

Joint Structured Models for Extraction from Overlapping Sources

Rahul Gupta
IIT Bombay, India
grahul@cse.iitb.ac.in

Sunita Sarawagi
IIT Bombay, India
sunita@cse.iitb.ac.in

Abstract

We consider the problem of jointly training structured models for extraction from sources whose instances enjoy partial overlap. This has important applications like user-driven ad-hoc information extraction on the web. Such applications present new challenges in terms of the number of sources and their arbitrary pattern of overlap not seen by earlier collective training schemes applied on two sources. We present an agreement-based learning framework and alternatives within it to trade-off tractability, robustness to noise, and extent of agreement. We provide a principled scheme to discover low-noise agreement sets in unlabeled data across the sources. Through extensive experiments over 58 real datasets, we establish that our method of additively rewarding agreement over maximal segments of text provides the best trade-offs, and also scores over alternatives such as collective inference, staged training, and multi-view learning.

1 Introduction

This paper addresses the problem of training multiple structured prediction models that share an output space but differ in their input data and feature space. Further, labeled data in each source is limited, but unlabeled data over the different sources overlap partially. This scenario is applicable in many text modeling tasks such as information extraction, dependency parsing, and word alignment. These tasks are increasingly being deployed in settings where supervision is limited but redundancy is abundant. A concrete motivation for our work comes from recent efforts to support rich forms of structured query-answering on the Web [5, 4]. A typical subtask here is building extraction models over multiple Web documents starting from a small seed of user-provided structured records.

Recently, many learning paradigms have been proposed to exploit the relatedness of multiple models. On one end of the spectrum we have collective inference [21, 3, 15, 8, 12] where each model is trained independently but prediction happens jointly to encourage agreement on overlapping content. On the other end are methods like multi-view learning [11, 9] and agreement-based learning [17, 18] that formulate a single objective to jointly train all models. Then there are methods in-between that train models sequentially or alternately [2, 5]. Our problem is different from traditional multi-view learning where multiple models are trained on different views of a *single* data source. However, by treating the different contexts in which each shared portion resides as a different view, we can apply multi-view learning to this problem. We elaborate on this and other alternatives in Section 4.

In agreement-based learning [17, 18] the goal is to train multiple models so as to maximize the likelihood of the labels agreeing on shared variables. However, these assume that all models need to agree on the same set of variables — this trivially holds for two sources where these methods have been applied. In our application the number of sources is often as large as 20. As number of sources increase, there is a bewildering number of ways in which they overlap. This makes it challenging to devise objectives that maximally exploit the overlap while accounting for noise in the agreement set and intractability of training. We are aware of no study where such issues are addressed in the context of jointly training more than two sources with partial overlap.

In this paper we propose an agreement-based model for training multiple structured models with arbitrary partial overlap among the sources. We propose several alternatives for enforcing agreement ranging from singleton variables, to groups of contiguous variables, to global models that lead to giant agreement graphs. For the task of information extraction, we present a strategy for selecting the unit of agreement that leads to a significant reduction in the noise in the agreement set compared to the existing naïve approach for choosing agreement sets. We present an extensive evaluation on 58 real-life collective extraction tasks covering a rich spectrum of

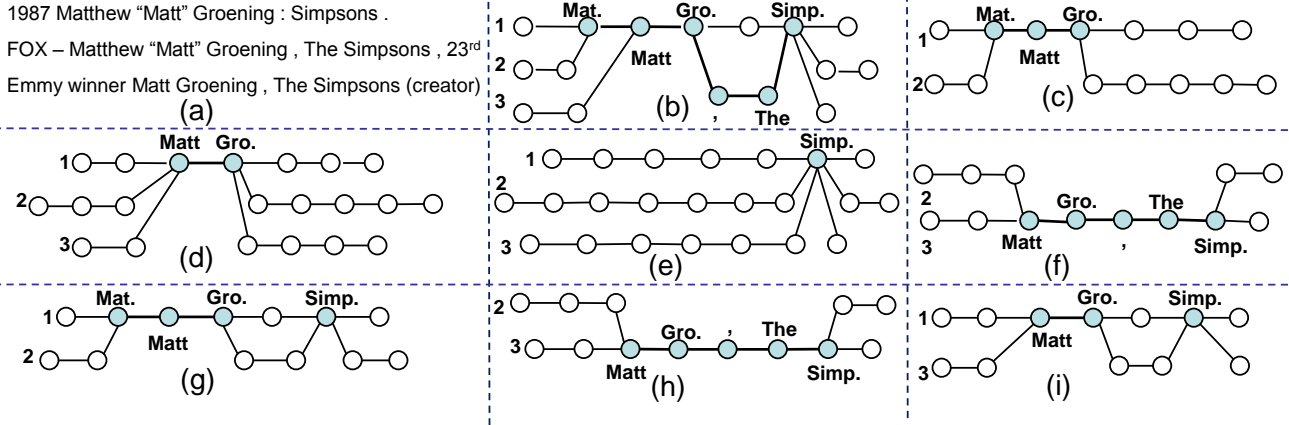


Figure 1: (a) Three samples sentences with $\mathcal{A}=\{(\text{Matthew Matt Groening}), (\text{Matt Groening}), (\text{Matt Groening , The Simpsons}), (\text{Simpsons})\}$ (b) The fused graph (c)-(f) Clique Agreement approximation (g)-(i) Instance Pair approximation.

data characteristics. This study reveals that agreement objectives that are additive over smaller components provide the best accuracy because of robustness against noise in the agreement term, while providing a tractable inference objective.

2 Collective training

Our goal is to collectively train S structured prediction models, where every source $s \in S$ comes with a small set L_s of labeled instances and many unlabeled instances U_s . The unlabeled instances from different sources share overlapping content and we seek to exploit this for better training. For extraction tasks, this overlap can be anywhere from the level of unigrams to non-contiguous segments, as illustrated in Figure 1(a). An overlap finding algorithm identifies such shared parts as an agreement set \mathcal{A} comprising of a set of cliques. Each clique $C \in \mathcal{A}$ contains a list of triples (s, i, r) indicating the source s , the instance $i \in s$, and the part $r \in i$ which has the same content as the other members of the clique. We do not make any assumption of mutual exclusion between cliques and in general each clique can span a variable number of members and token positions. In Section 3 we present our strategy for computing the agreement set.

As in standard structured learning, we define a probabilistic model $P_s(\mathbf{y}|\mathbf{x})$ for each source s using its feature vector $\mathbf{f}_s(\mathbf{x}, \mathbf{y})$ and parameters \mathbf{w}_s as:

$$P_s(\mathbf{y}|\mathbf{x}, \mathbf{w}_s) = \frac{1}{Z(\mathbf{x}, \mathbf{w}_s)} \exp(\mathbf{w}_s \cdot \mathbf{f}_s(\mathbf{x}, \mathbf{y})) \quad (1)$$

The traditional goal of training is to find a \mathbf{w}_s that maximizes the regularized likelihood of the labeled set in source s :

$$\text{LL}_s(L_s, \mathbf{w}_s) = \sum_{(\mathbf{x}, \mathbf{y}) \in L_s} \log P_s(\mathbf{y}|\mathbf{x}, \mathbf{w}_s) - \gamma R(\mathbf{w}_s) \quad (2)$$

We propose to augment this base objective with the likelihood that the S models agree on the labels of cliques in the agreement set \mathcal{A} . We first observe that the joint distribution over the labels \mathbf{Y} of all instances \mathbf{X} spanning all sources is:

$$P(\mathbf{Y}|\mathbf{X}, \mathbf{W}) \triangleq \prod_{s \in S} \prod_{i \in s} P(\mathbf{Y}_{si}|\mathbf{X}_{si}, \mathbf{w}_s) \quad (3)$$

where \mathbf{W} denotes $(\mathbf{w}_1, \dots, \mathbf{w}_S)$, \mathbf{X}_{si} represents instance i in source s , and \mathbf{Y}_{si} is a random variable for the structured output of this instance.

Now given an agreement set \mathcal{A} , consider the subset $\mathcal{Y}_{\mathcal{A}}$ of all possible labelings that are consistent with it:

$$\mathcal{Y}_{\mathcal{A}} \triangleq \{\mathbf{Y} : \forall C \in \mathcal{A}, (s, i, r), (s', i', r') \in C : \mathbf{Y}_{sir} = \mathbf{Y}_{s'i'r'}\} \quad (4)$$

The log likelihood of the agreement set is then:

$$\text{LL}(\mathcal{Y}_A, \mathbf{W}) \triangleq \log \Pr(\mathcal{Y}_A) = \log \sum_{\mathbf{Y} \in \mathcal{Y}_A} \prod_{s,i \in \mathcal{S}} P(\mathbf{Y}_{si} | \mathbf{X}_{si}, \mathbf{w}_s) \quad (5)$$

In the rest of the paper we use the short form (s, i) to denote the instance i in source s . Our goal now is to jointly train $\mathbf{w}_1, \dots, \mathbf{w}_S$ so as to maximize a weighted combination of the likelihoods of the labeled and agreement sets:

$$\max_{\mathbf{w}_1, \dots, \mathbf{w}_S} \sum_s \text{LL}(L_s, \mathbf{w}_s) + \lambda \text{LL}(\mathcal{Y}_A, \mathbf{W}) \quad (6)$$

2.1 Computing $\text{LL}(\mathcal{Y}_A, \mathbf{W})$

Using Equations 1, 3, and 5, we rewrite $\text{LL}(\mathcal{Y}_A, \mathbf{W})$ as

$$\text{LL}(\mathcal{Y}_A, \mathbf{W}) = \log \sum_{\mathbf{Y} \in \mathcal{Y}_A} \exp\left(\sum_{s,i} \mathbf{w}_s \mathbf{f}_s(\mathbf{Y}_{si}, \mathbf{w}_s)\right) - \sum_{s,i} \log Z(\mathbf{X}_{si}, \mathbf{w}_s) \quad (7)$$

The second part in this equation is the sum of the log partition function over individual instances which can be computed efficiently as long as the base models are tractable.

The first part is equal to the log partition function of a fused graphical model G_A constructed as follows: Initially, each instance (s, i) creates a graph G_{si} corresponding to its model $P_s(\cdot)$. For text tasks, this would typically be a chain model with a node for each token position. Next, for each clique $C \in \mathcal{A}$, and for each pair of triples $(s, i, r), (s', i', r') \in C$, we collapse the nodes of r in G_{si} with the corresponding nodes of r' in $G_{s'i'}$. In Figure 1(b) we show an example of such a fused graph created from the three instances of Figure 1(a) with four cliques in their agreement set.

Let K be the number of nodes in the final fused graph and z_1, \dots, z_K denote the node variables. Every node j in the initial graph G_{si} is now mapped to some final node $k \in 1, \dots, K$, and we denote this mapping by $\pi(s, i, j)$. The log-potential for a component \mathbf{c} in the fused graph is simply an aggregate of the log-potentials of the members \mathbf{c}' that collapsed onto it.

$$\theta_{\mathbf{c}}(\mathbf{z}_{\mathbf{c}}) \triangleq \sum_{(s,i,c'):\pi(s,i,c')=\mathbf{c}} \mathbf{w}_s \mathbf{f}_s(\mathbf{z}_{\mathbf{c}'}, \mathbf{X}_{si}, \mathbf{c}') \quad (8)$$

where we extend π to operate on node-sets as well. The above θ parameters now define a distribution over the fused variables z_1, \dots, z_K as follows:

$$P_A(\mathbf{z}|\theta) = \frac{1}{Z_A(\theta)} \exp\left(\sum_{\mathbf{c}} \theta_{\mathbf{c}}(\mathbf{z}_{\mathbf{c}})\right) \quad (9)$$

It is easy to see that the log partition function of this distribution is the same as the first term of Equation 7, so we can work with G_A instead. If the set of cliques in \mathcal{A} is such that the fused graph G_A has a small tree width, we can compute the log partition function $\log Z_A(\theta)$ efficiently. In other cases, we need to approximate the term in various ways. We discuss several such approximations in Section 2.3.

2.2 Training algorithm

The overall training objective of Equation 6 is not necessarily concave in \mathbf{w}_s because of the agreement term with sums within a log. As in [17, 11] it is easy to derive a variational approximation with extra variables to be solved using an EM algorithm. EM will give a local optima if the marginals of the P_A distribution can be computed exactly. Since this cannot be guaranteed for general fused graphs, we also explore the simpler approach of gradient ascent. In Section 5 we show that gradient ascent achieves better accuracy than EM while being faster. The gradient of $\text{LL}(\mathbf{Y}_A, \mathbf{W})$ is

$$\nabla \text{LL}(\mathbf{Y}_A, \mathbf{W}) = \sum_{s,i,c} \sum_{\mathbf{y}_c} (\mu_{\mathcal{A},\pi(s,i,c)}(\mathbf{y}_c) - \mu_{s,c}(\mathbf{y}_c | \mathbf{X}_{si})) \mathbf{f}_s(\mathbf{X}_{si}, \mathbf{y}_c, c)$$

where we use the notation $\mu_{s,c}, \mu_{\mathcal{A},c'}$ to denote the marginal probability at c of P_s and c' of P_A respectively. Note that the E-step of EM requires the computation of the same kind of marginal variables. These are computed using the same inference algorithms as used to compute the log partition function and we discuss the various options next.

2.3 Approximations

We explore two categories of approximations for training when $\text{LL}(\mathcal{A}, \mathbf{W})$ is intractable.

2.3.1 Partitioning the agreement set

The first category is based on approximating the $\text{Pr}(\mathcal{Y}_{\mathcal{A}})$ distribution with product of simpler distributions obtained by partitioning the set \mathcal{A} . We partition the agreement set \mathcal{A} into smaller subsets $\mathcal{A}_1, \dots, \mathcal{A}_R$ such that each $\text{Pr}(\mathcal{Y}_{\mathcal{A}_k})$ is easy to compute and $\bigcap_k \mathcal{Y}_{\mathcal{A}_k} = \mathcal{Y}_{\mathcal{A}}$. We then approximate $\text{Pr}(\mathcal{Y}_{\mathcal{A}})$ by $\prod_k \text{Pr}(\mathcal{Y}_{\mathcal{A}_k})$, thus replacing the corresponding log-likelihood term by $\sum_k \text{LL}(\mathcal{A}_k, \mathbf{W})$. We explore three such partitionings:

Clique Agreement In this case we have one partition per clique $C \in \mathcal{A}$. $G_{\mathcal{A}}$ now decomposes into several simpler graphs, where a simple graph has its nodes fused only at one clique. Figures 1(c)-(f) illustrate this decomposition for the fused model of Figure 1(b). The probability $\text{Pr}(\mathcal{Y}_{\{C\}})$ of agreement on members of a single clique C simplifies to

$$\text{Pr}(\mathcal{Y}_{\{C\}}) = \sum_{\mathbf{y} \in \mathbf{Y}_C} \prod_{(s,i,r) \in C} P_s(\mathbf{Y}_{sir} = \mathbf{y} | \mathbf{X}_{si}) \quad (10)$$

where \mathbf{Y}_C is set of all possible labelings for any member of C , and $P_s(\mathbf{Y}_{sir} = \mathbf{y})$ is the marginal probability of the part r taking the labeling \mathbf{y} under P_s .

This approximation is useful for two reasons. First, if the base models are sequences (e.g. in typical extraction tasks) and clique parts r are over contiguous positions in the sequence, the fused graph of $\text{Pr}(\mathcal{Y}_{\{C\}})$ is always a tree, such as the ones in Figures 1(c)-(f). Second, since for trees we can use sum-product to compute $\text{Pr}(\mathcal{Y}_{\{C\}})$ instead of Equation 10, we can now use arbitrarily long cliques, instead of choosing unigram cliques which is typically the norm in extraction applications.

Node Agreement We also consider a special case of the clique agreement approximation, called node agreement, in which each partition corresponds to agreement over a single variable as in Figure 1(e).

Instance Pair Agreement Another decomposition is based on picking pairs of instances and defining an agreement set on all cliques which they share. For the example in Figure 1, graphs marked (g),(h),and (i) demonstrate the fused graphs arising out of instance pair agreement. This scheme is expected to be useful when base models exhibit strong edge potentials. However, unlike for the above two decompositions, there is no guarantee that the fused graph is a tree (e.g. graph (g)). So, approximate inference may be required for some pairs.

2.3.2 Approximating $Z_{\mathcal{A}}(\theta)$

An alternate way to approximate $\text{LL}(\mathcal{Y}_{\mathcal{A}}, \mathbf{W})$ is to stick with the fused model but approximate the computation of $Z_{\mathcal{A}}(\theta)$. We consider two options:

Full BP In general, any available sum-product inference algorithm like Belief Propagation and their convergent tree reweighted versions [20, 14] can be used for approximating $Z_{\mathcal{A}}(\theta)$. However, these typically require multiple iterations and can be sometimes slow to converge.

OneStep TRW [17] propose a one-step approximation that reduces to a single step of the Tree reweighted (TRW) family of algorithms [14] where the roles of trees are played by individual instances. As in all TRW algorithms, this method guarantees that the log partition value it returns is an upper bound, but for maximization problems upper bounds are not very useful.

A downside of these approaches is that there is no guarantee that the approximation leads to a valid probability distribution. For example, we often observed that the approximate value of $Z_{\mathcal{A}}(\theta)$ was greater than $\sum_{s,i} \log Z(\mathbf{X}_{s,i}, \mathbf{w}_s)$ causing the probability of agreement to be greater than 1.

To summarize, we would ideally like to optimize the agreement-based objective in Equation 6 exactly by working with the equivalent fused graphical model of Equation 9. Due to intractability, we discussed various ways to decompose the agreement term or approximate the corresponding fused model. As we shall show in

Section 5, when there are noisy cliques in the agreement set, the tractable decompositions turn out to be much more robust than methods that approximate the fused model created from erroneous cliques.

3 Generating the agreement set

In this section we discuss our unsupervised strategy for finding agreement sets. But first we stress that the importance of this step cannot be overstated. As we show in Section 5, even the best collective training schemes are only as good as their agreement set. This has interesting parallels with other learning tasks e.g. semi-supervised learning, where recent work has shown the importance of creating good neighborhood graphs [13].

Traditional collective extraction methods have not focused on the process of finding quality agreement sets. These methods usually form a clique from arbitrary repetitions of unigrams [21, 8, 15]. This is inadequate because of two reasons. First, any strong first order dependencies cannot be transferred with only unigram cliques. Second, blindly marking repetitions of a token/n-gram as a clique can inject a lot of noise in the agreement set.

Instead we use a more principled strategy. We make the working assumption that significant content overlap among sources is caused by (approximate-)duplication of instances. So we assume that each instance has a hidden variable with value equal to its ‘canonical instance value’. Instances inside a source will have different values of this variable (as duplicates are rare inside a source), whereas these values will be shared across sources, thus forming clusters. Assume for now that these clusters are known. Given such a cluster of deemed duplicates, we find maximally long segments that repeat among the instances in the cluster, and add one clique per such segment to the agreement set. Segment repetitions outside the cluster are considered as false matches and ignored.

The task of optimally computing the clusters essentially reduces to the NP-hard multi-partite matching problem with suitably defined edge-weights. We tackle this by employing the following staged scheme: First, we order the sources using a natural criteria such as average pairwise similarity with the other sources. Each instance in the first source forms a singleton cluster. In stage s , we find a bipartite matching between source $s + 1$ and the clusters formed by the first s sources. An instance i in source $s + 1$ will be assigned to the cluster to which it is matched. Unmatched instances form new singleton clusters. The edge-weight between an instance i and a cluster is defined as the best similarity of i with any member instance of the cluster.

When our assumption of instance duplication does not hold, say when each instance is an arbitrary natural language sentence, the bipartite matching scores will be low and we revert to the conventional clique generation scheme. As we shall see in Section 5, our strategy generates much better agreement cliques in practice.

4 Relationship with other approaches

We now review various approaches relevant to collective training with partially overlapping sources. We omit Agreement-based learning as it has already been discussed in Sections 1 and 2.3.

4.1 Posterior regularization (PR)

The PR framework [10] trains a model with task-specific linear constraints on the posterior. The constrained optimization problem is solved via the EM algorithm on its variational form. PR has been shown to have interesting relationships with similar frameworks [16, 19, 6].

The aspect of PR most relevant to us is its application to multi-view learning [11]. Then the PR constraints translate to minimizing the Bhattacharayya distance between the various posteriors. This has two key differences with our setting. First their agreement set is at the level of full instances instead of arbitrary sub-parts. Moreover, their agreement set has no noise because the instances across views are known duplicates instead of assumed ones like in ours. The second and more interesting difference is that of the agreement term.

Assuming that we have only two sources s and s' , with only one shared clique \mathbf{c} , training the two models is the same as learning in the presence of two-views of \mathbf{c} . The agreement term under PR would be $\log \sum_{\mathbf{y}_{\mathbf{c}}} \sqrt{P_s(\mathbf{y}_{\mathbf{c}})P_{s'}(\mathbf{y}_{\mathbf{c}})}$, where $P_s(\mathbf{y}_{\mathbf{c}})$ is the marginal of \mathbf{c} . This is maximized when the two marginals are identical. In contrast, our agreement term of $\log \sum_{\mathbf{y}_{\mathbf{c}}} P_s(\mathbf{y}_{\mathbf{c}})P_{s'}(\mathbf{y}_{\mathbf{c}})$ is maximized when the marginals are identical *and peaked*. If both the base models are strong, their marginals will be almost peaked, resulting in little

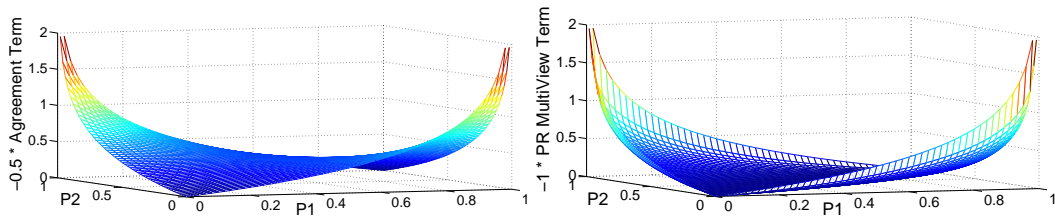


Figure 2: Comparison of the agreement (left) and multi-view (right) losses over two binomial posteriors

difference between the two terms. But a difference arises in the asymmetric case when one model is strong and peaked and the other is weak and flat. One possible maxima for the two-view term would be the strong model flattening out and becoming identical to the weaker one. As our agreement term is averse to flat marginals, it will avoid this maxima. Figure 4.1 illustrates the difference between the two terms for two binomial distributions.

Given such a relationship between the two terms, we compare the performance of our algorithms with the multi-view algorithm in Section 5. The multi-view objective will be optimized via EM because the gradient cannot be computed tractably for structured models.

4.2 Label transfer

Another set of approaches deal with transferring labeled data from one source to the other. One such inexpensive approach is an asymmetric staged strategy of training the model for a more confident source first, transferring its certain labels to the next source and so on. This requires a good ordering of sources to control error-cascades. In Section 5, we show that even with suitable heuristics, this scheme suffers from huge performance deviations. Similar label transfer ideas have been employed in training rule-based extraction wrappers [5].

More sophisticated methods in this class include CoBoosting [7], Co-Training [1], and the two-view Perceptron [2] that train two models in tandem by each model providing labeled data for the other. A detailed comparison of these models in [10] show that these methods are less robust than methods that jointly train all models.

4.3 Inference-only approaches

Another option is to only train the base models, and perform any corrections at runtime through collective inference. Such strategies have been used on a variety of NLP tasks [21, 8, 15]. These methods usually end up using cliques only over unigrams, with little focus on controlling their noise. The most common practice is marking arbitrary repetitions of a token as a clique. As we show in Section 5, our collective training algorithms are significantly better than collective inference, even with identical agreement sets. A prime limitation of inference-only approaches is that they cannot transfer the benefits of overlap to other instances which do not overlap.

5 Experimental evaluation

We present extensive experiments over several real datasets covering a rich diversity of data characteristics. Our first set of experiments seek to justify collective training by showing substantial benefits over base models, and alternatives like staged training and collective inference discussed in Sections 4.2 and 4.3 respectively. Second, we study our collective training approach in detail by comparing the accuracies of the various approximations made in Section 2.3. Third, we demonstrate the importance of choosing high quality agreement sets by comparing various set-generation schemes. Finally, we make a case that our simple gradient ascent algorithm is as accurate as existing traditional EM-based approaches [17, 11] while being considerably faster.

Datasets: We use a corpus of 58 real datasets, each comprising multiple HTML lists. All lists in a dataset contain semi-structured instances relevant to a dataset-specific relation e.g. University mottos, Caldecott medal winners, movies by James Cagney, Supreme court cases etc. The 58 datasets exhibit a wide spectrum of behavior in terms of their base accuracy, number of sources, number of cliques per instance, their noise levels,

	# Datasets	S	$ \mathcal{L} $	$ \mathcal{A} $	Instances	F1 Base	$ \mathcal{A} $ Noise
50F	2	9	4.0	23	75	55.23	0.10
50M	3	11	4.3	202	223	54.59	0.11
60F	4	6	3.5	147	409	67.53	0.07
60M	4	14	4.5	235	344	67.55	0.12
70F	3	9	5.3	146	346	76.67	0.35
70M	14	10	4.4	413	336	75.17	0.21
80F	9	14	4.0	172	575	85.75	0.04
80M	7	13	4.0	959	831	85.9	0.11
90F	6	10	4.0	154	440	94.95	0.04
90M	6	15	4.0	436	493	95.71	0.13
All	58	11	4.2	348	451	79.56	0.15
Std	0	5.8	1.1	500	432	12.24	0.14

Table 1: Properties of the datasets

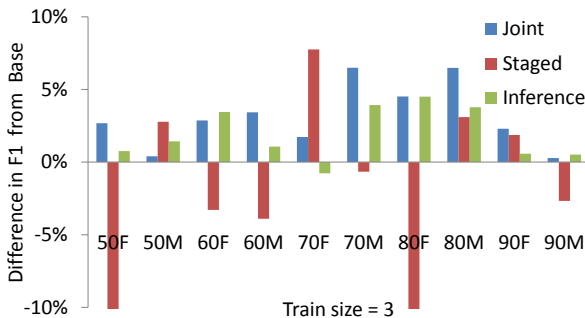


Figure 3: Comparing Clique Agreement against the Collective Inference and Staged approaches

and so on. For ease of presentation, we partition these 58 sources into ten groups by a paired criteria — base accuracy and relative size of the agreement sets. We create five bins for base accuracy values: 50–60, 60–70, and so on, and two bins for agreement set: “M” (many) when there are more than 0.5 cliques per instance and “F” (few) otherwise. Table 1 lists for each of the ten groups the number of datasets ($\#$), average number of sources (S), number of labels ($|\mathcal{L}|$), number of cliques ($|\mathcal{A}|$), instances, base F1 score, and noise in the agreement set \mathcal{A} . The last row in the table that lists the standard deviation of these values over all 58 sources illustrates the diversity of the dataset.

Task For each dataset, we mimic a user query by seeding with a handful of structured records. These are used to generate labeled data out of matching instances in each list of the dataset. The goal is to learn a robust model for each list and extract more instances from it. All comparisons are with 3 and 7 seed records only. Bigger training sets are not practical in this task as the seed structured records are provided through a manual query. All our numbers are averaged over five random selections of the seed training set. Our base model is a conditional random field trained using standard context features over the neighborhood of a word, along with class prior and edge features. Our ground truth consists of every token manually labeled with a relevant dataset-specific label. Using this ground truth, we denote a clique as pure if all its members agree on their true labels, and noisy otherwise. We measure model accuracy by the F1 score of the extracted entities. We set λ using a validation set.

5.1 Benefit of collective training

We first compare collective training, collective inference, and staged label-transfer methods with the base model, starting with only three labeled instances from the user. We chose the Clique agreement method (Clique) for collective training, and used the same agreement set for the collective training and collective inference.

Figure 3 shows the gains and losses of the three methods over the base model for each of ten groups. Collective

Data	Base	Agreement					PR
		Clique	Node	Pair	Full	TR1	EM
Train size = 3							
All	83.3	4.2	3.9	2.6	2.6	2.1	3.7
50F	55.2	2.7	3.5	2.9	2.9	3.5	1.0
50M	54.6	0.6	0.9	1.3	1.1	4.5	3.6
60F	66.9	2.9	2.6	0.8	0.6	1.5	1.5
60M	67.3	3.4	2.3	1.8	2.2	-0.1	3.4
70F	73.5	1.7	1.2	1.4	1.0	0.7	1.1
70M	76.1	6.5	5.8	3.8	4.5	3.7	6.9
80F	85.6	4.5	4.1	3.7	3.5	0.2	4.4
80M	86.6	6.5	6.0	3.8	3.4	3.6	4.5
90F	94.3	2.3	2.1	0.5	1.1	1.2	1.7
90M	96.1	0.3	0.6	-0.1	0.0	0.6	0.4
Train size = 7							
All	87.3	1.8	2.0	1.0	0.9	0.7	2.1
50F	52.5	4.2	5.6	4.8	4.8	4.2	3.4
50M	63.4	0.0	0.1	0.9	0.4	3.5	2.7
60F	76.2	1.9	1.5	0.1	0.1	0.5	1.2
60M	75.0	1.3	2.6	0.9	0.3	-1.5	2.4
70F	79.6	2.3	3.3	0.1	0.0	-0.8	2.5
70M	82.4	3.2	3.5	2.2	2.2	1.7	3.8
80F	90.3	1.5	1.4	1.1	1.2	0.6	1.6
80M	90.5	2.7	2.6	1.3	1.0	0.9	2.4
90F	96.6	0.6	0.2	0.1	0.2	0.2	0.8
90M	96.5	0.3	0.7	0.2	0.2	-0.2	0.5

Table 2: Comparing different training approximations in terms of F1 accuracy gain over the base model.

training clearly performs the best, and its gains are specially large for datasets whose base accuracy is in the 60–80% range, and which have big agreement sets. Overall, its F1 is 87.9% in contrast to the base accuracy of 83.7%. Even with a training size of seven, F1 improves from 87.4 to 89.2 (not shown in the figure). In contrast, the staged approach overall performs worse than Base and shows large swings in accuracy across datasets. It is highly sensitive to the ordering of sources, and the hard label-transfer often causes error-propagation to all downstream sources. Collective inference improves accuracy in a few cases but overall provides only a small gain of 0.3% beyond Base.

5.2 Comparing collective training objectives

We now compare the various approximations of the agreement term — Clique, Node, Instance Pair, Full (Full BP), and TR1 (OneStep TRW) as described in Section 2.3. Table 2 shows the gains in F1 for all the approaches over the base model.

Observe that Clique and Node agreement are two of the best performing methods. We explore two possible reasons for why they score over other approaches that fuse the influence of multiple cliques. One partial explanation is that 15% of the cliques in our agreement set are noisy. In such a case, fused methods would try hard at maximizing the likelihood of a wrongly fused graph. In contrast, the Clique and Node agreement models decompose over cliques, so they can choose to ignore the terms corresponding to erroneous cliques during optimization. A second reason common to all the losing approaches is the inexact nature of the optimization of their training objectives. To understand which of these is a plausible reason, we remove all noisy cliques using the ground truth and compare Clique and Instance Pair agreement. Accuracy improves by less than 0.6 F1 in both, and Clique continues to score over Instance Pair. This indicates that inexact gradient computation is perhaps a major reason why more complex fused approaches perform worse.

We also see that there is little difference between the Clique and Node agreement models. While one possible reason is the weakness of any first-order dependency in the true model, we find another interesting reason for this behavior. We note that in a general n-gram agreement clique, only a few positions might be erroneous. For

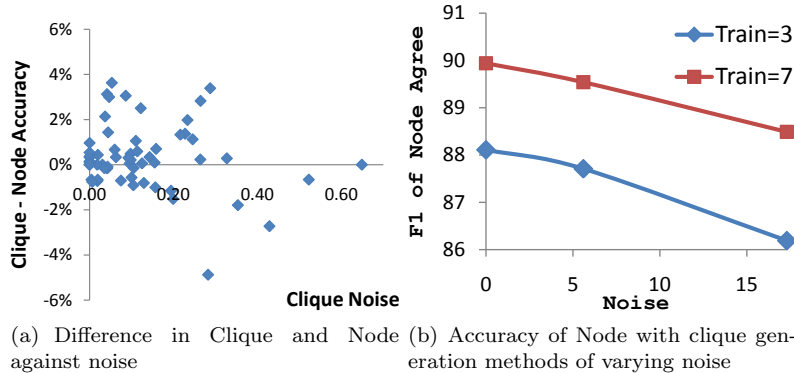


Figure 4: Effect of clique noise on collective training

example, the 15% noise measured at segment level reduces to 5.6% at position level. Since Node decomposes the clique over positions, it can ignore wrong positions during optimization and be more robust against noisy cliques. We corroborate this in Figure 4(a), where for each of the 58 datasets, we plot the difference in F1 of Clique and Node versus clique noise in the dataset’s agreement set. We observe that whenever Clique performs sufficiently worse than Node there is high noise in the cliques. In low noise settings, Clique is often much better than Node.

5.3 Noise in the Agreement Set

As discussed in Section 3, it is important to choose high quality agreement sets. Figure 4(b) shows the F1 scores of Node under three clique generation schemes of varying noise. The rightmost points are for the conventional practice of choosing arbitrary unigram repetitions as cliques and has a noise of 17.3% at position level. The middle point is our method of clique generation where we reduce the noise to 5.6% and the leftmost are ideal cliques with zero noise obtained by using the ground truth to remove all noisy unigrams. We find that our clique selection method enjoys accuracy very close to that with noise-free cliques and the accuracy with carelessly chosen cliques is much lower.

5.4 Comparison with EM-based approaches

In Section 4.1 we described how the PR framework [11] is applicable to our problem. We show its results in the last column of Table 2. The accuracy of PR is comparable to the Clique method showing that distance-based and likelihood terms serve similar goals in our setting. However, the PR approach is more than four times slower than our likelihood objective maximized using gradient ascent. The PR objective requires the EM algorithm for training. In typical feature-based structured models, the M-step tends to be expensive and it is best not wasted on working with fixed E-values. To evaluate the tradeoffs between EM and gradient-based training we also ran the EM algorithm of [17] whose gradient-based version we call TR1 in Table 2. We found the EM trainer (not shown) to have an F1 0.4% less than TR1 and also a factor of two slower.

6 Conclusion

We presented a framework for jointly training multiple extraction models exploiting partial content overlap across sources. Partial overlap opens up a slew of problems — choosing a noise-free agreement set, a training objective or its approximation, and an optimization algorithm. We showed that while decomposing the agreement term over cliques provides a tractable yet accurate method of agreement, it also turns out to be more robust against clique noise than methods that approximate the fused graph. We also presented a strategy for computing clean agreement sets that is far superior to the naïve alternative. Through extensive experiments on various real datasets we showed that our agreement-term decompositions on cliques and positions are more robust, accurate, and faster than alternatives like multi-view learning.

References

- [1] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *COLT*, 1998.
- [2] Ulf Brefeld, Christoph Büscher, and Tobias Scheffer. Multi-view hidden markov perceptrons. In *LWA*, pages 134–138, 2005.
- [3] Razvan Bunescu and Raymond J. Mooney. Collective information extraction with relational markov networks. In *ACL*, pages 439–446, 2004.
- [4] Michael J. Cafarella, Alon Halevy, Yang Zhang, Daisy Zhe Wang, and Eugene Wu. Webtables: Exploring the power of tables on the web. In *VLDB*, 2008.
- [5] Andrew Carlson, Justin Betteridge, Richard C. Wang, Estevam R. Hruschka Jr., and Tom M. Mitchell. Coupled semi-supervised learning for information extraction. In *WSDM*, 2010.
- [6] Ming-Wei Chang, Lev Ratinov, and Dan Roth. Guiding semi-supervision with constraint-driven learning. In *ACL*, pages 280–287, 2007.
- [7] Michael Collins and Yoram Singer. Unsupervised models for named entity classification. In *Proceedings of the SIGDAT - EMNLP*, 1999.
- [8] Jenny Rose Finkel, Trond Grenager, and Christopher D. Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In *ACL*, 2005.
- [9] Kuzman Ganchev, Jennifer Gillenwater, and Ben Taskar. Dependency grammar induction via bitext projection constraints. In *ACL*, pages 369–377, August 2009.
- [10] Kuzman Ganchev, Joao Graca, Jennifer Gillenwater, and Ben Taskar. Posterior regularization for structured latent variable models. Technical Report MS-CIS-09-16, University of Pennsylvania Department of Computer and Information Science, 2009.
- [11] Kuzman Ganchev, João Graça, John Blitzer, and Ben Taskar. Multi-view learning over structured and non-identical outputs. In *UAI*, 2008.
- [12] Rahul Gupta, Ajit A. Diwan, and Sunita Sarawagi. Efficient inference with cardinality-based clique potentials. In *ICML*, 2007.
- [13] T. Jebara, J. Wang, and S.F. Chang. Graph construction and b-matching for semi-supervised learning. In *ICML*, 2009.
- [14] Vladimir Kolmogorov and Martin J. Wainwright. On the optimality of tree-reweighted max-product message passing. In *UAI*, 2005.
- [15] Vijay Krishnan and Christopher D. Manning. An effective two-stage model for exploiting non-local dependencies in named entity recognition. In *ACL-COLING*, 2006.
- [16] Percy Liang, Michael I. Jordan, and Dan Klein. Learning from measurements in exponential families. In *ICML*, 2009.
- [17] Percy Liang, Dan Klein, and Michael I. Jordan. Agreement-based learning. In *NIPS*, 2008.
- [18] Percy Liang, Benjamin Taskar, and Dan Klein. Alignment by agreement. In *HLT-NAACL*, 2006.
- [19] Gideon Mann and Andrew McCallum. Generalized expectation criteria for semi-supervised learning of conditional random fields. In *ACL*, 2008.
- [20] Talya Meltzer, Amir Globerson, and Yair Weiss. Convergent message passing algorithms - a unifying view. In *UAI*, 2009.
- [21] Charles Sutton and Andrew McCallum. Collective segmentation and labeling of distant entities in information extraction. Technical Report TR # 04-49, University of Massachusetts, July 2004.