

Handling Large Structural Costs in Neural Networks with Slack Rescaling

Anonymous Authors¹

Abstract

A cost-sensitive learning objective is often important for a neural model to achieve good performance in many structured problems. An ideal objective should penalize an incorrect label by its structural discrepancy and a correct label by zero. Since this is non-convex and non-differentiable, one typically turns to a convex surrogate loss in practice (Tsochantaridis et al., 2004). A widely used surrogate is *margin rescaling* which promotes the score of the true label to be greater than the loss-augmented score of the best label. However, for problems with large structural costs, this formulation is not faithful to the ideal objective: even when all structures are correctly classified, the loss may remain high. Consequently, the model wastefully updates parameters on correct instances during training. In this work, we focus on an alternative cost-sensitive objective called *slack rescaling*. Unlike margin rescaling, slack rescaling is invariant to the absolute values of structural costs and ignores labels that are already well-separated. This can be seen as a tighter approximation to the ideal objective. Inference with slack rescaling is in general intractable, but we adapt recent development of an efficient inference algorithm to the deep network. We evaluate our approach on neural graph-based dependency parsing and report promising results.

1. Introduction

In many structured problems, it is often important to carefully account for structural properties of the label space to achieve good performance. This need arises in neural network modeling across many disciplines such as natural language processing (NLP), computer vision, and speech recognition (Bakir et al., 2007). As a concrete example, in

the graph-based neural dependency parser of Kiperwasser & Goldberg (2016) which is trained using structured hinge loss, the accuracy of the parser plunges from 91 to 79.4 when the training objective disregard structural costs.

Structured prediction concerns many interesting real world problems where there exists an internal structure within its output (Bakir et al., 2007). A few examples are dependency parsing in natural language processing, segmentation task in computer vision, and speech recognition task. In such examples, output label is consists of micro-labels, and its relationship within the micro-label is expressed by a structure, for instance, a tree or a sequence.

In the era of deep networks due to surprising success in various fields (Krizhevsky et al., 2012)(Goodfellow et al., 2014), building deep structured models to deal with the structured tasks is an important task, and it is well motivated by the success in many fields (Chen et al., 2015; Schwing & Urtasun, 2015).

Incorporating the structure in the model is the essence of the structured models. One of the most well-studied approaches is to incorporate the structure into the loss function. This is called cost-sensitive learning. Compare to zero one flat loss where loss is one if predicted label is different from the true label or zero if correct, in cost-sensitive learning, different loss incurs respect to the structure of the labels, i.e. the loss is calculated modeling how the labels are different. (Tsochantaridis et al., 2004) introduced convex surrogate loss for the cost-sensitive losses, called margin rescaling and slack rescaling.

Two formulations differ how the most violating label is calculated. While margin rescaling criteria is a sum of the feature score and the label score, the slack rescaling criteria is a product between the two. This enables margin rescaling to be computationally efficient since interplay between the two scores is not global, and can be decomposed respect to the substructure. However, margin rescaling can be dominated by one score alone. This is problematic for a large structured problem since even for the label which is well separated, it can be considered as a violating label. This results margin rescaling hard to optimize. In the other hand, criteria of slack rescaling is a product of the two scores, and it is difficult for the label to be considered as the most violating label if it has only one high score. Specifically,

^{*}Equal contribution ¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

labels with a margin larger than a constant margin of one cannot be the most violating label. However, the interplay between the two score is global, and does not decompose with respect to the substructure. Thus, loss augmented inference in slack rescaling formulation is computationally intensive, and intractable for large structures.

To reduce computational burden of loss augmented inference, Sarawagi & Gupta (2008) introduced a method utilizing margin rescaling inference as an oracle, and by calling this oracle iteratively to obtain an approximated slack rescaling argmax label. However, the drawback of the approach that makes the approach impractical is that it involves binary search over λ scalar, which can take tens of iterations. A dynamic programming approach is proposed in (Bauer et al., 2014) for a simple structures. Choi et al. (2016) improves upon (Sarawagi & Gupta, 2008) for efficient optimization and extends to finding exact optimum with a small modification of the margin rescaling oracle. In this paper, we propose an efficient algorithm using margin rescaling oracle iteratively. We demonstrate that slack rescaling can be done almost as same computational complexity as margin rescaling, calling oracle only 2 to 3 times in average in the later stage of the optimization.

We mainly consider the dependency parsing problem in NLP investigated in (Kiperwasser & Goldberg, 2016). We demonstrate that this slack rescaling approach can be utilized in deep network efficiently to achieve the higher performance. This paper is organized as follows: Firstly, we review two convex surrogate losses, margin rescaling and slack rescaling, and compare the pros and the cons of the two approaches. Secondly, we briefly describe how we adapt slack rescaling to the deep network. Thirdly, we visit our main focused application of dependency parsing. Fourthly, we show the empirical evaluation of our approach. We conclude with discussions.

2. Two Surrogate Losses

In this section, we review the two surrogate losses based on (Choi et al., 2016). Structural SVM (Tschantz et al., 2004; Taskar et al., 2003) with margin rescaling formulation is defined as

$$\begin{aligned} \min_{w, \xi} \quad & \frac{C}{2} \|w\|_2^2 + \frac{1}{n} \sum_i \xi_i \\ \text{s.t.} \quad & f_i(y_i) - f_i(y) \geq L(y, y_i) - \xi_i \quad \forall i, y \in \mathcal{Y} \\ & \xi_i \geq 0 \quad \forall i \end{aligned} \quad (1)$$

where C is the regularization constant. The slack rescaling formulation is

$$\begin{aligned} \min_{w, \xi} \quad & \frac{C}{2} \|w\|_2^2 + \frac{1}{n} \sum_i \xi_i \\ \text{s.t.} \quad & f_i(y_i) - f_i(y) \geq 1 - \frac{\xi_i}{L(y, y_i)} \quad \forall i, y \in \mathcal{Y} \\ & \xi_i \geq 0 \quad \forall i \end{aligned} \quad (2)$$

One important point is that surrogate losses are that they upper bound the task loss $\Delta(y, y_i)$ and they are convex. To see the difference directly, let $m(y) = f(y) - f(y_i)$ be the margin of label of y_i for a certain instance. Then, the loss occurred by label y in both formulations are

$$\text{Margin: } \xi_M(y) = L(y, y_i) + m(y) \quad (3)$$

$$\text{Slack: } \xi_S(y) = L(y, y_i) (1 + m(y)) \quad (4)$$

And the empirical loss for the instance is defined as max over the label

$$\text{Margin: } \xi_M = \max_y \xi_M(y) \quad (5)$$

$$\text{Slack: } \xi_S = \max_y \xi_S(y) \quad (6)$$

Since $\xi_M(y_i) = \xi_S(y_i) = 0$, the loss is non-negative. Two main differences of the two surrogate losses are tightness and constraint set of label \mathcal{Y} . To be specific, slack rescaling is tighter to the task loss $\Delta(y, y_i)$ than margin rescaling when the instance is not classified correctly, $m(y) > 0$, and the other way when the instance is classified not correctly. Also, slack rescaling considers a much smaller set of constraint set of \mathcal{Y} disregarding labels which are already well separated.

The first argument can be shown easily by showing following

$$\begin{aligned} \Delta(y, y_i) \leq \xi_S(y) \leq \xi_M(y) & \quad \text{if } m(y) > 0. \\ \Delta(y, y_i) \leq \xi_M(y) \leq \xi_S(y) & \quad \text{if } m(y) < 0. \end{aligned}$$

Since deep network is a powerful model, by increasing the power of deep network, we can expect all the training margin to be well classified. Then, the slack rescaling loss is much closer to the target loss than margin rescaling, especially for large structures with a large target loss $\max_{y, y_i} \Delta(y, y_i)$. Therefore, slack rescaling more preferable objective than the margin rescaling for deep network.

The other difference is the label set that they are considering. Since $\xi_M(y_i) = \xi_S(y_i) = 0$ and objective function takes maximum over the labels, we can disregard y such that $\xi_m(y) < 0$ or $\xi_s(y) < 0$ not changing the ξ_i . Writing down this label set explicitly,

$$\text{Margin: } \tilde{\mathcal{Y}}_m = \{y | L(y, y_i) < m(y)\} \quad (7)$$

$$\text{Slack: } \tilde{\mathcal{Y}}_s = \{y | 1 < m(y)\} \quad (8)$$

We can immediately see that $\tilde{\mathcal{Y}}_m \subseteq \tilde{\mathcal{Y}}_s$, and while slack rescaling removes all well separated labels from the consideration, margin rescaling is more conservative removing the labels by only removing labels separated with a

Algorithm 1 Slack rescaling argmax search

```

 $S \leftarrow \emptyset$ 
for  $t = 1, \dots, T$  do
     $\lambda \leftarrow \mathcal{O}(S)$ 
     $y_t = \arg \max h(y) + \lambda g(y)$ 
    if  $y_t \in S$  then
        return  $\arg \max_{y \in S} h(y)g(y)$ 
     $S \leftarrow S \cup \{y_t\}$ 

```

large margin. This implies that slack rescaling has less constraints to satisfy, and in turn easy to satisfy by removing unnecessary labels, and margin rescaling is hard to satisfy.

However, the main caveat for the slack rescaling is the computational cost. To see this, denote $h(y)$ as score function and error function as $g(y)$. For margin rescaling $h(y) = f(y) - f(y_i)$ and for slack rescaling $h(y) = 1 + f(y) - f(y_i)$, and $g(y) = \Delta(y, y_i)$ for both. Then, the slack rescaling objective is the product between the two, $h(y)g(y)$ and for the margin rescaling it is the sum, $h(y) + g(y)$. Since often for tractable structural problems, label decomposes over the substructure, and finding the maximum label in margin rescaling can benefit from the decomposition since the objective function also decomposes over the substructure. However, slack rescaling cannot benefit from the decomposition since the objective function is the product between the two, and it does not decompose over the substructure. Therefore, it is often intractable to do max slack rescaling label search for a large structure. We discuss this problem in the next section.

3. Adaptation to Deep Network

Here we describe the method to learn with slack rescaling formulation in the deep network. Moving from the margin rescaling formulation to slack rescaling formulation can be straightforward. This method can be widely applicable to the models that is using margin rescaling, changing the loss from (5) to (6). Since we update is mostly done via SGD, we only need to change argmax label of (5) to that of (6). This is the approach in (Choi et al., 2016) described in Algorithm 1. We can scale the error function by λ and find the argmax label in margin rescaling formulation iteratively. How to find the next λ depends on the algorithm. We use variant of Bisection search in (Choi et al., 2016). Then, we update the network according to the maximum label in slack rescaling formulation among the labels we found varying λ noted as S in the Algorithm 1.

Algorithm	Margin	Slack
Labeled attachment score	92.44	92.47
Unlabeled attachment score	93.89	94.02

Table 1. : Results on the Treebank dataset

4. Experiments

We experimented with the state-of-the-art neural dependency parser described in (Kiperwasser & Goldberg, 2016). It uses margin scaling formulation. In the problem, exploiting the structure is found crucial in for the high performance as described in the introduction. For the argmax search algorithm, we used a variant of the bisection search of (Choi et al., 2016), and by efficiently making use of the label cache, we only need to search argmax margin scaling label less than 3 times in the later epoch. We trained for a full dataset for 29 epoch report result on the dev set.

5. Discussion

We adapted slack rescaling approach to the deep neural network to achieve higher performance than the margin rescaling. We showed that a careful choosing of λ can lead to efficient slack rescaling so that it can be deployed in the large deep network. This implies that slack rescaling can be widely used for other structural tasks where the deep network is used.

References

- Bakir, Gükhan H., Hofmann, Thomas, Schölkopf, Bernhard, Smola, Alexander J., Taskar, Ben, and Vishwanathan, S. V. N. *Predicting Structured Data*. The MIT Press, 2007.
- Bauer, Alexander, Gornitz, N, Biegler, Franziska, Muller, K-R, and Kloft, Marius. Efficient algorithms for exact inference in sequence labeling svms. *Neural Networks and Learning Systems, IEEE Transactions on*, 25 (5):870–881, 2014.
- Chen, Liang-Chieh, Schwing, Alexander, Yuille, Alan, and Urtasun, Raquel. Learning deep structured models. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pp. 1785–1794, 2015.
- Choi, Heejin, Meshi, Ofer, and Srebro, Nathan. Fast and scalable structural svm with slack rescaling. In *Artificial Intelligence and Statistics*, pp. 667–675, 2016.
- Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.

330	Kiperwasser, Eliyahu and Goldberg, Yoav. Simple and ac-	385
331	curate dependency parsing using bidirectional lstm fea-	386
332	ture representations. <i>Transactions of the Association for</i>	387
333	<i>Computational Linguistics</i> , 4:313–327, 2016.	388
334		389
335	Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E.	390
336	Imagenet classification with deep convolutional neural	391
337	networks. In <i>Advances in neural information processing</i>	392
338	<i>systems</i> , pp. 1097–1105, 2012.	393
339		394
340	Sarawagi, Sunita and Gupta, Rahul. Accurate max-margin	395
341	training for structured output spaces. In <i>Proceedings of</i>	396
342	<i>the 25th international conference on Machine learning</i> ,	397
343	pp. 888–895. ACM, 2008.	398
344	Schwing, Alexander G and Urtasun, Raquel. Fully	399
345	connected deep structured networks. <i>arXiv preprint</i>	400
346	<i>arXiv:1503.02351</i> , 2015.	401
347		402
348	Taskar, B., Guestrin, C., and Koller, D. Max-margin	403
349	Markov networks. In <i>Advances in Neural Information</i>	404
350	<i>Processing Systems</i> . MIT Press, 2003.	405
351		406
352	Tsochantaridis, Ioannis, Hofmann, Thomas, Joachims,	407
353	Thorsten, and Altun, Yasemin. Support vector machine	408
354	learning for interdependent and structured output spaces.	409
355	In <i>Proceedings of the twenty-first international confer-</i>	410
356	<i>ence on Machine learning</i> , pp. 104. ACM, 2004.	411
357		412
358		413
359		414
360		415
361		416
362		417
363		418
364		419
365		420
366		421
367		422
368		423
369		424
370		425
371		426
372		427
373		428
374		429
375		430
376		431
377		432
378		433
379		434
380		435
381		436
382		437
383		438
384		439