# Cross Domain Regularization for Neural Ranking Models using Adversarial Learning

Daniel Cohen*
Center for Intelligent Information
Retrieval
University of Massachusetts Amherst
dcohen@cs.umass.edu

Bhaskar Mitra, Katja Hofmann
Microsoft AI & Research
bmitra@microsoft.com
katja.hofmann@microsoft.com

W. Bruce Croft
Center for Intelligent Information
Retrieval
University of Massachusetts Amherst
croft@cs.umass.edu

## ABSTRACT

Unlike traditional *learning to rank* models that depend on hand-crafted features, neural representation learning models learn higher level features for the ranking task by training on large datasets. Their ability to learn new features directly from the data, however, may come at a price. Without any special supervision, these models learn relationships that may hold only in the domain from which the training data is sampled, and generalize poorly to domains not observed during training. We study the effectiveness of adversarial learning as a cross domain regularizer in the context of the ranking task. We use an adversarial discriminator and train our neural ranking model on a small set of domains. The discriminator provides a negative feedback signal to discourage the model from learning domain specific representations. Our experiments show consistently better performance on held out domains in the presence of the adversarial discriminator—sometimes up to 30% on precision@1.

## 1 INTRODUCTION

Several neural ranking models have been proposed recently that estimate the relevance of a document to a query by considering the raw query-document text [14] or based on the patterns of exact query term matches in the document [5], or a combination of both [10]. These models typically learn to distinguish between the input feature distributions corresponding to a relevant and a less relevant query-document pair by observing a large number of relevant and non-relevant samples during training. Unlike traditional *learning to rank* (LTR) models that depend on hand-crafted features [8], these deep neural models learn higher level representations useful for the

target task directly from the data. Their ability to learn features from the training data is a powerful attribute that enables them to potentially discover new relationships not captured by hand-crafted features. However, as Mitra and Craswell [9] discuss, the ability to learn new features may come at the cost of poor generalization and performance on domains not observed during training. The model, for example, may observe that certain pairs of phrases—*e.g.*, "Theresa May" and "Prime Minister"—co-occur together more often than others in the training corpus. Or, the model may conclude that it is more important to learn a good representation for "Theresa May" than for "John Major" based on their relative frequency of occurrences in training queries. While these correlations and distributions are important if our goal is to achieve the best performance on a single domain, the model must learn to be more robust to them if we instead care about "out of box" performance on unseen domains, *e.g.*, older TREC collections [19]. In contrast, traditional retrieval models (*e.g.*BM25 [12]) and LTR models based on aggregated count based features—that make fewer distributional assumptions—typically exhibit more robust cross domain performances.

Our goal is to train deep neural ranking models that learn useful representations from the data without "overfitting" to the distributions of the training domains. Recently, adversarial learning has been shown to be an effective cross domain regularizer suitable for classification tasks [3, 17]. We adapt a similar strategy to force neural ranking models to learn more domain invariant representations. We train our neural ranking model on a small set of domains and evaluate its performance on held out domains. During training, we combine our ranking model with an adversarial discriminator that tries to predict the domain of the training sample based on the representations learned by the ranking model. The gradients from the adversarial components are reversed when backpropagating through the layers of the ranking model. This provides a negative feedback signal to the ranking model to discourage it from learning representations that may be significant only for specific domains. Our experiments show consistent improvements in ranking performance on held out domains from the proposed adversarial training—sometimes up to 30% improvement on precision@1.

## 2 RELATED WORK

Adversarial networks surfaced shortly after they were introduced in the generative adversarial network (GAN) model. Goodfellow et al. [4] present a generative model that learns a distribution $p_G(x)$ that matches a true distribution $p_{data}(x)$. The generative model receives training updates through a joint loss function shared with an adversarial network, the discriminator, that learns whether a sample is from $p_G(x)$ or $p_{data}(x)$ as a binary classification problem. The

(a) CosSim w/ adversarial discriminator

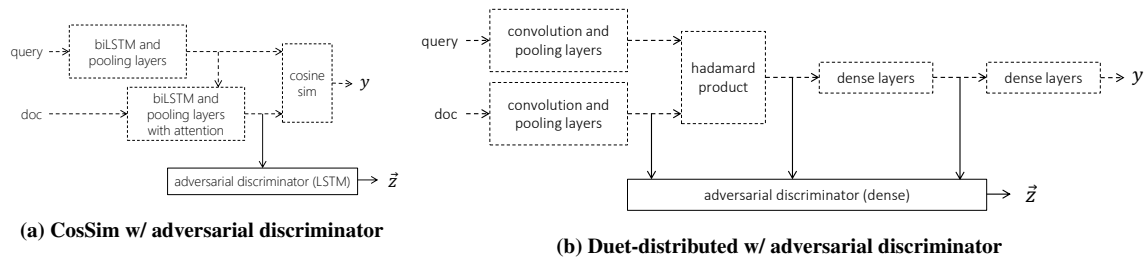(b) Duet-distributed w/ adversarial discriminator

Figure 1: Cross domain regularization of the two baseline models—CosSim and Duet-distributed—using an adversarial discriminator. The discriminator inspects the learned representations of the ranking model and provides a negative feedback signal for any representation that aids domain discrimination.

generator is penalized when the discriminator can successfully classify the sample origin, framing the relationship as a minimax game. While initially proposed for generating continuous data, Donahue et al. [2] extend this work by learning an encoder that maps the data to the latent space **z**. They show that this can learn useful features for image classification tasks without the need for supervised training. Tzeng et al. [18] first propose a form of domain agnostic representation via *domain confusion*, where the maximum mean discrepancy between the final layers of two identical networks over different domains is directly minimized. With the introduction of adversarial agents, Ganin et al. [3] approach the same task of domain agnostic representation by using an adversarial discriminator. The representation of the main network is forced away from a domain specific representation by reversing the gradient updates outside of the adversarial discriminator.

As previous methods used shared weights for both domains, Rozantsev et al. [13] expand on this work showing that unpairing a portion of the classification model, with only a small number of parameters shared prior to input into the final layers, can lead to effective adaptation in supervised and unsupervised settings. Recently, Tzeng et al. [17] have represented a number of past domain adaptation works in a unified framework, referred to as *Adversarial Discriminative Domain Adaptation*, that captures previous approaches as special cases and encompasses a GAN loss into the training of the classifier and adversarial discriminator. This methodology achieves robust domain agnostic models over computer vision collections.

## 3 CROSS DOMAIN REGULARIZATION USING ADVERSARIAL LEARNING

The motivation of the adversarial discriminator is to force the neural model to learn domain independent features that are useful to estimate relevance. Conventional neural ranking models are trained to only optimize for relevance evaluations, disregarding the nature of features learned internally. We propose using an adversarial agent to force the features learned by the ranking model to be domain agnostic by shifting the model parameters in the opposite direction to domain specific spaces on the manifold. This cross domain regularization via domain confusion [17] can be represented as a joint loss function:

$$
\begin{aligned}
\mathcal{L} = \; & \mathcal{L}_{\text{rel}}(q, doc_r, d_{nr}, \theta_D, \theta_{\text{rel}}) \\
& + \lambda \cdot \left( \mathcal{L}_{\text{adv}}(q, doc_r, \theta_D) + \mathcal{L}_{\text{adv}}(q, doc_{nr}, \theta_D) \right)
\end{aligned}
\tag{1}
$$

where $\mathcal{L}_{\text{rel}}$ is a relevance based loss function and $L_{\text{adv}}$ is the adversarial discriminator loss. $q, doc_r$, and $doc_{nr}$ are the query, the relevant document, and the non-relevant documents, respectively. Finally, $\theta_{\text{rel}}$ and $\theta_D$ are the parameters for the relevance and the adversarial models, respectively. $\lambda$ determines how strongly the domain confusion loss should impact the optimization process. We treat it as a hyper-parameter in our training regime. The ranking model is trained on a set of train domains $D_{\text{train}} = \{d_1, \ldots, d_k\}$ separate from the set of held out domains $D_{\text{test}} = \{d_{k+1}, \ldots, d_n\}$ on which it is evaluated.

The discriminator is a classifier that inspects the outputs of the hidden layers of the ranking model, and tries to predict the domain $d_{\text{true}} \in D_{\text{train}}$ of the training sample. The discriminator is trained using a standard cross-entropy loss.

$$
\mathcal{L}_{\text{adv}}(q, doc, \theta_D) = -\log\big(p(d_{\text{true}}|q, doc, \theta_D)\big)
\tag{2}
$$

$$
p(d_{\text{true}}|q, doc, \theta_D) = \frac{exp(z_{\text{true}})}{\sum_{j \in D_{\text{train}}} exp(z_j)}
\tag{3}
$$

Gradient updates are performed via backpropagation through all subsequent layers, including those belonging to the ranking model. However, as proposed by Ganin et al. [3], we utilize a gradient reversal layer. This layer transforms the standard gradient, $\frac{\delta L_{\text{adv}}}{\delta \theta}$ to its additive inverse, $-\frac{\delta \mathcal{L}_{\text{adv}}}{\delta \theta_{\text{rel}}}$. This results in $\theta_{\text{rel}}$ maximizing the domain identification loss, while still allowing $\theta_D$ to learn to discriminate domains. While not directly optimized, this can be viewed as modifying (1) via a sign change for $L_{\text{adv}}$.

*Passage Retrieval Models.* We evaluate our adversarial learning approach on the passage retrieval task. We employ the neural ranking model proposed by Tan et al. [16]—referred to as CosSim in the remaining sections—and the Duet model [10] as our baselines. Our focus in this paper is on learning domain agnostic text representations. Therefore, similar to Zamani et al. [20] we only consider the distributed sub-network of the Duet model.

The CosSim model is an LSTM-based interaction focused architecture. We train the CosSim model in the same manner as [16], with a margin of 0.2 over a hinge loss function. The Duet-distributed is trained by maximizing the log likelihood of the correct passage, as originally proposed in [10]. Similar to [11], we adapt the hyperparameters of the Duet model for passage retrieval. The output of

| source → target | Size | CosSim | | | | Duet-Dist. | | | |
| | | Original | | Adv | | Original | | Adv | |
| | | P@1 | MRR | P@1 | MRR | P@1 | MRR | P@1 | MRR |
|---|---|---|---|---|---|---|---|---|---|
| All→All | 142627 | **0.4229** | 0.6188 | 0.4213(-.3%) | **0.6214(+.4%)** | **0.4514** | **0.6136** | 0.4286(-5%)[†] | 0.6061(-1%)[†] |
| All*→Sports | 139000 | 0.3282 | 0.5194 | **0.4041(+23%)**[†] | **0.5925(+12%)**[†] | 0.2570 | 0.4567 | **0.3282(+28%)**[†] | **0.5011(+10%)**[†] |
| Sports→Sports | 3627 | 0.2146 | 0.5482 | - | - | 0.2415 | 0.3734 | - | - |
| All*→Home | 133372 | 0.3460 | 0.5275 | **0.3645(+5%)**[†] | **0.5433(+3%)**[†] | 0.3314 | 0.5285 | **0.3639(+10%)**[†] | **0.5457(+3%)**[†] |
| Home→Home | 9255 | 0.3014 | 0.5490 | - | - | 0.2477 | 0.4119 | - | - |
| All*→Politics | 138739 | 0.3100 | 0.5101 | **0.3580(+16%)**[†] | **0.5507(+8%)**[†] | 0.3400 | 0.5291 | **0.3516(+3%)**[†] | **0.5342(+3%)**[†] |
| Politics→Politics | 3888 | 0.2219 | 0.5234 | - | - | 0.2160 | 0.5388 | - | - |
| All*→Travel | 140150 | 0.2360 | 0.4486 | **0.2789(+18%)**[†] | **0.4723(+5%)**[†] | 0.2158 | 0.4196 | **0.2842(+32%)**[†] | **0.4532(+8%)**[†] |
| Travel→Travel | 2477 | 0.2263 | 0.5181 | - | - | 0.1895 | 0.3998 | - | - |

**Table 1: Performance across L4 topics, where metrics under each collections represents the performance of the model trained on the opposing two collections. All\* is the entire L4 collection with target topic removed. † represents significance against non adversarial model ($p < 0.05$, Wilcoxon test)**

the Hadamard product is significantly reduced by taking the max pooled representation, the query length is expanded to 20 from 8 tokens, and the max document length is reduced to 300 from the original 1000 tokens.

As opposed to past uses of adversarial approaches [3, 6, 17], ranking requires modeling an interaction between the query and the document. As shown in Figure 1a, the adversarial discriminator in our setting, therefore, inspects the joint query-document representation learned by the neural ranking models. For deeper architectures, such as the Duet-distributed, we allow the discriminator to inspect additional layers within the ranking model, as shown in Figure 1b.

## 4 EXPERIMENTS

### 4.1 Data

*L4.* We use Yahoo's Webscope L4 high quality "Manner" collection [15]. For evaluation and training, all answers that were not the highest voted were removed from the collection to reduce label noise during training and provide a better judgment of performance during evaluation. Training, development, and test sets were created from a 80-10-10 split. Telescoping is used to create answer pools for evaluation from the top 10 BM25 retrieved answers as in [1].

*InsuranceQA* In the InsuranceQA dataset, questions are created from real user submissions and the high quality answers come from insurance professionals. The dataset consists of 12,887 QA pairs for training, 1,000 pairs for validation, and two tests sets containing 1,800 pairs. For testing, each of the 1,800 QA pairs is evaluated with 499 randomly sampled candidate answers.

*WebAP* As both L4 and InsuranceQA are based on isolated passage retrieval for a directed question, we include the WebAP collection from Keikha et al. [7] to examine how well a model trained on isolated passages with specific questions can generalize to a more general passage retrieval task. The format of this collection consists of 82 TREC queries with a total of 8,027 answer passages in total. As only relevant answer passages are annotated in this collection, we create non-relevant documents by using a sliding window of random size. Evaluation is done over a telescoped list of top 100 BM25 retrieved documents.

### 4.2 Training

We experimented with two different training settings—updating the ranking model and the discriminator parameters alternately as proposed by Goodfellow et al. [4], and simultaneously. We also tried different values for $\lambda$. Based on our validation results, we choose to train the CosSim model with alternate updates and $\lambda = 1$. For the Duet-distributed model, we see best performance with simultaneous updates and $\lambda = 0.25$. All models were trained with PyTorch [1] and we implement early stopping based on the validation set.

### 4.3 Evaluation

We evaluate our proposed adversarial approach to cross domain regularization under two settings. Under the *cross topic* setup, we consider the 25 topics in the L4 dataset. We evaluate separately on four of these topics—Sports, Home, Politics, and Travel—each time training the corresponding models on the remaining 24 topics. For the *cross collection* setup, we consider all three collections introduced in Section 4.1. Similar to the cross topic setting, we evaluate our models on each collection individually while training on the remaining two. However, due to more pronounced differences in both size and distributions between these collections—as compared to the differences between the L4 topics—our basic adversarial approach had limited success on the cross collection task. Thus, we adopt two additional changes to our training regime: (i) we sample the training data from the training collections equally to avoid over-fitting to any single collection, and (ii) we feed training samples from the evaluation collection to the adversarial discriminator. We make sure that the training samples from the evaluation collection have no overlap with the test samples. In addition, we clarify that the ranking model receives no parameter updates from these training samples with respect to relevance judgments. These samples are only used to train the discriminator model's loss. This training setup may be appropriate when we want to train on some collections and evaluate on a different collection, where we can leverage the unlabeled documents from the target collection to at least guide the training of the adversarial component.

---

[1] https://github.com/pytorch/pytorch

| source → target | CosSim | | | | Duet-Dist. | | | |
|---|---|---|---|---|---|---|---|---|
| | Original | | Adv | | Original | | Adv | |
| | P@1 | MRR | P@1 | MRR | P@1 | MRR | P@1 | MRR |
| (InsuranceQA, L4)→ WebAP | 0.0901 | 0.2410 | **0.2500** | **0.3873** | 0.1250 | 0.4567 | **0.3286**[†] | **0.5011**[†] |
| (InsuranceQA, WebAP)→ L4 | 0.1120 | 0.2957 | **0.2424**[†] | **0.4335**[†] | 0.0758 | 0.1939 | **0.3908**[†] | **0.5642**[†] |
| (L4, WebAP)→ InsuranceQA | 0.1406 | 0.4267 | **0.1582** | **0.4717**[†] | 0.0489 | 0.1473 | **0.1622**[†] | **0.3059**[†] |

**Table 2: Performance across collections, where metrics under each collections represents the performance of the model trained on the opposing two collections. † represents significance against non adversarial model ($p < 0.05$, Wilcoxon test)**

## 5  RESULTS AND DISCUSSION

*Cross Topic.* Table 1 show the poor performance of the CosSim and Duet-distributed models on the four target topics when trained on the remaining collection. Notably, training on the topic specific data alone also performs poorly likely because of inadequate training data. However, in the presence of the adversarial discriminator both the models show significant improvement in performance on all held out topics. The improvements are somewhat bigger on the Duet-distributed baseline. We posit this is because the Duet-distributed model—with a deeper architecture—fits the training domain better at the cost of further loss in performance on the held out domains. Therefore, the adversarial learning has a stronger regularization opportunity on the Duet-distributed model.

*Cross Collection.* In similar vein as the cross topic evaluation, the incorporation of the adversarial signal significantly increases performance on the held out collections in Table 2. However, the difference in both size and distributional properties between these collections are far greater. Therefore, while the addition of the adversarial discriminator results in significant improvements—the absolute performance on the held out collections are still modest, even with adversarial regularization. We interpret these results as a reminder of the challenges in adapting these models to unseen domains.

## 6  CONCLUSION AND FUTURE WORK

The proposed adversarial approach to cross domain regularization shows significant performance improvements consistently under two evaluation settings (cross topic and cross collection) and over two different deep neural baselines. However, these improvements should be grounded in the realization that a model trained on large in-domain data is still likely to have a significant advantage over these models. Machine learning approaches to ad-hoc retrieval may need significantly more breakthroughs before achieving the level of robustness as some of the traditional retrieval models.

## 7  ACKNOWLEDGEMENTS

## REFERENCES

[1] Daniel Cohen and W. Bruce Croft. [n. d.]. End to End Long Short Term Memory Networks for Non-Factoid Question Answering. In *ICTIR '16*.
[2] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. 2016. Adversarial Feature Learning. *CoRR* abs/1605.09782 (2016).
[3] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. 2016. Domain-adversarial training of neural networks. *J. Mach. Learn. Res.* 17, 1 (2016), 2096–2030.
[4] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *NIPS 2014*. Curran Associates, Inc., 2672–2680. http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf
[5] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. 2016. A Deep Relevance Matching Model for Ad-hoc Retrieval. In *CIKM '16*. ACM, New York, NY, USA, 55–64. https://doi.org/10.1145/2983323.2983769
[6] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A. Efros, and Trevor Darrell. 2017. CyCADA: Cycle-Consistent Adversarial Domain Adaptation. *CoRR* abs/1711.03213 (2017).
[7] Mostafa Keikha, Jae Hyun Park, W. Bruce Croft, and Mark Sanderson. 2014. Retrieving Passages and Finding Answers. In *ADCS '14*. ACM, New York, NY, USA, Article 81, 4 pages. https://doi.org/10.1145/2682862.2682877
[8] Tie-Yan Liu et al. 2009. Learning to rank for information retrieval. *Foundations and Trends® in IR* 3, 3 (2009), 225–331.
[9] Bhaskar Mitra and Nick Craswell. 2018. An introduction to neural information retrieval. *Foundations and Trends® in IR (to appear)* (2018).
[10] Bhaskar Mitra, Fernando Diaz, and Nick Craswell. 2017. Learning to match using local and distributed representations of text for web search. In *WWW 17*. 1291–1299.
[11] Federico Nanni, Bhaskar Mitra, Matt Magnusson, and Laura Dietz. 2017. Benchmark for complex answer retrieval. In *Proc. ICTIR*. ACM, 293–296.
[12] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in IR* 3, 4 (2009), 333–389.
[13] Artem Rozantsev, Mathieu Salzmann, and Pascal Fua. 2016. Beyond Sharing Weights for Deep Domain Adaptation. *CoRR* abs/1603.06432 (2016).
[14] Aliaksei Severyn and Alessandro Moschitti. 2015. Learning to Rank Short Text Pairs with Convolutional Deep Neural Networks. In *SIGIR (SIGIR '15)*. ACM, New York, NY, USA, 373–382. https://doi.org/10.1145/2766462.2767738
[15] Mihai Surdeanu, Massimiliano Ciaramita, and Hugo Zaragoza. 2008. Learning to rank answers on large online QA collections. In *ACL:HLT*. 719–727.
[16] Ming Tan, Bing Xiang, and Bowen Zhou. 2015. LSTM-based Deep Learning Models for non-factoid answer selection. *CoRR* abs/1511.04108 (2015). http://arxiv.org/abs/1511.04108
[17] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. 2017. Adversarial discriminative domain adaptation. In *CVPR 17*, Vol. 1. 4.
[18] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. 2014. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474* (2014).
[19] Ellen M Voorhees, Donna K Harman, et al. 2005. *TREC: Experiment and evaluation in information retrieval*. Vol. 1. MIT press Cambridge.
[20] Hamed Zamani, Bhaskar Mitra, Xia Song, Nick Craswell, and Saurabh Tiwary. 2018. Neural ranking models with multiple document fields. In *Proc. WSDM*. ACM, 700–708.