

Sparse Communication via Mixed Distributions

António Farinhas¹ Wilker Aziz² Vlad Niculae³ André F. T. Martins^{1,4,5}

¹Instituto de Telecomunicações, Instituto Superior Técnico, Lisbon, Portugal ²ILLC, University of Amsterdam, The Netherlands
³IvI, University of Amsterdam, The Netherlands ⁴Unbabel, Lisbon, Portugal ⁵LUM LIS (Lisbon ELLIS Unit), Lisbon, Portugal



Motivation

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

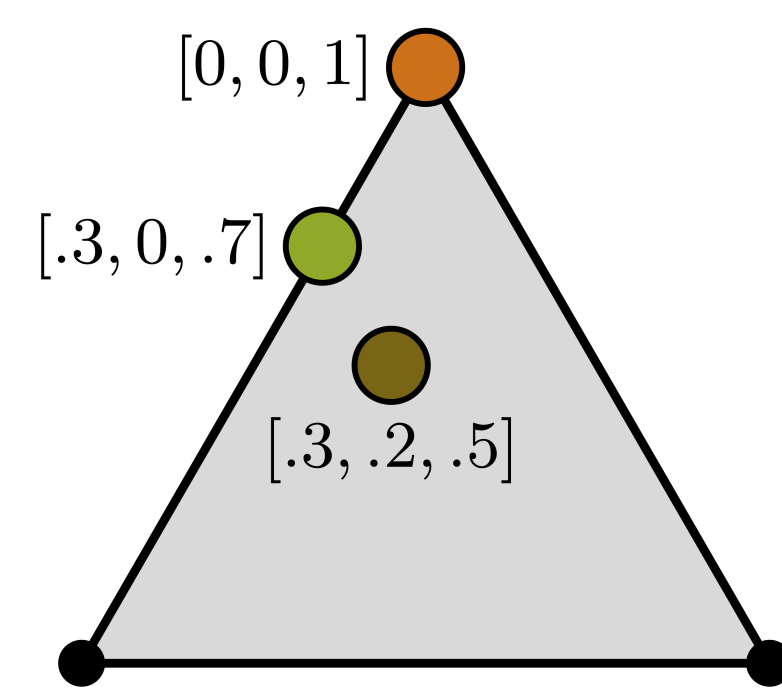
- Language is symbolic and mostly *discrete*
- Neural networks learn and use *continuous* representations

What happens in-between? Can sparsity help?

Transformations from \mathbb{R}^K to Δ_{K-1}

How can we convert a vector of real numbers $z \in \mathbb{R}^K$ (scores for the several symbols, often called *logits*) into a probability vector $y \in \Delta_{K-1}$?

- $y = \text{softmax}(z) \propto \exp(z)$
- $y = \lim_{\tau \rightarrow 0^+} \text{softmax}(z/\tau)$
- $y = \text{sparsemax}(z) := \arg \min_{y \in \Delta_{K-1}} \|y - z\|$



Densities over the simplex Δ_{K-1}

There are several works on defining distributions on the simplex Δ_{K-1}

- The Dirichlet, the Concrete (Maddison et al., 2017; Jang et al., 2017), and the Logistic-Normal (Atchison and Shen, 1980) are restricted to $\text{ri}(\Delta_{K-1})$
- The Hard Concrete Louizos et al. (2018) and rectified Gaussians (Palmer et al., 2017) are **mixed** discrete/continuous hybrids limited to $K = 2$

This work:

- Mathematical theory for handling mixed random variables
- Provide extensions to $K > 2$

Extending truncated densities to $K > 2$

We define probability densities with respect to the direct sum measure,

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

where μ_f is the $\dim(f)$ -dimensional Lebesgue measure for $\dim(f) > 0$, and the counting measure for $\dim(f) = 0$.

Mixed random variables

How to define probability densities?

- Define a probability mass function $P_F(f)$ on \mathcal{F}
- For each face $f \in \mathcal{F}$, define a probability density $p_{Y|F}(y | f)$ over $\text{ri}(f)$

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

$$\Pr\{y \in A\} = \int_A p_Y^\oplus(y) d\mu^\oplus = \sum_{f \in \mathcal{F}} P_F(f) \int_{A \cap \text{ri}(f)} p_{Y|F}(y | f) \quad (2)$$

Recovers both *discrete* and *continuous* distributions!

Information Theory for mixed random variables

See our paper for generalizations of information theoretic concepts such as entropy, Kullback-Leibler divergence, and mutual information.

The **entropy** of a r.v. X with respect to a measure μ is

$$H^\mu(X) = - \int_{\mathcal{X}} p_X(x) \log p_X(x) d\mu(x), \quad \text{with } \int_{\mathcal{X}} p_X(x) d\mu(x) = 1 \quad (3)$$

- \mathcal{X} finite, μ counting measure: **Shannon's discrete entropy**
- $\mathcal{X} \subseteq \mathbb{R}^k$ continuous, μ Lebesgue measure: **differential entropy**
- What if μ is the direct sum measure?

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

Average length of the optimal code where the sparsity pattern of $y \in \Delta_{K-1}$ must be encoded losslessly and where there is a predefined bit precision for the fractional entries of y .

The maximum entropy mixed distribution is written as a generalized Laguerre polynomial

- $\log_2(2 + 2^N)$ for $K = 2$, instead of $\log_2(2) = 1$ in the purely discrete case

Intrinsic characterization

Specify a mixture of distributions directly over the faces of Δ_{K-1} : P_F and $p_{Y|F}$ for each $f \in \mathcal{F}$.

Mixed Dirichlet (two parameters: $w \in \mathbb{R}^K$ and $\alpha \in \mathbb{R}_{>0}^K$)

- Sample a face $f \sim P_F(f) \propto \prod_{k \in f} \exp(w_k)$
- Sample $Y|F = f \sim \text{Dir}(\alpha|_f)$ where $\alpha|_f$ masks out entries of α not supported by f

Extrinsic characterization

Start with a distribution over the ambient space and project it to the simplex using *sparsemax*.

K-D Hard Concrete (generalization of the binary Hard Concrete for $K > 2$)

$$Y' \sim \text{Concrete}(z, \beta), \quad Y = \text{sparsemax}(\lambda Y'), \quad \text{with } \lambda \geq 1. \quad (5)$$

Gaussian-Sparsemax (generalization of a double-sided rectified Gaussian for $K > 2$)

$$N \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad Y = \text{sparsemax}(z + \Sigma^{1/2} N) \quad (6)$$

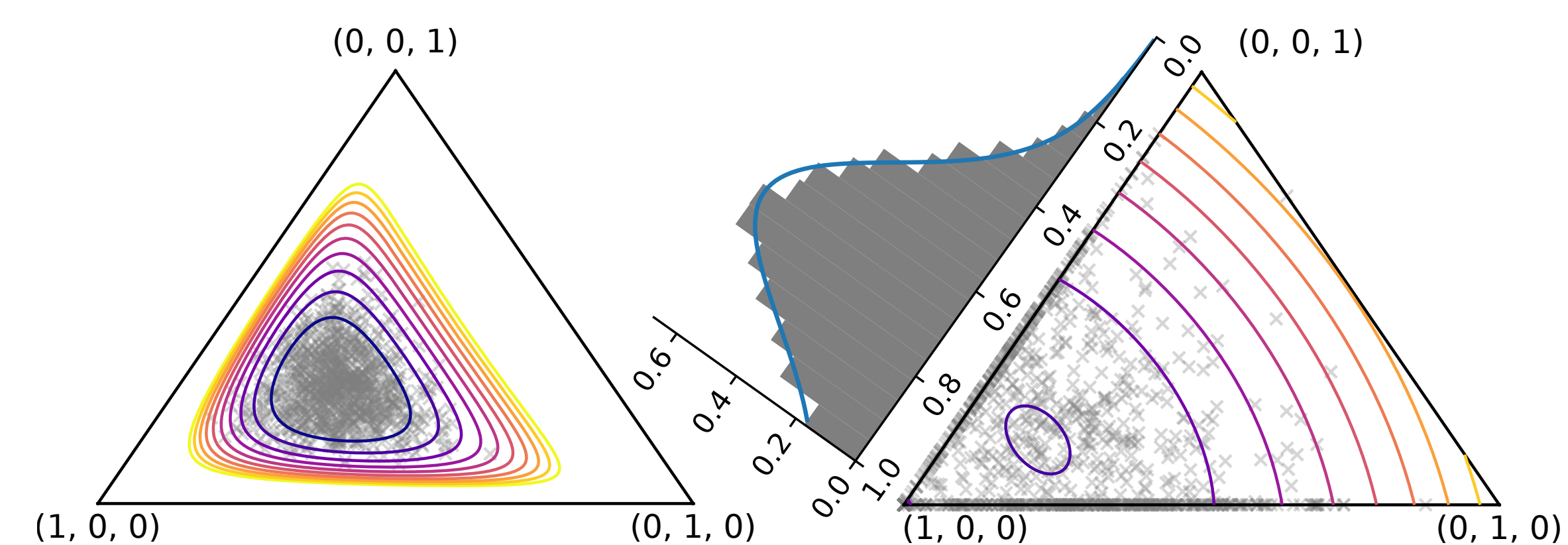


Figure: The Logistic-Normal (left), assigns zero probability to all faces but $\text{ri}(\Delta_{K-1})$. The Gaussian-Sparsemax (right) induces a distribution over the 1-dimensional edges (shown as a histogram), and assigns $\Pr\{(1, 0, 0)\} = .022$.

Can also be defined intrinsically

- Expressed via the orthant probability of multivariate Gaussians for $K > 2$
- Simple expression for $K = 2$; entropy and KL divergence in closed form

Experiments

Emergent communication game (inducing sparse communication between two agents)

- A *sender* sees an image and emits a single-symbol message from a fixed vocabulary
- A *receiver* reads the symbol and tries to identify the correct image out of a set of 16

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

(4) **Bit-Vector VAE on Fashion-MNIST** (studying the impact of the direct sum entropy)

- We consider 128 binary latent bits and maximize the ELBO

Method	Entropy	NLL	Sparsity (%) ↑
Binary Concrete	C ≈ 3.60	3.60	0
Gumbel-Softmax	D = 3.49	3.49	0
Gumbel-Softmax ST	D = 3.57	3.57	100
Hard Concrete	X ≈ 3.57	3.57	45.64
Gaussian-Sparsemax	X ≈ 3.53	3.53	82.82
Gaussian-Sparsemax	X = 3.49	3.49	73.83

Regression towards voting proportions (using the Mixed Dirichlet as a likelihood function)

- We use the UK election data; the observations are vectors of proportions over 5 parties
- Modeling simplex-valued data with the Dirichlet is tricky

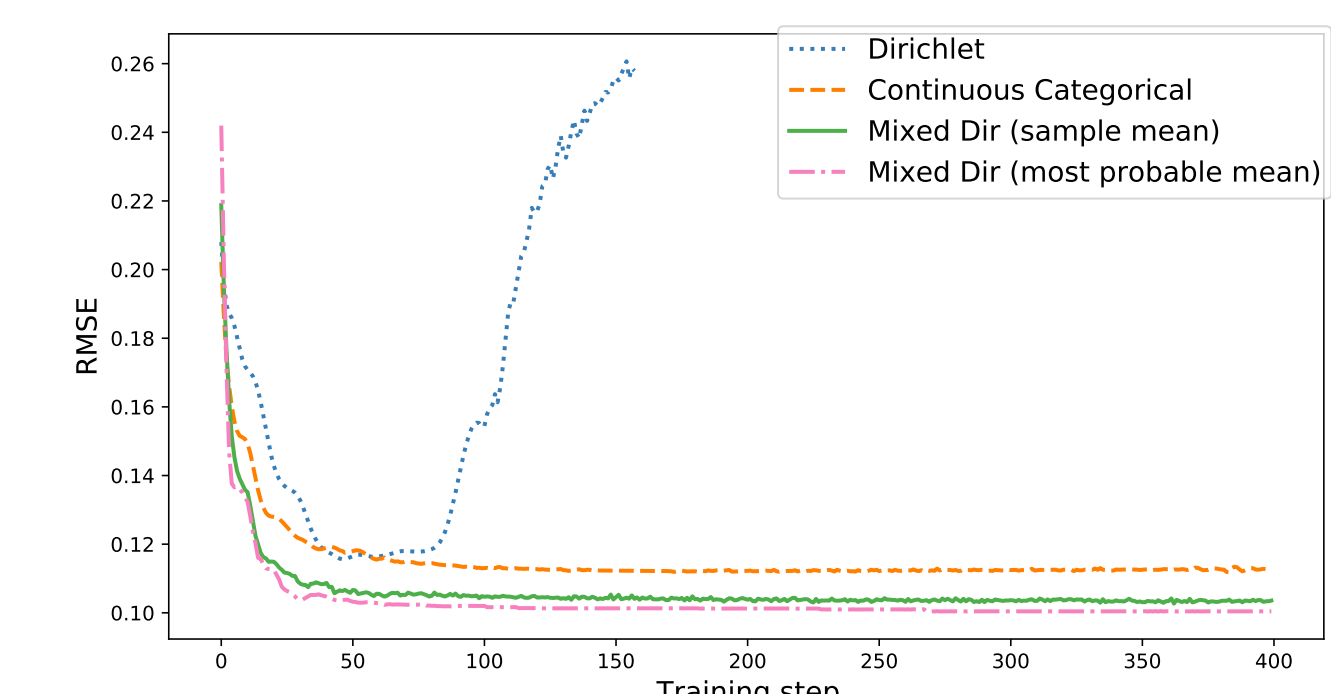


Figure: The Mixed Dirichlet addresses the pathologies of the Dirichlet in this setting, showing a slight advantage over the continuous categorical (Gordon-Rodriguez et al., 2020), likely due to the fact that Mixed Dirichlet samples are often sparse at test time.

Conclusions

- Mathematical framework for handling mixed random variables
- Direct sum measure as an alternative to the Lebesgue-Borel and the counting measures
- Generalizations of information theoretic concepts for mixed symbols
- Experiments on learning sparse representations and avoiding ill-defined log-likelihoods
- Future work:** More effective intrinsic parametrizations; mixed *structured* variables

Open-source code: <https://github.com/deep-spin/sparse-communication>

References

- Atchison, J. and Shen, S. M. (1980). Logistic-normal distributions: Some properties and uses. *Biometrika*.
- Gordon-Rodriguez, E., Loaiza-Ganem, G., and Cunningham, J. (2020). The continuous categorical: a novel simplex-valued exponential family. In *Proceedings of ICML*.
- Jang, E., Gu, S., and Poole, B. (2017). Categorical reparameterization with gumbel-softmax. In *Proceedings of ICLR*.
- Louizos, C., Welling, M., and Kingma, D. P. (2018). Learning sparse neural networks through l_0 regularization. In *Proceedings of ICLR*.
- Maddison, C. J., Mnih, A., and Teh, Y. W. (2017). The concrete distribution: A continuous relaxation of discrete random variables. In *Proceedings of ICLR*.
- 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*.