



# Latent Structure Models for NLP

André Martins

Instituto de Telecomunicações & IST & Unbabel

Tsvetomila Mihaylova

Instituto de Telecomunicações

Nikita Nangia

NYU

Vlad Niculae

Instituto de Telecomunicações

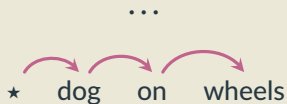
# I. Introduction

# Structured prediction and NLP

- **Structured prediction:** a machine learning framework for predicting structured, constrained, and interdependent outputs
- **NLP** deals with *structured* and *ambiguous* textual data:
  - machine translation
  - speech recognition
  - syntactic parsing
  - semantic parsing
  - information extraction
  - ...

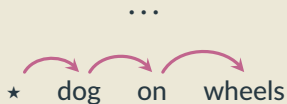
# Examples of structure in NLP

## Dependency parsing

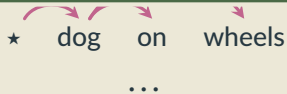


# Examples of structure in NLP

Dependency parsing



Exponentially many parse trees!  
Cannot enumerate.



# Examples of structure in NLP

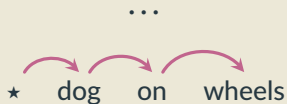
## POS tagging

VERB    PREP    NOUN  
dog    on    wheels

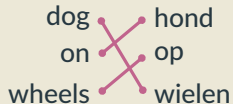
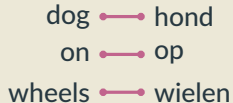
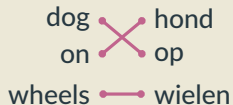
NOUN    PREP    NOUN  
dog    on    wheels

NOUN    DET    NOUN  
dog    on    wheels

## Dependency parsing



## Word alignments



# NLP 5 years ago:

## Structured prediction and pipelines



# NLP 5 years ago:

## Structured prediction and pipelines

- Big pipeline systems, connecting different structured predictors, trained separately
- **Advantages:** fast and simple to train, can rearrange pieces 😊



# NLP 5 years ago:

## Structured prediction and pipelines

- Big pipeline systems, connecting different structured predictors, trained separately
- **Advantages:** fast and simple to train, can rearrange pieces 😊
- **Disadvantage:** linguistic annotations required for each component 😓

# NLP 5 years ago:

## Structured prediction and pipelines

- Big pipeline systems, connecting different structured predictors, trained separately
- **Advantages:** fast and simple to train, can rearrange pieces 😊
- **Disadvantage:** linguistic annotations required for each component 😓
- **Bigger disadvantage:** error propagates through the pipeline 💩

# NLP today:

## End-to-end training



# NLP today:

## End-to-end training

- Forget pipelines—train everything from scratch!
- No more error propagation or linguistic annotations! 🎉

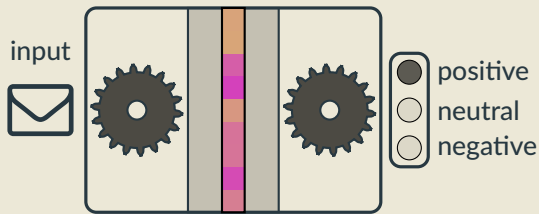
# NLP today:

## End-to-end training

- Forget pipelines—train everything from scratch!
- No more error propagation or linguistic annotations! 🎉
- Treat everything as *latent*! 🙌

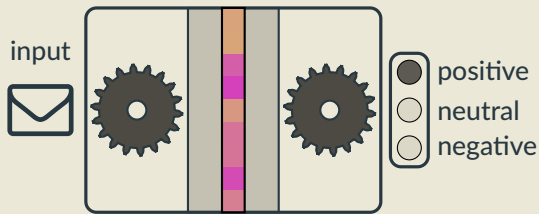
# Representation learning

- Uncover hidden representations useful for the *downstream task*.
- Neural networks are well-suited for this: *deep computation graphs*.



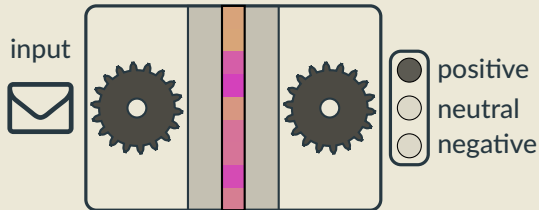
# Representation learning

- Uncover hidden representations useful for the *downstream task*.
- Neural networks are well-suited for this: *deep computation graphs*.
- Neural representations are unstructured, inscrutable.  
Language data has underlying structure!



# Latent structure models

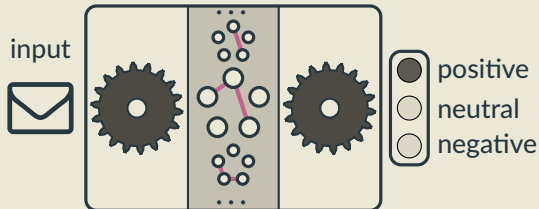
- Seek *structured* hidden representations instead!





# Latent structure models

- Seek *structured* hidden representations instead!



# Latent structure models aren't so new!

- They have a very long history in NLP:
  - IBM Models for SMT (latent word alignments) [Brown et al., 1993]
  - HMMs [Rabiner, 1989]
  - CRFs with hidden variables [Quattoni et al., 2007]
  - Latent PCFGs [Petrov and Klein, 2008, Cohen et al., 2012]
- Trained with EM, spectral learning, method of moments, ...
- Often, very strict assumptions (e.g. strong factorizations)
- Today, neural networks opened up some new possibilities!

# Why do we love latent structure models?

- The inferred latent variables can bring us some **interpretability**
- They offer a way of injecting prior knowledge as a **structured bias**
- Hopefully: Higher predictive power with fewer model parameters

# Why do we love latent structure models?

- The inferred latent variables can bring us some **interpretability**
- They offer a way of injecting prior knowledge as a **structured bias**
- Hopefully: Higher predictive power with fewer model parameters
  - smaller carbon footprint!

# What this tutorial is about:

- Discrete, combinatorial latent structures
- Often the structure is inspired by some linguistic intuition
- We'll cover both:
  - RL methods (structure built incrementally, reward coming from downstream task)
  - ... vs end-to-end differentiable approaches (global optimization, marginalization)
  - stochastic computation graphs
  - ... vs deterministic graphs.
- All plugged in *discriminative* neural models.

# This tutorial is *not* about:

- It's not about continuous latent variables
- It's not about deep generative learning
- We won't cover GANs, VAEs, etc.
- There are (very good) recent tutorials on deep variational models for NLP:
  - “Variational Inference and Deep Generative Models” (Schulz and Aziz, ACL 2018)
  - “Deep Latent-Variable Models for Natural Language” (Kim, Wiseman, Rush, EMNLP 2018)

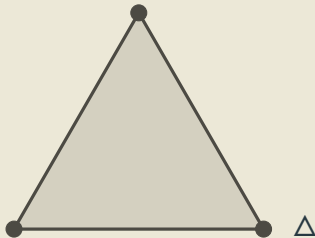
**Background**

# Unstructured vs structured

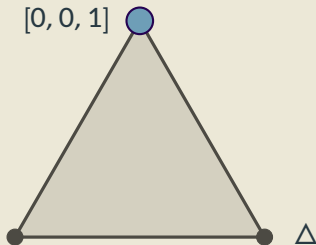
- To better explain the math, we'll often backtrack to *unstructured* models (where the latent variable is a categorical) before jumping to the *structured* ones



# The unstructured case: Probability simplex



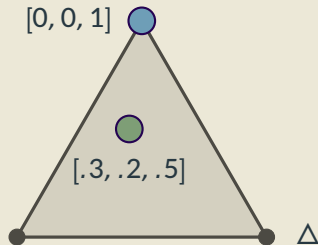
# The unstructured case: Probability simplex



- Each vertex is an *indicator vector*, representing one class:

$$\mathbf{z}_c = [0, \dots, 0, \underbrace{1}_{c^{\text{th}} \text{ position}}, 0, \dots, 0].$$

# The unstructured case: Probability simplex



- Each vertex is an *indicator vector*, representing one class:

$$\mathbf{z}_c = [0, \dots, 0, \underbrace{1}_{c^{\text{th}} \text{ position}}, 0, \dots, 0].$$

- Points inside are *probability vectors*, a convex combination of classes:

$$\mathbf{p} \geq \mathbf{0}, \quad \sum_c p_c = 1.$$

# What's the analogous of $\Delta$ for a structure?

- A structured object  $\mathbf{z}$  can be represented as a *bit vector*.

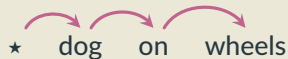
# What's the analogous of $\Delta$ for a structure?

- A structured object  $\mathbf{z}$  can be represented as a *bit vector*.
- Example:
  - a dependency tree can be represented a  $O(L^2)$  vector indexed by arcs
  - each entry is 1 iff the arc belongs to the tree
  - **structural constraints:** not all bit vectors represent valid trees!

# What's the analogous of $\Delta$ for a structure?

- A structured object  $\mathbf{z}$  can be represented as a *bit vector*.
- Example:
  - a dependency tree can be represented a  $O(L^2)$  vector indexed by arcs
  - each entry is 1 iff the arc belongs to the tree
  - **structural constraints:** not all bit vectors represent valid trees!

$$\mathbf{z}_1 = [1, 0, 0, 0, 1, 0, 0, 0, 1]$$



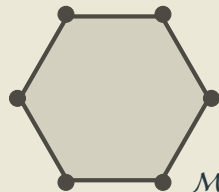
$$\mathbf{z}_2 = [0, 0, 1, 0, 0, 1, 1, 0, 0]$$



$$\mathbf{z}_3 = [1, 0, 0, 0, 1, 0, 0, 1, 0]$$

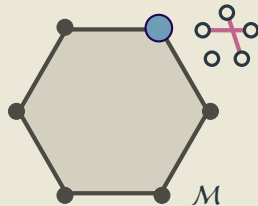


# The structured case: Marginal polytope



# The structured case: Marginal polytope

- Each vertex corresponds to one such *bit* vector  $\mathbf{z}$





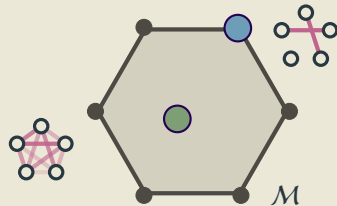
# The structured case: Marginal polytope

- Each vertex corresponds to one such *bit* vector  $\mathbf{z}$
- Points inside correspond to *marginal distributions*: convex combinations of structured objects

$$\boldsymbol{\mu} = \underbrace{p_1 \mathbf{z}_1 + \dots + p_N \mathbf{z}_N}_{\text{exponentially many terms}}, \quad \mathbf{p} \in \Delta.$$

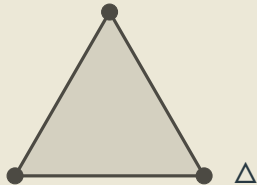
$$\begin{aligned} p_1 &= 0.2, & \mathbf{z}_1 &= [1, 0, 0, 0, 1, 0, 0, 0, 1] \\ p_2 &= 0.7, & \mathbf{z}_2 &= [0, 0, 1, 0, 0, 1, 1, 0, 0] \\ p_3 &= 0.1, & \mathbf{z}_3 &= [1, 0, 0, 0, 1, 0, 0, 1, 0] \end{aligned}$$

$$\Rightarrow \boldsymbol{\mu} = [.3, 0, .7, 0, .3, .7, .7, .1, .2].$$

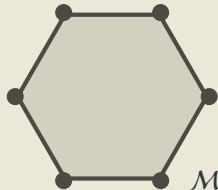


# Unstructured vs Structured

- Unstructured case: simplex  $\Delta$

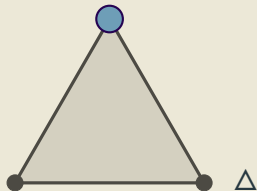


- Structured case: marginal polytope  $\mathcal{M}$

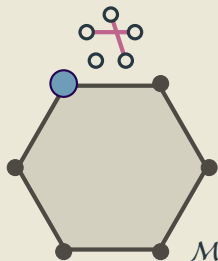


# Unstructured vs Structured

- Unstructured case: simplex  $\Delta$

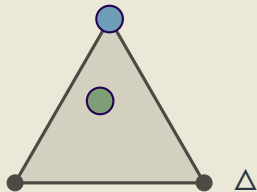


- Structured case: marginal polytope  $\mathcal{M}$

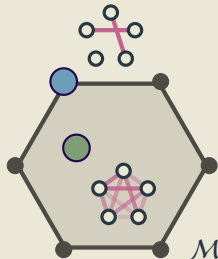


# Unstructured vs Structured

- Unstructured case: simplex  $\Delta$

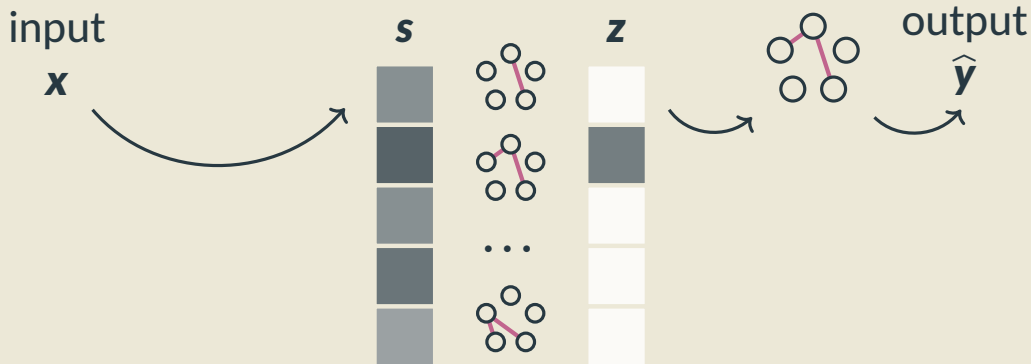


- Structured case: marginal polytope  $\mathcal{M}$



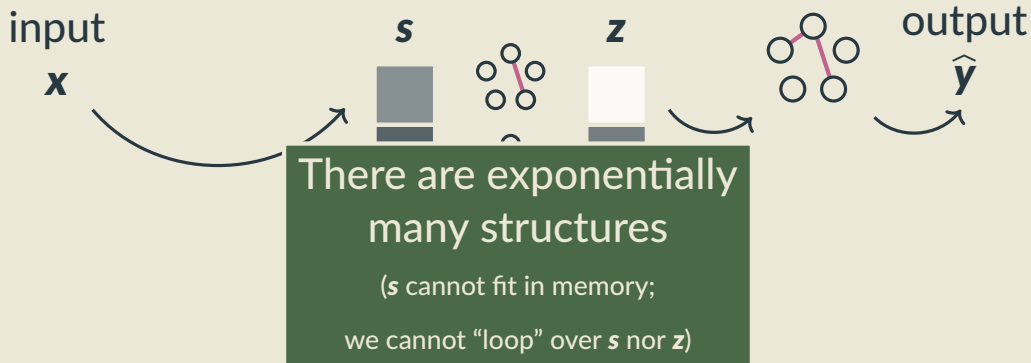
# Computing the most likely structure

is a very high-dimensional argmax

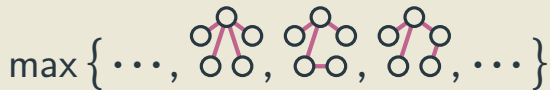
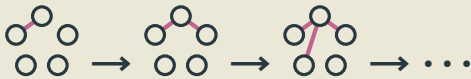


# Computing the most likely structure

is a very high-dimensional argmax



# Dealing with the combinatorial explosion



## 1. Incremental structures

- Build structure **greedily**, as sequence of discrete choices (e.g., shift-reduce).
- Scores (partial structure, action) tuples.
- **Advantages:** flexible, rich histories.
- **Disadvantages:** greedy, local decisions are suboptimal, error propagation.

## 2. Factorization into parts

- Optimizes **globally** (e.g. Viterbi, Chu-Liu-Edmonds, Kuhn-Munkres).
- Scores smaller parts.
- **Advantages:** optimal, elegant, can handle hard & global constraints.
- **Disadvantages:** strong assumptions.

# The challenge of discrete choices.

$$z = 1$$

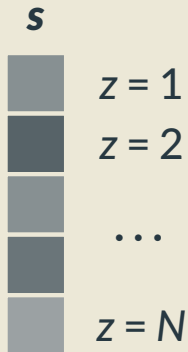
$$z = 2$$

...

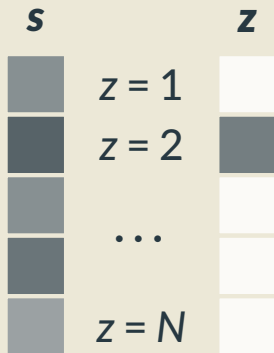
$$z = N$$



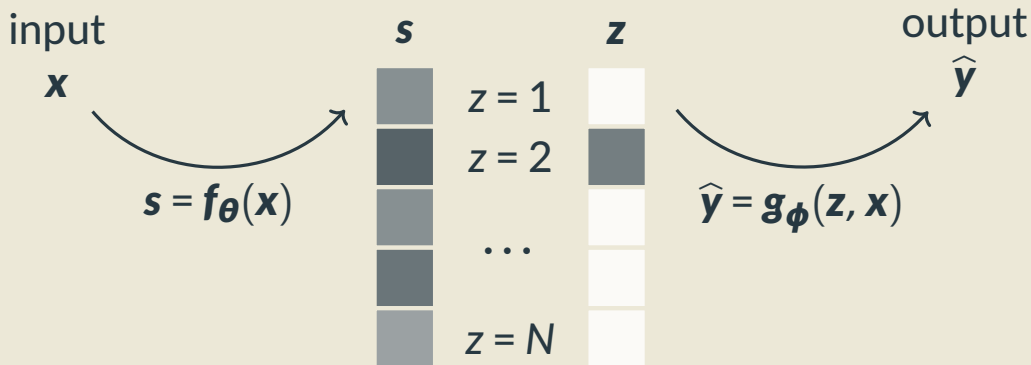
# The challenge of discrete choices.



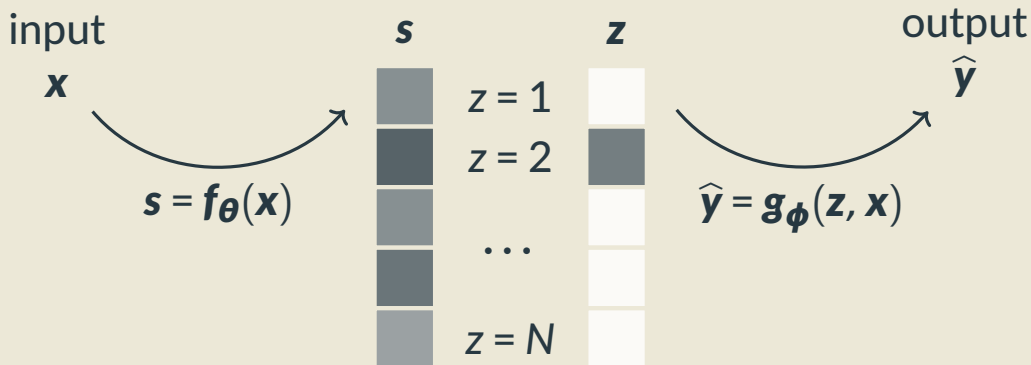
# The challenge of discrete choices.



# The challenge of discrete choices.

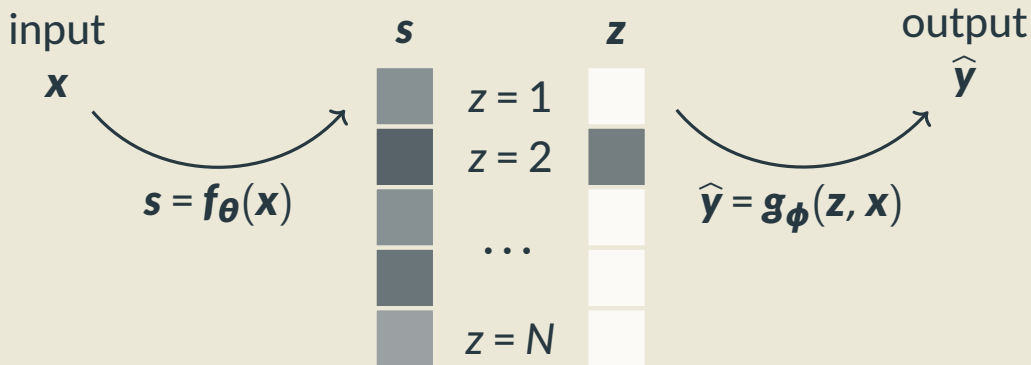


# The challenge of discrete choices.



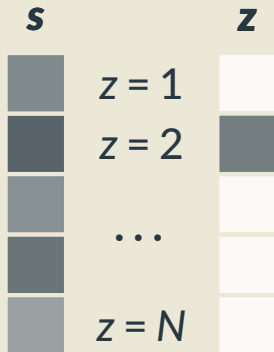
$$\frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{w}} = ?$$

# The challenge of discrete choices.



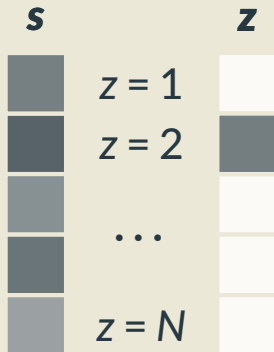
$$\frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{w}} = ? \quad \text{or, essentially,} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

# Discrete mappings are “flat”



$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

# Discrete mappings are “flat”



$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

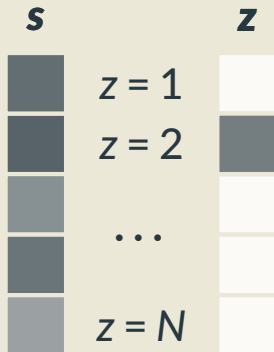
# Discrete mappings are “flat”



$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

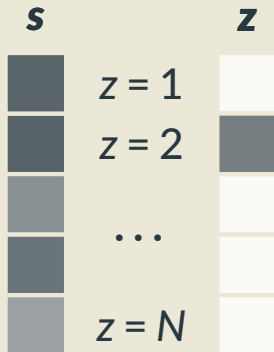


# Discrete mappings are “flat”



$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

# Discrete mappings are “flat”



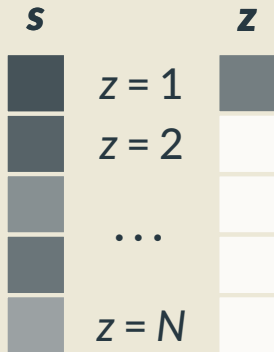
$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

# Discrete mappings are “flat”



$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

# Discrete mappings are “flat”



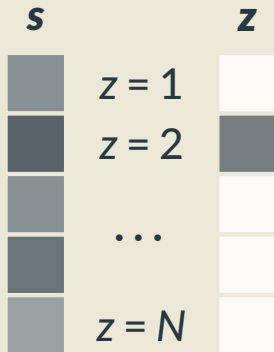
$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

# Discrete mappings are “flat”

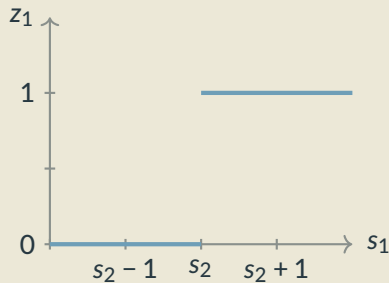


$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = ?$$

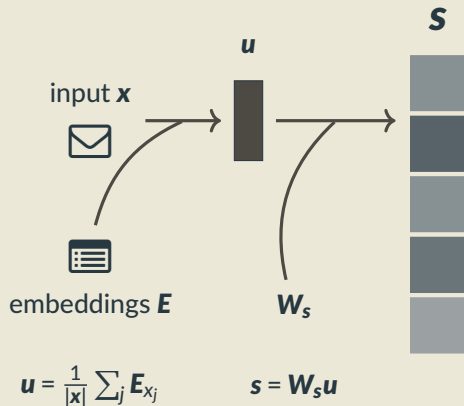
# Argmax



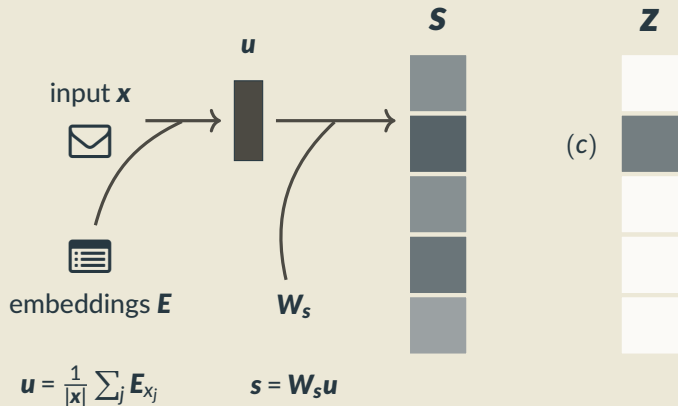
$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = \mathbf{0}$$



# Example: Regression with latent categorization



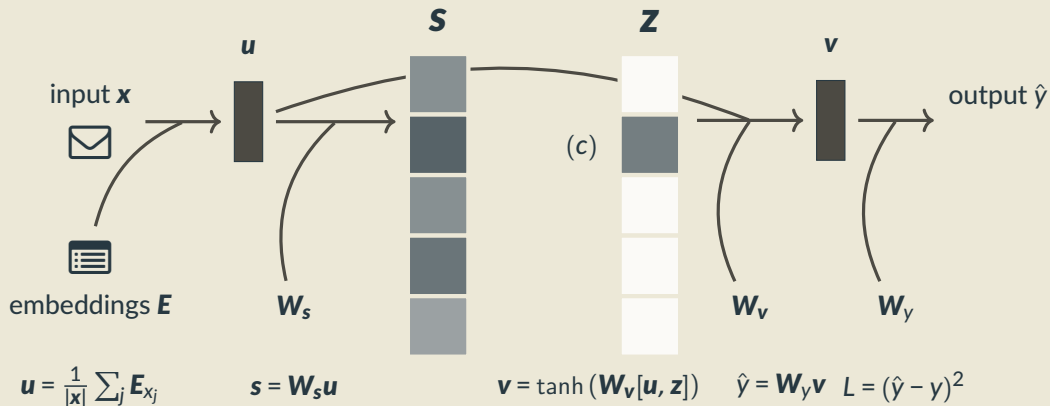
# Example: Regression with latent categorization



predict topic  $c$  ( $z = e_c$ )

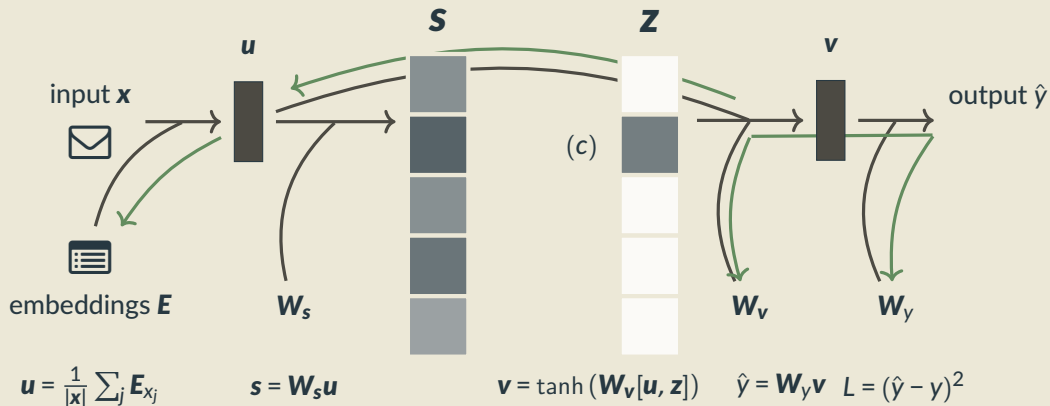


# Example: Regression with latent categorization

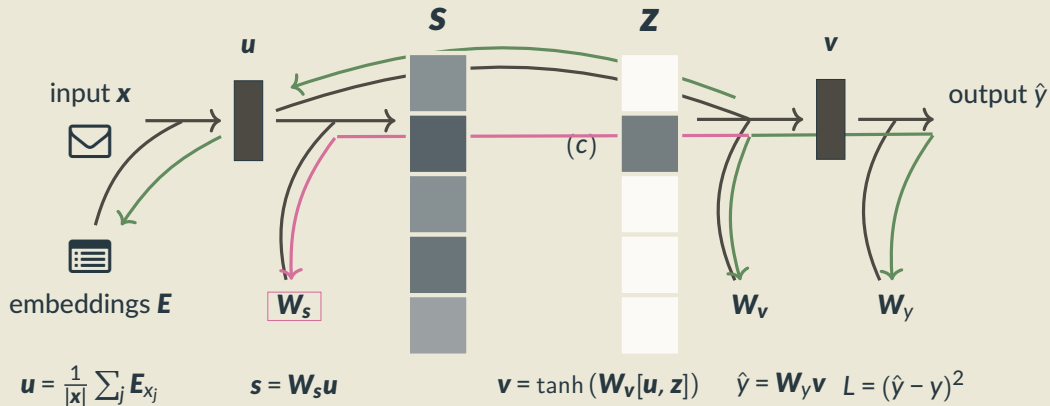


predict topic  $c$  ( $z = e_c$ )

# Example: Regression with latent categorization

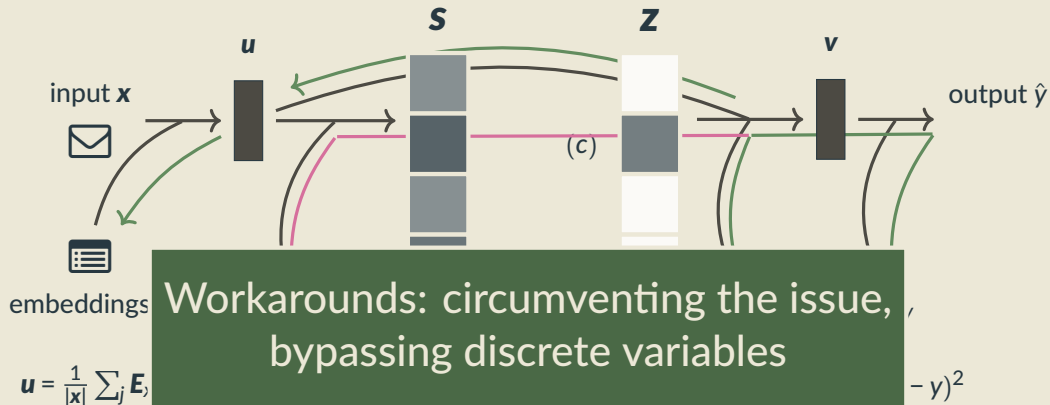


# Example: Regression with latent categorization

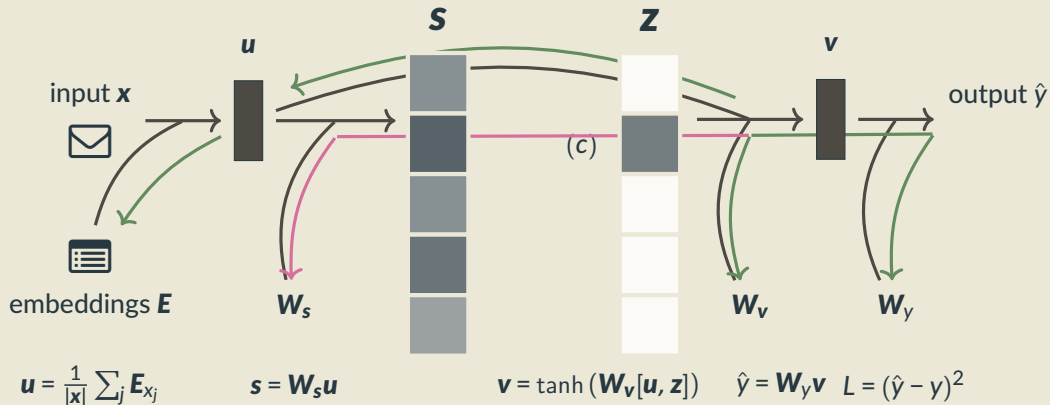


$$\frac{\partial L}{\partial W_s} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v} \frac{\partial v}{\partial z} \underbrace{\frac{\partial z}{\partial s}}_{\equiv 0} \frac{\partial s}{\partial W_s}$$

# Example: Regression with latent categorization

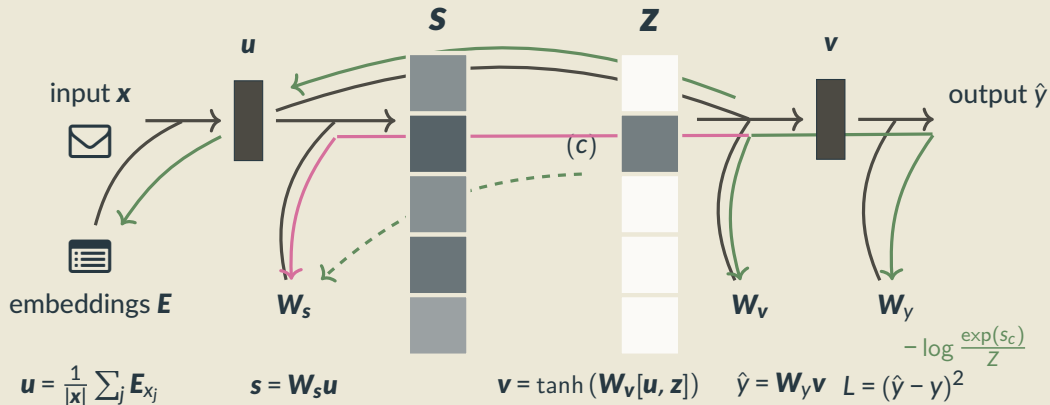


# Example: Regression with latent categorization



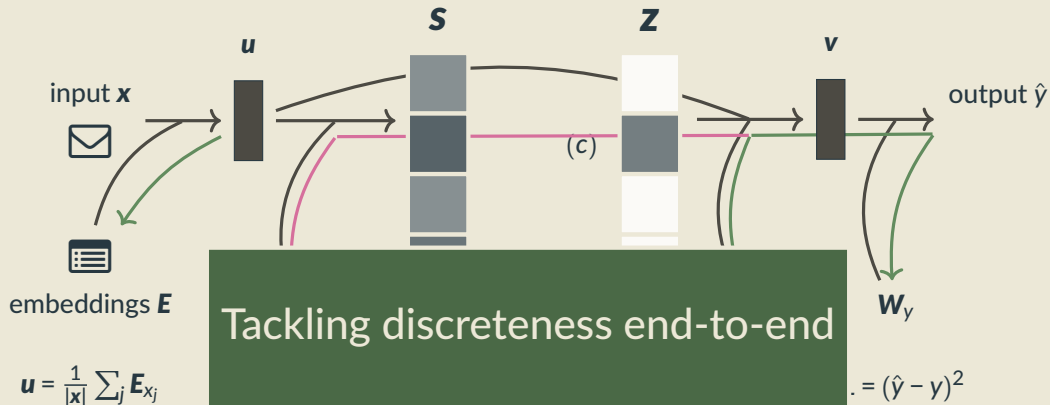
Option 1. Pretrain latent classifier  $W_s$

# Example: Regression with latent categorization

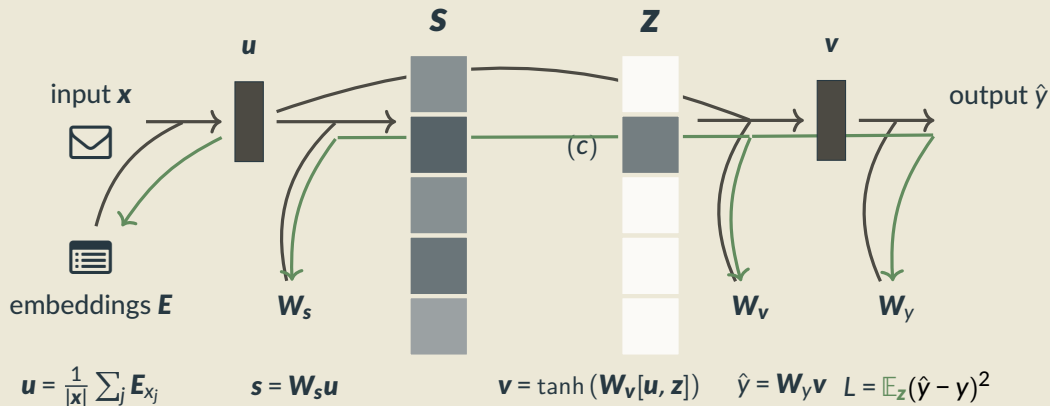


Option 2. Multi-task learning

# Example: Regression with latent categorization



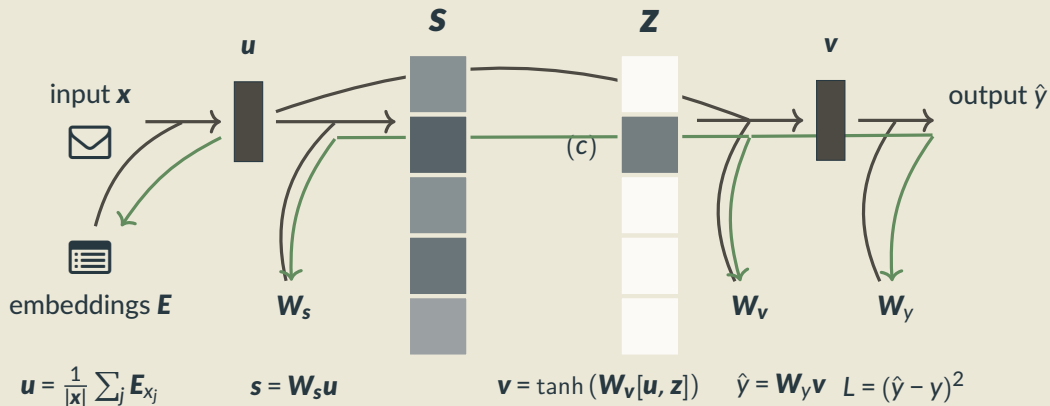
# Example: Regression with latent categorization



Option 3. Stochasticity!  $\frac{\partial \mathbb{E}_z(\hat{y}(z) - y)^2}{\partial W_s} \neq 0$

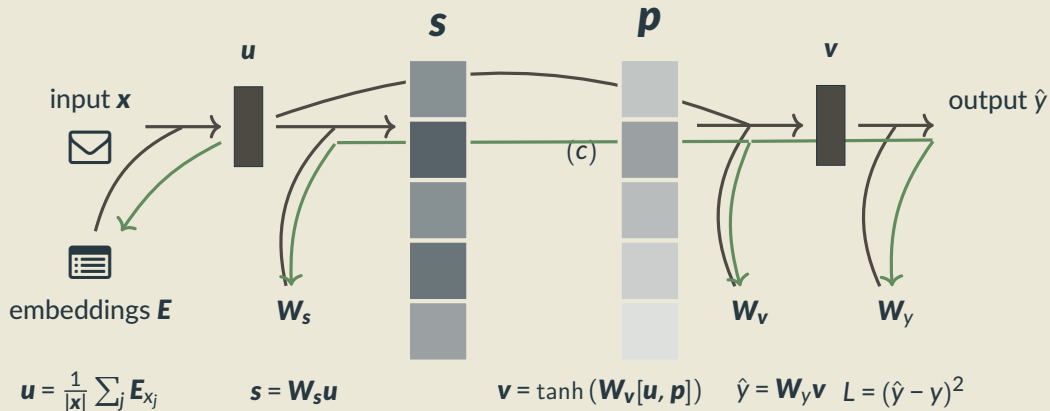


# Example: Regression with latent categorization



Option 4. Gradient surrogates (e.g. straight-through,  $\frac{\partial z}{\partial s} \leftarrow I$ )

# Example: Regression with latent categorization



Option 5. Continuous relaxation (e.g. softmax)

# Dealing with discrete latent variables

1. Pre-train external classifier
2. Multi-task learning
3. Stochastic latent variables
4. Gradient surrogates
5. Continuous relaxation

# Dealing with discrete latent variables

1. Pre-train external classifier
2. Multi-task learning
3. Stochastic latent variables (Part 2)
4. Gradient surrogates (Part 3)
5. Continuous relaxation (Part 4)

# Roadmap of the tutorial

- Part 1: Introduction ✓
- Part 2: Reinforcement learning
- Part 3: Gradient surrogates

## *Coffee Break*

- Part 4: End-to-end differentiable models
- Part 5: Conclusions

## **II. Reinforcement Learning Methods**

# Latent structure via marginalization

- Given a sentence-label pair  $(x, y)$  and its **known** parse tree  $\mathbf{z}$ ,

# Latent structure via marginalization

- Given a sentence-label pair  $(x, y)$  and its **known** parse tree  $\mathbf{z}$ , we can make a prediction  $\hat{y}(\mathbf{z}; x)$



# Latent structure via marginalization

- Given a sentence-label pair  $(x, y)$  and its **known** parse tree  $\mathbf{z}$ , we can make a prediction  $\hat{y}(\mathbf{z}; x)$  and incur a loss,

$$L(\hat{y}(\mathbf{z}; x), y)$$

# Latent structure via marginalization

- Given a sentence-label pair  $(x, y)$  and its **known** parse tree  $\mathbf{z}$ , we can make a prediction  $\hat{y}(\mathbf{z}; x)$  and incur a loss,

$$L(\hat{y}(\mathbf{z}; x), y) \text{ or simply } L(\mathbf{z})$$

# Latent structure via marginalization

- Given a sentence-label pair  $(x, y)$  and its **known** parse tree  $\mathbf{z}$ , we can make a prediction  $\hat{y}(\mathbf{z}; x)$  and incur a loss,

$$L(\hat{y}(\mathbf{z}; x), y) \text{ or simply } L(\mathbf{z})$$

# Latent structure via marginalization

- Given a sentence-label pair  $(x, y)$  and its **known** parse tree  $\mathbf{z}$ , we can make a prediction  $\hat{y}(\mathbf{z}; x)$  and incur a loss,

$$L(\hat{y}(\mathbf{z}; x), y) \text{ or simply } L(\mathbf{z})$$

- But we don't know  $\mathbf{z}$ !

# Latent structure via marginalization

- Given a sentence-label pair  $(x, y)$  and its **known** parse tree  $\mathbf{z}$ , we can make a prediction  $\hat{y}(\mathbf{z}; x)$  and incur a loss,

$$L(\hat{y}(\mathbf{z}; x), y) \text{ or simply } L(\mathbf{z})$$

- But we don't know  $\mathbf{z}$ !
- In this section:  
we jointly learn a structure prediction model  $\pi_{\theta}(\mathbf{z} \mid x)$

# Latent structure via marginalization

- Given a sentence-label pair  $(x, y)$  and its **known** parse tree  $\mathbf{z}$ , we can make a prediction  $\hat{y}(\mathbf{z}; x)$  and incur a loss,

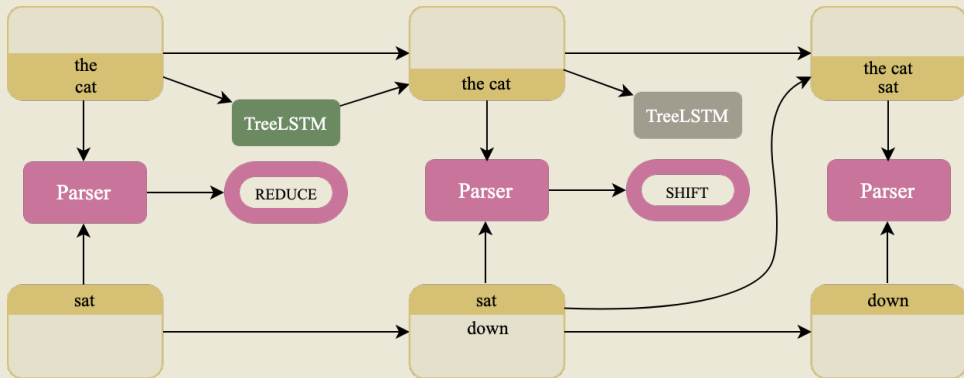
$$L(\hat{y}(\mathbf{z}; x), y) \text{ or simply } L(\mathbf{z})$$

- But we don't know  $\mathbf{z}$ !
- In this section:
  - we jointly learn a structure prediction model  $\pi_{\theta}(\mathbf{z} | x)$  by optimizing the **expected loss**,

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

**But first, supervised  
SPINN**

# Stack-augmented Parser-Interpreter Neural-Network





# Stack-augmented Parser-Interpreter Neural-Network

- Join learning: Combines a constituency parser and a sentence representation model.

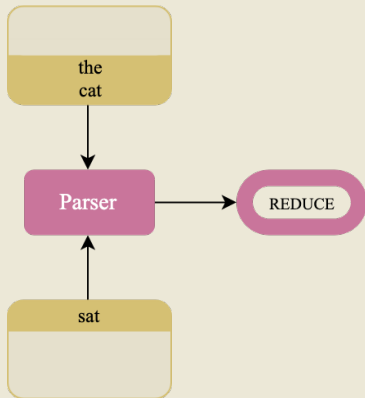
# Stack-augmented Parser-Interpreter Neural-Network

- Join learning: Combines a constituency parser and a sentence representation model.
- The parser is a transition-based **shift-reduce** parser. It looks at top two elements of stack and top element of the buffer.  $\rightarrow f_{\theta}(x)$

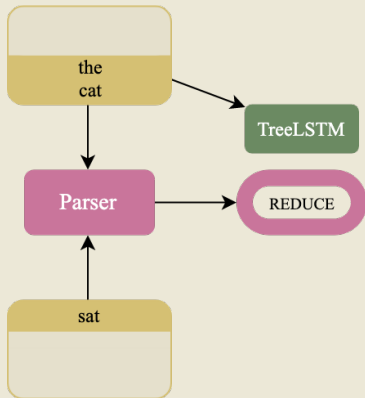
# Stack-augmented Parser-Interpreter Neural-Network

- Join learning: Combines a constituency parser and a sentence representation model.
- The parser is a transition-based **shift-reduce** parser. It looks at top two elements of stack and top element of the buffer.  $\rightarrow f_{\theta}(x)$
- **TreeLSTM** combines top two elements of the stack when the parser chooses the REDUCE action.

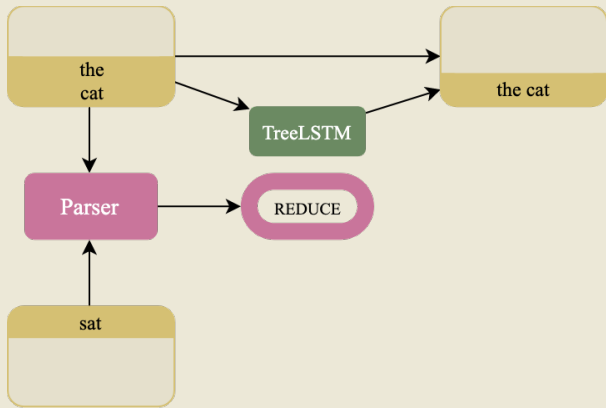
# Stack-augmented Parser-Interpreter Neural-Network



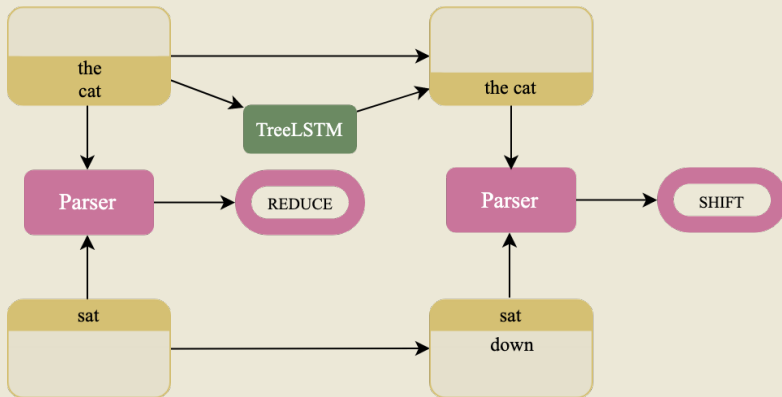
# Stack-augmented Parser-Interpreter Neural-Network



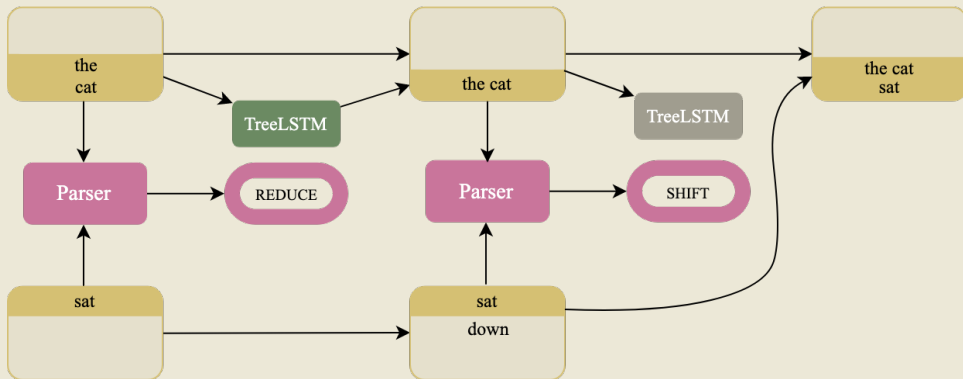
# Stack-augmented Parser-Interpreter Neural-Network



# Stack-augmented Parser-Interpreter Neural-Network

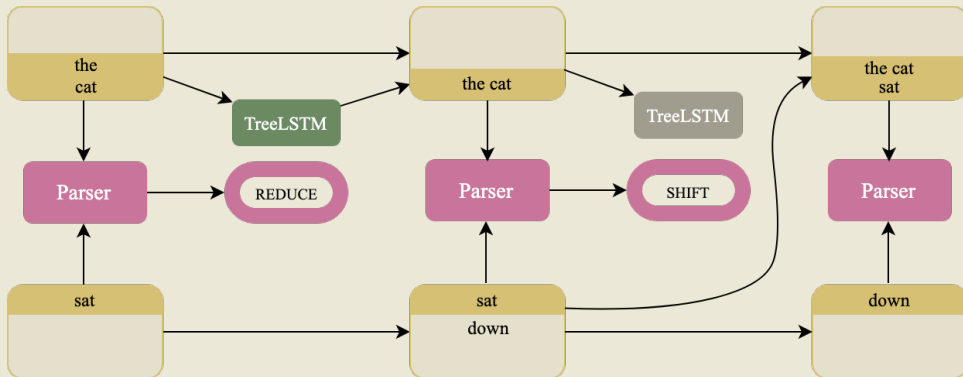


# Stack-augmented Parser-Interpreter Neural-Network





# Stack-augmented Parser-Interpreter Neural-Network



# Shift-Reduce parsing

We can write a shift-reduce style parse as a sequence of Bernoulli random variables,

$$\mathbf{z} = \{z_1, \dots, z_{2n-1}\}$$

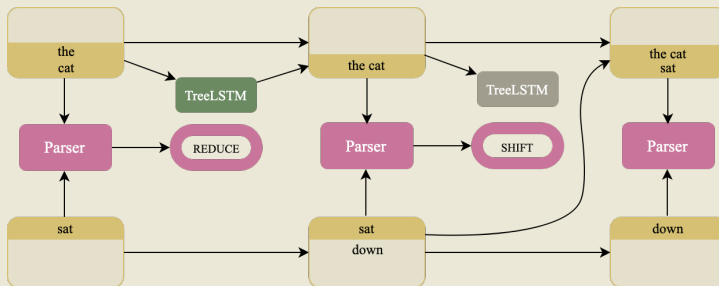
where,  $z_j \in \{0, 1\} \ \forall j \in [1, 2n - 1]$

# Shift-Reduce parsing

A sequence of Bernoulli trials but with conditional dependence,

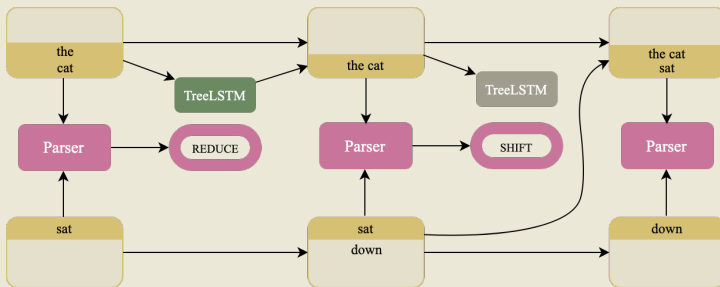
$$p(z_1, z_2, \dots, z_{2n-1}) = \prod_{j=1}^{2n-1} p(z_j \mid z_{<j})$$

# Latent structure learning with SPINN



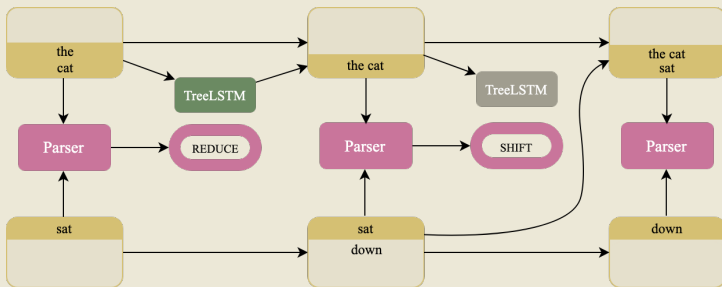
# Latent structure learning with SPINN

- But now, remove syntactic supervision from SPINN.



# Latent structure learning with SPINN

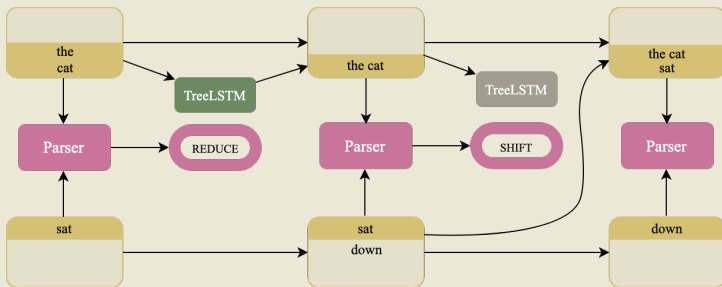
- But now, remove syntactic supervision from SPINN.



- We model the parse,  $\mathbf{z}$ , as a latent variable with our parser as the score function estimator,  $f_{\theta}(x)$ .

# Latent structure learning with SPINN

- But now, remove syntactic supervision from SPINN.



- We model the parse,  $\mathbf{z}$ , as a latent variable with our parser as the score function estimator,  $f_{\theta}(\mathbf{x})$ .
- With shift-reduce parsing, we're making discrete decisions  $\Rightarrow$  REINFORCE as a "natural" solution.

# Unsupervised SPINN



# Unsupervised SPINN

No syntactic supervision.

Only reward is from the downstream task.

We only get this reward after parsing the full sentence.

# SPINN with REINFORCE

Some basic terminology,

- The action space is  $z \in \{\text{SHIFT}, \text{REDUCE}\}$ , and  $\mathbf{z}$  is a sequence of actions.

# SPINN with REINFORCE

Some basic terminology,

- The action space is  $z \in \{\text{SHIFT}, \text{REDUCE}\}$ , and  $\mathbf{z}$  is a sequence of actions.
- Training parser network parameters,  $\mathbf{w}$  with REINFORCE

# SPINN with REINFORCE

Some basic terminology,

- The action space is  $z \in \{\text{SHIFT}, \text{REDUCE}\}$ , and  $\mathbf{z}$  is a sequence of actions.
- Training parser network parameters,  $\mathbf{w}$  with REINFORCE
- The state,  $\mathbf{s}$ , is the top two elements of the stack and the top element of the buffer.

# SPINN with REINFORCE

## Some basic terminology,

- The action space is  $\mathbf{z} \in \{\text{SHIFT}, \text{REDUCE}\}$ , and  $\mathbf{z}$  is a sequence of actions.
- Training parser network parameters,  $\mathbf{w}$  with REINFORCE
- The state,  $\mathbf{s}$ , is the top two elements of the stack and the top element of the buffer.
- Learning a policy network  $\pi(\mathbf{z} \mid \mathbf{s}; \mathbf{w})$

# SPINN with REINFORCE

## Some basic terminology,

- The action space is  $z \in \{\text{SHIFT}, \text{REDUCE}\}$ , and  $\mathbf{z}$  is a sequence of actions.
- Training parser network parameters,  $\mathbf{w}$  with REINFORCE
- The state,  $\mathbf{s}$ , is the top two elements of the stack and the top element of the buffer.
- Learning a policy network  $\pi(\mathbf{z} \mid \mathbf{s}; \mathbf{w})$
- Maximize the reward, where  $\mathcal{R}$  is performance on the downstream task like sentence classification.

# SPINN with REINFORCE

## Some basic terminology,

- The action space is  $z \in \{\text{SHIFT}, \text{REDUCE}\}$ , and  $\mathbf{z}$  is a sequence of actions.
- Training parser network parameters,  $\mathbf{w}$  with REINFORCE
- The state,  $\mathbf{s}$ , is the top two elements of the stack and the top element of the buffer
- Learning
- Maximizing the expected return

NOTE: Only a single reward at the end of parsing.

# Through the looking glass of REINFORCE

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})} [L(\mathbf{z})]$$



# Through the looking glass of REINFORCE

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z} | x)} [L(\mathbf{z})] = \nabla_{\boldsymbol{\theta}} \left[ \sum_{\mathbf{z}} L(\mathbf{z}) \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \right]$$

*(By definition of expectation. How to evaluate?)*

# Through the looking glass of REINFORCE

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z} | x)} [L(\mathbf{z})] = \nabla_{\boldsymbol{\theta}} \left[ \sum_{\mathbf{z}} L(\mathbf{z}) \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \right]$$

*(By definition of expectation. How to evaluate?)*

$$= \sum_{\mathbf{z}} L(\mathbf{z}) \nabla_{\boldsymbol{\theta}} \pi_{\boldsymbol{\theta}}(\mathbf{z} | x)$$

# Through the looking glass of REINFORCE

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z} | x)} [L(\mathbf{z})] = \nabla_{\boldsymbol{\theta}} \left[ \sum_{\mathbf{z}} L(\mathbf{z}) \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \right]$$

*(By definition of expectation. How to evaluate?)*

$$\begin{aligned} &= \sum_{\mathbf{z}} L(\mathbf{z}) \nabla_{\boldsymbol{\theta}} \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \\ &= \sum_{\mathbf{z}} L(\mathbf{z}) \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \end{aligned}$$

*(By Leibniz integral rule for log)*

# Through the looking glass of REINFORCE

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z} | x)} [L(\mathbf{z})] = \nabla_{\boldsymbol{\theta}} \left[ \sum_{\mathbf{z}} L(\mathbf{z}) \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \right]$$

*(By definition of expectation. How to evaluate?)*

$$\begin{aligned} &= \sum_{\mathbf{z}} L(\mathbf{z}) \nabla_{\boldsymbol{\theta}} \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \\ &= \sum_{\mathbf{z}} L(\mathbf{z}) \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{z} | x) \end{aligned}$$

*(By Leibniz integral rule for log)*

$$= \mathbb{E}_{\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z} | x)} [L(\mathbf{z}) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{z} | x)]$$

# SPINN with REINFORCE

Yogatama et al. (2017) uses REINFORCE to train SPINN!

# SPINN with REINFORCE

Yogatama et al. (2017) uses REINFORCE to train SPINN!

However, this vanilla implementation isn't very effective at learning syntax.

# SPINN with REINFORCE

Yogatama et al. (2017) uses REINFORCE to train SPINN!

However, this vanilla implementation isn't very effective at learning syntax.

This model fails to solve a simple toy problem.

# Toy problem: ListOps



[max 2 9 [min 4 7 ] 0 ]



# Toy problem: ListOps

Model	Accuracy		Self F1
	$\mu(\sigma)$	max	
LSTM	71.5 (1.5)	<b>74.4</b>	-
RL-SPINN	60.7 (2.6)	64.8	30.8
Random Trees	-	-	30.1

Model	F1 wrt.			Avg. Depth
	LB	RB	GT	
48D RL-SPINN	<b>64.5</b>	<b>16.0</b>	32.1	<b>14.6</b>
128D RL-SPINN	43.5	13.0	<b>71.1</b>	10.4
GT Trees	41.6	8.8	100.0	9.6
Random Trees	24.0	24.0	24.2	5.2

# Toy problem: ListOps

Model	Accuracy		Self F1
	$\mu(\sigma)$	max	
LSTM	71.5 (1.5)	<b>74.4</b>	-
RL-SPINN	60.7 (2.6)	64.8	30.8
Random Trees	-	-	30.1

But why?

	F1 wrt.			Avg. Depth
	LB	RB	GT	
	<b>4.5</b>	<b>16.0</b>	32.1	<b>14.6</b>
	3.5	13.0	<b>71.1</b>	10.4
GT Trees	41.6	8.8	100.0	9.6
Random Trees	24.0	24.0	24.2	5.2

# SPINN with REINFORCE

This system faces two big problems,

# SPINN with REINFORCE

This system faces two big problems,

1. High variance of gradients
2. Coadaptation

# High variance

- We have a single reward at the end of parsing.

# High variance

- We have a single reward at the end of parsing.
- We are sampling parses from very large search space!  
**Catalan number** of binary trees.

# High variance

- We have a single reward at the end of parsing.
- We are sampling parses from very large search space!  
**Catalan number** of binary trees.
- And the policy is stochastic.

# High variance

So, sometimes the policy lands in a “rewarding state”:

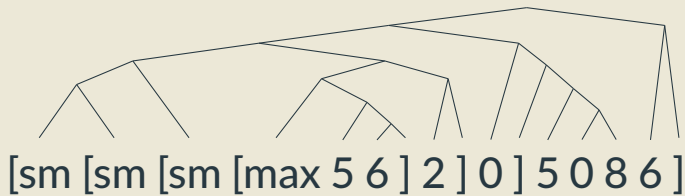


Figure: Truth: 7; Pred: 7



# High variance

Sometimes it doesn't:



Figure: Truth: 6; Pred: 5

# High variance

**Catalan parses** means we need many many samples to lower variance!

# High variance

**Catalan parses** means we need many many samples to lower variance!

Possible solutions,

1. Gradient normalization
2. Control variates, aka baselines

# Control variates

- A simple control variate: moving average of recent rewards

# Control variates

- A simple control variate: moving average of recent rewards
- Parameters are updated using the advantage which is the difference between the reward,  $\mathcal{R}$ , and the baseline prediction.

# Control variates

- A simple control variate: moving average of recent rewards
- Parameters are updated using the advantage which is the difference between the reward,  $\mathcal{R}$ , and the baseline prediction.

So,

$$\nabla \mathbb{E}_{\mathbf{z} \sim \pi(\mathbf{z})} = \mathbb{E}_{\mathbf{z} \sim \pi(\mathbf{z})} [(L(\mathbf{z}) - b(\mathbf{x})) \nabla \log \pi(\mathbf{z})]$$

# Control variates

- A simple control variate: moving average of recent rewards
- Parameters are updated using the advantage which is the difference between the reward,  $\mathcal{R}$ , and the baseline prediction.

So,

$$\nabla \mathbb{E}_{\mathbf{z} \sim \pi(\mathbf{z})} = \mathbb{E}_{\mathbf{z} \sim \pi(\mathbf{z})} [(L(\mathbf{z}) - b(\mathbf{x})) \nabla \log \pi(\mathbf{z})]$$

Which we can do because,

$$\sum_{\mathbf{z}} b(\mathbf{x}) \nabla \pi(\mathbf{z}) = b(\mathbf{x}) \sum_{\mathbf{z}} \nabla \pi(\mathbf{z}) = b(\mathbf{x}) \nabla 1 = 0$$

# Issues with SPINN with REINFORCE

This system faces two big problems,

1. High variance of gradients
2. Coadaptation



# Coadaptation problem

Learning composition function parameters  $\phi$  with backpropagation,  
and parser parameters  $\theta$  with REINFORCE.

# Coadaptation problem

Learning composition function parameters  $\phi$  with backpropagation, and parser parameters  $\theta$  with REINFORCE.

Generally,  $\phi$  will be learned more quickly than  $\theta$ , making it harder to explore the parsing search space and optimize for  $\theta$ .

# Coadaptation problem

Learning composition function parameters  $\phi$  with backpropagation, and parser parameters  $\theta$  with REINFORCE.

Generally,  $\phi$  will be learned more quickly than  $\theta$ , making it harder to explore the parsing search space and optimize for  $\theta$ .

Difference in variance of two gradient estimates.

# Coadaptation problem

Learning composition function parameters  $\phi$  with backpropagation, and parser parameters  $\theta$  with REINFORCE.

Generally,  $\phi$  will be learned more quickly than  $\theta$ ,

ma

Dif

Possible solution:

Proximal Policy Optimization (Schulman et al., 2017)

# Making REINFORCE+SPINN work

Havrylov et al. (2019) use,

1. Input dependent control variate
2. Gradient normalization
3. Proximal Policy Optimization

# Making REINFORCE+SPINN work

Havrylov et al. (2019) use,

1. Input dependent control variate
2. Gradient normalization
3. Proximal Policy Optimization

They solve ListOps!

# Should I? Shouldn't I?

- Low bias!

# Should I? Shouldn't I?

- Low bias!
- High variance 🙄



# Should I? Shouldn't I?

- Low bias!
- In a simple setting, with enough tricks, it can work! 😊
- High variance 😞

# Should I? Shouldn't I?

- Low bias!
- In a simple setting, with enough tricks, it can work! 😊
- High variance 😞
- Has not yet been very effective at learning English syntax.

# III. Gradient Surrogates

# So far:

- Tackled **expected loss** in a **stochastic computation graph**

$$\mathbb{E}_{\pi_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})}[L(\mathbf{z})]$$

# So far:

- Tackled **expected loss** in a **stochastic computation graph**

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Optimized with the **score function estimation** method.

## So far:

- Tackled **expected loss** in a **stochastic computation graph**

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Optimized with the **score function estimation** method.
- Struggled with variance & sampling.

## So far:

- Tackled **expected loss** in a **stochastic computation graph**

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Optimized with the **score function estimation** method.
- Struggled with variance & sampling.

## In this section:

- Consider the **deterministic alternative**:

## So far:

- Tackled **expected loss** in a **stochastic computation graph**

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Optimized with the **score function estimation** method.
- Struggled with variance & sampling.

## In this section:

- Consider the **deterministic alternative**:

pick “best” structure       $\hat{\mathbf{z}}(x) := \arg \max_{\mathbf{z} \in \mathcal{M}} \pi_{\theta}(\mathbf{z} | x)$



## So far:

- Tackled **expected loss** in a **stochastic computation graph**

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Optimized with the **score function estimation** method.
- Struggled with variance & sampling.

## In this section:

- Consider the **deterministic alternative**:

pick “best” structure

incur loss

$$\hat{\mathbf{z}}(x) := \arg \max_{\mathbf{z} \in \mathcal{M}} \pi_{\theta}(\mathbf{z} | x) \\ L(\hat{\mathbf{z}}(x))$$

## So far:

- Tackled **expected loss** in a **stochastic computation graph**

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Optimized with the **score function estimation** method.
- Struggled with variance & sampling.

## In this section:

- Consider the **deterministic alternative**:

pick “best” structure

$$\hat{\mathbf{z}}(x) := \arg \max_{\mathbf{z} \in \mathcal{M}} \pi_{\theta}(\mathbf{z} | x)$$

incur loss

$$L(\hat{\mathbf{z}}(x))$$

- 3A: try to optimize the deterministic loss directly

## So far:

- Tackled **expected loss** in a **stochastic computation graph**

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Optimized with the **score function estimation** method.
- Struggled with variance & sampling.

## In this section:

- Consider the **deterministic alternative**:

pick “best” structure

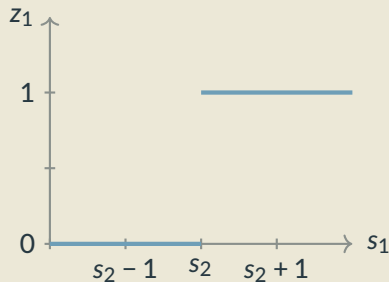
incur loss

$$\hat{\mathbf{z}}(x) := \arg \max_{\mathbf{z} \in \mathcal{M}} \pi_{\theta}(\mathbf{z} | x) \\ L(\hat{\mathbf{z}}(x))$$

- 3A: try to optimize the deterministic loss directly
- 3B: use this strategy to reduce variance in the stochastic model.

# Recap: The argmax problem

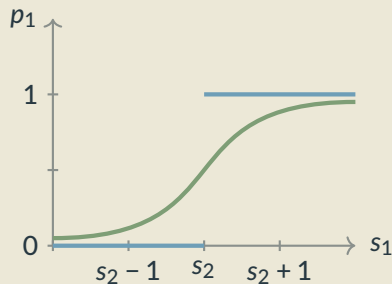
$$\mathbf{z} = \arg \max(\mathbf{s})$$



$$\frac{\partial \mathbf{z}}{\partial \mathbf{s}} = \mathbf{0}$$

# Softmax

$$p_j = \exp(s_j)/Z$$



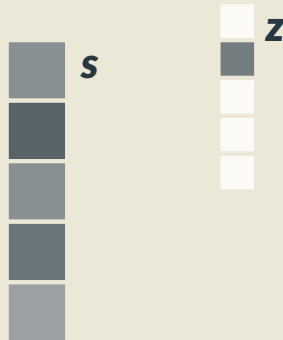
$$\frac{\partial \mathbf{p}}{\partial \mathbf{s}} = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top$$

# Straight-Through Estimator



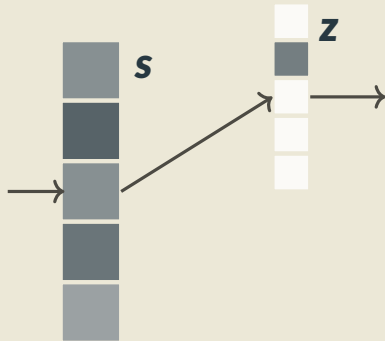
# Straight-Through Estimator

- Forward:  $\mathbf{z} = \arg \max(\mathbf{s})$



# Straight-Through Estimator

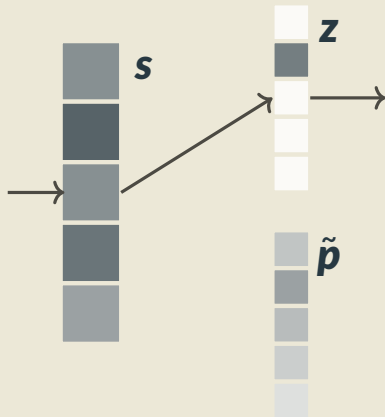
- Forward:  $\mathbf{z} = \arg \max(\mathbf{s})$





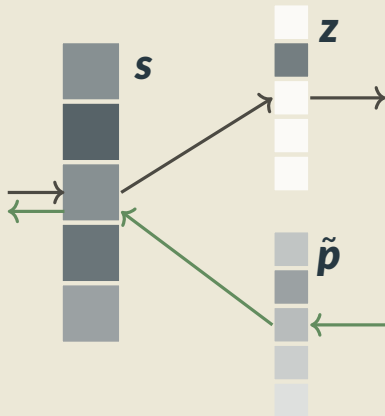
# Straight-Through Estimator

- Forward:  $\mathbf{z} = \arg \max(\mathbf{s})$
- Backward: pretend  $\mathbf{z}$  was some continuous  $\tilde{\mathbf{p}}$ ;  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} \neq \mathbf{0}$



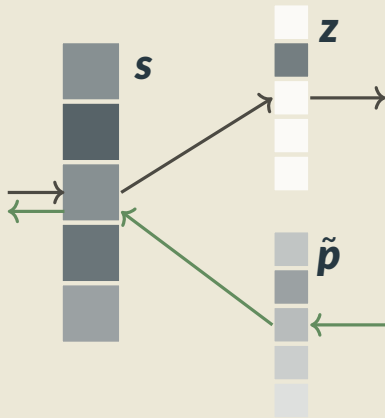
# Straight-Through Estimator

- Forward:  $\mathbf{z} = \arg \max(\mathbf{s})$
- Backward: pretend  $\mathbf{z}$  was some continuous  $\tilde{\mathbf{p}}$ ;  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} \neq \mathbf{0}$



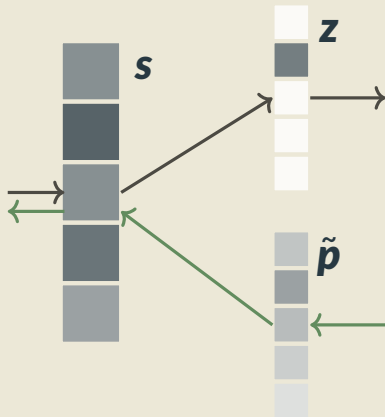
# Straight-Through Estimator

- Forward:  $\mathbf{z} = \arg \max(\mathbf{s})$
- Backward: pretend  $\mathbf{z}$  was some continuous  $\tilde{\mathbf{p}}$ ;  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} \neq \mathbf{0}$ 
  - simplest: identity,  $\tilde{\mathbf{p}}(\mathbf{s}) = \mathbf{s}$ ,  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} = \mathbf{I}$



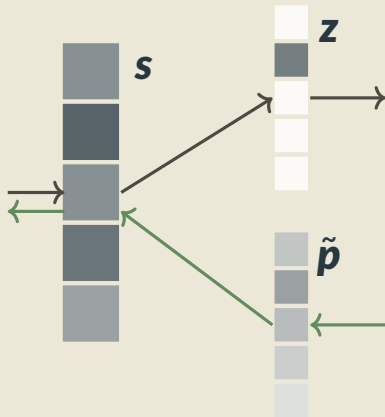
# Straight-Through Estimator

- Forward:  $\mathbf{z} = \arg \max(\mathbf{s})$
- Backward: pretend  $\mathbf{z}$  was some continuous  $\tilde{\mathbf{p}}$ ;  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} \neq \mathbf{0}$ 
  - simplest: identity,  $\tilde{\mathbf{p}}(\mathbf{s}) = \mathbf{s}$ ,  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} = \mathbf{I}$
  - others, e.g. softmax  $\tilde{\mathbf{p}}(\mathbf{s}) = \text{softmax}(\mathbf{s})$ ,  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} = \text{diag}(\tilde{\mathbf{p}}) - \tilde{\mathbf{p}}\tilde{\mathbf{p}}^T$



# Straight-Through Estimator

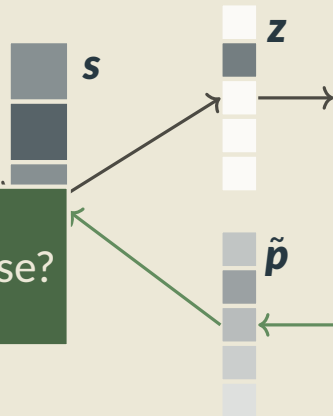
- Forward:  $\mathbf{z} = \arg \max(\mathbf{s})$
- Backward: pretend  $\mathbf{z}$  was some continuous  $\tilde{\mathbf{p}}$ ;  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} \neq \mathbf{0}$ 
  - simplest: identity,  $\tilde{\mathbf{p}}(\mathbf{s}) = \mathbf{s}$ ,  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} = \mathbf{I}$
  - others, e.g. softmax  $\tilde{\mathbf{p}}(\mathbf{s}) = \text{softmax}(\mathbf{s})$ ,  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} = \text{diag}(\tilde{\mathbf{p}}) - \tilde{\mathbf{p}}\tilde{\mathbf{p}}^T$
- More explanation in a while



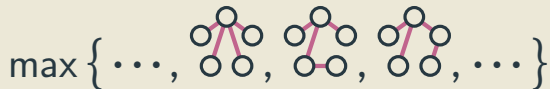
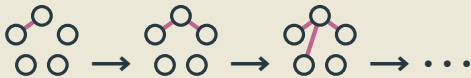
# Straight-Through Estimator

- Forward:  $\mathbf{z} = \arg \max(\mathbf{s})$
- Backward: pretend  $\mathbf{z}$  was some continuous  $\tilde{\mathbf{p}}$ ;  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} \neq \mathbf{0}$ 
  - simplest: identity,  $\tilde{\mathbf{p}}(\mathbf{s}) = \mathbf{s}$ ,  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} = \mathbf{I}$
  - others, e.g. softmax  $\tilde{\mathbf{p}}(\mathbf{s}) = \text{softmax}(\mathbf{s})$ ,  $\frac{\partial \tilde{\mathbf{p}}}{\partial \mathbf{s}} = \text{diag}(\tilde{\mathbf{p}}) - \tilde{\mathbf{p}}\tilde{\mathbf{p}}^T$
- More explanation

What about the structured case?



# Dealing with the combinatorial explosion



## 1. Incremental structures

- Build structure **greedily**, as sequence of discrete choices (e.g., shift-reduce).
- Scores (partial structure, action) tuples.
- **Advantages:** flexible, rich histories.
- **Disadvantages:** greedy, local decisions are suboptimal, error propagation.

## 2. Factorization into parts

- Optimizes **globally** (e.g. Viterbi, Chu-Liu-Edmonds, Kuhn-Munkres).
- Scores smaller parts.
- **Advantages:** optimal, elegant, can handle hard & global constraints.
- **Disadvantages:** strong assumptions.

# STE for incremental structures



# STE for incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)

# STE for incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.

# STE for incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.
- In this case, we just apply the straight-through estimator for each step.

# STE for incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.
- In this case, we just apply the straight-through estimator for each step.
- Forward: the **highest scoring action** for each step

# STE for incremental structures

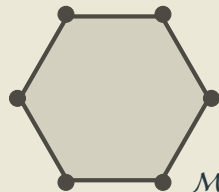
- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.
- In this case, we just apply the straight-through estimator for each step.
- Forward: the **highest scoring action** for each step
- Backward: pretend that we had used a **differentiable surrogate function**

# STE for incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.
- In this case, we just apply the straight-through estimator for each step.
- Forward: the **highest scoring action** for each step
- Backward: pretend that we had used a **differentiable surrogate function**

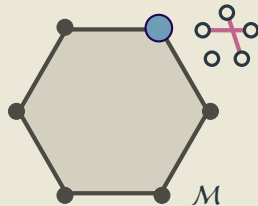
Example: Latent Tree Learning with Differentiable Parsers: Shift-Reduce Parsing and Chart Parsing [Maillard and Clark, 2018] (STE through beam search).

# The structured case: Marginal polytope



# The structured case: Marginal polytope

- Each vertex corresponds to one such *bit* vector  $\mathbf{z}$



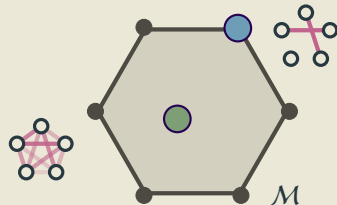


# The structured case: Marginal polytope

- Each vertex corresponds to one such *bit vector*  $\mathbf{z}$
- Points inside correspond to *marginal distributions*: convex combinations of structured objects

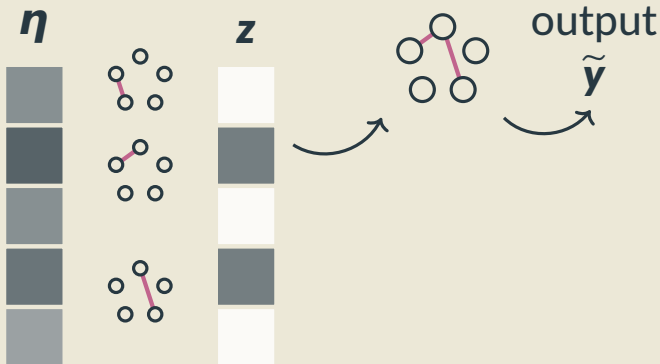
$$\boldsymbol{\mu} = \underbrace{p_1 \mathbf{z}_1 + \dots + p_N \mathbf{z}_N}_{\text{exponentially many terms}}, \quad \mathbf{p} \in \Delta.$$

$$\begin{aligned} p_1 &= 0.2, & \mathbf{z}_1 &= [1, 0, 0, 0, 1, 0, 0, 0, 1] \\ p_2 &= 0.7, & \mathbf{z}_2 &= [0, 0, 1, 0, 0, 1, 1, 0, 0] \\ p_3 &= 0.1, & \mathbf{z}_3 &= [1, 0, 0, 0, 1, 0, 0, 1, 0] \end{aligned} \quad \Rightarrow \quad \boldsymbol{\mu} = [.3, 0, .7, 0, .3, .7, .7, .1, .2].$$



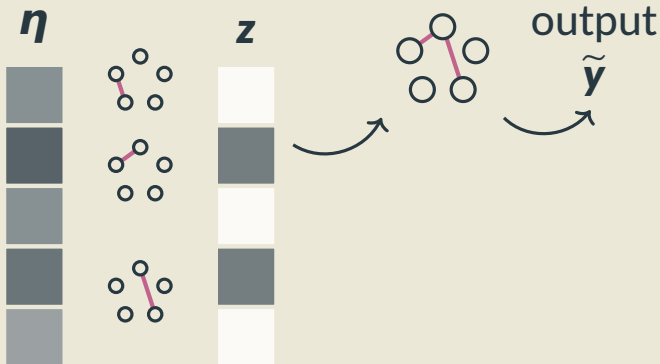
# STE for factorization into parts

- $\eta(i \rightarrow j)$ : score of arc  $i \rightarrow j$
- $z(i \rightarrow j)$ : is arc  $i \rightarrow j$  selected?



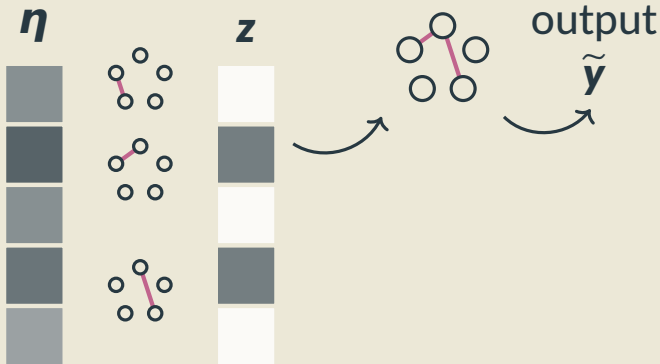
# STE for factorization into parts

- $\eta(i \rightarrow j)$ : score of arc  $i \rightarrow j$
- $z(i \rightarrow j)$ : is arc  $i \rightarrow j$  selected?
- Forward: structured argmax (task-specific algorithm!)



# STE for factorization into parts

- $\eta(i \rightarrow j)$ : score of arc  $i \rightarrow j$
- $z(i \rightarrow j)$ : is arc  $i \rightarrow j$  selected?
- Forward: structured argmax (task-specific algorithm!)
- Backward: identity  $\frac{\partial \tilde{\mu}}{\partial \eta} = I$



# Algorithms for specific structures

## Best structure (MAP)

**Sequence tagging**

Viterbi  
[Rabiner, 1989]

**Constituent trees**

CKY  
[Kasami, 1966, Younger, 1967]  
[Cocke and Schwartz, 1970]

**Temporal alignments**

DTW  
[Sakoe and Chiba, 1978]

**Dependency trees**

Max. Spanning Arborescence  
[Chu and Liu, 1965, Edmonds, 1967]

**Assignments**

Kuhn-Munkres  
[Kuhn, 1955, Jonker and Volgenant, 1987]

# Straight-Through Estimator

## Revisited

# Straight-Through Estimator

## Revisited

- In the forward pass: get the *argmax* for a structure.

# Straight-Through Estimator

## Revisited

- In the forward pass: get the *argmax* for a structure.
- if we had labels (multi-task learning),  $L = L_{\text{clf}}(\mathbf{z}; x, y) + L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}})$



# Straight-Through Estimator

## Revisited

- In the forward pass: get the *argmax* for a structure.
- if we had labels (multi-task learning),  $L = L_{\text{clf}}(\mathbf{z}; x, y) + L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}})$
- perceptron loss  $L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}}) = \mathbf{s}^T \mathbf{z} - \mathbf{s}^T \mathbf{z}^{\text{true}}$

# Straight-Through Estimator

## Revisited

- In the forward pass: get the *argmax* for a structure.
- if we had labels (multi-task learning),  $L = L_{\text{clf}}(\mathbf{z}; x, y) + L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}})$
- perceptron loss  $L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}}) = \mathbf{s}^{\text{T}} \mathbf{z} - \mathbf{s}^{\text{T}} \mathbf{z}^{\text{true}}$
- but we don't have labels, so we use the supervision from the downstream task:

# Straight-Through Estimator

## Revisited

- In the forward pass: get the *argmax* for a structure.
- if we had labels (multi-task learning),  $L = L_{\text{clf}}(\mathbf{z}; x, y) + L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}})$
- perceptron loss  $L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}}) = \mathbf{s}^T \mathbf{z} - \mathbf{s}^T \mathbf{z}^{\text{true}}$
- but we don't have labels, so we use the supervision from the downstream task:
- good guess for  $\mathbf{z}^{\text{true}}$ :

# Straight-Through Estimator

## Revisited

- In the forward pass: get the *argmax* for a structure.
- if we had labels (multi-task learning),  $L = L_{\text{clf}}(\mathbf{z}; x, y) + L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}})$
- perceptron loss  $L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}}) = \mathbf{s}^{\text{T}} \mathbf{z} - \mathbf{s}^{\text{T}} \mathbf{z}^{\text{true}}$
- but we don't have labels, so we use the supervision from the downstream task:
- good guess for  $\mathbf{z}^{\text{true}}$ :  $\arg \min_{\mathbf{p}} L_{\text{clf}}(\mathbf{p}, x, y)$  (unconstrained)

# Straight-Through Estimator

## Revisited

- In the forward pass: get the *argmax* for a structure.
- if we had labels (multi-task learning),  $L = L_{\text{clf}}(\mathbf{z}; x, y) + L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}})$
- perceptron loss  $L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}}) = \mathbf{s}^T \mathbf{z} - \mathbf{s}^T \mathbf{z}^{\text{true}}$
- but we don't have labels, so we use the supervision from the downstream task:
- good guess for  $\mathbf{z}^{\text{true}}$ :  $\arg \min_{\mathbf{p}} L_{\text{clf}}(\mathbf{p}, x, y)$  (unconstrained)
- one iteration of GD from  $\mathbf{z}$ :  $\mathbf{z}^{\text{true}} \leftarrow \tilde{\mathbf{p}} = \mathbf{z} - \nabla L_{\text{clf}}(\mathbf{z}, x, y)$

# Straight-Through Estimator

## Revisited

- In the forward pass: get the *argmax* for a structure.
- if we had labels (multi-task learning),  $L = L_{\text{clf}}(\mathbf{z}; x, y) + L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}})$
- perceptron loss  $L_{\text{hid}}(\mathbf{s}, \mathbf{z}^{\text{true}}) = \mathbf{s}^T \mathbf{z} - \mathbf{s}^T \mathbf{z}^{\text{true}}$
- but we don't have labels, so we use the supervision from the downstream task:
- good guess for  $\mathbf{z}^{\text{true}}$ :  $\arg \min_{\mathbf{p}} L_{\text{clf}}(\mathbf{p}, x, y)$  (unconstrained)
- one iteration of GD from  $\mathbf{z}$ :  $\mathbf{z}^{\text{true}} \leftarrow \tilde{\mathbf{p}} = \mathbf{z} - \nabla L_{\text{clf}}(\mathbf{z}, x, y)$
- perceptron loss yields:  $\frac{\partial L_{\text{hid}}}{\partial \mathbf{s}} = \mathbf{z} - (\mathbf{z} - \nabla L_{\text{clf}}(\mathbf{z}, x, y)) = \frac{\partial L_{\text{clf}}}{\partial \mathbf{z}}$

# Straight Through in the Structured Case

- Structured STE: perceptron update with fake annotation

$$\mathbf{z}^{\text{true}} = \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} L_{\text{clf}}(\boldsymbol{\mu}, x, y) \quad \approx \tilde{\boldsymbol{\mu}} = \mathbf{z} - \nabla L_{\text{clf}}(\mathbf{z}, x, y)$$

(one step of gradient descent)

# Straight Through in the Structured Case

- Structured STE: perceptron update with fake annotation

$$\mathbf{z}^{\text{true}} = \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} L_{\text{clf}}(\boldsymbol{\mu}, \mathbf{x}, y) \quad \approx \tilde{\boldsymbol{\mu}} = \mathbf{z} - \nabla L_{\text{clf}}(\mathbf{z}, \mathbf{x}, y)$$

(one step of gradient descent)

- SPIGOT takes into account the constraints; uses the fake annotation

$$\mathbf{z}^{\text{true}} = \arg \min_{\boldsymbol{\mu} \in \mathcal{M}} L_{\text{clf}}(\boldsymbol{\mu}, \mathbf{x}, y) \quad \approx \tilde{\boldsymbol{\mu}} = \text{Proj}_{\mathcal{M}}(\mathbf{z} - \nabla L_{\text{clf}}(\mathbf{z}, \mathbf{x}, y))$$

(one step of *projected* gradient descent!)



# Straight Through in the Structured Case

- Structured STE: perceptron update with fake annotation

$$\mathbf{z}^{\text{true}} = \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} L_{\text{clf}}(\boldsymbol{\mu}, \mathbf{x}, y) \quad \approx \tilde{\boldsymbol{\mu}} = \mathbf{z} - \nabla L_{\text{clf}}(\mathbf{z}, \mathbf{x}, y)$$

(one step of gradient descent)

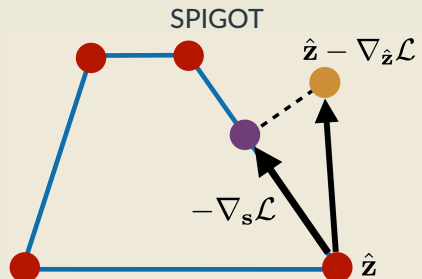
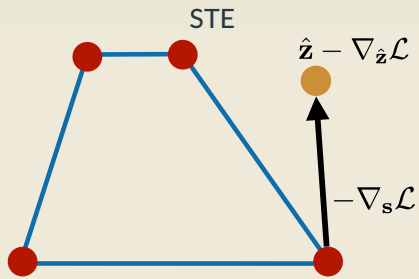
- SPIGOT takes into account the constraints; uses the fake annotation

$$\mathbf{z}^{\text{true}} = \arg \min_{\boldsymbol{\mu} \in \mathcal{M}} L_{\text{clf}}(\boldsymbol{\mu}, \mathbf{x}, y) \quad \approx \tilde{\boldsymbol{\mu}} = \text{Proj}_{\mathcal{M}}(\mathbf{z} - \nabla L_{\text{clf}}(\mathbf{z}, \mathbf{x}, y))$$

(one step of *projected* gradient descent!)

- We discuss a generic way to compute the projection in part 3.

# SPIGOT vs STE



# Summary: Straight-Through Estimator

# Summary: Straight-Through Estimator

We saw how to use the *Straight-Through Estimator* to allow learning models with *argmax* in the middle of the computation graph.

# Summary: Straight-Through Estimator

We saw how to use the *Straight-Through Estimator* to allow learning models with *argmax* in the middle of the computation graph.

We were optimizing  $L(\hat{\mathbf{z}}(x))$

# Summary: Straight-Through Estimator

We saw how to use the *Straight-Through Estimator* to allow learning models with *argmax* in the middle of the computation graph.

We were optimizing  $L(\hat{\mathbf{z}}(x))$

Now we will see how to apply STE for stochastic graphs, as an alternative approach of the score-function estimators.

# Stochastic node in the computation graph

# Stochastic node in the computation graph

Recall the stochastic objective:

$$\mathbb{E}_{\pi_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})}[L(\mathbf{z})]$$



# Stochastic node in the computation graph

Recall the stochastic objective:

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Score function methods (previous section).

# Stochastic node in the computation graph

Recall the stochastic objective:

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Score function methods (previous section). High variance. :(

# Stochastic node in the computation graph

Recall the stochastic objective:

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Score function methods (previous section). High variance. :(
- An alternative is using the *reparameterization trick* [Kingma and Welling, 2013].

# Categorical reparameterization

# Categorical reparameterization

- Sampling from a categorical value in the middle of the computation graph.

$$\mathbf{z} \sim \pi_{\theta}(\mathbf{z} \mid \mathbf{x}) \propto \exp \mathbf{s}_{\theta}(\mathbf{z} \mid \mathbf{x})$$



# Categorical reparameterization

- Sampling from a categorical value in the middle of the computation graph.  
 $\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathbf{x}) \propto \exp \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathbf{x})$
- What is the gradient of a sample  $\frac{\partial \mathbf{z}}{\partial \boldsymbol{\theta}}$ ?!

s

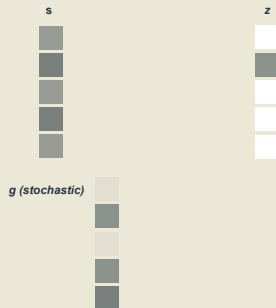


z



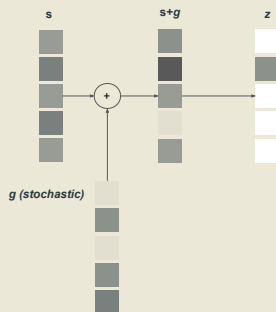
# Categorical reparameterization

- Sampling from a categorical value in the middle of the computation graph.  
 $\mathbf{z} \sim \pi_{\theta}(\mathbf{z} | \mathbf{x}) \propto \exp \mathbf{s}_{\theta}(\mathbf{z} | \mathbf{x})$
- What is the gradient of a sample  $\frac{\partial \mathbf{z}}{\partial \theta}$ ?!
- Reparameterization: Move the stochasticity out of the gradient path.



# Categorical reparameterization

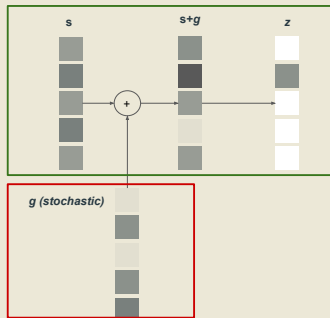
- Sampling from a categorical value in the middle of the computation graph.  
 $\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathbf{x}) \propto \exp \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathbf{x})$
- What is the gradient of a sample  $\frac{\partial \mathbf{z}}{\partial \boldsymbol{\theta}}$  ?!
- Reparameterization: Move the stochasticity out of the gradient path.
- Makes  $\mathbf{z}$  deterministic w.r.t.  $\mathbf{s}$ !





# Categorical reparameterization

- Sampling from a categorical value in the middle of the computation graph.  
 $\mathbf{z} \sim \pi_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathbf{x}) \propto \exp \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathbf{x})$
- What is the gradient of a sample  $\frac{\partial \mathbf{z}}{\partial \boldsymbol{\theta}} ?!$
- Reparameterization: Move the stochasticity out of the gradient path.
- Makes  $\mathbf{z}$  deterministic w.r.t.  $\mathbf{s}$ !

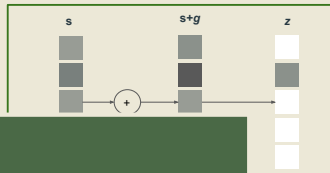


# Categorical reparameterization

- Sampling from a categorical value in the middle of the computation graph.

$$\mathbf{z} \sim \pi_{\theta}(\mathbf{z} | \mathbf{x}) \propto \exp \mathbf{s}_{\theta}(\mathbf{z} | \mathbf{x})$$

- What is the  $\xi$
- Reparameterize stochasticity
- Makes  $\mathbf{z}$  deterministic w.r.t.  $\mathbf{s}$ !



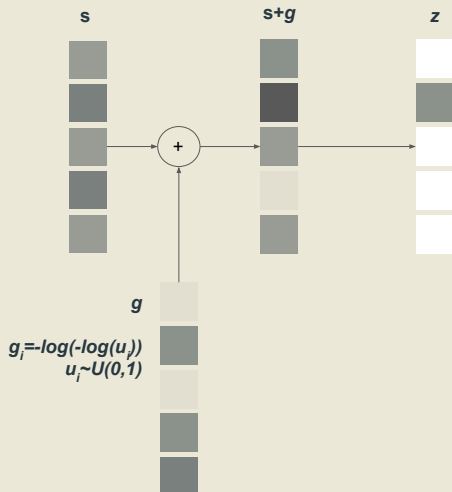
As a result:

Stochasticity is moved as an input.

We can backpropagate through the deterministic input to  $\mathbf{z}$ .



# Categorical reparameterization



# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

**1. Usually we sample  $z$ :**

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$



# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

## 2. An alternative way: Gumbel-max trick

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

## 2. An alternative way: Gumbel-max trick

- $u \sim \text{Uniform}(0, 1)$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

## 2. An alternative way: Gumbel-max trick

- $u \sim \text{Uniform}(0, 1)$
- $g_i = -\log(-\log(u))$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

## 2. An alternative way: Gumbel-max trick

- $u \sim \text{Uniform}(0, 1)$
- $g_i = -\log(-\log(u))$
- $\hat{z} = \arg \max_i (s_i + g_i)$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

The two methods are equivalent. *(Not obvious, but we will not prove it now.)*

## 2. An alternative way: Gumbel-max trick

- $u \sim \text{Uniform}(0, 1)$
- $g_i = -\log(-\log(u))$
- $\hat{z} = \arg \max_i (s_i + g_i)$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

The two methods are equivalent. (*Not obvious, but we will not prove it now.*)

Requires sampling from the Standard Gumbel Distribution  $G(0,1)$ .

## 2. An alternative way: Gumbel-max trick

- $u \sim \text{Uniform}(0, 1)$
- $g_i = -\log(-\log(u))$
- $\hat{z} = \arg \max_i (s_i + g_i)$



# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

The two methods are equivalent. *(Not obvious, but we will not prove it now.)*

Requires sampling from the Standard Gumbel Distribution  $G(0,1)$ .

*Derivation in this blog post from Ryan Adams:*

<http://lips.cs.princeton.edu/the-gumbel-max-trick-for-discrete-distributions/>

## 2. An alternative way: Gumbel-max trick

- $u \sim \text{Uniform}(0, 1)$
- $g_i = -\log(-\log(u))$
- $\hat{z} = \arg \max_i (s_i + g_i)$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_y \leq u < c_{y+1}$

The two methods are equivalent. *(Not obvious, but we will not prove it now.)*

Requires sampling from the Standard Gumbel Distribution  $G(0,1)$ .

*Derivation in this blog post from Ryan Adams:*

*<http://lips.cs.princeton.edu/the-gumbel-max-trick-for-discrete-distributions/>*

*More information in Tim Vieira's blog: <https://timvieira.github.io/blog/post/2014/07/31/gumbel-max-trick/>*

## 2. An alternative way: Gumbel-max trick

- $u \sim \text{Uniform}(0, 1)$
- $g_i = -\log(-\log(u))$
- $\hat{z} = \arg \max_i (s_i + g_i)$

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_{\hat{z}-1} < u \leq c_{\hat{z}}$

The t

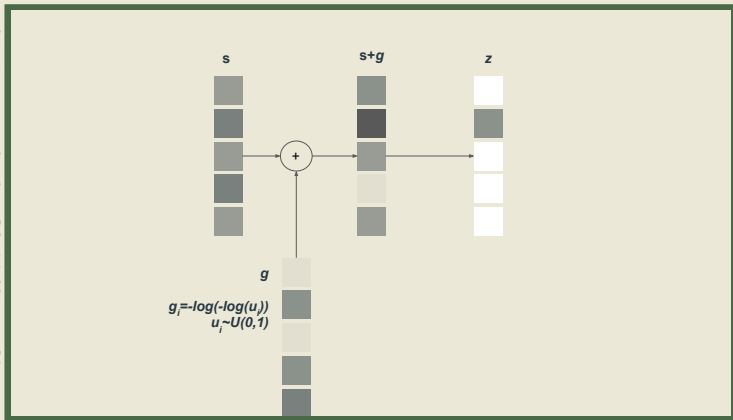
l

http:

More information in

## 2. An alternative way:

trick



now.)

).

ons/

l/gumbel-max-trick/

# Sampling from a categorical variable

We want to sample from a categorical variable with scores  $\mathbf{s}$  (class  $i$  has a score  $s_i$ )

## 1. Usually we sample $z$ :

- $p_i = \frac{\exp(s_i)}{\sum_{j=1}^K \exp(s_j)}$
- $c_i = \sum_{j \leq i} p_j$
- $u \sim \text{Uniform}(0, 1)$
- return  $\hat{z}$  s.t.  $c_{\hat{z}-1} < u \leq c_{\hat{z}}$

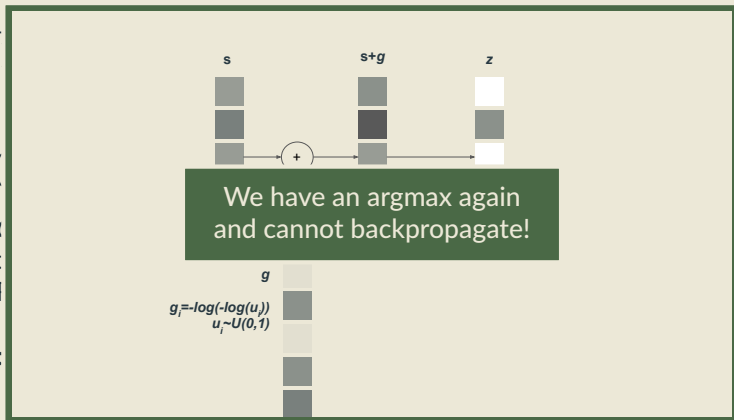
The t  
l

http:

More information in

## 2. An alternative way:

trick



now.)

).

ons/

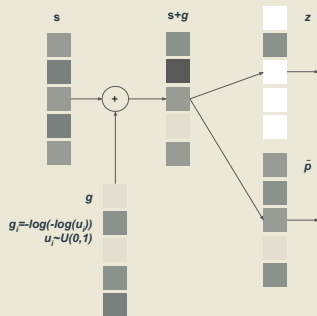
l/gumbel-max-trick/

# Straight-Through Gumbel Estimator

- Apply a variant of the Straight-Through Estimator to Gumbel-Max!

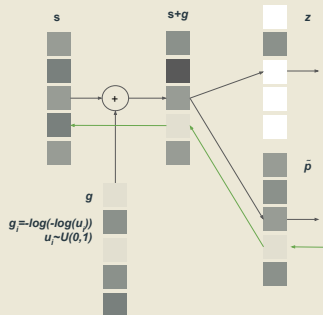
# Straight-Through Gumbel Estimator

- Apply a variant of the Straight-Through Estimator to Gumbel-Max!
- Forward:  $\mathbf{z} = \arg \max(\mathbf{s} + \mathbf{g})$



# Straight-Through Gumbel Estimator

- Apply a variant of the Straight-Through Estimator to Gumbel-Max!
- Forward:  $\mathbf{z} = \arg \max(\mathbf{s} + \mathbf{g})$
- Backward: pretend we had done  $\tilde{\mathbf{p}} = \text{softmax}(\mathbf{s} + \mathbf{g})$



# Straight-Through Gumbel Estimator

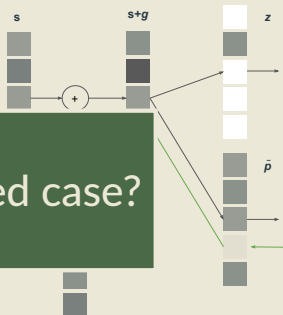
- Apply a variant of the Straight-Through Estimator to Gumbel-Max!

- Forward:  $\mathbf{z} = \arg$

- Backward: prete

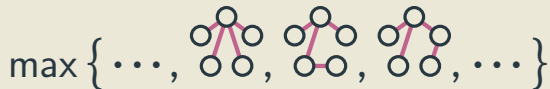
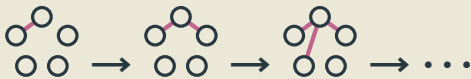
$$\tilde{\mathbf{p}} = \text{softmax}(\mathbf{s} + \mathbf{g})$$

What about the structured case?





# Dealing with the combinatorial explosion



## 1. Incremental structures

- Build structure **greedily**, as sequence of discrete choices (e.g., shift-reduce).
- Scores (partial structure, action) tuples.
- **Advantages:** flexible, rich histories.
- **Disadvantages:** greedy, local decisions are suboptimal, error propagation.

## 2. Factorization into parts

- Optimizes **globally** (e.g. Viterbi, Chu-Liu-Edmonds, Kuhn-Munkres).
- Scores smaller parts.
- **Advantages:** optimal, elegant, can handle hard & global constraints.
- **Disadvantages:** strong assumptions.

# Sampling from incremental structures

# Sampling from incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)

# Sampling from incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.

# Sampling from incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.
- Reparameterize the scores with Gumbel-Max - now we have a deterministic node.

# Sampling from incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.
- Reparameterize the scores with Gumbel-Max - now we have a deterministic node.
- Forward: the **argmax** from the reparameterized scores for each step

# Sampling from incremental structures

- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.
- Reparameterize the scores with Gumbel-Max - now we have a deterministic node.
- Forward: the **argmax** from the reparameterized scores for each step
- Backward: pretend that we had used a **differentiable surrogate functions**

# Sampling from incremental structures

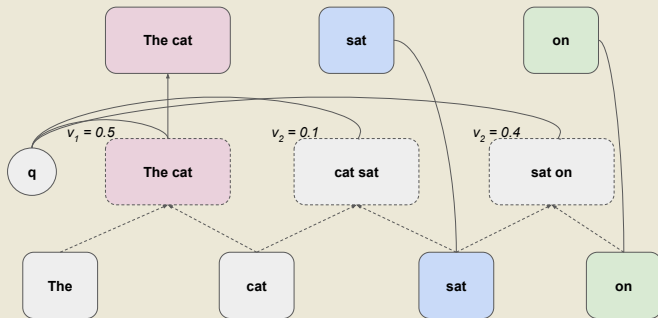
- Build a structure as a sequence of discrete choices (e.g., shift-reduce)
- Assigns a score to any (partial structure, action) tuple.
- Reparameterize the scores with Gumbel-Max - now we have a deterministic node.
- Forward: the **argmax** from the reparameterized scores for each step
- Backward: pretend that we had used a **differentiable surrogate functions**

Example: Gumbel Tree-LSTM [Choi et al., 2018].



# Example: Gumbel Tree-LSTM

- Building task-specific tree structures.
- Straight-Through Gumbel-Softmax at each step to select one arc.



# Sampling from factorized models

## Perturb-and-MAP

Reparameterize by **perturbing the arc scores**. (inexact!)

- Sample from the normal Gumbel distribution.
- Perturb the arc scores with the Gumbel noise.
- Compute MAP (task-specific algorithm).
- $\mathbf{g} \sim G(0, 1)$
- $\tilde{\boldsymbol{\eta}} = \boldsymbol{\eta} + \mathbf{g}$
- $\arg \max_{\mathbf{z} \in \mathcal{Z}} \tilde{\boldsymbol{\eta}}^T \mathbf{z}$

# Summary: Gradient surrogates

- Based on the **Straight-Through Estimator**.
- Can be used for stochastic or deterministic computation graphs.
- **Forward pass**: Get an argmax (might be structured).
- **Backpropagation**: use a function, which we hope is close to argmax.
- Examples:
  - Argmax for iterative structures and factorization into parts
  - Sampling from iterative structures and factorization into parts

# Gradient surrogates: Pros and cons

## Pros

- Do not suffer from the high variance problem of the score function methods.
- Allow for flexibility to select or sample a latent structured in the middle of the computation graph.
- Efficient computation.

## Cons

- The Gumbel sampling with Perturb-and-MAP is an approximation.
- Bias, due to function mismatch on the backpropagation (next section will address this problem.)

# Overview

$$\mathbb{E}_{\pi_{\boldsymbol{\theta}}(\mathbf{z}|x)}[L(\mathbf{z})]$$

$$L(\arg \max_{\mathbf{z}} \pi_{\boldsymbol{\theta}}(\mathbf{z} | x))$$

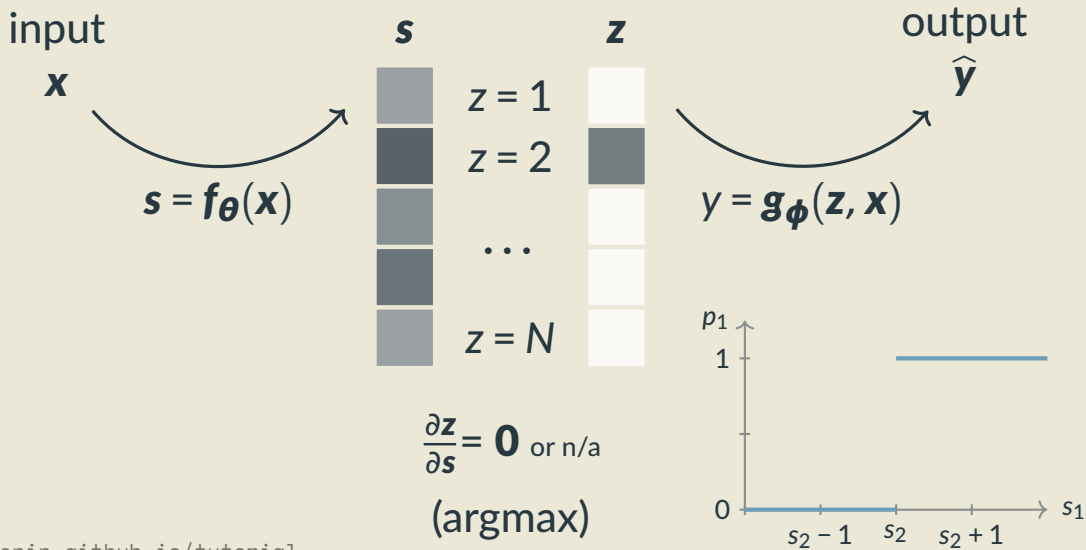
- Score Function Estimator
- Straight Through-Gumbel
- Perturb-and-Parse
- Straight Through (flavors)
- SPIGOT

# **IV. End-to-end differentiable methods**

# End-to-end differentiable methods

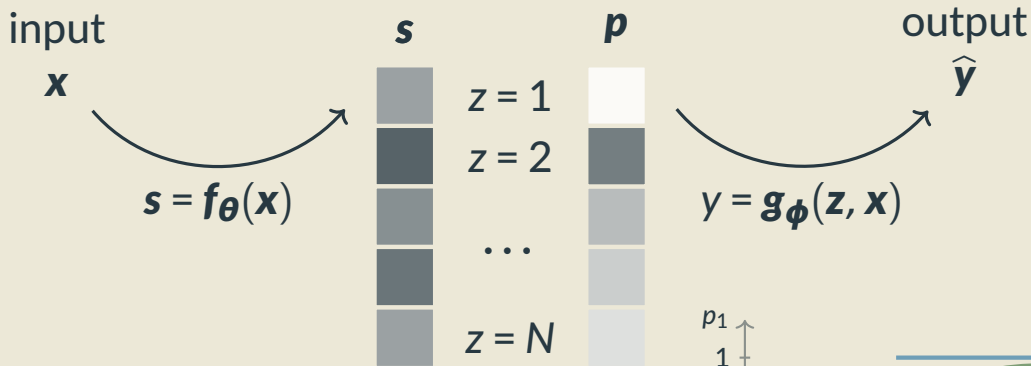
1. Digging into softmax
2. Alternatives to softmax
3. Generalizing to structured prediction
4. Stochasticity and global structures

# Recall: Discrete choices & differentiability





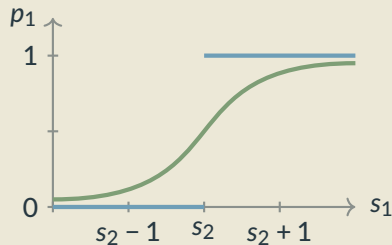
# One solution: smooth relaxation



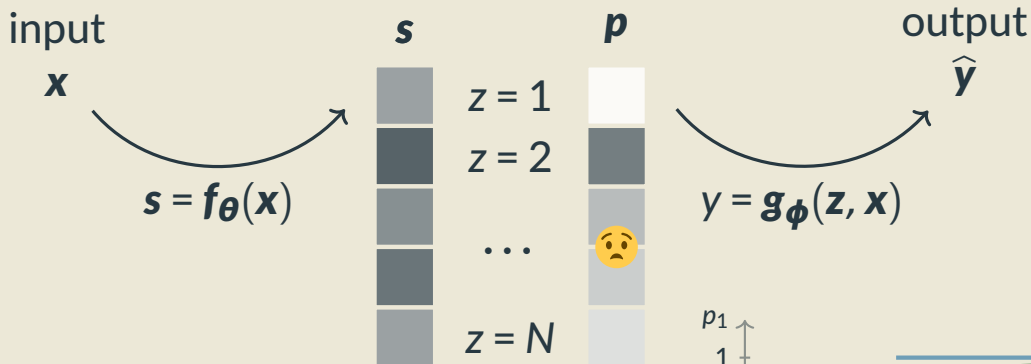
$\mathbf{p} = \text{softmax}(\mathbf{s}) = \mathbb{E}[\mathbf{z}]$ , i.e.  
replace  $\mathbb{E}[f(\mathbf{z})]$  with  $f(\mathbb{E}[\mathbf{z}])$

$$\frac{\partial \mathbf{p}}{\partial \mathbf{s}} = \text{😊}$$

(softmax)



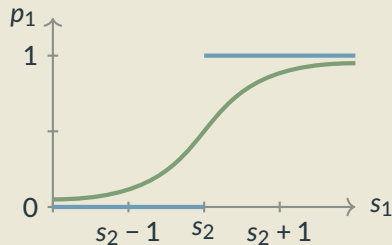
# One solution: smooth relaxation



$\mathbf{p} = \text{softmax}(\mathbf{s}) = \mathbb{E}[\mathbf{z}]$ , i.e.  
replace  $\mathbb{E}[f(\mathbf{z})]$  with  $f(\mathbb{E}[\mathbf{z}])$

$$\frac{\partial \mathbf{p}}{\partial \mathbf{s}} = \text{😊}$$

(softmax)



# Overview

$$\mathbb{E}_{\pi_{\boldsymbol{\theta}}(\mathbf{z}|x)}[L(\mathbf{z})]$$

$$L(\arg \max_{\mathbf{z}} \pi_{\boldsymbol{\theta}}(\mathbf{z} | x))$$

$$L(\mathbb{E}_{\pi_{\boldsymbol{\theta}}(\mathbf{z}|x)}[\mathbf{z}])$$

- Score Function Estimator
- Straight Through-Gumbel
- Perturb-and-Parse
- Straight Through (flavors)
- SPIGOT

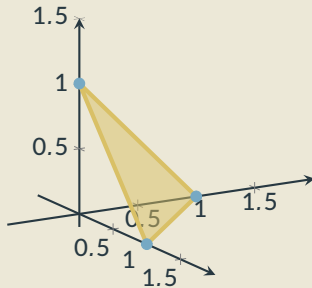
# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

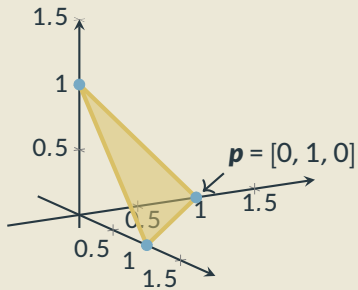
$\mathbf{p} \in \Delta$ : probability distribution over choices



# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

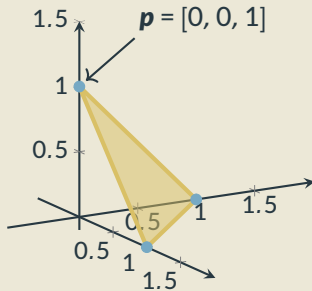
$\mathbf{p} \in \Delta$ : probability distribution over choices



# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

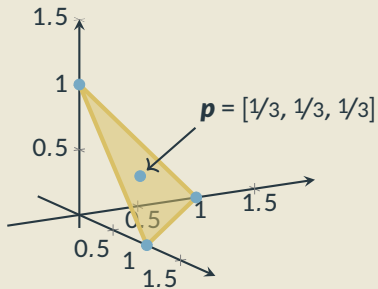
$\mathbf{p} \in \Delta$ : probability distribution over choices



# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

$\mathbf{p} \in \Delta$ : probability distribution over choices



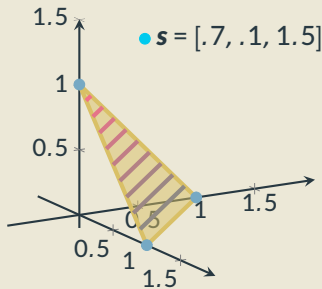


# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

$\mathbf{p} \in \Delta$ : probability distribution over choices

Expected score under  $\mathbf{p}$ :  $\mathbb{E}_{i \sim \mathbf{p}} s_i = \mathbf{p}^\top \mathbf{s}$



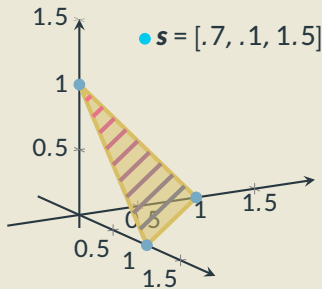
# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

$\mathbf{p} \in \Delta$ : probability distribution over choices

Expected score under  $\mathbf{p}$ :  $\mathbb{E}_{i \sim \mathbf{p}} s_i = \mathbf{p}^\top \mathbf{s}$

**argmax**



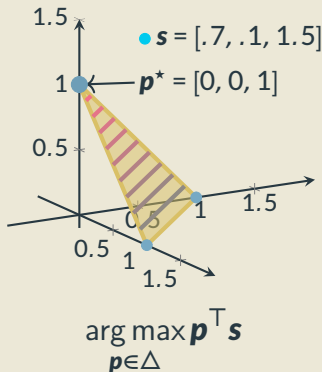
# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

$\mathbf{p} \in \Delta$ : probability distribution over choices

Expected score under  $\mathbf{p}$ :  $\mathbb{E}_{i \sim \mathbf{p}} s_i = \mathbf{p}^\top \mathbf{s}$

**argmax** maximizes **expected score**



# What is softmax?

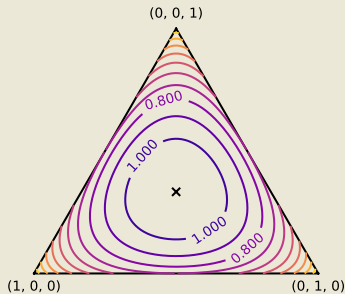
Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

$\mathbf{p} \in \Delta$ : probability distribution over choices

Expected score under  $\mathbf{p}$ :  $\mathbb{E}_{i \sim \mathbf{p}} s_i = \mathbf{p}^\top \mathbf{s}$

**argmax** maximizes **expected score**

Shannon entropy of  $\mathbf{p}$ :  $H(\mathbf{p}) = -\sum_i p_i \log p_i$



# What is softmax?

Often defined via  $p_i = \frac{\exp s_i}{\sum_j \exp s_j}$ , but where does it come from?

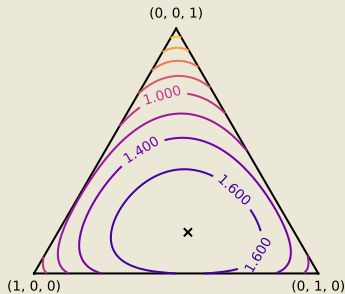
$\mathbf{p} \in \Delta$ : probability distribution over choices

Expected score under  $\mathbf{p}$ :  $\mathbb{E}_{i \sim \mathbf{p}} s_i = \mathbf{p}^\top \mathbf{s}$

**argmax** maximizes **expected score**

Shannon entropy of  $\mathbf{p}$ :  $H(\mathbf{p}) = -\sum_i p_i \log p_i$

**softmax** maximizes **expected score + entropy**:



$$\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$$

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ & \text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ & \text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = - \sum_j p_j s_j - p_j \log p_j - \mathbf{p}^\top \boldsymbol{\nu} + \tau(\mathbf{p}^\top \mathbf{1} - 1)$$



# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ & \text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -\sum_j p_j s_j - p_j \log p_j - \mathbf{p}^\top \boldsymbol{\nu} + \tau(\mathbf{p}^\top \mathbf{1} - 1)$$

Optimality conditions (KKT):

$$\begin{aligned} 0 &= \nabla_{p_i} \mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -s_i + \log p_i + 1 - \nu_i + \tau \\ \mathbf{p}^\top \boldsymbol{\nu} &= 0 \\ \mathbf{p} &\in \Delta \\ \boldsymbol{\nu} &\geq 0 \end{aligned}$$

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\log p_i = s_i + \nu_i - (\tau + 1)$$

$$\begin{aligned} &\text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ &\text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -\sum_j p_j s_j - p_j \log p_j - \mathbf{p}^\top \boldsymbol{\nu} + \tau(\mathbf{p}^\top \mathbf{1} - 1)$$

Optimality conditions (KKT):

$$0 = \nabla_{p_i} \mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -s_i + \log p_i + 1 - \nu_i + \tau$$

$$\mathbf{p}^\top \boldsymbol{\nu} = 0$$

$$\mathbf{p} \in \Delta$$

$$\boldsymbol{\nu} \geq 0$$

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ & \text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

$$\begin{aligned} \log p_i &= s_i + \nu_i - (\tau + 1) \\ \text{if } p_i &= 0, \text{ r.h.s. must be } -\infty, \\ \text{thus } p_i &> 0, \text{ so } \nu_i = 0. \end{aligned}$$

Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -\sum_j p_j s_j - p_j \log p_j - \mathbf{p}^\top \boldsymbol{\nu} + \tau(\mathbf{p}^\top \mathbf{1} - 1)$$

Optimality conditions (KKT):

$$0 = \nabla_{p_i} \mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -s_i + \log p_i + 1 - \nu_i + \tau$$

$$\mathbf{p}^\top \boldsymbol{\nu} = 0$$

$$\mathbf{p} \in \Delta$$

$$\boldsymbol{\nu} \geq 0$$

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ & \text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

$$\begin{aligned} \log p_i &= s_i + \nu_i - (\tau + 1) \\ \text{if } p_i &= 0, \text{ r.h.s. must be } -\infty, \\ \text{thus } p_i &> 0, \text{ so } \nu_i = 0. \end{aligned}$$

$$p_i = \exp(s_i) / \exp(\tau + 1) = \exp(s_i) / Z$$

Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -\sum_j p_j s_j - p_j \log p_j - \mathbf{p}^\top \boldsymbol{\nu} + \tau(\mathbf{p}^\top \mathbf{1} - 1)$$

Optimality conditions (KKT):

$$0 = \nabla_{p_i} \mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -s_i + \log p_i + 1 - \nu_i + \tau$$

$$\mathbf{p}^\top \boldsymbol{\nu} = 0$$

$$\mathbf{p} \in \Delta$$

$$\boldsymbol{\nu} \geq 0$$

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ & \text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -\sum_j p_j s_j - p_j \log p_j - \mathbf{p}^\top \boldsymbol{\nu} + \tau(\mathbf{p}^\top \mathbf{1} - 1)$$

Optimality conditions (KKT):

$$0 = \nabla_{p_i} \mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -s_i + \log p_i + 1 - \nu_i + \tau$$

$$\mathbf{p}^\top \boldsymbol{\nu} = 0$$

$$\mathbf{p} \in \Delta$$

$$\boldsymbol{\nu} \geq 0$$

$$\log p_i = s_i + \nu_i - (\tau + 1)$$

if  $p_i = 0$ , r.h.s. must be  $-\infty$ ,  
thus  $p_i > 0$ , so  $\nu_i = 0$ .

$$p_i = \exp(s_i) / \exp(\tau + 1) = \exp(s_i) / Z$$

Must find  $Z$  such that  $\sum_j p_j = 1$ .

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ & \text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -\sum_j p_j s_j - p_j \log p_j - \mathbf{p}^\top \boldsymbol{\nu} + \tau(\mathbf{p}^\top \mathbf{1} - 1)$$

Optimality conditions (KKT):

$$\begin{aligned} 0 &= \nabla_{p_i} \mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -s_i + \log p_i + 1 - \nu_i + \tau \\ \mathbf{p}^\top \boldsymbol{\nu} &= 0 \\ \mathbf{p} &\in \Delta \\ \boldsymbol{\nu} &\geq 0 \end{aligned}$$

$$\begin{aligned} \log p_i &= s_i + \nu_i - (\tau + 1) \\ \text{if } p_i &= 0, \text{ r.h.s. must be } -\infty, \\ \text{thus } p_i &> 0, \text{ so } \nu_i = 0. \end{aligned}$$

$$p_i = \exp(s_i) / \exp(\tau + 1) = \exp(s_i) / Z$$

Must find  $Z$  such that  $\sum_j p_j = 1$ .

Answer:  $Z = \sum_j \exp(s_j)$

# Variational form of softmax

**Proposition.** The unique solution to  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{s} + H(\mathbf{p})$  is given by  $p_j = \frac{\exp s_j}{\sum_i \exp s_i}$ .

Explicit form of the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_j p_j s_j - p_j \log p_j \\ & \text{subject to} && \mathbf{p} \geq 0, \mathbf{p}^\top \mathbf{1} = 1 \end{aligned}$$

Lagrangian:

$$\mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -\sum_j p_j s_j - p_j \log p_j - \mathbf{p}^\top \boldsymbol{\nu} + \tau(\mathbf{p}^\top \mathbf{1} - 1)$$

Optimality conditions (KKT):

$$0 = \nabla_{p_i} \mathcal{L}(\mathbf{p}, \boldsymbol{\nu}, \tau) = -s_i + \log p_i + 1 - \nu_i + \tau$$

$$\mathbf{p}^\top \boldsymbol{\nu} = 0$$

$$\mathbf{p} \in \Delta$$

$$\boldsymbol{\nu} \geq 0$$

$$\begin{aligned} \log p_i &= s_i + \nu_i - (\tau + 1) \\ \text{if } p_i &= 0, \text{ r.h.s. must be } -\infty, \\ \text{thus } p_i &> 0, \text{ so } \nu_i = 0. \end{aligned}$$

$$p_i = \exp(s_i) / \exp(\tau + 1) = \exp(s_i) / Z$$

Must find  $Z$  such that  $\sum_j p_j = 1$ .

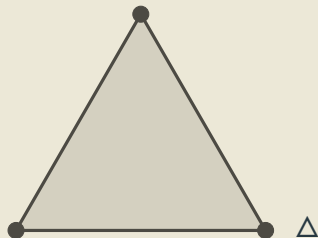
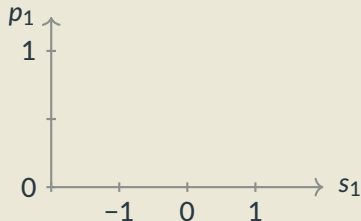
Answer:  $Z = \sum_j \exp(s_j)$

$$\text{So, } p_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}.$$

Classic result, e.g., [Boyd and Vandenberghe, 2004, Wainwright and Jordan, 2008]

# Generalizing softmax: Smoothed argmaxes

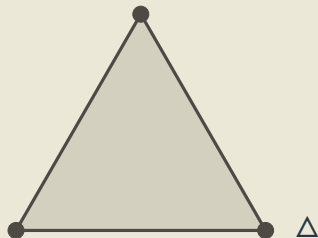
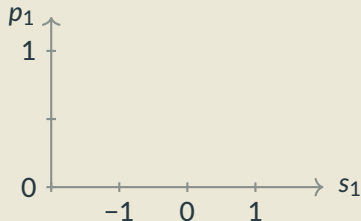
$$\hat{\mathbf{p}}_{\Omega}(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \mathbf{s} - \Omega(\mathbf{p})$$





# Generalizing softmax: Smoothed argmaxes

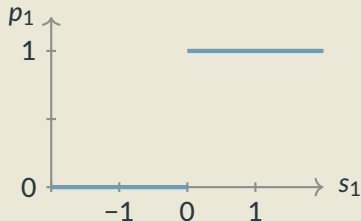
$$\hat{\mathbf{p}}_{\Omega}(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \mathbf{s} - \Omega(\mathbf{p})$$



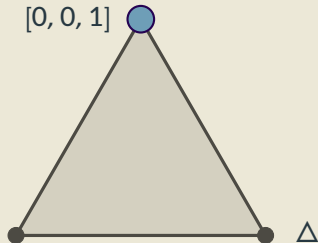
# Generalizing softmax: Smoothed argmaxes

$$\hat{\mathbf{p}}_{\Omega}(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \mathbf{s} - \Omega(\mathbf{p})$$

- argmax:  $\Omega(\mathbf{p}) = 0$



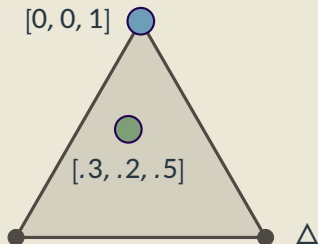
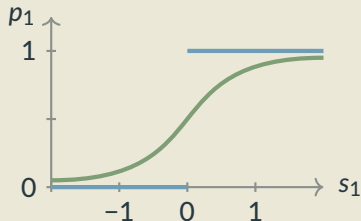
$[0, 0, 1]$



# Generalizing softmax: Smoothed argmaxes

$$\hat{\mathbf{p}}_{\Omega}(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \mathbf{s} - \Omega(\mathbf{p})$$

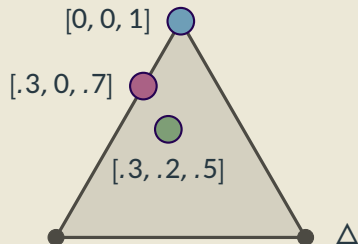
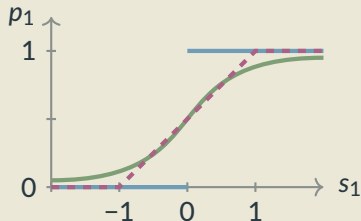
- argmax:  $\Omega(\mathbf{p}) = 0$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$



# Generalizing softmax: Smoothed argmaxes

$$\hat{\mathbf{p}}_{\Omega}(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \mathbf{s} - \Omega(\mathbf{p})$$

- argmax:  $\Omega(\mathbf{p}) = 0$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2$



# Generalizing softmax: Smoothed argmaxes

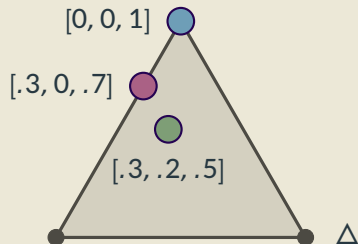
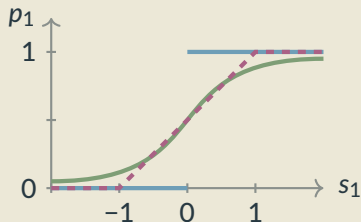
$$\hat{\mathbf{p}}_{\Omega}(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \mathbf{s} - \Omega(\mathbf{p})$$

- argmax:  $\Omega(\mathbf{p}) = 0$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2$
- $\alpha$ -entmax:  $\Omega(\mathbf{p}) = 1/\alpha(\alpha-1) \sum_j p_j^{\alpha}$

Generalized entropy interpolates in between [Tsallis, 1988]

Used in Sparse Seq2Seq: [Peters et al., 2019]

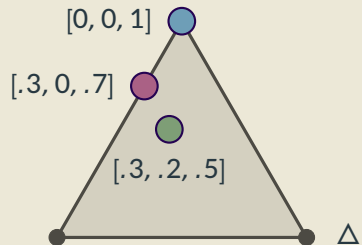
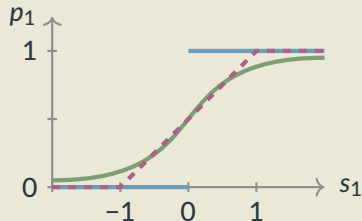
(Mon 13:50, poster session 2D)



# Generalizing softmax: Smoothed argmaxes

$$\hat{\mathbf{p}}_{\Omega}(\mathbf{s}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^{\top} \mathbf{s} - \Omega(\mathbf{p})$$

- argmax:  $\Omega(\mathbf{p}) = 0$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2$
- $\alpha$ -entmax:  $\Omega(\mathbf{p}) = 1/\alpha(\alpha-1) \sum_j p_j^{\alpha}$
- fusedmax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2 + \sum_j |p_j - p_{j-1}|$
- csparsesmax:  $\Omega(\mathbf{p}) = 1/2 \|\mathbf{p}\|_2^2 + \iota(\mathbf{a} \leq \mathbf{p} \leq \mathbf{b})$
- csoftmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j + \iota(\mathbf{a} \leq \mathbf{p} \leq \mathbf{b})$



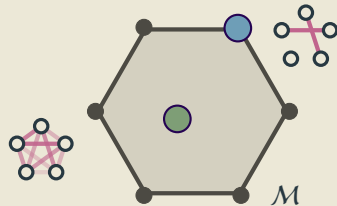
# The structured case: Marginal polytope

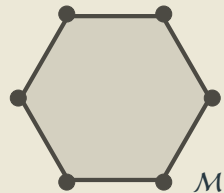
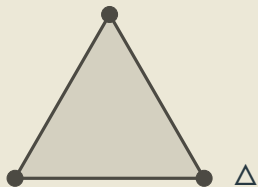
- Each vertex corresponds to one such *bit vector*  $\mathbf{z}$
- Points inside correspond to *marginal distributions*: convex combinations of structured objects

$$\boldsymbol{\mu} = \underbrace{p_1 \mathbf{z}_1 + \dots + p_N \mathbf{z}_N}_{\text{exponentially many terms}}, \quad \mathbf{p} \in \Delta.$$

$$\begin{aligned} p_1 &= 0.2, & \mathbf{z}_1 &= [1, 0, 0, 0, 1, 0, 0, 0, 1] \\ p_2 &= 0.7, & \mathbf{z}_2 &= [0, 0, 1, 0, 0, 1, 1, 0, 0] \\ p_3 &= 0.1, & \mathbf{z}_3 &= [1, 0, 0, 0, 1, 0, 0, 1, 0] \end{aligned}$$

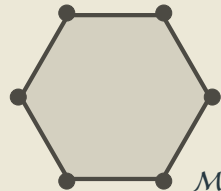
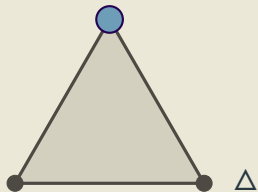
$$\Rightarrow \boldsymbol{\mu} = [.3, 0, .7, 0, .3, .7, .7, .1, .2].$$



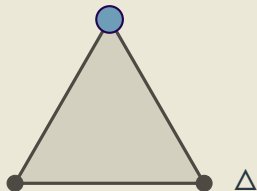




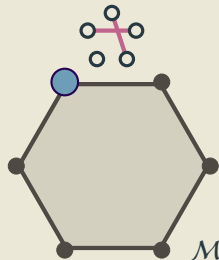
- $$\operatorname{argmax}_{\mathbf{p} \in \Delta} \operatorname{argmax} \mathbf{p}^T \mathbf{s}$$



•  $\text{argmax}_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s}$

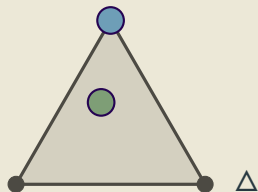


•  $\text{MAP} \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$

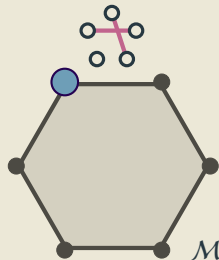


- **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s}$

- **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} + H(\mathbf{p})$

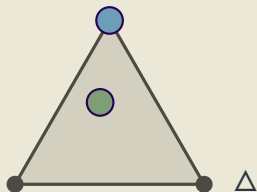


- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$



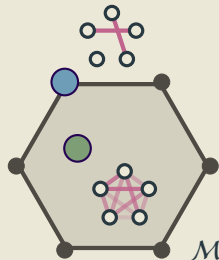
- **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s}$

- **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} + H(\mathbf{p})$



- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$

- **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$



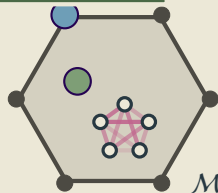
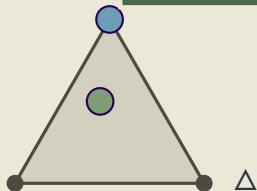
- **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s}$

- **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} + H(\mathbf{p})$

- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$

- **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$

Just like softmax relaxes argmax,  
marginals relax MAP **differentiably!**



- **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s}$

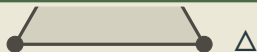
- **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} + H(\mathbf{p})$

- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$

- **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$

Just like softmax relaxes argmax,  
marginals relax MAP **differentiably!**

Unlike argmax/softmax, computation is not obvious!



# Algorithms for specific structures

	Best structure (MAP)	Marginals
Sequence tagging	Viterbi [Rabiner, 1989]	Forward-Backward [Rabiner, 1989]
Constituent trees	CKY [Kasami, 1966, Younger, 1967] [Cocke and Schwartz, 1970]	Probabilistic CKY
Temporal alignments	DTW [Sakoe and Chiba, 1978]	Soft-DTW [Cuturi and Blondel, 2017]
Dependency trees	Max. Spanning Arborescence [Chu and Liu, 1965, Edmonds, 1967]	Matrix-Tree [Kirchhoff, 1847]
Assignments	Kuhn-Munkres [Kuhn, 1955, Jonker and Volgenant, 1987]	#P-complete [Valiant, 1979, Taskar, 2004]

# Algorithms for specific structures

		Best structure (MAP)	Marginals
dyn. prog.	Sequence tagging	Viterbi [Rabiner, 1989]	Forward-Backward [Rabiner, 1989]
	Constituent trees	CKY [Kasami, 1966, Younger, 1967] [Cocke and Schwartz, 1970]	Probabilistic CKY
	Temporal alignments	DTW [Sakoe and Chiba, 1978]	Soft-DTW [Cuturi and Blondel, 2017]
	Dependency trees	Max. Spanning Arborescence [Chu and Liu, 1965, Edmonds, 1967]	Matrix-Tree [Kirchhoff, 1847]
	Assignments	Kuhn-Munkres [Kuhn, 1955, Jonker and Volgenant, 1987]	#P-complete [Valiant, 1979, Taskar, 2004]



# Derivatives of marginals 1: DP

**Dynamic programming:** marginals by **Forward-Backward**, **Inside-Outside**, etc.

# Derivatives of marginals 1: DP

**Dynamic programming:** marginals by **Forward-Backward, Inside-Outside**, etc.

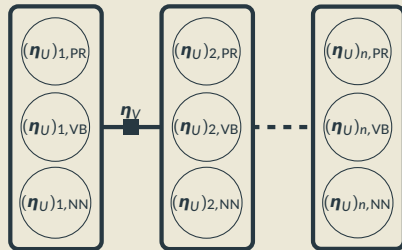
---

## Marginals in a sequence tagging model.

---

```
1 input:  $d$  tags,  $n$  tokens,  $\boldsymbol{\eta}_U \in \mathbb{R}^{n \times d}$ ,  $\boldsymbol{\eta}_V \in \mathbb{R}^{d \times d}$ 
2 initialize  $\boldsymbol{\alpha}_1 = \mathbf{0}$ ,  $\boldsymbol{\beta}_n = \mathbf{0}$ 
3 for  $i \in 2, \dots, n$  do                                # forward log-probabilities
4    $\alpha_{i,k} = \log \sum_{k'} \exp(\alpha_{i-1,k'} + (\boldsymbol{\eta}_U)_{i,k} + (\boldsymbol{\eta}_V)_{k',k})$  for all  $k$ 
5 for  $i \in n-1, \dots, 1$  do                                # backward log-probabilities
6    $\beta_{i,k} = \log \sum_{k'} \exp(\beta_{i+1,k'} + (\boldsymbol{\eta}_U)_{i+1,k'} + (\boldsymbol{\eta}_V)_{k,k'})$  for all  $k$ 
7  $Z = \sum_k \exp \alpha_{n,k}$                                     # partition function
8 return  $\boldsymbol{\mu} = \exp(\boldsymbol{\alpha} + \boldsymbol{\beta} - \log Z)$                 # marginals
```

---



# Derivatives of marginals 1: DP

**Dynamic programming:** marginals by **Forward-Backward, Inside-Outside**, etc.

- Alg. consists of differentiable ops: PyTorch autograd can handle it! (v. bad idea)

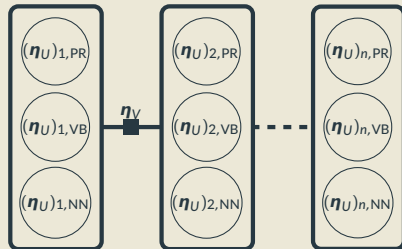
---

Marginals in a sequence tagging model.

---

```
1 input:  $d$  tags,  $n$  tokens,  $\boldsymbol{\eta}_U \in \mathbb{R}^{n \times d}$ ,  $\boldsymbol{\eta}_V \in \mathbb{R}^{d \times d}$ 
2 initialize  $\boldsymbol{\alpha}_1 = \mathbf{0}$ ,  $\boldsymbol{\beta}_n = \mathbf{0}$ 
3 for  $i \in 2, \dots, n$  do                                # forward log-probabilities
4    $\alpha_{i,k} = \log \sum_{k'} \exp(\alpha_{i-1,k'} + (\boldsymbol{\eta}_U)_{i,k} + (\boldsymbol{\eta}_V)_{k',k})$  for all  $k$ 
5 for  $i \in n-1, \dots, 1$  do                                # backward log-probabilities
6    $\beta_{i,k} = \log \sum_{k'} \exp(\beta_{i+1,k'} + (\boldsymbol{\eta}_U)_{i+1,k'} + (\boldsymbol{\eta}_V)_{k,k'})$  for all  $k$ 
7  $Z = \sum_k \exp \alpha_{n,k}$                                     # partition function
8 return  $\boldsymbol{\mu} = \exp(\boldsymbol{\alpha} + \boldsymbol{\beta} - \log Z)$                 # marginals
```

---



# Derivatives of marginals 1: DP

**Dynamic programming:** marginals by **Forward-Backward, Inside-Outside**, etc.

- Alg. consists of differentiable ops: PyTorch autograd can handle it! (v. bad idea)
- Better book-keeping: Li and Eisner [2009], Mensch and Blondel [2018]

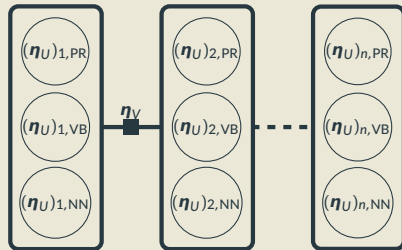
---

Marginals in a sequence tagging model.

---

```
1 input:  $d$  tags,  $n$  tokens,  $\boldsymbol{\eta}_U \in \mathbb{R}^{n \times d}$ ,  $\boldsymbol{\eta}_V \in \mathbb{R}^{d \times d}$ 
2 initialize  $\boldsymbol{\alpha}_1 = \mathbf{0}$ ,  $\boldsymbol{\beta}_n = \mathbf{0}$ 
3 for  $i \in 2, \dots, n$  do                                # forward log-probabilities
4    $\alpha_{i,k} = \log \sum_{k'} \exp(\alpha_{i-1,k'} + (\boldsymbol{\eta}_U)_{i,k} + (\boldsymbol{\eta}_V)_{k',k})$  for all  $k$ 
5 for  $i \in n-1, \dots, 1$  do                                # backward log-probabilities
6    $\beta_{i,k} = \log \sum_{k'} \exp(\beta_{i+1,k'} + (\boldsymbol{\eta}_U)_{i+1,k'} + (\boldsymbol{\eta}_V)_{k,k'})$  for all  $k$ 
7  $Z = \sum_k \exp \alpha_{n,k}$                                     # partition function
8 return  $\boldsymbol{\mu} = \exp(\boldsymbol{\alpha} + \boldsymbol{\beta} - \log Z)$                 # marginals
```

---



# Derivatives of marginals 1: DP

**Dynamic programming:** marginals by **Forward-Backward, Inside-Outside**, etc.

- Alg. consists of differentiable ops: PyTorch autograd can handle it! (v. bad idea)
- Better book-keeping: Li and Eisner [2009], Mensch and Blondel [2018]
- With circular dependencies, this breaks! Can get an approximation Stoyanov et al. [2011]

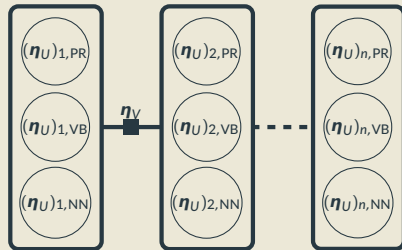
---

Marginals in a sequence tagging model.

---

```
1 input:  $d$  tags,  $n$  tokens,  $\boldsymbol{\eta}_U \in \mathbb{R}^{n \times d}$ ,  $\boldsymbol{\eta}_V \in \mathbb{R}^{d \times d}$ 
2 initialize  $\boldsymbol{\alpha}_1 = \mathbf{0}$ ,  $\boldsymbol{\beta}_n = \mathbf{0}$ 
3 for  $i \in 2, \dots, n$  do                                # forward log-probabilities
4    $\alpha_{i,k} = \log \sum_{k'} \exp(\alpha_{i-1,k'} + (\boldsymbol{\eta}_U)_{i,k} + (\boldsymbol{\eta}_V)_{k',k})$    for all  $k$ 
5 for  $i \in n-1, \dots, 1$  do                                # backward log-probabilities
6    $\beta_{i,k} = \log \sum_{k'} \exp(\beta_{i+1,k'} + (\boldsymbol{\eta}_U)_{i+1,k'} + (\boldsymbol{\eta}_V)_{k,k'})$    for all  $k$ 
7  $Z = \sum_k \exp \alpha_{n,k}$                                 # partition function
8 return  $\boldsymbol{\mu} = \exp(\boldsymbol{\alpha} + \boldsymbol{\beta} - \log Z)$                 # marginals
```

---



# Derivatives of marginals 2: Matrix-Tree

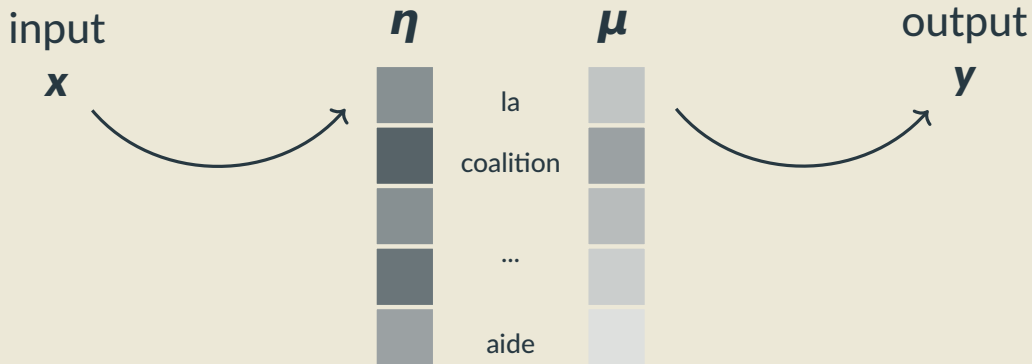
$\mathbf{L}(\mathbf{s})$ : Laplacian of the edge score graph

$$Z = \det \mathbf{L}(\mathbf{s})$$

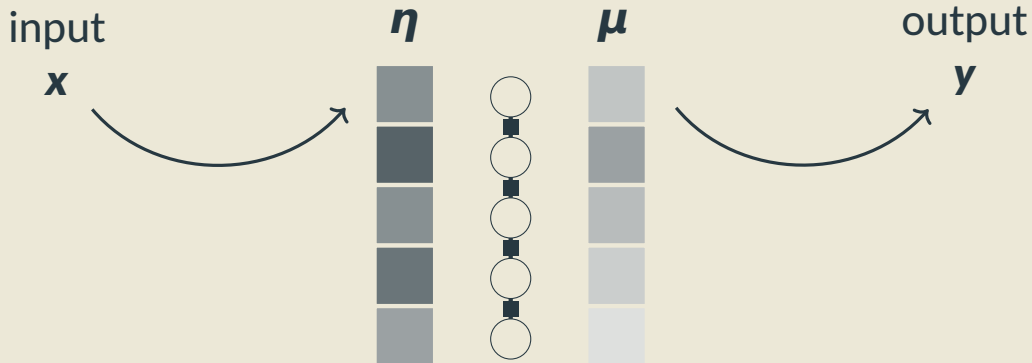
$$\boldsymbol{\mu} = \mathbf{L}(\mathbf{s})^{-1}$$

$$\nabla \boldsymbol{\mu} = \nabla \mathbf{L}^{-1} = \mathbf{L}^{-1} \left( \frac{\partial \mathbf{L}}{\partial \boldsymbol{\eta}} \right) \mathbf{L}^{-1}$$

# Structured Attention Networks

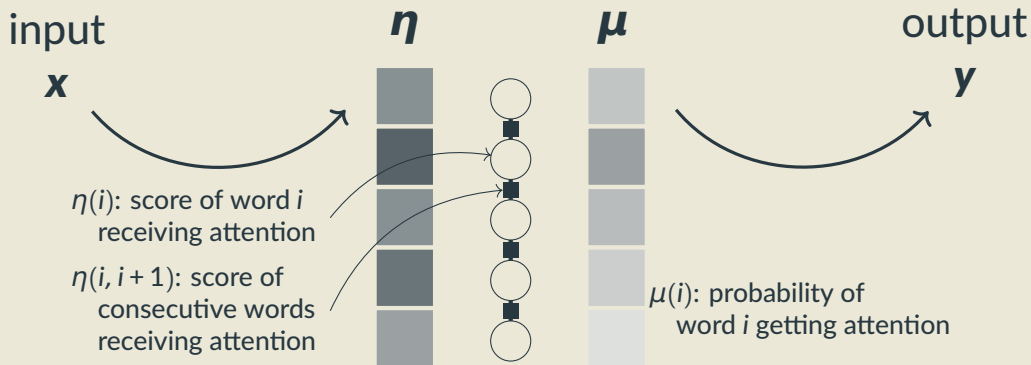


# Structured Attention Networks

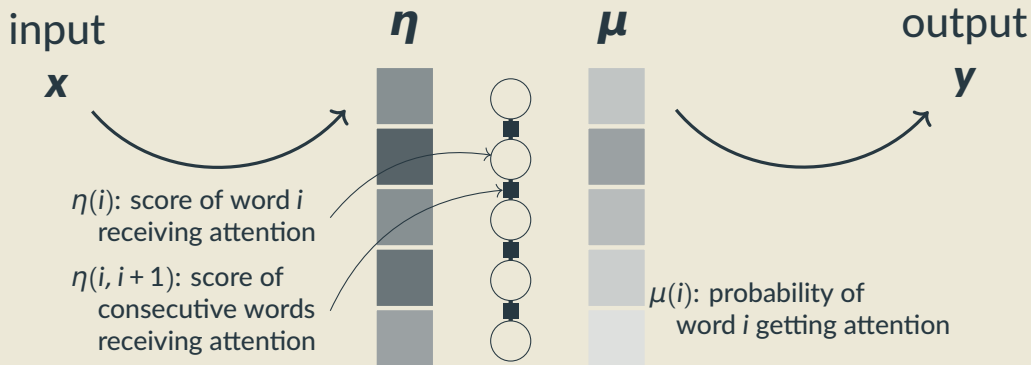




# Structured Attention Networks

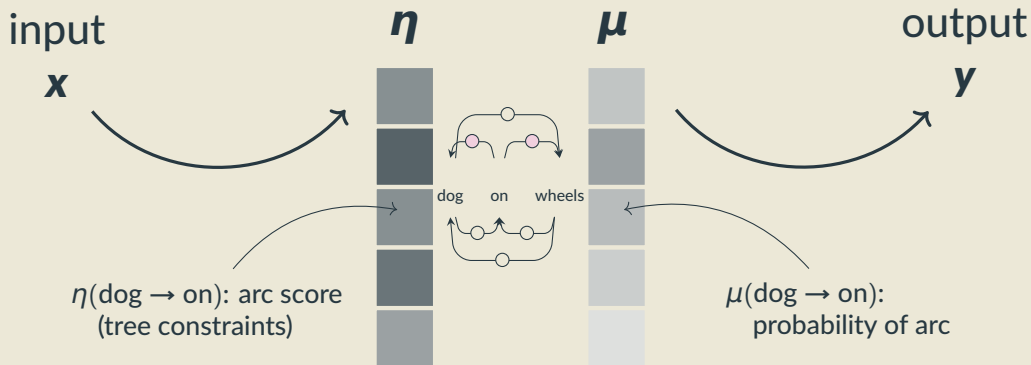


# Structured Attention Networks



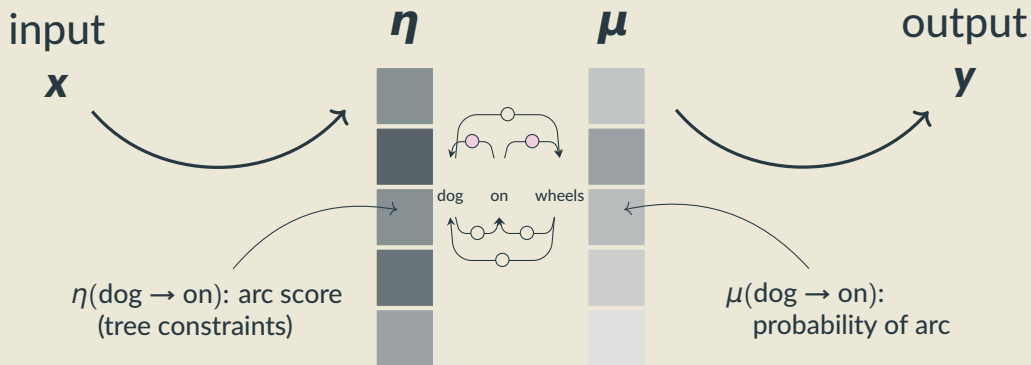
CRF marginals (from *forward-backward*) give attention weights  $\in (0, 1)$

# Structured Attention Networks



CRF marginals (from *forward-backward*) give attention weights  $\in (0, 1)$   
 Similar idea for projective dependency trees with *inside-outside*

# Structured Attention Networks

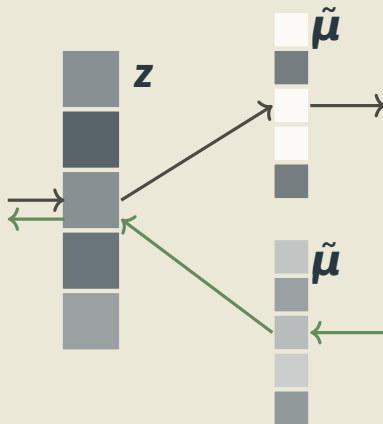


CRF marginals (from *forward-backward*) give attention weights  $\in (0, 1)$   
 Similar idea for projective dependency trees with *inside-outside*  
 and non-projective with the Matrix-Tree theorem [Liu and Lapata, 2018].

# Differentiable Perturb & Parse

## Extending Gumbel-Softmax to structured stochastic models

- Forward pass:  
sample structure  $z$  (approximately)
- Backward pass:  
pretend we did marginal inference



# Back-propagating through marginals

Pros:

# Back-propagating through marginals

Pros:

- Familiar algorithms for NLPers,

# Back-propagating through marginals

Pros:

- Familiar algorithms for NLPers,
- (Structured Attention Networks:) All computations exact.



# Back-propagating through marginals

Pros:

- Familiar algorithms for NLPers,
- (Structured Attention Networks:) All computations exact.

Cons:

- (Structured Attention Networks:) forward pass marginals are dense;  
(fixed by Perturb & MAP, at cost of rough approximation)

# Back-propagating through marginals

Pros:

- Familiar algorithms for NLPers,
- (Structured Attention Networks:) All computations exact.

Cons:

- (Structured Attention Networks:) forward pass marginals are dense;  
(fixed by Perturb & MAP, at cost of rough approximation)
- Efficient & numerically stable back-propagation through DPs is tricky;  
(somewhat alleviated by [Mensch and Blondel, 2018])

# Back-propagating through marginals

## Pros:

- Familiar algorithms for NLPers,
- (Structured Attention Networks:) All computations exact.

## Cons:

- (Structured Attention Networks:) forward pass marginals are dense;  
(fixed by Perturb & MAP, at cost of rough approximation)
- Efficient & numerically stable back-propagation through DPs is tricky;  
(somewhat alleviated by [Mensch and Blondel, 2018])
- Not applicable when marginals are unavailable.

# Back-propagating through marginals

Pros:

- Familiar algorithms for NLPers,
- (Structured Attention Networks:) All computations exact.

Cons:

- (Structured Attention Networks:) forward pass marginals are dense;  
(fixed by Perturb & MAP, at cost of rough approximation)
- Efficient & numerically stable back-propagation through DPs is tricky;  
(somewhat alleviated by [Mensch and Blondel, 2018])
- Not applicable when marginals are unavailable.
- Case-by-case algorithms required, can get tedious.

# Back-propagating through marginals

Pros:

- Familiar algorithms for NLPers,
- (Structured Attention Net

Cons:

- (Structured Attention Net (fixed by Perturb & MA
- Efficient & numerically stable (somewhat alleviated b
- Not applicable when mar
- Case-by-case algorithms

```

procedure BACKPROPINSIDEOUTSIDE( $\theta, p, \nabla_p^T$ )
  for  $s, t = 1, \dots, n, s \neq t$  do  $\triangleright$  Backpropagation uses the identity  $\nabla_p^T = (p \odot \nabla_p^T) / \nabla_p^{T+P}$ 
     $d[s, t] \leftarrow \log p[s, t] \odot \log \nabla_p^T[s, t]$   $\triangleright d \leftarrow \log(p \odot \nabla_p^T)$ 
     $\nabla_p^T, \nabla_p^T \leftarrow -\infty$   $\triangleright$  Initialize inside ( $\nabla_p^T$ ), outside ( $\nabla_p^T$ ) gradients, and log of  $\nabla_p^T$ 
    for  $s = 1, \dots, n - 1$  do  $\triangleright$  Backpropagate  $d$  to  $\nabla_p^T$  and  $\nabla_p^T$ 
      for  $t = s + 1, \dots, n$  do
         $\nabla_p^T[s, t, R, 0], \nabla_p^T[s, t, R, 1] \leftarrow d[s, t]$ 
         $\nabla_p^T[1, n, R, 1] \leftarrow -d[s, t]$ 
        if  $s > 1$  then
           $\nabla_p^T[s, t, L, 0], \nabla_p^T[s, t, L, 1] \leftarrow -d[s, t]$ 
           $\nabla_p^T[1, n, R, 1] \leftarrow -d[s, t]$ 
        for  $k = 1, \dots, n$  do  $\triangleright$  Backpropagate through outside vsp
          for  $s = 1, \dots, n - k$  do
             $t \leftarrow s + k$ 
             $v \leftarrow \nabla_p^T[s, t, R, 0] \odot \beta[s, t, R, 0]$   $\triangleright v, \gamma$  are temporary values
            for  $u = t, \dots, n$  do
               $\nabla_p^T[s, u, R, 0], \nabla_p^T[s, u, R, 1] \leftarrow v \odot \beta[s, u, R, 1] \odot \alpha[t, u, R, 1]$ 
            if  $s > 1$  then
               $v \leftarrow \nabla_p^T[s, t, L, 0] \odot \beta[s, t, L, 0]$ 
              for  $u = 1, \dots, s$  do
                 $\nabla_p^T[s, t, L, 1], \nabla_p^T[s, u, L, 1] \leftarrow v \odot \beta[s, u, L, 1] \odot \alpha[u, s, L, 1]$ 
               $v \leftarrow \nabla_p^T[s, t, L, 1] \odot \beta[s, t, L, 1]$ 
              for  $u = 1, \dots, s$  do
                 $\nabla_p^T[s, u, L, 1], \nabla_p^T[u, s, L, 0] \leftarrow v \odot \beta[s, u, L, 1] \odot \alpha[u, s, L, 1]$ 
            for  $u = 1, \dots, s - 1$  do
               $\gamma \leftarrow \beta[s, t, R, 0] \odot \alpha[u, s - 1, R, 1] \odot \theta_{s,t}$ 
               $\nabla_p^T[s, t, R, 0], \nabla_p^T[u, s - 1, R, 1], \log \nabla_p^T[u, t] \leftarrow v \odot \gamma$ 
               $\gamma \leftarrow \beta[s, t, L, 0] \odot \alpha[u, s - 1, L, 1] \odot \theta_{s,t}$ 
               $\nabla_p^T[s, t, L, 0], \nabla_p^T[u, s - 1, L, 1], \log \nabla_p^T[t, u] \leftarrow v \odot \gamma$ 
             $v \leftarrow \nabla_p^T[s, t, R, 1] \odot \beta[s, t, R, 1]$ 
            for  $u = 1, \dots, s$  do
               $\nabla_p^T[s, t, R, 1], \nabla_p^T[s, u, R, 0] \leftarrow v \odot \beta[s, u, R, 1] \odot \alpha[u, s, R, 0]$ 
            for  $u = t + 1, \dots, n$  do
               $\gamma \leftarrow \beta[s, u, R, 0] \odot \alpha[t + 1, u, L, 1] \odot \theta_{s,t}$ 
               $\nabla_p^T[s, u, R, 0], \nabla_p^T[t + 1, u, L, 1], \log \nabla_p^T[s, u] \leftarrow v \odot \gamma$ 
               $\gamma \leftarrow \beta[s, u, L, 0] \odot \alpha[t + 1, u, L, 1] \odot \theta_{s,t}$ 
               $\nabla_p^T[s, u, L, 0], \nabla_p^T[t + 1, u, L, 1], \log \nabla_p^T[u, t] \leftarrow v \odot \gamma$ 
          for  $k = n, \dots, 1$  do  $\triangleright$  Backpropagate through inside vsp
            for  $s = 1, \dots, n - k$  do
               $t \leftarrow s + k$ 
               $v \leftarrow \nabla_p^T[s, t, R, 1] \odot \alpha[s, t, R, 1]$ 
              for  $u = s + 1, \dots, t$  do
                 $\nabla_p^T[s, t, R, 0], \nabla_p^T[s, t, R, 1] \leftarrow v \odot \alpha[s, u, R, 0] \odot \alpha[s, t, R, 1]$ 
              if  $s > 1$  then
                 $v \leftarrow \nabla_p^T[s, t, L, 1] \odot \alpha[s, t, L, 1]$ 
                for  $u = s, \dots, t - 1$  do
                   $\nabla_p^T[s, u, L, 1], \nabla_p^T[s, t, L, 0] \leftarrow v \odot \alpha[s, u, L, 1] \odot \alpha[s, t, L, 0]$ 
                 $v \leftarrow \nabla_p^T[s, t, L, 0] \odot \alpha[s, t, L, 0]$ 
                for  $u = s, \dots, t - 1$  do
                   $\gamma \leftarrow \alpha[s, u, R, 1] \odot \alpha[u + 1, t, L, 1] \odot \theta_{s,t}$ 
                   $\nabla_p^T[s, u, R, 1], \nabla_p^T[s + 1, t, L, 1], \log \nabla_p^T[t, s] \leftarrow v \odot \gamma$ 
                 $v \leftarrow \nabla_p^T[s, t, R, 0] \odot \alpha[s, t, R, 0]$ 
                for  $u = s, \dots, t - 1$  do
                   $\gamma \leftarrow \alpha[s, u, R, 1] \odot \alpha[u + 1, t, L, 1] \odot \theta_{s,t}$ 
                   $\nabla_p^T[s, u, R, 1], \nabla_p^T[s + 1, t, L, 1], \log \nabla_p^T[s, t] \leftarrow v \odot \gamma$ 
            return sign exp log  $\nabla_p^T$   $\triangleright$  Exponentiate log gradient, multiply by sign, and return  $\nabla_p^T$ 

```

Figure 7: Backpropagation through the inside-outside algorithm to calculate the gradient with respect to the input potentials.  $\nabla_p^T$  denotes the Jacobian of  $\alpha$  with respect to  $\theta$  ( $\nabla_p^T$  is the gradient with respect to  $\theta$ ).  $\alpha, \beta \leftarrow \infty$  means  $\alpha = \alpha \odot \infty$  and  $\beta = \beta \odot \infty$ .

xact.

inals are dense;  
nation)

ugh DPs is tricky;

# Back-propagating through marginals

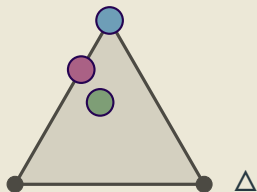
## Pros:

- Familiar algorithms for NLPers,
- (Structured Attention Networks:) All computations exact.

## Cons:

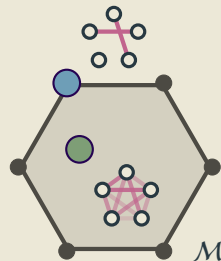
- (Structured Attention Networks:) forward pass marginals are dense;  
(fixed by Perturb & MAP, at cost of rough approximation)
- Efficient & numerically stable back-propagation through DPs is tricky;  
(somewhat alleviated by [Mensch and Blondel, 2018])
- Not applicable when marginals are unavailable.
- Case-by-case algorithms required, can get tedious.

- **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s}$
- **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} + H(\mathbf{p})$
- **sparsemax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} - 1/2 \|\mathbf{p}\|^2$

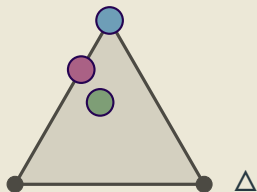


● **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$

● **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$



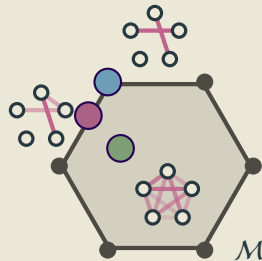
- **argmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s}$
- **softmax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} + H(\mathbf{p})$
- **sparsemax**  $\arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{s} - 1/2 \|\mathbf{p}\|^2$



- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta}$

- **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} + \tilde{H}(\boldsymbol{\mu})$

- **SparseMAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^T \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$





# SparseMAP solution

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

$$= \begin{array}{c} \circ \quad \circ \\ \diagup \quad \diagdown \\ \circ \quad \circ \end{array} = .6 \begin{array}{c} \circ \quad \circ \\ \diagup \quad \diagdown \\ \circ \quad \circ \end{array} + .4 \begin{array}{c} \circ \quad \circ \\ \diagup \quad \diagdown \\ \circ \quad \circ \end{array}$$

( $\boldsymbol{\mu}^*$  is unique, but may have multiple decompositions  $\boldsymbol{p}$ . Active Set recovers a sparse one.)

# Algorithms for SparseMAP

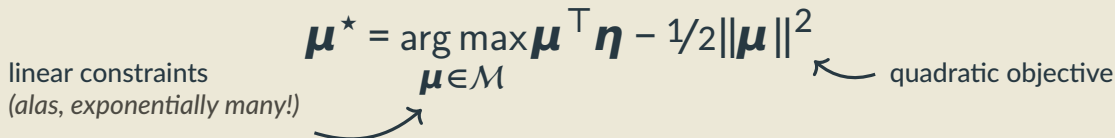
$$\boldsymbol{\mu}^{\star} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

linear constraints  
(*alas, exponentially many!*)

quadratic objective

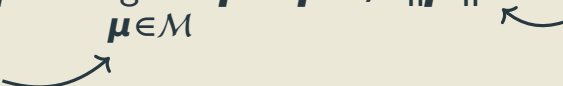


# Algorithms for SparseMAP

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

linear constraints  
(*alas, exponentially many!*)

quadratic objective



## Conditional Gradient

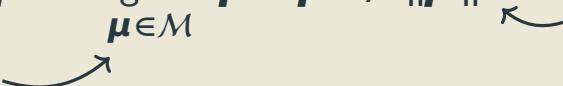
[Frank and Wolfe, 1956, Lacoste-Julien and Jaggi, 2015]

# Algorithms for SparseMAP

linear constraints  
(*alas, exponentially many!*)

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

quadratic objective



## Conditional Gradient

[Frank and Wolfe, 1956, Lacoste-Julien and Jaggi, 2015]

- select a new corner of  $\mathcal{M}$

# Algorithms for SparseMAP

linear constraints  
(*alas, exponentially many!*)

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

quadratic objective

## Conditional Gradient

[Frank and Wolfe, 1956, Lacoste-Julien and Jaggi, 2015]

- select a new corner of  $\mathcal{M}$

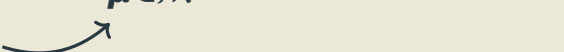
$$\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \underbrace{(\boldsymbol{\eta} - \boldsymbol{\mu}^{(t-1)})}_{\tilde{\boldsymbol{\eta}}}$$

# Algorithms for SparseMAP

linear constraints  
(*alas, exponentially many!*)

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

quadratic objective



## Conditional Gradient

[Frank and Wolfe, 1956, Lacoste-Julien and Jaggi, 2015]

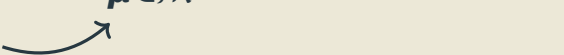
- select a new corner of  $\mathcal{M}$
- update the (sparse) coefficients of  $\boldsymbol{p}$ 
  - Update rules: vanilla, away-step, pairwise

# Algorithms for SparseMAP

linear constraints  
(*alas, exponentially many!*)

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

quadratic objective



## Conditional Gradient

[Frank and Wolfe, 1956, Lacoste-Julien and Jaggi, 2015]

- select a new corner of  $\mathcal{M}$
- update the (sparse) coefficients of  $\boldsymbol{p}$ 
  - Update rules: vanilla, away-step, pairwise
  - Quadratic objective: **Active Set**

a.k.a. Min-Norm Point, [Wolfe, 1976]

[Martins et al., 2015, Nocedal and Wright, 1999,

Vinyes and Obozinski, 2017]



# Algorithms for SparseMAP

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

linear constraints  
(*alas, exponentially many!*)

quadratic objective

## Conditional Gradient

[Frank and Wolfe, 1956, Lacoste-Julien and Jaggi, 2015]

- select a new corner
- update the (sparse)
  - Update rules: vanilla
  - Quadratic objective:

Active Set achieves  
**finite & linear** convergence!

a.k.a. Min-Norm Point, [Wolfe, 1976]

[Martins et al., 2015, Nocedal and Wright, 1999,

Vinyes and Obozinski, 2017]

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

linear constraints  
(*alas, exponentially many!*)

quadratic objective

## Conditional Gradient

[Frank and Wolfe, 1956, Lacoste-Julien and Jaggi, 2015]

- select a new corner of  $\mathcal{M}$
- update the (sparse) coefficients of  $\mathbf{p}$ 
  - Update rules: vanilla, away-step, pairwise
  - Quadratic objective: **Active Set**  
a.k.a. Min-Norm Point, [Wolfe, 1976]

[Martins et al., 2015, Nocedal and Wright, 1999,

Vinyes and Obozinski, 2017]

## Backward pass

$$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \text{ is sparse}$$

# Algorithms for SparseMAP

linear constraints  
(*alas, exponentially many!*)

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

quadratic objective

## Conditional Gradient

[Frank and Wolfe, 1956, Lacoste-Julien and Jaggi, 2015]

- select a new corner of  $\mathcal{M}$
- update the (sparse) coefficients of  $\mathbf{p}$ 
  - Update rules: vanilla, away-step, pairwise
  - Quadratic objective: **Active Set**  
a.k.a. Min-Norm Point, [Wolfe, 1976]

[Martins et al., 2015, Nocedal and Wright, 1999,

Vinyes and Obozinski, 2017]

## Backward pass

$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$  is sparse  
computing  $\left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}\right)^\top d\mathbf{y}$   
takes  $\mathcal{O}(\dim(\boldsymbol{\mu}) \text{nnz}(\mathbf{p}^*))$

# Algorithms for SparseMAP

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{\eta} - 1/2 \|\boldsymbol{\mu}\|^2$$

linear constraints  
(*alas, exponentially many!*)

quadratic objective

Condition

pass

Completely modular: just add MAP

[Frank and Wolfe, 1956]

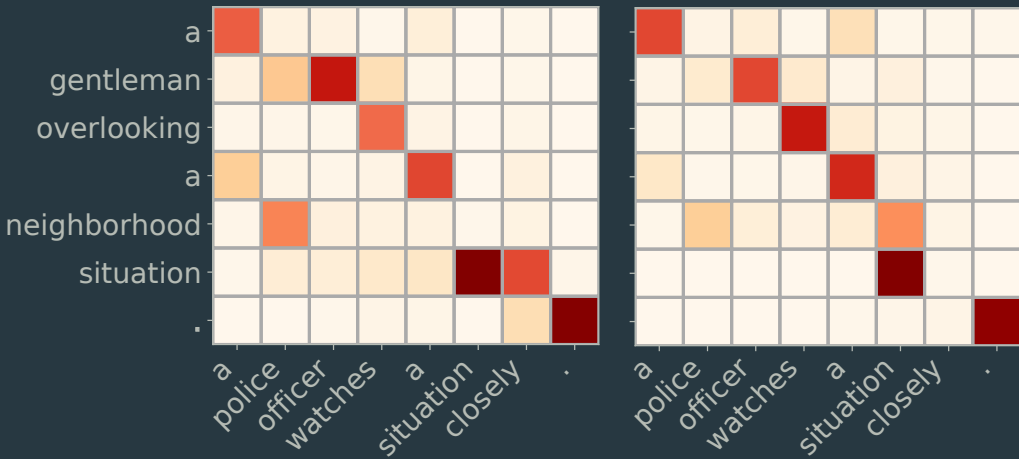
- select a new candidate  $\mathbf{c}$
- update the (sparse) coefficients or  $\mathbf{p}$ 
  - Update rules: vanilla, away-step, pairwise
  - Quadratic objective: **Active Set**

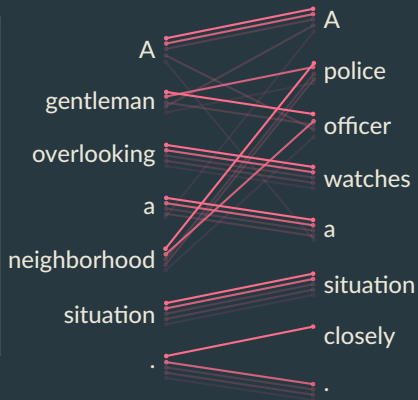
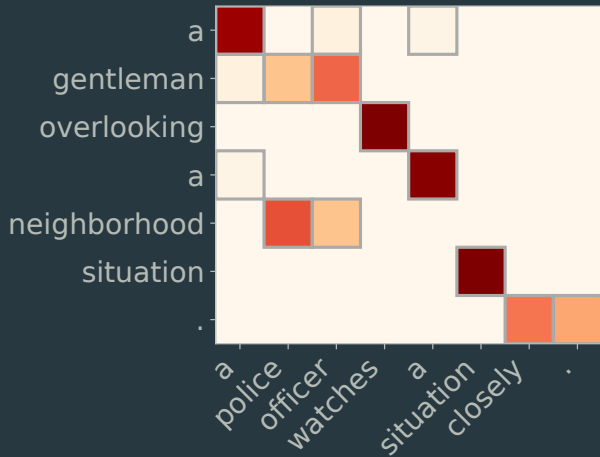
a.k.a. Min-Norm Point, [Wolfe, 1976]

[Martins et al., 2015, Nocedal and Wright, 1999,

Vinyes and Obozinski, 2017]

computing  $\left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}\right)^\top d\boldsymbol{\eta}$   
takes  $\mathcal{O}(\dim(\boldsymbol{\mu}) \text{nnz}(\mathbf{p}^*))$





# Overview

$$\mathbb{E}_{\pi_{\boldsymbol{\theta}}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Score Function Estimator
- Straight Through-Gumbel
- Perturb-and-Parse

$$L(\arg \max_{\mathbf{z}} \pi_{\boldsymbol{\theta}}(\mathbf{z} | x))$$

- Straight Through (flavors)
- SPIGOT

$$L(\mathbb{E}_{\pi_{\boldsymbol{\theta}}(\mathbf{z}|x)}[\mathbf{z}])$$

- Structured Attn. Networks
- SparseMAP

# Structured latent variables without sampling

$$p(y \mid x) = \sum_{\mathbf{z} \in \mathcal{Z}} p(y \mid \mathbf{z}, x) \pi(\mathbf{z} \mid x)$$




# Structured latent variables without sampling

$$p(y \mid x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\boldsymbol{\phi}}(y \mid \mathbf{z}, x) \pi_{\boldsymbol{\theta}}(\mathbf{z} \mid x)$$

# Structured latent variables without sampling

$$p(y \mid x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y \mid \mathbf{z}, x) \pi_{\theta}(\mathbf{z} \mid x)$$

e.g., a TreeLSTM defined by  $\mathbf{z}$



# Structured latent variables without sampling

$$p(y \mid x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\boldsymbol{\phi}}(y \mid \mathbf{z}, x) \pi_{\boldsymbol{\theta}}(\mathbf{z} \mid x)$$

e.g., a TreeLSTM defined by  $\mathbf{z}$

parsing model,  
using some score  $\boldsymbol{\theta}(\mathbf{z}; x)$

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

Exponentially large sum!

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

idea 1

idea 2

idea 3

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum_{h \in \mathcal{H}}$$

idea 1

idea 2

idea 3

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \theta}$$

idea 1

idea 2

idea 3

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \theta}$$

idea 1  $\pi_{\theta}(\mathbf{z} | x) = 1$  if  $\mathbf{z} = \mathbf{z}^*$  else 0

argmax

idea 2

idea 3



# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \theta}$$

idea 1  $\pi_{\theta}(\mathbf{z} | x) = 1$  if  $\mathbf{z} = \mathbf{z}^*$  else 0

argmax



idea 2

idea 3

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \theta}$$

idea 1  $\pi_{\theta}(\mathbf{z} | x) = 1$  if  $\mathbf{z} = \mathbf{z}^*$  else 0

argmax

idea 2

idea 3



# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \theta}$$

😊 ☹️

idea 1  $\pi_{\theta}(\mathbf{z} | x) = 1$  if  $\mathbf{z} = \mathbf{z}^*$  else 0

argmax

idea 2  $\pi_{\theta}(\mathbf{z} | x) \propto \exp(\text{score}_{\theta}(\mathbf{z}; x))$

softmax

idea 3

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \theta}$$

idea 1  $\pi_{\theta}(\mathbf{z} | x) = 1$  if  $\mathbf{z} = \mathbf{z}^*$  else 0

argmax



idea 2  $\pi_{\theta}(\mathbf{z} | x) \propto \exp(\text{score}_{\theta}(\mathbf{z}; x))$

softmax



idea 3

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \theta}$$

idea 1  $\pi_{\theta}(\mathbf{z} | x) = 1$  if  $\mathbf{z} = \mathbf{z}^*$  else 0

argmax



idea 2  $\pi_{\theta}(\mathbf{z} | x) \propto \exp(\text{score}_{\theta}(\mathbf{z}; x))$

softmax



idea 3

# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

idea 1  $\pi_{\theta}(\mathbf{z} | x) = 1$  if  $\mathbf{z} = \mathbf{z}^*$  else 0

idea 2  $\pi_{\theta}(\mathbf{z} | x) \propto \exp(\text{score}_{\theta}(\mathbf{z}; x))$

idea 3

argmax

softmax

SparseMAP

$$\sum_{h \in \mathcal{H}} \frac{\partial p(y | x)}{\partial \theta}$$



# Structured latent variables without sampling

sum over  
all possible trees

e.g., a TreeLSTM defined by  $\mathbf{z}$

$$p(y | x) = \sum_{\mathbf{z} \in \mathcal{Z}} p_{\phi}(y | \mathbf{z}, x) \pi_{\theta}(\mathbf{z} | x)$$

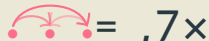
How to define  $\pi_{\theta}$ ?

parsing model,  
using some score  $\theta(\mathbf{z}; x)$

$$\sum \frac{\partial p(y | x)}{\partial \theta}$$


Different from SAN:  
classifier need not be defined outside  $\mathcal{Z}$ !  
Perturb-&-MAP sampling is not exact in general 😞

# Structured latent variables without sampling





# Structured latent variables without sampling

 $= .7x$

 $+ .3x$

 $+ 0x$   $+ \dots$

# Structured latent variables without sampling

$$\begin{aligned}
 & \text{Diagram 1} = .7 \times \text{Diagram 2} + .3 \times \text{Diagram 3} + 0 \times \text{Diagram 4} + \dots \\
 p(y \mid x) = & .7 \times p_{\phi}(y \mid \text{Diagram 1}) + .3 \times p_{\phi}(y \mid \text{Diagram 2})
 \end{aligned}$$

The diagrams consist of three red dots. Diagram 1 has a curved arrow from the first dot to the second, and a straight arrow from the second dot to the third. Diagram 2 has a curved arrow from the first dot to the second, and a curved arrow from the second dot to the third. Diagram 3 has a curved arrow from the first dot to the second, and a straight arrow from the second dot to the third. Diagram 4 has a curved arrow from the first dot to the second, and a curved arrow from the second dot to the third.

## Stanford Natural Language Inference (Accuracy)

[Kim et al., 2017]

Simple Attention 86.2

Structured Attention 86.8

[Liu and Lapata, 2018]

100D SAN - 86.8

Yogatama et al

100D RL-SPINN 80.5

[Choi et al., 2018]

100D ST Gumbel-Tree 82.6

300D - 85.6

600D - 86.0

[Corro and Titov, 2019b]

Latent Tree + 1 GCN - 85.2

Latent Tree + 2 GCN - 86.2

## Stanford Sentiment (Accuracy)

Socher et al

Bigram Naive Bayes 83.1

[Nicolae et al., 2018b]

TreeLSTM w/ CoreNLP 83.2

TreeLSTM w/ SparseMAP 84.7

[Corro and Titov, 2019b]

GCN w/ CoreNLP 83.8

GCN w/ Perturb-and-MAP 84.6

# V. Conclusions

# Is it syntax?!

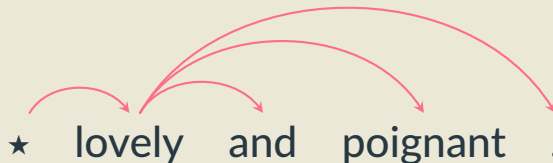
- Unlike e.g. unsupervised parsing, the structures we learn are guided by a **downstream objective** (typically discriminative).
- They don't typically resemble grammatical structure (yet) [Williams et al., 2018]  
(future work: more inductive biases and constraints?)

# Is it syntax?!

- Unlike e.g. unsupervised parsing, the structures we learn are guided by a **downstream objective** (typically discriminative).
- They don't typically resemble grammatical structure (yet) [Williams et al., 2018]  
(future work: more inductive biases and constraints?)
- Common to compare latent structures with parser outputs. But is it always desirable?

# Syntax vs. Composition Order

CoreNLP parse,  $p = 21.4\%$

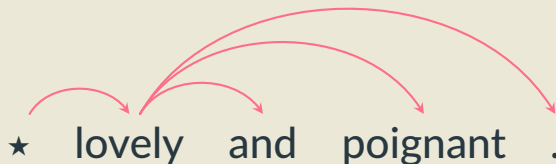


# Syntax vs. Composition Order

$p = 22.6\%$



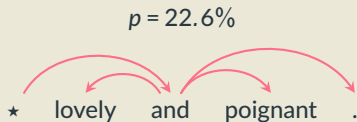
CoreNLP parse,  $p = 21.4\%$



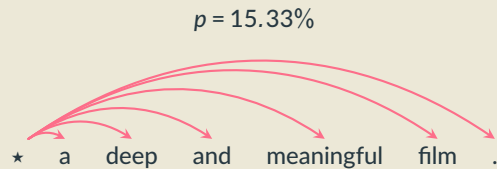
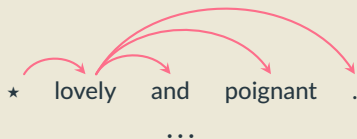
...



# Syntax vs. Composition Order



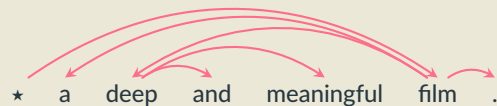
CoreNLP parse,  $p = 21.4\%$



$p = 15.27\%$



...  
CoreNLP parse,  $p = 0\%$



# Overview

$$\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[L(\mathbf{z})]$$

- Score Function Estimator
- Straight Through-Gumbel
- Perturb-and-Parse
- SparseMAP

## Model restrictions:

- $\text{dom } L$  may be only  $\mathcal{Z}$ ,
- $\nabla_{\mathbf{z}} f$  need not exist!

$$L(\arg \max_{\mathbf{z}} \pi_{\theta}(\mathbf{z} | x))$$

- Straight Through (flavors)
- SPIGOT

- $L(\mathbf{z})$  with  $\mathbf{z} \in \mathcal{Z}$  in forward
- needs (relaxed)  $\nabla_{\mathbf{z}} f$  in backward.

$$L(\mathbb{E}_{\pi_{\theta}(\mathbf{z}|x)}[\mathbf{z}])$$

- Structured Attn. Networks
- SparseMAP

- $L(\mathbf{z})$  must be relaxed and differentiable.
- (sparsity gets us closer to  $\mathcal{Z}$ ).

# Conclusions

- Latent structure models are desirable for interpretability, structural bias, and higher predictive power with fewer parameters.
- Stochastic latent variables can be dealt with RL or straight-through gradients.
- Deterministic argmax requires surrogate gradients (e.g. SPIGOT).
- Continuous relaxations of argmax include SANs and SparseMAP.
- Intuitively, some of these different methods are trying to do similar things or require the same building blocks (e.g. SPIGOT and SparseMAP).

# References I

- Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- Mathieu Blondel, André FT Martins, and Vlad Niculae. Learning classifiers with Fenchel-Young losses: Generalized entropies, margins, and algorithms. 2019.
- Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- Peter F Brown, Vincent J Della Pietra, Stephen A Della Pietra, and Robert L Mercer. The mathematics of statistical machine translation: Parameter estimation. *Computational linguistics*, 19(2):263–311, 1993.
- Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, Hui Jiang, and Diana Inkpen. Enhanced LSTM for natural language inference. In *Proc. of ACL*, 2017.
- Jihun Choi, Kang Min Yoo, and Sang-goo Lee. Learning to compose task-specific tree structures. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Yoeng-Jin Chu and Tseng-Hong Liu. On the shortest arborescence of a directed graph. *Science Sinica*, 14: 1396–1400, 1965.

# References II

- William John Cocke and Jacob T Schwartz. *Programming languages and their compilers*. Courant Institute of Mathematical Sciences., 1970.
- Shay B Cohen, Karl Stratos, Michael Collins, Dean P Foster, and Lyle Ungar. Spectral learning of latent-variable pcfgs. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 223–231, 2012.
- Caio Corro and Ivan Titov. Differentiable Perturb-and-Parse: Semi-Supervised Parsing with a Structured Variational Autoencoder. In *Proc. ICLR*, 2019a.
- Caio Corro and Ivan Titov. Learning latent trees with stochastic perturbations and differentiable dynamic programming. In *Proc. ACL*, 2019b. URL <https://caio-corro.fr/pdf/acl2019.pdf>.
- Marco Cuturi and Mathieu Blondel. Soft-dtw: a differentiable loss function for time-series. In *Proc. ICML*, 2017.
- Jack Edmonds. Optimum branchings. *J. Res. Nat. Bur. Stand.*, 71B:233–240, 1967.
- Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Nav. Res. Log.*, 3(1-2):95–110, 1956.
- Geoffrey Hinton. Neural networks for machine learning. In *Coursera video lectures*, 2012.

# References III

- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- Roy Jonker and Anton Volgenant. A shortest augmenting path algorithm for dense and sparse linear assignment problems. *Computing*, 38(4):325–340, 1987.
- Tadao Kasami. An efficient recognition and syntax-analysis algorithm for context-free languages. *Coordinated Science Laboratory Report no. R-257*, 1966.
- Yoon Kim, Carl Denton, Loung Hoang, and Alexander M Rush. Structured attention networks. In *Proc. of ICLR*, 2017.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Gustav Kirchhoff. Ueber die auflösung der gleichungen, auf welche man bei der untersuchung der linearen vertheilung galvanischer ströme geführt wird. *Annalen der Physik*, 148(12):497–508, 1847.
- Harold W Kuhn. The Hungarian method for the assignment problem. *Nav. Res. Log.*, 2(1-2):83–97, 1955.
- Simon Lacoste-Julien and Martin Jaggi. On the global linear convergence of Frank-Wolfe optimization variants. In *Proc. of NeurIPS*, 2015.

# References IV

- Zhifei Li and Jason Eisner. First-and second-order expectation semirings with applications to minimum-risk training on translation forests. In *Proc. of EMNLP*, 2009.
- Yang Liu and Mirella Lapata. Learning structured text representations. *TACL*, 6:63–75, 2018.
- Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- Jean Maillard and Stephen Clark. Latent tree learning with differentiable parsers: Shift-reduce parsing and chart parsing. *arXiv preprint arXiv:1806.00840*, 2018.
- Chaitanya Malaviya, Pedro Ferreira, and André F. T. Martins. Sparse and constrained attention for neural machine translation. In *Proc. of ACL*, 2018.
- André FT Martins and Ramón Fernandez Astudillo. From softmax to sparsemax: A sparse model of attention and multi-label classification. In *Proc. of ICML*, 2016.
- André FT Martins, Mário AT Figueiredo, Pedro MQ Aguiar, Noah A Smith, and Eric P Xing. AD3: Alternating directions dual decomposition for MAP inference in graphical models. *JMLR*, 16(1):495–545, 2015.

# References V

- Arthur Mensch and Mathieu Blondel. Differentiable dynamic programming for structured prediction and attention. In *Proc. of ICML*, 2018.
- Vlad Niculae and Mathieu Blondel. A regularized framework for sparse and structured neural attention. In *Proc. of NeurIPS*, 2017.
- Vlad Niculae, André FT Martins, Mathieu Blondel, and Claire Cardie. SparseMAP: Differentiable sparse structured inference. In *Proc. of ICML*, 2018a.
- Vlad Niculae, André FT Martins, and Claire Cardie. Towards dynamic computation graphs via sparse latent structure. In *Proc. of EMNLP*, 2018b.
- Jorge Nocedal and Stephen Wright. *Numerical Optimization*. Springer New York, 1999.
- George Papandreou and Alan L Yuille. Perturb-and-map random fields: Using discrete optimization to learn and sample from energy models. In *2011 International Conference on Computer Vision*, pages 193–200. IEEE, 2011.
- Hao Peng, Sam Thomson, and Noah A Smith. Backpropagating through structured argmax using a SPIGOT. In *Proc. ACL*, 2018.
- Ben Peters, Vlad Niculae, and André FT Martins. Sparse sequence-to-sequence models. 2019.



# References VI

- Slav Petrov and Dan Klein. Discriminative log-linear grammars with latent variables. In *Advances in neural information processing systems*, pages 1153–1160, 2008.
- Ariadna Quattoni, Sybor Wang, Louis-Philippe Morency, Michael Collins, and Trevor Darrell. Hidden conditional random fields. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (10):1848–1852, 2007.
- Lawrence R. Rabiner. A tutorial on Hidden Markov Models and selected applications in speech recognition. *P. IEEE*, 77(2):257–286, 1989.
- Hiroaki Sakoe and Seibi Chiba. Dynamic programming algorithm optimization for spoken word recognition. *IEEE Trans. on Acoustics, Speech, and Sig. Proc.*, 26:43–49, 1978.
- Veselin Stoyanov, Alexander Ropson, and Jason Eisner. Empirical risk minimization of graphical model parameters given approximate inference, decoding, and model structure. In *AISTATS*, 2011.
- Ben Taskar. *Learning structured prediction models: A large margin approach*. PhD thesis, Stanford University, 2004.
- Constantino Tsallis. Possible generalization of Boltzmann-Gibbs statistics. *Journal of Statistical Physics*, 52:479–487, 1988.
- Leslie G Valiant. The complexity of computing the permanent. *Theor. Comput. Sci.*, 8(2):189–201, 1979.

# References VII

- Marina Vinyes and Guillaume Obozinski. Fast column generation for atomic norm regularization. In *Proc. of AISTATS*, 2017.
- M. Wainwright and M. Jordan. *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers, 2008.
- Adina Williams, Andrew Drozdov, and Samuel R Bowman. Do latent tree learning models identify meaningful structure in sentences? *TACL*, 2018.
- Philip Wolfe. Finding the nearest point in a polytope. *Mathematical Programming*, 11(1):128–149, 1976.
- Daniel H Younger. Recognition and parsing of context-free languages in time  $n^3$ . *Information and Control*, 10(2): 189–208, 1967.