# Bringing machine learning & compositional semantics together: approaches

https://github.com/cgpotts/annualreview-complearning

Chris Potts
Stanford Linguistics

CS 244U: Natural language understanding

## Semantic parsing

$$\langle \boxed{u, t, r}, d \rangle$$

## Basic formulation

|       | Utterance                  | Logical form            |
|-------|----------------------------|-------------------------|
|       | seven minus five           | $(- \ 7 \ 5)$           |
|       | five minus seven           | $(- \ 5 \ 7)$           |
|       | three plus one             | $(- \ 7 \ 5)$           |
|       | minus three plus one       | $(+ \ \neg 3 \ 1)$      |
| Train | minus three plus one       | $\neg(+ \ 3 \ 1)$       |
|       | two minus two times two    | $(\times \ (- \ 2 \ 2) \ 2)$ |
|       | two minus two times two    | $(- \ 2 \ (\times \ 2 \ 2))$ |
|       | two plus three plus four   | $(+ \ 2 \ (+ \ 3 \ 4))$ |
|       | ⋮                          |                         |
|       | three minus one            | ?                       |
|       | three times one            | ?                       |
| Test  | minus six times four       | ?                       |
|       | one plus three plus five   | ?                       |
|       | ⋮                          |                         |

Table: Data requirements.

## Basic formulation

| | Utterance | Logical form |
|---|---|---|
| | seven minus five | (− 7 5) |
| | five minus seven | (− 5 7) |
| | three plus one | (− 7 5) |
| | minus three plus one | (+ ¬3 1) |
| Train | minus three plus one | ¬(+ 3 1) |
| | two minus two times two | (× (− 2 2) 2) |
| | two minus two times two | (− 2 (× 2 2)) |
| | two plus three plus four | (+ 2 (+ 3 4)) |
| | ⋮ | |
| | three minus one | ? |
| | three times one | ? |
| Test | minus six times four | ? |
| | one plus three plus five | ? |
| | ⋮ | |

Table: Data requirements.

| Syntax | Logical form |
|---|---|
| N → one | 1 |
| N → one | 2 |
| | ⋮ |
| N → two | 1 |
| N → two | 2 |
| | ⋮ |
| R → plus | + |
| R → plus | − |
| R → plus | × |
| R → minus | + |
| R → minus | − |
| R → minus | × |
| R → times | + |
| R → times | − |
| R → times | × |
| S → minus | ¬ |

N → S N ⌜S⌝⌜N⌝
N → N$_L$ R N$_R$ (⌜R⌝ ⌜N$_L$⌝ ⌜N$_R$⌝)

Table: Crude grammar.

# Learning framework

# Learning framework

1. Feature representations: $\phi(x, y) \in \mathbb{R}^d$

## Learning framework

1. Feature representations: $\phi(x, y) \in \mathbb{R}^d$
2. Scoring: $\text{Score}_{\mathbf{w}}(x, y) = \sum_{j=1}^{d} w_j \phi(x, y)_j$

## Learning framework

1. Feature representations: $\phi(x, y) \in \mathbb{R}^d$
2. Scoring: $\text{Score}_{\mathbf{w}}(x, y) = \sum_{j=1}^{d} w_j \phi(x, y)_j$
3. Multiclass hinge-loss objective function:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{(x,y) \in \mathcal{D}} \max_{y' \in \text{GEN}(x)} \left[ \text{Score}_{\mathbf{w}}(x, y') + c(y, y') \right] - \text{Score}_{\mathbf{w}}(x, y)$$

where $\mathcal{D}$ is a set of $(x, y)$ training examples and $c(a, b) = 1$ if $a \neq b$, else 0.

## Learning framework

1. Feature representations: $\phi(x, y) \in \mathbb{R}^d$
2. Scoring: $\text{Score}_{\mathbf{w}}(x, y) = \sum_{j=1}^{d} w_j \phi(x, y)_j$
3. Multiclass hinge-loss objective function:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{(x,y) \in \mathcal{D}} \max_{y' \in \text{Gen}(x)} [\text{Score}_{\mathbf{w}}(x, y') + c(y, y')] - \text{Score}_{\mathbf{w}}(x, y)$$

where $\mathcal{D}$ is a set of $(x, y)$ training examples and $c(a, b) = 1$ if $a \neq b$, else 0.
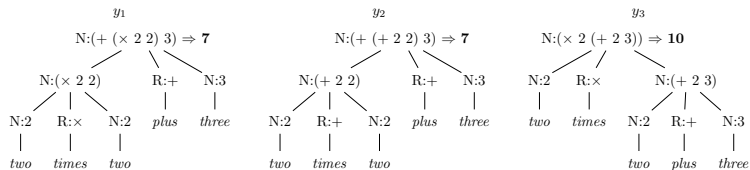
4. Optimization:
   StochasticGradientDescent($\mathcal{D}, T, \eta$)

   1   Initialize $\mathbf{w} \leftarrow \mathbf{0}$
   2   Repeat $T$ times
   3       **for** each $(x, y) \in \mathcal{D}$ (in random order)
   4            $\tilde{y} \leftarrow \arg\max_{y' \in \text{Gen}(x)} \text{Score}_{\mathbf{w}}(x, y') + c(y, y')$
   5            $\mathbf{w} \leftarrow \mathbf{w} + \eta(\phi(x, y) - \phi(x, \tilde{y}))$
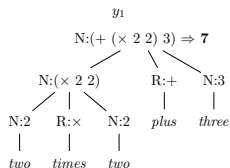   6   Return $\mathbf{w}$

# Example



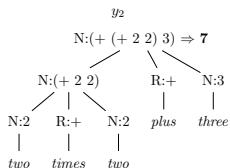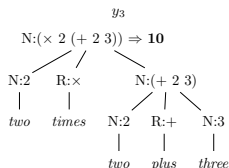(a) **Candidates** GEN($x$) for utterance $x = $ *two times two plus three*

# Example

(a) **Candidates** GEN($x$) for utterance $x = $ *two times two plus three*

## Example

**(a) Candidates** GEN($x$) for utterance $x = $ *two times two plus three*



**(b) Learning from logical forms (Section 4.1)**

## Derivational ambiguity

| Syntax | Logical form |
|---|---|
| N → one | 1 |
| N → two | 2 |
| | ⋮ |
| R → plus | + |
| R → minus | − |
| R → times | × |
| S → minus | ¬ |
| N → S N | ⌜S⌝⌜N⌝ |
| N → N_L R N_R | (⌜R⌝ ⌜N_L⌝ ⌜N_R⌝) |
| Q → n | (λf (f ⌜n⌝)) |
| N → U Q | (⌜Q⌝⌜U⌝) |

Table: Grammar with type-lifting.

Training instance: (*minus three*, ¬3)

$$N : \neg 3$$

```
        N : ¬3
       /      \
    U : ¬    N : 3
      |        |
    minus    three
```

$$N : ((\lambda f\ (f\ 3))\ \neg) \stackrel{\beta}{\Rightarrow} \neg 3$$

```
    U : ¬    Q : (λf (f 3))
      |             |
    minus         three
```

(Beta-conversion $\stackrel{\beta}{\Rightarrow}$ is the syntactic counterpart
of functional application.)

## Derivations as latent variables

- The training instances are $(u, r)$ pairs.

- Since $r$ might have multiple derivations, derivations are latent variables.

- Latent support vector machine objective:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{(x,r) \in \mathcal{D}} \max_{y' \in \text{Gen}(x)} [\text{Score}_{\mathbf{w}}(x, y') + c(r, \text{Root}(y'))] - \max_{y'' \in \text{Gen}(x,r)} \text{Score}_{\mathbf{w}}(x, y''),$$

where $\mathcal{D}$ is a set of (utterance, formula) pairs; $c(a, b) = 1$ if $a \neq b$, else 0; and $\text{Gen}(x, r) = \{y \in \text{Gen}(x) : \text{Root}(y) = r\}$

- Optimization:
  StochasticGradientDescent$(\mathcal{D}, T, \eta)$

```
1  Initialize w ← 0
2  Repeat T times
3      for each (x, r) ∈ D (in random order)
4          y ← arg max_{y''∈Gen(x,r)} Score_w(x, y'')
5          ỹ ← arg max_{y'∈Gen(x)} Score_w(x, y') + c(y, y')
6          w ← w + η(φ(x, y) − φ(x, ỹ))
7  Return w
```

## Learning from denotations

$$\langle \boxed{u, t, r, d} \rangle$$

## Motivations

### Semantic parsing

- What is the largest city in California?
- $\arg\max(\{c : \text{city}(c) \wedge \text{loc}(c, \text{CA})\}, \text{population})$

### Interpretive

- What is the largest city in California?
- Los Angeles.

Semantic parsing · · · · · · · · · · · · · · · · · · · · · · · · Derivations as latent variables · · · · · · · · · · · · · · · · · · · · · · · · Learning from denotations

○○○ · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · ○○ · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · ○●○○

## Basic formulation

|       | Utterance                | Denotation |
|-------|--------------------------|-----------:|
|       | seven minus five         | 2          |
|       | five minus seven         | −2         |
|       | three plus one           | 4          |
|       | minus three plus one     | −2         |
| Train | minus three plus one     | −4         |
|       | two minus two times two  | 0          |
|       | two minus two times two  | −2         |
|       | two plus three plus four | 9          |
|       | ⋮                        |            |
|       | three minus one          | ?          |
|       | three times one          | ?          |
| Test  | minus six times four     | ?          |
|       | one plus three plus five | ?          |
|       | ⋮                        |            |

Table: Data requirements.

## Basic formulation

| | Utterance | Denotation |
|---|---|---|
| | seven minus five | 2 |
| | five minus seven | −2 |
| | three plus one | 4 |
| | minus three plus one | −2 |
| Train | minus three plus one | −4 |
| | two minus two times two | 0 |
| | two minus two times two | −2 |
| | two plus three plus four | 9 |
| | ⋮ | |
| | three minus one | ? |
| | three times one | ? |
| Test | minus six times four | ? |
| | one plus three plus five | ? |
| | ⋮ | |

Table: Data requirements.

| Syntax | Logical form | Denotation |
|---|---|---|
| N → one | 1 | 1 |
| N → one | 2 | 2 |
| | ⋮ | |
| N → two | 1 | 1 |
| N → two | 2 | 2 |
| | ⋮ | |
| R → plus | + | addition |
| R → plus | − | subtraction |
| R → plus | × | multiplication |
| R → minus | + | addition |
| R → minus | − | subtraction |
| R → minus | × | multiplication |
| R → times | + | addition |
| R → times | − | subtraction |
| R → times | × | multiplication |
| S → minus | ¬ | negative |
| N → S N | $\ulcorner S \urcorner \ulcorner N \urcorner$ | $[\![ \ulcorner S \urcorner ]\!]([\![ \ulcorner N \urcorner ]\!])$ |
| N → $N_L$ R $N_R$ | $(\ulcorner R \urcorner \ulcorner \ulcorner N_L \urcorner \urcorner \ulcorner \ulcorner N_R \urcorner \urcorner)$ | $[\![ \ulcorner R \urcorner ]\!]([\![ \ulcorner N_L \urcorner ]\!], [\![ \ulcorner N_R \urcorner ]\!])$ |

Table: Crude grammar.

## Learning framework

Feature representations and scoring are as before.

1. Latent support vector machine objective:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{(x,d) \in \mathcal{D}} \max_{y' \in \text{Gen}(x)} [\text{Score}_{\mathbf{w}}(x, y') + c(d, [\![y']\!])] - \max_{y \in \text{Gen}(x,d)} \text{Score}_{\mathbf{w}}(x, y),$$

where $\text{Gen}(x, d) = \{y \in \text{Gen}(x) : [\![y]\!] = d\}$ is the set of logical forms that evaluate to denotation $d$.
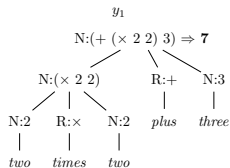
2. Optimization:
StochasticGradientDescent($\mathcal{D}, T, \eta$)

```
1   Initialize w ← 0
2   Repeat T times
3       for each (x, d) ∈ D (in random order)
4           y ← arg max_{y''∈GEN(x,d)} Score_w(x, y'')
5           ỹ ← arg max_{y'∈GEN(x)} Score_w(x, y') + c(y, y')
6           w ← w + η(φ(x, y) − φ(x, ỹ))
7   Return w
```
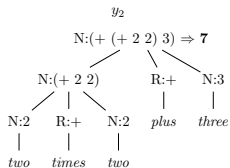
# Example

**(a) Candidates** GEN($x$) for utterance $x = two\ times\ two\ plus\ three$



$$\phi(x, y_1) = \begin{array}{|l|} \hline \text{R:}\times[times]:1 \\ \text{R:}+[plus]:1 \\ \text{top}[\text{R:}+]:1 \\ \hline \end{array}$$

$$\phi(x, y_2) = \begin{array}{|l|} \hline \text{R:}+[times]:1 \\ \text{R:}+[plus]:1 \\ \text{top}[\text{R:}+]:1 \\ \hline \end{array}$$

$$\phi(x, y_3) = \begin{array}{|l|} \hline \text{R:}\times[times]:1 \\ \text{R:}+[plus]:1 \\ \text{top}[\text{R:}\times]:1 \\ \hline \end{array}$$

## Example

**(a) Candidates** GEN($x$) for utterance $x = $ *two times two plus three*



$$\phi(x, y_1) = \begin{array}{|l|} \hline \text{R:}\times[times]:1 \\ \text{R:}+[plus]:1 \\ \text{top[R:}+]:1 \\ \hline \end{array}$$

$$\phi(x, y_2) = \begin{array}{|l|} \hline \text{R:}+[times]:1 \\ \text{R:}+[plus]:1 \\ \text{top[R:}+]:1 \\ \hline \end{array}$$

$$\phi(x, y_3) = \begin{array}{|l|} \hline \text{R:}\times[times]:1 \\ \text{R:}+[plus]:1 \\ \text{top[R:}\times]:1 \\ \hline \end{array}$$

**(c) Learning from denotations (Section 4.2)**

**Iteration 1**

$$\mathbf{w} = \begin{array}{|l|} \hline \text{R:}\times[times]:0 \\ \text{R:}+[times]:0 \\ \text{R:}+[plus]:0 \\ \text{top[R:}+]:0 \\ \text{top[R:}\times]:0 \\ \hline \end{array}$$

Scores:$[0, 0, 0]$
GEN($x, d$) $= \{y_1, y_2\}$
$y = y_1$ (tied with $y_2$)
$\tilde{y} = y_3$

$\Rightarrow$

**Iteration 2**

$$\mathbf{w} = \begin{array}{|l|} \hline \text{R:}\times[times]:0 \\ \text{R:}+[times]:0 \\ \text{R:}+[plus]:0 \\ \text{top[R:}+]:1 \\ \text{top[R:}\times]:\text{-}1 \\ \hline \end{array}$$

Scores:$[1, 1, -1]$
GEN($x, d$) $= \{y_1, y_2\}$
$y = y_1$ (tied with $y_2$)
$\tilde{y} = y_1$ (tied with $y_2$)