# Partial Or Complete, That's The Question

Qiang Ning<sup>1</sup>, Hangfeng He<sup>2</sup>, Chuchu Fan<sup>1</sup>, Dan Roth<sup>1,2</sup> <sup>1</sup>Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign <sup>2</sup>Department of Computer and Information Science, University of Pennsylvania

## Background

Many learning tasks are **structured**. Considering the structure in the learning/prediction phase leads to structured learning/prediction. **However**, the effect of structure on data collection has received less attention.



This paper aims to provide a better (both theoretical and empirical) understanding of **the role of structure in data collection for structured tasks**.

## Motivation

Structure: a set of variables that are not independent (mathematical definition shown in Sec. "Theory").
A common perception: "partial" data are of low quality and should be avoided in data collection.
This paper challenges it and provides a more principled understanding for "partial".



**Motivation**: if we annotate a structure only partially, the remaining part is constrained to a smaller space; that is, some "information" of the remaining part has been explained by existing annotations. This is a **benefit of "partial" that is overlooked previously** [\*SEM'18].

Approach: Partial Or Complete	
-------------------------------	--

To study the effect of structure on data collection, we investigate **two data collection paradigms**:

• **Complete**: complete as many structures as possible; the cost is leaving some structures empty

• Early stopping: evenly distribute the annotation budget to all structures; the cost is only getting partial annotations We train two systems accordingly and compare their performances on the same test set.



**Structure**: a vector of *d* random variables:  $Y = [Y_1, Y_2, ..., Y_d] \in C(\mathcal{L}^d) \subseteq \mathcal{L}^d$ , where  $\mathcal{L}$  is the label set. <u>Two simple cases:</u>

- When the variables are independent:  $\mathcal{C}(\mathcal{L}^d) = \mathcal{L}^d$
- When the structure is so "strong" that it requires all variables to share the same label:  $C(\mathcal{L}^d) = \{[\ell_1, \ell_1, \dots, \ell_1], [\ell_2, \ell_2, \dots, \ell_2], \dots, [\ell_{|\mathcal{L}|}, \ell_{|\mathcal{L}|}, \dots, \ell_{|\mathcal{L}|}\}$

# Complete annotation: All variables are labeled and a unique point in $C(\mathcal{L}^d)$ is determined.

**Partial annotation**: A subset of  $C(\mathcal{L}^d)$  is determined, which contains the "true" structure.

Let f<sub>k</sub> be the size of the feasible subset if k out of d variables are labeled.

 $\Box f_0 = \left| \mathcal{C}(\mathcal{L}^d) \right| \ge f_1 \ge f_2 \ge \dots \ge f_d = 1$ 

No annotation

#### Complete annotation

# $1 \rightarrow 2 \rightarrow 3 \rightarrow 4$

# Theory

# Chain structure in ranking problems:

- Linear chain with a transitivity constraint.
- When k out of d comparisons are given, the structure is a partial order.
- $f_k$ : # of linear extensions of these partial orders; need simulations to estimate  $I_k$  [ToC'91].

# Bipartite graph structure in assignment problems:

- Assign d agents to d' tasks (w.l.o.g, d < d') such that each agent handles exactly one task, and each task can only be handled by at most one agent.
- When k out of d agents are assigned, we need to assign the remaining d' k tasks to d k agents.

 $\Box \quad I_k = \log \frac{d'!}{(d'-k)!}$ 

# Sequence tagging problems:

- □ Shallow Parsing, NER, etc., are key examples.
- O cannot be immediately followed by I.
- No closed-form solution to  $I_k$ ; need dynamic programming simulations.

 $I_{k} - I_{k-1}$ 

# Algorithm: CoDL & SSPAN

**CoDL**: constraint-driven learning; can be seen as "structured self-learning" [ACL'07] **SSPAN**: structured self-learning with partial annotations; can be seen as a straightforward extension of CoDL [\*SEM'18]

## The SSPAN Algorithm

Input:  $\mathcal{T} = \{(x_i, y_i)\}_{i=1}^N, \mathcal{P} = \{(x_i, a_i)\}_{i=N+1}^{N+M}$ 1. Initialize  $\mathcal{H} = LEARN(\mathcal{T})$ 2. While convergence not reached, do 1)  $\tilde{\mathcal{P}} = \emptyset$ 2) Foreach  $(x_i, a_i) \in \mathcal{P}$  do *i.*  $\hat{y}_i = INFERENCE(x_i; \mathcal{H}), \text{ s.t.}$   $\hat{y}_i \in C(\mathcal{Y}^d)$  $\hat{y}_{i,j} = a_{i,j}, \forall a_{i,j} \neq \Box \checkmark$ 



**Define:**  $I_k \triangleq \log |C(\mathcal{L}^d)| - E[\log f_k]$ • Measures how much of  $C(\mathcal{L}^d)$  has been **disqualified** by k labels.

**Define:** A *k*-partial annotation  $A_k$  is a vector of random variables  $A_k = [A_{k,1}, A_{k,2}, ..., A_{k,d}] \in (\mathcal{L} \cup \Pi)^d$ , where  $\Pi$  is a special character for no label yet, s.t.

- $\square \ \sum_{i=1}^{d} \mathbb{I}(A_{k,i} \neq \sqcap) = k$
- $\square P(Y|A_k = a_k) = P(Y|Y_j = a_{k,j}, j \in \mathcal{J}), \text{ where } \mathcal{J} = \{j: a_{k,j} \neq \sqcap\}$
- $\Box$   $A_k$  means k variables in Y are correctly labeled

**Theorem:**  $I_k$  is the **mutual information** between Y and  $A_k$  when both Y and the k variables labeled in  $A_k$  follow uniform distributions.



20 30 40 50 60 70 80 90 100 Completeness k/d (%)

#### Implications of the curves above:

Diminishing return of new labels

- Better to annotate a structure partially than completely
- The slope may be an indicator for the strengths of structures



*ii.*  $\tilde{\mathcal{P}} = \tilde{\mathcal{P}} \cup \{ \boldsymbol{x}_i, \boldsymbol{\hat{y}}_i \}$ 3)  $\mathcal{H} = LEARN(\mathcal{T} + \tilde{\mathcal{P}})$ 3. Return  $\mathcal{H}$ 

SSPAN is conceptually a **hard-EM** algorithm for structured learning tasks. Without  $\rightleftharpoons$ , SSPAN goes back to CoDL.

## Conclusion

This paper provides a **unique view of structured annotations**: it is the reduction in the uncertainty of a target structure *Y*, by a random process *A* representing the annotation process.

- □ **Theoretically**, "partial" provides more information
- Empirically, "partial" leads to improved performance on three very different tasks

In general, we argue that any signal that has **non-zero mutual information** with Y can be viewed as "annotation".

#### <u>References</u>

**[ToC'91]** G. Brightwell and P. Winkler. Counting linear extensions is #p-complete. **[CoNLL'00]** TK Sang and S. Buchholz. Introduction to the CoNLL-2000 shared task: Chunking. **[NIPS'01]** V. Punyakanok and D. Roth. The use of classifiers in sequential inference. **[LREC'02]** P. Kingsbury and M. Palmer. From Treebank to PropBank. **[CoNLL'05]** X. Carreras and L. Marquez. Introduction to the CoNLL-2005 shared tasks: Semantic role labeling. **[ACL'07]** MW Chang, L. Ratinov, and D. Roth. Guiding semisupervision with constraint-driven learning. **[ACL'18a]** Q. Ning, Z. Feng, H. Wu, and D. Roth. Joint reasoning for temporal and causal relations. **[ACL'18b]** Q. Ning, H. Wu, and D. Roth. A multi-axis annotation scheme for event temporal relations. **[CVPR'18]** J. Choi, J. Krishnamurthy, A. Kembhavi, and A. Farhadi. Structured set matching networks for one-short part labeling. **[EMNLP'18]** Q. Ning, B. Zhou, Z. Feng, H. Peng, and D. Roth. CogCompTime: A tool for understanding time in natural language. **[LREC'18]** D. Khashabi et al. CogCompNLP: Your swiss army knife for NLP. **[\*SEM'18]** Q. Ning, Z. Yu, C. Fan, and D. Roth. Exploiting partially annotated data for temporal relation extraction.

#### **Acknowledgement**

The Allen Institute for Artificial Intelligence (allenai.org); the IBM-ILLINOIS Center for Cognitive Computing Systems Research (C3SR) - a research collaboration as part of the IBM AI Horizons Network; Contract HR0011-15-2-0025 with the US Defense Advanced Research Projects Agency (DARPA); the Army Research Laboratory (ARL).



![](_page_0_Picture_69.jpeg)

![](_page_0_Picture_70.jpeg)

![](_page_0_Picture_71.jpeg)