

# Learning with Sparse Latent Structure

#### Vlad Niculae

Instituto de Telecomunicações

Work with: André Martins, Claire Cardie, Mathieu Blondel

🖸 github.com/vene/sparsemap 🛛 💆 @vnfrombucharest

#### **Structured Prediction**

. . .







### **Structured Prediction**



# **Structured Prediction**







. . .

#### **Latent Structure Models**











#### \*freeze frame\*



с<sub>1</sub> с<sub>2</sub>

• • •

с<sub>N</sub>













$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$$



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$$



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$$



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$$



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$$



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$$



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$$



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$$



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \boldsymbol{0}$$

# Argmax vs. Softmax



$$\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \operatorname{diag}(\boldsymbol{p}) - \boldsymbol{p}\boldsymbol{p}^{\mathsf{T}}$$

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

# $\Delta = \{ \boldsymbol{p} \in \mathbb{R}^N : \, \boldsymbol{p} \ge \boldsymbol{0}, \, \, \boldsymbol{1}^\top \boldsymbol{p} = \boldsymbol{1} \}$











1.5

0.5

N = 2

1

N = 3

1.5







$$\max_{j} \boldsymbol{\theta}_{j} = \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^{\top} \boldsymbol{\theta}$$



$$\max_{j} \boldsymbol{\theta}_{j} = \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^{\top} \boldsymbol{\theta}$$



$$\max_{j} \boldsymbol{\theta}_{j} = \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^{\top} \boldsymbol{\theta}$$



$$\max_{j} \boldsymbol{\theta}_{j} = \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^{\top} \boldsymbol{\theta}$$


#### Variational Form of Argmax

$$\max_{j} \boldsymbol{\theta}_{j} = \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^{\top} \boldsymbol{\theta}$$



#### Variational Form of Argmax

$$\max_{j} \boldsymbol{\theta}_{j} = \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^{\top} \boldsymbol{\theta}$$



#### Variational Form of Argmax

$$\max_{j} \boldsymbol{\theta}_{j} = \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^{\top} \boldsymbol{\theta}$$





 $\theta_1$ 

 $\wedge$ 

$$\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta}) = \arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$

$$\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta}) = rg\max \boldsymbol{p}^{\mathsf{T}} \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$
  
 $\boldsymbol{p} \in \Delta$ 

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



$$\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta}) = \arg \max \boldsymbol{p}^{\mathsf{T}} \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$
  
 $\boldsymbol{p} \in \Delta$ 

- argmax:  $\Omega(\mathbf{p}) = \mathbf{0}$
- softmax:  $\Omega(\mathbf{p}) = \sum_{j} p_{j} \log p_{j}$



$$\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta}) = \arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$
  
 $\boldsymbol{p} \in \Delta$ 

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



(Martins and Astudillo, 2016)



softmax



sparsemax

sparsemax(
$$\boldsymbol{\theta}$$
) = arg max  $\boldsymbol{p}^{\top} \boldsymbol{\theta} - \frac{1}{2} \|\boldsymbol{p}\|_{2}^{2}$   
= arg min  $\|\boldsymbol{p} - \boldsymbol{\theta}\|_{2}^{2}$   
 $\boldsymbol{p} \in \Delta$ 

sparsemax(
$$\boldsymbol{\theta}$$
) = arg max  $\boldsymbol{p}^{\top}\boldsymbol{\theta} - \frac{1}{2} \|\boldsymbol{p}\|_{2}^{2}$   
= arg min  $\|\boldsymbol{p} - \boldsymbol{\theta}\|_{2}^{2}$   
=  $\boldsymbol{p} \in \Delta$ 

**Computation:** 

 $\boldsymbol{p}^{\star} = [\boldsymbol{\theta} - \tau \mathbf{1}]_{+}$  $\boldsymbol{\theta}_{i} > \boldsymbol{\theta}_{j} \Rightarrow p_{i} \ge p_{j}$ O(d) via partial sort

(Held et al., 1974; Brucker, 1984; Condat, 2016)

sparsemax(
$$\boldsymbol{\theta}$$
) = arg max  $\boldsymbol{p}^{\top}\boldsymbol{\theta} - \frac{1}{2}\|\boldsymbol{p}\|_{2}^{2}$   
 $\boldsymbol{p} \in \Delta$   
= arg min  $\|\boldsymbol{p} - \boldsymbol{\theta}\|_{2}^{2}$   
 $\boldsymbol{p} \in \Delta$   
mputation:  
=  $[\mathbf{0}, \mathbf{z}^{1}]$ 

$$\boldsymbol{p}^{\star} = [\boldsymbol{\theta} - \tau \mathbf{1}]_{+}$$
$$\boldsymbol{\theta}_{i} > \boldsymbol{\theta}_{j} \Rightarrow \boldsymbol{p}_{i} \ge \boldsymbol{p}_{j}$$
$$\mathcal{O}(d) \text{ via partial sort}$$

Cor

(Held et al., 1974; Brucker, 1984; Condat, 2016)

$$\begin{aligned} \mathbf{J}_{\text{sparsemax}} &= \text{diag}(\mathbf{s}) - \frac{1}{|\mathcal{S}|} \mathbf{s} \mathbf{s}^{\top} \\ \text{where } \mathcal{S} &= \{j : p_{j}^{\star} > 0\}, \\ s_{j} &= [\![j \in \mathcal{S}]\!] \end{aligned}$$

(Martins and Astudillo, 2016)

sparsemax(
$$\boldsymbol{\theta}$$
) = arg max  $\boldsymbol{p}^{\top} \boldsymbol{\theta} - 1/2 \|\boldsymbol{p}\|_{2}^{2}$   
= arg min  $\|\boldsymbol{p} - \boldsymbol{\theta}\|_{2}^{2}$   
**Computation:** Backward pass:  
 $\boldsymbol{p}^{\star} = [( argmin differentiation \\ \boldsymbol{\theta}_{i} > \boldsymbol{\theta}_{j} \\ (d) \text{ via} ]$ 
 $(Gould et al., 2016; Amos and Kolter, 2017)$ 
 $(g(s) - \frac{1}{|S|}ss^{\top} : p_{j}^{\star} > 0\}, \quad z \in S$ 

(Held et al., 1974; Brucker, 1984; Condat, 2016)

(Martins and Astudillo, 2016)



fusedmax ?!

$$\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta}) = \arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$
  
 $\boldsymbol{p} \in \Delta$ 

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



$$\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta}) = \arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$
  
 $\boldsymbol{p} \in \Delta$ 

- argmax:  $\Omega(\mathbf{p}) = \mathbf{0}$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax:  $\Omega(\mathbf{p}) = \frac{1}{2} ||\mathbf{p}||_2^2$ fusedmax:  $\Omega(\mathbf{p}) = \frac{1}{2} ||\mathbf{p}||_2^2 + \sum_j |p_j - p_{j-1}|$ csparsemax:  $\Omega(\mathbf{p}) = \frac{1}{2} ||\mathbf{p}||_2^2 + \iota(\mathbf{a} \le \mathbf{p} \le \mathbf{b})$



#### Fusedmax

$$fusedmax(\boldsymbol{\theta}) = \arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} - \frac{1}{2} \|\boldsymbol{p}\|_{2}^{2} - \sum_{2 \le j \le d} |p_{j} - p_{j-1}|$$
$$= \arg \min \|\boldsymbol{p} - \boldsymbol{\theta}\|_{2}^{2} + \sum_{2 \le j \le d} |p_{j} - p_{j-1}|$$
$$prox_{fused}(\boldsymbol{\theta}) = \arg \min \|\boldsymbol{p} - \boldsymbol{\theta}\|_{2}^{2} + \sum_{2 \le j \le d} |p_{j} - p_{j-1}|$$

**Proposition:** fusedmax( $\boldsymbol{\theta}$ ) = sparsemax(prox<sub>fused</sub>( $\boldsymbol{\theta}$ ))

(Niculae and Blondel, 2017)



(INICUIAE and BIONGEI, ZUIT)

$$\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta}) = \arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} - \Omega(\boldsymbol{p})$$
  
 $\boldsymbol{p} \in \Delta$ 

- argmax:  $\Omega(\mathbf{p}) = \mathbf{0}$
- softmax:  $\Omega(\mathbf{p}) = \sum_j p_j \log p_j$
- sparsemax:  $\Omega(\mathbf{p}) = \frac{1}{2} ||\mathbf{p}||_2^2$ fusedmax:  $\Omega(\mathbf{p}) = \frac{1}{2} ||\mathbf{p}||_2^2 + \sum_j |p_j - p_{j-1}|$ csparsemax:  $\Omega(\mathbf{p}) = \frac{1}{2} ||\mathbf{p}||_2^2 + \iota(\mathbf{a} \le \mathbf{p} \le \mathbf{b})$



finally

is essentially a (very high-dimensional) argmax



is essentially a (very high-dimensional) argmax



is essentially a (very high-dimensional) argmax



# Factorization Into Parts $\boldsymbol{\theta} = \mathbf{A}^{\top} \boldsymbol{\eta}$

# Factorization Into Parts $\theta = A^T \eta$

wheels dog on \*

∗→dog	[ 1	0	0		[ .1]	
on→dog	0	1	1		.2	
wheels→dog	0	0	0		1	
×→on	0	1	1		.3	
<b>A</b> = dog→on	1	0	0	 η =	.8	
wheels→on	0	0	0		.1	
*→wheels	0	0	0		3	
dog→wheels	0	1	0		.2	
on→wheels	1	0	1		1 ]	

# **Factorization Into Parts** $\boldsymbol{\theta} = \mathbf{A}^{\top} \boldsymbol{n}$

.1

.2

-.1

.3

.8

.1

-.3

.2

-.1

**n** =



1

0

0

0

1

0

0

0

1

∗→dog

on→dog

∗→on

\*→wheels

dog→wheels

on→wheels

wheels→dog

dog→on

wheels→on

**A** =

0

1 0

0

0

1

0

1

1 1

0

0 0

0

1

0

dog	~/	hond
on	$\checkmark$	ор
wheels		wielen

dog-hond	1		0	0			.1
dog—op	0		1	1			.2
dog—wielen	0		0	0			1
on-hond	0		0	0			.3
on—op	1		0	0		η=	.8
on-wielen	0		1	1			.1
eels-hond	0		1	0			3
eels—op	0		0	0			.2
eels-wielen	1		0	1			1
	dog-hond dog-op dog-wielen on-hond on-op on-wielen eels-hond eels-op eels-wielen	dog-hond1dog-op0dog-wielen0on-hond0on-op1on-wielen0eels-hond0eels-op0eels-wielen1	dog-hond1dog-op0dog-wielen0on-hond0on-op1on-wielen0eels-hond0eels-op0eels-wielen1	dog-hond         1         0           dog-op         0         1           dog-wielen         0         0           on-hond         0         0           on-op         1            on-wielen         0         1           eels-hond         0         1           eels-op         0         0           on-wielen         0         1	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c cccc} dog-hond & 1 & 0 & 0 \\ dog-op & 0 & 1 & 1 \\ dog-wielen & 0 & 0 & 0 \\ on-hond & 0 & 0 & 0 \\ on-op & 1 & \dots & 0 & 0 & \dots \\ on-wielen & 0 & 1 & 1 \\ eels-hond & 0 & 1 & 0 \\ eels-op & 0 & 0 & 0 \\ eels-wielen & 1 & 0 & 1 \\ \end{array}$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$











$$\mathcal{M} := \operatorname{conv} \left\{ \boldsymbol{a}_h : h \in \mathcal{H} \right\}$$
$$= \left\{ \boldsymbol{A} \boldsymbol{p} : \boldsymbol{p} \in \Delta \right\}$$





$$\mathcal{M} := \operatorname{conv} \left\{ \boldsymbol{a}_h : h \in \mathcal{H} \right\}$$
$$= \left\{ \boldsymbol{A} \boldsymbol{p} : \boldsymbol{p} \in \Delta \right\}$$
$$= \left\{ \mathbb{E}_{H \sim \boldsymbol{p}} \boldsymbol{a}_H : \boldsymbol{p} \in \Delta \right\}$$





















**MAP** arg max $\boldsymbol{\mu}^{\mathsf{T}}\boldsymbol{\eta}$  $\mu \in \mathcal{M}$ 

# *e.g.* dependency parsing → max. spanning tree matching → the Hungarian algorithm





• **argmax**  $\arg \max p^{\mathsf{T}} \boldsymbol{\theta}$  $p \in \Delta$ 

• softmax  $\arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} + H(\boldsymbol{p})$  $\boldsymbol{p} \in \Delta$ 







• **argmax** arg max  $p^{\top} \theta$  $p \in \Delta$ 

• softmax  $\arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} + H(\boldsymbol{p})$  $\boldsymbol{p} \in \Delta$ 










#### e.g. dependency parsing $\rightarrow$ the Matrix-Tree theorem

(Koo et al., 2007; D. A. Smith and N. A. Smith, 2007; McDonald and Satta, 2007)

As attention: (Liu and Lapata, 2018)







• **argmax** arg max 
$$p^{\top} \theta$$
  
 $p \in \Delta$ 

• softmax  $\arg \max \boldsymbol{p}^\top \boldsymbol{\theta} + H(\boldsymbol{p})$  $\boldsymbol{p} \in \Delta$ 

• sparsemax  $\arg \max \boldsymbol{p}^{\top} \boldsymbol{\theta} - \frac{1}{2} \|\boldsymbol{p}\|^2$  $\boldsymbol{p} \in \Delta$ 







(Niculae, Martins, Blondel, and Cardie, 2018)

2

• argmax 
$$\arg \max p^{T} \theta$$
  
 $p \in \Delta$ 
• softmax  $\arg \max p^{T} \theta + H(p)$   
 $p \in \Delta$ 
• softmax  $\arg \max p^{T} \theta + H(p)$   
 $p \in \Delta$ 
• sparsemax  $\arg \max p^{T} \theta - \frac{1}{2} \|p\|^{2}$ 
• sparseMAP  $\arg \max \mu^{T} \eta - \frac{1}{2} \|\mu\|$   
 $\mu \in \mathcal{M}$ 





#### **SparseMAP Solution**

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

$$= \overset{\circ}{0}\overset{\circ}{0}\overset{\circ}{0} = .6\overset{\circ}{0}\overset{\circ}{0}\overset{\circ}{0} + .4\overset{\circ}{0}\overset{\circ}{0}\overset{\circ}{0}$$

=  $\mathbf{A}\mathbf{p}^*$  with very sparse  $\mathbf{p}^* \in \Delta^N$ 

### **Algorithms for SparseMAP**

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

Algorithms for SparseMAP  

$$\mu^{\star} = \arg \max \mu^{\top} \eta - 1/2 \|\mu\|^{2}$$
Ilinear constraints  
(alas, exponentially many!) (alas, exponentially many!) (alas, exponentially many!)

Algorithms for SparseMAP  

$$\mu^{\star} = \arg \max \mu^{\top} \eta - 1/2 \|\mu\|^{2}$$
Inter constraints  
(alas, exponentially many!) (under a constraints) (under a constraints)

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

Algorithms for SparseMAP  

$$\mu^{\star} = \arg \max \mu^{\top} \eta - \frac{1}{2} \|\mu\|^{2}$$
(alas, exponentially many!) (alas, exponentially many!) (blue) (constraints) (constraints

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

#### • select a new corner of ${\mathcal M}$

Algorithms for SparseMAP  

$$\mu^{*} = \arg \max \mu^{\top} \eta - \frac{1}{2} \|\mu\|^{2}$$
(alas, exponentially many!) (alas, exponentially many!) (blue) (constraints) (constraints

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

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

$$\boldsymbol{a}_{\boldsymbol{y}^{\star}} = \arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \underbrace{(\boldsymbol{\eta} - \boldsymbol{\mu}^{(t-1)})}_{\widetilde{\boldsymbol{\eta}}}$$

# Algorithms for SparseMAP $\mu^{\star} = \arg \max \mu^{\top} \eta - 1/2 \|\mu\|^{2}$ Innear constraints (alas, exponentially many!) (alas, exponentially many!) (black)

#### **Conditional Gradient**

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

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

# Algorithms for SparseMAP $\mu^{\star} = \arg \max \mu^{\top} \eta - 1/2 \|\mu\|^{2}$ Inter constraints (alas, exponentially many!) (alas, exponentially many!) (blue) (blue) (classified of the second second

#### **Conditional Gradient**

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

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

#### • Quadratic objective: Active Set

(Nocedal and Wright, 1999, Ch. 16.4 & 16.5) (Wolfe, 1976; Vinyes and Obozinski, 2017)

Algorithms for SparseMAP  

$$\mu^{\star} = \arg \max \mu^{\top} \eta - 1/2 \|\mu\|^{2}$$
Inter constraints  
(alas, exponentially many!) (under a constraint) (und

(Frank and Wolfe, 1956; Lacost

- select a new corner
- update the (sparse)

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

- Update rules: vanilla, away-step, pairwise
- Quadratic objective: Active Set

(Nocedal and Wright, 1999, Ch. 16.4 & 16.5) (Wolfe, 1976; Vinyes and Obozinski, 2017)

# Algorithms for SparseMAP $\mu^{\star} = \arg \max \mu^{\top} \eta - 1/2 \|\mu\|^{2}$ Innear constraints (alas, exponentially many!) (alas, exponentially many!)

#### **Conditional Gradient**

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

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

#### • Quadratic objective: Active Set

(Nocedal and Wright, 1999, Ch. 16.4 & 16.5) (Wolfe, 1976; Vinyes and Obozinski, 2017)

#### **Backward pass**

Algorithms for SparseMAP  

$$\mu^{\star} = \arg \max \mu^{\top} \eta - \frac{1}{2} \|\mu\|^{2}$$
Intear constraints  
(alas, exponentially many!)

(Frank and Wolfe, 1956; Lacoste-Julien and Jaggi, 2015)

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

#### • Quadratic objective: Active Set

(Nocedal and Wright, 1999, Ch. 16.4 & 16.5) (Wolfe, 1976; Vinyes and Obozinski, 2017)

#### **Backward pass**

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



- Update rules: vanilla, away-step, pairwise
- Quadratic objective: Active Set

(Nocedal and Wright, 1999, Ch. 16.4 & 16.5) (Wolfe, 1976; Vinyes and Obozinski, 2017)  $\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \text{ is sparse} \\ \text{computing } \left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}\right)^{\mathsf{T}} \boldsymbol{d} \boldsymbol{y} \\ \text{takes } O(\dim(\boldsymbol{\mu}) \operatorname{nnz}(\boldsymbol{p}^*))$ 

NLI premise: A gentleman overlooking a neighborhood situation. hypothesis: A police officer watches a situation closely.



(Model: ESIM (Chen et al., 2017))

NLI premise: A gentleman overlooking a neighborhood situation. hypothesis: A police officer watches a situation closely.



(Model: ESIM (Chen et al., 2017))

NLI premise: A gentleman overlooking a neighborhood situation. hypothesis: A police officer watches a situation closely.



(Model: ESIM (Chen et al., 2017))

NLI premise: A gentleman overlooking a neighborhood situation. hypothesis: A police officer watches a situation closely.



(Proposed model: global matching)

### In code:

# Xp: (n\_prem x k) #  $Xh: (n_hypo \times k)$ Z = Xp a Xh.t() Up = softmax(Z, dim=1)Uh = softmax(Z, dim=0)Xp = cat([Xp, Up a Xh])Xh = cat([Xh, Uh.t() a Xp])

### In code:

# Xp: (n prem x k)# Xp: (n prem x k) # Xh: (n hupo x k) # Xh: (n hupo x k)Z = Xp a Xh.t()Z = Xp a Xh.tUp = softmax(Z, dim=1) $U = sparsemap_matching(Z)$ Uh = softmax(Z. dim=0)Xp = cat([Xp, Up a Xh]) Xp = cat([Xp, U a Xh]) $Xh = cat([Xh, Uh, t() a Xp]) \quad Xh = cat([Xh, U, t() a Xp])$ 





76.5%



а gentleman overlooking а neighborhood situation Poliofivations stored . POILOFICET DES STOTEN



## Dynamically inferring the computation graph

# So far: a structured hidden layer $\mathbb{E}_{H}[\boldsymbol{a}_{H}]$

Network must handle "soft" combinations of structures. Fine for attention, but can be limiting.













(Tai et al., 2015)



The bears eat the pretty ones

### Latent Dependency TreeLSTM

(Niculae, Martins, and Cardie, 2018)



### Latent Dependency TreeLSTM

(Niculae, Martins, and Cardie, 2018)

$$p(y|x) = \sum_{h \in \mathcal{H}} p(y \mid h, x) p(h \mid x)$$



input

Х
$$p(y \mid x) = \sum_{h \in \mathcal{H}} p (y \mid h, x) p (h \mid x)$$

$$p(y \mid x) = \sum_{h \in \mathcal{H}} p_{\phi}(y \mid h, x) p_{\pi}(h \mid x)$$

e.g., a TreeLSTM defined by h  $p(y \mid x) = \sum p_{\phi}^{\checkmark}(y \mid h, x) p_{\pi}(h \mid x)$ h∈H





## Exponentially large sum!



idea 1

idea 2

idea 3



idea 3



idea 3















## **SparseMAP**



## **SparseMAP**



## **SparseMAP**

# $p(y \mid x) = .7 \quad p_{\phi}(y \mid \widehat{f} \cdot \widehat{f} \cdot \widehat{f}) + .3 \quad p_{\phi}(y \mid \widehat{f} \cdot \widehat{f} \cdot \widehat{f}) + .3 \quad p_{\phi}(y \mid \widehat{f} \cdot \widehat{f} \cdot \widehat{f})$

85% -	
84% -	
83% -	
82% -	
81% -	

00	0	0/								
ou	17	°O	_							











### Sentiment classification (SST)





Sentence pair classification (P, H)  $p(y \mid P, H) = \sum_{h_P \in \mathcal{H}(P)} \sum_{h_H \in \mathcal{H}(H)} p_{\phi}(y \mid h_P, h_H) p_{\pi}(h_P \mid P) p_{\pi}(h_H \mid H)$ 



### **Reverse dictionary lookup**

given word description, predict word embedding (Hill et al., 2016) instead of p(y | x), we model  $\mathbb{E}_{p_m} g(x) = \sum_{h \in \mathcal{H}} g(x; h) p_m(h | x)$ 



### **Reverse dictionary lookup**

		(definitio	ns)		<i>·</i> ·		(concep	ots)	
accuracy@10	38% –				accuracy@10	38% -			
	36% –					36% -			
	34% -					34% -			
	32% -					32% -			
	30% —	LTR	Flat	Latent		30% —	LTR	Flat	Latent



### Natural Language Inference (SNLI)



#### **Reverse dictionary lookup**



### (concepts)



## Syntax vs. Composition Order

CoreNLP parse, p = 21.4%



## Syntax vs. Composition Order

p = 22.6%lovely and poignant  $\star$ CoreNLP parse, p = 21.4%

 $\star$  lovely and poignant .

• • •

## Syntax vs. Composition Order

*p* = 15.33%



## Conclusions

Differentiable & sparse structured inference

Generic, extensible algorithms

Interpretable structured attention

### Dynamically-inferred computation graphs



**O** github.com/vene/sparsemap 









# **Extra slides**

## Acknowledgements



This work was supported by the European Research Council (ERC StG DeepSPIN 758969) and by the Fundação para a Ciência e Tecnologia through contract UID/EEA/50008/2013.

Some icons by Dave Gandy and Freepik via flaticon.com.

## Danskin's Theorem

Let  $\phi : \mathbb{R}^d \times \mathcal{Z} \to \mathbb{R}, \mathcal{Z} \subset \mathbb{R}^d$  compact.  $\partial \max_{z \in \mathcal{Z}} \phi(x, z) = \operatorname{conv} \{ \nabla_x \phi(x, z^*) \mid z^* \in \arg\max_{z \in \mathcal{Z}} \phi(x, z) \}.$ 

Example: maximum of a vector

## Danskin's Theorem

Let 
$$\phi : \mathbb{R}^d \times \mathcal{Z} \to \mathbb{R}, \mathcal{Z} \subset \mathbb{R}^d$$
 compact.  
 $\partial \max_{\mathbf{z} \in \mathcal{Z}} \phi(\mathbf{x}, \mathbf{z}) = \operatorname{conv} \{ \nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{z}^*) \mid \mathbf{z}^* \in \arg\max_{\mathbf{z} \in \mathcal{Z}} \phi(\mathbf{x}, \mathbf{z}) \}.$ 

Example: maximum of a vector

$$\partial \max_{j \in [d]} \theta_j = \partial \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{\theta}$$
$$= \partial \max_{\boldsymbol{p} \in \Delta} \phi(\boldsymbol{p}, \boldsymbol{\theta})$$
$$= \operatorname{conv} \{ \nabla_{\boldsymbol{\theta}} \phi(\boldsymbol{p}^*, \boldsymbol{\theta})$$
$$= \operatorname{conv} \{ \boldsymbol{p}^* \}$$
# Danskin's Theorem

Let 
$$\phi : \mathbb{R}^d \times \mathbb{Z} \to \mathbb{R}, \mathbb{Z} \subset \mathbb{R}^d$$
 compact.  
 $\partial \max_{\mathbf{z} \in \mathbb{Z}} \phi(\mathbf{x}, \mathbf{z}) = \operatorname{conv} \{ \nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{z}^*) \mid \mathbf{z}^* \in \arg\max_{\mathbf{z} \in \mathbb{Z}} \phi(\mathbf{x}, \mathbf{z}) \}.$ 

#### Example: maximum of a vector



## **Example: Source Sentence with Three Words**



#### e.g., fertility constraints for NMT



constrained softmax: (Martins and Kreutzer, 2017) constrained sparsemax: (Malaviya et al., 2018)



#### **Structured Output Prediction**

rseMAP 
$$L_{\mathbf{A}}(\boldsymbol{\eta}, \boldsymbol{\bar{\mu}}) = \max_{\boldsymbol{\mu} \in \mathcal{M}} \{ \boldsymbol{\eta}^{\mathsf{T}} \boldsymbol{\mu} - \frac{1}{2} \| \boldsymbol{\mu} \|^2 \}$$
$$- \boldsymbol{\eta}^{\mathsf{T}} \boldsymbol{\bar{\mu}} + \frac{1}{2} \| \boldsymbol{\bar{\mu}} \|^2$$

Spa

Instance of a structured Fenchel-Young loss, like CRF, SVM, etc. (Blondel, Martins, and Niculae, 2019)

#### **Structured Output Prediction**

SparseMAP 
$$L_{\mathbf{A}}(\boldsymbol{\eta}, \bar{\boldsymbol{\mu}}) = \max_{\boldsymbol{\mu} \in \mathcal{M}} \{ \boldsymbol{\eta}^{\top} \boldsymbol{\mu} - \frac{1}{2} \| \boldsymbol{\mu} \|^2 \}$$
$$- \boldsymbol{\eta}^{\top} \bar{\boldsymbol{\mu}} + \frac{1}{2} \| \bar{\boldsymbol{\mu}} \|^2$$
cost-SparseMAP 
$$L_{\mathbf{A}}^{\rho}(\boldsymbol{\eta}, \bar{\boldsymbol{\mu}}) = \max_{\boldsymbol{\mu} \in \mathcal{M}} \{ \boldsymbol{\eta}^{\top} \boldsymbol{\mu} - \frac{1}{2} \| \boldsymbol{\mu} \|^2 + \rho(\boldsymbol{\mu}, \bar{\boldsymbol{\mu}}) \}$$
$$- \boldsymbol{\eta}^{\top} \bar{\boldsymbol{\mu}} + \frac{1}{2} \| \bar{\boldsymbol{\mu}} \|^2$$

Instance of a structured Fenchel-Young loss, like CRF, SVM, etc. (Blondel, Martins, and Niculae, 2019)





Universal Dependencies dataset

#### **Sparse Structured Output Prediction**

As models train, inference gets sparser!



# **Sparse Structured Output Prediction**

Inference captures linguistic ambiguity!



# **Sparse Structured Output Prediction**

Inference captures linguistic ambiguity!



### **References** I

- Amos, Brandon and J. Zico Kolter (2017). "OptNet: Differentiable optimization as a layer in neural networks". In: *Proc. of ICML*.
- Bertsekas, Dimitri P (1999). Nonlinear Programming. Athena Scientific Belmont.
- Blondel, Mathieu, André FT Martins, and Vlad Niculae (2019). "Learning with Fenchel-Young Losses". In: *preprint arXiv*:1901.02324.
- Brucker, Peter (1984). "An O(n) algorithm for quadratic knapsack problems". In: Operations Research Letters 3.3, pp. 163–166.
- Chen, Qian et al. (2017). "Enhanced LSTM for natural language inference". In: Proc. of ACL.
- Condat, Laurent (2016). "Fast projection onto the simplex and the  $l_1$  ball". In: Mathematical Programming 158.1-2, pp. 575–585.
- Danskin, John M (1966). "The theory of max-min, with applications". In: SIAM Journal on Applied Mathematics 14.4, pp. 641–664.
- Dantzig, George B, Alex Orden, and Philip Wolfe (1955). "The generalized simplex method for minimizing a linear form under linear inequality restraints". In: *Pacific Journal of Mathematics* 5.2, pp. 183–195.

# **References II**

Frank, Marguerite and Philip Wolfe (1956). "An algorithm for quadratic programming". In: *Nav. Res. Log.* 3.1-2, pp. 95–110.

Gould, Stephen et al. (2016). "On differentiating parameterized argmin and argmax problems with application to bi-level optimization". In: *preprint arXiv:1607.05447*.

- Held, Michael, Philip Wolfe, and Harlan P Crowder (1974). "Validation of subgradient optimization". In: *Mathematical Programming* 6.1, pp. 62–88.
- Hill, Felix et al. (2016). "Learning to understand phrases by embedding the dictionary". In: TACL 4.1, pp. 17–30.
- Kim, Yoon et al. (2017). "Structured attention networks". In: Proc. of ICLR.
- Koo, Terry et al. (2007). "Structured prediction models via the matrix-tree theorem". In: *Proc. of EMNLP*.
- Lacoste-Julien, Simon and Martin Jaggi (2015). "On the global linear convergence of Frank-Wolfe optimization variants". In: *Proc. of NeurIPS*.

#### **References III**

- Liu, Yang and Mirella Lapata (2018). "Learning structured text representations". In: TACL 6, pp. 63–75.
- Malaviya, Chaitanya, Pedro Ferreira, and André F. T. Martins (2018). "Sparse and constrained attention for neural machine translation". In: *Proc. of ACL*.
- Martins, André FT and Ramón Fernandez Astudillo (2016). "From softmax to sparsemax: A sparse model of attention and multi-label classification". In: *Proc. of ICML*.
- Martins, André FT and Julia Kreutzer (2017). "Learning What's Easy: Fully Differentiable Neural Easy-First Taggers". In: *Proc. of EMNLP*, pp. 349–362.
- McDonald, Ryan T and Giorgio Satta (2007). "On the complexity of non-projective data-driven dependency parsing". In: *Proc. of ICPT.*
- Niculae, Vlad and Mathieu Blondel (2017). "A regularized framework for sparse and structured neural attention". In: *Proc. of NeurIPS*.
- Niculae, Vlad, André FT Martins, Mathieu Blondel, et al. (2018). "SparseMAP: Differentiable sparse structured inference". In: *Proc. of ICML*.

#### **References IV**

- Niculae, Vlad, André FT Martins, and Claire Cardie (2018). "Towards dynamic computation graphs via sparse latent structure". In: *Proc. of EMNLP*.
- Nocedal, Jorge and Stephen Wright (1999). Numerical Optimization. Springer New York.
- Rabiner, Lawrence R. (1989). "A tutorial on Hidden Markov Models and selected applications in speech recognition". In: *P. IEEE* 77.2, pp. 257–286.
- Smith, David A and Noah A Smith (2007). "Probabilistic models of nonprojective dependency trees". In: *Proc. of EMNLP*.
- Tai, Kai Sheng, Richard Socher, and Christopher D Manning (2015). "Improved semantic representations from tree-structured Long Short-Term Memory networks". In: *Proc. of* ACL-IJCNLP.
- Taskar, Ben (2004). "Learning structured prediction models: A large margin approach". PhD thesis. Stanford University.
- Tibshirani, Robert et al. (2005). "Sparsity and smoothness via the fused lasso". In: Journal of the Royal Statistical Society: Series B (Statistical Methodology) 67.1, pp. 91–108.

#### **References V**

Valiant, Leslie G (1979). "The complexity of computing the permanent". In: *Theor. Comput. Sci.* 8.2, pp. 189–201.

Vinyes, Marina and Guillaume Obozinski (2017). "Fast column generation for atomic norm regularization". In: *Proc. of AISTATS*.

Wolfe, Philip (1976). "Finding the nearest point in a polytope". In: *Mathematical Programming* 11.1, pp. 128–149.