



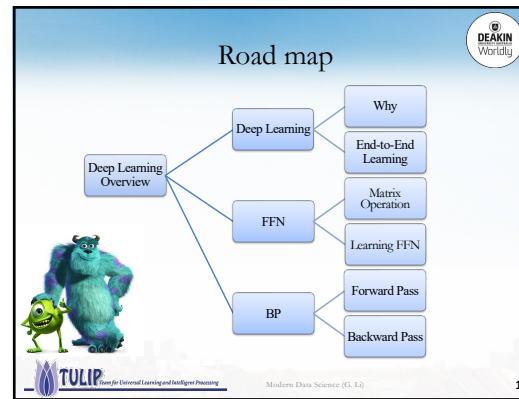
Lecture Notes on
Pattern Recognition

Module 09(A): Deep Learning Overview

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia





0



1

Deep Learning



- Deep Learning
- Why Deep Learning?
 - Modularization
 - End-to-End Learning

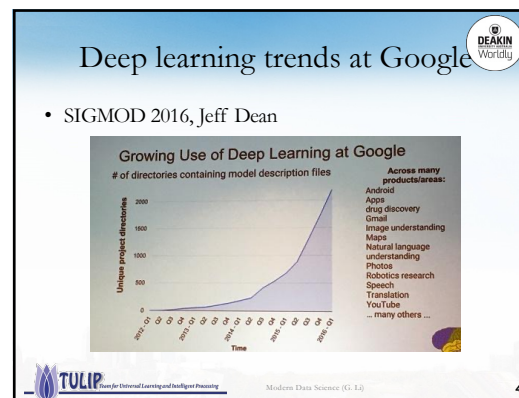
TULIP

Modern Data Science (G. Li)

2



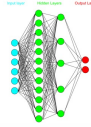
3



4

What is Deep Learning?

- Fast answer:
 - Fast answer: multilayer perceptron (aka deep neural networks) of the 1980s rebranded in 2006.
 - But early nets go stuck at 1-2 hidden layers.
- Slow answer:
 - distributed representation, multiple steps of computation, modelling the compositionality of the world, a better prior, advances in compute, data & optimization, neural architectures, etc.



TULIP

Modern Data Science (G. Li)

5

From ANN to Deep Learning

- 1958: Perceptron (linear model)
- 1969: Perceptron has limitation
- 1980s: Multi-layer perceptron
 - Do not have significant difference from DNN today
- 1986: Backpropagation
 - Usually more than 3 hidden layers is not helpful
- 1989: 1 hidden layer is “good enough”
 - why deep?



6

From ANN to Deep Learning

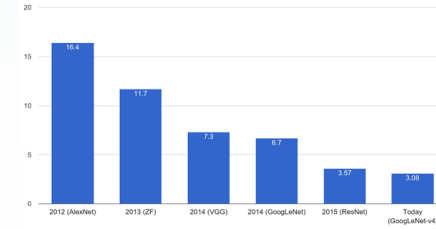
- 2006: RBM initialization
- 2009: GPU
- 2012: win ILSVRC image competition
- 2015.2:
 - Image recognition surpassing human-level performance
- 2016.3: Alpha GO beats Lee Sedol
- 2016.10:
 - Many system as good as humans



7

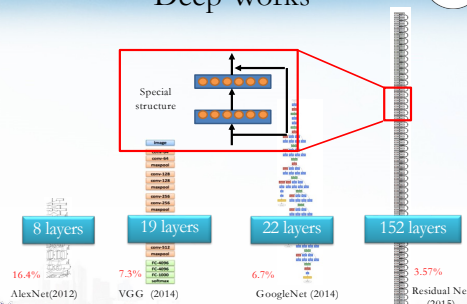
Deep Works

ImageNet Classification Error (Top 5)



8

Deep Works



9

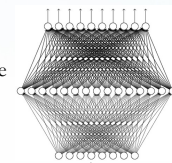
The Best of Machine Learning

- Strong/flexible priors:
 - Good features: Feature engineering
 - Data structure: HMM, CRF, MRF, Bayesian nets
 - Model structure, VC-dimension, regularization, sparsity:
 - SVM, compressed sensing
 - Manifold assumption, class/region separation:
 - Metric + semi-supervised learning
 - Factors of variation: PCA, ICA, FA
- Uncertainty quantification:
 - Bayesian, ensemble: RF, GBM
- Sharing statistical strength:
 - model reuse: transfer learning, domain adaption, multitask learning, lifelong learning

10

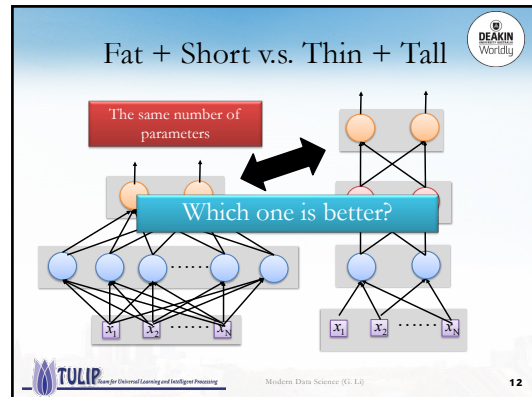
Universality Theorem

- Shallow network can represent any function.
- Given **enough** hidden neurons, Any continuous function $f: R^N \rightarrow R^M$ can be realized by a network with one hidden layer
- However, using deep structure is more effective.



Reference for the reason:
<http://neuralnetworksanddeeplearning.com/chap4.html>

11



12

Fat + Short v.s. Thin + Tall

Layer X Size	Word Error Rate (%)	Layer X Size	Word Error Rate (%)
1 X 2k	24.2		
2 X 2k	20.4		
3 X 2k	18.4		
4 X 2k	17.8		
5 X 2k	17.2	1 X 3772	22.5
7 X 2k	17.1	1 X 4634	22.6
		1 X 16k	22.1

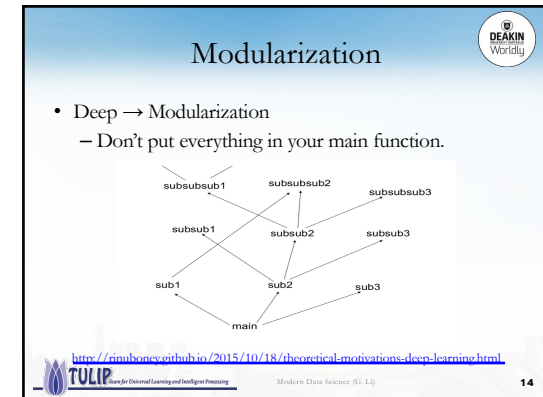
Why?

Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

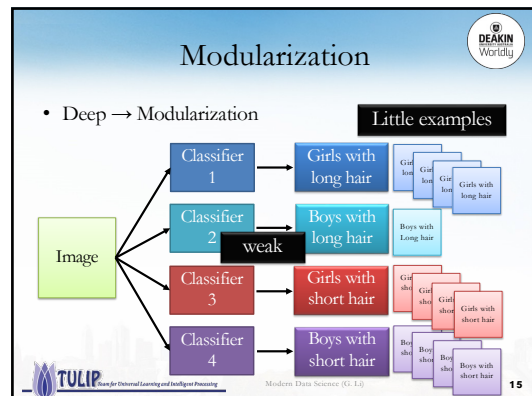
TULIP: Towards Universal Learning and Intelligent Processing

Modern Data Science (G. Li)

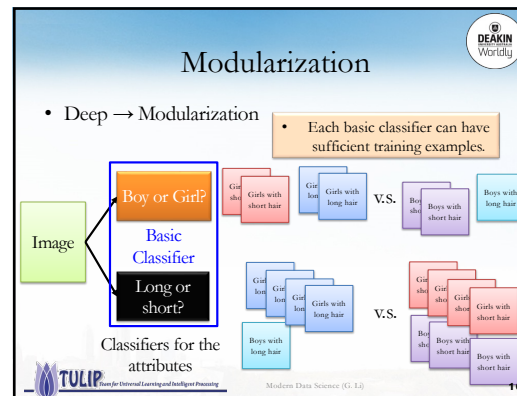
13



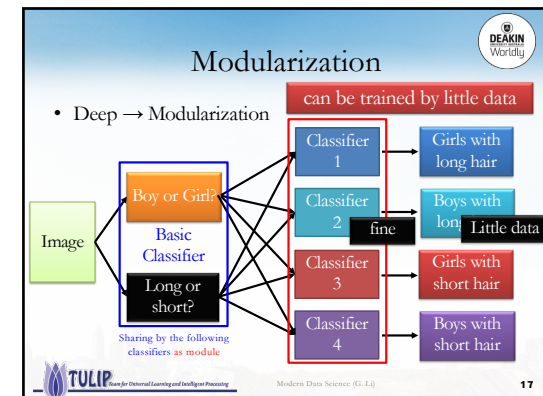
14



15



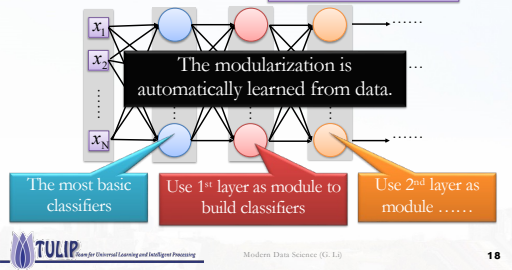
16



17

Modularization

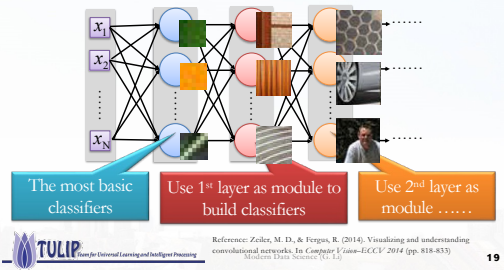
- Deep \rightarrow Modularization → Less training data?



18

Modularization

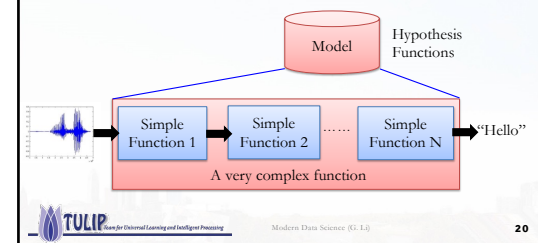
- Deep \rightarrow Modularization



19

End-to-end Learning

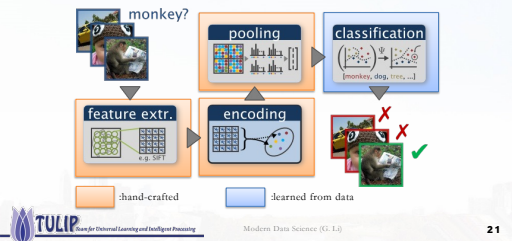
- End-to-end training:
 - What each function should do is learned automatically



20

End-to-end Learning

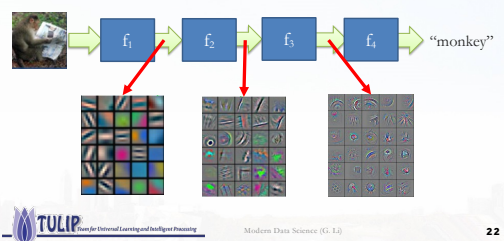
- Shallow Approach
 - Each box is a simple function in the production line:



21

End-to-end Learning

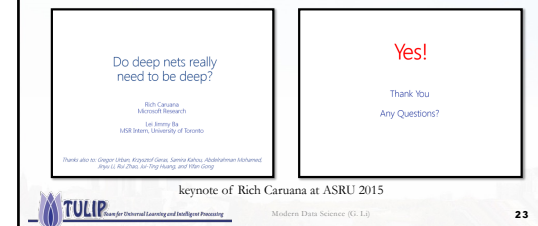
- Deep Learning
 - All functions are learned from data



22

To learn more ...

- Do Deep Nets Really Need To Be Deep? (by Rich Caruana)
 - <http://research.microsoft.com/apps/video/default.aspx?id=232373&r=1>



23

Tensorflow or MXNet?

	Languages	MultiGPU	Distributed	Mobile	Runtime Engine
Tensorflow	Python	Yes	No	Yes	?
MXNET	Python, R, Julia, Go	Yes	Yes	Yes	Yes

	Tensorflow (GoogleNet ES-1650/980)	Torch7	Caffe	MXNET
Time	940ms	172ms	170ms	180ms
Memory	all (OOM after 24)	2.1GB	2.2GB	1.6GB

Marlane Labs Carnegie Mellon University

Tensorflow vs. MXNET

K + TensorFlow

K mxnet

TULIP Tulip Institute for Data Science and Intelligent Processing

Modern Data Science (G. L.)

24

24

FFN

- FFN
- Matrix Operation
- Learning

TULIP Tulip Institute for Data Science and Intelligent Processing

Modern Data Science (G. L.)

25

25

Deep Learning is so simple

Step 1: define a set of function → Step 2: goodness of function → Step 3: pick the best function

TULIP Tulip Institute for Data Science and Intelligent Processing

Modern Data Science (G. L.)

26

26

Neural Network

Neural Network:

- Different connection leads to different network structures
- Network parameter θ : all the weights and biases in the "neurons"

TULIP Tulip Institute for Data Science and Intelligent Processing

Modern Data Science (G. L.)

27

27

Fully Connected Feedforward Network

Sigmoid Function

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

TULIP Tulip Institute for Data Science and Intelligent Processing

Modern Data Science (G. L.)

28

28

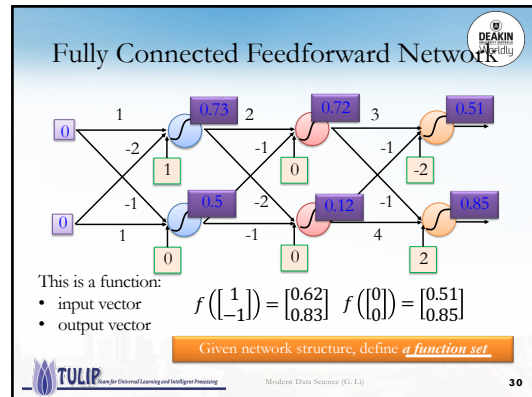
Fully Connected Feedforward Network

TULIP Tulip Institute for Data Science and Intelligent Processing

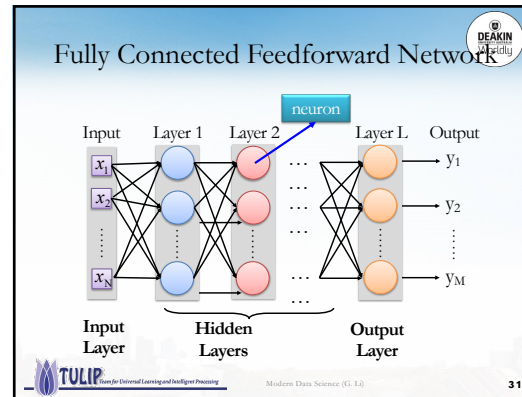
Modern Data Science (G. L.)

29

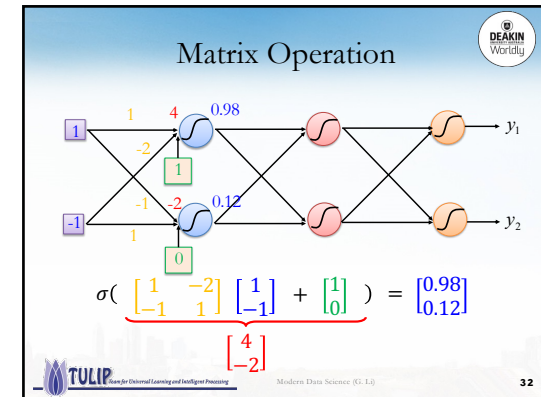
29



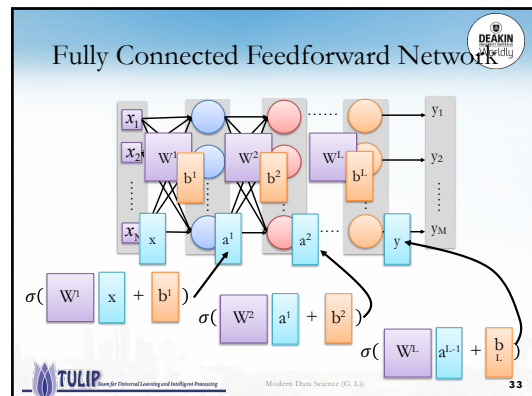
30



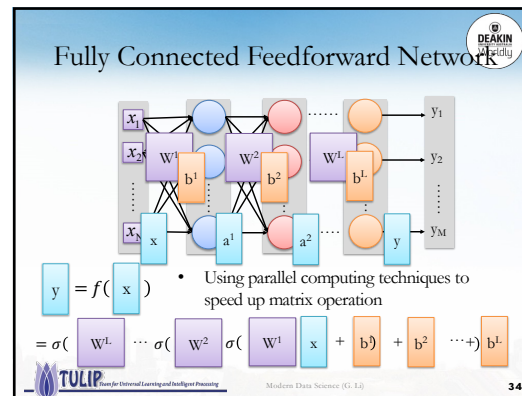
31



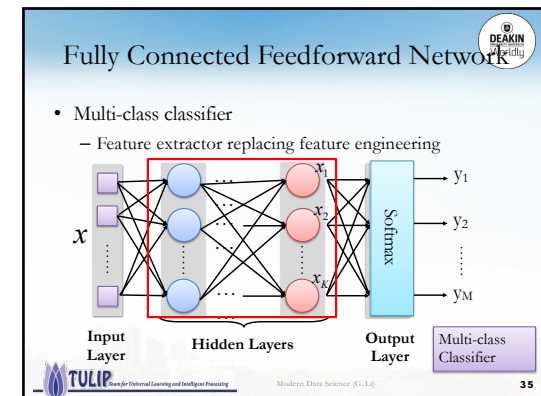
32



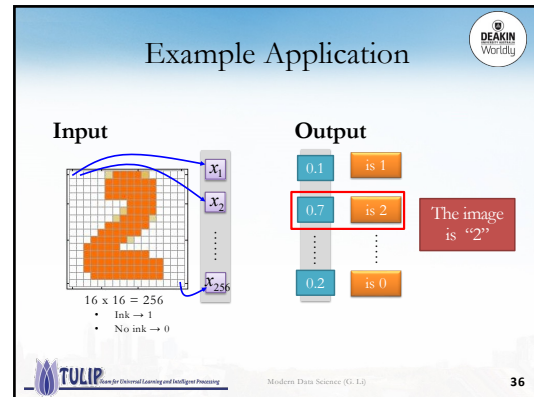
33



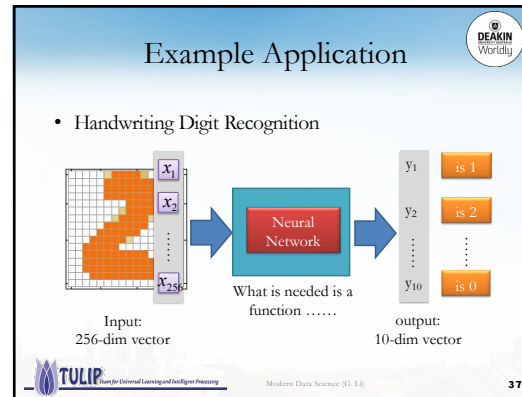
34



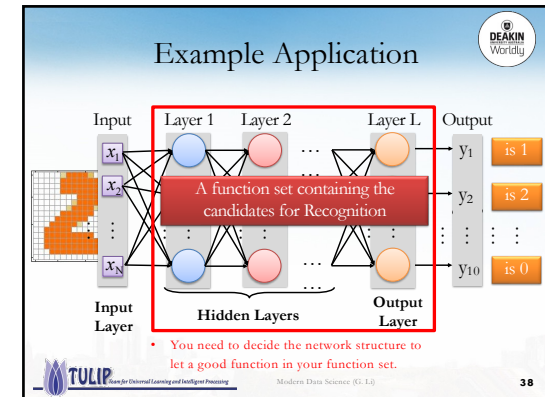
35



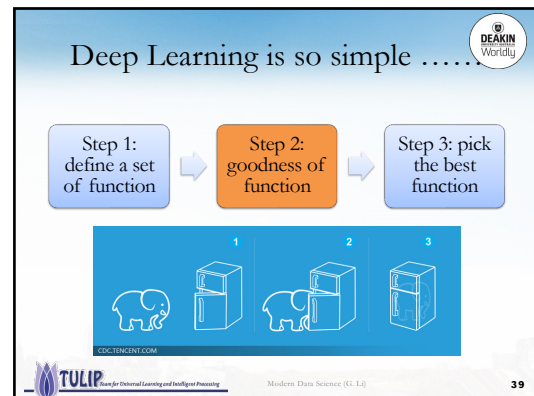
36



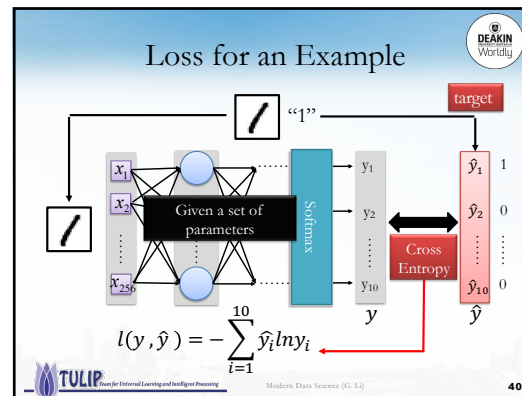
37



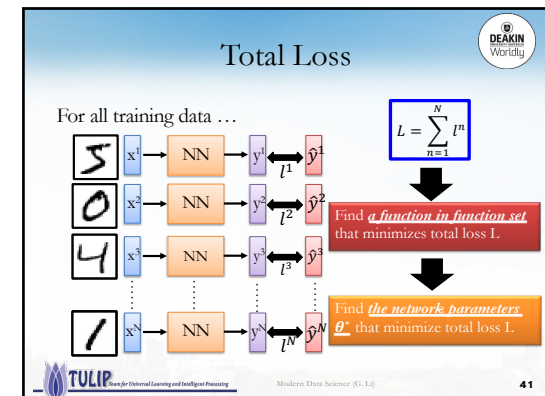
38



39



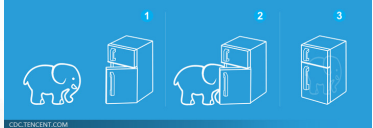
40



41

Deep Learning is so simple

Step 1: define a set of function → Step 2: goodness of function → Step 3: pick the best function

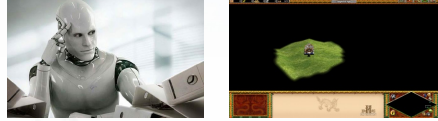


TULIP: Tools for Universal Learning and Intelligent Processing
Modern Data Science (G. L.)

42

Gradient Descent

- Gradient Descent is the “learning” of machines in deep learning
- Even alpha go using this approach.



TULIP: Tools for Universal Learning and Intelligent Processing
Modern Data Science (G. L.)

43

Gradient Descent

θ

w_1	0.2	Compute $\partial L / \partial w_1$ $-\mu \partial L / \partial w_1$	0.15
w_2	-0.1	Compute $\partial L / \partial w_2$ $-\mu \partial L / \partial w_2$	0.05
\vdots			
b_1	0.3	Compute $\partial L / \partial b_1$ $-\mu \partial L / \partial b_1$	0.2
\vdots			

$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial b_1} \\ \vdots \end{bmatrix}$$

gradient

TULIP: Tools for Universal Learning and Intelligent Processing
Modern Data Science (G. L.)

44

Gradient Descent

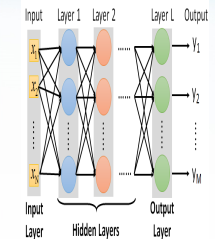
θ

w_1	0.2	Compute $\partial L / \partial w_1$ $-\mu \partial L / \partial w_1$	0.15	Compute $\partial L / \partial w_1$ $-\mu \partial L / \partial w_1$	0.09
w_2	-0.1	Compute $\partial L / \partial w_2$ $-\mu \partial L / \partial w_2$	0.05	Compute $\partial L / \partial w_2$ $-\mu \partial L / \partial w_2$	0.15
\vdots					
b_1	0.3	Compute $\partial L / \partial b_1$ $-\mu \partial L / \partial b_1$	0.2	Compute $\partial L / \partial b_1$ $-\mu \partial L / \partial b_1$	0.10
\vdots					

TULIP: Tools for Universal Learning and Intelligent Processing
Modern Data Science (G. L.)

45

Back Propagation



- Forward Pass
- Backward Pass

TULIP: Tools for Universal Learning and Intelligent Processing
Modern Data Science (G. L.)

46

Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

Starting Parameters $\theta^0 \rightarrow \theta^1 \rightarrow \theta^2 \rightarrow \dots$

$\nabla L(\theta)$

Compute $\nabla L(\theta^0)$	$\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$
Compute $\nabla L(\theta^1)$	$\theta^2 = \theta^1 - \eta \nabla L(\theta^1)$

Millions of parameters

To compute the gradients efficiently, we use backpropagation.

TULIP: Tools for Universal Learning and Intelligent Processing
Modern Data Science (G. L.)

47

Chain Rule

Case 1 $y = g(x)$ $z = h(y)$

$$\Delta x \rightarrow \Delta y \rightarrow \Delta z \quad \frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

Case 2 $x = g(s)$ $y = h(s)$ $z = k(x, y)$

$$\Delta s \rightarrow \Delta x \rightarrow \Delta z \quad \frac{dz}{ds} = \frac{\partial z}{\partial x} \frac{dx}{ds} + \frac{\partial z}{\partial y} \frac{dy}{ds}$$

TULIP Centre for Data Science and Intelligent Processing Modern Data Science (G1-L1) 48

48

Backpropagation

$L(\theta) = \sum_{n=1}^N l^n(\theta) \rightarrow \frac{\partial L(\theta)}{\partial w} = \sum_{n=1}^N \frac{\partial l^n(\theta)}{\partial w}$

TULIP Centre for Data Science and Intelligent Processing Modern Data Science (G1-L1) 49

49

Backpropagation

$z = x_1 w_1 + x_2 w_2 + b$

Forward pass:
Compute $\partial z / \partial w$ for all parameters

Backward pass:
Compute $\partial l / \partial z$ for all activation function inputs z
(Chain rule)

TULIP Centre for Data Science and Intelligent Processing Modern Data Science (G1-L1) 50

50

Backpropagation – Forward pass

Compute $\partial z / \partial w$ for all parameters

$z = x_1 w_1 + x_2 w_2 + b$

$\frac{\partial z}{\partial w_1} = ?$ x_1
 $\frac{\partial z}{\partial w_2} = ?$ x_2

The value of the input connected by the weight

TULIP Centre for Data Science and Intelligent Processing Modern Data Science (G1-L1) 51

51

Backpropagation – Forward pass

Compute $\partial z / \partial w$ for all parameters

$\frac{\partial z}{\partial w} = -1$ $\frac{\partial z}{\partial w} = 0.12$ $\frac{\partial z}{\partial w} = 0.11$

TULIP Centre for Data Science and Intelligent Processing Modern Data Science (G1-L1) 52

52

Backpropagation – Backward pass

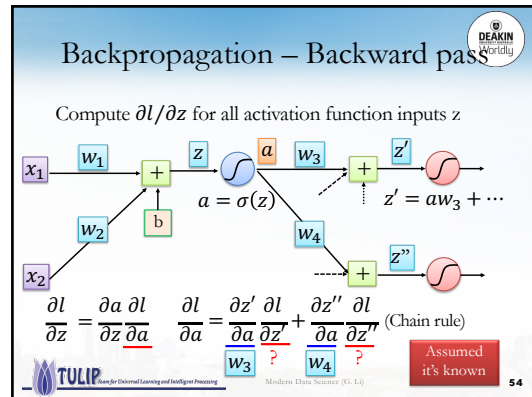
Compute $\partial l / \partial z$ for all activation function inputs z

$a = \sigma(z)$

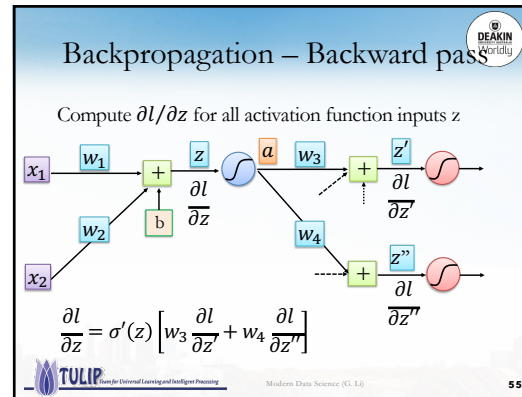
$\frac{\partial l}{\partial z} = \frac{\partial a}{\partial z} \frac{\partial l}{\partial a} \rightarrow \sigma'(z)$

TULIP Centre for Data Science and Intelligent Processing Modern Data Science (G1-L1) 53

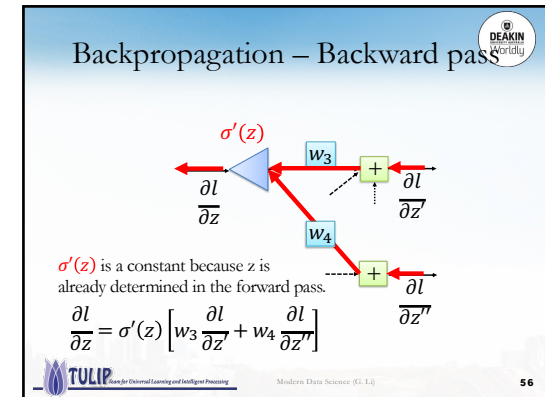
53



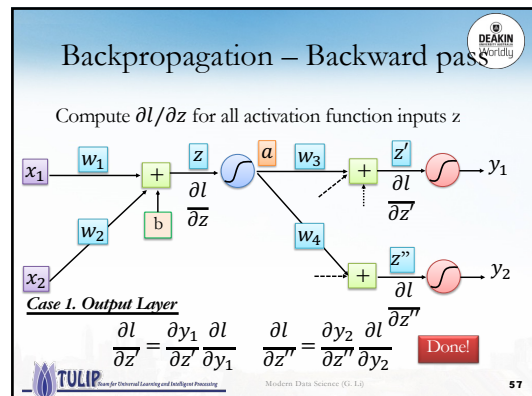
54



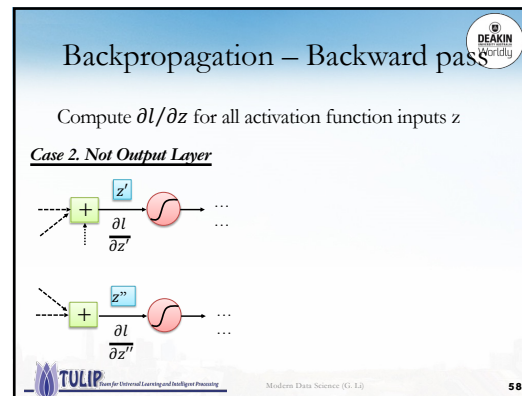
55



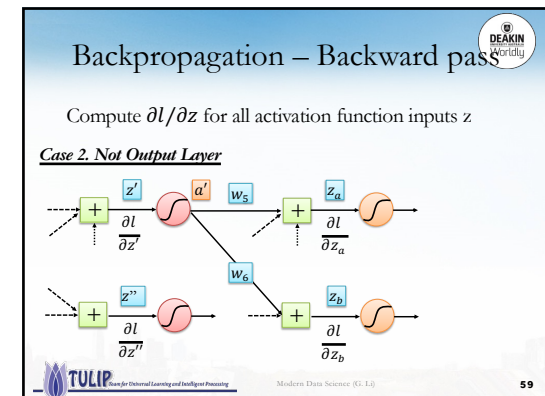
56



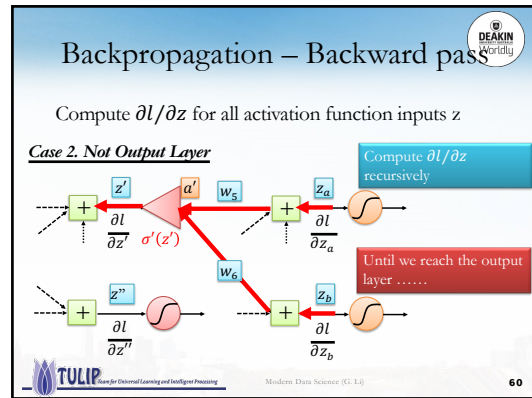
57



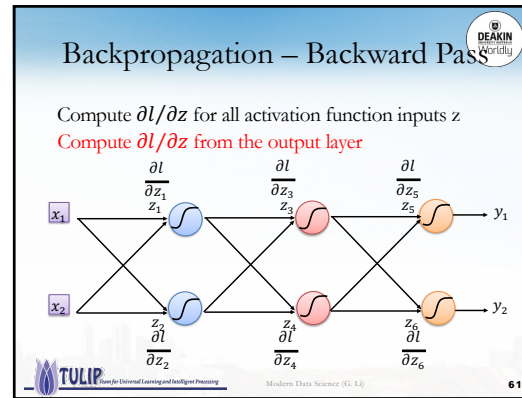
58



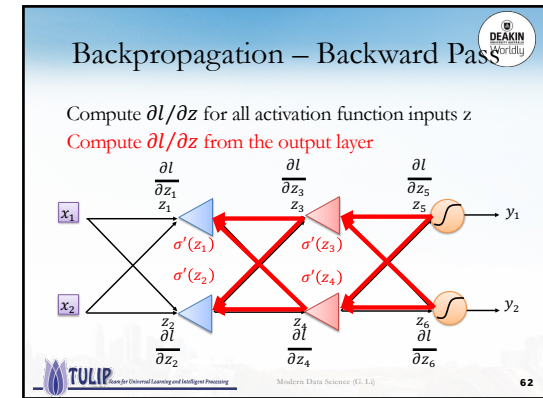
59



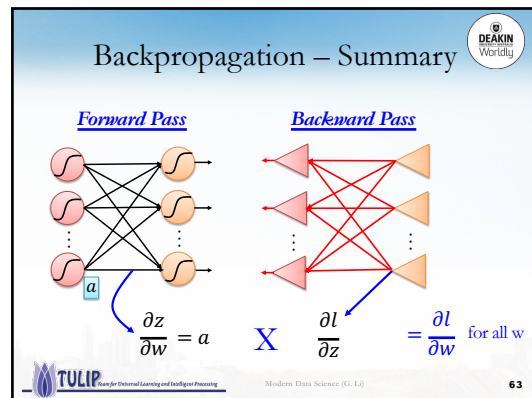
60



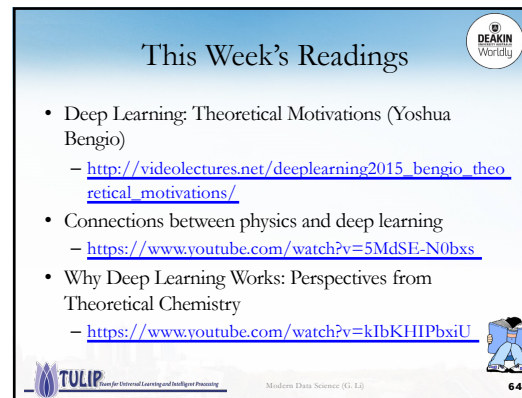
61



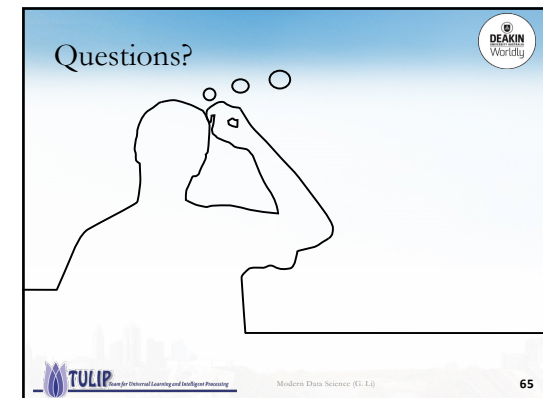
62



63



64



65