

A Bayesian Model of Multilingual Unsupervised Semantic Role Induction

Nikhil Garg

University of Geneva
Switzerland

nikgarg@gmail.com

James Henderson

University of Geneva
Switzerland

james.henderson@unige.ch

Abstract

We propose a Bayesian model of unsupervised semantic role induction in multiple languages, and use it to explore the usefulness of parallel corpora for this task. Our joint Bayesian model consists of individual models for each language plus additional latent variables that capture alignments between roles across languages. Because it is a generative Bayesian model, we can do evaluations in a variety of scenarios just by varying the inference procedure, without changing the model, thereby comparing the scenarios directly. We compare using only monolingual data, using a parallel corpus, using a parallel corpus with annotations in the other language, and using small amounts of annotation in the target language. We find that the biggest impact of adding a parallel corpus to training is actually the increase in mono-lingual data, with the alignments to another language resulting in small improvements, even with labeled data for the other language.

1 Introduction

Semantic Role Labeling (SRL) has emerged as an important task in Natural Language Processing (NLP) due to its applicability in information extraction, question answering, and other NLP tasks. SRL is the problem of finding predicate-argument structure in a sentence, as illustrated below:

[_{A0} Mike] has [_{PRED} written] [_{A1} a book] (*S1*)

Here, the predicate WRITE has two arguments: ‘Mike’ as *A0* or *the writer*, and ‘a book’ as *A1* or *the thing written*. The labels *A0* and *A1* correspond to the PropBank annotations (Palmer et al., 2005).

As the need for SRL arises in different domains and languages, the existing manually annotated

corpora become insufficient to build supervised systems. This has motivated work on unsupervised SRL (Lang and Lapata, 2011b; Titov and Klementiev, 2012a; Garg and Henderson, 2012). Previous work has indicated that unsupervised systems could benefit from the word alignment information in parallel text in two or more languages (Naseem et al., 2009; Snyder et al., 2009; Titov and Klementiev, 2012b). For example, consider the German translation of sentence *S1*:

[_{A0}Mike] hat [_{A1}ein Buch][_{PRED}geschrieben] (*S2*)

If sentences *S1* and *S2* have the word alignments: Mike-Mike, written-geschrieben, and book-Buch, the system might be able to predict *A1* for Buch, even if there is insufficient information in the monolingual German data to learn this assignment. Thus, in languages where the resources are sparse or not good enough, or the distributions are not informative, SRL systems could be made more accurate by using parallel data with resource rich or more amenable languages.

In this paper, we propose a joint Bayesian model for unsupervised semantic role induction in multiple languages. The model consists of individual Bayesian models for each language (Garg and Henderson, 2012), and *crosslingual latent variables* to incorporate soft role agreement between aligned constituents. This latent variable approach has been demonstrated to increase the performance in a multilingual unsupervised part-of-speech tagging model based on HMMs (Naseem et al., 2009). We investigate the application of this approach to unsupervised SRL, presenting the performance improvements obtained in different settings involving labeled and unlabeled data, and analyzing the annotation effort required to obtain similar gains using labeled data.

We begin by briefly describing the unsupervised SRL pipeline and the monolingual semantic role induction model we use, and then describe our multilingual model.

2 Unsupervised SRL Pipeline

As established in previous work (Gildea and Jurafsky, 2002; Pradhan et al., 2005), we use a standard unsupervised SRL setup, consisting of the following steps:

- 1. Syntactic Parsing** Off-the-shelf parsers can be used to syntactically parse a given sentence. We use a dependency parse because of its simplicity and easier comparison with the previous work in unsupervised SRL.
- 2. Predicate Identification** We select all the non-auxiliary verbs in a sentence as predicates.
- 3. Argument Identification** For a given predicate, this step classifies each constituent of the parse tree as a semantic argument or a non-argument. Heuristics based on syntactic features such as the dependency relation of a constituent to its head, path from the constituent to the predicate, etc. have been used in unsupervised SRL.
- 4. Argument Classification** Without access to semantic role labels, unsupervised SRL systems cast the problem as a clustering problem. Arguments of a predicate in all the sentences are divided into clusters such that each cluster corresponds to a single semantic role. The better this clustering is, the easier it becomes for a human to give it an actual semantic role label like *A0*, *A1*, etc. Our model assigns a role variable to every identified argument. This variable can take any value from 1 to N , where N is the number of semantic roles that we want to induce.

The task we model, unsupervised semantic role induction, is the step 4 of this pipeline.

3 Monolingual Model

We use the Bayesian model of Garg and Henderson (2012) as our base monolingual model. The semantic roles are predicate-specific. To model the role ordering and repetition preferences, the role inventory for each predicate is divided into *Primary* and *Secondary* roles as follows:

Primary Role (PR) Let there be a total of N roles (or clusters) for each predicate. Assign K of them as PRs $\{P_1, P_2, \dots, P_K\}$. Further, create 3 additional PRs: *START* denoting the start of the role sequence, *END* denoting its end, and *PRED* denoting the predicate. These $(K + 3)$ PRs are not allowed to repeat in a frame and their ordering defines the global role ordering.

Secondary Role (SR) The rest of the $(N - K)$ roles are called SRs $\{S_1, S_2, \dots, S_{N-K}\}$. Unlike PRs, they are not constrained to occur only once and only their ordering w.r.t. PRs is used in the probability model.

For example, the complete role sequence in a frame could be: $(START, P_3, S_1, S_1, PRED, P_2, S_5, END)$. The *ordering* is defined as the sequence of PRs, $(START, P_3, PRED, P_2, END)$. Each pair of consecutive PRs in an ordering is called an *interval*. Thus, $(P_3, PRED)$ is an interval that contains two SRs, S_1 and S_1 . An interval could also be empty, for instance $(START, P_3)$ contains no SRs. When we evaluate, these roles get mapped to gold roles. For instance, the PR P_2 could get mapped to a core role like *A0*, *A1*, etc. or to a modifier role like *AM-TMP*, *AM-MOD*, etc. Garg and Henderson (2012) reported that, in practice, PRs mostly get mapped to core roles and SRs to modifier roles, which conforms to the linguistic motivations for this distinction.

Figure 4 illustrates two copies of the monolingual model, on either side of the crosslingual latent variables. The generative process is as follows:

- 1. Predicate, Voice** The predicate p and its voice vc are treated as top-level visible variables.
- 2. Ordering (Generate PRs)** Select an ordered set of PRs from a multinomial distribution.

$$o \sim Multinomial(\theta_{p,vc}^{order})$$
- 3. Generate SRs** For each interval in the ordering o , a sequence of SRs is generated as:
 for each interval $I \in o$:
 draw an indicator $s \sim Binomial(\theta_{p,I,0}^{STOP})$
 while $s \neq STOP$:
 choose a SR $r \sim Multinomial(\theta_{p,I}^{SR})$
 draw an indicator $s \sim Binomial(\theta_{p,I,1}^{STOP})$
- 4. Generate Features** For each PR and SR, the features for that constituent are generated independently. To keep the model simple and comparable to previous unsupervised work, we only use three features: (i) dependency relation of the argument to its head, (ii) head word of the argument, and (iii) POS tag of the head word:
 for each generated role r :
 for each feature type f :
 choose a value $v_f \sim Multinomial(\theta_{p,r,f}^F)$

All the multinomial and binomial distributions have symmetric Dirichlet and beta priors respectively. Figure 1a gives the probability equations

$$P(\mathbf{r}, \mathbf{f} | p, vc) = \underbrace{P(o | p, vc)}_{o = \text{ordering}(\mathbf{r})} \underbrace{\prod_{r_i \in \mathbf{r} \cap PR} P(f_i | r_i, p)}_{\text{Primary Roles}} \underbrace{\prod_{I \in o} P(\mathbf{r}(I), \mathbf{f}(I) | I, p)}_{\text{Intervals}} \quad (1)$$

$$\text{where } P(\mathbf{r}(I), \mathbf{f}(I) | I, p) = \prod_{r_i \in \mathbf{r}(I)} \underbrace{P(\neg \text{stop} | I, p, \text{adj})}_{\text{generate indicator}} \underbrace{P(r_i | I, p)}_{\text{generate SR}} \underbrace{P(f_i | r_i, p)}_{\text{generate features}} \underbrace{P(\text{stop} | I, p, \text{adj})}_{\text{end of the interval}} \quad (2)$$

$$\text{and } P(f_i | r_i, p) = \prod_{t=1}^T P(f_{i,t} | r_i, p) \quad (3)$$

$$P(\mathbf{f} | p, vc) = \sum_{\mathbf{r}} P(\mathbf{r}, \mathbf{f} | p, vc) \quad (4)$$

(a) Probability equations for the monolingual model. Bold-faced variables denote a sequence of values. \mathbf{r} denotes the complete sequence of roles, and \mathbf{f} denotes the complete sequence of features. p and vc denote the predicate and its voice respectively. o denotes the ordering of PRs in the sequence \mathbf{r} and $\text{ordering}(\mathbf{r})$ is a function for computing this ordering. r_i and f_i denote the role and features at position i respectively, and $r(I)$ and $f(I)$ respectively denote the SR sequence and feature sequence in interval I . $f_{i,t}$ denotes the value of feature t at position i . $\text{adj} = 0$ for generating the first SR, and 1 for a subsequent one. Equation 1 gives the joint probability of the model and equation 4 gives the marginal probability of the observed features.

$$P(\mathbf{r}^{l1}, \mathbf{f}^{l1}, \mathbf{r}^{l2}, \mathbf{f}^{l2}, \mathbf{z} | p^{l1}, vc^{l1}, p^{l2}, vc^{l2}) = P(\mathbf{z}) \prod_{l \in \{l1, l2\}} P(\mathbf{r}^l, \mathbf{f}^l | \mathbf{z}, p^l, vc^l) \quad (5)$$

$$\approx P(\mathbf{z}) \prod_{l \in \{l1, l2\}} P(\mathbf{r}^l, \mathbf{f}^l | p^l, vc^l) \prod_{i,k: z_k \rightarrow r_i^l} P(r_i^l | z_k) \quad (6)$$

(a) Probability equations for the multilingual model. The superscript l denotes the variable for language l . \mathbf{z} denotes the common crosslingual latent variables for both languages. $z_k \rightarrow r_i^l$ denotes that the argument at position i in language l is connected to the crosslingual latent variable $\#k$.

Figure 3: Probability equations for the (a) monolingual and (b) multilingual model.

for the monolingual model. This formulation models the global role ordering and repetition preferences using PRs, and limited context for SRs using intervals. Ordering and repetition information was found to be helpful in supervised SRL as well (Punyakanok et al., 2004; Pradhan et al., 2005; Toutanova et al., 2008). More details, including the motivations behind this model, are in (Garg and Henderson, 2012).

4 Multilingual Model

The multilingual model uses word alignments between sentences in a parallel corpus to exploit role correspondences across languages. We make copies of the monolingual model for each language and add additional *crosslingual latent variables* (CLVs) to couple the monolingual models, capturing crosslingual semantic role patterns. Concretely, when training on parallel sentences, whenever the head words of the arguments are aligned, we add a CLV as a parent of the two corresponding role variables. Figure 4 illustrates this model. The generative process, as explained below, remains the same as the monolingual model for the most part, with the exception of aligned

roles which are now generated by both the monolingual process as well as the CLV.

1. Monolingual Data Given a parallel frame with the predicate pair $p1, p2$, generate two separate monolingual frames as in section 3.

2. Aligned Arguments For each aligned argument, first generate a crosslingual latent variable from a Chinese Restaurant Process (CRP). Then generate the two aligned roles:

for aligned arguments i, j :

draw a crosslingual latent variable:

$$z \sim CRP(\alpha_{p1, p2}^{CRP})$$

draw role for language $l1$:

$$r_i \sim \text{Multinomial}(\theta_{p1, p2, z, l1}^{align})$$

draw role for language $l2$:

$$r_j \sim \text{Multinomial}(\theta_{p1, p2, z, l2}^{align})$$

Every predicate-tuple has its own inventory of CLVs specific to that tuple. Each CLV z is a multi-valued variable where each value defines a distribution over role labels for each language (denoted by $\theta_{p1, p2, z, l}^{align}$ above). These distributions over labels are trained to be peaky, so that each value c for a CLV represents a correlation between the labels that c predicts in the two languages. For ex-

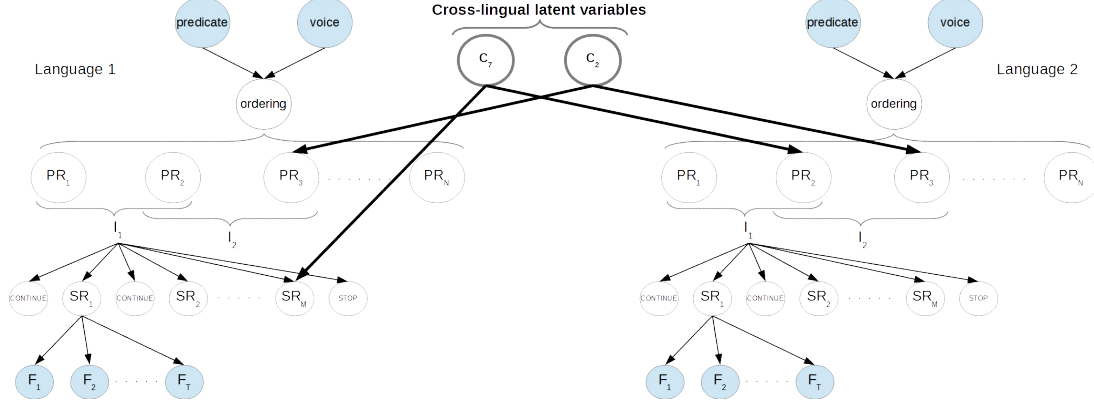


Figure 4: Multilingual model. The CLVs and their associated parameters are drawn in bold. PR_3 in language 1 is aligned to PR_3 in language 2 with the corresponding CLV taking the value c_2 , and SR_M is aligned to PR_2 with the CLV taking the value c_7 .

ample, a value c for the CLV z might give high probabilities to S_3 and S_8 in language 1, and to S_1 in language 2. If c is the only value for z that gives high probability to S_3 in language 1, and the monolingual model in language 1 decides to assign S_3 to the role for z , then z will predict S_1 in language 2, with high probability. We generate the CLVs via a Chinese Restaurant Process (Pitman, 2002), a non-parametric Bayesian model, which allows us to induce the number of CLVs for every predicate-tuple from the data. We continue to train on the non-parallel sentences using the respective monolingual models.

The multilingual model is deficient, since the aligned roles are being generated twice. Ideally, we would like to add the CLV as additional conditioning variables in the monolingual models. The new joint probability can be written as equation 5 (Figure 2a), which can be further decomposed following the decomposition of the monolingual model in Figure 1a. However, having this additional conditioning variable breaks the Dirichlet-multinomial conjugacy, which makes it intractable to marginalize out the parameters during inference. Hence, we use an approximation where we treat each of the aligned roles as being generated twice, once by the monolingual model and once by the corresponding CLV (equation 6).

This is the first work to incorporate the coupling of aligned arguments directly in a Bayesian SRL model. This makes it easier to see how to extend this model in a principled way to incorporate additional sources of information. First, the model scales gracefully to more than two languages. If there are a total of n languages, and there is an aligned argument in m of them, the

multilingual latent variable is connected to only those m aligned arguments.

Second, having one joint Bayesian model allows us to use the same model in various semi-supervised learning settings, just by fixing the annotated variables during training. Section 6.6 evaluates a setting where we have some labeled data in one language (called source), while no labeled data in the second language (called target). Note that this is different from a classic annotation projection setting (e.g. (Padó and Lapata, 2009)), where the role labels are mapped from source constituents to aligned target constituents.

5 Inference and Training

The inference problem consists of predicting the role labels and CLVs (the hidden variables) given the predicate, its voice, and syntactic features of all the identified arguments (the visible variables). We use a collapsed Gibbs-sampling based approach to generate samples for the hidden variables (model parameters are integrated out). The sample counts and the priors are then used to calculate the MAP estimate of the model parameters.

For the monolingual model, the role at a given position is sampled as:

$$P(r_i | \mathbf{r}_{-i}, \mathbf{f}, p, vc, D^-) \propto P(r_i, \mathbf{r}_{-i}, \mathbf{f} | p, vc, D^-) \\ = \int P(r_i, \mathbf{r}_{-i}, \mathbf{f} | \boldsymbol{\theta}, p, vc) P(\boldsymbol{\theta} | D^-) d\boldsymbol{\theta}$$

where the subscript $-i$ refers to all the variables except at position i , D^- refers to the variables in all the training instances except the current one, and $\boldsymbol{\theta}$ refers to all the model parameters. The above integral has a closed form solution due to Dirichlet-multinomial conjugacy.

For sampling roles in the multilingual model, we also need to consider the probabilities of roles being generated by the CLVs:

$$\begin{aligned} P(r_i | \mathbf{r}_{-i}, \mathbf{f}, p, vc, \mathbf{z}, D^-) &\propto P(r_i, \mathbf{r}_{-i}, \mathbf{f} | \mathbf{z}, p, vc, D^-) \\ &= \int P(r_i, \mathbf{r}_{-i}, \mathbf{f} | \boldsymbol{\theta}, \mathbf{z}, p, vc) P(\boldsymbol{\theta} | D^-) d\boldsymbol{\theta} \\ &= \int P(r_i, \mathbf{r}_{-i}, \mathbf{f} | \boldsymbol{\theta}, p, vc) \left(\prod_{j, k: z_k \rightarrow r_j} P(r_j | \boldsymbol{\theta}, z_k) \right) P(\boldsymbol{\theta} | D^-) d\boldsymbol{\theta} \end{aligned}$$

For sampling CLVs, we need to consider three factors: two corresponding to probabilities of generating the aligned roles, and the third one corresponding to selecting the CLV according to CRP.

$$P(z_k | r_i^{l1}, r_j^{l2}, D^{-,k}) \propto P(r_i^{l1} | z_k, D^{-,k}) P(r_j^{l2} | z_k, D^{-,k}) P(z_k | D^{-,k})$$

where the aligned roles r_i^{l1} and r_j^{l2} are connected to z_k , and $D^{-,k}$ refers to all the variables except z_k , r_i^{l1} , and r_j^{l2} .

We use the trained parameters to parse the monolingual data using the monolingual model. The crosslingual parameters are ignored even if they were used during training. Thus, the information coming from the CLVs acts as a regularizer for the monolingual models.

6 Experiments

6.1 Evaluation

Following the setting of Titov and Klementiev (2012b), we evaluate only on the arguments that were correctly identified, as the incorrectly identified arguments do not have any gold semantic labels. Evaluation is done using the metric proposed by Lang and Lapata (2011a), which has 3 components: (i) **Purity (PU)** measures how well an induced cluster corresponds to a single gold role, (ii) **Collocation (CO)** measures how well a gold role corresponds to a single induced cluster, and (iii) **F1** is the harmonic mean of PU and CO. For each predicate, let N denote the total number of argument instances, C_i the instances in the induced cluster i , and G_j the instances having label j in gold annotations. $PU = \frac{1}{N} \sum_i \max_j |C_i \cap G_j|$, $CO = \frac{1}{N} \sum_j \max_i |C_i \cap G_j|$, and $F1 = \frac{2 \cdot PU \cdot CO}{PU + CO}$. The score for each predicate is weighted by the number of its argument instances, and a weighted average is computed over all the predicates.

6.2 Baseline

We use the same baseline as used by Lang and Lapata (2011a) which has been shown to be difficult to outperform. This baseline assigns a semantic

role to a constituent based on its syntactic function, i.e. the dependency relation to its head. If there is a total of N clusters, $(N - 1)$ most frequent syntactic functions get a cluster each, and the rest are assigned to the N th cluster.

6.3 Closest Previous Work

This work is closely related to the cross-lingual unsupervised SRL work of Titov and Klementiev (2012b). Their model has separate monolingual models for each language and an extra penalty term which tries to maximize $P(r^{l1} | r^{l2})$ and $P(r^{l2} | r^{l1})$ i.e. for all the aligned arguments with role label r^{l1} in language 1, it tries to find a role label r^{l2} in language 2 such that the given proportion is maximized and vice versa. However, there is no efficient way to optimize the objective with this penalty term and the authors used an inference method similar to annotation projection. Further, the method does not scale naturally to more than two languages. Their algorithm first does monolingual inference in one language ignoring the penalty and then does the inference in the second language taking into account the penalty term. In contrast, our model adds the latent variables as a part of the model itself, and not an external penalty, which enables us to use the standard Bayesian learning methods such as sampling.

The monolingual model we use (Garg and Henderson, 2012) also has two main advantages over Titov and Klementiev (2012b). First, the former incorporates a global role ordering probability that is missing in the latter. Secondly, the latter defines *argument-keys* as a tuple of four syntactic features and all the arguments having the same argument-keys are assigned the same role. This kind of hard clustering is avoided in the former model where two constituents having the same set of features might get assigned different roles if they appear in different contexts.

6.4 Data

Following Titov and Klementiev (2012b), we run our experiments on the English (EN) and German (DE) sections of the CoNLL 2009 corpus (Hajič et al., 2009), and EN-DE section of the Europarl corpus (Koehn, 2005). We get about 40k EN and 36k DE sentences from the CoNLL 2009 training set, and about 1.5M parallel EN-DE sentences from Europarl. For appropriate comparison, we keep the same setting as in (Titov and Klementiev, 2012b) for automatic parses and ar-

gument identification, which we briefly describe here. The EN sentences are parsed syntactically using MaltParser (Nivre et al., 2007) and DE using LTH parser (Johansson and Nugues, 2008). All the non-auxiliary verbs are selected as predicates. In CoNLL data, this gives us about 3k EN and 500 DE predicates. The total number of predicate instances are 3.4M in EN (89k CoNLL + 3.3M Europarl) and 2.62M in DE (17k CoNLL + 2.6M Europarl). The arguments for EN are identified using the heuristics proposed by Lang and Lapata (2011a). However, we get an F1 score of 85.1% for argument identification on CoNLL 2009 EN data as opposed to 80.7% reported by Titov and Klementiev (2012b). This could be due to implementation differences, which unfortunately makes our EN results incomparable. For DE, the arguments are identified using the LTH system (Johansson and Nugues, 2008), which gives an F1 score of 86.5% on the CoNLL 2009 DE data. The word alignments for the EN-DE parallel Europarl corpus are computed using GIZA++ (Och and Ney, 2003). For high-precision, only the intersecting alignments in the two directions are kept. We define two semantic arguments as aligned if their head-words are aligned. In total we get 9.3M arguments for EN (240k CoNLL + 9.1M Europarl) and 4.43M for DE (32k CoNLL + 4.4M Europarl). Out of these, 0.76M arguments are aligned.

6.5 Main Results

Since the CoNLL annotations have 21 semantic roles in total, we use 21 roles in our model as well as the baseline. Following Garg and Henderson (2012), we set the number of PRs to 2 (excluding *START*, *END* and *PRED*), and SRs to 21-2=19. Table 1 shows the results.

In the first setting (Line 1), we train and test the monolingual model on the CoNLL data. We observe significant improvements in F1 score over the Baseline (Line 0) in both languages. Using the CoNLL 2009 dataset alone, Titov and Klementiev (2012b) report an F1 score of 80.9% (PU=86.8%, CO=75.7%) for German. Thus, our monolingual model outperforms their monolingual model in German. For English, they report an F1 score of 83.6% (PU=87.5%, CO=80.1%), but note that our English results are not directly comparable to theirs due to differences argument identification, as discussed in section 6.4. As their argument identification score is lower, perhaps their system

is discarding “difficult” arguments which leads to a higher clustering score.

In the second setting (Line 2), we use the additional monolingual Europarl (EP) data for training. We get equivalent results in English and a significant improvement in German compared to our previous setting (Line 1). The German dataset in CoNLL is quite small and benefits from the additional EP training data. In contrast, the English model is already quite good due to a relatively big dataset from CoNLL, and good accuracy syntactic parsers. Unfortunately, Titov and Klementiev (2012b) do not report results with this setting.

The third setting (Line 3) gives the results of our multilingual model, which adds the word alignments in the EP data. Comparing with Line 2, we get non-significant improvements in both languages. Titov and Klementiev (2012b) obtain an F1 score of 82.7% (PU=85.0%, CO=80.6%) for German, and 83.7% (PU=86.8%, CO=80.7%) for English. Thus, for German, our multilingual Bayesian model is able to capture the cross-lingual patterns at least as well as the external penalty term in (Titov and Klementiev, 2012b). We cannot compare the English results unfortunately due to differences in argument identification.

We also compared monolingual and bilingual training data using a setting that emulates the standard supervised setup of separate training and test data sets. We train only on the EP dataset and test on the CoNLL dataset. Lines 4 and 5 of Table 1 give the results. The multilingual model obtains small improvements in both languages, which confirms the results from the standard unsupervised setup, comparing lines 2 to 3.

These results indicate that little information can be learned about semantic roles from this parallel data setup. One possible explanation for this result is that the setup itself is inadequate. Given the definition of aligned arguments, only 8% of English arguments and 17% of German arguments are aligned. This plus our experiments suggest that improving the alignment model is a necessary step to making effective use of parallel data in multilingual SRI, for example by joint modeling with SRI. We leave this exploration to future work.

6.6 Multilingual Training with Labeled Data for One Language

Another motivation for jointly modeling SRL in multiple languages is the transfer of information

| | Dataset | | Model | English | | | German | | |
|---|--------------|---------|--------------|---------|-------|--------------|--------|-------|--------------|
| | Training | Testing | | PU | CO | F1 | PU | CO | F1 |
| 0 | CoNLL | CoNLL | Baseline | 78.23 | 79.46 | 78.84 | 83.09 | 79.32 | 81.16 |
| 1 | CoNLL | CoNLL | Monolingual | 76.29 | 83.13 | 79.56* | 82.54 | 81.94 | 82.24* |
| 2 | CoNLL+EP | CoNLL | Monolingual | 76.11 | 83.33 | 79.56 | 83.77 | 81.65 | 82.70* |
| 3 | 2×(CoNLL+EP) | CoNLL | Multilingual | 76.23 | 83.25 | 79.59 | 83.81 | 81.79 | 82.79 |
| 4 | EP | CoNLL | Monolingual | 73.26 | 80.60 | 76.76 | 83.72 | 81.28 | 82.48 |
| 5 | 2×EP | CoNLL | Multilingual | 73.07 | 81.24 | 76.94 | 83.59 | 81.50 | 82.54 |

Table 1: Main Results. A * denotes a significant improvement in the F1 score over the the previous line. We compute the significance using stratified shuffling and consider it significant if the p-value < 0.05.

| | Source | Target | English | | | German | | |
|---|--------------------|---------|---------|-------|--------------|--------|-------|--------------|
| | | | PU | CO | F1 | PU | CO | F1 |
| 1 | Multilingual Model | | 76.23 | 83.25 | 79.59 | 83.81 | 81.79 | 82.79 |
| 2 | English | German | NA | | | 83.83 | 81.83 | 82.82 |
| 3 | German | English | 76.26 | 83.37 | 79.66 | NA | | |

Table 2: Results for the Multilingual Model with using labeled data for the source language.

from a resource rich language to a resource poor language. We evaluated our model in a very general annotation transfer scenario, where we have a small labeled dataset for one language (source), and a large parallel unlabeled dataset for the source and another (target) language. We investigate whether this setting improves the parameter estimates for the target language. To this end, we clamp the role annotations of the source language in the CoNLL dataset using a predefined mapping¹, and do not sample them during training. This data gives us good parameters for the source language, which are used to sample the roles of the source language in the unlabeled Europarl data. The CLVs aim to capture this improvement and thereby improve sampling and parameter estimates for the target language. Table 2 shows the results of this experiment. We obtain small improvements in the target languages. As in the unsupervised setting, the small percentage of aligned roles probably limits the impact of the cross-lingual information.

6.7 Labeled Data in Monolingual Model

We explored the improvement in the monolingual model in a semi-supervised setting. To this end, we randomly selected $S\%$ of the sentences in the CoNLL dataset as “supervised sentences” and the rest $(100 - S)\%$ were kept unsupervised. Next, we clamped the role labels of the supervised sentences

¹A0 was mapped to the primary role P_1 , A1 to P_2 , and the rest were mapped to the secondary roles (S_1, \dots, S_{19}) in the order of their decreasing frequency.

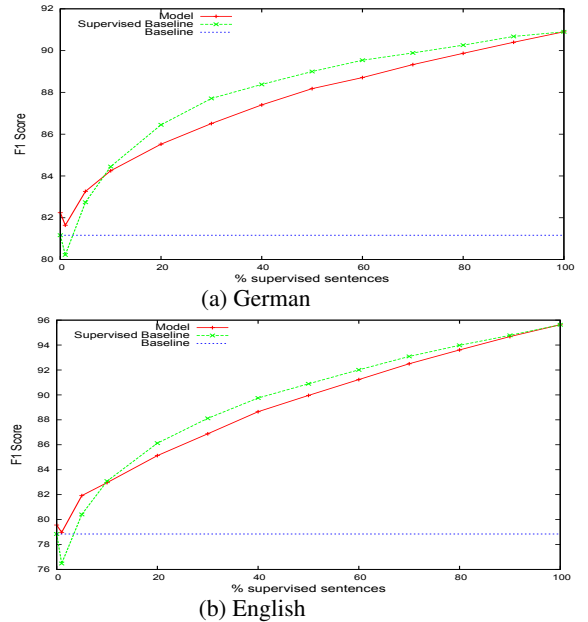


Figure 5: F1 with a portion of the data labeled.

using the predefined mapping from Section 6.6. Sampling was done on the unsupervised sentences as usual. We then measured the clustering performance using the trained parameters.²

To access the contribution of partial supervision better, we constructed a “supervised baseline” as follows. For predicates seen in the supervised sentences, a MAP estimate of the parameters was calculated using the predefined mapping. For the unseen predicates, the standard baseline was used.

Figures 5a and 5b show the performance varia-

²To account for the randomness in selecting the supervised sentences, the experiment was repeated 10 times and average of the performance numbers was taken.

tion with S . We make the following observations:

- In both languages, at around $S = 10$, the supervised baseline starts outperforming the semi-supervised model, which suggests that manually labeling about 10% of the sentences is a good enough alternative to our training procedure. Note that 10% amounts to about 3.6k sentences in German and 4k in English. We noticed that the proportion of seen predicates increases dramatically as we increase the proportion of supervised sentences. At 10% supervised sentences, the model has already seen 63% of predicates in German and 44% in English. This explains to some extent why only 10% labeled sentences are enough.
- For German, it takes about 3.5% or 1260 supervised sentences to have the same performance increase as 1.5M unlabeled sentences (Line 1 to Line 2 in Table 1). Adding about 180 more supervised sentences also covers the benefit obtained by alignments in the multilingual model (Line 2 to Line 3 in Table 1). There is no noticeable performance difference in English.

We also evaluated the performance variation on a completely unseen CoNLL test set. Since the test set is very small compared to the training set, the clustering evaluation is not as reliable. Nonetheless, we broadly obtained the same pattern.

7 Related Work

As discussed in section 6.3, our work is closely related to the crosslingual unsupervised SRL work of Titov and Klementiev (2012b). The idea of using *superlingual* latent variables to capture crosslingual information was proposed for POS tagging by Naseem et al. (2009), which we use here for SRL. In a semi-supervised setting, Padó and Lapata (2009) used a graph based approach to transfer semantic role annotations from English to German. Fürstenaue and Lapata (2009) used a graph alignment method to measure the semantic and syntactic similarity between dependency tree arguments of known and unknown verbs.

For monolingual unsupervised SRL, Swier and Stevenson (2004) presented the first work on a domain-general corpus, the British National Corpus, using 54 verbs taken from VerbNet. Garg and Henderson (2012) proposed a Bayesian model for this problem that we use here. Titov and

Klementiev (2012a) also proposed a closely related Bayesian model. Grenager and Manning (2006) proposed a generative model but their parameter space consisted of all possible linkings of syntactic constituents and semantic roles, which made unsupervised learning difficult and a separate language-specific rule based method had to be used to constrain this space. Other proposed models include an iterative split-merge algorithm (Lang and Lapata, 2011a) and a graph-partitioning based approach (Lang and Lapata, 2011b). Márquez et al. (2008) provide a good overview of the supervised SRL systems.

8 Conclusions

We propose a Bayesian model of semantic role induction (SRI) that uses crosslingual latent variables to capture role alignments in parallel corpora. The crosslingual latent variables capture correlations between roles in different languages, and regularize the parameter estimates of the monolingual models. Because this is a joint Bayesian model of multilingual SRI, we can apply the same model to a variety of training scenarios just by changing the inference procedure appropriately. We evaluate monolingual SRI with a large unlabeled dataset, bilingual SRI with a parallel corpus, bilingual SRI with annotations available for the source language, and monolingual SRI with a small labeled dataset. Increasing the amount of monolingual unlabeled data significantly improves SRI in German but not in English. Adding word alignments in parallel sentences results in small, non significant improvements, even if there is some labeled data available in the source language. This difficulty in showing the usefulness of parallel corpora for SRI may be due to the current assumptions about role alignments, which mean that only a small percentage of roles are aligned. Further analyses reveals that annotating small amounts of data can easily outperform the performance gains obtained by adding large unlabeled dataset as well as adding parallel corpora.

Future work includes training on different language pairs, on more than two languages, and with more inclusive models of role alignment.

Acknowledgments

This work was funded by the Swiss NSF grant 200021_125137 and EC FP7 grant PARLANCE.

References

- [Fürstenau and Lapata2009] H. Fürstenau and M. Lapata. 2009. Graph alignment for semi-supervised semantic role labeling. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1-Volume 1*, pages 11–20. Association for Computational Linguistics.
- [Garg and Henderson2012] N. Garg and J. Henderson. 2012. Unsupervised semantic role induction with global role ordering. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics.
- [Gildea and Jurafsky2002] D. Gildea and D. Jurafsky. 2002. Automatic labeling of semantic roles. *Computational Linguistics*, 28(3):245–288.
- [Grenager and Manning2006] T. Grenager and C.D. Manning. 2006. Unsupervised discovery of a statistical verb lexicon. In *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing*, pages 1–8. Association for Computational Linguistics.
- [Hajič et al.2009] J. Hajič, M. Ciaramita, R. Johansson, D. Kawahara, M.A. Martí, L. Màrquez, A. Meyers, J. Nivre, S. Padó, J. Štěpánek, et al. 2009. The CoNLL-2009 shared task: Syntactic and semantic dependencies in multiple languages. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning: Shared Task*, pages 1–18. Association for Computational Linguistics.
- [Johansson and Nugues2008] R. Johansson and P. Nugues. 2008. Dependency-based semantic role labeling of propbank. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 69–78. Association for Computational Linguistics.
- [Koehn2005] P. Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *MT summit*, volume 5.
- [Lang and Lapata2011a] J. Lang and M. Lapata. 2011a. Unsupervised semantic role induction via split-merge clustering. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics, Portland, Oregon*.
- [Lang and Lapata2011b] J. Lang and M. Lapata. 2011b. Unsupervised semantic role induction with graph partitioning. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 1320–1331, Edinburgh, Scotland, UK., July. Association for Computational Linguistics.
- [Màrquez et al.2008] L. Màrquez, X. Carreras, K.C. Litkowski, and S. Stevenson. 2008. Semantic role labeling: an introduction to the special issue. *Computational linguistics*, 34(2):145–159.
- [Naseem et al.2009] T. Naseem, B. Snyder, J. Eisenstein, and R. Barzilay. 2009. Multilingual part-of-speech tagging: Two unsupervised approaches. *Journal of Artificial Intelligence Research*, 36(1):341–385.
- [Nivre et al.2007] J. Nivre, J. Hall, J. Nilsson, A. Chanev, G. Eryigit, S. Kubler, S. Marinov, and E. Marsi. 2007. Maltparser: A language-independent system for data-driven dependency parsing. *Natural Language Engineering*, 13(2):95.
- [Och and Ney2003] F.J. Och and H. Ney. 2003. A systematic comparison of various statistical alignment models. *Computational linguistics*, 29(1):19–51.
- [Padó and Lapata2009] S. Padó and M. Lapata. 2009. Cross-lingual annotation projection for semantic roles. *Journal of Artificial Intelligence Research*, 36(1):307–340.
- [Palmer et al.2005] M. Palmer, D. Gildea, and P. Kingsbury. 2005. The proposition bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31(1):71–106.
- [Pitman2002] J. Pitman. 2002. Combinatorial stochastic processes. Technical report, Technical Report 621, Dept. Statistics, UC Berkeley, 2002. Lecture notes for St. Flour course.
- [Pradhan et al.2005] S. Pradhan, K. Hacioglu, V. Krugler, W. Ward, J.H. Martin, and D. Jurafsky. 2005. Support vector learning for semantic argument classification. *Machine Learning*, 60(1):11–39.
- [Punyakankok et al.2004] V. Punyakankok, D. Roth, W. Yih, and D. Zimak. 2004. Semantic role labeling via integer linear programming inference. In *Proceedings of the 20th international conference on Computational Linguistics*, page 1346. Association for Computational Linguistics.
- [Snyder et al.2009] B. Snyder, T. Naseem, and R. Barzilay. 2009. Unsupervised multilingual grammar induction. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1*, pages 73–81. Association for Computational Linguistics.
- [Swier and Stevenson2004] R. Swier and S. Stevenson. 2004. Unsupervised semantic role labelling. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 95–102.
- [Titov and Klementiev2012a] I. Titov and A. Klementiev. 2012a. A bayesian approach to unsupervised semantic role induction. In *Proceedings of the Conference of the European Chapter of the Association for Computational Linguistics*, page 12.

- [Titov and Klementiev2012b] I. Titov and A. Klementiev. 2012b. Crosslingual induction of semantic roles. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.
- [Toutanova et al.2008] K. Toutanova, A. Haghighi, and C.D. Manning. 2008. A global joint model for semantic role labeling. *Computational Linguistics*, 34(2):161–191.