

Lecture 7: Structured Prediction and Graphical Models

André Martins



instituto de
telecomunicações



Deep Structured Learning Course, Fall 2018

Announcements

- The project proposals have been graded. The deadline for the midterm report is November 28 and the final report is due January 2. The class presentations will be in January 9 and 16.
- The deadline for turning in Homework 2 is this Wednesday.
- Homework 3 is out.
- Due to dissertation season, we need to change the rooms for the November 14 and 21 classes. Check the webpage.

Today's Roadmap

Two lectures ago we talked about sequences and dynamic programming.
Today we'll extend this to **trees** and **graphs**:

- Structured prediction.
- Bayesian and Markov networks.
- Factor graphs.
- MAP Inference and Integer Linear Programming.
- Marginal Inference and Message-Passing Algorithms.
- Dual decomposition and linear programming relaxations.

Outline

① Structured Prediction and Factor Graphs

② Integer Linear Programming

③ Message-Passing Algorithms

Sum-Product

Max-Product

④ Dual Decomposition

Structured Prediction

- Input set \mathcal{X}
- For each $x \in \mathcal{X}$: a **large** set of candidate outputs $\mathcal{Y}(x)$
- A compatibility function $F_w(x, y)$ induced by a model w
(Linear model: $F_w(x, y) = w \cdot \phi(x, y)$)

Structured Prediction

- Input set \mathcal{X}
- For each $x \in \mathcal{X}$: a **large** set of candidate outputs $\mathcal{Y}(x)$
- A compatibility function $F_w(x, y)$ induced by a model w
(Linear model: $F_w(x, y) = w \cdot \phi(x, y)$)
- **Training problem:** learn the model w from data $\{\langle x_i, y_i \rangle\}_{i=1}^M$
- **Decoding problem (our focus):**

$$\hat{y} = \arg \max_{y \in \mathcal{Y}(x)} F_w(x, y)$$

Structured Prediction

- Input set \mathcal{X}
- For each $x \in \mathcal{X}$: a **large** set of candidate outputs $\mathcal{Y}(x)$
- A compatibility function $F_w(x, y)$ induced by a model w
(Linear model: $F_w(x, y) = w \cdot \phi(x, y)$)
- **Training problem:** learn the model w from data $\{\langle x_i, y_i \rangle\}_{i=1}^M$
- **Decoding problem (our focus):**

$$\hat{y} = \arg \max_{y \in \mathcal{Y}(x)} F_w(x, y)$$

- **Key assumption:** F_w decomposes into (overlapping) *parts*

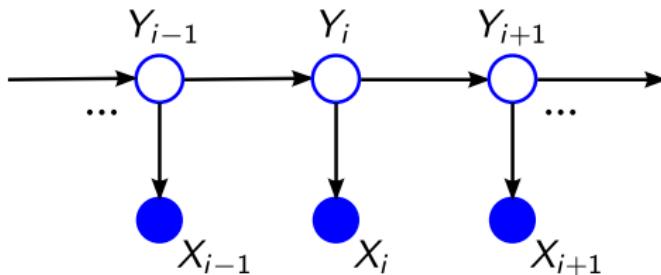
Three Important Questions

- Representation?
- Decoding/Inference?
- Learning the parameters?

Recap: Hidden Markov Models

F_w is a log-probability, factoring over emissions and transitions.

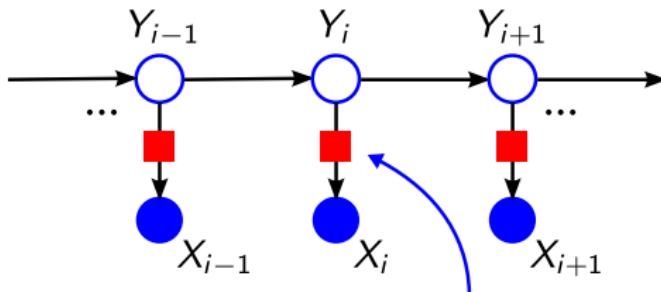
$$\mathbb{P}(x, y) = \prod_i \underbrace{\mathbb{P}(x_i | y_i)}_{\text{emissions}} \underbrace{\mathbb{P}(y_i | y_{i-1})}_{\text{transitions}}$$



Recap: Hidden Markov Models

F_w is a log-probability, factoring over emissions and transitions.

$$\mathbb{P}(x, y) = \prod_i \underbrace{\mathbb{P}(x_i | y_i)}_{\text{emissions}} \underbrace{\mathbb{P}(y_i | y_{i-1})}_{\text{transitions}}$$



$$\psi_i(y_i) := \mathbb{P}(x_i | y_i) \text{ (unary potentials)}$$

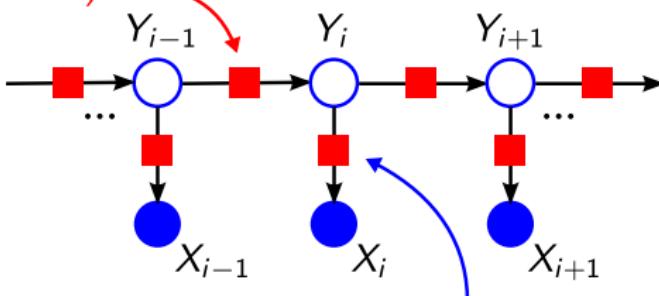
Recap: Hidden Markov Models

F_w is a log-probability, factoring over emissions and transitions.

$$\mathbb{P}(x, y) = \prod_i \underbrace{\mathbb{P}(x_i | y_i)}_{\text{emissions}} \underbrace{\mathbb{P}(y_i | y_{i-1})}_{\text{transitions}}$$

$$\psi_{i,i-1}(y_i, y_{i-1}) := \mathbb{P}(y_i | y_{i-1})$$

(pairwise potentials)



$$\psi_i(y_i) := \mathbb{P}(x_i | y_i) \text{ (unary potentials)}$$

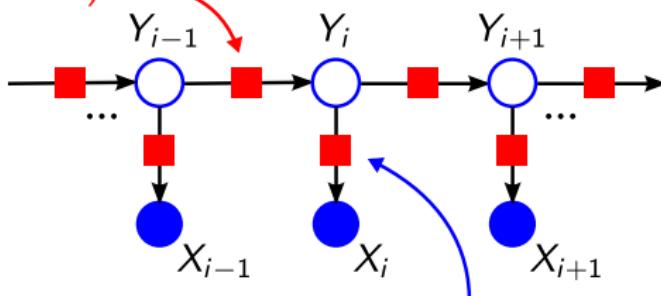
Recap: Hidden Markov Models

F_w is a log-probability, factoring over emissions and transitions.

$$\mathbb{P}(x, y) = \prod_i \underbrace{\mathbb{P}(x_i | y_i)}_{\text{emissions}} \underbrace{\mathbb{P}(y_i | y_{i-1})}_{\text{transitions}} = \prod_i \psi_i(y_i) \prod_i \psi_{i,i-1}(y_i, y_{i-1})$$

$$\psi_{i,i-1}(y_i, y_{i-1}) := \mathbb{P}(y_i | y_{i-1})$$

(pairwise potentials)



$$\psi_i(y_i) := \mathbb{P}(x_i | y_i) \text{ (unary potentials)}$$

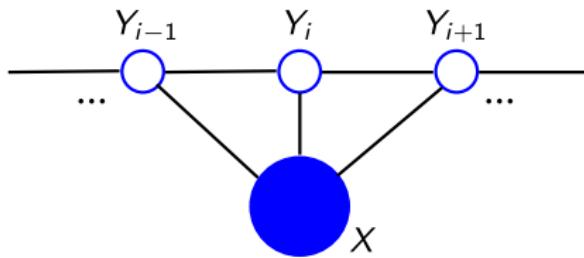
Recap: Hidden Markov Models

- Representation? **Directed sequence model.**
- Decoding/Inference? **Viterbi/forward-backward algorithms.**
- Learning the parameters? **Maximum likelihood (count and normalize).**

Recap: Conditional Random Fields

Same factorization, but globally normalized.

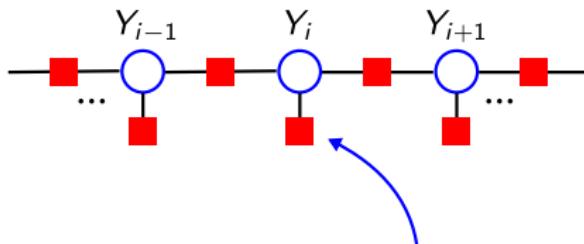
$$\mathbb{P}(y|x) = \frac{1}{Z(\mathbf{w}, x)} \exp \left(\sum_i \underbrace{\mathbf{w} \cdot \phi_i(x, y_i)}_{\text{nodes}} + \sum_i \underbrace{\mathbf{w} \cdot \phi_{i,i-1}(x, y_i, y_{i-1})}_{\text{edges}} \right)$$



Recap: Conditional Random Fields

Same factorization, but globally normalized.

$$\mathbb{P}(y|x) = \frac{1}{Z(\mathbf{w}, x)} \exp \left(\sum_i \underbrace{\mathbf{w} \cdot \phi_i(x, y_i)}_{\text{nodes}} + \sum_i \underbrace{\mathbf{w} \cdot \phi_{i,i-1}(x, y_i, y_{i-1})}_{\text{edges}} \right)$$



$$\psi_i(y_i) := \exp(\mathbf{w}^\top \mathbf{f}(x, y_i)) \text{ (unary potentials)}$$

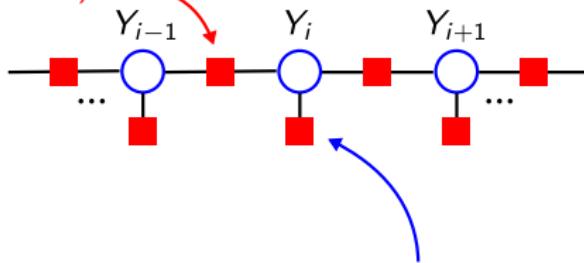
Recap: Conditional Random Fields

Same factorization, but globally normalized.

$$\mathbb{P}(y|x) = \frac{1}{Z(\mathbf{w}, x)} \exp \left(\sum_i \underbrace{\mathbf{w} \cdot \phi_i(x, y_i)}_{\text{nodes}} + \sum_i \underbrace{\mathbf{w} \cdot \phi_{i,i-1}(x, y_i, y_{i-1})}_{\text{edges}} \right)$$

$$\psi_{i,i-1}(y_i, y_{i-1}) := \exp(\mathbf{w}^\top \mathbf{f}(x, y_i, y_{i-1}))$$

(pairwise potentials)



$$\psi_i(y_i) := \exp(\mathbf{w}^\top \mathbf{f}(x, y_i)) \text{ (unary potentials)}$$

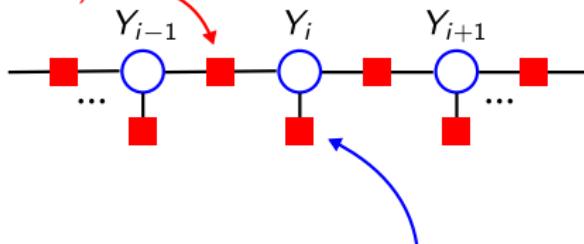
Recap: Conditional Random Fields

Same factorization, but globally normalized.

$$\begin{aligned}\mathbb{P}(y|x) &= \frac{1}{Z(\mathbf{w}, x)} \exp \left(\sum_i \underbrace{\mathbf{w} \cdot \phi_i(x, y_i)}_{\text{nodes}} + \sum_i \underbrace{\mathbf{w} \cdot \phi_{i,i-1}(x, y_i, y_{i-1})}_{\text{edges}} \right) \\ &\propto \prod_i \psi_i(y_i) \prod_i \psi_{i,i-1}(y_i, y_{i-1})\end{aligned}$$

$$\psi_{i,i-1}(y_i, y_{i-1}) := \exp(\mathbf{w}^\top \mathbf{f}(x, y_i, y_{i-1}))$$

(pairwise potentials)



$$\psi_i(y_i) := \exp(\mathbf{w}^\top \mathbf{f}(x, y_i)) \quad (\text{unary potentials})$$

Recap: Conditional Random Fields

- Representation? **Undirected sequence model.**
- Decoding/Inference? **Viterbi/forward-backward algorithms.**
- Learning the parameters? **Maximum conditional likelihood (convex optimization).**

Can we Generalize This to Trees?

Can we Generalize This to Trees?

Yes!

Let $\pi(i)$ denote the parent of node i in the tree.

- Tree-HMMs:

$$\mathbb{P}(x, y) = \prod_i \mathbb{P}(x_i | y_i) \mathbb{P}(y_i | y_{\pi(i)})$$

- Tree-CRFs:

$$\mathbb{P}(y|x) = \frac{1}{Z(\mathbf{w}, x)} \exp \left(\sum_i \mathbf{w} \cdot \phi_i(x, y_i) + \sum_i \mathbf{w} \cdot \phi_{i, \pi(i)}(x, y_i, y_{\pi(i)}) \right)$$

Dynamic programming still works!

(Just run it top-down from the root to the leaves.)

Can we Generalize This to Graphs?

- Sort of, but it's a bit more complicated.
- This leads us to **graphical models**.

Graphical Models

HMMs and CRFs are two instances of *graphical models*.

In general, graphical models come in two flavours:

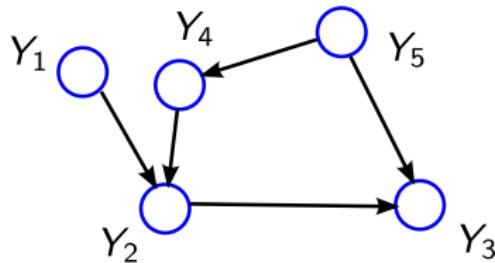
- Directed (Bayesian Networks)
- Undirected (Markov Networks)

Bayesian Networks

Useful to express **causality relations**.

Factors are **conditional probability tables**.

$$\mathbb{P}(y) = \mathbb{P}(y_1)\mathbb{P}(y_2|y_1, y_4)\mathbb{P}(y_3|y_2, y_5)\mathbb{P}(y_4|y_5)\mathbb{P}(y_5)$$



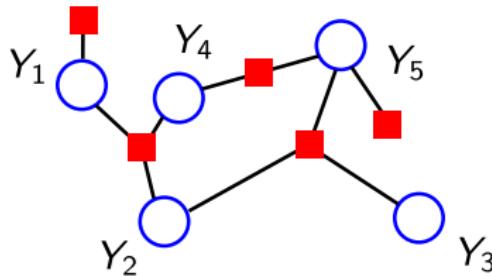
$$\mathbb{P}(y) = \prod_i \mathbb{P}(y_i|\text{parents}(y_i))$$

Bayesian Networks

Useful to express **causality relations**.

Factors are **conditional probability tables**.

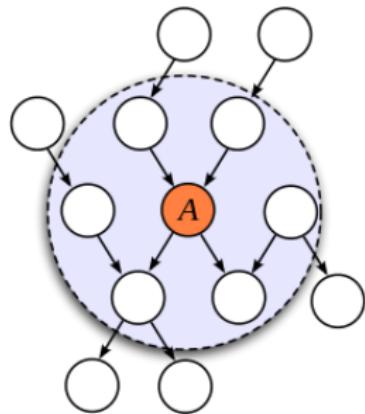
$$\mathbb{P}(y) = \mathbb{P}(y_1)\mathbb{P}(y_2|y_1, y_4)\mathbb{P}(y_3|y_2, y_5)\mathbb{P}(y_4|y_5)\mathbb{P}(y_5)$$



$$\mathbb{P}(y) = \prod_i \mathbb{P}(y_i|\text{parents}(y_i))$$

Markov Blanket

The minimum set of variables given which the target variable is **conditionally independent** of the rest of the variables.



$$\forall X \notin \{A\} \cup \text{MB}(A), \quad A \perp\!\!\!\perp X \mid \text{MB}(A).$$

This consists of:

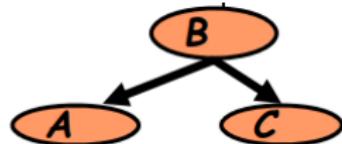
- The parents
- The children
- The co-parents!

Local Structures and Independencies

- Common parent

- Fixing B decouples A and C

"given the level of gene B, the levels of A and C are independent"



- Cascade

- Knowing B decouples A and C

"given the level of gene B, the level gene A provides no extra prediction value for the level of gene C"



- V-structure

- Knowing C couples A and B

because A can "explain away" B w.r.t. C

"If A correlates to C, then chance for B to also correlate to B will decrease"

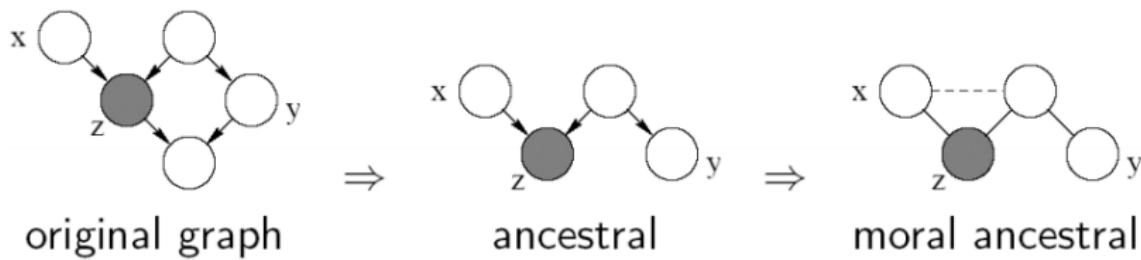


(Credits: Eric Xing)

D-Separation

How to check if variables x and y are **d-separated** (i.e. conditionally independent) given z ?

- ① Draw the **ancestor graph** (x , y , z , and all their ancestors)
- ② **Moralize** the graph (marry the parents) and replace all directed arcs by undirected edges
- ③ Remove the z and all its adjacent edges
- ④ Check if there is a path connecting x and y .



Summary of Bayesian Networks

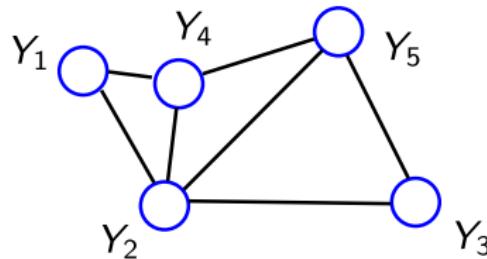
- The structure is a directed acyclic graph (DAG)
- A node is conditionally independent of every other node in the network given its Markov blanket
- Local conditional distributions and the DAG structure completely determine the joint distribution
- Generative process: sampling is easy (traverse the DAG in topological order)

Markov Networks

Useful to express **correlations between variables**.

Factors correspond to **cliques of the graph**.

$$\mathbb{P}(y) = \frac{1}{Z} \psi_{124}(y_1, y_2, y_4) \psi_{235}(y_2, y_3, y_5) \psi_{245}(y_2, y_4, y_5)$$



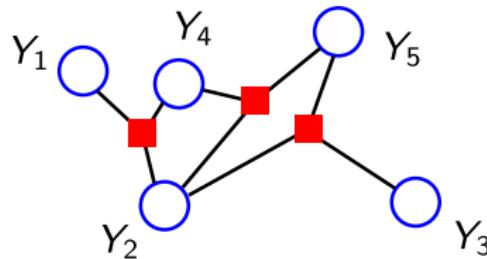
$$\mathbb{P}(y) \propto \prod_{s \in \text{cliques}(G)} \psi_s(y_s)$$

Markov Networks

Useful to express **correlations between variables**.

Factors correspond to **cliques of the graph**.

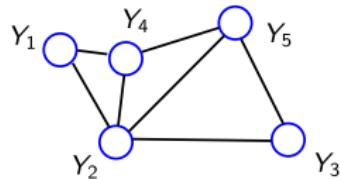
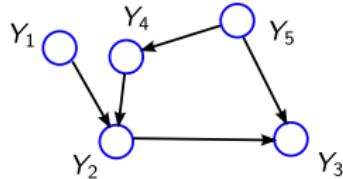
$$\mathbb{P}(y) = \frac{1}{Z} \psi_{124}(y_1, y_2, y_4) \psi_{235}(y_2, y_3, y_5) \psi_{245}(y_2, y_4, y_5)$$



$$\mathbb{P}(y) \propto \prod_{s \in \text{cliques}(G)} \psi_s(y_s)$$

Convert a Bayesian Net Into a Markov Net

$$\mathbb{P}(y) = \mathbb{P}(y_1)\mathbb{P}(y_2|y_1, y_4)\mathbb{P}(y_3|y_2, y_5)\mathbb{P}(y_4|y_5)\mathbb{P}(y_5) \quad \mathbb{P}(y) = \frac{1}{Z} \psi_{124}(y_1, y_2, y_4) \psi_{235}(y_2, y_3, y_5) \psi_{245}(y_2, y_4, y_5)$$



- **Moralize** the graph (marry the parents)
- Convert all direct arcs into undirected edges
- The resulting cliques are the factors of the distribution.

Summary of Markov Networks

- The structure is a undirected graph
- A node is conditionally independent of every other node in the network given its direct neighbors
- Local contingency functions (potentials) and the cliques in the graph completely determine the joint distribution
- No generative process: sampling may be hard

Conditional Independence

Graphical models are a great tool for modeling conditional independence
They link properties of the probability distribution with properties of the graph (reachability, D-separation, etc.)
Lots of literature about this: Pearl (1988); Lauritzen (1996); Koller and Friedman (2009)

MAP/Marginal Inference in a Graphical Model

We can always reduce a Bayes net to a Markov net, so let's focus on the latter.

- Variable elimination algorithm (generalizes forward-backward)
- Belief propagation algorithm (for tree structures)
- Junction tree algorithm (for general graphs)

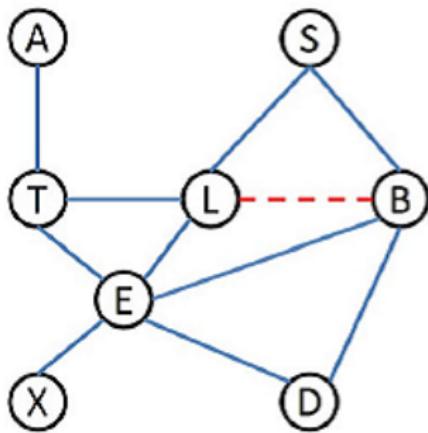
Junction Tree Algorithm

A generic **exact inference** algorithm for any graphical model!

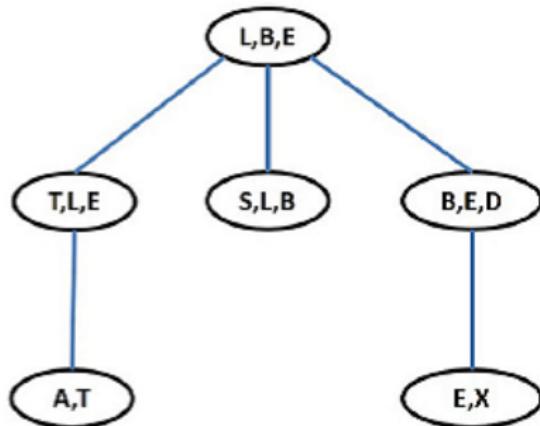
- ① Triangulate the graph to make it chordal
- ② Create super-nodes for the cliques of the graph and transform the graph into a **junction tree** (note: this may cause an exponential blow-up in the number of variables)
- ③ Run belief propagation on the junction tree.

Complexity: exponential in the size of the maximal clique (called the **treewidth** of the original graph)

Example of Junction Tree



(a) Triangulated graph

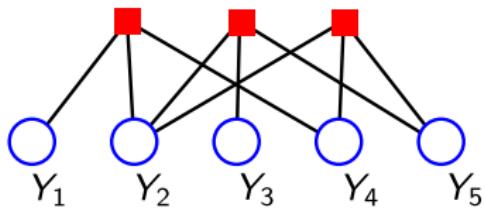
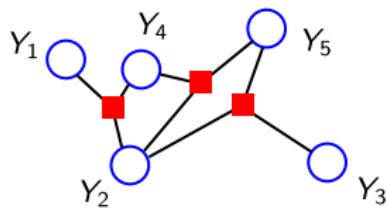


(b) Junction tree

An Intermediate Representation: Factor Graph

Sometimes it's more convenient to work with **factor graphs**, another graphical model representation:

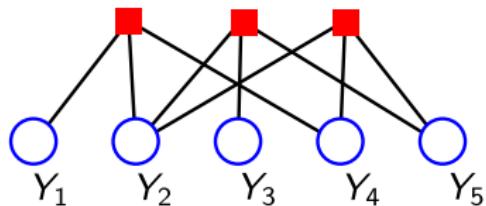
$$\mathbb{P}(y) = \frac{1}{Z} \psi_{124}(y_1, y_2, y_4) \psi_{235}(y_2, y_3, y_5) \psi_{245}(y_2, y_4, y_5)$$



Factor Graph

A bipartite graph with **variable nodes** and **factor nodes**

It makes *explicit* the factors of the distribution



$$\mathbb{P}(y) \propto \underbrace{\prod_i \psi_i(y_i)}_{\text{unitary potentials}} \times \underbrace{\prod_s \psi_s(y_s)}_{\text{higher-order potentials}}$$

With unary potentials only, all variables would be independent

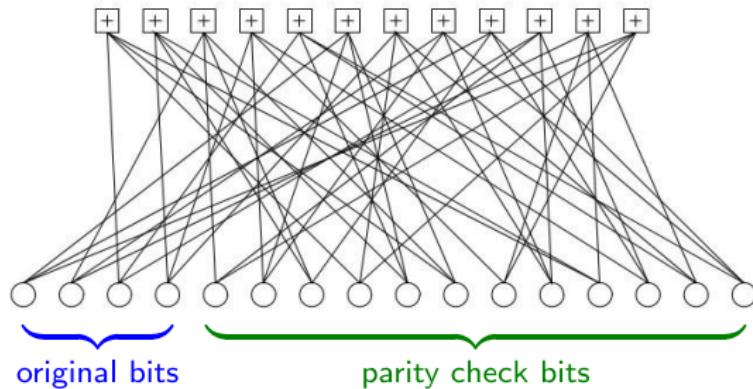
Higher-order potentials can model correlations, impose soft/hard constraints, etc.

Example: Low-Density Parity Check Codes

A message is transmitted through a noisy channel, corrupting some bits

Redundancy can help decoding the message, e.g. via additional parity check bits that can detect/correct errors (error-correcting codes)

High-level idea: increase redundancy to build more accurate decoders



(Adapted from MacKay 2003.)

Inference/Decoding

$$\mathbb{P}_\psi(y|x) = \frac{1}{Z(\psi, x)} \times \underbrace{\prod_i \psi_i(y_i)}_{\text{unary potentials}} \times \underbrace{\prod_s \psi_s(\mathbf{y}_s)}_{\text{higher-order potentials}}$$

Two decoding problems:

- **MAP decoding:** compute $\hat{y} = \arg \max_y \mathbb{P}_\psi(y|x)$
- **Marginal decoding:** compute every $\mathbb{P}_\psi(y_i|x)$ and $\mathbb{P}_\psi(\mathbf{y}_s|x)$; and evaluate the partition function $Z(\psi, x)$

Sometimes easy, in general intractable...

When is Decoding Easy?

- independent variables (trivial)
- sequence models (Viterbi, forward-backward)
- graphical models without cycles (variable elimination, belief propagation)
- graphical models with low treewidth (junction tree algorithm)

When is Decoding Easy?

- independent variables (trivial)
- sequence models (Viterbi, forward-backward)
- graphical models without cycles (variable elimination, belief propagation)
- graphical models with low treewidth (junction tree algorithm)

In general, for graphs with cycles, MAP decoding is NP-hard and marginal decoding is #P-hard

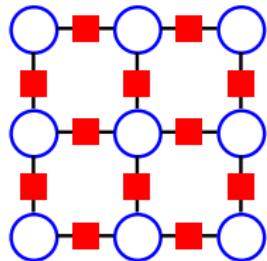
When is Decoding Easy?

- independent variables (trivial)
- sequence models (Viterbi, forward-backward)
- graphical models without cycles (variable elimination, belief propagation)
- graphical models with low treewidth (junction tree algorithm)

In general, for graphs with cycles, MAP decoding is NP-hard and marginal decoding is #P-hard

Note: tractability depends not only on the topology, but also on the *potentials*

Example: Ising and Potts Models



Ising/Potts grid



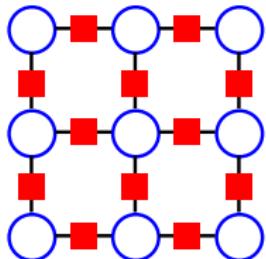
Ernst Ising, 1900–1998



Ren Potts, 1925–2005

All factors are pairwise, variables are binary (Ising) or multi-class (Potts)

Example: Ising and Potts Models



Ising/Potts grid



Ernst Ising, 1900–1998



Ren Potts, 1925–2005

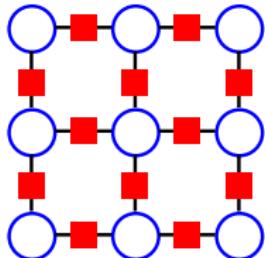
All factors are pairwise, variables are binary (Ising) or multi-class (Potts)

MAP decoding is tractable for *attractive* Ising models (i.e. Ising models with *supermodular* log-potentials):

$$\log \psi_{ij}(1, 1) + \log \psi_{ij}(0, 0) \geq \log \psi_{ij}(0, 1) + \log \psi_{ij}(1, 0)$$

Good approximations for *attractive* Potts models

Example: Ising and Potts Models



Ising/Potts grid



Ernst Ising, 1900–1998



Ren Potts, 1925–2005

All factors are pairwise, variables are binary (Ising) or multi-class (Potts)

MAP decoding is tractable for *attractive* Ising models (i.e. Ising models with *supermodular* log-potentials):

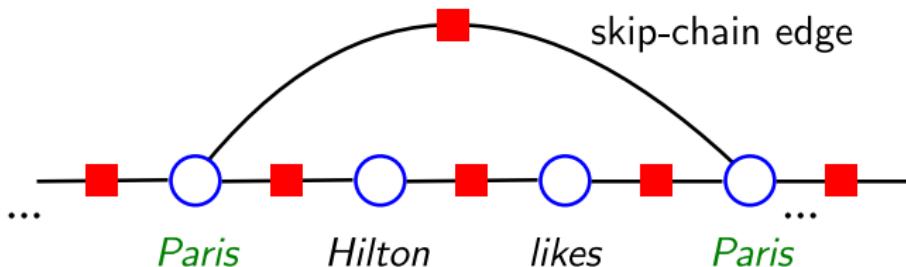
$$\log \psi_{ij}(1, 1) + \log \psi_{ij}(0, 0) \geq \log \psi_{ij}(0, 1) + \log \psi_{ij}(1, 0)$$

Good approximations for *attractive* Potts models

... but the general case is NP-hard and hard to approximate

Example: Skip-Chain CRFs

Skip-chain CRFs are useful to model long-range dependencies



Skip-chains introduce cycles, making decoding more expensive

We could write this information in the “state” and still decode with dynamic programming, but that would blow up the number of states

Beyond Graphical Models

Some NLP problems (e.g. parsing) require representations beyond graphical models

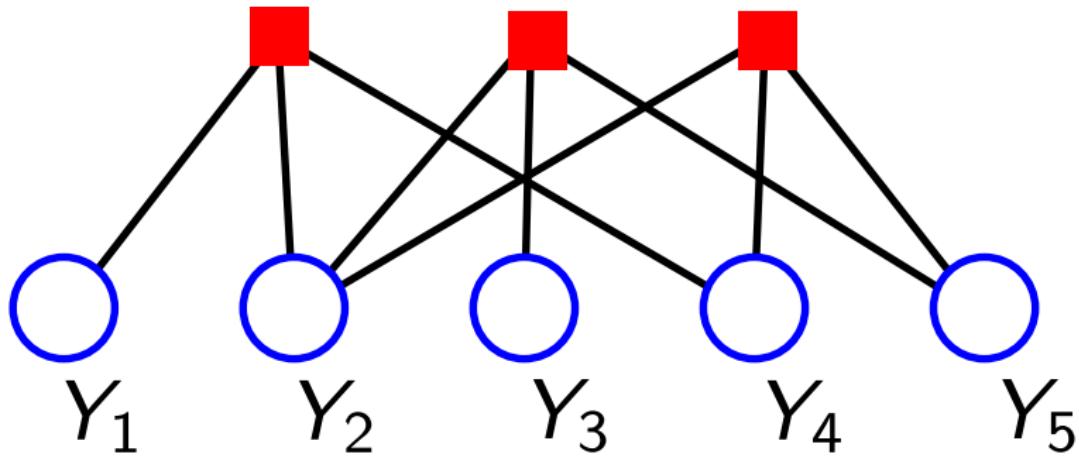
Dynamic programming algorithms (CKY, inside-outside) still work for those representations

Example: case-factor diagrams (McAllester et al., 2008)

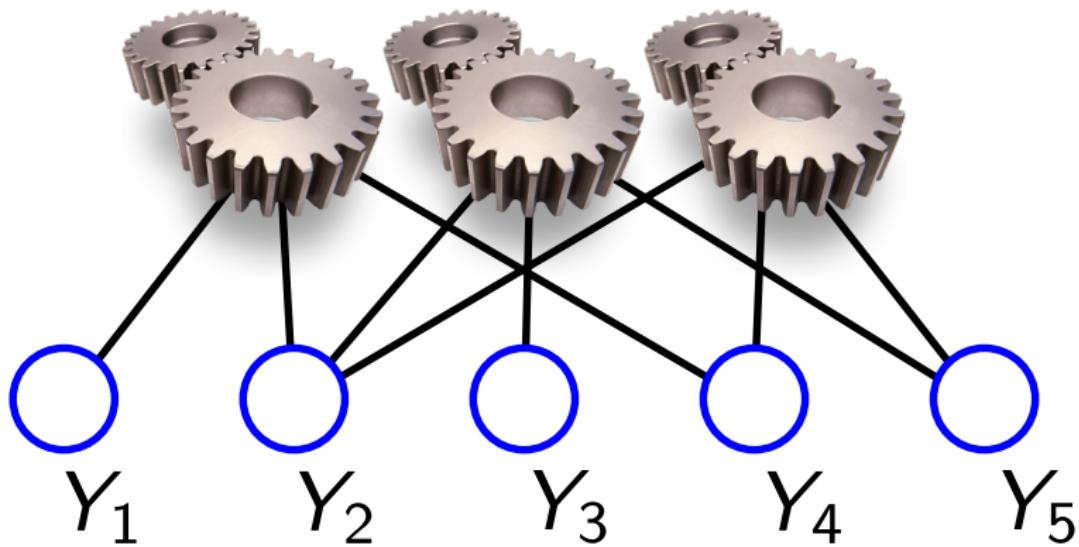
Other problems (e.g. matching, spanning trees) can be solved with combinatorial algorithms not related with dynamic programming

All these can still be represented as GMs by “generalizing” the notion of factor

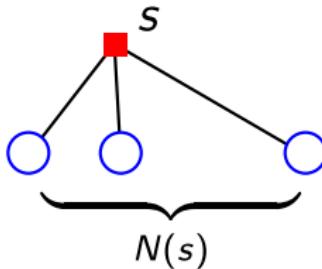
Factors as Machines



Factors as Machines



Three Kinds of Factors



Let $N(s)$ denote the set of variables that are *neighbors* of factor s .
(Its cardinality $|N(s)|$ is called the *degree* of s .)

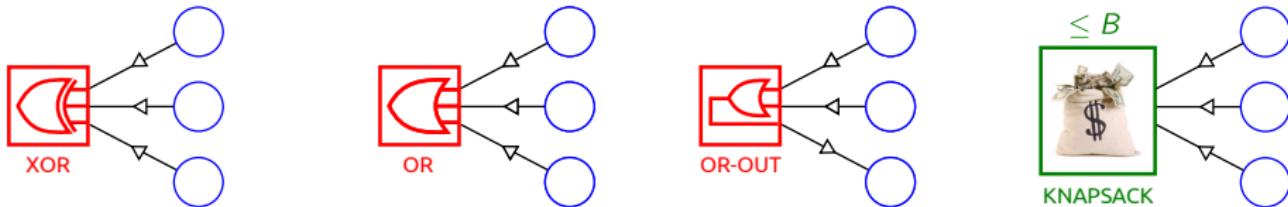
- ① **Dense factors:** $\psi_s(\mathbf{y}_s)$ has all $O(\exp(|N(s)|))$ degrees of freedom
- ② **Structured factors:** $\psi_s(\mathbf{y}_s)$ has internal structure
- ③ **Hard constraint factors:**

$$\psi_s(\mathbf{y}_s) := \begin{cases} 1, & \text{if } \mathbf{y}_s \in \mathcal{Y}_s \\ 0, & \text{otherwise.} \end{cases}$$

Examples of Structured Factors

- a factor for bipartite matching (Duchi et al., 2007)
- combining a sequential model (POS tagger) with a PCFG (Rush et al., 2010)
- combining CCG parsing and supertagging (Auli and Lopez, 2011)
- dependency parsing with head automata (Smith and Eisner, 2008; Koo et al., 2010)
- handling string-valued variables with factors that are finite state transducers (Dreyer and Eisner, 2009)
- inversion transduction grammar constraint (Burkett and Klein, 2012)

Examples of Hard Constraint Factors



Logic factors: can express arbitrary FOL constraints

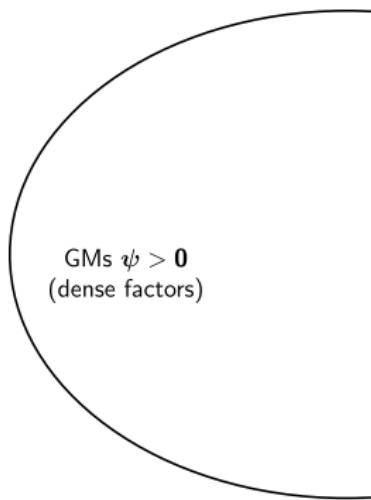
- Applications: Markov logic networks (Richardson and Domingos, 2006), constrained conditional models (Roth and Yih, 2004)

Knapsack factors: can express budget constraints

- Applications: summarization, diversity problems,...

(Martins et al., 2011b; Almeida and Martins, 2013; Martins et al., 2015)

GMs, $\psi \geq 0$
(dense, structured,
hard-constraint factors)



tractable problems

problems solvable with
dynamic programming

non-projective
dependency
grammars

projective
dependency
grammars

PCFGs

Approximate Decoding

What to do when exact decoding is intractable?

- Sampling methods (MCMC, etc.)
- Mean field algorithms
- LP relaxations
- Message-passing
- Dual decomposition

We'll highlight connections between several of these methods.

Approximate Decoding

What to do when exact decoding is intractable?

- Sampling methods (MCMC, etc.)
- Mean field algorithms
- LP relaxations
- Message-passing
- Dual decomposition

We'll highlight connections between several of these methods.

Global/Local Decoding

“Local” denotes independent problems within the scope of each factor

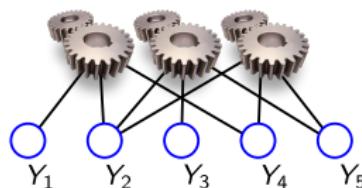
“Global” involves a global assignment of variables, consistent across factors

Key idea: “glue” the local evidence at the factors to obtain a global assignment

Our assumption: local decoding is easy, for every factor

We want to build a good (approximate) global decoder by invoking the local decoders.

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Outline

① Structured Prediction and Factor Graphs

② Integer Linear Programming

③ Message-Passing Algorithms

Sum-Product

Max-Product

④ Dual Decomposition

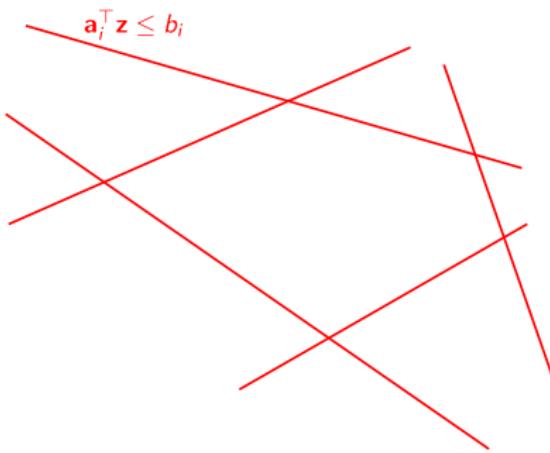
Linear Programming (Kantorovich, 1940; Dantzig, 1947)

$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N.$$

Linear objective

Linear constraints



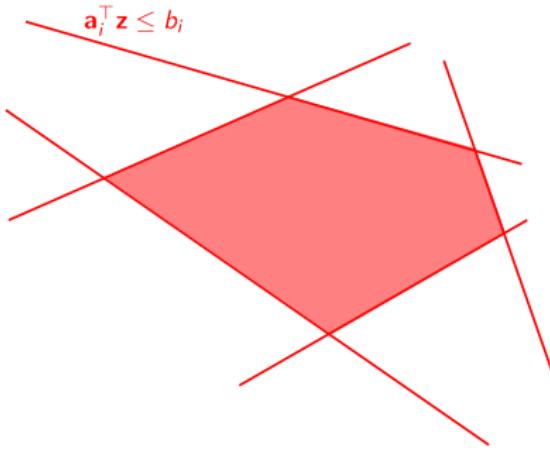
Linear Programming (Kantorovich, 1940; Dantzig, 1947)

$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N.$$

Linear objective

Linear constraints



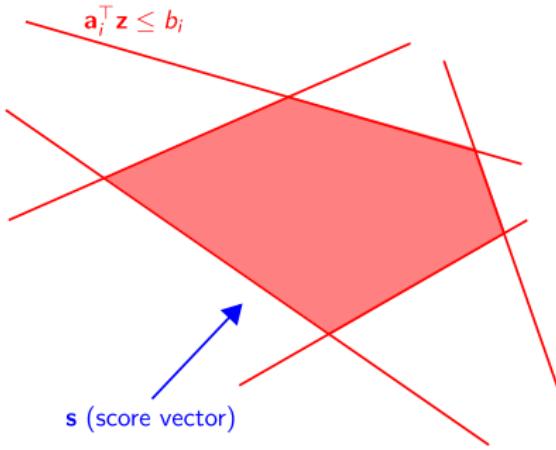
Linear Programming (Kantorovich, 1940; Dantzig, 1947)

$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N.$$

Linear objective

Linear constraints



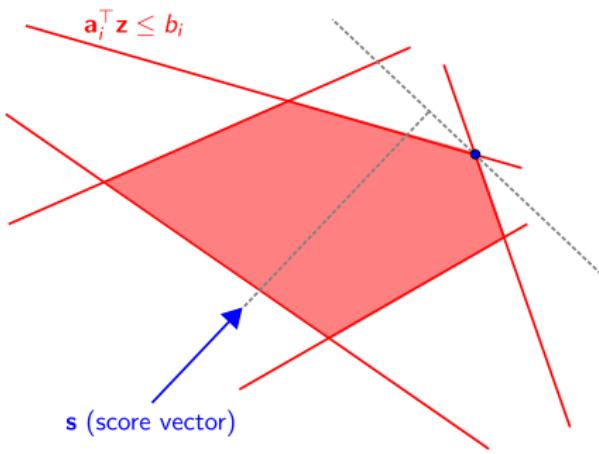
Linear Programming (Kantorovich, 1940; Dantzig, 1947)

$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N.$$

Linear objective

Linear constraints



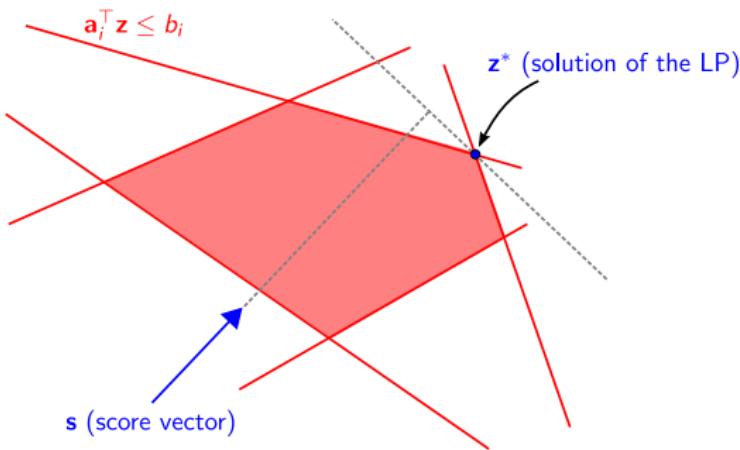
Linear Programming (Kantorovich, 1940; Dantzig, 1947)

$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N.$$

Linear objective

Linear constraints



Linear Programming (Kantorovich, 1940; Dantzig, 1947)

- If feasible and bounded, the solution is always attained at a vertex
- Can be solved in **polynomial time** (Khachiyan, 1980)
- Lots of off-the-shelf solvers (CPLEX, Gurobi, GLPK, LP_Solve, etc.)

Integer Linear Programming

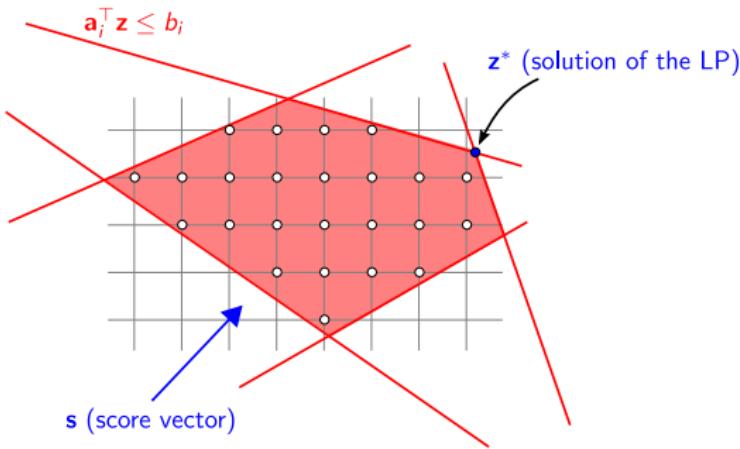
$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N,$$

z integer.

Linear objective

Linear constraints



Integer Linear Programming

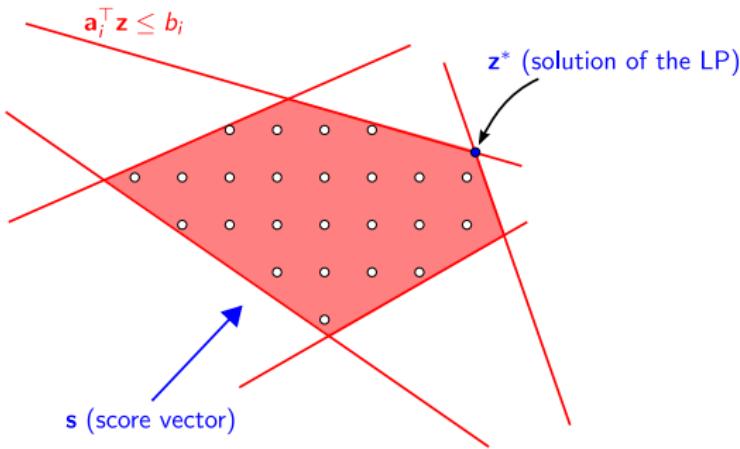
$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N,$$

z integer.

Linear objective

Linear constraints



Integer Linear Programming

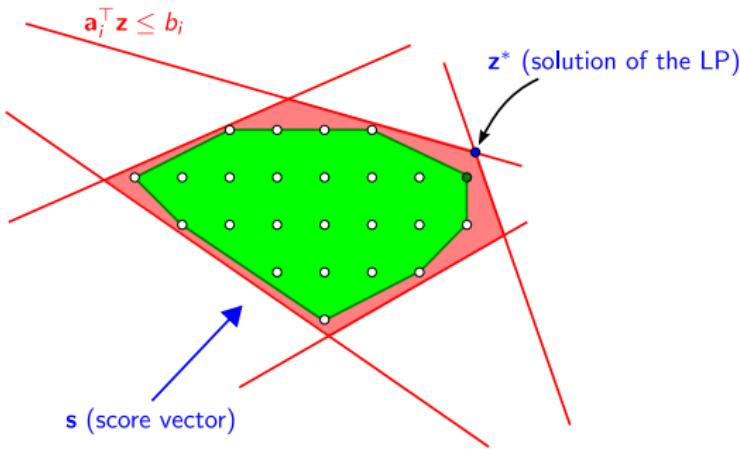
$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N,$$

z integer.

Linear objective

Linear constraints



Integer Linear Programming

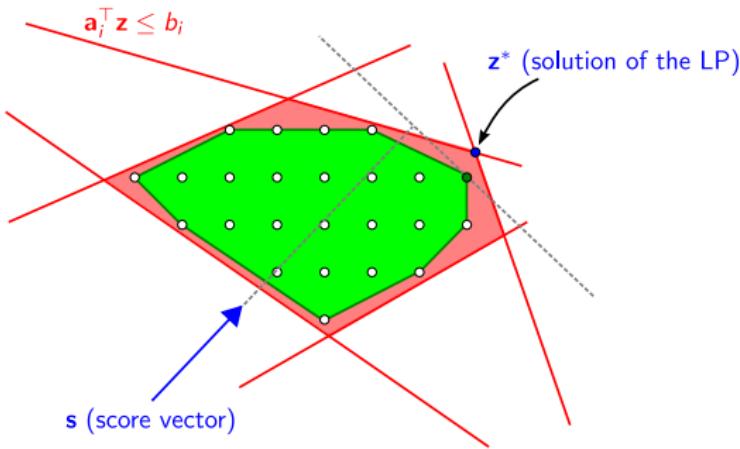
$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N,$$

z integer.

Linear objective

Linear constraints



Integer Linear Programming

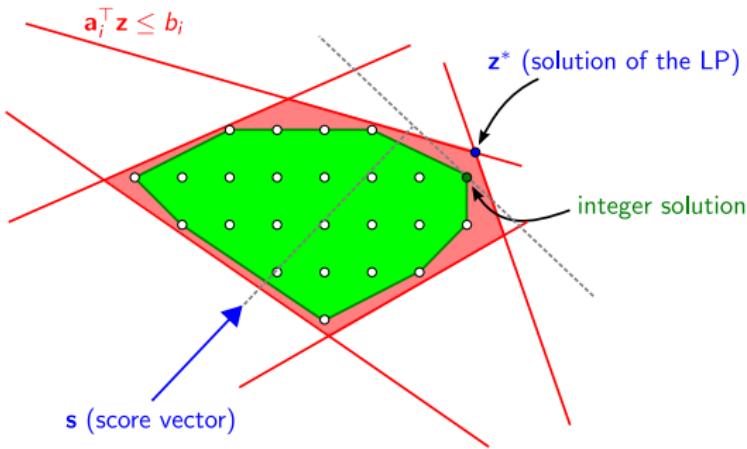
$$\max_z \quad s \cdot z$$

$$\text{s.t.} \quad a_i \cdot z \leq b_i, \quad i = 1, \dots, N,$$

z integer.

Linear objective

Linear constraints



Integer Linear Programming

In general, NP-hard (Karp, 1972)

Existing solvers are effective for small instances, but don't scale

LP relaxation: drops the integer constraints

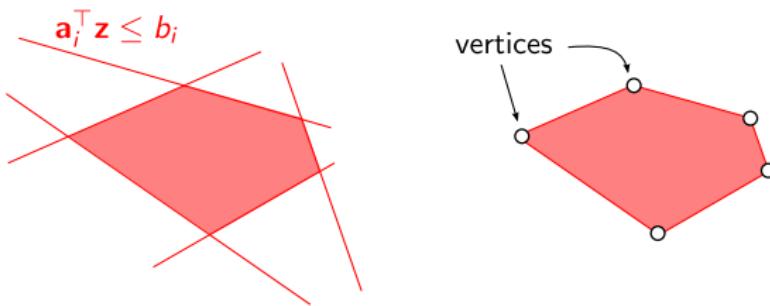
- Gives an upper bound of the solution of the ILP
- A common first step in exact algorithms (branch-and-bound, cutting plane, branch-and-cut)

Here's a very simple approximate algorithm:

- ① Solve the LP relaxation
- ② If the solution is integer, then it *is* the solution of the ILP
- ③ Otherwise, apply a rounding heuristic (problem-dependent)

Two Representations of Polytopes

Intersection of half-spaces (**H-representation**) or convex hull of a set of vertices (**V-representation**)

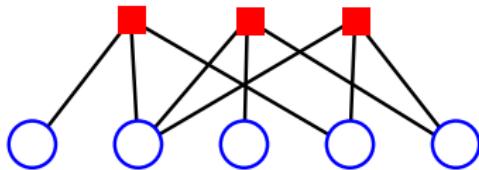


To call a solver, we need to specify a concise H-representation

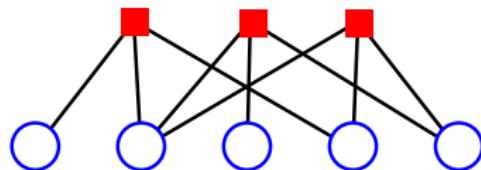
However, it may be difficult or impossible to obtain one if all we have is a V-representation

We next show how this relates to MAP decoding...

Structured Outputs as Bit-Vectors

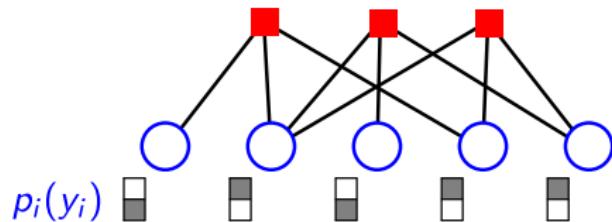


Structured Outputs as Bit-Vectors



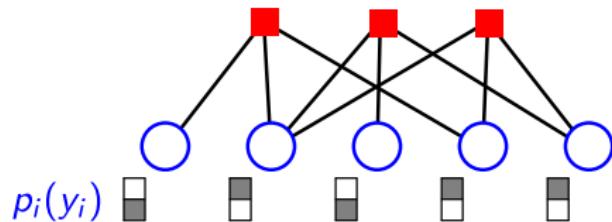
- One indicator $p_i(y_i)$ per each variable state

Structured Outputs as Bit-Vectors



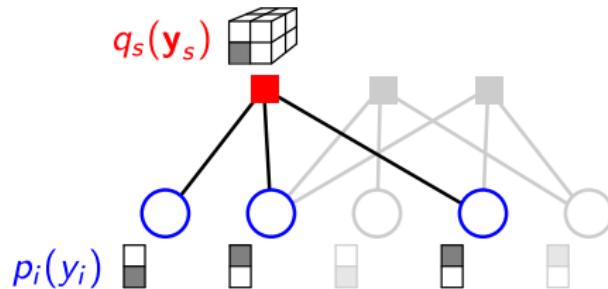
- One indicator $p_i(y_i)$ per each variable state

Structured Outputs as Bit-Vectors



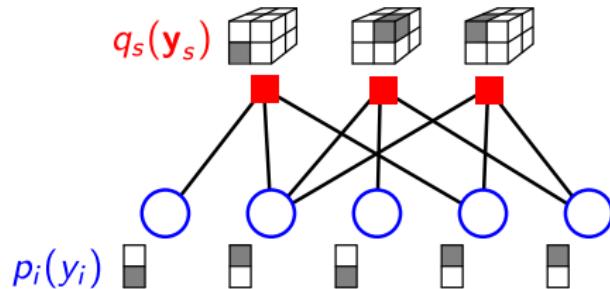
- One indicator $p_i(y_i)$ per each variable state
- One indicator $q_s(\mathbf{y}_s)$ per each factor configuration

Structured Outputs as Bit-Vectors



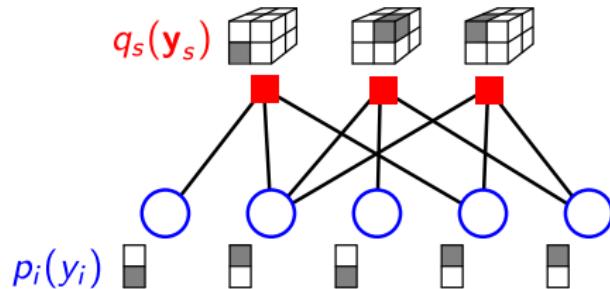
- One indicator $p_i(y_i)$ per each variable state
- One indicator $q_s(\mathbf{y}_s)$ per each factor configuration

Structured Outputs as Bit-Vectors



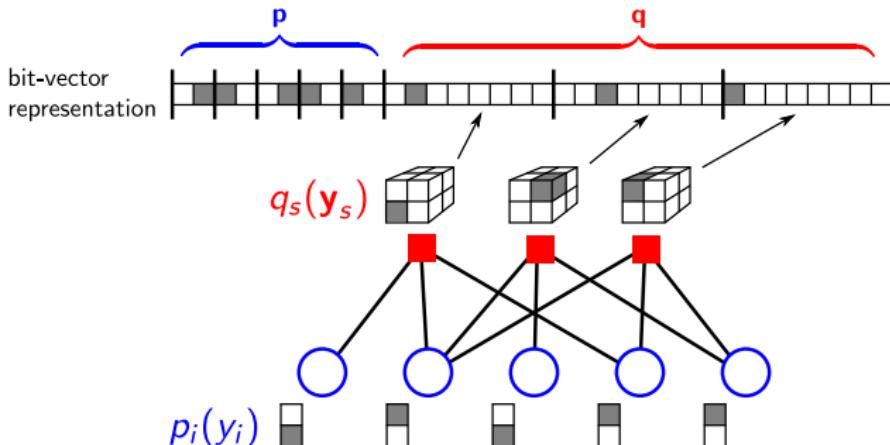
- One indicator $p_i(y_i)$ per each variable state
- One indicator $q_s(y_s)$ per each factor configuration

Structured Outputs as Bit-Vectors



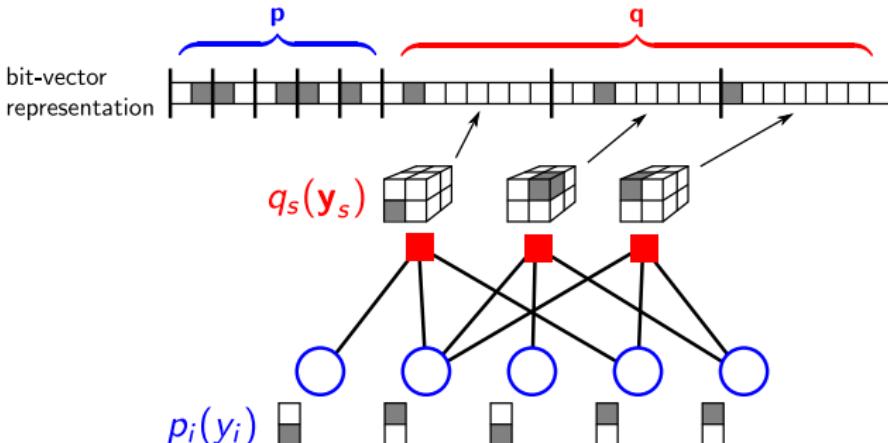
- One indicator $p_i(y_i)$ per each variable state
- One indicator $q_s(\mathbf{y}_s)$ per each factor configuration
- Overall: each global output $y \in \mathcal{Y}(x)$ is mapped to a bit-vector

Structured Outputs as Bit-Vectors



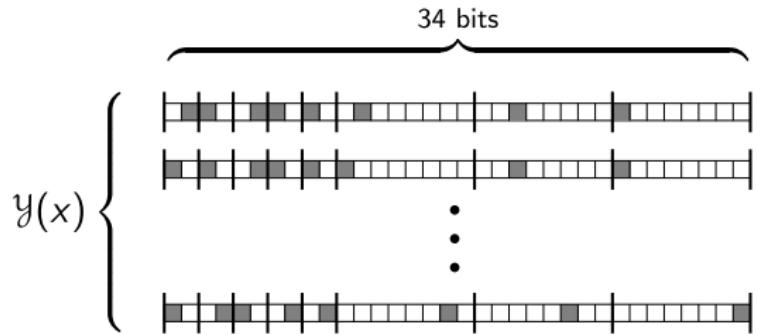
- One indicator $p_i(y_i)$ per each variable state
- One indicator $q_s(y_s)$ per each factor configuration
- Overall: each global output $y \in \mathcal{Y}(x)$ is mapped to a bit-vector

Structured Outputs as Bit-Vectors

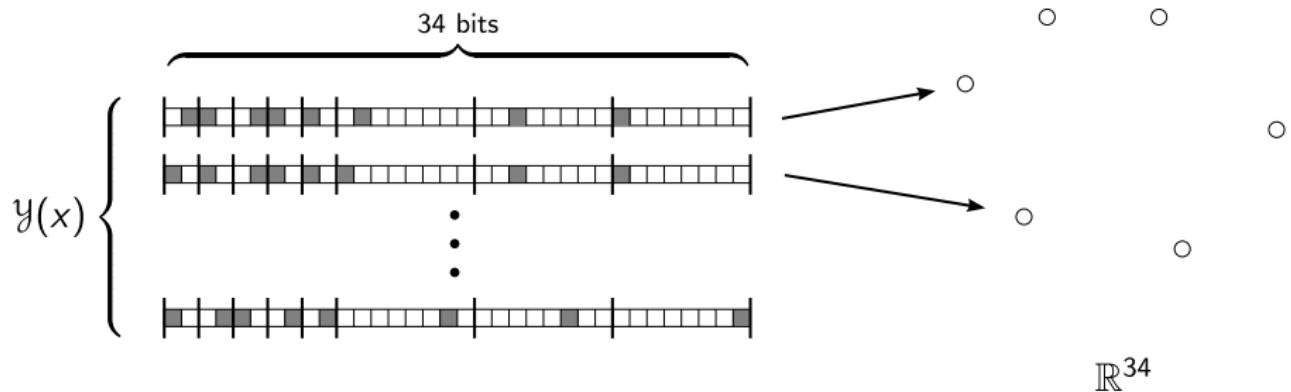


- One indicator $p_i(y_i)$ per each variable state
- One indicator $q_s(y_s)$ per each factor configuration
- Overall: each global output $y \in \mathcal{Y}(x)$ is mapped to a bit-vector
- **Note: not all bit vectors are valid (they must be consistent)**

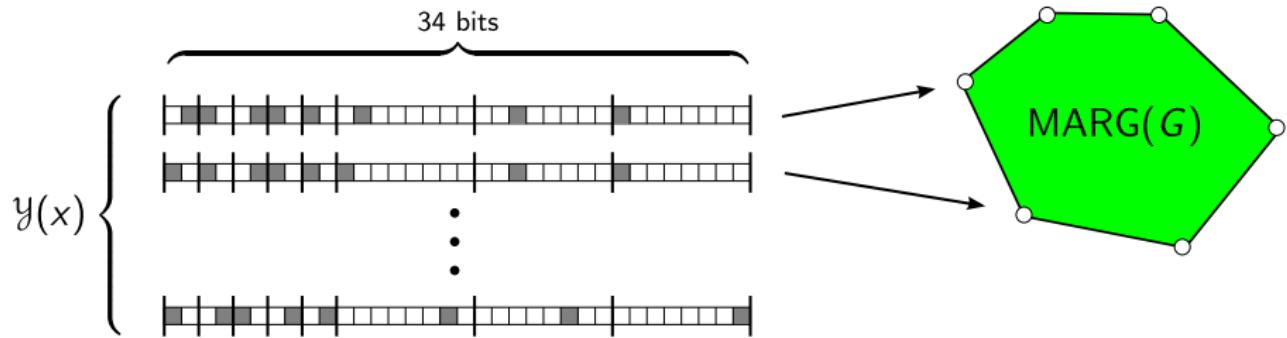
Marginal Polytope (Wainwright and Jordan, 2008)



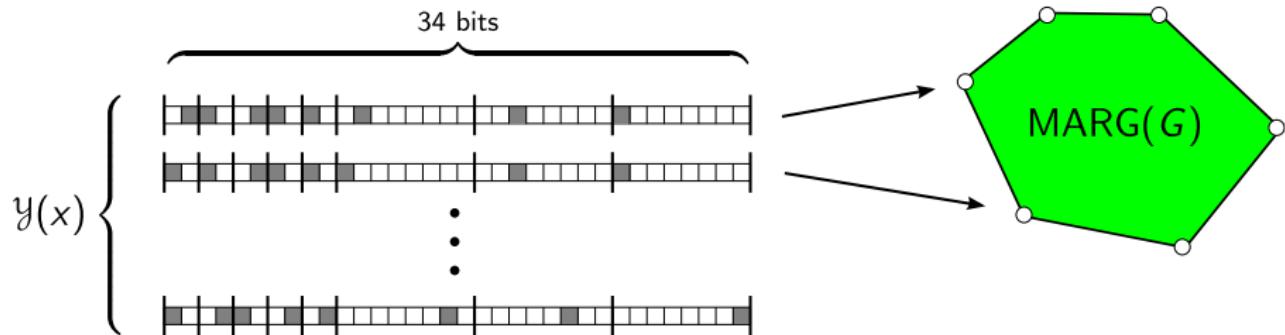
Marginal Polytope (Wainwright and Jordan, 2008)



Marginal Polytope (Wainwright and Jordan, 2008)

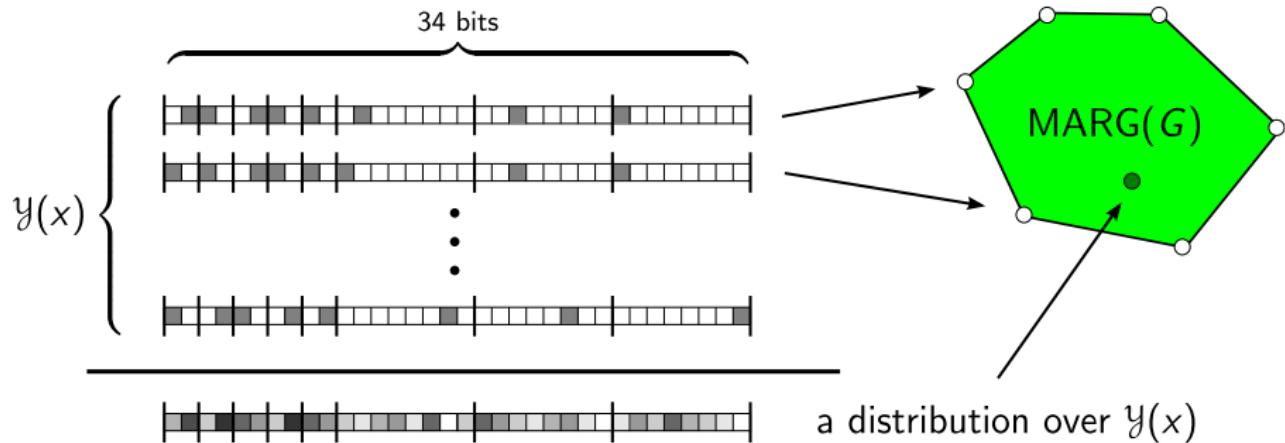


Marginal Polytope (Wainwright and Jordan, 2008)



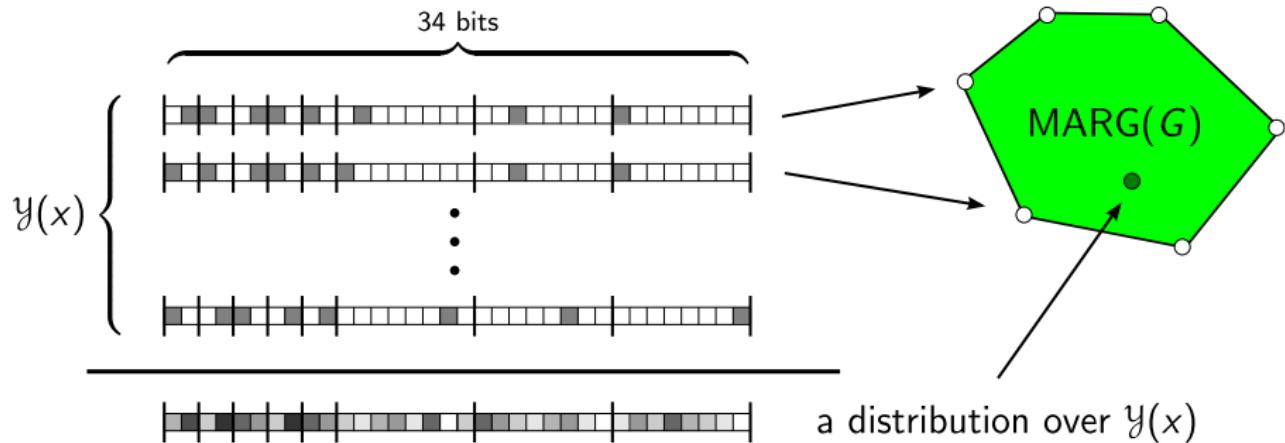
- Vertices of $\text{MARG}(G)$ correspond to outputs $y(x)$
- Points of $\text{MARG}(G)$ correspond to realizable marginals (more later)

Marginal Polytope (Wainwright and Jordan, 2008)



- Vertices of $\text{MARG}(G)$ correspond to outputs $\gamma(x)$
- Points of $\text{MARG}(G)$ correspond to realizable marginals (more later)

Marginal Polytope (Wainwright and Jordan, 2008)



- Vertices of $\text{MARG}(G)$ correspond to outputs $\gamma(x)$
- Points of $\text{MARG}(G)$ correspond to realizable marginals (more later)
- This is a V-representation, what about an H-representation?

H-Representation With Integer Constraints

In general, there's no concise H-representation for MARG(G)

... but we can represent its vertices if *integer constraints* are permitted:

$$\sum_{\mathbf{y}_s} q_s(\mathbf{y}_s) = 1, \quad q_s(\mathbf{y}_s) \geq \mathbf{0}, \quad \forall \mathbf{y}_s \in \mathcal{Y}_s \quad (\text{normalization})$$

$$p_i(y_i) = \sum_{\mathbf{y}_s \sim y_i} q_s(\mathbf{y}_s), \quad \forall i \in N(s) \quad (\text{marginalization})$$

q is integer **(integer constraints)**

This will open the door for formulating MAP decoding as an ILP.

MAP Decoding as an ILP

Recall the MAP decoding problem:

$$\begin{aligned}\hat{y} &= \arg \max_{y \in \mathcal{Y}(x)} P_\psi(y|x) \\ &= \arg \max_{y \in \mathcal{Y}(x)} \frac{1}{Z(\psi, x)} \prod_i \psi_i(y_i) \prod_s \psi_s(y_s) \\ &= \arg \max_{y \in \mathcal{Y}(x)} \sum_i \theta_i(y_i) + \sum_s \theta_s(y_s),\end{aligned}$$

where $\theta_i(y_i) := \log \psi_i(y_i)$ and $\theta_s(y_s) := \log \psi_s(y_s)$

We can rewrite this as an ILP:

$$\text{maximize } \sum_i \sum_{y_i} \theta_i(y_i) p_i(y_i) + \sum_s \sum_{y_s} \theta_s(y_s) q_s(y_s)$$

subject to $(p, q) \in \text{MARG}(G)$

Local Polytope

Obtained by relaxing the integer constraints

Regard p_i and q_s as probability distributions that must be **locally consistent**:

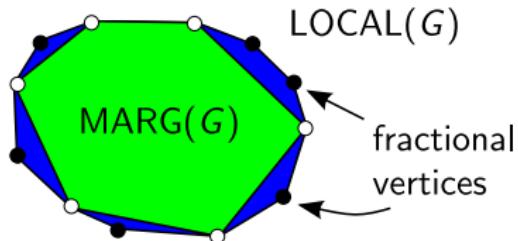
$$\sum_{y_s} q_s(y_s) = 1, \quad q_s(y_s) \geq 0, \quad \forall y_s \in \mathcal{Y}_s \quad (\text{normalization})$$

$$p_i(y_i) = \sum_{y_s \sim y_i} q_s(y_s), \quad \forall i \in N(s) \quad (\text{marginalization})$$

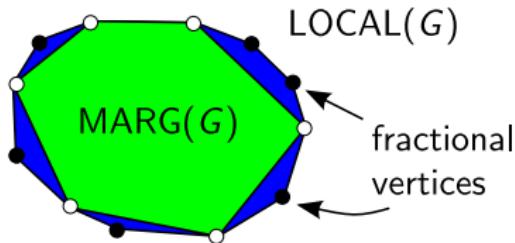
~~q is integer~~ **(integer constraints)**

The feasible points are *pseudo-marginals* (not necessarily valid marginals)

Local and Marginal Polytopes



Local and Marginal Polytopes



- $\text{LOCAL}(G)$ is an **outer bound** of $\text{MARG}(G)$
- It contains all the integer vertices of $\text{MARG}(G)$, plus spurious **fractional vertices**
- If the graph has no cycles, then $\text{LOCAL}(G) = \text{MARG}(G)$

LP-MAP Decoding

Solves a linear relaxation of MAP decoding, replacing $\text{MARG}(G)$ by $\text{LOCAL}(G)$:

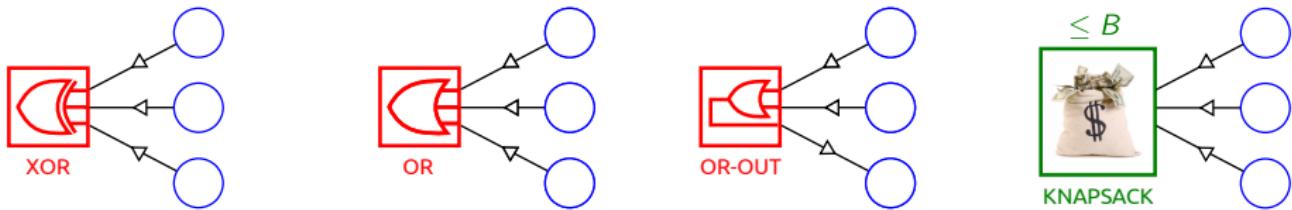
$$\text{maximize} \sum_i \sum_{y_i} \theta_i(y_i) p_i(y_i) + \sum_s \sum_{\mathbf{y}_s} \theta_s(\mathbf{y}_s) q_s(\mathbf{y}_s)$$

subject to $(p, q) \in \text{LOCAL}(G)$

If the solution is integer, we solved the problem exactly; otherwise, apply a rounding heuristic

Runtime is polynomial, but the procedure is only approximate.

What About Hard Constraint Factors?



Logic and knapsack/budget constraints can also be expressed *linearly*

Logic/Budget Constraints

Assume $z_1, z_2, \dots \in \{0, 1\}$ (binary variables)

Condition	Statement	Constraint
Implication	if z_1 then z_2	$z_1 \leq z_2$
Negation	z_1 iff $\neg z_2$	$z_1 = 1 - z_2$
OR	z_1 or z_2 or z_3	$z_1 + z_2 + z_3 \geq 1$
XOR	z_1 xor z_2 xor z_3	$z_1 + z_2 + z_3 = 1$
OR-OUT	$z_{12} = z_1 \vee z_2$	$z_{12} \geq z_1, z_{12} \geq z_2,$ $z_{12} \leq z_1 + z_2$
AND-OUT	$z_{12} = z_1 \wedge z_2$	$z_{12} \leq z_1, z_{12} \leq z_2,$ $z_{12} \geq z_1 + z_2 - 1$
Budget	at most B active units	$\sum_i z_i \leq B$
Knapsack	at most B total weight	$\sum_i w_i z_i \leq B$

More complex expressions via composition and De Morgan's laws

Summing Up ILPs

- MAP decoding can be expressed as an Integer Linear Program (ILP)
- Logic constraints can be incorporated easily
- Structured factors are harder (they need to be disassembled)
- The ILP can be relaxed for approximate decoding (LP-MAP)
- Geometrically: an outer bound of the *marginal polytope*
- The relaxation is tight if the graph G does not have cycles
- **Disadvantage: an off-the-shelf LP solver won't exploit the modularity of the problem**
- **Algorithms that exploit the structure of the LP will be the topic of the remaining sections**

Outline

① Structured Prediction and Factor Graphs

② Integer Linear Programming

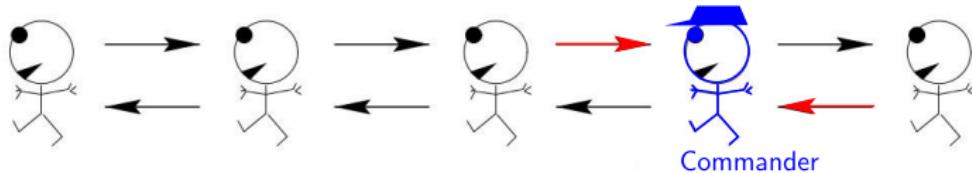
③ Message-Passing Algorithms

Sum-Product

Max-Product

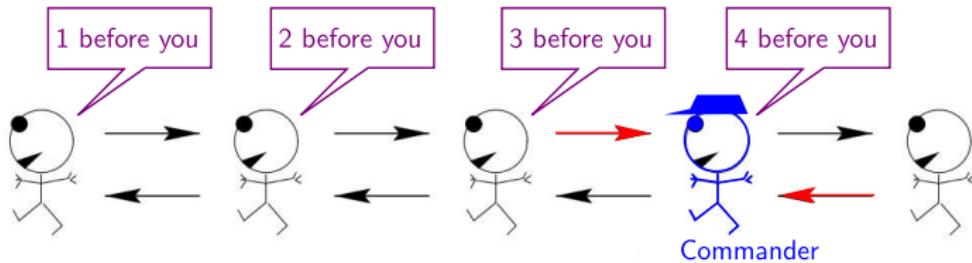
④ Dual Decomposition

Motivating Example: Counting Soldiers



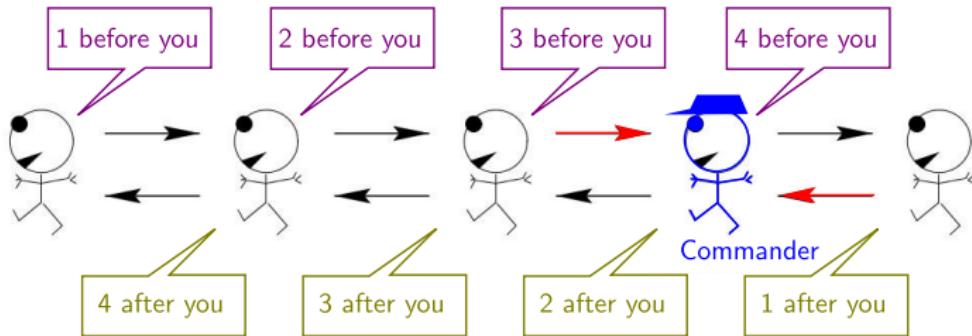
(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

Motivating Example: Counting Soldiers



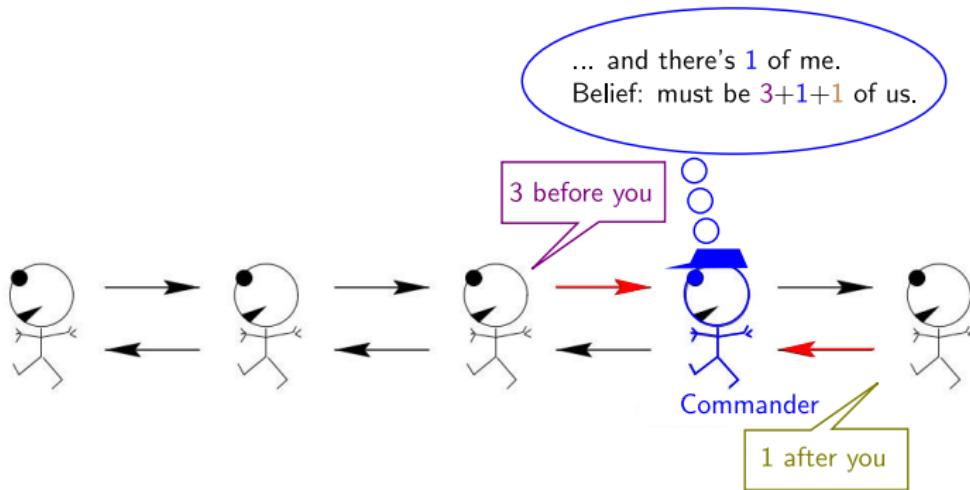
(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

Motivating Example: Counting Soldiers



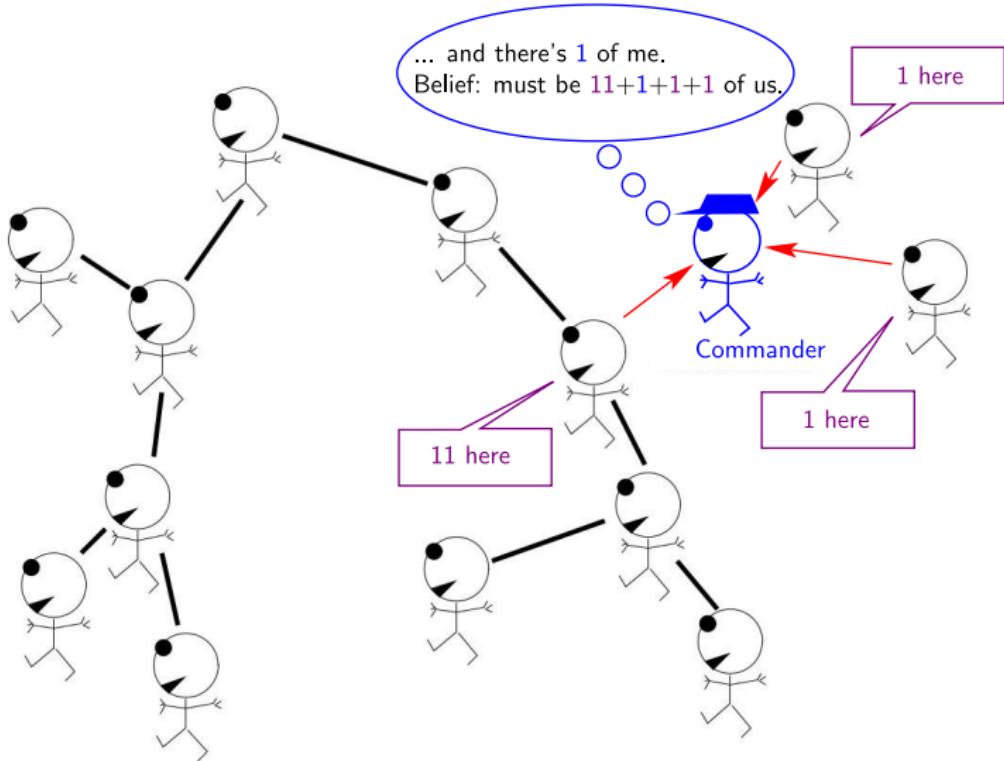
(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

Motivating Example: Counting Soldiers



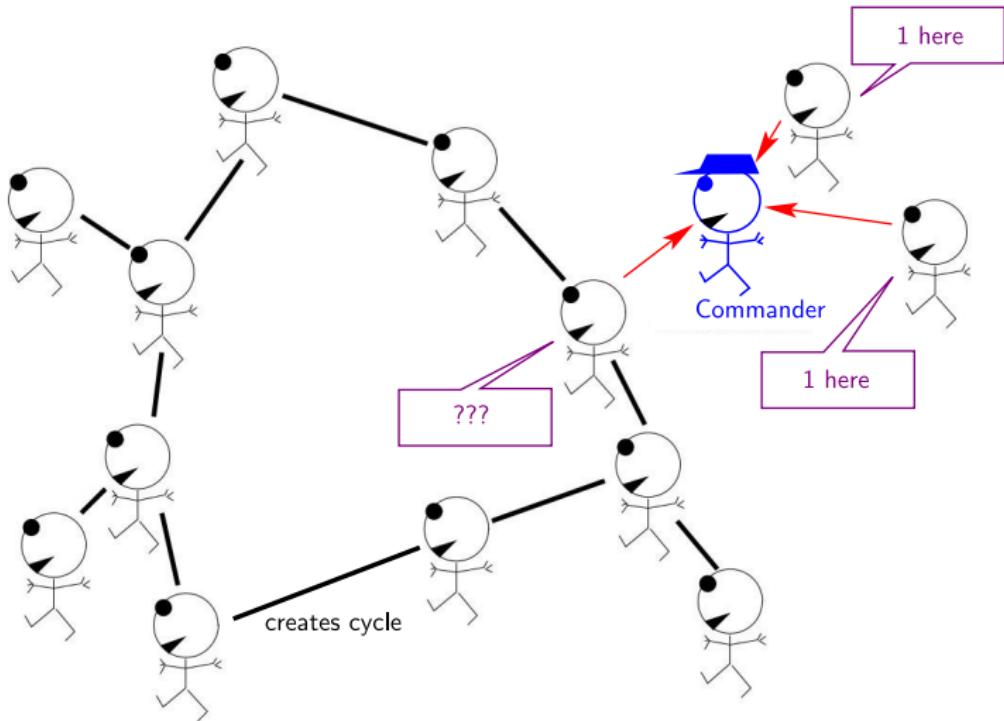
(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

Motivating Example: Counting Soldiers



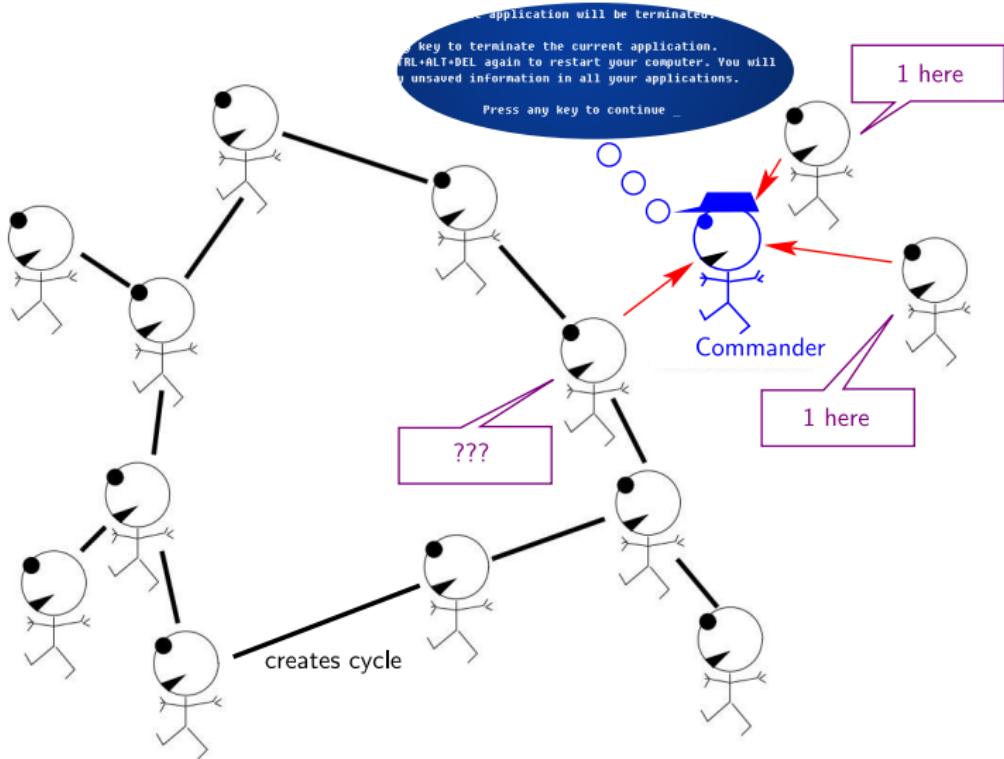
(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

Motivating Example: Counting Soldiers



(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

Motivating Example: Counting Soldiers



(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

Outline

① Structured Prediction and Factor Graphs

② Integer Linear Programming

③ Message-Passing Algorithms

Sum-Product

Max-Product

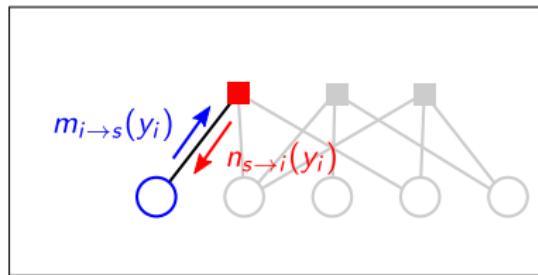
④ Dual Decomposition

Sum-Product Belief Propagation

Recall that $\mathbb{P}_\psi(y|x) := \frac{1}{Z(\psi,x)} \times \prod_i \psi_i(y_i) \times \prod_s \psi_s(\mathbf{y}_s)$

Alternate between computing two kinds of messages:

- **Variable-to-factor:** $m_{i \rightarrow s}(y_i) = \psi_i(y_i) \prod_{s' \in N(i) \setminus \{s\}} n_{s' \rightarrow i}(y_i)$
- **Factor-to-variable:** $n_{s \rightarrow i}(y_i) = \sum_{\mathbf{y}_s \sim y_i} \psi_s(\mathbf{y}_s) \prod_{j \in N(s) \setminus \{i\}} m_{j \rightarrow s}(y_j)$



Beliefs

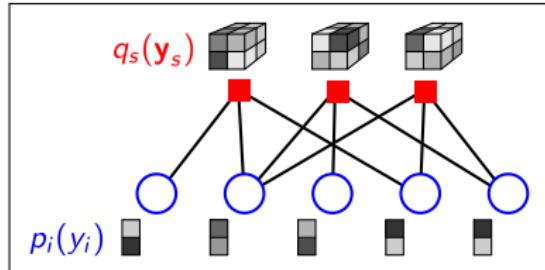
Given the messages, we compute local *beliefs*:

- **Variable beliefs:**

$$p_i(y_i) \propto \psi_i(y_i) \prod_{s \in N(i)} n_{s \rightarrow i}(y_i)$$

- **Factor beliefs:**

$$q_s(\mathbf{y}_s) \propto \psi_s(\mathbf{y}_s) \prod_{i \in N(s)} m_{i \rightarrow s}(y_i)$$



If the graph has no cycles, these beliefs converge to the true marginals

$$p_i(y_i) \rightarrow \mathbb{P}_{\psi}(y_i|x), \quad q_s(\mathbf{y}_s) \rightarrow \mathbb{P}_{\psi}(\mathbf{y}_s|x)$$

Otherwise: loopy BP (later)

Belief Propagation as Calibration

- Variable-to-factor messages:

$$m_{i \rightarrow s}(y_i) = \psi_i(y_i) \prod_{s' \in N(i) \setminus \{s\}} n_{s' \rightarrow i}(y_i) = \frac{p_i(y_i)}{n_{s \rightarrow i}(y_i)}$$

- Factor-to-variable messages:

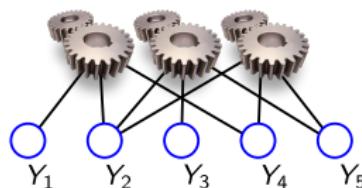
$$n_{s \rightarrow i}(y_i) = \sum_{\mathbf{y}_s \sim y_i} \psi_s(\mathbf{y}_s) \prod_{j \in N(s) \setminus \{i\}} m_{j \rightarrow s}(y_j) = \frac{\sum_{\mathbf{y}_s \sim y_i} q_s(\mathbf{y}_s)}{m_{i \rightarrow s}(y_i)}$$

- Calibration equations (attained at convergence):

$$p_i(y_i) = \sum_{\mathbf{y}_s \sim y_i} q_s(\mathbf{y}_s)$$

Punchline: to run sum-product BP, we only need local marginals

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Loopy Belief Propagation

What if the graph has cycles?

Loopy Belief Propagation

What if the graph has cycles?

We'll see that marginal decoding corresponds to optimizing a *free energy objective* over the *marginal polytope*

Sum-product “loopy” BP entails two approximations:

- ① Replaces $\text{MARG}(G)$ by $\text{LOCAL}(G)$
- ② Optimizes a *Bethe free energy* approximation

Step #1: Dual Parametrization

For any ψ , there are marginals \mathbf{p} , \mathbf{q} in $\text{MARG}(G)$ that parametrize \mathbb{P}_ψ

E.g. if the graph has no cycles:

$$\begin{aligned}\mathbb{P}_\psi(y|x) &= \frac{1}{Z(\psi, x)} \prod_i \psi_i(y_i) \times \prod_s \psi_s(\mathbf{y}_s) \\ &= \prod_i \mathbf{p}_i(y_i)^{1-|N(i)|} \times \prod_s \mathbf{q}_s(\mathbf{y}_s) \quad (* \text{ next slide}) \\ &:= \mathbb{P}_{\mathbf{p}, \mathbf{q}}(y|x)\end{aligned}$$

Therefore: a distribution can be represented as a point in $\text{MARG}(G)$

$\theta := \log(\psi)$ are called *canonical parameters*, and (\mathbf{p}, \mathbf{q}) *mean parameters*

(*) Derivation of Dual Parametrization

Assume a tree-shaped Bayes net (each variable i has a single parent π_i)

$$\begin{aligned}\mathbb{P}(y) &= \mathbb{P}(y_0) \prod_{i \neq 0} \mathbb{P}(y_i | y_{\pi_i}) \\&= \mathbb{P}(y_0) \prod_{i \neq 0} \frac{\mathbb{P}(y_i, y_{\pi_i})}{\mathbb{P}(y_{\pi_i})} \\&= \frac{\mathbb{P}(y_0) \prod_s \mathbb{P}(\mathbf{y}_s)}{\prod_j \mathbb{P}(y_j)^{|i:j=\pi_i|}} \\&= \frac{\mathbb{P}(y_0) \prod_s \mathbb{P}(\mathbf{y}_s)}{\mathbb{P}(y_0)^{|N(0)|} \prod_{j \neq 0} \mathbb{P}(y_j)^{|N(j)-1|}} \\&= \frac{\prod_s \mathbb{P}(\mathbf{y}_s)}{\prod_j \mathbb{P}(y_j)^{|N(j)|-1}} \\&= \prod_i p_i(y_i)^{1-|N(i)|} \times \prod_s q_s(\mathbf{y}_s).\end{aligned}$$

Step #2: Entropy and Log-Partition Function

Entropy of a probability distribution: $H(\mathbb{P}) = -\sum_y \mathbb{P}(y) \log \mathbb{P}(y)$

Definition: the Fenchel dual of a convex function $f : \mathbb{R}^D \rightarrow \mathbb{R} \cup \{+\infty\}$ is the convex function $f^* : \mathbb{R}^D \rightarrow \mathbb{R} \cup \{+\infty\}$ defined pointwise as
$$f^*(\mathbf{v}) := \sup_{\mathbf{u}} (\mathbf{v} \cdot \mathbf{u} - f(\mathbf{u}))$$

Step #2: Entropy and Log-Partition Function

Entropy of a probability distribution: $H(\mathbb{P}) = -\sum_y \mathbb{P}(y) \log \mathbb{P}(y)$

Definition: the Fenchel dual of a convex function $f : \mathbb{R}^D \rightarrow \mathbb{R} \cup \{+\infty\}$ is the convex function $f^* : \mathbb{R}^D \rightarrow \mathbb{R} \cup \{+\infty\}$ defined pointwise as
 $f^*(\mathbf{v}) := \sup_{\mathbf{u}} (\mathbf{v} \cdot \mathbf{u} - f(\mathbf{u}))$

Theorem (I): the **log-partition function** $\log Z(\theta)$ and the **negative entropy** $-H(\mathbb{P}_{\mathbf{p}, \mathbf{q}})$ are Fenchel dual:

$$\log Z(\theta) = \max_{(\mathbf{p}, \mathbf{q}) \in \text{MARG}(G)} \underbrace{\sum_i \theta_i \cdot \mathbf{p}_i + \sum_s \theta_s \cdot \mathbf{q}_s}_{\text{(negative) variational free energy}} + H(\mathbb{P}_{\mathbf{p}, \mathbf{q}}),$$

This underlies the well-known duality between maximum likelihood in log-linear models and maximum entropy.

Step #3: Loopy BP as Variational Inference

Theorem (II): The maximizers $(\mathbf{p}^*, \mathbf{q}^*)$ are the **true marginals** of \mathbb{P}_θ :

$$(\mathbf{p}^*, \mathbf{q}^*) = \arg \max_{(\mathbf{p}, \mathbf{q}) \in \text{MARG}(G)} \sum_i \theta_i \cdot \mathbf{p}_i + \sum_s \theta_s \cdot \mathbf{q}_s + H(\mathbb{P}_{\mathbf{p}, \mathbf{q}})$$

Step #3: Loopy BP as Variational Inference

Theorem (II): The maximizers $(\mathbf{p}^*, \mathbf{q}^*)$ are the **true marginals** of \mathbb{P}_θ :

$$(\mathbf{p}^*, \mathbf{q}^*) = \arg \max_{(\mathbf{p}, \mathbf{q}) \in \text{MARG}(G)} \sum_i \theta_i \cdot \mathbf{p}_i + \sum_s \theta_s \cdot \mathbf{q}_s + H(\mathbb{P}_{\mathbf{p}, \mathbf{q}})$$

Drawback: in general, $\text{MARG}(G)$ and $H(\mathbb{P}_{\mathbf{p}, \mathbf{q}})$ are both intractable

Step #3: Loopy BP as Variational Inference

Theorem (II): The maximizers $(\mathbf{p}^*, \mathbf{q}^*)$ are the **true marginals** of \mathbb{P}_{θ} :

$$(\mathbf{p}^*, \mathbf{q}^*) = \arg \max_{(\mathbf{p}, \mathbf{q}) \in \text{MARG}(G)} \sum_i \theta_i \cdot \mathbf{p}_i + \sum_s \theta_s \cdot \mathbf{q}_s + H(\mathbb{P}_{\mathbf{p}, \mathbf{q}})$$

Drawback: in general, $\text{MARG}(G)$ and $H(\mathbb{P}_{\mathbf{p}, \mathbf{q}})$ are both intractable

Yedidia et al. (2001) showed that loopy BP entails two approximations:

- ① Replace $\text{MARG}(G)$ by $\text{LOCAL}(G)$
- ② Approximate $H(\mathbb{P}_{\mathbf{p}, \mathbf{q}})$ by the Bethe entropy $H_{\text{Bethe}}(\mathbb{P}_{\mathbf{p}, \mathbf{q}})$

Both are exact when the graph does not have cycles

Bethe Entropy Approximation

Derived by “pretending” the graph has no cycles

We have seen

$$\mathbb{P}_\psi(y|x) \approx \prod_i p_i(y_i)^{1-|N(i)|} \times \prod_s q_s(y_s)$$

From which we can derive

$$\begin{aligned} H(\mathbb{P}_{\mathbf{p}, \mathbf{q}}) &\approx H_{\text{Bethe}}(\mathbb{P}_{\mathbf{p}, \mathbf{q}}) \\ &= \sum_i (1 - |N(i)|) H_i(\mathbf{p}_i) + \sum_s H_s(\mathbf{q}_s) \end{aligned}$$



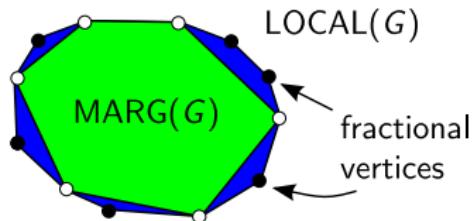
Hans Bethe, 1906–2005

A linear combination of local entropies:

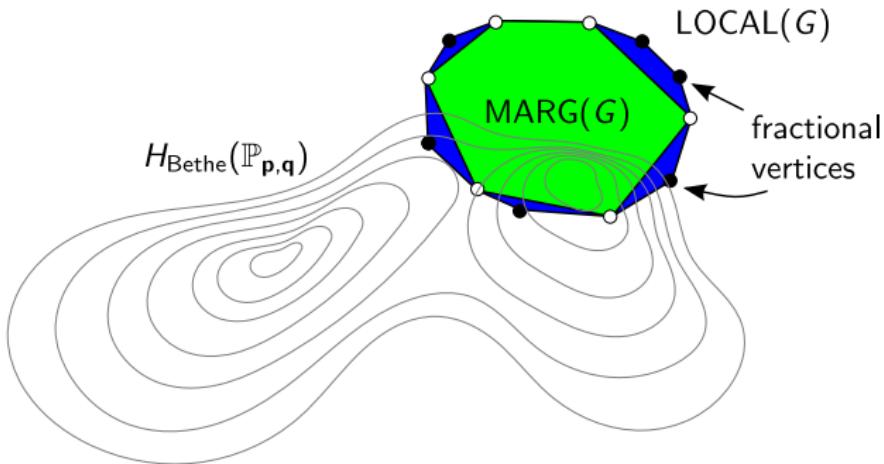
$$H_i(\mathbf{p}_i) = - \sum_{y_i} p_i(y_i) \log p_i(y_i), \quad H_s(\mathbf{q}_s) = - \sum_{y_s} q_s(y_s) \log q_s(y_s)$$

Not concave in general!

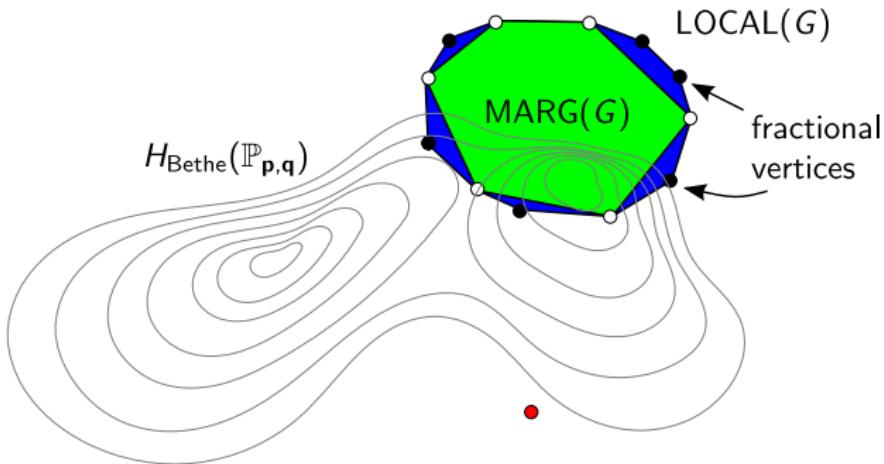
Geometric Illustration



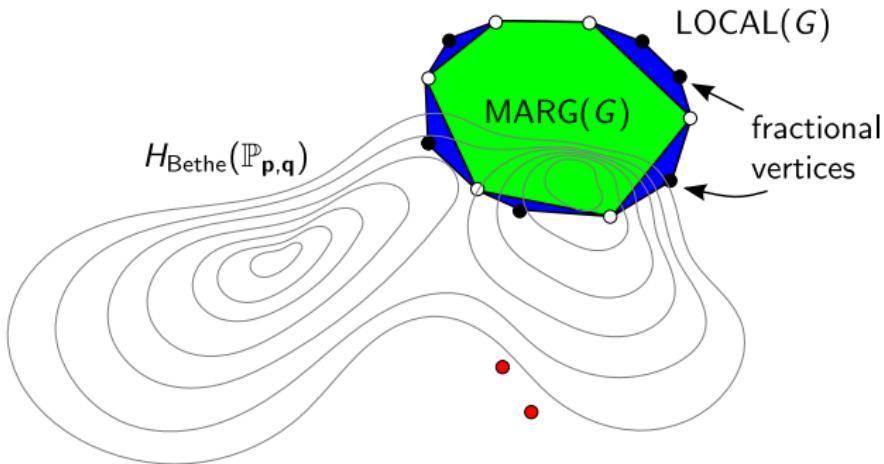
Geometric Illustration



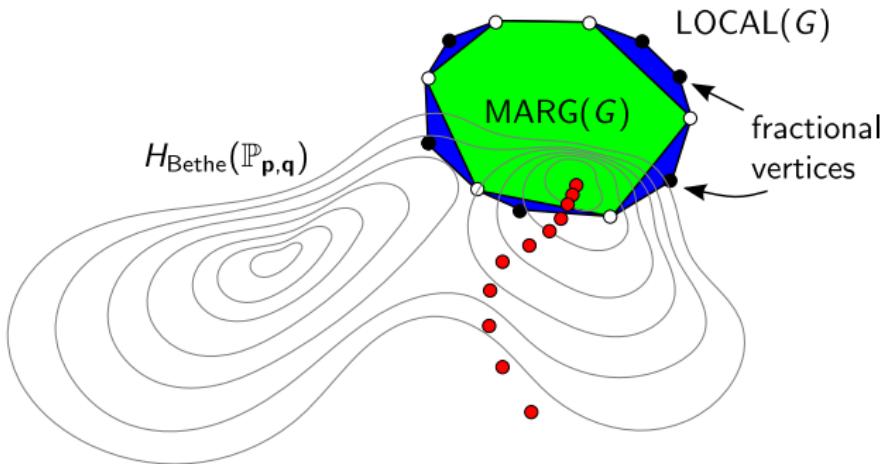
Geometric Illustration



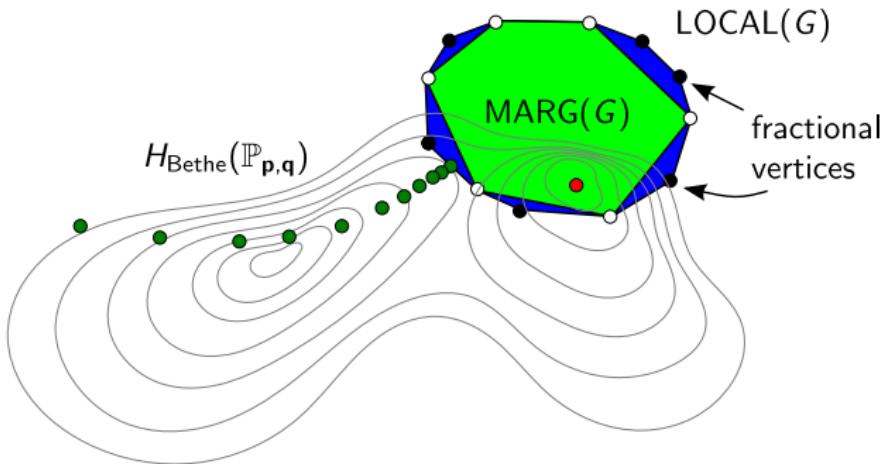
Geometric Illustration



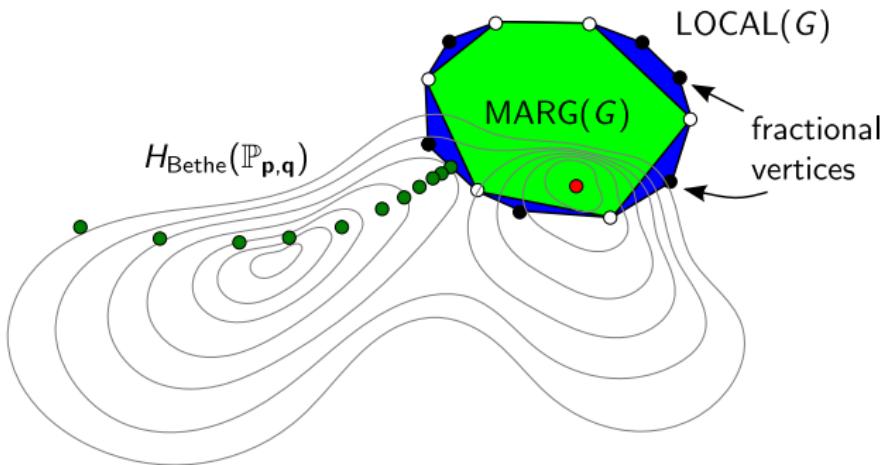
Geometric Illustration



Geometric Illustration



Geometric Illustration



If loopy BP converges, it reaches a stationary point of the approximate variational problem

$H_{Bethe}(\mathbb{P}_{p,q})$ is non-concave in general \Rightarrow local minima

Summary of Loopy BP

Advantages:

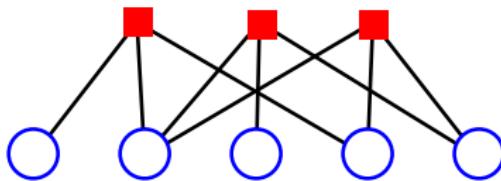
- Simple to implement
- Handles structured and logic factors (only need to compute local marginals)
- Often works well in practice (if cycles are not very influential)
- Often yields a reasonable approximation of $\log Z$ and H

Disadvantages:

- Doesn't give an upper/lower bound of $\log Z$ and H
- Entropy approximation is not concave (local minima)
- May not converge
- The final beliefs may not be realizable marginals

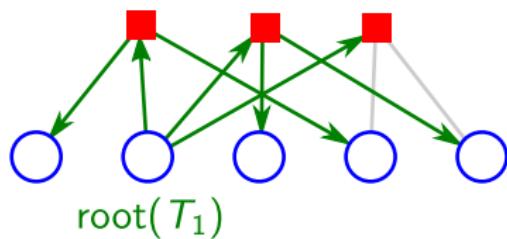
Tree Reweighted BP (Wainwright et al., 2005)

Key idea: cover the graph with a set of trees



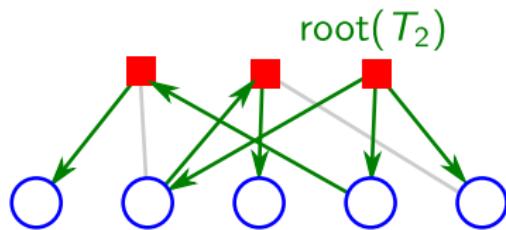
Tree Reweighted BP (Wainwright et al., 2005)

Key idea: cover the graph with a set of trees



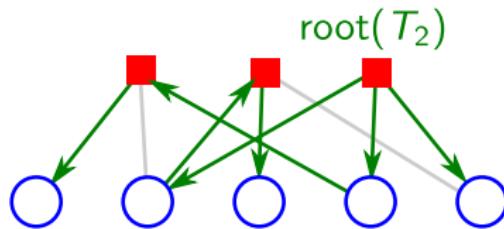
Tree Reweighted BP (Wainwright et al., 2005)

Key idea: cover the graph with a set of trees



Tree Reweighted BP (Wainwright et al., 2005)

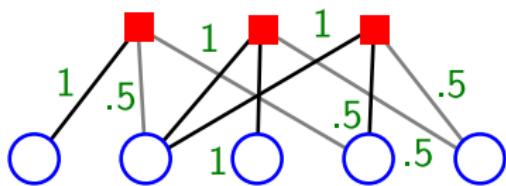
Key idea: cover the graph with a set of trees



Count the appearance probability $c_{is} > 0$ of each edge

Tree Reweighted BP (Wainwright et al., 2005)

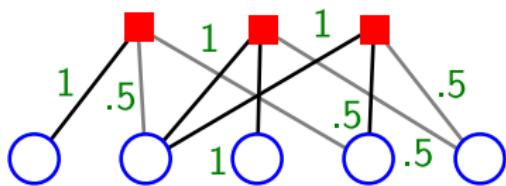
Key idea: cover the graph with a set of trees



Count the appearance probability $c_{is} > 0$ of each edge

Tree Reweighted BP (Wainwright et al., 2005)

Key idea: cover the graph with a set of trees



Count the appearance probability $c_{is} > 0$ of each edge

This results in a convex upper bound of $-H$ and $\log Z$:

$$H_{\text{TRBP}}(\mathbb{P}_{\mathbf{p}, \mathbf{q}}) = \sum_i (1 - \sum_{s \in N(i)} c_{is}) H_i(\mathbf{p}_i) + \sum_s H_s(\mathbf{q}_s)$$

(Note: if all $c_{is} = 1$ this would revert to the Bethe approximation)

TRBP Messages

- Variable-to-factor messages:

$$m_{i \rightarrow s}(y_i) = \frac{\psi_i(y_i) \prod_{s' \in N(i)} n_{s' \rightarrow i}^{c_{is'}}(y_i)}{n_{s \rightarrow i}(y_i)}$$

- Factor-to-variable messages:

$$n_{s \rightarrow i}(y_i) = \sum_{y_s \sim y_i} \frac{\psi_s(y_s) \prod_{j \in N(s)} m_{j \rightarrow s}^{c_{js}}(y_j)}{m_{i \rightarrow s}(y_i)}$$

- Variable beliefs:

$$p_i(y_i) \propto \psi_i(y_i) \prod_{s \in N(i)} n_{s \rightarrow i}^{c_{is}}(y_i)$$

- Factor beliefs:

$$q_s(y_s) \propto \psi_s(y_s) \prod_{i \in N(s)} m_{i \rightarrow s}^{c_{is}}(y_i)$$

Summary of TRBP

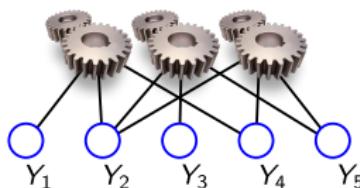
Advantages:

- Still simple to implement
- Entropy approximation is concave (no local minima)
- Gives an upper bound on $-H$ and $\log Z$
- Lots of knobs (the appearance probabilities)

Disadvantages:

- Lots of knobs (the appearance probabilities)
- Typically it's a very loose bound
- May not converge (but in practice always does, with dampening)
- The final beliefs may not be realizable marginals

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Norm-Product BP (Hazan and Shashua, 2010)

Subsumes loopy BP and TRBP

Relies on a convex approximation to the entropy using *counting numbers* $c_i \geq 0$ and $c_s > 0$ (in its simpler variant)

$$H_{\text{NPBP}}(\mathbb{P}_{\boldsymbol{p}, \boldsymbol{q}}) = \sum_i c_i H_i(\boldsymbol{p}_i) + \sum_s c_s H_s(\boldsymbol{q}_s)$$

Messages will become *norms*

Recall the definition of ℓ_p -norm: $\|\mathbf{x}\|_p = (\sum_i |x_i|^p)^{1/p}$

NPBP Messages

- Variable-to-factor messages:

$$m_{i \rightarrow s}(y_i) = \frac{\left(\psi_i(y_i) \prod_{s' \in N(i)} n_{s' \rightarrow i}(y_i) \right)^{c_s / (c_i + \sum_{s' \in N(i)} c'_s)}}{n_{s \rightarrow i}(y_i)}$$

- Factor-to-variable messages:

$$n_{s \rightarrow i}(y_i) = \left(\sum_{y_s \sim y_i} \left(\psi_s(y_s) \prod_{j \in N(s) \setminus \{i\}} m_{j \rightarrow s}(y_j) \right)^{1/c_s} \right)^{c_s}$$

- Variable beliefs:

$$p_i(y_i) \propto \left(\psi_i(y_i) \prod_{s \in N(i)} n_{s \rightarrow i}(y_i) \right)^{1 / (c_i + \sum_{s' \in N(i)} c'_s)}$$

- Factor beliefs:

$$q_s(y_s) \propto \left(\psi_s(y_s) \prod_{i \in N(s)} m_{i \rightarrow s}(y_i) \right)^{c_s}$$

Summary of NPB

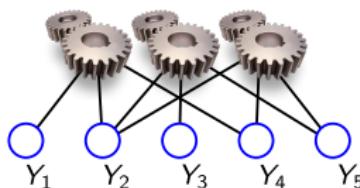
Advantages:

- Still simple to implement
- Entropy approximation is concave (no local minima)
- Always converges (primal-dual block ascent)
- Lots of knobs (the counting numbers)

Disadvantages:

- Lots of knobs (the counting numbers)
- Messages are not computed in parallel (otherwise, may not converge)
- The final beliefs may not be realizable marginals

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Outline

① Structured Prediction and Factor Graphs

② Integer Linear Programming

③ Message-Passing Algorithms

Sum-Product

Max-Product

④ Dual Decomposition

Zero-Limit Temperature

Define Z_ϵ where ϵ is a *temperature parameter*:

$$Z_\epsilon(\psi, x) = \left(\sum_{y \in \mathcal{Y}(x)} \prod_i \psi_i(y_i)^{1/\epsilon} \prod_s \psi_s(y_s)^{1/\epsilon} \right)^\epsilon$$

If $\epsilon = 1$, this becomes the partition function $Z(\psi, x)$

If $\epsilon \rightarrow 0$, this becomes the mode of $\mathbb{P}_\psi(y|x)$

Note that $Z_\epsilon(\psi, x) = Z(\psi^{1/\epsilon}, x)^\epsilon$ for any ϵ , i.e., **Z_ϵ can be computed by the same means as the partition function by scaling the potentials**

By choosing a small enough ϵ , any sum-product message-passing algorithm can be used to approximate the MAP

There is a trade-off between precision and numerical stability

Max-Product Belief Propagation

- For MAP decoding instead of marginal decoding
- Only change: factor-to-variable messages (max instead of \sum)

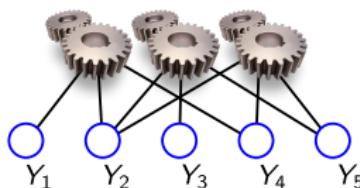
$$n_{s \rightarrow i}(y_i) = \max_{\mathbf{y}_s \sim y_i} \left(\psi_s(\mathbf{y}_s) \prod_{j \in N(s) \setminus \{i\}} m_{j \rightarrow s}(y_j) \right) = \frac{\max_{\mathbf{y}_s \sim y_i} q_s(\mathbf{y}_s)}{m_{i \rightarrow s}(y_i)}$$

- If the graph has no cycles, beliefs will converge to **max-marginals**:

$$p_i(y_i) \rightarrow \max_{y \sim y_i} \mathbb{P}_{\psi}(y|x), \quad q_s(\mathbf{y}_s) \rightarrow \max_{\mathbf{y} \sim \mathbf{y}_s} \mathbb{P}_{\psi}(y|x)$$

- Decoding the best max-marginal at each variable node gives the MAP
- With cycles: not guaranteed to converge, and even if it does, no relation with LP-MAP

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

TRW-S (Kolmogorov, 2006)

Same rationale as sum-product TRBP: cover the graph with spanning trees, and compute messages using edge appearance probabilities

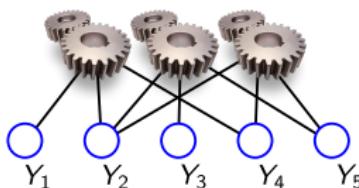
Only differences:

- Replace \sum with max
- Messages need to be computed *sequentially* for convergence

As max-product loopy BP, all is required is to compute *local max-marginals*

Under mild assumptions, gives the solution of LP-MAP

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Max-Product LP (Globerson and Jaakkola, 2008)

Derived by writing the dual of LP-MAP, and solving it with a **block coordinate descent algorithm**

The message updates need to be computed in a sequential schedule

Progress in the dual objective is monotonic

Drawback: since the dual is non-smooth, we may get stuck at a suboptimal point

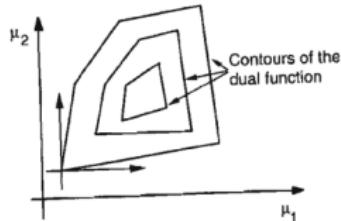


Figure 6.3.6. The basic difficulty with coordinate ascent for a nondifferentiable dual function. At some points it may be impossible to improve the dual function along any coordinate direction.

(From Bertsekas et al. (1999))

MPLP Messages

- Variable-to-factor messages:

$$m_{i \rightarrow s}(y_i) = \psi_i(y_i) \prod_{s' \in N(i) \setminus \{s\}} n_{s' \rightarrow i}(y_i)$$

- Factor-to-variable messages:

$$n_{s \rightarrow i}(y_i) = \frac{\max_{y_s \sim y_i} \left(\psi_s(y_s)^{1/|N(s)|} \prod_{j \in N(s)} m_{j \rightarrow s}(y_j)^{1/|N(s)|} \right)}{m_{i \rightarrow s}(y_i)}$$

Summary of MPLP

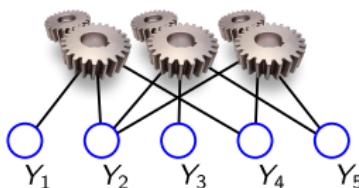
Advantages:

- Very simple to implement
- Handles structured and logic factors (only need to compute local max-marginals)
- Monotonically improves the dual
- No parameters to tune

Disadvantages:

- Can get stuck at a suboptimal solution (general problem with nonsmooth coordinate ascent)
- Messages are not computed in parallel (otherwise, may not converge)

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Summing Up Message-Passing

- BP algorithms and their variants can be used both for MAP and marginal decoding
- They need to compute local *marginals* (sum-product) or *max-marginals* (max-product)
- Always exact if the graph has no cycles; approximate otherwise
- They correspond to minimizing an energy approximation over the local polytope
- Some variants do convex approximations or compute upper bounds
- Two views of MAP decoding: (1) the near-zero temperature limit of marginal decoding; (2) a non-smooth optimization problem

Outline

① Structured Prediction and Factor Graphs

② Integer Linear Programming

③ Message-Passing Algorithms

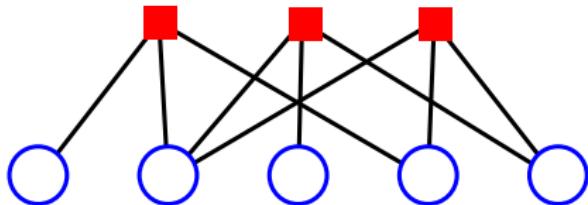
Sum-Product

Max-Product

④ Dual Decomposition

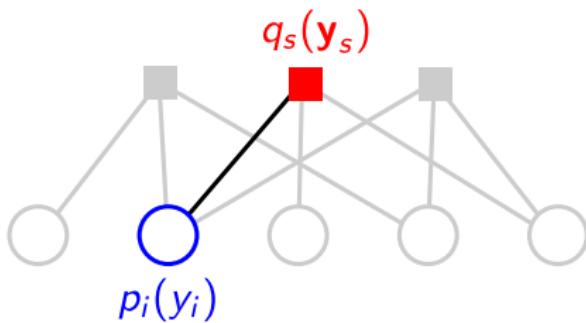
Dual Decomposition

- Old idea in optimization (Dantzig and Wolfe, 1960; Everett III, 1963)
- First proposed by Komodakis et al. (2007) in computer vision
- Introduced in NLP by Rush et al. (2010) for model combination
- Successful in syntax, semantics, MT: Koo et al. (2010); Chang and Collins (2011); Martins et al. (2011b); Almeida et al. (2014); Martins and Almeida (2014), and many others.



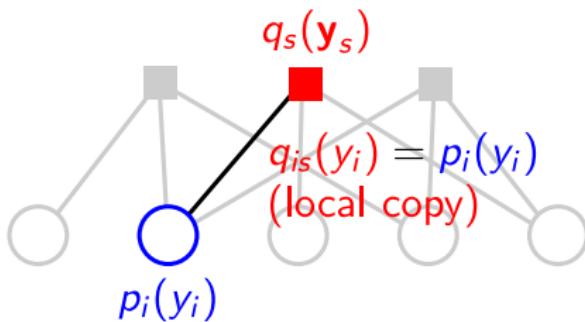
Dual Decomposition

- Old idea in optimization (Dantzig and Wolfe, 1960; Everett III, 1963)
- First proposed by Komodakis et al. (2007) in computer vision
- Introduced in NLP by Rush et al. (2010) for model combination
- Successful in syntax, semantics, MT: Koo et al. (2010); Chang and Collins (2011); Martins et al. (2011b); Almeida et al. (2014); Martins and Almeida (2014), and many others.



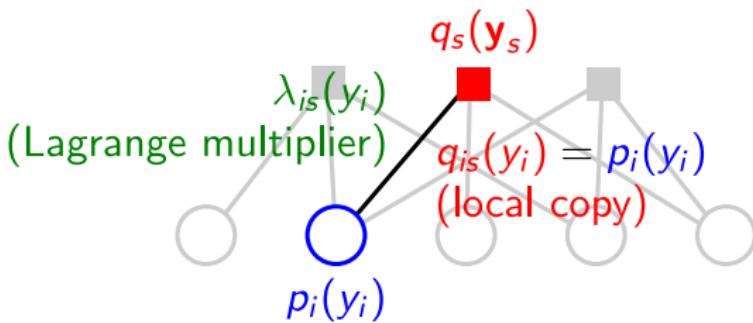
Dual Decomposition

- Old idea in optimization (Dantzig and Wolfe, 1960; Everett III, 1963)
- First proposed by Komodakis et al. (2007) in computer vision
- Introduced in NLP by Rush et al. (2010) for model combination
- Successful in syntax, semantics, MT: Koo et al. (2010); Chang and Collins (2011); Martins et al. (2011b); Almeida et al. (2014); Martins and Almeida (2014), and many others.



Dual Decomposition

- Old idea in optimization (Dantzig and Wolfe, 1960; Everett III, 1963)
- First proposed by Komodakis et al. (2007) in computer vision
- Introduced in NLP by Rush et al. (2010) for model combination
- Successful in syntax, semantics, MT: Koo et al. (2010); Chang and Collins (2011); Martins et al. (2011b); Almeida et al. (2014); Martins and Almeida (2014), and many others.



Recap: LP-MAP

Recall the LP-MAP problem:

$$\begin{aligned} & \text{maximize} \quad \sum_i \theta_i \cdot \mathbf{p}_i + \sum_s \theta_s \cdot \mathbf{q}_s \\ & \text{subject to} \quad \begin{cases} \mathbf{q}_s \in \Delta^{|\mathcal{Y}_s|}, \forall s \\ \mathbf{p}_i = \mathbf{M}_{is} \mathbf{q}_s, \forall i, s. \end{cases} \quad (\text{local polytope}) \end{aligned}$$

Matrix $\mathbf{M}_{is} \in \{0, 1\}^{|\mathcal{Y}_i| \times |\mathcal{Y}_s|}$ represents the constraints
 $p_i(y_i) = \sum_{y_s \sim y_i} q_s(y_s)$

We'll reformulate this problem by:

- ① Introducing copy variables $\mathbf{q}_{is} = \mathbf{p}_i$
- ② Defining $\theta_{is} := \theta_i / |N(i)|$

Reformulation of LP-MAP

The problem becomes:

$$\begin{aligned} & \text{maximize} \quad \sum_s \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} \theta_{is} \cdot \mathbf{q}_{is} \right) \\ & \text{subject to} \quad \begin{cases} \mathbf{q}_s \in \Delta^{|\mathcal{Y}_s|}, \quad \forall s \\ \mathbf{q}_{is} = \mathbf{M}_{is} \mathbf{q}_s, \quad \forall i, s \quad (\text{local polytope}) \\ \mathbf{q}_{is} = \mathbf{p}_i, \quad \forall i, s. \end{cases} \end{aligned}$$

By introducing Lagrange multipliers for the last constraints, we get the following **Lagrangian function**:

$$\mathcal{L}(\mathbf{p}, \mathbf{q}, \boldsymbol{\lambda}) = \sum_s \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} \theta_{is} \cdot \mathbf{q}_{is} \right) + \sum_{is} \boldsymbol{\lambda}_{is} \cdot (\mathbf{p}_i - \mathbf{q}_{is})$$

Dual of LP-MAP

The dual problem is

$$\text{minimize } \sum_s g_s(\lambda) \quad \text{subject to } \lambda \in \Lambda := \left\{ \lambda \mid \sum_{s \in N(i)} \lambda_{is} = \mathbf{0} \right\}$$

where the $g_s(\lambda)$ are **local subproblems**,

$$\begin{aligned} g_s(\lambda) &:= \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} \right) \\ &= \max_{\mathbf{y}_s \in \mathcal{Y}_s} \left(\theta_s(\mathbf{y}_s) + \sum_{i \in N(s)} (\theta_{is}(y_i) + \lambda_{is}(y_i)) \right) \end{aligned}$$

and $\bar{\mathbf{q}}_s \in \mathcal{Q}_s$ encodes the constraints $\begin{cases} \mathbf{q}_s \in \Delta^{|\mathcal{Y}_s|} \\ \mathbf{q}_{is} = \mathbf{M}_{is} \mathbf{q}_s, \forall i \in N(s). \end{cases}$

A subgradient can be computed by solving these local subproblems

Projected Subgradient (Komodakis et al., 2007)

initialize penalties λ to zero
repeat

until consensus (all $q_{is} = p_i$) or maximum number of iterations reached

Projected Subgradient (Komodakis et al., 2007)

initialize penalties λ to zero

repeat

for each factor s **do**

$$\bar{q}_s \leftarrow \arg \max_{\bar{q}_s \in \mathcal{Q}_s} \theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is}$$

end for

until consensus (all $\mathbf{q}_{is} = \mathbf{p}_i$) or maximum number of iterations reached

Projected Subgradient (Komodakis et al., 2007)

initialize penalties λ to zero

repeat

for each factor s **do**

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is}$$

end for

$$\mathbf{p}_i \leftarrow \frac{1}{|N(i)|} \sum_{s \in N(i)} \mathbf{q}_{is}$$

until consensus (all $\mathbf{q}_{is} = \mathbf{p}_i$) or maximum number of iterations reached

Projected Subgradient (Komodakis et al., 2007)

initialize penalties λ to zero

repeat

for each factor s **do**

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is}$$

end for

$$\mathbf{p}_i \leftarrow \frac{1}{|N(i)|} \sum_{s \in N(i)} \mathbf{q}_{is}$$

$$\lambda_{is} \leftarrow \lambda_{is} - \eta (\mathbf{q}_{is} - \mathbf{p}_i)$$

until consensus (all $\mathbf{q}_{is} = \mathbf{p}_i$) or maximum number of iterations reached

Projected Subgradient (Komodakis et al., 2007)

initialize penalties λ to zero

repeat

for each factor s **do**

$$\bar{q}_s \leftarrow \arg \max_{\bar{q}_s \in \mathcal{Q}_s} \theta_s \cdot q_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot q_{is}$$

end for

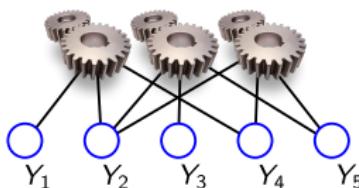
$$p_i \leftarrow \frac{1}{|N(i)|} \sum_{s \in N(i)} q_{is}$$

$$\lambda_{is} \leftarrow \lambda_{is} - \eta (q_{is} - p_i)$$

until consensus (all $q_{is} = p_i$) or maximum number of iterations reached

- Guaranteed to converge to an ϵ -accurate solution after at most $O(1/\epsilon^2)$ iterations
- **Problem:** too slow when there are many factors (Martins et al., 2011b)

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Two fundamental problems with the subgradient algorithm:

- ① The dual objective $\sum_s g_s(\lambda)$ is non-smooth
- ② Consensus is promoted only through updating λ (no memory about past updates)

Two fundamental problems with the subgradient algorithm:

- ① The dual objective $\sum_s g_s(\lambda)$ is non-smooth
- ② Consensus is promoted only through updating λ (no memory about past updates)

How can dual decomposition be accelerated?

Two fundamental problems with the subgradient algorithm:

- ① The dual objective $\sum_s g_s(\lambda)$ is non-smooth
- ② Consensus is promoted only through updating λ (no memory about past updates)

How can dual decomposition be accelerated?

- Jovic et al. (2010) smooth the objective and use gradient methods
- Martins et al. (2011a): augmented Lagrangian

Two fundamental problems with the subgradient algorithm:

- ① The dual objective $\sum_s g_s(\lambda)$ is non-smooth
- ② Consensus is promoted only through updating λ (no memory about past updates)

How can dual decomposition be accelerated?

- Jovic et al. (2010) smooth the objective and use gradient methods
- Martins et al. (2011a): augmented Lagrangian

Accelerated Gradient (Jojic et al., 2010)

Basic idea: make the dual objective smooth by adding an entropic perturbation with a near-zero ϵ temperature (also Johnson (2008))

The subproblems become local *marginal* computations instead of maximizations

With Nesterov's accelerated gradient method (Nesterov, 2005), the iteration bound goes from $O(1/\epsilon^2)$ to $O(1/\epsilon)$

Accelerated Gradient (Jojic et al., 2010)

Basic idea: make the dual objective smooth by adding an entropic perturbation with a near-zero ϵ temperature (also Johnson (2008))

The subproblems become local *marginal* computations instead of maximizations

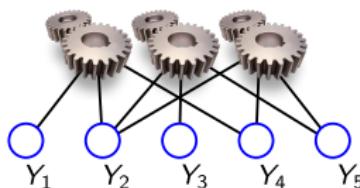
With Nesterov's accelerated gradient method (Nesterov, 2005), the iteration bound goes from $O(1/\epsilon^2)$ to $O(1/\epsilon)$

However: very sensitive to the temperature parameter

With low temperatures, may face numerical issues (in particular for some hard-constraint factors)

In practice, quite slow to take off (we'll see some plots later)

What Kind of Local Decoding Do We Need?



Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Two fundamental problems with the subgradient algorithm:

- ① The dual objective $\sum_s g_s(\lambda)$ is non-smooth
- ② Consensus is promoted only through updating λ (no memory about past updates)

How can dual decomposition be accelerated?

- Jovic et al. (2010) smooth the objective and use gradient methods ✓
- Martins et al. (2011a): augmented Lagrangian

Two fundamental problems with the subgradient algorithm:

- ① The dual objective $\sum_s g_s(\lambda)$ is non-smooth
- ② Consensus is promoted only through updating λ (no memory about past updates)

How can dual decomposition be accelerated?

- Jovic et al. (2010) smooth the objective and use gradient methods ✓
- Martins et al. (2011a): augmented Lagrangian

Alternating Directions Dual Decomposition (AD³)



Based on the **alternating direction method of multipliers (ADMM)**:

- an old method in optimization inspired by augmented Lagrangians (Gabay and Mercier, 1976; Glowinski and Marroco, 1975)
- a natural fit to consensus problems
- a natural “upgrade” of the subgradient algorithm (Boyd et al., 2011)

Augmented Lagrangian and ADMM

Basic idea: augment the Lagrangian function with a **quadratic penalty**

$$\begin{aligned}\mathcal{L}_\eta(\mathbf{p}, \mathbf{q}, \boldsymbol{\lambda}) = & \sum_s \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} \theta_{is} \cdot \mathbf{q}_{is} \right) + \sum_{is} \boldsymbol{\lambda}_{is} \cdot (\mathbf{p}_i - \mathbf{q}_{is}) \\ & - \frac{\eta}{2} \sum_{is} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2\end{aligned}$$

Method of multipliers (super-linear convergence):

- ① Maximize $\mathcal{L}_\eta(\mathbf{p}, \mathbf{q}, \boldsymbol{\lambda})$ jointly w.r.t. \mathbf{p} and \mathbf{q} (challenging)
- ② Multiplier update: $\boldsymbol{\lambda}_{is} \leftarrow \boldsymbol{\lambda}_{is} - \eta(\mathbf{q}_{is} - \mathbf{p}_i)$

Alternating direction method of multipliers: replace step 1 by separate maximizations (first w.r.t. \mathbf{q} , then \mathbf{p})

From Subgradient to AD³ (Martins et al., 2011a)

initialize penalties λ to zero

repeat

for each factor $s = 1, \dots, S$ **do**

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is}$$

end for

$$\mathbf{p}_i \leftarrow \frac{1}{|N(i)|} \sum_{s \in N(i)} \mathbf{q}_{is}$$

$$\lambda_{is} \leftarrow \lambda_{is} - \eta(\mathbf{q}_{is} - \mathbf{p}_i)$$

until consensus (all $\mathbf{q}_{is} = \mathbf{p}_i$) or maximum number of iterations reached

From Subgradient to AD³ (Martins et al., 2011a)

initialize penalties λ to zero

repeat

for each factor $s = 1, \dots, S$ **do**

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} - \frac{\eta}{2} \sum_{i \in N(s)} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2$$

end for

$$\mathbf{p}_i \leftarrow \frac{1}{|N(i)|} \sum_{s \in N(i)} \mathbf{q}_{is}$$

$$\lambda_{is} \leftarrow \lambda_{is} - \eta(\mathbf{q}_{is} - \mathbf{p}_i)$$

until consensus (all $\mathbf{q}_{is} = \mathbf{p}_i$) or maximum number of iterations reached

From Subgradient to AD³ (Martins et al., 2011a)

initialize penalties λ to zero

repeat

for each factor $s = 1, \dots, S$ **do**

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} - \frac{\eta}{2} \sum_{i \in N(s)} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2$$

end for

$$\mathbf{p}_i \leftarrow \frac{1}{|N(i)|} \sum_{s \in N(i)} \mathbf{q}_{is}$$

$$\lambda_{is} \leftarrow \lambda_{is} - \eta(\mathbf{q}_{is} - \mathbf{p}_i)$$

until consensus (all $\mathbf{q}_{is} = \mathbf{p}_i$) or maximum number of iterations reached

- **faster consensus:** regularize \mathbf{q} -step towards average votes in \mathbf{p}

From Subgradient to AD³ (Martins et al., 2011a)

initialize penalties λ to zero

repeat

for each factor $s = 1, \dots, S$ **do**

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} - \frac{\eta}{2} \sum_{i \in N(s)} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2$$

end for

$$\mathbf{p}_i \leftarrow \frac{1}{|N(i)|} \sum_{s \in N(i)} \mathbf{q}_{is}$$

$$\lambda_{is} \leftarrow \lambda_{is} - \eta(\mathbf{q}_{is} - \mathbf{p}_i)$$

until consensus (all $\mathbf{q}_{is} = \mathbf{p}_i$) or maximum number of iterations reached

- **faster consensus:** regularize \mathbf{q} -step towards average votes in \mathbf{p}
- **better stopping conditions:** keeps track of primal and dual residuals

Theoretical Guarantees of AD³

Convergent in primal and dual (Glowinski and Le Tallec, 1989)

Iteration bound: $O(1/\epsilon)$ (cf. $O(1/\epsilon^2)$ for projected subgradient)

Inexact AD³ subproblems: still convergent if residuals are summable
(Eckstein and Bertsekas, 1992)

Always dual feasible: can compute upper bounds and embed in
branch-and-bound toward *exact* decoding (Das et al., 2012)

Theoretical Guarantees of AD³

Convergent in primal and dual (Glowinski and Le Tallec, 1989)

Iteration bound: $O(1/\epsilon)$ (cf. $O(1/\epsilon^2)$ for projected subgradient)

Inexact AD³ subproblems: still convergent if residuals are summable
(Eckstein and Bertsekas, 1992)

Always dual feasible: can compute upper bounds and embed in
branch-and-bound toward *exact* decoding (Das et al., 2012)

**But: AD³ local subproblems are quadratic (more involved than in
projected subgradient)**

Theoretical Guarantees of AD³

Convergent in primal and dual (Glowinski and Le Tallec, 1989)

Iteration bound: $O(1/\epsilon)$ (cf. $O(1/\epsilon^2)$ for projected subgradient)

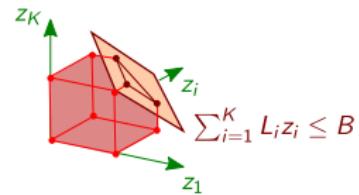
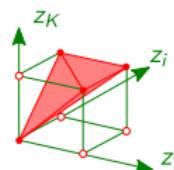
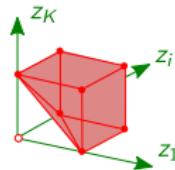
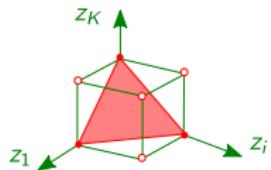
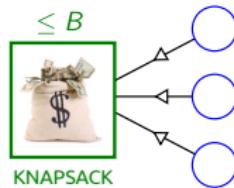
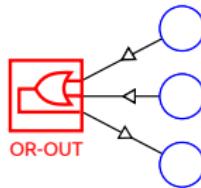
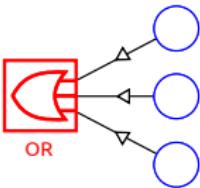
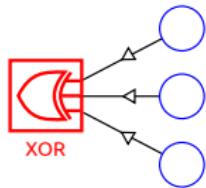
Inexact AD³ subproblems: still convergent if residuals are summable
(Eckstein and Bertsekas, 1992)

Always dual feasible: can compute upper bounds and embed in
branch-and-bound toward *exact* decoding (Das et al., 2012)

**But: AD³ local subproblems are quadratic (more involved than in
projected subgradient)**

Still—very easy and efficient for logic and knapsack factors!

Projecting onto Hard Constraint Polytopes



- Martins et al. (2011a): logic factors can be solved in $O(K)$ time
- **Almeida and Martins (2013): same for knapsack factors!**

Structured Factors

What about structured factors?

Structured Factors

What about structured factors?

Projected subgradient handles these quite well

- combinatorial machinery (Viterbi, Chu-Liu-Edmonds, Fulkerson-Ford, Floyd-Warshall,...)

We cannot solve the AD^3 subproblems with that machinery...

Structured Factors

What about structured factors?

Projected subgradient handles these quite well

- combinatorial machinery (Viterbi, Chu-Liu-Edmonds, Fulkerson-Ford, Floyd-Warshall,...)

We cannot solve the AD^3 subproblems with that machinery...

Or can we?

Structured Factors

What about structured factors?

Projected subgradient handles these quite well

- combinatorial machinery (Viterbi, Chu-Liu-Edmonds, Fulkerson-Ford, Floyd-Warshall,...)

We cannot solve the AD^3 subproblems with that machinery...

Or can we?

Active set method: seek the support of the solution by adding/removing components; very suitable for warm-starting (Nocedal and Wright, 1999)

An Active Set Method for the AD³ Subproblem

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} - \frac{\eta}{2} \sum_{i \in N(s)} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2 \right)$$

An Active Set Method for the AD³ Subproblem

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} - \frac{\eta}{2} \sum_{i \in N(s)} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2 \right)$$

Too many possible assignments: dimension of \mathbf{q}_s is $O(\exp(|N(s)|))$

An Active Set Method for the AD³ Subproblem

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} - \frac{\eta}{2} \sum_{i \in N(s)} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2 \right)$$

Too many possible assignments: dimension of \mathbf{q}_s is $O(\exp(|N(s)|))$

Key result: there's a sparse solution (only $O(|N(s)|)$ nonzeros)

An Active Set Method for the AD³ Subproblem

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} - \frac{\eta}{2} \sum_{i \in N(s)} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2 \right)$$

Too many possible assignments: dimension of \mathbf{q}_s is $O(\exp(|N(s)|))$

Key result: there's a sparse solution (only $O(|N(s)|)$ nonzeros)

Active set methods: seek the support of the solution by adding/removing components; very suitable for warm-starting (Nocedal and Wright, 1999)

Only requirement: a local-max oracle (as in projected subgradient)

An Active Set Method for the AD³ Subproblem

$$\bar{\mathbf{q}}_s \leftarrow \arg \max_{\bar{\mathbf{q}}_s \in \mathcal{Q}_s} \left(\theta_s \cdot \mathbf{q}_s + \sum_{i \in N(s)} (\theta_{is} + \lambda_{is}) \cdot \mathbf{q}_{is} - \frac{\eta}{2} \sum_{i \in N(s)} \|\mathbf{q}_{is} - \mathbf{p}_i\|^2 \right)$$

Too many possible assignments: dimension of \mathbf{q}_s is $O(\exp(|N(s)|))$

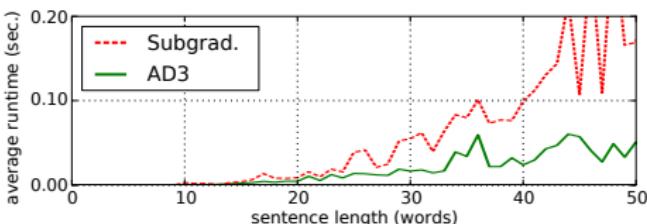
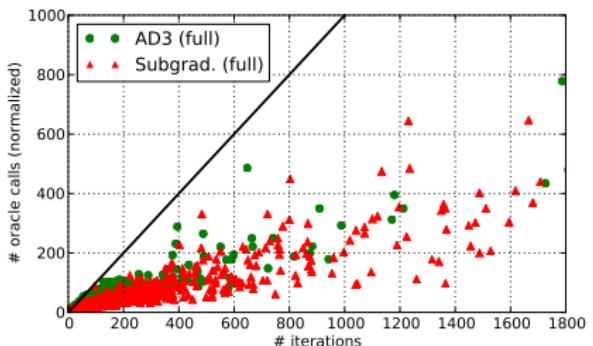
Key result: there's a sparse solution (only $O(|N(s)|)$ nonzeros)

Active set methods: seek the support of the solution by adding/removing components; very suitable for warm-starting (Nocedal and Wright, 1999)

Only requirement: a local-max oracle (as in projected subgradient)

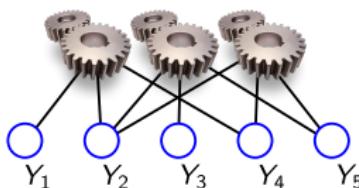
More info: Martins et al. (2015)

Runtime of AD³ vs PSDD (Parsing)



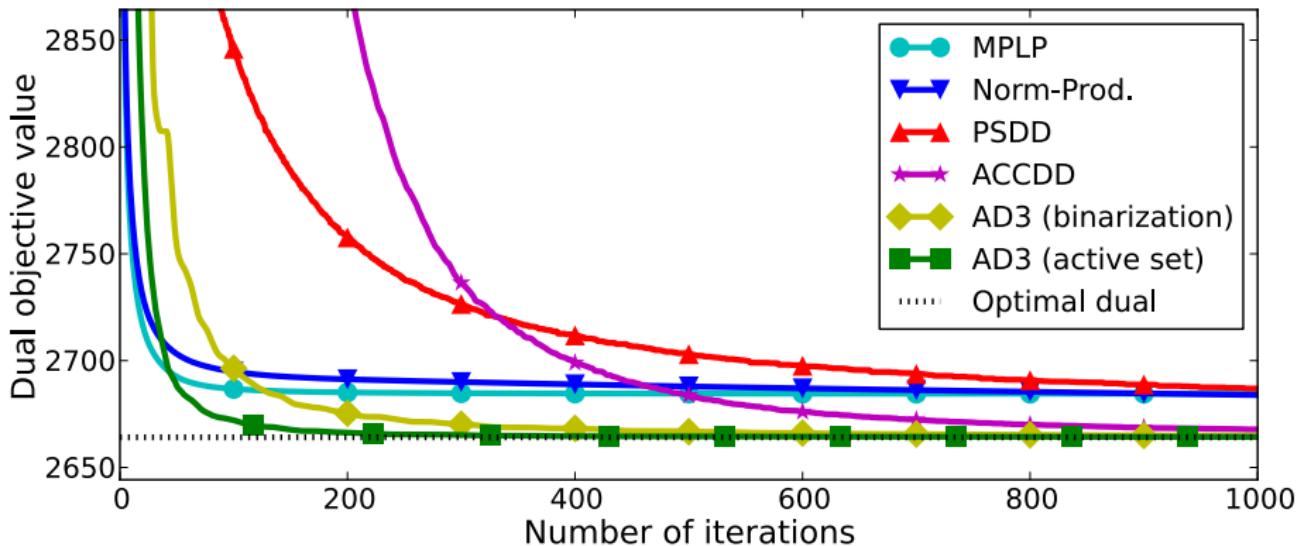
- **Caching** and **warm-starting** the subproblems reduces drastically the number of oracle calls—**huge speed-ups!!**
- **AD³ faster to achieve consensus** (due to the quadratic penalty)

What Kind of Local Decoding Do We Need?



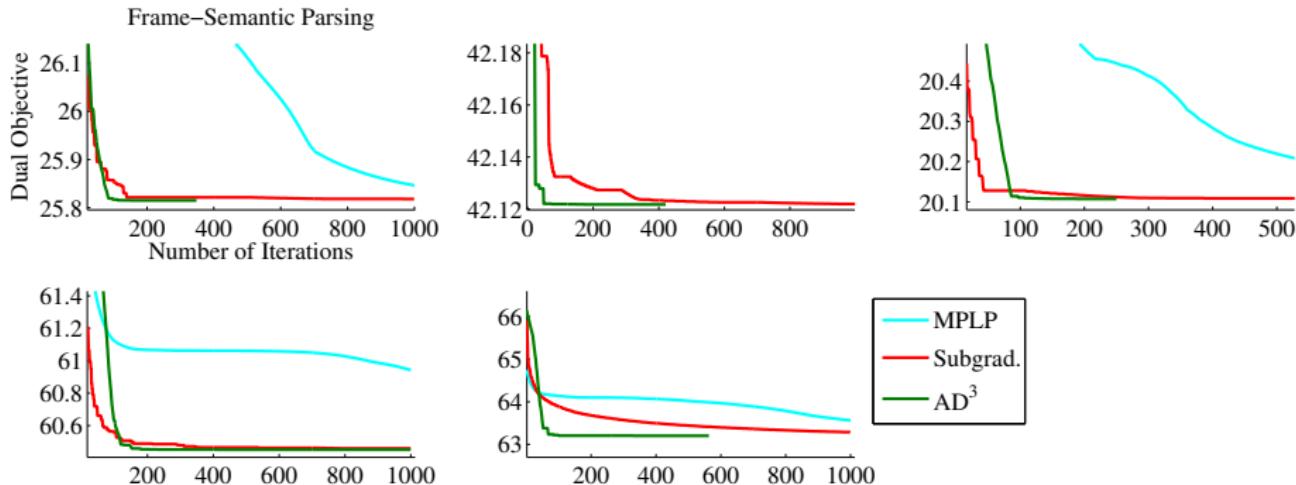
Algorithm	Local Operation
Sum-Prod. BP (Pearl, 1988)	marginals
TRBP (Wainwright et al., 2005)	marginals
Norm-Product BP (Hazan and Shashua, 2010)	marginals
Max-Prod. BP (Pearl, 1988)	max-marginals
TRW-S (Kolmogorov, 2006)	max-marginals
MPLP (Globerson and Jaakkola, 2008)	max-marginals
PSDD (Komodakis et al., 2007)	MAP
Accelerated DD (Jojic et al., 2010)	marginals
AD ³ (Martins et al., 2011a)	QP/MAP

Example: Potts Grid (20×20 , 8 states)



- A. Martins, M. Figueiredo, P. Aguiar, N. Smith, E. Xing (2014).
AD³: Alternating Directions Dual Decomposition for MAP Inference in Graphical Models.
Journal of Machine Learning Research (to appear).

Example: Frame-Semantic Parsing



- Embedded in a branch-and-bound procedure for exact decoding
- D. Das, A. Martins, N. Smith.
“An Exact DD Algorithm for Shallow Semantic Parsing with Constraints.”
*SEM Workshop, 2012.

Some Problems in Which AD³ Have Been Applied

- Dependency parsing (Martins et al., 2011b, 2013)
- Frame semantics (Das et al., 2012)
- Broad-coverage semantic parsing (Martins and Almeida, 2014)
- Compressive summarization (Almeida and Martins, 2013)
- Coreference resolution (Almeida et al., 2014)
- Argument mining with structured SVMs and RNNs (Niculae et al., 2017)
- Joint extraction of events and entities (Yang and Mitchell, 2016)
- Deep structured learning for facial expression intensity estimation (Walecki et al., 2017)

Some Problems in Which AD³ Have Been Applied

- Dependency parsing (Martins et al., 2011b, 2013)
- Frame semantics (Das et al., 2012)
- Broad-coverage semantic parsing (Martins and Almeida, 2014)
- Compressive summarization (Almeida and Martins, 2013)
- Coreference resolution (Almeida et al., 2014)
- Argument mining with structured SVMs and RNNs (Niculae et al., 2017)
- Joint extraction of events and entities (Yang and Mitchell, 2016)
- Deep structured learning for facial expression intensity estimation (Walecki et al., 2017)

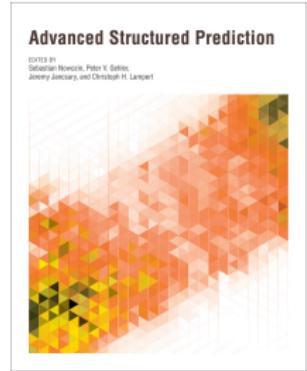
Some Problems in Which AD³ Have Been Applied

- Dependency parsing (Martins et al., 2011b, 2013)
- Frame semantics (Das et al., 2012)
- Broad-coverage semantic parsing (Martins and Almeida, 2014)
- Compressive summarization (Almeida and Martins, 2013)
- Coreference resolution (Almeida et al., 2014)
- Argument mining with structured SVMs and RNNs (Niculae et al., 2017)
- Joint extraction of events and entities (Yang and Mitchell, 2016)
- Deep structured learning for facial expression intensity estimation (Walecki et al., 2017)

... A great fit to many other applications!!

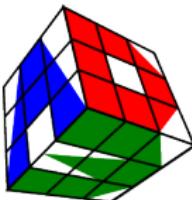
Literature Pointers

- André F. T. Martins.
“AD³: A Fast Decoder for Structured Prediction.”
Book chapter of *Advanced Structured Prediction*,
Sebastian Nowozin, Peter V. Gehler, Jeremy
Jancsary, and Christoph H. Lampert (Editors),
MIT Press, 2014.
- A. Martins, M. Figueiredo, P. Aguiar, N. Smith, E. Xing.
“AD3: Alternating Directions Dual Decomposition for MAP Inference
in Graphical Models.”
JMLR 2015.



More details: EMNLP 2014 tutorial on “LP Decoders for NLP.”

Try It Yourself: AD³ Toolkit



- Freely available at: <http://www.ark.cs.cmu.edu/AD3>
- Implemented in C++, includes a Python wrapper (thanks to Andy Mueller and Vlad Niculae)
- Implements MPLP, PSDD, AD³ for arbitrary factor graphs
- Many built-in factors: logic, knapsack, dense, and some structured factors
- You can implement your own factor (only need to write a local MAP decoder!)
- Toy examples included (parsing, coreference, Potts models)

Summing Up Dual Decomposition

- Dual decomposition is a general optimization technique that splits the dual into several subproblems (one per factor) that must agree on overlaps
- This can be used to solve LP-MAP
- We discussed three variants: subgradient (PSDD), accelerated gradient (ADD), and alternating directions (AD³)
- The algorithms are convergent and retrieve the true MAP if the graph has no cycles; they also give certificates when the solution of LP-MAP equals the MAP
- For PSDD and AD³ only local maximizations are necessary; ADD requires computing marginals

Conclusions

- Graphical models are a useful way of representing structural relations among random variables
- Many structured problems are NP-hard or expensive (constrained models, diversity, combination of structured models)
- Often they can be approximately decoded via Linear Programming (e.g., by relaxing an ILP)
- The structure inherent to these problems can be represented with a factor graph
- Message-passing and dual decomposition algorithms can solve these LPs efficiently, exploiting the structure of the graph
- Conceptually: approximate *global* decoding by invoking only *local* decoders (local maximizations, marginals, max-marginals, QPs, ...)
- AD³ is faster than the subgradient algorithm both in theory and in practice, and requires the same local decoders

Thank you!

Questions?



References I

- Almeida, M. B. and Martins, A. F. T. (2013). Fast and robust compressive summarization with dual decomposition and multi-task learning. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Almeida, M. S. C., Almeida, M. B., and Martins, A. F. T. (2014). A Joint Model For Quotation Attribution and Coreference Resolution. In *Proc. of Annual Meeting of the European Chapter of the Association for Computational Linguistics*.
- Auli, M. and Lopez, A. (2011). A Comparison of Loopy Belief Propagation and Dual Decomposition for Integrated CCG Supertagging and Parsing. In *Proc. of Annual Meeting of the Association for Computational Linguistics*.
- Bertsekas, D., Hager, W., and Mangasarian, O. (1999). *Nonlinear programming*. Athena Scientific.
- Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J. (2011). *Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers*. Now Publishers.
- Burkett, D. and Klein, D. (2012). Fast inference in phrase extraction models with belief propagation. In *Proc. of the North American Chapter of the Association for Computational Linguistics*, pages 29–38. Association for Computational Linguistics.
- Chang, Y.-W. and Collins, M. (2011). Exact decoding of phrase-based translation models through lagrangian relaxation. In *Proc. of Empirical Methods for Natural Language Processing*.
- Dantzig, G. and Wolfe, P. (1960). Decomposition principle for linear programs. *Operations Research*, 8(1):101–111.
- Dantzig, G. B. (1947). Maximization of a linear function of variables subject to linear inequalities. Published in T.C. Koopmans (ed.): *Activity Analysis of Production and Allocation*, New York-London 1951, pages 339–347.
- Das, D., Martins, A. F. T., and Smith, N. A. (2012). An Exact Dual Decomposition Algorithm for Shallow Semantic Parsing with Constraints. In *Proc. of First Joint Conference on Lexical and Computational Semantics (*SEM)*.
- Dreyer, M. and Eisner, J. (2009). Graphical models over multiple strings. In *Proc. of Empirical Methods in Natural Language Processing*, pages 101–110. Association for Computational Linguistics.
- Duchi, J., Tarlow, D., Elidan, G., and Koller, D. (2007). Using combinatorial optimization within max-product belief propagation. *Advances in Neural Information Processing Systems*, 19.
- Eckstein, J. and Bertsekas, D. (1992). On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators. *Mathematical Programming*, 55(1):293–318.

References II

- Everett III, H. (1963). Generalized Lagrange multiplier method for solving problems of optimum allocation of resources. *Operations Research*, 11(3):399–417.
- Gabay, D. and Mercier, B. (1976). A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers and Mathematics with Applications*, 2(1):17–40.
- Globerson, A. and Jaakkola, T. (2008). Fixing max-product: Convergent message passing algorithms for MAP LP-relaxations. *Neural Information Processing Systems*, 20.
- Glowinski, R. and Le Tallec, P. (1989). *Augmented Lagrangian and operator-splitting methods in nonlinear mechanics*. Society for Industrial Mathematics.
- Glowinski, R. and Marroco, A. (1975). Sur l'approximation, par éléments finis d'ordre un, et la résolution, par penalisation-dualité, d'une classe de problèmes de Dirichlet non linéaires. *Rev. Franc. Automat. Inform. Rech. Operat.*, 9:41–76.
- Hazan, T. and Shashua, A. (2010). Norm-product belief propagation: Primal-dual message-passing for approximate inference. *IEEE Transactions on Information Theory*, 56(12):6294–6316.
- Johnson, J. (2008). Equivalence of Entropy Regularization and Relative-Entropy Proximal Method. Unpublished manuscript.
- Jojic, V., Gould, S., and Koller, D. (2010). Accelerated dual decomposition for MAP inference. In *International Conference of Machine Learning*.
- Kantorovich, L. V. (1940). A new method of solving of some classes of extremal problems. In *Dokl. Akad. Nauk SSSR*, volume 28, pages 211–214.
- Karp, R. M. (1972). *Reducibility among combinatorial problems*. Springer.
- Khachiyan, L. G. (1980). Polynomial algorithms in linear programming. *USSR Computational Mathematics and Mathematical Physics*, 20(1):53–72.
- Koller, D. and Friedman, N. (2009). *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press.
- Kolmogorov, V. (2006). Convergent tree-reweighted message passing for energy minimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28:1568–1583.
- Komodakis, N., Paragios, N., and Tziritas, G. (2007). MRF optimization via dual decomposition: Message-passing revisited. In *Proc. of International Conference on Computer Vision*.

References III

- Koo, T., Rush, A. M., Collins, M., Jaakkola, T., and Sontag, D. (2010). Dual decomposition for parsing with non-projective head automata. In *Proc. of Empirical Methods for Natural Language Processing*.
- Lauritzen, S. (1996). *Graphical Models*. Clarendon Press, Oxford.
- MacKay, D. (2003). *Information Theory, Inference, and Learning Algorithms*, volume 7. Cambridge University Press.
- Martins, A. F. T., Almeida, M. B., and Smith, N. A. (2013). Turning on the turbo: Fast third-order non-projective turbo parsers. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Martins, A. F. T. and Almeida, M. S. C. (2014). Pribberam: A Turbo Semantic Parser with Second Order Features. In *Proc. of International Workshop on Semantic Evaluation (SemEval), task 8: Broad Coverage Semantic Dependency Parsing*.
- Martins, A. F. T., Figueiredo, M. A. T., Aguiar, P. M. Q., Smith, N. A., and Xing, E. P. (2011a). An Augmented Lagrangian Approach to Constrained MAP Inference. In *Proc. of International Conference on Machine Learning*.
- Martins, A. F. T., Figueiredo, M. A. T., Aguiar, P. M. Q., Smith, N. A., and Xing, E. P. (2015). AD³: Alternating Directions Dual Decomposition for MAP Inference in Graphical Models. *Journal of Machine Learning Research (to appear)*.
- Martins, A. F. T., Smith, N. A., Aguiar, P. M. Q., and Figueiredo, M. A. T. (2011b). Dual Decomposition with Many Overlapping Components. In *Proc. of Empirical Methods for Natural Language Processing*.
- McAllester, D., Collins, M., and Pereira, F. (2008). Case-factor diagrams for structured probabilistic modeling. *Journal of Computer and System Sciences*, 74(1):84–96.
- Nesterov, Y. (2005). Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1):127–152.
- Niculae, V., Park, J., and Cardie, C. (2017). Argument mining with structured svms and rnns. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Nocedal, J. and Wright, S. (1999). *Numerical optimization*. Springer verlag.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann.
- Richardson, M. and Domingos, P. (2006). Markov logic networks. *Machine Learning*, 62(1):107–136.
- Roth, D. and Yih, W. (2004). A linear programming formulation for global inference in natural language tasks. In *International Conference on Natural Language Learning*.

References IV

- Rush, A., Sontag, D., Collins, M., and Jaakkola, T. (2010). On dual decomposition and linear programming relaxations for natural language processing. In *Proc. of Empirical Methods for Natural Language Processing*.
- Smith, D. and Eisner, J. (2008). Dependency parsing by belief propagation. In *Proc. of Empirical Methods for Natural Language Processing*.
- Wainwright, M. and Jordan, M. (2008). *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers.
- Wainwright, M. J., Jaakkola, T., and Willsky, A. (2005). A new class of upper bounds on the log partition function. *IEEE Transactions on Information Theory*, 51(7):2313–2335.
- Walecki, R., Rudovic, O., Pavlovic, V., Schuller, B., and Pantic, M. (2017). Deep structured learning for facial expression intensity estimation. In *Proceedings of IEEE Int'l Conf. on Computer Vision and Pattern Recognition*, Honolulu, Hawaii.
- Yang, B. and Mitchell, T. (2016). Joint extraction of events and entities within a document context. In *Proc. of the North American Chapter of the Association for Computational Linguistics*.
- Yedidia, J. S., Freeman, W. T., and Weiss, Y. (2001). Generalized belief propagation. In *Neural Information Processing Systems*.