

### From Symbolic to Neural Approaches to NLP - Case Studies of Machine Reading and Dialogue

#### **Jianfeng Gao**

Thanks for the slides by **Bill Dolan, Michel Galley**, **Lihong Li, Yi-Min Wang et al.** Joint work with many Microsoft colleagues and interns (see the list of collaborators) Microsoft AI & Research Nov. 11, 2017, Dalian, China

### Outline

- Part 1: The transition of NLP to neural approaches
  - Deep learning leads to paradigm shift in NLP
  - The powers of deep learning
  - Cast study of Deep Semantic Similarity Models
- Part 2: Neural machine reading models for question answering
- Part 3: Deep reinforcement learning for task-completion dialogue

### Traditional definition of NLP: the branch of AI

- Deal with analyzing, understanding and generating the languages that humans use naturally (natural language)
- Study knowledge of language at different levels
  - Phonetics and Phonology the study of linguistic sounds
  - Morphology the study of the meaning of components of words
  - Syntax the study of the structural relationships between words
  - Semantics the study of meaning
  - Discourse they study of linguistic units larger than a single utterance

### Traditional NLP component stack



- Natural language understand (NLU): parsing (speech) input to semantic meaning and update the system state
- 2. Application reasoning and execution: take the next action based on state
- **3. Natural language generation (NLG):** generating (speech) response from action

### *Pragmatic* definition: building computer systems

- Process large text corpora, turning information into knowledge
  - Text classification
  - Information retrieval and extraction
  - Machine reading comprehension and question answering
  - ...
- Enable human-computer interactions, making knowledge accessible to humans in the most natural way
  - Dialogue and conversational agents
  - Machine translation
  - ...

### Challenge of NLP: the diversity of natural language

Many-to-many mapping btw *symbolic* language and *semantic* meaning

#### Ambiguity

Example: I made her duck.

- I cooked waterfowl for her.
- I cooked waterfowl belonging to her.
- I created the plaster duck she owns.
- I caused her to quickly lower her head or body.
- I waved my magic wand and turned her into undifferentiated waterfowl.

#### Paraphrase

Example: How long is the X river?

- The Mississippi River is 3,734 km (2,320 mi) long.
- ...is a short river, some 4.5 miles (7.2 km) in length
- The total length of the river is 2,145 kilometers.
- ... at the estimated length of 5,464 km (3,395 mi)...
- ... has a meander length of 444 miles (715 km)...
- ... Bali's longest river, measuring approximately 75 kilometers from source to mouth.
- The ... mainstem is 2.75 miles (4.43 km) long although total distance from headwater source tributaries to the sea is 14 miles (23 km).

### Deep Learning (DL) leads to a paradigm shift in NLP: from symbolic to neural approaches

#### Traditional symbolic approaches

- Discrete, symbolic space
- Human comprehensible
  - easy to debug
- Computationally inefficient
  - Sensitive to ambiguity/paraphrase
  - Cascaded models prone to error propagation and require careful feature engineering





#### Neural approaches

- Continuous, neural space
- Human incomprehensible
  - hard to debug
- Computationally efficient
  - Robust to ambiguity/paraphrase
  - E2E learning leads to better performance and simplified systems



"film", "award" film-genre/films-in-this-genre film/cinematography cinematographer/film award-honor/honored-for netflix-title/netflix-genres director/film award-honor/honored-for

### E2E approaches based on DL

#### Discrete, symbolic space

- Human comprehensible
- Input: *x*
- Output: *y*







## The powers of deep learning

- 1. End-to-end Learning
  - Simplifies systems, reduces effort for feature engineering and localization
- 2. Strong Representation Power
  - Due to novel DNN architectures and learning algorithms; leads to high accuracy in many tasks
- 3. Semantic Representation Learning
  - Leads to a paradigm shift in NLP/IR: from symbolic to neural computation
- 4. New Applications and Experience
  - E.g., link language to real-world signals such as images and machine state
- 5. Deep Reinforcement Learning
  - Makes it possible to build intelligent agents for real-world applications such as goaloriented dialogue

### State of the art results on NLP application-level tasks

Task	Test set	Metric	Best non- neural	Best neural	Source
Machine Translation	Enu-deu newstest16	BLEU	31.4	34.8	http://matrix.statmt.org
	Deu-enu newstest16	BLEU	35.9	39.9	http://matrix.statmt.org
Sentiment Analysis	Stanford sentiment bank	5-class Accuracy	71.0	80.7	Socher+ 13
Question Answering	WebQuestions test set	F1	39.9	52.5	<u>Yih+ 15</u>
Entity Linking	Bing Query Entity Linking set	AUC	72.3	78.2	<u>Gao+ 14b</u>
Image Captioning	COCO 2015 challenge	Turing test pass%	25.5	32.2	<u>Fang+ 15</u>
Sentence compression	Google 10K dataset	F1	0.75	0.82	<u>Fillipova+ 15</u>
Response Generation	Sordoni dataset	BLEU-4	3.98	5.82	<u>Li+ 16a</u>

### Deep Semantic Similarity Model (DSSM)

- Compute semantic similarity between two text strings X and Y
  - Map X and Y to feature vectors in a latent semantic space via deep neural net
  - Compute the cosine similarity between the feature vectors

Tasks	X	Υ	Ref
Web search	Search query	Web document	Huang+ 13; Shen+ 14; Palangi+ 16
Entity linking	Entity mention and context	Entity and its corresponding page	<u>Gao+ 14b</u>
Online recommendation	Doc in reading	Interesting things / other docs	<u>Gao+ 14b</u>
Image captioning	Image	Text	<u>Fang+ 15</u>
Machine translation	Sentence in language A	Translations in language B	<u>Gao+ 14a</u>
Question answering	Question	Answer	<u>Yih+ 15</u>

### 3. Semantic Representation Learning



### DSSM for entity linking

ray of light

#### The Einstein Theory of Relativity

(1) The perihelion of Mercury shows a discrepancy which has long puzzled astronomers. This discrepancy is fully accounted for by Einstein. At the time when he published his theory, this was its only experimental verification.

(2) Modern physicists were willing to suppose that light might be subject to gravitation—i.e., that a ray of light passing near a great mass like the sun might be deflected to the extent to which a particle moving with the same velocity would be deflected according to the orthodox theory of gravitation. But Einstein's theory required that the light should be deflected just twice as much as this. The matter could only be tested during an eclipse among a number of bright stars. Fortunately a peculiarly favourable eclipse occurred last year. The results of the observations

#### Ray of Light (Experiment)



#### Ray of Light (Song)



Ray of Light is the seventh studio album by American singersongwriter Madonna, released on March 3,

 1998 by Maverick Records. After giving birth to her daughter Lourdes, Madonna started working on her new album with producers Babyface, Patrick Leonard an...

 Release date
 Mar 3, 1998

 Artist
 Madonna

 Awards
 Grammy Award for B...



### DSSM: Compute Similarity in Semantic Space

Relevance measured by cosine similarity

Word sequence

 $X_t$ 



**Learning:** maximize the similarity between X (source) and Y (target)

**Representation:** use DNN to extract abstract semantic features, f or g is a

- Multi-Layer Perceptron (MLP) if text is a bag of words [<u>Huang+ 13</u>]
- Convolutional Neural Network (CNN) if text is a bag of chunks [<u>Shen+ 14</u>]
- Recurrent Neural Network (RNN) if text is a sequence of words [<u>Palangi+ 16</u>]

### DSSM: Compute Similarity in Semantic Space

Relevance measured by cosine similarity

Semantic layer	h
Max pooling layer	v
Convolutional layer	$C_t$
Word hashing layer	$f_t$
Word sequence	$X_t$



**Learning:** maximize the similarity between X (source) and Y (target)

**Representation:** use DNN to extract abstract semantic representations

**Convolutional and Max-pooling layer:** identify key words/concepts in X and Y

**Word hashing:** use sub-word unit (e.g., letter *n*-gram) as raw input to handle very large vocabulary

### Learning DSSM from Labeled X-Y Pairs



• Map X and Y into the same semantic space via deep neural net

### Learning DSSM from Labeled X-Y Pairs



• Positive Y are closer to X than negative Y in that space

### Learning DSSM from Labeled X-Y Pairs

- Consider a query X and two docs  $Y^+$  and  $Y^-$ 
  - Assume  $Y^+$  is more relevant than  $Y^-$  with respect to X
- $sim_{\theta}(X, Y)$  is the cosine similarity of X and Y in semantic space, mapped by DSSM parameterized by  $\theta$

• 
$$\Delta = \operatorname{sim}_{\theta}(X, Y^+) - \operatorname{sim}_{\theta}(X, Y^-)$$

- We want to maximize  $\Delta$
- $Loss(\Delta; \boldsymbol{\theta}) = \log(1 + \exp(-\gamma \Delta))$
- Optimize θ using mini-batch SGD on GPU



### New Applications and Experience

Neural approaches allow language models to be grounded in the world, i.e., link language to real-world signals such as images, machine state, sensor data from biomedical devices.



Output of a neural conversation model trained on 250K Twitter conversations sparked by a tweeted photo

### Social Bots [MSR Data-Driven Conversation]

- The success of Xiaolce (小冰)
- Problem setting and evaluation
  - Maximize the user engagement by automatically generating
  - enjoyable and useful conversations
- Learning a neural conversation engine
  - A data driven engine trained on social chitchat data [Sordoni+ 15; Li+ 16a]
  - Persona based models and speaker-role based models [<u>Li+ 16b</u>; Luan+ 17]
  - Image-grounded models [Mostafazadeh+ 17]
  - Knowledge-grounded models [Ghazvininejad+ 17]





### Outline

- Part 1: The transition of NLP to neural approaches
- Part 2: Neural machine reading models for question answering
  - MindNet: a case study of symbolic approaches
  - Neural approaches to MRC and QA
  - ReasoNet: a case study of neural approaches
  - Ongoing research: visualize the reasoning process in neural space
- Part 3: Deep reinforcement learning for task-completion dialogue

### Question Answering (QA) on Knowledge Base



#### Large-scale knowledge graphs

- Properties of billions of entities
- Plus relations among them

An QA Example:

**Question:** what is Obama's citizenship?

- Query parsing: (Obama, Citizenship,?)
- Identify and infer over relevant subgraphs: (Obama, BornIn, Hawaii) (Hawaii, PartOf, USA)
- correlating semantically relevant relations: BornIn ~ Citizenship

#### Answer: USA

### Symbolic approaches to QA: production system

https://en.wikipedia.org/wiki/Production\_system\_(computer\_science)

- Production rules
  - condition—action pairs
  - Represent (world) knowledge as a graph
- Working memory
  - Contains a description of the current state of the world in a reasoning process
- Recognizer-act controller
  - Update working memory by searching and firing a production rule
- A case study: MSR MindNet [Dolan+ 93; <u>Richardson+ 98</u>]

### Case study of Question Answering with MindNet

- Build a MindNet graph from:
  - Text of dictionaries
  - Target corpus, e.g. an encyclopedia (Encarta 98)
- Build a dependency graph from query
- Model QA as a graph matching procedure
  - Heuristic fuzzy matching for synonyms, named entities, wh-words, etc.
  - Some common sense reasoning (e.g. dates, math)
- Generate answer string from matched subgraph
  - Including well-formed answers that didn't occur in original corpus



### Fuzzy Match against MindNet

Input LF:



Who assassinated Abraham Lincoln?

American actor <u>John Wilkes Booth</u>, who was a violent backer of the South during the Civil War, <u>shot Abraham Lincoln</u> at Ford's Theater in Washington, D.C., on April 14, 1865.

### Generate output string



"John Wilkes Booth shot Abraham Lincoln"

## Worked beautifully!

- Just not very often...
- What went wrong?
  - One major reason: paraphrase alternations
    - The Mississippi River is 3,734 km (2,320 mi) long.
    - ...is nearly 86 km long...
    - ...is a short river, some 4.5 miles (7.2 km) in length

"How long is the X river?"

- The total length of the river is 2,145 kilometres (1,333 mi).
- ... at the estimated length of 5,464 km (3,395 mi)...
- ... is a 25-mile (40 km) tributary of ...
- ... has a meander length of 444 miles (715 km)...
- ... Bali's longest river, measuring approximately 75 kilometers from source to mouth.
- The ... mainstem is 2.75 miles (4.43 km) long although total distance from headwater source tributaries to the sea is 14 miles (23 km).

#### Symbolic Space

- **Knowledge Representation** 
  - *Explicitly* store a BIG but incomplete knowledge graph (KG)
  - Words, relations, templates
  - High-dim, discrete, sparse vectors
- Inference
  - Slow on a big KG
  - Keyword/template matching is sensitive to paraphrase alternations
- Human comprehensible but not computationally efficient

#### Squire Trelawney, Dr. Livesey, nails, and the sabre cut acro Is a poultry and the rest of these gentleme one cheek, a dirty, livid white. I having asked me to write down remember him looking round the the whole particulars about Treas-ure Island, from the beginning cover and whistling to himself as he did so, and then breaking out to the end, keeping nothing back but the bearings of the island, and that only because there is still Fifteen men on the dead man's Fifteen men on the des treasure not yet lifted, I take up my pen in the year of grace 17-and go back to the time when my voice that seemed to have been father kept the Admiral Benbow tuned and broken at the capstan inn and the brown old seaman with the sabre cut first took up his bars. Then he rapped on the door with a bit of stick like a handspike lodging under our roof. that he carried, and when my fa I remember him as if it were ther appeared, called roughly for a glass of rum. This, when it was esterday, as he came plodding to the inn door, his sea-chest brought to him, he drank slowly, following behind him in a handlike a connoisseur, lingering on the taste and still looking about him at the cliffs and up at our barrow; a tall, strong, heavy, nut-brown man, his tarry pigtail falling over the shoulder of his simboard. I led blue coat, his hands ragged d scarred, with black, broken at length; 'and a pleasant sittyated



#### **Neural Space**

- **Knowledge Representation** 
  - Implicitly store entities and structure of KG in a *compact* way that is more generalizable
  - Semantic concepts/classes
  - Low-dim, cont., dense vectors shaped by KG
- Inference
  - Fast on compact memory
  - Semantic matching is robust to paraphrase alternations
- **Computationally efficient but not human** comprehensible yet



"film", "award" film-genre/films-in-this-genre film/cinematography cinematographer/film award-honor/honored-for netflix-title/netflix-genres director/film award-honor/honored-for

### From symbolic to neural computation



### Case study: ReasoNet with Shared Memory



- Production Rules → Shared memory encodes task-specific knowledge
  - Long-term memory: encode KB for answering all questions in QA on KB
  - Short-term memory: encode the passage(s) which contains the answer of a question in QA on Text
- Working memory → Hidden state S<sub>t</sub> Contains a description of the current state of the world in a reasoning process
- Recognizer-act controller → Search controller performs multi-step inference to update S<sub>t</sub> of a question using knowledge in shared memory
- Input/output modules are task-specific

### KB relation paths in symbolic vs. neural spaces



### Search controller for KB QA



### Joint learning of Shared Memory and Search Controller



### Joint learning of Shared Memory and Search Controller



# Shared Memory: long-term memory to store learned knowledge, like human brain

- Knowledge is learned via performing tasks, e.g., update memory to answer new questions
- New knowledge is *implicitly* stored in memory cells via gradient update
- Semantically relevant relations/entities can be compactly represented using similar vectors.



#### The Knowledge Base Question Answering Results on WN18 and FB15K

Model	Additional Information	WN1	8	FB15k	
		Hits@10(%)	MR	Hits@10(%)	MR
SE (Bordes et al., 2011)	NO	80.5	985	39.8	162
Unstructured (Bordes et al., 2014)	NO	38.2	304	6.3	979
TransE (Bordes et al., 2013)	NO	89.2	251	47.1	125
TransH (Wang et al., 2014)	NO	86.7	303	64.4	87
TransR (Lin et al., 2015b)	NO	92.0	225	68.7	77
CTransR (Lin et al., 2015b)	NO	92.3	218	70.2	75
KG2E (He et al., 2015)	NO	93.2	348	74.0	59
TransD (Ji et al., 2015)	NO	92.2	212	77.3	91
TATEC (García-Durán et al., 2015)	NO	-	-	76.7	58
NTN (Socher et al., 2013)	NO	66.1	-	41.4	-
DISTMULT (Yang et al., 2014)	NO	94.2	-	57.7	-
STransE (Nguyen et al., 2016)	NO	94.7 (93)	244 (206)	79.7	69
RTransE (García-Durán et al., 2015)	Path	-	-	76.2	50
PTransE (Lin et al., 2015a)	Path	-	-	84.6	58
NLFeat (Toutanova et al., 2015)	Node + Link Features	94.3	-	87.0	-
Random Walk (Wei et al., 2016)	Path	94.8	-	74.7	-
ReasoNet (Shen+ 16a)	NO	95.3	249	92.7	38

### Visualization of Reasoning in MRC Models

Translate Natural Language to Image
Multi-step Image Editing via Dialogue





"There is a blue vase by the foot of the sofa"



## First Step: Let there be Color, and Shape!

Task 1: Text-guided Image Colorization





Color Image

#### Task 2: Text-guided Image Segmentation



**RGB-depth Image** 

"At the middle of the kitchen lies a blue chair.
The upper cabinet has a black microwave.
There is a black garbage can on the left of the blue chair.
There is a white garbage can on the right of the blue chair.
On its left there is a red fire distinguisher."





Segmentation

"The flower has purple petals with a white stamen"

BW Image

### Multi-Modal ReasoNet

The flower has purple petals with a white stamen



### Outline

- Part 1: The transition of NLP to neural approaches
- Part 2: Neural machine reading models for question answering
- Part 3: Deep reinforcement learning for task-completion dialogue
  - Dialogue as RL
  - Case study 1: InfoBot with end-to-end learning RL
  - Case study 2: Composite task completion bot with Hierarchical RL
  - Ongoing research: subgoal discovery for hierarchical RL

## Multi-turn (goal-oriented) dialogue



4Z

## (Deep) Reinforcement Learning for Dialogue



Application	State	Action	Reward
Task Completion Bots (Movies, Restaurants,)	User input + Context	Dialog act + slot_value	Task success rate # of turns
Info Bots (Q&A bot over KB, Web etc.)	Question + Context	Clarification questions, Answers	Relevance of answer # of turns
Social Bot (Xiaolce)	Conversation history	Response	Engagement(?)

### A user simulator for RL and evaluation



- Robustness: automatic action selection based on uncertainty by RL
- Flexibility: allow user-initiated behaviors
- Reproducibility: a R&D setting that allows consistent comparisons of competing methods

### InfoBot as an interactive search engine

- Problem setting
  - User is looking for a piece of information from one or more tables/KBs
  - System must iteratively ask for user constraints ("slots") to retrieve the answer
- A general rule-based approach
  - Given current beliefs, ask for slot with maximum uncertainty
  - Works well in most cases but,
    - Has no notion of what the user is likely to be looking for or likely to know
    - No principled way to deal with errors/uncertainty in language understanding

### InfoBot as an interactive search engine



Agent

### Deep Reinforcement Learning



Agent

### End-to-End Learning [<u>Dhuwan+17</u>]



——Reinforcement Learning

### **Dual Exploration**

- Agent should explore actions as well as KB outputs
  - Share similarities with RL Neural Turing Machines (NTU)
- Optimizing expected return

$$J( heta) = E_{a \sim \pi, I \sim p_T} \left[ \sum_{h=0}^{H} \gamma^h r_h \right]$$

• via REINFORCE

$$abla_ heta J( heta) = E_{a \sim \pi, I \sim p_{\mathcal{T}}} \left[ \left( 
abla_ heta \log p(I) + \sum_{h=0}^H 
abla_ heta \log \pi(a_h) 
ight) \sum_{k=0}^H \gamma^k r_k 
ight]$$

### Result on IMDB using KB-InfoBot w/ simulated users



Agent	Success Rate	Avg Turns	Avg Reward
Rule-Soft	0.76	3.94	0.83
RL-Hard	0.75	3.07	0.86
RL-Soft	0.80	3.37	0.98
E2E-RL	0.83	3.27	1.10



### Results on real users



# Composite task completion bot with Hierarchical RL [Peng+ 17]



### A hierarchical policy learner





Similar to HAM [Parr & Russell 98] and hierarchical DQN [Kulkarni+ 16]

### Results on simulated and real users



### Subgoal discovery for HRL:



#### divided and conquer

Figure 3: Subgoals for the landmarks problem (Sutton et al., 1999). Though the solution with subgoals may not be optimal, having the subgoals could usually reduce the search space, and potentially accelerate the learning efficiency.

### The 4-room game



Figure 7: Termination probability visualization for the 4-room experiment. Each time the agent travels from the upper-left corner cell to the lower-right corner cell. The visualization shows the termination probabilities of the RNN generative models in the HRL training after the sequence segmentation process. Darker colors mean higher probabilities.

### Summary

- The transition of NLP to neural approaches
- Neural approaches to MRC and QA
  - Knowledge representation and search in neural space
  - A case study: ReasoNet w/ long-term memory
  - Ongoing research: visualize the reasoning process in neural space
  - Learn more at <u>Deep Learning for Machine Reading Comprehension</u>
- An intelligent, human-like, open-domain conversational system
  - Dialogue as RL
  - Case study 1: InfoBot with end-to-end learning RL
  - Case study 2: Composite task completion bot with Hierarchical RL
  - Ongoing research: subgoal discovery for hierarchical RL
  - Learn more at <u>Deep RL for goal-oriented dialogues</u>

### **Contact Information:**

www.microsoft.com/en-us/research/people/jfgao/

### **Collaborators:**

Faisal Ahmed, Chris Brockett, Asli Celikyilmaz, Ming-Wei Chang, Weizhu Chen, Yun-Nung Chen, Li Deng, Bhuwan Dhingra, Bill Dolan, Michel Galley, Marjan Ghazvininejad, Xiaodong He, Po-Sen Huang, Sungjin Lee, Jiwei Li, Lihong Li, Xiujun Li, Zachary Lipton, Xiaodong Liu, Rangan Majumder, Nasrin Mostafazadeh, Baolin Peng, Mir Rosenberg, Yelong Shen, Alessandro Sordoni, Saurabh Tiwary, Lucy Vanderwende, Luan Yi, Scott Yih et al.