# SEMI-SUPERVISED STRUCTURED PREDICTION WITH NEURAL CRF AUTOENCODER

*Xiao Zhang[1], Yong Jiang[2], Hao Peng[1], Kewei Tu[2] and Dan Goldwaswer[1]*
Purdue University[1], ShanghaiTech University[2]
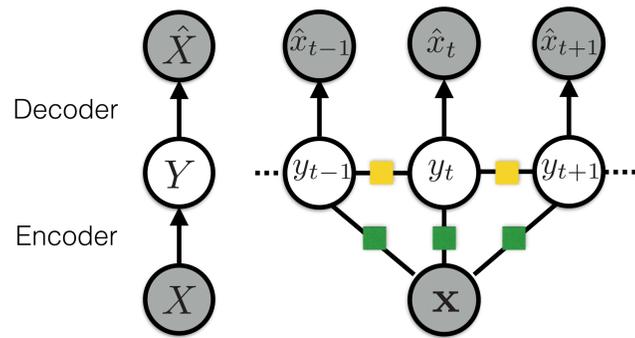
emnlp 2017

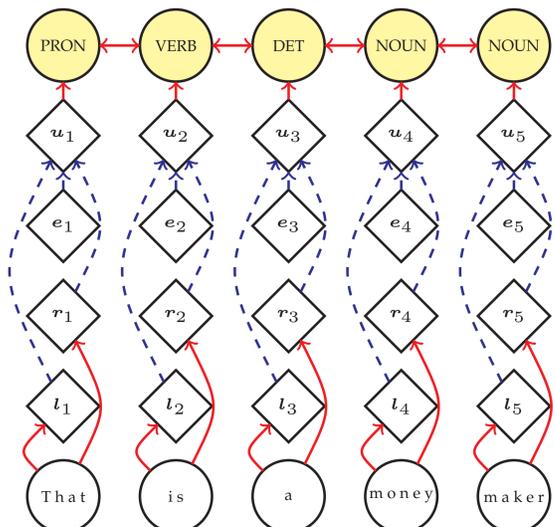PURDUE UNIVERSITY

## MODEL



**From autoencoder (left) to NCRF-AE (right).**

- A generalized autoencoder (AE) models the conditional distribution $P(\hat{X}|X)$.
- A CRF autoencoder (CRF-AE) extends AE by modeling a sequence $P(\hat{x}|x)$, using a CRF as an encoder and a generative model as a decoder.
- A neural CRF autoencoder (NCRF-AE) extends CRF-AE by using deep neural networks as feature extractors in the CRF encoder with the mixed EM algorithm.

$$P_{\Theta,\Lambda}(\hat{x}|x) = \sum_{y} P_{\Theta,\Lambda}(\hat{x}, y|x)$$

$$= \sum_{y} \underbrace{P_{\Theta}(\hat{x}|y)}_{\text{Decoder}} \underbrace{P_{\Lambda}(y|x)}_{\text{Encoder}}.$$

**Neural CRF encoder:**



$$P_{\Lambda}(y|x) = e^{\Phi(x,y)}/Z, Z = \sum_{\tilde{y}} e^{\Phi(x,\tilde{y})};$$

$$\Phi(x,y) = \sum_{t} \phi(x, y_t) + \psi(y_{t-1}, y_t).$$

## ABSTRACT

We propose an end-to-end neural CRF autoencoder (NCRF-AE) model for semi-supervised learning of sequential structured prediction problems. Our NCRF-AE consists of two parts: an encoder which is a CRF model enhanced by deep neural networks, and a decoder which is a generative model trying to reconstruct the input. Our model has a unified structure with different loss functions for labeled and unlabeled data with shared parameters. We developed a variation of the EM algorithm for optimizing both the encoder and the decoder simultaneously by decoupling their parameters. Our experimental results over the Part-of-Speech (POS) tagging task on eight different languages show that the NCRF-AE model can outperform competitive systems in both supervised and semi-supervised scenarios.

## LEARNING

**Unified learning framework** Loss functions for labeled and unlabeled data with *shared parameters*:
Labeled data: $loss_l = -\log P_{\Theta,\Lambda}(\hat{x}, y|x)$; Unlabeled data: $loss_u = -\log P_{\Theta,\Lambda}(\hat{x}|x)$.

**Parameter learning using EM** *Decoupled parameters* update in the E and M steps respectively:

E: $\sum_i \log P(\hat{x}_i|x_i) \geq \sum_i \sum_{y_i} Q(y_i) \log \frac{P(\hat{x}_i, y_i|x_i)}{Q(y_i)}$; M: $\arg\max_{\Theta^{(t)}} \sum_{y \to x} \log \theta_{y \to x}^{(t)} E_{y \sim Q}[C(y, x)]$ $s.t. \sum_x \theta_{y \to x}^{(t)} = 1$.

**Decoding using Viterbi** $y^* = \arg\max_y P_{\Theta,\Lambda}(\hat{x}, y|x)$.

---

**Algorithm 1** Obtain Expected Count ($T_e$)

**Require:** the expected count table $T_e$
1: **for** an unlabeled data example $x_i$ **do**
2:   Compute the forward messages: $\alpha(y, t)$   $\forall y, t$. ▷ $t$ is the position in a sequence.
3:   Compute the backward messages: $\beta(y, t)$   $\forall y, t$.
4:   Calculate the expected count for each $x$ in $x_i$: $P(y_t|x_t) \propto \alpha(y, t) \times \beta(y, t)$.
5:   $T_e(x_t, y_t) \leftarrow T_e(x_t, y_t) + P(y_t|x_t)$   ▷ $T_e$ is the expected count table.
6: **end for**

---

**Algorithm 2** Mixed Expectation-Maximization

1: Initialize expected count table $T_e$ using labeled data $\{x, y\}_i^l$ and use it as $\Theta^{(0)}$ in the decoder.
2: Initialize $\Lambda^{(0)}$ in the encoder randomly.
3: **for** $t$ in $epochs$ **do**
4:   Train the encoder on labeled data $\{x, y\}^l$ and unlabeled data $\{x\}^u$ to update $\Lambda^{(t-1)}$ to $\Lambda^{(t)}$.
5:   Re-initialize expected count table $T_e$ with **0**s.
6:   Use labeled data $\{x, y\}^l$ to calculate real counts and update $T_e$.
7:   Use unlabeled data $\{x\}^u$ to compute the expected counts with parameters $\Lambda^{(t)}$ and $\Theta^{(t-1)}$ and update $T_e$.
8:   Obtain $\Theta^{(t)}$ globally and analytically based on $T_e$.
9: **end for**

## RESULTS

| Models (Supervised) | English | French | German | Italian | Russian | Spanish | Indonesian | Croatian |
|---|---|---|---|---|---|---|---|---|
| HMM | 86.28% | 91.23% | 85.59% | 92.03% | 79.82% | 91.31% | 89.40% | 86.98% |
| CRF | 89.96% | 93.40% | 86.83% | 94.07% | 83.38% | 91.47% | 88.63% | 86.90% |
| LSTM | 90.50% | 94.16% | 88.40% | 94.96% | 84.87% | 93.17% | 89.42% | 88.95% |
| NCRF | 91.52% | 95.07% | 90.27% | 96.20% | 93.37% | 93.34% | 92.32% | 93.85% |
| NCRF-AE | **92.50%** | **95.28%** | **90.50%** | **96.64%** | **93.60%** | **93.86%** | **93.96%** | **94.32%** |

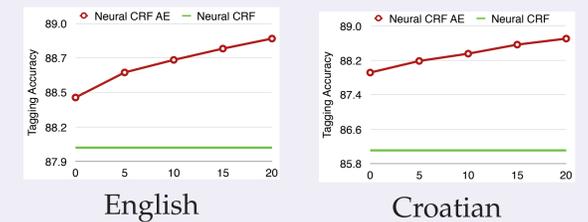| Models (Semi-supervised) | English | French | German | Italian | Russian | Spanish | Indonesian | Croatian |
|---|---|---|---|---|---|---|---|---|
| NCRF-AE (Only Labeled) | 88.41% | 93.69% | 90.75% | 92.17% | 87.82% | 91.70% | 89.06% | 87.92% |
| HMM-EM | 79.92% | 88.15% | 77.01% | 84.57% | 72.96% | 86.77% | 83.61% | 77.20% |
| NCRF-AE (Hard EM) | 86.79% | 92.83% | 89.78% | 90.68% | 86.39% | 91.30% | 88.86% | 86.55% |
| NCRF-AE | **89.43%** | **93.89%** | **90.99%** | **92.85%** | **88.93%** | **92.17%** | **89.41%** | **89.14%** |

## EXPERIMENTS

We evaluated our model on the POS tagging task, in both the supervised and semi-supervised learning settings, over 8 different languages from the UD (Universal Dependencies) 1.4 dataset.

### Error analysis: an example form the test set

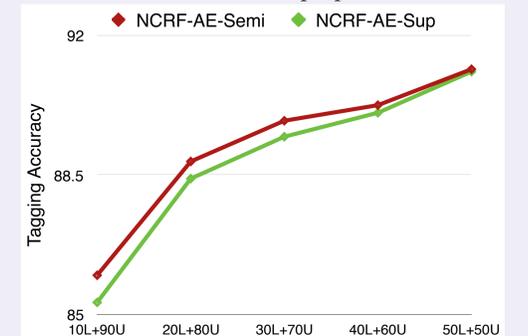| Text | Google | is | a | nice | search | engine | . |
|---|---|---|---|---|---|---|---|
| **Gold** | PROPN | VERB | DET | ADJ | NOUN | NOUN | PUNCT |
| **NCRF-AE** | PROPN | VERB | DET | ADJ | NOUN | NOUN | PUNCT |
| **NCRF** | NOUN | VERB | DET | ADJ | NOUN | NOUN | PUNCT |
| **LSTM** | PROPN | VERB | DET | ADJ | VERB | NOUN | PUNCT |

### Semi-supervised learning effect

UD POS tagging accuracy versus increasing proportion of unlabeled sequences using 20% labeled data. The green straight line is the performance of the neural CRF trained over the labeled data.



English          Croatian

### Varying sizes of labeled data on English

We gradually increased the proportion of labeled data, and in accordance decreased the proportion of unlabeled.



## FUTURE WORK

- Use embeddings for POS tags to compute both the transition score and the generative decoder score.
- Add a prior for predicted labels and cast it into the variational inference framework.