

# From Sparse Modeling to Sparse Communication

André Martins



Erice International School on Complexity

Ettore Majorana Foundation and Centre for Scientific Culture, Erice, Italy

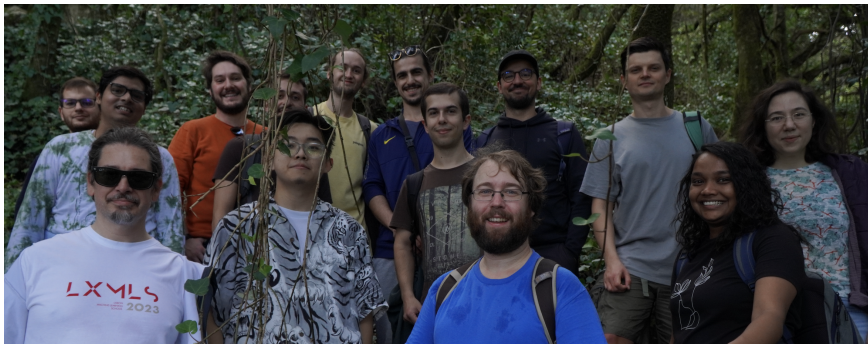
April 27, 2024

# Our Amazing Team



**SARDINE:** Structure AwaRe moDelling for Natural Language

# Our Amazing Team



**SARDINE:** Structure AwaRe moDElling for Natural Language

# DeepSPIN & DECOLLAGE



- ERC starting grant (2018–23) and consolidator grant (2023–28)
- Goal: put together *deep learning* and *structured prediction* for *natural language processing*
- More details: <https://deep-spin.github.io>



## From Sparse Modeling ...

- Mostly used with linear models, lots of work in the 2000s
- Main idea: embed a sparse regularizer (e.g.  $\ell_1$ -norm) in the learning objective
- Irrelevant features get zero weight and can be discarded
- Extensions to structured sparsity (group-lasso, fused-lasso, etc.)

## ... to Sparse Communication:

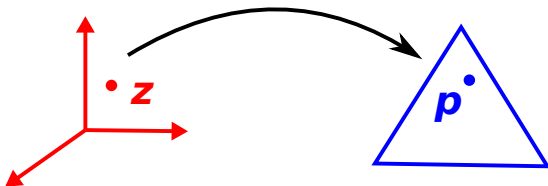
- Mostly used with neural networks, most work after 2015
- Main idea: sparse neuron activations (biological plausibility)
- Predictions are triggered by a few neurons only (input-dependent)
- Example: ReLUs, dropout, sparse attention mechanisms

# This Talk

An inventory of transformations that capture **sparsity** and **structure**:

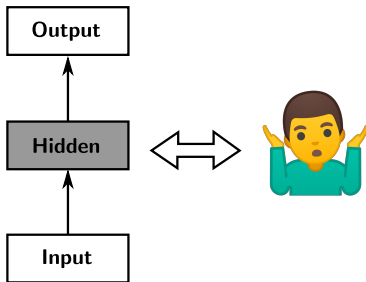
- All differentiable (efficient forward and backward propagation)
- Can be used at hidden (attention) or output layers (loss)
- Can make a bridge between the **continuous** and **discrete** worlds
- Effective in several natural language processing tasks.

Building block:

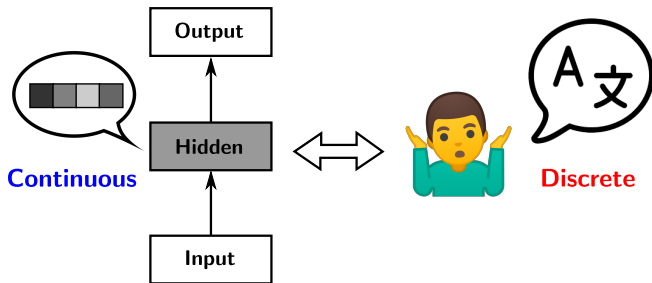


Sparse transformations from the Euclidean space to the simplex  $\triangle$ .

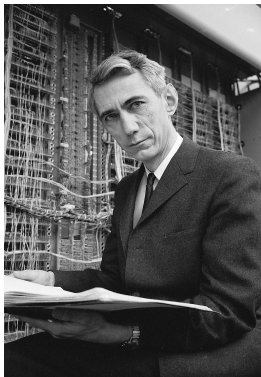
# Machine-Human Communication



# Machine-Human Communication







# The Bell System Technical Journal

Vol. XXVII

July, 1948

No. 3

---

## A Mathematical Theory of Communication

By C. E. SHANNON

INTRODUCTION

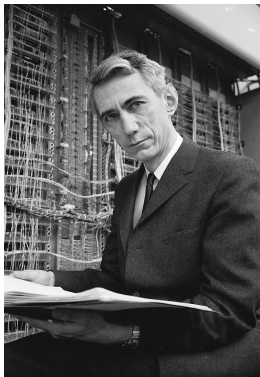
PART I: DISCRETE NOISELESS SYSTEMS

PART II: THE DISCRETE CHANNEL WITH NOISE

PART III: MATHEMATICAL PRELIMINARIES

PART IV: THE CONTINUOUS CHANNEL

PART V: THE RATE FOR A CONTINUOUS SOURCE



# The Bell System Technical Journal

Vol. XXVII

July, 1948

No. 3

## A Mathematical Theory of Communication

By C. E. SHANNON

INTRODUCTION

PART I: DISCRETE NOISELESS SYSTEMS

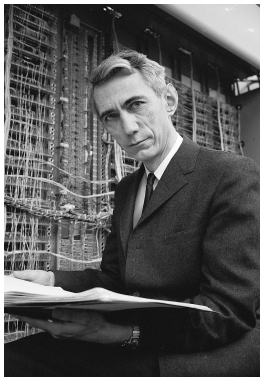
PART II: THE DISCRETE CHANNEL WITH NOISE

PART III: MATHEMATICAL PRELIMINARIES

PART IV: THE CONTINUOUS CHANNEL

PART V: THE RATE FOR A CONTINUOUS SOURCE





# The Bell System Technical Journal

Vol. XXVII

July, 1948

No. 3

## A Mathematical Theory of Communication

By C. E. SHANNON

### INTRODUCTION

PART I: DISCRETE NOISELESS SYSTEMS

PART II: THE DISCRETE CHANNEL WITH NOISE

PART III: MATHEMATICAL PRELIMINARIES

PART IV: THE CONTINUOUS CHANNEL

PART V: THE RATE FOR A CONTINUOUS SOURCE

$\Sigma$

$\int$

$$\Sigma \text{ vs. } \int$$

Commonly we have to opt between **discrete** or **continuous** models:

- Language is symbolic and *discrete*
- Neural networks use (and learn) *continuous* representations

We should look at what happens in-between!

**Sparsity** might help with this, but...

# $\Sigma$ vs. $\int$

Commonly we have to opt between **discrete** or **continuous** models:

- Language is symbolic and *discrete*
- Neural networks use (and learn) *continuous* representations

We should look at what happens in-between!

**Sparsity** might help with this, but...

... sparse probabilities are understudied and often excluded from theory:

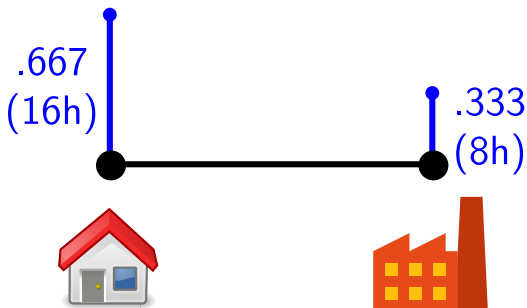
- Hammersley-Clifford theorem in graphical models
- Pitman-Koopman-Darmois theorem (sufficient statistics and exponential families)
- Log-likelihood is  $-\infty$  if estimated probability is 0.

## Motivating Example: John's Life

John splits his day as follows: he works 8h/day, and stays home 16h/day.

## Motivating Example: John's Life

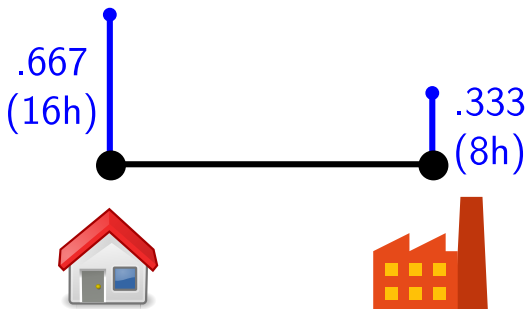
John splits his day as follows: he works 8h/day, and stays home 16h/day.



## Motivating Example: John's Life

John splits his day as follows: he works 8h/day, and stays home 15h/day.

He is in transit 1h/day to commute to work and back.

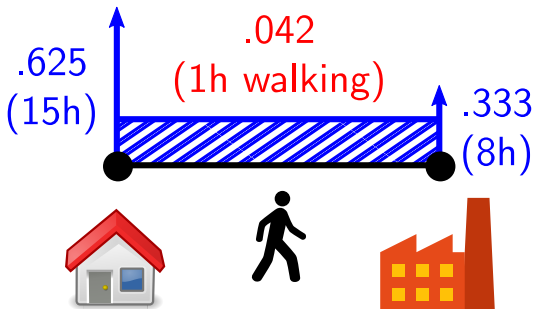




## Motivating Example: John's Life

John splits his day as follows: he works 8h/day, and stays home 15h/day.

He is in transit 1h/day to commute to work and back.



## Motivating Example: John's Life

John splits his day as follows: he works 8h/day, and stays home 15h/day.

He is in transit 1h/day to commute to work and back.

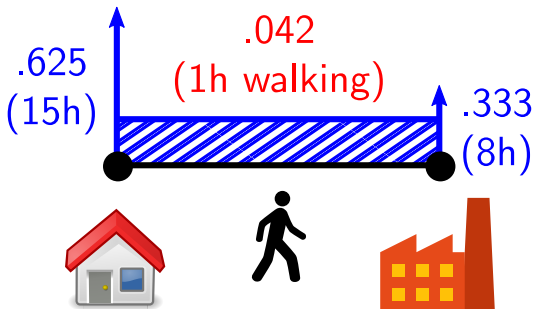


Is John's location a *discrete* or *continuous* random variable?

## Motivating Example: John's Life

John splits his day as follows: he works 8h/day, and stays home 15h/day.

He is in transit 1h/day to commute to work and back.



Is John's location a *discrete* or *continuous* random variable? It's **mixed**.

# Outline

- 1 Sparse Transformations
- 2 Fenchel-Young Losses
- 3 Sparse Hopfield Networks
- 4 Mixed Distributions
- 5 Conclusions

## Recap: Softmax and Argmax

Softmax exponentiates and normalizes:

$$\text{softmax}(\mathbf{z}) = \frac{\exp(\mathbf{z})}{\sum_{k=1}^K \exp(z_k)}$$

- **Fully dense:**  $\text{softmax}(\mathbf{z}) > 0, \forall \mathbf{z}$
- Used both as a loss function (cross-entropy) and for attention.

## Recap: Softmax and Argmax

Softmax exponentiates and normalizes:

$$\text{softmax}(\mathbf{z}) = \frac{\exp(\mathbf{z})}{\sum_{k=1}^K \exp(z_k)}$$

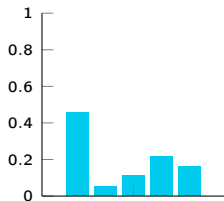
- **Fully dense:**  $\text{softmax}(\mathbf{z}) > 0, \forall \mathbf{z}$
- Used both as a loss function (cross-entropy) and for attention.

Argmax can be written as:

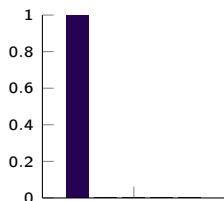
$$\begin{aligned} \text{argmax}(\mathbf{z}) &:= \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^\top \mathbf{p} \\ &= \lim_{\tau \rightarrow 0^+} \text{softmax}(\mathbf{z}/\tau) \quad (\text{temperature trick}) \end{aligned}$$

- Retrieves a **one-hot vector** for the highest scored index.

softmax( $\mathbf{z}$ )



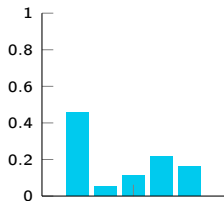
argmax( $\mathbf{z}$ )



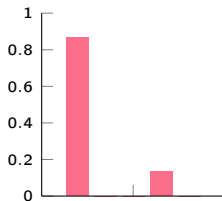
(Same  $\mathbf{z} = [1.0716, -1.1221, -0.3288, 0.3368, 0.0425]$ )

- Argmax is an extreme case of sparsity, but it is **discontinuous**.
- Is there a **sparse** and **differentiable** alternative?

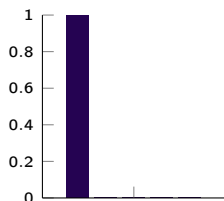
softmax( $\mathbf{z}$ )



sparsemax( $\mathbf{z}$ )



argmax( $\mathbf{z}$ )



(Same  $\mathbf{z} = [1.0716, -1.1221, -0.3288, 0.3368, 0.0425]$ )

- Argmax is an extreme case of sparsity, but it is **discontinuous**.
- Is there a **sparse** and **differentiable** alternative?



# Sparsemax (Martins and Astudillo, 2016, ICML)

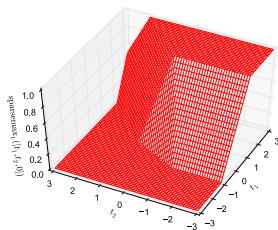
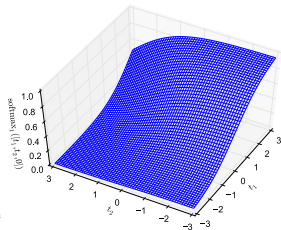
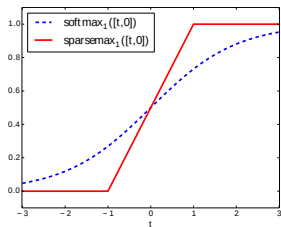
Euclidean projection of  $\mathbf{z}$  onto the probability simplex  $\Delta$ :

$$\begin{aligned}\text{sparsemax}(\mathbf{z}) &:= \arg \min_{\mathbf{p} \in \Delta} \|\mathbf{p} - \mathbf{z}\|^2 \\ &= \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^\top \mathbf{p} - \frac{1}{2} \|\mathbf{p}\|^2.\end{aligned}$$

- Likely to hit the boundary of the simplex, in which case  $\text{sparsemax}(\mathbf{z})$  becomes sparse (hence the name)
- End-to-end differentiable
- Forward pass:  $O(K \log K)$  or  $O(K)$ , (almost) as fast as softmax
- Backprop: sublinear, **better than softmax!**

# Sparsemax in 2D and 3D

(Martins and Astudillo, 2016, ICML)



- Sparsemax is piecewise linear, but asymptotically similar to softmax.

For convex  $\Omega$ , define the  $\Omega$ -regularized argmax transformation:

$$\operatorname{argmax}_{\Omega}(\mathbf{z}) := \operatorname{argmax}_{\mathbf{p} \in \Delta} \mathbf{z}^{\top} \mathbf{p} - \Omega(\mathbf{p})$$

- **Argmax** corresponds to **no regularization**,  $\Omega \equiv 0$
- **Softmax** amounts to **entropic regularization**,  $\Omega(\mathbf{p}) = \sum_{i=1}^K p_i \log p_i$
- **Sparsemax** amounts to  $\ell_2$ -regularization,  $\Omega(\mathbf{p}) = \frac{1}{2} \|\mathbf{p}\|^2$

Is there something in-between?

# Entmax (Peters et al., 2019a, ACL)

Parametrized by  $\alpha \geq 0$ :

$$\Omega_{\alpha}(\mathbf{p}) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \left( \sum_{i=1}^K p_i^{\alpha} - 1 \right) & \text{if } \alpha \neq 1 \\ \sum_{i=1}^K p_i \log p_i & \text{if } \alpha = 1. \end{cases}$$



Related to **Tsallis generalized entropies** (Tsallis, 1988).

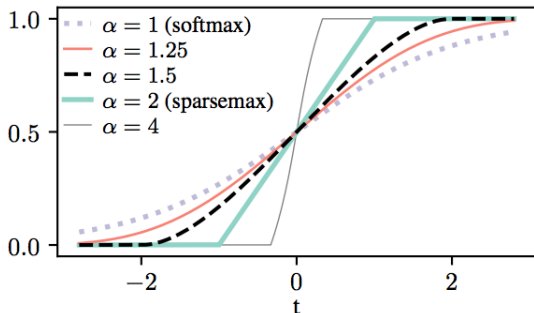
- **Argmax** corresponds to  $\alpha \rightarrow \infty$
- **Softmax** amounts to  $\alpha \rightarrow 1$
- **Sparsemax** amounts to  $\alpha = 2$ .



**Key result:** always sparse for  $\alpha > 1$ , sparsity increases with  $\alpha$

- Forward pass for general  $\alpha$  can be done with a bisection algorithm
- Backward pass runs in sublinear time.

## Entmax in 2D (Peters et al., 2019a, ACL)



$\alpha = 1.5$  is a sweet spot!

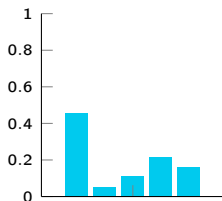
- Efficient exact algorithm (nearly as fast as softmax), smooth, and good empirical performance.

Pytorch code: <https://github.com/deep-spin/entmax>

# Sparse Transformations (Peters et al., 2019a, ACL)

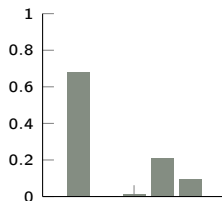
$\alpha = 1$

softmax( $\mathbf{z}$ )



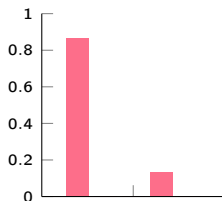
$\alpha = 1.5$

1.5-entmax( $\mathbf{z}$ )



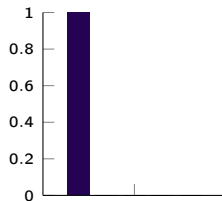
$\alpha = 2$

sparsemax( $\mathbf{z}$ )



$\alpha = \infty$

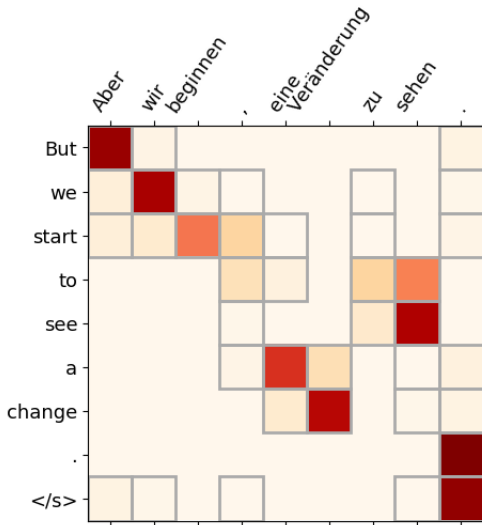
argmax( $\mathbf{z}$ )



(Same  $\mathbf{z} = [1.0716, -1.1221, -0.3288, 0.3368, 0.0425]$ )

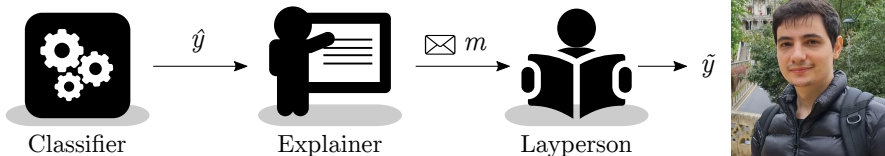
# Example: Sparse Attention for Machine Translation

- **Selects** source words when generating a target word (sparse alignments)
- Better interpretability
- Can also model fertility: **constrained** sparsemax (Malaviya et al., 2018, ACL)
- Can also learn  $\alpha$  (**adaptively sparse** transformers): (Correia et al., 2019, EMNLP)



# Example: Sparse Attention for Explainability

(Treviso and Martins, 2020, BlackboxNLP)



- A classifier makes a prediction
- An “explainer” (embedded or not in the classifier) generates a sparse **message** that explains the classifier’s decision
- The layperson receives the message and tries to guess the classifier’s prediction (also called *simulatability*, *forward simulation/prediction*)
- **Communication success rate**: how often the two predictions match?
- Follow-up: Scaffold Maximizing Training (Fernandes et al., 2022, NeurIPS)



# LP-SparseMAP (Niculae et al., 2018, ICML) (Niculae and Martins, 2020, ICML)

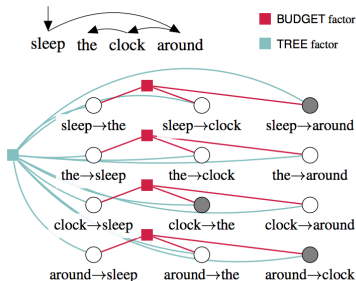
Generalizes sparsemax to **structures**.

Works both as output and hidden layer.

Can handle logic variables and constraints through a **factor graph**.

Returns **sparse** and **differentiable** combination of structures.

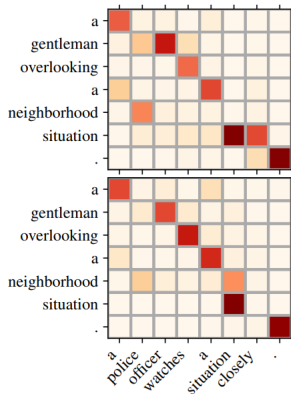
Efficient forward/backprop (requires only a MAP oracle).



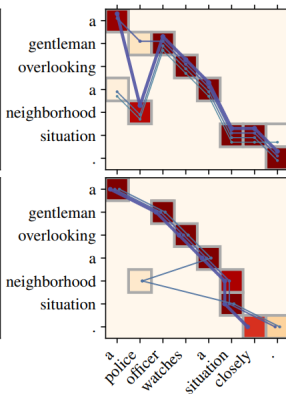
```
fg = TorchFactorGraph()
u = fg.variable_from(arc_scores)
fg.add(DepTree(u))
for k in range(n):
    fg.add(Budget(u[:, k], budget=5))
fg.solve()
```

# Example: Latent Structured Alignments in SNLI

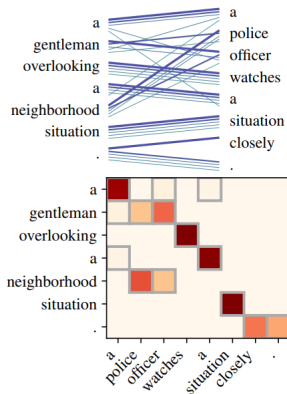
(Niculae et al., 2018)



(a) softmax



(b) sequence

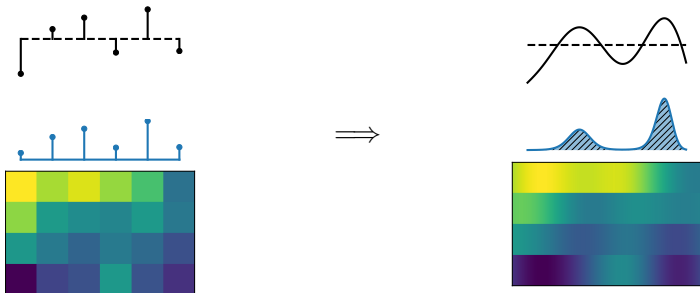
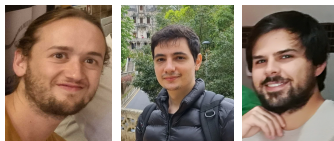


(c) matching

# Sparse and Continuous Attention

(Martins et al., 2020a, 2022a, NeurIPS, JMLR)

- So far: attention over a **finite set** (words, pixel regions, etc.)
- We generalize attention to *arbitrary sets*, possibly continuous.
- Applications: VQA; long-range  $\infty$ -former (Martins et al., 2022b, ACL)



# From Discrete to Continuous Attention

(Martins et al., 2020a, NeurIPS)

(Bahdanau et al., 2015, ICLR)

Finite set  $S = \{1, \dots, K\}$

# From Discrete to Continuous Attention

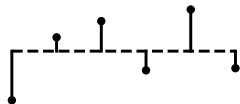
(Martins et al., 2020a, NeurIPS)

(Bahdanau et al., 2015, ICLR)

Finite set  $S = \{1, \dots, K\}$

Three ingredients:

- Score vector  $\mathbf{z} \in \mathbb{R}^K$



# From Discrete to Continuous Attention

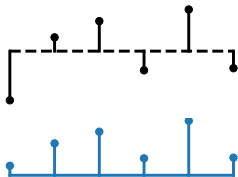
(Martins et al., 2020a, NeurIPS)

(Bahdanau et al., 2015, ICLR)

Finite set  $S = \{1, \dots, K\}$

Three ingredients:

- Score vector  $\mathbf{z} \in \mathbb{R}^K$
- Transformation from  $\mathbf{z}$  to probability vector  $\mathbf{p} \in \Delta^K$



# From Discrete to Continuous Attention

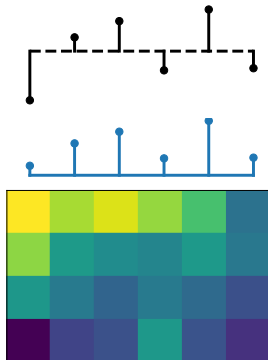
(Martins et al., 2020a, NeurIPS)

(Bahdanau et al., 2015, ICLR)

Finite set  $S = \{1, \dots, K\}$

Three ingredients:

- Score vector  $\mathbf{z} \in \mathbb{R}^K$
- Transformation from  $\mathbf{z}$  to probability vector  $\mathbf{p} \in \Delta^K$
- Value matrix  $V \in \mathbb{R}^{D \times K}$



# From Discrete to Continuous Attention

(Martins et al., 2020a, NeurIPS)

(Bahdanau et al., 2015, ICLR)

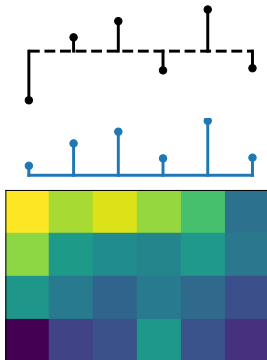
Finite set  $S = \{1, \dots, K\}$

Three ingredients:

- Score vector  $\mathbf{z} \in \mathbb{R}^K$
- Transformation from  $\mathbf{z}$  to probability vector  $\mathbf{p} \in \Delta^K$
- Value matrix  $V \in \mathbb{R}^{D \times K}$

Output:

- **Weighted average**  $V\mathbf{p} \in \mathbb{R}^D$





# From Discrete to Continuous Attention

(Martins et al., 2020a, NeurIPS)

Our work:

Measure space  $S$  (e.g. continuous)

Three ingredients:

- Score vector  $\mathbf{z} \in \mathbb{R}^K$
- Transformation from  $\mathbf{z}$  to probability vector  $\mathbf{p} \in \Delta^K$
- Value matrix  $V \in \mathbb{R}^{D \times K}$

Output:

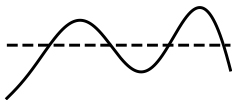
- Weighted average  $V\mathbf{p} \in \mathbb{R}^D$

# From Discrete to Continuous Attention

(Martins et al., 2020a, NeurIPS)

Our work:

Measure space  $S$  (e.g. continuous)



Three ingredients:

- *Score function*  $f : S \rightarrow \mathbb{R}$
- Transformation from  $\mathbf{z}$  to probability vector  $\mathbf{p} \in \Delta^K$
- Value matrix  $V \in \mathbb{R}^{D \times K}$

Output:

- Weighted average  $V\mathbf{p} \in \mathbb{R}^D$

# From Discrete to Continuous Attention

(Martins et al., 2020a, NeurIPS)

Our work:

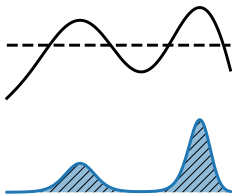
Measure space  $S$  (e.g. continuous)

Three ingredients:

- Score function  $f : S \rightarrow \mathbb{R}$
- Transformation from  $f$  to density  $p : S \rightarrow \mathbb{R}_+$ ,  $\int_S p = 1$
- Value matrix  $V \in \mathbb{R}^{D \times K}$

Output:

- Weighted average  $V\mathbf{p} \in \mathbb{R}^D$



# From Discrete to Continuous Attention

(Martins et al., 2020a, NeurIPS)

Our work:

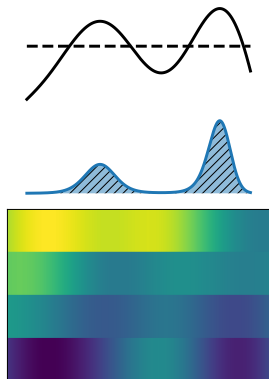
Measure space  $S$  (e.g. continuous)

Three ingredients:

- Score function  $f : S \rightarrow \mathbb{R}$
- Transformation from  $f$  to density  $p : S \rightarrow \mathbb{R}_+$ ,  $\int_S p = 1$
- Value function  $V : S \rightarrow \mathbb{R}^D$

Output:

- Weighted average  $V \mathbf{p} \in \mathbb{R}^D$



# From Discrete to Continuous Attention

(Martins et al., 2020a, NeurIPS)

Our work:

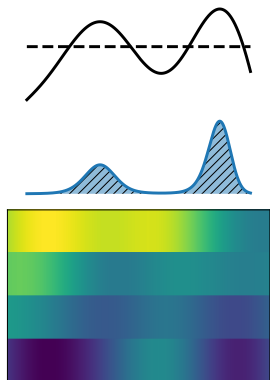
Measure space  $S$  (e.g. continuous)

Three ingredients:

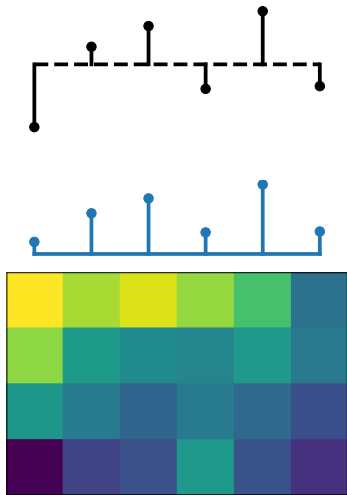
- Score function  $f : S \rightarrow \mathbb{R}$
- Transformation from  $f$  to density  $p : S \rightarrow \mathbb{R}_+$ ,  $\int_S p = 1$
- Value function  $V : S \rightarrow \mathbb{R}^D$

Output:

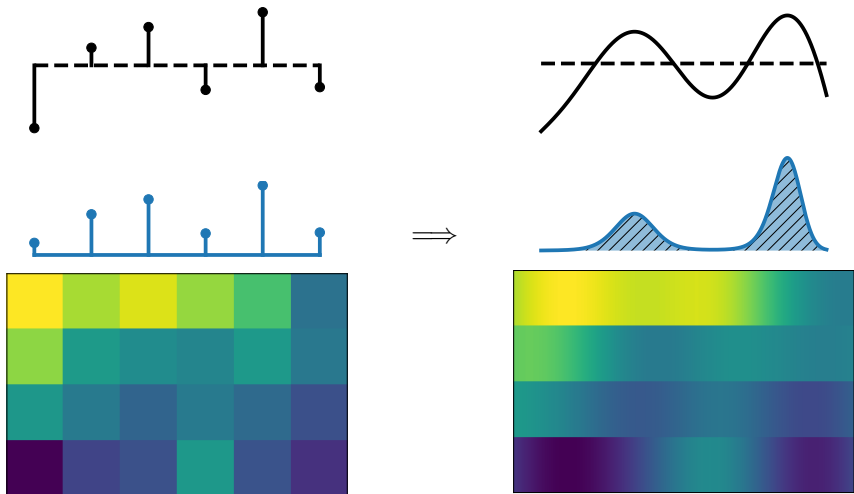
- $\mathbb{E}_p[V(t)] = \int_S p(t)V(t) \in \mathbb{R}^D$



# From Discrete to Continuous Attention



# From Discrete to Continuous Attention

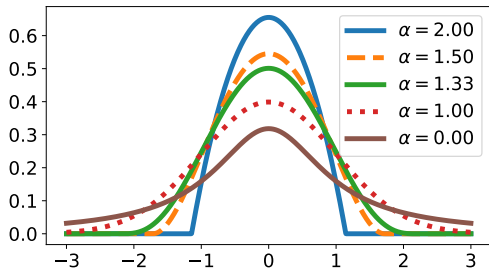


# Sparse and Continuous

How to generalize the concept of **sparsity** to non-finite (e.g. continuous) domains  $S$ ?

- A density with base measure  $\mu$  is **sparse** iff  $\mu(S \setminus \text{supp}(p)) > 0$ .

Examples with  $S = \mathbb{R}$ :  **$q$ -Gaussians for  $\alpha = 2 - q > 1$**



$q$	$\alpha$	
0	2	Epanechnikov
$2/3$	$4/3$	triweight
$1/2$	$3/2$	biweight
1	1	Gaussian
2	0	Cauchy

These can be generalized to  $S = \mathbb{R}^K$  (later)



# $\Omega$ -Regularized Prediction Map ( $\Omega$ -RPM)

Transforms **score function**  $f$  into **probability density**  $p \equiv \hat{p}_\Omega[f]$ :

$$\hat{p}_\Omega[f] = \operatorname{argmax}_p \mathbb{E}_p[f(t)] - \Omega(p), \quad \Omega \text{ convex regularizer.}$$

# $\Omega$ -Regularized Prediction Map ( $\Omega$ -RPM)

Transforms **score function**  $f$  into **probability density**  $p \equiv \hat{p}_\Omega[f]$ :

$$\hat{p}_\Omega[f] = \operatorname{argmax}_p \mathbb{E}_p[f(t)] - \Omega(p), \quad \Omega \text{ convex regularizer.}$$

$-\Omega_\alpha$  **Tsallis  $\alpha$ -entropy**  $\implies$   $\alpha$ -**entmax** (deformed exponential family):

$$\hat{p}_{\Omega_\alpha}[f](t) = \begin{cases} \exp(f(t) - \tau) & \text{if } \alpha = 1 \\ (1 + (\alpha - 1)(f(t) - \tau))_+^{\frac{1}{\alpha-1}} & \text{if } \alpha \neq 1 \end{cases}$$

This is the  $q$ -**exponential function**, with  $q = 2 - \alpha$ .

**Particular cases:** (continuous) softmax ( $\alpha = 1$ ) and sparsemax ( $\alpha = 2$ ).

Blondel et al. (2020a, JMLR), Martins and Astudillo (2016, ICML), Peters et al. (2019b, ACL)

Tsallis (1988, "Possible Generalization of Boltzmann-Gibbs Statistics", J. of Stat. Phys.)

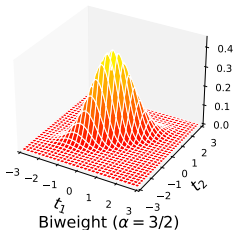
# Example: Multivariate $q$ -Gaussians

Quadratic score function:  $f(t) = -\frac{1}{2}(t - \mu)^\top \Sigma^{-1}(t - \mu)$

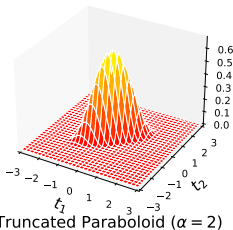
Nice properties:

- They're instances of **elliptical distributions**
- Efficient to sample from (Beta + sphere trick)
- FY loss (Bregman) and Wasserstein distances computable in closed form
- For  $\alpha = 2 - q = \frac{k}{k-1}$  with  $k \in \mathbb{Z}$ , attention and its gradient have closed form for the 1D case.

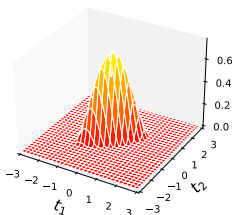
Gaussian ( $\alpha = 1$ )



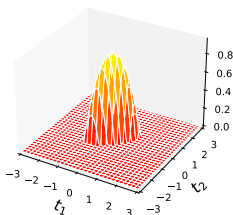
Triweight ( $\alpha = 4/3$ )



Biweight ( $\alpha = 3/2$ )



Truncated Paraboloid ( $\alpha = 2$ )



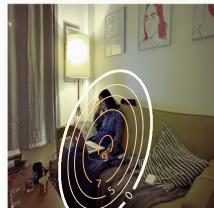
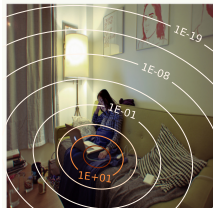
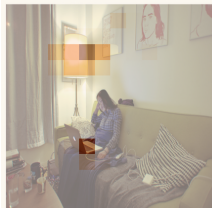
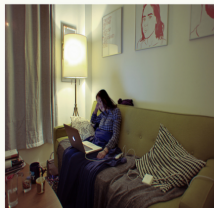
# Example: Visual Question Answering

What is the woman looking at?

tv

computer

computer

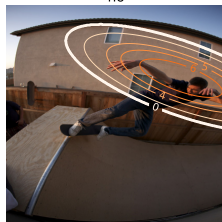


Is the man wearing a hat?

yes

no

no



(original image)

(discrete attention)

(continuous softmax)

(continuous sparsemax)

# Outline

- 1 Sparse Transformations
- 2 Fenchel-Young Losses**
- 3 Sparse Hopfield Networks
- 4 Mixed Distributions
- 5 Conclusions

# Loss Functions

Define the training objective to fit the model to the data.

Assess compatibility between:

- **groundtruth**  $\mathbf{y}$  (e.g. response variable)
- **model output**  $\mathbf{z}$  (e.g. last layer of a neural network).

Examples:

- **squared loss** in regression ( $\mathbf{y} \in \mathbb{R}^K, \mathbf{z} \in \mathbb{R}^K$ ):

$$L(\mathbf{z}, \mathbf{y}) = \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|^2$$

- **cross-entropy loss** in logistic regression ( $\mathbf{y} \in \Delta, \mathbf{z} \in \mathbb{R}^K$ ):

$$\begin{aligned} L(\mathbf{z}, \mathbf{y}) &= - \sum_i y_i \log[\text{softmax}(\mathbf{z})]_i \\ &= - \sum_i y_i z_i + \log \sum_i \exp(z_i). \end{aligned}$$

# Entmax Losses

- Entmax can also be used as a loss in the **output layer** (to replace logistic/cross-entropy loss)
- However, not expressed as a log-likelihood (which could lead to  $\log(0)$  problems due to sparsity)
- Instead, we build a entmax loss inspired by **Fenchel-Young losses**.

## Recap: $\Omega$ -Regularized Argmax (Nicolae and Blondel, 2017, NeurIPS)

For convex  $\Omega$ , define the  $\Omega$ -regularized argmax transformation:

$$\operatorname{argmax}_{\Omega}(\mathbf{z}) := \operatorname{arg\,max}_{\mathbf{p} \in \Delta} \mathbf{z}^{\top} \mathbf{p} - \Omega(\mathbf{p})$$

- **Argmax** corresponds to **no regularization**,  $\Omega \equiv 0$
- **Softmax** amounts to **entropic regularization**,  $\Omega(\mathbf{p}) = \sum_{i=1}^K p_i \log p_i$
- **Sparsemax** amounts to  $\ell_2$ -regularization,  $\Omega(\mathbf{p}) = \frac{1}{2} \|\mathbf{p}\|^2$

All these are particular cases of  $\alpha$ -entmax (Peters et al., 2019a, ACL).



# Fenchel-Young Losses (Blondel et al., 2020b, JMLR)

Assess compatibility between **groundtruth**  $\mathbf{y} \in \Delta$  and **scores**  $\mathbf{z} \in \mathbb{R}^K$

Convex conjugate  $\Omega^*(\mathbf{z}) := \max_{\mathbf{p} \in \Delta} \mathbf{z}^\top \mathbf{p} - \Omega(\mathbf{p})$

$$L_\Omega(\mathbf{z}, \mathbf{y}) := \Omega^*(\mathbf{z}) + \Omega(\mathbf{y}) - \mathbf{z}^\top \mathbf{y}$$

Recover **cross-entropy loss**:  $\Omega(\mathbf{p}) = \sum_i p_i \log p_i \Rightarrow \Omega^*(\mathbf{z}) = \log \sum_i \exp(z_i)$ .

# Fenchel-Young Losses (Blondel et al., 2020b, JMLR)

Assess compatibility between **groundtruth**  $\mathbf{y} \in \Delta$  and **scores**  $\mathbf{z} \in \mathbb{R}^K$

Convex conjugate  $\Omega^*(\mathbf{z}) := \max_{\mathbf{p} \in \Delta} \mathbf{z}^\top \mathbf{p} - \Omega(\mathbf{p})$

$$L_\Omega(\mathbf{z}, \mathbf{y}) := \Omega^*(\mathbf{z}) + \Omega(\mathbf{y}) - \mathbf{z}^\top \mathbf{y}$$

Recover **cross-entropy loss**:  $\Omega(\mathbf{p}) = \sum_i p_i \log p_i \Rightarrow \Omega^*(\mathbf{z}) = \log \sum_i \exp(z_i)$ .

**Properties:**

- $L_\Omega(\mathbf{z}, \mathbf{y}) \geq 0$  (automatic from **Fenchel-Young inequality**)
- $L_\Omega(\mathbf{z}, \mathbf{y}) = 0$  iff  $\mathbf{y} = \operatorname{argmax}_\Omega(\mathbf{z})$
- $L_\Omega$  is convex and differentiable with  $\nabla L_\Omega(\mathbf{z}, \mathbf{y}) = \operatorname{argmax}_\Omega(\mathbf{z}) - \mathbf{y}$

Also called “mixed-type Bregman divergences” (Amari, 2016).

## Definition: Loss Margin

Some loss functions (e.g. the **hinge loss** in SVMs) are associated to the concept of **margin**.

A loss function  $L(\mathbf{z}, \mathbf{y})$  has a **margin** if there is finite  $m \geq 0$  such that

$$\forall i \in [K], \quad L(\mathbf{z}, \mathbf{e}_i) = 0 \Leftrightarrow z_i - \max_{j \neq i} z_j \geq m.$$

The smallest such  $m$  is called the margin of  $L$ .

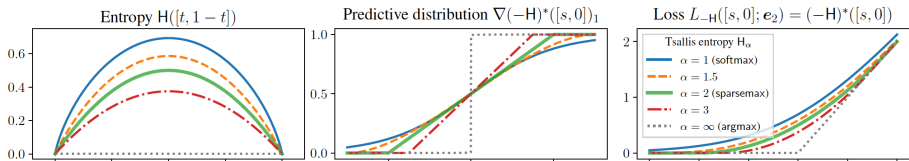
If  $L_\Omega$  is a Fenchel-Young loss, this condition is equivalent to

$$\operatorname{argmax}_{\Omega}(\mathbf{z}) = \mathbf{e}_i.$$

Corollary: cross-entropy loss does **not** have a margin.

# Entmax Transformations and Losses

(Blondel et al., 2020b, JMLR)



- Key result: for all  $\alpha > 1$ ,  $\alpha$ -entmax transformations are **sparse** and lead to losses with **margins**!
- The **margin**  $m$  is related to the **slope** of the entropy in the simplex corners! ( $m = \frac{1}{\alpha-1}$  for entmax losses.)
- See paper for details!

Pytorch code: <https://github.com/deep-spin/entmax>

# Example: Machine Translation

(Peters et al., 2019a, ACL) (Peters and Martins, 2021, NAACL)

<b>This</b>	92.9%	<b>is another</b>	<b>view</b>	49.8%	<b>at</b>	95.7%	<b>the tree of life .</b>
So	5.9%		<b>look</b>	27.1%	<b>on</b>	5.9%	
And	1.3%		glimpse	19.9%	,	1.3%	
Here	<0.1%		kind	2.0%			
			looking	0.9%			
			way	0.2%			
			vision	<0.1%			
			gaze	<0.1%			



(Source: "Dies ist ein weiterer Blick auf den Baum des Lebens.")

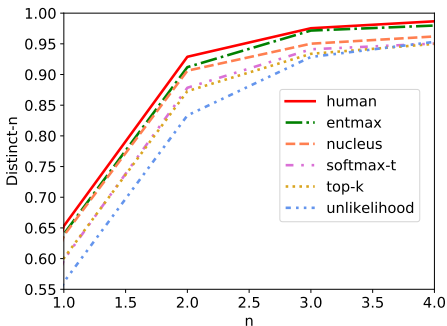
- Only a few words get non-zero probability at each time step
- Auto-completion when several words in a row have probability 1
- Useful for predictive translation.

# Entmax Sampling (Martins et al., 2020b, EMNLP)

Use the entmax loss for training language models.

At test time, **sample** from this sparse distribution.

Better quality with less repetitions than other methods:



# Outline

- 1 Sparse Transformations
- 2 Fenchel-Young Losses
- 3 Sparse Hopfield Networks**
- 4 Mixed Distributions
- 5 Conclusions

# What are Hopfield Networks? (Amari, 1972; Hopfield, 1982)

A model of **recurrent neural network**:

- Named after John Hopfield, an American physicist and neuroscientist.
- Designed for associative memory and associative recall.
- Stores and recalls patterns, making it useful for pattern recognition.
- Inspired by the human brain's associative memory.
  - Our brain stores and retrieves information not through explicit memory addresses but by associating content with memories.
  - **Content-based recall**: in the brain, seeing or hearing a partial cue can trigger the recall of associated memories.



# Hopfield Networks (Amari, 1972; Hopfield, 1982)

A model of **associative memory**:

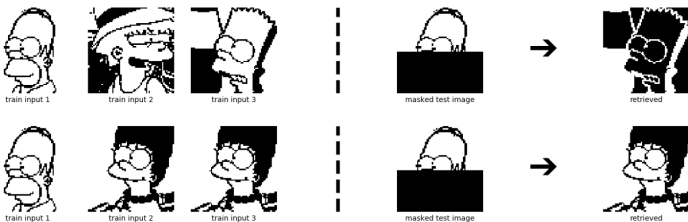
- Memory patterns  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \{\pm 1\}^{N \times D}$ , query  $\mathbf{q} \in \{\pm 1\}^D$
- Energy  $E(\mathbf{q}) = -\frac{1}{2} \|\mathbf{X}\mathbf{q}\|^2$
- Hopfield dynamics  $\mathbf{q}_{t+1} = \text{sign}(\mathbf{X}^\top \mathbf{X}\mathbf{q}_t)$
- Memory patterns are **attractors** (but many spurious attractors)
- Memory capacity is only  $N \lesssim 0.138D = \mathcal{O}(D)$ .

# Hopfield Networks



<https://ml-jku.github.io/hopfield-layers/>

What if we store more than one pattern?



<https://ml-jku.github.io/hopfield-layers/>

# Hopfield Networks

What if we try with even more?



There are alternative energies with much better capacity  
(Krotov and Hopfield, 2016; Demircigil et al., 2017; Ramsauer et al., 2020).

Caveat: needs to store patterns **explicitly**.

# Dense Associative Memories

(Krotov and Hopfield, 2016; Demircigil et al., 2017)

- Krotov et al. proposed a new energy function:

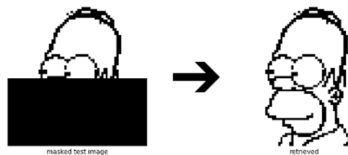
$$E(\mathbf{q}) = -F(\mathbf{X}\mathbf{q}) \qquad F(\mathbf{x}) = \begin{cases} \mathbf{x}^n & \text{if } \mathbf{x} \geq 0 \\ 0 & \text{if } \mathbf{x} < 0 \end{cases}$$

- Hopfield dynamics  $\mathbf{q}_{t+1} = \text{sign}(\mathbf{X}^\top \text{spow}(\mathbf{X}\mathbf{q}_t, n - 1))$
- Demircigil et al. further expanded the energy function:

$$E(\mathbf{q}) = -\exp(\mathbf{X}\mathbf{q})$$

- Hopfield dynamics  $\mathbf{q}_{t+1} = \text{sign}(\mathbf{X}^\top \exp(\mathbf{X}\mathbf{q}_t))$
- Memory capacity is now  $\mathcal{O}(\exp(D))$

# Example



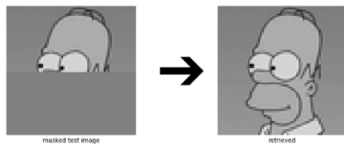
# Modern Hopfield Networks (Ramsauer et al., 2020)

- Operates on **continuous space**,  $\mathbf{X} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{q} \in \mathbb{R}^D$
- Energy:

$$E(\mathbf{q}) = -\beta^{-1} \log \sum_{i=1}^N \exp(\beta \mathbf{x}_i^\top \mathbf{q}) + \frac{1}{2} \|\mathbf{q}\|^2.$$

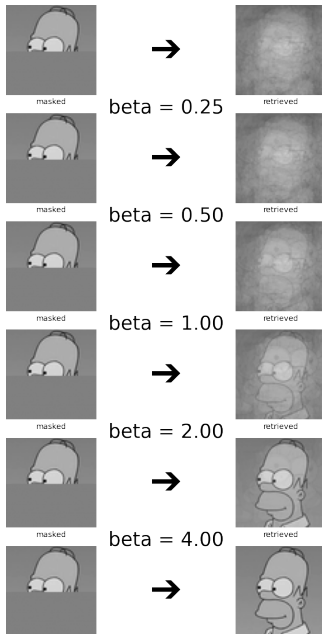
- Hopfield dynamics  $\mathbf{q}_{t+1} = \mathbf{X}^\top \text{softmax}(\beta \mathbf{X} \mathbf{q}_t)$
- Similar to self-attention in transformers!
- Memory patterns are **close** to attractors (but there can be some spurious attractors)
- Memory capacity is  $\mathcal{O}(\exp(D))$  (but retrieval is only **approximate**)

# Modern Hopfield Networks



If some stored patterns are similar to each other, then a **metastable state** near the similar patterns appears.

# Sensitivity to Temperature





# Research Questions

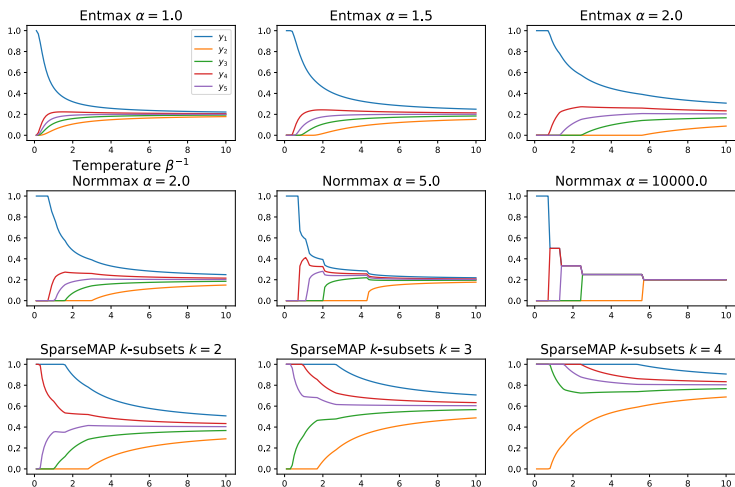
- Can we design high capacity, continuous-space Hopfield networks with **exact** retrieval?
- Can we make them less sensitive to temperature?
- Can we extend them to handle structure?

# Research Questions

- Can we design high capacity, continuous-space Hopfield networks with **exact** retrieval?
- Can we make them less sensitive to temperature?
- Can we extend them to handle structure?

Yes, if we use sparse transformations and Fenchel-Young losses!

# Sparse and Structured Transformations



Regularization path of sparse and structured transformations.  
Shown is  $\operatorname{argmax}_{\Omega}(\beta \mathbf{z})$  as a function of the temperature  $\beta^{-1}$  where  
 $\mathbf{z} = [1.0716, -1.1221, -0.3288, 0.3368, 0.0425]^T$ .

- Hopfield-Fenchel-Young energy, induced by convex  $\Omega$ :

$$\begin{aligned} E(\mathbf{q}) &= -\beta^{-1} \Omega^*(\beta \mathbf{X} \mathbf{q}) + \frac{1}{2} \|\mathbf{q}\|^2 \\ &= -L_{\beta^{-1} \Omega}(\mathbf{X} \mathbf{q}; \mathbf{u}) + \frac{1}{2} \|\mathbf{q} - \mathbf{X}^\top \mathbf{u}\|^2 + \text{const.} \end{aligned}$$

- Includes MHNs as a particular case
- Hopfield dynamics  $\mathbf{q}_{t+1} = \mathbf{X}^\top \operatorname{argmax}_{\Omega}(\beta \mathbf{X} \mathbf{q})$
- If  $\operatorname{argmax}_{\Omega}$  is a sparse transformation, memory patterns are **exactly** attractors (but there can be some spurious attractors)
- Memory capacity is still  $\mathcal{O}(\exp(D))$  (but retrieval can be exact)

# Exact Convergence to Single Pattern

Define separation of pattern  $\mathbf{x}_i$  from data (Ramsauer et al., 2020):

$$\Delta_i = \mathbf{x}_i^\top \mathbf{x}_i - \max_{j \neq i} \mathbf{x}_i^\top \mathbf{x}_j.$$

Assume  $L_\Omega$  is a FY loss with margin  $m$ , and let  $\mathbf{x}_i$  be a memory pattern outside the convex hull of the other memory patterns. Then,

- $\mathbf{x}_i$  is a stationary point of the energy iff  $\Delta_i \geq \frac{1}{m\beta}$ .
- if the initial query satisfies  $\mathbf{q}_0^\top (\mathbf{x}_i - \mathbf{x}_j) \geq \frac{1}{m\beta}$  for all  $j \neq i$ , then the update rule converges to  $\mathbf{x}_i$  exactly in one iteration.
- if the patterns are normalized and  $\Delta_i \geq \frac{1}{m\beta} + 2M\epsilon$ , then any  $\mathbf{q}_0$   $\epsilon$ -close to  $\mathbf{x}_i$  ( $\|\mathbf{q}_0 - \mathbf{x}_i\| \leq \epsilon$ ) will converge to  $\mathbf{x}_i$  in one iteration.

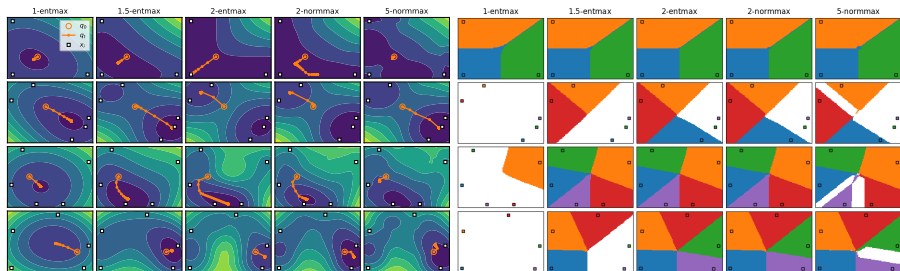
Margins:  $m = 1$  for normmax and  $m = \frac{1}{\alpha-1}$  for  $\alpha$ -entmax.

# Storage Capacity with Exact Retrieval

Assume patterns are randomly placed on the sphere with uniform distribution. Then, with probability  $1 - p$ , the HFY network can store and exactly retrieve  $N = \mathcal{O}(\sqrt{p}\zeta^{\frac{D-1}{2}})$  patterns in one iteration under a  $\epsilon$ -perturbation if

$$\epsilon \leq \frac{M}{2} \left( 1 - \cos \frac{1}{\zeta} \right) - \frac{m}{2\beta M}.$$

# Example: Hopfield Dynamics and Basis of Attraction



As  $\alpha$  increases:

- $\alpha$ -entmax converges more often to a single pattern.
- $\alpha$ -normmax tends to converge towards an attractor which is a uniform average of some patterns.

# Structured Hopfield Networks

Similar guarantees for **structured sparse transformations**:

- Hopfield dynamics  $\mathbf{q}_{t+1} = \mathbf{X}^\top \text{SparseMAP}(\beta \mathbf{X} \mathbf{q}_t)$

Examples of structural constraints:

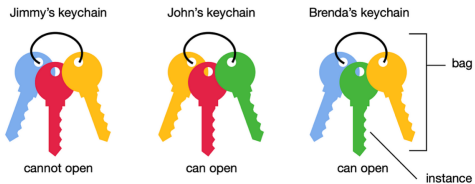
- **$k$ -subsets:**
  - Retrieve subsets of  $k$  patterns, e.g., to take into account a  $k$ -ary relation among patterns or to perform top- $k$  retrieval.
- **sequential  $k$ -subsets:**
  - Promote consecutive memory items to be both (or none) retrieved.

Other structures (trees, graphs, matchings, ...) are possible.



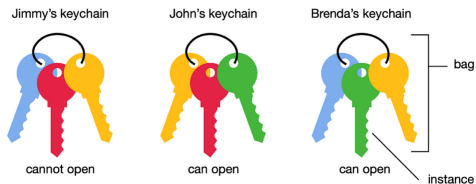
# Example: Multiple Instance Learning

MIL



# Example: Multiple Instance Learning

MIL



$K$ -MIL

At least  $K$  positive instances

# Example: Multiple Instance Learning

Methods	MNIST				MIL benchmarks		
	$K=1$	$K=2$	$K=3$	$K=5$	Fox	Tiger	Elephant
1-entmax (softmax)	98.4 ± 0.2	94.6 ± 0.5	91.1 ± 0.5	89.0 ± 0.3	66.4 ± 2.0	87.1 ± 1.6	92.6 ± 0.6
1.5-entmax	97.6 ± 0.8	96.0 ± 0.9	90.4 ± 1.1	92.4 ± 1.4	66.3 ± 2.0	87.3 ± 1.5	92.4 ± 1.0
2.0-entmax (sparsemax)	97.9 ± 0.2	96.7 ± 0.5	92.9 ± 0.9	91.6 ± 1.0	66.1 ± 0.6	<b>87.7 ± 1.4</b>	91.8 ± 0.6
2.0-normmax	97.9 ± 0.3	96.6 ± 0.6	93.9 ± 0.7	92.4 ± 0.7	66.1 ± 2.5	86.4 ± 0.8	92.4 ± 0.7
5.0-normmax	98.2 ± 0.5	97.2 ± 0.3	95.8 ± 0.4	93.2 ± 0.5	66.4 ± 2.3	85.5 ± 0.6	93.0 ± 0.7
SparseMAP, $k = 2$	<b>98.7 ± 0.3</b>	<b>97.9 ± 0.2</b>	92.2 ± 0.5	92.2 ± 0.5	66.8 ± 2.7	85.4 ± 0.6	<b>93.1 ± 1.0</b>
SparseMAP, $k = 3$	97.2 ± 0.4	96.2 ± 1.2	<b>97.3 ± 0.4</b>	92.0 ± 0.6	<b>68.0 ± 2.6</b>	85.0 ± 0.5	90.8 ± 0.7
SparseMAP, $k = 5$	97.5 ± 0.9	97.4 ± 0.5	94.5 ± 0.8	<b>96.2 ± 1.1</b>	65.5 ± 1.9	79.8 ± 1.2	89.9 ± 0.6

- Normmax consistently performs well across datasets, likely because of its adaptability to near-uniform metastable states of varying sizes.
- When  $K > 1$ , the  $k$ -subsets method works best with  $k = K$ .

See Santos et al. (2024) for more experiments (e.g. text rationalization).

# Outline

- 1 Sparse Transformations
- 2 Fenchel-Young Losses
- 3 Sparse Hopfield Networks
- 4 Mixed Distributions**
- 5 Conclusions

# Mixed Distributions (Farinhas et al., 2022, ICLR)

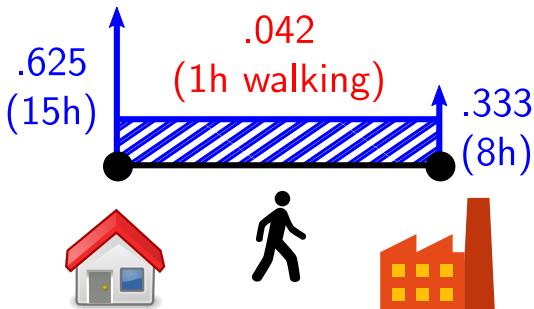


- We saw how to obtain sparse probability distributions.
- How can we use them to bridge the gap between *discrete* and *continuous* domains?
- We'll see how next.

## Back to John's Life

John splits his day as follows: he works 8h/day, and stays home 15h/day.

He is **in transit 1h/day** to commute to work and back.



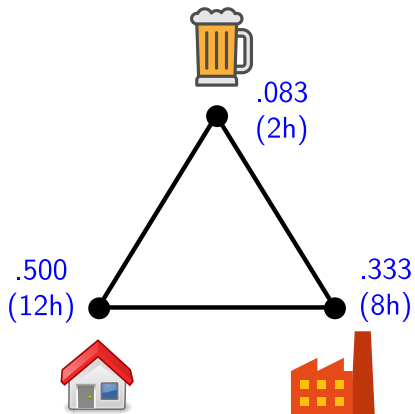
## Back to John's Life

John splits his day as follows: he works 8h/day, and stays home 15h/day.  
He is **in transit 1h/day** to commute to work and back.



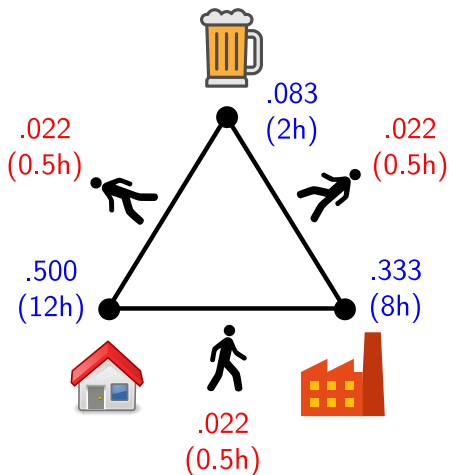
That's a sad life!

After work, John spends 2h in the pub with friends.

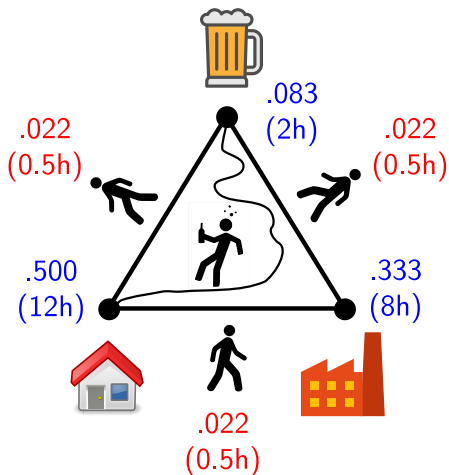




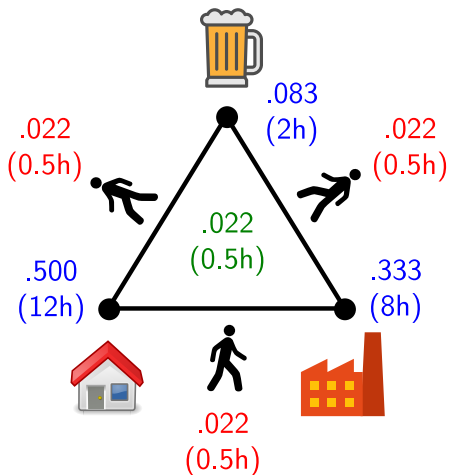
After work, John spends 2h in the pub with friends.



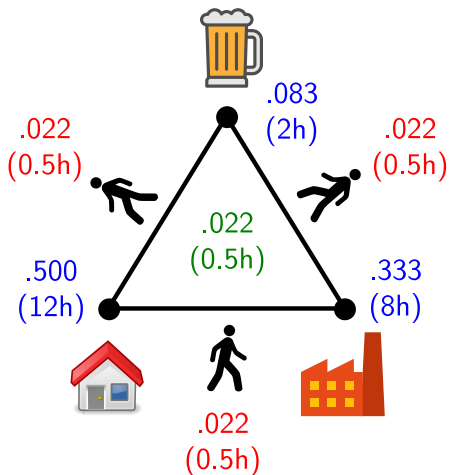
After work, John spends 2h in the pub with friends.



After work, John spends 2h in the pub with friends.



After work, John spends 2h in the pub with friends.



We need a way to represent this probability mass in vertices, edges, face.

# Densities over the simplex $\triangle$

We denote by  $\text{ri}(\triangle)$  the **relative interior** of  $\triangle$ .

Common densities on the simplex:

- Dirichlet distribution
- Logistic-Normal (a.k.a. Gaussian-Softmax)
- Concrete (a.k.a. Gumbel-Softmax)

None of these place any probability mass on the boundary  $\triangle \setminus \text{ri}(\triangle)$ .

# Truncated Densities in the Binary Case ( $K = 2$ )

When  $K = 2$ , the simplex is isomorphic to unit interval,  $\Delta_1 \simeq [0, 1]$ .

A point in  $\Delta_1$  can be represented as  $\mathbf{y} = [y, 1 - y]$ .

**Truncated** densities have been proposed for  $K = 2$ :

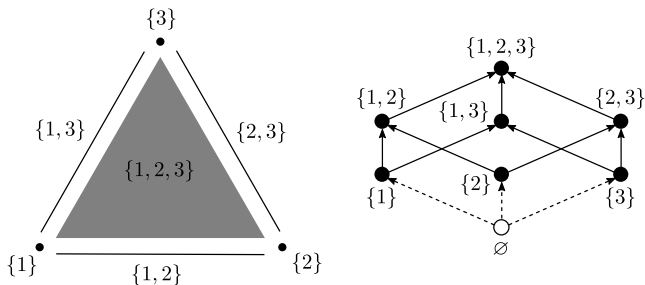
- Binary Hard Concrete (Louizos et al., 2018)
- Rectified Gaussian (Hinton and Ghahramani, 1997; Palmer et al., 2017)

We generalize them for  $K > 2$ .

# Our Approach: Face Stratification

How to extend these “truncated densities” to  $K > 2$ ?

Our solution relies on the **face lattice** of the simplex:



0-faces are vertices, 1-faces are edges, ..., the  $(K - 1)$ -face is  $\triangle$  itself.

We define a **direct sum measure** on the stratified  $\triangle$  and define probability densities w.r.t. this base measure.

# Mixed Random Variables (Farinhas et al., 2022, ICLR)

Discrete RVs assign probability only to **0-faces** (vertices of  $\triangle$ ).

Continuous RVs assign probability only to the **maximal face** ( $\text{ri}(\triangle)$ ).

**Mixed RVs generalize both:** can assign probability to **all faces** of  $\triangle$ .



# Mixed Random Variables (Farinhas et al., 2022, ICLR)

Discrete RVs assign probability only to **0-faces** (vertices of  $\Delta$ ).

Continuous RVs assign probability only to the **maximal face** ( $\text{ri}(\Delta)$ ).

**Mixed RVs generalize both:** can assign probability to **all faces** of  $\Delta$ .

They can be defined via:

- Their **face probability mass function**  $P_F(f) = \Pr\{\mathbf{y} \in \text{ri}(f)\}$ ,  $f \in \mathcal{F}$ .
- Their **face-conditional densities**  $p_{Y|F}(\mathbf{y} | f)$ , for  $f \in \mathcal{F}$ ,  $\mathbf{y} \in \text{ri}(f)$ .

The probability of a set  $A \subseteq \Delta$  is given by:

$$\Pr\{\mathbf{y} \in A\} = \sum_{f \in \mathcal{F}} P_F(f) \int_{A \cap \text{ri}(f)} p_{Y|F}(\mathbf{y} | f).$$

# Extrinsic vs Intrinsic (Farinhas et al., 2022, ICLR)

Two ways of characterizing mixed RVs:

- **Extrinsic characterizat**on: start with a distribution over  $\mathbb{R}^K$  and then apply a deterministic transformation to project it to  $\triangle$
- **Intrinsic characterizat**on: specify a mixture of distributions directly over the faces of  $\triangle$ , by specifying  $P_F$  and  $p_{Y|F}$  for each  $f \in \mathcal{F}$

# $K$ -D Hard Concrete (Farinhas et al., 2022, ICLR)

Uses an **extrinsic characterization**, via “stretch-and-project.”

Generative story:

$$Y \sim \text{HardConcrete}(\mathbf{z}, \lambda, \tau) \quad \Leftrightarrow \quad \begin{aligned} Y' &\sim \text{Concrete}(\mathbf{z}, \lambda) \\ Y &= \text{sparsemax}(\tau Y'), \quad \text{with } \tau \geq 1. \end{aligned}$$

- Recovers the binary Hard Concrete for  $K = 2$
- The larger  $\tau$ , the higher the tendency to hit a non-maximal face of the simplex and induce sparsity.

# Gaussian-Sparsemax (Farinhas et al., 2022, ICLR)

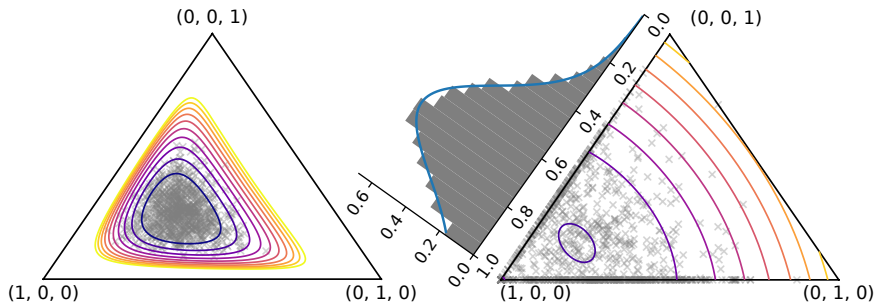
Uses an **extrinsic characterization**, by sampling from a Gaussian and projecting.

Generative story:

$$Y \sim \text{GaussianSparsemax}(\mathbf{z}, \Sigma) \Leftrightarrow \begin{aligned} N &\sim \mathcal{N}(0, I) \\ Y &= \text{sparsemax}(\mathbf{z} + \Sigma^{1/2} N). \end{aligned}$$

- Sparsemax counterpart of the Logistic-Normal.
- Can assign nonzero probability mass to the boundary of the simplex.
- When  $K = 2$ , we recover the double-sided rectified Gaussian.
- For  $K > 2$ , an intrinsic representation can be expressed via the orthant probability of multivariate Gaussians.

# Logistic-Normal vs Gaussian-Sparsemax (Farinhas et al., 2022, ICLR)



Logistic-Normal (left) assigns zero probability to all faces but  $\text{ri}(\Delta)$

Gaussian-Sparsemax (right) is a **mixed distribution**: it assigns probability to the *full* simplex, including its boundary.

# Information Theory of Mixed Random Variables

(Farinhas et al., 2022, ICLR)

“Direct sum” entropy using  $\mu^\oplus$  as the base measure:

$$\begin{aligned} H^\oplus(Y) &:= H(F) + H(Y | F) \\ &= \underbrace{-\sum_{f \in \mathcal{F}} P_F(f) \log P_F(f)}_{\text{discrete entropy}} + \sum_{f \in \mathcal{F}} P_F(f) \underbrace{\left( -\int_{\mathbf{y}} p_{Y|F}(\mathbf{y} | f) \log p_{Y|F}(\mathbf{y} | f) \right)}_{\text{differential entropy}}. \end{aligned}$$

- Average length of the optimal code where  $f$  must be encoded **losslessly** and where  $\mathbf{y}|_f$  has a predefined bit precision  $N$
- Max-ent is written as a generalized Laguerre polynomial (see paper)
  - e.g.  $\log_2(2 + 2^N)$  for  $K = 2$  (vs.  $\log_2(2) = 1$  in the purely discrete case)
- KL divergence and mutual information defined similarly.

# Experiment: Emergent Communication

The first agent needs to communicate a **code** to the second agent that represents a given image.

Given the code, the second agent needs to identify the correct image among **16** possibilities. (Random guess is  $1/16 = 6.25\%$ .)

Success average and standard error over 10 runs:

Method	Success (%)	Nonzeros ↓
Gumbel-Softmax	78.84 ±8.07	256
Gumbel-Softmax ST	49.96 ±9.51	1
<i>K</i> -D Hard Concrete	76.07 ±7.76	21.43 ±17.56
Gaussian-Sparsemax	<b>80.88</b> ±0.50	1.57 ±0.02

(See paper for more experiments with VAEs on FashionMNIST and MNIST.)

# Outline

- 1 Sparse Transformations
- 2 Fenchel-Young Losses
- 3 Sparse Hopfield Networks
- 4 Mixed Distributions
- 5 Conclusions**



# Conclusions

- Transformations from real numbers to distributions are ubiquitous
- We introduced new transformations that handle **sparsity**, **constraints**, and **structure**
- All are differentiable and their gradients are efficient to compute
- Can be used as hidden layers or as output layers (Fenchel-Young losses)
- Mixed distributions are in-between the discrete and continuous worlds
- Sparse communication potentially useful as a path for explainability.

# Thank you!



# References I

- Amari, S.-I. (1972). Learning patterns and pattern sequences by self-organizing nets of threshold elements. *IEEE Transactions on computers*, 100(11):1197–1206.
- Amari, S.-i. (2016). *Information geometry and its applications*, volume 194. Springer.
- Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In *Proc. ICLR*.
- Blondel, M., Martins, A. F., and Niculae, V. (2020a). Learning with fenchel-young losses. *Journal of Machine Learning Research*, 21(35):1–69.
- Blondel, M., Martins, A. F. T., and Niculae, V. (2020b). Learning with fenchel-young losses. *Journal of Machine Learning Research*, 21(35):1–69.
- Correia, G., Niculae, V., and Martins, A. F. T. (2019). Adaptively sparse transformers. In *Proceedings of the Empirical Methods for Natural Language Processing*.
- Demircigil, M., Heusel, J., Löwe, M., Upgang, S., and Vermet, F. (2017). On a model of associative memory with huge storage capacity. *Journal of Statistical Physics*, 168:288–299.
- Farinhas, A., Aziz, W., Niculae, V., and Martins, A. F. (2022). Sparse communication via mixed distributions. In *Proc. of ICLR*.
- Fernandes, P., Treviso, M., Pruthi, D., Martins, A. F., and Neubig, G. (2022). Learning to scaffold: Optimizing model explanations for teaching. *arXiv preprint arXiv:2204.10810*.
- Hinton, G. E. and Ghahramani, Z. (1997). Generative models for discovering sparse distributed representations. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 352(1358):1177–1190.

## References II

- Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558.
- Hu, J. Y.-C., Yang, D., Wu, D., Xu, C., Chen, B.-Y., and Liu, H. (2023). On sparse modern hopfield model. In *Advances in Neural Information Processing Systems*.
- Krotov, D. and Hopfield, J. J. (2016). Dense associative memory for pattern recognition. *Advances in neural information processing systems*, 29.
- Louizos, C., Welling, M., and Kingma, D. P. (2018). Learning sparse neural networks through  $l_0$  regularization. In *International Conference on Learning Representations*.
- Malaviya, C., Ferreira, P., and Martins, A. F. T. (2018). Sparse and Constrained Attention for Neural Machine Translation. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Martins, A. and Astudillo, R. (2016). From softmax to sparsemax: A sparse model of attention and multi-label classification. In Balcan, M. F. and Weinberger, K. Q., editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1614–1623, New York, New York, USA. PMLR.
- Martins, A., Farinhas, A., Treviso, M., Niculae, V., Aguiar, P., and Figueiredo, M. (2020a). Sparse and continuous attention mechanisms. *Advances in Neural Information Processing Systems*, 33.
- Martins, A. F. T., Treviso, M., Farinhas, A., Aguiar, P. M., Figueiredo, M. A., Blondel, M., and Niculae, V. (2022a). Sparse continuous distributions and Fenchel-Young losses. *Journal of Machine Learning Research (to appear)*, 23(257):1–74.
- Martins, P. H., Marinho, Z., and Martins, A. F. (2020b). Sparse text generation. In *Empirical Methods for Natural Language Processing*.

## References III

- Martins, P. H., Marinho, Z., and Martins, A. F. T. (2022b).  $\infty$ -former: Infinite memory transformer. In *Proc. of NAACL-HLT*.
- Niculae, V. and Blondel, M. (2017). A regularized framework for sparse and structured neural attention. *arXiv preprint arXiv:1705.07704*.
- Niculae, V. and Martins, A. F. (2020). Lp-sparsemap: Differentiable relaxed optimization for sparse structured prediction. In *International Conference on Machine Learning*.
- Niculae, V., Martins, A. F. T., Blondel, M., and Cardie, C. (2018). SparseMAP: Differentiable Sparse Structured Inference. In *Proc. of the International Conference on Machine Learning*.
- Palmer, A. W., Hill, A. J., and Scheduling, S. J. (2017). Methods for stochastic collection and replenishment (scar) optimisation for persistent autonomy. *Robotics and Autonomous Systems*, 87:51–65.
- Peters, B. and Martins, A. F. (2021). Smoothing and shrinking the sparse seq2seq search space. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2642–2654.
- Peters, B., Niculae, V., and Martins, A. F. T. (2019a). Sparse sequence-to-sequence models. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*.
- Peters, B., Niculae, V., and Martins, A. F. T. (2019b). Sparse sequence-to-sequence models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1504–1519, Florence, Italy. Association for Computational Linguistics.
- Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Adler, T., Gruber, L., Holzleitner, M., Pavlović, M., Sandve, G. K., et al. (2020). Hopfield networks is all you need. *arXiv preprint arXiv:2008.02217*.

# References IV

- Santos, S., Niculae, V., McNamee, D., and Martins, A. F. (2024). Sparse and structured hopfield networks. *arXiv preprint arXiv:2402.13725*.
- Treviso, M. and Martins, A. F. (2020). The explanation game: Towards prediction explainability through sparse communication. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 107–118.
- Tsallis, C. (1988). Possible generalization of boltzmann-gibbs statistics. *Journal of Statistical Physics*, 52:479–487.