

Simplex-to-Euclidean Bijections for Categorical Flow Matching

Bernardo Williams¹, Victor M. Yeom-Song^{2,3}, Marcelo Hartmann¹, Arto Klami¹

¹University of Helsinki

²Aalto University

³ELLIS Institute Finland

The Problem: Discrete Data Generation



Binary images ($K = 2$ categories)

ACCCGCGACAGACACAAA
CTGGCGCTGCCCTGTGGGA
CGCGGACACTCGATCTCGG
AGCTCATTGCACGCTGCTG
CGGTTTCAGGTCGGGGAACG
CCAGCTGCGCAGCGGAAG
CACGATGGGCCAGCTAATG
GTCGGGGCATGGACCAATT
GGCGTCGACGGGGGCGGGA
TCGGGGAAGGTATATAAGC
CAGCCGGGGCGGCCGGGGC

DNA sequences ($K = 4$
categories)

the founder of modern anarchist theory in what is property proudhon answers with the famous accusation property is theft in this work he opposed the institution of decreed property propri t where owners have complete rights to use and abuse their property as they wish such as exploiting workers for profit in its place proudhon supported what he called possession individuals can have limited rights to use resources capital and goods in accordance with principles of equality and justice proudhon s vision of anarchy which he called mutualism mutuellisme involved an exchange economy where individuals and groups could trade the products of their labor using labor notes which represented the amount of working time

Text ($K = 27$ categories)

The Problem: Discrete data generation

- ▶ K Categorical data are vertices of the simplex:

$$c \in \{1, \dots, K\},$$

$$\Delta^D := \{\mathbf{x} \in \mathbb{R}_+^K : \sum_{k=1}^K x_k = 1\}$$

$$e_c := (0, \dots, \underbrace{1}_{c\text{-th entry}}, \dots, 0) \in \Delta^D$$

- ▶ **Problem:** The simplex has boundaries and **non-Euclidean** geometry e.g. Fisher.

Goal. Reuse *Euclidean generative models* while respecting Simplex geometry.

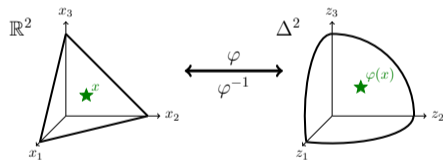


Figure: The Fisher information geometry is the sphere.

Method: Move the Modeling Problem to Euclidean Space

1. Stochastic interpolation

$$\mathbf{x} = \lambda \mathbf{e}_c + (1 - \lambda) \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \text{Dirichlet}(\boldsymbol{\alpha}).$$

2. Apply a smooth bijection

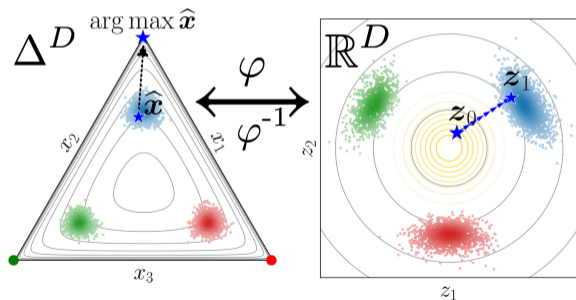
$$\mathbf{z}_1 = \varphi(\mathbf{x}), \quad \varphi : \Delta^D \rightarrow \mathbb{R}^D.$$

3. Train flow matching

$$\mathbf{z}_t = (1 - t)\mathbf{z}_0 + t\mathbf{z}_1, \quad \dot{\mathbf{z}}_t = \mathbf{z}_1 - \mathbf{z}_0,$$

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}[\|\mathbf{v}_\theta(t, \mathbf{z}_t) - \dot{\mathbf{z}}_t\|^2].$$

Sampling. Draw $\mathbf{z}_0 \sim p_0$, solve ODE $\dot{\mathbf{z}}_t = \mathbf{v}_\theta(t, \mathbf{z}_t)$, map back with $\hat{\mathbf{x}} = \varphi^{-1}(\mathbf{z}_1)$, then use $\arg \max_k \hat{\mathbf{x}}_k$.



Method details

Bijection 1: Isometric Logratio transform (ILR) H Helmert matrix.

$$\varphi : \Delta^D \rightarrow \mathbb{R}^D, \mathbf{x} \mapsto \mathbf{z} = \mathbf{H} \log \mathbf{x},$$

$$\varphi^{-1} : \mathbb{R}^D \rightarrow \Delta^D, \mathbf{z} \mapsto \mathbf{x} = \text{softmax}(\mathbf{H}^\top \mathbf{z}).$$

Isometric to $(\mathbb{R}^D, \|\cdot\|_2)$.

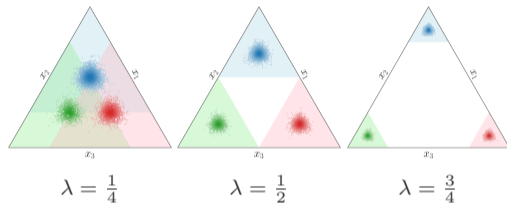
Bijection 2: Stick-Breaking (SB)
(coordinates are centered before)

$$\varphi : \overset{\circ}{\Delta}^D \rightarrow \mathbb{R}^D, z_k = \log \left(x_k / \left(1 - \sum_{i=1}^k x_i \right) \right),$$

$$\varphi^{-1} : \mathbb{R}^D \rightarrow \overset{\circ}{\Delta}^D, x_k = \sigma(z_k) \prod_{i=1}^{k-1} (1 - \sigma(z_i)).$$

Exact recovery:

Interpolation: $\mathbf{x} = \lambda \mathbf{e}_c + (1 - \lambda)\boldsymbol{\epsilon}$,

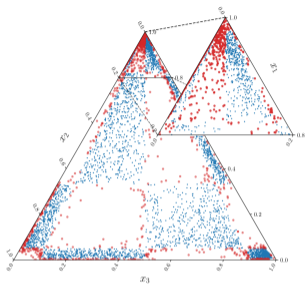


if $\lambda \geq 1/2$ recover original category

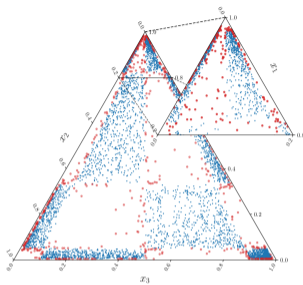
$$c = \arg \max_k x_k.$$

^oJ. J. Egozcue et al., *Isometric logratio transformations for compositional data analysis*, 2003.

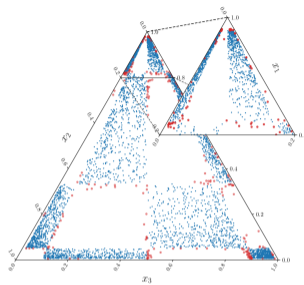
Empirical results



(a) LinearFM, 25.9%



(b) SFM, 12.9%



(c) Ours, 5.4%

Dataset	Metric	Baseline	Ours
BMNIST	FID	SFM 4.62	ILR 4.36
DNA	SP-MSE	SFM .0258	SB .0214
Text8	NLL	SFM 6.85	ILR 6.81

Pattern. Our approach beats simplex-based methods on discrete-data generation.

Takeaway

FM- $\overset{\circ}{\Delta}^D$ makes discrete-data generation look like continuous flow matching.

- ▶ Stochastic interpolation turns vertices into recoverable interior points.
- ▶ Smooth bijections move the simplex interior to \mathbb{R}^D .
- ▶ Existing Euclidean generative models can be reused directly.
- ▶ Examples include diffusion models, consistency models, flow maps, and mean flows.

See you at Poster 187!



[Link to the paper](#)