# Efficient Computation of Expectations under Spanning Tree Distributions

September 8, 2021

# Dependency Trees

- ▶ Classical representation of text
- ▶ Similar to phrase structure grammars
  - ▶ Phrase structure: Labels groups of words
  - ▶ Dependency: Labels relationship between pairs of words
- ▶ Useful if a language has free word order
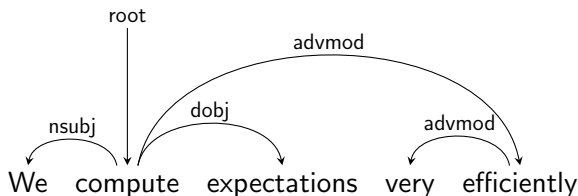- ▶ Spanning trees



Figure: Example of a dependency tree

# Efficient Computation of Expectations under Spanning Tree Distributions

- A framework for computing expectations (of decomposable functions) over spanning trees
  - Unifies previous algorithms
  - Additionally show large asymptotic and empirical speed improvements over particular past implementations
- Relies on connections between moments and derivatives
  - Uses automatic differentiation for easy to implement and efficient algorithms
- Fun demonstration of a general 'inference with automatic differentiation' recipe
  - Compute partition function
  - Use AD to compute expectations

# Outline

Goal is to compute expectations over spanning trees

- ▶ Background: Distributions over (spanning) trees
- ▶ Method: Connecting expectations to derivatives
    - ▶ Use properties of our choice of tree distribution and decomposable functions
    - ▶ Stitch together into efficient algorithms with automatic differentiation
- ▶ Computational complexity

# Distributions over Trees

Assuming fixed length sentence, with $N$ nodes:

▶ Weighted edges

$$(i \xrightarrow{w_{ij}} j) \in \mathcal{E}$$

▶ Trees weights

$$w(d) := \prod_{(i \to j) \in d} w_{ij}$$

▶ Tree probability obtained via normalization

$$p(d) := \frac{w(d)}{Z},$$

where

$$Z := \sum_{d \in \mathcal{D}} w(d) = \sum_{d \in \mathcal{D}} \prod_{(i \to j) \in d} w_{ij}$$

▶ Computation of $Z$ is tractable despite exponentially many trees in $\mathcal{D}$

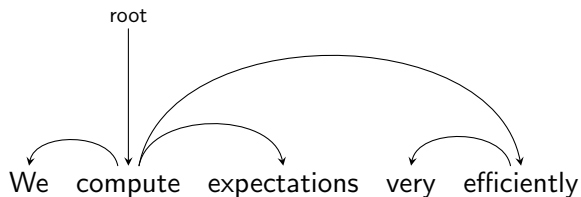# Distributions over spanning trees



Figure: Example of a dependency tree

# Matrix-Tree Theorem (MTT)

- Compute partition function $Z$ using the graph Laplacian $L \in \mathbb{R}^{N \times N}$: $Z = |L|$.
    - Simple (undirected, unweights) graphs: $L = D - A$, where $D$ is the degree matrix and $A$ the adjacency
    - Weighted directed graphs:

$$L_{ij} := \begin{cases} \sum_{i' \in \mathcal{N} \smallsetminus \{j\}} w_{i'j} & \textbf{if } i = j \\ -w_{ij} & \textbf{otherwise} \end{cases}$$

- Determinant can be computed in $O(N^3)$ time
- Flexibility in choice of $L$ depending on example, care is needed
- Results are general for choice of Laplacian $L$

# Expectations

- Main goal is to compute expectations / totals

$$\mathbb{E}_d\big[f(d)\big] := \sum_{d \in \mathcal{D}} p(d)f(d) = \frac{1}{Z}\sum_{d \in \mathcal{D}} w(d)f(d) = \frac{1}{Z}\overline{f}$$

- Need to consider every unique tree with nonzero mass unless $f : \mathcal{D} \to \mathbb{R}^F$ decomposes
- Consider two families of decomposable $f$
  - Additively decomposable: Shannon entropy, KL
  - Second-order additively decomposable: Gradient of entropy, covariance

# Intuition for decompositions

- Recall: Tree distribution is a distribution over a binary vector of size $O(N^2)$
- Compute expectations by
  - Storing the relevant parts of trees and their marginal probabilities
  - Applying $f$ to those parts
- Consider different levels of decomposability of $f$
  - Does not decompose: exponential num trees
  - Decomposes over pairs of edges: $O(N^4)$ evals
  - Decomposes over edges: $O(N^2)$ evals
- Combine $p(\mathrm{part})f(\mathrm{part})$ to compute expectation (unnormalized total using $p(\mathrm{part}) = \frac{1}{Z}\tilde{w}_{\mathrm{part}}$)

# (Unnormalized) marginals

▶ Unary marginals (prob of one edge appearing in a tree)

$$p((i \rightarrow j) \in d) = \frac{1}{Z} \sum_{d \in \mathcal{D}_{ij}} w(d) = \frac{1}{Z} \widetilde{w_{ij}}$$

▶ Binary marginals (prob of two edges appearing together in a tree)

$$p((i \rightarrow j), (k \rightarrow l) \in d) = \frac{1}{Z} \sum_{d \in \mathcal{D}_{ij,kl}} w(d) = \frac{1}{Z} \widetilde{w_{ij,kl}}$$

▶ Want to show that
   ▶ Can decompose totals into sums over marginals (enumerate parts of trees)
   ▶ Can compute marginals using automatic differentiation (cheap gradient principle, simple implementation)
   ▶ Avoid storing full probability tables (requires low-level tricks, not our focus / see paper for details)

# Additively decomposable functions

Allows us to decompose $r : \mathcal{D} \to \mathbb{R}^R$ into functions of edges and combine with (unnormalized) unary marginals (over those tree edges)

$$r(d) = \sum_{(i \to j) \in d} r_{ij}, \tag{1}$$

with $r_{ij} \in \mathbb{R}^R$

# Additively decomposable functions: Marginals

Marginals are expectations

$$p((i \to j) \in d) = \mathbb{E}_{p(d)} \left[ 1((i \to j) \in d) \right]_{ij} = [\sum_d p(d) r(d)]_{ij}$$

Set

$$r(d) = 1((i \to j) \in d)$$

Proof

$$1((i \to j) \in d) = \sum_{(k \to l) \in d} 1(kl = ij)$$

$$\Rightarrow r_{ij} = 1(kl = ij)$$

# Additively decomposable functions: Shannon Entropy

Shannon entropy is an expectation

$$\mathbb{E}_{p(d)} \left[ - \log p(d) \right]$$

Set

$$r(d) = - \log p(d)$$

Proof

$$- \log p(d) = - \log(\frac{1}{Z} \prod_{(i \to j) \in d} w_{ij})$$

$$= \log Z - \sum_{(i \to j) \in d} \log w_{ij}$$

$$\Rightarrow r_{ij} = \frac{1}{N} \log Z - \log w_{ij}$$

# First order totals: Recipe

- Rewrite unary marginal as gradient
- Rewrite total as sum over unary parts

# First order totals: Marginal

For any edge $(i \rightarrow j)$,

$$\widetilde{w_{ij}} := \sum_{d \in \mathcal{D}_{ij}} w(d) = \frac{\partial Z}{\partial w_{ij}} w_{ij}, \qquad (2)$$

where $\mathcal{D}_{ij} := \{d \in \mathcal{D} \mid (i \rightarrow j) \in d\}$

▶ Significantly easier to implement via AD than manually

# First order totals: Marginal derivation

$$
\begin{aligned}
\widetilde{w_{ij}} &= \sum_{d \in \mathcal{D}_{ij}} w(d) && \text{(Defn)} \\
&= \sum_{d \in \mathcal{D}_{ij}} \prod_{(i' \to j') \in d} w_{i'j'} && \text{(Defn of } w(d)) \\
&= w_{ij} \sum_{d \in \mathcal{D}_{ij}} \prod_{\substack{(i' \to j') \in \\ d \smallsetminus \{(i \to j)\}}} w_{i'j'} && \text{(Pull out } w_{ij}) \\
&= w_{ij} \frac{\partial}{\partial w_{ij}} \sum_{d \in \mathcal{D}} \prod_{(i' \to j') \in d} w_{i'j'} && \text{(Sum-prod deriv trick)} \\
&= \frac{\partial Z}{\partial w_{ij}} w_{ij} && \text{(Defn of } Z)
\end{aligned}
$$

## First order totals: Total

For any additively decomposable function $r: \mathcal{D} \mapsto \mathbb{R}^R$,

$$\bar{r} = \sum_{(i \to j) \in \mathcal{E}} \widetilde{w_{ij}} r_{ij} = \sum_{(i \to j) \in \mathcal{E}} \frac{\partial Z}{\partial w_{ij}} w_{ij} r_{ij} \tag{3}$$

Proof.

$$
\begin{aligned}
\bar{r} &= \sum_{d \in \mathcal{D}} w(d) r(d) && \text{(Defn)} \\
&= \sum_{d \in \mathcal{D}} w(d) \sum_{(i \to j) \in d} r_{ij} && \text{(Defn of } r(d)) \\
&= \sum_{d \in \mathcal{D}} \sum_{(i \to j) \in d} w(d) r_{ij} && \text{(Distributive prop)} \\
&= \sum_{(i \to j) \in \mathcal{E}} \sum_{d \in \mathcal{D}_{ij}} w(d) r_{ij} && \text{(Commutative prop)} \\
&= \sum_{(i \to j) \in \mathcal{E}} \widetilde{w_{ij}} r_{ij} && \text{(Defn of } \widetilde{w_{ij}})
\end{aligned}
$$

# First order totals: Algorithm

1: **def** $\mathrm{T}_1\Big(w \colon \mathcal{E} \mapsto \mathbb{R}, r \colon \mathcal{E} \mapsto \mathbb{R}^R\Big)$ :

2:     ▷ *Compute first-order total; requires* $\mathcal{O}\Big(N^3 R'\Big)$ *time,* $\mathcal{O}\Big(N^2+R\Big)$ *space.*

3:     Compute all $\widetilde{w_{ij}}$ via AD in $\mathcal{O}(N^3)^1$

4:     $\overline{r} \leftarrow \displaystyle\sum_{(i \to j) \in \mathcal{E}} \widetilde{w_{ij}} r_{ij}$         ▷ $\mathcal{O}\Big(N^2 R'\Big)$

5:     **return** $\overline{r}$

---

[1]Constant factor of $|L|$ runtime with AD

# Second-order additively decomposable functions

- Allows us to decompose into functions of pairs of edges and combine with (unnormalized) binary marginals

$$t(d) = r(d)s(d)^\top \tag{4}$$

- Outer product of two additively decomposable functions, $r \colon \mathcal{D} \mapsto \mathbb{R}^R$ and $s \colon \mathcal{D} \mapsto \mathbb{R}^S$
- $t(d) \in \mathbb{R}^{R \times S}$ is generally a matrix.

# Second-order additively decomposable functions: Binary marginals

Binary marginals are expectations

$$p((i \to j), (k \to l) \in d) = \mathbb{E}_{p(d)} \left[ 1((i \to j), (k \to l) \in d) \right]_{ij,kl}$$
$$= \left[ \sum_d p(d) r(d) s(d)^\top \right]_{ij,kl}$$

Set
$$r(d) = 1((i \to j) \in d), \qquad s(d) = 1((k \to l) \in d)$$

which we know are additively decomposable

# Second-order totals: Recipe

- Rewrite binary marginal as Hessian of partition function
- Rewrite grad of intermediate total as Hessian-vector product
- Rewrite second-order total as Jacobian-matrix product or Hessian-matrix product

# Second-order totals: Marginal

Unnormalized binary marginals

$$\widetilde{w_{ij,kl}} := \sum_{d \in \mathcal{D}_{ij,kl}} w(d) = \frac{\partial^2 Z}{\partial w_{ij} \, \partial w_{kl}} w_{ij} w_{kl} \tag{5}$$

where $\mathcal{D}_{ij,kl} := \{d \in \mathcal{D} \mid (i \rightarrow j) \in d, (k \rightarrow l) \in d\}$

# Second-order totals: Marginal

Same tricks as first-order marginals

$$\widetilde{w_{ij,kl}} = \sum_{d \in \mathcal{D}_{ij,kl}} w(d)$$

$$= \sum_{d \in \mathcal{D}_{ij,kl}} \prod_{(k' \to l') \in d} w_{k'l'}$$

$$= w_{ij} w_{kl} \frac{\partial^2}{\partial w_{ij} \partial w_{kl}} \sum_{d \in \mathcal{D}} \prod_{(i' \to j') \in d} w_{i'j'}$$

$$= \frac{\partial^2 Z}{\partial w_{ij} \partial w_{kl}} w_{ij} w_{kl}$$

## Second-order totals: Grad of intermediate total

Write $\nabla \bar{r}$ in terms of Hessian $\nabla^2 Z$, as a step towards writing $\bar{t}$ in terms of $\nabla^2 Z$

$$w_{ij} \frac{\partial \bar{r}}{\partial w_{ij}}$$

$$= w_{ij} \frac{\partial}{\partial w_{ij}} \left( \sum_{(k \to l) \in \mathcal{E}} \frac{\partial Z}{\partial w_{kl}} w_{kl} r_{kl} \right)$$

$$= w_{ij} \frac{\partial Z}{\partial w_{ij}} r_{ij} + w_{ij} \sum_{(k \to l) \in \mathcal{E}} \frac{\partial^2 Z}{\partial w_{ij} \partial w_{kl}} w_{kl} r_{kl}$$

$$= \widetilde{w_{ij}} r_{ij} + \sum_{(k \to l) \in \mathcal{E}} \widetilde{w_{ij,kl}} r_{kl}$$

# Second-order totals: Total

We can both

- ▶ Write total $\bar{t}$ in terms of Jacobian $\frac{\partial \bar{r}}{\partial w_{ij}}$ or gradients $\nabla \bar{r}_n$
    - ▶ Run backward for each dimension or $\bar{r}$: $O(RN^3 + R^2N^2)$
- ▶ Write total $\bar{t}$ in terms of Hessian $\nabla^2 Z$
    - ▶ Hessian is large $O(N^4)$ entries
    - ▶ This takes $O(N^4)$ time, but can be optimized to $O(N^2 + RS)$ space

## Second-order totals: Total

$$\overline{t} = \sum_{d \in \mathcal{D}} w(d) r(d) s(d)^\top$$

$$= \sum_{d \in \mathcal{D}} w(d) r(d) \sum_{(i \to j) \in d} s_{ij}^\top$$

$$= \sum_{d \in \mathcal{D}} \sum_{(i \to j) \in d} w(d) r(d) s_{ij}^\top$$

$$= \sum_{(i \to j) \in \mathcal{E}} \sum_{d \in \mathcal{D}_{ij}} w(d) r(d) s_{ij}^\top$$

$$= \sum_{(i \to j) \in \mathcal{E}} w_{ij} \frac{\partial}{\partial w_{ij}} \left( \sum_{d \in \mathcal{D}} w(d) r(d) \right) s_{ij}^\top$$

$$= \sum_{(i \to j) \in \mathcal{E}} w_{ij} \frac{\partial \overline{r}}{\partial w_{ij}} s_{ij}^\top$$

## Second-order totals: Jacobian version

1: **def** $\mathrm{T}_2^{\mathrm{v}}\Big(w\colon \mathcal{E} \mapsto \mathbb{R}, r\colon \mathcal{E} \mapsto \mathbb{R}^R, s\colon \mathcal{E} \mapsto \mathbb{R}^S\Big)$ :

2:   ▷ *Compute second-order total with gradient-vector products; requires* $\mathcal{O}(R(N^3+N^2R'+N^2S'))$ *time,* $\mathcal{O}\Big(N^2R+R\,S\Big)$ *space.*

3:   **for** $n = 1 \ldots R$ :        ▷ $\mathcal{O}\Big(R(N^3 + N^2R')\Big)$

4:     Compute $\nabla \overline{r}_n$ using reverse-mode AD on $[\mathrm{T}_1(w, r)]_n$

5:   ▷ *Requires* $\mathcal{O}\Big(N^2RS'\Big)$

6:   **return** $\displaystyle\sum_{(i \to j) \in \mathcal{E}} \frac{\partial \overline{r}}{\partial w_{ij}} w_{ij} s_{ij}^{\top}$

▶ Recommended approach

# Second-order totals: Hessian version

1: **def** $\mathrm{T}_2^{\mathrm{h}}\Big(w\colon \mathcal{E} \mapsto \mathbb{R}, r\colon \mathcal{E} \mapsto \mathbb{R}^R, s\colon \mathcal{E} \mapsto \mathbb{R}^S\Big):$

2: ▷ *Compute second-order total by materializing Hessian; requires* $\mathcal{O}\Big(N^4 R'S'\Big)$ *time,* $\mathcal{O}\Big(N^2 + R\,S\Big)$ *space.*

3: Compute all $\widetilde{w_{ij}}$ via $\nabla Z$

4: Compute all $\widetilde{w_{ij,kl}}$ via $\nabla^2 Z$ (can compute this efficiently)

5: ▷ *Requires* $\mathcal{O}\Big(N^4 R'S'\Big)$

6: **return** $\displaystyle\sum_{(i \to j) \in \mathcal{E}} \widetilde{w_{ij}}\, r_{ij} s_{ij}{}^\top + \sum_{(k \to l) \in \mathcal{E}} \widetilde{w_{ij,kl}}\, r_{ij} s_{kl}{}^\top$

- ▶ Better than previous algo at dealing with sparsity, i.e. if $R$ and $S$ are really big and sparse
- ▶ Can come up with a 3rd algo that has the best of both worlds by computing the second term involving the Hessian-matrix product more efficiently

# Refs I

[1] Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In *International Conference on Machine Learning*, pages 2113–2122, 2015.

[2] Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back-propagation for bilevel optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 1723–1732, 2019.