

Sparse Sequence-to-Sequence Models

Ben Peters Instituto de Telecomunicações

→ **Vlad Niculae** IT

André Martins IT & Unbabel

Sequence-to-Sequence With Attention

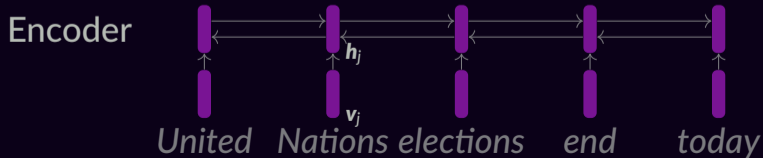
United Nations elections end today

Sequence-to-Sequence With Attention

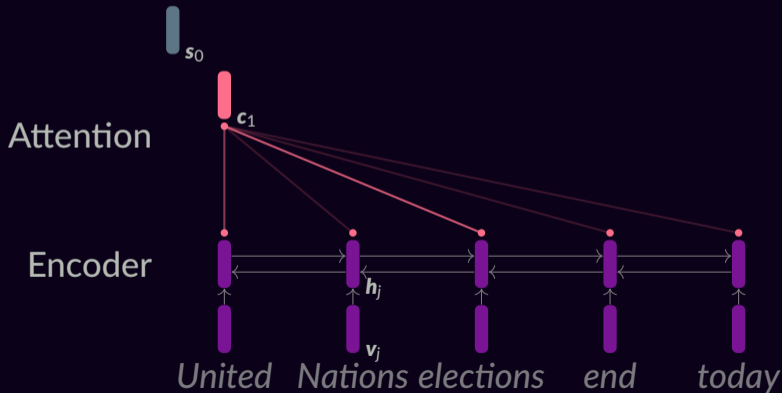
Encoder



Sequence-to-Sequence With Attention



Sequence-to-Sequence With Attention



attention weights
computed with
softmax:

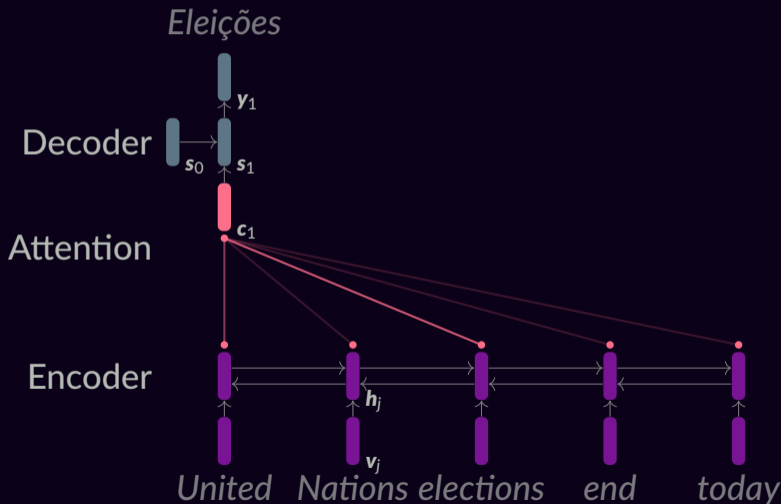
for some decoder state s_t ,
compute contextually
weighted average of input c_t :

$$z_j = s_t^T \mathbf{W}^{(a)} h_j$$

$$\pi_j = \text{softmax}_j(\mathbf{z})$$

$$c_t = \sum_j \pi_j h_j$$

Sequence-to-Sequence With Attention



predictive probability
(also using *softmax*!)

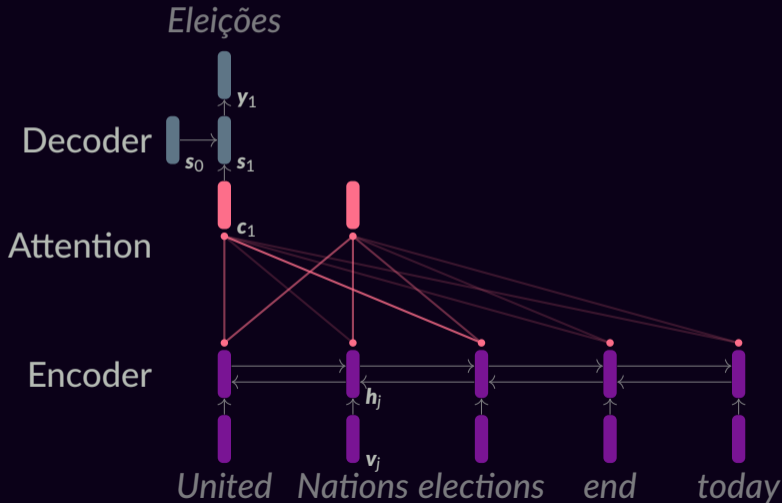
$$\mathbf{u}_t = \tanh(\mathbf{W}^{(u)}[\mathbf{s}_t; \mathbf{c}_t])$$

$$P(y_t | y_{1:t-1}, \mathbf{x}) = \text{softmax}(\mathbf{V}\mathbf{u}_t)$$

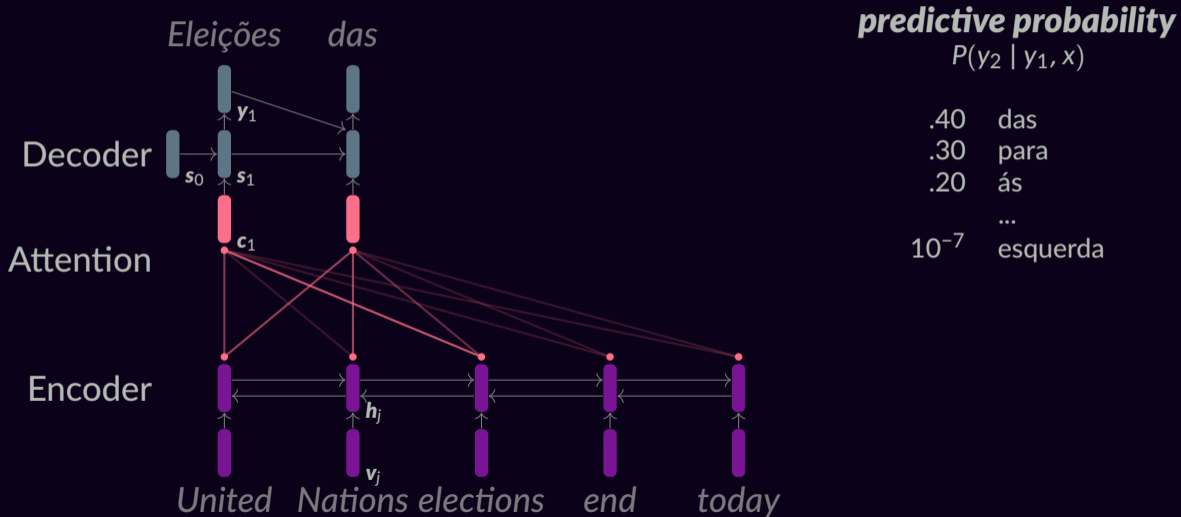
$P(y_1 | \mathbf{x})$

.70	Eleições
.11	Os
.10	As
.09	Nações
...	...
10^{-6}	Bucarest

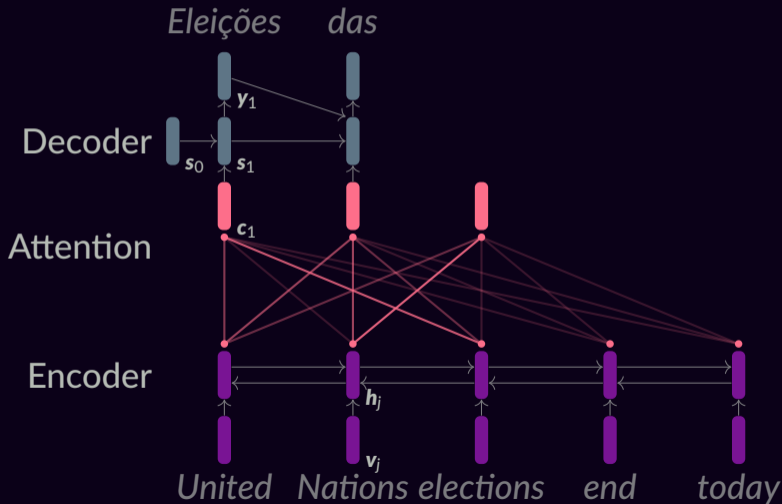
Sequence-to-Sequence With Attention



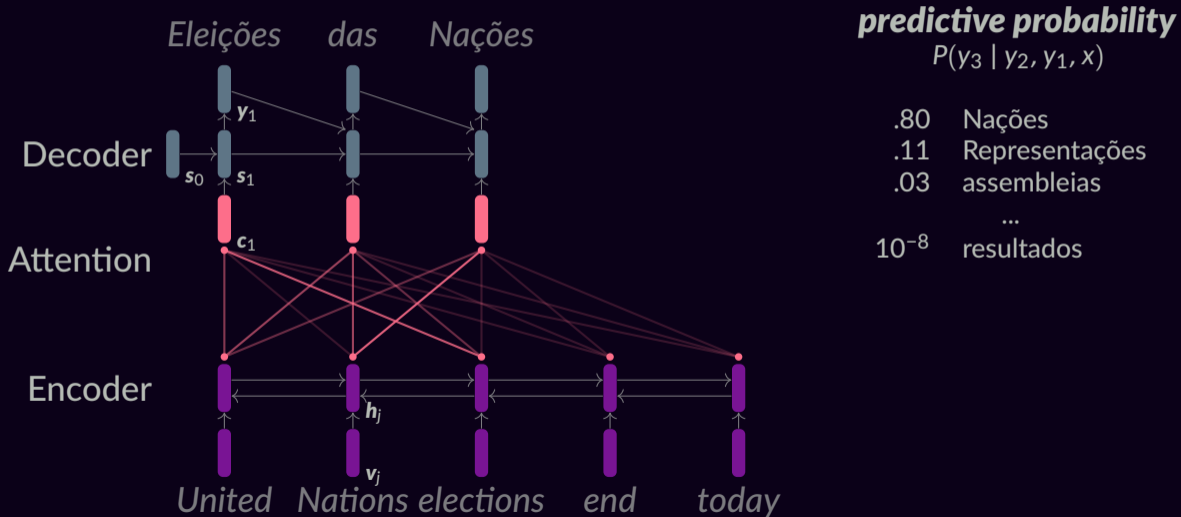
Sequence-to-Sequence With Attention



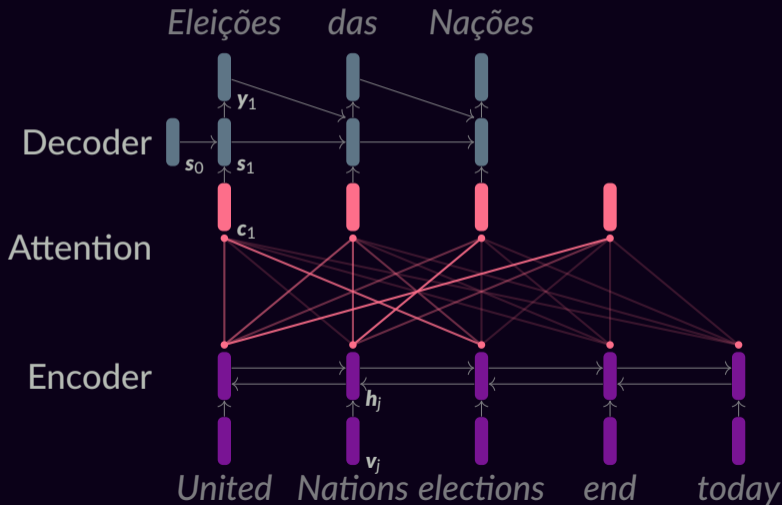
Sequence-to-Sequence With Attention



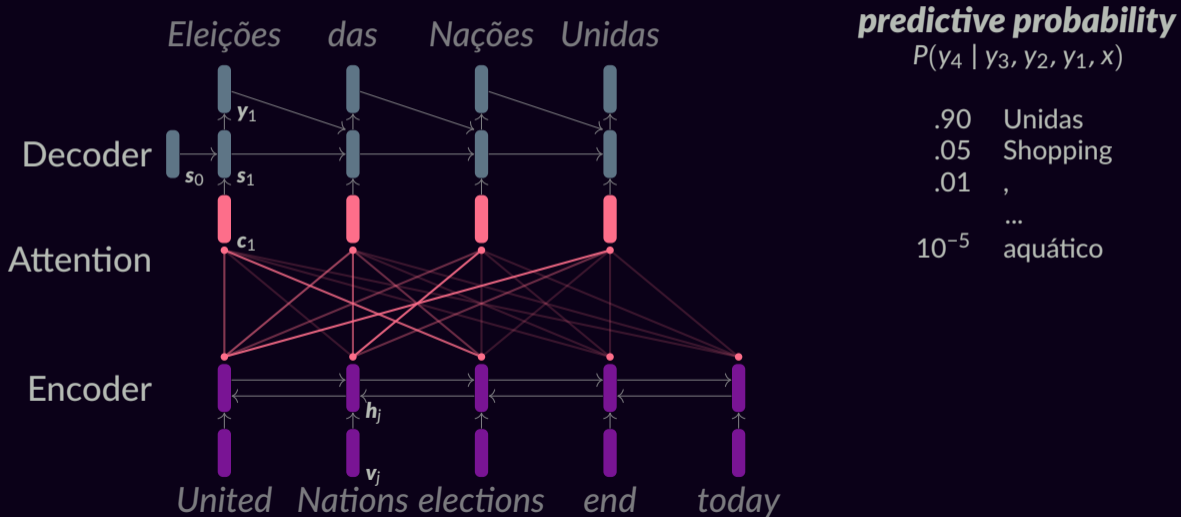
Sequence-to-Sequence With Attention



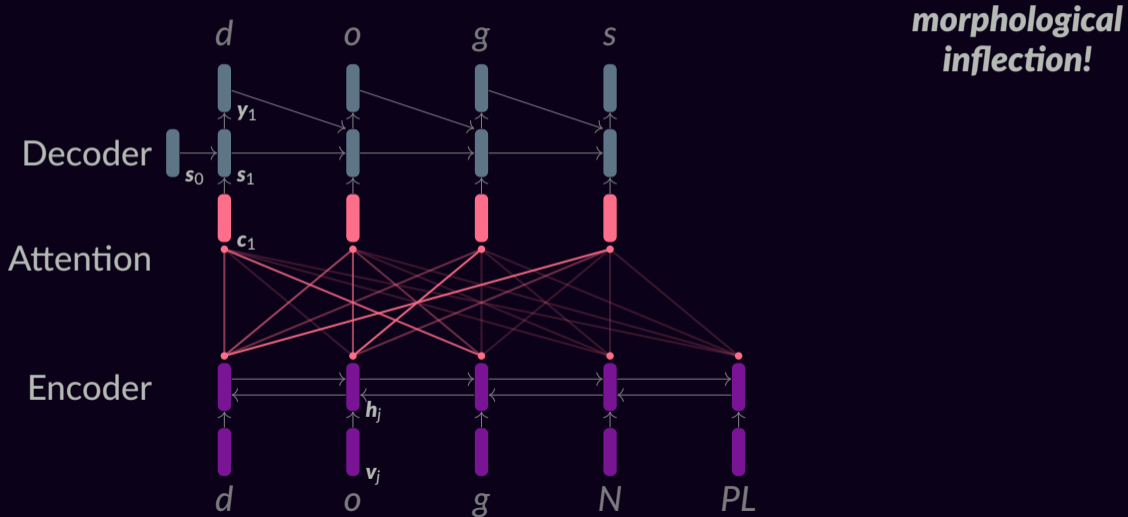
Sequence-to-Sequence With Attention



Sequence-to-Sequence With Attention



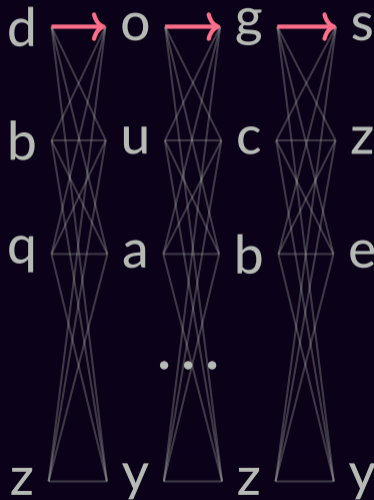
Sequence-to-Sequence With Attention



The Space of Outputs

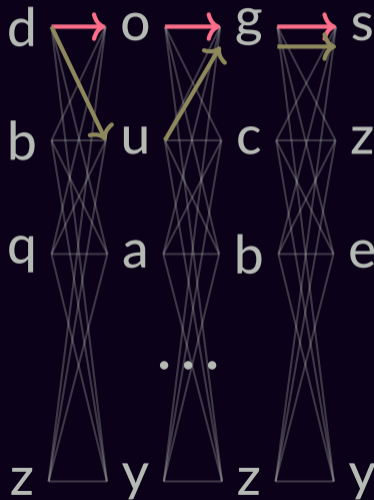


The Space of Outputs



$$p(\cdot) = 0.60$$

The Space of Outputs



$$p(\cdot) = 0.60$$

$$p(\cdot) = 0.13$$

The Space of Outputs



$$p(\cdot) = 0.60$$

$$p(\cdot) = 0.13$$

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

The Space of Outputs: Made Sparse!



$$p(\cdot) = 0.70$$

$$p(\cdot) = 0.20$$

$$p(\cdot) = 0 \text{ !!}$$

Softmax plays two roles in seq2seq:

attention weights

for some decoder state \mathbf{s}_t , compute contextually weighted average of input \mathbf{c}_t :

$$\mathbf{z}_j = \mathbf{s}_t^\top \mathbf{W}^{(a)} \mathbf{h}_j$$

$$\pi_j = \text{softmax}_j(\mathbf{z})$$

$$\mathbf{c}_t = \sum_j \pi_j \mathbf{h}_j$$

output probabilities

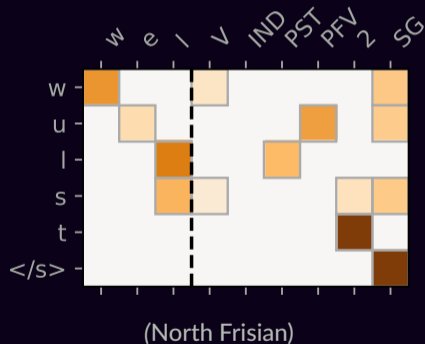
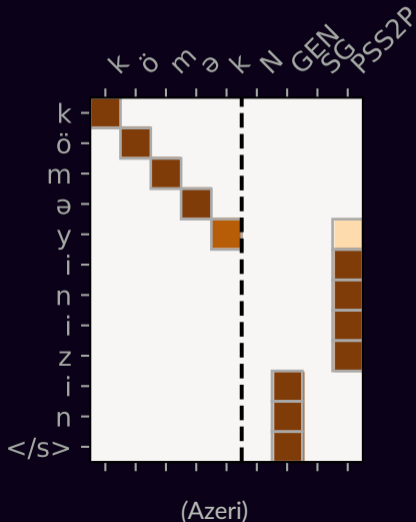
predict the probability of the next word:

$$\mathbf{u}_t = \tanh(\mathbf{W}^{(u)}[\mathbf{s}_t; \mathbf{c}_t])$$

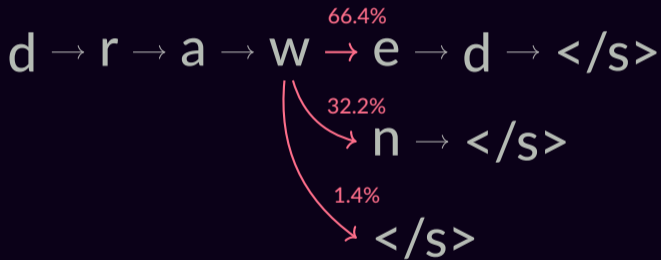
$$P(y_t | y_{1:t-1}, \mathbf{x}) = \text{softmax}(\mathbf{V}\mathbf{u}_t)$$

Our work: replace softmax
with a *family of* new sparsity-inducing alternatives

Sparse Attention Weights / Alignments



Sparse Predictive Probabilities



What is softmax?

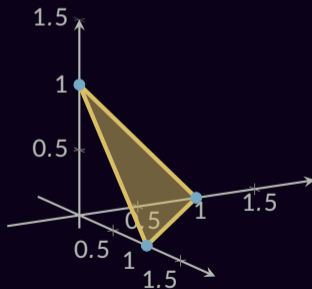
Often defined via $p_i := \frac{\exp z_i}{\sum_j \exp z_j}$, but where does it come from?

What is softmax?

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

$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

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

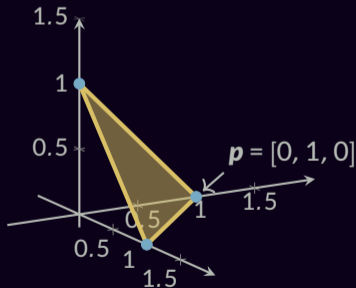


What is softmax?

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

$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

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

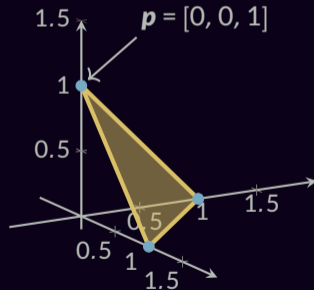


What is softmax?

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

$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

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

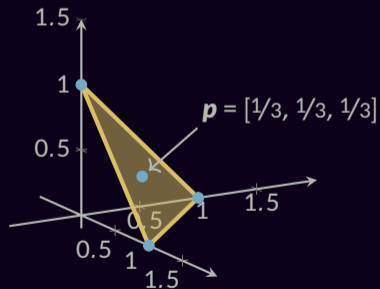


What is softmax?

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

$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

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



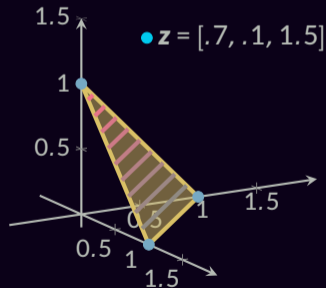
What is softmax?

Often defined via $p_j := \frac{\exp z_j}{\sum_j \exp z_j}$, but where does it come from?

$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

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

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



What is softmax?

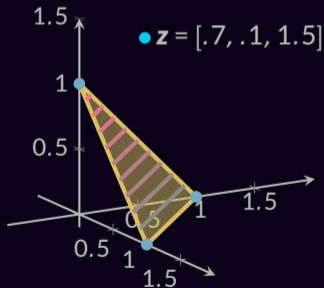
Often defined via $p_i := \frac{\exp z_i}{\sum_j \exp z_j}$, but where does it come from?

$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

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

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

argmax



What is softmax?

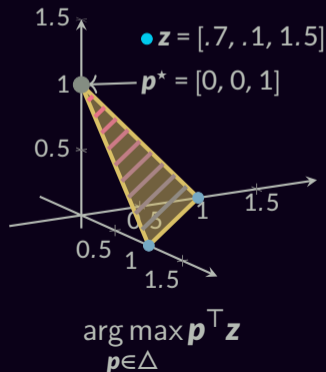
Often defined via $p_j := \frac{\exp z_j}{\sum_j \exp z_j}$, but where does it come from?

$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

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

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

argmax maximizes **expected score**



What is softmax?

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

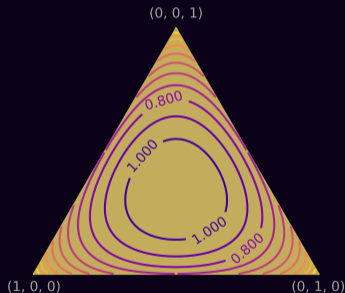
$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

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

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

argmax maximizes **expected score**

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



What is softmax?

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

$$\Delta := \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0}, \sum_j p_j = 1\}$$

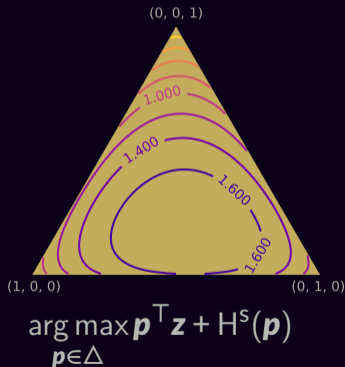
$\mathbf{p} \in \Delta$: probability distribution over choices

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

argmax maximizes **expected score**

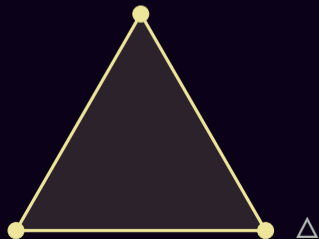
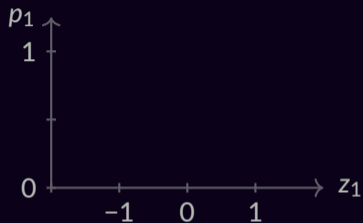
Shannon entropy of \mathbf{p} : $H^s(\mathbf{p}) := -\sum_i p_i \log p_i$

softmax maximizes **expected score + entropy**:



Generalizing Softmax Using Entropies

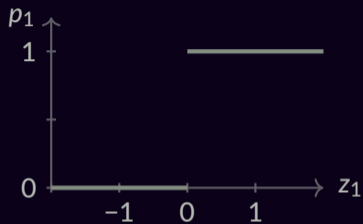
$$\boldsymbol{\pi}_H(\mathbf{z}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{z} + H(\mathbf{p})$$



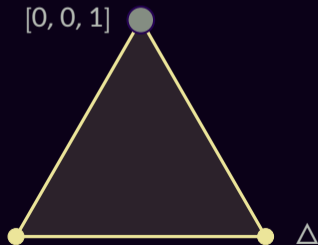
Generalizing Softmax Using Entropies

$$\boldsymbol{\pi}_H(\mathbf{z}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{z} + H(\mathbf{p})$$

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



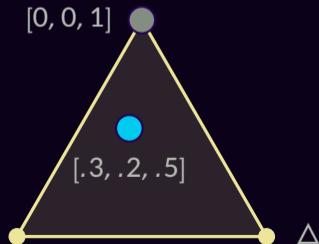
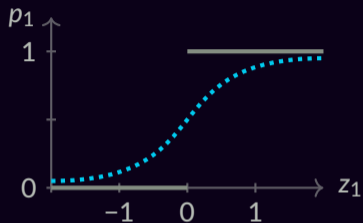
$[0, 0, 1]$



Generalizing Softmax Using Entropies

$$\boldsymbol{\pi}_H(\mathbf{z}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{z} + H(\mathbf{p})$$

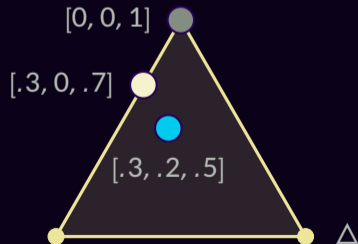
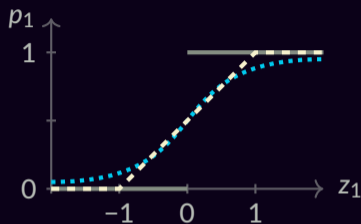
- argmax: $H^0(\mathbf{p}) = 0$
- softmax: $H^s(\mathbf{p}) = -\sum_j p_j \log p_j$



Generalizing Softmax Using Entropies

$$\boldsymbol{\pi}_H(\mathbf{z}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{z} + H(\mathbf{p})$$

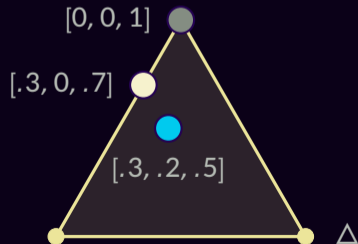
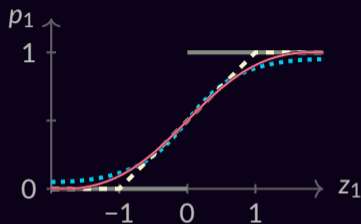
- argmax: $H^0(\mathbf{p}) = 0$
- softmax: $H^s(\mathbf{p}) = -\sum_j p_j \log p_j$
- sparsemax: $H^g(\mathbf{p}) = 1/2 \sum_j p_j (1 - p_j)$



Generalizing Softmax Using Entropies

$$\boldsymbol{\pi}_H(\mathbf{z}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{z} + H(\mathbf{p})$$

- argmax: $H^0(\mathbf{p}) = 0$
- softmax: $H^s(\mathbf{p}) = -\sum_j p_j \log p_j$
- sparsemax: $H^g(\mathbf{p}) = 1/2 \sum_j p_j(1 - p_j)$
- α -entmax: $H_\alpha^t(\mathbf{p}) = \frac{1}{\alpha(\alpha-1)} \sum_j (p_j - p_j^\alpha)$



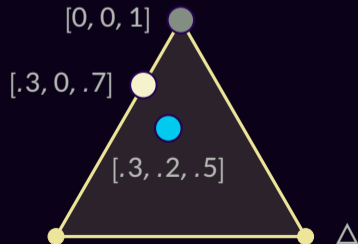
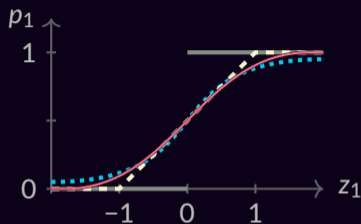
Generalizing Softmax Using Entropies

$$\boldsymbol{\pi}_H(\mathbf{z}) = \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{z} + H(\mathbf{p})$$

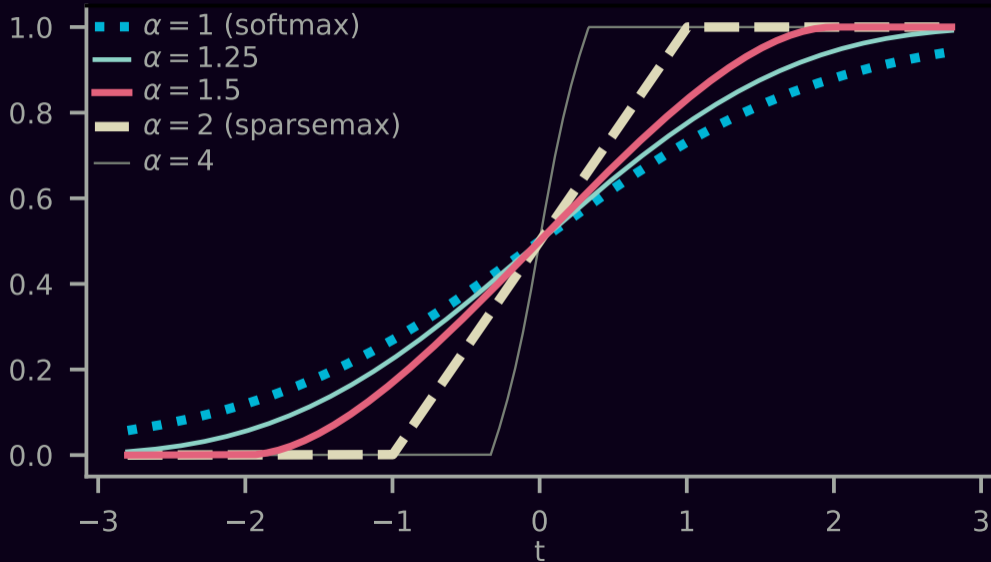
- argmax: $H^0(\mathbf{p}) = 0$
- softmax: $H^s(\mathbf{p}) = -\sum_j p_j \log p_j$
- sparsemax: $H^g(\mathbf{p}) = 1/2 \sum_j p_j(1 - p_j)$
- α -entmax: $H_\alpha^t(\mathbf{p}) = \frac{1}{\alpha(\alpha-1)} \sum_j (p_j - p_j^\alpha)$

Tsallis α -entropy (Tsallis, 1988).

Depicted: $\alpha = 1.5$. Uncovers softmax ($\alpha \rightarrow 1$) and sparsemax ($\alpha = 2$).



$$\pi_\alpha([t, 0])_1$$



Computing α -entmax

$$\boldsymbol{\pi}_{H_{\alpha}^t}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{p}^T \mathbf{z} + H_{\alpha}^t(\mathbf{p})$$

Solution has the form:

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

Algorithms:

bisection

- approximate; bracket $\tau \in [\tau_{lo}, \tau_{hi}]$
- gain 1 bit per $O(d)$ iteration
- float32 has 23 mantissa bits

sort-based

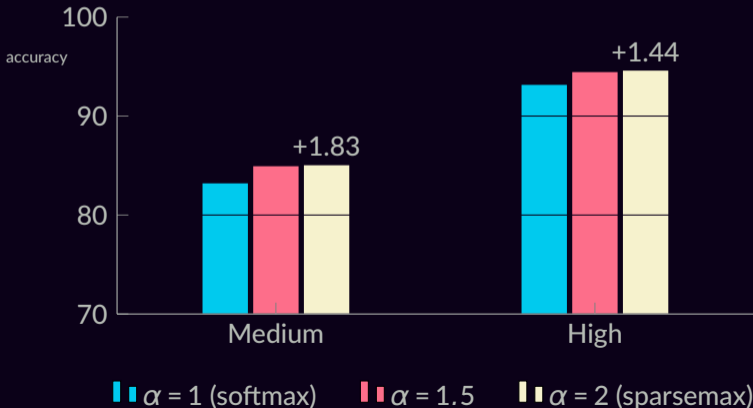
- exact algorithm, $O(d \log d)$
- available only for $\alpha \in \{1.5, 2\}$
- huge speed-up from partial sorting when expecting sparse solutions

Morphological inflection

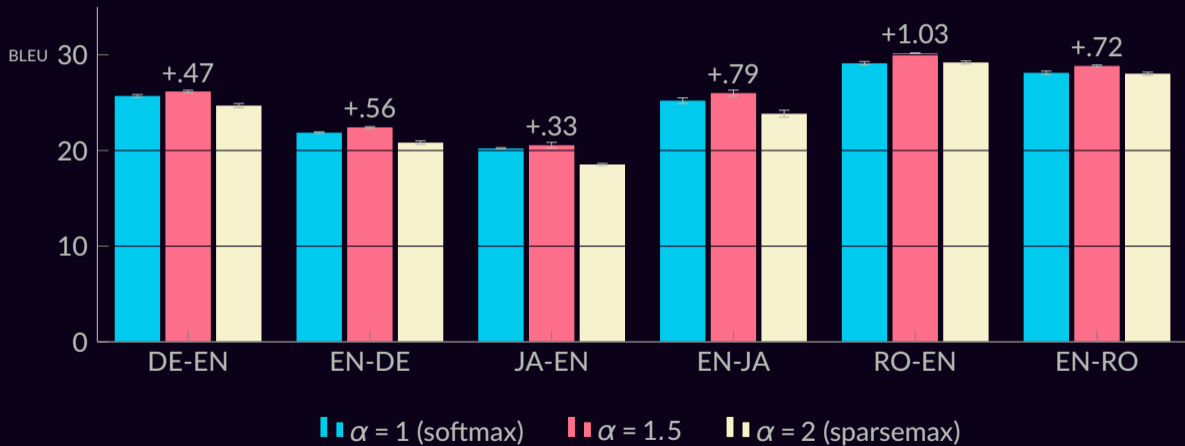
SIGMORPHON 2018. Shared multi-lingual model.

Medium: 1k pairs per language, 102 languages.

High: 10k pairs per language, 86 languages.

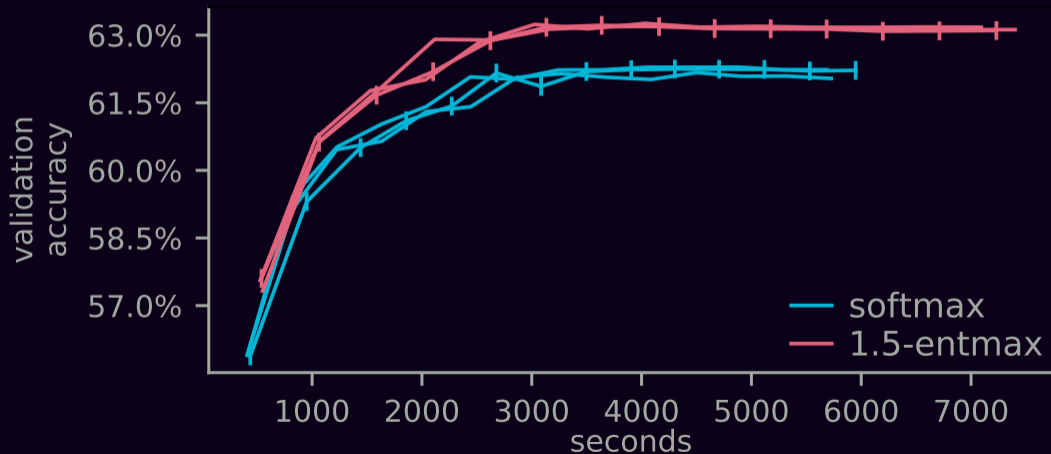


Neural Machine Translation



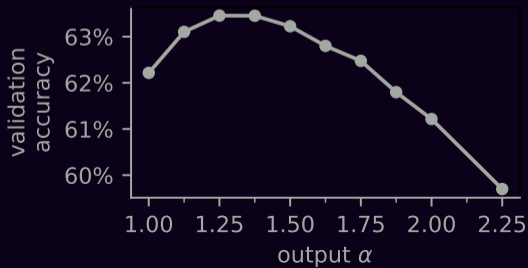
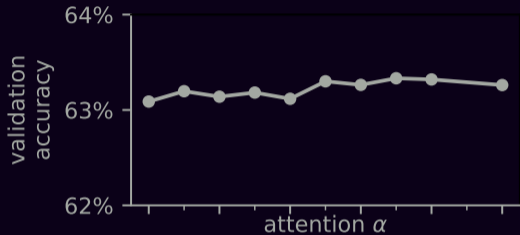
Sparse Mappings Don't Slow Down Training

Training timing on three DE-EN runs.

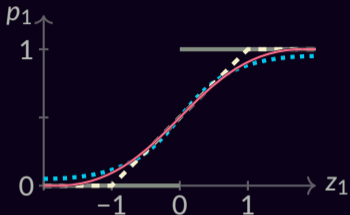


Impact of Fine Tuning α

Grid search on DE-EN.

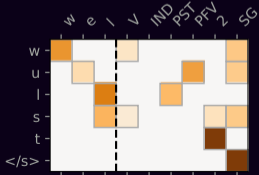


Sparse Seq2Seq: Conclusions

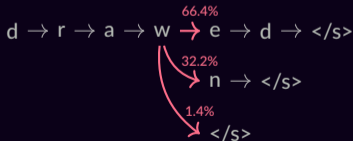


New family of sparse mappings α -entmax, algorithms for efficient forward & backward passes.

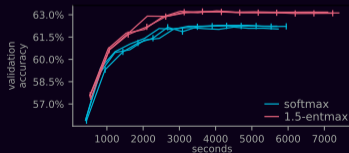
sparse attention weights



sparse output space



performance improvements



Acknowledgements



This work was supported by the European Research Council (ERC StG DeepSPIN 758969) and by the Fundação para a Ciência e Tecnologia through contract UID/EEA/50008/2019 and CMUPERI/TIC/0046/2014 (GoLocal).

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

References I

- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2015). “Neural machine translation by jointly learning to align and translate”. In: *Proc. of ICLR*.
- Blondel, Mathieu, André FT Martins, and Vlad Niculae (2019). “Learning classifiers with Fenchel-Young losses: Generalized entropies, margins, and algorithms”. In: *Proc. AISTATS*.
- 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*.
- Niculae, Vlad and Mathieu Blondel (2017). “A regularized framework for sparse and structured neural attention”. In: *Proc. of NeurIPS*.
- Peters, Ben and André FT Martins (2019). “IT-IST at the SIGMORPHON 2019 shared task: Sparse two-headed models for inflection”. In: *Proc. SIGMORPHON*.
- Tsallis, Constantino (1988). “Possible generalization of Boltzmann-Gibbs statistics”. In: *Journal of Statistical Physics* 52, pp. 479–487.