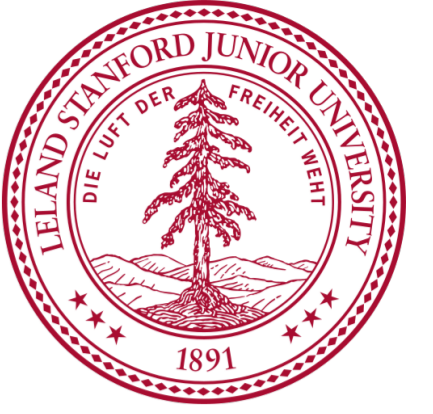


# Learning with Relaxed Supervision

Jacob Steinhardt Percy Liang

{jsteinhardt, pliang}@cs.stanford.edu



## Intractable Supervision

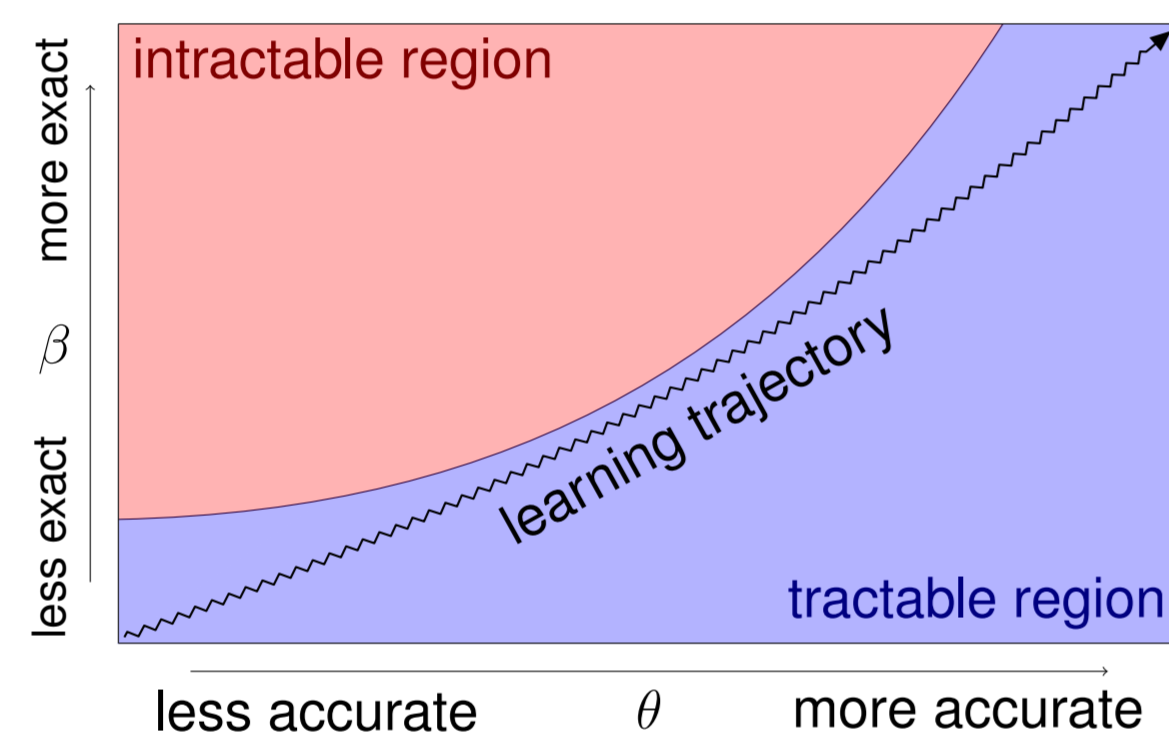
For weakly-supervised tasks, inference can be intractable:

input  $x$ : What is the largest city in California?  
 latent  $z$ :  $\text{argmax}(\lambda x.\text{CITY}(x) \wedge \text{LOC}(x, \text{CA}), \lambda x.\text{POPULATION}(x))$   
 output  $y$ : Los Angeles

Computing  $p(z | x, y)$  requires inverting arbitrary logical forms!

- Still want to exploit likely statistical relationships (CITY and Los Angeles)
- Need a way to relax the supervision so we can learn tractably.
- Want to maintain good statistical properties (asymptotic consistency).

## Our Approach



- Start with intractable supervision  $q_\infty(y | z)$
- Replace with family of relaxed supervision functions  $q_\beta(y | z)$
- Derive constraints on  $(\theta, \beta)$  that ensure tractability of inference
- Optimize likelihood within the tractable region

Intuition:

- Supervision is intractable if too harsh relative to model accuracy.
- Initially need very forgiving supervision, can eventually incorporate full supervision (done adaptively over course of optimization).

## The Relaxation

Assume relationship between  $z$  and  $y$  given by constraints  $\mathbb{S}_j, j = 1, \dots, k$  (think machine translation, checking that each word of the output is correct).

Relaxation based on weighted count of constraint violations:

$$p_\beta(y | z) \propto \exp\left(-\sum_{j=1}^k \beta_j (1 - \mathbb{S}_j(z, y))\right) \quad (\dagger)$$

When  $\beta = 0$ ,  $p_\beta$  is uniform; when  $\beta = \infty$ ,  $p_\beta$  is original supervision.

**Challenges:** normalization constant of  $p_\beta$ ; ensuring tractable inference.

## Framework

Assumptions:

- $x \rightarrow z \rightarrow y$ , where  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is observed and  $z \in \mathcal{Z}$  is unobserved.
- Parameterized family  $p_\theta(z | x)$ .
- $z \rightarrow y$  is a known deterministic function  $y = f(z)$ .

Hence, letting  $\mathbb{S}(z, y) \in \{0, 1\}$  denote the constraint  $[y = f(z)]$ :

$$p_\theta(y | x) = \sum_z \mathbb{S}(z, y) p_\theta(z | x).$$

Goal: decompose  $\mathbb{S}$  into smaller components  $\mathbb{S}_j$ .

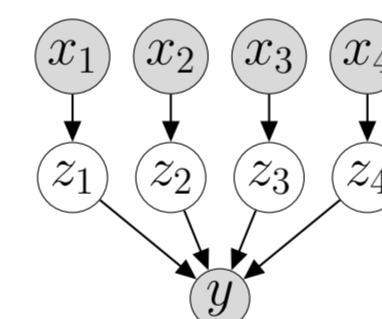
- Define projections  $\pi_j: \mathcal{Y} \rightarrow \mathcal{Y}_j$ .
- Projected constraint:  $\mathbb{S}_j(z, y) \stackrel{\text{def}}{=} [\pi_j(f(z)) = \pi_j(y)]$ .
- If  $\pi_1 \times \dots \times \pi_k$  is one-to-one, then can decompose  $\mathbb{S}$  as  $\mathbb{S} = \bigwedge_{j=1}^k \mathbb{S}_j$ .

## Example Decompositions

Translation from Unordered Supervision

Goal: infer substitution cipher Model  $p_\theta(z | x)$ : soft substitutions

input  $x$ : a b a a  
 latent  $z$ : d c d d  
 output  $y$ : {c: 1, d: 3}  
 (cipher: {a  $\mapsto$  d, b  $\mapsto$  c, ...})



Supervision:  $y = \text{multiset}(z)$

Decomposition ( $y$  and  $z$  match if all counts match):

$$[y = \text{multiset}(z)] \Leftrightarrow \bigwedge_{j=1}^V [\text{count}(z, j) = \text{count}(y, j)]$$

Conjunctive Semantic Parsing

Side information: predicates  $Q_1, \dots, Q_m$ .

- e.g.  $Q_6 = [\text{DOG}] = \text{set of all dogs}$

input  $x$ : brown dog (input utterance)  
 latent  $z$ : ( $Q_{11}, Q_6$ ) (set of all brown objects, set of all dogs)  
 output  $y$ :  $Q_{11} \cap Q_6$  (denotation, observed as a set)

For  $z = (Q_{j_1}, \dots, Q_{j_L})$ , define the denotation  $[[z]] = Q_{j_1} \cap \dots \cap Q_{j_L}$ .

Decomposition ( $y$  and  $z$  match if contained in same predicates):

$$[y = [[z]]] \Leftrightarrow \bigwedge_{j=1}^m \mathbb{I}[[z] \subseteq Q_j] = \mathbb{I}[y \subseteq Q_j]$$

## Theory

**Lemma (normalization constant).** For any  $z$ , the log-normalization constant of  $p_\beta(y | z)$  is bounded above by

$$A(\beta) \stackrel{\text{def}}{=} \sum_{j=1}^k \log(1 + (|\mathcal{Y}_j| - 1) \exp(-\beta_j)).$$

**Lemma (asymptotic consistency).** Suppose that we use  $A(\beta)$  above as a surrogate normalization constant for  $p_\beta$ . Then, the MLE of  $(\theta, \beta)$  asymptotically recovers the true model parameters.

## Tractability Constraints

Typical expression for gradient (for some features  $\phi(x, z, y)$ ):

$$\nabla \log p_{\theta, \beta}(y | x) = \underbrace{\mathbb{E}_{z|x, y}[\phi(x, z, y)]}_{\text{model + supervision}} - \underbrace{\mathbb{E}_{z|x}[\phi(x, z, y)]}_{\text{model}}.$$

To learn, need to sample  $p_{\theta, \beta}(z | x, y) \propto p_\theta(z | x) \exp(\beta^\top \mathbb{S}_{1:k}(z, y))$  (see  $(\dagger)$ ).

- For large  $\beta$ , this is as intractable as the original supervision.
- Need a way to constrain  $\beta$  to yield tractable inference.

Inference algorithm: rejection sampling.

- Sample from  $p_\theta(z | x)$ , accept with probability  $p_\beta(y | z)$ .

Constrain expected number of rejections based on computational budget  $\tau$ :

$$\begin{array}{ll} \text{minimize} & \mathbb{E}_{x, y}[-\log p_{\theta, \beta}(y | x)] \quad (\mathcal{L}) \\ \text{subject to} & \mathbb{E}_{x, y}[\text{Rejections}(x, y)] \leq \tau \quad (\mathcal{C}) \end{array}$$

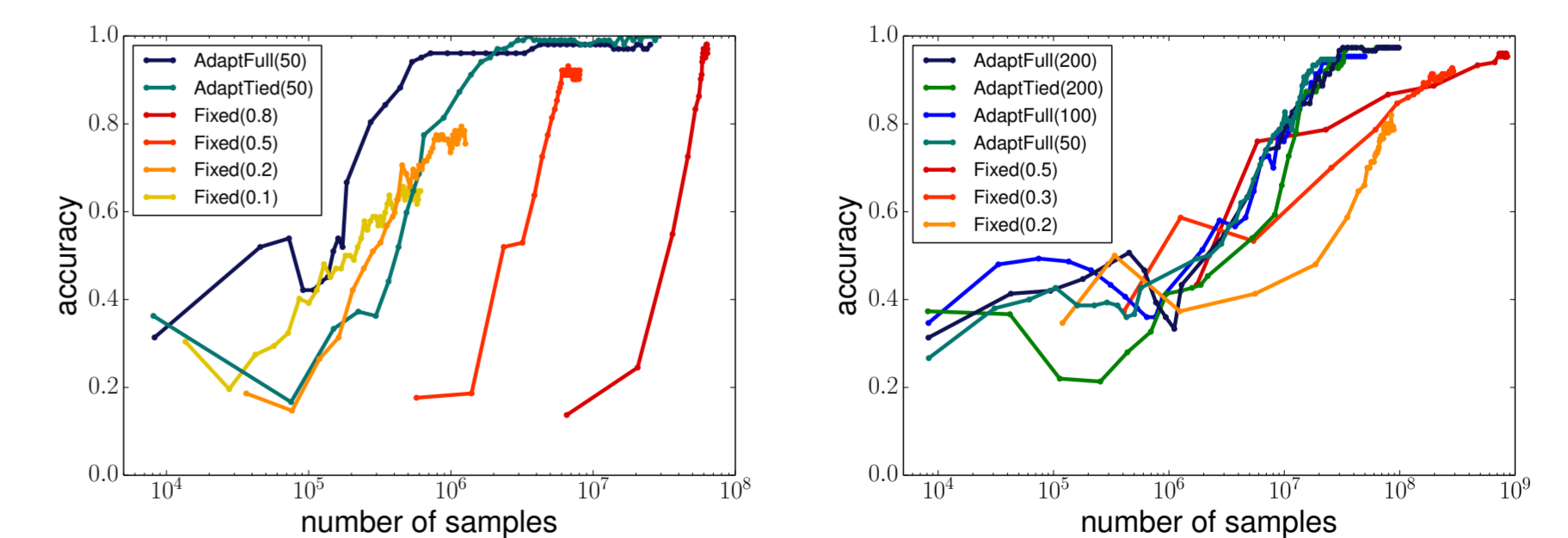
Amazingly,  $(\mathcal{C})$  is well-behaved enough to admit an EM-like procedure for constrained optimization! (See paper for full details.)

## Experiments

Implemented our relaxed supervision algorithm on both the unordered translation and conjunctive semantic parsing tasks.

Compared a fixed value of relaxation  $\beta$  (FIXED) to optimizing  $\beta$  subject to our tractability constraints (ADAPT).

Our tractability constraints improve efficiency by orders of magnitude while also improving accuracy:



(a) unordered translation (b) conjunctive semantic parsing

Reproducible experiments on CodaLab: [worksheets.codalab.org](https://worksheets.codalab.org)