



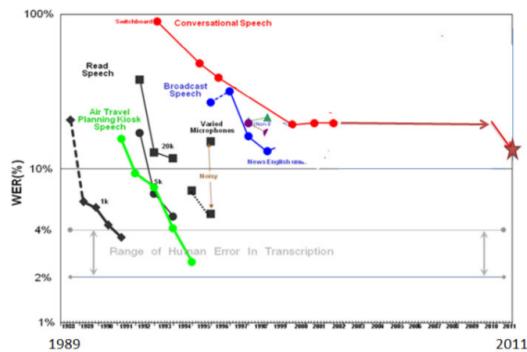
Lecture 14.1: Deep Learning



CS221 / Summer 2019 / Jia

[figure from Li Deng]

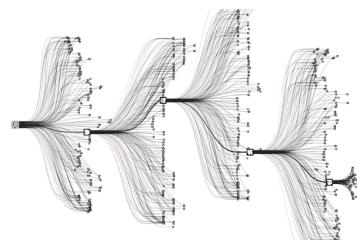
Speech recognition (2009-2011)



- Steep drop in WER due to deep learning
- IBM, Google, Microsoft all switched over from GMM-HMM

CS221 / Summer 2019 / Jia

Go (2016)

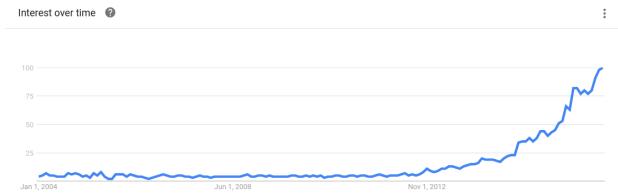


- Defeated world champion Le Sedol 4-1
- Simple architecture (in contrast, DeepBlue was search + hand-crafted heuristics)
- 2017: AlphaGoZero does not require human expert games as supervision

CS221 / Summer 2019 / Jia

Google Trends

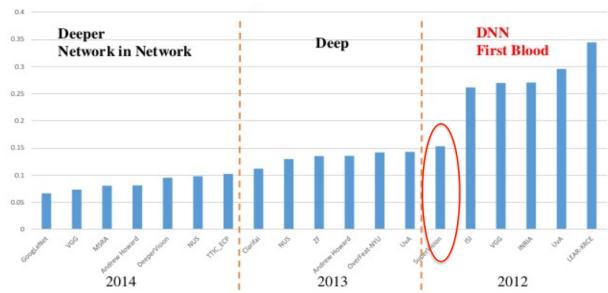
Query: deep learning



CS221 / Summer 2019 / Jia

[Krizhevsky et al., 2012, a.k.a. AlexNet]

Object recognition (2012)

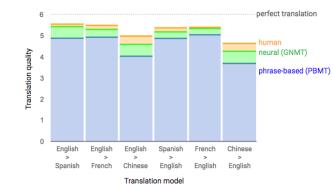


- Landslide win in ILSVRC object recognition competition
- Computer vision community switched to CNNs
- Simpler than hand-engineered features (SIFT)

CS221 / Summer 2019 / Jia

Machine translation (2016)

Input sentence:	Translation (PBMT):	Translation (GNMT):	Translation (human):
李克強今日將啟動中國國家主席和加拿大總理特魯多的首次訪加兩國總理首次對話。	Li Keqiang will start the annual dialogue mechanism with Prime Minister Trudeau of Canada.	Li Keqiang will initiate the annual dialogue mechanism with the Canadian Prime Minister. Li Keqiang and hold the first annual dialogue between the two premiers.	Li Keqiang will start the annual dialogue mechanism with Prime Minister Trudeau of Canada.



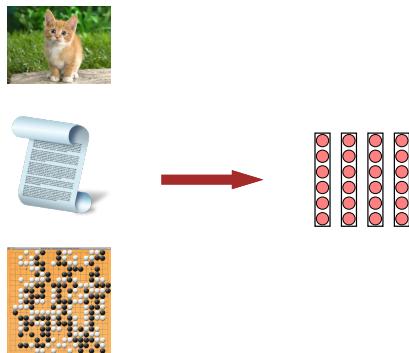
- Decisive wins have taken longer to achieve in NLP (words are meaningful in a way that pixels are not)
- Current state-of-the-art in machine translation
- Simpler architecture (throw out word alignment, phrases tables, language models)

CS221 / Summer 2019 / Jia

5

What is deep learning?

A family of techniques for learning compositional vector representations of complex data.



Roadmap

Feedforward neural networks

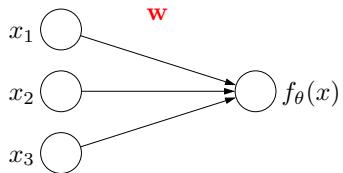
Convolutional neural networks

Recurrent neural networks

Unsupervised and transfer learning

Final remarks

Review: linear predictors

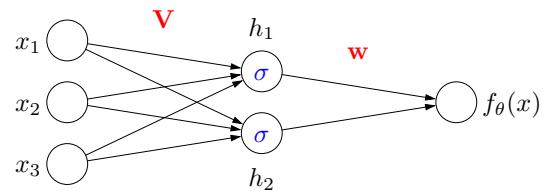


Output:

$$f_\theta(x) = \mathbf{w} \cdot x$$

Parameters: $\theta = \mathbf{w}$

Review: neural networks



Intermediate hidden units:

$$h_j(x) = \sigma(\mathbf{v}_j \cdot x) \quad \sigma(z) = (1 + e^{-z})^{-1}$$

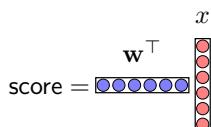
Output:

$$f_\theta(x) = \mathbf{w} \cdot \mathbf{h}(x)$$

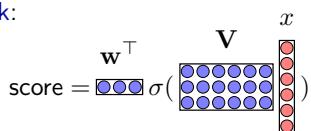
Parameters: $\theta = (\mathbf{V}, \mathbf{w})$

Deep neural networks

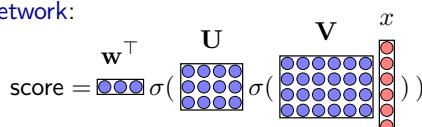
1-layer neural network:



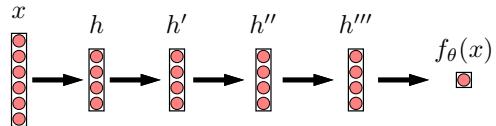
2-layer neural network:



3-layer neural network:



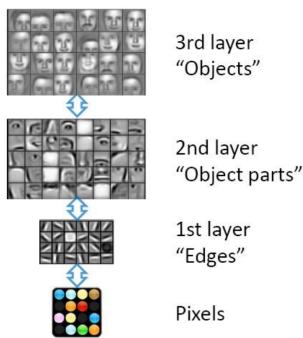
Depth



Intuitions:

- Hierarchical feature representations
- Formal theory/understanding is still incomplete

What's learned?



Review: optimization

Regression:

$$\text{Loss}(x, y, \theta) = (f_\theta(x) - y)^2$$

Key idea: minimize training loss

$$\text{TrainLoss}(\theta) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \text{Loss}(x, y, \theta)$$

$$\min_{\theta \in \mathbb{R}^d} \text{TrainLoss}(\theta)$$

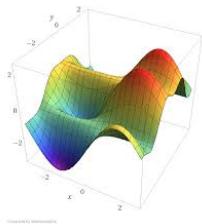
Algorithm: stochastic gradient descent

For $t = 1, \dots, T$:

For $(x, y) \in \mathcal{D}_{\text{train}}$:

$$\theta \leftarrow \theta - \eta_t \nabla_{\theta} \text{Loss}(x, y, \theta)$$

Training



- Non-convex optimization
- No theoretical guarantees that it works
- Before 2000s, empirically very difficult to get working

What's different today

Computation (time/memory)



Information (data)



Other tricks: adaptive step sizes (ADAM), batch normalization, more hidden units, Dropout, pretraining

Summary

- Deep networks learn hierarchical representations of data
- Train via SGD, use backpropagation to compute gradients
- Non-convex optimization, but works empirically given enough compute and data

Roadmap

Feedforward neural networks

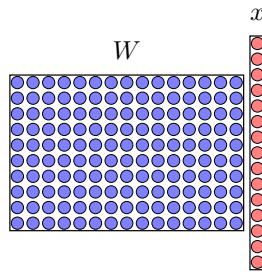
Convolutional neural networks

Recurrent neural networks

Unsupervised and transfer learning

Final remarks

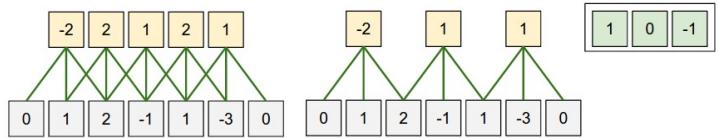
Motivation



- **Observation:** images are not arbitrary vectors
- **Goal:** leverage spatial structure of images (translation invariance)

[figure from Andrej Karpathy]

Prior knowledge



- **Local connectivity:** each hidden unit operates on a local image patch (3 instead of 7 connections per hidden unit)
- **Parameter sharing:** processing of each image patch is same (3 parameters instead of $3 \cdot 5$)
- **Intuition:** try to match a pattern in image

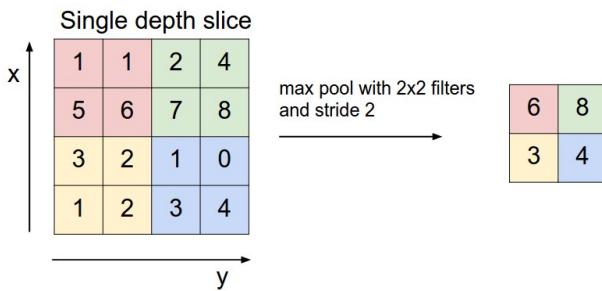
[Andrej Karpathy's demo]

CS221 / Summer 2019 / Jia

18 CS221 / Summer 2019 / Jia

19

Max-pooling



- Intuition: test if there exists a pattern in neighborhood
- Reduce computation, prevent overfitting

CS221 / Summer 2019 / Jia

20 CS221 / Summer 2019 / Jia

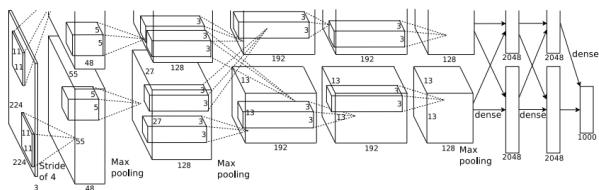
21

History of CNNs

- 1980: Neocognitron, a.k.a. convolutional neural networks (Fukushima)
- 1989: backpropagation on convolutional neural networks (LeCun)

Krizhevsky et al., 2012]

AlexNet

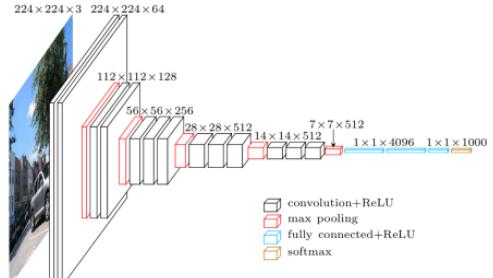


- **Non-linearity:** use ReLU ($\max(z, 0)$) instead of logistic
- **Data augmentation:** translate, horizontal reflection, vary intensity, dropout (guard against overfitting)
- **Computation:** parallelize across two GPUs (6 days)
- **Results on ImageNet:** 16.4% error (next best was 25.8%)

[Simonyan/Zisserman, 2014]

[image credit: Davi Frossard]

VGGNet



- **Architecture:** deeper but smaller filters; uniform
- **Computation:** 4 GPUs for 2-3 weeks
- **Results on ImageNet:** 7.3% error (AlexNet: 16.4%)

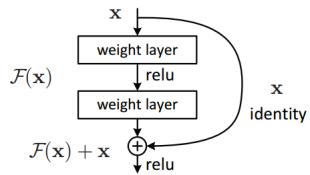
CS221 / Summer 2019 / Jia

22 CS221 / Summer 2019 / Jia

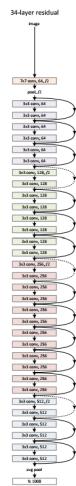
23

Residual networks

$$x \mapsto \sigma(Wx) + x$$



- Key idea: make it easy to learn the identity (good inductive bias)
- Enables training 152 layer networks
- Results on ImageNet: 3.6% error



Summary

- Key idea: locality of connections, capture spatial structure
- Filters have parameter sharing; most parameters in last fully connected layers
- Depth really matters
- Applications to text, Go, drug design, etc.

Roadmap

Feedforward neural networks

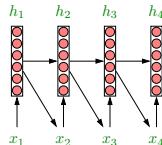
Convolutional neural networks

Recurrent neural networks

Unsupervised and transfer learning

Final remarks

Recurrent neural networks



$$h_1 = \text{Encode}(x_1)$$

$$x_2 \sim \text{Decode}(h_1)$$

$$h_2 = \text{Encode}(h_1, x_2)$$

$$x_3 \sim \text{Decode}(h_2)$$

$$h_3 = \text{Encode}(h_2, x_3)$$

$$x_4 \sim \text{Decode}(h_3)$$

$$h_4 = \text{Encode}(h_3, x_4)$$

Update context vector:

$$h_t = \text{Encode}(h_{t-1}, x_t)$$

Predict next character:

$$x_{t+1} = \text{Decode}(h_t)$$

context h_t compresses x_1, \dots, x_t

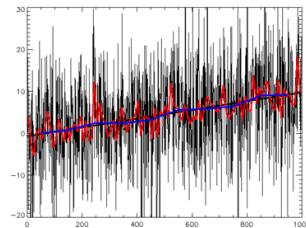
Motivation: modeling sequences

Sentences:

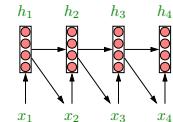
$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \quad x_7 \quad x_8 \quad x_9 \quad x_{10} \quad x_{11} \quad x_{12}$

Paris Talks Set Stage for Action as Risks to the Climate Rise

Time series:



Simple recurrent network



$$\text{Encode}(h_{t-1}, x_t) = \sigma(V h_{t-1} + W)$$

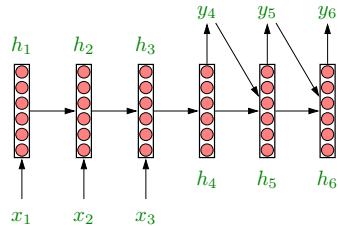
$$\text{Decode}(h_t) \sim \text{softmax}(W' h_t) = p(x_{t+1})$$

Sequence-to-sequence model

Motivation: machine translation

x: Je crains l'homme de un seul livre.

y: Fear the man of one book.



Read in a sentence first, output according to RNN:

$$h_t = \text{Encode}(h_{t-1}, x_t \text{ or } y_{t-1}), \quad y_t = \text{Decode}(h_t)$$

Attention-based models

Motivation: long sentences — compress to finite dimensional vector?

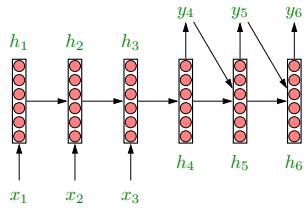
Eine Folge von Ereignissen bewirkte, dass aus Beethovens Studienreise nach Wien ein dauerhafter und endgültiger Aufenthalt wurde. Kurz nach Beethovens Ankunft, am 18. Dezember 1792, starb sein Vater. 1794 besetzten französische Truppen das Rheinland, und der kurfürstliche Hof musste fliehen.



Key idea: attention

Learn to look back at your notes.

Attention-based models



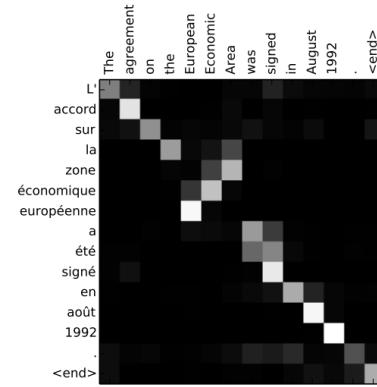
Distribution over input positions:

$$\alpha_t = \text{softmax}([\text{Attend}(h_1, h_{t-1}), \dots, \text{Attend}(h_L, h_{t-1})])$$

Generate with attended input:

$$h_t = \text{Encode}(h_{t-1}, y_{t-1}, \sum_{j=1}^L \alpha_j h_j)$$

Machine translation



Summary

- Recurrent neural networks: model sequences (non-linear version of Kalman filter or HMM)
- LSTMs mitigate the vanishing gradient problem
- Attention-based models: when only part of input is relevant at a time
- Newer models with "external memory": memory networks, neural Turing machines



Roadmap

Feedforward neural networks

Convolutional neural networks

Recurrent neural networks

Unsupervised and transfer learning

Final remarks

Motivation

- Deep neural networks require lot of data
- Sometimes not very much labeled data for task of interest
- ...but there is more data for other tasks
- ...and there is plenty of unlabeled data (text, images, videos)

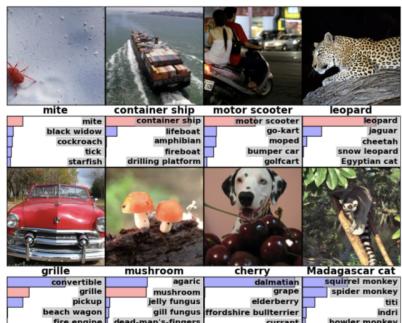
Humans are very good at learning without direct supervision

Transfer learning

Key idea: reuse features that were useful on another task

- Train model on large **source** dataset
- Compute features extracted by trained model
- Train model on small **target** dataset using source model features

Source task: Object recognition



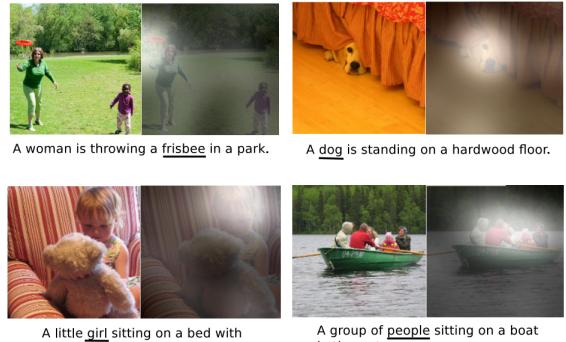
- Map images x to object types y
- Lots of training data (from ImageNet)

Unsupervised pre-training

labeled

unlabeled

Target task: Image captioning



- Map images x to **image captions** y
- Less training data, use ImageNet features (+ fine-tuning)

Reading comprehension (SQuAD)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail.... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called "showers".

What causes precipitation to fall?
gravity

100K examples

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?
graupel

Where do water droplets collide with ice crystals to form precipitation?
within a cloud

Raw text

Stanford University (officially Leland Stanford Junior University)^[10] (colloquially "the Farm") is a private research university in Stanford, California. Stanford is known for its academic strength, wealth, proximity to Silicon Valley, and ranking as one of the world's top universities.^{[11][12][13][14][15]}
 The university was founded in 1895 by Leland and Jane Stanford in memory of their only child, Leland Stanford Jr., who had died of typhoid fever at age 15 the previous year. Stanford was a U.S. Senator and former Governor of California who made his fortune as a railroad tycoon. The school admitted its first students on October 1, 1891,^[27] as a coeducational and non-denominational institution.
 Stanford University struggled financially after the death of Leland Stanford in 1893 and again after much of the campus was damaged by the 1906 San Francisco earthquake.^[16] Following World War II, Provost Frederic Terman supported faculty and graduates' entrepreneurship to build self-sufficient local industry in what would later be known as Silicon Valley.^[17] The university is also one of the top fundraising institutions in the country, becoming the first school to raise more than a billion dollars in a year.^[18]
 The university is organized around three traditional schools consisting of 40 academic departments at the undergraduate and graduate level and four professional schools that focus on graduate programs in Law, Medicine, Education and Business. Stanford's undergraduate program is one of the top three most selective in the United States by acceptance rate.^{[19][20][21][22]} Students compete in 36 varsity sports, and the university is one of two private institutions in the Division I FBS Pac-12 Conference. It has gained 117 NCAA team championships,^[24] the most for a university. Stanford athletes have won 512 individual championships,^[25] and Stanford has won the NACDA Directors' Cup for 23 consecutive years, beginning in 1994–1995.^[26] In addition, Stanford students and alumni have won 270 Olympic medals including 139 gold medals.^[27]

...

3.3 billion words



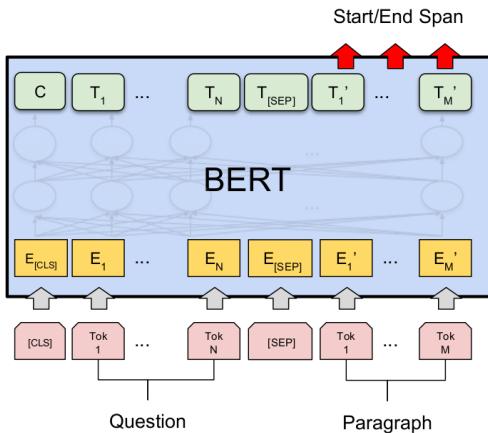
BERT

Paris Talks ___ Stage for _____ as Risks to ___ Climate Rise
 ↓
 Paris Talks Set Stage for Action as Risks to the Climate Rise

- Tasks: fill in words, predict whether is next sentence

- Trained on 3.3B words, 4 days on 64 TPUs

BERT



Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar et al. '16)	82.304	91.221
1 Oct 05, 2018	BERT (ensemble) Google A.I.	87.433	93.160
2 Oct 05, 2018	BERT (single model) Google A.I.	85.083	91.835
2 Sep 09, 2018	nlnet (ensemble) Microsoft Research Asia	85.356	91.202
2 Sep 26, 2018	nlnet (ensemble) Microsoft Research Asia	85.954	91.677
3 Jul 11, 2018	QANet (ensemble) Google Brain & CMU	84.454	90.490
4 Jul 08, 2018	r-net (ensemble) Microsoft Research Asia	84.003	90.147
5 Mar 19, 2018	QANet (ensemble) Google Brain & CMU	83.877	89.737
5 Sep 09, 2018	nlnet (single model) Microsoft Research Asia	83.468	90.133
5 Jun 20, 2018	MARS (ensemble) YUANFUDAO research NLP	83.982	89.796
6 Sep 01, 2018	MARS (single model) YUANFUDAO research NLP	83.185	89.547

Unsupervised learning

- Principle: make up prediction tasks (e.g., x given x or context)
- Hard task → pressure to learn something
- Loss minimization using SGD
- Discriminatively fine tune: initialize feedforward neural network and backpropagate to optimize task accuracy



Roadmap

Feedforward neural networks

Convolutional neural networks

Recurrent neural networks

Unsupervised and transfer learning

Final remarks

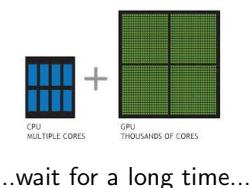
Getting things to work

Better optimization algorithms: SGD, SGD+momentum, AdaGrad, AdaDelta, momentum, Nesterov, Adam

Tricks: initialization, gradient clipping, batch normalization, dropout

More hyperparameter tuning: step sizes, architectures

Better hardware: GPUs, TPUs



...wait for a long time...

Theory: why does it work?

Two questions:

- Approximation: why are neural networks good hypothesis classes?
- Optimization: why can SGD optimize a high-dimensional non-convex problem?

Partial answers:

- 1-layer neural networks can approximate any continuous function on compact set [Cybenko, 1989; Barron, 1993]
- Use statistical physics to analyze loss surfaces [Choromanska et al., 2014]

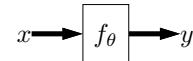
Summary

Phenomena	Ideas
Fixed vectors	Feedforward NNs
Spatial structure	convolutional NNs
Sequence	recurrent NNs LSTMs
Sequence-to-sequence	encoder-decoder attention-based models
Transfer/Unsupervised	any auxiliary task

Extensibility: able to compose modules



Learning programs: think about analogy with a computer



Outlook