# Differentiable dynamic programming for structured prediction and attention

Arthur Mensch
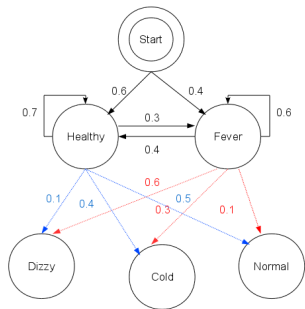
École Normale Supérieure
Département de Mathématiques Appliquées
Paris, France

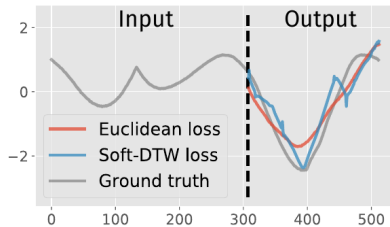January 29, 2019

# Dynamic programming in machine learning



(Fig: Wikipedia)

Belief propagation
Viterbi algorithm



Input    Output

— Euclidean loss
— Soft-DTW loss
— Ground truth

(Fig: Cuturi & Blondel)

Dynamic time warping



Value iteration

# New layers for deep predictive modelling

**Dynamic programming**
- Function evaluation on a big space
- From evaluation on smaller space (divide and conquer)
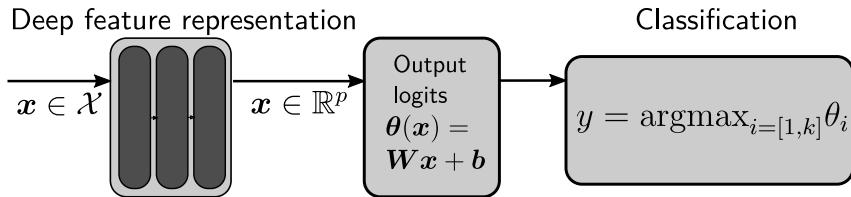
**Modern machine learning**
- Modular models, trainable with gradient descent

**This talk: Principled differentiable dynamic programming layers**
- Application to complex predictive models (e.g. attention mechanisms)
- Sparse output probabilities in structured setting

# Modular structured prediction: potentials + linear programming

**Classification:** $\mathcal{Y} = [1, k]$

Deep feature representation

Classification



$\boldsymbol{x} \in \mathcal{X}$

$\boldsymbol{x} \in \mathbb{R}^p$

Output logits
$\boldsymbol{\theta}(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b}$

$y = \mathrm{argmax}_{i=[1,k]} \theta_i$

# Modular structured prediction: potentials + linear programming

**Classification:** $\mathcal{Y} = [1, k]$

Deep feature representation

Classification

$\boldsymbol{x} \in \mathcal{X}$    $\boldsymbol{x} \in \mathbb{R}^p$

Output logits
$\boldsymbol{\theta}(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b}$

$y = \mathrm{argmax}_{i=[1,k]} \theta_i$

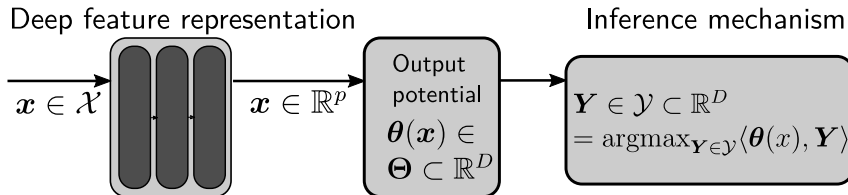**Structured output ?** $\mathcal{Y} \subset \mathbb{R}^D$ (edges of a polytope), *e. g.* a tag sequence

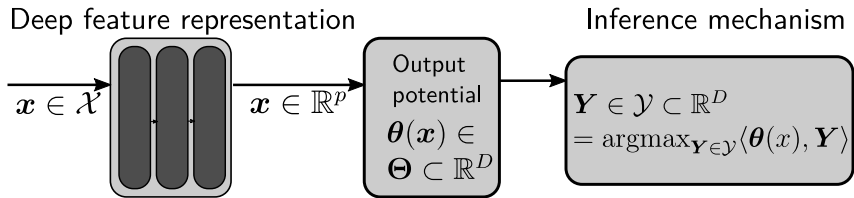# Modular structured prediction: potentials + linear programming

**Classification:** $\mathcal{Y} = [1, k]$

Deep feature representation

Classification

$\boldsymbol{x} \in \mathcal{X}$ → $\boldsymbol{x} \in \mathbb{R}^p$ →

Output logits
$\boldsymbol{\theta}(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b}$

→ $y = \mathrm{argmax}_{i=[1,k]} \theta_i$

**Structured output ?** $\mathcal{Y} \subset \mathbb{R}^D$ (edges of a polytope), *e.g.* a tag sequence

Deep feature representation

Inference mechanism

$\boldsymbol{x} \in \mathcal{X}$ → $\boldsymbol{x} \in \mathbb{R}^p$ →

Output potential
$\boldsymbol{\theta}(\boldsymbol{x}) \in \boldsymbol{\Theta} \subset \mathbb{R}^D$

→ $\boldsymbol{Y} \in \mathcal{Y} \subset \mathbb{R}^D$
$= \mathrm{argmax}_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle$

# Training

**Structure prediction:**

**Structure prediction:** *Structured perceptron loss*



Deep feature representation

Inference mechanism

$\boldsymbol{x} \in \mathcal{X}$

$\boldsymbol{x} \in \mathbb{R}^p$

Output potential

$\boldsymbol{\theta}(\boldsymbol{x}) \in$

$\boldsymbol{\Theta} \subset \mathbb{R}^D$

Gradient

$\boldsymbol{Y} \in \mathcal{Y} \subset \mathbb{R}^D$

$= \operatorname{argmax}_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle$

Gradient

$\boldsymbol{Y}_{\mathsf{true}} \in \mathcal{Y}$

$\max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle$

$- \langle \boldsymbol{\theta}(x), \boldsymbol{Y}_{\mathsf{true}} \rangle$

Loss

# Training

**Structure prediction:** *Structured perceptron loss*



Deep feature representation

Inference mechanism

$\boldsymbol{x} \in \mathcal{X}$

$\boldsymbol{x} \in \mathbb{R}^p$

Output potential $\boldsymbol{\theta}(\boldsymbol{x}) \in \boldsymbol{\Theta} \subset \mathbb{R}^D$

Gradient

$\boldsymbol{Y} \in \mathcal{Y} \subset \mathbb{R}^D = \operatorname{argmax}_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle$

$\boldsymbol{Y}_{\text{true}} \in \mathcal{Y} \cdots\cdots\triangleright$

$\max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle - \langle \boldsymbol{\theta}(x), \boldsymbol{Y}_{\text{true}} \rangle$

Gradient

Loss

Backprop through the max operator.

# Training

**Structure prediction:** *Structured perceptron loss*



Deep feature representation

Inference mechanism

$\boldsymbol{x} \in \mathcal{X}$

$\boldsymbol{x} \in \mathbb{R}^p$

Output potential

$\boldsymbol{\theta}(\boldsymbol{x}) \in \boldsymbol{\Theta} \subset \mathbb{R}^D$

Gradient

$\boldsymbol{Y} \in \mathcal{Y} \subset \mathbb{R}^D$
$= \mathrm{argmax}_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle$

Gradient

$\boldsymbol{Y}_{\text{true}} \in \mathcal{Y}$ ·······▸ $\max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle - \langle \boldsymbol{\theta}(x), \boldsymbol{Y}_{\text{true}} \rangle$ Loss

Backprop through the max operator. Not differentiable everywhere.

# Structured prediction as an inner layer

**Example:** Attention mechanisms,[1] where $c$ are the attention weights.



Deep feature representation

Inference mechanism

$x \in \mathcal{X}$

$x \in \mathbb{R}^p$

Output potential $\theta(x) \in \Theta \subset \mathbb{R}^D$

$c \in \mathcal{C} \subset \mathbb{R}^D$
$= \operatorname{argmax}_{c \in \mathcal{C}} \langle \theta(x), c \rangle$

Inference layer may be in the middle of the model

$Y \in \mathcal{Y} \subset \mathbb{R}^D$

---

[1] Dzmitry Bahdanau et al. (2015). "Neural Machine Translation by Jointly Learning to Align and Translate". In: *Proc. of ICLR*.

# Structured prediction as an inner layer

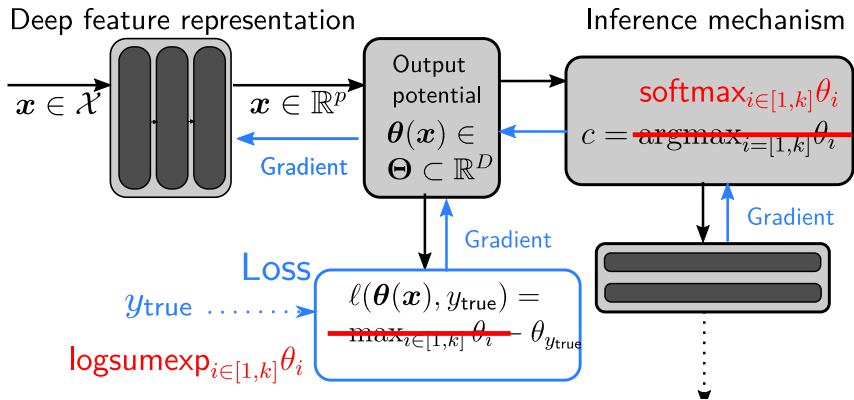**Example:** Attention mechanisms,[1] where $\boldsymbol{c}$ are the attention weights.



Deep feature representation

$\boldsymbol{x} \in \mathcal{X}$

$\boldsymbol{x} \in \mathbb{R}^p$

Gradient

Output potential

$\boldsymbol{\theta}(\boldsymbol{x}) \in$ $\boldsymbol{\Theta} \subset \mathbb{R}^D$

Gradient

Inference mechanism

$\boldsymbol{c} \in \mathcal{C} \subset \mathbb{R}^D$ $= \arg\max_{\boldsymbol{c} \in \mathcal{C}} \langle \boldsymbol{\theta}(x), \boldsymbol{c} \rangle$

Gradient

$\boldsymbol{Y}_{\text{true}} \in \mathcal{Y} \cdots\!\!\rightarrow \Delta(\boldsymbol{Y}, \boldsymbol{Y}_{\text{true}}) \cdots\!\!\rightarrow \boldsymbol{Y} \in \mathcal{Y} \subset \mathbb{R}^D$

We need to backpropagate through the argmax.

[1] Dzmitry Bahdanau et al. (2015). "Neural Machine Translation by Jointly Learning to Align and Translate". In: *Proc. of ICLR*.

# Structured prediction as an inner layer

**Example:** Attention mechanisms,[1] where $\boldsymbol{c}$ are the attention weights.



Deep feature representation

$\boldsymbol{x} \in \mathcal{X}$   $\boldsymbol{x} \in \mathbb{R}^p$

Output potential $\boldsymbol{\theta}(\boldsymbol{x}) \in \boldsymbol{\Theta} \subset \mathbb{R}^D$

Gradient

Inference mechanism

$\boldsymbol{c} \in \mathcal{C} \subset \mathbb{R}^D$
$= \operatorname{argmax}_{\boldsymbol{c} \in \mathcal{C}} \langle \boldsymbol{\theta}(x), \boldsymbol{c} \rangle$

Gradient

Gradient

$\boldsymbol{Y}_{\text{true}} \in \mathcal{Y} \cdots\cdots\blacktriangleright \Delta(\boldsymbol{Y}, \boldsymbol{Y}_{\text{true}}) \cdots\cdots\blacktriangleright \boldsymbol{Y} \in \mathcal{Y} \subset \mathbb{R}^D$

We need to backpropagate through the argmax. Zero gradient.

[1] Dzmitry Bahdanau et al. (2015). "Neural Machine Translation by Jointly Learning to Align and Translate". In: *Proc. of ICLR*.

# Gradient from regularization: from max to softmax

# Gradient from regularization: from max to softmax



Multinomial loss, softmax attention: **differentiable layers**

# Questions and goal

- From **max** to **softmax**: Where does this comes from and can we use different smoothing techniques ?

# Questions and goal

- From **max** to **softmax**: Where does this comes from and can we use different smoothing techniques ?
- How to smooth a wide class of **structured prediction** LP problems?

$$\max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle \qquad\qquad \boldsymbol{Y} \in \mathcal{Y} \subset \mathbb{R}^D = \operatorname*{argmax}_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle$$

# Questions and goal

- From **max** to **softmax**: Where does this comes from and can we use different smoothing techniques ?
- How to smooth a wide class of **structured prediction** LP problems?

$$\max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle \qquad \boldsymbol{Y} \in \mathcal{Y} \subset \mathbb{R}^D = \underset{\boldsymbol{Y} \in \mathcal{Y}}{\mathrm{argmax}} \langle \boldsymbol{\theta}(x), \boldsymbol{Y} \rangle$$

Focus on inference mechanisms that relies on **dynamic programming**

- Smooth **max** layers for new structured losses
- Differentiable **argmax** layers for test and inner inference mechanisms

# Contributions

**Generic framework for differentiable structured prediction:**

- Regularizing the max operators with strongly convex penalties.
- May output **sparse** continuous outputs

**Applications:**

- End-to-end audio to score alignment
- Named entity recognition with sparse predictions
- Block sparse attention mechanisms

## <u>Extends and ground in theory</u>[2,3,4,5]

[2] Yann LeCun et al. (2006). "A tutorial on energy-based learning". In: *Predicting structured data* 1.0.

[3] Guillaume Lample et al. (2016). "Neural Architectures for Named Entity Recognition". In: *Proc. of NAACL*, pp. 260–270.

[4] Yoon Kim et al. (2017). "Structured Attention Networks". In: *Proc. of ICLR*.

[5] Marco Cuturi and Mathieu Blondel (2017). "Soft-DTW: a Differentiable Loss Function for Time-Series". In: *Proc. of ICML*, pp. 894–903.

# Dynamic programming

Dynamic programming solve the structure prediction problem

$$\text{LP}(\boldsymbol{\theta}) \triangleq \max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}, \boldsymbol{Y} \rangle$$

by splitting the combinatorial set $\mathcal{Y} \subset \mathbb{R}^D$ into **sets of smaller dimensions**

- Compute $\text{LP}(\boldsymbol{\theta})$ in linear time $\mathcal{O}(D)$ vs exponential naive resolution

# Dynamic programming

Dynamic programming solve the structure prediction problem

$$\text{LP}(\boldsymbol{\theta}) \triangleq \max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}, \boldsymbol{Y} \rangle$$

by splitting the combinatorial set $\mathcal{Y} \subset \mathbb{R}^D$ into **sets of smaller dimensions**

- Compute $\text{LP}(\boldsymbol{\theta})$ in linear time $\mathcal{O}(D)$ vs exponential naive resolution

Also provide the **argmax** in $\mathcal{O}(D)$:

$$\underset{\boldsymbol{Y} \in \mathcal{Y}}{\text{argmax}} \langle \boldsymbol{\theta}, \boldsymbol{Y} \rangle$$

# Dynamic programming

Dynamic programming solve the structure prediction problem

$$\mathsf{LP}(\boldsymbol{\theta}) \triangleq \max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}, \boldsymbol{Y} \rangle$$

by splitting the combinatorial set $\mathcal{Y} \subset \mathbb{R}^D$ into **sets of smaller dimensions**

- Compute $\mathsf{LP}(\boldsymbol{\theta})$ in linear time $\mathcal{O}(D)$ vs exponential naive resolution

Also provide the **argmax** in $\mathcal{O}(D)$:

$$\underset{\boldsymbol{Y} \in \mathcal{Y}}{\mathrm{argmax}} \langle \boldsymbol{\theta}, \boldsymbol{Y} \rangle$$

**Examples:**

- Viterbi algorithm for infering tag sequences
- Dynamic time warping algorithm for infering alignment matrices

# Dynamic programming as best path in a DAG

## Directed acyclic graph

- $G = (\mathcal{N}, \mathcal{E})$, with 1 root and 1 leaf, nodes numbered in topo. order $[1, N]$
- Edge $(i, j)$ has weight $\theta_{i,j}$ — $j$ parent, $i$ child. $\boldsymbol{\theta} \in \mathbb{R}^{n \times n}$ incidence matrix
- Path $\boldsymbol{Y} \in \mathcal{Y} \subset \{0, 1\}^{N \times N}$: $y_{i,j} = 1$ iff $(i, j)$ is taken

**Single path value:** $\langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$

**Highest score among all paths**

$$\text{LP}(\boldsymbol{\theta}) = \max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$$

# Maximum value computation (finding the max)

- **Max value from** $1$ **to** $i$

$$v_i(\boldsymbol{\theta}) = \max_{j \in \mathcal{P}_i} \theta_{i,j} + v_j(\boldsymbol{\theta})$$

- One pass over the graph

$$(v_1 = 0, v_2, \dots, v_n \triangleq \mathrm{DP}(\boldsymbol{\theta}))$$

= **Bellman equation**



Value computation

# Maximum value computation (finding the max)

- **Max value from $1$ to $i$**

$$v_i(\boldsymbol{\theta}) = \max_{j \in \mathcal{P}_i} \theta_{i,j} + v_j(\boldsymbol{\theta})$$

- One pass over the graph

$$(v_1 = 0, v_2, \ldots, v_n \triangleq \mathrm{DP}(\boldsymbol{\theta}))$$

= **Bellman equation**



$\overrightarrow{\text{Value computation}}$

**The DP recursion solves the linear problem (Bellman, 1958)**

$$\mathrm{DP}(\boldsymbol{\theta}) = \mathrm{LP}(\boldsymbol{\theta}) = \max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$$

# Best path computation (finding the argmax)

What if we want to find the LP solution (a.k.a. perform inference ?)

# Best path computation (finding the argmax)

What if we want to find the LP solution (a.k.a. perform inference ?)

**The argmax is computable using backpropagation** = backtracking

$$\partial\mathrm{DP}(\boldsymbol{\theta}) = \partial_{\boldsymbol{\theta}}(\boldsymbol{\theta} \to \max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle)) = \mathrm{conv}(\operatorname*{argmax}_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle)$$

- When the argmax is unique: $\partial_{\boldsymbol{\theta}}\mathrm{DP}(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$

# Best path computation (finding the argmax)

What if we want to find the LP solution (a.k.a. perform inference ?)

**The argmax is computable using backpropagation** = backtracking

---

**Danskin theorem (Danskin, 1966)**

$$\partial \mathrm{DP}(\boldsymbol{\theta}) = \partial_{\boldsymbol{\theta}}(\boldsymbol{\theta} \to \max_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle)) = \mathrm{conv}(\underset{\boldsymbol{Y} \in \mathcal{Y}}{\mathrm{argmax}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle)$$

---

- When the argmax is unique: $\partial_{\boldsymbol{\theta}} \mathrm{DP}(\boldsymbol{\theta}) = \mathrm{argmax}_{\boldsymbol{Y} \in \mathcal{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$

---

**Dynamic programming layers**

- **Max** layer: $\boldsymbol{\theta} \to \mathrm{DP}(\boldsymbol{\theta}) = \max_{\boldsymbol{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$
- **Argmax** layer: $\boldsymbol{\theta} \to \partial_{\boldsymbol{\theta}} \mathrm{DP}(\boldsymbol{\theta}) \sim \mathrm{argmax}_{\boldsymbol{Y}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$

# Example: Linear conditional random field

$(x_1, \ldots, x_T)$ observation, $(y_1, \ldots, y_T) \in [S]^T$ states. $\mathbf{Y} \in \mathcal{Y} \in \{0,1\}^{S \times S \times T}$
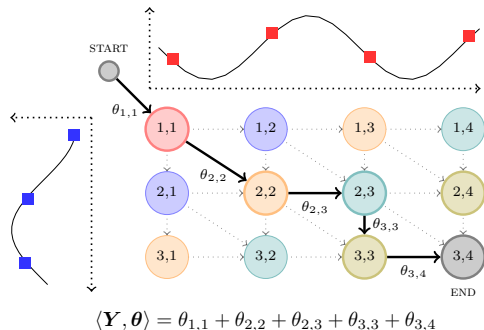
$$\mathbf{y} = \operatorname*{argmax}_{\mathbf{y} \in \mathcal{Y}} \sum_{t=1}^{T} \theta_t(y_t, y_{t-1}, x_t) = \operatorname*{argmax}_{\mathbf{Y} \in \mathcal{Y}} \langle \boldsymbol{\theta}, \mathbf{Y} \rangle$$



$$\langle \mathbf{Y}, \boldsymbol{\theta} \rangle = \theta_{1,3,1} + \theta_{2,1,3} + \theta_{3,2,1}$$

$\mathbf{Y}$ computed with dynamic programming = **Viterbi algorithm.**

# Example: Dynamic time warping



$$\langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle = \theta_{1,1} + \theta_{2,2} + \theta_{2,3} + \theta_{3,3} + \theta_{3,4}$$

**Elastic matching**

- Two time-series $\boldsymbol{A}$, $\boldsymbol{B}$
- Distance matrix: $\boldsymbol{\theta}_{i,j} = \|a_i - b_i\|_2^2$

**Alignment matrices**

- $(1,1) \to (N_A, N_B)$
- $\downarrow, \to, \searrow$ moves

Best alignment: $\boldsymbol{Y}(\boldsymbol{A}, \boldsymbol{B}) = \underset{\boldsymbol{Y} \in \mathcal{Y}}{\mathrm{argmax}} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$

DTW distance: $d(\boldsymbol{A}, \boldsymbol{B}) = \underset{\boldsymbol{Y} \in \mathcal{Y}}{\max} \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$

Computable by dynamic programming

- $\mathcal{Y}$ set of alignment matrices
- $\boldsymbol{\theta}$ distance matrix

# Regularizing dynamic programming

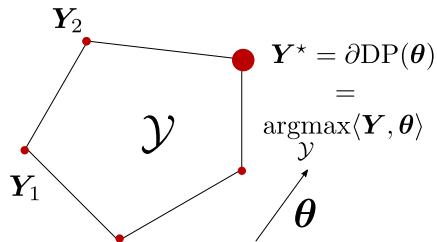## Obstacles to end-to-end training

- Max layer $\theta \rightarrow DP(\theta)$ is not differentiable everywhere
- Argmax layer $\theta \rightarrow \partial DP(\theta)$ is piecewise constant / not defined

# Regularizing dynamic programming

**Obstacles to end-to-end training**

- Max layer $\boldsymbol{\theta} \to \mathrm{DP}(\boldsymbol{\theta})$ is not differentiable everywhere
- Argmax layer $\boldsymbol{\theta} \to \partial\mathrm{DP}(\boldsymbol{\theta})$ is piecewise constant / not defined

**Culprit is the Bellman recursion**

$$\boldsymbol{x} \in \mathbb{R}^d \to \max(\boldsymbol{x}) \in \mathbb{R}$$

- Not differentiable everywhere
- Piecewise linear (null Hessian)



Hard geometry

# Regularizing dynamic programming

**Obstacles to end-to-end training**
- Max layer $\boldsymbol{\theta} \to \mathrm{DP}(\boldsymbol{\theta})$ is not differentiable everywhere
- Argmax layer $\boldsymbol{\theta} \to \partial\mathrm{DP}(\boldsymbol{\theta})$ is piecewise constant / not defined

**Culprit is the Bellman recursion**

$$\boldsymbol{x} \in \mathbb{R}^d \to \max(\boldsymbol{x}) \in \mathbb{R}$$

- Not differentiable everywhere
- Piecewise linear (null Hessian)



$\boldsymbol{Y}_2$

$\boldsymbol{Y}^\star = \partial\mathrm{DP}(\boldsymbol{\theta})$
$=$
$\underset{\mathcal{Y}}{\mathrm{argmax}}\langle \boldsymbol{Y}, \boldsymbol{\theta}\rangle$

$\mathcal{Y}$

$\boldsymbol{Y}_1$

$\boldsymbol{\theta}$

Hard geometry

Solution: smooth the maximum operator

# Max smoothing

$\Omega : \mathbb{R} \to \mathbb{R}$ strongly-convex function. $\boldsymbol{x} \in \mathbb{R}^d$. $\Delta^d$: $d$-dim simplex.

### Smoothed max operator (Moreau, 1965; Nesterov, 2005)

$$\max{}_{\Omega}(\boldsymbol{x}) = \max{}_{\boldsymbol{y} \in \Delta^d} \langle \boldsymbol{x}, \boldsymbol{y} \rangle - \sum_{i=1}^{d} \Omega(y_i)$$

# Max smoothing

$\Omega : \mathbb{R} \to \mathbb{R}$ strongly-convex function. $\boldsymbol{x} \in \mathbb{R}^d$. $\Delta^d$: $d$-dim simplex.

## Smoothed max operator (Moreau, 1965; Nesterov, 2005)

$$\max\nolimits_\Omega(\boldsymbol{x}) = \max\nolimits_{\boldsymbol{y} \in \Delta^d} \langle \boldsymbol{x}, \boldsymbol{y} \rangle - \sum\nolimits_{i=1}^d \Omega(y_i)$$

**Properties:**

- Consistent smoothing: $\max_0(\boldsymbol{x}) = \max(\boldsymbol{x})$
- Twice differentiable almost everywhere with non-zero Hessian

# Examples of regularization

**Shannon entropy:** $\Omega(x) = x \log(x) \longrightarrow$ *Softmax* operator

$$\max{}_\Omega(\boldsymbol{x}) = \log(Z), \text{ where } Z = \sum_j \exp(x_j)$$

$$\nabla \max{}_\Omega(\boldsymbol{x}) = (\exp(x_i)/Z)_{i \in \mathbb{R}^d}$$

# Examples of regularization

**Shannon entropy:** $\Omega(x) = x \log(x) \longrightarrow$ *Softmax* operator

$$\max\nolimits_{\Omega}(\boldsymbol{x}) = \log(Z), \text{ where } Z = \sum\nolimits_j \exp(x_j)$$

$$\nabla \max\nolimits_{\Omega}(\boldsymbol{x}) = (\exp(x_i)/Z)_{i \in \mathbb{R}^d}$$

$\ell_2^2$ **norm:** $\Omega(x) = x^2 \longrightarrow$ *Sparsemax* (Martins and Astudillo, 2016)

$\nabla \max\nolimits_{\Omega}(\boldsymbol{x}) = \operatorname{argmin}_{p \in \triangle^d} \|\boldsymbol{x} - \boldsymbol{p}\|_2^2$   Sparse: eucl. projection on simplex

# Dynamic programming regularization

**1. Smooth max:** $\max_\Omega(\boldsymbol{x}) = \max_{\boldsymbol{y} \in \Delta^d} \langle \boldsymbol{x}, \boldsymbol{y} \rangle - \sum_{i=1}^{d} \Omega(y_i)$

**2. Bellman recursion:** $v_i = \max_{j \in \mathcal{P}_i} \theta_{i,j} + v_j, \quad \mathrm{DP}(\boldsymbol{\Theta}) \triangleq v_N$

# Dynamic programming regularization

1. **Smooth max:** $\max_{\Omega}(x) = \max_{y \in \Delta^d} \langle x, y \rangle - \sum_{i=1}^{d} \Omega(y_i)$

2. **Bellman recursion:** $v_i = \max_{j \in \mathcal{P}_i} \theta_{i,j} + v_j, \quad \mathrm{DP}(\Theta) \triangleq v_N$

**Bottom-up construction**

For all $i \in [N]$:

$$v_i(\boldsymbol{\theta}) = \max_{\Omega}(\theta_{i,j} + v_j)_{j \in \mathcal{P}_i}$$

$$\mathrm{DP}_{\Omega}(\boldsymbol{\theta}) \triangleq v_N(\boldsymbol{\theta})$$



$v_1$ $v_2$ $v_5$ $v_7$

$\theta_{21}$ $\theta_{52}$ $\theta_{75}$

$\theta_{31}$ $\theta_{53}$ $v_6$ $\theta_{76}$

$\theta_{41}$ $v_3$

$\theta_{64}$

$v_4$

$\overrightarrow{\text{Value computation}}$

# Regularized best-path: $\nabla \mathrm{DP}_\Omega(\boldsymbol{\theta})$

**From max to smoothed max:**

$$\boldsymbol{Y}(\boldsymbol{\theta}) = \partial \mathrm{DP}(\boldsymbol{\theta}) \implies \boldsymbol{Y}_\Omega(\boldsymbol{\theta}) \triangleq \nabla \mathrm{DP}_\Omega(\boldsymbol{\theta})$$

# Regularized best-path: $\nabla DP_{\Omega}(\boldsymbol{\theta})$

**From max to smoothed max:**

$$\boldsymbol{Y}(\boldsymbol{\theta}) = \partial DP(\boldsymbol{\theta}) \implies \boldsymbol{Y}_{\Omega}(\boldsymbol{\theta}) \triangleq \nabla DP_{\Omega}(\boldsymbol{\theta})$$

Computed with backpropagation

**Requirements:** Gradients of Bellman equations

$$\boldsymbol{q}_i = \nabla \max{}_{\Omega}(\theta_{i,j} + v_j)_{j \in \mathcal{P}_i}$$



Gradient computation

# Differentiable dynamic programming layers

# Differentiable DP properties

**Usable for loss design**: $\theta \to \mathrm{DP}_\Omega(\theta)$ is convex, bounds $\mathrm{DP}(\theta)$

# Differentiable DP properties

**Usable for loss design**: $\theta \rightarrow DP_\Omega(\theta)$ is convex, bounds $DP(\theta)$

**Is local regularization equivalent to global regularization ?**

$$LP_\Omega(\theta) \triangleq \max_\Omega \left(\langle Y, \theta \rangle\right)_{Y \in \mathcal{Y}} = \max_{p \in \triangle^D} \left\langle p, \left(\langle Y, \theta \rangle\right)_{y \in \mathcal{Y}} \right\rangle - \Omega(p)$$

**Theorem**:
- $DP_\Omega(\theta) = LP_\Omega(\theta)$ **if and only if** $\Omega(p) = -\gamma \sum_i p_i \log p_i$
- $DP_\Omega(\theta)$ is the CRF (Lafferty et al., 2001) log-partition

$$DP_\Omega(\theta) = \log\left(\sum_{Y \in \mathcal{Y}} \exp(\langle \theta, Y \rangle)\right)$$

New motivation for Shannon reg. But $\ell_2^2$ has other interesting properties.
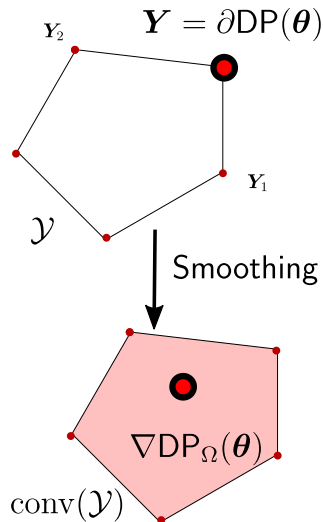
# Relaxed gradient properties



**Probabilistic interpretation**

Backprop defines a distribution $\mathcal{D}_\Omega$ on the set of paths $\mathcal{Y}$

$$\nabla \mathrm{DP}_\Omega(\boldsymbol{\theta}) = \mathbb{E}_{\mathcal{D}_\Omega}[\boldsymbol{Y}] \in \mathrm{conv}(\mathcal{Y})$$

$\Rightarrow$ Probability of path $\boldsymbol{Y}$: $p_{\boldsymbol{\theta},\Omega}(\boldsymbol{Y})$

# Relaxed gradient properties



$Y = \partial \mathrm{DP}(\boldsymbol{\theta})$

**Probabilistic interpretation**

Backprop defines a distribution $\mathcal{D}_\Omega$ on the set of paths $\mathcal{Y}$

$$\nabla \mathrm{DP}_\Omega(\boldsymbol{\theta}) = \mathbb{E}_{\mathcal{D}_\Omega}[\boldsymbol{Y}] \in \mathrm{conv}(\mathcal{Y})$$

$\Rightarrow$ Probability of path $\boldsymbol{Y}$: $p_{\boldsymbol{\theta},\Omega}(\boldsymbol{Y})$

- **Shannon:** Gibbs distribution: $p_{\boldsymbol{\theta},\Omega}(\boldsymbol{Y}) \propto \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$
- $\ell_2^2$: $\mathcal{D}_\Omega$ has a small support $\rightarrow \nabla \mathrm{DP}_\Omega(\boldsymbol{\theta})$ is **sparse**

# Backpropagating through $\nabla DP_\Omega(\Theta)$

**Regularized best-path layer:** $\theta \in \mathbb{R}^{N \times N} \to \nabla DP_\Omega(\theta)$

**Jacobian ?** $\nabla \nabla DP_\Omega(\Theta) = \nabla^2 DP_\Omega(\Theta) = $ Hessian

### Hessian vector-product

$$\nabla(\nabla DP_\Omega(\Theta))Z = \nabla^2 DP_\Omega(\Theta)Z, \qquad Z \in \mathbb{R}^{n \times n} \quad \text{direction}$$

Computable in $\mathcal{O}(|\mathcal{E}|)$: reverse-on-forward differentiation

# Differentiable dynamic programming layers

**Highest-score layer, forward-pass**

$$\boldsymbol{\theta} \in \mathbb{R}^{N \times N} \to \mathrm{DP}_{\Omega}(\boldsymbol{\theta})$$

**Highest score layer, backward pass**
**Best path layer, forward-pass**

$$\boldsymbol{\theta} \in \mathbb{R}^{N \times N} \to \nabla \mathrm{DP}_{\Omega}(\boldsymbol{\theta})$$

# Differentiable dynamic programming layers

**Highest-score layer, forward-pass**

$$\boldsymbol{\theta} \in \mathbb{R}^{N \times N} \to \mathrm{DP}_\Omega(\boldsymbol{\theta})$$

**Highest score layer, backward pass**
**Best path layer, forward-pass**

$$\boldsymbol{\theta} \in \mathbb{R}^{N \times N} \to \nabla\mathrm{DP}_\Omega(\boldsymbol{\theta})$$

**Best-path layer: backward pass**

$$\boldsymbol{\theta}, \boldsymbol{Z} \in \mathbb{R}^{N \times N} \times \mathbb{R}^{N \times N} \to \nabla^2\mathrm{DP}_\Omega(\boldsymbol{\theta})\boldsymbol{Z}$$

# Differentiable dynamic programming layers

**Highest-score layer, forward-pass**

$$\boldsymbol{\theta} \in \mathbb{R}^{N \times N} \to \mathrm{DP}_\Omega(\boldsymbol{\theta})$$

**Highest score layer, backward pass**
**Best path layer, forward-pass**

$$\boldsymbol{\theta} \in \mathbb{R}^{N \times N} \to \nabla\mathrm{DP}_\Omega(\boldsymbol{\theta})$$

**Best-path layer: backward pass**

$$\boldsymbol{\theta}, \boldsymbol{Z} \in \mathbb{R}^{N \times N} \times \mathbb{R}^{N \times N} \to \nabla^2\mathrm{DP}_\Omega(\boldsymbol{\theta})\boldsymbol{Z}$$

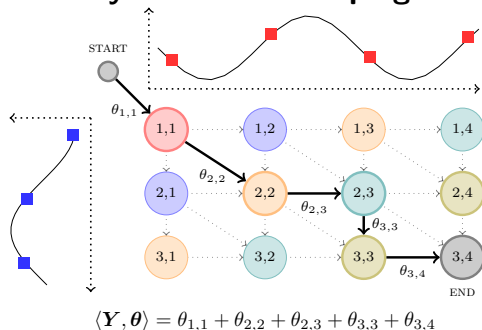- Sparse/dense output with $\ell_2$/entropy regularization
- Total computational cost: $\mathcal{O}(|\mathcal{E}|)$

# Applications



**Viterbi**

the    boat    sank

$$\langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle = \theta_{1,3,1} + \theta_{2,1,3} + \theta_{3,2,1}$$

**Dynamic time warping**

$$\langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle = \theta_{1,1} + \theta_{2,2} + \theta_{2,3} + \theta_{3,3} + \theta_{3,4}$$
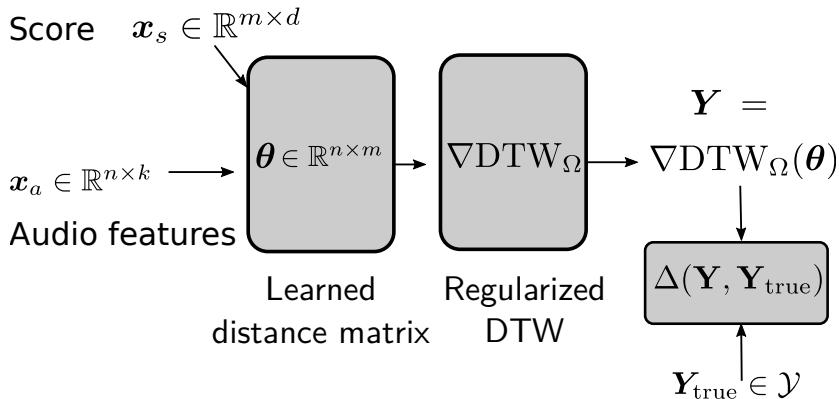
$$\nabla \mathsf{Vit}_\Omega : \mathbb{R}^{T \times S \times S} \rightarrow \mathbb{R}^{T \times S \times S}$$

$$\nabla \mathsf{DTW}_\Omega : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$$

# Audio-to-score alignment

- **Input data:** audio sequence $\boldsymbol{x}_a \in \mathbb{R}^{n \times k}$, one-hot key sequence $\boldsymbol{x}_s \in \mathbb{R}^{m \times d}$
- **Labels:** Alignment $\boldsymbol{Y}_{\text{true}} \in \mathcal{Y} \subset \mathbb{R}^{n \times m}$



Score $\quad \boldsymbol{x}_s \in \mathbb{R}^{m \times d}$

$\boldsymbol{x}_a \in \mathbb{R}^{n \times k}$

Audio features

$\boldsymbol{\theta} \in \mathbb{R}^{n \times m}$

Learned
distance matrix

$\nabla \text{DTW}_\Omega$

Regularized
DTW

$\boldsymbol{Y} = \nabla \text{DTW}_\Omega(\boldsymbol{\theta})$

$\Delta(\mathbf{Y}, \mathbf{Y}_{\text{true}})$

$\boldsymbol{Y}_{\text{true}} \in \mathcal{Y}$

# Metric learning experiment

**Learn the distance matrix:**

- Baseline: multinomial classification, audio-frame to score key
- Our model: end-to-end training of a linear model with final soft-DTW layer
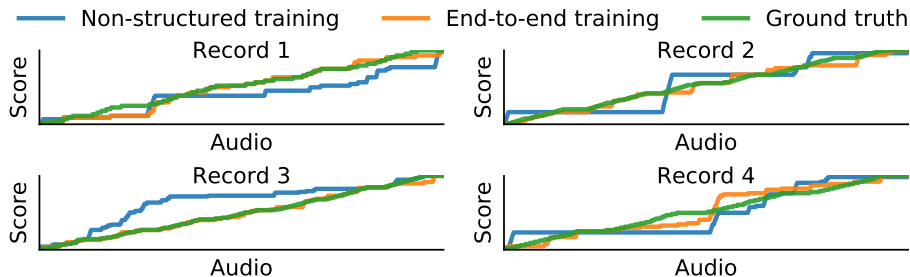
**Data:** Supervised dataset: 10 annotated Bach quatuors (**Bach10**)

**Validation:**

- Leave-one-out prediction
- At test time: Hard DTW on the learned distance matrix
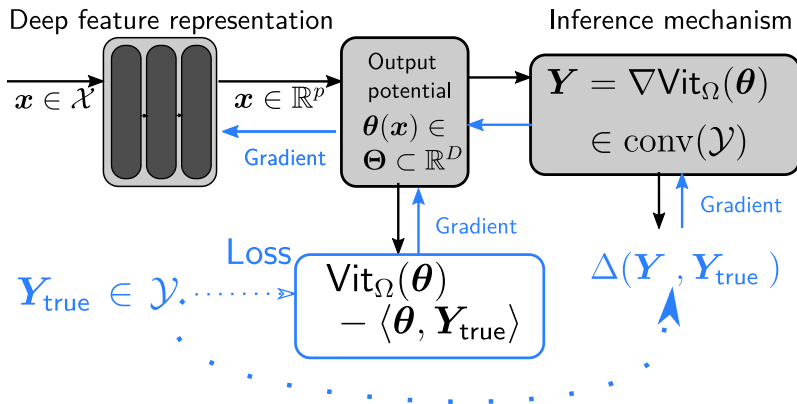- RMSE between predicted onsets

# Results

| RMSE | Test set | Train set |
|------|----------|-----------|
| End-to-end training | $\mathbf{1.26 \pm 0.64}$ | $\mathbf{0.17 \pm 0.01}$ |
| Non-structure training | $3.70 \pm 2.85$ | $1.80 \pm 0.14$ |
| Random | $14.64 \pm 2.63$ | $14.64 \pm 0.29$ |

# Named entity recognition

- **Input data:** Sentences $\boldsymbol{x}$ of length $T$
- **Labels $\boldsymbol{Y}$:** {Begin/Inside/Outside}{Person/Org./Loc./Misc.}



Deep feature representation

Inference mechanism

$\boldsymbol{x} \in \mathcal{X}$

$\boldsymbol{x} \in \mathbb{R}^p$

Output potential
$\boldsymbol{\theta}(\boldsymbol{x}) \in$
$\Theta \subset \mathbb{R}^D$

$\boldsymbol{Y} = \nabla\mathsf{Vit}_{\Omega}(\boldsymbol{\theta})$
$\in \mathrm{conv}(\mathcal{Y})$

Gradient

Gradient

Loss

Gradient

$\boldsymbol{Y}_{\mathsf{true}} \in \mathcal{Y}$

$\mathsf{Vit}_{\Omega}(\boldsymbol{\theta})$
$- \langle \boldsymbol{\theta}, \boldsymbol{Y}_{\mathsf{true}} \rangle$

$\Delta(\boldsymbol{Y}, \boldsymbol{Y}_{\mathsf{true}})$

# K-best set predictions in named entity recognition

$\Omega = \ell_2^2$: $\nabla \mathrm{DP}_\Omega(\boldsymbol{\theta}(x))$ k-best subset of $\boldsymbol{Y}$ such that $p_{\boldsymbol{\theta}, \Omega}(\boldsymbol{Y}) \neq 0$

# Quantitative comparison of losses

- **Potential-convex loss:** $\ell_\Omega(\boldsymbol{\theta}, \boldsymbol{Y}) = \mathrm{DP}_\Omega(\boldsymbol{\theta}) - \langle \boldsymbol{Y}, \boldsymbol{\theta} \rangle$
- **Cost-sensitive loss:** $\ell_\Omega(\boldsymbol{\theta}, \boldsymbol{Y}) = \Delta(\nabla \mathrm{DP}_\Omega(\boldsymbol{\theta}), \boldsymbol{Y})\rangle.$

| $\Omega$ | Loss | English | Spanish | German | Dutch |
|---|---|---|---|---|---|
| Negentropy | Convex loss | 90.80 | **86.68** | 77.35 | **87.56** |
| | Cost-sensitive | 90.47 | 86.20 | **77.56** | 87.37 |
| $\ell_2^2$ | Convex loss | **90.86** | 85.51 | 76.01 | 86.58 |
| | Cost-sensitive | 89.49 | 84.07 | 76.91 | 85.90 |
| Lample et al., 2016[6] | | *90.96* | *85.75* | *78.76* | *81.74* |

- $\ell_2^2$ reg. achieves comparable accuracy with more interpretable predictions
- Training directly from potential-derived losses is slightly better

---

[6] Guillaume Lample et al. (2016). "Neural Architectures for Named Entity Recognition". In: *Proc. of NAACL*, pp. 260–270.

# Structured attention — Neural machine transation

- Compute an attention vector $c$: 2 state linear-chain CRF
- $c = \mathbb{E}[z]$, $z_i = 1$ if attention, $z_i = 0$ if no-attention
- Use $\text{Vit}_\Omega$, with sparse marginal computation $\Omega = \ell_2^2$.
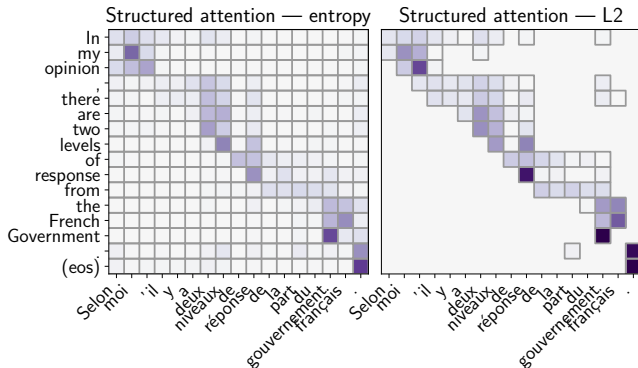


Structured attention — entropy    Structured attention — L2

Similar BLEU scores WMT14 1M

| Attention model | fr→en | en→fr |
|---|---|---|
| Softmax | **27.96** | **28.08** |
| CRF + entropy | **27.96** | 27.98 |
| CRF + $\ell_2^2$ reg. | 27.21 | 27.28 |

# Structured attention — Neural machine transation

- Compute an attention vector $c$: 2 state linear-chain CRF
- $c = \mathbb{E}[z]$, $z_i = 1$ if attention, $z_i = 0$ if no-attention
- Use $\mathrm{Vit}_\Omega$, with sparse marginal computation $\Omega = \ell_2^2$.



Structured attention — entropy    Structured attention — L2

Similar BLEU scores WMT14 1M

| Attention model | fr→en | en→fr |
|---|---|---|
| Softmax | **27.96** | **28.08** |
| CRF + entropy | **27.96** | 27.98 |
| CRF + $\ell_2^2$ reg. | 27.21 | 27.28 |

## Block sparse attention

# Conclusion

**General framework to use DP algorithms in arbitrary networks**

- Efficient and stable algorithms
- Flexibility of regularization (sparse output)

**Experiments:** $\ell_2$/entropy have similar performance

- $\ell_2^2$: More interpretable outputs / k-best sets with sparsity

*PyTorch* package *didyprog* available (fast custom Viterbi and DTW layer)

# Conclusion

**General framework to use DP algorithms in arbitrary networks**

- Efficient and stable algorithms
- Flexibility of regularization (sparse output)

**Experiments:** $\ell_2$/entropy have similar performance

- $\ell_2^2$: More interpretable outputs / k-best sets with sparsity

*PyTorch* package *didyprog* available (fast custom Viterbi and DTW layer)



Arthur Mensch and Mathieu Blondel (2018). "Differentiable Dynamic Programming for Structured Prediction and Attention". In: *Proceedings of the International Conference on Machine Learning*

# Conclusion

**General framework to use DP algorithms in arbitrary networks**

- Efficient and stable algorithms
- Flexibility of regularization (sparse output)

**Experiments:** $\ell_2$/entropy have similar performance

- $\ell_2^2$: More interpretable outputs / k-best sets with sparsity

*PyTorch* package *didyprog* available (fast custom Viterbi and DTW layer)



Arthur Mensch and Mathieu Blondel (2018). "Differentiable Dynamic Programming for Structured Prediction and Attention". In: *Proceedings of the International Conference on Machine Learning*

**Related work at Google:** Framwork formalizes differentiable beam search (Goyal et al., 2017), similar effor in reinforcement learning (Haarnoja et al., 2018)

# Bibliography I

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2015). "Neural Machine Translation by Jointly Learning to Align and Translate". In: *Proc. of ICLR*.

Bellman, Richard (1958). "On a routing problem". In: *Quarterly of applied mathematics* 16.1, pp. 87–90.

Cuturi, Marco and Mathieu Blondel (2017). "Soft-DTW: a Differentiable Loss Function for Time-Series". In: *Proc. of ICML*, pp. 894–903.

Danskin, John M (1966). "The theory of max-min, with applications". In: *SIAM Journal on Applied Mathematics* 14.4, pp. 641–664.

Goyal, Kartik et al. (July 2017). "A Continuous Relaxation of Beam Search for End-to-End Training of Neural Sequence Models". In: *arXiv:1708.00111 [cs]*. arXiv: 1708.00111 [cs].

# Bibliography II

Haarnoja, Tuomas et al. (2018). "Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor". en. In: *Proceedings of the International Conference on Machine Learning*, p. 10.

Kim, Yoon et al. (2017). "Structured Attention Networks". In: *Proc. of ICLR*.

Lafferty, John, Andrew McCallum, and Fernando CN Pereira (2001). "Conditional random fields: Probabilistic models for segmenting and labeling sequence data". In: *Proc. of ICML*, pp. 282–289.

Lample, Guillaume et al. (2016). "Neural Architectures for Named Entity Recognition". In: *Proc. of NAACL*, pp. 260–270.

LeCun, Yann et al. (2006). "A tutorial on energy-based learning". In: *Predicting structured data* 1.0.

# Bibliography III

Martins, André F.T. and Ramón Fernandez Astudillo (2016). "From softmax to sparsemax: A sparse model of attention and multi-label classification". In: *Proc. of ICML*, pp. 1614–1623.

Mensch, Arthur and Mathieu Blondel (2018). "Differentiable Dynamic Programming for Structured Prediction and Attention". In: *Proceedings of the International Conference on Machine Learning.*

Moreau, Jean-Jacques (1965). "Proximité et dualité dans un espace hilbertien". In: *Bullet de la Société Mathémathique de France* 93.2, pp. 273–299.

Nesterov, Yurii (2005). "Smooth minimization of non-smooth functions". In: *Mathematical Programming* 103.1, pp. 127–152.