



1

Nature-inspired computing

- Nature has always served as a source of inspiration for engineers and scientists
- The best problem solver known in nature is:
 - **the (human) brain** that created “the wheel, New York, wars and so on” (after Douglas Adams’ Hitch-Hikers Guide)
 - **the evolution mechanism** that created the human brain (after Darwin’s Origin of Species)
- Answer 1 → neurocomputing
 - Today
- Answer 2 → evolutionary computing
 - Weeks 7 + 8

2 / 39

2

1

Outline

- Artificial Neural Networks (ANNs)
 - Neural networks
 - Logistic regression – revision
 - Different learning tasks
 - ANNs
 - Cost function
 - Gradient descent
 - Backpropagation
- Deep Learning

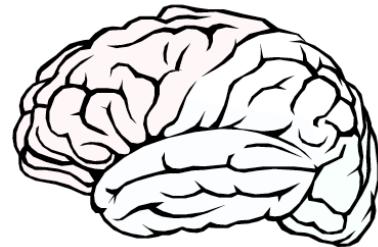
3

Supervised Learning using Artificial Neural Networks

4

Neural Networks

- Human brain consists of around 86 billion interconnected “neurons”
- Each neuron is connected to about 1000+ neighbouring neurons
- Each neuron responds to stimuli and passes output to other neurons
- In ML, artificial neural networks are loosely modelled on the human brain



5

Biological Networks of Neurons

- A neuron receives inputs (electrochemical signals) from its neighbours
- If enough inputs are received at the same time, it gets *activated*
- Once activated, a neuron fires, giving an output that may (in turn) activate connected neurons
- The behaviour of a biological network of neurons is dependent on connection types and activation strength (*synaptic function*)

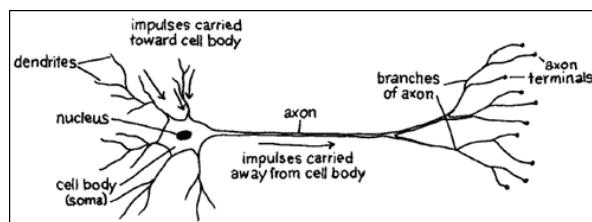
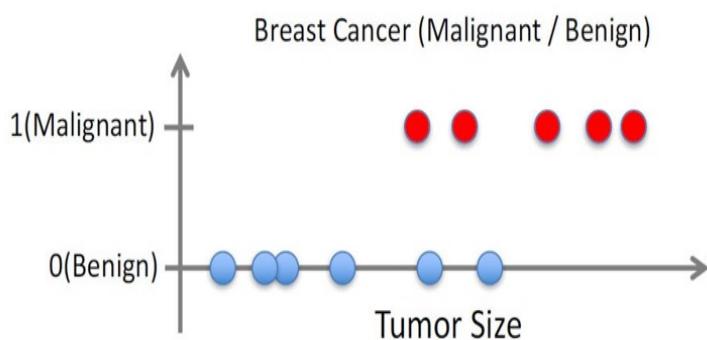


Image Source: <https://science.education.nih.gov/supplements/webversions/BrainAddiction/guide/lesson2-1.html>

6

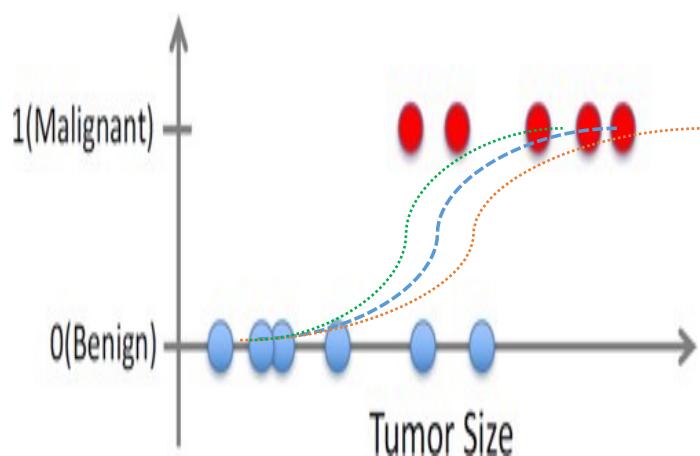
Review - Logistic regression

- It is actually a classification technique!



Source: Andrew Ng/Stanford Online

7



- Note that the blue model is not the only one. There are many. We still face the problem of which function is the “right” one

- → minimize the cost function using Gradient Descent (GD) or Maximum Likelihood Estimation (MLE)

8

4

Logistic regression - mathematically

- $\hat{y} = f_{W,b}(x) = \sigma(g_{W,b}(x)) = \frac{1}{1+e^{-g_{W,b}(x)}}$
- More precisely: $\text{Prob}(y = 1|x) = \frac{1}{1+e^{-g_{W,b}(x)}}$
- E.g., $\text{Prob}(y = \text{Malignant}|\text{size}) = \frac{1}{1+e^{-g_{W,b}(\text{size})}}$
- where $g_{w,b}(\text{size}) = b + w * \text{size}$

9

Logistic regression - mathematically

- $\hat{y} = f_{W,b}(x) = \sigma(g_{W,b}(x)) = \frac{1}{1+e^{-g_{W,b}(x)}}$
- Of course, for multivariate problem:

$$g_{W,b}(x_1, \dots, x_n) = b + w_1 * x_1 + \dots + w_n * x_n$$

$$W \text{ is a vector } W = (w_1, \dots, w_n)$$

Some people write: $g_W(x_1, \dots, x_n) = w_0 + w_1 * x_1 + \dots + w_n * x_n = Wx$

by assuming: $x_0 = 1$ and the vector $W = (w_0, w_1, \dots, w_n)$

We are supposed to write $W^T x$ but we can keep it simple

10

Logistic regression - mathematically

- $\hat{y} = f_W(x) = \sigma(g_W(x)) = \frac{1}{1+e^{-g_W(x)}}$

- The cost function:

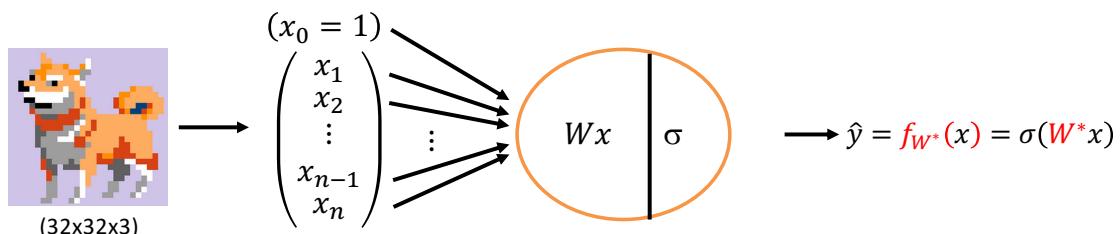
$$J(W) = \frac{1}{m} \sum_{i=1}^m y^i \log \hat{y}^i + (1 - y^i) \log(1 - \hat{y}^i)$$

- The model f_W can then be trained by minimizing $J(W)$ using the Gradient Descent (GD) algorithm.

11

Logistic Regression for Image classification

Task 1: Determine if a dog is in the image: $\begin{cases} 1 & \text{if there is a dog} \\ 0 & \text{if there isn't a dog} \end{cases}$



General steps in a Gradient Descent (GD) algorithm:

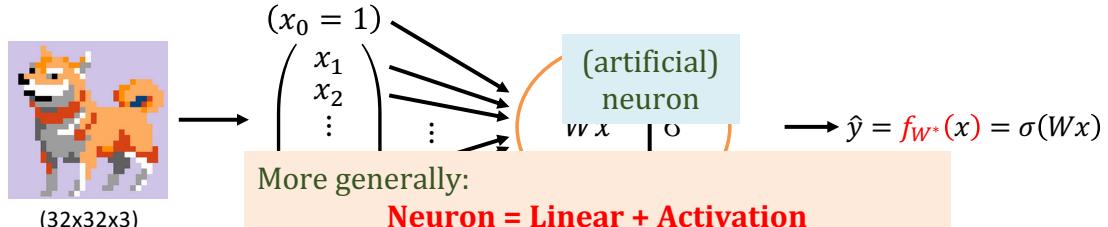
- Initialise (randomly) W
- Repeat the updates $W := W - \alpha \frac{\partial J(W)}{\partial W}$ until convergence to a minimum

Use the trained model f_{W*} for inference: Given an input x , $\hat{y} = f_{W*}(x)$

12

Logistic Regression for Image classification

Task 1: Determine if a dog is in the image: $\begin{cases} 1 & \text{if there is a dog} \\ 0 & \text{if there isn't a dog} \end{cases}$



General steps (& there are many **activation functions**: Sigmoid, step, linear, tanH, Gaussian, ReLU, etc.)

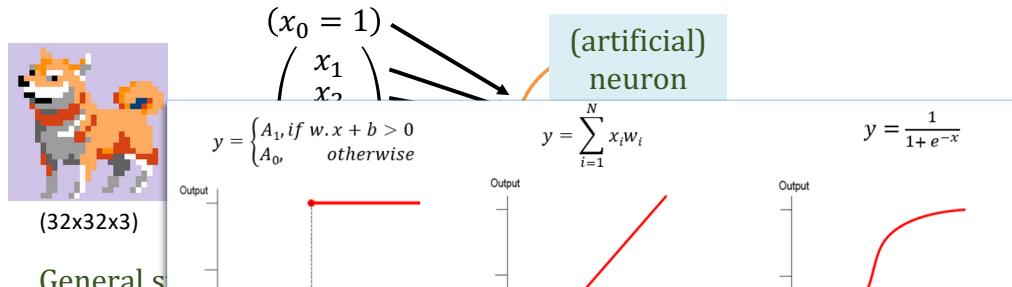
- Initialise (W)
- Repeat the updates $W := W - \alpha \frac{\partial J(W)}{\partial W}$ until convergence to a minimum

Use the trained model f_{W^*} for inference: Given an input x , $\hat{y} = f_{W^*}(x)$

13

Logistic Regression for Image classification

Task 1: Determine if a dog is in the image: $\begin{cases} 1 & \text{if there is a dog} \\ 0 & \text{if there isn't a dog} \end{cases}$

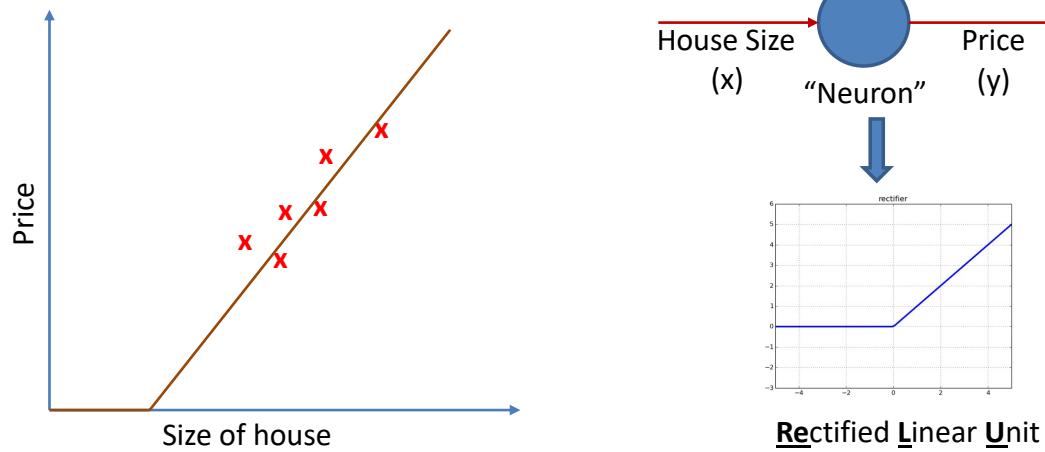


- Initialise (W)
- Find the threshold (b)
- Use W and b to predict

Learn more about activation functions - <https://medium.com/the-theory-of-everything/understanding-activation-functions-in-neural-networks-9491262884e0>

14

Price Prediction - Single Neuron



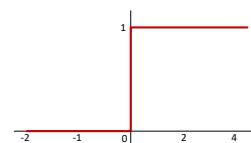
15

Perceptron

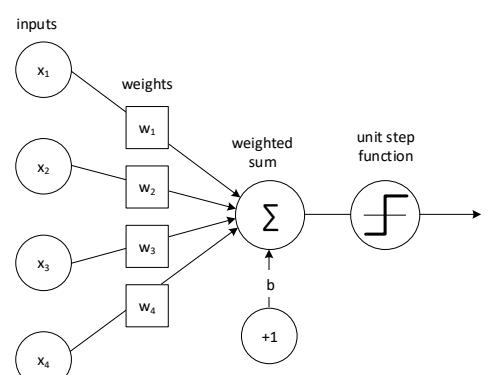
- Simplest neural network possible
- Single neuron that performs binary classification
- Uses a *step activation function*
 - Also referred to as *Heaviside Step Function*

$$f(x) = \begin{cases} 1, & \text{if } w \cdot x + b > t \\ 0, & \text{otherwise} \end{cases}$$

where $w \cdot x$ is the dot product $\sum_{i=1}^N x_i w_i$ and t is the threshold



Step Activation Function



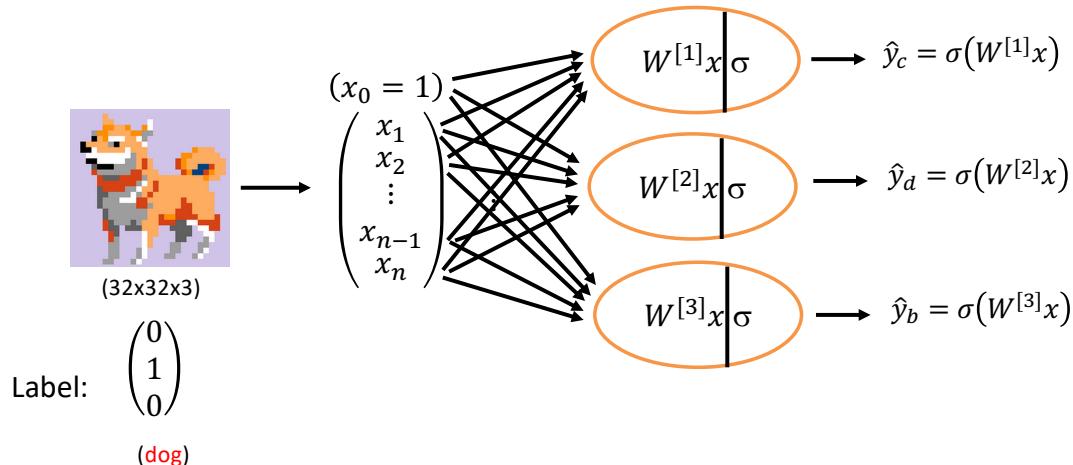
16

What about Multinomial Classification? (aka. Multiclass classification)

SWINBURNE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

Task 2: Classifying images into cat/dog/bird?



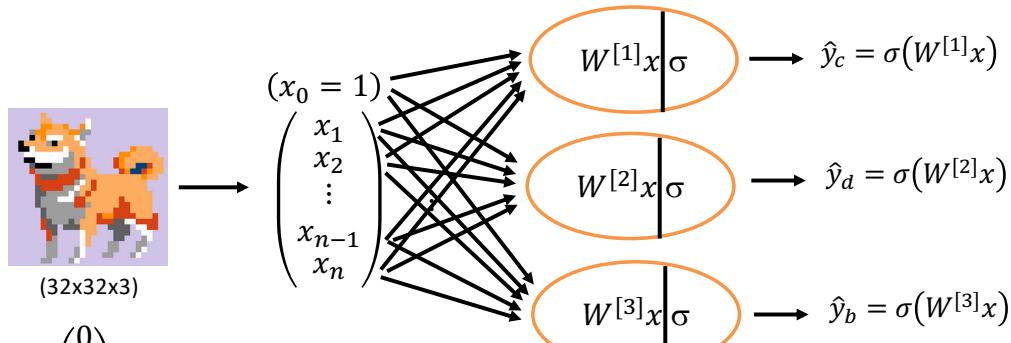
17

What about Multinomial Classification? (aka. Multiclass classification)

SWIN
BUR
NE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

Task 2: Classifying images into cat/dog/bird?



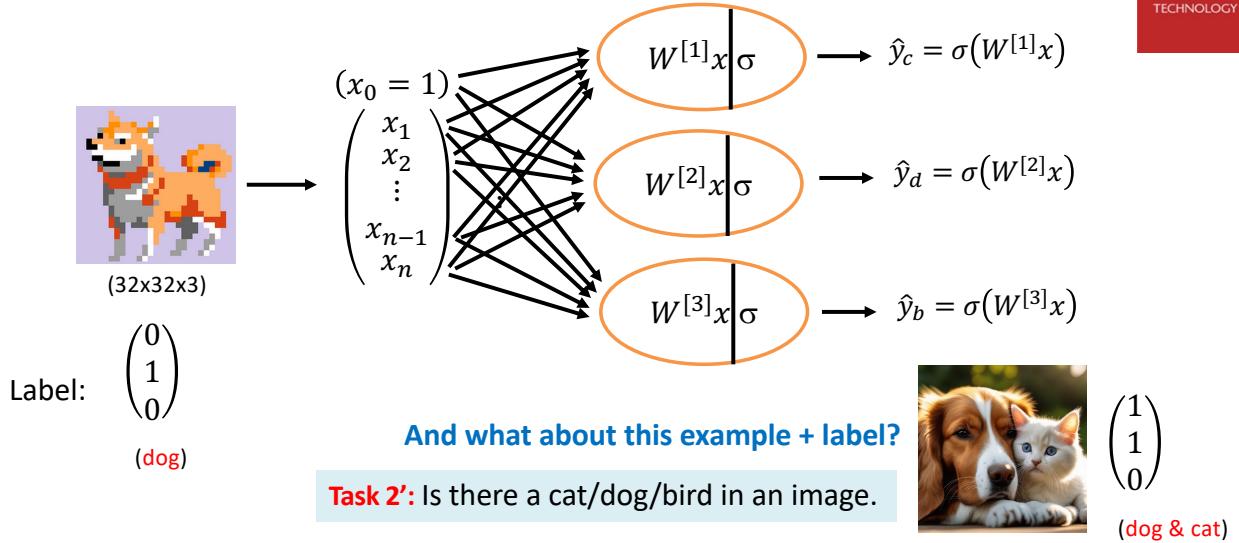
Cost function:

$$J(W) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^3 (y^i \log(\sigma(W^{[j]}x^i)) + (1-y^i) \log(1-\sigma(W^{[j]}x^i)))$$

18

What about Multinomial Classification? (aka. Multiclass classification)

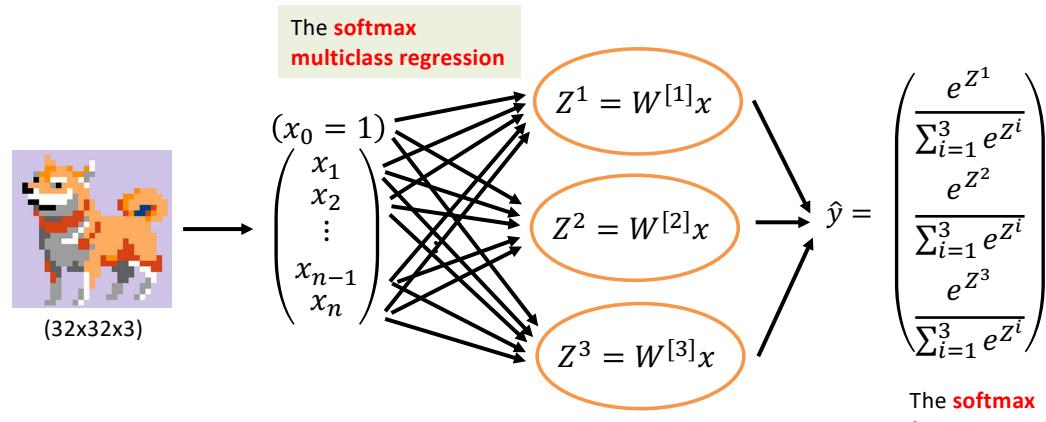
Task 2: Classifying images into cat/dog/bird?



19

Multinomial Classification under constraints

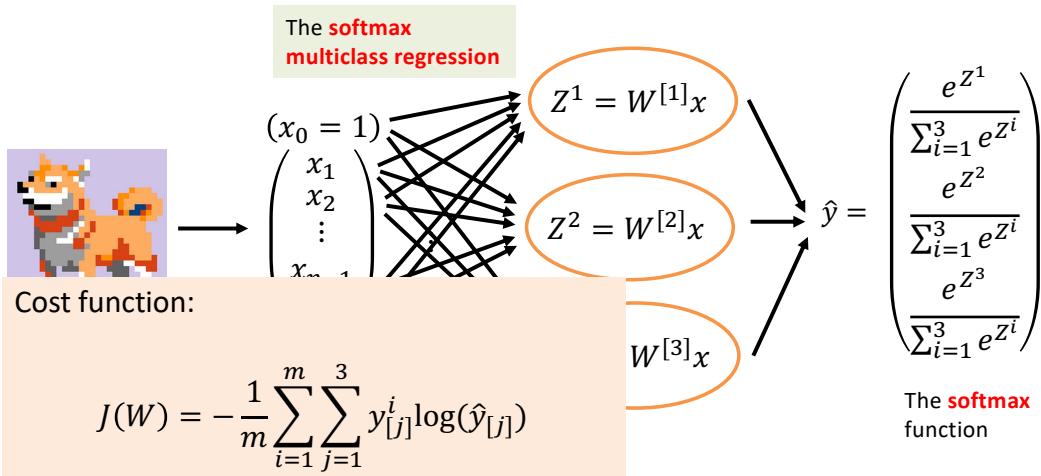
Task 3: Which animal (a cat, dog or bird) is in an image (and at most one of them)?



20

Multinomial Classification under constraints

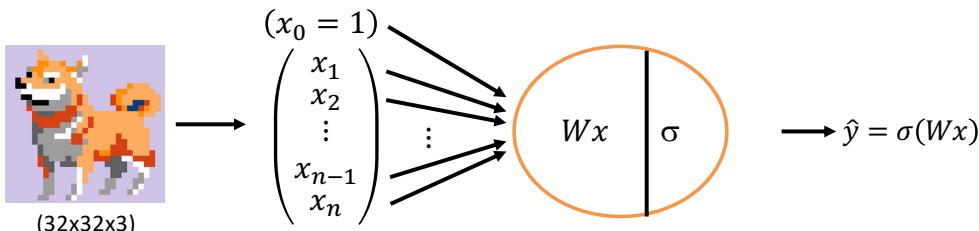
Task 3: Which animal (a cat, dog or bird) is in an image (and at most one of them)?



21

Logistic Regression for Image classification

Task 1: Determine if a dog is in the image: $\begin{cases} 1 & \text{if there is a dog} \\ 0 & \text{if there isn't a dog} \end{cases}$



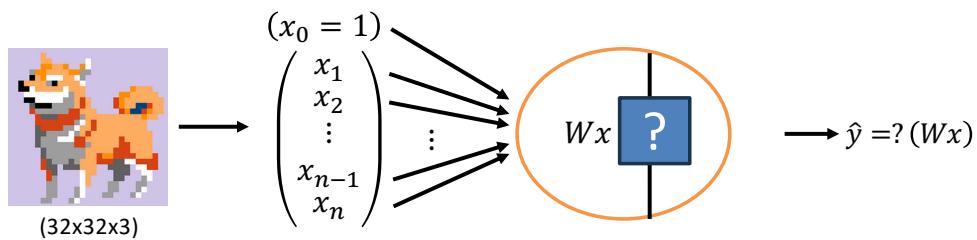
General steps in a Logistic Regression algorithm:

- Initialise (randomly) W
- Find the optimal W^* ($W = W - \alpha \frac{\partial J(W)}{\partial W}$ over a number of iterations)
- Use W^* for inference: Given an input x , $\hat{y} = \sigma(W^*x)$

22

Machine Learning for Image processing

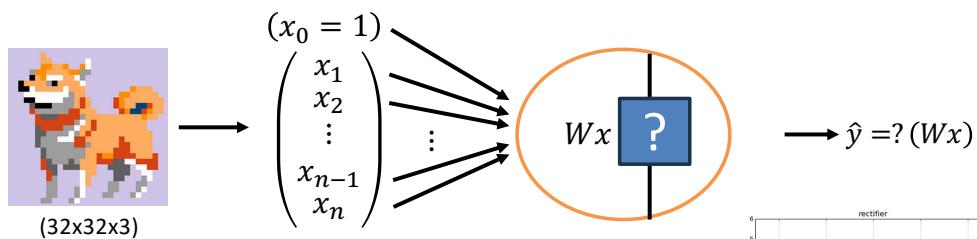
Task 1': Determine the age of the dog is in the image?



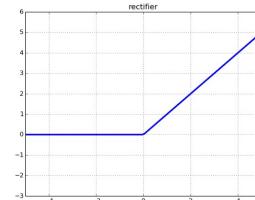
23

Machine Learning for Image processing

Task 1': Determine the age of the dog is in the image?



Remember the ReLU activation function?



Rectified Linear Unit

24

Here come the big question...

- Would it suffice to have just linear + activation (with the right activation function, of course)
 - Classification? 
 - Regression? 
- Then should we complicate our lives with neurons and neural networks?

25

Things are not always simple

- There can be hundreds/thousands/ even millions of people we need to recognise
- There are millions of species of animals



26

Things are not always simple

- There are thousands of languages used by people around the world
 - Different syntaxes and grammars
 - Different rules



27

Things are not always simple

- Even in a single language (e.g., English), voice recognition systems have to deal with many different accents



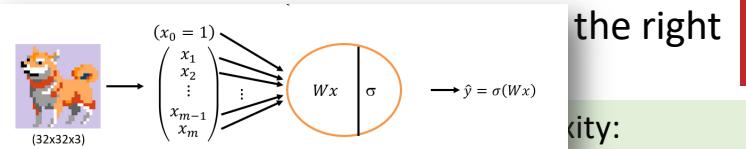
28

Here come the big question...

- Would it suffice to activation function

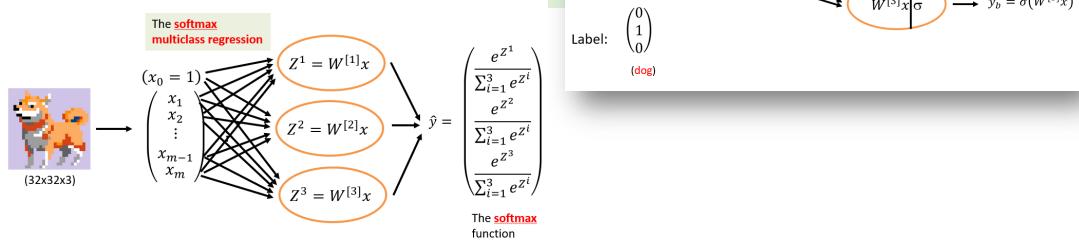
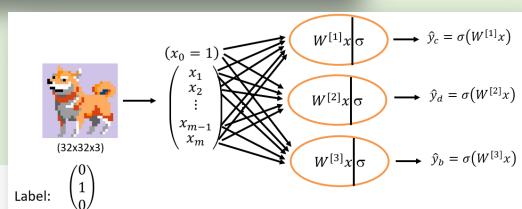
We have seen different

- Task 1 (and 1')
- Task 2 (and 2')
- Task 3



the right

city:



29

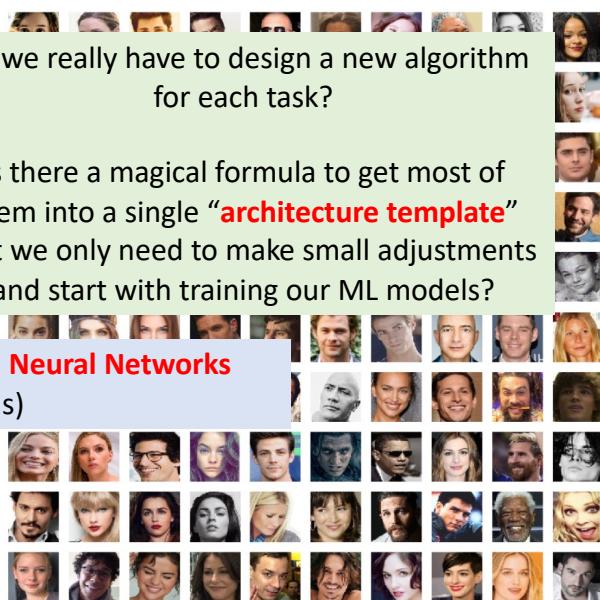
These can be your task 100th, 157th, ...

- There can be hundreds/thousands/ even millions of people we need to recognise
- There are millions of species

Do we really have to design a new algorithm for each task?

Is there a magical formula to get most of them into a single “**architecture template**” that we only need to make small adjustments and start with training our ML models?

Welcome to **Artificial Neural Networks**
(ANNs)



30

These can be your task 100th, 157th, ...

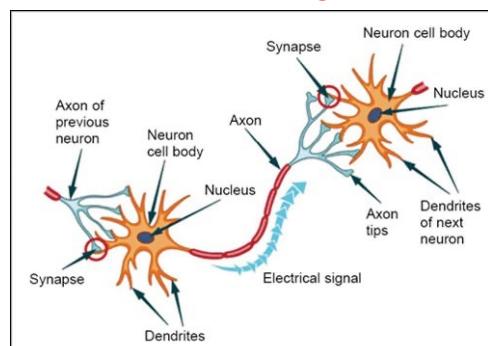
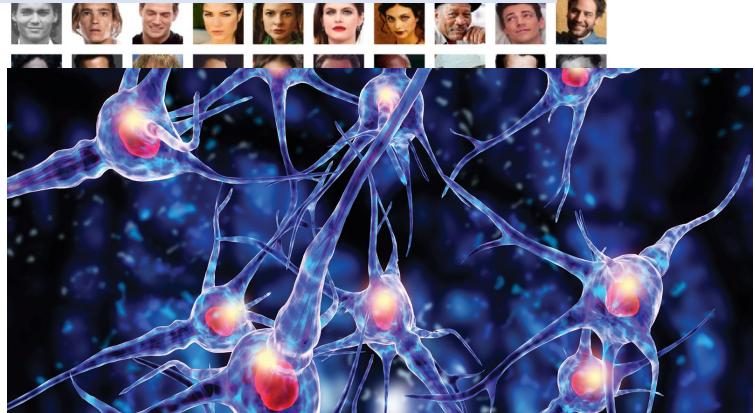
SWIN
BUR
NE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

- There can be hundreds/thousands/ even millions of people we need to recognise



Welcome to Artificial Neural Networks (ANNs)

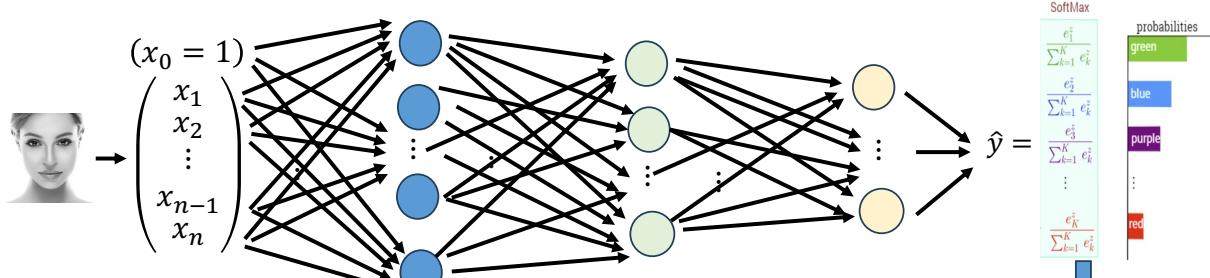
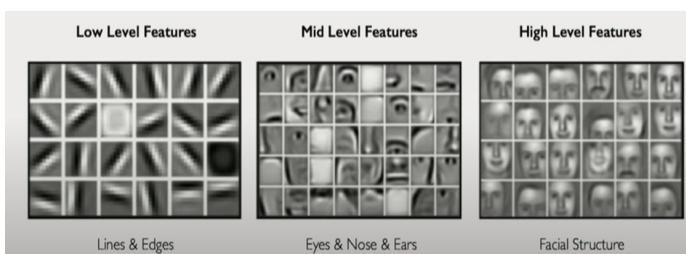


31

Welcome to Artificial Neural Networks (ANNs)

SWIN
BUR
NE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

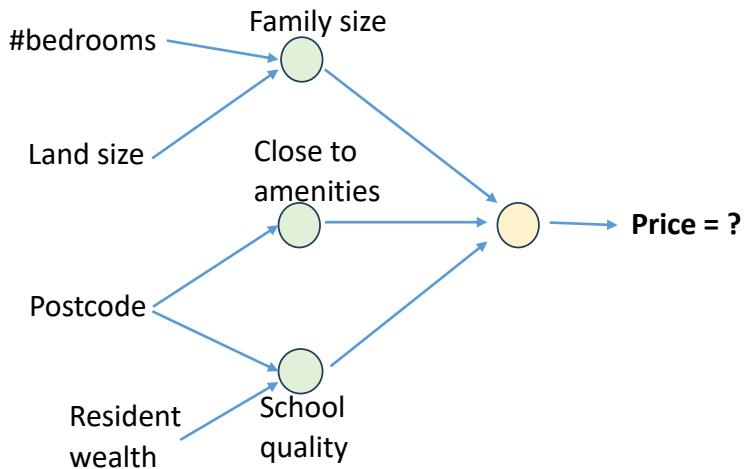


32

16

(Automatic) Feature Engineering

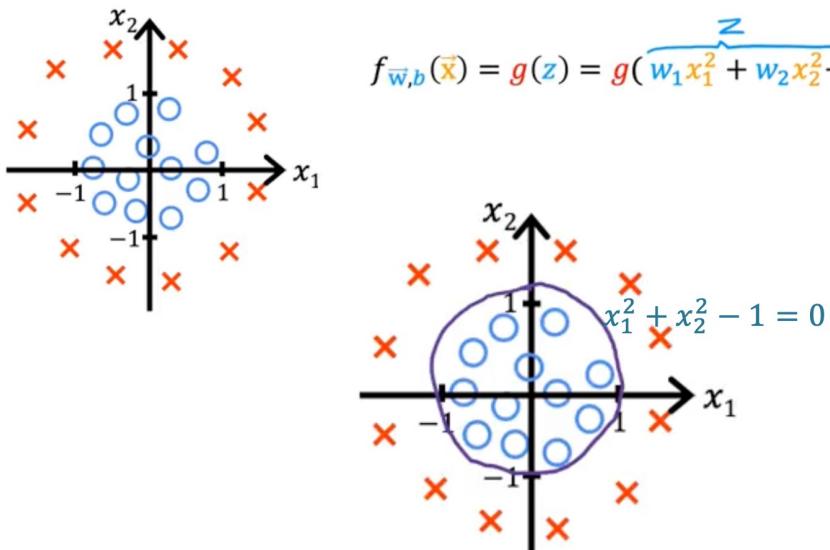
Task: Predict the house price.



33

(Automatic) Feature Engineering

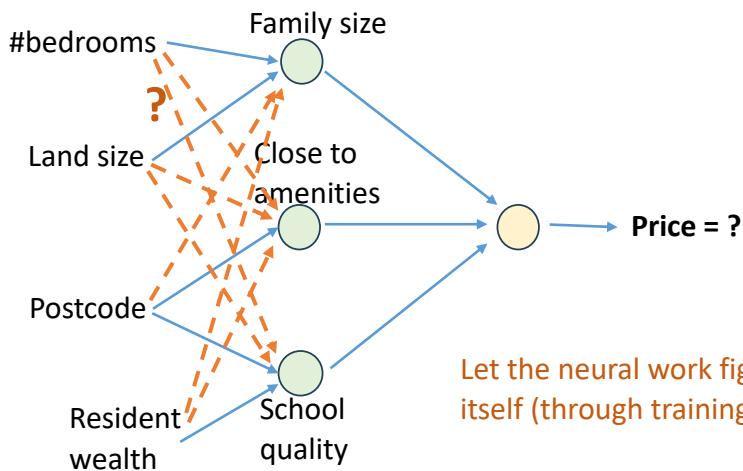
$$f_{\bar{w}, b}(\vec{x}) = g(z) = g(w_1 x_1^2 + w_2 x_2^2 + b)$$



34

(Automatic) Feature Engineering

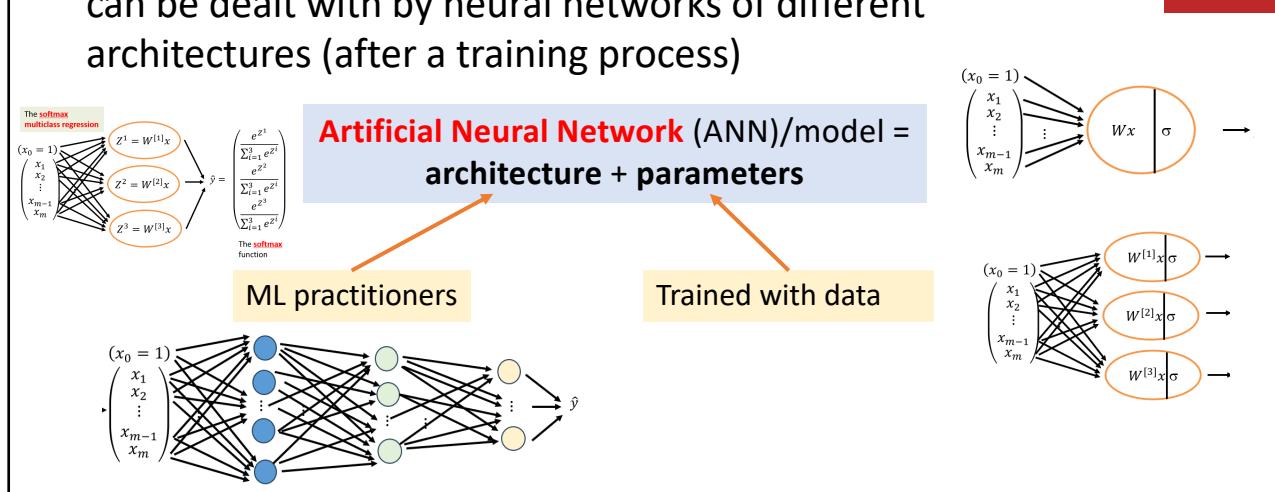
Task: Predict the house price.



35

Important observation

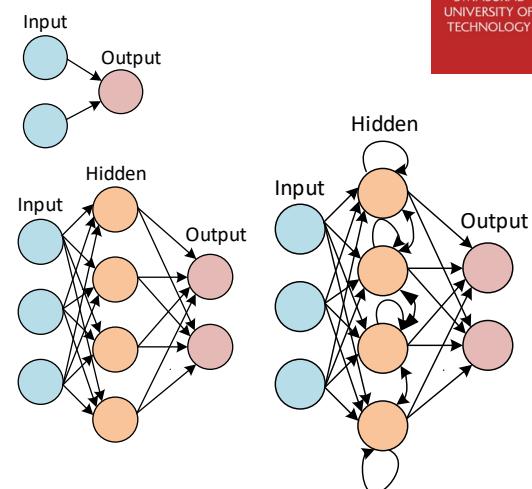
- There can be many ML problems/tasks but (most of) them can be dealt with by neural networks of different architectures (after a training process)



36

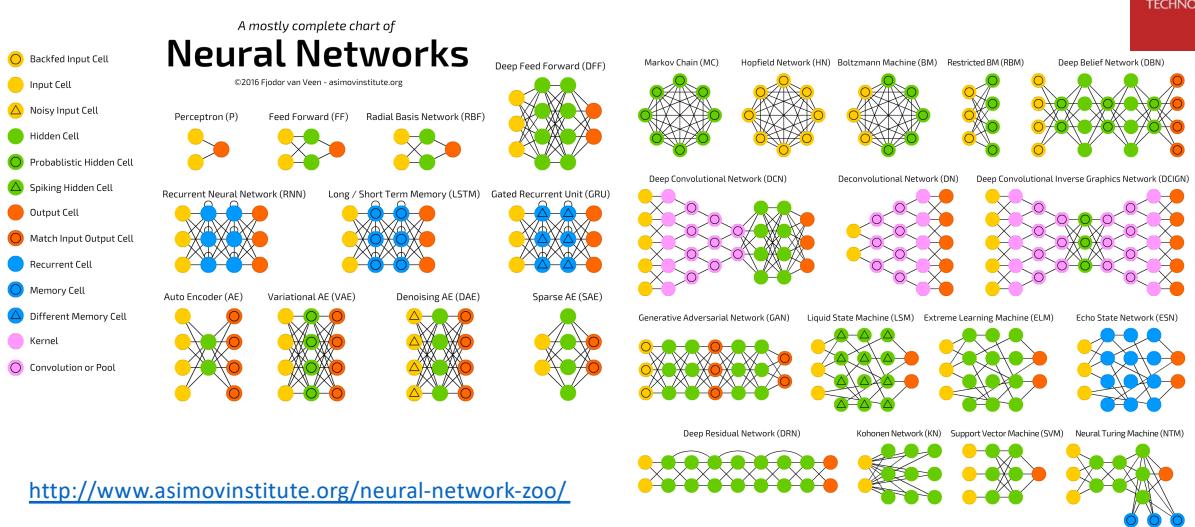
Complex Neural Network - Architectures

- There are 2 main categories of neural network structures:
 - **Acyclic/Feed-Forward Networks** – no loops, no internal state
 - Single-layer perceptron
 - Multi-layer perceptron
 - **Cyclic/Recurrent Networks** – feeds its outputs back into its own inputs, directed cycles, supports short-term memory



37

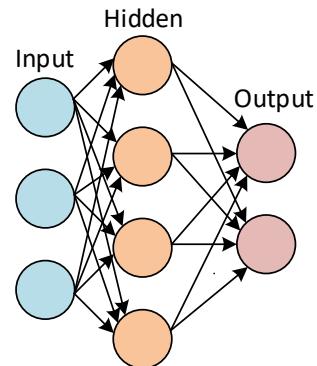
Types of Neural Networks



38

Multiple-Layer Neural Networks

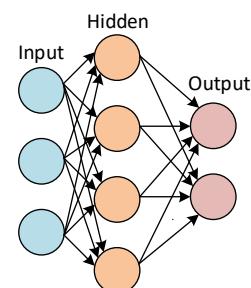
- A complex neural network involves a set of neurons (e.g. perceptrons) interconnected in layers to solve more complex problems
 - Input layer neurons connect to hidden layer neurons
 - Hidden layer neurons connect to output layer neurons
 - Output layer can have multiple outputs
- Multilayer neural networks can classify a range of functions including non-linearly separable ones



39

Learning in Multiple-Layer Neural Networks

- Multilayer NN learn the same way as perceptrons
- **Forward pass** – compute the values from input to output
- **Backward pass** – propagate errors backwards to adjust the connection weights
- Error occurs when expected output != actual output
- **Cost function** – a measure of “how good or bad” a NN did w.r.t its given training sample and the expected output



40

Gradient Descent for Artificial Neural Networks

41

Logistic Regression for Image classification

Task 1: Determine if a dog is in the image: $\begin{cases} 1 & \text{if there is a dog} \\ 0 & \text{if there isn't a dog} \end{cases}$

It's unlikely giving us a good classifier:

The flattened single feature vector representing the m pixels loses the spatial relationships within the image

→ Require effective feature engineering to extract meaningful features (e.g., edges, textures, colour histograms, etc.)

(32x32x3)

$\backslash x_n /$

General steps in a Gradient Descent (GD) algorithm:

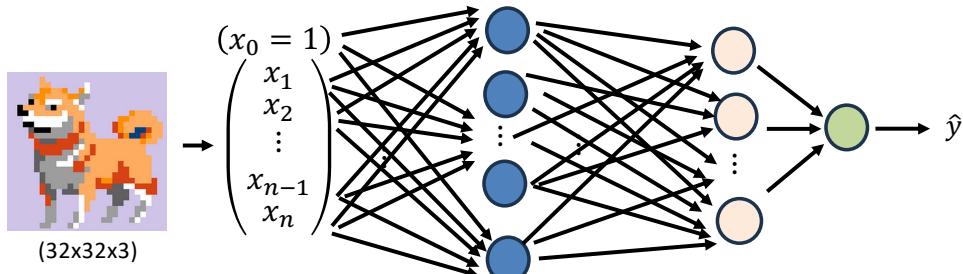
- Initialise (randomly) W
- Repeat the updates $W := W - \alpha \frac{\partial J(W)}{\partial W}$ until convergence to a minimum

Use the trained model f_{W^*} for inference: Given an input x , $\hat{y} = f_{W^*}(x)$

42

ANN for Image classification

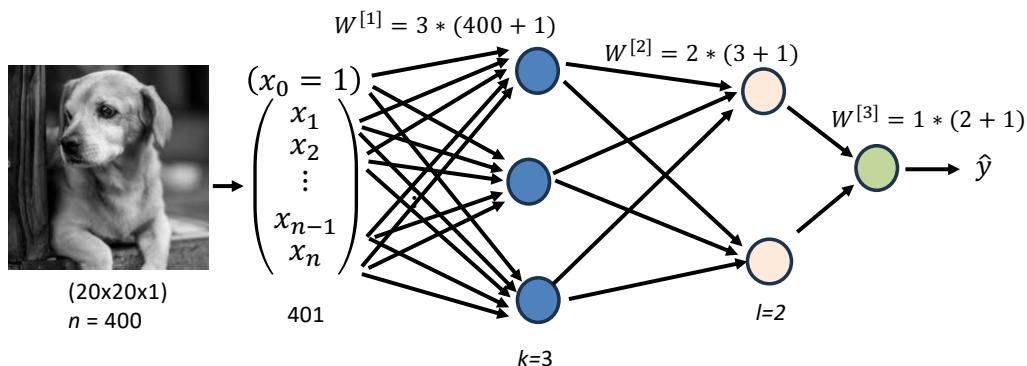
Task 1: Determine if a dog is in the image: $\begin{cases} 1 & \text{if there is a dog} \\ 0 & \text{if there isn't a dog} \end{cases}$



43

ANN for Image classification

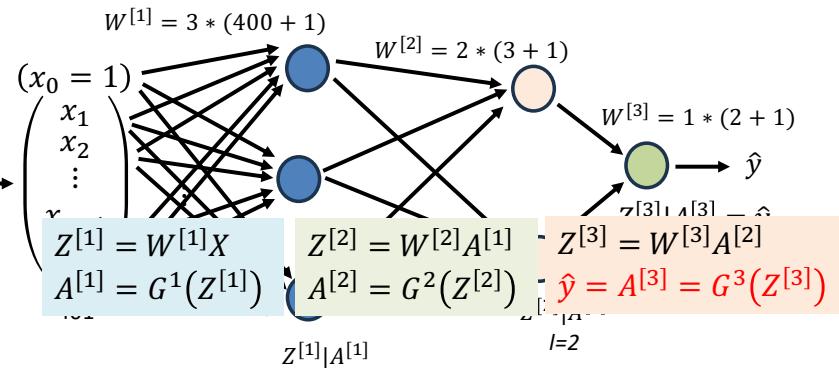
Task 1: Determine if a dog is in the image: $\begin{cases} 1 & \text{if there is a dog} \\ 0 & \text{if there isn't a dog} \end{cases}$



44

ANN for Image classification

Task 1: Determine if a dog is in the image: $\begin{cases} 1 & \text{if there is a dog} \\ 0 & \text{if there isn't a dog} \end{cases}$

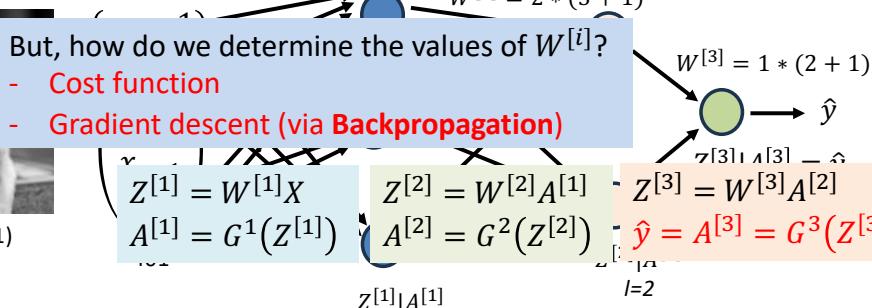


Note: G^i are the activation functions used at the neurons (and can be different at different neurons)

45

ANN for Image classification

If we have the model (i.e., have determined the values of $W^{[i]}$'s), then the steps below enable us to use the model for inference



Note: G^i are the activation functions used at the neurons (and can be different at different neurons)

46

ANN – Cost function

$$J(\mathbf{W}) = J(W^{[1]}, W^{[2]}, W^{[3]}) \\ = \sum_{i=1}^m loss^i(\mathbf{W}) = \sum_{i=1}^m loss(\hat{y}^{(i)}, y^{(i)}; \mathbf{W})$$

For instance, for the image classification problem:

$$loss(\hat{y}^{(i)}, y^{(i)}; \mathbf{W}) = -y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

That is, we want to compute

$$W^{[l]} := W^{[l]} - \alpha \frac{\partial J(\mathbf{W})}{\partial W^{[l]}} \text{ for every layer } l = 1, 2, 3$$

- Repeat the updates $\mathbf{W} := \mathbf{W} - \alpha \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$ until convergence to a minimum

47

ANN – Backpropagation

Remember: $Z^{[3]} = W^{[3]}A^{[2]}$ $J(\mathbf{W}) = \sum_{i=1}^m loss(\hat{y}^{(i)}, y^{(i)}; \mathbf{W})$

Chain rule:

$$\frac{\partial J(\mathbf{W})}{\partial W^{[3]}} = \boxed{\frac{\partial J(\mathbf{W})}{\partial A^{[3]}}} \cdot \boxed{\frac{\partial A^{[3]}}{\partial Z^{[3]}}} \cdot \boxed{\frac{\partial Z^{[3]}}{\partial W^{[3]}}}$$

Similarly:

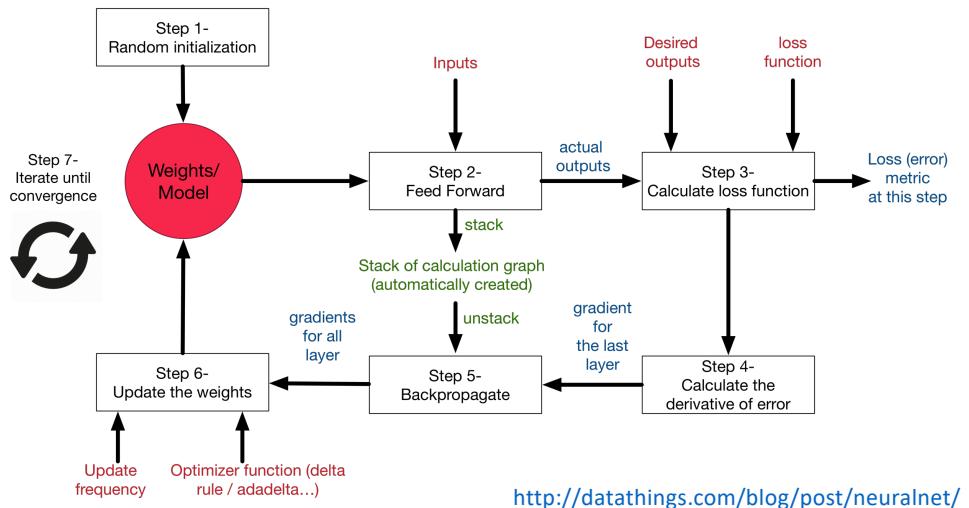
$$Z^{[2]} = W^{[2]}A^{[1]} \quad Z^{[1]} = W^{[1]}X \\ A^{[2]} = G^2(Z^{[2]}) \quad A^{[1]} = G^1(Z^{[1]})$$

$$\frac{\partial J(\mathbf{W})}{\partial W^{[2]}} = \boxed{\frac{\partial J(\mathbf{W})}{\partial Z^{[3]}}} \cdot \boxed{\frac{\partial Z^{[3]}}{\partial A^{[2]}}} \cdot \boxed{\frac{\partial A^{[2]}}{\partial Z^{[2]}}} \cdot \boxed{\frac{\partial Z^{[2]}}{\partial W^{[2]}}}$$

$$\frac{\partial J(\mathbf{W})}{\partial W^{[1]}} = \boxed{\frac{\partial J(\mathbf{W})}{\partial Z^{[2]}}} \cdot \boxed{\frac{\partial Z^{[2]}}{\partial A^{[1]}}} \cdot \boxed{\frac{\partial A^{[1]}}{\partial Z^{[1]}}} \cdot \boxed{\frac{\partial Z^{[1]}}{\partial W^{[1]}}}$$

48

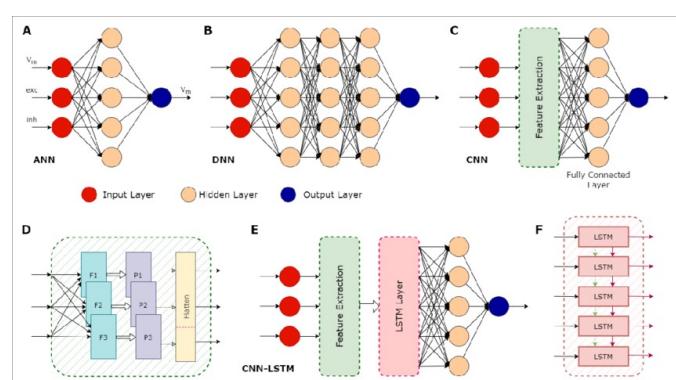
Learning in Multiple-Layer Neural Networks



49

Why DL now?

- Better algorithms & understanding
- Computing power (GPUs, TPUs, ...)
- Data with labels
- Open-source tools and models



Source: <https://www.ip-paris.fr/education/masters/mention-mathematiques-appliquees-statistique/master-year-2-data-science>

50

Why DL now?

- Better algorithms & understanding
- Computing power (GPUs, TPUs, ...)
- Data with labels
- Open-source tools and models



GPU and TPU

Source: <https://www.ip-paris.fr/education/masters/mention-mathematiques-appliquees-statistique/master-year-2-data-science>

51

Why DL now?

- Better algorithms & understanding
- Computing power (GPUs, TPUs, ...)
- Big data
- Open-source tools and models



Source: <https://www.ip-paris.fr/education/masters/mention-mathematiques-appliquees-statistique/master-year-2-data-science>

52

Why DL now?

- Better algorithms & understanding
- Computing power (GPUs, TPUs, ...)
- Data with labels
- Open-source tools and models

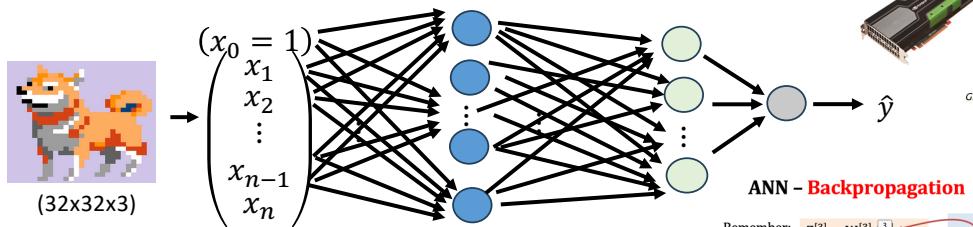


Source: <https://www.ip-paris.fr/education/masters/mention-mathematiques-appliquees-statistique/master-year-2-data-science>

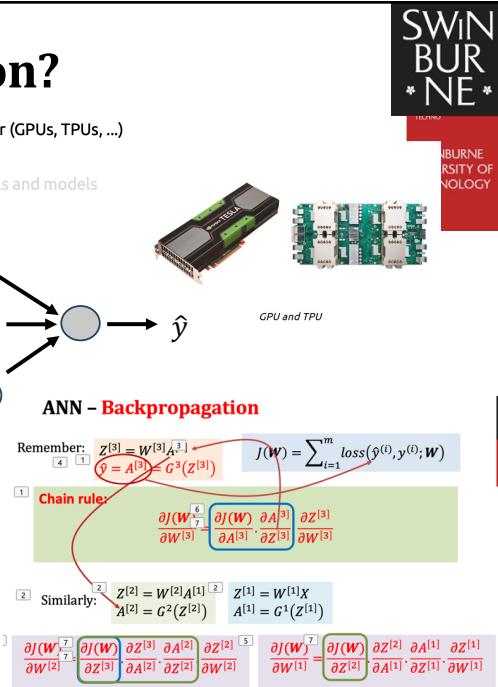
53

Why GPU/TPU? Parallelization?

- Computing power (GPUs, TPUs, ...)
- Data with labels
- Open-source tools and models

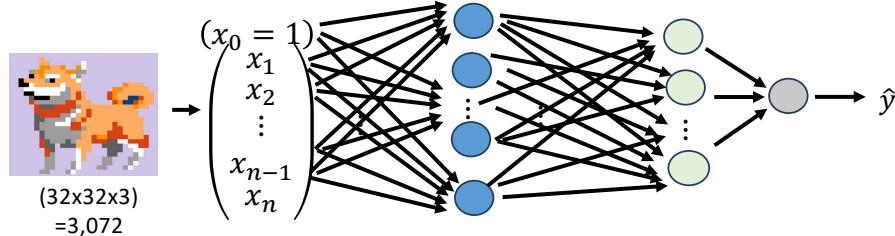


$$J(\mathbf{W}) = \sum_{i=1}^m \text{loss}(\hat{y}^{(i)}, y^{(i)}; \mathbf{W})$$



54

What about batch, epoch, SGD?



$$J(\mathbf{W}) = \sum_{i=1}^m loss(\hat{y}^{(i)}, y^{(i)}; \mathbf{W}) \quad m=10,000,000??$$

55

Wh



(32x3
=3,C

$J(\mathbf{W}) =$

Aspect

Batch Gradient Descent

Stochastic Gradient Descent (SGD)

Data Processing

Uses the whole training dataset to compute the gradient.

Uses a single training sample to compute the gradient.

Convergence Speed

Slower, takes longer to converge.

Faster, converges quicker due to frequent updates.

Convergence Accuracy

More accurate, gives precise gradient estimates.

Less accurate due to noisy gradient estimates.

Computational and Memory Requirements

Requires significant computation and memory.

Requires less computation and memory.

Optimization of Non-Convex Functions

Can get stuck in local minima.

Can escape local minima and find the global minimum.

Suitability for Large Datasets

Not ideal for very large datasets due to slow computation.

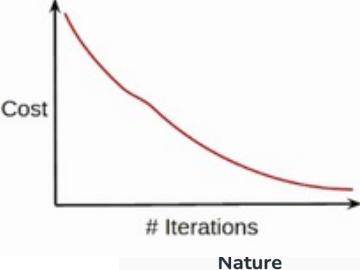
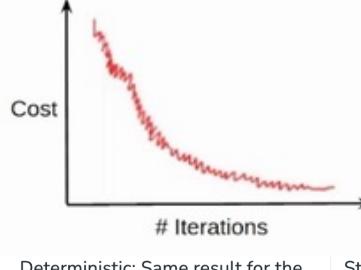
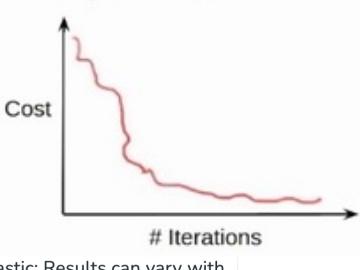
Can handle large datasets effectively.

Nature

Deterministic: Same result for the same initial conditions.

Stochastic: Results can vary with different initial conditions.

56

What	Aspect	Batch Gradient Descent	Stochastic Gradient Descent (SGD)	Swinburne University of Technology
	Data Processing	Uses the whole training dataset to compute the gradient.	Uses a single training sample to compute the gradient.	
	Convergence Speed	Slower, takes longer to converge.	Faster, converges quicker due to frequent updates.	
	Convergence Accuracy	More accurate, gives precise gradient	Less accurate due to noisy	
<ul style="list-style-type: none"> Cost function reduces smoothly  <p>Nature</p>	<ul style="list-style-type: none"> Lot of variations in cost function  <p>Deterministic: Same result for the same initial conditions.</p>	<ul style="list-style-type: none"> Smaller cost function as compared to SGD  <p>Stochastic: Results can vary with different initial conditions.</p>		

57

Developing a Neural Network

- Needs a fair degree of mathematical competence
- Not simple (also very time consuming)
- Considerations include:
 - Identification of features in the problem domain worth mapping into inputs/outputs (Problem Representation)
 - Determining the number of hidden layers and nodes
 - Setting up initial parameters – weights, thresholds, increments
 - What learning method to use?
- No rules exist, based on experience (and solid understanding of problem domain)

58

Weakness of Neural Networks

- Neural networks are *opaque* – it is very hard to check that the weights after training are sensible
- A neural network can never explain how it arrived at a particular conclusion
- Loss of human control, lack of trust since AI cannot explain all decisions and actions, or take responsibility – **Black Box AI**
- Some work is being done on combating the opaqueness problem
- Partnership on AI (<https://www.partnershiponai.org>) is a step towards safe and trustworthy AI technologies – supported by tech giants including Google, IBM and Microsoft

59

Machine Learning as a Service

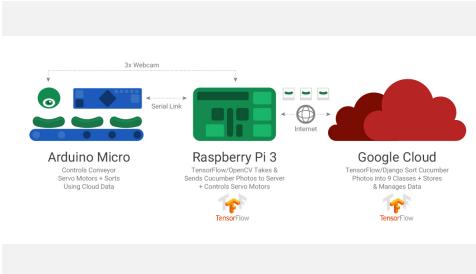
- Machine learning is no longer for experts only
- Software developers can use cutting-edge, commercially usable machine learning frameworks & related libraries:
 - Scikit-learn
 - Google Tensorflow, Keras
 - PyTorch
 - Microsoft Cognitive Toolkit
 - Hugging Face Transformers
 - FastAI
 - NumPy, Pandas, Matplotlib, Seaborn, etc.
- Major cloud providers including Amazon, Google and Microsoft offer machine learning as a service

60

Cucumber Farming & TensorFlow

SWIN
BUR
* NE *

SWINBURNE
UNIVERSITY OF
TECHNOLOGY



- Cucumber classification into 9 different classes using TensorFlow based on size, thickness, colour, texture, small scratches, whether they are crooked, whether they have prickles, etc.

<https://cloud.google.com/blog/products/ai-machine-learning/how-a-japanese-cucumber-farmer-is-using-deep-learning-and-tensorflow>