

# COMP34711

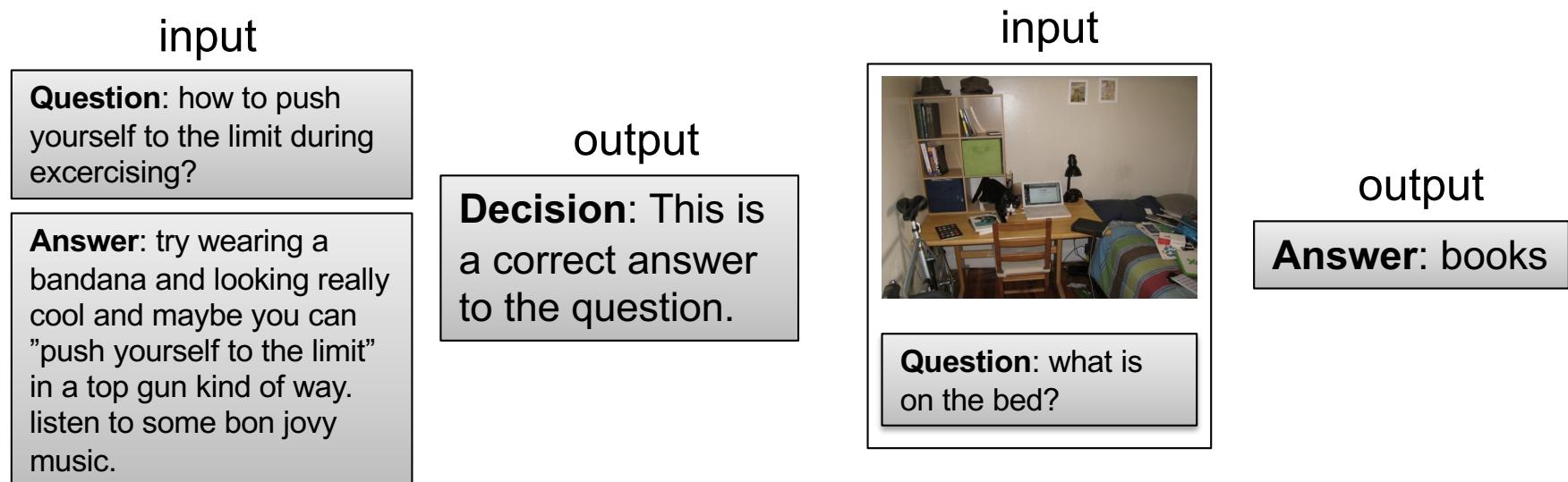
## Deep Learning for NLP I

Tingting Mu  
Department of Computer Science

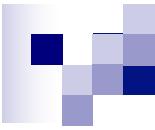
Slides are based on CS224d: Deep Learning for Natural Language Processing

# Deep Learning

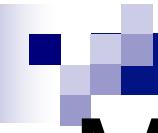
- Deep learning (DL) and neural network (NN) are often used interchangeably. The subtle difference is that DL emphasises networks with higher number of layers.
- DL mostly refers to NN based techniques for building end-to-end systems, which take raw objects as the input.



- Deep learning is the state of the art technique for many NLP tasks.  
Browse state-of-the-art in NLP: <https://paperswithcode.com/area/natural-language-processing>



## ■ *What can neural network do?*



# Model Sequence Probability

- Neural networks can be used to model probability distribution of a sequence.

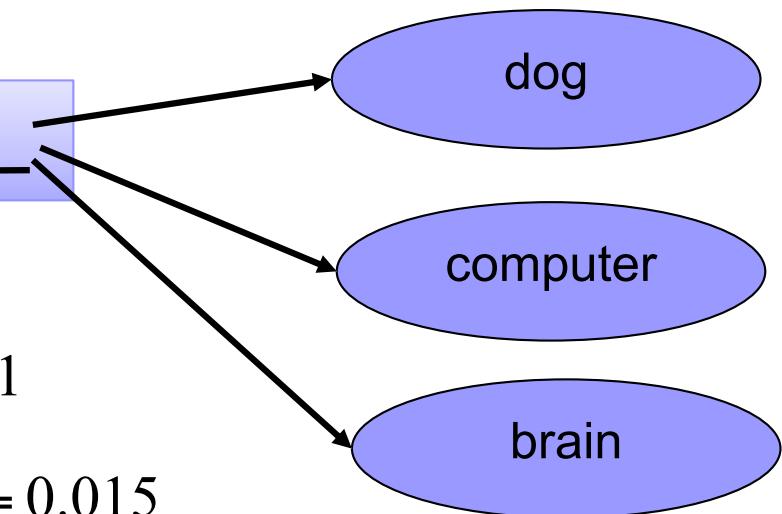
$$p(x_k | x_1, x_2, \dots, x_{k-1})$$

$$p(x_1, x_2, \dots, x_{k-1}, x_k)$$

# Example: Generate Text

- For example:

The student is typing with a \_\_\_\_\_?

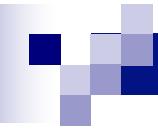


$$p(\text{dog} | \text{the,student,is,typing,with,a}) = 0.001$$

$$p(\text{computer} | \text{the,student,is,typing,with,a}) = 0.015$$

$$p(\text{brain} | \text{the,student,is,typing,with,a}) = 0.003$$

“Computer” is chosen with highest probability.



# Sequence Representation and Classification

- Neural networks can be used to learn a representation vector of a sequence, and use it to classify the sequence.

$$f(x_1 \rightarrow x_2 \rightarrow \dots, x_{k-1} \rightarrow x_k) = \text{class}$$

# Example: Sentence Classification

- Sentiment analysis: Classify the sentiment of a sentence.

$f(\text{overall I enjoyed the movie a lot}) = \text{positive class}$

$f(\text{the quality of this knife is not great}) = \text{negative class}$

- Spam filtering.

$f(\text{from micro soft team... online security team}) = \text{spam class}$

**From:** "Microsoft Team" <[inegon06@netscape.com](mailto:inegon06@netscape.com)>  
**Date:** 7 October 2015 at 17:53:10 BST  
**To:** <[customerservice@outlook.com](mailto:customerservice@outlook.com)>  
**Subject:** Avoid Suspension 2015!!!



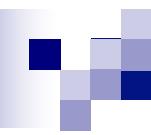
Dear Subscriber,

Your Microsoft account has been compromised. You must update it immediately or your account will be closed.

[Click here](#) to update

Sincerely,

Microsoft Online Security Team



# Sequence Labelling

- Neural networks can be used to learn a representation vector for each state (element) in a sequence, and use it to predict the class label for each state.

$$f(x_1 \rightarrow x_2 \rightarrow \dots, x_{k-1} \rightarrow x_k) \\ = \text{class}_1 \rightarrow \text{class}_2 \rightarrow \dots, \text{class}_{k-1} \rightarrow \text{class}_k$$

# Example: POS Tagging

- Part-of-speech tagging: Assign a word to a priorly defined lexical class.

*f(the startled cat knocked over the vase)*  
= (DT JJ NN VBN IN DT NN)

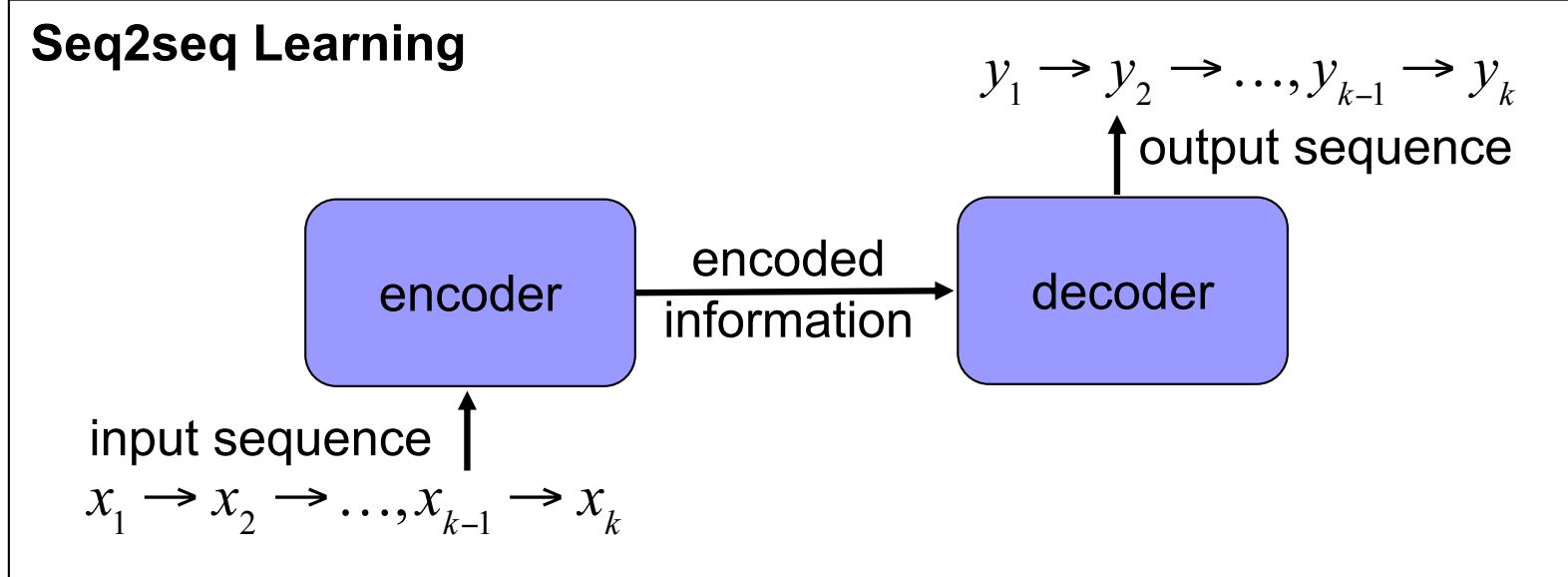
- Named entity recognition: Assign a word to a name class (e.g. in a name class or in no name class).

*f(Fred showed Sue Mengqiu Huang's new painting)*  
= (PER O PER PER PER O O O)

# Sequence to Sequence Learning

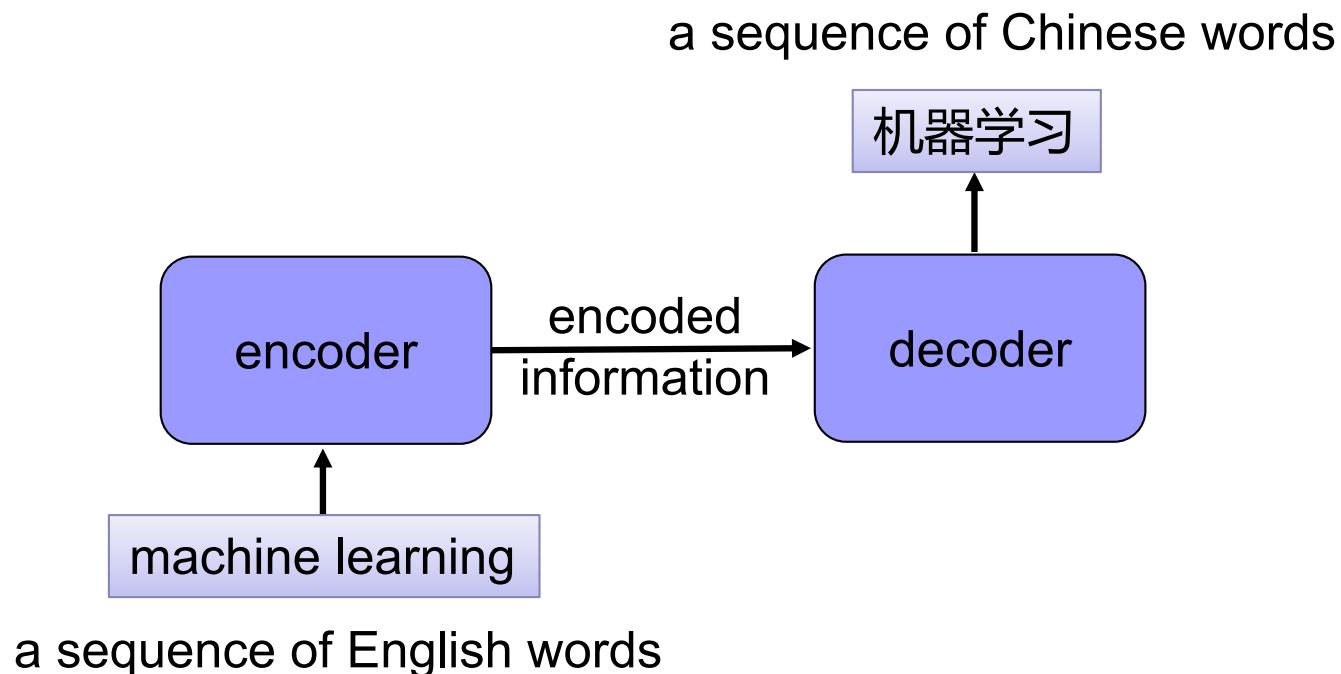
- Neural networks can be used to **encode** information in an input sequence (seq), and **decode** it to generate an output sequence (2seq).

$$f(x_1 \rightarrow x_2 \rightarrow \dots, x_{k-1} \rightarrow x_k) = y_1 \rightarrow y_2 \rightarrow \dots, y_{k-1} \rightarrow y_k$$



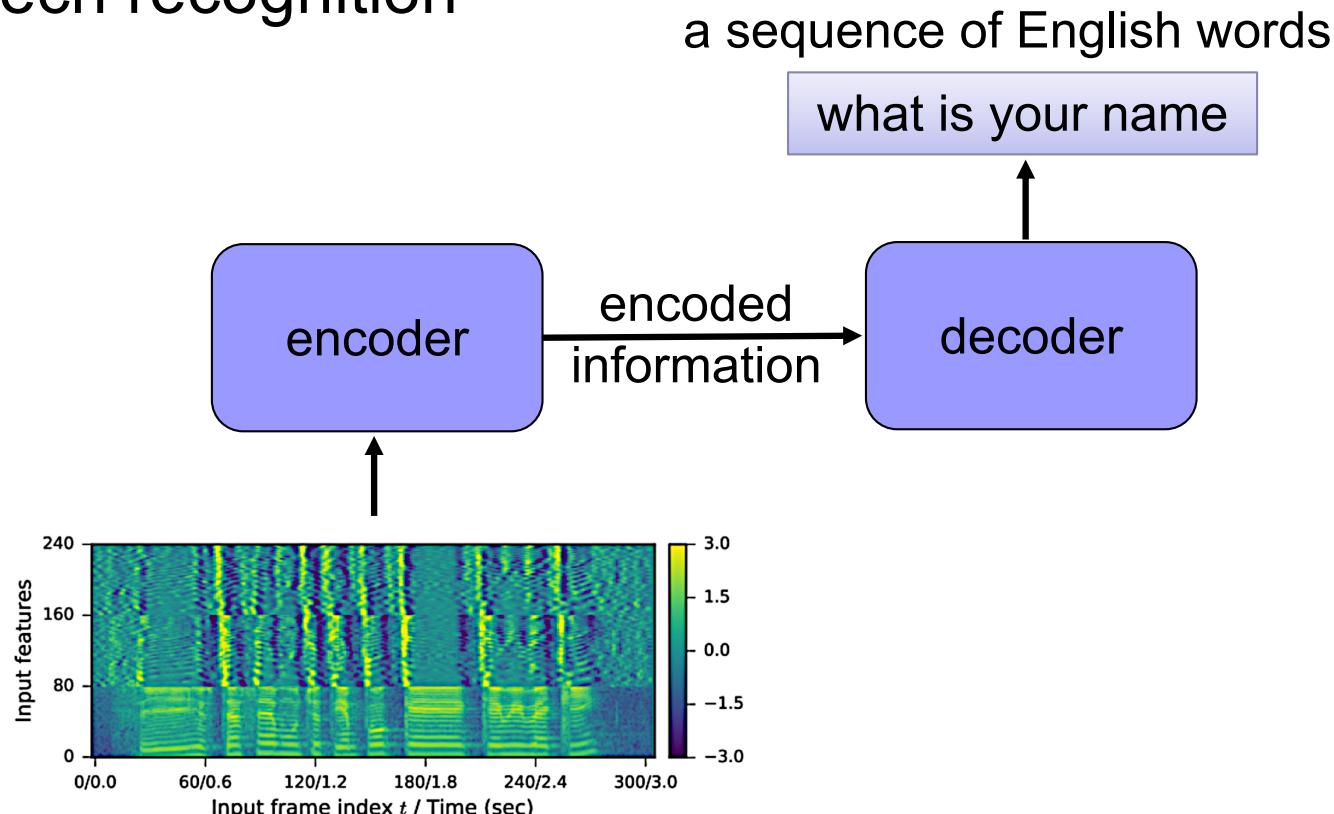
# Examples

## ■ Machine translation



# Examples

## ■ Speech recognition



a sequence of speech feature vectors

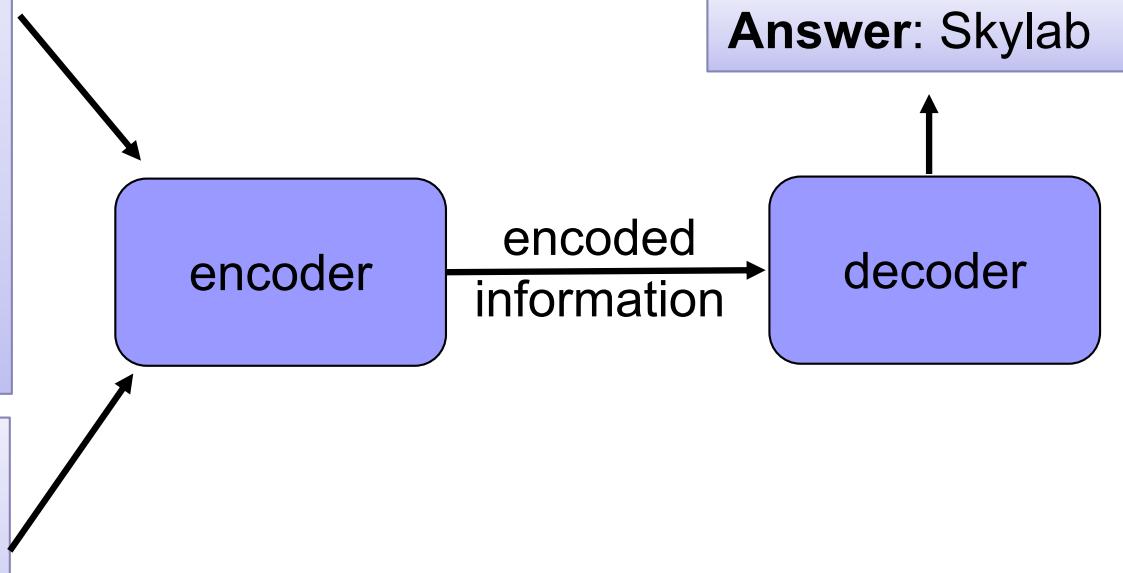
Feature figure is from <https://www.arxiv-vanity.com/papers/1703.08581/>

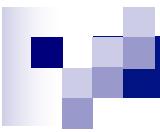
# Examples

- Question answering

**Context:** Apollo ran from 1961... vehicles were also used for an Apollo Applications Program, which consisted of Skylab, a space station that supported three manned missions in 1973–74, and the Apollo–Soyuz...

**Question:** What space station supported three manned missions in 1973-1974?



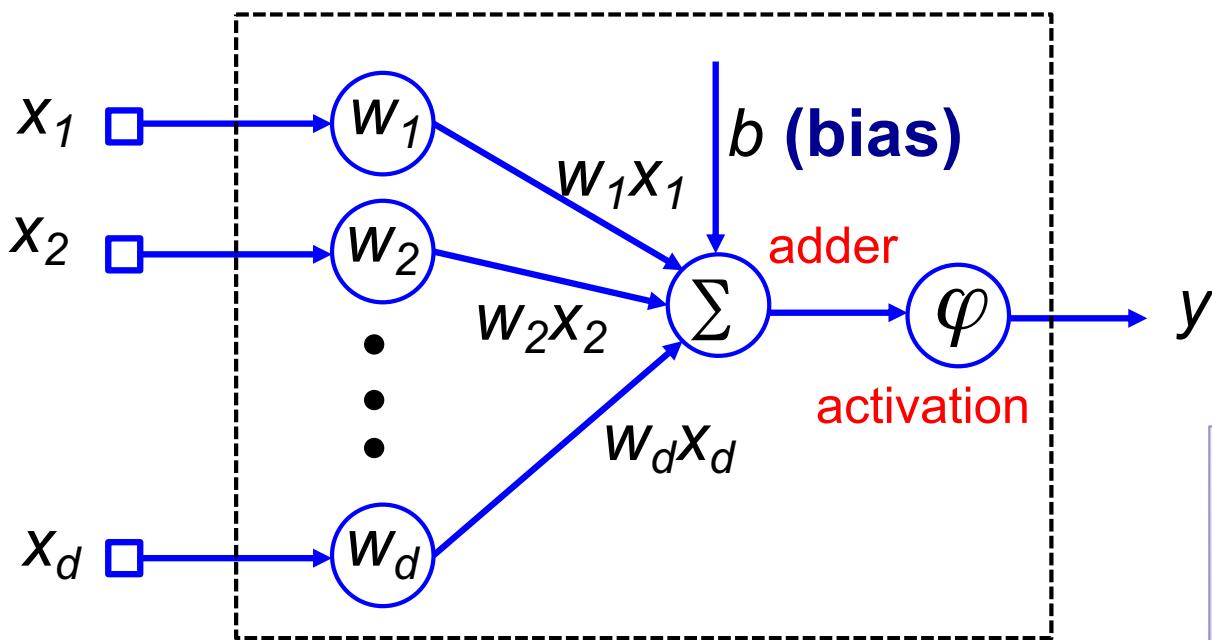


## ■ ***Basic Neural Network Architectures***

*This part is taught in your year 2 Machine Learning course.  
Revisit the knowledge.*

# Single Neuron

- An artificial neuron: multiple inputs  $[x_1, x_2, \dots, x_d]$  and one output  $y$ .



**neuron**    
$$y = \varphi \left( \sum_{i=1}^d w_i x_i + b \right)$$

Given  $d$  input values, a single neuron has  $d+1$  parameters. Training is the process of finding the best values of these  $d+1$  parameters.

# Typical Activation Functions

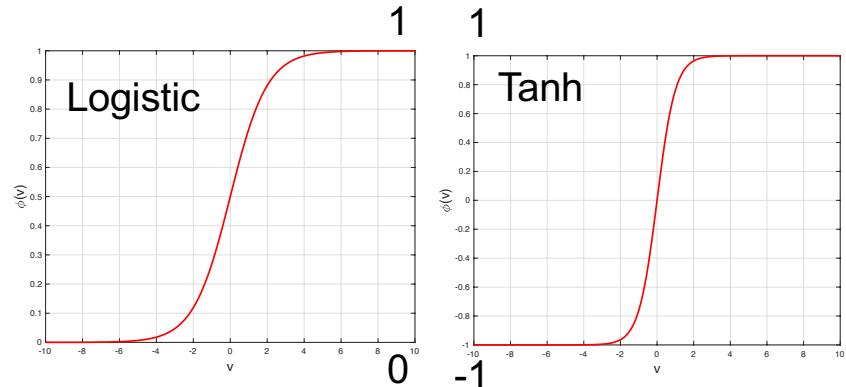
- Sigmoid function (“S”-shaped curve):

Logistic function:

$$\varphi(v) = \frac{1}{1 + \exp(-v)} \in (0, 1)$$

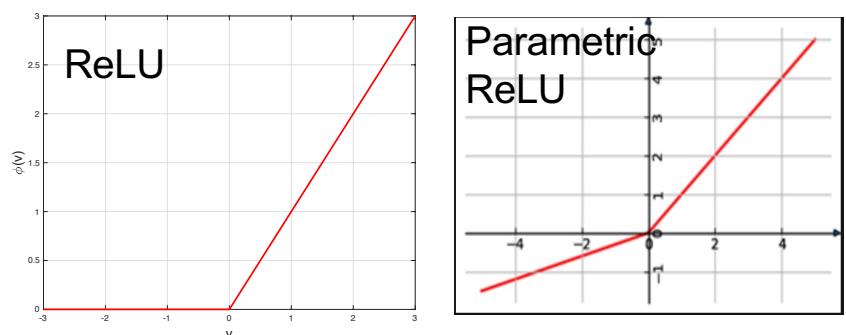
Hyperbolic tangent:

$$\varphi(v) = \tanh(v) = \frac{\exp(2v) - 1}{\exp(2v) + 1} \in (-1, +1)$$



- Rectified linear unit (ReLU):

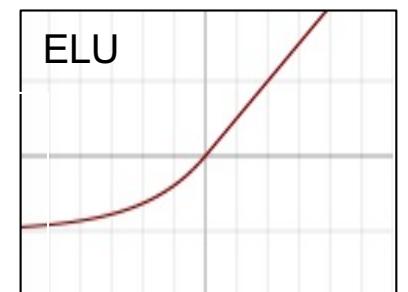
$$\varphi(v) = \begin{cases} v & \text{if } v \geq 0 \\ 0 & \text{if } v < 0 \end{cases}$$



- Parametric ReLU (a=0.01 Leaky Relu) ■ ELU

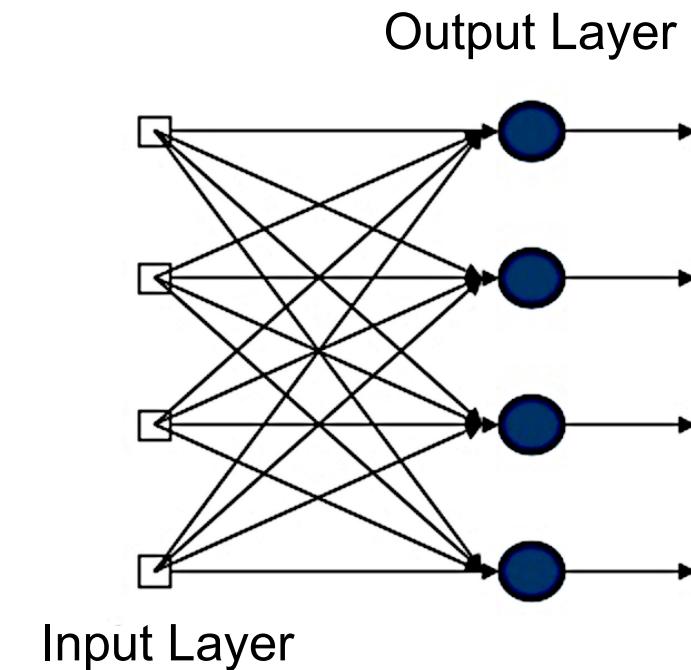
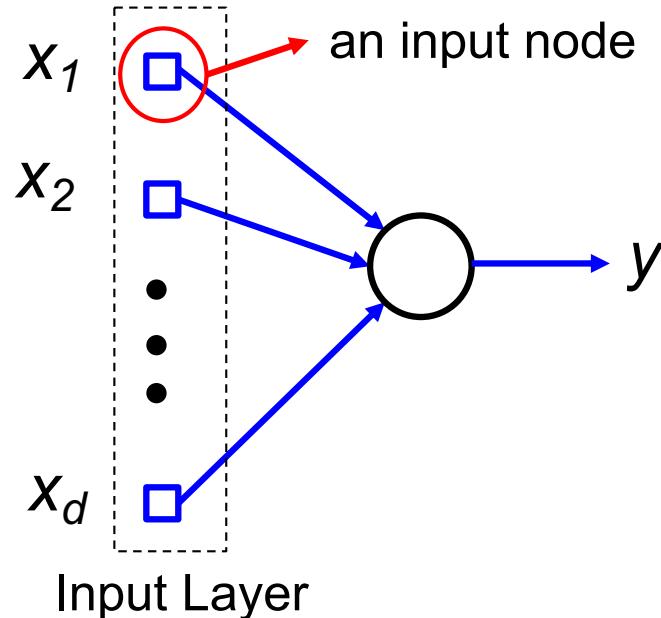
$$\varphi(v) = \begin{cases} v & \text{if } v \geq 0 \\ av & \text{if } v < 0 \end{cases}$$

$$\varphi(v) = \begin{cases} v & \text{if } v \geq 0 \\ a(e^v - 1) & \text{if } v < 0 \end{cases}$$



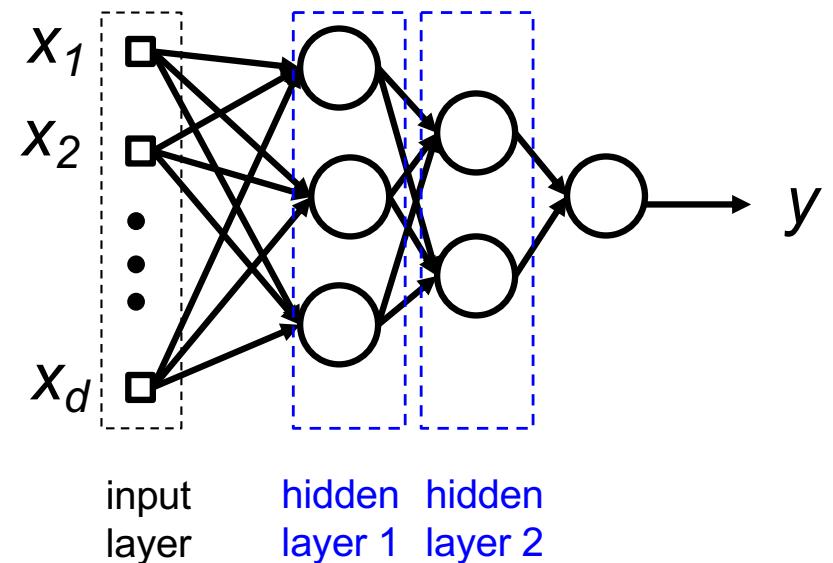
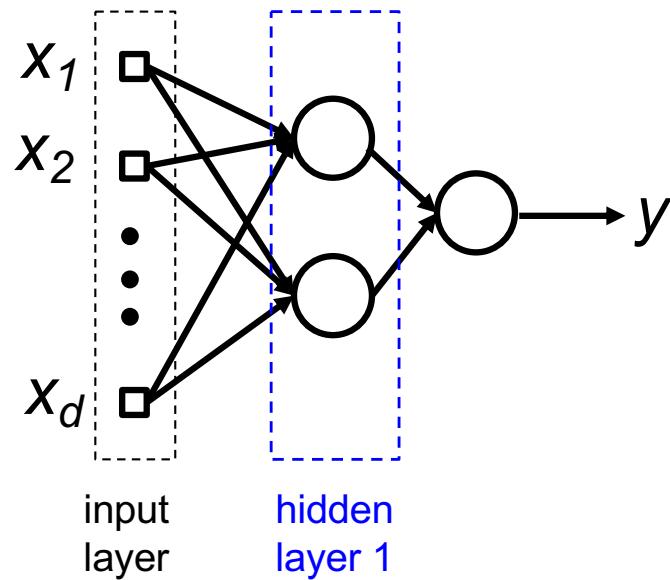
# Single Layer Perceptron

- A single layer perceptron has one input layer and one output layer.



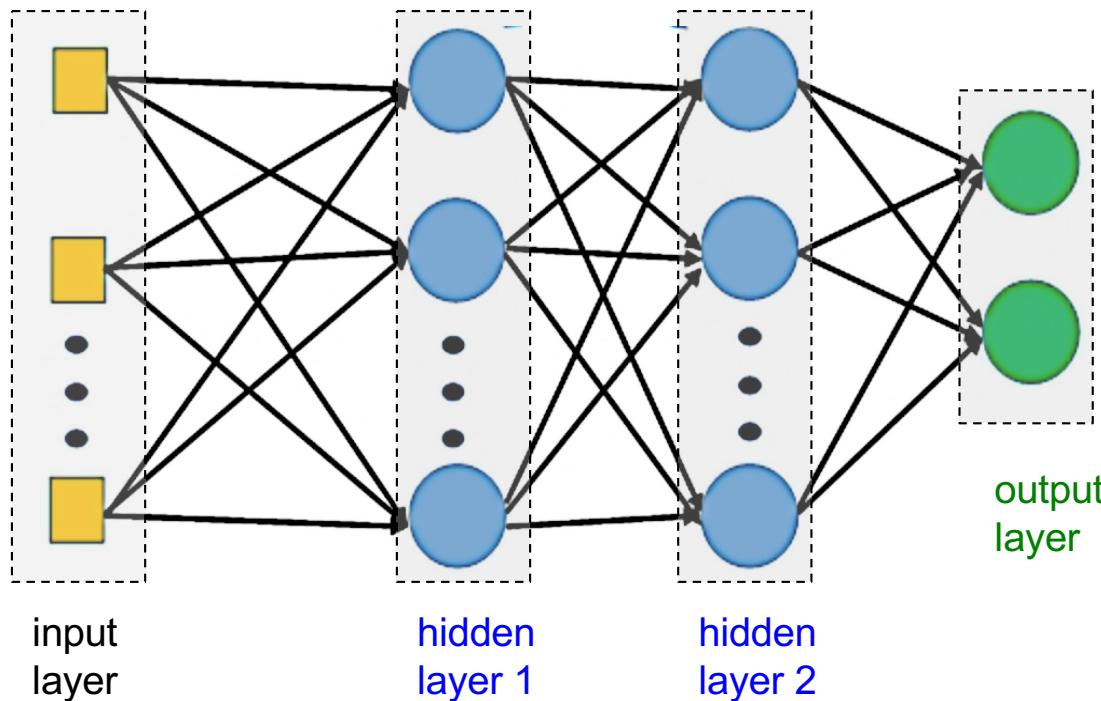
# Adding Hidden Layers!

- The presence of hidden layers allows to formulate more complex functions.



# Multilayer Perceptron

- A **multilayer perceptron** (MLP), also called **feedforward artificial neural network**, consists of at least three layers of nodes: input, hidden (at least one) and output layers.
- In a **fully connected** neural network, each neuron is connected to all the neurons in the previous layer through non-zero weight.

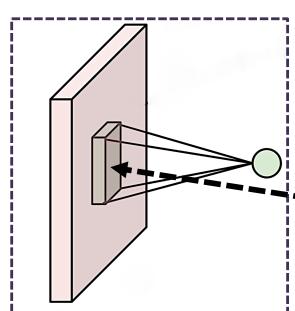


This is a powerful tool that can help you to create a complex multi-input and multi-output function.

# Convolutional Neural Network (CNN)

Green cube is a convolutional filter:

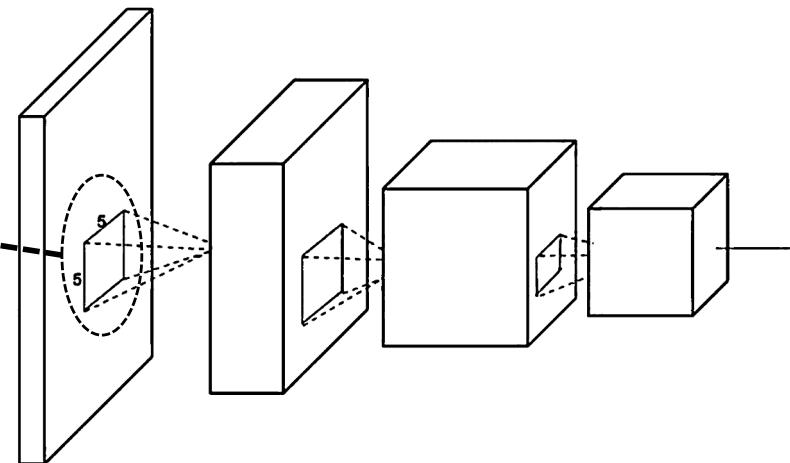
a (width) x b (height) x d (depth)



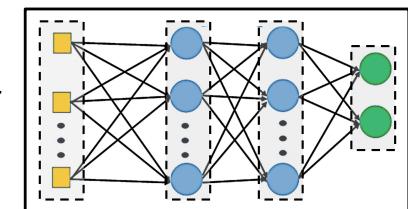
Red cube is a set of data matrices:

m (width) x n (height) x d (depth)

Convolutional and Pooling Layers



Fully Connected Layers



Main CNN recipe:

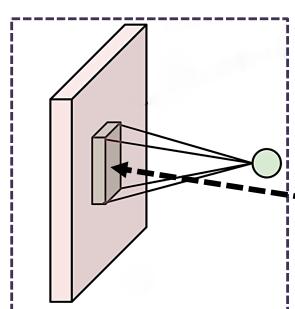
- **2D/3D Neurons:** Neurons are arranged in 2 dimensions (width and height) or 3 dimensions (width, height and depth).
- Convolutional and pooling layers: local connection, weight sharing, down-sampling.
- Fully connected layers in the end.

# Convolutional Neural Network (CNN)

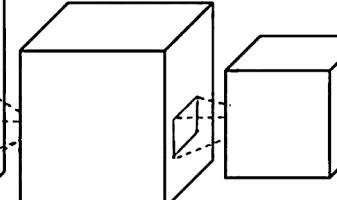
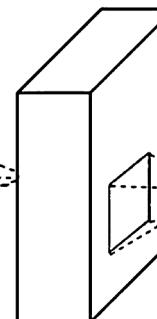
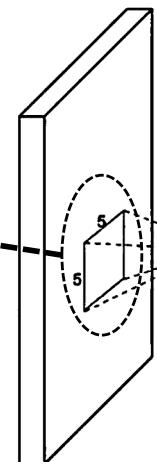
Green cube is a convolutional filter:

a (width) x b (height) x d (depth)

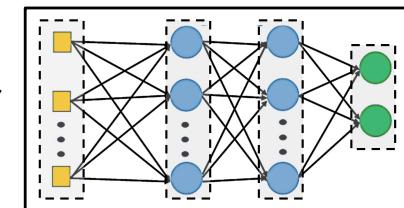
Convolutional and Pooling Layers



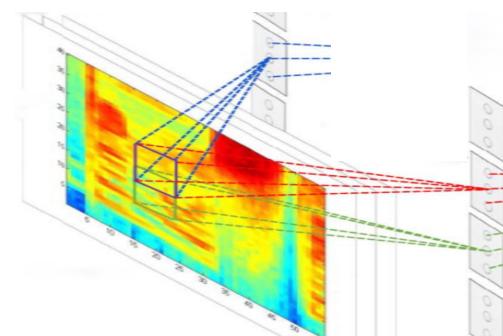
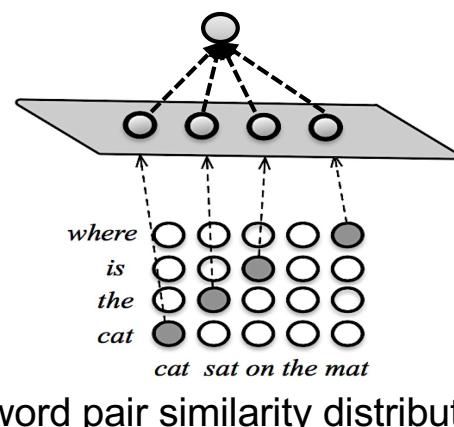
Red cube is a set of data matrices:  
m (width) x n (height) x d (depth)



Fully Connected Layers



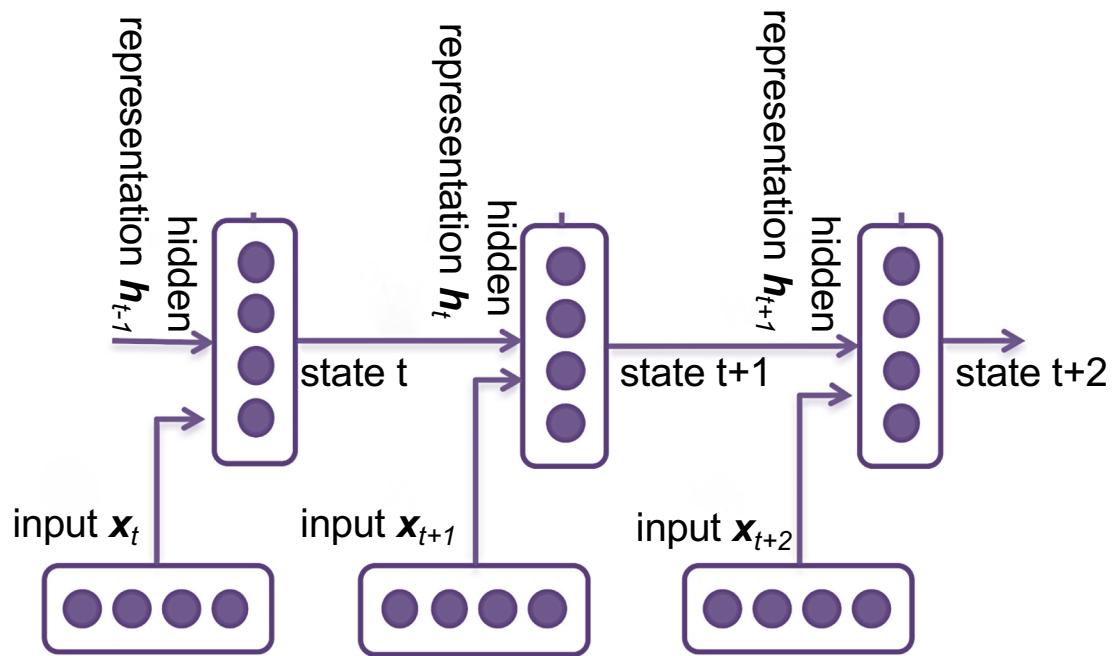
CNNs are good at capturing data patterns existing in a **3-D data cube** or a **2-D data matrix**, widely used for processing images, videos, text and signals.



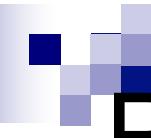
Abdel-Hamid et al. ,2014'

Energy pattern of audio signals over  
different frequency bands

# Recurrent Neural Network (RNN)



RNNs are good at capturing data patterns and encoding information in sequential data.

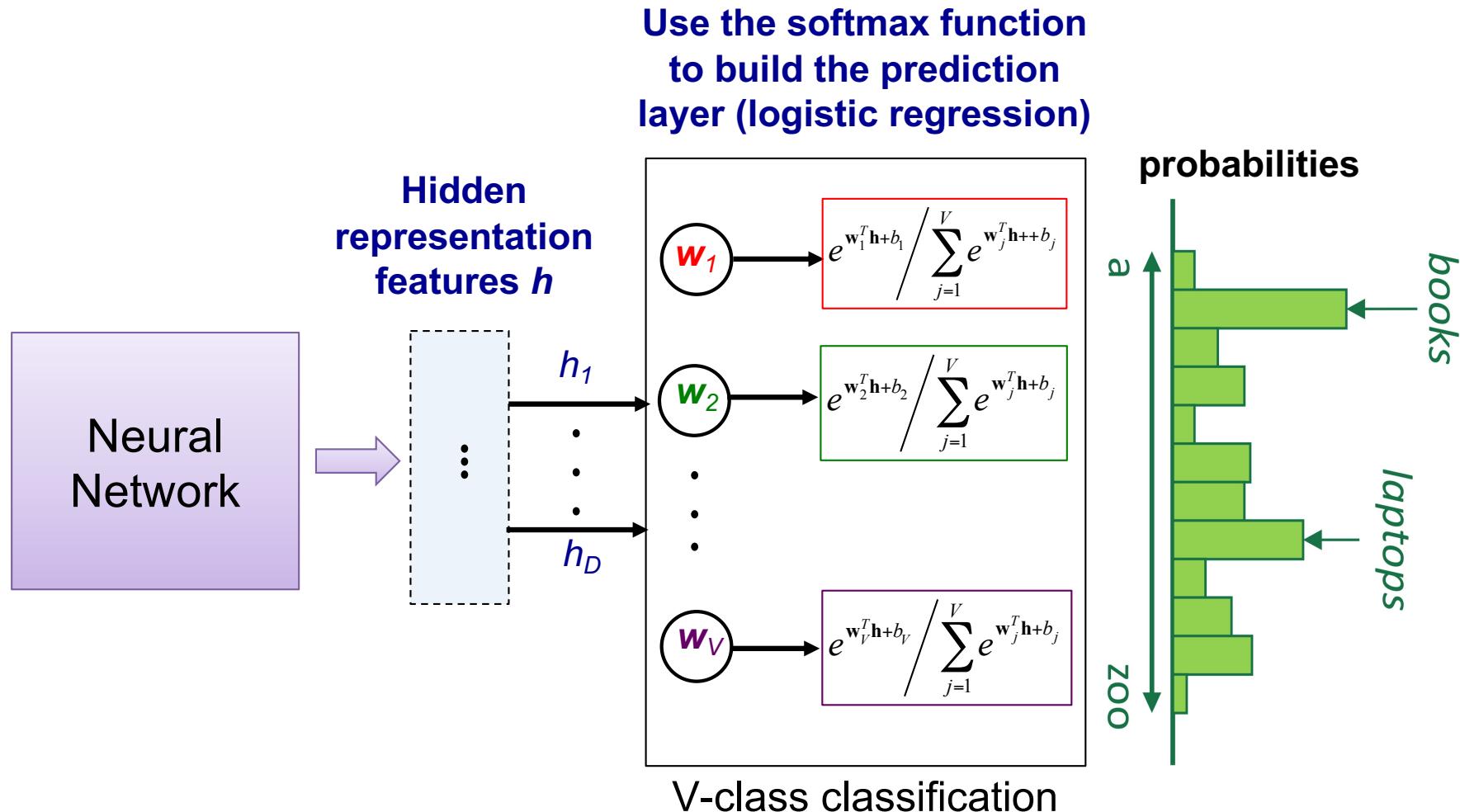


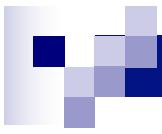
# Prediction Layer

- Most NLP tasks can be decomposed into a series of **classification** tasks.
- The hidden representation vectors returned by neural networks are usually used as the input of a **logistic regression** model (a classifier). This is referred to as the prediction layer.
- Training neural network weights aims at minimising a **classification error loss** calculated using your training data.
- Cross entropy loss is a common option.

# Prediction Layer

- Here is an example of the prediction layer.





*Thank you!*