

# Lagrangian Relaxation for MAP Inference

SPFLODD

October 8, 2013

# Outline

- An elegant example of a relaxation to TSP
- A common problem in NLP: finding consensus
- Basic Lagrangian relaxation
- Solving the problem with subgradient
- AD<sup>3</sup>: an alternative approach to decomposition and optimization using the augmented Lagrangian

# Traveling Salesman Problem

- Given: a graph  $(V, E)$  with edge weight function  $\theta$
- Tour: a subset of  $E$  corresponding to a path that starts and ends in the same place, and visits every other node exactly once.
- TSP: Find the maximum-scoring tour.
  - NP-hard

$$\max_{y \in \mathcal{Y}_{\text{tour}}} \sum_{e \in E} y_e \theta_e$$

# Another Problem

- 1-tree: a tree on edges for  $\{2, \dots, |V|\}$ , plus two edges from  $E$  that link the tree to vertex 1.
  - All tours are 1-trees.
  - All 1-trees where every vertex has degree 2 are tours.
  - Easy to solve.

# Held and Karp (1971)

$$\mathcal{Y}_{\text{tour}} = \left\{ y : y \in \mathcal{Y}_{1\text{-tree}} \wedge \forall i \in \{1, \dots, |V|\}, \sum_{e:i \in e} y_e = 2 \right\}$$

$$\max_{y \in \mathcal{Y}_{\text{tour}}} \sum_{e \in E} y_e \theta_e$$

transforming the constraints

$$\max_{y \in \mathcal{Y}_{1\text{-tree}}} \sum_{e \in E} y_e \theta_e \text{ s.t. } \forall i, \sum_{e:i \in e} y_e = 2$$

Lagrangian dual

$$L(u) = \max_{y \in \mathcal{Y}_{1\text{-tree}}} \sum_{e \in E} y_e \theta_e + \sum_{i=1}^{|V|} u_i \left( \sum_{e:i \in e} y_e - 2 \right)$$

# LR Algorithm for TSP

1. Initialize  $u^{(0)} = 0$
2. Repeat for  $k = 1, 2, \dots$ :

$$y^{(k)} \leftarrow \arg \max_{y \in \mathcal{Y}_{1\text{-tree}}} \sum_{e \in E} y_e \theta_e + \sum_{i=1}^{|V|} u_i^{(k-1)} \left( \sum_{e: i \in e} y_e - 2 \right)$$

$$\forall i, u_i^{(k)} \leftarrow u_i^{(k-1)} - \delta_k \left( \sum_{e: i \in e} y_e - 2 \right)$$

If this converges to a solution that satisfies the constraints, it is a solution to the TSP.

# Lagrangian Relaxation, More Generally

- Assume a linear scoring function that is “hard” to maximize.

$$\max_{\mathbf{y} \in \mathcal{Y}} \boldsymbol{\theta}^\top \mathbf{y}$$

- Rewrite the problem as something easier, with linear constraints (relaxation):  
 $\mathcal{Y} = \{\mathbf{y} \in \mathcal{Y}' : \mathbf{A}\mathbf{y} = \mathbf{b}\}$

$$\max_{\mathbf{y} \in \mathcal{Y}'} \boldsymbol{\theta}^\top \mathbf{y}$$

$$\text{s.t. } \mathbf{A}\mathbf{y} = \mathbf{b}$$

- Tackle the dual problem:

$$\min_{\mathbf{u}} \max_{\mathbf{y} \in \mathcal{Y}'} \boldsymbol{\theta}^\top \mathbf{y} + \mathbf{u}^\top (\mathbf{A}\mathbf{y} - \mathbf{b})$$

# Theory

- The dual function (of  $\mathbf{u}$ ) upper bounds the MAP problem.
- A subgradient algorithm can be applied to minimize the dual; it will converge in the limit.
- If the solution to the dual problem satisfies the constraints, it is also a solution to the primal (relaxed) problem ( $\mathcal{Y}'$ ).
  - If the relaxation is *tight*, we also have a solution to the original primal problem ( $\mathcal{Y}$ ).



# Dual Decomposition (A Special Case of LR)

- Assume the objective decomposes into two parts, coupled only through the linear constraints:

$$\begin{aligned} \max_{\mathbf{y} \in \mathcal{Y}, \mathbf{z} \in \mathcal{Z}} \quad & \boldsymbol{\theta}^\top \mathbf{y} + \boldsymbol{\psi}^\top \mathbf{z} \\ \text{s.t.} \quad & \mathbf{A}\mathbf{y} + \mathbf{C}\mathbf{z} = \mathbf{b} \end{aligned}$$

- The relaxation:

$$\max_{\mathbf{y} \in \mathcal{Y}, \mathbf{z} \in \mathcal{Z}} \boldsymbol{\theta}^\top \mathbf{y} + \boldsymbol{\psi}^\top \mathbf{z} \equiv \left( \max_{\mathbf{y} \in \mathcal{Y}} \boldsymbol{\theta}^\top \mathbf{y}, \max_{\mathbf{z} \in \mathcal{Z}} \boldsymbol{\psi}^\top \mathbf{z} \right)$$

# Dual Decomposition

$$\min_{\mathbf{u}} \max_{\mathbf{y} \in \mathcal{Y}, \mathbf{z} \in \mathcal{Z}} \boldsymbol{\theta}^\top \mathbf{y} + \boldsymbol{\psi}^\top \mathbf{z} + \mathbf{u}^\top (\mathbf{A}\mathbf{y} + \mathbf{C}\mathbf{z} - \mathbf{b})$$

1. Initialize  $\mathbf{u}^{(0)} = \mathbf{0}$
2. Repeat for  $k = 1, 2, \dots$ :

$$\mathbf{y}^{(k)} \leftarrow \max_{\mathbf{y} \in \mathcal{Y}} \boldsymbol{\theta}^\top \mathbf{y} + \mathbf{u}^{(k-1)\top} \mathbf{A}\mathbf{y}$$

$$\mathbf{z}^{(k)} \leftarrow \max_{\mathbf{z} \in \mathcal{Z}} \boldsymbol{\psi}^\top \mathbf{z} + \mathbf{u}^{(k-1)\top} \mathbf{C}\mathbf{z}$$

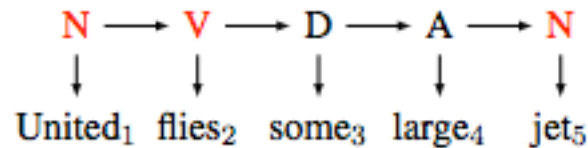
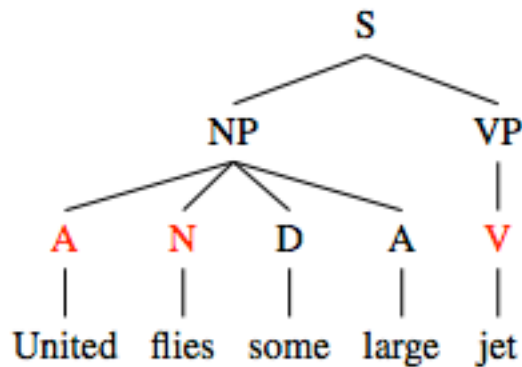
$$\mathbf{u}^{(k)} \leftarrow \mathbf{u}^{(k-1)} - \delta_k \left( \mathbf{A}\mathbf{y}^{(k)} + \mathbf{C}\mathbf{z}^{(k)} - \mathbf{b} \right)$$

# Consensus Problems in NLP

- Key example:
  - Find the jointly-best parse (under a WCFG) and sequence labeling (under an HMM); see Rush et al. (2010)
- Other examples:
  - Finding a lexicalized phrase structure parse that is jointly-best under a WCFG and a dependency model (Rush et al., 2010)
  - Decoding in phrase-based translation (Chang and Collins, 2011).

# Example Run (k = 1)

$$\forall i \in \{1, \dots, n\}, \forall N \in \mathcal{N}, \mathbf{y}[N, i, i] = \mathbf{z}[N, i]$$



$$\mathbf{u}[N, i]^{(1)} = \mathbf{u}[N, i]^{(0)} - \delta_k \left( \mathbf{y}[N, i, i]^{(1)} - \mathbf{z}[N, i]^{(1)} \right)$$

$$\mathbf{u}[A, 1] = -1$$

$$\mathbf{u}[N, 2] = -1$$

$$\mathbf{u}[V, 5] = -1$$

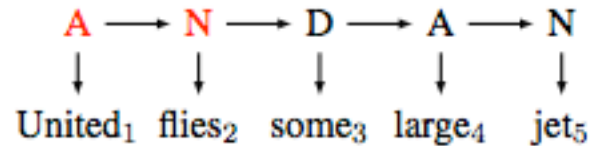
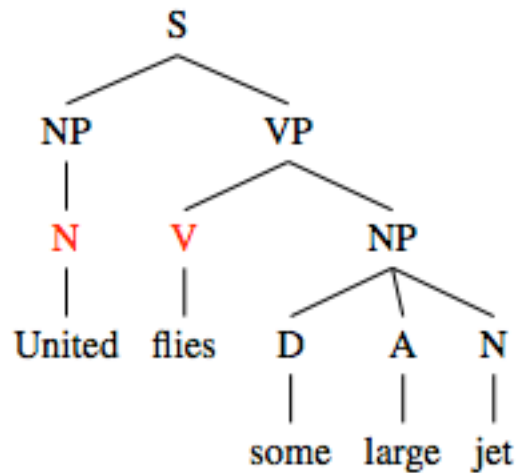
$$\mathbf{u}[N, 1] = 1$$

$$\mathbf{u}[V, 2] = 1$$

$$\mathbf{u}[N, 5] = 1$$

# Example Run (k = 2)

$$\forall i \in \{1, \dots, n\}, \forall N \in \mathcal{N}, \mathbf{y}[N, i, i] = \mathbf{z}[N, i]$$



$$\mathbf{u}[N, i]^{(2)} = \mathbf{u}[N, i]^{(1)} - \delta_k \left( \mathbf{y}[N, i, i]^{(2)} - \mathbf{z}[N, i]^{(2)} \right)$$

$$\downarrow \mathbf{u}[N, 1]$$

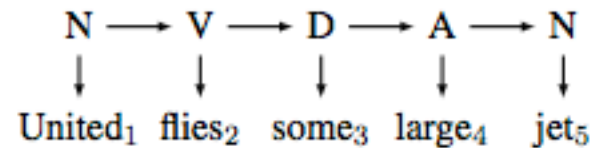
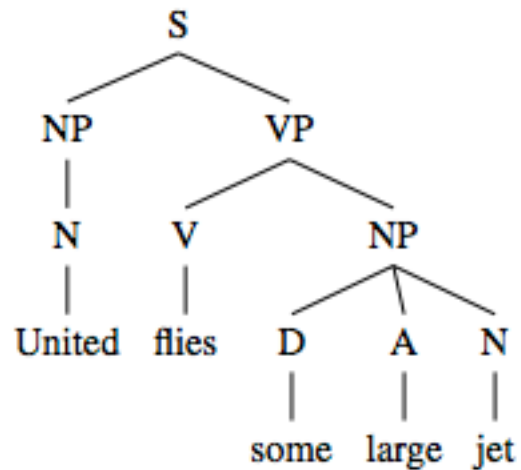
$$\downarrow \mathbf{u}[V, 1]$$

$$\uparrow \mathbf{u}[A, 1]$$

$$\uparrow \mathbf{u}[N, 1]$$

# Example Run (k = 3)

$$\forall i \in \{1, \dots, n\}, \forall N \in \mathcal{N}, \mathbf{y}[N, i, i] = \mathbf{z}[N, i]$$



# “Certificate”

- Proof that we have solved the original problem: constraints hold.
  - This is easy to check given  $\mathbf{y}$  and  $\mathbf{z}$ .
- In published NLP papers so far, this happens most of the time (better than 98%).

# What can go wrong?

- It can take many iterations to converge.
- Oscillation between different solutions; failure to agree.
  - Suggested solution: add more variables for “bigger parts” and enforce agreement among them with more constraints.



# What does this have to do with ILP?

- The linear constraints are expressed in terms of an integer-vector representation of the output space.
  - Just like when we treated decoding as an ILP.
- The subproblems *could* be expressed as ILPs, though we'd prefer to use poly-time combinatorial algorithms to solve them if we can.

# Consensus Problems, Revisited

- What if we just have a hard combinatorial optimization problem?
  - There isn't always a straightforward decomposition.
- Martins et al. (2011): shatter a decoding problem into *many* “small” subproblems (instead of two “big” ones).
  - Instead of dynamic programming as a subroutine, LP relaxations of “small” subproblems.
  - Extra LP relaxation step.

# Martins' Alternative Formulation

- Original problem: 
$$\max_{\mathbf{y}_1 \in \mathcal{Y}_1, \dots, \mathbf{y}_S \in \mathcal{Y}_S, \mathbf{w} \in \mathbb{R}^D} \sum_{s=1}^S \boldsymbol{\theta}_s^\top \mathbf{y}_s$$
$$\text{s.t. } \forall s, \mathbf{A}_s \mathbf{w} = \mathbf{y}_s$$
- Convex relaxation: 
$$\max_{\mathbf{y}_1 \in \text{conv}(\mathcal{Y}_1), \dots, \mathbf{y}_S \in \text{conv}(\mathcal{Y}_S), \mathbf{w} \in \mathbb{R}^D} \sum_{s=1}^S \boldsymbol{\theta}_s^\top \mathbf{y}_s$$
$$\text{s.t. } \forall s, \mathbf{A}_s \mathbf{w} = \mathbf{y}_s$$
- Dual:

$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_S} \max_{\mathbf{y}_1 \in \text{conv}(\mathcal{Y}_1), \dots, \mathbf{y}_S \in \text{conv}(\mathcal{Y}_S), \mathbf{w} \in \mathbb{R}^D} \sum_{s=1}^S \boldsymbol{\theta}_s^\top \mathbf{y}_s + \sum_s \mathbf{u}_s^\top (\mathbf{y}_s - \mathbf{A}_s \mathbf{w})$$

# Augmented Lagrangian (Hestenes, 1969; Powell, 1969)

$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_S} \max_{\mathbf{y}_1 \in \text{conv}(\mathcal{Y}_1), \dots, \mathbf{y}_S \in \text{conv}(\mathcal{Y}_S), \mathbf{w} \in \mathbb{R}^D} \sum_{s=1}^S \boldsymbol{\theta}_s^\top \mathbf{y}_s + \sum_s \mathbf{u}_s^\top (\mathbf{y}_s - \mathbf{A}_s \mathbf{w}) + \frac{\rho}{2} \sum_s \|\mathbf{y}_s - \mathbf{A}_s \mathbf{w}\|_2^2$$



$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_S} \max_{\mathbf{y}_1 \in \text{conv}(\mathcal{Y}_1), \dots, \mathbf{y}_S \in \text{conv}(\mathcal{Y}_S), \mathbf{w} \in \mathbb{R}^D} \sum_{s=1}^S \boldsymbol{\theta}_s^\top \mathbf{y}_s + \sum_s \mathbf{u}_s^\top (\mathbf{y}_s - \mathbf{A}_s \mathbf{w})$$

# Alternating Directions Method of Multipliers

(Gabay and Mercier, 1976; Glowinski and Marroco, 1975)

## Dual Decomposition (AD<sup>3</sup>)

- Alternate between updating  $\mathbf{y}$  and  $\mathbf{w}$ :

$$\forall s, \mathbf{y}_s \leftarrow \arg \max_{\mathbf{y}_s \in \text{conv}(\mathcal{Y}_s)} \boldsymbol{\theta}_s^\top \mathbf{y}_s + \mathbf{u}_s^\top \mathbf{y}_s + \frac{\rho}{2} \|\mathbf{y}_s - \mathbf{A}_s \mathbf{w}\|_2^2$$

$$\mathbf{w} \leftarrow \arg \max_{\mathbf{w}} \sum_s \mathbf{u}_s^\top \mathbf{A}_s \mathbf{w} + \frac{\rho}{2} \sum_s \|\mathbf{y}_s - \mathbf{A}_s \mathbf{w}\|_2^2$$

- Subgradient step for dual variables  $\mathbf{u}$  is similar to before:

$$\forall s, \mathbf{u}_s^{(k)} \leftarrow \mathbf{u}_s^{(k-1)} - \delta_k (\mathbf{y}_s - \mathbf{A}_s \mathbf{w})$$

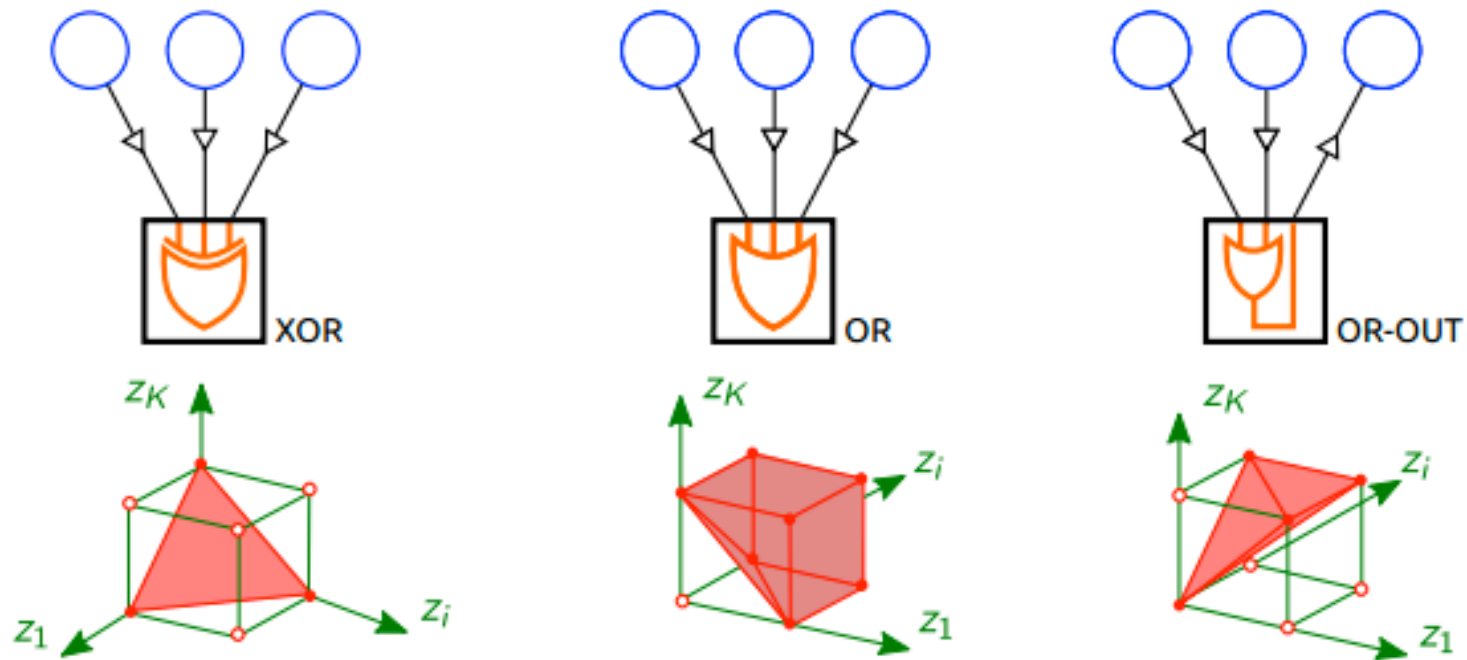
# Massive Decomposition

- Most extreme: every factor (MN) or “part” is a separate subproblem.

$$\forall s, \mathbf{y}_s \leftarrow \arg \max_{\mathbf{y}_s \in \text{conv}(\mathcal{Y}_s)} \boldsymbol{\theta}_s^\top \mathbf{y}_s + \mathbf{u}_s^\top \mathbf{y}_s + \frac{\rho}{2} \|\mathbf{y}_s - \mathbf{A}_s \mathbf{w}\|_2^2$$

- Some kinds of MN factors can be solved very efficiently ...

# XOR, OR, OR-with-Output Solvable in $O(K \log K)$



# AD<sup>3</sup> and “Big” Subproblems?

- Return to Rush and Collins’ example.
  - One subproblem is “WCFG” and one is “HMM tagger.”

$$\forall s, \mathbf{y}_s \leftarrow \arg \max_{\mathbf{y}_s \in \text{conv}(\mathcal{Y}_s)} \boldsymbol{\theta}_s^\top \mathbf{y}_s + \mathbf{u}_s^\top \mathbf{y}_s + \frac{\rho}{2} \|\mathbf{y}_s - \mathbf{A}_s \mathbf{w}\|_2^2$$

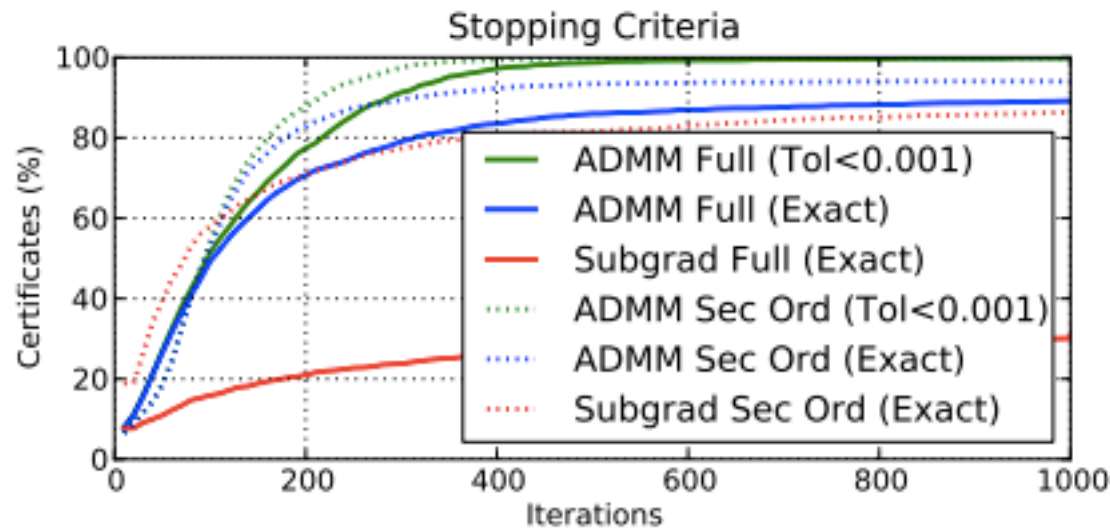
- In dependency parsing, “max arborescence” might be a subproblem.
  - Why can’t we use AD<sup>3</sup>?



# Pros and Cons

- Con: Subproblems are now *quadratic*.
  - Linear decoders as subroutines?
- Con: Fractional solutions.
- Pro: Better stopping criteria: residuals.
  - Primal residuals measure amount by which primal constraints are violated.
  - Dual residuals measure amount by which dual optimality is violated.
- Pro: Certificates as before (for each  $s$ ,  $\mathbf{A}_s \mathbf{w} = \mathbf{y}_s$ )

# Convergence of AD<sup>3</sup> vs. Subgradient



Dependency parsing:

- ADMM = AD<sup>3</sup>
- Sec Ord = Second order model for which subgradient optimization is possible
- Full = second order model with all-siblings, directed paths, and non-projective arcs

# Take-Home Messages

- Dual decomposition is useful for consensus problems.
  - Subgradient DD when there are a few subproblems with good specialized solvers.
  - AD<sup>3</sup> when you've got a big problem with lots of hard and soft constraints. (There is a library.)
- Attractive guarantees (cf. beam search).
- Only MAP inference.

# References

- “A tutorial on dual decomposition and Lagrangian relaxation for inference in natural language processing,” by A. Rush and M. Collins, *JAIR* 45:305-362, 2013.
- “Alternating directions dual decomposition” by A. Martins et al., arXiv 1212.6550.