30 Questions in Neural Machine Translation

Jia Xu

Graduate Center & Hunter College, CUNY

MT Tutorial @ JSALT'19

³Άνδρα μοι ἔννεπε, Μοῦσα, πολύτροπον... Odyssey, Homer, 700 BC

^{*}Άνδρα μοι ἔννεπε, Μοῦσα, πολύτροπον... Odyssey, Homer, 700 BC



ancient Greek

³Άνδρα μοι ἔννεπε, Μοῦσα, πολύτροπον... Odyssey, Homer, 700 BC



ancient Greek

²Άνδρα μοι ἔννεπε, Μοῦσα, πολύτροπον... Odyssey, Homer, 700 BC

Virum mihi, Camena, insece versutum...



ancient Greek



Latin



2000 year old translation of <u>Sing in me, Muse, and through me tell the story of the man of twists and turns...</u>



2000 year old translation of Sing in me, Muse, and through me tell the story of the man of twists and turns...

translating between human languages is an old problem





































nowadays machine translation is applied in society, science, arts, commerce and finance, literature, military, ...







machine translation:

automatically translate from one human language to another

machine translation evaluations & resources

- competitions with data resources and baseline platforms, e.g.
 - NIST: 2002- [https://www.nist.gov/itl/iad/mig/openmt15-evaluation]
 - WMT: 2006- [http://www.statmt.org]
 - IWSLT: 2004- [https://workshop2019.iwslt.org]
- projects, e.g.
 - GALE, TC-Star, EuroMatrix, BOLT, and more and more
- datasets, e.g.
 - LDC: [https://www.ldc.upenn.edu]

machine translation development



year

machine translation development



year

1934	1954	1966 1968	1982	1993	2003	2005	2016	2019



itary use			,				> 500	M commerci	al users
rule	-based M ⁻	Г	exam	ple-	MT Ł	based of	on machir	ne learnir	ng
direct tran	direct transfer-based interlingua				statist word -based	ical M phrase -basec	F (SMT) syntax based	neural (NM	MT T)
1934	1954	1966 1968	1982	1993	2	.003	2005	2016	2019

Google translates over 100 billion words a day

litary use							> 500	M commerci	al users
rule-	based M	Г	exan	nple-	MT Ł	based c	on machir	ne learnir	ng È
direct trans	direct transfer-based interlingua				statist word -based	ical M7 phrase -based	(SMT) syntax -based	neural (NM	MT T)
1934	1954	1966 1968	1982	1993	2	.003	2005	2016	2019

problem of each approach:

lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretability
--------------------	-----------------------	------------------------	--------------------	-----------------	-----------------------------	------------------------------	--------------------------------

Google translates over 100 billion words a day

features

interpretability

milita	ary use		、				> 500 M	commercia	lusers
	rul	e-based M	T	example-	MT	based on	machine	learnin	g
	direct tra	insfer-based i	nterlingua	based MT	statist word -based	cical MT (phrase -based	SMT) syntax -based	neural (NMT	MT -)
	934	1954	1966 1968	1982 1993	3	2003 20	05	2016	2019
orol	olem [•] of e	ach approa	ch:						
	lack	language	difficult	lack	lack	need syntax	human defined	robus	tness

generality

context

structure

dependent

to define

generality



milita	ry use				,		> 500 M	commercial users
	rul	e-based M ⁻	Г	example-	MT	based on	machine	elearning
	direct transfer-based interlingua			based MT	statist word -based	cical MT (phrase -based	SMT) syntax -based	neural MT (NMT)
	934	1954	1966-1968	1982 1993	3	2003 20	05	2016 2019
prot	olem of ea	ach approa	ch:					
	lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretabili

milita	iry use				, 		> 500 M	commercial users	
	rul	e-based M ⁻	Г	example-	MT	based on	machine	e learning	
	direct tra	nsfer-based i	nterlingua	based MT	statist word -based	tical MT (phrase -based	SMT) syntax -based	neural MT (NMT)	
1934 1954 1966 1968 1982 1993 2003 2005 2016 2019									
	lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretabilit	

milita	ary use				, <u></u>		> 500 M	commercial users			
	rul	e-based M	Г	example-	MT	based on	machine	elearning			
	direct tra	nsfer-based i	nterlingua	based MT	statist word -based	tical MT (phrase -based	SMT) syntax -based	neural MT (NMT)			
- I Drot	1934 1954 1966 1968 1982 1993 2003 2005 2016 2019 roblem of each approach:										
	lack generality	lack language difficult nerality dependent to define		lack generality	lack context	need syntax structure	human defined features	robustness interpretability			

milita	ary use				, <u></u>		> 500 M	commercial users
	rul	e-based M	Т	example-	MT	based on	machine	elearning
	direct tra	insfer-based i	nterlingua	based MT	statist word -based	cical MT (phrase -based	SMT) syntax -based	neural MT (NMT)
orol	934 olem of ea	1954 ach approa	1966 1968 ach:	1982 1993	3	2003 20	05	2016 2019
	lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretabili

milita	iry use				,		> 500 M	commercial users	
	rul	e-based M	Т	example-	MT	based on	machine	e learning	
	direct tra	Insfer-based i	nterlingua	based MT	statist word -based	cical MT (phrase -based	SMT) syntax -based	neural MT (NMT)	
- I Droł	934 olem of ea	1954 ach approa	1966 1968 .ch:	1982 1993	3	2003 20	05	2016 2019	
	lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretabilit	

milita	ary use				,		> 500 M	commercia	al users
	rul	e-based M	T	example-	MT	based on	machine	learnir	g
	direct tra	ansfer-based i	nterlingua	based MT	statist word -based	statistical MT (SMT) word phrase syntax -based -based -based			МТ Г)
- 	934	1954	1966 1968	1982 1993	3	2003 20	05	2016	2019
prot	Diem of e	acn approa							
	lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robus interpre	tness etability

Google translates over 100 billion words a day

litary use							> 500	M commerci	al users
rule-	based M	Г	exan	nple-	MT Ł	based c	on machir	ne learnir	ng È
direct trans	direct transfer-based interlingua				statist word -based	ical M7 phrase -based	(SMT) syntax -based	neural (NM	MT T)
1934	1954	1966 1968	1982	1993	2	.003	2005	2016	2019

problem of each approach:

lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretability
--------------------	-----------------------	------------------------	--------------------	-----------------	-----------------------------	------------------------------	--------------------------------

Google translates over 100 billion words a day

ilitary use					> 500	M commerci	al users	
rule-	based MT	exan	nple-	MT based on machine learning				
direct trans	direct transfer-based interlingua			statistical word phr -based -ba	neural MT (NMT)			
1934	1954 1966 196	8 1982	1993	2003	2005	2016	2019	

problem of each approach:

lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretability
--------------------	-----------------------	------------------------	--------------------	-----------------	-----------------------------	------------------------------	--------------------------------

Question #1: how to enhance NMT robustness?

Google translates over 100 billion words a day

litary use							> 500	M commerci	al users
rule-	rule-based MT				MT based on machine learning				
direct trans	direct transfer-based interlingua			sed IT	statistical MT (SMT) word phrase syntax -based -based -based			neural MT (NMT)	
1934	1954	1966 1968	1982	1993	2	.003	2005	2016	2019

problem of each approach:

lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretability
--------------------	-----------------------	------------------------	--------------------	-----------------	-----------------------------	------------------------------	--------------------------------
Google translates over 100 billion words a day

litary use							> 500	M commerci	al users				
rule-	rule-based MT				rule-based MT example-				MT Ł	based c	on machir	ne learnir	ng È
direct trans	sfer-based in	nterlingua	based MT		statist word -based	ical M7 phrase -based	(SMT) syntax -based	neural (NM	MT T)				
1934	1954	1966 1968	1982	1993	2	.003	2005	2016	2019				

problem of each approach:

lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretability
--------------------	-----------------------	------------------------	--------------------	-----------------	-----------------------------	------------------------------	--------------------------------

Google translates over 100 billion words a day

litary use							> 500 N	1 commercia	al users				
rule	rule-based MT				rule-based MT example-				MT b	ased o	n machir	ne learnir	g
direct tran	isfer-based in	nterlingua	based MT		statist word -based	ical MT phrase -based	(SMT) syntax -based	neural (NM	МТ Г)				
1934	1954	1966 1968	1982	1993	2	003	2005	2016	2019				

problem of each approach:

lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretability
--------------------	-----------------------	------------------------	--------------------	-----------------	-----------------------------	------------------------------	--------------------------------

Question #2: how to increase interpretability?

Google translates over 100 billion words a day

litary use							> 500	M commerci	al users				
rule-	rule-based MT				rule-based MT example-				MT Ł	based c	on machir	ne learnir	ng È
direct trans	sfer-based in	nterlingua	based MT		statist word -based	ical M7 phrase -based	(SMT) syntax -based	neural (NM	MT T)				
1934	1954	1966 1968	1982	1993	2	.003	2005	2016	2019				

problem of each approach:

lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretability
--------------------	-----------------------	------------------------	--------------------	-----------------	-----------------------------	------------------------------	--------------------------------

Google translates over 100 billion words a day

litary use							> 500	M commerci	al users				
rule-	rule-based MT				rule-based MT example-				MT Ł	based c	on machir	ne learnir	ng È
direct trans	sfer-based in	nterlingua	based MT		statist word -based	ical M7 phrase -based	(SMT) syntax -based	neural (NM	MT T)				
1934	1954	1966 1968	1982	1993	2	.003	2005	2016	2019				

problem of each approach:

lack generality	language dependent	difficult to define	lack generality	lack context	need syntax structure	human defined features	robustness interpretability
--------------------	-----------------------	------------------------	--------------------	-----------------	-----------------------------	------------------------------	--------------------------------

> 500 M commercial users
MT based on machine learning
statistical MT (SMT) neural MT
word phrase syntax
-based -based -based

 > 500 M commercial users
 MT based on machine learning
 statistical MT (SMT) neural MT word phrase syntax (NMT)
 -based -based -based

noisy channel model

> 500 M commercial users

 MT based on machine learning

 statistical MT (SMT)
 neural MT

 word
 phrase
 syntax

 -based
 -based
 -based

noisy channel model

[Shannon, 1948] Information Theory

 > 500 M commercial users
 MT based on machine learning
 statistical MT (SMT) neural MT (NMT)
 word phrase syntax -based -based

noisy channel model

[Shannon, 1948] Information Theory

- input: source sentence (observation) f
- output: target sentence (decision) e
- Bayes decision rule

 > 500 M commercial users
 MT based on machine learning
 statistical MT (SMT) neural MT (NMT)
 word phrase syntax -based -based

noisy channel model

[Shannon, 1948] Information Theory

- input: source sentence (observation) f
- output: target sentence (decision) e
- Bayes decision rule

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$

 > 500 M commercial users
 MT based on machine learning
 statistical MT (SMT) neural MT (NMT)
 word phrase syntax -based -based

noisy channel model

[Shannon, 1948] Information Theory

- input: source sentence (observation) f
- output: target sentence (decision) e
- Bayes decision rule

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$

$$\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$$
$$= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$$















machine translation











introduced in the previous tutorial







R









word pair semantic distance preserved



word pair semantic distance preserved



Question #3: better text embedding/representation? BERT, ELMO, GloVec, FastText, ...



word pair semantic distance preserved



word pair semantic distance preserved














low distortion embedding space



Can we learn only from distance in low distortion embedding space?



Can we learn only from distance in low distortion embedding space?

No, in NLP we need language models



Can we learn only from distance in low distortion embedding space?

No, in NLP we need language models

setting: multi-class classification + metric structure + 2 experts we cannot combine them to a better one by querying them [PYX, 16]

machine translation components

what are

language models

machine translation components

what are language models



machine translation components







$$\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$$
$$= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$$



```
\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}
= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}
                                        language
                                          model
```

evaluating language models: perplexity

let e_i be a word in the document that contains N words

PPL (perplexity) is measured as

$$\log PPL = -\frac{1}{N} \sum_{i=1}^{N} \log P(e_i | h_i)$$

 h_i is the history of word e_i

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$

- $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$
- let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ $= Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule \dot{i} estimated as $\approx P(e_i|e_1,\cdots,e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ $= Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule \dot{i} estimated as $\approx P(e_i|e_1,\cdots,e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1, \cdots, e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1, \cdots, e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule \dot{i} estimated as statistical $\rightarrow \approx P(e_i | \cdot, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule \dot{i} estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

previously

 $\hat{e} = \operatorname{argmax}_{e} \{ Pr(e|f) \}$ $= \operatorname{argmax}_{e} \{ Pr(f|e) \cdot Pr(e) \}$ let e_1^I be a sentence of a sequence of words $e_1, e_2 \cdots e_I$ $Pr(e_1^I) = Pr(e_1, e_2, \cdots, e_I)$ = $Pr(e_i|e_1,\cdots,e_{i-1})$ chain rule i estimated as statistical $\rightarrow \approx P(e_i | e_1, \cdots, e_{i-n+1})$ relative frequency

long-term memory for language models

long-term memory for language models

- n-gram LM making prediction on fixed windows
- past n words my not be sufficient to capture the context

long-term memory for language models

- n-gram LM making prediction on fixed windows
- past n words my not be sufficient to capture the context

RNNs are capable of conditioning the model on all previous words

neural language model with RNN



- x_t : input word vector at time t
- W: weights matrix to condition t
- h_{t-1} : output of the non-linear function at the previous time step
- σ : the non-linearity function
Gated Recurrent Units (GRU)

- problem of vanishing gradients makes RNNs hard to train for long-term dependency
- use more complex units for activation



another type of complex activation unit



Question #4: contextual memory in language model



another type of complex activation unit



another type of complex activation unit



Question #5: affective neuron activation function



another type of complex activation unit



another type of complex activation unit









NMT I: sequence-to-sequence with RNN

[Sutskever, 93]



[Sutskever, 93]

$$\frac{1}{4} \left[-P("W") - P("X") - P("Y") - P("Z") \right]$$



[Sutskever, 93]



[Sutskever, 93]



[Sutskever, 93]



[Sutskever, 93]



back propagation operates "end-to-end"

maximize the log probability of a correct translation given the source sentence

Question #6: better training criterion? Maximum Likelihood, squared error, MAP, cross-entropy, minimum risk, ..

[Sutskever, 93]



[Sutskever, 93]



[Sutskever, 93]



back propagation operates "end-to-end"

maximize the log probability of a correct translation given the source sentence

Question #7: better training algorithm? error back propagation, contrastive estimation, ...

[Sutskever, 93]



[Sutskever, 93]



NMT II: encoder & decoder with attention

[Luong et.al., 15]



NMT III: multiple models with CNN

[Gehring, 16]



NMT III: multiple models with CNN

[Gehring, 16]











greedy search



method: take most probable word in each step problem: no way to undo decisions

• W____

- WX____
- WXZ____ (no way back!)

exhaustive search



ideally: find a translation that maximize

$$P(y_1, \cdots, y_{T'} | x_1, \cdots, x_T) = \prod_{t=1}^{T'} P(y_t | x, y_1 \cdots, y_{t-1})$$

thod: compute all possible sequences y

problem: expensive

me

each step tracking V (vocabulary) words complexity $O(V^T)$

beam search



$$\operatorname{score}(y_1, \cdots, y_t) = \log P_{LM}(y_1, \cdots, y_t | x)$$
$$= \sum_{i=1}^t \log P_{LM}(y_i | y_1 \cdots, y_{i-1}, x)$$

method: on each search step, keep track of the k most probable (higher score) partial translations problem: no guarantee for optimal solution

efficient!
























Question #8: more efficient or controlled search? binary NMT, constraint

ensemble



random initialization or outputs from different iterations









human evaluation is expensive, develop automatic evaluation criteria hypothesis: ``Montreal is a city .'' reference: ``Montreal , a giant playground .''

- method to measure error rates:
 - edit distance: insertion, deletion, substitution
 - word error rate: normalized edit distance
 - HTER: Human Targeted Translation Error Rate

human evaluation is expensive, develop automatic evaluation criteria

hypothesis: ``Montreal is a city .'' reference: ``Montreal , a giant playground .'' substitution#=1

- method to measure error rates:
 - edit distance: insertion, deletion, substitution
 - word error rate: normalized edit distance
 - HTER: Human Targeted Translation Error Rate

human evaluation is expensive, develop automatic evaluation criteria

hypothesis: ``Montreal is a city .'' reference: ``Montreal , a giant playground .'' substitution#=1+1

- method to measure error rates:
 - edit distance: insertion, deletion, substitution
 - word error rate: normalized edit distance
 - HTER: Human Targeted Translation Error Rate

human evaluation is expensive, develop automatic evaluation criteria

hypothesis: ``Montreal is a city ." reference: ``Montreal , a giant playground ." substitution#=I+I deletion#=I

- method to measure error rates:
 - edit distance: insertion, deletion, substitution
 - word error rate: normalized edit distance
 - HTER: Human Targeted Translation Error Rate

human evaluation is expensive, develop automatic evaluation criteria

hypothesis: ``Montreal is a city ." reference: ``Montreal , a giant <u>playground</u> ." substitution#=1+1 deletion#=1 insertion#=0

- method to measure error rates:
 - edit distance: insertion, deletion, substitution
 - word error rate: normalized edit distance
 - HTER: Human Targeted Translation Error Rate

human evaluation is expensive, develop automatic evaluation criteria

hypothesis: ``Montreal is a city ." reference: ``Montreal , a giant <u>playground</u> ." substitution#=I+I deletion#=I insertion#=0 edit distance#=I+I+I=3

- method to measure error rates:
 - edit distance: insertion, deletion, substitution
 - word error rate: normalized edit distance
 - HTER: Human Targeted Translation Error Rate

human evaluation is expensive, develop automatic evaluation criteria

hypothesis: ``Montreal is a city ." reference: ``Montreal , a giant <u>playground</u> ." substitution#=1+1 deletion#=1 insertion#=0 edit distance#=1+1+1=3

- method to measure error rates:
 - edit distance: insertion, deletion, substitution
 - word error rate: normalized edit distance 3/6=0.5
 - HTER: Human Targeted Translation Error Rate

BLEU (Bilingual Evaluation Understudy) hypothesis: ``Montreal is a city .'' reference: ``Montreal , a giant playground .''

- method to measure accuracy
- most well-cited evaluation criterion so far
 - precision of I-gram through 4-gram
 - brevity penalty
- mainly relies on n-gram coverage

BLEU (Bilingual Evaluation Understudy)

hypothesis: ``Montreal is a city .'' reference: ``Montreal , a giant playground .''

- method to measure accuracy
- most well-cited evaluation criterion so far
 - precision of I-gram through 4-gram
 - brevity penalty
- mainly relies on n-gram coverage

BLEU (Bilingual Evaluation Understudy)

hypothesis: ``Montreal is a city .'' reference: ``Montreal , a giant playground .''

- method to measure accuracy
- most well-cited evaluation criterion so far
 - precision of I-gram through 4-gram
 - brevity penalty
- mainly relies on n-gram coverage

BLEU (Bilingual Evaluation Understudy)

hypothesis: ``Montreal is a city . '' reference: ``Montreal , a giant playground . ''

- method to measure accuracy
- most well-cited evaluation criterion so far
 - precision of I-gram through 4-gram
 - brevity penalty
- mainly relies on n-gram coverage

BLEU (Bilingual Evaluation Understudy)

hypothesis: ``Montreal is a city . " reference: ``Montreal , a giant playground . "

- method to measure accuracy
- most well-cited evaluation criterion so far
 - precision of I-gram through 4-gram
 - brevity penalty
- mainly relies on n-gram coverage

l-gram#=3 2-gram#=0 3-gram#=0 4-gram#=0

BLEU (Bilingual Evaluation Understudy)

hypothesis: ``Montreal is a city . " reference: ``Montreal , a giant playground . "

- method to measure accuracy
- most well-cited evaluation criterion so far
 - precision of I-gram through 4-gram
 - brevity penalty
- mainly relies on n-gram coverage

Question #9: higher correlation with human judgement? rich literature

BLEU (Bilingual Evaluation Understudy)

hypothesis: ``Montreal is a city . " reference: ``Montreal , a giant playground . "

- method to measure accuracy
- most well-cited evaluation criterion so far
 - precision of I-gram through 4-gram
 - brevity penalty



Question #10: better quality estimation?

Question #9: higher correlation with human judgement? rich literature









- tokenization: separate words from punctuation marks, typically based on rules
- text normalization
- word segmentation
- sentence segmentation
- domain classification

- tokenization: separate words from punctuation marks, typically based on rules
- text normalization

Question #11: text normalization

- word segmentation
- sentence segmentation
- domain classification

- tokenization: separate words from punctuation marks, ullettypically based on rules
- text normalization ightarrow
- ullet



word segmentation Question #12: better subword?

- sentence segmentation ullet
- domain classification ullet

- tokenization: separate words from punctuation marks, typically based on rules
- text normalization Question #11: text normalization
 word segmentation Question #12: better subword?
 sentence segmentation Question #13: monolingual and bilingual sentence segmentation
 - domain classification

- tokenization: separate words from punctuation marks, typically based on rules
- text normalization Question #11: text normalization
 word segmentation Question #12: better subword?
 sentence segmentation Question #13: monolingual and bilingual sentence segmentation
 domain classification Question #14: domain adaptation
text normalization

formal text: ``are you coming to the class tomorrow?" informal text: ``r u cuming 2 class tomr?"

- bad translation: style, domain change, noise e.g. mis-spelling
- goal: translate different lexical variations
 - add noise to training: [Michell, et.al., 19]
 - word clustering [Khan et. al., 19]



明天来上课吗?

- possible segmentation boundaries: 2^{k-1}
 - k: number of characters
- n-gram approach:

n-gram approach

$$\hat{f}_1^J(c_1^K) = \operatorname*{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}}^{k_j} | c_{k_{j-2-n+1}}^{k_{j-1-n}}, ..., c_{k_{j-2}+1}^{k_{j-1}})$$

$$= \operatorname{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}+1}^{k_j})$$

明天来上课吗?

- possible segmentation boundaries: 2^{k-1}
 - k: number of characters
- n-gram approach:

n-gram approach

$$\hat{f}_1^J(c_1^K) = \operatorname*{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}}^{k_j} | c_{k_{j-2-n+1}}^{k_{j-1-n}}, ..., c_{k_{j-2}+1}^{k_{j-1}})$$

$$= \operatorname{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}+1}^{k_j})$$

明天来上课吗?

- possible segmentation boundaries: 2^{k-1}
 - k: number of characters
- n-gram approach:

n-gram approach

$$\hat{f}_1^J(c_1^K) = \operatorname*{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}}^{k_j} | c_{k_{j-2-n+1}}^{k_{j-1-n}}, ..., c_{k_{j-2}+1}^{k_{j-1}})$$

$$= \operatorname{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}+1}^{k_j})$$

明天来上课吗?

- possible segmentation boundaries: 2^{k-1}
 - k: number of characters
- n-gram approach:

n-gram approach

$$\hat{f}_1^J(c_1^K) = \operatorname*{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}}^{k_j} | c_{k_{j-2-n+1}}^{k_{j-1-n}}, ..., c_{k_{j-2}+1}^{k_{j-1}})$$

$$= \operatorname{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}+1}^{k_j})$$

明天来上课吗?

- possible segmentation boundaries: 2^{k-1}
 - k: number of characters
- n-gram approach:

n-gram approach

$$\hat{f}_1^J(c_1^K) = \operatorname*{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}}^{k_j} | c_{k_{j-2-n+1}}^{k_{j-1-n}}, ..., c_{k_{j-2}+1}^{k_{j-1}})$$

$$= \operatorname{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}+1}^{k_j})$$

明天来上课吗?

- possible segmentation boundaries: 2^{k-1}
 - k: number of characters
- n-gram approach:

n-gram approach

$$\hat{f}_1^J(c_1^K) = \operatorname*{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}}^{k_j} | c_{k_{j-2-n+1}}^{k_{j-1-n}}, ..., c_{k_{j-2}+1}^{k_{j-1}})$$

$$= \operatorname{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}+1}^{k_j})$$

明天来上课吗?

- possible segmentation boundaries: 2^{k-1}
 - k: number of characters
- n-gram approach:

n-gram approach

$$\hat{f}_1^J(c_1^K) = \operatorname*{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}}^{k_j} | c_{k_{j-2-n+1}}^{k_{j-1-n}}, ..., c_{k_{j-2}+1}^{k_{j-1}})$$

$$= \operatorname{argmax}_{k_1^J,J} \prod_{j=1}^J \Pr(c_{k_{j-1}+1}^{k_j})$$



character based statistical MT

明天来上课吗?

• Gibbs sampling: joint model word alignment and word segmentation $\operatorname{argmax}_{f_1^J} \{ P(f_1^{J\lambda_1} P(e_1^J, b_1 I | f_1^J)^{\lambda_2} P(f_1^J, a_1^J | e_1^I)^{\lambda_3}$



character based statistical MT

明天来上课吗?

• Gibbs sampling: joint model word alignment and word segmentation $\operatorname{argmax}_{f_1^J} \{ P(f_1^{J\lambda_1} P(e_1^J, b_1 I | f_1^J)^{\lambda_2} P(f_1^J, a_1^J | e_1^I)^{\lambda_3} \}$

statistical



character based statistical MT

明天来上课吗?

• Gibbs sampling: joint model word alignment and word segmentation $\operatorname{argmax}_{f_1^J} \{ P(f_1^{J\lambda_1} P(e_1^J, b_1 I | f_1^J)^{\lambda_2} P(f_1^J, a_1^J | e_1^I)^{\lambda_3} \}$

statistical



Question #15: what can we borrow from statistical MT?

character based neural MT

- integrate with neural network framework
 - [Ling, et.al., 15], [Cherry, et.al. 18], [Lee, et.al., 18]...



- efficient and affective, NMT independent
- [Gage, 1994], [Sennrich, et.al., 16]
 - aaabdaaabac
 - ZabdZabac (Z=aa)
 - ZYdZYac (Y=ab; Z=aa)
 - XdXac (X=ZY;Y=ab; Z=aa)

- efficient and affective, NMT independent
- [Gage, 1994], [Sennrich, et.al., 16]
 - aaabdaaabac
 - ZabdZabac (Z=aa)
 - ZYdZYac (Y=ab; Z=aa)
 - XdXac (X=ZY;Y=ab; Z=aa)

Question #16: better subword e.g. with morphological knowledge?

- efficient and affective, NMT independent
- [Gage, 1994], [Sennrich, et.al., 16]
 - aaabdaaabac
 - ZabdZabac (Z=aa)
 - ZYdZYac (Y=ab; Z=aa)
 - XdXac (X=ZY;Y=ab; Z=aa)

Question #16: better subword e.g. with morphological knowledge?

Question #17: unseen words?

- efficient and affective, NMT independent
- [Gage, 1994], [Sennrich, et.al., 16]
 - aaabdaaabac
 - ZabdZabac (Z=aa)
 - ZYdZYac (Y=ab; Z=aa)
 - XdXac (X=ZY;Y=ab; Z=aa)

Question #16: better subword e.g. with morphological knowledge?

Question #17: unseen words?

Question #18: named entities?

- efficient and affective, NMT independent
- [Gage, 1994], [Sennrich, et.al., 16]
 - aaabdaaabac
 - ZabdZabac (Z=aa)
 - ZYdZYac (Y=ab; Z=aa)
 - XdXac (X=ZY;Y=ab; Z=aa)

















- German to French translation is good
- French to German translation is bad
- use German to French MT system to translate German monolingual data e.g. [Sennrich, et.al., 16]
- add the synthetic data into French to German MT training
 - usually greatly improves translation quality due to data augmentation

- German to French translation is good
- French to German translation is bad
- use German to French MT system to translate German monolingual data e.g. [Sennrich, et.al., 16]
- add the synthetic data into French to German MT training
 - usually greatly improves translation quality due to data augmentation

Question #20: back translation

- German to French translation is good
- French to German translation is bad
- use German to French MT system to translate German monolingual data e.g. [Sennrich, et.al., 16]
- add the synthetic data into French to German MT training
 - usually greatly improves translation quality due to data augmentation

- German to French translation is good
- French to German translation is bad
- use German to French MT system to translate German monolingual data e.g. [Sennrich, et.al., 16]
- add the synthetic data into French to German MT training
 - usually greatly improves translation quality due to data augmentation

- lack of German French parallel training data
- rich data of German English, and French English
- generate German French parallel data using German English and French - English MT systems
- add the synthetic data into French to German MT training
 - usually greatly improves translation quality

- lack of German French parallel training data
- rich data of German English, and French English
- generate German French parallel data using German English and French - English MT systems
- add the synthetic data into French to German MT training
 - usually greatly improves translation quality

Question #21: pivot translation

- lack of German French parallel training data
- rich data of German English, and French English
- generate German French parallel data using German English and French - English MT systems
- add the synthetic data into French to German MT training
 - usually greatly improves translation quality

- lack of German French parallel training data
- rich data of German English, and French English
- generate German French parallel data using German English and French - English MT systems
- add the synthetic data into French to German MT training
 - usually greatly improves translation quality

multiple languages

- observation in WMT'19: adding Hindi English parallel data improves Gujarati - English translation
- how about many other languages? [Firat, et.al., 16], [Johnson, et.al., 17]



multiple languages

- observation in WMT'19: adding Hindi English parallel data improves Gujarati - English translation
- how about many other languages? [Firat, et.al., 16], [Johnson, et.al., 17]



Question #22: multi-lingual and zero resource

multiple languages

- observation in WMT'19: adding Hindi English parallel data improves Gujarati - English translation
- how about many other languages? [Firat, et.al., 16], [Johnson, et.al., 17]


multiple languages

- observation in WMT'19: adding Hindi English parallel data improves Gujarati - English translation
- how about many other languages? [Firat, et.al., 16], [Johnson, et.al., 17]



the concept of interlingua





the concept of interlingua interlingua analysis generation transfer target

direct translation

.....



the concept of interlingua





the concept of interlingua











metaalgorithm

boosting

improves prediction accuracy; non-parallelizable



Natürliche Lebensräume wurden zerstört. Dies ist eine ihrer Hauptaufgaben. Das kann so nicht weitergehen. Die Aussprache ist geschlossen. Natural habitats were destroyed. This is a major task. That cannot continue. That concludes the debate.



Natürliche Lebensräume wurden zerstört.	Natural habitats were destroyed.
Dies ist eine ihrer Hauptaufgaben.	This is a major task.
Das kann so nicht weitergehen.	I hat cannot continue.
Die Aussprache ist geschlossen.	That concludes the debate.



Natürliche Lebensräume wurden zerstört.	Natural habitats were destroyed.
Dies ist eine ihrer Hauptaufgaben.	This is a major task.
Das kann so nicht weitergehen.	l hat cannot continue.
Die Aussprache ist geschlossen.	That concludes the debate.





Natürliche Lebensräume wurden zerstört.	Natural habitats were destroyed.
Dies ist eine ihrer Hauptaufgaben.	This is a major task.
Das kann so nicht weitergehen	Lhat cannot continue
Die Aussprache ist geschlossen.	That concludes the debate.







Natürliche Lebensräume wurden zerstört.	Natural habitats were destroyed.
Dies ist eine ihrer Hauptaufgaben.	This is a major task.
Das kann so nicht weitergehen	Lhat cannot continue
Die Aussprache ist geschlossen.	That concludes the debate.





Natürliche Lebensräume wurden zerstört.	Natural habitats were destroyed.
Dies ist eine ihrer Hauptaufgaben.	This is a major task.
Das kann so nicht weitergehen.	Lhat cannot continue
Die Aussprache ist geschlossen.	That concludes the debate.





Natürliche Lebensräume wurden zerstört.	Natural habitats were destroyed.
Dies ist eine ihrer Hauptaufgaben.	This is a major task.
Das kann so nicht weitergehen.	Lhat cannot continue
Die Aussprache ist geschlossen.	That concludes the debate.







Natürliche Lebensräume wurden zerstört.	Natural habitats were destroyed.
Dies ist eine ihrer Hauptaufgaben.	This is a major task.
Das kann so nicht weitergehen	Lhat cannot continue
Die Aussprache ist geschlossen.	That concludes the debate.





















each sample is a parallel sentence

Naturliche Lebensraume wurden zerstort.	Natural habitats were destroyed.
Dies ist eine ihrer Hauptaufgaben.	This is a major task.
Das kann so nicht weitergehen.	I hat cannot continue.
Die Aussprache ist geschlossen.	That concludes the debate.





























now we are at this point of the execution of Design-Bagging
















bootstrapping with combinatorial design

N=9, m=3, b=6 Question #24: enhance bootstrapping







metaalgorithm





































30 Questions

Question #1: how to enhance NMT robustness?

Question #2: how to increase interpretability?

Question #3: better text embedding/representation? BERT, ELMO, GloVec, FastText, ...

Question #4: contextual memory in language model

Question #5: affective neuron activation function

Question #6: better training criterion? Maximum Likelihood, squared error, MAP, cross-entropy, minimum risk, ..

Question #7: better training algorithm? error back propagation, contrastive estimation, ...

Question #8: more efficient or controlled search? binary NMT, constraint

Question #9: higher correlation with human judgement? rich literature

Question #10: better quality estimation?

Question #11: text normalization

Question #12: better subword?

Question #13: monolingual and bilingual sentence segmentation

Question #14: domain adaptation

Question #15: what can we borrow from statistical MT?

Question #16: better subword e.g. with morphology?

Question #17: unseen words?

Question #18: named entities?

Question #19: higher quality in unsupervised MT?

Question #20: back translation

Question #21: pivot translation

Question #22: multi-lingual and zero resource

Question #23: interlingua exists? [Lu, et.al., 18]

Question #24: syntax in NMT

Question #25: multi-modal in NMT

Question #26: translation retrieval

Question #27: probing task

Question #28: word alignment

Question #29: MT using quantum information

Question #30: error analysis

questions?