

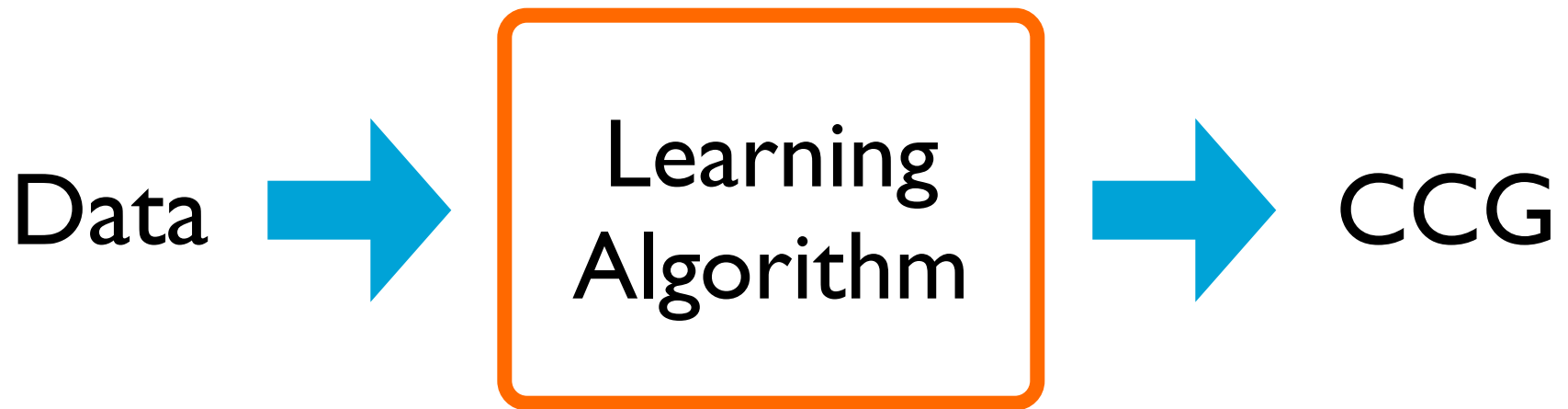
# Semantic Parsing with Combinatory Categorical Grammars

Yoav Artzi, Nicholas FitzGerald and Luke Zettlemoyer  
University of Washington

ACL 2013 Tutorial  
Sofia, Bulgaria



# Learning



- What kind of data/supervision we can use?
- What do we need to learn?

# Supervised Data

show	me	flights	to	Boston
$S/N$		$N$	$PP/NP$	$NP$
$\lambda f.f$		$\lambda x.flight(x)$	$\lambda y.\lambda x.to(x, y)$	$BOSTON$
			$PP$	$\lambda x.to(x, BOSTON)$
			$N \setminus N$	$\lambda f.\lambda x.f(x) \wedge to(x, BOSTON)$
		$N$	$\lambda x.flight(x) \wedge to(x, BOSTON)$	
		$S$	$\lambda x.flight(x) \wedge to(x, BOSTON)$	

# Supervised Data

show	me	flights	to	Boston
$S/N$		$N$	$PP/NP$	$NP$
$\lambda f.f$		$\lambda x.flight(x)$	$\lambda y.\lambda x.to(x, y)$	$BOSTON$
				$\rightarrow$
			$PP$	
			$\lambda x.to(x, BOSTON)$	
				$\rightarrow$
			$N \setminus N$	
			$\lambda f.\lambda x.f(x) \wedge to(x, BOSTON)$	
				$\leftarrow$
			$N$	
			$\lambda x.flight(x) \wedge to(x, BOSTON)$	
				$\rightarrow$
			$S$	
			$\lambda x.flight(x) \wedge to(x, BOSTON)$	

Latent

# Supervised Data

Supervised learning is done from pairs  
of sentences and logical forms

Show me flights to Boston

$\lambda x. flight(x) \wedge to(x, BOSTON)$

I need a flight from baltimore to seattle

$\lambda x. flight(x) \wedge from(x, BALTIMORE) \wedge to(x, SEATTLE)$

what ground transportation is available in san francisco

$\lambda x. ground\_transport(x) \wedge to\_city(x, SF)$

# Weak Supervision

- Logical form is latent
- “Labeling” requires less expertise
- Labels don’t uniquely determine correct logical forms
- Learning requires executing logical forms within a system and evaluating the result

# Weak Supervision

## Learning from Query Answers

What is the largest state that borders Texas?

*New Mexico*

# Weak Supervision

## Learning from Query Answers

What is the largest state that borders Texas?

*New Mexico*

$\text{argmax}(\lambda x. \text{state}(x)$   
 $\wedge \text{border}(x, TX), \lambda y. \text{size}(y))$

$\text{argmax}(\lambda x. \text{river}(x)$   
 $\wedge \text{in}(x, TX), \lambda y. \text{size}(y))$



# Weak Supervision

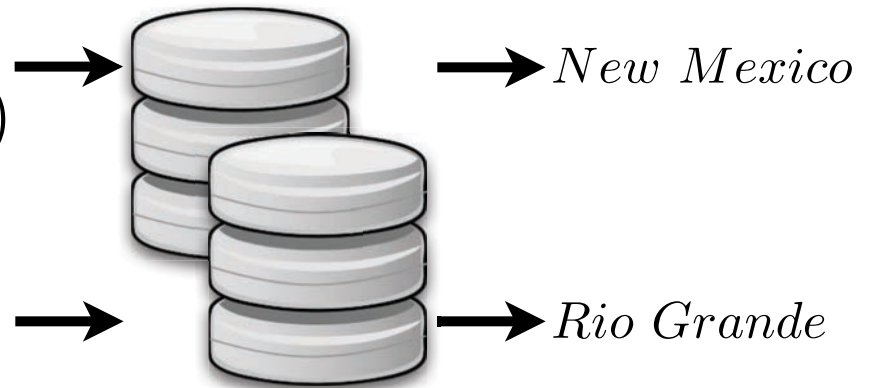
## Learning from Query Answers

What is the largest state that borders Texas?

*New Mexico*

$\text{argmax}(\lambda x. \text{state}(x)$   
 $\wedge \text{border}(x, TX), \lambda y. \text{size}(y))$

$\text{argmax}(\lambda x. \text{river}(x)$   
 $\wedge \text{in}(x, TX), \lambda y. \text{size}(y))$



# Weak Supervision

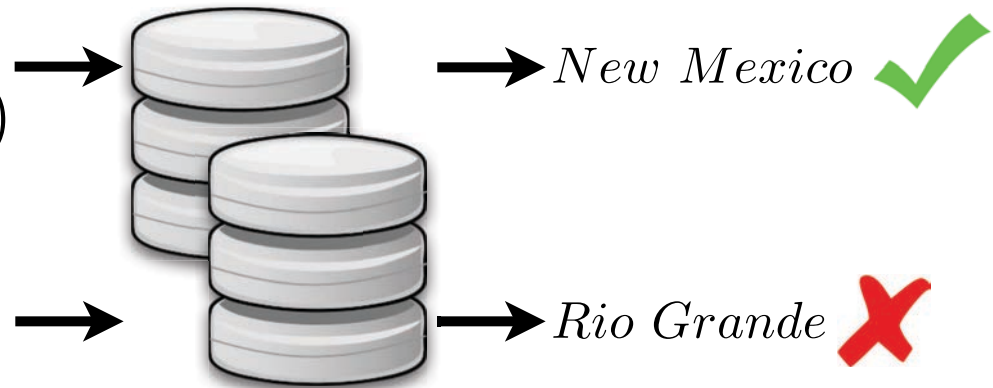
## Learning from Query Answers

What is the largest state that borders Texas?

*New Mexico*

$\text{argmax}(\lambda x. \text{state}(x) \wedge \text{border}(x, TX), \lambda y. \text{size}(y))$

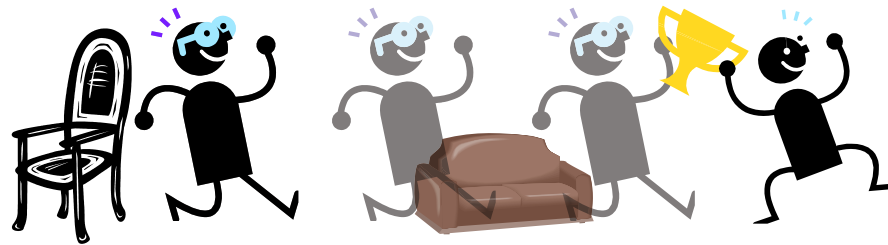
$\text{argmax}(\lambda x. \text{river}(x) \wedge \text{in}(x, TX), \lambda y. \text{size}(y))$



# Weak Supervision

## Learning from Demonstrations

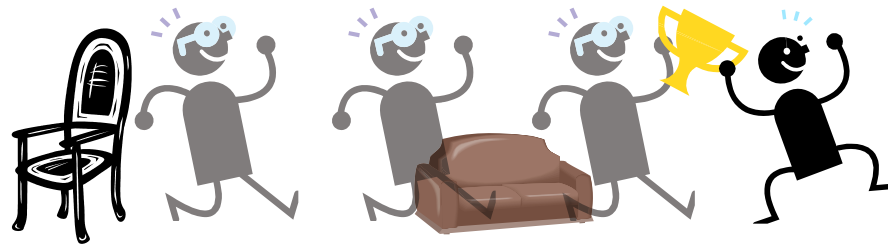
at the chair, move forward three steps past the sofa



# Weak Supervision

## Learning from Demonstrations

at the chair, move forward three steps past the sofa



Some examples from other domains:

- Sentences and labeled game states [Goldwasser and Roth 2011]
- Sentences and sets of physical objects [Matuszek et al. 2012]

Parsing

Learning

Modeling

- Structured perceptron
- A unified learning algorithm
- Supervised learning
- Weak supervision

# Structured Perceptron

- Simple additive updates
  - Only requires efficient decoding ( $\text{argmax}$ )
  - Closely related to maxent and other feature rich models
  - Provably finds linear separator in finite updates, if one exists
- Challenge: learning with hidden variables

# Structured Perceptron

**Data:**  $\{(x_i, y_i) : i = 1 \dots n\}$

For  $t = 1 \dots T$  :

[iterate epochs]

For  $i = 1 \dots n$ :

[iterate examples]

$$y^* \leftarrow \arg \max_y \langle \theta, \Phi(x_i, y) \rangle$$

[predict]

If  $y^* \neq y_i$ :

[check]

$$\theta \leftarrow \theta + \Phi(x_i, y_i) - \Phi(x_i, y^*)$$

[update]

# One Derivation of the Perceptron

Log-linear model: 
$$p(y|x) = \frac{e^{w \cdot f(x,y)}}{\sum_{y'} e^{w \cdot f(x,y')}}$$

Step 1: Differentiate, to maximize data log-likelihood

$$update = \sum_i f(x_i, y_i) - E_{p(y|x_i)} f(x_i, y)$$

Step 2: Use online, stochastic gradient updates, for example  $i$ :

$$update_i = f(x_i, y_i) - E_{p(y|x_i)} f(x_i, y)$$

Step 3: Replace expectations with maxes (Viterbi approx.)

$$update_i = f(x_i, y_i) - f(x_i, y^*) \text{ where } y^* = \arg \max_y w \cdot f(x_i, y)$$



# The Perceptron with Hidden Variables

Log-linear

model:  $p(y|x) = \sum_h p(y, h|x) \quad p(y, h|x) = \frac{e^{w \cdot f(x, h, y)}}{\sum_{y', h'} e^{w \cdot f(x, h', y')}}$

Step 1: Differentiate marginal, to maximize data log-likelihood

$$update = \sum_i E_{p(h|y_i, x_i)} [f(x_i, h, y_i)] - E_{p(y, h|x_i)} [f(x_i, h, y)]$$

Step 2: Use online, stochastic gradient updates, for example  $i$ :

$$update_i = E_{p(y_i, h|x_i)} [f(x_i, h, y_i)] - E_{p(y, h|x_i)} [f(x_i, h, y)]$$

Step 3: Replace expectations with maxes (Viterbi approx.)

$$update_i = f(x_i, h', y_i) - f(x_i, h^*, y^*) \text{ where}$$

$$y^*, h^* = \arg \max_{y, h} w \cdot f(x_i, h, y) \quad \text{and} \quad h' = \arg \max_h w \cdot f(x_i, h, y_i)$$

# Hidden Variable Perceptron

**Data:**  $\{(x_i, y_i) : i = 1 \dots n\}$

For  $t = 1 \dots T$  : [iterate epochs]

For  $i = 1 \dots n$ : [iterate examples]

$y^*, h^* \leftarrow \arg \max_{y, h} \langle \theta, \Phi(x_i, h, y) \rangle$  [predict]

If  $y^* \neq y_i$ : [check]

$h' \leftarrow \arg \max_h \langle \theta, \Phi(x_i, h, y_i) \rangle$  [predict hidden]

$\theta \leftarrow \theta + \Phi(x_i, h', y_i) - \Phi(x_i, h^*, y^*)$  [update]

# Hidden Variable Perceptron

- No known convergence guarantees
  - Log-linear version is non-convex
- Simple and easy to implement
  - Works well with careful initialization
- Modifications for semantic parsing
  - Lots of different hidden information
  - Can add a margin constraint, do probabilistic version, etc.

# Learning Choices

## Validation Function

$$\mathcal{V} : \mathcal{Y} \rightarrow \{t, f\}$$

- Indicates correctness of a parse  $y$
- Varying  $\mathcal{V}$  allows for differing forms of supervision

## Lexical Generation Procedure

$$GENLEX(x, \mathcal{V}; \Lambda, \theta)$$

- Given:
  - sentence  $x$
  - validation function  $\mathcal{V}$
  - lexicon  $\Lambda$
  - parameters  $\theta$
- Produce an overly general set of lexical entries

Initialize  $\theta$  using  $\Lambda_0$  ,  $\Lambda \leftarrow \Lambda_0$

For  $t = 1 \dots T, i = 1 \dots n$  :

**Step 1:** (Lexical generation)

- a. Set  $\lambda_G \leftarrow GENLEX(x_i, \mathcal{V}_i; \Lambda, \theta)$ ,  
 $\lambda \leftarrow \Lambda \cup \lambda_G$
- b. Let  $Y$  be the  $k$  highest scoring parses from  $GEN(x_i; \lambda)$
- c. Select lexical entries from the highest scoring valid parses:  
$$\lambda_i \leftarrow \bigcup_{y \in MAXV_i(Y; \theta)} LEX(y)$$
- d. Update lexicon:  $\Lambda \leftarrow \Lambda \cup \lambda_i$

**Step 2:** (Update parameters)

**Output:** Parameters  $\theta$  and lexicon  $\Lambda$

# Unification-based $GENLEX(x, z; \Lambda, \theta)$

I want a flight to Boston  
 $\lambda x.flight(x) \wedge to(x, BOS)$

1. Find highest scoring correct parse
2. Find splits that most increases score
3. Return new lexical entries

Iteration 2

