

# Improving Joint Training of Inference Networks and Structured Prediction Energy Networks

Lifu Tu<sup>1</sup> Yuanzhe Pang<sup>2</sup> Kevin Gimpel<sup>1</sup>

<sup>1</sup>Toyota Technological Institute at Chicago

<sup>2</sup>New York University

4th Workshop on Structured Prediction for NLP, 2020

# Table of Contents

- 1 Motivation
- 2 Inference Networks (InfNets) [Tu and Gimpel (2018, 2019)]
- 3 Joint Training of SPENs and Two InfNets
- 4 Conclusions

[Finkel et al, 2005]

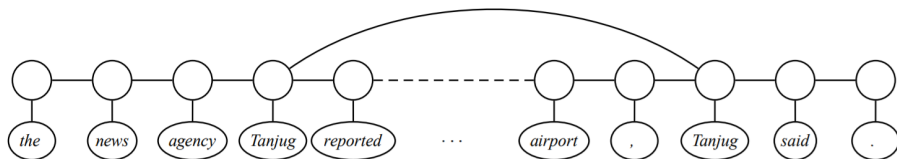


Figure: An example from CoNLL 2003 Named Entity Recognition

Enable label consistency



# Motivation

Structured prediction:

- Capture Label dependency
- Avoid repetition in text generation
- ...

However,

- Hard to capture long dependency between structure outputs.
  - ▶ Energy-based models
- Inference for energy-based model is computational challenging!
  - ▶ Intractable
  - ▶ Exact inference/gradient descent inference is slow

# Table of Contents

- 1 Motivation
- 2 Inference Networks (InfNets) [Tu and Gimpel (2018, 2019)]
- 3 Joint Training of SPENs and Two InfNets
- 4 Conclusions

# Energy Function and Inference Network

## Definition

An energy function [LeCun et al., 2006; Belanger and McCallum, 2016]  $E_{\Theta} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  parametrized by  $\Theta$  that uses a functional architecture to compute a scalar energy for an input/output pair.

## Test-time inference

At test time, for a given input  $\mathbf{x}$ , prediction is done by choosing the output with the lowest energy.

$$\hat{\mathbf{y}} = \arg \min_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y}).$$

## Inference Networks [Tu and Gimpel, 2018]

A test-time inference network  $A_{\Psi} : \mathcal{X} \rightarrow \mathcal{Y}_R$  is parameterized by  $\Psi$  and trained with the goal that

$$A_{\Psi}(\mathbf{x}) \approx \arg \min_{\mathbf{y} \in \mathcal{Y}_R(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y}).$$

# Training Objectives

## SPEN Loss

Belanger and McCallum (2016) use a structured hinge loss for training SPENs

$$\min_{\Theta} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} \left[ \underbrace{\max_{\mathbf{y} \in \mathcal{Y}_R(\mathbf{x})} (\Delta(\mathbf{y}, \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, \mathbf{y}) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i))}_{\text{cost-augmented inference}} \right]_+$$



# Training Objectives

## SPEN Loss

Belanger and McCallum (2016) use a structured hinge loss for training SPENs

$$\min_{\Theta} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} \left[ \underbrace{\max_{\mathbf{y} \in \mathcal{Y}_R(\mathbf{x})} (\Delta(\mathbf{y}, \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, \mathbf{y}) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i))}_{\text{cost-augmented inference}} \right]_+$$

## Our Objective [Tu and Gimpel, 2018]

We parametrize an inference network using  $\Phi$ , alternately optimize  $\Theta$  and  $\Phi$  (like adversarial training):

$$\min_{\Theta} \max_{\Phi} \sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} [\Delta(F_{\Phi}(\mathbf{x}_i), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, F_{\Phi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i)]_+$$

# Pipeline

Recall:  $\Theta$  is params for energy function,  $\Phi$  for **cost-augmented** InfNet,  $\Psi$  for **test-time** InfNet.

Step 1:  $\hat{\Theta}, \hat{\Phi} = \min_{\Theta} \max_{\Phi}$

$$\sum_{\langle \mathbf{x}_i, \mathbf{y}_i \rangle \in \mathcal{D}} [\Delta(F_{\Phi}(\mathbf{x}_i), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, F_{\Phi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i)]_+$$

Update  $\Phi$  to yield output with low energy and high cost.

Step 2:  $\hat{\Psi} = \arg \min_{\Psi} E_{\Theta}(\mathbf{x}, A_{\Psi}(\mathbf{x}))$  where  $A_{\Psi}$  is initialized by trained  $F_{\Phi}$ .

# Table of Contents

- 1 Motivation
- 2 Inference Networks (InfNets) [Tu and Gimpel (2018, 2019)]
- 3 Joint Training of SPENs and Two InfNets**
- 4 Conclusions

# New Objective

## Objective

$$\hat{\Theta}, \hat{\Phi} = \min_{\Theta} \max_{\Phi} \sum_i \underbrace{[\Delta(F_{\Phi}(\mathbf{x}), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, F_{\Phi}(\mathbf{x})) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i)]_+}_{\text{margin-rescaled loss}}$$

Cost-augmented inference:  $F_{\Phi} \approx \arg \min_{\mathbf{y}'} (E_{\Theta}(\mathbf{x}, \mathbf{y}') - \Delta(\mathbf{y}', \mathbf{y}))$ ,

Test-time inference:  $A_{\Psi} \approx \arg \min_{\mathbf{y}'} E_{\Theta}(\mathbf{x}, \mathbf{y}')$ .

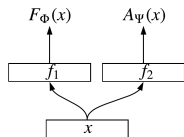
# New Objective

$$\hat{\Theta}, \hat{\Phi}, \hat{\Psi} = \min_{\Theta} \max_{\Phi, \Psi} \sum_i \underbrace{[\Delta(F_{\Phi}(\mathbf{x}), \mathbf{y}_i) - E_{\Theta}(\mathbf{x}_i, F_{\Phi}(\mathbf{x})) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i)]_+}_{\text{margin-rescaled loss}}$$
$$+ \lambda \underbrace{[-E_{\Theta}(\mathbf{x}_i, A_{\Psi}(\mathbf{x}_i)) + E_{\Theta}(\mathbf{x}_i, \mathbf{y}_i)]_+}_{\text{perceptron loss}}$$

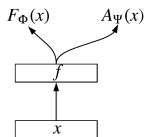
Cost-augmented inference:  $F_{\Phi} \approx \arg \min_{\mathbf{y}'} (E_{\Theta}(\mathbf{x}, \mathbf{y}') - \Delta(\mathbf{y}', \mathbf{y}))$ ,

Test-time inference:  $A_{\Psi} \approx \arg \min_{\mathbf{y}'} E_{\Theta}(\mathbf{x}, \mathbf{y}')$ .

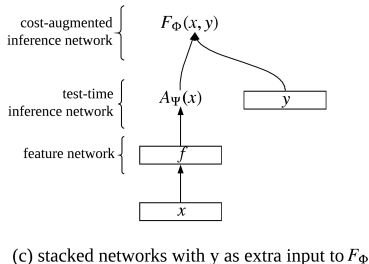
# Joint parametrization for cost-augmented ( $F_\Phi$ ) and test-time ( $A_\Psi$ ) inference networks



(a) separated networks



(b) shared feature networks



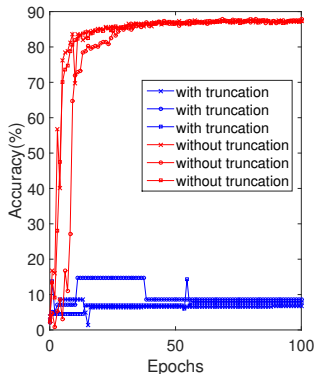
(c) stacked networks with  $y$  as extra input to  $F_\Phi$

Cost-augmented inference:  $F_\Phi \approx \arg \min_{y'} (E_\Theta(\mathbf{x}, \mathbf{y}') - \Delta(\mathbf{y}', \mathbf{y}))$ ,

Test-time inference:  $A_\Psi \approx \arg \min_{y'} E_\Theta(\mathbf{x}, \mathbf{y}')$ .

# Removing Zero Truncation

For inference network objective:  $\max_{\Psi} [h_{\Psi}]_{+} \rightarrow \max_{\Psi} h_{\Psi}$

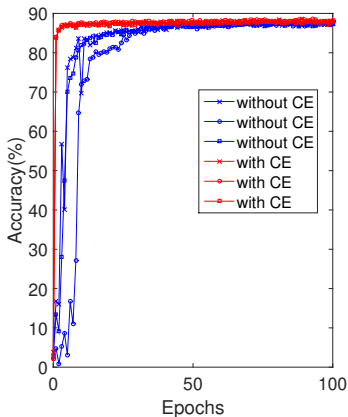


training trajectories without  
CE

Without truncation, the inference network can work well even without any stabilization terms.

# Local Cross Entropy Loss

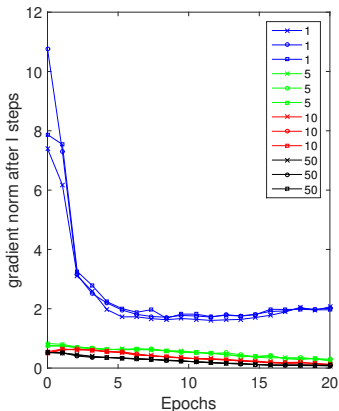
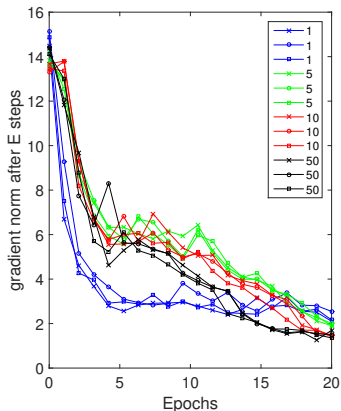
Including  $\sum_{t=1}^{|\mathbf{y}|} \text{CE}(\mathbf{y}_t, A(\mathbf{x})_t)$  or not when without truncation?



The local loss helps speed up convergence and improve accuracy



# Multiple Inference Network Update Steps



Multiple steps of inner loop optimization help inference network maintain near its optimal solution for the given input.

# Experimental Setup

- Energy function: BLSTM-CRF
- Inference network architecture: BLSTM
- Two sequence labeling tasks: Twitter POS tagging (POS) and CoNLL 2003 Named Entity Recognition (NER)

# Comparison of Loss Functions

	POS	NER
baseline: margin-rescaled	89.3	85.2
separated	89.4	85.0
new losses: shared	85.6	85.6
stacked	<b>89.8</b>	85.6

New losses outperform the baseline.

# Comparing cost-augmented ( $F_\phi$ ) and test-time ( $A_\psi$ ) inference networks

		POS	NER
		$A_\psi - F_\phi$	$A_\psi - F_\phi$
margin-rescaled		0.2	0
separated		2.2	0.4
combined	shared	1.9	0.5
	stacked	<b>2.6</b>	<b>1.7</b>

Stacked parameterization shows largest difference between  $F_\phi$  and  $A_\psi$

# Qualitative Analysis

test-time ( $A_\psi$ )	cost-augmented ( $F_\phi$ )
common noun	proper noun
proper noun	common noun
common noun	adjective
proper noun	proper noun + possessive
adverb	adjective
preposition	adverb
adverb	preposition
verb	common noun
adjective	verb
common noun	verb

Table: Top 10 most frequent output differences between  $A_\psi$  and  $F_\phi$ .

$F_\phi$  tends to output tags that are highly confusable with those output by  $A_\psi$ !

# Table of Contents

- 1 Motivation
- 2 Inference Networks (InfNets) [Tu and Gimpel (2018, 2019)]
- 3 Joint Training of SPENs and Two InfNets
- 4 Conclusions

# Conclusions

- SPENs are powerful but learning and inference are hard (due to gradient descent for inference)
- Inference networks can make it easier and more efficient to use SPENs
- Separating inference networks for the two inference problems (cost-augmented and test-time inference) improves accuracy and leads to complementary functionality

Q&A