Overview for today

- Natural Language Processing with NNs [~15m]
 Supervised models
- Unsupervised Learning [~45m]

• Memory in Neural Nets [~30m]

Natural Language Processing

Slides from:



Jason Weston Tomas Mikolov Antoine Bordes Wojciech Zaremba

NLP

- Many different problems
 - Language modeling
 - Machine translation
 - Q&A
- Recent attempts to address with neural nets

 Yet to achieve same dramatic gains as vision/speech

Language modeling

- Natural language is a sequence of sequences
- Some sentences are more likely than others:
 - "How are you ?" has a high probability
 - "How banana you ? " has a low probability

Neural Network Language Models



Bengio, Y., Schwenk, H., Sencal, J. S., Morin, F., & Gauvain, J. L. (2006). Neural probabilistic language models. In Innovations in Machine Learning (pp. 137-186). Springer Berlin Heidelberg.

[Slide: Antoine Border & Jason Weston, EMNLP Tutorial 2014]

Recurrent Neural Network Language Models

Key idea: *input to predict next word is current word plus context fed-back from previous word (i.e. remembers the past with recurrent connection).*



Figure: Recurrent neural network based LM

Recurrent neural network based language model. Mikolov et al., Interspeech, '10.

[Slide: Antoine Border & Jason Weston, EMNLP Tutorial 2014]

Recurrent neural networks - schema



Backpropagation through time

- The intuition is that we unfold the RNN in time
- We obtain deep neural network with shared weights U and W



[Slide: Thomas Mikolov, COLING 2014]

Backpropagation through time

- We train the unfolded RNN using normal backpropagation + SGD
- In practice, we limit the number of unfolding steps to 5 – 10
- It is computationally more efficient to propagate gradients after few training examples (batch mode)

Tomas Mikolov, COLING 2014



[Slide: Thomas Mikolov, COLING 2014]

NNLMS vs. RNNS: Penn Treebank Results (Mikolov)

Model	Weight	PPL
3-gram with Good-Turing smoothing (GT3)	0	165.2
5-gram with Kneser-Ney smoothing (KN5)	0	141.2
5-gram with Kneser-Ney smoothing + cache	0.0792	125.7
Maximum entropy model	0	142.1
Random clusterings LM	0	170.1
Random forest LM	0.1057	131.9
Structured LM	0.0196	146.1
Within and across sentence boundary LM	0.0838	116.6
Log-bilinear LM	0	144.5
Feedforward NNLM	0	140.2
Syntactical NNLM	0.0828	131.3
Combination of static RNNLMs	0.3231	102.1
Combination of adaptive RNNLMs	0.3058	101.0
ALL	1	83.5

Recent uses of NNLMs and RNNs to improve machine translation: Fast and Robust NN Joint Models for Machine Translation, Devlin et al, ACL '14. Also Kalchbrenner '13, Sutskever et al., '14., Cho et al., '14.

[Slide: Antoine Border & Jason Weston, EMNLP Tutorial 2014]

Language modelling – RNN samples

the meaning of life is that only if an end would be of the whole supplier. widespread rules are regarded as the companies of refuses to deliver. in balance of the nation's information and loan growth associated with the carrier thrifts are in the process of slowing the seed and commercial paper.

More depth gives more power



LSTM - Long Short Term Memory

[Hochreiter and Schmidhuber, Neural Computation 1997]

- Ad-hoc way of modelling long dependencies
- Many alternative ways of modelling it
- Next hidden state is modification of previous hidden state (so information doesn't decay too fast).



For simple explanation, see [Recurrent Neural Network Regularization, Wojciech Zaremba, Ilya Sutskever, Oriol Vinyals, arXiv 1409.2329, 2014]

RNN-LSTMs for Machine Translation



Sequence to Sequence Learning with Neural Networks, Ilya Sutskever, Oriol Vinyals, Quoc Le, NIPS 2014

Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation, Kyunghyun Cho, Bart van Merrienboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Yoshua Bengio, EMNLP 2014

Visualizing Internal Representation

t-SNE projection of network state at end of input sentence



Ilya Sutskever, Oriol Vinyals, Quoc Le, NIPS 2014

Translation - examples

• FR: Les avionneurs se querellent au sujet de la largeur des sièges alors que de grosses commandes sont en jeu

• Google Translate: Aircraft manufacturers are quarreling about the seat width as large orders are at stake

• LSTM: Aircraft manufacturers are concerned about the width of seats while large orders are at stake

• Ground Truth: Jet makers feud over seat width with big orders at stake

[Sequence to Sequence Learning with Neural Networks, Ilya Sutskever, Oriol Vinyals, Quoc Le, NIPS 2014]

Image Captioning: Vision + NLP

- Generate short text descriptions of image, given just picture.
 - Use Convnet to extract image features
- RNN or LSTM model takes image features as input, generates text



Many recent works on this:

٠

- Baidu/UCLA: Explain Images with Multimodal Recurrent Neural Networks
- Toronto: Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models
- Berkeley: Long-term Recurrent Convolutional Networks for Visual Recognition and Description
- Google: Show and Tell: A Neural Image Caption Generator
- Stanford: Deep Visual-Semantic Alignments for Generating Image Description
- UML/UT: Translating Videos to Natural Language Using Deep Recurrent Neural Networks
- Microsoft/CMU: Learning a Recurrent Visual Representation for Image Caption Generation
- Microsoft: From Captions to Visual Concepts and Back

Image Captioning Examples



[men (0.59)] [group (0.66)] [woman (0.64)] [people (0.89)] [holding (0.60)] [playing (0.61)] [tennis (0.69)] [court (0.51)] [standing (0.59)] [skis (0.58)] [street (0.52)] [man (0.77)] [skateboard (0.67)]

a group of people standing next to each other people stand outside a large ad for gap featuring a young boy



[person (0.55)] [street (0.53)] [holding (0.55)] [group (0.63)] [slope (0.51)] [standing (0.62)] [snow (0.91)] [skis (0.74)] [player (0.54)] [people (0.85)] [men (0.57)] [skiing (0.51)] [skateboard (0.89)] [riding (0.75)] [tennis (0.74)] [trick (0.53)] [skate (0.52)] [woman (0.52)] [man (0.86)] [down (0.61)]

a group of people riding skis down a snow covered slope a guy on a skate board on the side of a ramp



g in the direction of the pigeons



a baby elephant standing next to each other on a field elephants are playing together in a shallow watering hole

From Captions to Visual Concepts and Back, Hao Fang* Saurabh Gupta* Forrest Iandola* Rupesh K. Srivastava*, Li Deng Piotr Dollar, Jianfeng Gao Xiaodong He, Margaret Mitchell John C. Platt, C. Lawrence Zitnick, Geoffrey Zweig, CVPR 2015.

Unsupervised Learning

Motivation

• Most successes obtained with supervised models, e.g. Convnets



• Unsupervised learning methods less successful

• But likely to be very important in long-term

Historical Note

- Deep Learning revival started in ~2006
 Hinton & Salakhudinov Science paper on RBMs
- Unsupervised Learning was focus from 2006-2012

 In ~2012 great results in vision, speech with supervised methods appeared

 Less interest in unsupervised learning

Arguments for Unsupervised Learning

- Want to be able to exploit unlabeled data
 Vast amount of it often available
 - Essentially free
- Good regularizer for supervised learning – Helps generalization
 - Transfer learning
 - Zero / one-shot learning

Another Argument for Unsupervised Learning

When we're learning to see, nobody's telling us what the right answers are — we just look. Every so often, your mother says "that's a dog", but that's very little information.

You'd be lucky if you got a few bits of information — even one bit per second — that way. The brain's visual system has 10^{14} neural connections. And you only live for 10^9 seconds.

So it's no use learning one bit per second. You need more like 10^5 bits per second. And there's only one place you can get that much information: from the input itself.

— Geoffrey Hinton, 1996

Taxonomy of Approaches

- Autoencoder (most unsupervised Deep Learning methods)
 - RBMs / DBMs
 - Denoising autoencoders
 - Predictive sparse decomposition
- Decoder-only
 - Sparse coding
 - Deconvolutional Nets
- Encoder-only
 - Implicit supervision, e.g. from video
- Adversarial Networks

Loss involves some kind of reconstruction error

Auto-Encoder





Auto-Encoder Example 2

• Predictive Sparse Decomposition [Ranzato et al., '07]



Auto-Encoder Example 2

• Predictive Sparse Decomposition [Kavukcuoglu et al., '09]



Stacked Auto-Encoders



Training phase 2: Supervised Fine-Tuning



Effects of Pre-Training

• From [Hinton & Salakhudinov, Science 2006]



See also: Why Does Unsupervised Pre-training Help Deep Learning? Dumitru Erhan, Yoshua Bengio ,Aaron Courville, Pierre-Antoine Manzagol PIERRE-Pascal Vincent, Sammy Bengio, JMLR 2010

Deep Boltzmann Machines



Shape Boltzmann Machine



Figure 2. **Undirected models of shape:** (a) 1D slice of a Markov Random Field. (b) Restricted Boltzmann Machine in 1D. (c) Deep Boltzmann Machine in 1D. (d) 1D slice of a Shape Boltzmann Machine. (e) Shape Boltzmann Machine in 2D.



"The Shape Boltzmann Machine: a Strong Model of Object Shape", Ali Eslami, Nicolas Heess and John Winn, CVPR 2012

Variational Auto-Encoder

• [Kingma & Welling, ICLR 2014]



[Slide: Ian Goodfellow, Deep Learning workshop, ICML 2015]

Decoder-Only Models

- Examples:
 - Sparse coding
 - Deconvolutional Networks [Zeiler & Fergus, '10]
- No encoder to compute features
- So need to perform optimization
 Can be relatively fast

Sparse Coding (Patch-based)

• Over-complete linear decomposition of input y using dictionary D

$$= 0.3 \times + 0.5 \times + 0.2 \times$$

9

$$C(y,D) = \underset{z}{\operatorname{argmin}} \ \frac{\lambda}{2} \|Dz - y\|_{2}^{2} + |z|_{1}$$



Dictionary D

- *l*₁ regularization yields solutions with few non-zero elements
- Output is sparse vector: $z = [0, 0.3, 0, \dots, 0.5, \dots, 0.2, \dots, 0]$
Deconvolutional Network Layer

 Convolutional form of sparse coding [Zeiler & Fergus, CVPR 2010]. Also Kavukcuoglu et al. NIPS 2010



Overall Architecture (2 layers)



Generative Models using Convnets

- Learning to Generate Chairs with Convolutional Neural Networks, Alexey Dosovitskiy, Jost Tobias Springenberg and Thomas Brox, 1411.5928, 2014
- Supervised training of convnet to draw chairs



Some other interesting generative models

- "Generative Image Modeling Using Spatial LSTMs", L. Theis and M. Bethge, arXiv 1506.03478, 2015
- "Texture synthesis and the controlled generation of natural stimuli using convolutional neural networks", Leon A. Gatys, Alexander S. Ecker, Matthias Bethge, . arXiv:1505.07376, 2015



Encoder-Only Models

• In vision setting, essentially a convnet trained without explicit class labels

- Learn invariances
 - Unsupervised feature learning by augmenting single images, Alexey Dosovitskiy, Jost Tobias Springenberg and Thomas Brox, NIPS 2014
- Learn from video
 - Unsupervised Learning of Visual Representations using Videos Xiaolong Wang, Abhinav Gupta, arXiv 1505.00687, 2015

Unsupervised Learning of Transformations

[Unsupervised feature learning by augmenting single images, Alexey Dosovitskiy, Jost Tobias Springenberg and Thomas Brox, NIPS 2014]

- Take patches from images
- For each patch, make lots of peturbed versions



- Treat each patch + peturbed copies as a separate classs
- Train supervised convnet

	-		•	-
	STL-10	CIFAR-10-reduced	CIFAR-10	Caltech-101
K-means [6]	60.1 ± 1	70.7 ± 0.7	82.0	
Multi-way local pooling [5]		—		77.3 ± 0.6
Slowness on videos [25]	61.0	—		74.6
Receptive field learning [16]		_	$[83.11]^1$	75.3 ± 0.7
Hierarchical Matching Pursuit (HMP) [3]	64.5 ± 1	—		
Multipath HMP [4]		_		82.5 ± 0.5
Sum-Product Networks [8]	62.3 ± 1	—	$[83.96]^1$	
View-Invariant K-means [15]	63.7	72.6 ± 0.7	81.9	
This paper	67.4 ± 0.6	69.3 ± 0.4	77.5	$76.6 \pm 0.7^{\ 2}$

Unsupervised Learning from Video

 Unsupervised Learning of Visual Representations using Videos, Xiaolong Wang, Abhinav Gupta, arXiv 1505.00687, 2015



Generative Adversarial Networks

[Generative Adversarial Nets, Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio, NIPS 2014]



Generative Adversarial Networks

[Generative Adversarial Nets, Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio, NIPS 2014]

• Minimax value function:

```
\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]
                                       Discriminator's
                                                                         Discriminator's
        Discriminator
           pushes up
                                            ability to
                                                                              ability to
                                    recognize data as
                                                                              recognize
Generator
                                           being real
                                                                             generator
   pushes
                                                                        samples as being
    down
                  [Slide: Ian Goodfellow, Deep Learning workshop,
                                                                                   fake
                  ICML 2015]
```

Generative Adversarial Networks

[Generative Adversarial Nets, Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio, NIPS 2014]



[Slide: Ian Goodfellow, Deep Learning workshop, ICML 2015]

Adversarial Network Samples









CIFAR-10 (convolutional)

Adversarial Network using Laplacian Pyramid

• [Denton + Chintala, et al. arXiv 1506.05751, 2015]



Adversarial Network using Laplacian Pyramid

• [Denton + Chintala, et al. arXiv 1506.05751, 2015]



Memory in Neural Networks

Sainbayar Sukhbaatar



Introduction

• Recently, there has been lot of interest in incorporating memory and attention to neural networks

- Memory Networks, NTM, Learning to attend ...

- Neural networks are not good at remembering things, especially when input is large but only part of it is relevant
- Adding external memory and learning to attend on important part is key

Outline

- Implicit Internal memory – RNN, LSTM
- Explicit External memory – MemNN, NTM
- Attention models

- MT, Speech, Image, Pointer Network

Implicit Internal Memory

Internal state of the model can be used for memory
 Recurrent Neural Networks (RNNs)



- Computation and memory is mixed
 - Complex computation requires many layers of nonlinearity
 - But some information is lost with each non-linearity
 - Gradient vanishing, Catastrophic forgetting

Ways to Prevent Forgetting in RNNs

- Split state into fast and slow changing parts: structurally constrained recurrent nets (Mikolov et al., 2014)
 - Fast changing part is good for computation
 - Slow changing part is good for storing information
- Gated units for internal state
 - Control when to forget/write using gates
 - Long-short term memory (LSTM) (see Graves, 2013)
 - Simpler Gated Recurrent Unit (GRU) (Cho et al., 2014)
- Other problems
 - Memory capacity is fixed and limited by the dimension of state vector (computation is $O(N^2)$ where N is memory capacity)
 - Vulnerable to distractions in inputs
 - Restricted to sequential inputs

Stack memory for RNN (Joulin et al., 2014)

- Added a stack module to RNN, which can hold a list of vectors
- Action on stack: push, pop and no-op
- More powerful with multiple stacks
- Stack are updated in continuous manner → differentiable
 → trainable by backpropagation + search
- Applied to counting, memorization, binary addition



External Global Memory

- Separate memory from computation
 - Add separate memory module for storage
 - Memory contains list/set of items



- Main module can read and write to the memory
- Advantage: long-term, scalable, flexible

Selective Addressing is Key for Memory

- Often, you only want to interact with few items in memory at once
 - Memory needs some addressing mechanism
- Memory addressing types
 - Soft or hard addressing
 - Soft addressing can be trained by backpropagation
 - Hard addressing is not differentiable (e.g. can be trained with reinforcement learning or additional training signal for where to attend)
 - Context and Location based addressing
 - When input is ordered in some way, location based addressing is useful
 - Location addressing is same as context if location is embedded in the context (e.g. MemN2N)

Memory Networks (Weston et al., 2014)

- Neural network with large external memory
- Writes everything to the memory, but reads only relative information
- Hard addressing: max of the inner product between then internal state and memory contents
- Location based addressing: can compare two memory items by their relative location
- Can perform multiple memory lookups (hops) before producing an output
- Requires additional training signals for training hard addressing
- Applied to toy and large-scale QA tasks



End-to-end Memory Networks (Sukhbaatar et al., 2015)

- Soft addressing: replaced hard max with **softmax**
- End-to-end training: softmax is differentiable →
 can train with backpropagation
- Location addressing: location/time is embedded into the context (special words for "Time=4")
- Applied to toy QA and language modeling

End-to-end Memory Networks (Sukhbaatar et al., 2015)



Single Memory Lookup

End-to-end Memory Networks (Sukhbaatar et al., 2015)



Single Memory Lookup

Multiple Memory Lookup

RNN viewpoint of End-to-End MemNN



Inputs are fed to RNN one-by-one in order. RNN has only one chance to look at a certain input symbol.

Place all inputs in the memory. Let the model decide which part it reads next.

Attention during memory lookups

Samples from toy QA tasks (bAbI dataset)

Result

Story (1: 1 supporting fact)	Support	Hop 1	Hop 2	Hop 3	Story (2: 2 supporting facts)
Daniel went to the bathroom.		0.00	0.00	0.03	John dropped the milk.
Mary travelled to the hallway.		0.00	0.00	0.00	John took the milk there.
John went to the bedroom.		0.37	0.02	0.00	Sandra went back to the bathroom.
John travelled to the bathroom.	yes	0.60	0.98	0.96	John moved to the hallway.
Mary went to the office.	-	0.01	0.00	0.00	Mary went back to the bedroom.
Where is John? Answer: bathroom	Predict	ion: bath	nroom		Where is the milk? Answer: hall
					1
Story (16: basic induction)	Support	Hop 1	Hop 2	Hop 3	Story (18: size reasoning)
Brian is a frog.	ves	0.00	0 08	0.00	The suitease is bigger than the choir
	,	0.00	0.30	0.00	
Lily is gray.	,	0.07	0.00	0.00	The box is bigger than the chocolate
Lily is gray. Brian is yellow.	yes	0.07 0.07	0.00	0.00	The box is bigger than the chocolate The chest is bigger than the chocolate
Lily is gray. Brian is yellow. Julius is green.	yes	0.07 0.07 0.06	0.00 0.00 0.00	0.00 0.00 1.00 0.00	The box is bigger than the chocolat The chest is bigger than the chocolat The chest fits inside the container.
Lily is gray. Brian is yellow. Julius is green. Greg is a frog.	yes yes	0.07 0.07 0.06 0.76	0.00 0.00 0.00 0.02	0.00 0.00 1.00 0.00 0.00	The box is bigger than the chocolat The chest is bigger than the chocolat The chest fits inside the container. The chest fits inside the box.

ohn dropped the milk.		0.06	0.00	0.00
ohn took the milk there.	yes	0.88	1.00	0.00
andra went back to the bathroom.		0.00	0.00	0.00
ohn moved to the hallway.	yes	0.00	0.00	1.00
ary went back to the bedroom.		0.00	0.00	0.00
/here is the milk? Answer: hallway	Predictio	n: hallwa	у	
tory (18: size reasoning)	Support	Hop 1	Hop 2	Hop 3
tory (18: size reasoning) ne suitcase is bigger than the chest.	Support yes	Hop 1 0.00	Hop 2 0.88	Hop 3 0.00
tory (18: size reasoning) ne suitcase is bigger than the chest. ne box is bigger than the chocolate.	Support yes	Hop 1 0.00 0.04	Hop 2 0.88 0.05	Hop 3 0.00 0.10
tory (18: size reasoning) ne suitcase is bigger than the chest. ne box is bigger than the chocolate. ne chest is bigger than the chocolate.	Support yes yes	Hop 1 0.00 0.04 0.17	Hop 2 0.88 0.05 0.07	Hop 3 0.00 0.10 0.90
tory (18: size reasoning) ne suitcase is bigger than the chest. ne box is bigger than the chocolate. ne chest is bigger than the chocolate. ne chest fits inside the container.	Support yes yes	Hop 1 0.00 0.04 0.17 0.00	Hop 2 0.88 0.05 0.07 0.00	Hop 3 0.00 0.10 0.90 0.00
tory (18: size reasoning) ne suitcase is bigger than the chest. ne box is bigger than the chocolate. ne chest is bigger than the chocolate. ne chest fits inside the container. ne chest fits inside the box.	Support yes yes	Hop 1 0.00 0.04 0.17 0.00 0.00	Hop 2 0.88 0.05 0.07 0.00 0.00	Hop 3 0.00 0.10 0.90 0.00 0.00

Support Hop 1 Hop 2 Hop 3

	Test error	Failed tasks
MemNN	6.7%	4
LSTM	51%	20
MemN2N	12.4%	11

Neural Turing Machine (Graves et al., 2014)

- Learns how to write to the memory
- Soft addressing \rightarrow backpropagation training
- Location addressing: small continuous shift of attention
- Complex addressing mechanism: need to sharpen after convolution
- Controller can be LSTM-RNN or feed-forward neural network
- Applied to learn algorithms such as sort, associative recall and copy.
- Hard addressing with reinforcement learning (Zaremba et al., 2015)



RNNsearch: Attention in Machine Translation (Bahdanau et al., 2015)

- RNN based encoder and decoder model
- Decoder can look at past encoder states using soft attention
- Attention mechanism is implement by a small neural network
 - It takes the current decoder state and a past encoder state and outputs a score. Then the all scores are fed to softmax to get attention weights
- Applied to machine translation. Significant improvement in translation of longer sentences



Image caption generation with attention (Xu et al., 2015)

- Encoder: lower convolutional layer of a deep ConvNet (because need spatial information)
- Decoder: LSTM RNN with soft spatial attention
 - Decoder state and encoder state at single location are fed to small NN to get score at that location
- Network attends to the object when it is generating a word for it
- Also hard attention is tried with reinforcement learning



A woman is throwing a <u>frisbee</u> in a park.



A dog is standing on a hardwood floor.



A <u>stop</u> sign is on a road with a mountain in the background.



A little <u>girl</u> sitting on a bed with a teddy bear.



A group of <u>people</u> sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Video description generation (Yao et al., 2015)



+Local+Global: A man and a woman are talking on the road

Ref: A man and a woman ride a motorcycle



Ref: A woman is frying food

(bottom: ground truth)

L. Yao, A. Torabi, K. Cho, N. Ballas, C. Pal, H. Larochelle, and A. Courville, "Describing videos by exploiting temporal structure," *arXiv: 1502.08029*, 2015.

Location-aware attention for speech (Chorowski et al., 2015)

- RNN based encoder-decoder model with attention (similar to RNNsearch)
- Location based addressing: previous attention weights are used as feature for the current attention (good when subsequent attention locations are highly correlated)
- Improvement with sharpening and smoothing of memory addressing



Pointer Network: attention as an output (Vinyals et al., 2015)

- RNN based encoder-decoder model for discrete optimization problems
- Decoder can attend to previous encoder states (similar to RNNsearch, content based soft attention by a small NN)
- Rather than fixed output classes, attention weights determine output
- Input to the most attended encoder state becomes an output
 → can output any sequence of inputs



Resources

• EMNLP 2014 tutorial

- http://emnlp2014.org/tutorials.html#embedding

- CVPR2014 deep learning tutorial

 <u>https://sites.google.com/site/deeplearningcvpr2014/</u>
- ICML 2013 deep learning tutorial
 - http://www.cs.nyu.edu/~yann/talks/lecun-ranzatoicml2013.pdf

Software

- Caffe (http://caffe.berkeleyvision.org/)
 Vision-centric
- Torch (http://torch.ch/)
 - Lua-based library for Deep Learning
 - Currently used by FAIR and Google Deep Mind
- Theano (http://deeplearning.net/software/theano/)
 - Automatic differentiation
 - Python-based
Thanks!

Facebook AI Research colleagues & NYU PhD students:







Yaniv Taigman



Soumith Chintala Emily Denton



Jason Weston Tomas Mikolov Ronan Collobert

Marc'Aurelio Ranzato

Sainbayar Sukhbaatar

FAIR Overview

Facebook Al Research

- Toward Artificial Intelligence (AI), with Machine Learning.
- Established Dec 2013 (1.5 year old)
 - initiative of CEO and CTO
 - Iead by Yann Lecun







FAIR Overview

Facebook AI Research

- ~35 researcher scientists
 - Machine Learning, Computer Vision and Natural Language Processing
- ~15 research engineers
 - Software support, prototyping, interaction with product teams...
- Locations:
 - New York City
 - Menlo Park (HQ)
 - Paris



FAIR Mission

Facebook AI Research

- Advance the state of the art of AI
 - Publish research in best conferences and journals
 - Open-source code release
- Produce software tools for AI research and applications
- Help FB products to leverage advances in AI
 - Software prototyping, architecting, interaction with product teams...



FAIR Impact

Machine Learning @ FB



Computer Vision

- Face detection and identification
- Object detection, scene classification
- Video classification



- Natural Language
 - Tag prediction for search, feed ranking, ad targeting
- Computational Advertising
 - Ads targeting
 - User interest modeling

Huge Scale Deployment of Machine Learning

- 1.4 billion monthly active users
- 850 million daily active users (1 in 7 people on Earth)
- More images uploaded than any other website
 - 400M+ new Facebook photos/day (no labels)
 - 60M+ Instagram images/day (most with hashtags)
 - ~ 500 Billion photos total
- Face and Object recognition models applied to <u>every</u> image
- 5M video uploads/day & growing rapidly
 - More video playback than YouTube

We are hiring!

Internships

https://www.facebook.com/careers/department?
 dept=grad&req=a0IA000000CzCGuMAN

-

Postdoc positions

- Ex-postdocs now faculty at Berkeley, Harvard
- Full-time positions

https://research.facebook.com/ai



Memory References

- A. Graves. Generating sequences with recurrent neural networks. arXiv preprint: 1308.0850, 2013
- T. Mikolov, A. Joulin, S. Chopra, M. Mathieu, and M. A. Ranzato. Learning longer memory in recurrent neural networks. arXiv preprint:1412.7753, 2014
- A. Joulin, and T. Mikolov. Inferring Algorithmic Patterns with Stack-Augmented Recurrent Nets. arXiv preprint:1503.01007, 2015
- K. Cho, B. van Merriënboer, D. Bahdanau, and Y. Bengio. On the properties of neural machine translation: Encoder-decoder approaches. arXiv preprint:1409.1259, 2014
- J. Weston, S. Chopra, and A. Bordes. Memory networks. In International Conference on Learning Representations (ICLR), 2015
- S. Sukhbaatar, A. Szlam, J. Weston, and R. Fergus. End-To-End Memory Networks. arXiv preprint:1503.08895, 2015

Memory References

- A. Graves, G. Wayne, and I. Danihelka. Neural turing machines. arXiv preprint: 1410.5401, 2014
- W. Zaremba, and I. Sutskever. Reinforcement Learning Neural Turing Machines. arXiv preprint:1505.00521, 2015
- D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In International Conference on Learning Representations (ICLR), 2015
- K. Xu, J. L. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. ICML, 2015
- L. Yao, A. Torabi, K. Cho, N. Ballas, C. Pal, H. Larochelle, and A. Courville. Describing videos by exploiting temporal structure. arXiv preprint: 1502.08029, 2015
- J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio. Attention-based models for speech recognition. arXiv preprint: 1506.07503, 2015
- O. Vinyals, M. Fortunato, and N. Jaitly. Pointer networks. arXiv preprint:1506.03134, 2015

Neuroscience of memory

- hippocampus
 - Densely connected
 - Vital for new memory formation
 - From few days to few years
 - Place / grid cells
- Neo-cortex
 - Can keep memory much longer

Memory types

- Short-term memory (working memory)
 Limited capacity
- Long term memory
 - Explicit / Declarative
 - Semantic memory
 - Episodic memory
 - Implicit
 - Procedural memory
 - Priming