

# Moment Distributionally Robust Tree Structured Prediction

Yeshu Li <sup>1</sup>   Danyal Saeed <sup>1</sup>   Xinhua Zhang <sup>1</sup>   Brian D. Ziebart <sup>1</sup>  
Kevin Gimpel <sup>2</sup>

<sup>1</sup>Department of Computer Science  
University of Illinois at Chicago

<sup>2</sup>Toyota Technological Institute at Chicago

Nov. 2022



# Statistical Learning: Test vs Train

- The ultimate goal is to perform well on test data, in terms of a performance metric  $L$ , typically piecewise-constant and discontinuous, usually intractable to optimize on training data

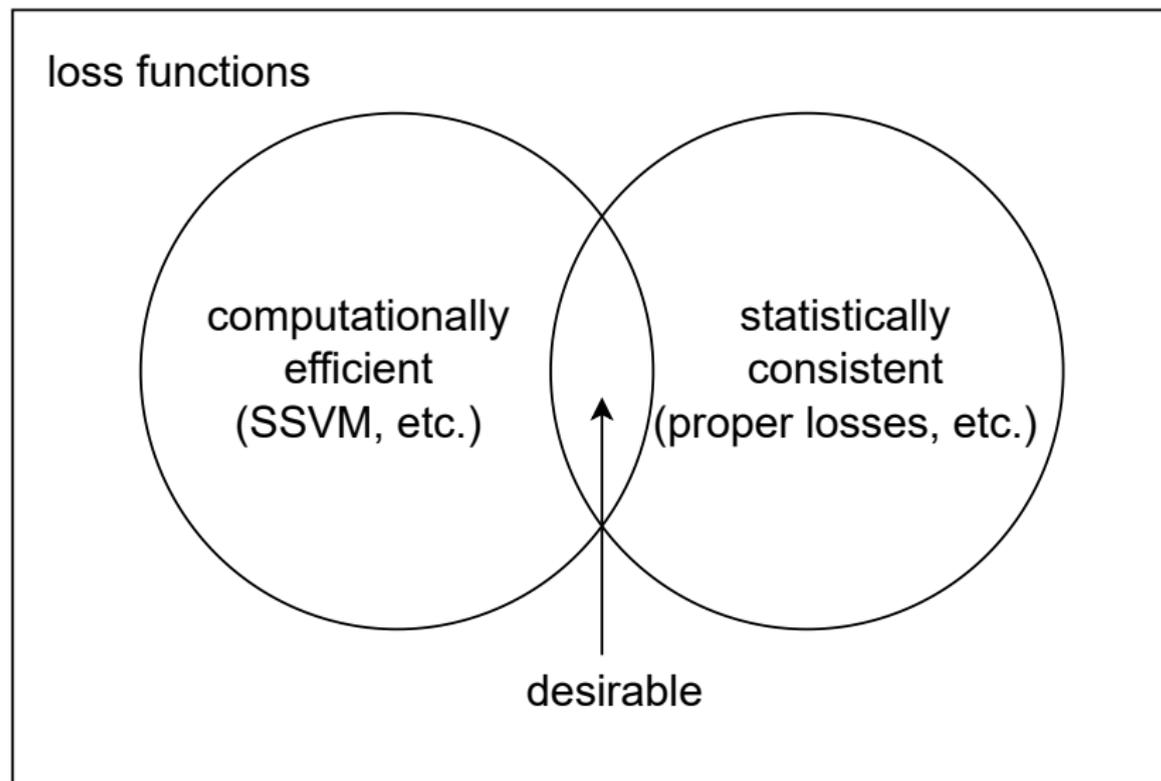
$$\min_{f \in \mathcal{F}} \sum_{i=1}^{n_{\text{train}}} L(f(\mathbf{x}_i), \mathbf{y}_i)$$

- Resort to surrogate loss and mappings for easier optimization

$$\min_{\theta} \sum_{i=1}^{n_{\text{train}}} S(g_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

- Might lead to discrepancy between training and testing objectives

# Design of Losses



# Fisher Consistency

## Definition (Fisher consistency)

Given a space  $\mathcal{H}$  of all measurable functions  $\mathcal{X} \mapsto \mathcal{P}(\mathcal{Y})$ , a space  $\mathcal{G}$  of all measurable functions  $\mathcal{X} \mapsto \mathcal{T}$ , and a surrogate loss function  $S : \mathcal{T} \times \mathcal{P}(\mathcal{Y}) \mapsto \mathbb{R}_+$ , we say that the surrogate loss  $S$  is Fisher consistent with respect to the loss  $L : \mathcal{P}(\mathcal{Y}) \times \mathcal{P}(\mathcal{Y}) \mapsto \mathbb{R}_+$  if there exists a function  $d : \mathcal{T} \mapsto \mathcal{P}(\mathcal{Y})$  such that for every probability distribution  $\mathbb{P}$  on  $\mathcal{X} \times \mathcal{Y}$ , every minimizer  $g^*$  of the surrogate risk reaches Bayes optimal risk:

$$R_{\mathbb{P}}^S(g^*) = \min_{g \in \mathcal{G}} R_{\mathbb{P}}^S(g) \implies R_{\mathbb{P}}^L(d \circ g^*) = \min_{h \in \mathcal{H}} R_{\mathbb{P}}^L(h).$$

## Existing Methods

- Empirical risk minimization ( $K$ -best lists (Smith and Eisner, 2006), automatic differentiation (Stoyanov and Eisner, 2012))

$$\min_{\theta} \sum_{i=1}^m \sum_{\mathbf{y} \in \mathcal{Y}} L(\mathbf{y}_i, \mathbf{y}) \Pr_{\theta}(\mathbf{y} | \mathbf{x}_i) = \sum_{i=1}^m \sum_{\mathbf{y} \in \mathcal{Y}} L(\mathbf{y}_i, \mathbf{y}) \frac{\exp(\text{score}_{\theta}(\mathbf{x}_i, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\text{score}_{\theta}(\mathbf{x}_i, \mathbf{y}'))}$$

- Conditional log-likelihood (deep learning (Dozat and Manning, 2017), matrix tree theorem (Koo et al., 2007))

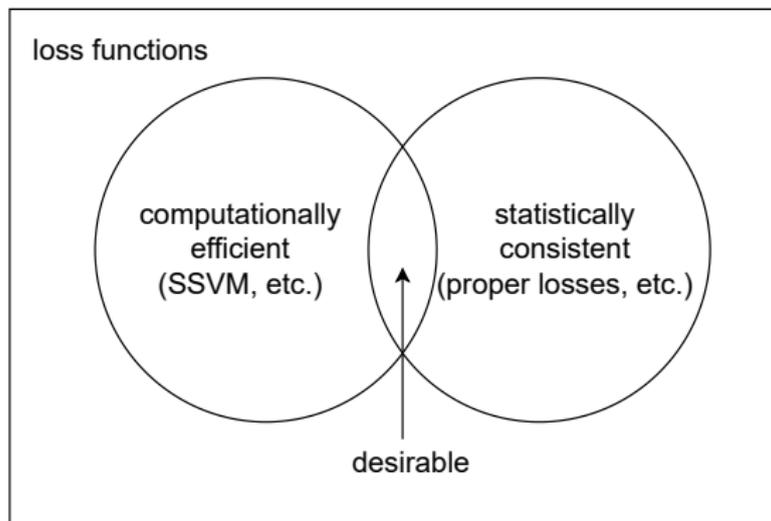
$$\min_{\theta} \sum_{i=1}^m -\text{score}_{\theta}(\mathbf{x}_i, \mathbf{y}_i) + \log \sum_{\mathbf{y}} \exp(\text{score}_{\theta}(\mathbf{x}_i, \mathbf{y}))$$

- Max-margin (Structured SVM (Taskar et al., 2004))

$$\min_{\theta} \sum_{i=1}^m -\text{score}_{\theta}(\mathbf{x}_i, \mathbf{y}_i) + \max_{\mathbf{y}} (\text{score}_{\theta}(\mathbf{x}_i, \mathbf{y}) + L(\mathbf{y}_i, \mathbf{y}))$$

# Motivation

- Existing methods are either **NOT Fisher consistent** or **NOT convex**
- Existing methods rely on instinctive **regularization** to combat overfitting, without good probabilistic interpretation

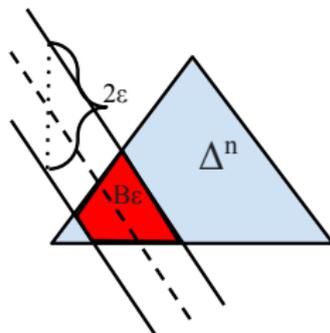


# Proposed Method (Primal)

Distributionally robust optimization based on moment divergence

$$\inf_{\tilde{\mathbb{P}}} \sup_{\check{\mathbb{P}} \in B_\varepsilon(\tilde{\mathbb{P}})} \mathbb{E}_{\check{\mathbb{P}}_{\mathbf{X}}, \hat{\mathbb{P}}_{\check{\mathbf{Y}}|\mathbf{X}}, \check{\mathbb{P}}_{\check{\mathbf{Y}}|\mathbf{X}}} L(\hat{\mathbf{Y}}, \check{\mathbf{Y}})$$

- $B_\varepsilon(\tilde{\mathbb{P}}) := \{\check{\mathbb{P}} : \|\mathbb{E}_{\check{\mathbb{P}}_{\mathbf{X}, \check{\mathbf{Y}}}} \phi(\mathbf{X}, \check{\mathbf{Y}}) - \mathbb{E}_{\tilde{\mathbb{P}}_{\mathbf{X}, \mathbf{Y}}} \phi(\mathbf{X}, \mathbf{Y})\| \leq \varepsilon, \check{\mathbb{P}}_{\mathbf{X}} = \tilde{\mathbb{P}}_{\mathbf{X}}\}$  is a **convex compact ambiguity set** quantifying the uncertainty about the underlying true distribution
- $\phi(\mathbf{x}, \mathbf{y}) := \sum_{f \in \mathcal{F}} \phi(\mathbf{x}, \mathbf{y}_f)$  is a feature mapping decomposable over factors  $f$



# Proposed Method (Dual)

## Proposition (dual norm regularization)

The distributionally robust tree structured prediction problem based on moment divergence can be rewritten as

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{\mathbb{P}_{\mathbf{X}, \mathbf{Y}}^{\text{emp}}} \underbrace{\min_{\mathbb{P}} \max_{\mathbb{Q}} \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{Y}}|\mathbf{X}}, \mathbb{Q}_{\check{\mathbf{Y}}|\mathbf{X}}} \ell(\hat{\mathbf{Y}}, \check{\mathbf{Y}}) + \boldsymbol{\theta}^{\top} (\phi(\mathbf{X}, \check{\mathbf{Y}}) - \phi(\mathbf{X}, \mathbf{Y})) + \varepsilon \|\boldsymbol{\theta}\|_*}_{\ell_{\text{adv}}(\boldsymbol{\theta}, (\mathbf{X}, \mathbf{Y}))},$$

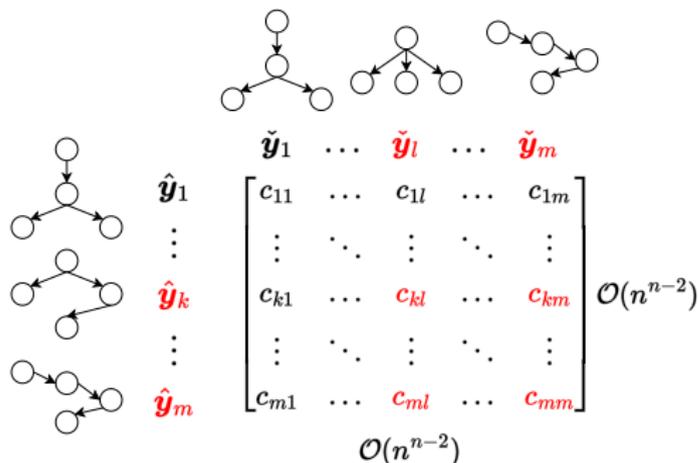
where  $\boldsymbol{\theta} \in \mathbb{R}^d$  is the vector of Lagrangian multipliers and  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$ .

# Algorithm: Game Theory

- The inner minimax problems are independent conditioned on  $\mathbf{X}$

$$\min_{\theta} \mathbb{E}_{\mathbb{P}_{\mathbf{X}, \mathbf{Y}}^{\text{emp}}} \min_{\mathbb{P}} \max_{\mathbb{Q}} \mathbb{E}_{\mathbb{P}_{\hat{\mathbf{Y}}|\mathbf{X}}, \mathbb{Q}_{\check{\mathbf{Y}}|\mathbf{X}}} \text{payoff}(\hat{\mathbf{Y}}, \check{\mathbf{Y}})$$

- Use constraint generation to find a Nash equilibrium  $\mathbb{P}_{\hat{\mathbf{Y}}|\mathbf{X}}^*, \mathbb{Q}_{\check{\mathbf{Y}}|\mathbf{X}}^*$
- No convergence guarantees



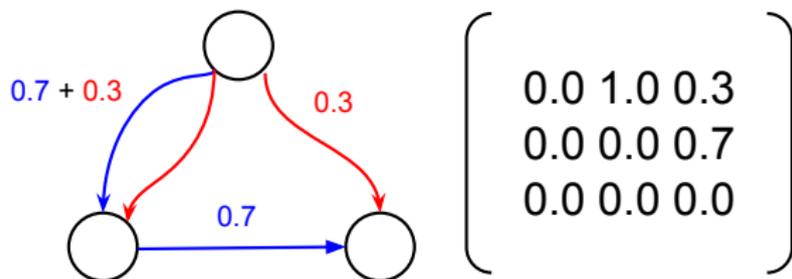
## Algorithm: Marginal Distribution

- Assuming additive  $L$ , rearrange the order of optimization variables

$$\max_{\mathbf{q}^{(i)} \in \mathcal{A}_{\text{arb}}} \min_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m \min_{\mathbf{p} \in \mathcal{A}_{\text{arb}}} (\mathbf{q}^{(i)} - \mathbf{p}_{\text{emp}}^{(i)})^{\top} \Phi^{(i)} \boldsymbol{\theta} - \langle \mathbf{p}, \mathbf{q}^{(i)} \rangle$$

$$+ \frac{\mu}{2} \|\mathbf{p}\|_2^2 - \frac{\mu}{2} \|\mathbf{q}^{(i)}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

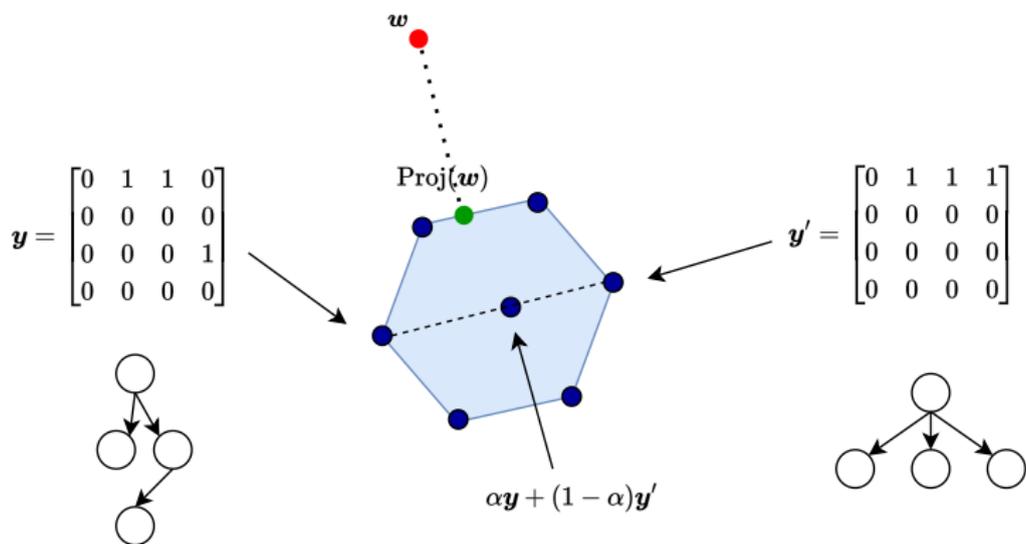
- Converges to the global optimum
- Requires an oracle for projection onto the arborescence (directed spanning tree) polytopes



# Projection on the Arborescence Polytope

Constrained quadratic programming

$$\min_{\mathbf{x} \in \mathcal{A}_{\text{arb}}} f(\mathbf{x}) := \|\mathbf{x} - \mathbf{w}\|_2^2$$



# Projection on the Arborescence Polytope

## Proposed solutions

- Frank-Wolfe (Frank and Wolfe, 1956)
  - Based on minimum weight directed spanning trees
  - Low per-iteration cost
  - Sub-linear convergence rate  $\mathcal{O}(\frac{1}{\epsilon})$
- Alternating Direction Method of Multipliers (Boyd et al., 2011)
  - Based on a compact representation of the first-order arborescence polytope (Friesen, 2019)
  - Higher per-iteration cost
  - Linear convergence rate  $\mathcal{O}(\log \frac{1}{\epsilon})$

# Excess True Risk Bound

## Theorem

Given  $m$  samples, a non-negative loss  $\ell(\cdot, \cdot)$  such that  $|\ell(\cdot, \cdot)| \leq K$ , a feature function  $\phi(\cdot, \cdot)$  such that  $\|\phi(\cdot, \cdot)\| \leq B$ , a positive ambiguity level  $\varepsilon > 0$ , then, for any  $\rho \in (0, 1]$ , with a probability at least  $1 - \rho$ , the following excess true worst-case risk bound holds:

$$\max_{\mathbb{Q} \in \mathcal{B}(\mathbb{P}^{\text{true}})} R_{\mathbb{Q}}^L(\theta_{\text{emp}}^*) - \max_{\mathbb{Q} \in \mathcal{B}(\mathbb{P}^{\text{true}})} R_{\mathbb{Q}}^L(\theta_{\text{true}}^*) \leq \frac{4KB}{\varepsilon\sqrt{m}} \left( 1 + \frac{3}{2} \sqrt{\frac{\ln(4/\rho)}{2}} \right),$$

where  $\theta_{\text{emp}}^*$  and  $\theta_{\text{true}}^*$  are the optimal parameters learned under  $\mathbb{P}^{\text{emp}}$  and  $\mathbb{P}^{\text{true}}$  respectively. The original risk of  $\theta$  under  $\mathbb{Q}$  is

$R_{\mathbb{Q}}^L(\theta) := \mathbb{E}_{\mathbb{Q}_{\mathbf{X}, \mathbf{Y}}, \mathbb{P}_{\hat{\mathbf{Y}}|\mathbf{X}}^{\theta}} \ell(\hat{\mathbf{Y}}, \mathbf{Y})$  with Bayes prediction

$\mathbb{P}_{\mathbf{Y}|\mathbf{X}}^{\theta} \in \arg \min_{\mathbb{P}} \max_{\mathbb{Q}} \mathbb{E}_{\mathbb{Q}_{\check{\mathbf{Y}}|\mathbf{X}} \mathbb{P}_{\check{\mathbf{Y}}|\mathbf{X}}} \ell(\hat{\mathbf{Y}}, \check{\mathbf{Y}}) + \theta^T \phi(\mathbf{x}, \check{\mathbf{Y}}).$

# Fisher Consistency

Corollary of Proposition C.2 in Nowak et al. (2020)

When  $\varepsilon = 0$ ,  $\ell_{\text{adv}}$  is Fisher consistent with respect to  $\ell$ . Namely,  $\mathbb{P}_{\hat{Y}|\mathbf{X}}^{\theta_{\text{true}}^*}$  is the probabilistic prediction made by the Bayes optimal decision rule.

# Results

- Datasets: PTB, CTB, UD
- Baseline: Deep Biaffine Attention (Dozat and Manning, 2017)
- Marginal: use marginal probabilities with full gradient
- Stochastic: use marginal probabilities with mini-batch
- Game: constraint generation (double oracle)

Table 1: Comparison of mean unlabeled attachment score (UAS) and execution time under different training set sizes. Time refers to the CPU time taken to finish one gradient descent step. Statistically significant differences compared to *BiAF* are marked with † (paired t-test,  $p < 0.05$ ). The best UAS are highlighted in bold.

Method	Time (s)	PTB				CTB				UD Dutch			
		m = 10	50	100	1000	m = 10	50	100	1000	m = 10	50	100	1000
BiAF (baseline)	0.34	93.48	<b>96.87</b>	<b>96.95</b>	<b>97.16</b>	88.45	90.89	91.15	<b>91.70</b>	90.86	93.80	94.15	94.98
Marginal (ours)	0.28	94.51†	96.81†	96.92	97.12	89.19†	91.03†	91.27	91.67	<b>92.41†</b>	94.22†	94.50†	<b>95.15†</b>
Stochastic (ours)	2.72	<b>94.62†</b>	96.81	96.93	97.14	<b>89.27†</b>	91.03†	<b>91.27</b>	91.66	92.40†	94.23†	94.47	95.14†
Game (ours)	7.25	94.51†	96.86	96.92	97.08†	89.22†	<b>91.06†</b>	91.22	91.57†	92.32†	<b>94.34†</b>	<b>94.59†</b>	95.01

# More Results

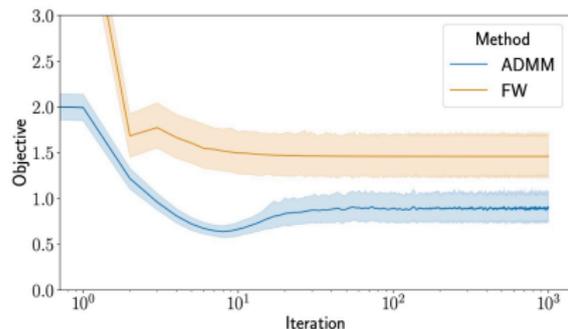


Figure 1: Convergence of Alternating Direction Method of Multipliers (ADMM) and Frank-Wolfe (FW) for random points with 95% confidence intervals.

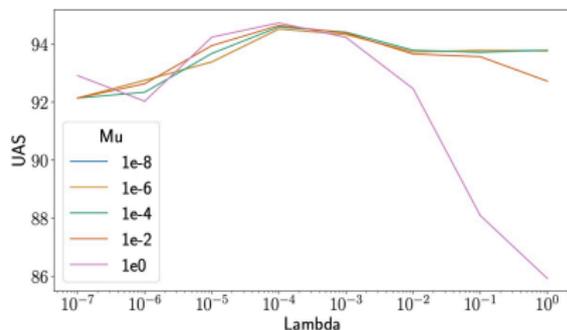
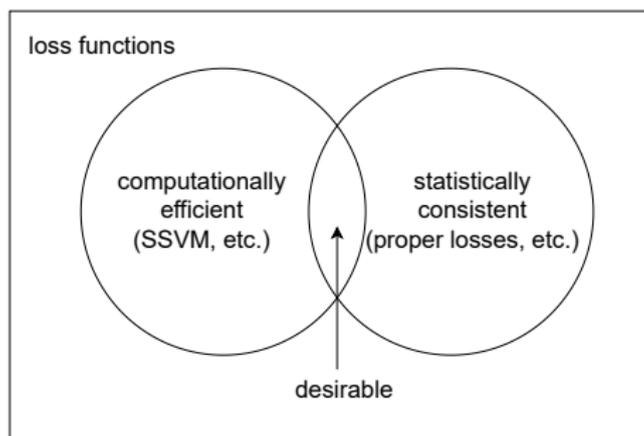


Figure 2: The best unlabeled attachment score (UAS) with the Marginal algorithm as  $\mu$  and  $\lambda$  vary in logarithmic scales.

# Conclusion

- Tree structured prediction from first principles in DRO
  - Dependency, directed, undirected, higher-order trees
- Generalization bounds and Fisher consistency
- Efficient projection oracles on arborescence polytopes
- Code available at <https://github.com/DanielLeee/drtreesp>



# Future Work

- General structured prediction tasks
- Other ambiguity sets in distributionally robust optimization
- End-to-end representation learning (code available)

$$\frac{\partial}{\partial \psi(\mathbf{x})} \ell_{\text{adv}} \in \frac{1}{B} \sum_{i=1}^B (\mathbf{q}^{(i)*} - \mathbf{p}_{\text{emp}}^{(i)*})$$

# References

- Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J., et al. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122.
- Dozat, T. and Manning, C. D. (2017). Deep biaffine attention for neural dependency parsing. In *International Conference on Learning Representations*.
- Frank, M. and Wolfe, P. (1956). An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110.
- Friesen, M. (2019). *Extended formulations for higher order polytopes in combinatorial optimization*. PhD thesis, Otto von Guericke University Magdeburg.
- Koo, T., Globerson, A., Carreras Pérez, X., and Collins, M. (2007). Structured prediction models via the matrix-tree theorem. In *Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 141–150.
- Nowak, A., Bach, F., and Rudi, A. (2020). Consistent structured prediction with max-min margin markov networks. In *International Conference on Machine Learning*, pages 7381–7391. PMLR.
- Smith, D. A. and Eisner, J. (2006). Minimum risk annealing for training log-linear models. In *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*, pages 787–794.
- Stoyanov, V. and Eisner, J. (2012). Minimum-risk training of approximate crf-based nlp systems. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 120–130.
- Taskar, B., Klein, D., Collins, M., Koller, D., and Manning, C. D. (2004). Max-margin parsing. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 1–8.

Thank you!

Q & A