

# Universal Schema for Entity Type Prediction

Limin Yao  
Dept. of Computer Science  
University of Massachusetts  
Amherst, MA, United States  
lmyao@cs.umass.edu

Sebastian Riedel  
Dept. of Computer Science  
University College London  
London, United Kingdom  
s.riedel@ucl.ac.uk

Andrew McCallum  
Dept. of Computer Science  
University of Massachusetts  
Amherst, MA, United States  
mccallum@cs.umass.edu

## ABSTRACT

Categorizing entities by their type is useful in many applications, such as knowledge base construction, relation extraction and query intent prediction. Fine-grained entity type ontologies are especially valuable, but typically difficult to design because of endless quandaries about level of detail and boundary cases. Automatically classifying entities by type is challenging as well, usually involving hand-labeling data and training a supervised predictor. This paper presents a *universal schema* approach to fine-grained entity type prediction. The set of types is taken as the union of textual surface patterns (e.g. appositives) and pre-defined types from available databases (e.g. Freebase)—yielding not tens or hundreds of types, but tens of thousands of entity types, such as financier, criminologist, and musical trio. We robustly learn mutual implicature among this large union by matrix completion using embeddings learned from probabilistic matrix factorization, thus avoiding the need for hand-labeled data. Experimental results demonstrate more than 30% reduction in error versus a traditional classification approach on predicting fine-grained entities types.

## 1. INTRODUCTION

Classifying entities into different categories is a common task in many NLP systems. In some cases, such as knowledge base construction, entity types may be a prominent user-visible feature [6, 18]. In others, such as relation extraction [36, 27] or query intent discovery [1] entity types are hidden variables included to improve accuracy on the target task. Occasionally the ontology of entity types is coarse, such as the four types in the CoNLL-2003 shared task (person, organization, location and miscellaneous), but often finer-grained ontologies are more useful. For example, specializations of people, including politician, scientist, and athlete are defined in [16, 17, 12]. Others are even more detailed; for instance the Unified Medical Language System (UMLS) defines an ontology of 987,321 biomedical concepts.

Defining such ontologies is a significant challenge, often giv-

ing rise to debates about desired granularity and subtle questions about boundary cases. These difficulties appear both when the assignment of entities to types is exclusive and when it is one-to-many.

Once the ontology is defined, the problem of building the automated classification system remains. The most common approach is supervised training from a set of entity mentions labeled into the ontology [16, 34]. However labeling such data is difficult—especially with fine-grained ontologies. Furthermore, when the ontology evolves or expands (as it often does), the data labeling must be re-visited. Even when used as hidden variables, the set of entity types may warrant adjustment because an ontology tuned to the task at hand typically performs better—for example the authors of [24] show that the entity types in the WordNet ontology [14] is not as effective as those derived from automatic clustering for the task of learning selectional preferences. Unsupervised clustering may also be employed to derive entity types [35, 13, 24], but the resulting types often have peculiar, undesirable boundary and granularity choices.

This paper presents an approach to fine-grained entity type classification that avoids the need to hand-design an ontology, avoids the need for labeled data, and avoids the boundary difficulties that arise from forcing our semantics into finite, pre-defined, somewhat arbitrary “boxes”. We accomplish this by adopting the *universal schema* approach, which is previously applied to relation extraction [26, 37], and extending it to entity types. In “universal schema”, our types are the *union* of all available types from all input sources, including multiple pre-existing ontologies and naturally-occurring textual surface-form expressions that indicate entity type, such as appositives, isa-expressions, or even adjectival or verb phrases. For example, “James Cameron” may appear as a *person/director* in Freebase, as a *PERSON* in TAC/KBP, and as a *movie-mogul, Canadian citizen, and jerk* in various appositives in available text. Rather than five, fifty or five-hundred entity types, this universal schema approach typically yields tens-of-thousands of entity types (particularly from the textual surface forms). It does not force the natural diversity and ambiguity of the original input types into a smaller set of types.

The key characteristic of universal schema is that it models directed implicature among the many candidate types of an entity by casting the problem as a large matrix completion task. Each row in the matrix corresponds to an entity;

each column an entity type; some cells of the matrix are observed and marked true; many are unobserved; it is the job of matrix completion to “fill in” the matrix, marking the unobserved cells as either true or false. For example, although we may not have directly observed that “Barack Obama” is a *leader*, our model will infer it by having observed that he is a *president* and *commander-in-chief*—doing so by leveraging various patterns of co-occurrences among these types in other entities. Similarly it will infer that he is not a *movie-mogul* or *masterpiece*. As in our previous work, we perform this matrix completion task using probabilistic matrix factorization—efficiently estimating vector embeddings for both entities and types by online stochastic gradient descent optimization. The confidence in an entity’s type assignment is determined by the dot-product of the corresponding embeddings, mapped through a logistic function.

Having so many entity types, including types appearing in natural language, allows users to query our system in natural language. That is, rather than having to learn an idiosyncratic ontology, users can ask about entity types in their own vocabulary, and we will most likely already have a column to match. The large number of entity types does make evaluation a challenge. We cannot evaluate every cell in the matrix. Thus we choose to evaluate a subset of the columns (entity types) on a closed set of entities which we have annotated. Here our approach achieves a 15% absolute increase in F1 versus the traditional classification method. Furthermore, although it does not leverage the diversity of universal schema, we also compare against a baseline method for Freebase type prediction [4]; here we achieve similar results as the baseline. In spite of the large number of types, training our system is still efficient, taking approximately 6 hours on one machine for 100 iterations, 100 components on about 503K entities and 16K types.

## 2. FACTORIZATION MODELS

We present a matrix factorization model to collectively learn semantic implication among unary relations and predict new relations for entities. We fill a matrix  $E \times R$  with unary relation instances, where  $E$  corresponds to entities and  $R$  to unary relations. Assume we index an entity with  $e$  and a relation with  $r$ . Each matrix cell is a binary variable, denoted as  $x_{e,r}$ . The variable is 1 when relation  $r$  holds for entity  $e$ , and 0 otherwise. For example, observing “directed by Marzieh Meshkini”, we fill the corresponding cell (Marzieh Meshkini, directed by X) with 1.

In our matrix factorization approach we associate each entity and relation with latent vectors  $\mathbf{a}_e$  and  $\mathbf{v}_r$  in a  $K$ -dimensional space, respectively. The dot-product  $\theta_{e,r} = \sum_c a_{e,c} v_{r,c}$  of these vectors for a given entity  $e$  and unary relation  $r$  then becomes the natural parameters of a Bernoulli distribution that generates the observed binary data. That is, the probability of  $x_{e,r} = 1$  is given by  $\sigma(\theta_{e,r})$  where  $\sigma$  is the sigmoid function. This model corresponds to an instantiation of generalized PCA [9].

To learn the low dimensional latent vectors we maximize the log likelihood of the observed cells under the probabilistic model above. Notice that in our training data we only observe positive cells and have no accurate data on which re-

lations do *not* hold for an entity. However, learning requires negative training data. We address this issue by sampling negative relations for an entity based on their frequencies in the whole dataset. In our experiments we use stochastic gradient optimization to effectively deal with the large scale of our matrices. In each iteration, we traverse random permutations of all training cells, randomly sample some negative cells for each training cell, and update the corresponding  $\mathbf{a}_e$  and  $\mathbf{v}_r$  vectors for the positive and negative cells based on their gradients (omitted here for brevity).

## 3. EXPERIMENTS

We extract unary relations from New York Times data for the years 1990 to 2007 [28]. We preprocess the documents by performing NER tagging [15] and dependency parsing [23]. Following [33], we extract dependency paths originating from a (named) entity mention as unary relations. Specifically, we traverse from the head token of the entity mention to the root of the dependency tree. Whenever we come across a content word (nouns, adjectives etc.), the current (lexicalized) path from the entity mention to this content word node is used as one unary relation. We stop when approaching a verb or a clause boundary. Additionally, when a verb is encountered, other modifiers of the verb are also included in the path. For example, we can have “X buy share”, “X roll over”. This yields many relations that could serve as entity types, including appositive structures. For example, the unary relation “X, a magnate” can define “magnate” as the corresponding entity type.

In training our matrix factorization model, we set the regularizers  $\lambda$ s for both entity and unary relation vectors as 0.02. We use 100 components, and run 100 iterations using stochastic gradient optimization. We experiment with different configurations and the current one results in the best performance.

Our task is to predict missing cells in the matrix. In the following, we design experiments to measure the accuracies of these predictions.

### 3.1 Pattern-based Evaluation

Our universal schema approach typically yields tens-of-thousands of entity types, particularly from the textual surface patterns. In this section, we demonstrate our predictions on these unary relations. Since there is no ground truth for these patterns (other than a subsample of positive-only cells), we cannot easily evaluate them. Instead, we query some of them and list the top ranked entities according to our model and the baseline. The set of queried unary relations in our experiments consists of patterns based on appositives. For example, “X, a scientist”, “X, an actor”, and “X, a band”. Intuitively these relations are most directly corresponding to entity types, *scientist*, *actor*, and *band* in this case. We ask human annotators to annotate each returned entity for a given pattern-based relation. The annotators are provided with sentences in which the entities are mentioned.

Entities from NYT articles are split into training and test set: 355,942 entities vs 147,359 entities. In total the input matrix has about 500K rows, 16K columns. When training the model we hide all query patterns in the test set. We compare our approach against a binary classifier that con-

Query	Univ	ME
politician	<b>0.738</b>	0.448
scientist	<b>0.499</b>	0.354
magnate	0.433	<b>0.460</b>
band	0.413	<b>0.427</b>
reporter	<b>0.437</b>	0.330
actor	<b>0.649</b>	0.518
player	<b>0.840</b>	0.711
magazine	<b>0.845</b>	0.675
Micro	<b>0.701</b>	0.557
Macro	<b>0.607</b>	0.490

**Table 1: F1 measure on 8 patterns of different approaches. Our approach (Univ) achieves significant better performance on almost of all these patterns.**

siders entities co-occurring with the query pattern as positive examples and all others as negative examples. As the classification model we use maximum entropy (ME).

Entities in the test set are selected for annotation by the following rules. Entities are ranked with respect to each pattern by our system and the baseline system separately. Top 100 entities of each target pattern ranked by each system are shown to the annotators. We also acquire annotations from Freebase. For example, we consider all entities that have labels *politician*, *us.congressperson*, *us.senator*, *us.vicepresident* as instances of our target pattern “X, a politician”. Similarly we obtain entities for target patterns “X, a player”, and “X, an actor”. Mappings from Freebase labels to these three patterns are from [20]. In total, we have annotations for 14,991 entities in the test set.

We measure precision, recall and F1 for entities in this set. For each entity, patterns with probabilities above a threshold are considered as true. The threshold is 0.5. In scenarios where an entity has no patterns above the threshold, the top ranked one is selected. This may lower the precision of each system, however, it does increase the recall and F1 score for both our approach and the baseline.

Table 1 lists the F1 measures for each pattern. We can see that our approach performs significantly better than the baseline on 6 patterns. We perform slightly worse on two patterns. On micro average, we gain about 15% in F1 score. When analyzing the errors made by the baseline system we see most problems when there are no patterns above the threshold. In these cases the baseline’s top ranked patterns (now with score under the threshold) are mostly incorrect. However, for our model the top ranked patterns, when under the threshold, are still often correct.

To interpret the embeddings of the unary relations, we perform hierarchical agglomerative clustering on relation vectors. Some example patterns that are in the same cluster with each query pattern are listed in Table 2. We can tell that our approach can learn diverse and accurate patterns that are indicative of the target patterns.

### 3.2 Closed Set Evaluation

System	Univ	ME	UW
Micro	0.515	0.501	0.553
Macro	0.172	0.120	0.180

**Table 3: F1 scores of different approaches.**

Our framework can also incorporate entity types from ontologies as columns. To make comparison against traditional approaches for entity type classification, we evaluate our predictions on pre-defined entity types as well. Specifically, we label entities occurring in a set of held-out with labels from a predefined set of types in Freebase. In this dataset all entities are annotated with all possible relations exhaustively, and this enables us to measure precision, recall and F1. We choose to compare against UW [20], in which the authors employ a multi-class multi-label classifier for fine-grained entity recognition. They test their model on a dataset of 18 documents and about 430 sentences. In this data set, entity mentions are annotated and each mention is labeled with all possible entity types. Our approach is only concerned about the types of entities (sets of entity mentions), not entity mentions. We therefore collapse entity mentions that have the same surface string into single entities. As the UW works on an entity mention basis, we aggregate their output for entity mentions of the same entity. Again we also compare against maximum entropy (ME) classifiers that are trained on labels of entities instead of mentions. Here we have one binary classifier per unary relation, and hence the same entity can have several relations.

Notice that our approach (Univ) and the ME baseline use the same training data as that is used in UW [20]. The input matrix has about 623K entities and 724K columns, 12M positive cells and 36M sampled negative cells. It takes about 10 hours to train the model for 100 iterations using 100 components.

Table 3 shows the F1 scores of different systems. The macro average numbers are small due to that many types get 0.0 in F1. Our approach is better than ME, slightly worse than UW. Looking closely at different entity types, we find that our approach is better at predicting “sports.league”, “athlete”, “person/coach”, “education/department”; It is worse at predicting “location/city”, “location/province”, “location/country”, and “news.agency”. Our model has higher recall but sometimes suffers from lower precision for fine-grained types that have fewer training instances, such as “news.agency”.

## 4. RELATED WORK

**Universal Schema.** The authors previously introduced *universal schema* for relation extraction [26, 37]. There the rows correspond to entity-pairs and columns correspond to relation types. A special three-part parameterization for matrix factorization is employed to complete the matrix. In this paper we extend the universal schema approach to entity types. This represents a first step towards future work in joint entity-type and relation-type prediction, both with universal schema and matrix factorization.

**Low Dimensional Embedding.** Learning low dimensional embeddings of high-dimensional NLP data has been

Query	Patterns
politician	legislator, official, politician like, vote, campaign of, criticism from, ally, accuse of, defeat by, election, lash out, lawmaker, whip, endorse, oppose by, run against, re-elected, opponent, conservative
scientist	criminologist, biologist, psychologist, sociologist, professor, researcher, neuroscientist, co-author with, ecologist, expert, physicist
band	tune, album, song by, act like, 's singer, hit for, 's song, country, concert, singer of, wing with, music of, tour with, musician like, trio, duo, blues, sound of, rock, folk, pop, recording by,

**Table 2: Top similar patterns to the target queries. We list words on the patterns for simplicity, adding a placeholder X when necessary.**

of both long-standing [5, 2, 10, 3] and rising interest [32, 21]. Much of this work has been for the embedding of individual words [5, 2, 3, 21]. Some has been for structured natural language processing, such as part-of-speech tags, chunks, named entity tags and semantic roles [10], or for parsing [32]. Some recent work uses tensor factorization for embedding relations from triples of semantic role labeling (subj-verb-obj) [19], from a structured knowledge base [22] or from WordNet relations [8]. Our work is the first of which to predict “open domain” universal schema for relations or entity types by leveraging natural language inputs.

**Entailment.** Our work is also related to semantic inference over text [11], in which given a hypothesis and textual evidence, the system must predict whether or not the text entails the hypothesis. In our framework, observed cells of the matrix are the evidence, and newly predicted cells are the hypotheses. Szpektor and Dagan [33] aim to discover implications among unary patterns for predicting new unary facts. We have a similar goal here. They concentrate on verb-triggered patterns, whereas we focus on patterns that define entity types, including noun-triggered patterns (such as appositives) and verb-triggered patterns. They employ distributional similarity where we use matrix factorization. Other work addresses semantic inference over relation instances, for example, by learning rules that conjoin textual patterns extracted by OpenIE [29, 30].

**Fine-Grained Entity Type Classification.** Classifying entities into large ontologies is a commonly tackled task and is widely acknowledged as useful, not only for knowledge base construction, but also query log prediction [25]. Some researchers have explored entity type classification specifically for categories of people [16, 17, 12]. Others have defined ontologies with a wider variety of entity types, but not implemented methods to automatically classify instances into the ontology [31]. Large-scale knowledge bases, such as Freebase and its fine-grained entity types, have significant collections of entities that can be used for training traditional classification methods by distant supervision [20]. Others have also performed entity type classification with a multi-label classifier in a hierarchy of types [38]. The main differences to our approach are (1) that we use matrix factorization rather than classification as the framework for our model, and more importantly (2) we do not restrict ourselves to predefined entity types, instead leveraging the wide diversity of naturally available data. Even when a pre-existing knowledge base can provide supervision for a classifier, the

resulting entity type classifier is still limited by the types envisioned by the creators of the knowledge base ontology. Furthermore, note that even when the goal is merely classification into a specific ontology, matrix factorization’s striving to predict many other text-based entity types provides a kind of multi-task learning [7] that can be beneficial.

## 5. CONCLUSION

This paper has presented *universal schema* for assigning entities into multiple of over 16,000 entity types on NYT data. We use the term “universal” because the set of types is formed by the union of textual surface patterns and multiple input entity type ontologies. We find that our approach, based on matrix factorization, reduces error by more than 30% in comparison with a traditional classifier on a set of entities from the widely-diverse textual appositive-derived entity types.

There are significant opportunities for future work. We have begun to integrate this paper’s entity type model with our previous relation type model, and we expect further increases in accuracy to arise from learning these two jointly. We also plan to explore new strategies for obtaining negative training signal and for integrating more observed features of the entity mentions.

## 6. REFERENCES

- [1] K. Balog and R. Neumayer. Hierarchical target type identification for entity-oriented queries. In *Proceedings of CIKM*, 2012.
- [2] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin. A neural probabilistic language model. *Journal of Machine Learning Research*, 3:1137–1155, 2003.
- [3] J. Blitzer, K. Q. Weinberger, L. K. Saul, and F. C. N. Pereira. Hierarchical distributed representations for statistical language modeling. In *Proceedings of NIPS*, 2004.
- [4] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD ’08: Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250, New York, NY, USA, 2008. ACM.
- [5] P. F. Brown, P. V. deSouza, R. L. Mercer, V. J. D. Pietra, and J. C. Lai. Class-based n-gram models of natural language. *Computational Linguistics*, 18(4):467–479, 1992.

- [6] A. Carlson, J. Betteridge, R. Wang, E. Hruschka, and T. Mitchell. Coupled semi-supervised learning for information extraction. In *Third ACM International Conference on Web Search and Data Mining (WSDM '10)*, 2010.
- [7] R. A. Caruana. Multitask learning: A knowledge-based source of inductive bias. In *Proceedings of ICML*, 1993.
- [8] D. Chen, R. Socher, C. D. Manning, and A. Y. Ng. Learning new facts from knowledge bases with neural tensor networks and semantic word vectors. In *Proceedings of workshop at ICLR*, 2013.
- [9] M. Collins, S. Dasgupta, and R. E. Schapire. A generalization of principal component analysis to the exponential family. In *Proceedings of NIPS*, 2001.
- [10] R. Collobert and J. Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of ICML*, 2008.
- [11] I. Dagan, O. Glickman, and B. Magnini. The pascal recognising textual entailment challenge. In *Proceedings of the PASCAL Challenges Workshop on Recognising Textual Entailment*, 2005.
- [12] A. Ekbal, E. Sourjikova, A. Frank, and S. P. Ponzetto. Assessing the challenge of fine-grained named entity recognition and classification. In *Named Entities Workshop, ACL*, 2010.
- [13] M. Elsner, E. Charniak, and M. Johnson. Structured generative models for unsupervised named-entity clustering. In *Proceedings of NAACL*, 2009.
- [14] C. Fellbaum. *WordNet: An Electronic Lexical Database*. Bradford Books, 1998.
- [15] J. R. Finkel, T. Grenager, and C. Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL '05)*, pages 363–370, June 2005.
- [16] M. Fleischman and E. Hovy. Fine grained classification of named entities. In *Proceedings of Coling*, 2002.
- [17] C. Giuliano and A. Gliozzo. Instance-based ontology population exploiting named-entity substitution. In *Proceedings of Coling*, 2008.
- [18] J. Hoffart, F. M. Suchanek, K. Berberich, and G. Weikum. Yago2: A spatially and temporally enhanced knowledge base from wikipedia. *Artificial Intelligence*, 2012.
- [19] R. Jenatton, N. L. Roux, A. Bordes, and G. Obozinski. A latent factor model for highly multi-relational data. In *Proceedings of NIPS*, 2012.
- [20] X. Ling and D. S. Weld. Fine-grained entity recognition. In *Proceedings of AAAI*, 2012.
- [21] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representation in vector space. In *Proceedings of workshop at ICLR*, 2013.
- [22] M. Nickel, V. Tresp, and H.-P. Kriegel. Factorizing yago: Scalable machine learning for linked data. In *Proceedings of WWW*, 2012.
- [23] J. Nivre, J. Hall, and J. Nilsson. Memory-based dependency parsing. In *Proceedings of CoNLL*, pages 49–56, 2004.
- [24] P. Pantel, R. Bhagat, B. Coppola, T. Chklovski, and E. Hovy. ISP: Learning Inferential Selectional Preferences. In *Proceedings of NAACL HLT*, 2007.
- [25] P. Pantel, T. Lin, and M. Gamon. Mining entity types from query logs via user intent modeling. In *Proceedings of ACL*, 2012.
- [26] S. Riedel, L. Yao, A. McCallum, and B. Marlin. Relation extraction with matrix factorization and universal schemas. In *Proceedings of NAACL*, 2013.
- [27] D. Roth and W. tau Yih. Global inference for entity and relation identification via a linear programming formulation. 2007.
- [28] E. Sandhaus. *The New York Times Annotated Corpus*. Linguistic Data Consortium, Philadelphia, 2008.
- [29] S. Schoenmackers, O. Etzioni, and D. S. Weld. Scaling textual inference to the web. In *EMNLP '08: Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 79–88, Morristown, NJ, USA, 2008. Association for Computational Linguistics.
- [30] S. Schoenmackers, O. Etzioni, D. S. Weld, and J. Davis. Learning first-order horn clauses from web text. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, EMNLP '10*, pages 1088–1098, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.
- [31] S. Sekine. Extended named entity ontology with attribute information. In *LREC*, 2008.
- [32] R. Socher, C. D. Manning, and A. Y. Ng. Learning continuous phrase representations and syntactic parsing with recursive neural networks. In *NIPS*, 2010.
- [33] I. Szpektor and I. Dagan. Learning entailment rules for unary templates. In *Proceedings of Coling*, 2008.
- [34] H. Tanev and B. Magnini. Weakly supervised approaches for ontology population. In *Proceedings of 11st Conference of the European Chapter of the Association for Computational Linguistics*, 2006.
- [35] L. Yao, A. Haghighi, S. Riedel, and A. McCallum. Structured relation discovery using generative models. In *Proceedings of EMNLP*, 2011.
- [36] L. Yao, S. Riedel, and A. McCallum. Collective cross-document relation extraction without labelled data. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 1013–1023, Cambridge, MA, October 2010. Association for Computational Linguistics.
- [37] L. Yao, S. Riedel, and A. McCallum. Probabilistic databases of universal schema. In *Proceedings of the AKBC-WEKEX Workshop at NAACL*, 2012.
- [38] M. A. Yosef, S. Bauer, J. Hoffart, M. Spaniol, and G. Weikum. Hyena: Hierarchical type classification for entity names. In *Proceedings of Coling*, 2012.