Probabilistic Graphical Models: Lagrangian Relaxation Algorithms for Natural Language Processing

Alexander M. Rush

(based on joint work with Michael Collins, Tommi Jaakkola, Terry Koo, David Sontag)

Uncertainty in language

natural language is notoriusly ambiguous, even in toy sentences

United flies some large jet

problem becomes very difficult in real-world sentences

Lin made the first start of his surprising N.B.A. career Monday just as the Knicks lost Amar'e Stoudemire and Carmelo Anthony, turning the roster upside down and setting the scene for another strange and wondrous evening at Madison Square Garden.

NLP challenges

Lin made the first start of his surprising N.B.A. career Monday just as the Knicks lost Amar'e Stoudemire and Carmelo Anthony, turning the roster upside down and setting the scene for another strange and wondrous evening at Madison Square Garden.

- who lost Stoudemire?
- what is turning the roster upside down?
- are "strange and wonderous" paired?
- what about "down and setting"?

NLP applications



New! Click the words above to view alternate translations. Dismiss

This lecture

- 1. introduce models for natural language processing in the context of decoding in graphical models
- 2. describe research in Lagrangian relaxation methods for inference in natural language problems.

Decoding in NLP

focus: decoding as a combinatorial optimization problem

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

where f is a scoring function and ${\mathcal Y}$ is a set of structures also known as the MAP problem

$$f(y;x) = p(y|x)$$

Algorithms for exact decoding

coming lectures will explore algorithms for decoding

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

for some problems, simple combinatorial algorithms are exact

• dynamic programming - e.g. tree-structured MRF's



- min cut for certain types of weight structures
- minimum spanning tree certain NLP problems

can depend on structure of the graphical model

Tagging



Parsing





Decoding complexity

issue: simple combinatorial algorithms do not scale to richer models

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

need decoding algorithms for complex natural language tasks

motivation:

- richer model structure often leads to improved accuracy
- exact decoding for complex models tends to be intractable

Phrase-based translation



Word alignment



Decoding tasks

high complexity

- combined parsing and part-of-speech tagging (Rush et al., 2010)
- "loopy" HMM part-of-speech tagging
- syntactic machine translation (Rush and Collins, 2011)

NP-Hard

- symmetric HMM alignment (DeNero and Macherey, 2011)
- phrase-based translation (Chang and Collins, 2011)
- higher-order non-projective dependency parsing (Koo et al., 2010)

in practice:

- approximate decoding methods (coarse-to-fine, beam search, cube pruning, gibbs sampling, belief propagation)
- approximate models (mean field, variational models)

Motivation

cannot hope to find exact algorithms (particularly when NP-Hard) **aim:** develop decoding algorithms with formal guarantees

method:

- derive fast algorithms that provide certificates of optimality
- show that for practical instances, these algorithms often yield exact solutions
- provide strategies for improving solutions or finding approximate solutions when no certificate is found

Lagrangian relaxation helps us develop algorithms of this form

Lagrangian relaxation

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
- ...

Lagrangian relaxation algorithms

Simple - uses basic combinatorial algorithms

Efficient - faster than solving exact decoding problems

Strong guarantees

- gives a certificate of optimality when exact
- direct connections to linear programming relaxations

Outline

- 1. background on exact NLP decoding
- 2. worked algorithm for combined parsing and tagging
- 3. important theorems and formal derivation
- 4. more examples from parsing and alignment

1. Background: Part-of-speech tagging



Graphical model formulation

given:

- a sentence of length n and a tag set ${\mathcal T}$
- one variable for each word, takes values in $\ensuremath{\mathcal{T}}$
- edge potentials heta(i-1,i,t',t) for all $i\in {\it n},\,t,t'\in {\cal T}$

example:



United₁ flies₂ some₃ large₄ jet_5

 $\mathcal{T} = \{A, D, N, V\}$

note: for probabilistic HMM $\theta(i-1, i, t', t) = \log(p(w_i|t)p(t|t'))$

Decoding problem

let $\mathcal{Y} = \mathcal{T}^n$ be the set of assignments to the variables decoding problem: $\arg \max_{y \in \mathcal{Y}} f(y)$ $f(y) = \sum_{i=2}^n \theta(i-1, i, y(i-1), y(i))$

example:

Ν

United₁ flies₂ some₃ large₄ jet_5

Combinatorial solution: Viterbi algorithm

Set $c(i, t) = -\infty$ for all $i \neq 1, t \in \mathcal{T}$ Set c(1, t) = 0 for all $t \in \mathcal{T}$ For i = 2 to n For t = 1 to \mathcal{T} $c(i,t) \leftarrow \max_{t' \in \mathcal{T}} c(i-1,t') + \theta(i-1,i,t',t)$ **Return** $\max_{t' \in \mathcal{T}} c(n, t')$

Note: run time $O(n|\mathcal{T}|^2)$, can greatly reduce $|\mathcal{T}|$ in practice special case of the belief propagation (next week)













Parsing



Parse decoding

parsing has potentials:

- $\theta(A \rightarrow B \ C)$ for each internal rules
- $\theta(A \rightarrow w)$ for pre-terminal (bottom) rules

not easily represented as a graphical model

• complicated independence structure



luckily, similar decoding algorithms to Viterbi

- CKY algorithm for context-free parsing $O(G^3n^3)$
- Eisner algorithm for dependency parsing $O(n^3)$

2. Worked example

aim: walk through a Lagrangian relaxation algorithm for combined parsing and part-of-speech tagging

- introduce formal notation for parsing and tagging
- give assumptions necessary for decoding
- step through a run of the Lagrangian relaxation algorithm

Combined parsing and part-of-speech tagging



goal: find parse tree that optimizes

$$score(S \rightarrow NP VP) + score(VP \rightarrow V NP) +$$

...+ $score(N \rightarrow V) + score(N \rightarrow United) + ...$

Constituency parsing

notation:

- ${\mathcal Y}$ is set of constituency parses for input
- $y \in \mathcal{Y}$ is a valid parse
- f(y) scores a parse tree

goal:

$$rg\max_{y\in\mathcal{Y}}f(y)$$

example: a context-free grammar for constituency parsing



Part-of-speech tagging

notation:

- ${\mathcal Z}$ is set of tag sequences for input
- $z \in \mathcal{Z}$ is a valid tag sequence
- g(z) scores of a tag sequence

goal:

 $\arg\max_{z\in\mathcal{Z}}g(z)$

example: an HMM for part-of speech tagging



Identifying tags

notation: identify the tag labels selected by each model

- y(i, t) = 1 when parse y selects tag t at position i
- z(i, t) = 1 when tag sequence z selects tag t at position i

example: a parse and tagging with y(4, A) = 1 and z(4, A) = 1



Combined optimization

goal:

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

such that for all $i = 1 \dots n$, $t \in \mathcal{T}$,

$$y(i,t)=z(i,t)$$

i.e. find the best parse and tagging pair that agree on tag labels equivalent formulation:

$$\arg\max_{y\in\mathcal{Y}}f(y)+g(l(y))$$

where $I: \mathcal{Y} \rightarrow \mathcal{Z}$ extracts the tag sequence from a parse tree

Exact method: Dynamic programming intersection

can solve by solving the product of the two models

example:

- parsing model is a context-free grammar
- tagging model is a first-order HMM
- can solve as CFG and finite-state automata intersection



some large jet
Intersected parsing and tagging complexity

let *G* be the number of grammar non-terminals

parsing CFG require $O(G^3n^3)$ time with rules VP \rightarrow V NP



with intersection $O(G^3n^3|\mathcal{T}|^3)$ with rules $VP_{N,V} \rightarrow V_{N,V} NP_{V,N}$

becomes $O(G^3n^3|\mathcal{T}|^6)$ time for second-order HMM

Tagging assumption

assumption: optimization with *u* can be solved efficiently

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

example: HMM with scores with edge potentials θ

$$g(z) = \sum_{i=2}^{n} \theta(i-1, i, z(i-1), z(i))$$

where

$$\begin{aligned} \arg \max_{z \in \mathcal{Z}} & g(z) - \sum_{i,t} u(i,t) z(i,t) = \\ \arg \max_{z \in \mathcal{Z}} & \sum_{i=2}^{n} (\theta(i-1,i,z(i-1),z(i)) - u(i,z(i))) - u(1,z(1)) \end{aligned}$$

(Modified) Viterbi algorithm

Set
$$c(i, t) = -\infty$$
 for all $i \neq 1, t \in \mathcal{T}$
Set $c(1, t) = u(1, t)$ for all $t \in \mathcal{T}$
For $i = 2$ to n
For $t = 1$ to \mathcal{T}
 $c(i, t) \leftarrow \max_{t' \in \mathcal{T}} c(i - 1, t') + \theta(i - 1, i, t', t) - u(i, t)$
Return $\max_{t' \in \mathcal{T}} c(n, t')$

Note: same algorithm as before with modified potentials

Parsing assumption

assumption: optimization with *u* can be solved efficiently

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

example: CFG with rule scoring function θ

$$f(y) = \sum_{X \to Y \ Z \in y} \theta(X \to Y \ Z) + \sum_{(i,X) \in y} \theta(X \to w_i)$$

where

$$\begin{aligned} \arg \max_{y \in \mathcal{Y}} \quad f(y) + \sum_{i,t} u(i,t)y(i,t) = \\ \arg \max_{y \in \mathcal{Y}} \quad \sum_{X \to Y \ Z \in y} \theta(X \to Y \ Z) + \sum_{(i,X) \in y} (\theta(X \to w_i) + u(i,X)) \end{aligned}$$

Lagrangian relaxation algorithm

Set
$$u^{(1)}(i, t) = 0$$
 for all $i, t \in \mathcal{T}$
For $k = 1$ to K
 $y^{(k)} \leftarrow \arg \max_{y \in \mathcal{Y}} f(y) + \sum_{i,t} u^{(k)}(i, t)y(i, t)$ [Parsing]
 $z^{(k)} \leftarrow \arg \max_{z \in \mathcal{Z}} g(z) - \sum_{i,t} u^{(k)}(i, t)z(i, t)$ [Tagging]
If $y^{(k)}(i, t) = z^{(k)}(i, t)$ for all i, t Return $(y^{(k)}, z^{(k)})$
Else $u^{(k+1)}(i, t) \leftarrow u^{(k)}(i, t) - \alpha_k(y^{(k)}(i, t) - z^{(k)}(i, t))$

CKY Parsing

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

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

Viterbi Decoding

United₁ flies₂ some₃ large₄ jet₅

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

Key



Viterbi Decoding

United₁ flies₂ some₃ large₄ jet₅

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

Key

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



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

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

Key



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

Key

Penalties u(i, t) = 0 for all 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

CKY Parsing

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

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

Viterbi Decoding

United₁ flies₂ some₃ large₄ jet₅

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

Key

CKY Parsing		S				Penal	ties
	NP		VP			u(i,t)=0	for all <i>i</i> ,t
	N	V		NP		Iteration	1
	United	l	D	Ă	N	u(1,A)	-1
	onned	mes	i somo	largo	 	u(1, N)	1
			Some	laige	jet	u(2, N)	-1
$y^* = a$	$\operatorname{rgmax}_{y\in\mathcal{Y}}(f)$	(y) + 2	y u(i, t)y(i,t))		u(2, V)	1
		'	i,t			u(5, V)	-1
Viterbi Decoding						u(5, N)	1

 $United_1 \ flies_2 \ some_3 \ large_4 \ jet_5$

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

Key



Penalties

u(i, t) = 0 for all i, t $\frac{\text{Iteration } 1}{u(1, A)} - 1$

1

1

Ν

1

iet

u(z, N)	-1
u(2, V)	1
u(5, V)	-1
u(5, N)	1

u(1, N)

..(2 M)

Viterbi Decoding



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

Key



y(i, t) = 1 if y contains tag t at position i

CKY Parsing

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

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

Viterbi Decoding

United₁ flies₂ some₃ large₄ jet₅

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

Key

CKY Parsing		S				Penalt	ies
	NP		VP			u(i,t)=0 for	or all <i>i</i> ,t
	I N	V		NP		Iteration 1	
	United	l	D	A	N	u(1,A)	-1
	onited	mes	i somo		 iot	u(1, N)	1
			some	large	jet	u(2, N)	-1
y * =	$= \arg \max_{y \in \mathcal{Y}} (f$	f(y) + 2) [u(i, t])y(i,t))		u(2, V)	1
			1,1			u(5, V)	-1
Viterbi Decodi	ng					u(5, N)	1
						Iteration 2	
I	United ₁ flies ₂	some ₃	$large_4$	jet ₅		u(5, V)	-1
						u(5, N)	1
<i>z</i> * =	$= rg\max_{z\in\mathcal{Z}}(g)$	$(z) - \sum_{i}$	$\sum_{t} u(i, t)$	z(i,t))			
Kev							

 $\begin{array}{rcl} y \\ f(y) & \leftarrow \ \mathsf{CFG} \\ \mathcal{Y} & \leftarrow \ \mathsf{Parse \ Trees} \\ y(i,t) = 1 & \mathsf{if} \\ \end{array} \begin{array}{rcl} y \ \mathsf{contains \ tag \ } t \ \mathsf{at \ position \ } i \end{array} \begin{array}{rcl} g(z) & \leftarrow \ \mathsf{HMM} \\ \mathcal{Z} & \leftarrow \ \mathsf{Taggings} \\ \end{array}$



y(i, t) = 1 if y contains tag t at position i

CKY Parsing	S					Penalt	ies
	NP	VP				u(i,t) = 0 for	or all <i>i</i> ,t
	N V	N	ĪP			Iteration 1	
	United flies	D	À N	J		u(1,A)	-1
		- I some		.+		u(1, N)	1
*						u(2, N)	-1
$y^{\star} =$	$\arg \max_{y \in \mathcal{Y}} (f(y) +$	$\sum_{i,j} u(i,t)y(i)$	(, t))			u(2, V)	1
		I,t				u(5, V)	-1
Viterbi Decoding	g					u(5, N)	1
	$N \longrightarrow V \longrightarrow D$	$\longrightarrow A \longrightarrow N$					
	\downarrow \downarrow \downarrow	↓ ↓				Iteration 2	
Ur	nited ₁ flies ₂ some	e ₃ large ₄ jet ₅				u(5, V)	-1
						u(5, N)	1
$z^* = z$	$\arg \max(g(z) -$	$\sum u(i,t)z(i$, t))				
	z∈Z	i,t				Conver	ged
Key					<i>y</i> *	$= \arg \max_{y \in \mathcal{Y}} f$	(y) + g(y)
f(y)	← CFG ← Parse Trees		g(z) 7	↓	HMM Taggings		
y(i,t) = 1	if y contains ta	g t at position	i	~	198811182		

Main theorem

theorem: if at any iteration, for all *i*, $t \in T$

$$y^{(k)}(i,t) = z^{(k)}(i,t)$$

then $(y^{(k)}, z^{(k)})$ is the global optimum

proof: focus of the next section

Convergence



2. Formal properties

 $\ensuremath{\operatorname{aim:}}$ formal derivation of the algorithm given in the previous section

- derive Lagrangian dual
- prove three properties
 - upper bound
 - convergence
 - optimality
- describe subgradient method

Lagrangian

goal:

$$rg\max_{y\in\mathcal{Y},z\in\mathcal{Z}}f(y)+g(z)$$
 such that $y(i,t)=z(i,t)$

Lagrangian:

$$L(u, y, z) = f(y) + g(z) + \sum_{i,t} u(i, t) (y(i, t) - z(i, t))$$

redistribute terms

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

Lagrangian dual

Lagrangian:

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

Lagrangian dual:

$$L(u) = \max_{y \in \mathcal{Y}, z \in \mathcal{Z}} L(u, y, z)$$

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

Theorem 1. Upper bound

define:

• y^*, z^* is the optimal combined parsing and tagging solution with $y^*(i, t) = z^*(i, t)$ for all i, t

theorem: for any value of u

$$L(u) \geq f(y^*) + g(z^*)$$

L(u) provides an upper bound on the score of the optimal solution **note:** upper bound may be useful as input to branch and bound or A* search

Theorem 1. Upper bound (proof)

theorem: for any value of u, $L(u) \ge f(y^*) + g(z^*)$ proof:

$$L(u) = \max_{y \in \mathcal{Y}, z \in \mathcal{Z}} L(u, y, z)$$
(1)

$$\geq \max_{y \in \mathcal{Y}, z \in \mathcal{Z}: y = z} L(u, y, z)$$
(2)

$$= \max_{y \in \mathcal{Y}, z \in \mathcal{Z}: y=z} f(y) + g(z)$$
(3)

$$= f(y^*) + g(z^*)$$
 (4)

Formal algorithm (reminder)

Set
$$u^{(1)}(i, t) = 0$$
 for all $i, t \in \mathcal{T}$
For $k = 1$ to K
 $y^{(k)} \leftarrow \arg \max_{y \in \mathcal{Y}} f(y) + \sum_{i,t} u^{(k)}(i, t)y(i, t)$ [Parsing]
 $z^{(k)} \leftarrow \arg \max_{z \in \mathcal{Z}} g(z) - \sum_{i,t} u^{(k)}(i, t)z(i, t)$ [Tagging]
If $y^{(k)}(i, t) = z^{(k)}(i, t)$ for all i, t Return $(y^{(k)}, z^{(k)})$
Else $u^{(k+1)}(i, t) \leftarrow u^{(k)}(i, t) - \alpha_k(y^{(k)}(i, t) - z^{(k)}(i, t))$

Theorem 2. Convergence

notation:

- $u^{(k+1)}(i,t) \leftarrow u^{(k)}(i,t) + lpha_k(y^{(k)}(i,t) z^{(k)}(i,t))$ is update
- $u^{(k)}$ is the penalty vector at iteration k
- $\alpha_k > 0$ is the update rate at iteration k

theorem: for any sequence $\alpha^1, \alpha^2, \alpha^3, \ldots$ such that

$$\lim_{t\to\infty}\alpha^t=0 \quad \text{and} \quad \sum_{t=1}^\infty \alpha^t=\infty,$$

we have

$$\lim_{t\to\infty}L(u^t)=\min_u L(u)$$

i.e. the algorithm converges to the tightest possible upper bound **proof:** by subgradient convergence (next section)

Dual solutions

define:

• for any value of *u*

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

and

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

• y_u and z_u are the dual solutions for a given u

Theorem 3. Optimality

theorem: if there exists u such that

$$y_u(i,t)=z_u(i,t)$$

for all i, t then

$$f(y_u) + g(z_u) = f(y^*) + g(z^*)$$

i.e. if the dual solutions agree, we have an optimal solution

 (y_u, z_u)

Theorem 3. Optimality (proof)

theorem: if u such that $y_u(i, t) = z_u(i, t)$ for all i, t then

$$f(y_u) + g(z_u) = f(y^*) + g(z^*)$$

proof: by the definitions of y_u and z_u

$$L(u) = f(y_u) + g(z_u) + \sum_{i,t} u(i,t)(y_u(i,t) - z_u(i,t))$$

= $f(y_u) + g(z_u)$

since $L(u) \ge f(y^*) + g(z^*)$ for all values of u

$$f(y_u) + g(z_u) \geq f(y^*) + g(z^*)$$

but y^* and z^* are optimal

$$f(y_u) + g(z_u) \leq f(y^*) + g(z^*)$$

Dual optimization

Lagrangian dual:

$$L(u) = \max_{y \in \mathcal{Y}, z \in \mathcal{Z}} L(u, y, z)$$

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

goal: dual problem is to find the tightest upper bound

 $\min_{u} L(u)$

Dual subgradient

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

properties:

- L(u) is convex in u (no local minima)
- L(u) is not differentiable (because of max operator)

handle non-differentiability by using subgradient descent

define: 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 algorithm

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

recall, y_u and z_u are the argmax's of the two terms subgradient:

$$g_u(i,t) = y_u(i,t) - z_u(i,t)$$

subgradient descent: move along the subgradient

$$u'(i,t) = u(i,t) - \alpha \left(y_u(i,t) - z_u(i,t) \right)$$

guaranteed to find a minimum with conditions given earlier for $\boldsymbol{\alpha}$

4. More examples

aim: demonstrate similar algorithms that can be applied to other decoding applications

- context-free parsing combined with dependency parsing
- combined translation alignment

Combined constituency and dependency parsing (Rush et al., 2010)

setup: assume separate models trained for constituency and dependency parsing

problem: find constituency parse that maximizes the sum of the two models

example:

· combine lexicalized CFG with second-order dependency parser

Lexicalized constituency parsing

notation:

- ${\mathcal Y}$ is set of lexicalized constituency parses for input
- $y \in \mathcal{Y}$ is a valid parse
- f(y) scores a parse tree

goal:

$$rg\max_{y\in\mathcal{Y}}f(y)$$

example: a lexicalized context-free grammar


Dependency parsing

define:

- $\ensuremath{\mathcal{Z}}$ is set of dependency parses for input
- $z \in \mathcal{Z}$ is a valid dependency parse
- g(z) scores a dependency parse

example:



Identifying dependencies

notation: identify the dependencies selected by each model

- y(i,j) = 1 when word *i* modifies of word *j* in constituency parse *y*
- z(i,j) = 1 when word *i* modifies of word *j* in dependency parse *z*

example: a constituency and dependency parse with y(3,5) = 1 and z(3,5) = 1



V

Combined optimization

goal:

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

such that for all $i = 1 \dots n$, $j = 0 \dots n$,

$$y(i,j) = z(i,j)$$

CKY Parsing

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

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

Dependency Parsing

 $*_0$ United₁ flies₂ some₃ large₄ jet₅

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

Key



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

Dependency Parsing

*0 United1 flies2 some3 large4 jet5

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

Key



Dependency Parsing



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

Key

Penalties u(i,j) = 0 for all 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



Penalties		
u(i,j) = 0 for	r all <i>i,j</i>	
Iteration 1		
<i>u</i> (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

CKY Parsing

Penalties

ι	u(i,j)=0	for	all	i,j
	Iteration	1		
	u(2,3)		-1	_
	u(5,3)		1	

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

Dependency Parsing

 $*_0$ United₁ flies₂ some₃ large₄ jet₅

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

Key



Penalties

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

$$\frac{\text{Iteration 1}}{u(2,3) - 1}$$

$$u(5,3) = 1$$

Dependency Parsing

 $*_0$ United₁ flies₂ some₃ large₄ jet₅

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

Key



u(i,j) = 0 for all i,j $\frac{\text{Iteration } 1}{u(2,3) - 1}$ 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



Dependency Parsing



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

Key

Penalties

u(i,j) = 0 for	r all <i>i</i> ,j
Iteration 1	
u(2,3)	-1
u(5,3)	1

Converged

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

Convergence



Integrated Constituency and Dependency Parsing: Accuracy



F_1 Score

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

Corpus-level tagging

setup: given a corpus of sentences and a trained sentence-level tagging model

problem: find best tagging for each sentence, while at the same time enforcing inter-sentence soft constraints

example:

- test-time decoding with a trigram tagger
- constraint that each word type prefer a single POS tag

Corpus-level tagging



Sentence-level decoding

notation:

- \mathcal{Y}_i is set of tag sequences for input sentence i
- $\mathcal{Y} = \mathcal{Y}_1 imes \ldots imes \mathcal{Y}_m$ is set of tag sequences for the input corpus
- $Y \in \mathcal{Y}$ is a valid tag sequence for the corpus
- $F(Y) = \sum_{i} f(Y_i)$ is the score for tagging the whole corpus

goal:

 $\arg\max_{Y\in\mathcal{Y}}F(Y)$

example: decode each sentence with a trigram tagger



Inter-sentence constraints

notation:

- ${\mathcal Z}$ is set of possible assignments of tags to word types
- $z \in \mathcal{Z}$ is a valid tag assignment
- g(z) is a scoring function for assignments to word types

example: an MRF model that encourages words of the same type to choose the same tag



 $g(z_1) > g(z_2)$

Identifying word tags

notation: identify the tag labels selected by each model

- Y_s(i, t) = 1 when the tagger for sentence s at position i selects tag t
- z(s, i, t) = 1 when the constraint assigns at sentence s position i the tag t

example: a parse and tagging with $Y_1(5, N) = 1$ and z(1, 5, N) = 1



Combined optimization

goal:

$$\arg \max_{Y \in \mathcal{Y}, z \in \mathcal{Z}} F(Y) + g(z)$$

such that for all $s = 1 \dots m$, $i = 1 \dots n$, $t \in \mathcal{T}$,

$$Y_s(i,t)=z(s,i,t)$$

Tagging

Penalties u(s, i, t) = 0 for all s, i, t

MRF

Key

Penalties u(s, i, t) = 0 for all s, i, t



MRF

Key

Penalties u(s, i, t) = 0 for all s, i, t



language

language language

Key

Penalties u(s, i, t) = 0 for all s, i, t



 \leftarrow Sentence-level tagging

 \mathcal{Y}

 $Y_s(i, t) = 1$

if

MRF g(z)4 Z Inter-sentence constraints \leftarrow sentence s has tag t at position i



 $Y_s(i,t) = 1$ if sentence s has tag t at position i

Inter-sentence constraints \leftarrow

Penalties

s, i, t) = 0 for	r all <i>s,i,t</i>
Iteration 1	
u(1, 5, N)	-1
u(1, 5, A)	1
u(3, 1, N)	-1
u(3, 1, A)	1

Penalties

u(s, i, t) = 0 for	all	s,i,t
Iteration 1		
u(1, 5, N)	-1	_
u(1, 5, A)	1	
u(3, 1, N)	-1	
u(3, 1, A)	1	

Tagging

MRF

Key

Penalties



u(s, i, t) = 0 for	r all s	s,i,t
Iteration 1		
u(1, 5, N)	-1	
u(1, 5, A)	1	
u(3, 1, N)	-1	
u(3, 1, A)	1	

MRF

Key

U (N)

language

(R)

now

N

beings

u(s, i, t) = 0 for all s, i, t $\underbrace{\frac{\text{Iteration } 1}{u(1, 5, N) - 1}}_{u(1, 5, A) - 1}$ u(3, 1, N) - 1 u(3, 1, A) - 1

Penalties

MRF

Tagging



(P)

my

A

language

Ρ

us

first

N)

arts

Ν

human

language l

V

is

(v)

studies

V

makes

(N)

English

P

He

N

Language

language language

Key

Tagging V (P) A (N) N English is my first language (v)P A N) (R)He studies language arts now Ν N V Ρ N makes beings Language us human Ν Ν Ν Ν

u(s, i, t) = 0 for all s, i, tIteration 1 u(1, 5, N)-1 u(1, 5, A)1 u(3, 1, N) -1

Penalties

u(3, 1, A)1

MRF



Key

F(Y)MRF Tagging model g(z)⇐ \Leftarrow \mathcal{Y} Sentence-level tagging 7. \Leftarrow Inter-sentence constraints \leftarrow $Y_{s}(i, t) = 1$ if sentence s has tag t at position i



Sentence-level tagging \Leftarrow $Y_{s}(i, t) = 1$ sentence s has tag t at position iif

 \mathcal{Y}

Inter-sentence constraints \leftarrow

Tagging u(s, i, t) = 0 for all s, i, t(v)P Iteration 1 (N) (A) N u(1, 5, N)-1 English is my first language u(1, 5, A)1 (v)P N) N) (R)u(3, 1, N) -1u(3, 1, A)1 He studies language arts now Iteration 2 (v)P N N N u(1, 5, N)-1 makes beings Language human us u(1, 5, A)1 **MRF** u(3, 1, N)-1 u(3, 1, A)1 u(2, 3, N)1 u(2, 3, A)-1

Penalties

- Key

Penalties Tagging u(s, i, t) = 0 for all s, i, tV Iteration 1 (Р) Â (N) Ν u(1, 5, N)-1 English is mv first language u(1, 5, A)1 (v)P N N) (R)u(3, 1, N) -1u(3, 1, A)1 He studies language arts now Iteration 2 ٧. Ν N Ρ Ν u(1, 5, N)-1 makes beings Language human us u(1, 5, A)1 **MRF** u(3, 1, N)-1 Ν u(3, 1, A)1 u(2, 3, N)1 Ν Ν Ν u(2, 3, A)-1 language language language Key F(Y)MRF Tagging model g(z) \Leftarrow \leq \mathcal{Y} Sentence-level tagging 7. ⇐ Inter-sentence constraints

 \leftarrow

 $Y_s(i,t) = 1$ sentence s has tag t at position iif

Combined alignment (DeNero and Macherey, 2011)

setup: assume separate models trained for English-to-French and French-to-English alignment

problem: find an alignment that maximizes the score of both models

example:

• HMM models for both directional alignments (assume correct alignment is one-to-one for simplicity)

Alignment



Graphical model formulation

given:

- French sentence of length n and English sentence of length m
- one variable for each English word, takes values in $\{1, \ldots, n\}$
- edge potentials heta(i-1,i,f',f) for all $i\in {\it n},\,f,f'\in\{1,\ldots,{\it n}\}$

example:



English-to-French alignment

define:

- $\mathcal Y$ is set of all possible English-to-French alignments
- $y \in \mathcal{Y}$ is a valid alignment
- f(y) scores of the alignment

example: HMM alignment


French-to-English alignment

define:

- \mathcal{Z} is set of all possible French-to-English alignments
- $z\in\mathcal{Z}$ is a valid alignment
- g(z) scores of an alignment

example: HMM alignment



Identifying word alignments

notation: identify the tag labels selected by each model

- y(i,j) = 1 when e-to-f alignment y selects French word i to align with English word j
- z(i, j) = 1 when f-to-e alignment z selects French word i to align with English word j

example: two HMM alignment models with y(6,5) = 1 and z(6,5) = 1



Combined optimization

goal:

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

such that for all $i = 1 \dots n$, $j = 1 \dots n$,

$$y(i,j) = z(i,j)$$

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

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

French-to-English

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

Key



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

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

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

Key



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

Key

f(y)	\Leftarrow	HMM Alignment	g(z)	\Leftarrow	HMM Alignment
\mathcal{Y}^{-}	\Leftarrow	English-to-French model	\mathcal{Z}	\Leftarrow	French-to-English model
y(i, j) = 1	if	French word <i>i</i> aligns to English word <i>i</i>			



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



u(i,j) = 0 for all i,jIteration 1 u(3,2)-1 u(2,2)1 u(2,3) -1 u(3,3) 1

⇐ HMM Alignment f(y) \mathcal{Y} $g(z) \leftarrow \text{HMM Alignment}$ ← English-to-French model \mathcal{Z} \leftarrow French-to-English model y(i, j) = 1 if French word *i* aligns to English word *j*

Penalties

u(i,j)=0	for all <i>i</i> ,j
Iteration	1
u(3,2)	-1
u(2,2)	1
u(2,3)	-1
u(3,3)	1

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

French-to-English

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

Key



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

o-English

French-to-English

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

Key



u(i, j	i) = 0	for	all	i,j
lte	ration	1		
u(:	3, 2)		-1	_
u(:	2,2)		1	
u(1	2,3)		-1	
u(:	3, 3)		1	

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

laid a fourrure rouge

French-to-English



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

Key



u(i,j)=0	for all <i>i,j</i>
Iteration	1
u(3,2)	-1
u(2,2)	1
u(2,3)	-1
u(3,3)	1

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

laid a fourrure rouge

French-to-English



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

Key



u(i,j)=0	for all <i>i,j</i>
Iteration	1
u(3,2)	-1
u(2,2)	1
u(2,3)	-1
u(3,3)	1

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

laid a fourrure rouge

French-to-English



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

Key

MAP problem in Markov random fields



given: binary variables $x_1 \dots x_n$

goal: MAP problem

$$\arg\max_{x_1...x_n}\sum_{(i,j)\in E}f_{i,j}(x_i,x_j)$$

where each $f_{i,j}(x_i, x_j)$ is a local potential for variables x_i, x_j

Dual decomposition for MRFs (Komodakis et al., 2010)

goal:

$$\arg\max_{x_1...x_n}\sum_{(i,j)\in E}f_{i,j}(x_i,x_j)$$

equivalent formulation:

$$\arg\max_{x_1\ldots x_n, y_1\ldots y_n}\sum_{(i,j)\in\mathcal{T}_1}f_{i,j}'(x_i,x_j)+\sum_{(i,j)\in\mathcal{T}_2}f_{i,j}'(y_i,y_j)$$

such that for $i = 1 \dots n$,

$$x_i = y_i$$

Lagrangian:

$$L(u, x, y) = \sum_{(i,j)\in T_1} f'_{i,j}(x_i, x_j) + \sum_{(i,j)\in T_2} f'_{i,j}(y_i, y_j) + \sum_i u_i(x_i - y_i)$$

4. Practical issues

aim: overview of practical dual decomposition techniques

- tracking the progress of the algorithm
- choice of update rate α_k
- lazy update of dual solutions
- extracting solutions if algorithm does not converge

Optimization tracking

at each stage of the algorithm there are several useful values

track:

- $y^{(k)}$, $z^{(k)}$ are current dual solutions
- $L(u^{(k)})$ is the current dual value
- $y^{(k)}$, $I(y^{(k)})$ is a potential primal feasible solution
- $f(y^{(k)}) + g(l(y^{(k)}))$ is the potential primal value

Tracking example



example run from syntactic machine translation (later in talk)

• current primal

$$f(y^{(k)}) + g(l(y^{(k)}))$$

current dual

 $L(u^{(k)})$

Optimization progress

useful signals:

- $L(u^{(k)}) L(u^{(k-1)})$ is the dual change (may be positive)
- $\min_{k} L(u^{(k)})$ is the best dual value (tightest upper bound)
- $\max_{k} f(y^{(k)}) + g(l(y^{(k)}))$ is the best primal value

the optimal value must be between the best dual and primal values

Progress example



Update rate

choice of α_k has important practical consequences

- α_k too high causes dual value to fluctuate
- α_k too low means slow progress



Update rate

practical: find a rate that is robust to varying inputs

- α_k = c (constant rate) can be very fast, but hard to find constant that works for all problems
- $\alpha_k = \frac{c}{k}$ (decreasing rate) often cuts rate too aggressively, lowers value even when making progress
- rate based on dual progress
 - $\alpha_k = \frac{c}{t+1}$ where t < k is number of iterations where dual value increased
 - robust in practice, reduces rate when dual value is fluctuating

Lazy decoding

idea: don't recompute $y^{(k)}$ or $z^{(k)}$ from scratch each iteration

lazy decoding: if subgradient $u^{(k)}$ is sparse, then $y^{(k)}$ may be very easy to compute from $y^{(k-1)}$

use:

- helpful if y or z factor naturally into independent components
- can be important for fast decompositions

Lazy decoding example





recall corpus-level tagging example

at this iteration, only sentence 2 receives a weight update

with lazy decoding

$$Y_1^{(k)} \leftarrow Y_1^{(k-1)} \\ Y_3^{(k)} \leftarrow Y_3^{(k-1)}$$

Lazy decoding results

lazy decoding is critical for the efficiency of some applications



recomputation statistics for non-projective dependency parsing

Approximate solution

upon agreement the solution is exact, but this may not occur otherwise, there is an easy way to find an approximate solution **choose:** the structure $y^{(k')}$ where

$$k' = \arg\max_k f(y^{(k)}) + g(I(y^{(k)}))$$

is the iteration with the best primal score

guarantee: the solution $y^{k'}$ is non-optimal by at most

$$(\min_{k} L(u^{k})) - (f(y^{(k')}) + g(l(y^{(k')})))$$

there are other methods to estimate solutions, for instance by averaging solutions (see ?)

Choosing best solution



non-exact example from syntactic translation

best approximate primal solution occurs at iteration 63

Early stopping results

early stopping results for constituency and dependency parsing



Early stopping results

early stopping results for non-projective dependency parsing



define:

- ► source-language sentence words $x_1, ..., x_N$
- phrase translation p = (s, e, t)
- translation derivation $y = p_1, \ldots, p_L$

example:

x_1	<i>x</i> ₂	<i>x</i> 3	<i>x</i> ₄	<i>x</i> 5	<i>x</i> 6
das	muss	unsere	sorge	gleichermaßen	sein

define:

- source-language sentence words x_1, \ldots, x_N
- phrase translation p = (s, e, t)
- translation derivation $y = p_1, \ldots, p_L$

example:



 $y = \{(1, 2, \text{this must}),$

define:

- source-language sentence words x_1, \ldots, x_N
- phrase translation p = (s, e, t)
- translation derivation $y = p_1, \ldots, p_L$

example:



 $y = \{(1, 2, \text{this must}), (5, 5, \text{also}), \}$

define:

- source-language sentence words x_1, \ldots, x_N
- phrase translation p = (s, e, t)
- translation derivation $y = p_1, \ldots, p_L$

example:



 $y = \{(1, 2, \text{this must}), (5, 5, \text{also}), (6, 6, \text{be}), \}$

define:

- source-language sentence words x_1, \ldots, x_N
- phrase translation p = (s, e, t)
- translation derivation $y = p_1, \ldots, p_L$

example:



 $y = \{(1, 2, \text{this must}), (5, 5, \text{also}), (6, 6, \text{be}), (3, 4, \text{our concern})\}$

define:

- source-language sentence words x_1, \ldots, x_N
- phrase translation p = (s, e, t)
- translation derivation $y = p_1, \ldots, p_L$

example:



 $y = \{(1, 2, \text{this must}), (5, 5, \text{also}), (6, 6, \text{be}), (3, 4, \text{our concern})\}$

define:

- source-language sentence words x_1, \ldots, x_N
- phrase translation p = (s, e, t)
- translation derivation $y = p_1, \ldots, p_L$

example:



 $y = \{(1, 2, \text{this must}), (5, 5, \text{also}), (6, 6, \text{be}), (3, 4, \text{our concern})\}$
derivation:

 $y = \{(1, 2, \text{this must}), (5, 5, \text{also}), (6, 6, \text{be}), (3, 4, \text{our concern})\}$



$$f(y) = h(e(y)) + \sum_{k=1}^{L} g(p_k) + \sum_{k=1}^{L-1} \eta |t(p_k) + 1 - s(p_{k+1})|$$

- language model score h
- phrase translation score g
- distortion penalty η

derivation:

 $y = \{(1, 2, \text{this must}), (5, 5, \text{also}), (6, 6, \text{be}), (3, 4, \text{our concern})\}$



$$f(y) = h(e(y)) + \sum_{k=1}^{L} g(p_k) + \sum_{k=1}^{L-1} \eta |t(p_k) + 1 - s(p_{k+1})|$$

- ► language model score *h*
- phrase translation score g
- distortion penalty η

derivation:

 $y = \{(1, 2, \text{this must}), (5, 5, \text{also}), (6, 6, \text{be}), (3, 4, \text{our concern})\}$



$$f(y) = h(e(y)) + \sum_{k=1}^{L} g(p_k) + \sum_{k=1}^{L-1} \eta |t(p_k) + 1 - s(p_{k+1})|$$

- language model score h
- phrase translation score g
- distortion penalty η

derivation:

 $y = \{(1, 2, \text{this must}), (5, 5, 2 \text{ lso}), (6, 6, \text{be}), (3, 4, \text{our concern})\}$



$$f(y) = h(e(y)) + \sum_{k=1}^{L} g(p_k) + \sum_{k=1}^{L-1} \eta |t(p_k) + 1 - s(p_{k+1})|$$

- language model score h
- phrase translation score g
- distortion penalty η

 \mathcal{Y}' : only requires the total number of words translated to be N

 $\mathcal{Y}' = \{y : \sum_{i=1}^{N} y(i) = N \text{ and the distortion limit } d \text{ is satisfied}\}$

example:



 \mathcal{Y}' : only requires the total number of words translated to be N

 $\mathcal{Y}' = \{y : \sum_{i=1}^{N} y(i) = N \text{ and the distortion limit } d \text{ is satisfied}\}$

example:



 \mathcal{Y}' : only requires the total number of words translated to be N

 $\mathcal{Y}' = \{y : \sum_{i=1}^{N} y(i) = N \text{ and the distortion limit } d \text{ is satisfied}\}$

example:



 \mathcal{Y}' : only requires the total number of words translated to be N

 $\mathcal{Y}' = \{y : \sum_{i=1}^{N} y(i) = N \text{ and the distortion limit } d \text{ is satisfied}\}$

example:





$$\mathcal{Y} = \{ y : y(i) = 1 \forall i = 1 \dots N \} \qquad \qquad \boxed{1 \ 1 \dots 1}$$

rewrite:





$$\mathcal{Y}' = \{ y : \sum_{i=1}^{N} y(i) = N \}$$



$$\mathcal{Y}' = \{ y : \sum_{i=1}^{N} y(i) = N \}$$

original: $\arg \max f(y)$ $y \in \mathcal{Y}$ exact DP is NP-hard $\mathcal{Y} = \{ y : y(i) = 1 \ \forall i = 1 \dots N \}$ $1 || 1 | \dots | 1 |$ rewrite: $rg\max_{y\in\mathcal{Y}'}f(y)$ such that $y(i) = 1 \; \forall i = 1 \dots N$ can be solved efficiently by DP 1 0 $\mathcal{Y}' = \{ \boldsymbol{y} : \sum_{i=1}^{N} \boldsymbol{y}(i) = \boldsymbol{N} \}$ sum to N



Iteration 1:

▶ update
$$u(i)$$
: $u(i) \leftarrow u(i) - \alpha(y(i) - 1)$
 $\alpha = 1$
 $u(i) \ 0 \ 0 \ 0 \ 0 \ 0$
 $y(i)$
 $x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6$
das muss unsere sorge gleichermaßen sein

Iteration 1:

▶ update
$$u(i)$$
: $u(i) \leftarrow u(i) - \alpha(y(i) - 1)$
 $\alpha = 1$



Iteration 1:

▶ update u(i): $u(i) \leftarrow u(i) - \alpha(y(i) - 1)$ $\alpha = 1$



Iteration 2:

▶ update
$$u(i)$$
: $u(i) \leftarrow u(i) - \alpha(y(i) - 1)$
 $\alpha = 0.5$
 $u(i) \ 1 \quad 0 \quad -1 \quad -1 \quad 1 \quad 0$
 $y(i)$
 $x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6$
das muss unsere sorge gleichermaßen sein

Iteration 2:

▶ update
$$u(i)$$
: $u(i) \leftarrow u(i) - \alpha(y(i) - 1)$
 $\alpha = 0.5$



Iteration 2:

► update u(i): $u(i) \leftarrow u(i) - \alpha(y(i) - 1)$ $\alpha = 0.5$



Iteration 3:

▶ update
$$u(i)$$
: $u(i) \leftarrow u(i) - \alpha(y(i) - 1)$
 $\alpha = 0.5$
 $u(i) \ 1 \ -0.5 \ -0.5 \ -0.5 \ 0.5 \ 0$
 $y(i)$
 $x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6$
das muss unsere sorge gleichermaßen sein

Iteration 3:

▶ update
$$u(i)$$
: $u(i) \leftarrow u(i) - \alpha(y(i) - 1)$
 $\alpha = 0.5$



In some cases, we never reach y(i) = 1 for $i = 1 \dots N$

If dual L(u) is not decreasing fast enough run for 10 more iterations count number of times each constraint is violated add 3 most often violated constraints

Iteration 41:



Iteration 42:



Iteration 43:



Iteration 44:



Iteration 50:



Iteration 51:



Add 2 hard constraints (x_5, x_6) to the dynamic program

Iteration 51:



Add 2 hard constraints (x_5, x_6) to the dynamic program

Experiments: German to English

- Europarl data: German to English
- Test on 1,824 sentences with length 1-50 words
- Converged: 1,818 sentences (99.67%)

Experiments: Number of Iterations



Experiments: Number of Hard Constraints Required



Experiments: Mean Time in Seconds

# words	1-10	11-20	21-30	31-40	41-50	All	
mean	0.8	10.9	57.2	203.4	679.9	120.9	
median	0.7	8.9	48.3	169.7	484.0	35.2	

Comparison to ILP Decoding

	(sec.)	(sec.)
1-10	275.2	132.9
11-15	2,707.8	1,138.5
16-20	20,583.1	3,692.6

Higher-order non-projective dependency parsing

setup: given a model for higher-order non-projective dependency parsing (sibling features)

problem: find non-projective dependency parse that maximizes the score of this model

difficulty:

- model is NP-hard to decode
- complexity of the model comes from enforcing combinatorial constraints

strategy: design a decomposition that separates combinatorial constraints from direct implementation of the scoring function

Non-Projective Dependency Parsing



Important problem in many languages.

Problem is NP-Hard for all but the simplest models.
Dual Decomposition

A classical 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



Non-Projective Dependency Parsing



- Starts at the root symbol *
- Each word has a exactly one parent word
- Produces a tree structure (no cycles)
- Dependencies can cross



f(y) =



 $f(y) = score(head = *_0, mod = saw_2)$



 $f(y) = score(head = *_0, mod = saw_2) + score(saw_2, John_1)$



 $f(y) = score(head = *_0, mod = saw_2) + score(saw_2, John_1) + score(saw_2, movie_4)$



 $f(y) = score(head = *_0, mod = saw_2) + score(saw_2, John_1)$

 $+score(saw_2, movie_4) + score(saw_2, today_5)$



 $f(y) = score(head = *_0, mod = saw_2) + score(saw_2, John_1)$ $+ score(saw_2, movie_4) + score(saw_2, today_5)$ $+ score(movie_4, a_3) + \dots$



e.g. $score(*_0, saw_2) = \log p(saw_2|*_0)$ (generative model)



 $f(y) = score(head =*_0, mod =saw_2) + score(saw_2, John_1)$ $+ score(saw_2, movie_4) + score(saw_2, today_5)$ $+ score(movie_4, a_3) + \dots$

e.g. $score(*_0, saw_2) = \log p(saw_2|*_0)$ (generative model) or $score(*_0, saw_2) = w \cdot \phi(saw_2, *_0)$ (CRF/perceptron model)



 $f(y) = score(head =*_0, mod =saw_2) + score(saw_2, John_1)$ $+ score(saw_2, movie_4) + score(saw_2, today_5)$ $+ score(movie_4, a_3) + \dots$

e.g. $score(*_0, saw_2) = \log p(saw_2|*_0)$ (generative model) or $score(*_0, saw_2) = w \cdot \phi(saw_2, *_0)$ (CRF/perceptron model)

 $y^* = \arg \max_y f(y) \Leftarrow \text{Minimum Spanning Tree Algorithm}$



f(y) =



 $f(y) = score(head = *_0, prev = NULL, mod = saw_2)$



 $f(y) = score(head = *_0, prev = NULL, mod = saw_2)$

+*score*(saw₂, NULL, John₁)



 $f(y) = score(head = *_0, prev = NULL, mod = saw_2)$

 $+score(saw_2, NULL, John_1) + score(saw_2, NULL, movie_4)$



 $f(y) = score(head = *_0, prev = \text{NULL}, mod = \text{saw}_2)$ +score(saw_2, NULL, John_1) +score(saw_2, NULL, movie_4) +score(saw_2, movie_4, today_5) + ...



 $\begin{aligned} f(y) &= \textit{score}(\textit{head} = *_0, \textit{prev} = \text{NULL}, \textit{mod} = \text{saw}_2) \\ &+ \textit{score}(\text{saw}_2, \text{NULL}, \text{John}_1) + \textit{score}(\text{saw}_2, \text{NULL}, \text{movie}_4) \\ &+ \textit{score}(\text{saw}_2, \text{movie}_4, \textbf{today}_5) + \dots \end{aligned}$

e.g. $score(saw_2, movie_4, today_5) = \log p(today_5|saw_2, movie_4)$



 $f(y) = score(head = *_0, prev = \text{NULL}, mod = \text{saw}_2)$ +score(saw_2, NULL, John_1) +score(saw_2, NULL, movie_4) +score(saw_2, movie_4, today_5) + ...

e.g. $score(saw_2, movie_4, today_5) = \log p(today_5|saw_2, movie_4)$ or $score(saw_2, movie_4, today_5) = w \cdot \phi(saw_2, movie_4, today_5)$



 $f(y) = score(head = *_0, prev = \text{NULL}, mod = \text{saw}_2)$ +score(saw_2, NULL, John_1) +score(saw_2, NULL, movie_4) +score(saw_2, movie_4, today_5) + ...

e.g. $score(saw_2, movie_4, today_5) = \log p(today_5|saw_2, movie_4)$ or $score(saw_2, movie_4, today_5) = w \cdot \phi(saw_2, movie_4, today_5)$

$$y^* = \arg \max_y f(y) \Leftarrow \mathsf{NP-Hard}$$

Thought Experiment: Individual Decoding

 a_0 John₁ saw₂ a_3 movie₄ today₅ that₆ he₇ liked₈

Thought Experiment: Individual Decoding



 $score(saw_2, NULL, John_1) + score(saw_2, NULL, movie_4) + score(saw_2, movie_4, today_5)$

Thought Experiment: Individual Decoding



 $score(saw_2, NULL, John_1) + score(saw_2, NULL, movie_4) + score(saw_2, movie_4, today_5)$

 $score(saw_2, NULL, John_1) + score(saw_2, NULL, that_6)$



 $score(saw_2, NULL, John_1) + score(saw_2, NULL, movie_4) + score(saw_2, movie_4, today_5)$

 $score(saw_2, NULL, John_1) + score(saw_2, NULL, that_6)$

 $score(saw_2, NULL, a_3) + score(saw_2, a_3, he_7)$





Under Sibling Model, can solve for each word with Viterbi decoding.



Idea: Do individual decoding for each head word using dynamic programming.



Idea: Do individual decoding for each head word using dynamic programming.



Idea: Do individual decoding for each head word using dynamic programming.



Idea: Do individual decoding for each head word using dynamic programming.



Idea: Do individual decoding for each head word using dynamic programming.



Idea: Do individual decoding for each head word using dynamic programming.



Idea: Do individual decoding for each head word using dynamic programming.



Idea: Do individual decoding for each head word using dynamic programming.



Idea: Do individual decoding for each head word using dynamic programming.

If we're lucky, we'll end up with a valid final tree.

But we might violate some constraints.

Dual Decomposition Structure

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

Dual Decomposition Structure

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

Rewrite as
$$\underset{z \in \mathcal{Z}, y \in \mathcal{Y}}{\operatorname{argmax}} f(z) + g(y)$$

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

Rewrite as argmax
$$f(z) + g(y)$$

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

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









Set penalty weights equal to 0 for all edges.

For k = 1 to K

Set penalty weights equal to 0 for all edges.

For k = 1 to K

 $z^{(k)} \leftarrow \mathsf{Decode} \ (f(z) + \mathrm{penalty})$ by Individual Decoding

Set penalty weights equal to 0 for all edges.

For k = 1 to K $z^{(k)} \leftarrow \text{Decode } (f(z) + \text{penalty})$ by Individual Decoding $y^{(k)} \leftarrow \text{Decode } (g(y) - \text{penalty})$ by Minimum Spanning Tree

Set penalty weights equal to 0 for all edges.

For k = 1 to K $z^{(k)} \leftarrow \text{Decode } (f(z) + \text{penalty})$ by Individual Decoding $y^{(k)} \leftarrow \text{Decode } (g(y) - \text{penalty})$ by Minimum Spanning Tree If $y^{(k)}(i,j) = z^{(k)}(i,j)$ for all i,j Return $(y^{(k)}, z^{(k)})$

Set penalty weights equal to 0 for all edges.

For k = 1 to K $z^{(k)} \leftarrow \text{Decode } (f(z) + \text{penalty})$ by Individual Decoding $y^{(k)} \leftarrow \text{Decode } (g(y) - \text{penalty})$ by Minimum Spanning Tree If $y^{(k)}(i,j) = z^{(k)}(i,j)$ for all i,j Return $(y^{(k)}, z^{(k)})$

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

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

 $*_0$ John $_1$ saw $_2$ a $_3$ movie $_4$ today $_5$ that $_6$ he $_7$ liked $_8$ $z^* = rg\max_{z\in\mathcal{Z}}(f(z)+\sum_{i,j}u(i,j)z(i,j))$

Minimum Spanning Tree

Kev

 $*_0$ John₁ saw₂ a₃ movie₄ today₅ that₆ he₇ liked₈

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

$$f(z) \quad \Leftarrow \quad \text{Sibling Model} \qquad g(y) \quad \Leftarrow \quad \text{Arc-Factored Model}$$

$$\mathcal{Z} \quad \Leftarrow \quad \text{No Constraints} \qquad \mathcal{Y} \quad \Leftarrow \quad \text{Tree Constraints}$$

y(i,j) = 1 if y contains dependency i,j



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

Minimum Spanning Tree

Key

 $*_0$ John₁ saw₂ a₃ movie₄ today₅ that₆ he₇ liked₈

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

$$f(z) \quad \Leftarrow \quad \text{Sibling Model} \qquad g(y) \quad \Leftarrow \quad \text{Arc-Factored Model}$$

$$Z \quad \leftarrow \quad \text{No Constraints} \qquad \qquad \mathcal{Y} \quad \Leftarrow \quad \text{Tree Constraints}$$

y(i,j) = 1 if y contains dependency i, j

Individual Decoding u(i, j) = 0 for all i, j*0 today₅ that₆ liked₈ John₁ movie₄ he₇ saw₂ a₃ $z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum_{i,i} u(i,j)z(i,j))$ Minimum Spanning Tree *0 .John₁ movie₄ today₅ that₆ he₇ liked₈ sawo as $y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,i} u(i,j)y(i,j))$ Key g(y)⇐ Arc-Factored Model

Penalties

⇐ Sibling Model⇐ No Constraints f(z)Z ν ⇐ Tree Constraints y(i, j) = 1 if y contains dependency i, j





Minimum Spanning Tree



$$\mathcal{Z} \qquad \leftarrow \qquad \text{No Constraints} \qquad \qquad \mathcal{Y} \qquad \leftarrow \qquad \text{Tree Constraints} \\ \mathcal{Y}(i,j) = 1 \quad \text{if} \qquad y \text{ contains dependency } i,j$$



Minimum Spanning Tree

Kev

 $*_0$ John₁ saw₂ a₃ movie₄ today₅ that₆ he₇ liked₈

$$\begin{split} y^* &= \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j)) \\ f(z) &\Leftarrow \text{Sibling Model} \\ \mathcal{Z} &\Leftarrow \text{No Constraints} \\ \mathcal{Y} &\Leftarrow \text{Tree Constraints} \end{split}$$

y(i,j) = 1 if y contains dependency i, j



-1

-1

1

1

Minimum Spanning Tree



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

Key

f(z) \mathcal{Z} Sibling Model \Leftarrow g(y)Arc-Factored Model \Leftarrow No Constraints \mathcal{Y} Tree Constraints ⇐ y(i,j) = 1y contains dependency i, jif



Individual Decoding

Penalties

i,j

	u(i,j) = 0 for	or all
	Iteration 1	
	u(8,1)	-1
· · · · · · · · · · · · · · · · · · ·	<i>u</i> (4,6)	-1
\circ_0 John $_1$ saw $_2$ a $_3$ movie $_4$ today $_5$ that $_6$ he $_7$ liked $_8$	u(2,6)	1
$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum u(i,j)z(i,j))$	u(8,7)	1
i,j	Iteration 2	
Minimum Spanning Tree	u(8,1)	-1
	<i>u</i> (4,6)	-2
	<i>u</i> (2,6)	2
	u(8,7)	1
${}^{*}_{0}$ John $_{1}$ saw $_{2}$ a $_{3}$ movie $_{4}$ today $_{5}$ that $_{6}$ he $_{7}$ liked $_{8}$		
$y^* = rg\max_{y \in \mathcal{Y}}(g(y) - \sum_{i,j} u(i,j)y(i,j))$		
Key $f(z) \leftarrow \text{Sibling Model} g(y) \leftarrow \text{Arc-Fact}$	ored Model	

f(z) \mathcal{Z} \Leftarrow No Constraints \mathcal{Y} ← Tree Constraints

y(i,j) = 1 if y contains dependency i, j

Individual Decoding	Penalties
	u(i,j) = 0 for all i,j
	Iteration 1
	u(8,1) -1
	<i>u</i> (4,6) -1
\circ_0 John $_1$ saw $_2$ a $_3$ movie $_4$ today $_5$ that $_6$ he $_7$ liked $_8$	u(2,6) 1
$z^* = \arg \max_{z \in \mathcal{Z}} (f(z) + \sum u(i,j)z(i,j))$	<i>u</i> (8,7) 1
i,j	Iteration 2
Minimum Spanning Tree	u(8,1) -1
	<i>u</i> (4,6) -2
	u(2,6) 2
	u(8,7) 1
$*_0$ John ₁ saw ₂ a ₃ movie ₄ today ₅ that ₆ he ₇ liked ₈	
$y^* = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$	
Key $f(z) \iff$ Sibling Model $g(y) \iff$ Arc-Fa $\mathcal{Z} \iff$ No Constraints $\mathcal{Y} \iff$ Tree C $y(i,j) = 1$ if y contains dependency i,j	actored Model Constraints



Individual Decoding	Penalties
	u(i,j) = 0 for all i,j
	Iteration 1
	u(8,1) -1
$*_0$ John ₁ saw ₂ a ₃ movie ₄ today ₅ that ₆ he ₇ liked ₈	u(4,6) -1
	u(2,6) 1
$z^* = rg\max_{z \in \mathcal{Z}} (f(z) + \sum u(i,j)z(i,j))$	u(8,7) 1
i,j	Iteration 2
Minimum Spanning Tree	u(8,1) -1
	u(4,6) -2
	u(2,6) 2
*o lohn, sawo az movier todaya thata her likedo	u(8,7) 1
*	Converged
$y^{\star} = \arg \max_{y \in \mathcal{Y}} (g(y) - \sum_{i,j} u(i,j)y(i,j))$	$y^* = rg\max_{y\in\mathcal{Y}} f(y) + g(y)$
Key $f(z) \notin$ Sibling Model $g(y) \notin$ Arc-Factore $\mathcal{Z} \notin$ No Constraints $\mathcal{Y} \notin$ Tree Const $y(i,j) = 1$ if y contains dependency i,j	ed Model raints

Guarantees

Theorem

If at any iteration $y^{(k)} = z^{(k)}$, then $(y^{(k)}, z^{(k)})$ is the global optimum.

In experiments, we find the global optimum on 98% of examples.

Guarantees

Theorem

If at any iteration $y^{(k)} = z^{(k)}$, then $(y^{(k)}, z^{(k)})$ is the global optimum.

In experiments, we find the global optimum on 98% of examples.

If we do not converge to a match, we can still return an approximate solution (more in the paper).

Extensions



Head Automata (Eisner, 2000)

Generalization of Sibling models

Allow arbitrary automata as local scoring function.

Experiments

Properties:

- Exactness
- Parsing Speed
- Parsing Accuracy
- Comparison to Individual Decoding
- ► Comparison to LP/ILP

Training:

Averaged Perceptron (more details in paper)

Experiments on:

- CoNLL Datasets
- English Penn Treebank
- Czech Dependency Treebank

How often do we exactly solve the problem?



 Percentage of examples where the dual decomposition finds an exact solution.

Parsing Speed



- Number of sentences parsed per second
- Comparable to dynamic programming for projective parsing

Accuracy

	Arc-Factored	Prev Best	Grandparent
Dan	89.7	91.5	91.8
Dut	82.3	85.6	85.8
Por	90.7	92.1	93.0
Slo	82.4	85.6	86.2
Swe	88.9	90.6	91.4
Tur	75.7	76.4	77.6
Eng	90.1		92.5
Cze	84.4		87.3

Prev Best - Best reported results for CoNLL-X data set, includes

- Approximate search (McDonald and Pereira, 2006)
- Loop belief propagation (Smith and Eisner, 2008)
- (Integer) Linear Programming (Martins et al., 2009)

Comparison to Subproblems



 F_1 for dependency accuracy

Comparison to $\ensuremath{\mathsf{LP}}\xspace/\ensuremath{\mathsf{ILP}}\xspace$

Martins et al.(2009): Proposes two representations of non-projective dependency parsing as a linear programming relaxation as well as an exact ILP.

Use an LP/ILP Solver for decoding

We compare:

- Accuracy
- Exactness
- Speed

Both LP and dual decomposition methods use the same model, features, and weights w.

Comparison to LP/ILP: Accuracy



All decoding methods have comparable accuracy

Comparison to LP/ILP: Exactness and Speed



Percentage with exact solution

Sentences per second

Summary

presented Lagrangian relaxation as a method for decoding in NLP

formal guarantees

- gives certificate or approximate solution
- can improve approximate solutions by tightening relaxation

efficient algorithms

- uses fast combinatorial algorithms
- can improve speed with lazy decoding

widely applicable

 demonstrated algorithms for a wide range of NLP tasks (parsing, tagging, alignment, mt decoding)

References I

- Y. Chang and M. Collins. Exact Decoding of Phrase-based Translation Models through Lagrangian Relaxation. In *To appear proc. of EMNLP*, 2011.
- J. DeNero and K. Macherey. Model-Based Aligner Combination Using Dual Decomposition. In *Proc. ACL*, 2011.
- J. Duchi, D. Tarlow, G. Elidan, and D. Koller. Using Combinatorial Optimization within Max-Product Belief Propagation. In *NIPS*, pages 369–376, 2007.
- D. Klein and C.D. Manning. Factored A* Search for Models over Sequences and Trees. In *Proc IJCAI*, volume 18, pages 1246–1251. Citeseer, 2003.
- N. Komodakis, N. Paragios, and G. Tziritas. Mrf energy minimization and beyond via dual decomposition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2010. ISSN 0162-8828.

References II

- Terry Koo, Alexander M. Rush, Michael Collins, Tommi Jaakkola, and David Sontag. Dual decomposition for parsing with non-projective head automata. In *EMNLP*, 2010. URL http://www.aclweb.org/anthology/D10-1125.
- B.H. Korte and J. Vygen. *Combinatorial Optimization: Theory and Algorithms*. Springer Verlag, 2008.
- A.M. Rush and M. Collins. Exact Decoding of Syntactic Translation Models through Lagrangian Relaxation. In *Proc. ACL*, 2011.
- A.M. Rush, D. Sontag, M. Collins, and T. Jaakkola. On Dual Decomposition and Linear Programming Relaxations for Natural Language Processing. In *Proc. EMNLP*, 2010.
- D.A. Smith and J. Eisner. Dependency Parsing by Belief Propagation. In *Proc. EMNLP*, pages 145–156, 2008. URL http://www.aclweb.org/anthology/D08-1016.