



Summary

We develop flows for discrete distributions.

• Discrete autoregressive flows enable multiple levels of autoregressivity.

Ex. Bidirectional language models that can generate data.

• Discrete bipartite flows enable flexible models with parallel generation.

Ex. Non-autoregressive text models with an exact likelihood.



Bipartite Flows.

$$y_{\geq d} = x_{\geq d} \cdot \sigma(x_{< d}) + \mu(x_{< q})$$
$$\left|\frac{\partial f}{\partial x}\right| = \prod \sigma(x_{< d})_i$$

Discrete Flows Invertible Generative Models of Discrete Data

Dustin Tran¹ Keyon Vafa^{1,2} Kumar Krishna Agrawal¹ Laurent Dinh¹ Ben Poole¹ **dustinvtran**

Discrete Change of Variables







(b) Factorized Base

Let x be a discrete random variable and y = f(x) where f is some function of **x**. The induced probability mass function of **y** is:

$$p(\mathbf{y} = y) = \sum_{x \in f^{-1}(y)} p(\mathbf{x} = x)$$

For an invertible function f, this simplifies to

$$p(\mathbf{y}=y)=p(\mathbf{x}=f^{-1}(y)).$$

Discrete Flows

Use location-scale transformation on the modulo integer space: $\cdot \mathbf{x}_d$) mod K.

$$\mathbf{y}_d = (\boldsymbol{\mu}_d + \boldsymbol{\sigma}_d)$$

 σ_d and μ_d are autoregressive or bipartite functions of y in 0,1,...,K-1 and 1,...,K-1.

For the flow to be invertible, σ_d and K must be coprime (inverse uses Euclid's algorithm). For example: mask noninvertible values for a given K; or make K prime.

Training Discrete Flows

The maximum likelihood objective per datapoint is

 $\log p(\mathbf{y}) = \log p(f^{-1}(\mathbf{y})).$

Gradient descent with respect to base distribution parameters is straightforward.

Gradient descent with respect to flow parameters requires backpropagation through the discrete-output function. We use the straight-through gradient estimator.

¹Google Brain ²Columbia University





(c) 1 Flow

4.5 4.0 3.5 3.0 2.5 2.0 1.5 1.0 0 1 2 3 Num				
3-layer LSTM Ziegler and Ru Ziegler and Ru Bipartite flow				
LSTM (Cooijr 64-layer Transforme Bipartite flow (4 Bipartite flow (8 Bipartite flow (8				
Autor				
D = 2, K = 2 D = 5, K = 5 D = 5, K = 10 D = 10, K = 5				
D = 2, K = 2 D = 5, K = 5 D = 5, K = 10 D = 10, K = 5 Full-Rank Disc Bipartite flows g				
D = 2, K = 2 D = 5, K = 5 D = 5, K = 10 D = 10, K = 5 Full-Rank Disc Bipartite flows g				
D = 2, K = 2 $D = 5, K = 5$ $D = 5, K = 10$ $D = 10, K = 5$ Full-Rank Disc Bipartite flows generative flows ge				
D = 2, K = 2 $D = 5, K = 5$ $D = 5, K = 10$ $D = 10, K = 5$ Full-Rank Disc Bipartite flows of Bipartite flows of Gradient bias • Devising more References "Normalizing Flo Papamakarios of				
D = 2, K = 2 $D = 5, K = 5$ $D = 5, K = 10$ $D = 10, K = 5$ Full-Rank Disc Bipartite flows generative flo				



Experiments								
H=256 L=2 H=512 L=2 H=512 L=2 with the second secon	th_scale	Texta flows nona 100x	8. LSTM per flow. Bipartite a get best utoregressive results and speedup.					
Test NLL (bpc) Generation								
(Merity et al., 201	l <mark>8</mark>)	1.18	3.8 min					
ush (2019) (AF/SC	(F)	1.4	-6 -					
ush (2019) (IAF/S	CF)	1.6	-					
		1.3	8 0.17 sec					
1ans+2016)	bpc 1.43	Gen. 19.8s	Penn Tree Bank. 2-3 layer Transformer per					
(Al-Rfou+2018)	1.13	35.5s	flow. Bipartite flows get					
flows, w/ σ)	1.60	0.15s	best nonautoregressive					
flows, w/o σ)	1.29	0.16s	results 1000x speedun					
flows, w/ σ)	1.23	0.16s						

pregressive Base	Autoregressive Flow	Factorized Base	Bipartite Flow
0.9	0.9	1.3	1.0
7.7	7.6	8.0	7.9
10.7	10.3	11.5	10.7
15.9	15.7	16.6	16.0

:rete. Discrete autoregressive flows are best. get similar performance as autoregressive base.

Limitations

permutations only.

over many flows and many classes.

re flexible transformations. RNG/compression?

lows for Probabilistic Modeling and Inference." G. et al., 2019.

ete Flows and Lossless Compression." E. : al., 2019.