

# From Sparse Modeling to Sparse Communication

André Martins



IC Colloquium EPFL, Lausanne, October 3, 2022

# Our Amazing Team



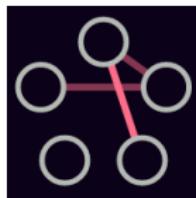
**SARDINE:** Structure AwaRe moDellNg for Natural LanguagE

# Our Amazing Team



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

# DeepSPIN



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



**Unbabel**



instituto de  
telecomunicações



**TÉCNICO**  
LISBOA

## 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:

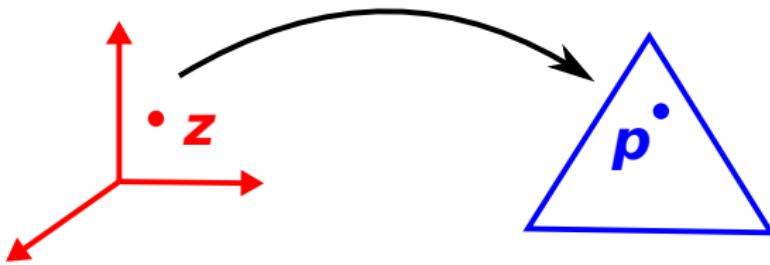
- Mostly used with neural networks, most work after 2015
- Main idea: sparse neuron activations (biological plausibility)
- Predictions are triggered by a few neurons only (input-dependent)
- Example: ReLUs, dropout, sparse attention mechanisms

# This Talk

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

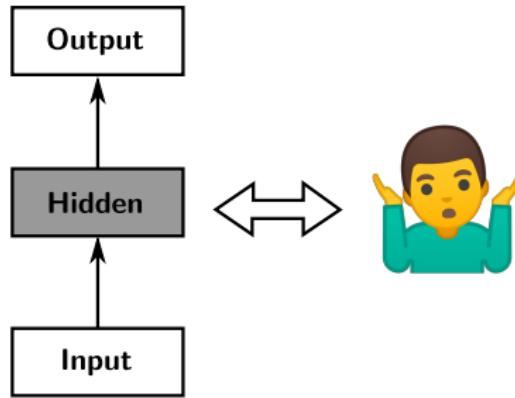
- All differentiable (efficient forward and backward propagation)
- Can be used at hidden (attention) or output layers (loss)
- Can make a bridge between the **continuous** and **discrete** worlds
- Effective in several natural language processing tasks.

Building block:

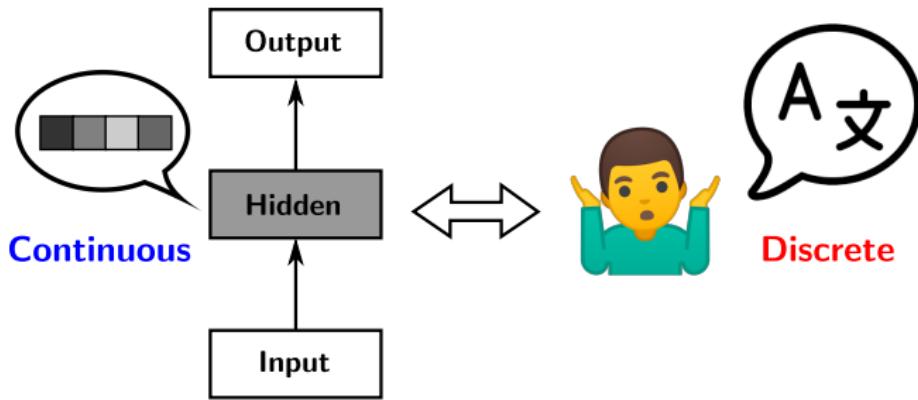


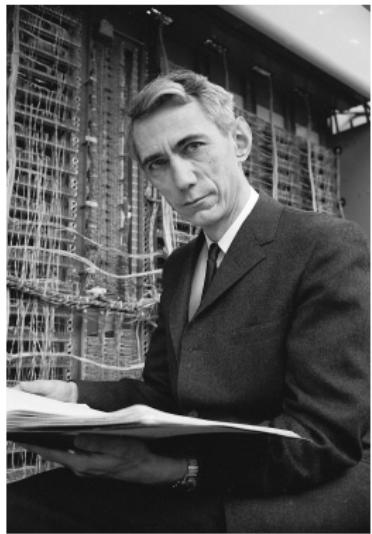
Sparse transformations from the Euclidean space to the simplex  $\Delta$ .

# 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

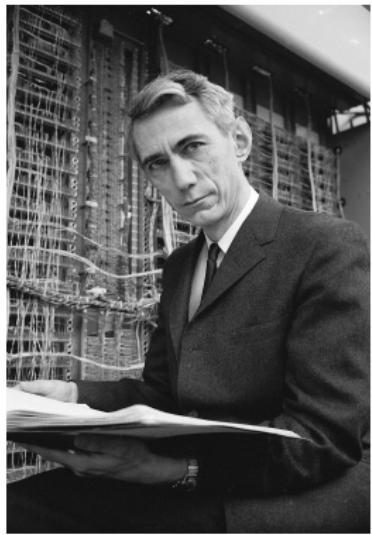
### PART I: DISCRETE NOISELESS SYSTEMS

### 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

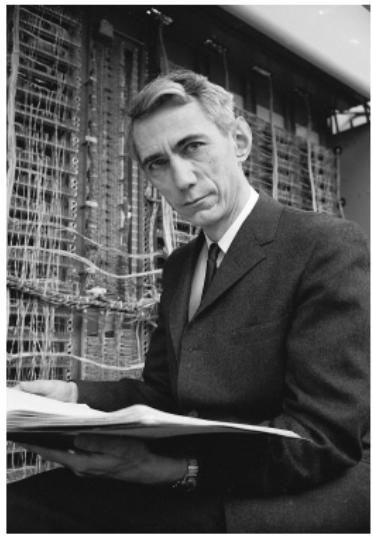
- PART I: DISCRETE NOISELESS SYSTEMS**
- PART II: THE DISCRETE CHANNEL WITH NOISE**

$\Sigma$

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

$\Sigma$

PART III: MATHEMATICAL PRELIMINARIES

PART IV: THE CONTINUOUS CHANNEL

$\int$

PART V: THE RATE FOR A CONTINUOUS SOURCE

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

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...

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

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...

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

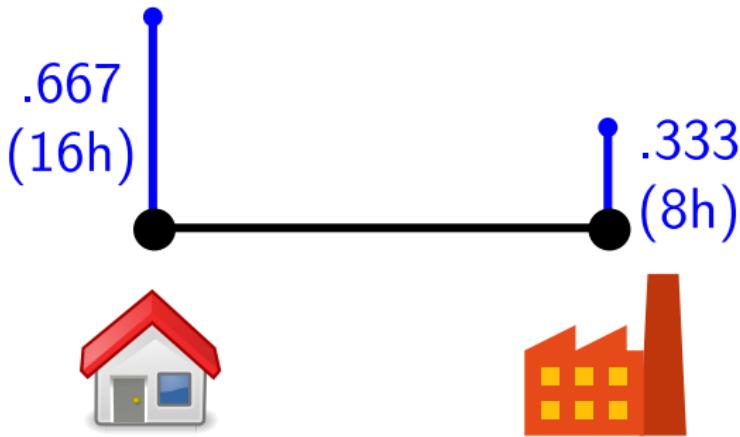
- Hammersley-Clifford theorem in graphical models
- Pitman-Koopman-Darmois theorem (sufficient statistics and exponential families)
- Log-likelihood is  $-\infty$  if estimated probability is 0.

## 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

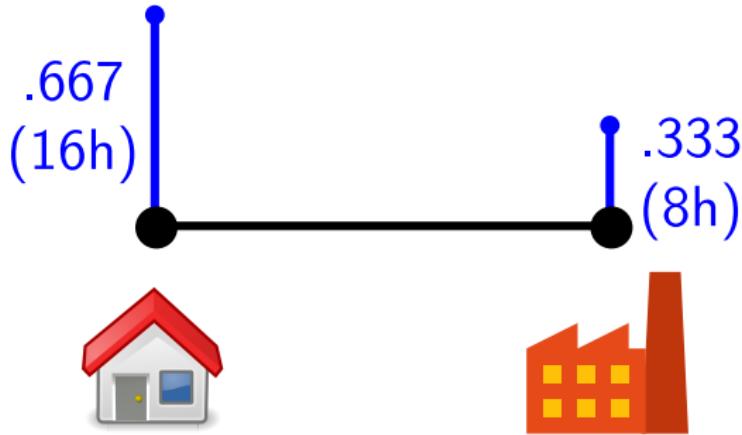
John splits his day as follows: he works 8h/day, and stays home 16h/day.



## Motivating Example: John's Life

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

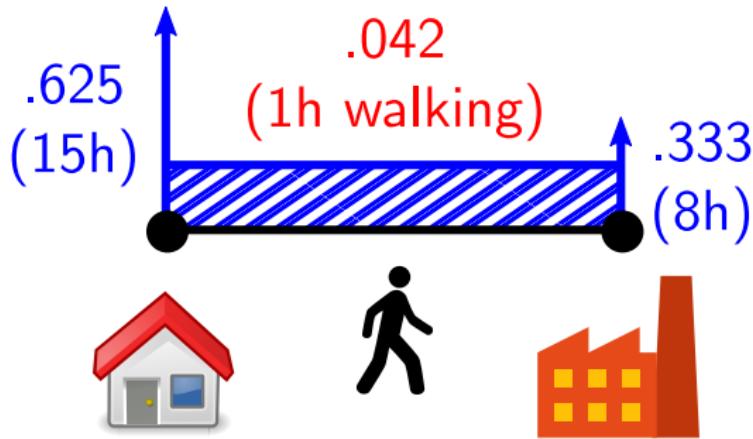
He is **in transit 1h/day** to commute to work and back.



## Motivating Example: John's Life

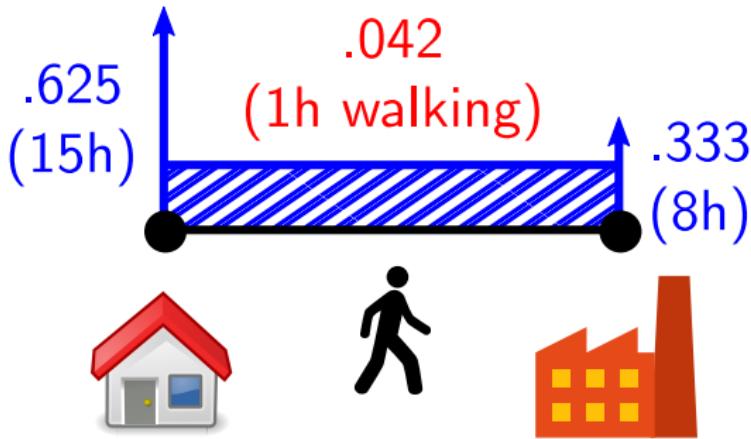
John splits his day as follows: he works 8h/day, and stays home **15h/day**.

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



## Motivating Example: John's Life

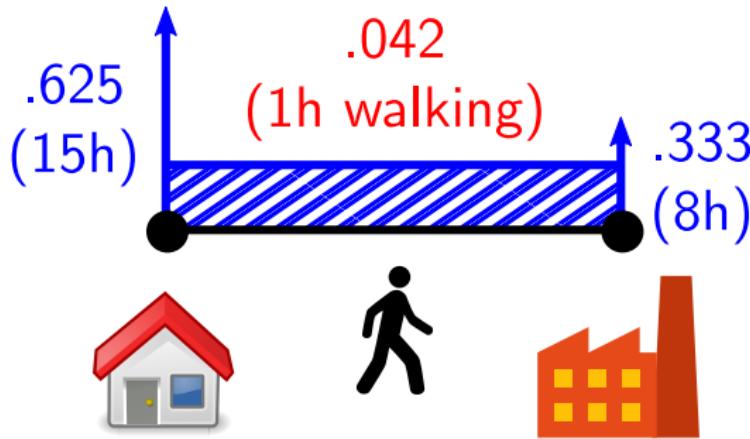
John splits his day as follows: he works 8h/day, and stays home **15h/day**.  
He is **in transit 1h/day** to commute to work and back.



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 Mixed Distributions

4 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:

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

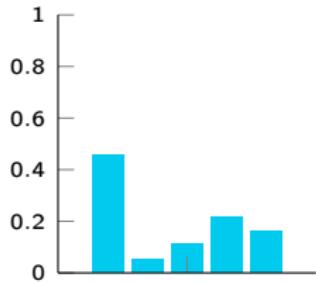
- Fully dense:  $\text{softmax}(\mathbf{z}) > 0, \forall \mathbf{z}$
- Used both as a loss function (cross-entropy) and for attention.

Argmax can be written as:

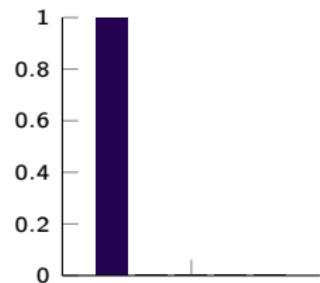
$$\begin{aligned}\text{argmax}(\mathbf{z}) &:= \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^\top \mathbf{p} \\ &= \lim_{\tau \rightarrow 0^+} \text{softmax}(\mathbf{z}/\tau) \quad (\text{temperature trick})\end{aligned}$$

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

$\text{softmax}(\mathbf{z})$



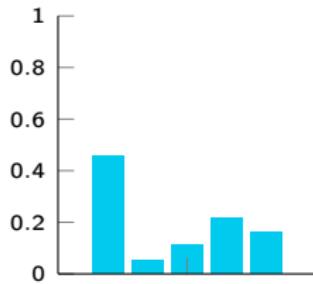
$\text{argmax}(\mathbf{z})$



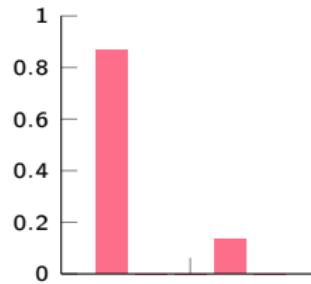
(Same  $\mathbf{z} = [1.0716, -1.1221, -0.3288, 0.3368, 0.0425]$ )

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

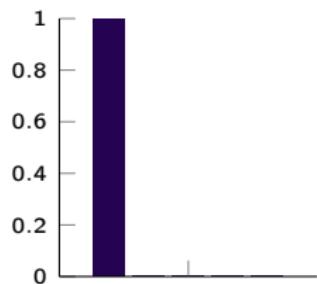
$\text{softmax}(\mathbf{z})$



$\text{sparsemax}(\mathbf{z})$



$\text{argmax}(\mathbf{z})$



(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)

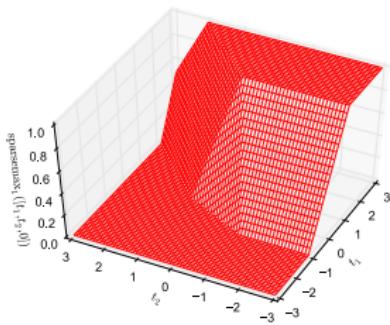
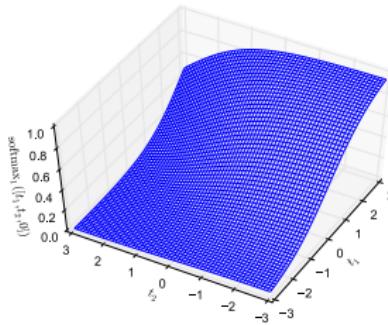
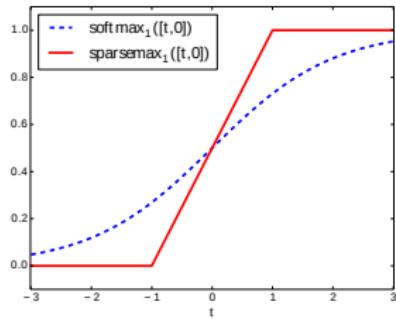
Euclidean projection of  $\mathbf{z}$  onto the probability simplex  $\Delta$ :

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

- 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.

## $\Omega$ -Regularized Argmax (Niculae and Blondel, 2017, NeurIPS)

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

$$\text{argmax}_{\Omega}(z) := \arg \max_{\mathbf{p} \in \Delta} z^T \mathbf{p} - \Omega(\mathbf{p})$$

- 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., 2019, ACL)

Parametrized by  $\alpha \geq 0$ :

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



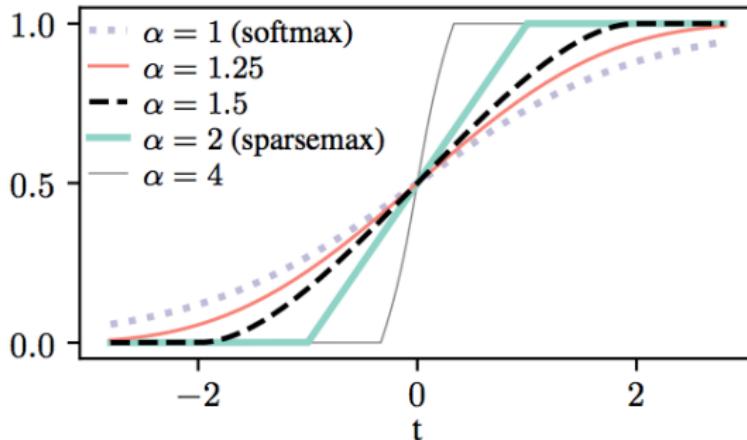
Related to **Tsallis generalized entropies** ([Tsallis, 1988](#)).

- **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., 2019, ACL)

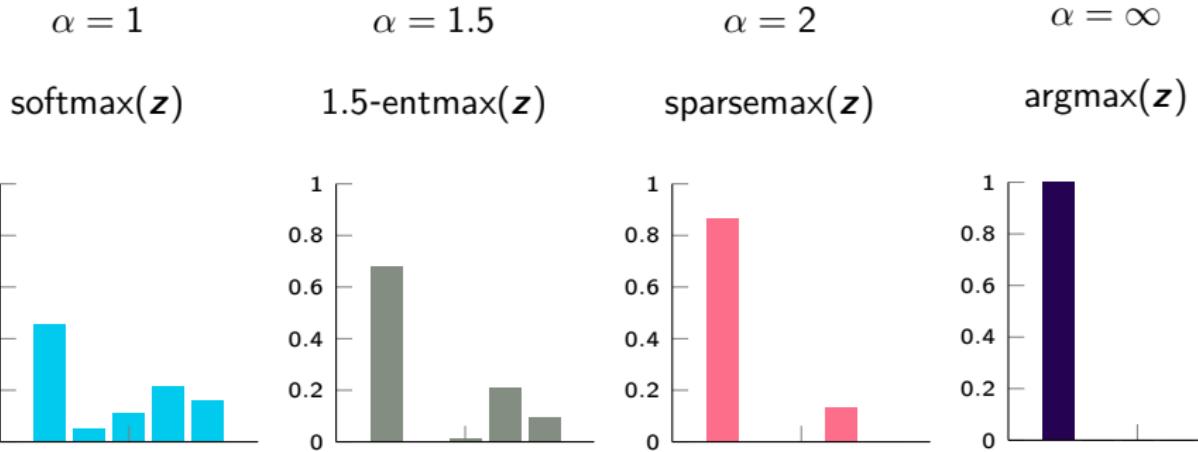


$\alpha = 1.5$  is a sweet spot!

- 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., 2019, ACL)

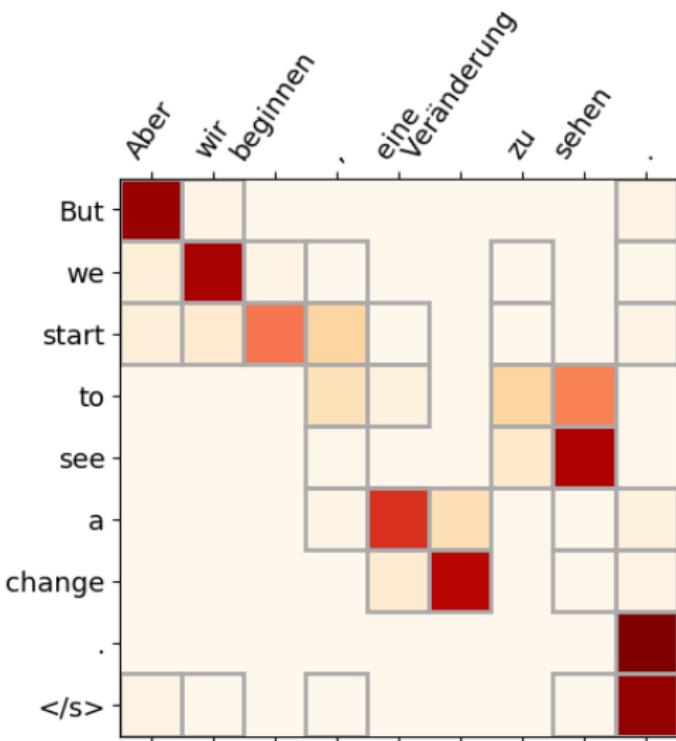


(Same  $\mathbf{z} = [1.0716, -1.1221, -0.3288, 0.3368, 0.0425]$ )

# Example: Sparse Attention for Machine Translation

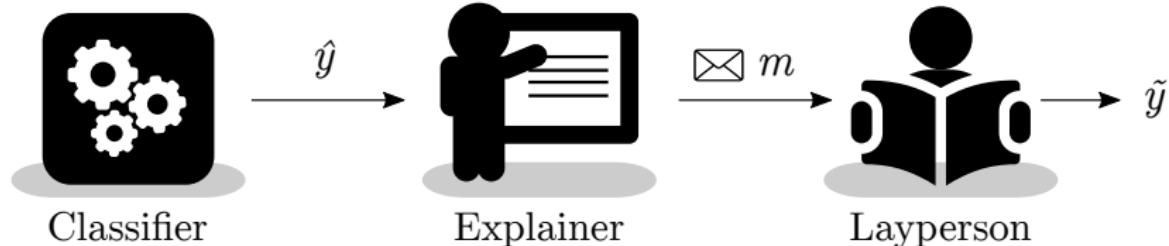
(Peters et al., 2019, ACL)

- 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

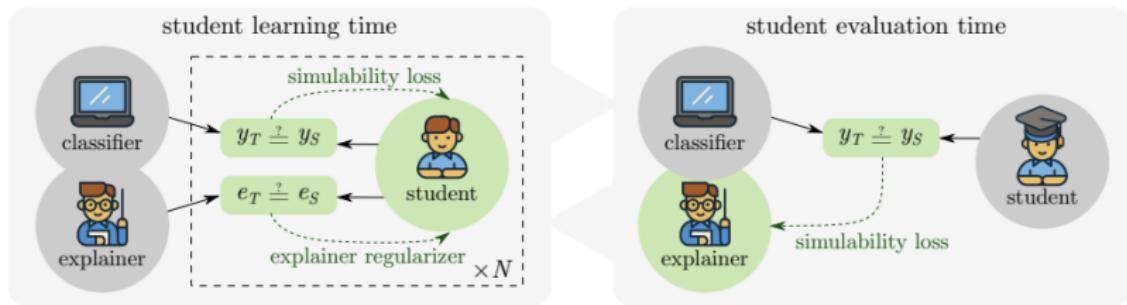
(Treviso and Martins, 2020, BlackboxNLP)



- 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 *simulability*, *forward simulation/prediction*)
- **Communication success rate:** how often the two predictions match?

# New: Scaffold Maximizing Training (SMaT)

(Fernandes et al., 2022, NeurIPS)



- Similar to above, but uses **bilevel optimization** (a la meta-learning) to learn to optimize model explanations for teaching.
- The teacher explainer is a linear combination of the classifier's attention layers and is used to regularize the student (Pruthi et al., 2022)
- The performance of the student in held-out data (simulability loss) is used as an objective for the teacher explainer.

## Other Related Transformations

Constrained softmax ([Martins and Kreutzer, 2017, EMNLP](#)),

Constrained sparsemax ([Malaviya et al., 2018, ACL](#)):

- Allows placing a **budget** on how much attention a word can receive
- Useful to model **fertility** in machine translation

Fusedmax ([Niculae and Blondel, 2017, NeurIPS](#)):

- Can promote **structured sparsity** (contiguous selection)

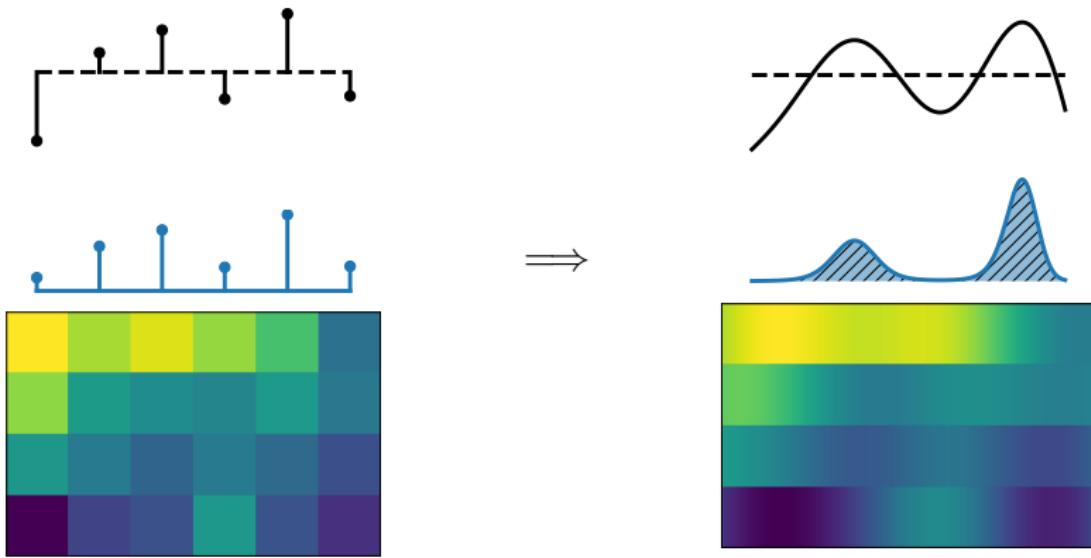
(LP-)SparseMAP [Niculae et al. \(2018, ICML\)](#), [Niculae and Martins \(2020, ICML\)](#):

- Extends sparsemax to **sparse structured prediction**.
- Can be used as hidden differentiable layer or output layer.
- Works with arbitrary factor graph (e.g. logic constraints).

# Sparse and Continuous Attention

(Martins et al., 2020a, NeurIPS)

- 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., 2022, ACL)



# Outline

1 Sparse Transformations

2 Fenchel-Young Losses

3 Mixed Distributions

4 Conclusions

# 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 (Niculae and Blondel, 2017, NeurIPS)

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

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

- **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., 2019, ACL).

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

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

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

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

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

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

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

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

## Properties:

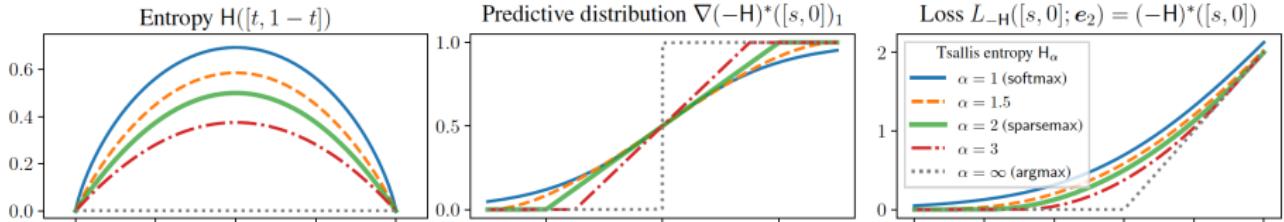
- $L_\Omega(\mathbf{z}, \mathbf{q}) \geq 0$  (automatic from Fenchel-Young inequality)
- $L_\Omega(\mathbf{z}, \mathbf{q}) = 0$  iff  $\mathbf{q} = \text{argmax } \Omega(\mathbf{z})$
- $L_\Omega$  is convex and differentiable with  $\nabla L_\Omega(\mathbf{z}, \mathbf{q}) = \text{argmax } \Omega(\mathbf{z})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

Recovers cross-entropy loss, sparsemax loss, and many other known losses

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

# Entmax Transformations and Losses

(Blondel et al., 2020, JMLR)



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

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

# Example: Machine Translation

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

<b>This</b>	92.9%	<b>is another</b>	view	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.")

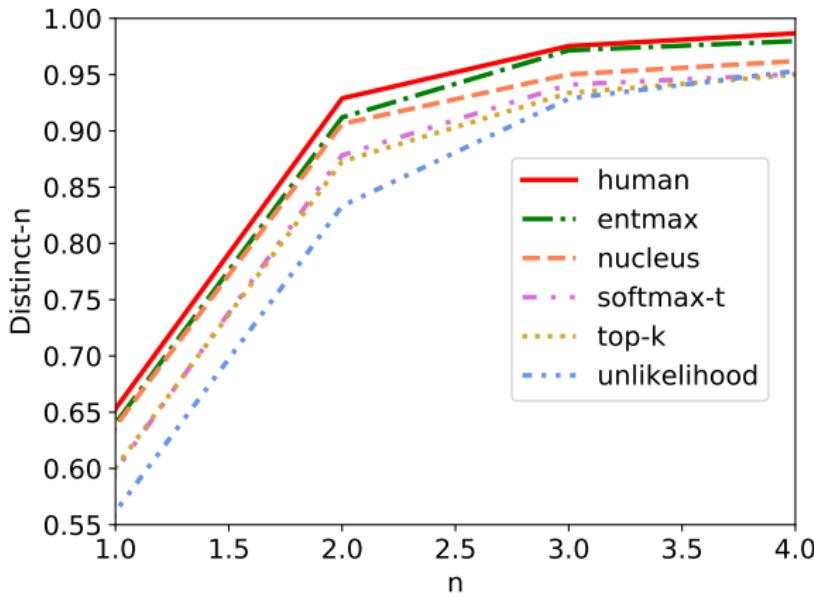
- Only a few words get non-zero probability at each time step
- Auto-completion when several words in a row have probability 1
- Useful for predictive translation.

# 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 Mixed Distributions

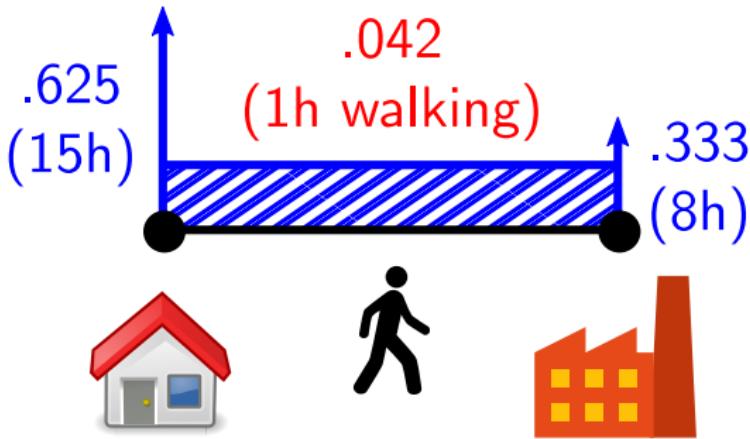
4 Conclusions

# 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.

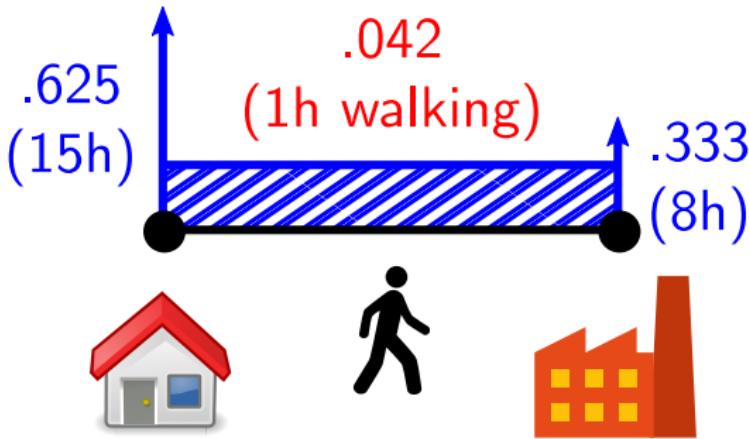
## Back to John's Life

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.



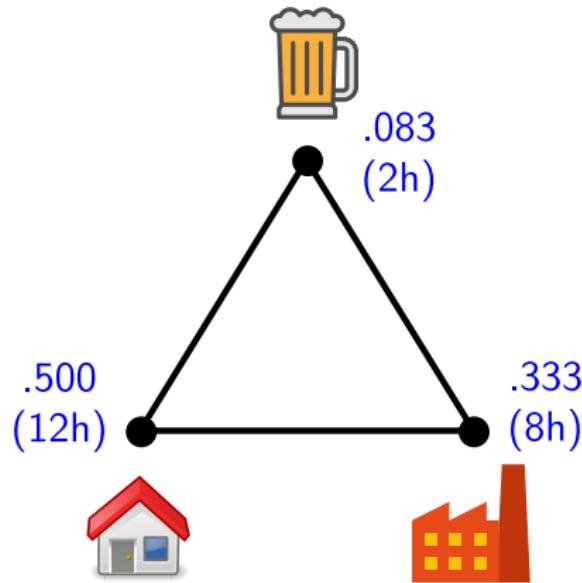
## Back to John's Life

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.

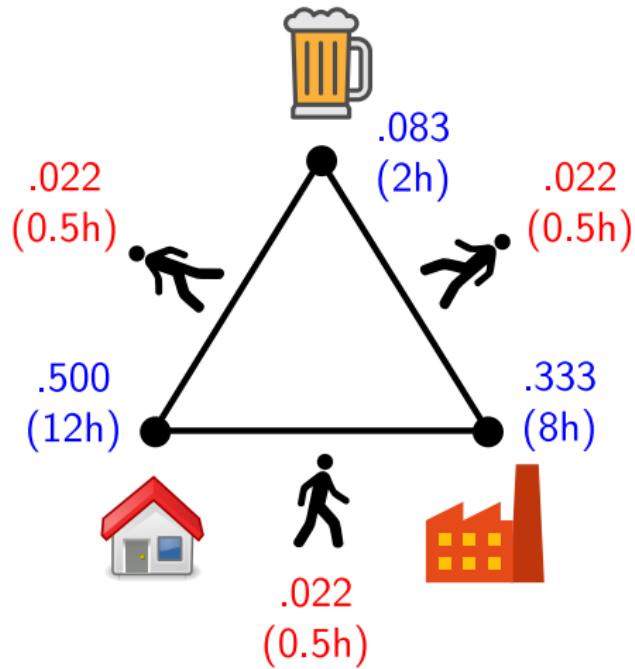


That's a sad life!

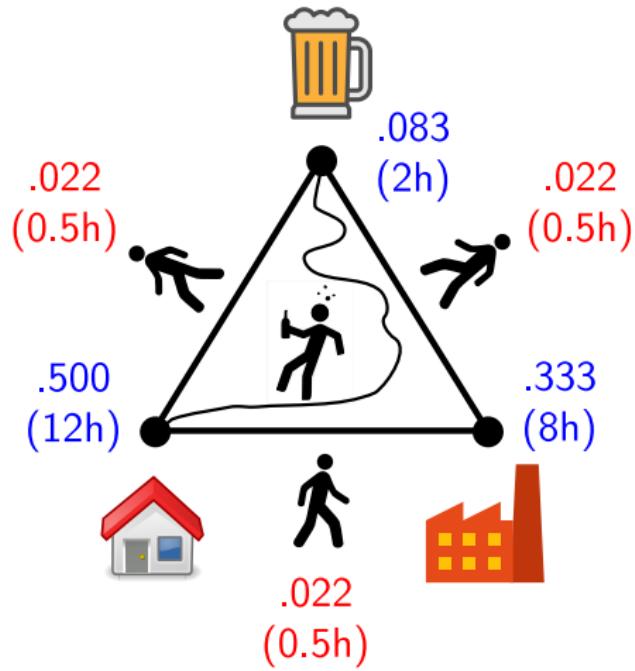
After work, John spends 2h in the pub with friends.



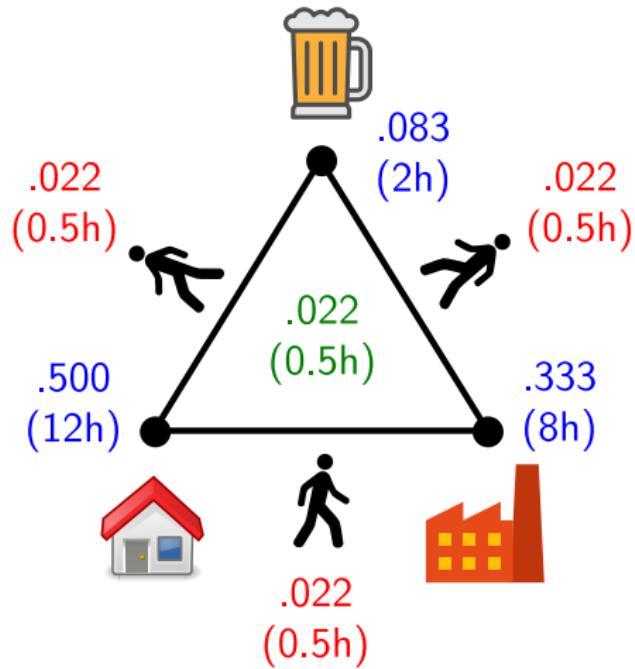
After work, John spends 2h in the pub with friends.



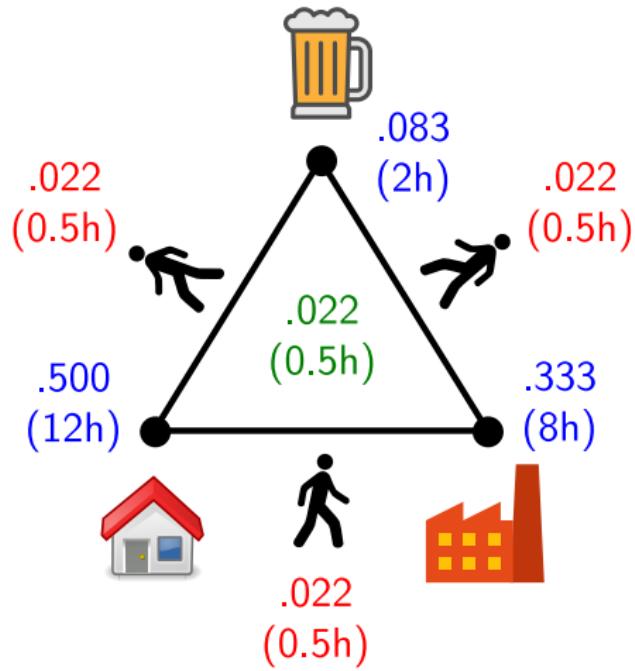
After work, John spends 2h in the pub with friends.



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 $\Delta_{K-1}$

We denote by  $\text{ri}(\Delta_{K-1})$  the **relative interior** of  $\Delta_{K-1}$ .

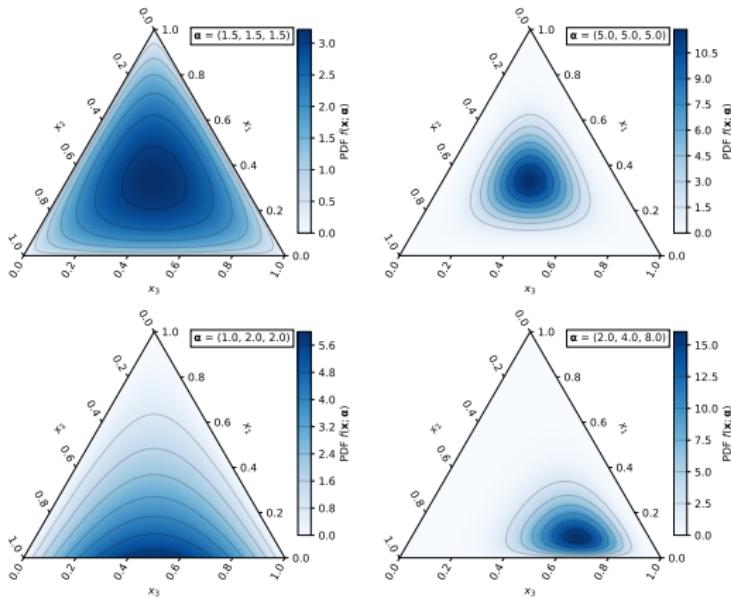
Common densities on the simplex:

- 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  
 $\Delta_{K-1} \setminus \text{ri}(\Delta_{K-1})$ .

# Dirichlet Distribution

$$Y \sim \text{Dirichlet}(\alpha) \Leftrightarrow p_Y(y; \alpha) \propto \prod_{k=1}^K y_k^{\alpha_k - 1}, \quad \alpha > 0.$$

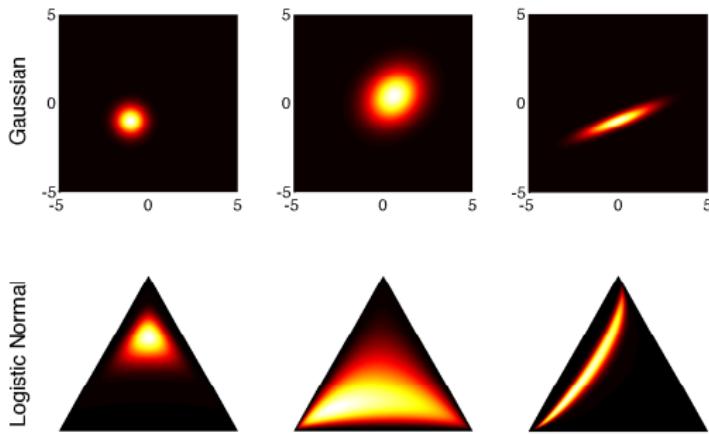


# Logistic Normal (a.k.a. Gaussian-Softmax)

(Atchison and Shen, 1980)

Generative story:

$$Y \sim \text{LogisticNormal}(z, \Sigma) \Leftrightarrow N \sim \mathcal{N}(0, I) \\ Y = \text{softmax}(z + \Sigma^{\frac{1}{2}} N).$$



# Concrete (a.k.a. Gumbel-Softmax)

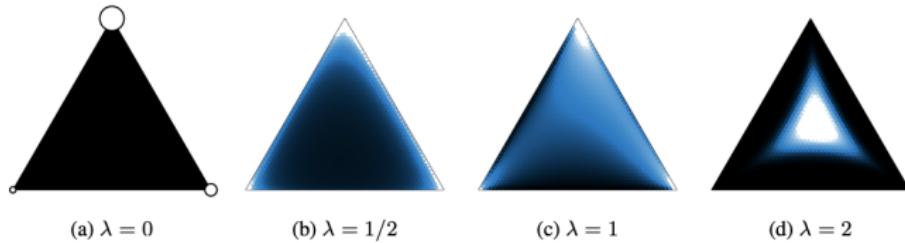
(Maddison et al., 2017; Jang et al., 2017)

Continuous relaxation of a categorical.

Approaches categorical as  $\lambda \rightarrow 0^+$  (Luce, 1959; Papandreou and Yuille, 2011).

Generative story:

$$Y \sim \text{Concrete}(z, \lambda) \iff \begin{aligned} G_k &\sim \text{Gumbel}(0, 1) \\ Y &= \text{softmax}(\lambda^{-1}(z + G)). \end{aligned}$$



## 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  $\mathbf{y} = [y, 1 - y]$ .

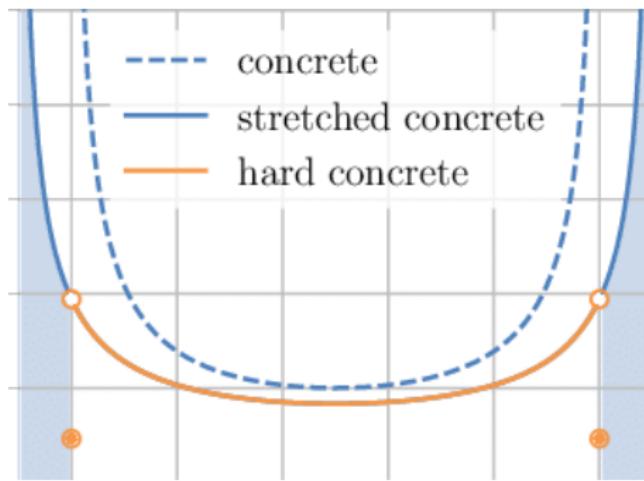
Truncated densities have been proposed for  $K = 2$ :

- Binary Hard Concrete
- Rectified Gaussian

# Binary Hard Concrete

(Louizos et al., 2018)

- Stretches the Concrete and applies a “hard” sigmoid transformation to place point masses at 0 and 1.
- Similar to spike-and-slab (Mitchell and Beauchamp, 1988; Ishwaran et al., 2005).



# Rectified Gaussian

(Hinton and Ghahramani, 1997; Palmer et al., 2017)

- Applies a “hard” sigmoid transformation to a univariate Gaussian.

$$p_Y(y) = \mathcal{N}(y; z, \sigma^2) + \frac{1 - \operatorname{erf}\left(\frac{z}{\sqrt{2}\sigma}\right)}{2} \delta_0(y) + \frac{1 + \operatorname{erf}\left(\frac{z-1}{\sqrt{2}\sigma}\right)}{2} \delta_1(y).$$

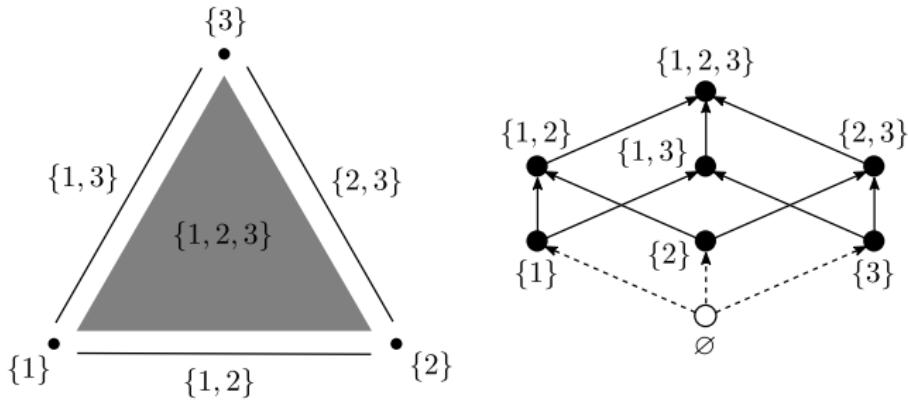
Extending such distributions to the multivariate case ( $K > 2$ ) is non-trivial:

- Combinatorially many multiple order Diracs would be needed
- Dirac deltas have  $-\infty$  differential entropy.

# 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, etc.

There is one  $(K - 1)$ -face: the simplex  $\Delta_{K-1}$  itself.

## Direct Sum Measure (Farinhas et al., 2022, ICLR)

Let  $\mathcal{F}$  denote the set of proper faces of  $\Delta_{K-1}$ ; we have  $|\mathcal{F}| = 2^K - 1$ .

We define a **direct sum measure**  $\mu^\oplus$  on  $\Delta_{K-1}$  as a sum of Lebesgue measures on each non-vertex face, and a counting measure on the vertices:

$$\mu^\oplus(A) = \sum_{f \in \mathcal{F}} \mu_f(A \cap \text{ri}(f)), \quad A \subseteq \Delta_{K-1}.$$

We 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  $\Delta_{K-1}$ ).

Continuous RVs assign probability only to the **maximal face** ( $\text{ri}(\Delta_{K-1})$ ).

**Mixed RVs generalize both:** can assign probability to **all faces** of  $\Delta_{K-1}$ .

# Mixed Random Variables

(Farinhas et al., 2022, ICLR)

Discrete RVs assign probability only to **0-faces** (vertices of  $\Delta_{K-1}$ ).

Continuous RVs assign probability only to the **maximal face** ( $\text{ri}(\Delta_{K-1})$ ).

**Mixed RVs generalize both:** can assign probability to **all faces** of  $\Delta_{K-1}$ .

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**  $p_{Y|F}(\mathbf{y} \mid f)$ , for  $f \in \mathcal{F}, \mathbf{y} \in \text{ri}(f)$ .

The probability of a set  $A \subseteq \Delta_{K-1}$  is given by:

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

Two ways of characterizing mixed RVs:

- **Extrinsic characterizaton**: start with a distribution over  $\mathbb{R}^K$  and then apply a deterministic transformation to project it to  $\Delta_{K-1}$
- **Intrinsic characterizaton**: specify a mixture of distributions directly over the faces of  $\Delta_{K-1}$ , by specifying  $P_F$  and  $p_{Y|F}$  for each  $f \in \mathcal{F}$

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

Generative story:

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

- 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.

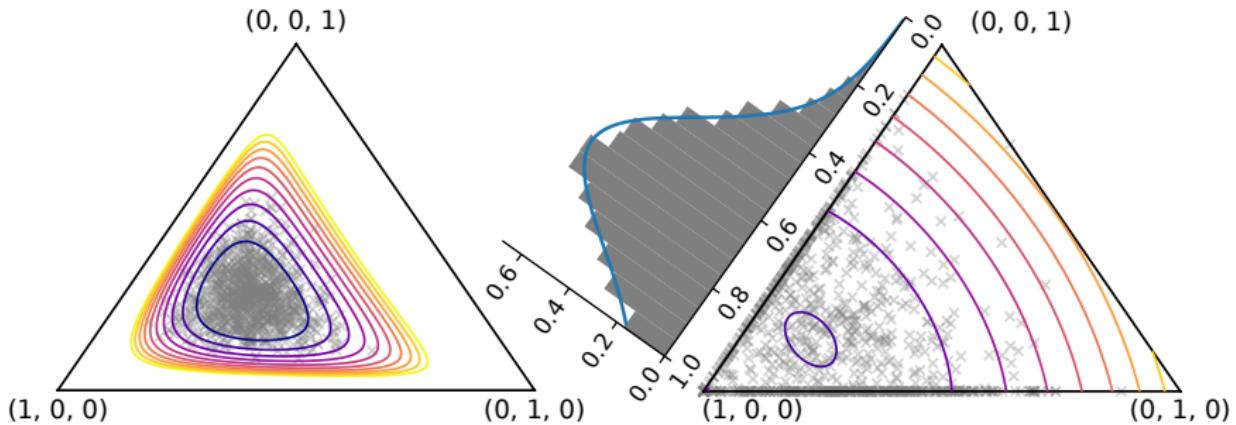
Uses an **extrinsic characterization**, by sampling from a Gaussian and projecting.

Generative story:

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

- 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  $\text{ri}(\Delta_{K-1})$

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

Uses an **intrinsic characterization**.

- Uses two parameters:  $\mathbf{w} \in \mathbb{R}^K$  and  $\alpha \in \mathbb{R}_{>0}^K$
- First, sample a face  $f \sim P_F(f) \propto \prod_{k \in f} w_k$ , where  $\mathbf{w} \in \mathbb{R}^K$
- Then, sample  $Y|F = f \sim \text{Dir}(\alpha|_f)$ , where  $\alpha|_f$  “masks out” entries of  $\alpha$  not supported by  $f$ .
- Sampling  $f$  can be done in  $\mathcal{O}(K)$  with dynamic programming.

# Information Theory of Mixed Random Variables

(Farinhas et al., 2022, ICLR)

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

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

- Average length of the optimal code where  $f$  must be encoded **losslessly** and where  $y|_f$  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 $\pm 8.07$	256
Gumbel-Softmax ST	49.96 $\pm 9.51$	1
K-D Hard Concrete	76.07 $\pm 7.76$	21.43 $\pm 17.56$
Gaussian-Sparsemax	<b>80.88</b> $\pm 0.50$	1.57 $\pm 0.02$

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

# Outline

1 Sparse Transformations

2 Fenchel-Young Losses

3 Mixed Distributions

4 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
- Examples: Gaussian-Sparsemax, Gumbel-Sparsemax, Mixed Dirichlet
- Sparse communication potentially useful as a path for explainability.

# Thank You!

DeepSPIN (“Deep Structured Prediction in NLP”)

- ERC starting grant, started in 2018
- Topics: deep learning, structured prediction, NLP
- More details: <https://deep-spin.github.io>



**Unbabel**



instituto de  
telecomunicações



TÉCNICO  
LISBOA

# References I

- Amari, S.-i. (2016). *Information geometry and its applications*, volume 194. Springer.
- Atchison, J. and Shen, S. M. (1980). Logistic-normal distributions: Some properties and uses. *Biometrika*, 67(2):261–272.
- Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., and Samek, W. (2015). On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):130–140.
- Bastings, J., Aziz, W., and Titov, I. (2019). Interpretable neural predictions with differentiable binary variables. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2963–2977, Florence, Italy. Association for Computational Linguistics.
- Blondel, M., Martins, A. F. T., and Niculae, V. (2020). Learning with fenchel-young losses. *Journal of Machine Learning Research*, 21(35):1–69.
- Church, K. W. and Hanks, P. (1990). Word association norms, mutual information, and lexicography. *Computational linguistics*, 16(1):22–29.
- Correia, G., Niculae, V., and Martins, A. F. T. (2019). Adaptively sparse transformers. In *Proceedings of the Empirical Methods for Natural Language Processing*.
- Farinhas, A., Aziz, W., Niculae, V., and Martins, A. F. (2022). Sparse communication via mixed distributions. In *Proc. of ICLR*.
- Feng, S., Wallace, E., Grissom II, A., Iyyer, M., Rodriguez, P., and Boyd-Graber, J. (2018). Pathologies of neural models make interpretations difficult. In *Proc. EMNLP*, pages 3719–3728.
- 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*.

## References II

- Guyon, I., Weston, J., Barnhill, S., and Vapnik, V. (2002). Gene selection for cancer classification using support vector machines. *Machine learning*, 46(1-3):389–422.
- 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.
- Huber, P. J. (1964). Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*, 35(1):73–101.
- Ishwaran, H., Rao, J. S., et al. (2005). Spike and slab variable selection: frequentist and bayesian strategies. *Annals of Statistics*, 33(2):730–773.
- Jang, E., Gu, S., and Poole, B. (2017). Categorical reparameterization with gumbel-softmax. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- Kohavi, R. and John, G. H. (1997). Wrappers for feature subset selection. *Artificial intelligence*, 97(1-2):273–324.
- Lei, T., Barzilay, R., and Jaakkola, T. (2016). Rationalizing neural predictions. *preprint arXiv:1606.04155*.
- Li, J., Chen, X., Hovy, E., and Jurafsky, D. (2016a). Visualizing and understanding neural models in nlp. In *Proc. NAACL-HLT*, pages 681–691.
- Li, J., Monroe, W., and Jurafsky, D. (2016b). Understanding neural networks through representation erasure. *arXiv preprint arXiv:1612.08220*.
- Louizos, C., Welling, M., and Kingma, D. P. (2018). Learning sparse neural networks through  $l_0$  regularization. In *International Conference on Learning Representations*.

## References III

- Luce, R. D. (1959). *Individual choice behavior: A theoretical analysis*. New York: Wiley, 1959.
- Maddison, C. J., Mnih, A., and Teh, Y. W. (2017). The concrete distribution: A continuous relaxation of discrete random variables. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- 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 Computation 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. and Kreutzer, J. (2017). Fully differentiable neural easy-first taggers. In *Proc. of Empirical Methods for Natural Language Processing*.
- Martins, P. H., Marinho, Z., and Martins, A. F. (2020b). Sparse text generation. In *Empirical Methods for Natural Language Processing*.
- Martins, P. H., Marinho, Z., and Martins, A. F. T. (2022).  $\infty$ -former: Infinite memory transformer. In *Proc. of NAACL-HLT*.
- Mitchell, T. J. and Beauchamp, J. J. (1988). Bayesian variable selection in linear regression. *Journal of the American Statistical Association*, 83(404):1023–1032.

## References IV

- 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 Scheding, S. J. (2017). Methods for stochastic collection and replenishment (scar) optimisation for persistent autonomy. *Robotics and Autonomous Systems*, 87:51–65.
- Papandreou, G. and Yuille, A. L. (2011). Perturb-and-map random fields: Using discrete optimization to learn and sample from energy models. In *2011 International Conference on Computer Vision*, pages 193–200.
- 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. (2019). Sparse sequence-to-sequence models. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*.
- Pruthi, D., Bansal, R., Dhingra, B., Soares, L. B., Collins, M., Lipton, Z. C., Neubig, G., and Cohen, W. W. (2022). Evaluating explanations: How much do explanations from the teacher aid students? *Transactions of the Association for Computational Linguistics*, 10:359–375.
- Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). Why should i trust you?: Explaining the predictions of any classifier. In *Proc. ACM SIGKDD*, pages 1135–1144. ACM.

# References V

- Serrano, S. and Smith, N. A. (2019). Is attention interpretable? In *Proc. ACL*.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society*, 58(1):267–288.
- 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.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., and Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. In *Proc. ICLR*, pages 2048–2057.
- Zhang, T. (2004). Statistical Behavior and Consistency of Classification Methods Based on Convex Risk Minimization. *Annals of Statistics*, pages 56–85.
- Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the royal statistical society*, 67(2):301–320.