# On Dual Decomposition
# and Linear Programming Relaxations
# for Natural Language Processing

Alexander M. Rush, David Sontag,
Michael Collins, and Tommi Jaakkola

# Dynamic Programming

Dynamic programming is a dominant technique in NLP.

- ▶ Fast
- ▶ Exact
- ▶ Easy to implement

Examples:

- ▶ Viterbi algorithm for hidden Markov models
- ▶ CKY algorithm for weighted context-free grammars

$$y^* = \arg\max_y f(y) \leftarrow \text{Decoding}$$
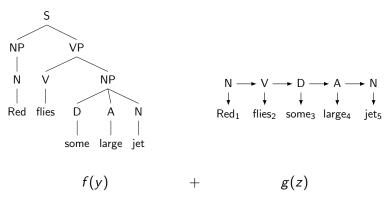
# Model Complexity

Unfortunately, dynamic programming algorithms do not scale well with model complexity.

As our models become complex, these algorithms can explode in terms of computational or implementational complexity.

Integration:
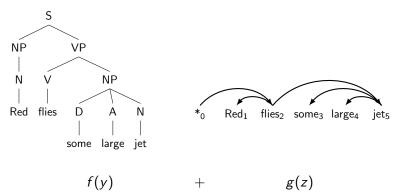
- $f \leftarrow$ Easy
- $g \leftarrow$ Easy
- $f + g \leftarrow$ Hard

# Integration (1)



$$f(y) \qquad + \qquad g(z)$$

▶ Classical problem in NLP.

▶ The dynamic programming intersection is prohibitively slow and complicated to implement.

# Integration (2)



$$f(y) \quad\quad + \quad\quad g(z)$$

- Important for improving parsing accuracy.

- The dynamic programming intersection is slow and complicated to implement.

# Dual Decomposition

A general technique for constructing decoding algorithms

Solve complicated models

$$y^* = \arg\max_y f(y)$$

by decomposing into smaller problems.

Upshot: Can utilize a toolbox of combinatorial algorithms.

- ▶ Dynamic programming
- ▶ Minimum spanning tree
- ▶ Shortest path
- ▶ Min-Cut
- ▶ ...

# Dual Decomposition Algorithms

Simple - Uses basic dynamic programming algorithms

Efficient - Faster than full dynamic programming intersections

Strong Guarantees - Gives a certificate of optimality when exact

In experiments, we find the global optimum on 99% of examples.

Widely Applicable - Similar techniques extend to other problems

Algorithm

Experiments

LP Relaxations

# Integrated Parsing and Tagging



HMM



CFG

# Integrated Parsing and Tagging



HMM

Dual Decomposition

CFG

# HMM for Tagging

$$N \longrightarrow V \longrightarrow D \longrightarrow A \longrightarrow N$$

$$\downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow$$

$$\text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

▶ Let $\mathcal{Z}$ be the set of all valid taggings of a sentence and $g(z)$ be a scoring function.

e.g. $g(z) = \log p(\mathrm{Red}_1|\mathrm{N}) + \log p(\mathrm{V}|\mathrm{N}) + ...$

# HMM for Tagging

$$N \longrightarrow V \longrightarrow D \longrightarrow A \longrightarrow N$$

$$\downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow$$

$$\text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

- Let $\mathcal{Z}$ be the set of all valid taggings of a sentence and $g(z)$ be a scoring function.

e.g. $g(z) = \log p(\text{Red}_1|\text{N}) + \log p(\text{V}|\text{N}) + ...$

$$z^* = \arg\max_{z \in \mathcal{Z}} g(z) \leftarrow \text{Viterbi decoding}$$

# CFG for Parsing



- Let $\mathcal{Y}$ be the set of all valid parse trees for a sentence and $f(y)$ be a scoring function.

e.g. $f(y) = \log p(\text{S} \rightarrow \text{NP VP}|\text{S}) + \log p(\text{NP} \rightarrow \text{N}|\text{NP}) + ...$

# CFG for Parsing



▶ Let $\mathcal{Y}$ be the set of all valid parse trees for a sentence and $f(y)$ be a scoring function.

e.g. $f(y) = \log p(\mathrm{S} \to \mathrm{NP\ VP}|\mathrm{S}) + \log p(\mathrm{NP} \to \mathrm{N}|\mathrm{NP}) + ...$

$$y^* = \arg\max_{y \in \mathcal{Y}} f(y) \leftarrow \text{CKY Algorithm}$$

Problem Definition

- Find parse tree that optimizes

$$score(\mathrm{S} \rightarrow \mathrm{NP\ VP}) + score(\mathrm{VP} \rightarrow \mathrm{V\ NP}) +$$

$$... + score(\mathrm{Red_1}, \mathrm{N}) + score(\mathrm{V}, \mathrm{N}) + ...$$

- Conventional Approach (Bar Hillel et al., 1961)
  - Replace rules like $S \rightarrow NP\ VP$

    with rules like $S_{N,N} \rightarrow NP_{N,V}\ VP_{V,N}$

    Painful. $O(t^6)$ increase in complexity for trigram tagging.

# The Integrated Parsing and Tagging Problem

Find $\underset{y \in \mathcal{Y},\ z \in \mathcal{Z}}{\text{argmax}}$ $f(y) + g(z)$

such that for all $i, t$, $y(i, t) = z(i, t)$

Where $y(i, t) = 1$ if parse includes tag $t$ at position $i$

$z(i, t) = 1$ if tagging includes tag $t$ at position $i$

$$\text{Find} \quad \underset{y \in \mathcal{Y}, \ z \in \mathcal{Z}}{\text{argmax}} \quad f(y) \ + \ g(z)$$

Trees

such that for all $i, t, \ y(i, t) = z(i, t)$

Where $y(i, t) = 1$ if parse includes tag $t$ at position $i$

$z(i, t) = 1$ if tagging includes tag $t$ at position $i$

Find $\quad \underset{y \in \mathcal{Y},\ z \in \mathcal{Z}}{\text{argmax}} \quad f(y) \; + \; g(z)$

Trees $\qquad$ Taggings

such that for all $i, t, \ \ y(i, t) = z(i, t)$

Where $y(i, t) = 1$ if parse includes tag $t$ at position $i$

$z(i, t) = 1$ if tagging includes tag $t$ at position $i$

# The Integrated Parsing and Tagging Problem

Find      argmax      $f(y) + g(z)$

$y \in \mathcal{Y}, z \in \mathcal{Z}$

such that for all $i, t,\; y(i, t) = z(i, t)$

Where $y(i, t) = 1$ if parse includes tag $t$ at position $i$

$z(i, t) = 1$ if tagging includes tag $t$ at position $i$

Find    argmax    $f(y) \ + \ g(z)$

$y \in \mathcal{Y}, \ z \in \mathcal{Z}$

Trees    Taggings

CFG    HMM

such that for all $i, t, \ y(i, t) = z(i, t)$

Where $y(i, t) = 1$ if parse includes tag $t$ at position $i$

$z(i, t) = 1$ if tagging includes tag $t$ at position $i$

Find  argmax  $f(y) + g(z)$

CFG    HMM

$y \in \mathcal{Y},\ z \in \mathcal{Z}$

Trees    Taggings

such that for all $i, t,\ \ y(i, t) = z(i, t)$

Constraints

Where $y(i, t) = 1$ if parse includes tag $t$ at position $i$

$z(i, t) = 1$ if tagging includes tag $t$ at position $i$

Set penalty weights equal to 0 for the tag at each position.

**For** $k = 1$ **to** $K$

# Algorithm Sketch

Set penalty weights equal to 0 for the tag at each position.

**For** $k = 1$ **to** $K$

$\quad y^{(k)} \leftarrow$ Decode $(f(y) + \text{penalty})$ by CKY Algorithm

Set penalty weights equal to 0 for the tag at each position.

**For** $k = 1$ **to** $K$

$y^{(k)} \leftarrow$ Decode $(f(y) + \mathrm{penalty})$ by CKY Algorithm

$z^{(k)} \leftarrow$ Decode $(g(z) - \mathrm{penalty})$ by Viterbi Decoding

# Algorithm Sketch

Set penalty weights equal to 0 for the tag at each position.

**For** $k = 1$ **to** $K$

$y^{(k)} \leftarrow$ Decode $(f(y) + \text{penalty})$ by CKY Algorithm

$z^{(k)} \leftarrow$ Decode $(g(z) - \text{penalty})$ by Viterbi Decoding

**If** $y^{(k)}(i, t) = z^{(k)}(i, t)$ for all $i, t$ **Return** $(y^{(k)}, z^{(k)})$

Set penalty weights equal to 0 for the tag at each position.

**For** $k = 1$ **to** $K$

$y^{(k)} \leftarrow$ Decode $(f(y) + \text{penalty})$ by CKY Algorithm

$z^{(k)} \leftarrow$ Decode $(g(z) - \text{penalty})$ by Viterbi Decoding

**If** $y^{(k)}(i, t) = z^{(k)}(i, t)$ for all $i, t$ **Return** $(y^{(k)}, z^{(k)})$

**Else** Update penalty weights based on $y^{(k)}(i, t) - z^{(k)}(i, t)$

# CKY Parsing

$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

## Viterbi Decoding

$$\text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

## Key

| | | | | | |
|---|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ | Taggings |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ | | | |

# CKY Parsing

$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

# Viterbi Decoding

$$\text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

# Key

| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

# CKY Parsing

```
                    S
         NP                   VP
   A     N     D      A       V
   |     |     |      |       |
  Red  flies some  large    jet
```

$$y^* = \arg \max_{y \in \mathcal{Y}} (f(y) + \sum_{i,t} u(i, t) y(i, t))$$

# Viterbi Decoding

```
   N ──→ V ──→ D ──→ A ──→ N
   ↓      ↓      ↓      ↓      ↓
 Red₁  flies₂ some₃ large₄ jet₅
```

$$z^* = \arg \max_{z \in \mathcal{Z}} (g(z) - \sum_{i,t} u(i, t) z(i, t))$$

# Key

| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i, t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

# CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

# Viterbi Decoding



$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

## Penalties

$u(i,t) = 0$ for all $i,t$

## Key

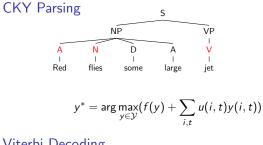| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

# CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}\left(f(y) + \sum_{i,t} u(i,t)y(i,t)\right)$$

# Viterbi Decoding



$$z^* = \arg\max_{z \in \mathcal{Z}}\left(g(z) - \sum_{i,t} u(i,t)z(i,t)\right)$$

# Key

| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

# Penalties

$u(i,t) = 0$ for all $i, t$

| Iteration 1 | |
|---|---|
| $u(1, A)$ | -1 |
| $u(1, N)$ | 1 |
| $u(2, N)$ | -1 |
| $u(2, V)$ | 1 |
| $u(5, V)$ | -1 |
| $u(5, N)$ | 1 |

## CKY Parsing

$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

## Viterbi Decoding

$$\text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

## Penalties

$u(i,t) = 0$ for all $i,t$

| Iteration 1 | |
|---|---|
| $u(1,A)$ | -1 |
| $u(1,N)$ | 1 |
| $u(2,N)$ | -1 |
| $u(2,V)$ | 1 |
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

## Key

| | | | | | |
|---|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ | Taggings |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ | | | |

# CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

# Viterbi Decoding

$$\text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

# Key

| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

# Penalties

$u(i,t) = 0$ for all $i,t$

| Iteration 1 | |
|---|---|
| $u(1,A)$ | -1 |
| $u(1,N)$ | 1 |
| $u(2,N)$ | -1 |
| $u(2,V)$ | 1 |
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

## CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}\left(f(y) + \sum_{i,t} u(i,t)y(i,t)\right)$$

## Viterbi Decoding



$$z^* = \arg\max_{z \in \mathcal{Z}}\left(g(z) - \sum_{i,t} u(i,t)z(i,t)\right)$$

## Penalties

$u(i,t) = 0$ for all $i,t$

| Iteration 1 | |
|---|---|
| $u(1,A)$ | -1 |
| $u(1,N)$ | 1 |
| $u(2,N)$ | -1 |
| $u(2,V)$ | 1 |
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

## Key

| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

# CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

# Viterbi Decoding



$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

# Key

| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

# Penalties

$u(i,t) = 0$ for all $i,t$

| Iteration 1 | |
|---|---|
| $u(1, A)$ | -1 |
| $u(1, N)$ | 1 |
| $u(2, N)$ | -1 |
| $u(2, V)$ | 1 |
| $u(5, V)$ | -1 |
| $u(5, N)$ | 1 |

| Iteration 2 | |
|---|---|
| $u(5, V)$ | -1 |
| $u(5, N)$ | 1 |

# CKY Parsing

$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

# Viterbi Decoding

Red$_1$  flies$_2$  some$_3$  large$_4$  jet$_5$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

# Penalties

$u(i,t) = 0$ for all $i,t$

| Iteration 1 | |
|---|---|
| $u(1,A)$ | -1 |
| $u(1,N)$ | 1 |
| $u(2,N)$ | -1 |
| $u(2,V)$ | 1 |
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

| Iteration 2 | |
|---|---|
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

# Key

| | | | | | |
|---|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ | Taggings |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ | | | |

# CKY Parsing

```
                    S
         ┌──────────┴──────────┐
        NP                    VP
         │          ┌──────────┴──────────┐
         N          V                    NP
         │          │          ┌──────────┼──────────┐
        Red       flies        D          A          N
                               │          │          │
                             some       large       jet
```

$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

# Viterbi Decoding

$\text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

# Penalties

$u(i,t) = 0$ for all $i,t$

| Iteration 1 | |
|---|---|
| $u(1,A)$ | -1 |
| $u(1,N)$ | 1 |
| $u(2,N)$ | -1 |
| $u(2,V)$ | 1 |
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

| Iteration 2 | |
|---|---|
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

# Key

| | | | | | |
|---|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ | Taggings |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ | | | |

## CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

## Viterbi Decoding



$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

## Key

| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

## Penalties

$u(i,t) = 0$ for all $i, t$

| Iteration 1 | |
|---|---|
| $u(1, A)$ | -1 |
| $u(1, N)$ | 1 |
| $u(2, N)$ | -1 |
| $u(2, V)$ | 1 |
| $u(5, V)$ | -1 |
| $u(5, N)$ | 1 |

| Iteration 2 | |
|---|---|
| $u(5, V)$ | -1 |
| $u(5, N)$ | 1 |

## CKY Parsing



$$y^* = \arg \max_{y \in \mathcal{Y}}(f(y) + \sum_{i,t} u(i,t)y(i,t))$$

## Viterbi Decoding



$$z^* = \arg \max_{z \in \mathcal{Z}}(g(z) - \sum_{i,t} u(i,t)z(i,t))$$

### Key

| | | |
|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees |
| $y(i,t) = 1$ | if | $y$ contains tag $t$ at position $i$ |

| | | |
|---|---|---|
| $g(z)$ | $\Leftarrow$ | HMM |
| $\mathcal{Z}$ | $\Leftarrow$ | Taggings |

## Penalties

$u(i,t) = 0$ for all $i,t$

| Iteration 1 | |
|---|---|
| $u(1,A)$ | -1 |
| $u(1,N)$ | 1 |
| $u(2,N)$ | -1 |
| $u(2,V)$ | 1 |
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

| Iteration 2 | |
|---|---|
| $u(5,V)$ | -1 |
| $u(5,N)$ | 1 |

### Converged

$$y^* = \arg \max_{y \in \mathcal{Y}} f(y) + g(y)$$

**Theorem**
If at any iteration $y^{(k)}(i, t) = z^{(k)}(i, t)$ for all $i, t$, then $(y^{(k)}, z^{(k)})$ is the global optimum.

In experiments, we find the global optimum on 99% of examples.

Theorem

If at any iteration $y^{(k)}(i,t) = z^{(k)}(i,t)$ for all $i, t$, then $(y^{(k)}, z^{(k)})$ is the global optimum.
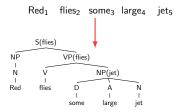
In experiments, we find the global optimum on 99% of examples.

If we do not converge to a match, we can still get a result (more in paper).
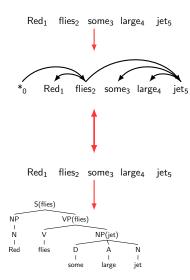
# Integrated CFG and Dependency Parsing

$Red_1$   $flies_2$   $some_3$   $large_4$   $jet_5$

Dependency Model

$*_0$   $Red_1$   $flies_2$   $some_3$   $large_4$   $jet_5$

$Red_1$   $flies_2$   $some_3$   $large_4$   $jet_5$

Lexicalized CFG

```
            S(flies)
      NP          VP(flies)
       |
       N        V         NP(jet)
       |        |
      Red     flies    D      A      N
                       |      |      |
                     some   large   jet
```

# Integrated CFG and Dependency Parsing



Dependency Model

Dual Decomposition

Lexicalized CFG

# Dependency Parsing



$*_0$    $\text{Red}_1$    $\text{flies}_2$    $\text{some}_3$    $\text{large}_4$    $\text{jet}_5$

- Let $\mathcal{Z}$ be the set of all valid dependency parses of a sentence and $g(z)$ be a scoring function.

e.g. $g(z) = \log p(\text{some}_3 | \text{jet}_5, \text{large}_4) + ...$

# Dependency Parsing



- Let $\mathcal{Z}$ be the set of all valid dependency parses of a sentence and $g(z)$ be a scoring function.

e.g. $g(z) = \log p(\text{some}_3 | \text{jet}_5, \text{large}_4) + ...$

$$z^* = \arg\max_{z \in \mathcal{Z}} g(z) \leftarrow \text{Eisner (2000) algorithm}$$

# Lexicalized PCFG



- Let $\mathcal{Y}$ be the set of all valid dependency parses of a sentence and $f(y)$ be a scoring function.

e.g. $f(y) = \log p(\text{S(flies)} \rightarrow \text{NP(Red) VP(flies)}|\text{S(flies)}) + ...$

# Lexicalized PCFG



- Let $\mathcal{Y}$ be the set of all valid dependency parses of a sentence and $f(y)$ be a scoring function.

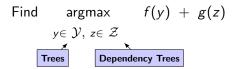e.g. $f(y) = \log p(\text{S(flies)} \rightarrow \text{NP(Red) VP(flies)}|\text{S(flies)}) + ...$

$$y^* = \arg \max_{y \in \mathcal{Y}} f(y) \leftarrow \text{Modified CKY algorithm}$$

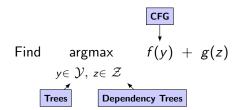$$\text{Find} \quad \underset{y \in \mathcal{Y}, \ z \in \mathcal{Z}}{\operatorname{argmax}} \quad f(y) \ + \ g(z)$$
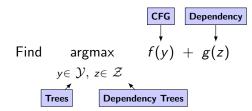
such that for all $i, j, \ y(i,j) = z(i,j)$

Where $y(i,j) = 1$ if parse includes dependency from word $i$ to $j$

$z(i,j) = 1$ if parse includes dependency from word $i$ to $j$

$$\text{Find} \quad \underset{y \in \mathcal{Y}, \ z \in \mathcal{Z}}{\text{argmax}} \quad f(y) \ + \ g(z)$$

Trees

such that for all $i, j, \ \ y(i,j) = z(i,j)$

Where $y(i,j) = 1$ if parse includes dependency from word $i$ to $j$

$z(i,j) = 1$ if parse includes dependency from word $i$ to $j$

$$\text{Find} \quad \underset{y \in \mathcal{Y}, \; z \in \mathcal{Z}}{\text{argmax}} \quad f(y) \; + \; g(z)$$

Trees     Dependency Trees

such that for all $i, j, \; y(i,j) = z(i,j)$

Where $y(i,j) = 1$ if parse includes dependency from word $i$ to $j$

$z(i,j) = 1$ if parse includes dependency from word $i$ to $j$

# The Integrated Constituency and Dependency Parsing Problem

$$\text{Find} \quad \underset{y \in \mathcal{Y}, \ z \in \mathcal{Z}}{\text{argmax}} \quad f(y) + g(z)$$

CFG

Trees

Dependency Trees

such that for all $i, j, \ y(i,j) = z(i,j)$

Where $y(i,j) = 1$ if parse includes dependency from word $i$ to $j$

$z(i,j) = 1$ if parse includes dependency from word $i$ to $j$

# The Integrated Constituency and Dependency Parsing Problem

Find $\quad$ argmax $\quad f(y) + g(z)$

$$y \in \mathcal{Y}, \ z \in \mathcal{Z}$$

CFG | Dependency

Trees | Dependency Trees

such that for all $i, j, \ y(i,j) = z(i,j)$

Where $y(i,j) = 1$ if parse includes dependency from word $i$ to $j$

$z(i,j) = 1$ if parse includes dependency from word $i$ to $j$

Find  argmax  $f(y) + g(z)$

CFG    Dependency

$y \in \mathcal{Y}, \ z \in \mathcal{Z}$

Trees    Dependency Trees

such that for all $i, j, \ y(i,j) = z(i,j)$

Constraints

Where $y(i,j) = 1$ if parse includes dependency from word $i$ to $j$

$z(i,j) = 1$ if parse includes dependency from word $i$ to $j$

## CKY Parsing

$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,j} u(i,j)y(i,j))$$

## Dependency Parsing

$$*_0 \quad \text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,j} u(i,j)z(i,j))$$

## Key

| | | | | | |
|---|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ | Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ | Dependency Trees |
| $y(i,j) = 1$ | if | $y$ contains dependency $i,j$ | | | |

## CKY Parsing

$$u(i,j) = 0 \text{ for all } i,j$$



$$y^* = \arg\max_{y \in \mathcal{Y}}\left(f(y) + \sum_{i,j} u(i,j)y(i,j)\right)$$

## Dependency Parsing

$$*_0 \quad \text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

$$z^* = \arg\max_{z \in \mathcal{Z}}\left(g(z) - \sum_{i,j} u(i,j)z(i,j)\right)$$

## Key

| | | | | |
|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ $\Leftarrow$ | Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ $\Leftarrow$ | Dependency Trees |
| $y(i,j) = 1$ | if | $y$ contains dependency $i,j$ | | |

# CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,j} u(i,j)y(i,j))$$

# Penalties

$u(i,j) = 0$ for all $i,j$

# Dependency Parsing



$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,j} u(i,j)z(i,j))$$

# Key

| | | | | | |
|---|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ | Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ | Dependency Trees |
| $y(i,j) = 1$ | if | $y$ contains dependency $i,j$ | | | |

## CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,j} u(i,j)y(i,j))$$

## Dependency Parsing



$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,j} u(i,j)z(i,j))$$

## Key

| | | | | | |
|---|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ | Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ | Dependency Trees |
| $y(i,j) = 1$ | if | $y$ contains dependency $i,j$ | | | |

# CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,j} u(i,j)y(i,j))$$

# Dependency Parsing



$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,j} u(i,j)z(i,j))$$

## Penalties

$u(i,j) = 0$ for all $i,j$

| Iteration 1 | |
| --- | --- |
| $u(2,3)$ | -1 |
| $u(5,3)$ | 1 |

## Key

| | | | | |
| --- | --- | --- | --- | --- |
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ Dependency Trees |
| $y(i,j) = 1$ | if | $y$ contains dependency $i,j$ | | |

CKY Parsing

$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,j} u(i,j)y(i,j))$$

Dependency Parsing

$$*_0 \quad Red_1 \quad flies_2 \quad some_3 \quad large_4 \quad jet_5$$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,j} u(i,j)z(i,j))$$

Key

| | | | | |
| --- | --- | --- | --- | --- |
| $f(y)$ | $\Leftarrow$ | CFG | $g(z) \Leftarrow$ | Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z} \quad \Leftarrow$ | Dependency Trees |
| $y(i,j) = 1$ | if | $y$ contains dependency $i,j$ | | |

# CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,j} u(i,j)y(i,j))$$

## Penalties

$u(i,j) = 0$ for all $i,j$

| Iteration 1 | |
|---|---|
| $u(2,3)$ | -1 |
| $u(5,3)$ | 1 |

# Dependency Parsing

$$*_0 \quad \text{Red}_1 \quad \text{flies}_2 \quad \text{some}_3 \quad \text{large}_4 \quad \text{jet}_5$$

$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,j} u(i,j)z(i,j))$$

## Key

| | | | | |
|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ CFG | | $g(z)$ | $\Leftarrow$ Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ Parse Trees | | $\mathcal{Z}$ | $\Leftarrow$ Dependency Trees |
| $y(i,j) = 1$ | if $y$ contains dependency $i,j$ | | | |

# CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}} (f(y) + \sum_{i,j} u(i,j) y(i,j))$$

## Penalties

$u(i,j) = 0$ for all $i,j$

| Iteration 1 | |
|---|---|
| $u(2,3)$ | -1 |
| $u(5,3)$ | 1 |

# Dependency Parsing



$$z^* = \arg\max_{z \in \mathcal{Z}} (g(z) - \sum_{i,j} u(i,j) z(i,j))$$

## Key

| | | | | | |
|---|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ | Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ | Dependency Trees |
| $y(i,j) = 1$ | if | $y$ contains dependency $i,j$ | | | |

## CKY Parsing



$$y^* = \arg\max_{y \in \mathcal{Y}}(f(y) + \sum_{i,j} u(i,j)y(i,j))$$

## Dependency Parsing



$$z^* = \arg\max_{z \in \mathcal{Z}}(g(z) - \sum_{i,j} u(i,j)z(i,j))$$

## Penalties

$u(i,j) = 0$ for all $i,j$

| Iteration 1 | |
|---|---|
| $u(2,3)$ | -1 |
| $u(5,3)$ | 1 |

**Converged**

$$y^* = \arg\max_{y \in \mathcal{Y}} f(y) + g(y)$$

## Key

| | | | | |
|---|---|---|---|---|
| $f(y)$ | $\Leftarrow$ | CFG | $g(z)$ | $\Leftarrow$ Dependency Model |
| $\mathcal{Y}$ | $\Leftarrow$ | Parse Trees | $\mathcal{Z}$ | $\Leftarrow$ Dependency Trees |
| $y(i,j) = 1$ | if | $y$ contains dependency $i,j$ | | |

Algorithm

Experiments

LP Relaxations

# Experiment

Properties:

- ► Exactness
- ► Parsing Accuracy

Experiments on:

- ► English Penn Treebank

Models

- ► Collins (1997) Model 1
- ► Semi-Supervised Dependency Parser (Koo, 2008)
- ► Trigram Tagger (Toutanova, 2000)

# How quickly do the models converge?



Integrated Dependency Parsing

Integrated POS Tagging

# Integrated Constituency and Dependency Parsing: Accuracy



$F_1$ Score

- ▶ Collins (1997) Model 1
- ▶ Fixed, First-best Dependencies from Koo (2008)
- ▶ Dual Decomposition

## Integrated Parsing and Tagging: Accuracy



$F_1$ Score

- ▶ Fixed, First-Best Tags From Toutanova (2000)
- ▶ Dual Decomposition

Algorithm

Experiments

LP Relaxations

Theorem

- ▶ If the dual decomposition algorithm converges, then $(y^{(k)}, z^{(k)})$ is the global optimum.

Questions

- ▶ What problem is dual decomposition solving?
- ▶ How come the algorithm doesn't always converge?

Dual decomposition searches over a <span style="color:red">linear programming relaxation</span> of the original problem.

# Convex Hulls for CKY

A parse tree can be represented as a binary vector $y \in \mathcal{Y}$.
$y(A \to B\ C, i, j, k) = 1$ if rule $A \to B\ C$ is used at span $i, j, k$.
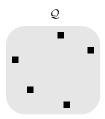


Parsing

- ▶ If $f$ is linear, $\arg \max_{y \in conv(\mathcal{Y})} f(y)$ is a linear program.
- ▶ The best point in an LP is a vertex. So CKY solves this LP.

# Convex Hulls for CKY

A parse tree can be represented as a binary vector $y \in \mathcal{Y}$.
$y(A \to B\ C, i, j, k) = 1$ if rule $A \to B\ C$ is used at span $i, j, k$.



Parsing

- If $f$ is linear, $\arg\max\limits_{y \in conv(\mathcal{Y})} f(y)$ is a linear program.
- The best point in an LP is a vertex. So CKY solves this LP.

# Convex Hulls for CKY

A parse tree can be represented as a binary vector $y \in \mathcal{Y}$.
$y(A \rightarrow B\ C, i, j, k) = 1$ if rule $A \rightarrow B\ C$ is used at span $i, j, k$.



conv($\mathcal{Y}$)

Parsing

- If $f$ is linear, $\arg\max\limits_{y \in conv(\mathcal{Y})} f(y)$ is a linear program.
- The best point in an LP is a vertex. So CKY solves this LP.

# Convex Hulls for CKY

A parse tree can be represented as a binary vector $y \in \mathcal{Y}$.
$y(A \rightarrow B\ C, i, j, k) = 1$ if rule $A \rightarrow B\ C$ is used at span $i, j, k$.



Parsing

- If $f$ is linear, $\arg\max\limits_{y \in conv(\mathcal{Y})} f(y)$ is a linear program.
- The best point in an LP is a vertex. So CKY solves this LP.

## Combined Problem

$$\mathcal{Q} = \{(y, z) \colon y \in \mathcal{Y}, z \in \mathcal{Z},$$
$$y(i, t) = z(i, t) \text{ for all } (i, t)\}$$



$\mathcal{Q}$

# Combined Problem
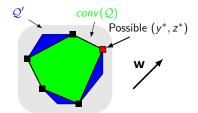
$$\mathcal{Q} = \{(y, z) \colon y \in \mathcal{Y}, z \in \mathcal{Z},$$
$$y(i, t) = z(i, t) \text{ for all } (i, t)\}$$

$\mathcal{Q}$

# Combined Problem

$$\mathcal{Q} = \{(y, z) \colon y \in \mathcal{Y}, z \in \mathcal{Z},$$
$$y(i, t) = z(i, t) \text{ for all } (i, t)\}$$

conv($\mathcal{Q}$)

$$\mathcal{Q} = \{(y, z) \colon y \in \mathcal{Y}, z \in \mathcal{Z},$$
$$y(i, t) = z(i, t) \text{ for all } (i, t)\}$$

$\mathcal{Q}'$
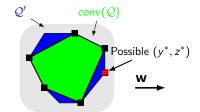


$$\mathcal{Q}' = \{(\mu, \nu) \colon \mu \in \text{conv}(\mathcal{Y}), \nu \in \text{conv}(\mathcal{Z}),$$
$$\mu(i, t) = \nu(i, t) \text{ for all } (i, t)\}$$

Dual decomposition searches over $\mathcal{Q}'$

## Combined Problem

$$\mathcal{Q} = \{(y, z) \colon y \in \mathcal{Y}, z \in \mathcal{Z},$$
$$y(i, t) = z(i, t) \text{ for all } (i, t)\}$$



$$\mathcal{Q}' = \{(\mu, \nu) \colon \mu \in \mathrm{conv}(\mathcal{Y}), \nu \in \mathrm{conv}(\mathcal{Z}),$$
$$\mu(i, t) = \nu(i, t) \text{ for all } (i, t)\}$$

Dual decomposition searches over $\mathcal{Q}'$

# Combined Problem

$$\mathcal{Q} = \{(y, z) : y \in \mathcal{Y}, z \in \mathcal{Z},$$
$$y(i, t) = z(i, t) \text{ for all } (i, t)\}$$



$$\mathcal{Q}' = \{(\mu, \nu) : \mu \in \text{conv}(\mathcal{Y}), \nu \in \text{conv}(\mathcal{Z}),$$
$$\mu(i, t) = \nu(i, t) \text{ for all } (i, t)\}$$

Dual decomposition searches over $\mathcal{Q}'$

Depending on the weight vector, $(y^*, z^*) \in \mathcal{Q}'$ could be in $\mathcal{Q}$ or in the strict outer bound.

# Combined Problem

$$\mathcal{Q} = \{(y, z) : y \in \mathcal{Y}, z \in \mathcal{Z},$$
$$y(i, t) = z(i, t) \text{ for all } (i, t)\}$$



$$\mathcal{Q}' = \{(\mu, \nu) : \mu \in \text{conv}(\mathcal{Y}), \nu \in \text{conv}(\mathcal{Z}),$$
$$\mu(i, t) = \nu(i, t) \text{ for all } (i, t)\}$$

Dual decomposition searches over $\mathcal{Q}'$

Depending on the weight vector, $(y^*, z^*) \in \mathcal{Q}'$ could be in $\mathcal{Q}$ or in the strict outer bound.

# Are there points strictly in the outer bound?



$\mathcal{Q}'$

Possible $(y^*, z^*)$?

**Taggings**   0.5x   A → A → A   + 0.5x   A → B → B
              ↓    ↓    ↓         ↓    ↓    ↓
             $w_1$  $w_2$  $w_3$        $w_1$  $w_2$  $w_3$

Best result can be a
fractional solution.

Convex
combination of
these structures.

**Parses**   0.5 x

```
      X                      X
     / \                    / \
    A   X        + 0.5 x   A   X
    |  / \                 |  / \
   w_1 A  B               w_1 B  A
       |  |                   |  |
      w_2 w_3               w_2 w_3
```

# Summary

A Dual Decomposition algorithm for integrated decoding

Simple - Uses only simple, off-the-shelf dynamic programming algorithms to solve a harder problem.

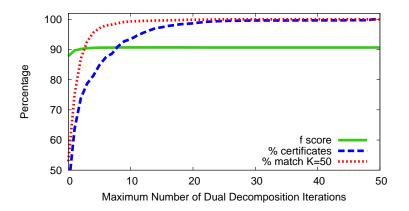Efficient - Faster than classical methods for dynamic programming intersection.

Strong Guarantees - Solves a linear programming relaxation which gives a certificate of optimality.

Finds the exact solution on 99% of the examples.

Widely Applicable - Similar techniques extend to other problems

Appendix

Iterative Progress

## Deriving the Algorithm

Goal:
$$y^* = \arg\max_{y \in \mathcal{Y}} f(y)$$

Rewrite:
$$\arg\max_{z \in \mathcal{Z}, y \in \mathcal{Y}} f(z) + g(y)$$

$$\text{s.t. } z(i,j) = y(i,j) \text{ for all } i,j$$

Lagrangian: $L(u, y, z) = f(z) + g(y) + \sum_{i,j} u(i,j) \left( y(i,j) - z(i,j) \right)$

# Deriving the Algorithm

Goal:
$$y^* = \arg\max_{y \in \mathcal{Y}} f(y)$$

Rewrite:
$$\arg\max_{z \in \mathcal{Z}, y \in \mathcal{Y}} f(z) + g(y)$$

s.t. $z(i,j) = y(i,j)$ for all $i, j$

Lagrangian: $L(u, y, z) = f(z) + g(y) + \sum_{i,j} u(i,j)\left(y(i,j) - z(i,j)\right)$

The dual problem is to find $\min_u L(u)$ where

$$L(u) = \max_{y \in \mathcal{Y}, z \in \mathcal{Z}} L(u, y, z) = \max_{z \in \mathcal{Z}} \left( f(z) + \sum_{i,j} u(i,j)z(i,j) \right)$$

$$+ \max_{y \in \mathcal{Y}} \left( g(y) - \sum_{i,j} u(i,j)y(i,j) \right)$$

Dual is an upper bound: $L(u) \geq f(z^*) + g(y^*)$ for any $u$

# A Subgradient Algorithm for Minimizing $L(u)$

$$L(u) \;=\; \max_{z \in \mathcal{Z}} \left( f(z) + \sum_{i,j} u(i,j) y(i,j) \right) + \max_{y \in \mathcal{Y}} \left( g(y) - \sum_{i,j} u(i,j) z(i,j) \right)$$

$L(u)$ is convex, but not differentiable. A *subgradient* of $L(u)$ at $u$ is a vector $g_u$ such that for all $v$,

$$L(v) \geq L(u) + g_u \cdot (v - u)$$

Subgradient methods use updates $u' = u - \alpha g_u$

In fact, for our $L(u)$, $g_u(i,j) = z^*(i,j) - y^*(i,j)$

# Related Work

- Methods that use general purpose linear programming or integer linear programming solvers (Martins et al. 2009; Riedel and Clarke 2006; Roth and Yih 2005)

- Dual decomposition/Lagrangian relaxation in combinatorial optimization (Dantzig and Wolfe, 1960; Held and Karp, 1970; Fisher 1981)

- Dual decomposition for inference in MRFs (Komodakis et al., 2007; Wainwright et al., 2005)

- Methods that incorporate combinatorial solvers within loopy belief propagation (Duchi et al. 2007; Smith and Eisner 2008)