Hierarchical MT Training using Max-Violation Perceptron

Kai Zhao Liang Huang City University of New York



Haitao Mi Abe Ittycheriah IBM T.J. Watson Research Center



MT is hard





Discriminative Training for SMT

- discriminative training is not very successful yet in MT
 - although dominant in parsing/tagging
 - can use arbitrary, overlapping, lexicalized features
- most efforts on MT training tune feature weights on the small dev set (~Ik sents) not the training set!
 - as a result can only use ~10 dense features (MERT)
 - or ~I0k rather impoverished features (MIRA/PRO)
- Liang et al ('06) train on the training set but not successful

training set (>100k sentences)



- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron

training example

那人咬了狗_x

the man bit the dog y

(Liang et al '06)

- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron

training example

during online learning...

the man bit the dog y

(Liang et al '06)

- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron

training example







the man bit the dog y

- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron

during online learning...

training example



- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron





during online learning...



- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron



- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron



- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron



- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron



(Liang et al '06)

reward correct penalize wrong

- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron



$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(x, d^*) - \Phi(x, \hat{d})$$
reward penalize
correct wrong

- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron



- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron

- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron

- with latent variables (hidden derivations)
- structured perceptron => latent-variable structured perceptron

- Phrase-based translation suffers from distortion limit
 - can only use a **small** portion of bitext (low forced decoding reachability)
 - translation quality is often slightly worse than hierarchical models

reachability % (in # of words) on FBIS

- Phrase-based translation suffers from distortion limit
 - can only use a **small** portion of bitext (low forced decoding reachability)
 - translation quality is often slightly worse than hierarchical models
- Hiero handles reordering better
 - potentially more sentences to train
 - learn more sparse features

reachability % (in # of words) on FBIS

- Phrase-based translation suffers from distortion limit
 - can only use a **small** portion of bitext (low forced decoding reachability)
 - translation quality is often slightly worse than hierarchical models
- Hiero handles reordering better
 - potentially more sentences to train
 - learn more sparse features
- Challenge: generalize max-violation latent perceptron to hypergraphs?

reachability % (in # of words) on FBIS

7

Reference Bush held talks with Sharon

- Reference derivation rank too low in the beam and gets pruned
- We call it a **Violation** iff. score(viterbi) score(reference) > 0

Fixing Search Error: Max-Violation

- standard perceptron does not guarantee violation
 - w/ pruning, the correct derivation might score higher at the ends
 - called "invalid" update b/c it doesn't fix the search error
- max-violation: update at where the violation is maximum
 - "worst-mistake" in the search
 - learns more and faster
- "violation-fixing perceptron" (Huang et al 2012)

Fixing Search Error: Max-Violation

- standard perceptron does not guarantee violation
 - w/ pruning, the correct derivation might score higher at the ends
 - called "invalid" update b/c it doesn't fix the search error
- max-violation: update at where the violation is maximum
 - "worst-mistake" in the search
 - learns more and faster

10

standard update

(no guarantee!)

Fixing Search Error: Max-Violation

- standard perceptron does not guarantee violation
 - w/ pruning, the correct derivation might score higher at the ends
 - called "invalid" update b/c it doesn't fix the search error
- max-violation: update at where the violation is maximum
 - "worst-mistake" in the search
 - learns more and faster
- "violation-fixing perceptron" (Huang et al 2012)

 Update at where the violation b/w the Viterbi and the reference is maximum (Zhang et al. '13)

 Update at where the violation b/w the Viterbi and the reference is maximum (Zhang et al. '13)

- multiple reference derivations at one node
 - at one node, pick the best reference
 - globally choose max-violated best reference

- multiple reference derivations at one node
 - at one node, pick the best reference
 - globally choose max-violated best reference

- multiple reference derivations at one node
 - at one node, pick the best reference
 - globally choose max-violated best reference

Latent-Variable Max-Violation Perceptron

Experiments Setting

train	dev	test
IWSLT09	IWSLT04	IWSLT05
30k short sentences	16 references	16 references
FBIS	NIST06	NIST08
240k sentences	4 references	4 references

- 18 dense features from cdec
- budgeted sparse features based on Word-Edge features
 - atomic features: C/E boundary words; C boundary characters
 - complex features
 - combination of atomic features within limited budget
- we remove all I-count rules in rule extraction
- trigram LM trained from target side

Reachability on FBIS

% reachability (in # of words)

Contribution of Sparse Features (on IWSLT09)

• BLEU scores on IWSLT09 (16 refs) and FBIS (4 refs)

• BLEU scores on IWSLT09 (16 refs) and FBIS (4 refs)

Conclusion

- A latent-variable violation-fixing perceptron framework for general structured prediction problems with inexact search over hypergraphs
- Compared with PBMT, it can use **more** training sentences.
- Compared with MERT/PRO, it is simpler in theory and practice, and achieves better translations

