# Differentiable Adaptive Sparsity For Neural Networks
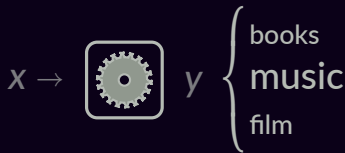
**Vlad Niculae**

Instituto de Telecomunicações

# Choosing Between *K* Options

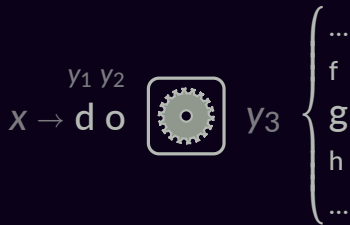A building block in many ML tasks!

multi-class classification

$x \rightarrow$ ⚙ $y$ 
$\begin{cases} \text{books} \\ \textbf{music} \\ \text{film} \end{cases}$

# **Choosing Between _K_ Options**
### A building block in many ML tasks!

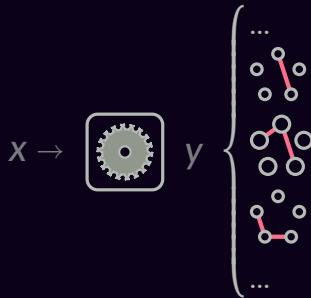multi-class classification

sequence generation

# Choosing Between *K* Options
## A building block in many ML tasks!

multi-class classification

sequence generation

structured output prediction

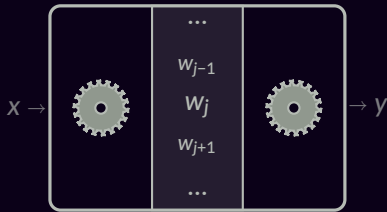# Choosing Between *K* Options

A building block in many ML tasks!

multi-class classification
sequence generation
structured output prediction $\left.\vphantom{\begin{array}{c}1\\2\\3\end{array}}\right\}$ output

# Choosing Between *K* Options

A building block in many ML tasks!

multi-class classification
sequence generation
structured output prediction } output

neural attention

# Choosing Between *K* Options

A building block in many ML tasks!

multi-class classification

sequence generation

structured output prediction

} output

neural attention

structured hidden layers

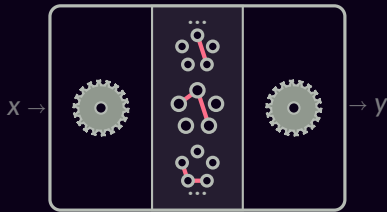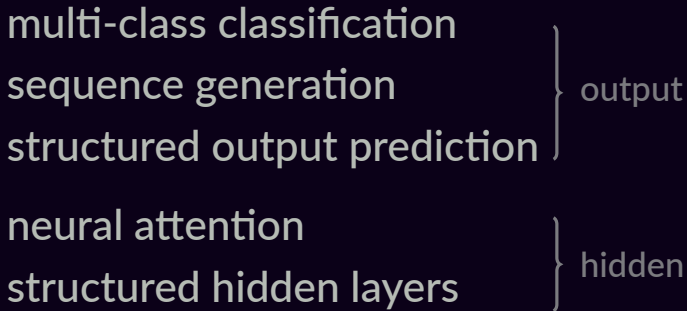# Choosing Between *K* Options
A building block in many ML tasks!

multi-class classification
sequence generation
structured output prediction

} output

neural attention
structured hidden layers

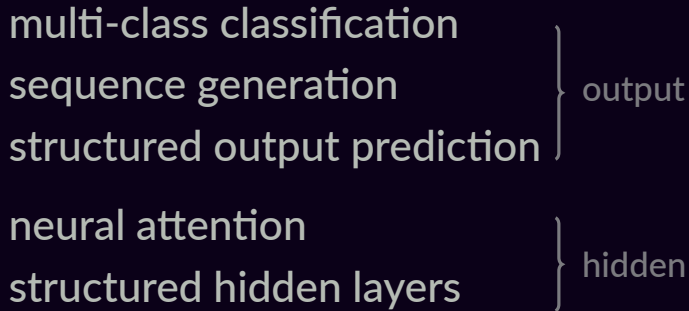} hidden

# Choosing Between *K* Options

A building block in many ML tasks!

multi-class classification
sequence generation
structured output prediction
} output

neural attention
structured hidden layers
} hidden

Deterministic sparse & structured mappings and losses via a general, constructive framework.

# Outline

# Perceptron & Argmax

$\theta$

$x$    $f_w(x)$

$c_1$
$c_2$
$\ldots$

$c_k$

# Perceptron & Argmax

$$p := \text{argmax}(\boldsymbol{\theta})$$

$\boldsymbol{\theta}$ $\qquad$ $\boldsymbol{p}$

$c_1$

$x$

$c_2$

$f_w(x)$

$\ldots$

$c_k$

- very sparse predictions

# Perceptron & Argmax



$$p := \text{argmax}(\boldsymbol{\theta})$$
$$L(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) = \langle \boldsymbol{\theta}, \mathbf{p} \rangle - \langle \boldsymbol{\theta}, \mathbf{y}^{\text{true}} \rangle$$
$$\partial_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) \ni \mathbf{p} - \mathbf{y}^{\text{true}}$$

- very sparse predictions
- famous update rule

# Perceptron & Argmax

$$p := \operatorname{argmax}(\boldsymbol{\theta})$$
$$L(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) = \langle \boldsymbol{\theta}, p \rangle - \langle \boldsymbol{\theta}, \mathbf{y}^{\text{true}} \rangle$$
$$\partial_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) \ni p - \mathbf{y}^{\text{true}}$$



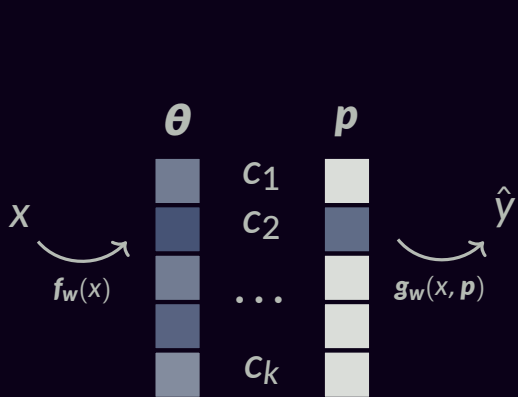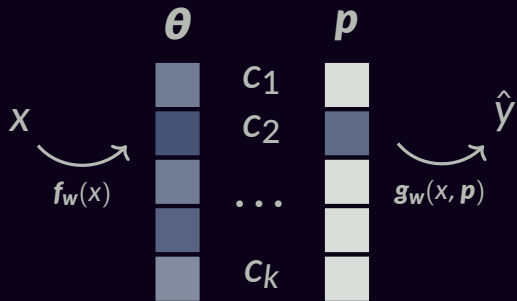- very sparse predictions
- famous update rule

# Perceptron & Argmax

$$p := \mathrm{argmax}(\boldsymbol{\theta})$$
$$L(\boldsymbol{\theta}; \mathbf{y}^{\mathrm{true}}) = \langle \boldsymbol{\theta}, \boldsymbol{p} \rangle - \langle \boldsymbol{\theta}, \mathbf{y}^{\mathrm{true}} \rangle$$
$$\partial_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{y}^{\mathrm{true}}) \ni \boldsymbol{p} - \mathbf{y}^{\mathrm{true}}$$



$\boldsymbol{\theta}$ $\boldsymbol{p}$

$c_1$

$x$ $c_2$ $\hat{y}$

$f_w(x)$ ... $g_w(x, \boldsymbol{p})$

$c_k$

- very sparse predictions
- famous update rule
- can't use as hidden layer: $\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \mathbf{0}$ *a.e.*

# Perceptron & Argmax

$$p := \text{argmax}(\boldsymbol{\theta})$$
$$L(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) = \langle \boldsymbol{\theta}, \boldsymbol{p} \rangle - \langle \boldsymbol{\theta}, \mathbf{y}^{\text{true}} \rangle$$
$$\partial_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) \ni \boldsymbol{p} - \mathbf{y}^{\text{true}}$$

$\boldsymbol{\theta}$  $\boldsymbol{p}$

$c_1$

$x$  $c_2$  $\hat{y}$

$f_w(x)$  $\cdots$  $g_w(x, \boldsymbol{p})$

$c_k$

- very sparse predictions
- famous update rule
- can't use as hidden layer: $\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \mathbf{0}$ *a.e.*

# Logistic Regression & Softmax

$$\boldsymbol{p} := \mathrm{softmax}(\boldsymbol{\theta})$$



$\boldsymbol{\theta}$     $\boldsymbol{p}$

$c_1$

$c_2$

$x$   $f_{\boldsymbol{w}}(x)$

$\cdots$

$c_k$

- dense predictive distribution (Gibbs)

# Logistic Regression & Softmax



$$\boldsymbol{p} := \mathrm{softmax}(\boldsymbol{\theta})$$

$$L(\boldsymbol{\theta}; \mathbf{y}^{\mathrm{true}}) = \log \sum_j \exp \theta_j - \langle \boldsymbol{\theta}, \mathbf{y}^{\mathrm{true}} \rangle$$

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{y}^{\mathrm{true}}) = \boldsymbol{p} - \mathbf{y}^{\mathrm{true}}$$

- dense predictive distribution (Gibbs)
- loss gradient:
  *expected* − *observed* statistics

# Logistic Regression & Softmax

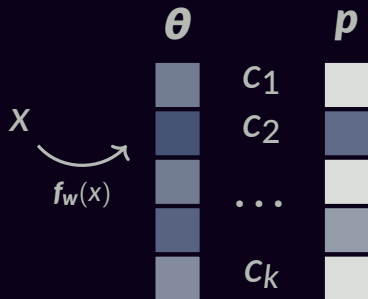$$\boldsymbol{p} := \text{softmax}(\boldsymbol{\theta})$$

$$L(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) = \log \sum_j \exp \theta_j - \langle \boldsymbol{\theta}, \mathbf{y}^{\text{true}} \rangle$$

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) = \boldsymbol{p} - \mathbf{y}^{\text{true}}$$

$\boldsymbol{\theta}$    $\boldsymbol{p}$

$c_1$

$c_2$

$x$    $\hat{y}$

$f_w(x)$    $g_w(x, \boldsymbol{p})$

$\ldots$

$c_k$

- **dense** predictive distribution (Gibbs)
- loss gradient:
  *expected − observed* statistics
- soft hidden layers: $\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}} = \text{diag}(\boldsymbol{p}) - \boldsymbol{p}\boldsymbol{p}^\top$
  (neural attention)

# Sparsemax

$$p := \mathrm{sparsemax}(\boldsymbol{\theta}) = \mathrm{proj}_{\triangle}(\boldsymbol{\theta})$$

$\boldsymbol{\theta}$  $p$
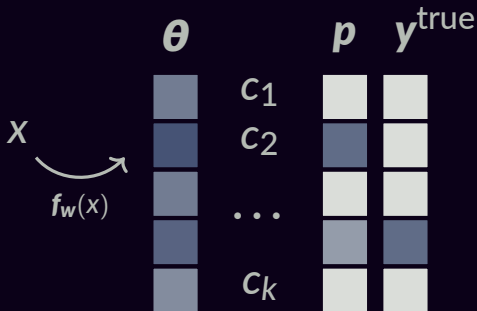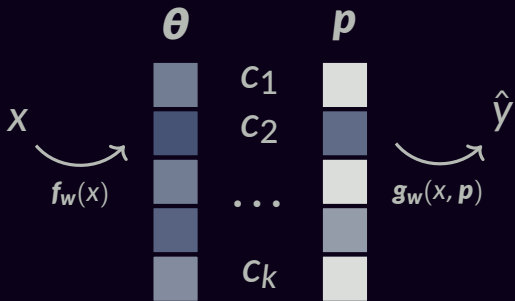
$c_1$

$c_2$

$x$

$f_w(x)$

$\dots$

$c_k$

- **sparse** predictive distribution

# Sparsemax

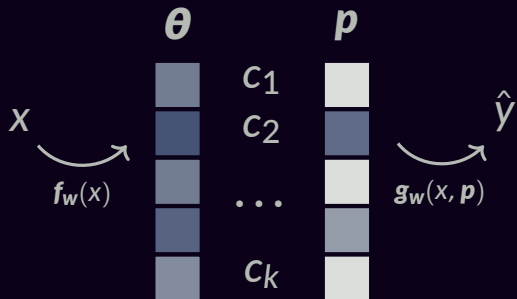$$p := \text{sparsemax}(\boldsymbol{\theta}) = \text{proj}_{\triangle}(\boldsymbol{\theta})$$

$$L(\boldsymbol{\theta}, \mathbf{y}^{\text{true}}) = \; ?$$

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \mathbf{y}^{\text{true}}) = \boldsymbol{p} - \mathbf{y}^{\text{true}}$$

$\boldsymbol{\theta}$  $\boldsymbol{p}$  $\mathbf{y}^{\text{true}}$

$c_1$

$x$

$f_w(x)$

$c_2$

$\cdots$

$c_k$

- **sparse** predictive distribution
- reverse-engineer loss from gradient
  *expected − observed* statistics

# Sparsemax

$$\boldsymbol{p} := \text{sparsemax}(\boldsymbol{\theta}) = \text{proj}_\triangle(\boldsymbol{\theta})$$

$$L(\boldsymbol{\theta}, \boldsymbol{y}^{\text{true}}) = \ ?$$

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \boldsymbol{y}^{\text{true}}) = \boldsymbol{p} - \boldsymbol{y}^{\text{true}}$$



$\boldsymbol{\theta}$  $\boldsymbol{p}$

$c_1$

$c_2$

$x$  $\hat{y}$

$\boldsymbol{f_w}(x)$  . . .  $\boldsymbol{g_w}(x, \boldsymbol{p})$

$c_k$

- **sparse** predictive distribution
- reverse-engineer loss from gradient
  *expected* – *observed* statistics
- sparse attention by deriving $\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}}$

# Sparsemax

(Martins and Astudillo, 2016)

$$\boldsymbol{p} := \mathrm{sparsemax}(\boldsymbol{\theta}) = \mathrm{proj}_\triangle(\boldsymbol{\theta})$$

$$L(\boldsymbol{\theta}, \boldsymbol{y}^{\mathrm{true}}) = \ ?$$

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \boldsymbol{y}^{\mathrm{true}}) = \boldsymbol{p} - \boldsymbol{y}^{\mathrm{true}}$$

$x$

$\boldsymbol{\theta}$     $\boldsymbol{p}$

$c_1$
$c_2$
$\ldots$
$c_k$

$f_w(x)$     $g_w(x, \boldsymbol{p})$

$\hat{y}$

- **sparse** predictive distribution
- reverse-engineer loss from gradient
  *expected – observed* statistics
- sparse attention by deriving $\frac{\partial \boldsymbol{p}}{\partial \boldsymbol{\theta}}$

*where do softmax-like functions come from?*

# A Softmax Origin Story 🦸🏽‍♀️

# A Softmax Origin Story 🦸

*First, some background.*

The simplex $\triangle := \{ \boldsymbol{p} \in \mathbb{R}^k \ : \ \boldsymbol{p} \geq \boldsymbol{0}, \ \sum_j p_j = 1 \}$

# A Softmax Origin Story 🦸

*First, some background.*

The simplex $\triangle := \{ \boldsymbol{p} \in \mathbb{R}^k \ : \ \boldsymbol{p} \geq \boldsymbol{0}, \ \sum_j p_j = 1 \}$



$\boldsymbol{p} = [0, 1, 0]$

# A Softmax Origin Story 🦸🏽‍♀️

*First, some background.*

The simplex $\triangle := \{\, \boldsymbol{p} \in \mathbb{R}^k \ : \ \boldsymbol{p} \geq \boldsymbol{0}, \ \ \sum_j p_j = 1 \,\}$



$\boldsymbol{p} = [0, 0, 1]$

# A Softmax Origin Story 🦸

*First, some background.*

The simplex $\triangle := \{ \boldsymbol{p} \in \mathbb{R}^k \ : \ \boldsymbol{p} \geq \boldsymbol{0}, \ \sum_j p_j = 1 \}$



$\boldsymbol{p} = [\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$

# A Softmax Origin Story 🦸🏽‍♀️

*First, some background.*

The simplex $\triangle := \{ \boldsymbol{p} \in \mathbb{R}^k \; : \; \boldsymbol{p} \geq \boldsymbol{0}, \; \sum_j p_j = 1 \}$

Extended value functions $f : \mathbb{R}^k \to \mathbb{R} \cup \{ \infty \}$

# A Softmax Origin Story 🦸

*First, some background.*

The simplex $\triangle := \{ \boldsymbol{p} \in \mathbb{R}^k \ : \ \boldsymbol{p} \geq \boldsymbol{0}, \ \sum_j p_j = 1 \}$

Extended value functions $f : \mathbb{R}^k \to \mathbb{R} \cup \{\infty\}$

$\mathrm{dom}(f) :=$ points where $f$ is finite

# A Softmax Origin Story 🦸

*First, some background.*

The simplex $\triangle := \{ \boldsymbol{p} \in \mathbb{R}^k \ : \ \boldsymbol{p} \geq \boldsymbol{0}, \ \sum_j p_j = 1 \}$

Extended value functions $f : \mathbb{R}^k \to \mathbb{R} \cup \{\infty\}$

$\mathrm{dom}(f) :=$ points where $f$ is finite

Indicator function: $\iota_{\mathcal{S}}(\boldsymbol{x}) = \begin{cases} 0, & \boldsymbol{x} \in \mathcal{S} \\ \infty, & \boldsymbol{x} \notin \mathcal{S} \end{cases}$

# A Softmax Origin Story 🦸

*First, some background.*

The simplex $\triangle := \{ \boldsymbol{p} \in \mathbb{R}^k \ : \ \boldsymbol{p} \geq \boldsymbol{0}, \ \sum_j p_j = 1 \}$

Extended value functions $f : \mathbb{R}^k \to \mathbb{R} \cup \{\infty\}$

$\mathrm{dom}(f) :=$ points where $f$ is finite

Indicator function: $\iota_{\mathcal{S}}(\boldsymbol{x}) = \begin{cases} 0, & \boldsymbol{x} \in \mathcal{S} \\ \infty, & \boldsymbol{x} \notin \mathcal{S} \end{cases}$

$(f + \iota_{\mathcal{S}}$ is $f$ restricted to $\mathcal{S})$

# A Softmax Origin Story 🦸

*First, some background.*

The simplex $\triangle := \{\boldsymbol{p} \in \mathbb{R}^k \; : \; \boldsymbol{p} \geq \boldsymbol{0}, \; \sum_j p_j = 1\}$

Extended value functions $f : \mathbb{R}^k \to \mathbb{R} \cup \{\infty\}$

$\mathrm{dom}(f) :=$ points where $f$ is finite

Indicator function: $\iota_{\mathcal{S}}(\boldsymbol{x}) = \begin{cases} 0, & \boldsymbol{x} \in \mathcal{S} \\ \infty, & \boldsymbol{x} \notin \mathcal{S} \end{cases}$

($f + \iota_{\mathcal{S}}$ is $f$ restricted to $\mathcal{S}$)

Fenchel conjugate of $f : \mathbb{R}^k \to \mathbb{R} \cup \{\infty\}$ :

$$f^*(\boldsymbol{x}) := \sup_{\boldsymbol{p} \in \mathrm{dom}(f)} \langle \boldsymbol{p}, \boldsymbol{x} \rangle - f(\boldsymbol{p})$$

# A Softmax Origin Story 🦸

Let $\Omega = \iota_\triangle$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle$$

# A Softmax Origin Story 🦸🏾‍♀️

Let $\Omega = \iota_\triangle$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle$$

$$\partial \Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle$$

# A Softmax Origin Story 🦸

Let $\Omega = \iota_\triangle$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad\qquad = \max(\boldsymbol{\theta})$$

$$\partial \Omega^*(\boldsymbol{\theta}) = \underset{\boldsymbol{p} \in \triangle}{\operatorname{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad\qquad \ni \operatorname{argmax}(\boldsymbol{\theta})$$

# A Softmax Origin Story 🦸

Let $\Omega = \iota_{\triangle}$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad = \max(\boldsymbol{\theta})$$

$$\partial\Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad \ni \operatorname{argmax}(\boldsymbol{\theta})$$

# A Softmax Origin Story 🦸

Let $\Omega = \iota_\triangle$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p}\in\triangle}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle \qquad = \max(\boldsymbol{\theta})$$

$$\partial\Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p}\in\triangle}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle \qquad \ni \operatorname{argmax}(\boldsymbol{\theta})$$



$\boldsymbol{\theta} = [.7, .1, 1.5]$

$\boldsymbol{p}^\star = [0, 0, 1]$

$$\operatorname*{argmax}_{\boldsymbol{p}\in\triangle}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle = \{\boldsymbol{p}^\star\}$$

# A Softmax Origin Story 🦸🏽‍♀️

Let $\Omega = \iota_\triangle$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad = \max(\boldsymbol{\theta})$$

$$\partial\Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad \ni \operatorname{argmax}(\boldsymbol{\theta})$$

Shannon entropy of $\boldsymbol{p}$ $\quad$ $H_1(\boldsymbol{p}) := -\sum_j p_j \log p_j$

# A Softmax Origin Story 🦸

Let $\Omega = \iota_\triangle$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad = \max(\boldsymbol{\theta})$$

$$\partial \Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad \ni \operatorname{argmax}(\boldsymbol{\theta})$$

**Shannon entropy of $\boldsymbol{p}$** $\quad H_1(\boldsymbol{p}) := -\sum_j p_j \log p_j$

Let $\Omega = -H_1(\boldsymbol{p}) + \iota_\triangle$. Then,

# A Softmax Origin Story 🦸

Let $\Omega = \iota_\triangle$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p}\in\triangle}\langle \boldsymbol{p}, \boldsymbol{\theta}\rangle \qquad = \max(\boldsymbol{\theta})$$

$$\partial\Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p}\in\triangle}\langle \boldsymbol{p}, \boldsymbol{\theta}\rangle \qquad \ni \operatorname{argmax}(\boldsymbol{\theta})$$

**Shannon entropy of $\boldsymbol{p}$** $\quad H_1(\boldsymbol{p}) := -\sum_j p_j \log p_j$

Let $\Omega = -H_1(\boldsymbol{p}) + \iota_\triangle$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p}\in\triangle}\langle \boldsymbol{p}, \boldsymbol{\theta}\rangle + H_1(\boldsymbol{p}) \qquad = \operatorname{logsumexp}(\boldsymbol{\theta})$$

$$\nabla\Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p}\in\triangle}\langle \boldsymbol{p}, \boldsymbol{\theta}\rangle + H_1(\boldsymbol{p}) = \operatorname{softmax}(\boldsymbol{\theta})$$

# A Softmax Origin Story 🦸🏽‍♀️

Let $\Omega = \iota_{\triangle}$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad\qquad = \max(\boldsymbol{\theta})$$

$$\partial\Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle \qquad\qquad \ni \operatorname{argmax}(\boldsymbol{\theta})$$

**Shannon entropy of $\boldsymbol{p}$**  $\quad H_1(\boldsymbol{p}) := -\sum_j p_j \log p_j$

Let $\Omega = -H_1(\boldsymbol{p}) + \iota_{\triangle}$. Then,

$$\Omega^*(\boldsymbol{\theta}) = \max_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle + H_1(\boldsymbol{p}) \qquad = \operatorname{logsumexp}(\boldsymbol{\theta})$$

$$\nabla\Omega^*(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle + H_1(\boldsymbol{p}) = \operatorname{softmax}(\boldsymbol{\theta})$$

**Softmax is an entropy-regularized argmax!**



(0, 0, 1)

1.000
1.400
1.600
1.600

(1, 0, 0)       (0, 1, 0)

# Outline

# Regularized Prediction Functions

## A family of softmax-like mappings

$$\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta}) = \underset{\boldsymbol{p} \in \text{dom}(\Omega)}{\text{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{p}) = \nabla \Omega^*(\boldsymbol{\theta})$$

# Regularized Prediction Functions

## A family of softmax-like mappings

$$\boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p} \in \mathrm{dom}(\Omega)} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{p}) \quad = \nabla\Omega^*(\boldsymbol{\theta})$$

Let $\mathrm{dom}(\Omega) = \triangle$. We recover

# Regularized Prediction Functions

## A family of softmax-like mappings

$$\pi_{\Omega}(\boldsymbol{\theta}) = \underset{\boldsymbol{p} \in \text{dom}(\Omega)}{\text{argmax}} \; \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{p}) \; = \nabla \Omega^*(\boldsymbol{\theta})$$

Let $\text{dom}(\Omega) = \triangle$. We recover

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

- softmax: $\Omega(\boldsymbol{p}) = -H_1(\boldsymbol{p}) = \sum_j p_j \log p_j$

# Regularized Prediction Functions
## A family of softmax-like mappings

$$\boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) = \underset{\boldsymbol{p}\in\text{dom}(\Omega)}{\text{argmax}} \; \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{p}) \; = \nabla\Omega^*(\boldsymbol{\theta})$$

Let $\text{dom}(\Omega) = \triangle$. We recover

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

- **softmax:** $\Omega(\boldsymbol{p}) = -H_1(\boldsymbol{p}) = \sum_j p_j \log p_j$

- **sparsemax:** $\Omega(\boldsymbol{p}) = -H_2(\boldsymbol{p}) = \frac{1}{2}\sum_j p_j(p_j - 1)$

# Regularized Prediction Functions

Regularization brings:

# Regularized Prediction Functions

Regularization brings:

- improved uncertainty handling:
  (predictions become hedged bets)

# Regularized Prediction Functions

Regularization brings:

- improved uncertainty handling:
  (predictions become hedged bets)

- smoothing effect (Nesterov, 2005; Kakade et al., 2009)
  $\Omega$ strongly convex $\Rightarrow \Omega^*$ smooth,
  $\Rightarrow \pi_\Omega$ differentiable almost everywhere

[0, 0, 1]

[.3, 0, .7]

[.3, .2, .5]

...

$w_{j-1}$

$w_j$

$w_{j+1}$

...

# Regularized Prediction Functions

Regularization brings:

- improved uncertainty handling:
  (predictions become hedged bets)

- smoothing effect (Nesterov, 2005; Kakade et al., 2009)
  $\Omega$ strongly convex $\Rightarrow \Omega^*$ smooth,
  $\Rightarrow \boldsymbol{\pi}_\Omega$ differentiable almost everywhere

- ability to add inductive bias



[0, 0, 1]
[.3, 0, .7]
[.3, .2, .5]



...
$w_{j-1}$
$w_j$
$w_{j+1}$
...

fusedmax: $\Omega(\boldsymbol{p}) = -\mathrm{H}_2(\boldsymbol{p}) + \sum_{j=1}^{k} |p_i - p_{i-1}|$

# Outline

$$\text{perceptron} \iff \text{argmax}$$
$$\text{logistic regression} \iff \text{softmax}$$

# What motivates this connection?

# Fenchel-Young Losses

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) := \Omega^*(\boldsymbol{\theta}) + \Omega(\mathbf{y}^{\text{true}}) - \langle \boldsymbol{\theta}, \mathbf{y}^{\text{true}} \rangle$$

$\Omega$:    a regularizer

$\mathbf{y}^{\text{true}} \in \text{dom}(\Omega)$:    target   (e.g. $\boldsymbol{e}_k$)

$\boldsymbol{\theta} \in \mathbb{R}^d$:    prediction scores

# Fenchel-Young Losses

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) := \Omega^*(\boldsymbol{\theta}) + \Omega(\mathbf{y}^{\text{true}}) - \langle \boldsymbol{\theta}, \mathbf{y}^{\text{true}} \rangle$$

$\Omega$:   a regularizer

$\mathbf{y}^{\text{true}} \in \text{dom}(\Omega)$:   target   (e.g. $\boldsymbol{e}_k$)

$\boldsymbol{\theta} \in \mathbb{R}^d$:   prediction scores

Based on the FY inequality:

$$\Omega^\star(\boldsymbol{\theta}) + \Omega(\boldsymbol{p}) \geq \langle \boldsymbol{\theta}, \boldsymbol{p} \rangle$$

**Properties:**

1. Non-negativity:

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) \geq 0$$

# Fenchel-Young Losses

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) := \Omega^*(\boldsymbol{\theta}) + \Omega(\mathbf{y}^{\text{true}}) - \langle \boldsymbol{\theta}, \mathbf{y}^{\text{true}} \rangle$$

$\Omega$:  a regularizer
$\mathbf{y}^{\text{true}} \in \text{dom}(\Omega)$:  target  (e.g. $\boldsymbol{e}_k$)
$\boldsymbol{\theta} \in \mathbb{R}^d$:  prediction scores

Based on the FY inequality:

$$\Omega^\star(\boldsymbol{\theta}) + \Omega(\boldsymbol{p}) \geq \langle \boldsymbol{\theta}, \boldsymbol{p} \rangle$$

**Properties:**

1. Non-negativity:

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) \geq 0$$

2. Zero loss:

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) = 0 \iff \boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) = \mathbf{y}^{\text{true}}$$

The natural loss for the mapping $\boldsymbol{\pi}_\Omega$.

# Fenchel-Young Losses

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) := \Omega^*(\boldsymbol{\theta}) + \Omega(\mathbf{y}^{\text{true}}) - \langle \boldsymbol{\theta}, \mathbf{y}^{\text{true}} \rangle$$

$\Omega$:   a regularizer

$\mathbf{y}^{\text{true}} \in \text{dom}(\Omega)$:   target   (e.g. $\boldsymbol{e}_k$)

$\boldsymbol{\theta} \in \mathbb{R}^d$:   prediction scores

Based on the FY inequality:

$$\Omega^\star(\boldsymbol{\theta}) + \Omega(\boldsymbol{p}) \geq \langle \boldsymbol{\theta}, \boldsymbol{p} \rangle$$

The natural loss for the mapping $\boldsymbol{\pi}_\Omega$.

**Properties:**

1. Non-negativity:

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) \geq 0$$

2. Zero loss:

$$L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) = 0 \iff \boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) = \mathbf{y}^{\text{true}}$$

3. Convex and differentiable:

$$\nabla_{\boldsymbol{\theta}} L_\Omega(\boldsymbol{\theta}; \mathbf{y}^{\text{true}}) = \boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) - \mathbf{y}^{\text{true}}$$

# Well-Known Fenchel-Young Losses

| | $\mathrm{dom}(\Omega)$ | $\Omega(\boldsymbol{p})$ | $\boldsymbol{\pi}_\Omega(\boldsymbol{\theta})$ |
|---|---|---|---|
| Perceptron | $\triangle^k$ | $0$ | $\mathrm{argmax}(\boldsymbol{\theta})$ |
| Logistic Regression | $\triangle^k$ | $- \mathrm{H}_1(\boldsymbol{p})$ | $\mathrm{softmax}(\boldsymbol{\theta})$ |
| Sparsemax | $\triangle^k$ | $- \mathrm{H}_2(\boldsymbol{p})$ | $\mathrm{sparsemax}(\boldsymbol{\theta})$ |

# Well-Known Fenchel-Young Losses

| | $\mathrm{dom}(\Omega)$ | $\Omega(\boldsymbol{p})$ | $\boldsymbol{\pi}_\Omega(\boldsymbol{\theta})$ |
|---|---|---|---|
| Perceptron | $\triangle^k$ | $0$ | $\mathrm{argmax}(\boldsymbol{\theta})$ |
| Logistic Regression | $\triangle^k$ | $-\mathrm{H}_1(\boldsymbol{p})$ | $\mathrm{softmax}(\boldsymbol{\theta})$ |
| Sparsemax | $\triangle^k$ | $-\mathrm{H}_2(\boldsymbol{p})$ | $\mathrm{sparsemax}(\boldsymbol{\theta})$ |
| Hinge (SVM) | $\triangle^k$ | $\langle \boldsymbol{p}, \boldsymbol{y}^{\mathrm{true}} - \boldsymbol{1}\rangle$ | $\mathrm{argmax}(\boldsymbol{1} - \boldsymbol{y}^{\mathrm{true}} + \boldsymbol{\theta})$ |

# Well-Known Fenchel-Young Losses

| | $\mathrm{dom}(\Omega)$ | $\Omega(\boldsymbol{p})$ | $\boldsymbol{\pi}_\Omega(\boldsymbol{\theta})$ |
|---|---|---|---|
| Perceptron | $\triangle^k$ | 0 | $\mathrm{argmax}(\boldsymbol{\theta})$ |
| Logistic Regression | $\triangle^k$ | $-\,\mathrm{H}_1(\boldsymbol{p})$ | $\mathrm{softmax}(\boldsymbol{\theta})$ |
| Sparsemax | $\triangle^k$ | $-\,\mathrm{H}_2(\boldsymbol{p})$ | $\mathrm{sparsemax}(\boldsymbol{\theta})$ |
| Hinge (SVM) | $\triangle^k$ | $\langle \boldsymbol{p}, \boldsymbol{y}^{\mathrm{true}} - \mathbf{1} \rangle$ | $\mathrm{argmax}(\mathbf{1} - \boldsymbol{y}^{\mathrm{true}} + \boldsymbol{\theta})$ |
| Squared | $\mathbb{R}^k$ | $\frac{1}{2}\|\boldsymbol{p}\|^2$ | $\boldsymbol{\theta}$ |

# Well-Known Fenchel-Young Losses

|  | $\mathrm{dom}(\Omega)$ | $\Omega(\boldsymbol{p})$ | $\boldsymbol{\pi}_{\Omega}(\boldsymbol{\theta})$ |
|---|---|---|---|
| Perceptron | $\triangle^k$ | $0$ | $\mathrm{argmax}(\boldsymbol{\theta})$ |
| Logistic Regression | $\triangle^k$ | $-\mathrm{H}_1(\boldsymbol{p})$ | $\mathrm{softmax}(\boldsymbol{\theta})$ |
| Sparsemax | $\triangle^k$ | $-\mathrm{H}_2(\boldsymbol{p})$ | $\mathrm{sparsemax}(\boldsymbol{\theta})$ |
| Hinge (SVM) | $\triangle^k$ | $\langle \boldsymbol{p}, \boldsymbol{y}^{\mathrm{true}} - \mathbf{1} \rangle$ | $\mathrm{argmax}(\mathbf{1} - \boldsymbol{y}^{\mathrm{true}} + \boldsymbol{\theta})$ |
| Squared | $\mathbb{R}^k$ | $\frac{1}{2}\|\boldsymbol{p}\|^2$ | $\boldsymbol{\theta}$ |
| One-vs-all | $[0, 1]^k$ | $-\sum_j \mathrm{H}_1([p_i, 1 - p_i])$ | $\mathrm{sigmoid}(\boldsymbol{\theta})$ |

# Well-Known Fenchel-Young Losses

| | $\mathrm{dom}(\Omega)$ | $\Omega(\boldsymbol{p})$ | $\boldsymbol{\pi}_\Omega(\boldsymbol{\theta})$ |
|---|---|---|---|
| Perceptron | $\triangle^k$ | $0$ | $\mathrm{argmax}(\boldsymbol{\theta})$ |
| Logistic Regression | $\triangle^k$ | $-\mathrm{H}_1(\boldsymbol{p})$ | $\mathrm{softmax}(\boldsymbol{\theta})$ |
| Sparsemax | $\triangle^k$ | $-\mathrm{H}_2(\boldsymbol{p})$ | $\mathrm{sparsemax}(\boldsymbol{\theta})$ |
| Hinge (SVM) | $\triangle^k$ | $\langle \boldsymbol{p}, \boldsymbol{y}^{\mathrm{true}} - \mathbf{1} \rangle$ | $\mathrm{argmax}(\mathbf{1} - \boldsymbol{y}^{\mathrm{true}} + \boldsymbol{\theta})$ |
| Squared | $\mathbb{R}^k$ | $\frac{1}{2}\|\boldsymbol{p}\|^2$ | $\boldsymbol{\theta}$ |
| One-vs-all | $[0, 1]^k$ | $-\sum_j \mathrm{H}_1([p_i, 1 - p_i])$ | $\mathrm{sigmoid}(\boldsymbol{\theta})$ |
| ... and more! | | | |

# Generalized Entropies

A function $H(\boldsymbol{p})$ quantifying uncertainty in $\boldsymbol{p} \in \triangle^k$:

1. $H(\boldsymbol{p}) = 0$ if $\boldsymbol{p} \in \{\boldsymbol{e}_k\}$
2. $H$ strictly concave
3. $H(\boldsymbol{p}) = H(\boldsymbol{Pp})$
   (permutation-invariant)



Tsallis entropies, Rényi entropies, norm entropies, etc.

# Tsallis Entropies

$$H_\alpha(\boldsymbol{p}) = \frac{1}{\alpha(\alpha - 1)} \sum_j (p_j - p_j^\alpha)$$

$\alpha \to 1$     Shannon
$\alpha = 2$     Gini
$\alpha \to \infty$   0

# Tsallis Entropies

$$\mathsf{H}_\alpha(\boldsymbol{p}) = \frac{1}{\alpha(\alpha-1)} \sum_j (p_j - p_j^\alpha)$$

$\alpha \to 1$    Shannon
$\alpha = 2$    Gini
$\alpha \to \infty$    0

generate Tsallis $\alpha$-entmax mappings & losses!

$$\boldsymbol{\pi}_{-H_\alpha}([t, 0])_1$$

Legend:
- $\alpha = 1$ (softmax)
- $\alpha = 1.25$
- $\alpha = 1.5$
- $\alpha = 2$ (sparsemax)
- $\alpha = 4$

$$\boldsymbol{\pi}_{-H_\alpha}([t, 0])_1$$



Legend:
- $\alpha = 1$ (softmax)
- $\alpha = 1.25$
- $\alpha = 1.5$
- $\alpha = 2$ (sparsemax)
- $\alpha = 4$

$$\pi_{-H_\alpha}([t, 0])_1$$

- ⋯ $\alpha = 1$ (softmax)
- — $\alpha = 1.25$
- — $\alpha = 1.5$
- --- $\alpha = 2$ (sparsemax)
- $\alpha = 4$

$$\boldsymbol{\pi}_{-\mathrm{H}_\alpha}\big([t, 0]\big)_1$$

- $\alpha = 1$ (softmax)
- $\alpha = 1.25$
- $\alpha = 1.5$
- $\alpha = 2$ (sparsemax)
- $\alpha = 4$

# Properties of $\alpha$-entmax Mappings & Losses

$\boldsymbol{\pi}_{-H_\alpha}$ is sparse for $\alpha > 1$

(Novel general condition:
$\boldsymbol{\pi}_\Omega$ is sparse iff. $\partial\Omega(\boldsymbol{p}) \neq \varnothing$ for any $\boldsymbol{p} \in \triangle$)

# Properties of $\alpha$-entmax Mappings & Losses

$\boldsymbol{\pi}_{-H_\alpha}$ is sparse for $\alpha > 1$

(Novel general condition:
$\boldsymbol{\pi}_\Omega$ is sparse iff. $\partial\Omega(\boldsymbol{p}) \neq \varnothing$ for any $\boldsymbol{p} \in \triangle$)

$L_{-H_\alpha}$ has the margin property:

$$\theta_k \geq \underbrace{\frac{1}{\alpha-1}}_{m} + \max_{j \neq k} \theta_j \Rightarrow L_{-H_\alpha}(\boldsymbol{\theta}; \boldsymbol{e}_k) = 0$$

(Equivalence result between sparsity and margins)

# Computing $\alpha$-entmax

$$\boldsymbol{\pi}_{-\mathrm{H}_\alpha}(\boldsymbol{\theta}) := \underset{\boldsymbol{p} \in \triangle}{\mathrm{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle + \mathrm{H}_\alpha(\boldsymbol{p})$$

# Computing $\alpha$-entmax

$$\boldsymbol{\pi}_{-\mathrm{H}_\alpha}(\boldsymbol{\theta}) := \operatorname*{argmax}_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle + \mathrm{H}_\alpha(\boldsymbol{p})$$

Solution has the form:

$$\boldsymbol{\pi}_{-\mathrm{H}_\alpha}(\boldsymbol{\theta}) = [(\alpha - 1)\boldsymbol{\theta} - \tau\mathbf{1}]_+^{1/\alpha - 1}$$

# Computing $\alpha$-entmax

$$\boldsymbol{\pi}_{-\mathsf{H}_\alpha}(\boldsymbol{\theta}) := \underset{\boldsymbol{p} \in \triangle}{\mathrm{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle + \mathsf{H}_\alpha(\boldsymbol{p})$$

Solution has the form:

$$\boldsymbol{\pi}_{-\mathsf{H}_\alpha}(\boldsymbol{\theta}) = [(\alpha - 1)\boldsymbol{\theta} - \tau\mathbf{1}]_+^{1/\alpha-1}$$

**Algorithms:**

***bisection***

- approximate; bracket $\tau \in [\tau_{\mathsf{lo}}, \tau_{\mathsf{hi}}]$

- gain 1 bit per $O(d)$ iteration

- `float32` has 23 mantissa bits

# Computing $\alpha$-entmax

$$\boldsymbol{\pi}_{-H_\alpha}(\boldsymbol{\theta}) := \underset{\boldsymbol{p}\in\triangle}{\operatorname{argmax}}\langle\boldsymbol{p}, \boldsymbol{\theta}\rangle + H_\alpha(\boldsymbol{p})$$

Solution has the form:

$$\boldsymbol{\pi}_{-H_\alpha}(\boldsymbol{\theta}) = [(\alpha - 1)\boldsymbol{\theta} - \tau\mathbf{1}]_+^{1/\alpha-1}$$

**Algorithms:**

| *bisection* | *sort-based* |
| --- | --- |
| • approximate; bracket $\tau \in [\tau_{lo}, \tau_{hi}]$ | • exact algorithm, $O(d\log d)$ |
| • gain 1 bit per $O(d)$ iteration | • available only for $\alpha \in \{1.5, 2\}$ |
| • `float32` has 23 mantissa bits | • For $\alpha = 2$, known since Held et al. (1974)! |

# Backward Pass

## (general result)

$$\boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) = \underset{\boldsymbol{p}\in\triangle}{\operatorname{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{p}) \quad = \nabla\Omega^*(\boldsymbol{\theta})$$

$\boldsymbol{p} = \boldsymbol{\pi}_\Omega(\boldsymbol{\theta})$

$\boldsymbol{\theta}$

$$\boldsymbol{J} = \frac{\partial \boldsymbol{\pi}_\Omega}{\partial \boldsymbol{\theta}}$$

# Backward Pass
## (general result)

$$\pi_\Omega(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{p} \in \triangle} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{p}) \; = \nabla\Omega^*(\boldsymbol{\theta})$$

- $\boldsymbol{J}$ symmetric ($= \nabla\nabla\Omega^*$).

$\boldsymbol{p} = \pi_\Omega(\boldsymbol{\theta})$

$\boldsymbol{\theta}$

# **Backward Pass**
## (general result)

$$\boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) = \underset{\boldsymbol{p} \in \triangle}{\text{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{p}) \quad = \nabla\Omega^*(\boldsymbol{\theta})$$

- $\boldsymbol{J}$ symmetric (=$\nabla\nabla\Omega^*$).
- $(\boldsymbol{J})_{ij} = 0$ if $p_i = 0$ or $p_j = 0$.

# Backward Pass

## (general result)

$$\boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) = \underset{\boldsymbol{p} \in \triangle}{\mathrm{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle - \Omega(\boldsymbol{p}) \quad = \nabla\Omega^*(\boldsymbol{\theta})$$

- $\boldsymbol{J}$ symmetric (=$\nabla\nabla\Omega^*$).

- $(\boldsymbol{J})_{ij} = 0$ if $p_i = 0$ or $p_j = 0$.

- Let $(\boldsymbol{H})_{ij} = \frac{\partial^2\Omega}{\partial p_i \partial p_j}(\boldsymbol{p})$ for nonzero $i, j$.

  Let $\boldsymbol{S} = \boldsymbol{H}^{-1}$ and $\boldsymbol{s} = \boldsymbol{1S}$.
  Then, $\bar{\boldsymbol{J}} = \boldsymbol{S} - \frac{1}{\langle \boldsymbol{1}, \boldsymbol{s} \rangle}\, \boldsymbol{ss}^\top$.



$p = \boldsymbol{\pi}_\Omega(\boldsymbol{\theta})$

# **Backward Pass**
## (general result)

$$\boldsymbol{\pi}_\Omega(\boldsymbol{\theta}) = \underset{\boldsymbol{p}\in\triangle}{\mathrm{argmax}}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle - \Omega(\boldsymbol{p}) \quad = \nabla\Omega^*(\boldsymbol{\theta})$$

- $\boldsymbol{J}$ symmetric $(=\nabla\nabla\Omega^*)$.
- $(\boldsymbol{J})_{ij} = 0$ if $p_i = 0$ or $p_j = 0$.
- Let $(\boldsymbol{H})_{ij} = \frac{\partial^2\Omega}{\partial p_i \partial p_j}(\boldsymbol{p})$ for nonzero $i, j$.

  Let $\boldsymbol{S} = \boldsymbol{H}^{-1}$ and $\boldsymbol{s} = \mathbf{1}\boldsymbol{S}$.
  Then, $\bar{\boldsymbol{J}} = \boldsymbol{S} - \frac{1}{\langle\mathbf{1},\boldsymbol{s}\rangle}\,\boldsymbol{s}\boldsymbol{s}^\top$.
- For $-H_\alpha$, $\boldsymbol{S} = \mathrm{diag}(\bar{\boldsymbol{p}}^{2-\alpha})$.



$\boldsymbol{p} = \boldsymbol{\pi}_\Omega(\boldsymbol{\theta})$

$\boldsymbol{\theta}$

# Outline

# Sequence-to-Sequence With Attention

*United   Nations elections    end      today*

# Sequence-to-Sequence With Attention

Encoder

$v_j$

*United*   *Nations*  *elections*   *end*      *today*

# Sequence-to-Sequence With Attention

Encoder

$h_j$

$v_j$

*United*  *Nations*  *elections*  *end*  *today*

# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)



Attention

Encoder

$s_0$

$c_1$

$h_j$

$v_j$

*United  Nations  elections  end  today*

**attention weights**
computed with
*softmax*:

for some decoder state $s_t$,
compute contextually
weighted average of input $c_t$:

$$\theta_j = s_t^\top W^{(a)} h_j$$

$$p = \text{softmax}(\theta)$$

$$c_t = \sum_j p_j h_j$$

# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)



**predictive probability**
(also using *softmax*!)

$$u_t = \tanh(W^{(u)}[s_t; c_t])$$
$$P(y_t \mid y_{1:t-1}, x) = \text{softmax}(Vu_t)$$

$P(y_1 \mid x)$

| | |
|---|---|
| .70 | Eleições |
| .11 | Os |
| .10 | As |
| .09 | Nações |
| | ... |
| $10^{-6}$ | Amsterdam |

Decoder

Attention

Encoder

*Eleições*

$y_1$

$s_0$ $s_1$

$c_1$

$h_j$

$v_j$

*United   Nations  elections   end   today*

# Sequence-to-Sequence With Attention

# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)

# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)

# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)

*Eleições*    *das*    *Nações*

Decoder

$y_1$   $s_0$   $s_1$

Attention   $c_1$

Encoder   $h_j$   $v_j$

*United*   *Nations elections*   *end*   *today*

**predictive probability**

$P(y_3 \mid y_2, y_1, x)$

| | |
|---|---|
| .80 | Nações |
| .11 | Representações |
| .03 | assembleias |
| | ... |
| $10^{-8}$ | resultados |

# Sequence-to-Sequence With Attention

# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)



**predictive probability**

$P(y_4 \mid y_3, y_2, y_1, x)$

| | |
|---|---|
| .90 | Unidas |
| .05 | Shopping |
| .01 | , |
| | ... |
| $10^{-5}$ | aquático |

# Sequence-to-Sequence With Attention

(Bahdanau et al., 2015)



*morphological inflection!*

# The Space of Outputs

# The Space of Outputs



$p(\cdot) = 0.60$

# The Space of Outputs



$p(\cdot) = 0.60$

$p(\cdot) = 0.13$

# The Space of Outputs



$p(\bullet) = 0.60$

$p(\bullet) = 0.13$

$p(\bullet) = 10^{-4}$

# The Space of Outputs: Made Sparse!

$p(\cdot) = 0.70$

$p(\cdot) = 0.20$

$p(\cdot) = 0$ ‼

# Morphological Inflection

(Peters, Niculae, and Martins, 2019)

SIGMORPHON 2018 data, shared multi-lingual model.

Accuracy

- $\alpha = 1$ (softmax)
- $\alpha = 1.5$
- $\alpha = 2$ (sparsemax)

d → r → a → w → e → d → </s>   66.4%

n → </s>   32.2%

</s>   1.4%

+1.83

Neural Machine Translation

(Peters, Niculae, and Martins, 2019)

# Sparse Mappings Don't Slow Down Training

(Peters, Niculae, and Martins, 2019)

Training timing on three DE-EN runs.
Ticks = passes over data.

validation accuracy

63.0%
61.5%
60.0%
58.5%
57.0%

1000  2000  3000  4000  5000  6000  7000
seconds

— softmax
— 1.5-entmax

# Impact of Fine Tuning $\alpha$

Grid search on DE-EN.

# Outline

# Sparse Transformers

# Adaptively Sparse Transformers

Transformers have 6 × 4 × 3 attention heads:
maybe *not all* should be sparse.

# Adaptively Sparse Transformers

Transformers have 6 × 4 × 3 attention heads:
maybe *not all* should be sparse.

Let each attention head learn its $\alpha$!

$$\frac{\partial \pi_{-H_{\alpha}}}{\partial \alpha}$$

# Neural Machine Translation

(Correia, Niculae, and Martins, 2019)

Trajectories of $\alpha$ During Training

(Correia, Niculae, and Martins, 2019)

decoder, layer 1, head 8
encoder, layer 1, head 3
encoder, layer 1, head 4
encoder, layer 2, head 8
encoder, layer 6, head 2

# Previous Position Head

(Correia, Niculae, and Martins, 2019)

softmax   1.5-entmax   $\alpha$-entmax

Learned $\alpha$ = 1.91.

# Outline

# Structured Prediction

# Structured Prediction

# Structured Prediction

# Factorization Into Parts

$$\theta = A^\top \eta$$

# Factorization Into Parts

$$\boldsymbol{\theta} = \boldsymbol{A}^\top \boldsymbol{\eta}$$

# Factorization Into Parts

$$\boldsymbol{\theta} = \boldsymbol{A}^{\top} \boldsymbol{\eta}$$

$\mathcal{M}$

$$\mathcal{M} := \text{conv}\left\{ \boldsymbol{a}_h : h \in \mathcal{H} \right\}$$

$\mathcal{M}$

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



$\mathcal{M}$

$\mathcal{M}$

- **argmax** $\underset{p \in \triangle}{\text{argmax}} \langle p, \theta \rangle$

$\mathcal{M}$

$$\underset{\substack{\textbf{argmax}}}{}\ \underset{p \in \triangle}{\text{argmax}} \langle p, \theta \rangle \qquad \qquad \textbf{MAP}\ \underset{\mu \in \mathcal{M}}{\text{argmax}} \langle \mu, \eta \rangle$$

$\mathcal{M}$

**argmax** $\underset{p \in \triangle}{\text{argmax}} \langle p, \theta \rangle$

**MAP** $\underset{\mu \in \mathcal{M}}{\text{argmax}} \langle \mu, \eta \rangle$

**softmax** $\underset{p \in \triangle}{\text{argmax}} \langle p, \theta \rangle + H(p)$

**argmax** $\underset{\boldsymbol{p} \in \triangle}{\mathrm{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle$

**MAP** $\underset{\boldsymbol{\mu} \in \mathcal{M}}{\mathrm{argmax}} \langle \boldsymbol{\mu}, \boldsymbol{\eta} \rangle$

**softmax** $\underset{\boldsymbol{p} \in \triangle}{\mathrm{argmax}} \langle \boldsymbol{p}, \boldsymbol{\theta} \rangle + \mathsf{H}(\boldsymbol{p})$

**marginals** $\underset{\boldsymbol{\mu} \in \mathcal{M}}{\mathrm{argmax}} \langle \boldsymbol{\mu}, \boldsymbol{\eta} \rangle + \widetilde{\mathsf{H}}(\boldsymbol{\mu})$

$\mathcal{M}$

# Algorithms for specific structures

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

- **argmax** $\underset{\boldsymbol{p}\in\triangle}{\operatorname{argmax}}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle$

- **softmax** $\underset{\boldsymbol{p}\in\triangle}{\operatorname{argmax}}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle + \mathrm{H}(\boldsymbol{p})$

- **sparsemax** $\underset{\boldsymbol{p}\in\triangle}{\operatorname{argmax}}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle - \tfrac{1}{2}\|\boldsymbol{p}\|^2$

- **MAP** $\underset{\boldsymbol{\mu}\in\mathcal{M}}{\operatorname{argmax}}\langle\boldsymbol{\mu},\boldsymbol{\eta}\rangle$

- **marginals** $\underset{\boldsymbol{\mu}\in\mathcal{M}}{\operatorname{argmax}}\langle\boldsymbol{\mu},\boldsymbol{\eta}\rangle + \widetilde{\mathrm{H}}(\boldsymbol{\mu})$

- **argmax** $\underset{\boldsymbol{p}\in\triangle}{\mathrm{argmax}}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle$

- **softmax** $\underset{\boldsymbol{p}\in\triangle}{\mathrm{argmax}}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle + \mathrm{H}(\boldsymbol{p})$

- **sparsemax** $\underset{\boldsymbol{p}\in\triangle}{\mathrm{argmax}}\langle\boldsymbol{p},\boldsymbol{\theta}\rangle - \tfrac{1}{2}\|\boldsymbol{p}\|^2$

- **MAP** $\underset{\boldsymbol{\mu}\in\mathcal{M}}{\mathrm{argmax}}\langle\boldsymbol{\mu},\boldsymbol{\eta}\rangle$

- **marginals** $\underset{\boldsymbol{\mu}\in\mathcal{M}}{\mathrm{argmax}}\langle\boldsymbol{\mu},\boldsymbol{\eta}\rangle + \widetilde{\mathrm{H}}(\boldsymbol{\mu})$

- **SparseMAP** $\underset{\boldsymbol{\mu}\in\mathcal{M}}{\mathrm{argmax}}\langle\boldsymbol{\mu},\boldsymbol{\eta}\rangle - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$

# Generic Algorithm for SparseMAP

$$\boldsymbol{\mu}^{\star} = \operatorname*{argmax}_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

# Generic Algorithm for SparseMAP

linear constraints
*(alas, exponentially many!)*

$$\boldsymbol{\mu}^{\star} = \underset{\boldsymbol{\mu} \in \mathcal{M}}{\mathrm{argmax}}\, \boldsymbol{\mu}^{\top}\boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

quadratic objective

# Generic Algorithm for SparseMAP

linear constraints
*(alas, exponentially many!)*

$$\boldsymbol{\mu}^\star = \underset{\boldsymbol{\mu} \in \mathcal{M}}{\text{argmax}} \, \boldsymbol{\mu}^\top \boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

quadratic objective

### Conditional Gradient

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

# Generic Algorithm for SparseMAP

linear constraints
*(alas, exponentially many!)*

$$\boldsymbol{\mu}^{\star} = \underset{\boldsymbol{\mu} \in \mathcal{M}}{\operatorname{argmax}} \, \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - \frac{1}{2} \|\boldsymbol{\mu}\|^{2}$$

quadratic objective

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

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

# Generic Algorithm for SparseMAP

linear constraints
*(alas, exponentially many!)*

$$\boldsymbol{\mu}^{\star} = \underset{\boldsymbol{\mu} \in \mathcal{M}}{\operatorname{argmax}} \, \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

quadratic objective

### Conditional Gradient

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

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

$$\boldsymbol{a}_{y^{\star}} = \underset{\boldsymbol{\mu} \in \mathcal{M}}{\operatorname{argmax}} \, \boldsymbol{\mu}^{\top} \underbrace{(\boldsymbol{\eta} - \boldsymbol{\mu}^{(t-1)})}_{\tilde{\boldsymbol{\eta}}}$$

# Generic Algorithm for SparseMAP

linear constraints
*(alas, exponentially many!)*

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

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

# Generic Algorithm for SparseMAP

linear constraints
*(alas, exponentially many!)*

$$\boldsymbol{\mu}^\star = \underset{\boldsymbol{\mu} \in \mathcal{M}}{\mathrm{argmax}}\, \boldsymbol{\mu}^\top \boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

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

# Generic Algorithm for SparseMAP

$$\boldsymbol{\mu}^{\star} = \operatorname*{argmax}_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

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: van

  - Quadratic objective: **Active Set**
    (Nocedal and Wright, 1999, Ch. 16.4 & 16.5)
    (Wolfe, 1976; Vinyes and Obozinski, 2017)

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

# Generic Algorithm for SparseMAP

linear constraints
*(alas, exponentially many!)*

$$\boldsymbol{\mu}^\star = \underset{\boldsymbol{\mu} \in \mathcal{M}}{\operatorname{argmax}} \, \boldsymbol{\mu}^\top \boldsymbol{\eta} - \tfrac{1}{2}\|\boldsymbol{\mu}\|^2$$

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

### Backward pass

$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$ is sparse;
precomputed in forward pass!

# Generic Algorithm for SparseMAP

linear constraints
*(alas, exponentially many!)*

$$\boldsymbol{\mu}^{\star} = \operatorname*{argmax}_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

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

## Backward pass

$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$ is sparse;
precomputed in forward pass!

# Generic Algorithm for SparseMAP

$$\boldsymbol{\mu}^{\star} = \operatorname*{argmax}_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^{\top} \boldsymbol{\eta} - \frac{1}{2}\|\boldsymbol{\mu}\|^2$$

linear constraints
*(alas, exponentially many!)*

quadratic objective

**Condit**  Completely modular: just add MAP  **pass**

(Frank and Wolfe, 1956

- select a new c

- 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)

$\frac{\cdot}{\partial \boldsymbol{\eta}}$ is sparse;
precomputed in forward pass!

Sparse Structured Attention for Alignments

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

NLI

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

input

(P, H)

output

A          A
gentleman   police
overlooking  officer
...         ...
situation   closely

entails

contradicts

neutral

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

# Sparse Structured Attention for Alignments

NLI

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

input

(P, H)

output

A          A
gentleman  police
overlooking officer
...        ...
situation  closely

entails
contradicts
neutral

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

# Sparse Structured Attention for Alignments

# Sparse Structured Output Prediction

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

Dependency Parsing, Universal Dependencies

Structured SVM · CRF · SparseMAP

# Sparse Structured Output Prediction
## Training

# Sparse Structured Output Prediction
## Validation: 25% unambiguous, 66% ≤ 5

# Sparse Structured Output Prediction

## Inference captures linguistic ambiguity!

.24

★ the broccoli looks browned around the edges .

.76

**Summary:** Fenchel-Young losses and mappings, a framework for:

*insight into sparsity & margins*

*sparse attention weights*

*sparse output space*



**Next steps:** sparsity in stochastic and generative models.

# Extra slides

# Acknowledgements

*softmax*

*sparsemax*

*1.5-entmax*

# Expressions for Margins

- Main result: $L_{-H}(\boldsymbol{\theta}, \boldsymbol{e}_k)$ has margin $m$ iff. $m\boldsymbol{e}_k \in \partial(-H)(\boldsymbol{e}_k)$.
- If H twice-differentiable, $m_H = \nabla_j H(\boldsymbol{e}_k) - \nabla_k H(\boldsymbol{e}_k)$.
- If $H = \sum_j h(p_j)$ separable, $m_H = h'(0) - h'(1)$.

# Relation With Bregman Divergences

- Bregman divergences are defined in primal space: $B_\Omega : \text{dom}\,\Omega \times \text{dom}\,\Omega \to \mathbb{R}_+$

$$B_\Omega(\boldsymbol{y}||\boldsymbol{p}) := \Omega(\boldsymbol{y}) - \Omega(\boldsymbol{p}) = \langle \nabla\Omega(\boldsymbol{p}), \boldsymbol{y} - \boldsymbol{p} \rangle$$

- FY losses are in **mixed** space: $L_\Omega : \text{dom}(\Omega^\star) \times \text{dom}(\Omega) \to \mathbb{R}_+$
- Denoting $\boldsymbol{\theta} = \nabla\Omega(\boldsymbol{p})$ gives $B_\Omega(\boldsymbol{y}||\boldsymbol{p}) = L_\Omega(\boldsymbol{\theta}; \boldsymbol{y})$.
- However, starting from $\boldsymbol{\theta}$, $B_\Omega(\boldsymbol{y}|| \boldsymbol{\pi}_\Omega(\boldsymbol{\theta}))$ not always convex. ("link function" approach).

# Danskin's Theorem

Let $\phi : \mathbb{R}^k \times \mathcal{Z} \to \mathbb{R}$, $\mathcal{Z} \subset \mathbb{R}^k$ compact.

$$\partial \max_{\boldsymbol{z} \in \mathcal{Z}} \phi(\boldsymbol{x}, \boldsymbol{z}) = \operatorname{conv} \{ \nabla_{\boldsymbol{x}} \phi(\boldsymbol{x}, \boldsymbol{z}^\star) \mid \boldsymbol{z}^\star \in \operatorname*{argmax}_{\boldsymbol{z} \in \mathcal{Z}} \phi(\boldsymbol{x}, \boldsymbol{z}) \}.$$

**Example: maximum of a vector**

# Danskin's Theorem

Let $\phi : \mathbb{R}^k \times \mathcal{Z} \to \mathbb{R}$, $\mathcal{Z} \subset \mathbb{R}^k$ compact.

$$\partial \max_{\boldsymbol{z} \in \mathcal{Z}} \phi(\boldsymbol{x}, \boldsymbol{z}) = \text{conv} \{ \nabla_{\boldsymbol{x}} \phi(\boldsymbol{x}, \boldsymbol{z}^\star) \mid \boldsymbol{z}^\star \in \underset{\boldsymbol{z} \in \mathcal{Z}}{\text{argmax}} \, \phi(\boldsymbol{x}, \boldsymbol{z}) \}.$$

**Example: maximum of a vector**

$$\partial \max_{j \in [d]} \theta_j = \partial \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{\theta}$$

$$= \partial \max_{\boldsymbol{p} \in \Delta} \phi(\boldsymbol{p}, \boldsymbol{\theta})$$

$$= \text{conv} \{ \nabla_{\boldsymbol{\theta}} \phi(\boldsymbol{p}^\star, \boldsymbol{\theta}) \}$$

$$= \text{conv} \{ \boldsymbol{p}^\star \}$$

# Danskin's Theorem

Let $\phi : \mathbb{R}^k \times \mathcal{Z} \to \mathbb{R}$, $\mathcal{Z} \subset \mathbb{R}^k$ compact.

$$\partial \max_{\boldsymbol{z} \in \mathcal{Z}} \phi(\boldsymbol{x}, \boldsymbol{z}) = \text{conv} \{ \nabla_{\boldsymbol{x}} \phi(\boldsymbol{x}, \boldsymbol{z}^\star) \mid \boldsymbol{z}^\star \in \underset{\boldsymbol{z} \in \mathcal{Z}}{\text{argmax}} \, \phi(\boldsymbol{x}, \boldsymbol{z}) \}.$$

## Example: maximum of a vector

$$\begin{aligned}
\partial \max_{j \in [d]} \theta_j &= \partial \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{\theta} \\
&= \partial \max_{\boldsymbol{p} \in \Delta} \phi(\boldsymbol{p}, \boldsymbol{\theta}) \\
&= \text{conv} \{ \nabla_{\boldsymbol{\theta}} \phi(\boldsymbol{p}^\star, \boldsymbol{\theta}) \} \\
&= \text{conv} \{ \boldsymbol{p}^\star \}
\end{aligned}$$



$\boldsymbol{\theta} = [t, 0]$

$\max_j \theta_j$

$\{ g_1 \mid \boldsymbol{g} \in \partial \max_j \theta_j \}$

# Example: Source Sentence with Three Words



softmax · sparsemax · csoftmax · csparsemax

(0.52, 0.35, 0.13) · (0.7, 0.3, 0) · (0.52, 0.35, 0.13) · (0.7, 0.3, 0)

(0.36, 0.44, 0.2) · (0.4, 0.6, 0) · (0.36, 0.44, 0.2) · (0.3, 0.7, 0)

(0.18, 0.27, 0.55) · (0, 0.15, 0.85) · (0.12, 0.21, 0.67) · (0, 0, 1)

Fertilities

# e.g., fertility constraints for NMT



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

# References I

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2015). "Neural machine translation by jointly learning to align and translate". In: *Proc. of ICLR*.

Baker, James K (1979). "Trainable grammars for speech recognition". In: *The Journal of the Acoustical Society of America* 65.S1, S132–S132.

Bertsekas, Dimitri P (1999). *Nonlinear Programming*. Athena Scientific Belmont.

Blondel, Mathieu, André FT Martins, and Vlad Niculae (2019). "Learning classifiers with Fenchel-Young losses: Generalized entropies, margins, and algorithms". In: *Proc. AISTATS*.

Boyd, Stephen and Lieven Vandenberghe (2004). *Convex Optimization*. Cambridge University Press.

Chen, Qian, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, Hui Jiang, and Diana Inkpen (2017). "Enhanced LSTM for natural language inference". In: *Proc. of ACL*.

Chu, Yoeng-Jin and Tseng-Hong Liu (1965). "On the Shortest Arborescence of a Directed Graph". In: *Science Sinica* 14, pp. 1396–1400.

Cocke, William John and Jacob T Schwartz (1970). *Programming languages and their compilers*. Courant Institute of Mathematical Sciences.

Correia, Gonçalo M., Vlad Niculae, and André FT Martins (2019). "Adaptively Sparse Transformers". In: *Proc. EMNLP*.

Cuturi, Marco and Mathieu Blondel (2017). "Soft-DTW: a differentiable loss function for time-series". In: *Proc. of ICML*.

# References II

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.

DeGroot, Morris H (1962). "Uncertainty, information, and sequential experiments". In: *The Annals of Mathematical Statistics*, pp. 404–419.

Edmonds, Jack (1967). "Optimum branchings". In: *J. Res. Nat. Bur. Stand.* 71B, pp. 233–240.

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

Grünwald, Peter D and A Philip Dawid (2004). "Game theory, maximum entropy, minimum discrepancy and robust Bayesian decision theory". In: *Annals of Statistics*, pp. 1367–1433.

Held, Michael, Philip Wolfe, and Harlan P Crowder (1974). "Validation of subgradient optimization". In: *Mathematical Programming* 6.1, pp. 62–88.

Jonker, Roy and Anton Volgenant (1987). "A shortest augmenting path algorithm for dense and sparse linear assignment problems". In: *Computing* 38.4, pp. 325–340.

Kakade, Sham, Shai Shalev-Shwartz, and Ambuj Tewari (2009). "On the duality of strong convexity and strong smoothness: Learning applications and matrix regularization". In: *Tech Report.*

# References III

Kasami, Tadao (1966). "An efficient recognition and syntax-analysis algorithm for context-free languages". In: *Coordinated Science Laboratory Report no. R-257*.

Kirchhoff, Gustav (1847). "Ueber die Auflösung der Gleichungen, auf welche man bei der Untersuchung der linearen Vertheilung galvanischer Ströme geführt wird". In: *Annalen der Physik* 148.12, pp. 497–508.

Kuhn, Harold W (1955). "The Hungarian method for the assignment problem". In: *Nav. Res. Log.* 2.1-2, pp. 83–97.

Lacoste-Julien, Simon and Martin Jaggi (2015). "On the global linear convergence of Frank-Wolfe optimization variants". In: *Proc. of NeurIPS*.

Malaviya, Chaitanya, Pedro Ferreira, and André FT 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.

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

Niculae, Vlad and Mathieu Blondel (2017). "A regularized framework for sparse and structured neural attention". In: *Proc. of NeurIPS*.

# References IV

Niculae, Vlad, André FT Martins, Mathieu Blondel, and Claire Cardie (2018). "SparseMAP: Differentiable sparse structured inference". In: *Proc. of ICML.*

Nocedal, Jorge and Stephen Wright (1999). *Numerical Optimization*. Springer New York.

Peters, Ben, Vlad Niculae, and André FT Martins (2019). "Sparse sequence-to-sequence models". In: *Proc. ACL.*

Rabiner, Lawrence R. (1989). "A tutorial on Hidden Markov Models and selected applications in speech recognition". In: *P. IEEE* 77.2, pp. 257–286.

Sakoe, Hiroaki and Seibi Chiba (1978). "Dynamic programming algorithm optimization for spoken word recognition". In: *IEEE Trans. on Acoustics, Speech, and Sig. Proc.* 26, pp. 43–49.

Taskar, Ben (2004). "Learning structured prediction models: A large margin approach". PhD thesis. Stanford University.

Tsallis, Constantino (1988). "Possible generalization of Boltzmann-Gibbs statistics". In: *Journal of Statistical Physics* 52, pp. 479–487.

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

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin (2017). "Attention Is All You Need". In: *Proc. of NeurIPS.*

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

# References V

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

📄 Younger, Daniel H (1967). "Recognition and parsing of context-free languages in time $n^3$". In: *Information and Control* 10.2, pp. 189–208.