

# Weakly Supervised Label Smoothing

Gustavo Penha and Claudia Hauff

TU Delft

{g.penha-1,c.hauff}@tudelft.nl

**Abstract.** We study Label Smoothing (LS), a widely used regularization technique, in the context of neural learning to rank (L2R) models. LS combines the ground-truth labels with a uniform distribution, encouraging the model to be less confident in its predictions. We analyze the relationship between the non-relevant documents—specifically how they are sampled—and the effectiveness of LS, discussing how LS can be capturing “*hidden similarity knowledge*” between the relevant and non-relevant document classes. We further analyze LS by testing if a curriculum-learning approach, i.e., starting with LS and after a number of iterations using only ground-truth labels, is beneficial. Inspired by our investigation of LS in the context of neural L2R models, we propose a novel technique called *Weakly Supervised Label Smoothing* (WSLS) that takes advantage of the retrieval scores of the negative sampled documents as a weak supervision signal in the process of modifying the ground-truth labels. WSLS is simple to implement, requiring no modification to the neural ranker architecture. Our experiments across three retrieval tasks—passage retrieval, similar question retrieval and conversation response ranking—show that WSLS for pointwise BERT-based rankers leads to consistent effectiveness gains. The source code is available at <https://anonymous.open.science/r/dac85d48-6f71-4261-a7d8-040da6021c52/>.

## 1 Introduction

Neural Learning to Rank (L2R) models are traditionally trained using large amounts of strongly labeled data, i.e., human generated relevance judgements. For example, in ad hoc retrieval each instance is comprised of a query, a document and a relevance judgment. All the other documents in the collection that were not labeled as (non-)relevant for the query, while not specified explicitly, can be viewed as non-relevant for the query. Since utilizing an entire corpus for training a L2R model is practically infeasible, the typical procedure is to rely on the top- $k$  ranked documents for a query obtained from an efficient (but less effective) retrieval model such as BM25. While research has shown that the negative sampler (NS), i.e. the technique to select documents to use as negative samples for a query, matters a great deal in the effectiveness of the learned ranker [1,14,10,23,2] there has been no work on how to make use of the *scores* of the NS, which are currently ignored in the training of L2R models—only the content of the documents are employed.

In this work we first aim to understand, in the realm of neural L2R<sup>1</sup>, a widely used and successful [28,27,22] regularization technique called Label Smoothing [21] (LS), that penalizes the divergence between the predictions and a uniform distribution. We begin by looking into how the choice of NS impacts LS, since in the binary relevance prediction problem LS penalizes the model less than normal training when predicting a negative document as relevant and vice versa. We also analyze whether it is beneficial to use a curriculum-learning inspired procedure for the hyper-parameter that controls the LS strength as shown by recent work on understanding LS in other domains [4,24]. This initial exploration to understand LS leads to the following research question: **RQ1** *Is label smoothing an effective regularizer for neural L2R (and if so, under what conditions)?* Our experimental results on three different retrieval tasks reveal that LS is indeed an effective regularization technique for neural L2R, specifically when **(a)** there is similarity between the relevant and the non-relevant sampled documents, i.e. when we use BM25 as the NS technique, and **(b)** a curriculum-like approach is used to control the strength of the smoothing.

Inspired by our findings, we propose the Weakly Supervised Label Smoothing (WSLS) technique which exploits the NS retrieval scores, as opposed to LS where all labels are smoothed equally, for training neural L2R models. Instead of interpolating the ground-truth label distribution with a uniform distribution (as done in LS), we interpolate it with the NS score distribution. WSLS has two benefits compared to using the ground-truth labels: (a) it regularizes the neural ranker by penalizing overconfident predictions and (b) it provides additional supervision signal through weak supervision [3] for the negative sampled documents. WSLS is simple to implement, and requires no modification to the neural ranker architecture, but only to the labels using weak supervision scores that are readily available. Our experiments to answer our second research question (**RQ2** *Is WSLS more effective than LS for training neural L2R models?*) reveal that WSLS is a better way of smoothing the labels by providing additional weak supervision obtained from the negative sampling procedure. We reach relative gains of 0.5% in effectiveness across tasks.

## 2 Background: Label Smoothing (LS)

Given an input instance  $x$  (a query and document combination), two classes ( $k = 0$  means not relevant and  $k = 1$  relevant, and thus here  $K = 2$ ), a ground truth distribution  $q(k | x)$  and predictions from the neural L2R model  $p(k | x) = \frac{\exp(z_k)}{\sum_{i=1}^K \exp(z_i)}$ , where  $z_i$  are the logits, we can use the cross entropy loss for training:  $\ell = -\sum_{k=0}^K \log(p(k)q(k))$ , where  $q(k) = \delta_{k,y}$ , and  $\delta_{k,y}$  is Dirac delta (equals 1 for  $z = y$  and 0 otherwise). Maximizing the log-likelihood of the correct label is approached if the logit corresponding to the ground-truth label is much greater than all other logits:  $z_y \gg z_k$  for all  $k \neq y$ . This encourages the model

<sup>1</sup> Binary relevance prediction is quite different from other domains such as image classification and language modelling which employ up to thousands of distinct classes.

to be overconfident in its predictions, which might not generalize well. Label smoothing [21] is a regularization mechanism to encourage the model to be less confident. Given a distribution  $u(k)$ , *independent* of the training example  $x$ , and a smoothing parameter  $\epsilon$ , for a training example with ground-truth label  $y$ , we replace the label distribution  $q(k | x) = \delta_{k,y}$  with  $q'(k | x) = (1 - \epsilon)\delta_{k,y} + \epsilon u(k)$ . In LS the uniform distribution is employed, i.e.  $u(k) = 1/K$ .

While LS is a widely used technique to regularize models, the reasons underlying its successes [28,27,22] and failures [12,19] remain unclear. Müller et. al [17] showed that while LS impairs teacher models to do knowledge distillation<sup>2</sup> it improves the models' calibration, i.e. how representative the predictions are with respect to the true likelihood of correctness [7].

**Curriculum Learning for Label Smoothing (T-LS)** Xu et. al [24] argued that given the empirical evidence of LS ineffectiveness in certain cases, it is natural to combine LS with the ground-truth labels during training in a two-stage training procedure and thus proposed T-LS: starts training with LS, i.e.  $\epsilon > 0$ , and after  $X$  training instances use normal training, i.e.  $\epsilon = 0$  (the unmodified ground-truth labels are used). Similarly, Dogan et. al [4] proposed to move from a distribution of labels smoothed by the similarity between label classes towards the ground-truth labels with a curriculum leaning procedure. In this paper we resort to T-LS<sup>3</sup> [24] to test whether a curriculum learning inspired approach for  $\epsilon$  is required or not in the training of neural L2R models.

### 3 Weakly Supervised Label Smoothing (WSLS)

We propose to replace the uniform distribution  $u(k)$  that is independent of the example  $x$ , with a weakly supervised function  $w(k | x)$ , which is readily available for documents with label 0 as part of the negative sampling procedure of L2R, at no additional cost: the negative sampler (NS) score<sup>4</sup>. Specifically,  $q'(k | x) = (1 - \epsilon)\delta_{k,y} + \epsilon NS(k | x)$ , where  $NS(k | x)$  is the negative sampling procedure score for instance  $x$  and label class  $k$ . If we use *BM25* to retrieve negative samples<sup>5</sup>, then for  $k = 0$  we have  $q'(k | x) = (1 - \epsilon)\delta_{k,y} + \epsilon BM25(x)$  and when  $k = 1$  we fall back to LS since we have strong labeled data:  $q'(k | x) = (1 - \epsilon)\delta_{k,y} + \epsilon \frac{1}{K}$ . In the same way we can induce a curriculum learning procedure for LS resulting in T-LS (see §2), we can do it for WSLS, for which we refer to as T-WSLS<sup>6</sup>.

## 4 Experimental Setup

**Tasks & Datasets:** In order to evaluate our research questions, we resort to the three following retrieval tasks: passage retrieval using the 2020 Deep Learning

<sup>2</sup> Knowledge distillation [8] is the process of using the predictions of a teacher model with higher complexity/size as the ground-truth distribution for a weaker model.

<sup>3</sup> Initial experiments where we decreased  $\epsilon$  linearly [4] were as effective as T-LS [24].

<sup>4</sup> When building the training col. of triplets for DL TREC 2019/20, for each query there is 1 relevant passage; the non-relevant passages are retrieved using BM25.

<sup>5</sup> Since the BM25 scores are not between 0 and 1 we apply min-max scaling.

<sup>6</sup> This bears resemblance to strategies for combining strong and weak labeled data [20].

track of TREC (TREC-DL) dataset<sup>7</sup>, similar question retrieval with the Quora Question Pairs [9] (QQP) dataset and conversation response ranking with the MANtIS [18] dataset. We use them due to the large amount of labeled examples (required for training neural ranking models) and diversity of tasks.

**Implementation details & Evaluation:** We use BERT-based ranking as a strong neural L2R baseline. We follow previous research [15] and fine-tune BERT using the [CLS] token to predict binary relevance—the query and the document are concatenated using the [SEP] token and used as input—using the cross-entropy loss and Adam optimizer [11] with  $lr = 5^{-6}$  and  $\epsilon = 1^{-8}$ . We train with a batch size of 32 and fine-tune the models for 50000 training instances. We train and test each model 5 times using different random seeds with 10 total candidate documents by query<sup>8</sup>. We resort to a standard evaluation metric in conversation response ranking [25,6]: recall at position  $K$  with  $n$  candidates<sup>9</sup>:  $R_n@K$ . Since all tasks here are concerned with re-ranking  $R_n@K$  is a suitable sampled metric<sup>10</sup> to compare models on how high the relevant documents are ranked when having only  $n$  candidates. We resort to a robust and widely used NS to obtain such candidates: BM25. We refer to using the query as input to BM25 and selecting the top  $n - 1$  ranked documents as  $NS_{BM25}$ . We also use random sampling ( $NS_{random}$ )—which samples candidate documents from the whole collection with the same probability and thus brings documents that are quite different from the relevant one—to better understand LS.

**Table 1.** Average  $R_{10}@1$  and the standard deviation results of 5 runs with different random seeds for BERT with label smoothing (w. LS) and BERT with two-stage label smoothing (w. T-LS) for different negative samplers during training ( $NS_{BM25}$  and  $NS_{random}$ ) and  $\epsilon = 0.2$  for the development set. Bold indicate the highest values for each dataset and  $\blacktriangle/\blacktriangledown$  superscripts indicate significant gains and losses respectively over the baseline (BERT) using paired Student’s t-test with confidence level of 0.95.

	$NS_{BM25}$			$NS_{random}$		
	TREC-DL	QQP	MANtIS	TREC-DL	QQP	MANtIS
BERT	0.568±.00	0.581±.03	0.612±.01	<b>0.385±.01</b>	0.444±.01	<b>0.350±.01</b>
w. LS	0.564±.01 $\blacktriangledown$	0.593±.01 $\blacktriangle$	0.612±.01	0.304±.05 $\blacktriangledown$	0.440±.03 $\blacktriangledown$	0.348±.01 $\blacktriangledown$
w. T-LS	<b>0.570±.01<math>\blacktriangle</math></b>	<b>0.598±.01<math>\blacktriangle</math></b>	0.612±.01	0.382±.02 $\blacktriangledown$	0.444±.01	0.345±.01 $\blacktriangledown$

## 5 Results

**Effectiveness of Label Smoothing for Neural Ranking (RQ1)** Table 1 displays the dev. set results<sup>11</sup> for the LS and T-LS techniques when changing

<sup>7</sup> Since the official test split is not available we split the dev. set into two.

<sup>8</sup> So for example in the passage retrieval task there are 10 passages and only one is relevant for each query and for similar question retrieval there are 10 questions.

<sup>9</sup> For example  $R_{10}@1$  indicates the number of relevant documents found at the first position when the model has to rank 10 candidates.

<sup>10</sup> A sampled metric uses a sample of documents instead of the whole collection [13].

<sup>11</sup> Since we do not do any hyper-parameter tuning for RQ1, we resort to the dev. set to avoid overusing the test set.

the NS. The results reveal that when training BERT with  $NS_{\text{random}}$  to sample negative documents, it is not effective to use any type of label smoothing. In fact there is a consistent and statistically significant decrease in the effectiveness compared to BERT. In contrast, when we sample documents to train with  $NS_{\text{BM25}}$  we observe that there are significant gains to train BERT with T-LS, with the exception of MANTIS where there is no statistical difference. When we compare LS with T-LS, we see that it is indeed beneficial to use a curriculum-learning approach for label smoothing (T-LS), which indicates that being more permissive of the mistakes in the first half of training is effective—this is in line with results obtained in other domains [4,24]. **This answers our first RQ positively: label smoothing is an effective regularization technique to train neural L2R models, with gains of 1% of  $R_{10}@1$  compared to standard training (BERT) on average across three different retrieval tasks when (a) using  $NS_{\text{BM25}}$  and (b) a curriculum learning approach for LS.**

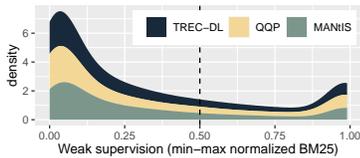
We hypothesize that label smoothing is effective for training neural L2R models if the negative documents are similar to the relevant documents for the query. Our results when changing from  $NS_{\text{BM25}}$  to  $NS_{\text{random}}$  support this hypothesis. Intuitively, if the negative document is random and thus very dissimilar to the query, using a label smoothing regularizer will penalize the model less for this mistake, which might hinder learning. When using label smoothing with a negative document that was sampled using BM25, we are penalizing the model less for choosing a document that is similar to the query in terms of exact matching words. In this way we are teaching the model the similarity between the classes relevant and non-relevant<sup>12</sup> by means of documents that are closer to the classification frontier. Our findings align with [16]: training with topically similar (but non-relevant) documents—as opposed to random documents—allows the model to better discriminate between documents provided by an earlier retrieval stage.

**Effectiveness of Weakly Supervised Label Smoothing (RQ2)** Before we dive into the effectiveness of T-WSLS<sup>13</sup> (RQ2), we investigate the distribution of the normalized weak supervision scores from  $NS_{\text{BM25}}$  in Figure 1. There is a high density for low scores indicating that only a few of the sampled documents receive scores close to the maximum of the list (0.99 score after min-max scaling) and most of them are closer to the minimum (0.00). This is very different from the uniform distribution used by T-LS (dashed vertical line), which does not change according to the sample, and with two classes ( $K = 2$ ) is equal to 0.5, whereas the mean of the weak supervision distribution is 0.33. This suggests that the optimal  $\epsilon$  for T-WSLS is different from T-LS.

Based on this observation, we test different values of  $\epsilon$  on the dev. set in order to tune this hyper-parameter and use it on the test set. Figure 2 displays the effect of  $\epsilon$  on the effectiveness of the proposed approach. The highest  $R_{10}@1$

<sup>12</sup> A similar reasoning can be found in recent work which discusses that the similarity between classes on the wrong responses, i.e. “*hidden similarity knowledge*” [8], is helpful for learning better neural networks [26,4,5].

<sup>13</sup> Based on RQ1 results we use the two-stage approaches here (T-LS and T-WSLS).

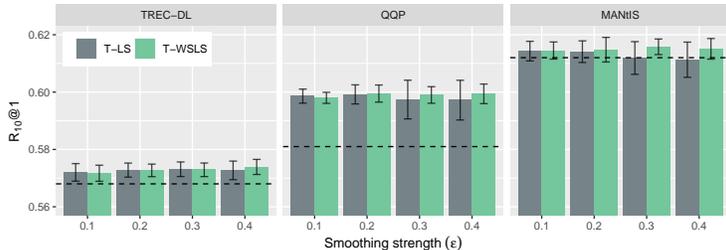


**Fig. 1.** Stacked and smoothed weak supervision distributions used for WLS from the min-max normalized scores of  $NS_{BM25}$ . The dashed vertical line indicates the distribution used by LS (uniform with  $K = 2$ ).

	TREC-DL	QQP	MANTIS
BERT	$0.599 \pm .00$	$0.595 \pm .01$	$0.609 \pm .01$
w. T-LS	$0.601 \pm .00^*$	$0.596 \pm .01$	$0.607 \pm .01$
w. T-WLS	<b><math>0.604 \pm .00^{\Delta}</math></b>	<b><math>0.598 \pm .01^{\Delta}</math></b>	$0.609 \pm .01^{\Delta}$

**Table 2.** Average  $R_{10}@1$  and the standard deviation results of 5 runs with different random seeds for the test set.  $\blacktriangle/\blacktriangledown$  and  $\triangle/\triangledown$  superscripts indicate significant gains and losses over the baselines (BERT) and (BERT w. T-LS) respectively using paired Student’s t-test with confidence level of 0.95 and Bonferroni correction.

values are observed for T-WLS: 0.574 (+1% over the baseline w/o T-WLS) for TREC-DL when  $\epsilon = 0.4$ , 0.600 (+3.2%) for QQP when  $\epsilon = 0.2$  and 0.6151 (+0.5%) for MANTIS when  $\epsilon = 0.4$ . When we apply the best models (for both T-LS and T-WLS) found using the dev. set on the test set, we see in Table 2 that BERT w. T-WLS outperforms both BERT and BERT w. LS with statistical significance (with the exception of MANTIS where there is no difference). **This answers RQ2 indicating that WLS is indeed more effective than LS with statistically significant gains on all tasks against T-LS and with an average of 0.5% improvement over BERT.**



**Fig. 2.** T-LS and T-WLS sensitivity to the hyperparameter  $\epsilon$  for the dev. set. Error bars indicate the 95% confidence intervals for  $R_{10}@1$  over 5 runs with different random seeds. Dashed horizontal lines indicate the baseline w/o label smoothing ( $\epsilon = 0$ ).

## 6 Conclusion

We studied LS in the context of neural L2R models. Our findings indicate that LS is effective when there is similarity between relevant and non-relevant documents and that using curriculum learning for the strength of the regularization is effective. We proposed a technique that combines the weak supervision scores of negative sampled documents with label smoothing (WLS) which outperforms LS on different retrieval tasks. In future work we will explore WLS in a wider range of retrieval models and tasks.

## References

1. Aslam, J.A., Kanoulas, E., Pavlu, V., Savev, S., Yilmaz, E.: Document selection methodologies for efficient and effective learning-to-rank. In: Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval. pp. 468–475 (2009)
2. Cohen, D., Jordan, S.M., Croft, W.B.: Learning a better negative sampling policy with deep neural networks for search. In: Proceedings of the 2019 ACM SIGIR International Conference on Theory of Information Retrieval. pp. 19–26 (2019)
3. Dehghani, M., Zamani, H., Severyn, A., Kamps, J., Croft, W.B.: Neural ranking models with weak supervision. In: Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 65–74 (2017)
4. Dogan, U., Deshmukh, A.A., Machura, M., Igel, C.: Label-similarity curriculum learning. arXiv preprint arXiv:1911.06902 (2019)
5. Furlanello, T., Lipton, Z.C., Tschannen, M., Itti, L., Anandkumar, A.: Born again neural networks. arXiv preprint arXiv:1805.04770 (2018)
6. Gu, J.C., Li, T., Liu, Q., Zhu, X., Ling, Z.H., Su, Z., Wei, S.: Speaker-aware bert for multi-turn response selection in retrieval-based chatbots. arXiv preprint arXiv:2004.03588 (2020)
7. Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. arXiv preprint arXiv:1706.04599 (2017)
8. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015)
9. Iyer, S., Dandekar, N., Csernai, K.: First quora dataset release: Question pairs (2017)
10. Karpukhin, V., Oğuz, B., Min, S., Wu, L., Edunov, S., Chen, D., Yih, W.t.: Dense passage retrieval for open-domain question answering. arXiv preprint arXiv:2004.04906 (2020)
11. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
12. Kornblith, S., Shlens, J., Le, Q.V.: Do better imagenet models transfer better? In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2661–2671 (2019)
13. Krichene, W., Rendle, S.: On sampled metrics for item recommendation. In: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. pp. 1748–1757 (2020)
14. Li, J., Tao, C., Feng, Y., Zhao, D., Yan, R., et al.: Sampling matters! an empirical study of negative sampling strategies for learning of matching models in retrieval-based dialogue systems. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 1291–1296 (2019)
15. Lin, J., Nogueira, R., Yates, A.: Pretrained transformers for text ranking: Bert and beyond. arXiv preprint arXiv:2010.06467 (2020)
16. Mitra, B., Diaz, F., Craswell, N.: Learning to match using local and distributed representations of text for web search. In: Proceedings of the 26th International Conference on World Wide Web. pp. 1291–1299 (2017)
17. Müller, R., Kornblith, S., Hinton, G.E.: When does label smoothing help? In: Advances in Neural Information Processing Systems. pp. 4694–4703 (2019)

18. Penha, G., Balan, A., Hauff, C.: Introducing mantis: a novel multi-domain information seeking dialogues dataset. arXiv preprint arXiv:1912.04639 (2019)
19. Seo, J.W., Jung, H.G., Lee, S.W.: Self-augmentation: Generalizing deep networks to unseen classes for few-shot learning. arXiv preprint arXiv:2004.00251 (2020)
20. Shnarch, E., Alzate, C., Dankin, L., Gleize, M., Hou, Y., Choshen, L., Aharonov, R., Slonim, N.: Will it blend? blending weak and strong labeled data in a neural network for argumentation mining. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). pp. 599–605 (2018)
21. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2818–2826 (2016)
22. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)
23. Xiong, L., Xiong, C., Li, Y., Tang, K.F., Liu, J., Bennett, P., Ahmed, J., Overwijk, A.: Approximate nearest neighbor negative contrastive learning for dense text retrieval. arXiv preprint arXiv:2007.00808 (2020)
24. Xu, Y., Xu, Y., Qian, Q., Li, H., Jin, R.: Towards understanding label smoothing. arXiv preprint arXiv:2006.11653 (2020)
25. Yuan, C., Zhou, W., Li, M., Lv, S., Zhu, F., Han, J., Hu, S.: Multi-hop selector network for multi-turn response selection in retrieval-based chatbots. In: EMNLP. pp. 111–120 (2019)
26. Yuan, L., Tay, F.E., Li, G., Wang, T., Feng, J.: Revisiting knowledge distillation via label smoothing regularization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3903–3911 (2020)
27. Zeyer, A., Irie, K., Schlüter, R., Ney, H.: Improved training of end-to-end attention models for speech recognition. arXiv preprint arXiv:1805.03294 (2018)
28. Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8697–8710 (2018)