

Rethinking Negative Sampling for Handling Missing Entity Annotations

Anonymous ACL submission

Abstract

Negative sampling is highly effective in handling missing annotations for named entity recognition (NER). One of our contributions is an analysis on how it makes sense through introducing two insightful concepts: missampling and uncertainty. Empirical studies show low missampling rate and high uncertainty are both essential for achieving promising performances with negative sampling. Based on the sparsity of named entities, we also theoretically derive a lower bound for the probability of zero missampling rate, which is only relevant to sentence length. The other contribution is an adaptive and weighted sampling distribution that further improves negative sampling via our former analysis. Experiments on synthetic datasets and well-annotated datasets (e.g., CoNLL-2003) show that our proposed approach benefits negative sampling in terms of F1 score and loss convergence. Besides, models with improved negative sampling have achieved new state-of-the-art results on real-world datasets (e.g., EC).

1 Introduction

With powerful neural networks and abundant well-labeled corpora, named entity recognition (NER) models have achieved promising performances (Huang et al., 2015; Ma and Hovy, 2016; Akbik et al., 2018; Li et al., 2020a). However, in many scenarios, available training data is low-quality, which means a portion of named entities are absent in annotations. Fig. 1 depicts a sentence and its incomplete annotations. Fine-grained NER (Ling and Weld, 2012) is a typical case. Its training data is mainly obtained through applying weak supervision to unlabeled corpora. Past works (Shang et al., 2018b; Li et al., 2021) find missing annotations impact NER models and refer this to *unlabeled entity problem*.

Recently, Li et al. (2021) find it’s the misguidance of unlabeled entities to NER models in train-

Japan	,	co-hosts	of	the	World Cup	in	2002	
<u>LOC</u>					<u>MISC</u>		<u>TIME</u>	
and	ranked	20th	in	the	world	by	FIFA	,
		<u>QUANTITY</u>					<u>ORG</u>	
are	favourites	to	regain	their	title	here	.	

Figure 1: A toy example to show *unlabeled entity problem*. The phrases underlined with dashed lines are the named entities neglected by annotators.

ing that causes their poor performances. To eliminate this adverse impact, they propose a simple yet effective approach based on negative sampling. Compared with its counterparts (Li and Liu, 2005; Tsuboi et al., 2008; Shang et al., 2018b; Peng et al., 2019), this method is of high flexibility, without relying on external resources, heuristics, etc.

While negative sampling has handled missing annotations well, there is no systematic study on how it works, especially what potential factors are involved. From a number of experiments, we find missampling and uncertainty both worth receiving attention. Missampling means that some unlabeled entities are mistakenly drawn into the set of training negatives by negative sampling. To quantitatively describe this, we define missampling rate, the proportion of unlabeled entities in sampled negatives, for a sentence. Uncertainty indicates how hard a sampled negative is for NER models to recognize, and we use entropy to estimate it. Empirical studies show low missampling rate and high uncertainty are both indispensable for effectively applying negative sampling. Besides, based on the observation that entities are commonly sparse, we provide a lower bound for the probability of zero missampling rate with theoretical proof, which is only related to sentence length.

Originally, Li et al. (2021) adopt uniform sampling distribution for negative sampling. Inspired by former findings, we introduce a weighted sampling distribution to displace the uniform one,

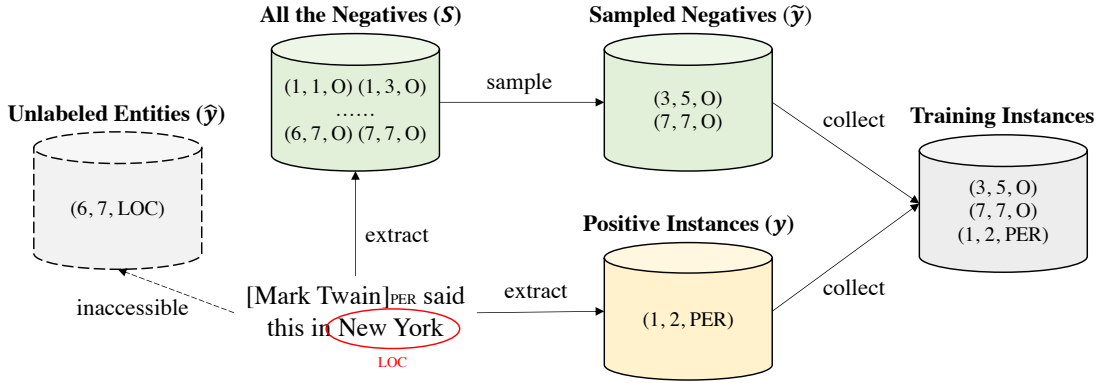


Figure 2: An example to depict how negative sampling collects training negatives given an annotated sentence. The phrase marked by a red circle is an unlabeled entity.

073 which takes missampling and uncertainty into account. Our distribution is purely computed from
 074 the predictions of an NER model. This means it coevolves with the model throughout the training
 075 process. The adaptive property of our method is appealing since it doesn't rely on manual annotations
 076 or additional models to indicate valuable negatives.
 077

078 We have conducted extensive experiments to verify the effectiveness of our weighed sampling distribution.
 079 Results on synthetic datasets and well-annotated datasets (e.g., OntoNotes 5.0) show that
 080 weighted sampling distribution improves negative sampling in performances and loss convergence.
 081 Notably, with improved negative sampling, our NER models have established new state-of-the-art
 082 performances on real-world datasets, like EC (Yang et al., 2018).
 083
 084
 085
 086
 087
 088
 089

090 2 Preliminaries

091 2.1 Unlabeled Entity Problem

092 Given an n -length sentence, $\mathbf{x} = [x_1, x_2, \dots, x_n]$, an annotator (e.g., human) will mark a set of named
 093 entities from it as $\mathbf{y} = \{y_1, y_2, \dots, y_m\}$. n is sequence length and m is set size. Every entity,
 094 y_k , of the set, \mathbf{y} , is denoted as a tuple, (i_k, j_k, l_k) . (i_k, j_k) is the span of the entity that corresponds
 095 to the phrase, $\mathbf{x}_{i_k, j_k} = [x_{i_k}, x_{i_k+1}, \dots, x_{j_k}]$, and l_k is its label. *Unlabeled entity problem* occurs
 096 when some ground truth named entities, $\hat{\mathbf{y}}$, are missed by annotators, which means they are not
 097 contained in the labeled entity collection, \mathbf{y} . In distantly supervised NER (Mintz et al., 2009; Ren
 098 et al., 2015; Fries et al., 2017), this is resulted from the limited coverage of external resources, such as
 099 predefined ontology. In other situations (e.g., fine-grained NER where manual annotation is extremely
 100
 101
 102
 103
 104
 105
 106
 107

108 hard), the cause may be the negligence of human
 109 annotators.

110 Take Fig. 2 as an example. The set of labeled
 111 entities is $\mathbf{y} = [(1, 2, \text{PER})]$, that of unlabeled
 112 entities is $\hat{\mathbf{y}} = \{(6, 7, \text{LOC})\}$, and that of ground-
 113 truth entities is $\mathbf{y} \cup \hat{\mathbf{y}}$.

114 Let \mathcal{S} denote the set that includes all spans of a
 115 sentence, \mathbf{x} , except the ones of annotated named en-
 116 tities, \mathbf{y} . Every span in this set is labeled with "O",
 117 indicating that it's a possible negative. A standard
 118 training strategy for NER models is to minimize
 119 the loss on annotated positives, \mathbf{y} , and all nega-
 120 tive candidates, \mathcal{S} . Unfortunately, since \mathcal{S} might
 121 contain unlabeled entities in $\hat{\mathbf{y}}$, NER models are
 122 seriously misguided in training. To address this
 123 problem, (Li et al., 2021) propose to circumvent
 124 unlabeled entities with negative sampling.

125 2.2 Training with Negative Sampling

126 The core idea is to uniformly sample a few negative
 127 candidates, $\tilde{\mathbf{y}}$, from \mathcal{S} for reliably training NER
 128 models. Under this scheme, the training instances
 129 contain sampled negatives, $\tilde{\mathbf{y}}$, and positives from
 130 annotated entities, \mathbf{y} . With them, $\mathbf{y} \cup \tilde{\mathbf{y}}$, a cross-
 131 entropy loss is incurred as

$$132 \mathcal{J} = \sum_{(i,j,l) \in \mathbf{y} \cup \tilde{\mathbf{y}}} -\log P(l | \mathbf{x}_{i,j}; \theta). \quad (1)$$

133 $P(l | \mathbf{x}_{i,j}; \theta)$ is the probability that the ground truth
 134 label of the span, (i, j) , is l and θ represents the
 135 parameters of a model. Following Li et al. (2021),
 136 our NER models are all span-based, which treat a
 137 span, instead of a single token, as the basic unit for
 138 labeling.

139 Negative sampling is probable to avoid mod-
 140 els being exposed to unlabeled entities. As Fig. 2
 141 shows, the false negative, $(6, 7, \text{O})$, is not involved

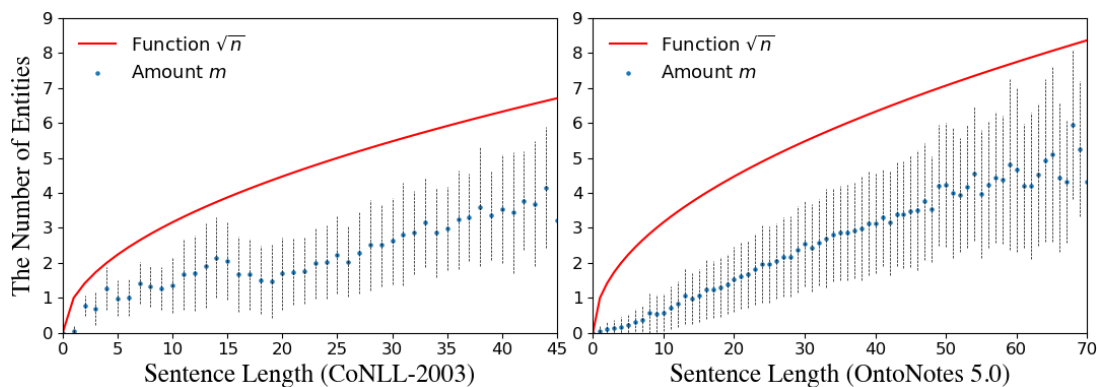


Figure 3: The comparisons between changes of entity number and square root curve.

Average γ	0.76%	1.52%	4.11%
F1 Score	89.86	87.35	83.11

Table 1: The effects of γ on F1.

in training. Li et al. (2021) have empirically confirmed the effectiveness of negative sampling in handling unlabeled entities. However, there is no systematic study to explain how it works, and what factors are relevant.

3 Analyzing Negative Sampling

We analyze how negative sampling leads NER models that suffer from missing entity annotations to promising results from two angles: missampling and uncertainty.

3.1 Missampling Rate

3.1.1 Definition

Missampling rate, γ , is defined as, for a sentence, the proportion of unlabeled entities contained in sampled negatives, $\tilde{\mathbf{y}}$. Formally, it’s computed as

$$\gamma = 1 - \frac{\#\{(i, j, l) \mid (i, j, l) \in \hat{\mathbf{y}}; (i, j, O) \notin \tilde{\mathbf{y}}\}}{\#\tilde{\mathbf{y}}},$$

where $\#$ is an operation that measures the size of an unordered set.

The missampling rate γ reflects the quality of training instances, $\mathbf{y} \cup \tilde{\mathbf{y}}$. A lower averaged rate over the whole dataset means that the NER model meets fewer unlabeled entities in training. Intuitively, this leads to higher F1 scores since there is less misguidance from missing annotations to the model. Hence, missampling is an essential factor for analysis.

3.1.2 Missampling Affects Performance

We design a simulation experiment to empirically verify the above intuition. Like Li et al. (2021), we build synthetic datasets as follows. We start from a well-labeled dataset, i.e., CoNLL-2003 (Sang and De Meulder, 2003), and then mimic *unlabeled entity problem* by randomly masking manually annotated entities with a fixed probability p (e.g., 0.7). In this way, we can obtain unlabeled entities, $\hat{\mathbf{y}}$, and annotated entities, \mathbf{y} , for every sentence, \mathbf{x} .

We can obtain different pairs of a missampling rate and an F1 score through running a negative sampling based model on different synthetic datasets. Table 1 demonstrates several cases, and we can see the trend that lower missampling rates lead to better performances. Therefore, we conclude that missampling affects the effectiveness of negative sampling.

3.1.3 Theoretical Guarantee

We also theoretically prove that negative sampling is very robust to unlabeled entities based on a natural property of named entities.

Entity Sparsity. Unlike other sequence labeling tasks, such as syntactic chunking (Sang and Buchholz, 2000) and part-of-speech tagging (Schmid, 1994), named entities (i.e., non-“O” segments) are commonly sparse in NER datasets.

Fig. 3 depicts some statistics of two common NER datasets, CoNLL-2003 and OntoNotes 5.0. The blue points are the averaged number of entities for sentences of fixed lengths. Every point stands on the center of a dashed line, whose length is the 1.6 variance of the entity numbers. The red curves are the square roots of sentence lengths. To avoid being influenced by “rare events” we erase the points supported by too few cases (i.e., 20).

From the above figure, we can see that the number of ground truth named entities (i.e., unlabeled entities, $\hat{\mathbf{y}}$, and annotated ones, \mathbf{y}) in a sentence is generally smaller than the square root of sentence length, \sqrt{n} . Empirically, we have $\#\mathbf{y} + \#\hat{\mathbf{y}} \leq \sqrt{n}$.

Theorem 1. For a n -length sentence \mathbf{x} , assume $\tilde{\mathbf{y}}$ is the set of sampled negatives with size $\lceil \lambda n \rceil$ ($0 < \lambda < 1$) via negative sampling. If the premise of entity sparsity holds, then the probability of zero missampling rate, i.e., $\gamma = 0$, is bounded.

Proof. Since $\tilde{\mathbf{y}}$ is uniformly sampled from \mathcal{S} without replacement, the probability q that $\gamma = 0$ for a single sentence \mathbf{x} can be formulated as

$$q = \prod_{0 \leq i < \lceil \lambda n \rceil} \left(1 - \frac{\#\hat{\mathbf{y}}}{\frac{n(n+1)}{2} - m - i} \right),$$

where $m = \#\mathbf{y}$. The i -th product term is the probability that, at the i -th sampling turn, the i -th sampled candidate doesn't belong to unlabeled entity set, $\hat{\mathbf{y}}$.

Then we can derive the following inequalities:

$$\begin{aligned} q &\geq \prod_{0 \leq i < \lceil \lambda n \rceil} \left(1 - \frac{\sqrt{n} - m}{\frac{n(n+1)}{2} - m - i} \right) \\ &\geq \prod_{0 \leq i < \lceil \lambda n \rceil} \left(1 - \frac{\sqrt{n}}{\frac{n(n+1)}{2} - i} \right) \\ &> \left(1 - \frac{2\sqrt{n}}{n(n-1) + 2} \right)^{\lceil \lambda n \rceil} \end{aligned}$$

The first inequality holds because of the assumption; the second one holds because $\frac{\sqrt{n} - m}{\frac{n(n+1)}{2} - m - i}$ is monotonically decreases as m increases, and $m \geq 0$; the last inequality hold since $\frac{\sqrt{n}}{\frac{n(n+1)}{2} - i}$ increases with decreasing i , $i < \lceil \lambda n \rceil$, and $\lceil \lambda n \rceil \leq n$.

Because $(1+a)^b \geq 1+ba$ for $a \geq -1 \cap b \geq 1$ and $\lceil \lambda n \rceil < \lambda n + 1$, we have

$$\begin{aligned} q &> \left(1 - \frac{2\sqrt{n}}{n(n-1) + 2} \right)^{\lceil \lambda n \rceil} \\ &\geq 1 - \frac{2(\lambda n + 1)\sqrt{n}}{n(n-1) + 2} \\ &> 1 - \frac{4\lambda\sqrt{n}}{n-1} \end{aligned}$$

The right-most term monotonically increases with the sentence length n , and thus the probability of zero missampling rate for every sentence has a lower bound. \square

H	Top- k	Middle- k	Bottom- k
F1 Score	88.82	87.72	85.56

Table 2: The effects of H on F1.

This theorem shows that missampling rates for standard negative sampling are controllable, and implies why negative sampling succeeds in handling missing annotations.

3.2 Uncertainty

3.2.1 Definition

Assume $P_o(l | \mathbf{x}_{i,j})$ is an oracle model that accurately estimates a label distribution over every span (i, j) . The uncertainty is defined as the entropy of this distribution:

$$\begin{aligned} H(L | \mathbf{X} = \mathbf{x}_{i,j}) = \\ \sum_{l \in L} -P_o(l | \mathbf{x}_{i,j}) \log P_o(l | \mathbf{x}_{i,j}), \end{aligned}$$

where L and \mathbf{X} represent the label space and a span, $\mathbf{x}_{i,j}$, respectively.

Note that the oracle model $P_o(l | \mathbf{x}_{i,j})$ is generally unreachable. the common practice is to additionally train a model $P(l | \mathbf{x}_{i,j}; \theta)$ (see Sec. 2.2) to approximate it. Besides, the approximate model is learned on held-out training data to avoid overconfident estimation.

Uncertainties essentially measure how difficult a case is for models to make a decision (Jurado et al., 2015). In active learning, uncertainty is used to mine hard unlabeled instances for human annotator (Settles, 2009). In our scenario, we suspect that the uncertainty of sampled negatives plays an important role in our training with negative sampling.

3.2.2 Uncertainty Affects Performances

We design an empirical experiment to verify our hypothesis. Specifically, we first randomly and equally split the entire training data with masked entities into two parts, and the first part is used to train an oracle model P_o . For every sentence \mathbf{x} in the second part, we then sample three subsets from \mathcal{S} as training negatives: the first subset denoted by $\tilde{\mathbf{y}}^t$ corresponding to the top- k uncertainties, and the second denoted by $\tilde{\mathbf{y}}^m$ corresponding to middle- k uncertainties, and the third denoted by $\tilde{\mathbf{y}}^b$ corresponding to the bottom- k uncertainties, with $k = \lceil \lambda n \rceil$. Since missampling affects F1 scores as aforementioned, we eliminate the effect on missampling rate by setting $\gamma = 0$ when constructing both subsets, i.e., neither subset contains

any spans included in $\hat{\mathbf{y}}$. Finally, we respectively train three models on top of three negative subsets according to Eq. 1, and report their performances on test data in Table 2. We can see that the model trained on $\tilde{\mathbf{y}}^t$ achieves the best performance, which validates our hypothesis.

4 Improving Negative Sampling

The previous section shows that the effectiveness of negative sampling is dependent on two factors: missampling and uncertainty. As a result, if we had considered both quantities when sampling negatives, we should see larger improvements from final models. In this section, we propose an adaptive and weighted sampling distribution based on these two factors.

Unfortunately, since missampling rate is defined on top of the unlabeled entities $\hat{\mathbf{y}}$ which is unknown in practice, it is not straightforward to apply missampling for improving negative sampling. Therefore, we assume that an oracle model, $z_{i,j,l} = P_o(l | \mathbf{x}_{i,j})$, exists, which is likely to predict the ground-truth label for every span $\mathbf{x}_{i,j}$. Then we define a score $v_{i,j}$ as the difference between the score $z_{i,j,O}$ and the maximum label score on the span (i, j) :

$$v_{i,j} = z_{i,j,O} - \max_{l \in \mathcal{L}} z_{i,j,l}. \quad (2)$$

Intuitively, if $v_{i,j}$ is high, then $z_{i,j,O}$ is high and $\max_{l \in \mathcal{L}} z_{i,j,l}$ is low. In other words, $\mathbf{x}_{i,j}$ is likely to be with ‘‘O’’ label and thus the missampling rate should be small. Hence sampling such a span as a negative won’t hurt NER models. Note that $\max_{l \in \mathcal{L}} z_{i,j,l}$ in the right hand acts as normalization, making $v_{i,j}$ comparable among different spans (i, j) .

We also define an uncertainty score, $u_{i,j}$, as the entropy of the label distribution for a span:

$$\begin{aligned} u_{i,j} &= H(L | \mathbf{X} = \mathbf{x}_{i,j}) \\ &= - \sum_{l \in \mathcal{L}} z_{i,j,l} \log z_{i,j,l}. \end{aligned} \quad (3)$$

As discussed in Sec. 3.2.2, training a NER model with the negatives of higher uncertainty scores, $u_{i,j}$, brings better performances.

Based on $v_{i,j}$ and $u_{i,j}$, we design the following weighted sampling distribution to displace the uniform one when sampling k negatives from \mathcal{S}

without replacement:

$$\begin{cases} r_{i,j} = u_{i,j} * (1 + v_{i,j})^\mu \\ e_{i,j} = \frac{\exp(r_{i,j}/T)}{\sum_{(i',j',O) \in \mathcal{S}} \exp(r_{i',j'}/T)} \end{cases}, \quad (4)$$

where $T \geq 1$ is a temperature to control the smoothness of sampling distribution. $\mu \geq 1$ is to make a trade-off between $v_{i,j}$ and $u_{i,j}$: a high μ will ensure a low missampling rate while a low μ will ensure a high uncertainty score.

To make our approach practical for use, we should specify how to approximate the oracle model, $P_o(l | \mathbf{x}_{i,j})$. In the simulation experiment in Sec. 3.2.1, the oracle model is a fixed model via standard negative sampling which is learned on held-out training data. It’s natural to use such a fixed model to approximate the oracle model here. However, this will cause a side-effect that our approach is not self-contained due to its dependence on an external model.

Consequently, we consider an adaptive style: directly using the NER model, $P(l | \mathbf{x}_{i,j}; \theta)$, itself as the oracle model whose parameter θ is learned during the training process. Under this scheme, T is scheduled as $\sqrt{C - c}$, where C is the number of training epochs and $0 \leq c < C$ is the current epoch number. Since the NER model $P(l | \mathbf{x}_{i,j}; \theta)$ is not accurate in early epochs of training, a more uniform sampling distribution (i.e., higher T) is safer for sampling negatives.

Finally, we get a weighted sampling distribution with the NER model, $P(l | \mathbf{x}_{i,j}; \theta)$, adaptively approximating the oracle model. Our training procedure is the same as that of vanilla negative sampling (see Fig. 2), except for sampling distribution.

5 Experiments

To evaluate our proposed variant (i.e., negative sampling w/ weighted sampling distribution), we have conducted extensive experiments on under-annotated cases: synthetic datasets and real-world datasets. We also validate its superiority in well-annotated scenarios.

5.1 Settings

The well-annotated datasets are CoNLL-2003 and OntoNotes 5.0. CoNLL-2003 contains 22137 sentences and is split into 14987, 3466, and 3684 sentences for training set, development set, and test set, respectively. OntoNotes 5.0 contains 76714

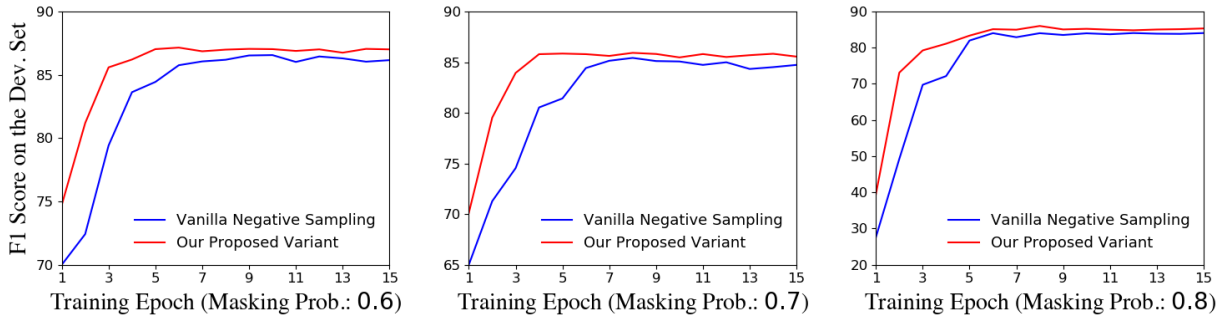


Figure 4: The changes of F1 scores with training epochs on some synthetic datasets.

Masking Prob.	CoNLL-2003		OntoNotes 5.0	
	Vanilla Neg. Sampling	Our Variant	Vanilla Neg. Sampling	Our Variant
0.5	89.22	89.51	88.17	88.31
0.6	87.65	88.03	87.53	88.02
0.7	86.24	86.97	86.42	86.85
0.8	78.84	82.05	85.02	86.12
0.9	51.47	60.57	74.26	80.55

Table 3: The comparisons of F1 scores on synthetic datasets.

sentences from a wide variety of sources. We follow the same format and partition as in Luo et al. (2020). The construction of synthetic datasets is based on well-annotated datasets and has been already described in Sec. 3.

Following prior works (Nooralahzadeh et al., 2019; Li et al., 2021), we adopt EC and NEWS as the real-world datasets. Both of them are collected by Yang et al. (2018). The data contains 2400 sentences annotated by human and is divided into three portions: 1200 for training set, 400 for development set, and 800 for test set. Yang et al. (2018) build an entity dictionary of size 927 and apply distant supervision on a raw corpus to get extra 2500 training cases. NEWS is constructed from MSRA (Levov, 2006). Training set is of size 3000, development set is of size 3328, and test set is of size 3186 are all sampled from MSRA. Yang et al. (2018) collect an entity dictionary of size 71664 and perform distant supervision on the remaining data to obtain extra 3722 cases for training. Both EC and NEWS contain massive incomplete annotations. NER models trained on them suffer from *unlabeled entity problem*.

We adopt the same configurations for all the datasets. The dimensions of scoring layers are 256. L2 regularization and dropout ratio are 10^{-5} and 0.4, respectively. We set $\mu = 8$. This setting is obtained via grid search. We use Adam (Kingma and Ba, 2014) to optimize models. Our models run on

GeForce RTX 2080T. At test time, we convert the predictions from our models into IOB format and use conllevl¹ script to compute the F1 score. In all the experiments, the improvements of our models over the baselines are statistically significant with a rejection probability lower than 0.01.

5.2 Results on Under-annotated Scenarios

We show how NER models with our proposed approach perform on two types of datasets: synthetic datasets (e.g., CoNLL-2003) and real-world datasets (e.g., EC). Synthetic datasets offer us a chance to qualitatively analyze how our approach reacts to changing mask probabilities. For example, we will show that weighted sampling distribution is beneficial in fast loss convergence. Real-world datasets provide more appropriate cases to evaluate NER models, since missing annotations are caused by limited knowledge resources, rather than intentional masking.

5.2.1 Results on Synthetic Datasets

Fig. 4 shows the changes of F1 scores from vanilla negative sampling and our proposed variant with training epochs. The synthetic datasets are constructed from OntoNotes 5.0. We can see that, compared with vanilla negative sampling, our proposed variant obtains far better performances on

¹<https://www.clips.uantwerpen.be/conll2000/chunking/conllevl.txt>.

Method		EC	NEWS
Partial CRF (Yang et al., 2018)		60.08	78.38
Positive-unlabeled (PU) Learning (Peng et al., 2019)		61.22	77.98
Weighted Partial CRF (Jie et al., 2019)		61.75	78.64
-----		-----	-----
BERT-MRC (Li et al., 2020a)		55.72	74.55
BERT-Biaffine Model (Yu et al., 2020)		55.99	74.57
-----		-----	-----
Li et al. (2021)	Vanilla Negative Sampling	66.17	85.39
	w/o BERT, w/ BiLSTM	64.68	82.11
This Work	Our Proposed Variant	67.03	86.15
	w/o BERT, w/ BiLSTM	65.81	83.79

Table 4: The experiment results on two real-world datasets.

Method	CoNLL-2003	OntoNotes 5.0
Flair Embedding (Akbik et al., 2018)	93.09	89.3
HCR w/ BERT (Luo et al., 2020)	93.37	90.30
BERT-MRC (Li et al., 2020a)	93.04	91.11
BERT-Biaffine Model (Yu et al., 2020)	93.5	91.3
-----	-----	-----
Vanilla Negative Sampling (Li et al., 2021)	93.42	90.59
Our Proposed Variant	93.68	91.17

Table 5: The experiment results on well-annotated datasets.

the first few epochs and converges much faster. These results clearly verify the superiority of our weighted sampling distribution.

Table 3 compares vanilla negative sampling with our proposed variant in terms of F1 score. We can draw two conclusions. Firstly, our approach greatly improves the effectiveness of negative sampling. For example, when masking probability p is 0.8, we increase the F1 scores by 4.07% on CoNLL-2003 and 1.29% on OntoNotes 5.0. Secondly, our variant is still robust when *unlabeled entity problem* is very serious. Setting masking probability p from 0.5 to 0.9, our performance on OntoNotes 5.0 only drops by 8.79%. By contrast, it's 32.33% for vanilla negative sampling.

5.2.2 Results on Real-world Datasets

Real-world datasets contain a high percentage of partial annotations caused by distant supervision. Hence, the models trained on them are faced with serious *unlabeled entity problem*.

Table 4 diagrams the results. The F1 scores of negative sampling and Partial CRF are from their papers. We have additionally reported the results of PU Learning², Weighted Partial CRF³,

BERT-MRC⁴, and BERT-Biaffine Model⁵, using their codes. We can draw three conclusions from the table. Firstly, we can see that BERT-MRC and BERT-Biaffine Model both perform poorly on real-world datasets. This manifests the huge adverse impacts of unlabeled entities on models. Secondly, our variant has achieved new state-of-the-art results on the two datasets. Our scores outnumber those of vanilla negative sampling by 1.30% and 0.89% on them. Thirdly, to make fair comparisons, we also report the results of using Bi-LSTM, instead of BERT, as the sentence encoder. This version still notably surpasses prior methods on the two datasets. For example, compared with Weighted Partial CRF, our improvements are 6.57% on EC and 6.55% on NEWS.

5.3 Results on Well-annotated Scenarios

As a by-product, we also evaluate the effectiveness of the proposed method on the well-annotated datasets CoNLL-2003 and OntoNotes 5.0. As shown in Table 5, we have achieved excellent performances on well-annotated datasets. The F1 scores of baselines are copied from Li et al. (2021). With our weighted sampling distribution, the results of negative sampling are improved by 0.28%

²<https://github.com/v-mipeng/LexiconNER>.

³https://github.com/allanj/ner_incomplete_annotation.

⁴<https://github.com/ShannonAI/mrc-for-flat-nested-ner>.

⁵<https://github.com/juntaoy/biaffine-ner>.

on CoNLL-2003 and 0.64% on OntoNotes 5.0. Our model even outperforms BERT-Biaffine Model by 0.19% on CoNLL-2003. Compared with a strong baseline, Flair Embedding, our improvements of F1 scores are 0.63% and 2.09% on the two datasets. These results further verify the effectiveness of the proposed sampling distribution.

The comparison here is in fact unfair for our model, because negative sampling only utilizes a small part of negatives, $\lceil \lambda n \rceil$ rather than $\frac{n(n+1)}{2} - m$ (see Sec. 2 for the details of these numbers). We also have tried using all the negatives for training our model, and found the resulting performances significantly outnumber those of baselines. The purpose of Table 5 is to confirm that negative sampling even works well for situations with complete entity annotations.

6 Related Work

A number of NER models (Lample et al., 2016; Akbik et al., 2018; Clark et al., 2018; Li et al., 2020b; Yu et al., 2020) based on end-to-end neural networks and well-labeled data have achieved promising performances. A representative work is Bi-LSTM CRF (Huang et al., 2015). However, in many situations (e.g., distantly supervised NER), these seemingly perfect models severely suffer from *unlabeled entity problem*, where massive named entities are not annotated in training data. There are some techniques developed by earlier works to mitigate this issue. Fuzzy CRF and AutoNER (Shang et al., 2018b) allow NER models to learn from high-quality phrases that might be potential named entities. Mining these phrases demands external resources (Shang et al., 2018a), which is not flexible for practical usage. Moreover, there is no guarantee that unlabeled entities are fully covered by these phrases. PU Learning (Peng et al., 2019; Mayhew et al., 2019) adopts a weighted training loss and assigns low weights to false negative instances. This approach is limited by requiring prior information or heuristics. Partial CRF (Yang et al., 2018; Jie et al., 2019) is an extension of CRF, which marginalizes the loss over all candidates that are compatible with the incomplete annotation. While being theoretically attractive, this approach still needs a portion of well-annotated data to obtain true negatives, which limits its use in real-world applications. For example, in fine-grained NER (Ling and Weld, 2012), all the training data are produced through weak supervision, and its

manual annotation is very difficult, so obtaining enough high-quality data is not practical.

Recently, Li et al. (2021) find that unlabeled entities severely misguide the NER models during training. Based on this observation, they introduce a simple yet effective approach using negative sampling. It’s much more flexible than other methods, without resorting to external resources, heuristics, etc. However, Li et al. (2021) haven’t well explained why negative sampling works and there are weaknesses in their principle analysis. In this paper, we first show two factors that affect how negative sampling avoids NER models from being impacted by missing annotations. Notably, a theoretical guarantee is provided for the zero missampling rate. Then, we propose weighted sampling distribution to further improve negative sampling based on our former findings.

7 Conclusion

In this work, we have made two contributions. On the one hand, we analyze why negative sampling succeeds in handling *unlabeled entity problem* from two perspectives: missampling and uncertainty. Empirical studies show both low missampling rates and high uncertainties are essential for applying negative sampling. Based on entity sparsity, we also provide a theoretical lower bound for the probability of zero missampling rate. On the other hand, we propose an adaptive and weighted sampling distribution that takes missampling and uncertainty into account. We have conducted extensive experiments to verify whether this further improves the effectiveness of negative sampling. Results on synthetic datasets and well-annotated datasets show that our approach benefits in performances and loss convergence. With improved negative sampling, our NER models also have achieved new state-of-the-art results on real-world datasets.

References

- Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. *Contextual string embeddings for sequence labeling*. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1638–1649, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Kevin Clark, Minh-Thang Luong, Christopher D. Manning, and Quoc Le. 2018. *Semi-supervised sequence modeling with cross-view training*. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1914–

679 Erik F Sang and Sabine Buchholz. 2000. Introduc-
680 tion to the conll-2000 shared task: Chunking. *arXiv*
681 *preprint cs/0009008*.

682 Erik F Sang and Fien De Meulder. 2003. Intro-
683 duction to the conll-2003 shared task: Language-
684 independent named entity recognition. *arXiv*
685 *preprint cs/0306050*.

686 Helmut Schmid. 1994. [Part-of-speech tagging with](#)
687 [neural networks](#). In *COLING 1994 Volume 1: The*
688 *15th International Conference on Computational*
689 *Linguistics*.

690 Burr Settles. 2009. Active learning literature survey.

691 Jingbo Shang, Jialu Liu, Meng Jiang, Xiang Ren,
692 Clare R Voss, and Jiawei Han. 2018a. Automated
693 phrase mining from massive text corpora. *IEEE*
694 *Transactions on Knowledge and Data Engineering*,
695 30(10):1825–1837.

696 Jingbo Shang, Liyuan Liu, Xiaotao Gu, Xiang Ren,
697 Teng Ren, and Jiawei Han. 2018b. [Learning](#)
698 [named entity tagger using domain-specific dictio-](#)
699 [nary](#). In *Proceedings of the 2018 Conference on*
700 *Empirical Methods in Natural Language Processing*,
701 pages 2054–2064, Brussels, Belgium. Association
702 for Computational Linguistics.

703 Yuta Tsuboi, Hisashi Kashima, Shinsuke Mori, Hiroki
704 Oda, and Yuji Matsumoto. 2008. [Training condi-](#)
705 [tional random fields using incomplete annotations](#).
706 In *Proceedings of the 22nd International Conference*
707 *on Computational Linguistics (Coling 2008)*, pages
708 897–904, Manchester, UK. Coling 2008 Organizing
709 Committee.

710 Yaosheng Yang, Wenliang Chen, Zhenghua Li,
711 Zhengqiu He, and Min Zhang. 2018. [Distantly su-](#)
712 [pervised NER with partial annotation learning and](#)
713 [reinforcement learning](#). In *Proceedings of the 27th*
714 *International Conference on Computational Linguis-*
715 *tics*, pages 2159–2169, Santa Fe, New Mexico, USA.
716 Association for Computational Linguistics.

717 Juntao Yu, Bernd Bohnet, and Massimo Poesio. 2020.
718 [Named entity recognition as dependency parsing](#). In
719 *Proceedings of the 58th Annual Meeting of the Asso-*
720 *ciation for Computational Linguistics*, pages 6470–
721 6476, Online. Association for Computational Lin-
722 guistics.