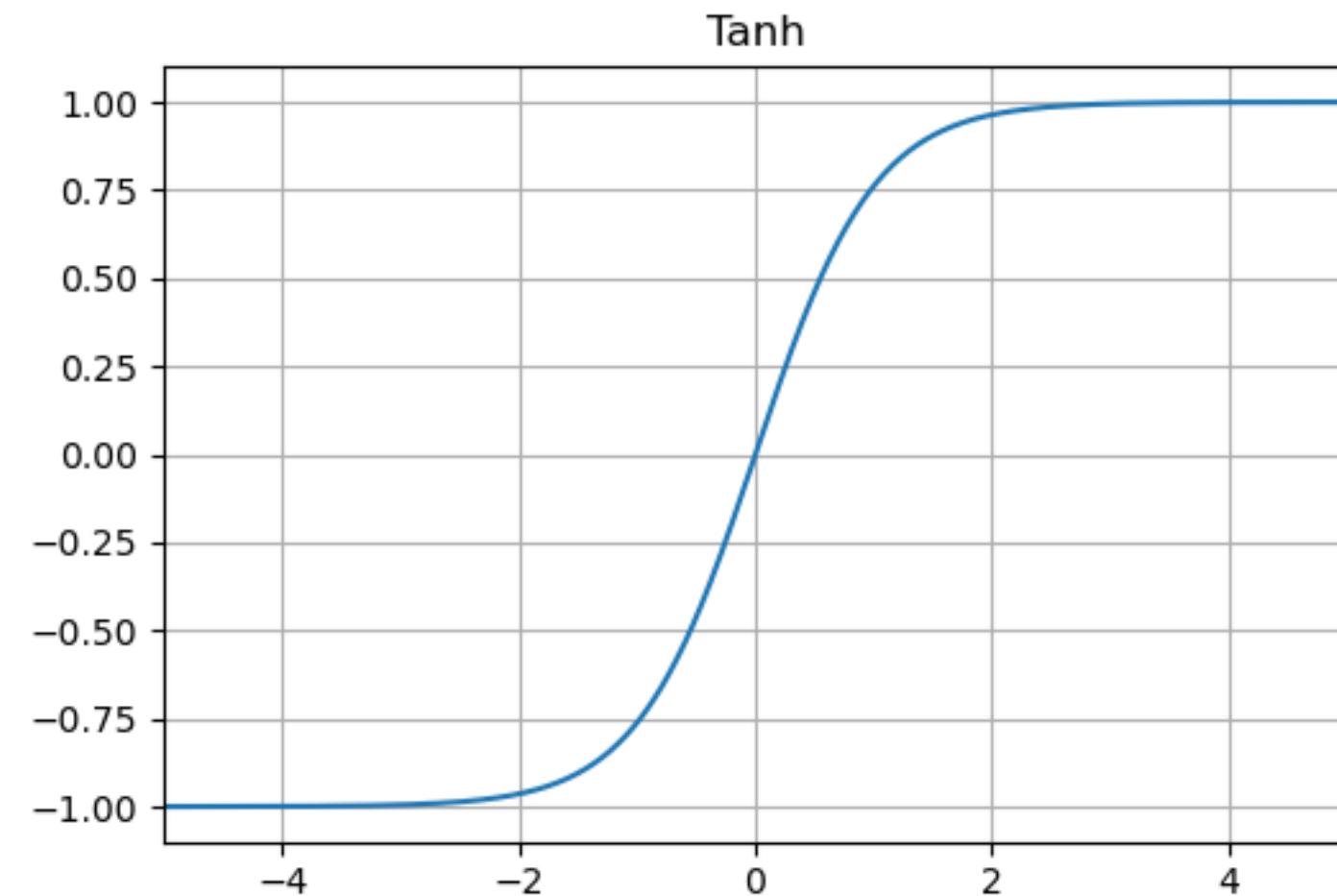
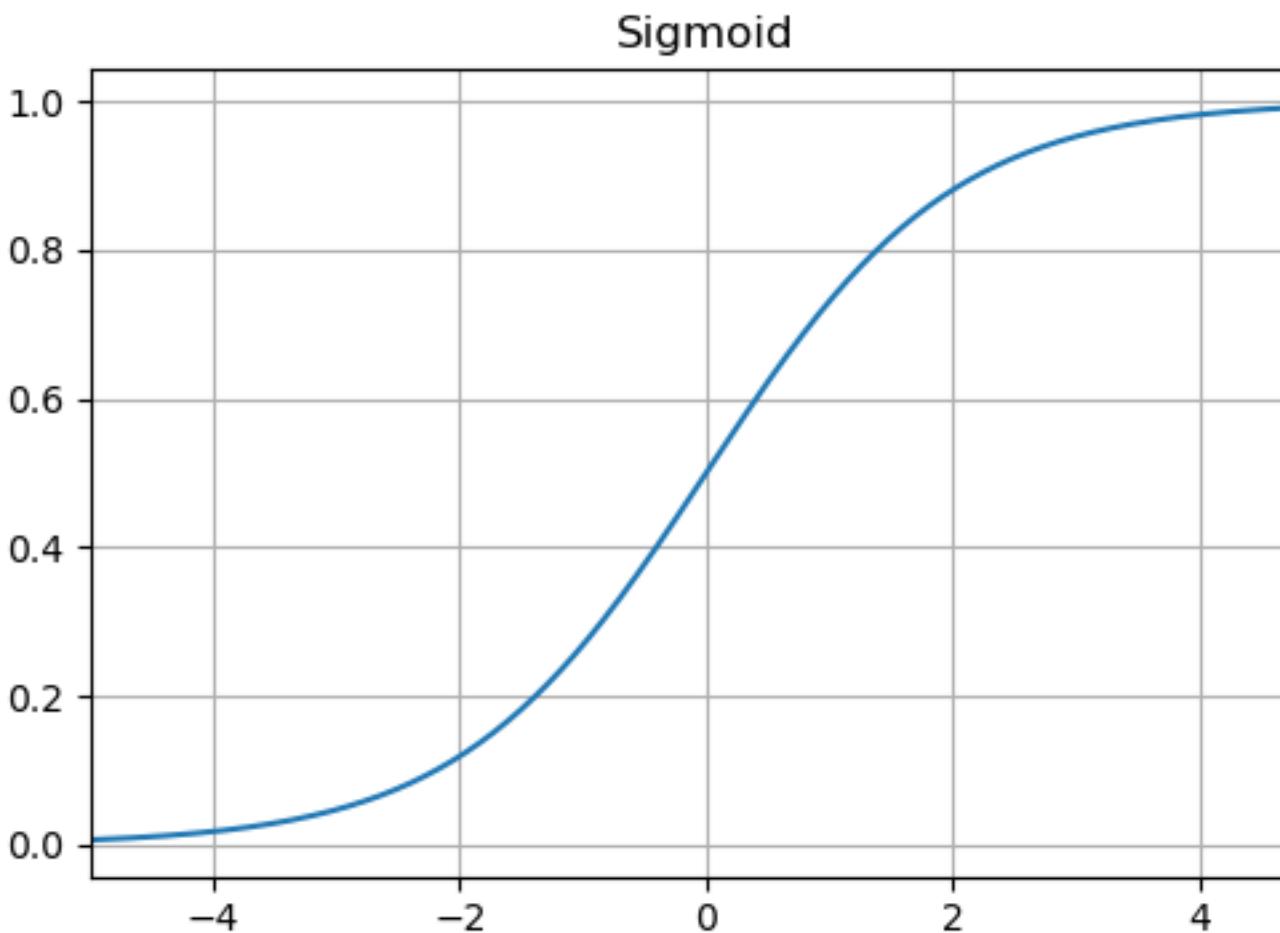
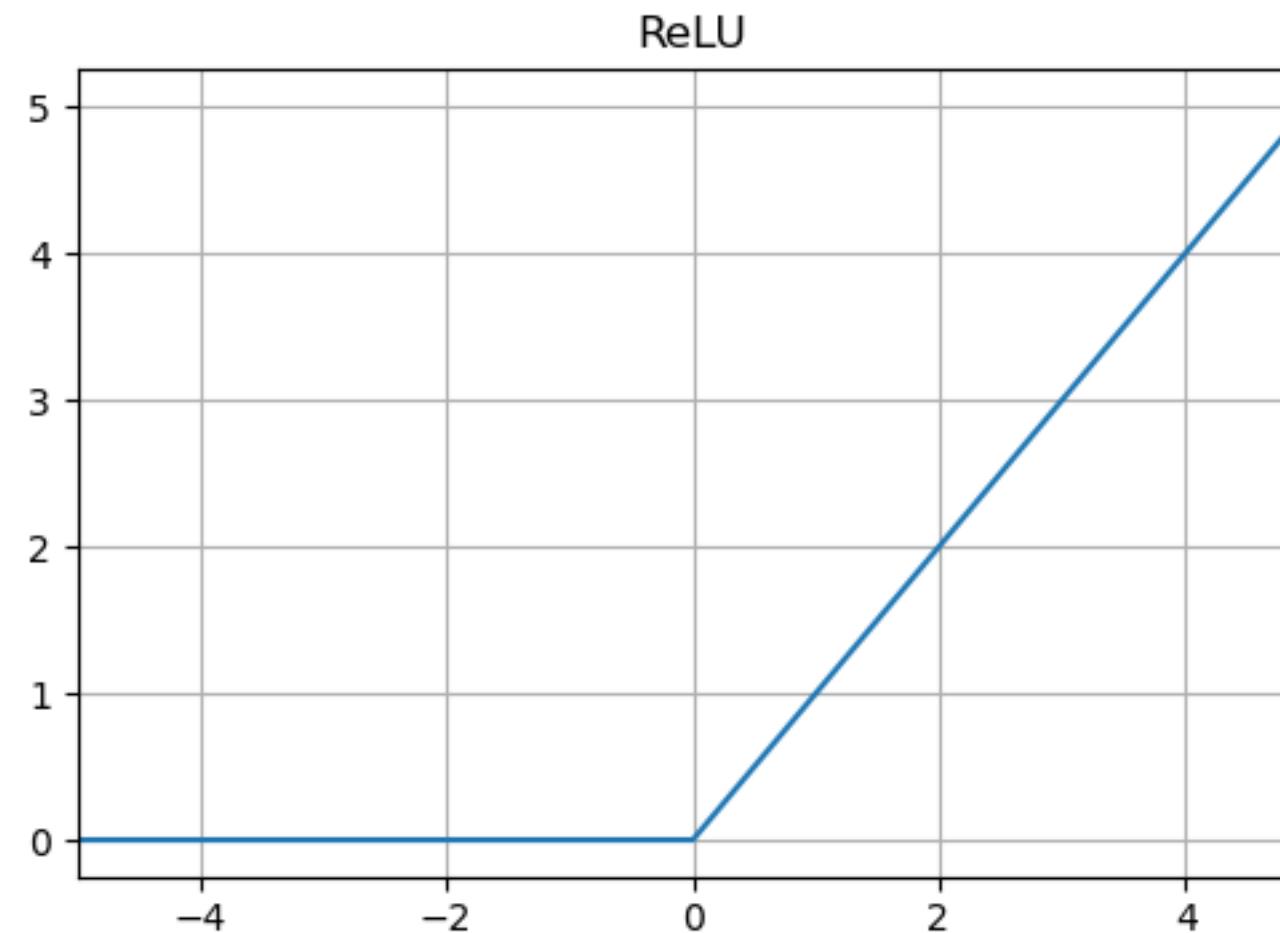
A perceptron is a function that sums its weighted inputs



$$y = f \left(\sum_{i=1}^n w_i x_i + b \right)$$



Activation Functions





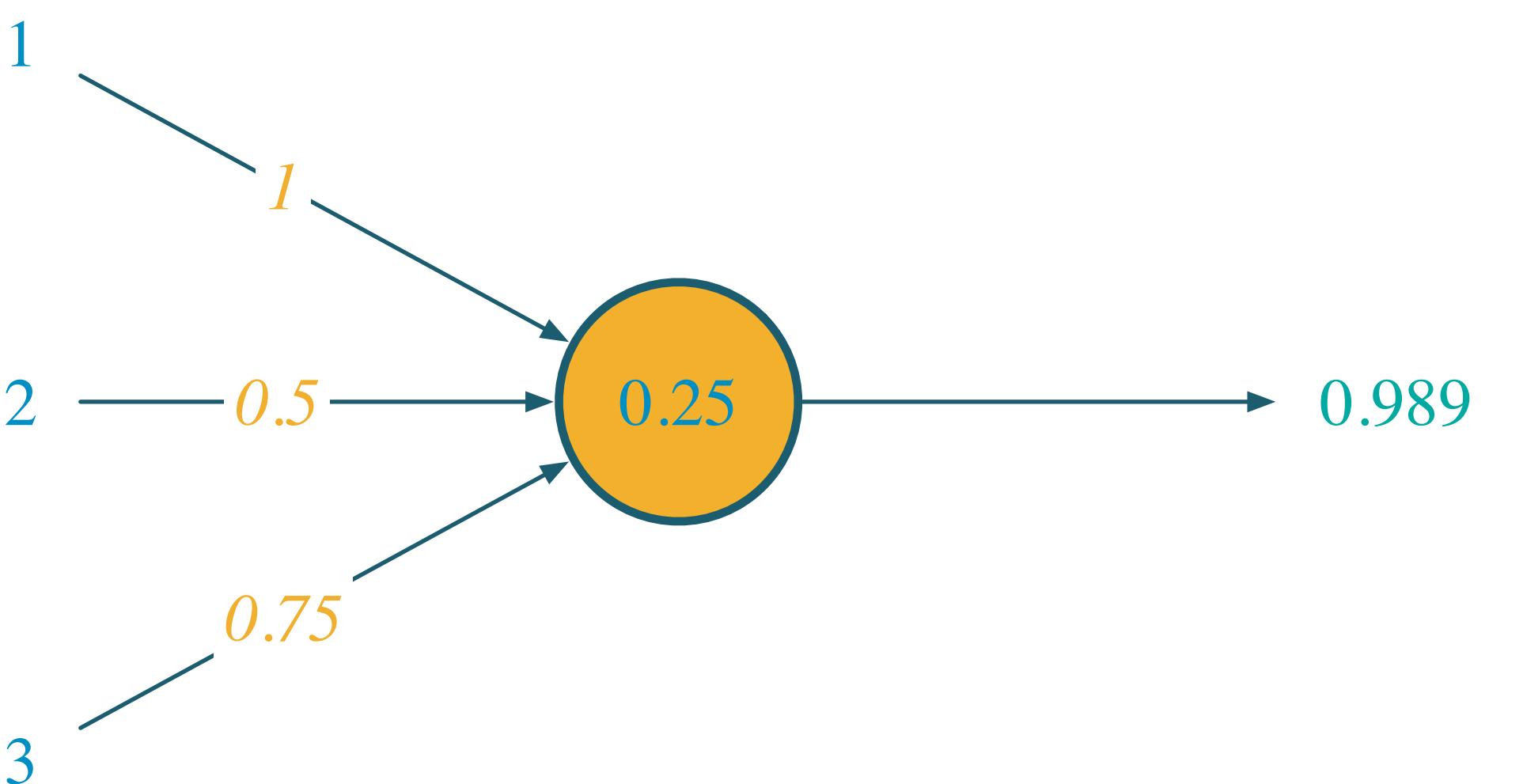
Examples

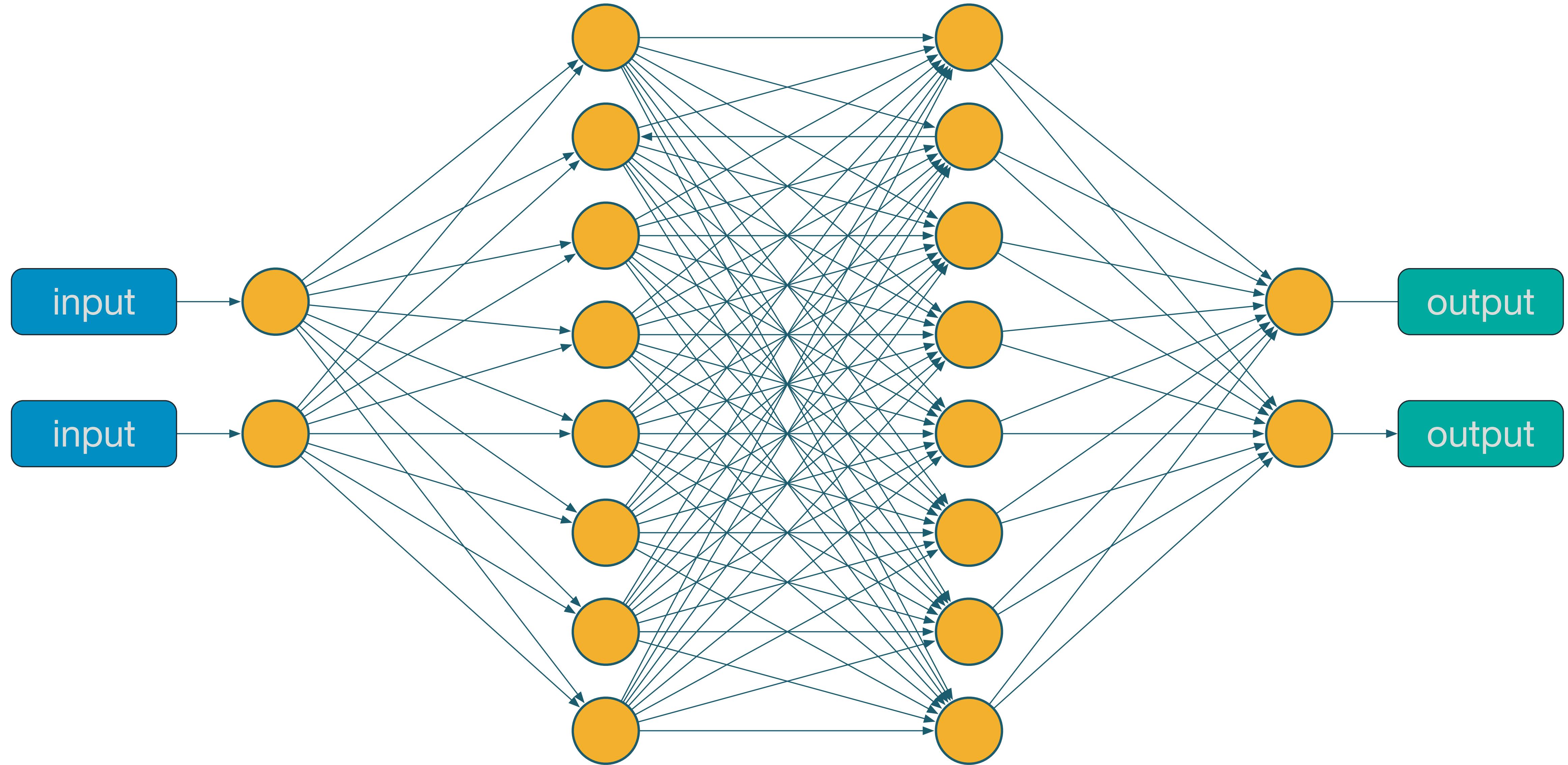
3 weighted inputs

Bias $b = 0.25$

Activation function is sigmoid

Resulting value is 0.989







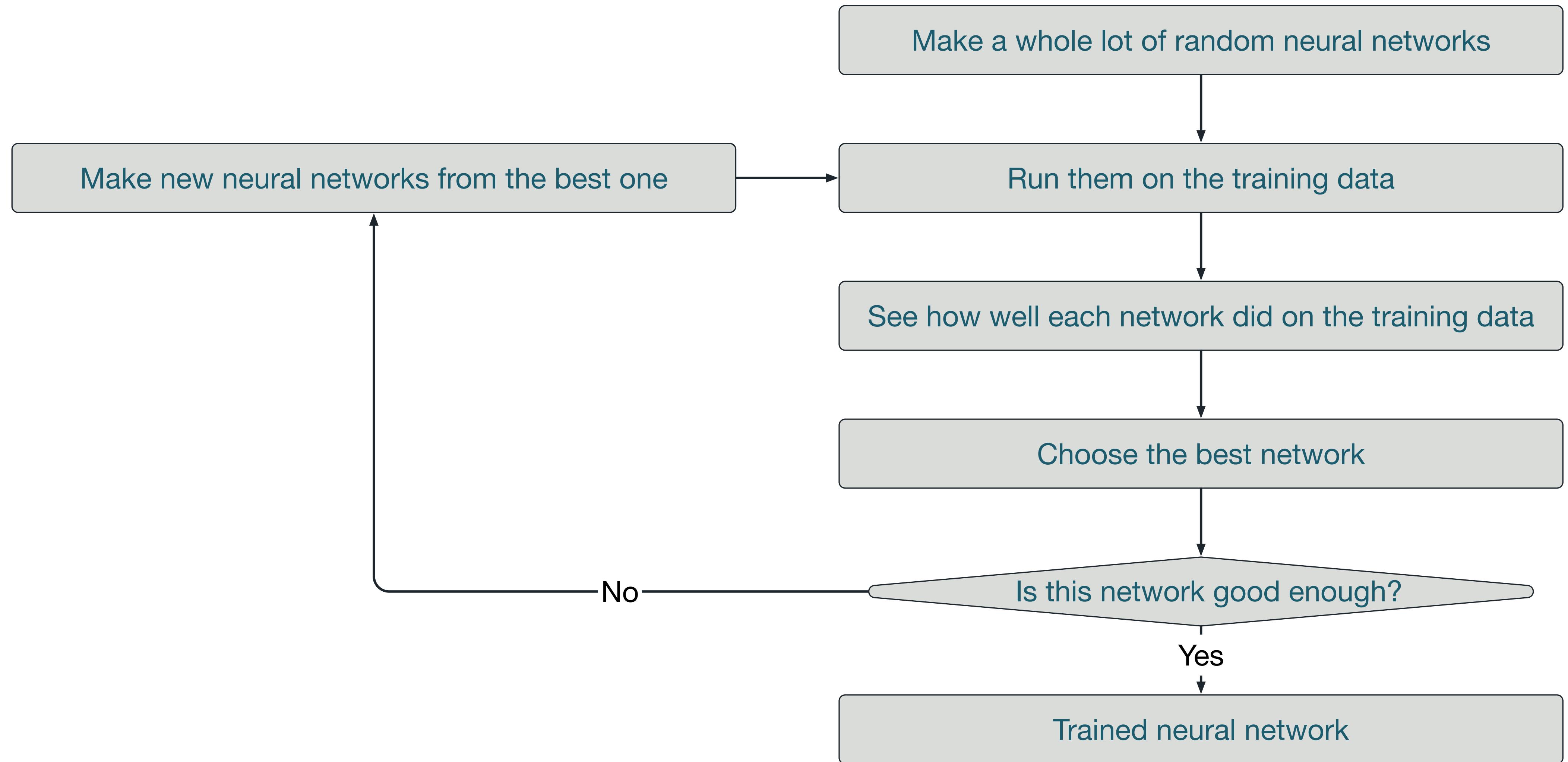
Training Neural Networks 1

To train a neural network, we use training data

Training data are simply data where we know the inputs and also the outputs we want to get

Training involves randomly changing weights and then hoping the new weights are better than the old ones

There is a risk of overfitting





Tweak the perceptron weights until the neural network gives the right answer

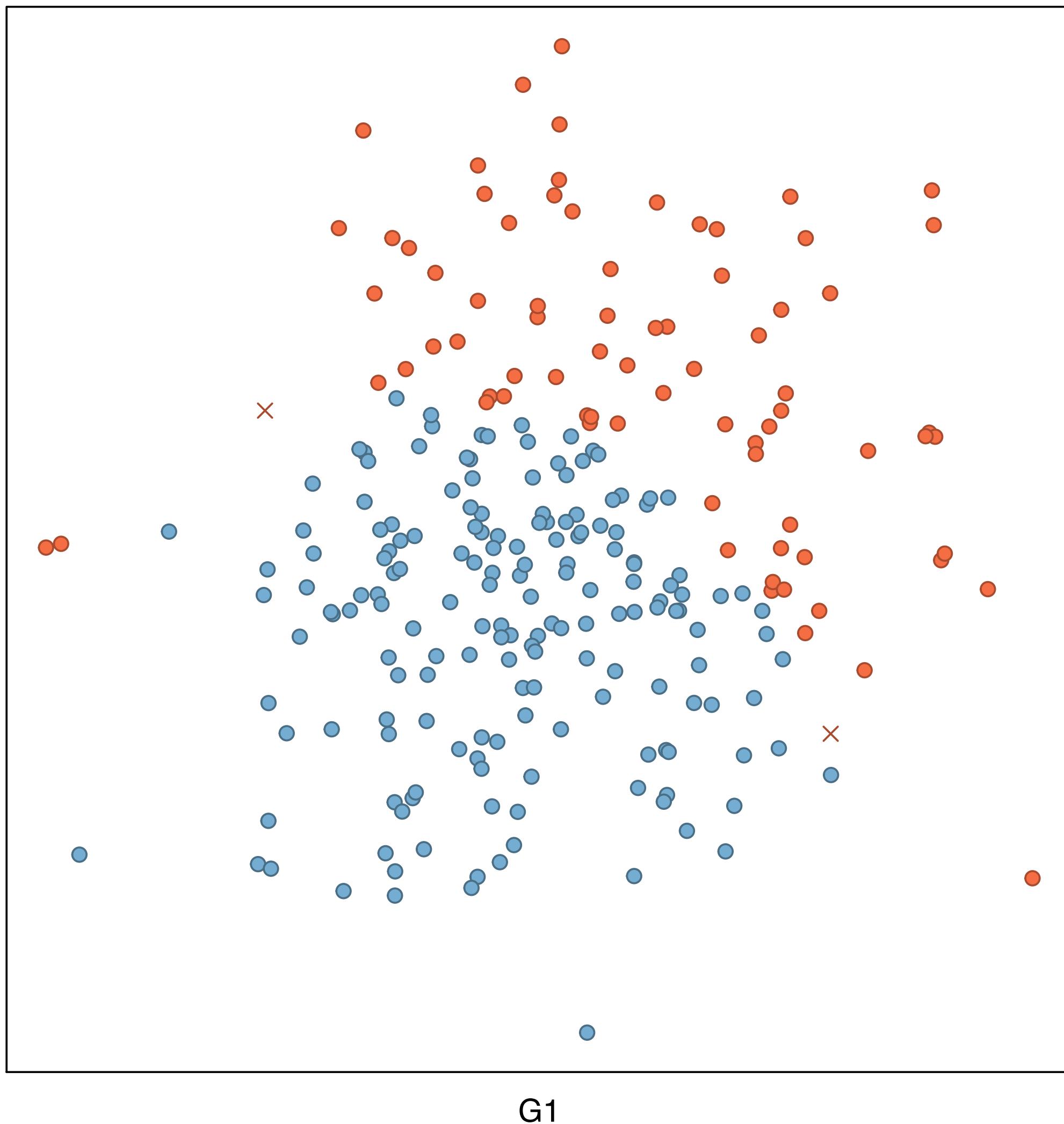
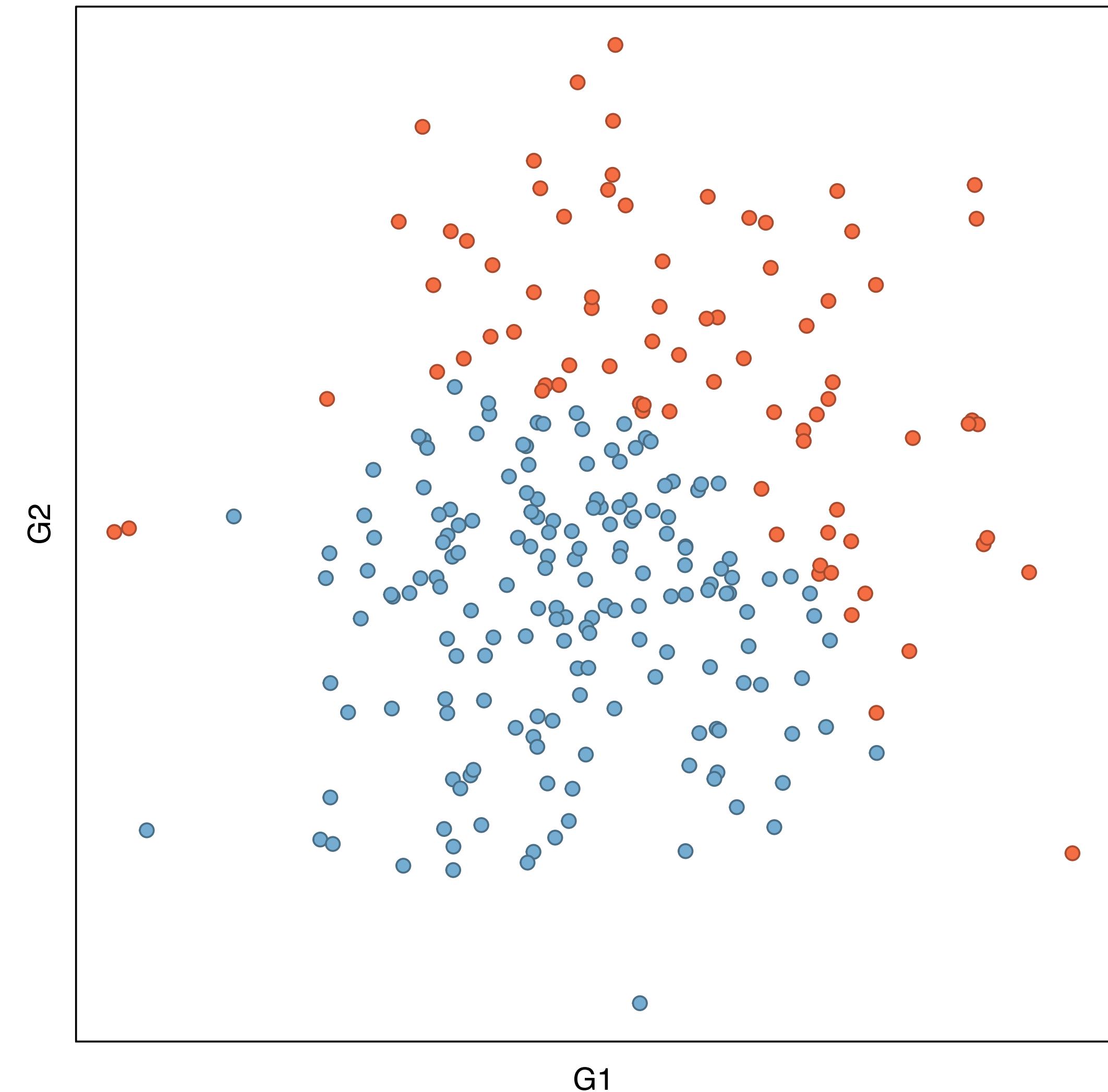
The trouble is, there are *a lot of weights*

In our example, there are 96 different weights, each of which can be any number

As we don't control what the neural network decides to learn, it is crucial that we:

- Use appropriate training data
- Use enough (lots) training data

It is very easy to underestimate how many training data points we need

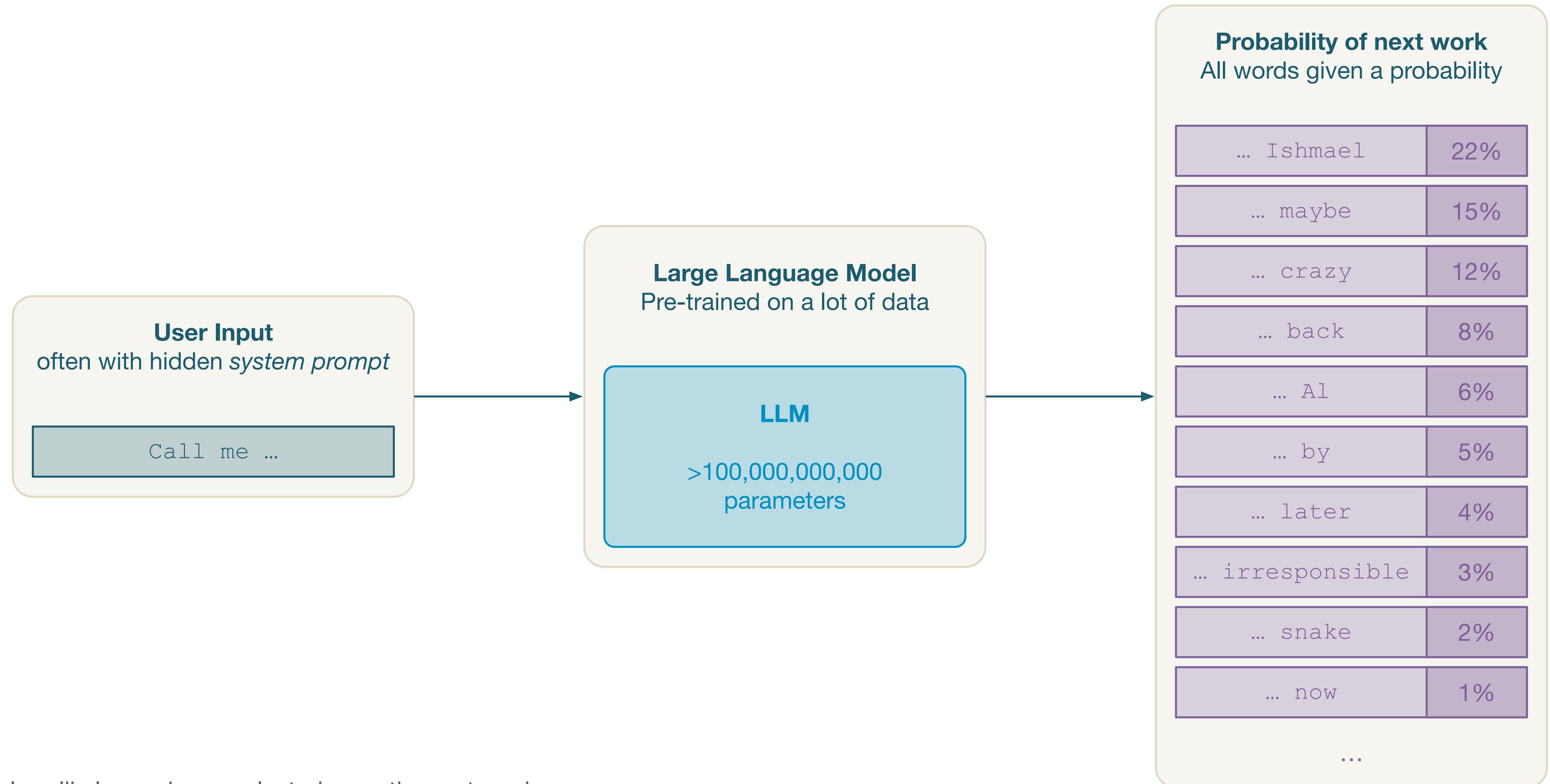


We train our neural network on 10,000 training points, then make it guess on our real data

After all that, it still gets two points wrong!

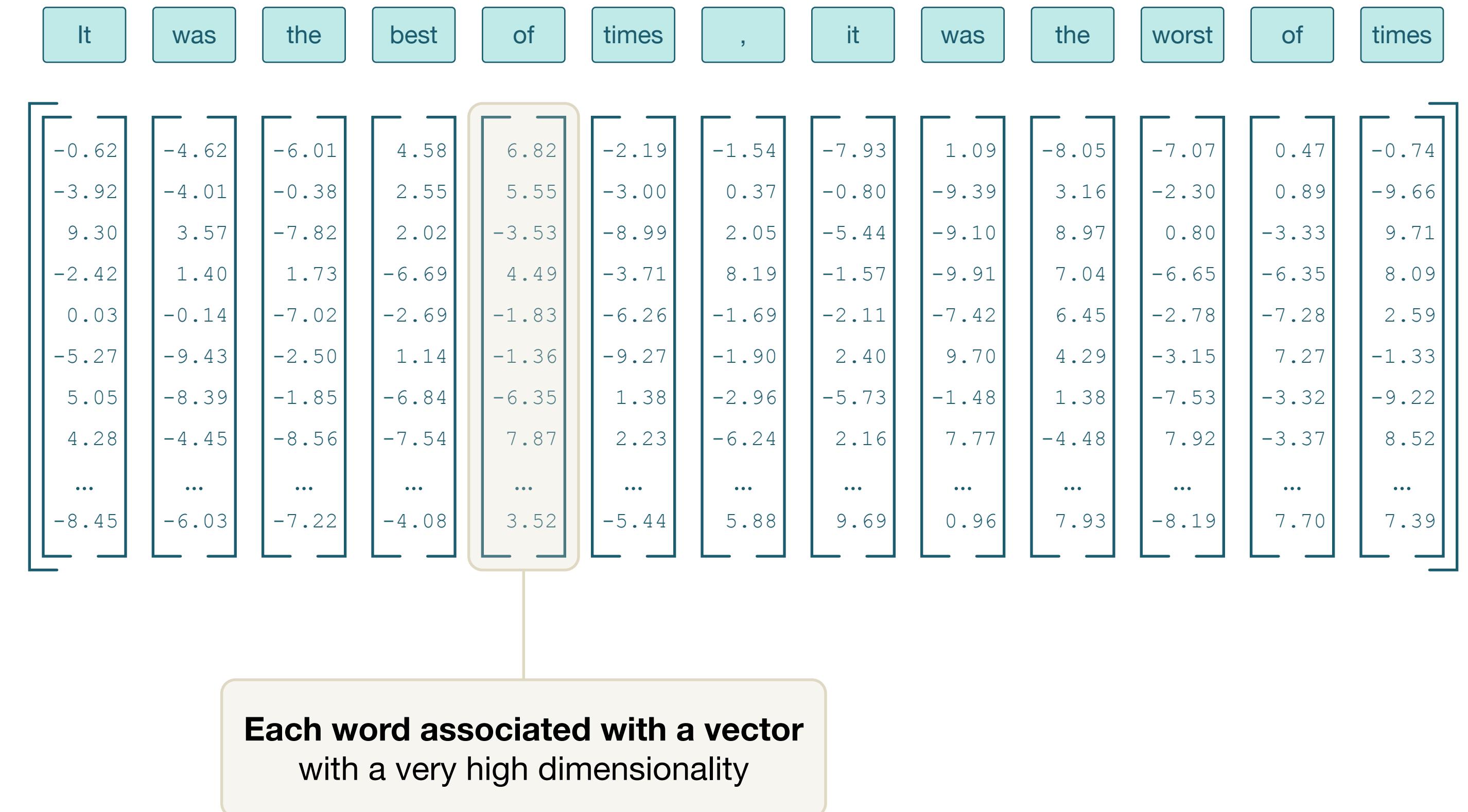


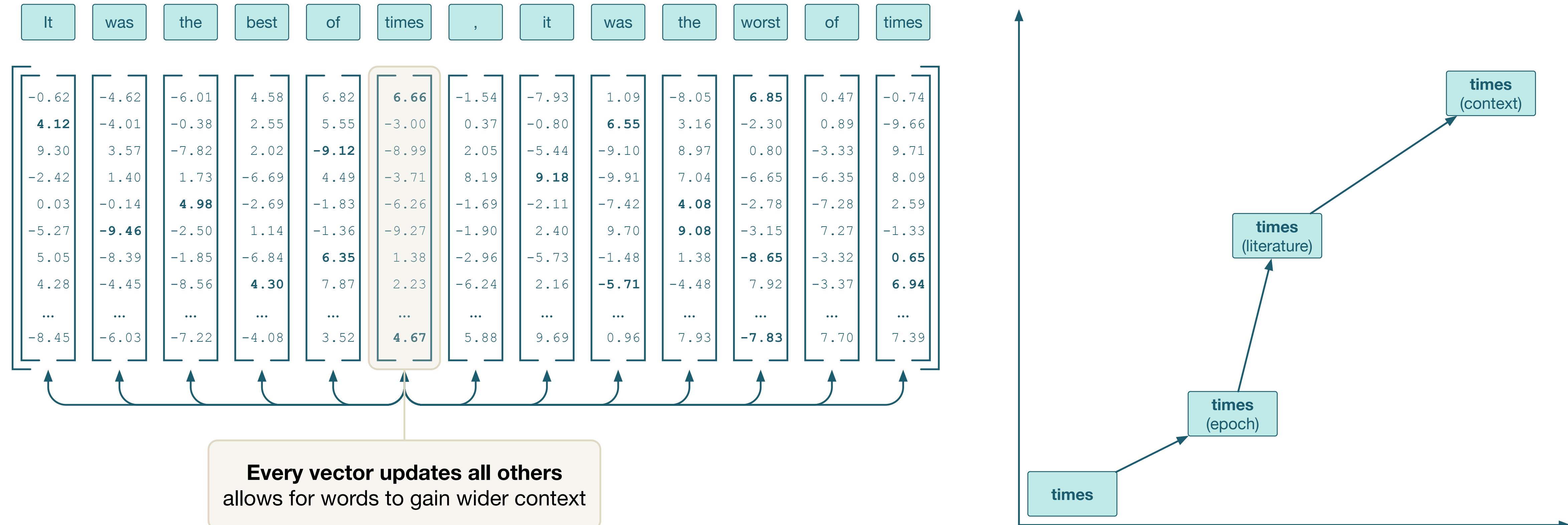
Large Language Models





Transformers



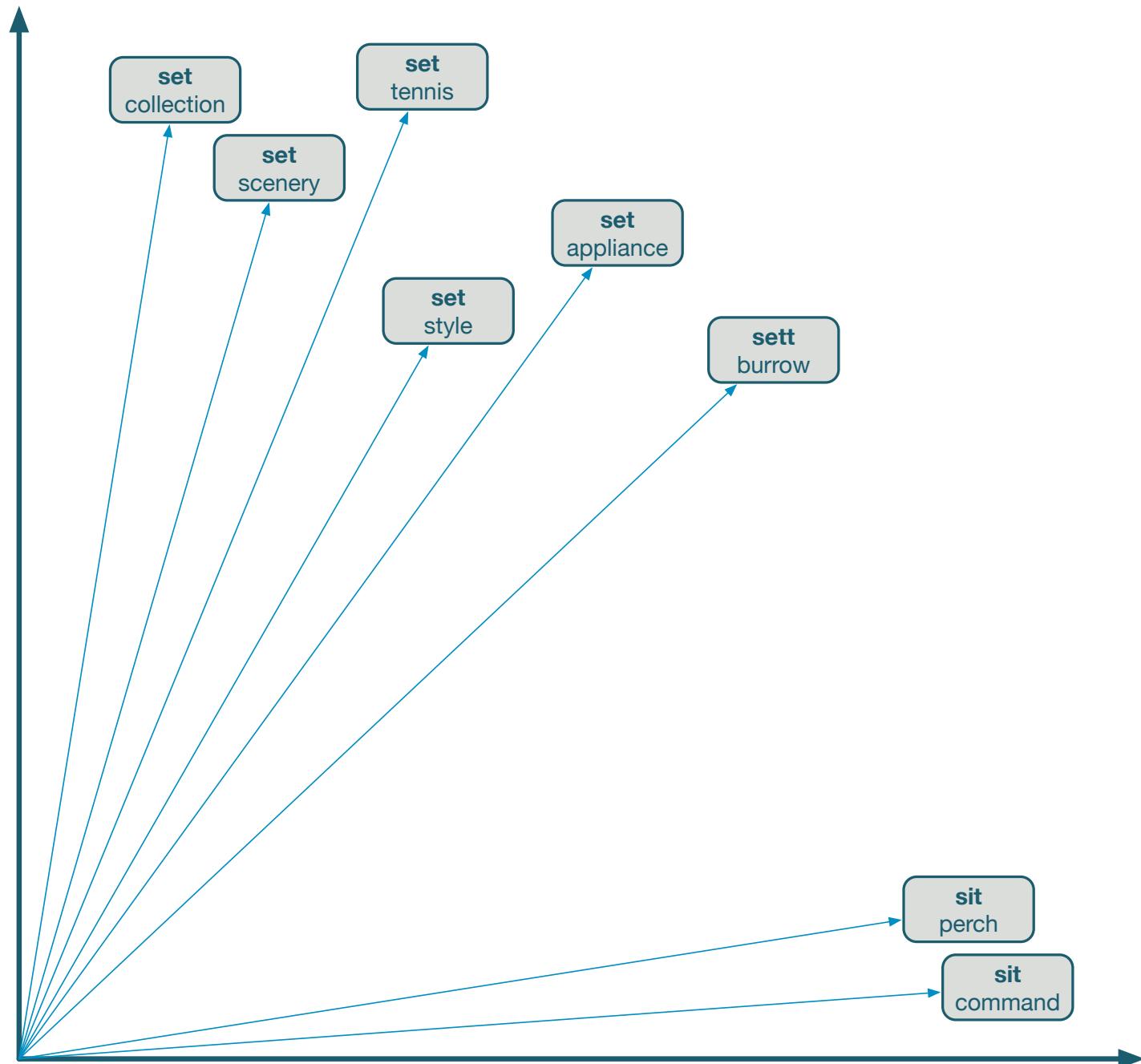




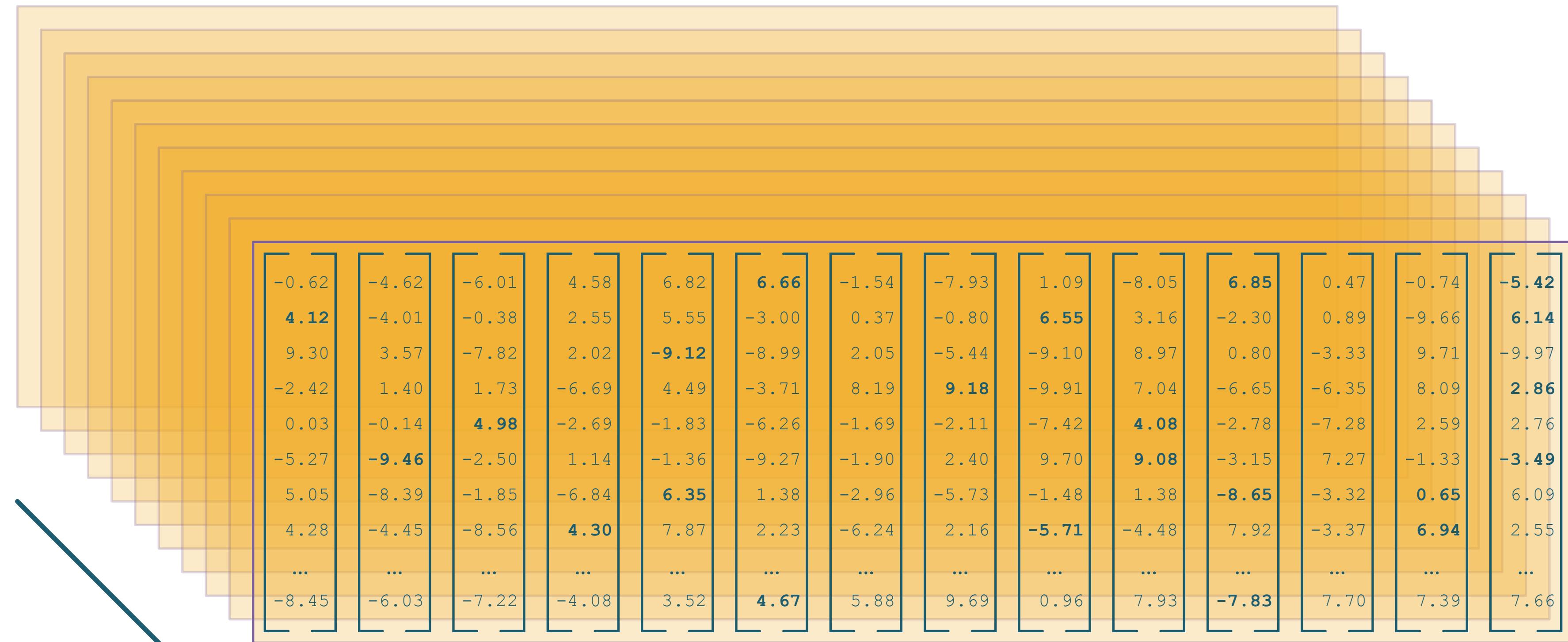
The mapping of a word to a vector is an *Embedding*

Similar terms will appear close in higher-dimension space

Embeddings *are* mathematical vectors



$$E(\text{man}) - E(\text{woman}) + E(\text{uncle}) \approx E(\text{aunt})$$



Vectors fed into multiple layers
each layer takes input from the last

It	was	the	best	of	times	,	it	was	the	worst	of	times	?
-0.62	-4.62	-6.01	4.58	6.82	-2.19	-1.54	-7.93	1.09	-8.05	-7.07	0.47	-0.74	1.83
-3.92	-4.01	-0.38	2.55	5.55	-3.00	0.37	-0.80	-9.39	3.16	-2.30	0.89	-9.66	-8.24
9.30	3.57	-7.82	2.02	-3.53	-8.99	2.05	-5.44	-9.10	8.97	0.80	-3.33	9.71	-9.97
-2.42	1.40	1.73	-6.69	4.49	-3.71	8.19	-1.57	-9.91	7.04	-6.65	-6.35	8.09	4.31
0.03	-0.14	-7.02	-2.69	-1.83	-6.26	-1.69	-2.11	-7.42	6.45	-2.78	-7.28	2.59	2.76
-5.27	-9.43	-2.50	1.14	-1.36	-9.27	-1.90	2.40	9.70	4.29	-3.15	7.27	-1.33	-2.41
5.05	-8.39	-1.85	-6.84	-6.35	1.38	-2.96	-5.73	-1.48	1.38	-7.53	-3.32	-9.22	6.09
4.28	-4.45	-8.56	-7.54	7.87	2.23	-6.24	2.16	7.77	-4.48	7.92	-3.37	8.52	2.55
...
-8.45	-6.03	-7.22	-4.08	3.52	-5.44	5.88	9.69	0.96	7.93	-8.19	7.70	7.39	7.66

Final vector used for prediction
after multiple rounds

Probability of next word

All words given a probability

... ,	85%
... .	8%
... ;	3%
... -	1%
... in	<1%
... for	<1%
... and	<1%
... during	<1%
... but	<1%
... when	<1%

...

