

# Distantly-Supervised Named Entity Recognition with Noise-Robust Learning and Language Model Augmented Self-Training

---

EMNLP 2021

발표: Baek, Hyeongryeol

## Abstract

1. Task: Distantly-supervised Named Entity Recognition
2. Contribution
  - a. Noise robust learning
  - b. Self-training methods

## NER

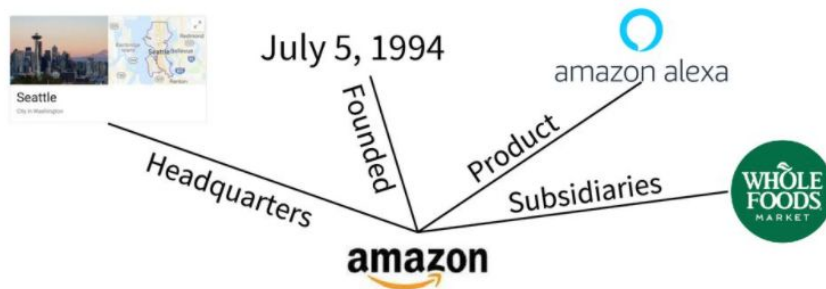
1. Named Entity: a named entity is a real-world object, such as **a person, location, organization, product, etc.**, that can be denoted with a proper name.
2. Distantly Labeled: matching **entity mentions** in the target corpus with **typed entities** in **external gazetteers** or **knowledge bases**.
  - a. gazetteers: **a collection of common entity names**, e.g., Random Name, US First Names Database , Word Lists, etc. for the type PER\*
  - b. Knowledge base (KB): A knowledge base is any system in which **knowledge is stored**, maintained, and accessed.

## KB example

- Predicate: **술어**

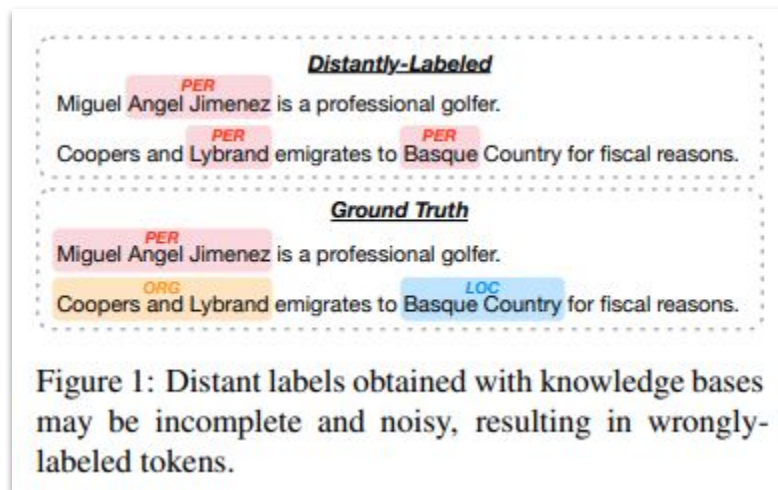
e.g. In the sentence "We went to the airport", "**went to the airport**" is the predicate.

- **Triples** as (subject, predicate, object)  
(**Amazon**, /organization/company/headquarters, **Seattle**)  
(**Amazon**, /organization/company/founded, **July 05, 1994**)  
(**Amazon**, /organization/company/product, **Amazon Alexa**)  
(**Amazon**, /organization/company/subsidiary, **Whole Food Market**)



## Distantly-labeled data

- Problem: incomplete and noisy entity labels -> yielding deteriorated performance



## Contribution

1. Noise-robust learning scheme
  - a. Introducing noise-robust loss function
  - b. Removing a noisy label
2. Proposing unsupervised contextualized augmentation approach
  - a. Novel self-training methods

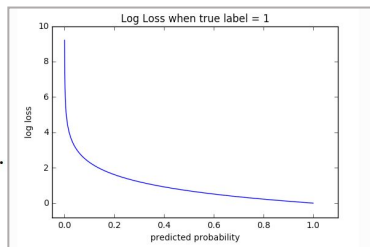
## 1. Noise-robust learning scheme

- a. Introducing noise-robust loss function: giving **less weights** to tokens on which the model prediction is **less consistent with the given labels**

$$\mathcal{L}_{\text{CE}} = - \sum_{i=1}^n \log f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta}),$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{CE}} = - \sum_{i=1}^n \frac{\nabla_{\boldsymbol{\theta}} f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta})}{f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta})}.$$

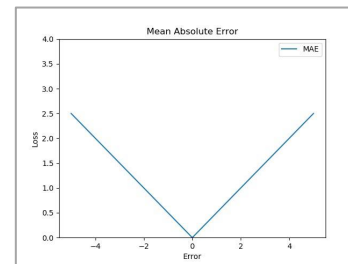
$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta_i} &= -\sum_k y_k \frac{\partial \log(P(y_k))}{\partial \theta_i} \\ &= -\sum_k y_k \frac{\frac{\partial \log(P(y_k))}{\partial P(y_k)} \times \frac{\partial P(y_k)}{\partial \theta_i}}{\frac{\partial \log(P(y_k))}{\partial P(y_k)}} \\ &= -\sum_k y_k \frac{1}{P(y_k)} \times \frac{\partial P(y_k)}{\partial \theta_i} \\ &\quad \left\{ \begin{array}{l} P_i(1-P_i) \text{ (i=k)} \\ -P_i P_i \text{ (i \neq k)} \end{array} \right. \end{aligned}$$



CE  
: Less consistent with  
given  $y_i$  is weighed  
more  $\rightarrow$  sensitive to  
noisy labels

$$\mathcal{L}_{\text{MAE}} = \sum_{i=1}^n (1 - f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta})),$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{MAE}} = - \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta}).$$



MAE  
: noise-tolerant but  
worsening convergence

<Cross entropy loss vs Mean absolute error>

## 1. Noise-robust learning scheme

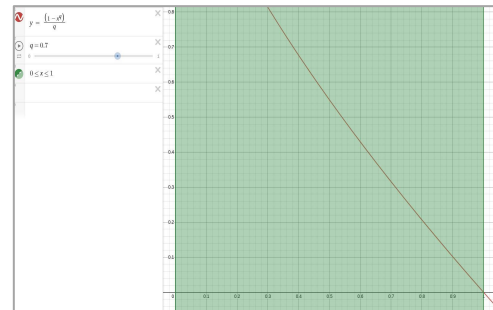
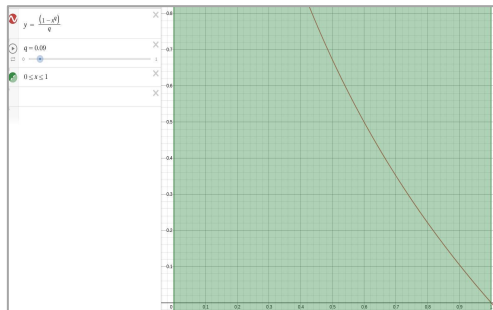
a. Introducing noise-robust loss function: giving **less weights** to tokens on which the model prediction is **less consistent with the given labels**

- GCE vs CE: more noise robust
- GCE vs MAE: giving more attention to difficult tokens

$$\mathcal{L}_{\text{GCE}} = \sum_{i=1}^n \frac{1 - f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta})^q}{q}, \quad (3)$$

where  $0 < q < 1$  is a hyperparameter: When  $q \rightarrow 1$ ,  $\mathcal{L}_{\text{GCE}}$  approximates  $\mathcal{L}_{\text{MAE}}$ ; when  $q \rightarrow 0$ ,  $\mathcal{L}_{\text{GCE}}$  approximates  $\mathcal{L}_{\text{CE}}$  (using L'Hôpital's rule; see Appendix A for the derivation). The gradient is computed as:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\text{GCE}} = - \sum_{i=1}^n \frac{\nabla_{\boldsymbol{\theta}} f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta})}{f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta})^{1-q}}. \quad (4)$$



<Generalized Cross entropy>



## 1. Noise-robust learning scheme

- b. Removing a noisy label: excluding tokens with lower prediction probability than predefined threshold (  $f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta}) \leq \tau$  where  $\tau$  is a threshold value )

$$\mathcal{L}_{\text{GCE}} = \sum_{i=1}^n w_i \frac{1 - f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta})^q}{q}, \quad (5)$$

where  $w_i = 1$  at the start of training and is periodically updated once every several batches as  $w_i = \mathbb{1}(f_{i,y_i}(\mathbf{x}; \boldsymbol{\theta}) > \tau)$ , where  $\mathbb{1}(\cdot)$  is the indicator function.

<Final GCE>

## 1. Noise-robust learning scheme

- + Model ensemble: multiple models' prediction are likely to be consistent on clean data. Training Model\_ENS using  $K$  ( $K=5$ ) models' prediction.

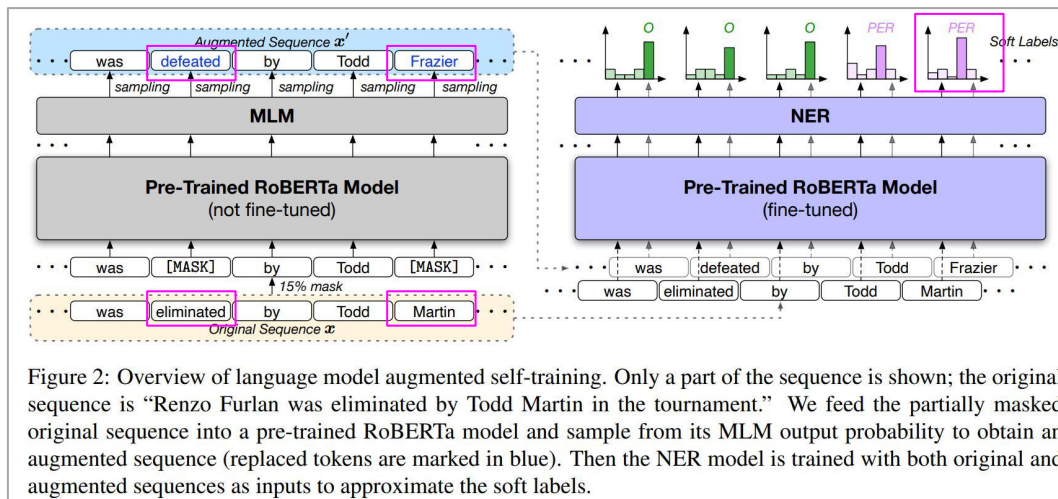
$$\mathcal{L}_{\text{ENS}} = \sum_{i=1}^n \text{KL} \left( \bar{f}_i(\mathbf{x}; \{\boldsymbol{\theta}_k\}_{k=1}^K) \parallel f_i(\mathbf{x}; \boldsymbol{\theta}_{\text{ENS}}) \right), \quad (6)$$

where  $\bar{f}_i(\mathbf{x}; \{\boldsymbol{\theta}_k\}_{k=1}^K) = \frac{1}{K} \sum_{k=1}^K f_i(\mathbf{x}; \boldsymbol{\theta}_k)$  is the  $K$  models' averaged prediction, and we find that  $K = 5$  is sufficient to provide stable ensemble model performance.

## 2. Proposing unsupervised contextualized augmentation approach

### a. Novel self-training methods

- Augmented seq.: token-level perturbation; but likely to be preserving entity type



## 2. Proposing unsupervised contextualized augmentation approach

### a. Novel self-training methods

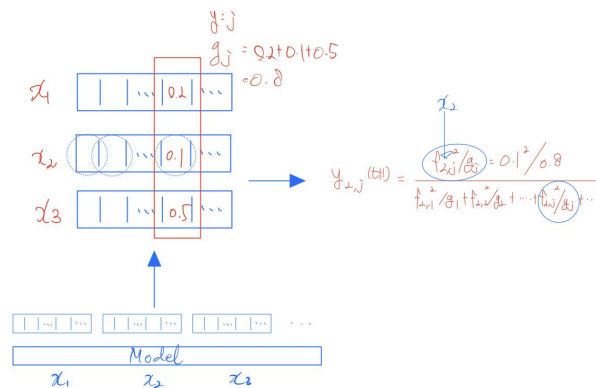
- Soft Labels: enhancing high-confidence predictions while demote low-confidence ones

$$y_{i,j}^{(t+1)} = \frac{f_{i,j}(\mathbf{x}; \theta^{(t)})^2 / g_j}{\sum_{j'} \left( f_{i,j'}(\mathbf{x}; \theta^{(t)})^2 / g_{j'} \right)}, \quad (8)$$

$$g_j = \sum_i f_{i,j}(\mathbf{x}; \theta^{(t)}).$$

Then the model  $\theta^{(t+1)}$  at the next iteration is updated by approximating the soft labels with both the original sequence and the augmented sequence as inputs, via the KL divergence loss:

$$\mathcal{L}_{ST} = \sum_{i=1}^n \text{KL} \left( y_i^{(t+1)} \parallel f_i(\mathbf{x}; \theta^{(t+1)}) \right) + \sum_{i=1}^n \text{KL} \left( y_i^{(t+1)} \parallel f_i(\mathbf{x}'; \theta^{(t+1)}) \right). \quad (9)$$



## 2. Proposing unsupervised contextualized augmentation approach

### a. Novel self-training methods

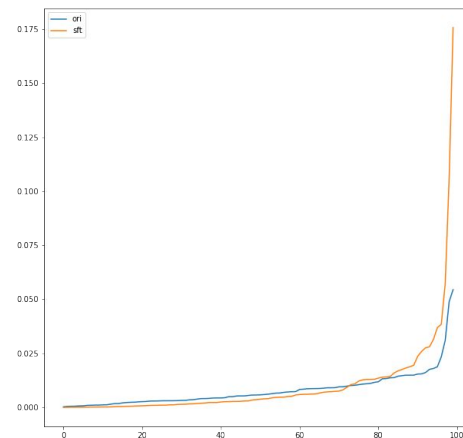
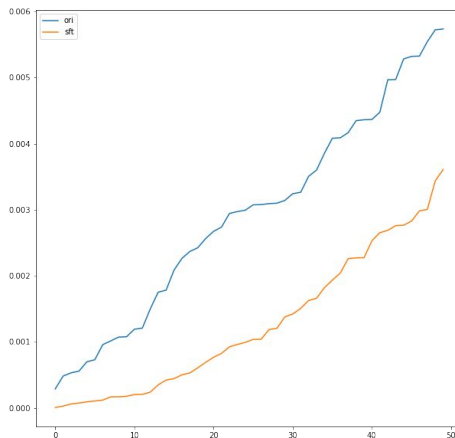
- Soft Labels: enhancing high-confidence predictions while demote low-confidence ones

$$y_{i,j}^{(t+1)} = \frac{f_{i,j}(\mathbf{x}; \boldsymbol{\theta}^{(t)})^2 / g_j}{\sum_{j'} (f_{i,j'}(\mathbf{x}; \boldsymbol{\theta}^{(t)})^2 / g_{j'})}, \quad (8)$$

$$g_j = \sum_i f_{i,j}(\mathbf{x}; \boldsymbol{\theta}^{(t)}).$$

Then the model  $\boldsymbol{\theta}^{(t+1)}$  at the next iteration is updated by approximating the soft labels with both the original sequence and the augmented sequence as inputs, via the KL divergence loss:

$$\mathcal{L}_{ST} = \sum_{i=1}^n \text{KL} \left( y_i^{(t+1)} \parallel f_i(\mathbf{x}; \boldsymbol{\theta}^{(t+1)}) \right) + \sum_{i=1}^n \text{KL} \left( y_i^{(t+1)} \parallel f_i(\mathbf{x}'; \boldsymbol{\theta}^{(t+1)}) \right). \quad (9)$$



## Experiments

Ablations	Pre.	Rec.	F1
<b>RoSTER</b>	0.859	0.849	0.854
w/o GCE	0.817	0.843	0.830
w/o NR	0.830	0.836	0.833
w/o ST	0.844	0.812	0.828

Table 3: Ablation study on CoNLL03 dataset. We compare our full method with ablations (see texts for the abbreviation meanings).

Ablations	Mean (Std.) F1
<b>w. ensemble</b>	0.828 (0.009)
<b>w/o ensemble</b>	0.817 (0.025)

Table 4: Mean and standard deviation (std.) F1 scores of 5 runs (before self-training) with and without model ensemble on CoNLL03 dataset.

## Model Components

: Generalized cross entropy, noisy label removal, self-training, ensemble

**Distant Match:** Shanghai-Ek [Chor]<sub>PER</sub> is jointly owned by the Shanghai Automobile Corporation and [Ek Chor]<sub>PER</sub> China Motorcycle.

**Ground Truth:** [Shanghai-Ek Chor]<sub>ORG</sub> is jointly owned by the [Shanghai Automobile Corporation]<sub>ORG</sub> and [Ek Chor China Motorcycle]<sub>ORG</sub>.

**AutoNER:** Shanghai-Ek [Chor]<sub>PER</sub> is jointly owned by the Shanghai Automobile Corporation and [Ek Chor]<sub>PER</sub> [China]<sub>LOC</sub> Motorcycle.

**BOND:** [Shanghai-Ek Chor]<sub>PER</sub> is jointly owned by the [Shanghai]<sub>LOC</sub> [Automobile Corporation]<sub>ORG</sub> and [Ek Chor]<sub>PER</sub> [China Motorcycle]<sub>ORG</sub>.

**RoSTER:** [Shanghai-Ek Chor]<sub>ORG</sub> is jointly owned by the [Shanghai Automobile Corporation]<sub>ORG</sub> and [Ek Chor China Motorcycle]<sub>ORG</sub>.

Table 6: Case study with **RoSTER** and baselines. The sentence is from CoNLL03.

## Conclusion

- studying the distantly-supervised NER problem without using any human annotations but only distantly-labeled data
- proposing a **noise-robust learning scheme**, consisting of a new loss function and a noisy label removal step.
- proposing a **self-training method** that guides model refinement with its own **high-confidence predictions** and enforces the model to make consistent predictions on original and augmented sequences generated by PLMs