

Brain Intelligence and Artificial Intelligence

人脑智能与机器智能

Lecture 5 – A brief intro to AI & Hands-on

Quanying Liu (刘泉影)

SUSTech, BME department

Email: liuqy@sustech.edu.cn

Lecture 5 – Intro to AI & hands on

- **A general introduction to AI**
 - 3 key components
 - The network **Architecture**
- **AI learning: Gradient Descent (GD) to minimize a *loss function***
 - What is Gradient Descent?
 - Gradient Descent to train deep NNs → Error Backpropagation
- **Hands-on** (pytorch), thanks to three TAs
- **Error Back-propagation (BP)** in fully-connected NN
 - Backpropagation – forward pass
 - Backpropagation – backward pass
 - BP in the brain

3 Key Components

1. Data

- Without labels: a data point (x), a dataset (X)
- With labels: a pair of data (x, y), a dataset (X, Y)

2. Model

- **Network**: the high-level concept, a function (input, output)
- **Loss function** (objective function) : the goal of machine learning
- **Parameters**: the variables to be learned/estimated

3. Optimization algorithm

- The algorithm to minimize the loss function (or maximize the objective function)

Benchmark datasets in computer vision (机器视觉)



- 14 million images and 1000 categories.
- Largest database of labeled images.



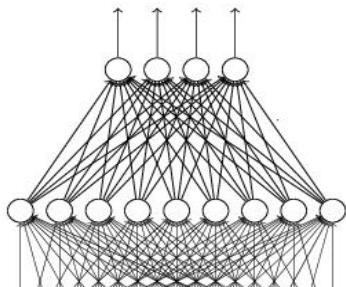
- Images in Fish category.
- Captures variations of fish.

The network architecture

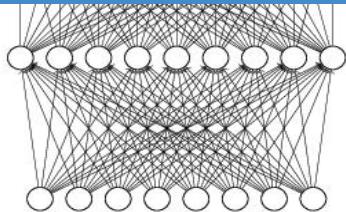
- MLP (fully-connected)
- CNN (convolution)
- GNN (graph)
- Transformer (attention)
- ...

The architecture of CNN

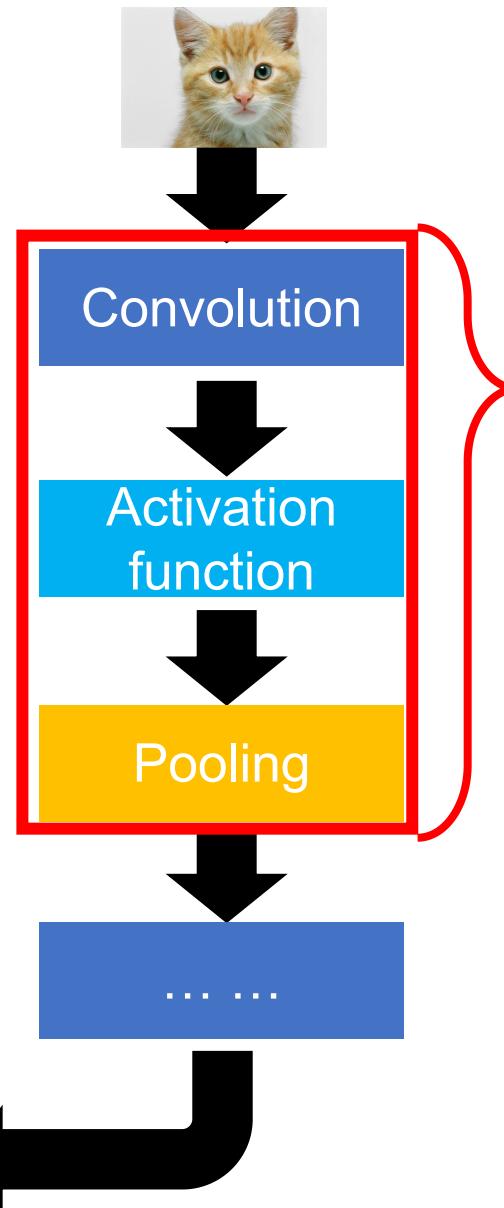
cat / dog



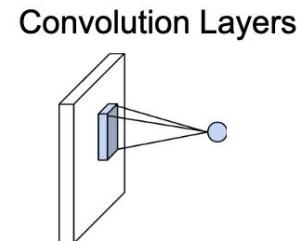
Fully Connected
Feedforward network



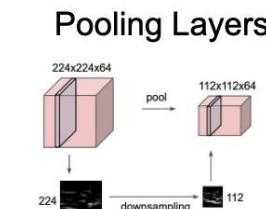
Flatten



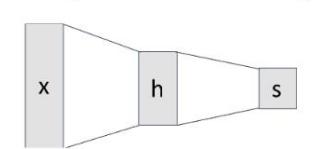
A block involves some core **components**.
These blocks can repeat many times.



Convolution Layers

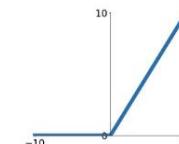


Pooling Layers



Fully-Connected Layers

Activation Function

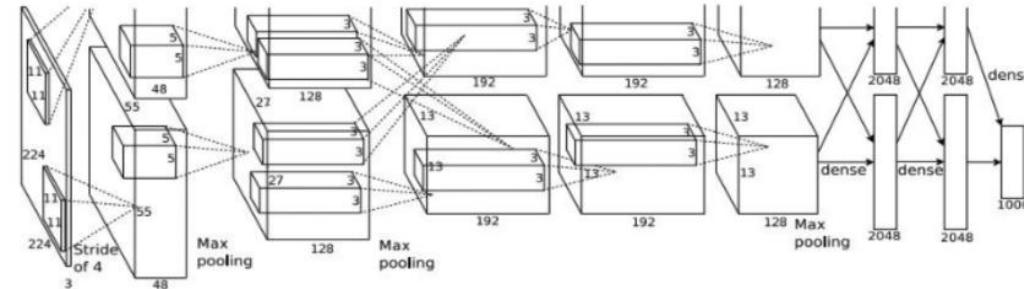


Normalization

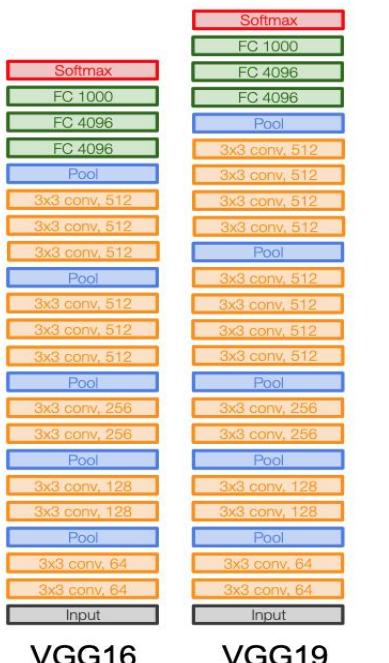
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

The architecture of CNN

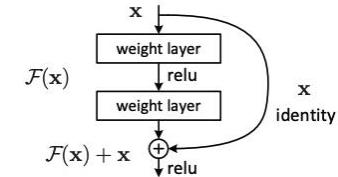
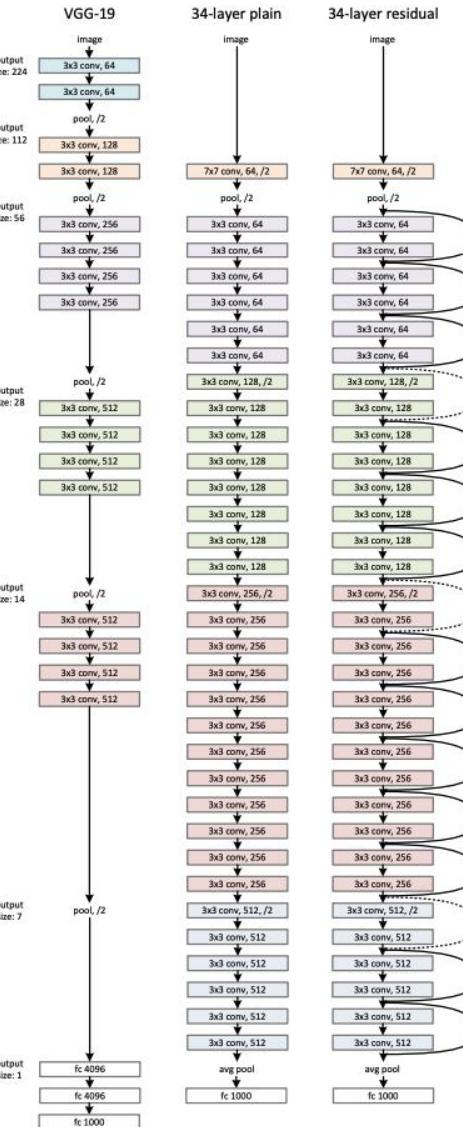
AlexNet



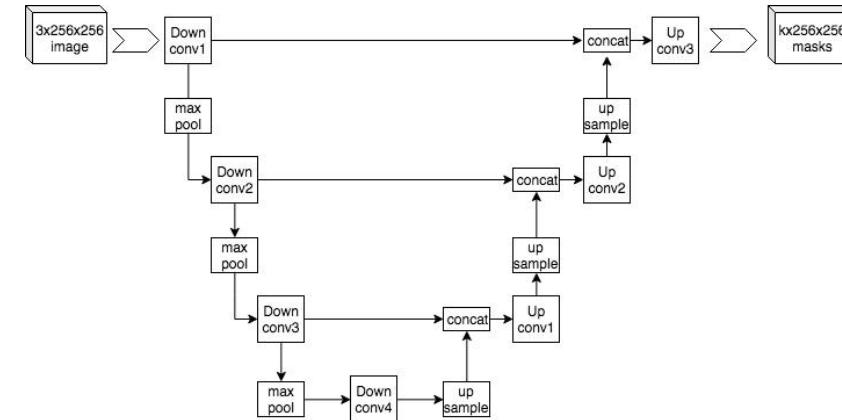
VGGNet



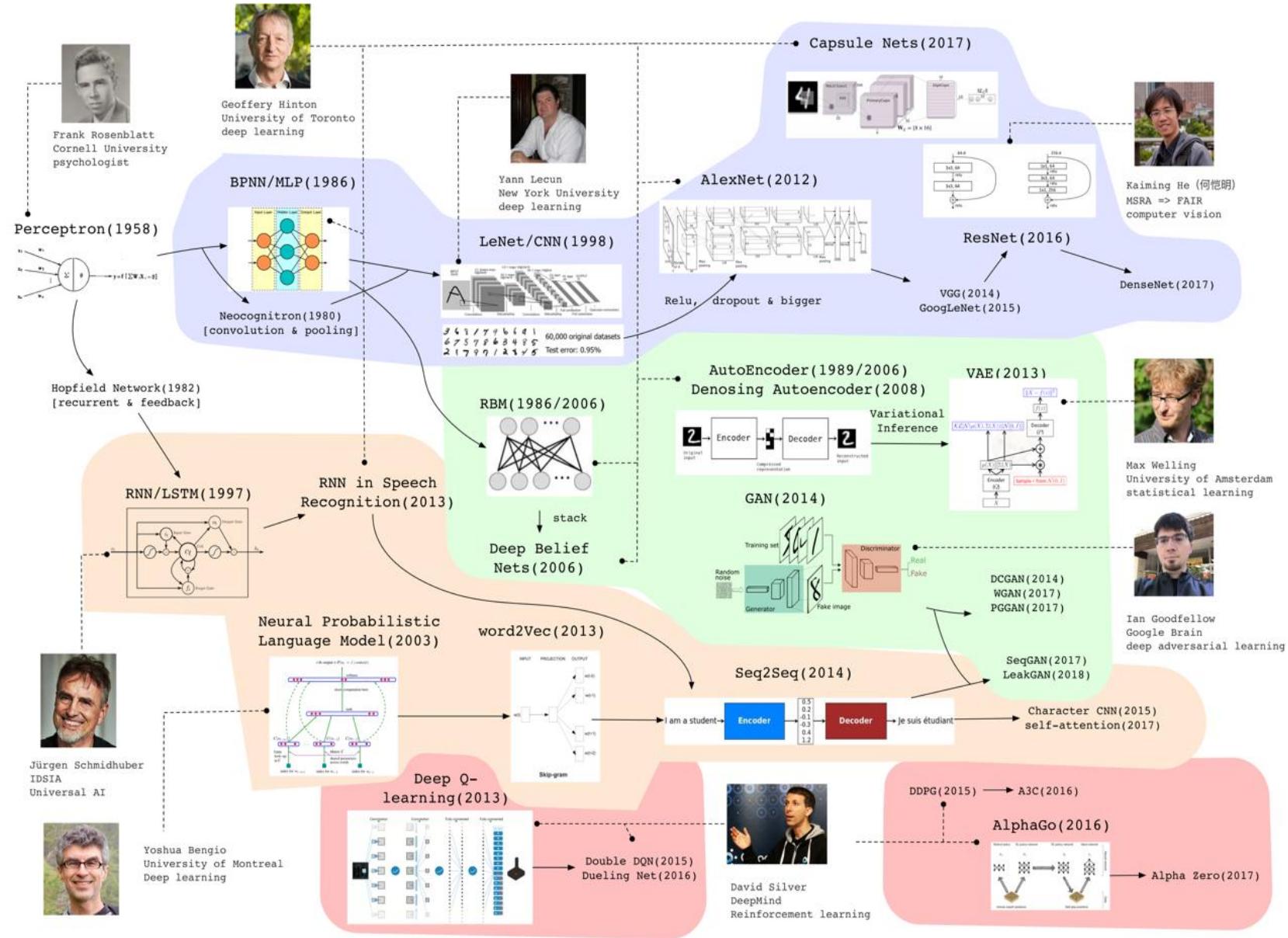
AlexNet



U-Net



《2020 中国人工智能白皮书》



Computer Vision and CNN

- 1979, *Neocognitron* by Fukushima
- 1986, *Backpropogation MLP* by Hinton
- 1998, *LeNet-5* by LeCun
- 2012, *AlexNet* by Hinton
- 2016, *ResNet* by 何恺明

Generative Models

- 1986–2006, *RBM* by Hinton
- 1989/2006, *AutoEncoder* by Hinton
- 2014, *GAN* by Goodfellow

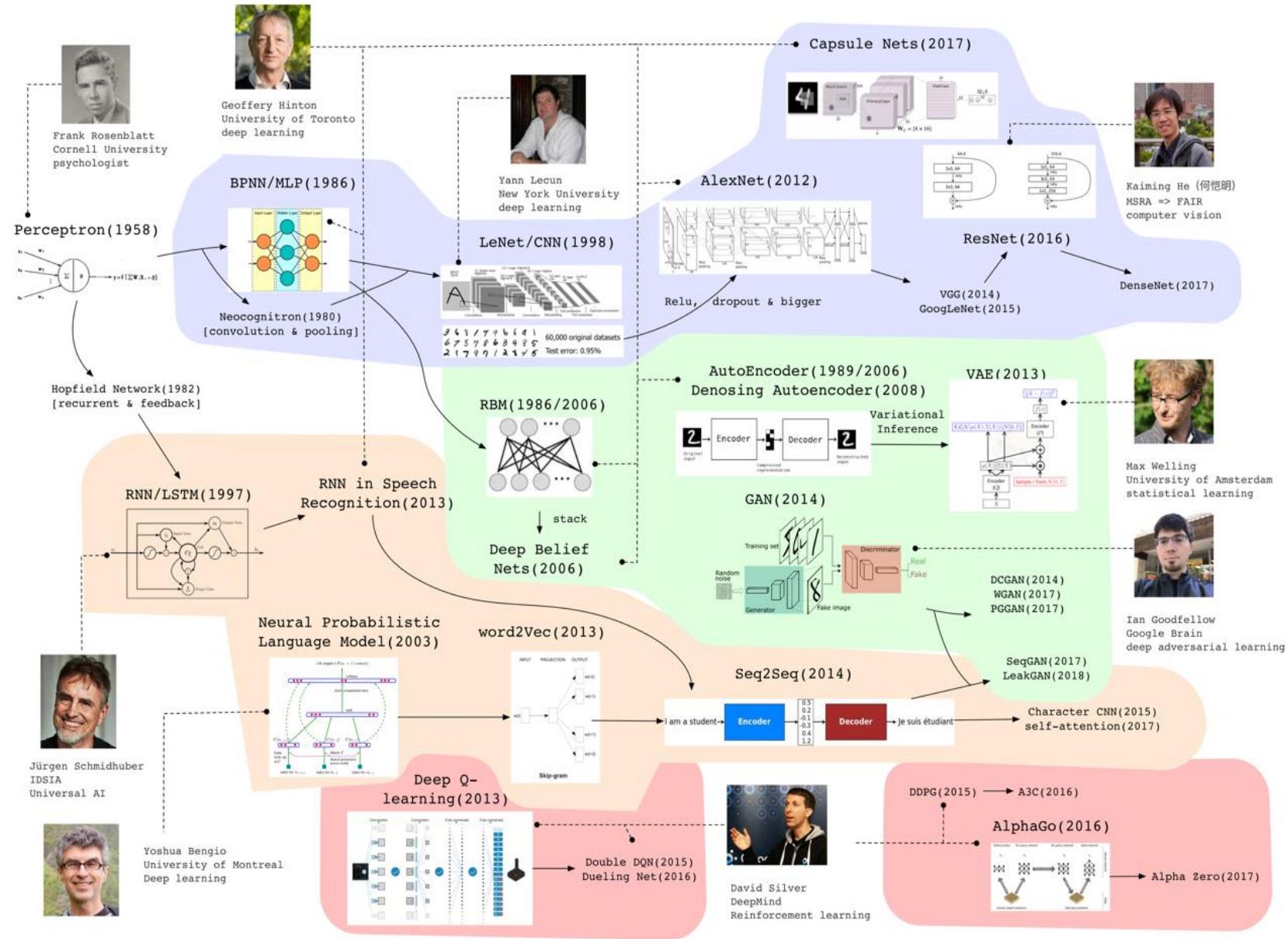
Sequence Models

- 1982, *Hopfield Network* by Hopfield
- 1997, *LSTM* by Schmidhuber
- 2013, *RNN* by Hinton
- 2017, *transformer* by google

Reinforcement Learning

- 2013, *Deep Q-learning* by Silver
- 2016, *AlphaGo* by deepmind

《2020 中国人工智能白皮书》



Computer Vision and CNN

- 1979, *Neocognitron* by Fukushima
- 1986, *Backpropagation MLP* by Hinton
- 1998, *LeNet-5* by LeCun
- 2012, *AlexNet* by Hinton
- 2016, *ResNet* by 何恺明
- 2020, *ViT* by google

Generative Models

- 1986–2006, *RBM* by Hinton
- 1989/2006, *AutoEncoder* by Hinton
- 2014, *GAN* by Goodfellow
- 2020, *Diffusion model* by Ho & Ahbeel

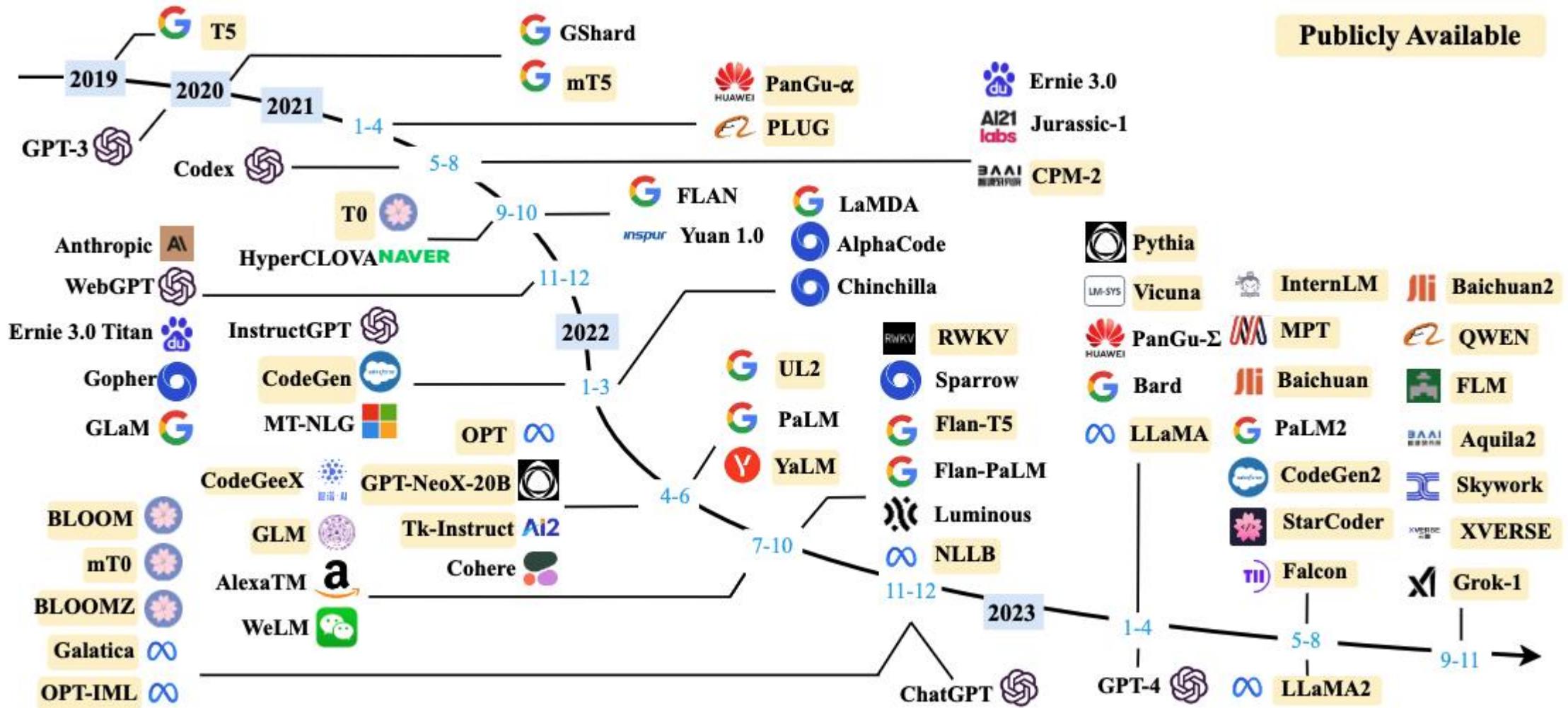
Sequence Models

- 1982, *Hopfield Network* by Hopfield
- 1997, *LSTM* by Schmidhuber
- 2013, *RNN* by Hinton
- 2017, *transformer* by google
- 2022, *chatGPT* by openAI

Reinforcement Learning

- 2013, *Deep Q-learning* by Silver
- 2016, *AlphaGo* by deepmind
- 2020, *AlphaFold* by deepmind

大模型时代



The Rise and Rise of A.I. Large Language Models (LLMs) & their associated bots like ChatGPT

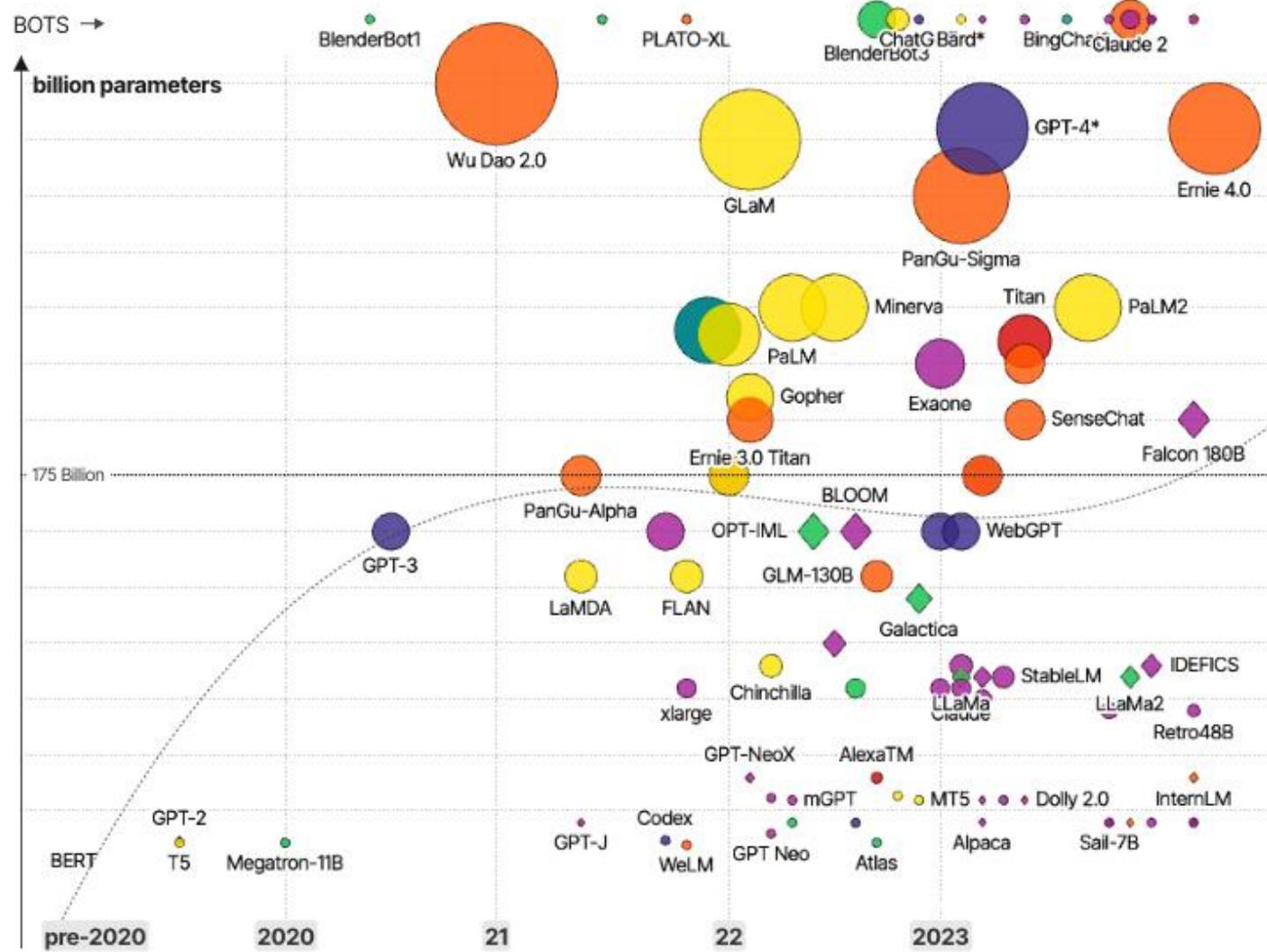


size = no. of parameters



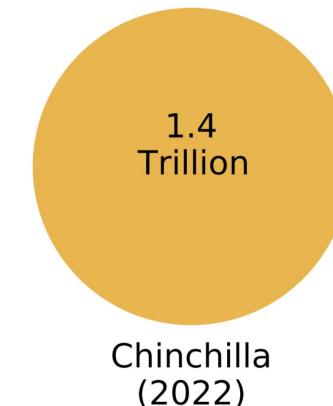
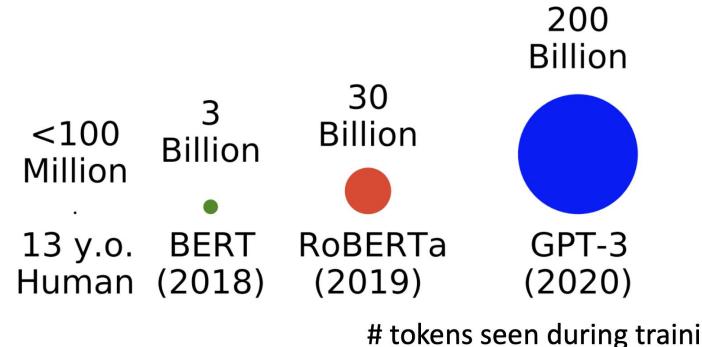
open-access

- Amazon-owned
- Chinese
- Google
- Meta / Facebook
- Microsoft
- OpenAI
- Other



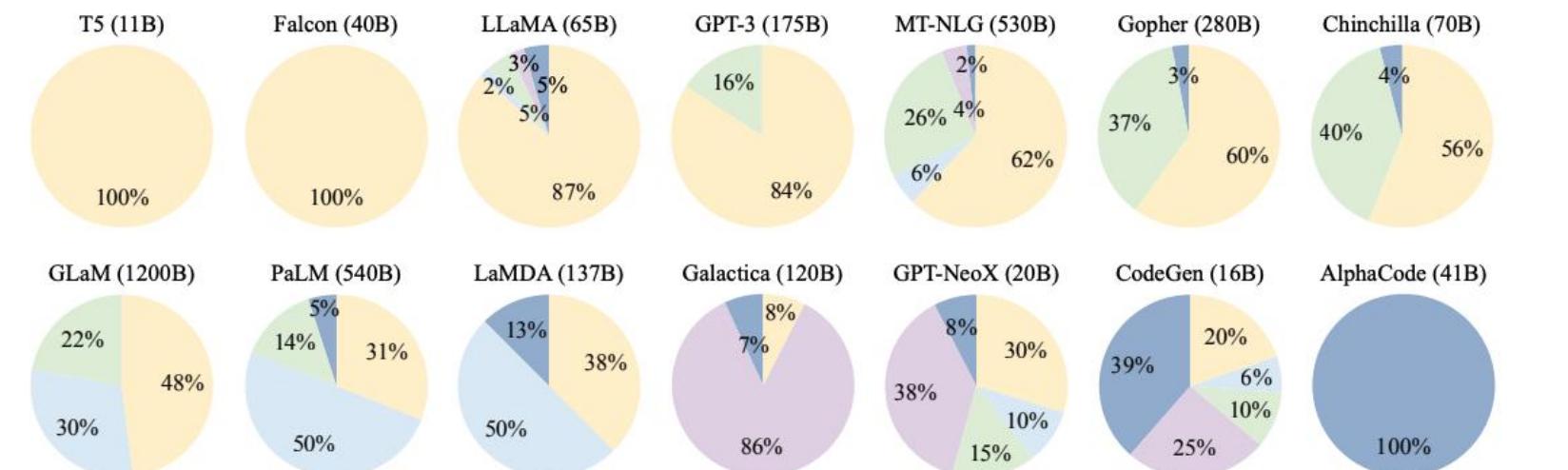
AI模型的训练数据越来越大

<https://babylm.github.io/>



训练数据越来越大

数据种类越来越丰富



Webpages

Conversation Data

Books & News

Scientific Data

Code

C4 (800G, 2019), OpenWebText (38G, 2023), Wikipedia (21G, 2023)

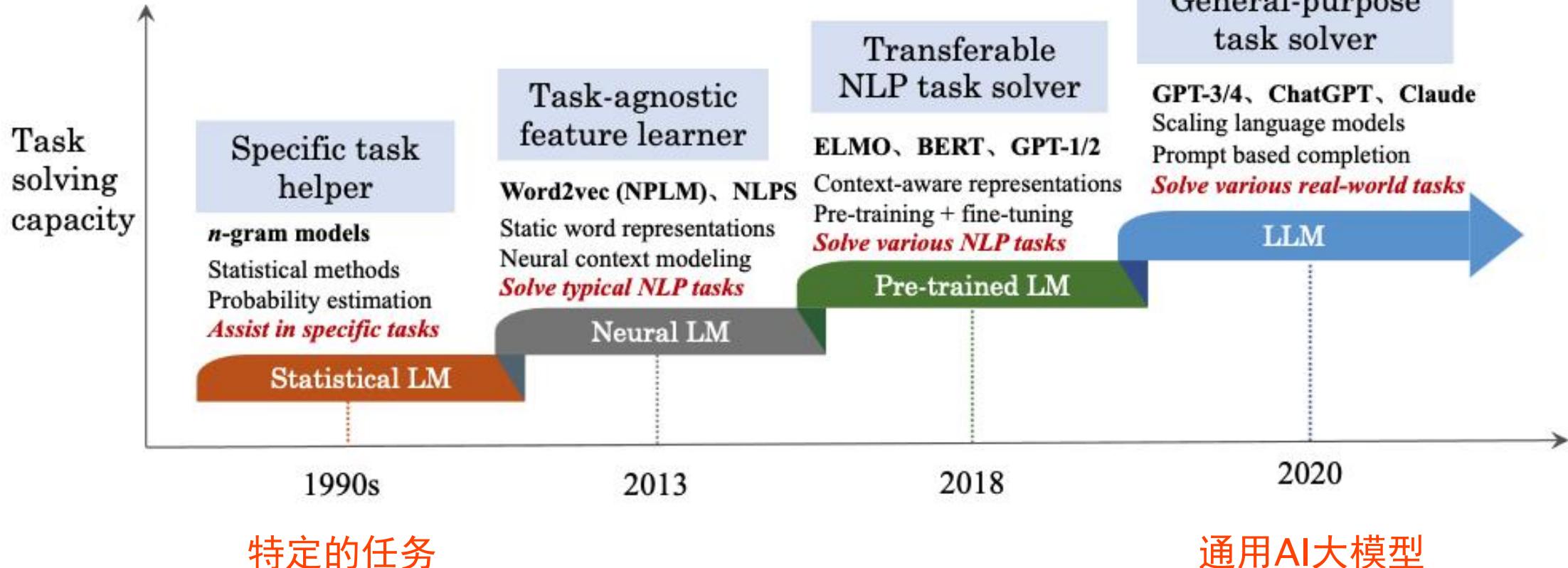
the Pile - StackExchange (41G, 2020)

BookCorpus (5G, 2015), Gutenberg (-, 2021), CC-Stories-R (31G, 2019), CC-NEWES (78G, 2019), REALNEWS (120G, 2019)

the Pile - ArXiv (72G, 2020), the Pile - PubMed Abstracts (25G, 2020)

BigQuery (-, 2023), the Pile - GitHub (61G, 2020)

AI模型能解决的任务越来越丰富



How to learn?

Supervised Learning

Given labels

e.g. classification task, regression task

Unsupervised Learning

unknown labels

e.g. clustering task, data reconstruction task

Self-supervised Learning

No labels

But the data itself can be the label: mask, next-word prediction

Tasks fall into two categories: pretext tasks, downstream tasks

The goal of AI: minimize the loss function

Goal: to solve an optimization problem

$$\theta^* = \operatorname{argmin}_{\theta} L(\theta)$$

L : loss function

θ : parameters

How to solve it?

Take linear regression as an example:

- 1) Analytical solution
- 2) Gradient descent

$$\beta^{(j+1)} \leftarrow \beta^{(j)} - \eta \frac{d\mathcal{L}(\beta)}{d\beta}$$

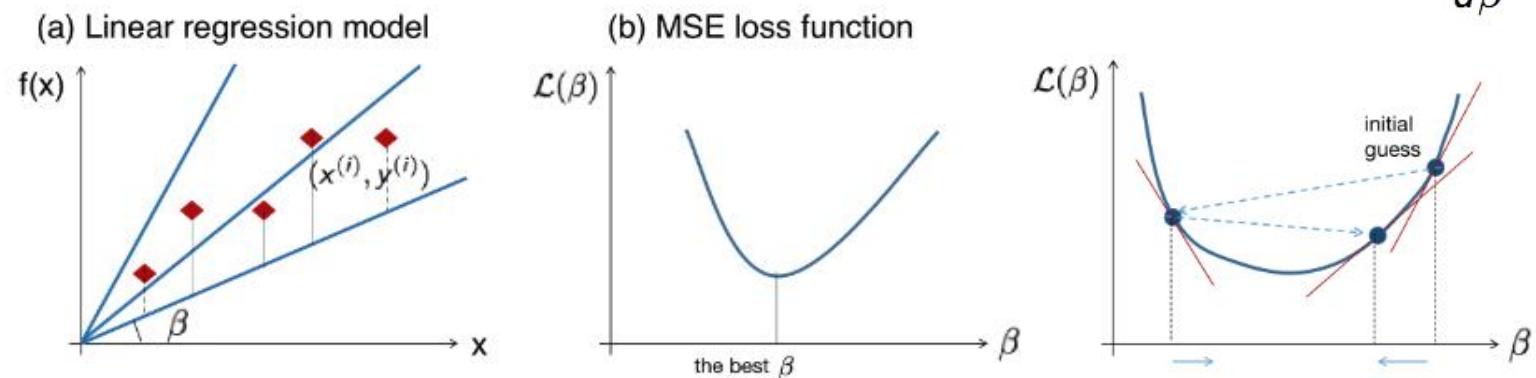
The simplest linear regression model: $y = \beta x$

Parameter: the slope β

Loss function:

$$\mathcal{L}(\beta) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \beta x^{(i)})^2$$

$$\beta^{(j+1)} \leftarrow \beta^{(j)} - \eta \frac{d\mathcal{L}(\beta)}{d\beta}$$



Gradient Descent (GD)

Goal: to solve an optimization problem

$$\theta^* = \operatorname{argmin}_{\theta} L(\theta)$$

L : loss function

θ : parameters

Suppose that θ has two variables $\{\theta_1, \theta_2\}$

Randomly start at $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$

$$\theta^1 = \begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta^0)/\partial \theta_1 \\ \partial L(\theta^0)/\partial \theta_2 \end{bmatrix} \rightarrow$$

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

$$\theta^2 = \begin{bmatrix} \theta_1^2 \\ \theta_2^2 \end{bmatrix} = \begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta^1)/\partial \theta_1 \\ \partial L(\theta^1)/\partial \theta_2 \end{bmatrix} \rightarrow$$

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

.....

until converge to θ^*

Gradient:

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta)/\partial \theta_1 \\ \partial L(\theta)/\partial \theta_2 \end{bmatrix}$$

Learning rate: η

An example to calculate the gradient of loss function

Loss function: $L(\theta) = \frac{1}{2}(\theta - y)^2$

$$L(\theta) = \frac{1}{2} \left([\theta_1 - y_1] \right)^2$$

Gradient:

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta)/\partial \theta_1 \\ \partial L(\theta)/\partial \theta_2 \end{bmatrix} = ?$$

$$\nabla L \left(\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \right) = \begin{bmatrix} \theta_1 - y_1 \\ \theta_2 - y_2 \end{bmatrix}$$

Gradient Descent to train Neural Networks

Network parameters

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

Initiate
Parameters

$$\theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow \dots$$

Compute $\nabla L(\theta^0)$ *Compute $\nabla L(\theta^1)$*

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

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta)/\partial w_1 \\ \partial L(\theta)/\partial w_2 \\ \vdots \\ \partial L(\theta)/\partial b_1 \\ \partial L(\theta)/\partial b_2 \\ \vdots \end{bmatrix}$$

ANNs have Millions of parameters

To compute the gradients efficiently in ANNs,
we use the error ***back-propagation***.

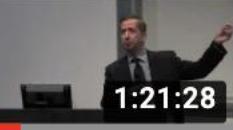
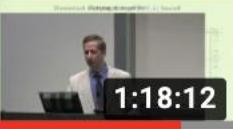
Some valuable materials for self-learning

highly recommended by 刘泉影

Caltech Machine Learning course

CS 156, by Prof. Yaser Abu-Mostafa

<https://www.youtube.com/watch?v=mbyG85GZ0PI&list=PLD63A284B7615313A>

 1:21:28	Lecture 01 - The Learning Problem caltech	 1:13:31	Lecture 07 - The VC Dimension caltech	 1:26:12	Lecture 13 - Validation caltech
 1:16:49	Lecture 02 - Is Learning Feasible? caltech	 1:16:51	Lecture 08 - Bias-Variance Tradeoff caltech	 1:14:16	Lecture 14 - Support Vector Machines caltech
 1:19:44	Lecture 03 -The Linear Model I caltech	 1:27:14	Lecture 09 - The Linear Model II caltech	 1:18:19	Lecture 15 - Kernel Methods caltech
 1:18:22	Lecture 04 - Error and Noise caltech	 1:25:16	Lecture 10 - Neural Networks caltech	 1:22:08	Lecture 16 - Radial Basis Functions caltech
 1:16:58	Lecture 05 - Training Versus Testing caltech	 1:19:49	Lecture 11 - Overfitting caltech	 1:16:18	Lecture 17 - Three Learning Principles caltech
 1:18:12	Lecture 06 - Theory of Generalization caltech	 1:15:14	Lecture 12 - Regularization caltech	 1:09:28	Lecture 18 - Epilogue caltech

Stanford Machine Learning course

CS 229, by Andrew Ng

<https://www.bilibili.com/video/BV19e411W7ga/>

III P1 1. Lecture 1 – Welcome

1:15:20

P2 2. Lecture 2 – Linear Regression and G... 1:18:17

P3 3. Lecture 3 – Locally Weighted & Logi... 1:19:35

P4 4. Lecture 4 – Perceptron & Generaliz... 1:22:02

P5 5. Lecture 5 – GDA & Naive Bayes 1:18:52

P6 6. Lecture 6 – Support Vector Machines 1:20:57

P7 7. Lecture 7 – Kernels 1:20:25

P8 8. Lecture 8 – Data Splits, Models & C... 1:23:26

P9 9. Lecture 9 – Approx_Estimation Erro... 1:26:03

P10 10. Lecture 10 – Decision Trees and E... 1:20:41

P11 11. Lecture 11 – Introduction to Neural ... 1:20:14

P12 12. Lecture 12 – Backprop & Improvin... 1:16:38

P13 13. Lecture 13 – Debugging ML Model... 1:18:55

P14 14. Lecture 14 – Expectation-Maximiz... 1:20:32

P15 15. Lecture 15 – EM Algorithm & Fact... 1:19:48

P16 16. Lecture 16 – Independent Compon... 1:18:10

P17 17. Lecture 17 – MDPs & Value_Policy I... 1:19:15

P18 18. Lecture 18 – Continous State MDP... 1:20:15

P19 19. Lecture 19 – Reward Model & Line... 1:21:07

P20 20. Lecture 20 – RL Debugging and ... 1:12:43

B站大学 【机器学习】白板推导

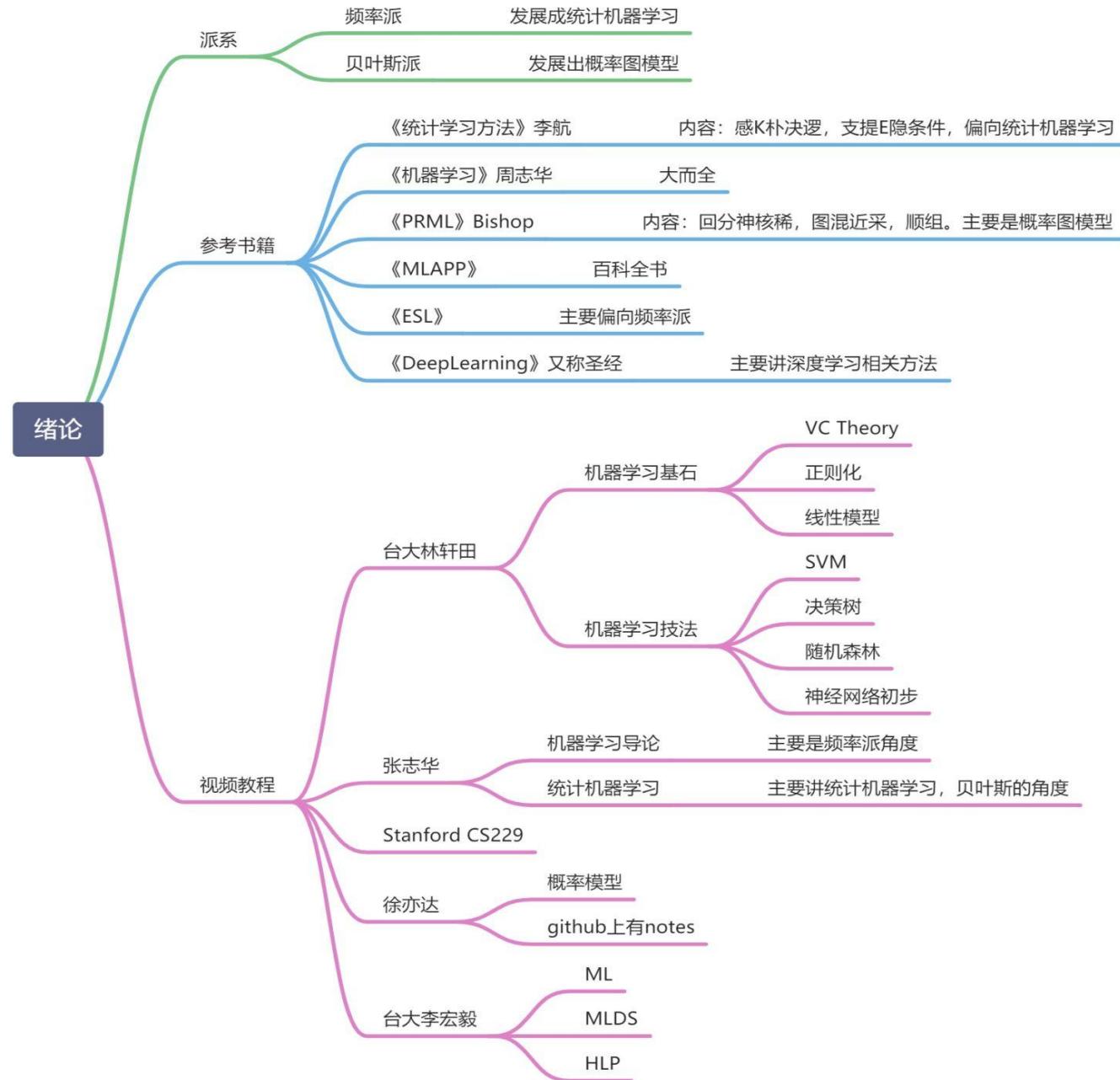
by shuhai008

<https://www.bilibili.com/video/BV1aE411o7qd>

纯理论推导



机器学习白板推导



李飞飞 Stanford CS231n

2021 年

-  [A brief history of computer vision](#)
-  [Deep Learning for Computer Vision Lecture 1 - Overview](#)
-  [Lecture 2: Image Classification with Linear Classifiers](#)
-  [Lecture 3: Regularization and Optimization](#)
-  [Lecture 4: Neural Networks and Backpropagation](#)
-  [Lecture 5: Image Classification with CNNs](#)
-  [Lecture 6: CNN Architectures](#)
-  [Lecture 7: Training Neural Networks](#)
-  [Lecture 8: Visualizing and Understanding](#)
-  [Lecture 9: Object Detection and Image Segmentation](#)
-  [Lecture 10: Recurrent Neural Networks](#)
-  [Lecture 11: Attention and Transformers](#)
-  [Lecture 12: Video Understanding](#)
-  [Lecture 13: Generative Models](#)
-  [Lecture 14: Self-supervised Learning](#)
-  [CS231N: Low-Level Vision](#)
-  [Lecture 16: 3D Vision](#)

最新的课程没用找到录屏

B站上只有2017年

<https://www.bilibili.com/video/BV1nJ411z7fe/>

2019年的课件

<https://www.bilibili.com/video/BV1hi4y1t7kF/>

李飞飞 – Stanford CS231n (Deep Learning for Computer Vision)

2022年 <http://cs231n.stanford.edu/2022/schedule.html>

--- Deep Learning Basics

03/31 Lecture 2: Image Classification with Linear Classifiers
The data-driven approach
K-nearest neighbor
Linear Classifiers

Image Classification Problem
Linear Classification

Algebraic / Visual / Geometric viewpoints
SVM and Softmax loss
[slides]

04/01 Python / Numpy Review Session
[Colab] [Tutorial]
⌚ 1:30-2:30pm PT

04/05 Lecture 3: Regularization and Optimization
Regularization
Stochastic Gradient Descent
Momentum, AdaGrad, Adam
Learning rate schedules
[slides]

04/07 Lecture 4: Neural Networks and Backpropagation
Multi-layer Perceptron
Backpropagation
[slides]

Backprop
Linear backprop example
Suggested Readings:
1. Why Momentum Really Works

--- Perceiving and Understanding the Visual World

04/12 Lecture 5: Image Classification with CNNs
History
Higher-level representations, image features
Convolution and pooling
[slides]

Convolutional Networks

04/14 Lecture 6: CNN Architectures
Batch Normalization
Transfer learning
AlexNet, VGG, GoogLeNet, ResNet
[slides]

AlexNet, VGGNet, GoogLeNet, ResNet

04/19 Lecture 7: Training Neural Networks
Activation functions
Data processing
Weight initialization

Neural Networks, Parts 1, 2, 3
Suggested Readings:
1. Stochastic Gradient Descent Tricks
2. Efficient Backprop

2023年

ht --- Deep Learning Basics

04/06 Lecture 2: Image Classification with Linear Classifiers
The data-driven approach
K-nearest neighbor
Linear Classifiers
Algebraic / Visual / Geometric viewpoints
SVM and Softmax loss
[slides]

04/07 Python / Numpy Review Session
[Colab] [Tutorial]
⌚ 1:30-2:20pm PT

04/11 Lecture 3: Regularization and Optimization
Regularization
Stochastic Gradient Descent
Momentum, AdaGrad, Adam
Learning rate schedules
[slides]

04/13 Lecture 4: Neural Networks and Backpropagation
Multi-layer Perceptron
Backpropagation
[slides]

--- Perceiving and Understanding the Visual World

04/18 Lecture 5: Image Classification with CNNs
History
Higher-level representations, image features
Convolution and pooling
[slides]

04/20 Lecture 6: CNN Architectures
Batch Normalization
Transfer learning
AlexNet, VGG, GoogLeNet, ResNet
[slides]

04/25 Lecture 7: Training Neural Networks
Activation functions
Data processing
Weight initialization

Neural Networks, Parts 1, 2, 3
Suggested Readings:
1. Stochastic Gradient Descent Tricks
2. Efficient Backprop

李飞飞 – Stanford CS231n (Deep Learning for Computer Vision)

2022年 <http://cs231n.stanford.edu/2022/schedule.html>

04/21 Lecture 8: Visualizing and Understanding
Feature visualization and inversion
Adversarial examples
DeepDream and style transfer
[\[slides\]](#)

04/26 Lecture 9: Object Detection and Image Segmentation
Single-stage detectors
Two-stage detectors
Semantic/Instance/Panoptic segmentation
[\[slides\]](#)

04/28 Lecture 10: Recurrent Neural Networks
RNN, LSTM, GRU
Language modeling
Image captioning
Sequence-to-sequence
[\[slides\]](#)

05/03 Lecture 11: Attention and Transformers
Self-Attention
Transformers
[\[slides\]](#)

05/05 Lecture 12: Video Understanding
Video classification
3D CNNs
Two-stream networks
Multimodal video understanding
[\[slides\]](#)

FCN, R-CNN, Fast R-CNN, Faster R-CNN, YOLO

Suggested Readings:
1. [DL book RNN chapter](#)
2. [Understanding LSTM Networks](#)

Suggested Readings:
1. Attention is All You Need [[Original Transformers Paper](#)]
2. Attention? Attention [[Blog by Lilian Weng](#)]
3. The Illustrated Transformer [[Blog by Jay Alammar](#)]
4. ViT: Transformers for Image Recognition [[Paper](#)] [[Blog](#)] [[Video](#)]
5. DETR: End-to-End Object Detection with Transformers [[Paper](#)] [[Blog](#)] [[Video](#)]

2023年

<http://cs231n.stanford.edu/2023/schedule.html>

04/27 Lecture 8: Recurrent Neural Networks
RNN, LSTM, GRU
Language modeling
Image captioning
Sequence-to-sequence
[\[slides\]](#)

05/02 Lecture 9: Attention and Transformers
Self-Attention
Transformers
[\[slides\]](#)

05/04 Lecture 10: Video Understanding
Video classification
3D CNNs
Two-stream networks
Multimodal video understanding
[\[slides\]](#)

05/09 Lecture 11: Object Detection and Image Segmentation
Single-stage detectors
Two-stage detectors
Semantic/Instance/Panoptic segmentation
[\[slides\]](#)

Suggested Readings:
1. [Attention is All You Need \[Original Transformers Paper\]](#)
2. Attention? Attention [[Blog by Lilian Weng](#)]
3. The Illustrated Transformer [[Blog by Jay Alammar](#)]
4. ViT: Transformers for Image Recognition [[Paper](#)] [[Blog](#)] [[Video](#)]
5. DETR: End-to-End Object Detection with Transformers [[Paper](#)] [[Blog](#)] [[Video](#)]

FCN, R-CNN, Fast R-CNN, Faster R-CNN, YOLO

05/11 Lecture 12: Visualizing and Understanding
Feature visualization and inversion
Adversarial examples
DeepDream and style transfer
[\[slides\]](#)

李飞飞 – Stanford CS231n (Deep Learning for Computer Vision)

2022年 <http://cs231n.stanford.edu/2022/schedule.html>

--- Reconstructing and Interacting with the Visual World	
05/12	Lecture 13: Generative Models Supervised vs. Unsupervised learning Pixel RNN, Pixel CNN Variational Autoencoders Generative Adversarial Networks [slides]
05/17	Lecture 14: Self-supervised Learning Pretext tasks Contrastive learning Multisensory supervision [slides]
05/19	Lecture 15: Low-Level Vision (Guest Lecture by Prof. Jia Deng from Princeton University) Optical flow Depth estimation Stereo vision [slides]
05/24	

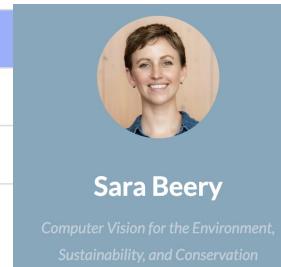
2023年	
http://cs231n.stanford.edu/schedule.html	
05/16	In-Class Midterm ⌚ 12:00-1:20pm
05/18	Lecture 13: Self-supervised Learning Pretext tasks Contrastive learning Multisensory supervision [slides] Suggested Readings: 1. Lilian Weng Blog Post 2. DINO: Emerging Properties in Self-Supervised Vision Transformers [Paper] [Blog] [Video]
05/23	Lecture 14: Robot Learning Deep Reinforcement Learning Model Learning Robotic Manipulation [slides]
05/25	Lecture 15: Generative Models (Guest Lecture by Dr. Ruiqi Gao from Google DeepMind) Generative Adversarial Network Diffusion models Autoregressive models [slides]
05/30	Lecture 16: 3D Vision 3D shape representations Shape reconstruction Neural implicit representations [slides]

李飞飞 – Stanford CS231n (Deep Learning for Computer Vision)

2022年 <http://cs231n.stanford.edu/2022/schedule.html>

--- Reconstructing and Interacting with the Visual World		
05/12	Lecture 13: Generative Models Supervised vs. Unsupervised learning Pixel RNN, Pixel CNN Variational Autoencoders Generative Adversarial Networks [slides]	Suggested Readings: 1. Image GPT: Generative Pretraining Pixels [Paper] [Blog]
05/17	Lecture 14: Self-supervised Learning Pretext tasks Contrastive learning Multisensory supervision [slides]	Suggested Readings: 1. Lilian Weng Blog Post 2. DINO: Emerging Properties in Self-Supervised Vision Transformers [Paper] [Blog] [Video]
05/19	Lecture 15: Low-Level Vision (Guest Lecture by Prof. Jia Deng from Princeton University) Optical flow Depth estimation Stereo vision [slides]	Optical flow Depth estimation Stereo vision [slides]
05/24	Lecture 16: 3D Vision 3D shape representations Shape reconstruction Neural implicit representations [slides]	
--- Human-Centered Applications and Implications		
05/26	Lecture 17: Human-Centered Artificial Intelligence AI & healthcare	
05/31	Lecture 18: Fairness in Visual Recognition (Guest Lecture by Prof. Olga Russakovsky from Princeton University)	

2023年
<http://cs231n.stanford.edu/schedule.html>

--- Generative and Interactive Visual Intelligence		
05/16	In-Class Midterm	⌚ 12:00-1:20pm
05/18	Lecture 13: Self-supervised Learning Pretext tasks Contrastive learning Multisensory supervision [slides]	Suggested Readings: 1. Lilian Weng Blog Post 2. DINO: Emerging Properties in Self-Supervised Vision Transformers [Paper] [Blog] [Video]
05/23	Lecture 14: Robot Learning Deep Reinforcement Learning Model Learning Robotic Manipulation [slides]	 
05/25	Lecture 15: Generative Models (Guest Lecture by Dr. Ruiqi Gao from Google DeepMind) Generative Adversarial Network Diffusion models Autoregressive models [slides]	
05/30	Lecture 16: 3D Vision 3D shape representations Shape reconstruction Neural implicit representations [slides]	
--- Human-Centered Applications and Implications		
06/01	Lecture 17: Human-Centered Artificial Intelligence	
06/06	Lecture 18: Guest Lecture by Prof. Sara Beery from MIT	<p>Sara Beery Computer Vision for the Environment, Sustainability, and Conservation</p>

机器人+强化学习+大模型
通往具身智能

Christ Manning – Stanford CS224n (Natural Language Processing with Deep Learning)

2023年 <https://web.stanford.edu/class/cs224n/index.html#schedule>

B站 2021年课程的视频合集 <https://www.bilibili.com/video/BV18Y411p79k>

Tue Jan 10 Word Vectors (*by John Hewitt*)
Week 1 [slides] [notes]

Gensim word vectors example:
[code] [preview]

Suggested Readings:

1. [Efficient Estimation of Word Representations in Vector Space](#) (original word2vec paper)
2. [Distributed Representations of Words and Phrases and their Compositionality](#) (negative sampling paper)

Thu Jan 12 Word Vectors, Word Window Classification, Language Models
[slides] [notes]

Suggested Readings:

1. [GloVe: Global Vectors for Word Representation](#) (original GloVe paper)
2. [Improving Distributional Similarity with Lessons Learned from Word Embeddings](#)
3. [Evaluation methods for unsupervised word embeddings](#)

Additional Readings:

1. [A Latent Variable Model Approach to PMI-based Word Embeddings](#)
2. [Linear Algebraic Structure of Word Senses, with Applications to Polysemy](#)
3. [On the Dimensionality of Word Embedding](#)

Christ Manning – Stanford CS224n (Natural Language Processing with Deep Learning)

2023年 <https://web.stanford.edu/class/cs224n/index.html#schedule>

Tue Jan 17 Backprop and Neural
Networks
Week 2 [slides] [notes]

Suggested Readings:

1. [matrix calculus notes](#)
2. [Review of differential calculus](#)
3. [CS231n notes on network architectures](#)
4. [CS231n notes on backprop](#)
5. [Derivatives, Backpropagation, and Vectorization](#)
6. [Learning Representations by Backpropagating Errors](#)
(seminal Rumelhart et al. backpropagation paper)

Additional Readings:

1. [Yes you should understand backprop](#)
2. [Natural Language Processing \(Almost\) from Scratch](#)

Thu Jan 19 Dependency Parsing
[slides] [notes]
[slides (annotated)]

Suggested Readings:

1. [Incrementality in Deterministic Dependency Parsing](#)
2. [A Fast and Accurate Dependency Parser using Neural Networks](#)
3. [Dependency Parsing](#)
4. [Globally Normalized Transition-Based Neural Networks](#)
5. [Universal Stanford Dependencies: A cross-linguistic typology](#)
6. [Universal Dependencies website](#)
7. [Jurafsky & Martin Chapter 14](#)

Christ Manning – Stanford CS224n (Natural Language Processing with Deep Learning)

Tue Jan 24 Week 3	Recurrent Neural Networks and Language Models [slides] [notes (lectures 5 and 6)]	Suggested Readings: <ol style="list-style-type: none">1. N-gram Language Models (textbook chapter)2. The Unreasonable Effectiveness of Recurrent Neural Networks (blog post overview)3. Sequence Modeling: Recurrent and Recursive Neural Nets (Sections 10.1 and 10.2)4. On Chomsky and the Two Cultures of Statistical Learning5. Sequence Modeling: Recurrent and Recursive Neural Nets (Sections 10.3, 10.5, 10.7-10.12)6. Learning long-term dependencies with gradient descent is difficult (one of the original vanishing gradient papers)7. On the difficulty of training Recurrent Neural Networks (proof of vanishing gradient problem)8. Vanishing Gradients Jupyter Notebook (demo for feedforward networks)9. Understanding LSTM Networks (blog post overview)
	Tue Jan 31 Week 4	Final Projects: Custom and Default; Practical Tips [slides] [notes]
Thu Feb 2 Week 5	Self-Attention and Transformers (<i>by John Hewitt</i>) [slides] [notes]	Suggested Readings: <ol style="list-style-type: none">1. Default Project Handout2. Attention Is All You Need3. The Illustrated Transformer4. Transformer (Google AI blog post)5. Layer Normalization6. Image Transformer7. Music Transformer: Generating music with long-term structure
	Tue Feb 7	Pretraining (<i>by John Hewitt</i>) [slides]
Thu Feb 9	Natural Language Generation (<i>by Xiang Lisa Li</i>) [slides]	Suggested Readings: <ol style="list-style-type: none">1. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding2. Contextual Word Representations: A Contextual Introduction3. The Illustrated BERT, ELMo, and co.4. Martin & Jurafsky Chapter on Transfer Learning
		Suggested Readings: <ol style="list-style-type: none">1. The Curious Case of Neural Text Degeneration2. Get To The Point: Summarization with Pointer-Generator Networks3. Hierarchical Neural Story Generation4. How NOT To Evaluate Your Dialogue System

Christ Manning – Stanford CS224n (Natural Language Processing with Deep Learning)

Tue Feb 14	Prompting, Reinforcement Learning from Human Feedback (by Jesse Mu) [slides]	Suggested Readings: 1. Language Models are Few-Shot Learners 2. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models 3. Finetuned Language Models Are Zero-Shot Learners 4. Learning to summarize from human feedback
Thu Feb 16	Question Answering [slides]	Suggested readings: 1. SQuAD: 100,000+ Questions for Machine Comprehension of Text 2. Bidirectional Attention Flow for Machine Comprehension 3. Reading Wikipedia to Answer Open-Domain Questions 4. Latent Retrieval for Weakly Supervised Open Domain Question Answering 5. Dense Passage Retrieval for Open-Domain Question Answering 6. Learning Dense Representations of Phrases at Scale
Tue Feb 21	ConvNets, Tree Recursive Neural Networks and Constituency Parsing [slides]	Suggested readings: 1. Convolutional Neural Networks for Sentence Classification 2. Improving neural networks by preventing co-adaptation of feature detectors 3. A Convolutional Neural Network for Modelling Sentences 4. Parsing with Compositional Vector Grammars. 5. Constituency Parsing with a Self-Attentive Encoder
Thu Feb 23	Insights between NLP and Linguistics (by Isabel Papadimitriou) [slides]	

Tue Feb 28	Code Generation (by Gabriel Poesia) [slides]	Suggested readings: 1. Program Synthesis with Large Language Models 2. Competition-level code generation with AlphaCode 3. Evaluating Large Language Models Trained on Code
Wed Mar 1	Training Large Language Models (by John Hewitt)	⌚ 3:30pm - 4:20pm Skilling Auditorium
Thu Mar 2	Multimodal Deep Learning (by Douwe Kiela) [slides]	
Tue Mar 7	Coreference Resolution [slides]	Suggested readings: 1. Coreference Resolution Chapter from Jurafsky and Martin 2. End-to-end Neural Coreference Resolution
Thu Mar 9	Analysis and Interpretability Basics (by John Hewitt) [slides]	
Fri Mar 10	Latex Tutorial (by Rishi Desai)	⌚ 3:30pm - 4:20pm Skilling Auditorium
Tue Mar 14	Model Interpretability and Editing (by Been Kim) [slides]	
Thu Mar 16	Final Project Emergency Assistance (no lecture)	Extra project office hours available during usual lecture time, see Ed.
Sat Mar 18		
Monday Mar 20	Poster Session	⌚ 5pm-9pm [More details] Location: Tressider Oak Lounge

B站大学 跟李沐学AI

<https://space.bilibili.com/1567748478>

【完结】动手学深度学习 PyTorch版 76

The thumbnail shows a grid of video thumbnails for the completed series. It includes a summary video, optimization algorithms, object detection competition results, and fine-tuning BERT. Each thumbnail has a play button, duration, and date.

视频标题	时长	发布日期
【完结】73 – 课程总结和进阶学习【动手学深度学习v2】	27:55	2021-8-25
72 优化算法【动手学深度学习v2】	43:17	2021-8-24
71 目标检测竞赛总结【动手学深度学习v2】	23:58	2021-8-22
70 BERT微调【动手学深度学习v2】	31:26	2021-8-19

合集 · 【更新中】AI 论文精读 53

如何读论文

如何读论文【论文精读·1】

35.6万 2021-10-6



9年后重读深度学习奠基之作之一：AlexNet【论文精读·2】

23.8万 2021-10-14



AlexNet论文逐段精读【论文精读】

22.8万 2021-10-15



撑起计算机视觉半边天的ResNet【论文精读】

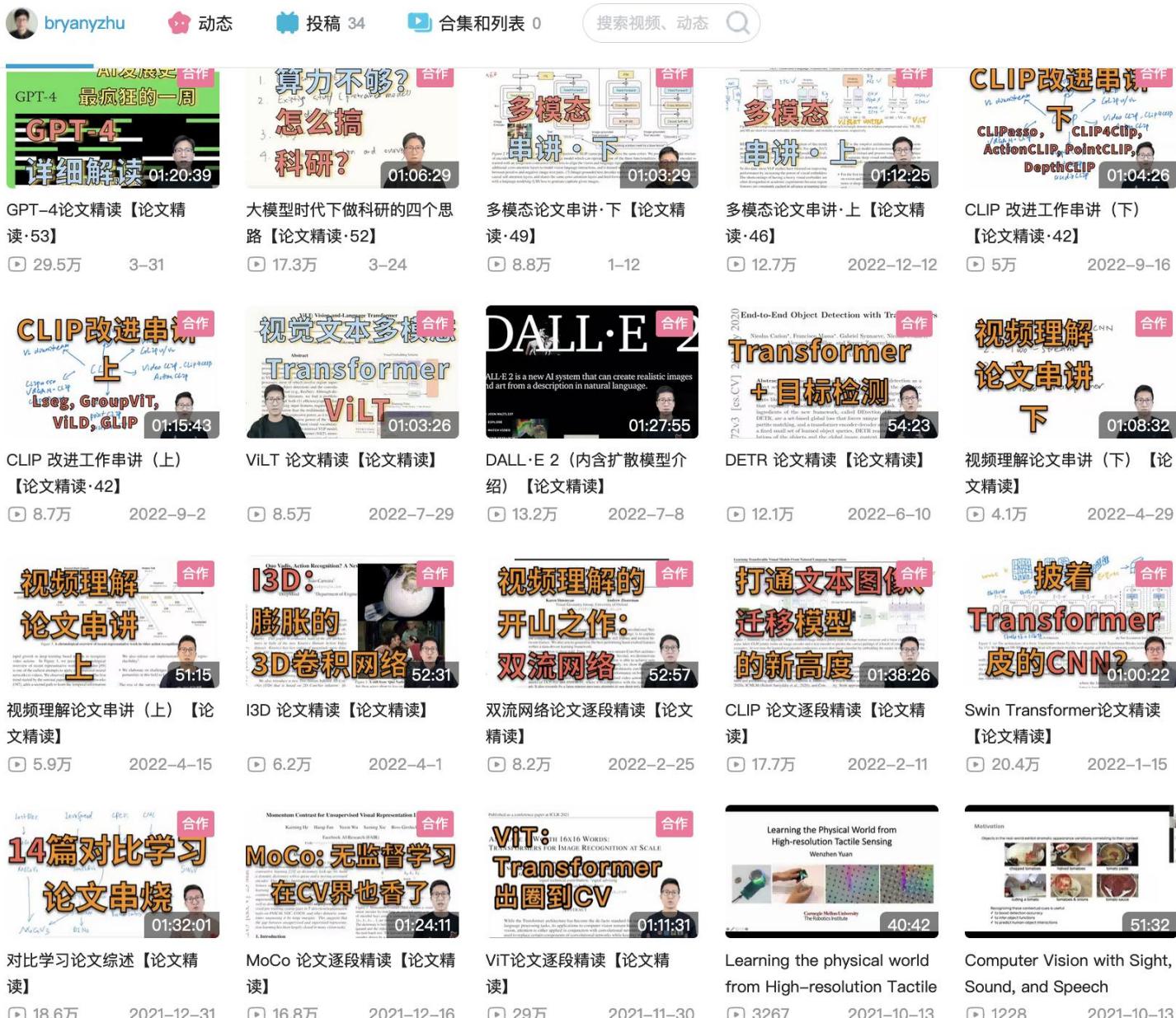
17.6万 2021-10-21

B站大学 bryanyzhu

<https://space.bilibili.com/1511278611>

朱老师的视频以CV为主

后来转到多模态模型



B站大学 清华NLP刘知远团队大模型公开课

<https://www.bilibili.com/video/BV1UG411p7zv>

		视频选集 (130/137) 	自动连播 
P31	3-11 Transformer结构--Transformer优...	P42 4-2 Prompt-Learning和Delta-Tuning--背景与介绍	05:5
P32	3-12 预训练语言模型--语言建模概述	P43 4-3 Prompt-Learning--基本组成与应用	07:2
P33	3-13 预训练语言模型--PLM介绍	P44 4-4 Prompt-Learning--PTM选取	
P34	3-14 预训练语言模型--MLM任务的应用	P45 4-5 Prompt-Learning--Template构造	
P35	3-15 预训练语言模型--前沿大模型介绍	P46 4-6 Prompt-Learning--Verbalizer构建	
P36	3-16 Transformers教程--Introduction	P47 4-7 Prompt-Learning--训练新范式	
P37	3-17 Transformers教程--使用Transformers	P48 4-8 Prompt-Learning--应用	
P38	3-18 Transformers教程--Tokenization	P49 4-9 Prompt-Learning--总结	07:5
P39	3-19 Transformers教程--常用API介绍	P50 4-10 Delta-Tuning--背景与介绍	14:1
P40	3-20 Transformers教程--Demo讲解	P51 4-11 Delta-Tuning--增量式tuning	01:1
		P128 9-6 大模型中的神经元--背景介绍 P129 9-7 大模型中的神经元--激活情况分析 P130 9-8 大模型中的神经元--转换模型架构 P131 9-9 大模型神经元的应用--学到的特定... P132 9-10 大模型神经元的应用--作为迁移指标 P133 9-11 大模型神经元的应用--表示情感 P134 9-12 大模型认知能力--介绍 P135 9-13 大模型认知能力--下游任务实例 P136 9-14 大模型认知能力--挑战与限制	12:1 24:1 16:3 13:3 07:5 14:1 01:1

AI Hands-on (pytorch)

By 楼可心、王淞、钱富元

Error Back-Propagation

Recall Calculus: Chain Rule

Case 1

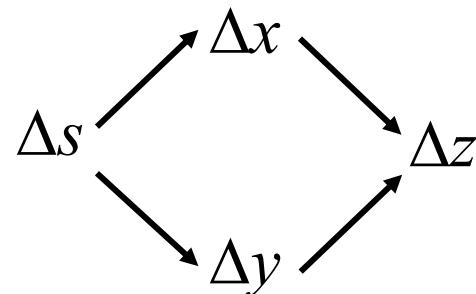
$$y = g(x) \quad z = h(y)$$

$$\Delta x \rightarrow \Delta y \rightarrow \Delta z$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

Case 2

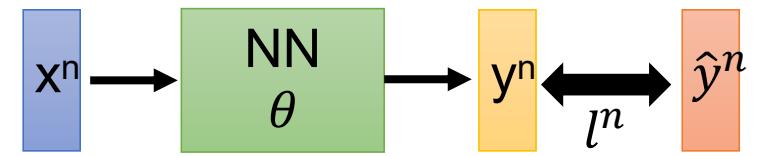
$$x = g(s) \quad y = h(s) \quad z = k(x, y)$$



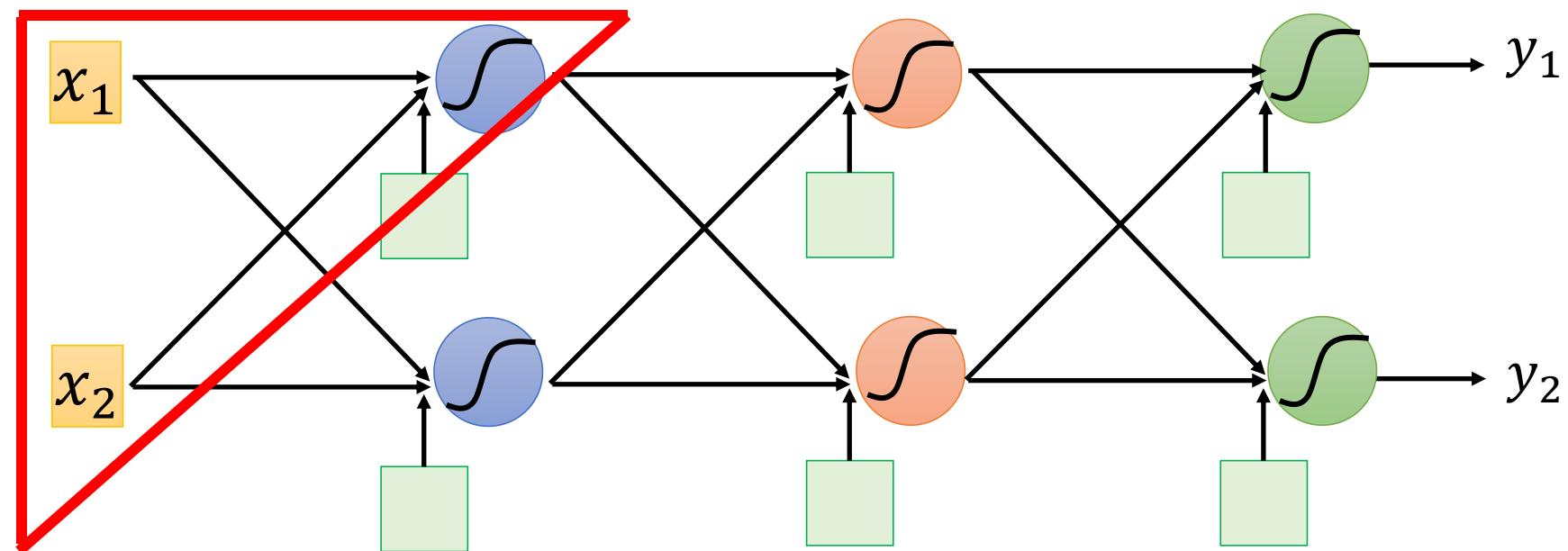
$$\frac{dz}{ds} = \frac{\partial z}{\partial x} \frac{dx}{ds} + \frac{\partial z}{\partial y} \frac{dy}{ds}$$

Error Backpropagation

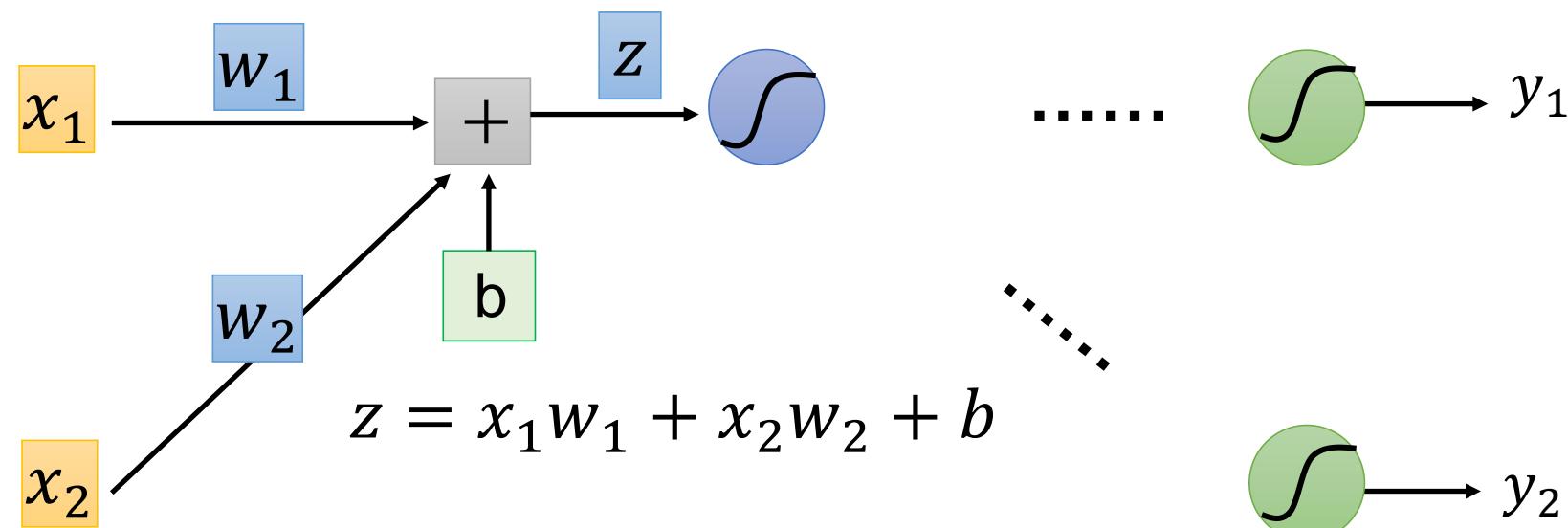
Backpropagation computes how slightly changing each synapse strength (**weight**) would change the network's **loss**, using the chain rule.



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



BP in fully-connected neural networks (fcNN)



Forward pass:

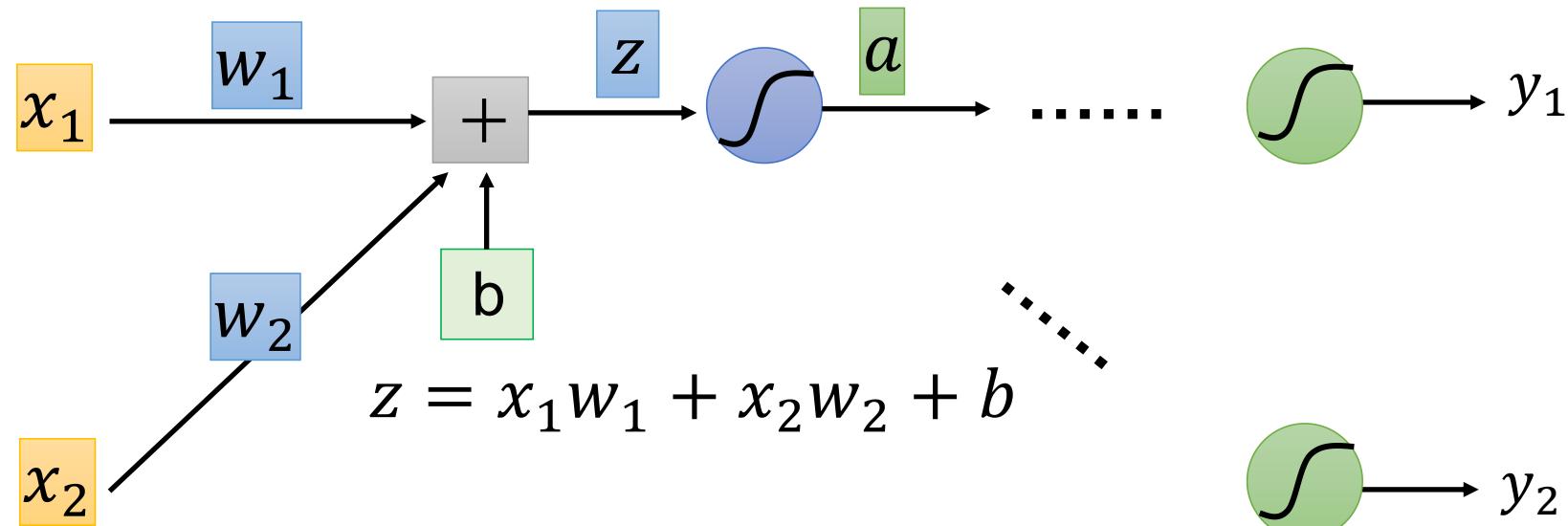
$\frac{\partial l}{\partial w} = ? \quad \frac{\partial l}{\partial z} \frac{\partial z}{\partial w}$ Compute $\frac{\partial z}{\partial w}$ for all parameters

Backward pass:

(Chain rule) Compute $\frac{\partial l}{\partial z}$ for all activation function inputs z

Backpropagation – Forward pass

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

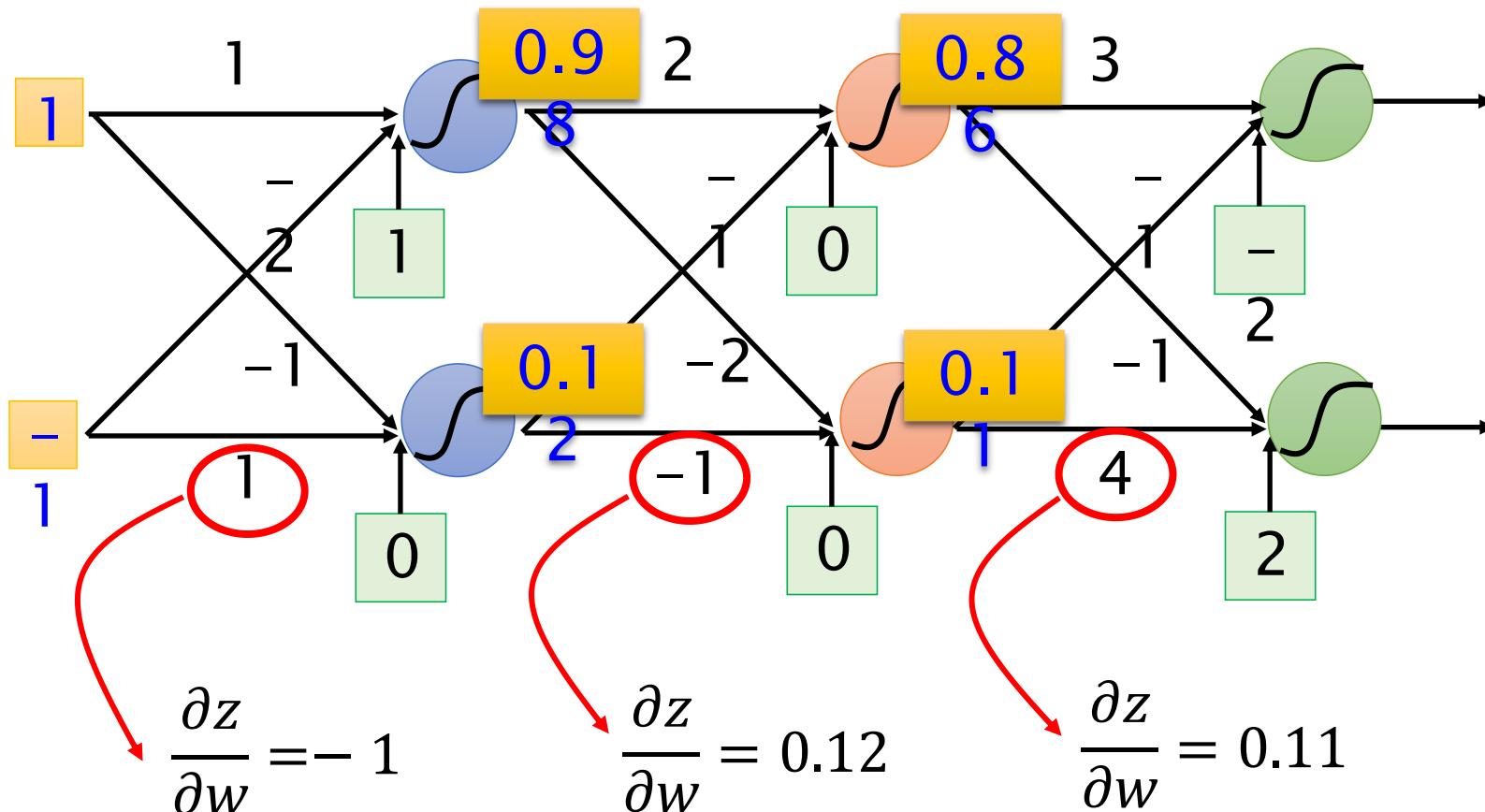


$$\begin{aligned}\partial z / \partial w_1 &= ? x_1 \\ \partial z / \partial w_2 &= ? x_2\end{aligned}$$

The value of the **input** connected by the weight

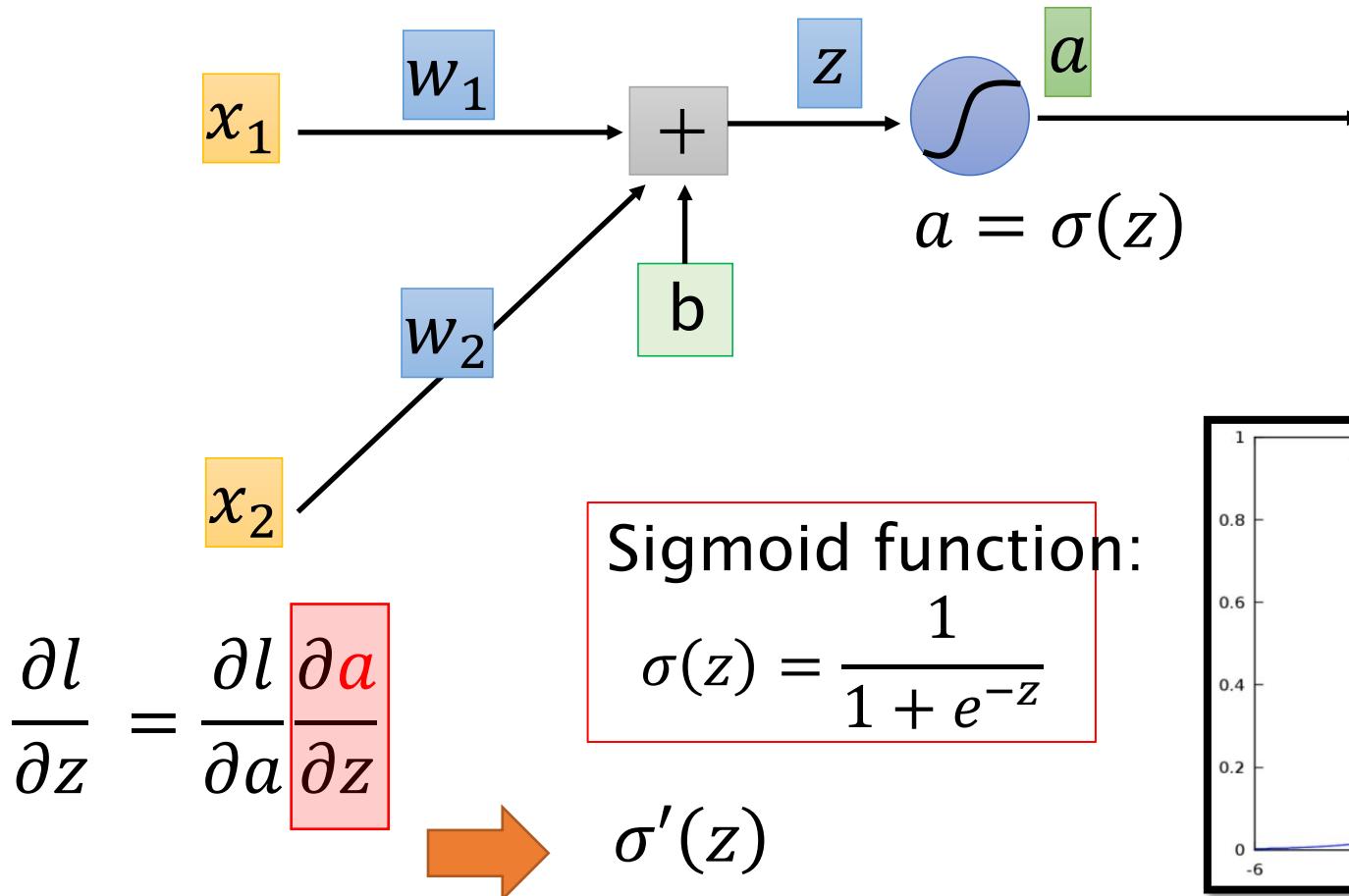
Backpropagation – Forward pass

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



Backpropagation – Backward pass

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

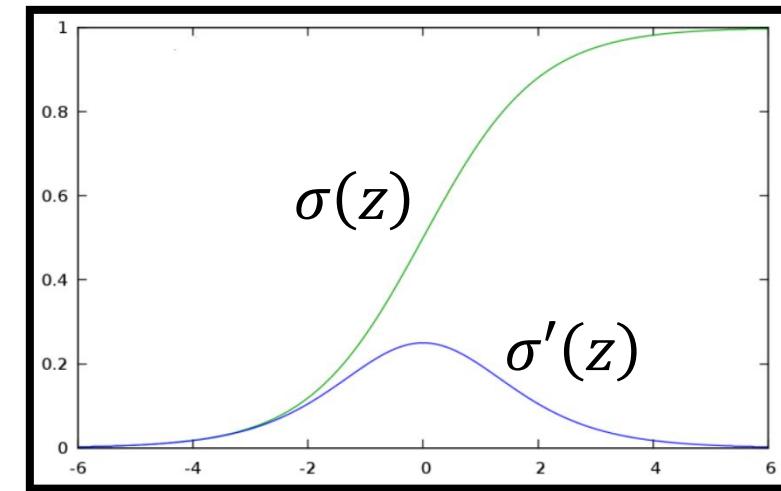
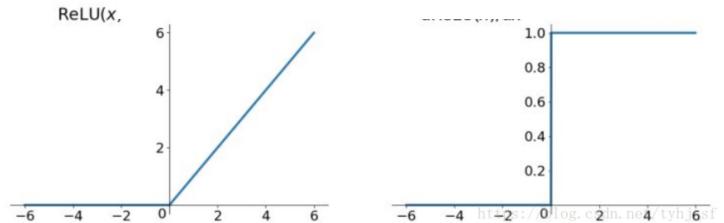


Sigmoid function:

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

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\text{Relu} = \max(0, x)$$



Backpropagation – Backward pass

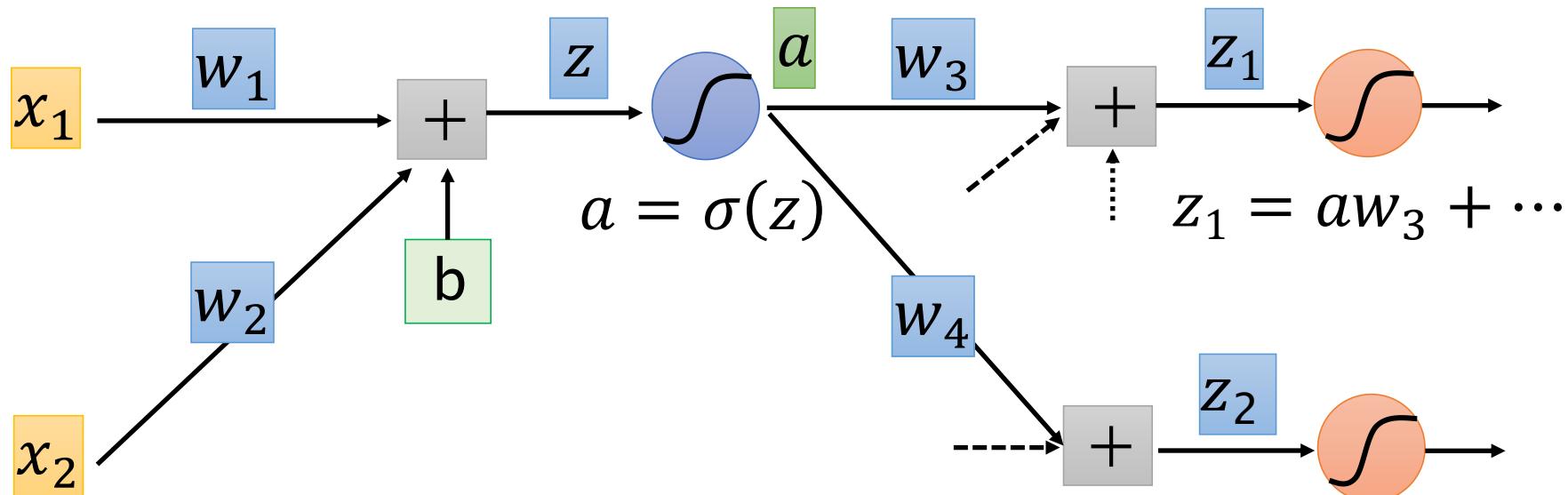
Activation function and its derivative

- ReLU: $h(a) = \max\{0, a\}$ $h'(a) = \begin{cases} 0, & \text{for } a < 0; \\ 1, & \text{for } a > 0 \end{cases}$
- Leaky ReLU: $h(a) = \begin{cases} \alpha a, & \text{for } a \leq 0 \\ a, & \text{for } a > 0 \end{cases}$
 $h'(a) = \begin{cases} \alpha, & \text{for } a \leq 0 \\ 1, & \text{for } a > 0 \end{cases}$
- Sigmoid function: $h(a) = \frac{1}{1+e^{-a}}$ $h'(a) = h(a)(1 - h(a))$
- Softmax function: $h(\mathbf{a}) = \frac{e^{\mathbf{a}_k}}{\sum_{k=1}^K e^{\mathbf{a}_k}}$
 $\frac{\partial h(\mathbf{a}_i)}{\partial \mathbf{a}_j} = h(\mathbf{a}_i)(\delta_{ij} - h(\mathbf{a}_j)),$ where δ_{ij} is called Kronecker δ
- TanH: $h(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$ $h'(a) = 1 - h^2(a)$

Backpropagation – Backward pass

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

z



$$\frac{\partial l}{\partial z} = \frac{\partial l}{\partial a} \frac{\partial a}{\partial z}$$

$\sigma'(z)$

$$\frac{\partial l}{\partial a} = \frac{\partial l}{\partial z_1} \frac{\partial z_1}{\partial a} + \frac{\partial l}{\partial z_2} \frac{\partial z_2}{\partial a}$$

(Chain rule)

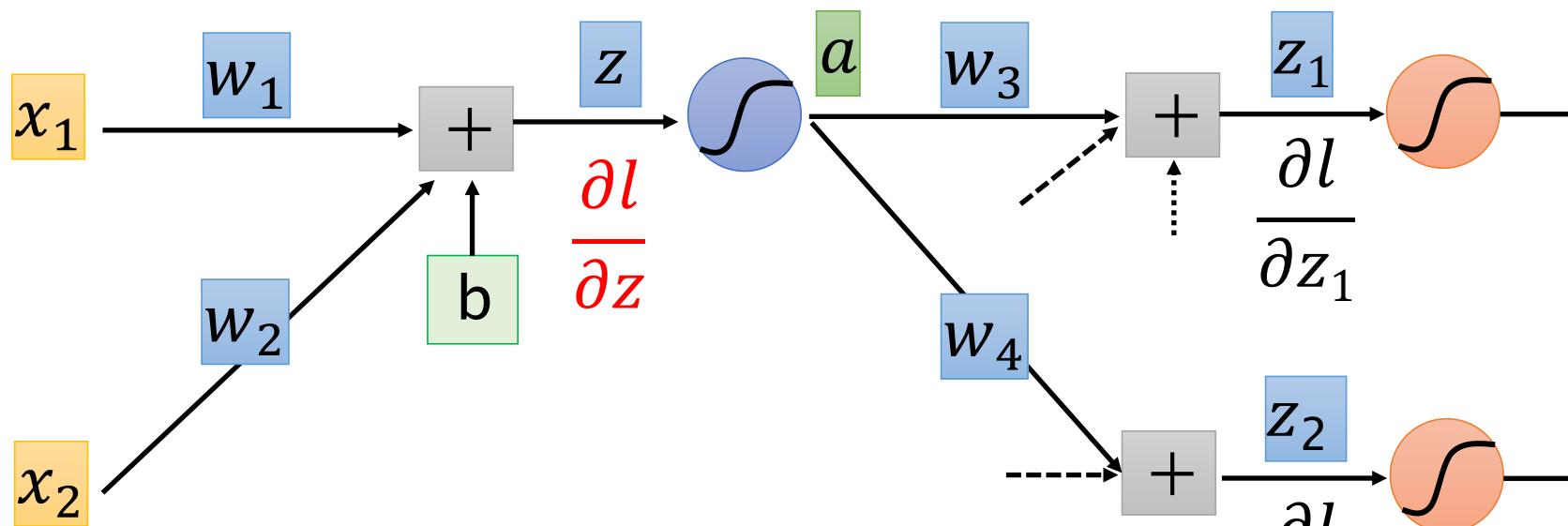
? w_3 ? w_4

Assumed it's known

Backpropagation – Backward pass

Compute $\frac{\partial l}{\partial z}$ for all activation function inputs

z

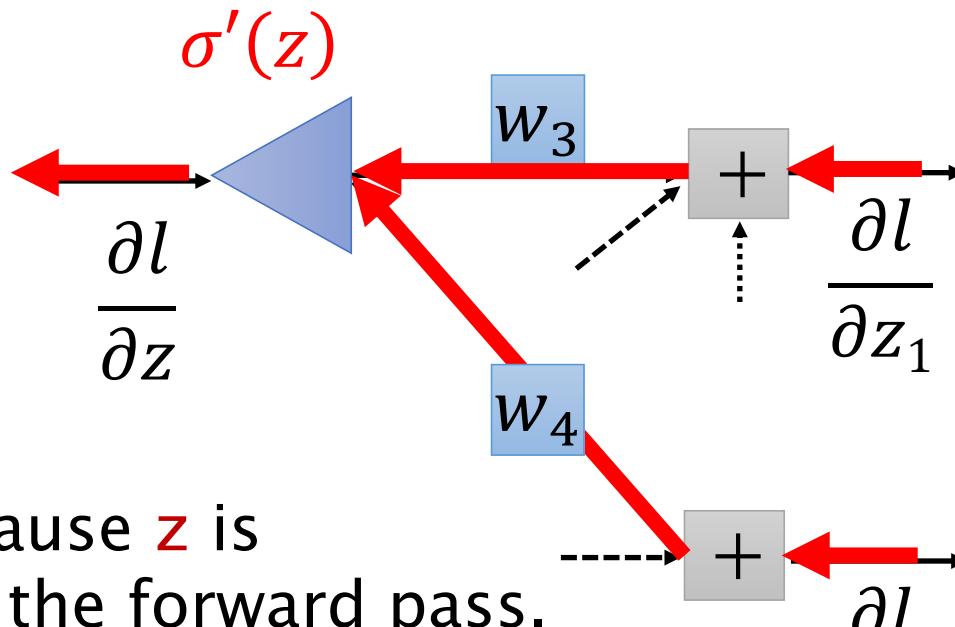


$$\frac{\partial l}{\partial z} = \sigma'(z) \left[w_3 \frac{\partial l}{\partial z_1} + w_4 \frac{\partial l}{\partial z_2} \right]$$

Backpropagation – Backward pass

Compute $\frac{\partial l}{\partial z}$ for all activation function inputs

z



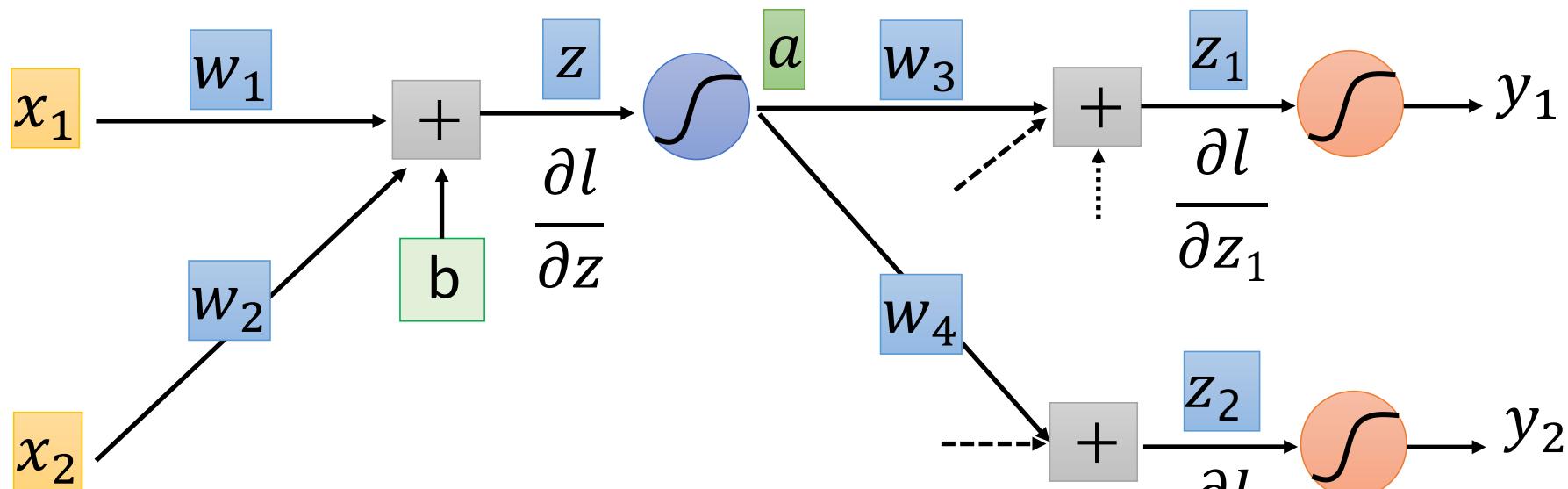
$\sigma'(z)$ is a constant because z is already determined in the forward pass.

$$\frac{\partial l}{\partial z} = \sigma'(z) \left[w_3 \frac{\partial l}{\partial z_1} + w_4 \frac{\partial l}{\partial z_2} \right]$$

Backpropagation – Backward pass

Compute $\frac{\partial l}{\partial z}$ for all activation function inputs

z



Case 1. Output

Layer

$$\frac{\partial l}{\partial z_1} = \frac{\partial l}{\partial y_1} \frac{\partial y_1}{\partial z_1}$$

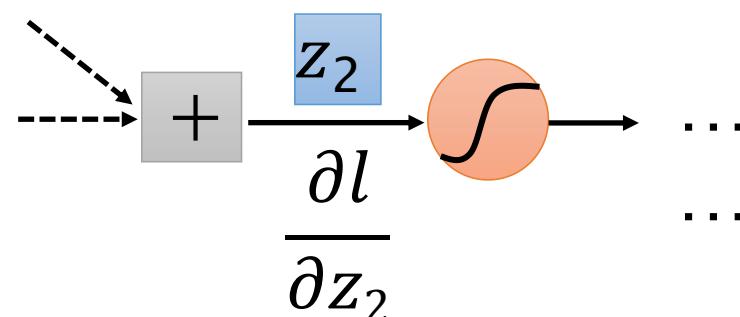
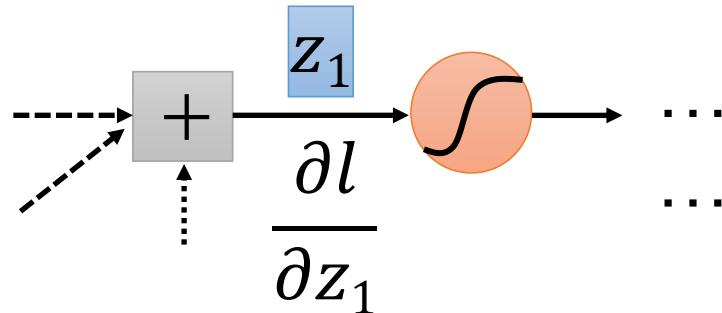
$$\frac{\partial l}{\partial z_2} = \frac{\partial l}{\partial y_2} \frac{\partial y_2}{\partial z_2}$$

Done!

Backpropagation – Backward pass

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

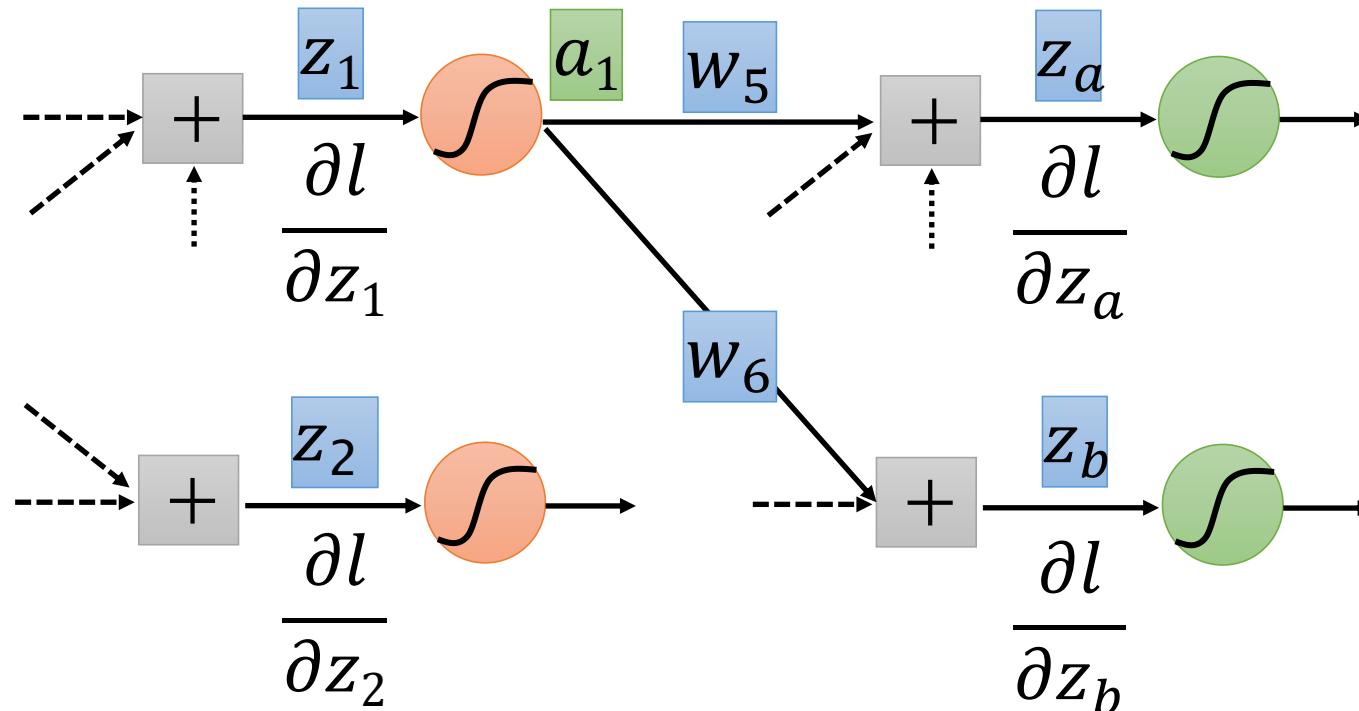
Case 2. Not Output Layer



Backpropagation – Backward pass

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

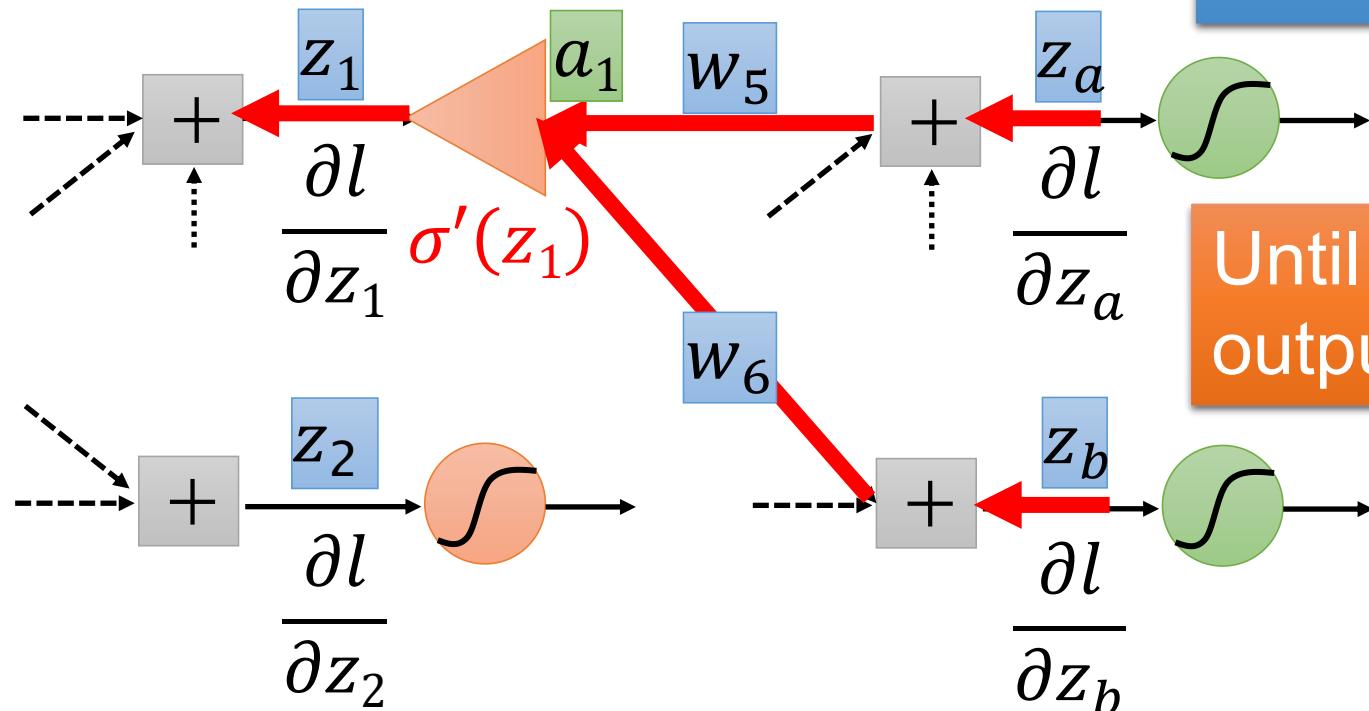
Case 2. Not Output Layer



Backpropagation – Backward pass

Compute $\frac{\partial l}{\partial z}$ for all activation function inputs z

Case 2. Not Output Layer



Compute $\frac{\partial l}{\partial z}$ recursively

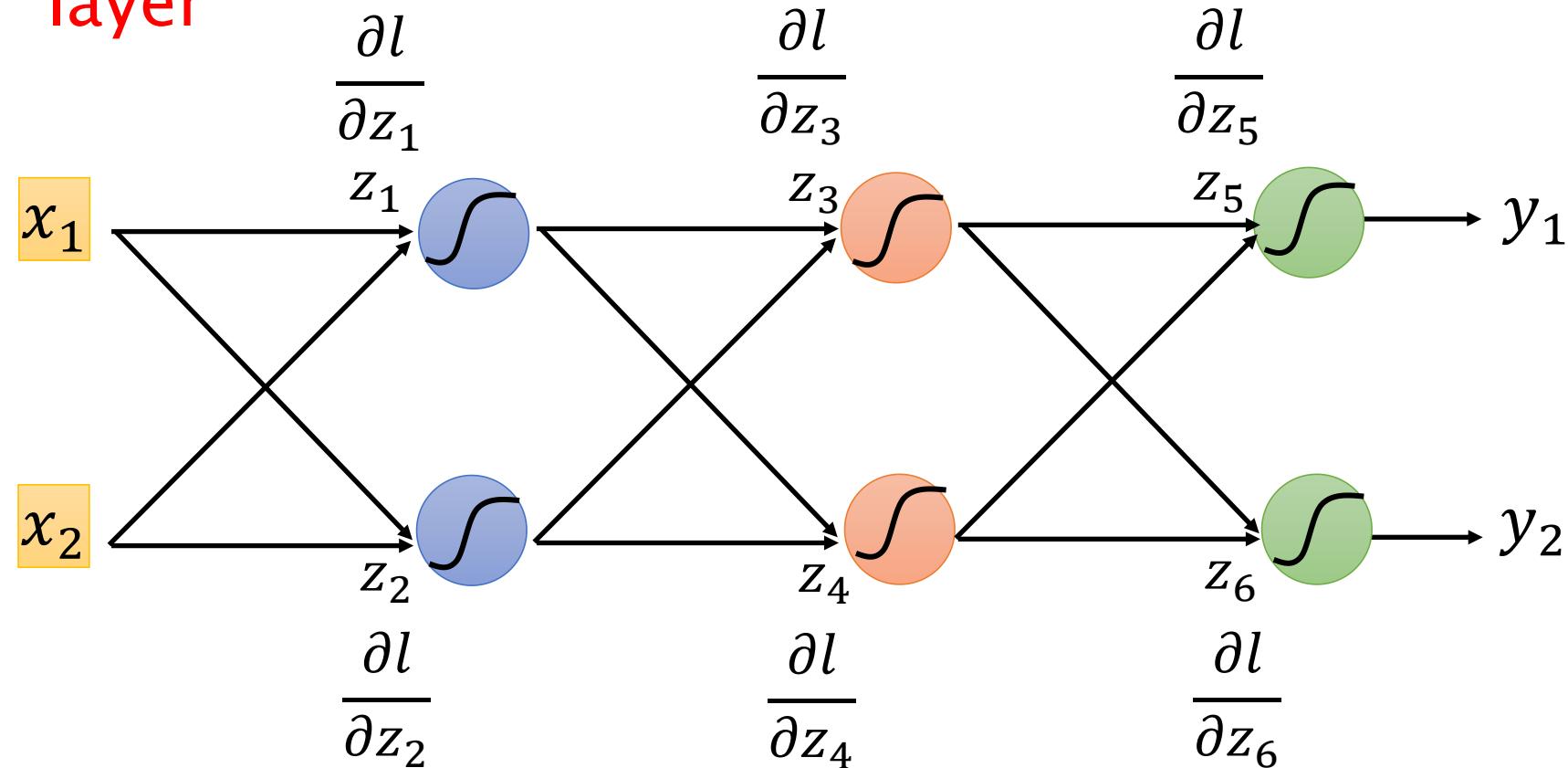
Until we reach the output layer

Backpropagation – Backward pass

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

Compute $\partial l / \partial z$ from the output layer

$$\frac{\partial l}{\partial z} = \sigma'(z) \left[w_3 \frac{\partial l}{\partial z_1} + w_4 \frac{\partial l}{\partial z_2} \right]$$

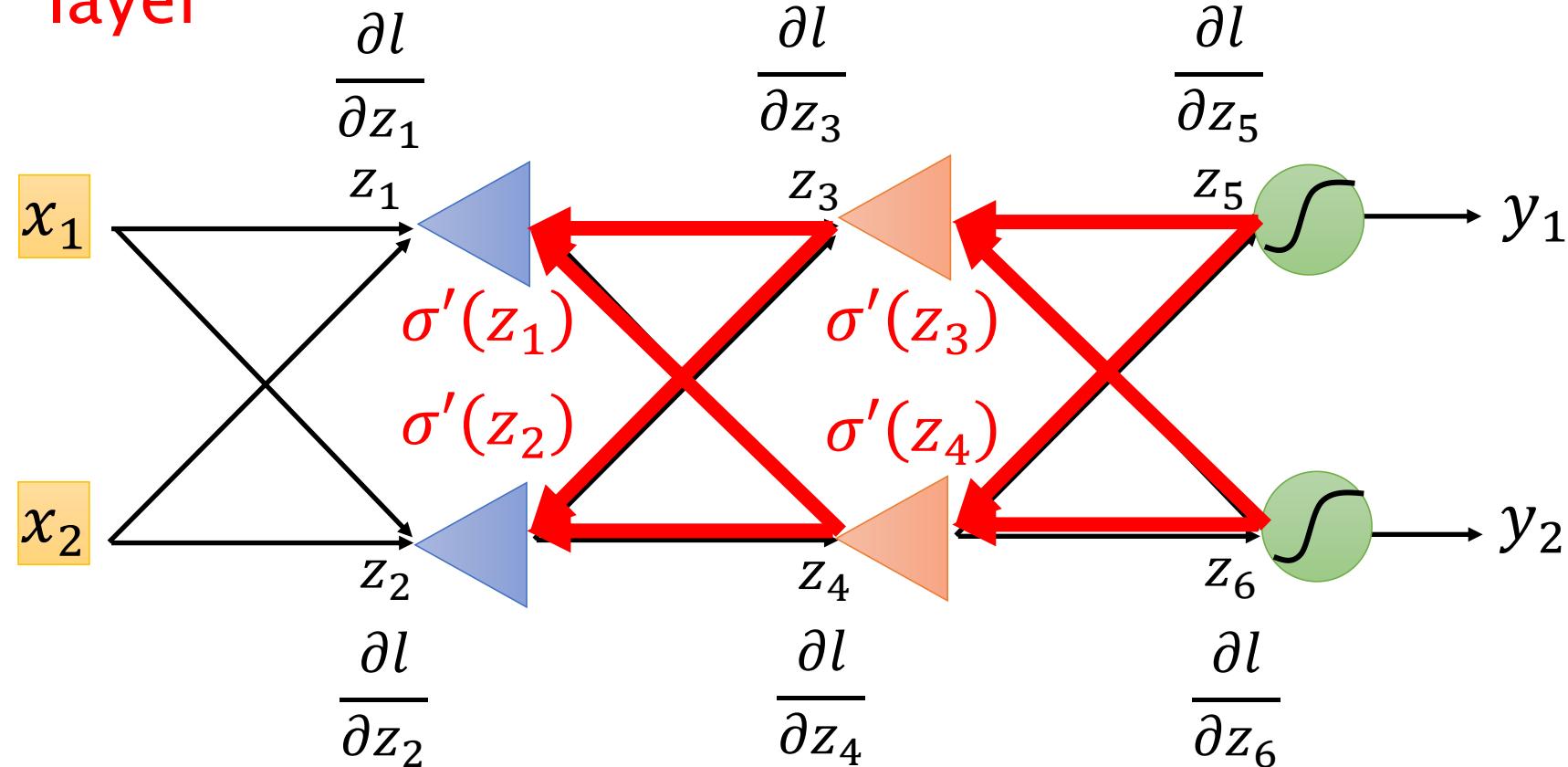


Backpropagation – Backward pass

Compute $\frac{\partial l}{\partial z}$ for all activation function inputs z

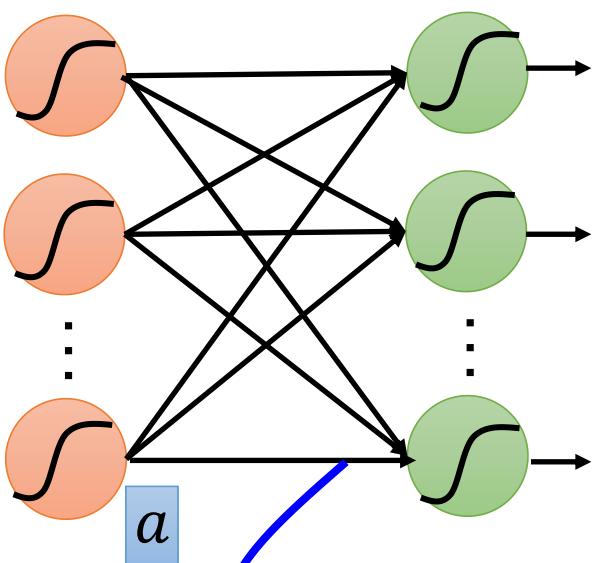
Compute $\frac{\partial l}{\partial z}$ from the output layer

$$\frac{\partial l}{\partial z} = \sigma'(z) \left[w_3 \frac{\partial l}{\partial z_1} + w_4 \frac{\partial l}{\partial z_2} \right]$$



Backpropagation – Overview

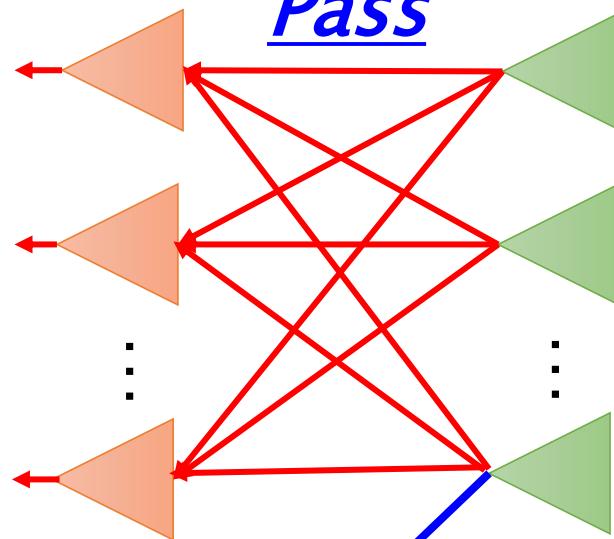
Forward Pass



- Initiate parameters all w^0
- Go forward pass to get all a^0
- $\frac{\partial z}{\partial w_{i,j}^0} = a_{i,j}^0$

$$\frac{\partial z}{\partial w} = a$$

Backward Pass



X

$$\frac{\partial l}{\partial z} = \frac{\partial l}{\partial w}$$

for all w

$$\frac{\partial l}{\partial z_{i,j}} = \sigma'(z_{i,j}) \left[\sum w_{i,j+1} \frac{\partial l}{\partial z_{i,j+1}} \right]$$

$$\frac{\partial l}{\partial z_{1,n}} = \frac{\partial l}{\partial y_{1,n}} \frac{\partial y_{1,n}}{\partial z_{1,n}}$$

Backpropagation (BP) in the Brain?

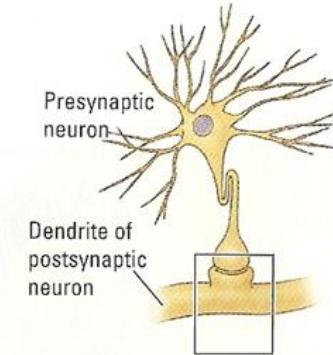
- There is **no** direct evidence that the brain uses a backprop-like algorithm for learning. → Most of researchers think BP in brain is biologically implausible.
- **Difficulties** in implementing BP in the Brain:
 1. Backprop demands **synaptic symmetry (the same weights)** in the forward and backward paths → how to design synaptic connections in forward and backward neural circuits?
 2. Error signals are **signed** and potentially **extreme-valued**. → how to convey signed and extreme-valued errors in real neuron spikes?
 3. Feedback in brains **alters** neural activity. → In NNs, feedback delivers error signals that do **not** influence the activity states of neurons produced by feedforward propagation. But it does change neural activity in the brain.

Learning in the brain (Hebbian learning)



“When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A’s efficiency, as one of the cells firing B, is increased.”

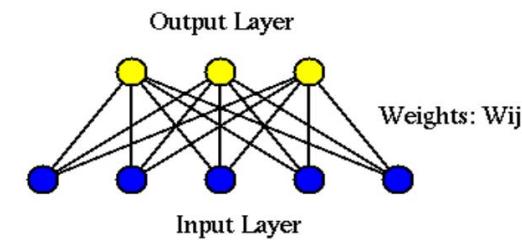
D. O. Hebb, *Organization of Behavior*, 1949



D. O. Hebb

In other words: “**Cells that fire together wire together.**”

Mathematically, this is often written as: $\Delta w_{ij} = \varepsilon x_i x_j$



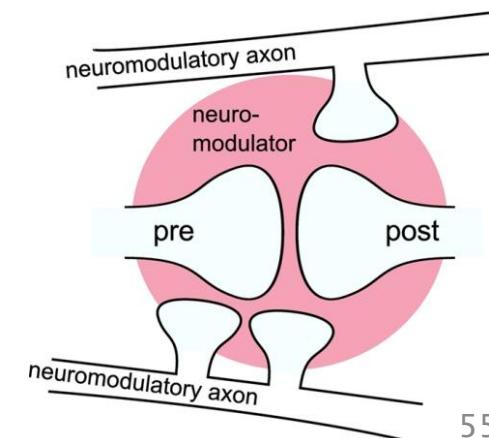
More complex and sophisticated ideas have been under continual exploration for over a half a century, including:

Reward-modulated learning (reinforcement learning)

Competitive learning

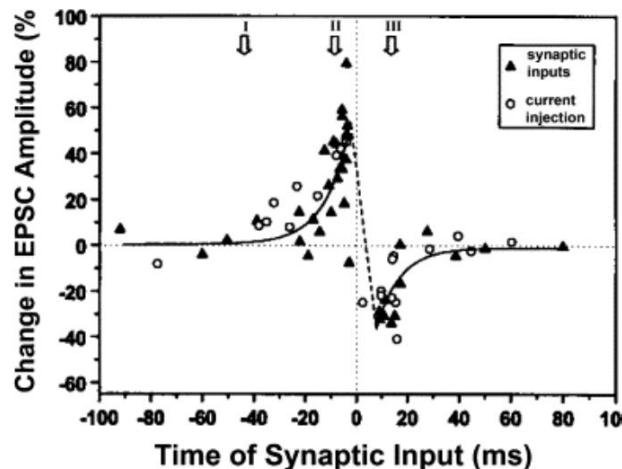
Error correcting learning

Spike-time dependent plasticity



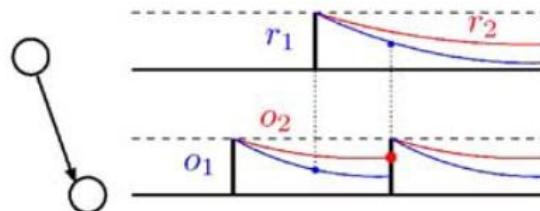
大脑中的前馈学习 及 神经递质调控学习规则

双向STDP学习规则（依赖放电时刻的突触可塑性）



Bi and Poo 1998 J. Neurosci.

三向STDP: 脉冲神经元模型



BCM学习率: 放电率模型

$$dw/dt = A_3^+ \tau_+ \tau_y r_{pre} r \left(r - \frac{\bar{r}(t)}{r_0} \bar{r}(t) \right)$$

三向STDP与BCM放电率学习法则的对应关系

Zenke et. al., 2013 PloS Comp. Biol.;
Pfister and Gerstner 2006 J. Neurosci.

$$\tau_- do_1/dt = -o_1$$

$$\tau_y do_2/dt = -o_2$$

$$\tau_+ dp_1/dt = -p_1$$

$t = t^{post}$ 时：

$$p_1 = p_1 + 1$$

$$\Delta w_{LTP} = p_1 (t - \epsilon) A_3^+$$

$t = t^{pre}$ 时：

$$o_1 = o_1 + 1$$

$$\Delta w_{LTD} = -o_1 A_2^-$$

$$o_2 = o_2 + 1$$

放电率平衡:

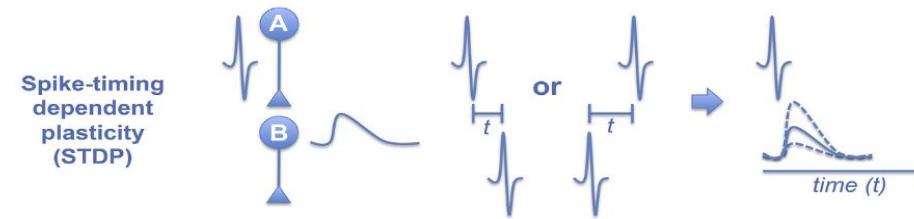
$$A_2^-(t) = \frac{\tau_+ \tau_y \bar{r}(t)^2}{\tau_- r_0} A_3^+$$

神经递质与任务反馈和大脑状态息息相关，携带全局信用信息。

加入神经调质可解决基于前馈学习（STDP 或 BCM）的“信用分配”在时空上的局域性问题，以及任务相关性问题。

乙酰胆碱，去甲肾上腺素，多巴胺，血清素等神经调质的受体的激活和抑制可以调控：

- **学习强度，开关**（调整 A_3^+ , A_2^- 即 r_0 ）
- **学习时间窗口大小**（调整 τ_- 和 τ_y, τ_+ ）
- **脉冲时刻依赖关系** (t_{pre}, t_{post} 和 LTD, LTP 的关系)



同时，这些神经调质的浓度和释放时间与 t_{pre}, t_{post} 可以决定调质的具体调控方向。

因此，为解决局部信用分配问题，需建立前馈学习中相关参数受神经调质调控的动力学模型

神经调质如何调控STDP

Brzozko et. al., 2019 Neuron

大脑与反向传播(Backpropagation and the brain)--上篇

Original NCC lab 神经计算与控制实验室 2020-05-11 22:37

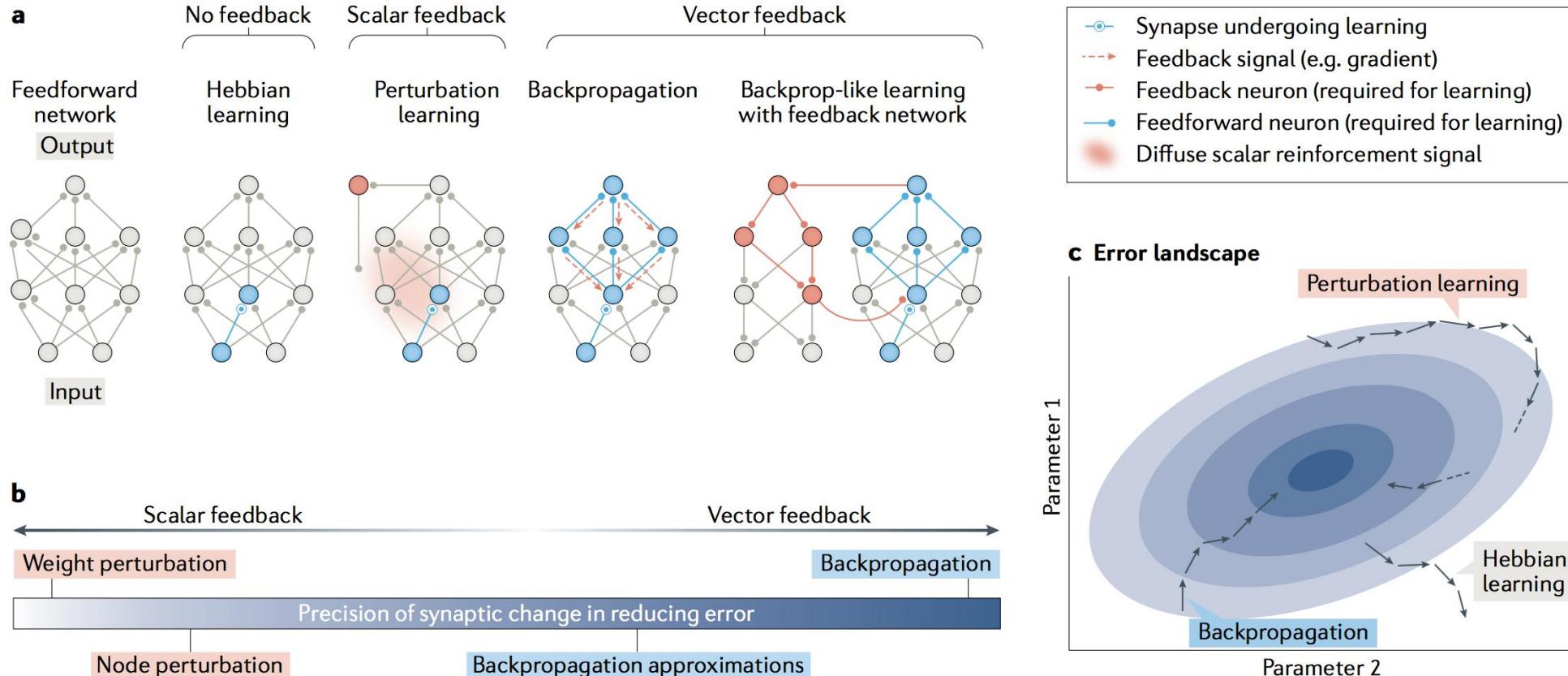


Scan to Follow

<https://mp.weixin.qq.com/s/gDfDYWGE DK9uxkABuAQ6jg>

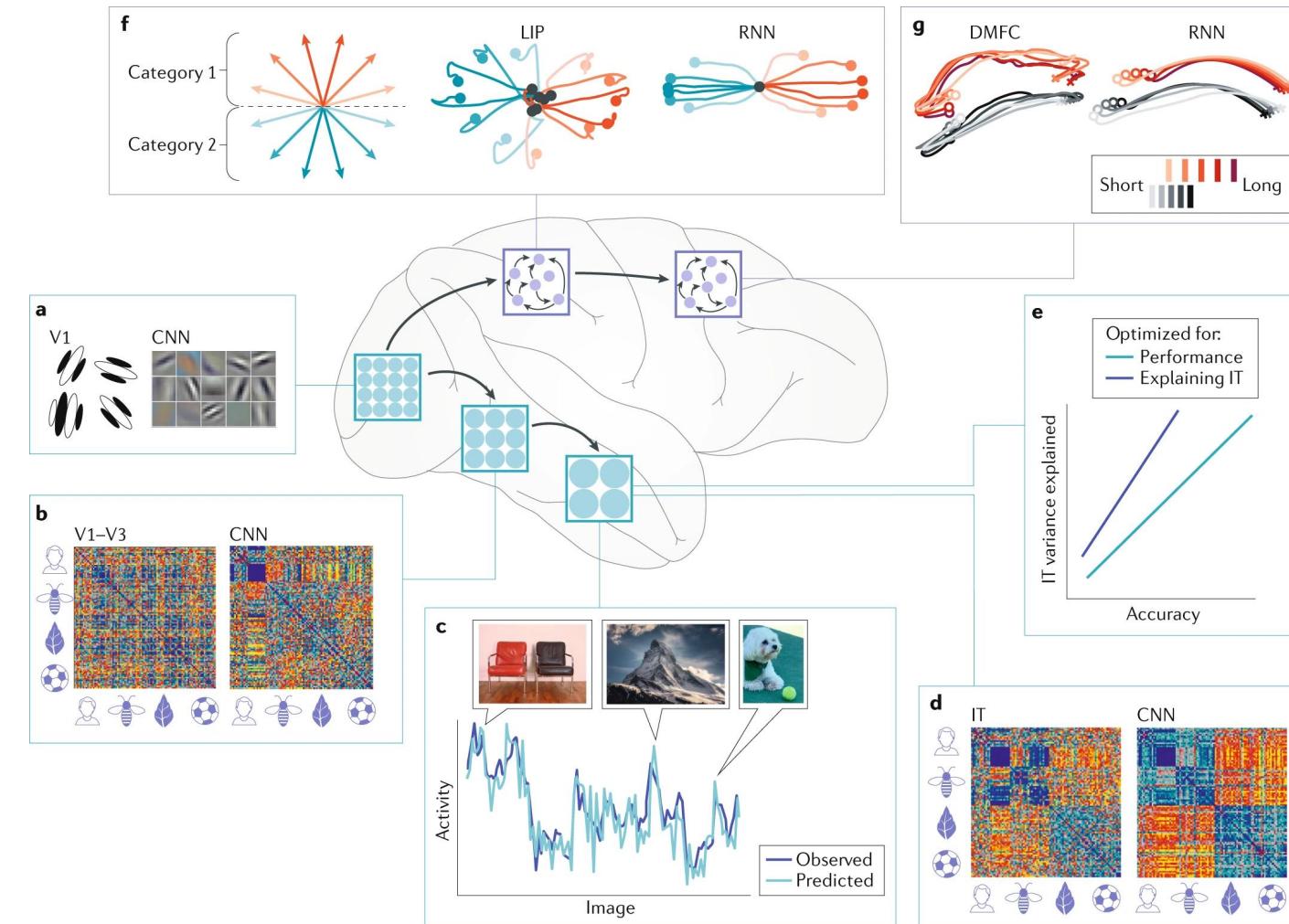
Timothy, Santoro, Marrs, Akerman, Hinton (2020) Backpropagation and the brain, *Nature*

近日 Hinton 等人的研究认为，尽管人脑和神经网络存在明显的物理差异，但大脑仍具有执行反向传播的核心原理的能力。我们对其原文进行了完整的翻译。



AI与BI的相似与差异

- ✓ **组成单元相似:** 都有神经元
- ✓ **网络结构相似:** 都有层级结构
- ✓ **信息表征相似:** 人工神经元活动与大脑神经元活动相关
- ✓ **功能实现相似:** 能完成同样的任务，例如物体识别、语义理解、工作记忆等等



AI与BI的相似与差异

- ✓ **组成单元相似:** 都有神经元
- ✓ **网络结构相似:** 都有层级结构
- ✓ **信息表征相似:** 人工神经元活动与大脑神经元活动相关
- ✓ **功能实现相似:** 能完成同样的任务, 例如物体识别、语义理解、工作记忆等等

- ❖ **学习方式:** 局部反馈 vs 全局误差
- ❖ **所需数据量:** 少量 vs 海量
- ❖ **能耗:** 20瓦 vs 1287兆瓦(单次训练 GPT3)
- ❖ **能力:** 因果推理、知识泛化

