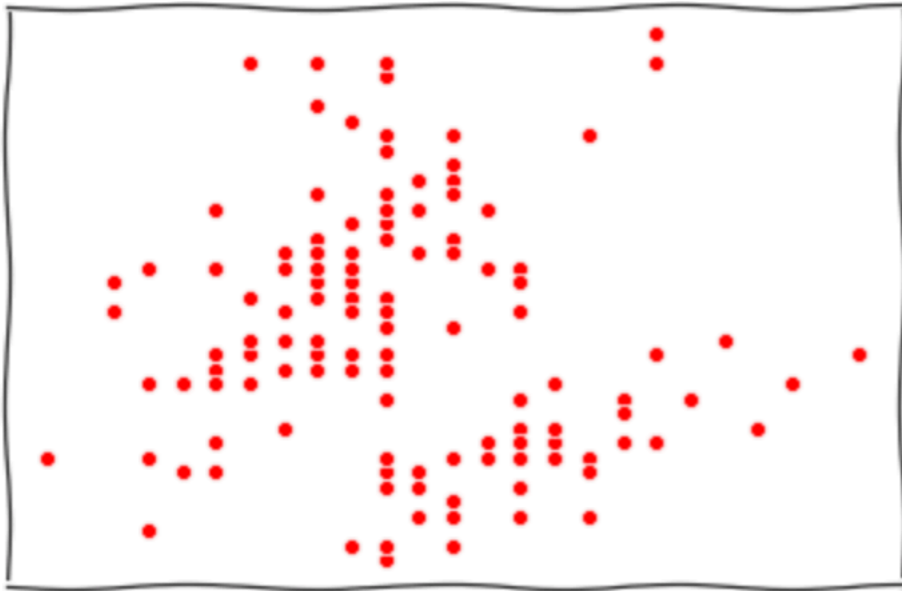# Lec 20: Triangular Flows

Yaoliang Yu

July 16, 2020
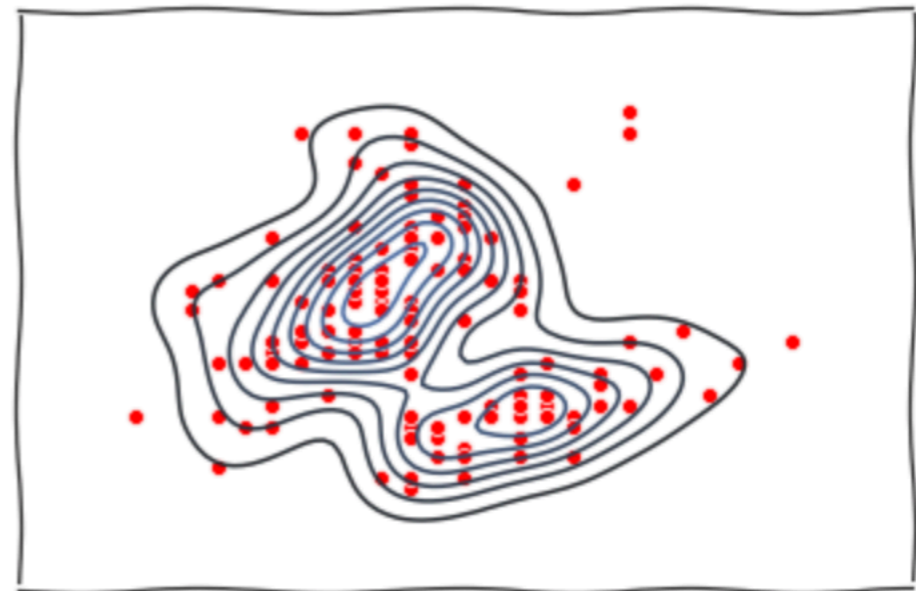
UNIVERSITY OF WATERLOO

# density estimation



**data** $= \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n\}$
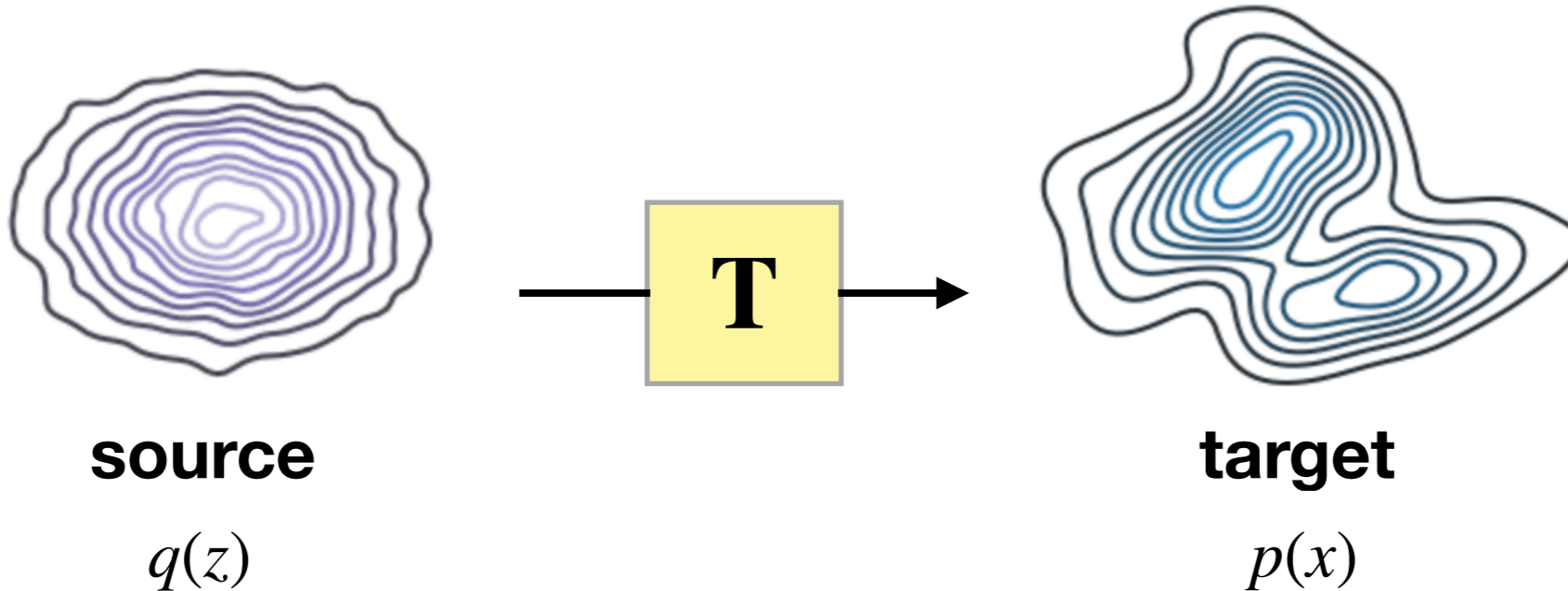
**estimate** $p(\mathbf{x})$

**If we can generate, then we can classify**

# Realistic

# Overview

find *deterministic* maps from source to target density



**source**

$q(z)$

**target**

$p(x)$

learn *bijective* & *differentiable* transformations

→ change of variables gives target density

computation of inverse and Jacobian must be *cheap*

→ always possible via triangular maps

$$x_i = \mathbf{T}(z_i)$$

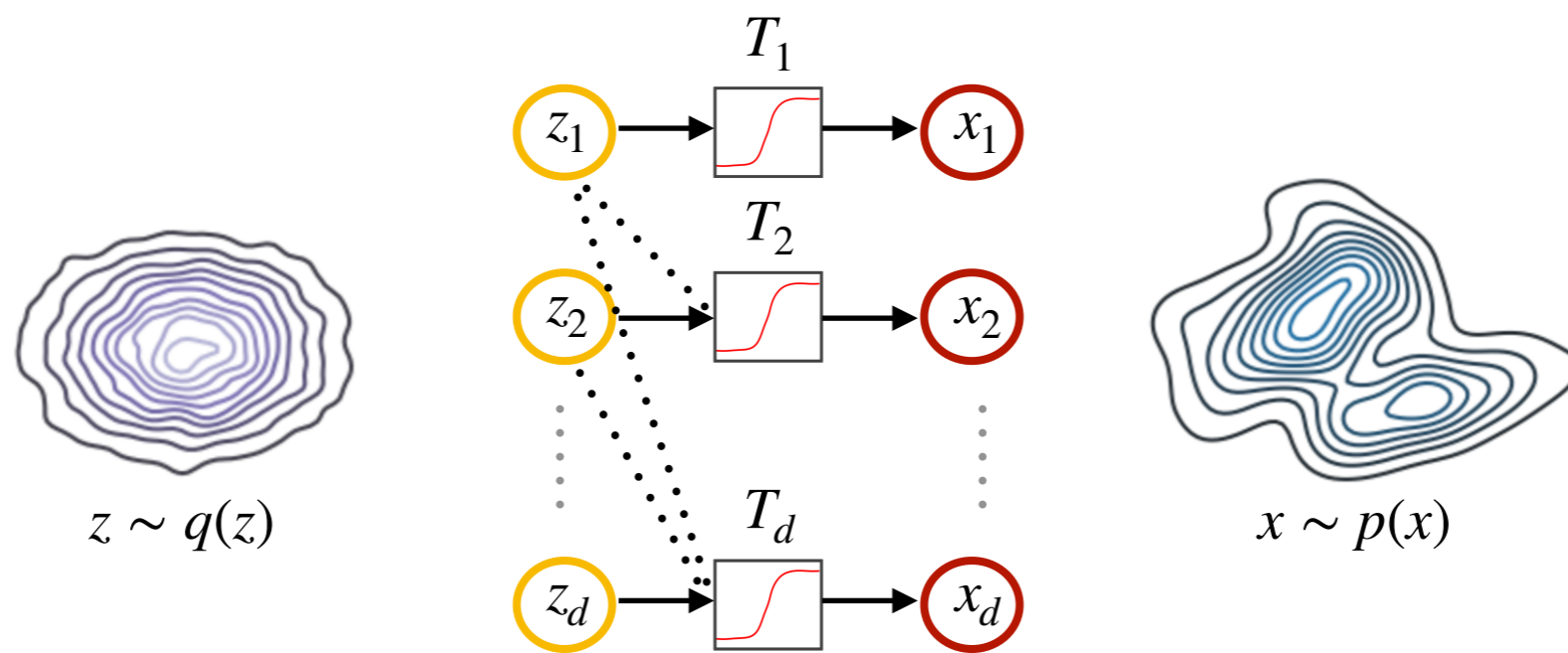$$(T_{\#}q)(x_i) = q(\mathbf{T}^{-1}(x_i)) \cdot \left| \mathbf{T}'\left( \mathbf{T}^{-1}(x_i) \right) \right|$$

$$\max_{\mathbf{T}} \mathscr{L}(\mathbf{T}) := \max_{\mathbf{T}} \prod_{i=1}^{n} (T_{\#}q)(x_i)$$

**unifying framework**

density estimation via increasing triangular maps

$T_1$

$z_1 \rightarrow \boxed{\int} \rightarrow x_1$

$T_2$

$z_2 \rightarrow \boxed{\int} \rightarrow x_2$

$T_d$

$z_d \rightarrow \boxed{\int} \rightarrow x_d$

$z \sim q(z)$

$x \sim p(x)$

# In a nutshell

**Given simulator for sampling from a normal distribution**

**How to simulate samples from a chi^2 distribution?**

# increasing triangular maps

$$\mathbf{T} : \mathbb{R}^d \to \mathbb{R}^d$$

$$x_1 = T_1(z_1)$$
$$x_2 = T_2(z_1, z_2)$$
$$x_3 = T_3(z_1, z_2, z_3)$$
$$\vdots$$
$$x_d = T_d(z_1, z_2, z_3, \ldots, z_d)$$

$$\nabla_{\mathbf{z}} \mathbf{T} = \begin{bmatrix} \dfrac{\partial T_1}{\partial z_1} & 0 & \ldots & 0 \\ \dfrac{\partial T_2}{\partial z_1} & \dfrac{\partial T_2}{\partial z_2} & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{\partial T_d}{\partial z_1} & \dfrac{\partial T_d}{\partial z_2} & \ldots & \dfrac{\partial T_d}{\partial z_d} \end{bmatrix}$$

**triangular :** $T_j$ **is a function of** $z_1, z_2, \ldots, z_j$

**increasing :** $T_j$ **is increasing w.r.t** $z_j$

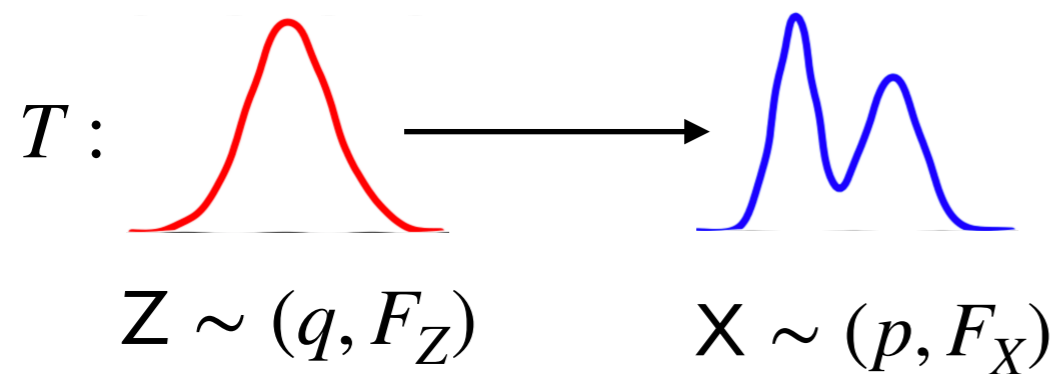$$\frac{\partial T_j}{\partial z_j} > 0$$

**triangular maps**

inverse and Jacobian are *easy* to compute

Theorem (paraphrase) : there always exists a unique* increasing triangular map that transforms a source density to a target density
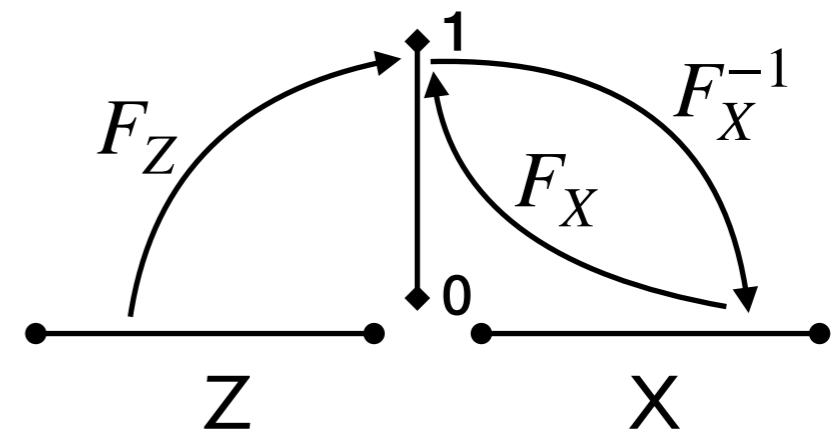
*Bogachev, V. et. al. Triangular Transformation of Measures, Sbornik: Mathematics, 2005*

*** for a fixed ordering**

# examples

increasing rearrangement

$$T := F_X^{-1} \circ F_Z$$



$T :$

$Z \sim (q, F_Z)$   $X \sim (p, F_X)$

Knothe-Rosenblatt transformation

iterative application of increasing rearrangement

$T :$

$Z \sim (q, F)$   $X \sim (p, G)$

$$T_1(z_1) := G_1^{-1} \circ F_1(z_1)$$

$$T_2(z_2, z_1) := G_{2|1}^{-1} \circ F_{2|1}(z_2)$$

*Villani, C. Optimal Transport: Old and New, vol 338, Springer, 2008*

# More Examples



$$TT^\top = \Sigma$$

$$Z \sim \mathcal{N}(0,I) \quad \longrightarrow \quad TZ =: X \sim \mathcal{N}(0,\Sigma)$$

Unique increasing triangular T = chol($\Sigma$)

# Maximum Likelihood revisited



**T**

$T_1$
$z_1 \rightarrow \boxed{} \rightarrow x_1$

$T_2$
$z_2 \rightarrow \boxed{} \rightarrow x_2$

$T_d$
$z_d \rightarrow \boxed{} \rightarrow x_d$

**source**

$q(\mathbf{z})$

$\mathbf{D}\Big( \quad , \quad \Big)$

**KL**

**Push-forward**

$(\mathbf{T}_{\#}q)(\mathbf{x})$

$\|$

$q(\mathbf{T}^{-1}(\mathbf{x})) \cdot \left| \mathbf{T}'\big(\mathbf{T}^{-1}(\mathbf{x})\big) \right|^{-1}$

**target**

$p(\mathbf{x})$

**learn T by maximizing likelihood**

$$\min_{\mathbf{T}} \ \sum_{i=1}^{n} \left[ -\log q\big(\mathbf{T}^{-1}(\mathbf{x}_i)\big) + \sum_{j} \log \partial_j T_j\big(\mathbf{T}^{-1}(\mathbf{x}_i)\big) \right]$$

*Marzouk, Y. et.al. Sampling via Measure Transport: An Introduction, Springer, 2016*
*[JSY]. Sum-of-squares Polynomial Flow. ICML, 2019*
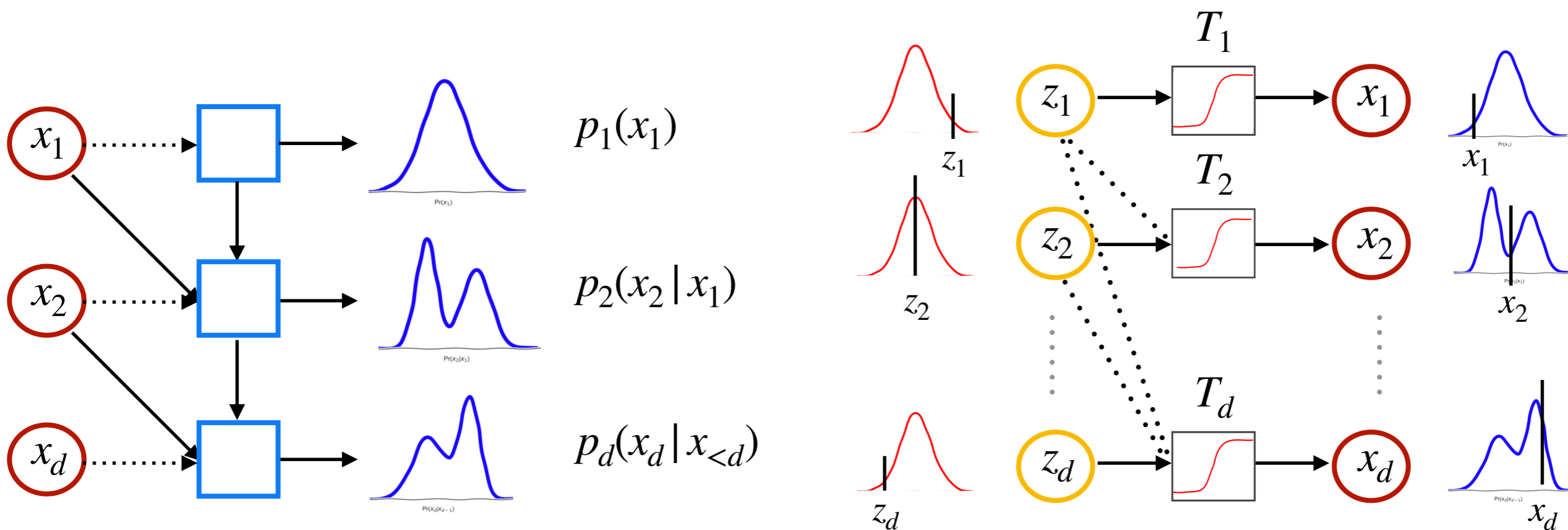
**explicitly** **evaluating** $q_\theta(\mathbf{x})$

# flow models as triangular maps

study commonalities & differences of flow based methods

# autoregressive models

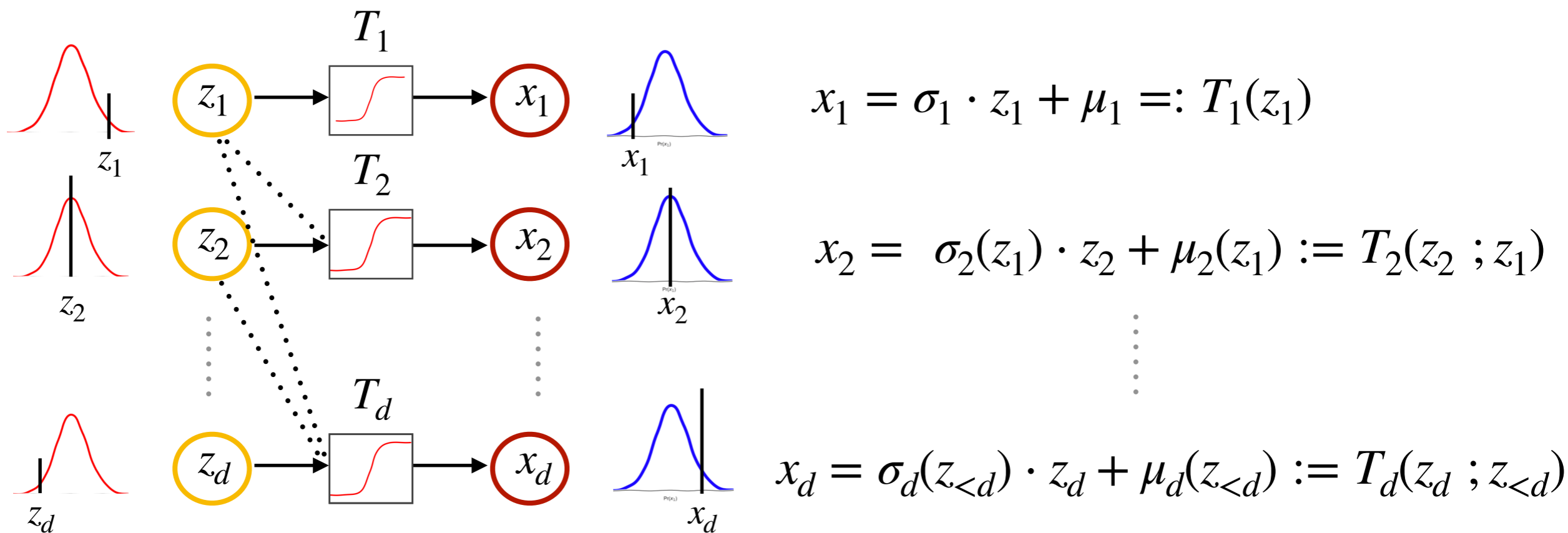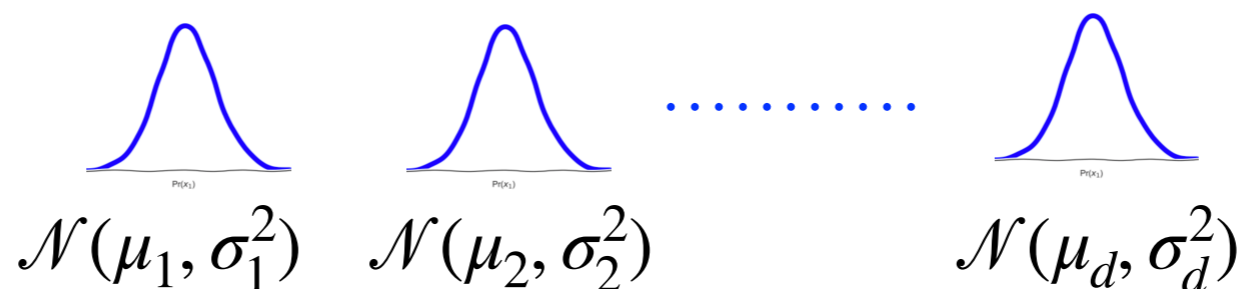$$p(x) = p_1(x_1) \cdot p_2(x_2 \,|\, x_1) \cdot \ldots \cdot p_d(x_d \,|\, x_{<d})$$



choosing a conditional implicitly fixes a family of triangular maps
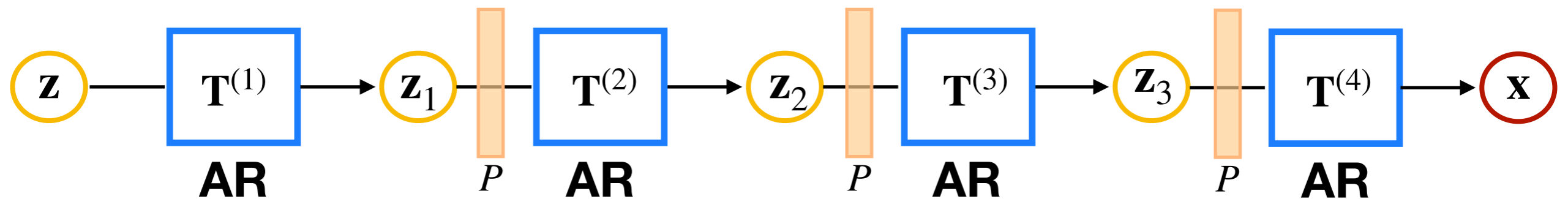
$$x_j = T_j \left( z_j; \ \theta_j(z_{<j}) \right)$$

*Larochelle, H and Murray, I. The Neural Autoregressive Distribution Estimator, AISTATS, 2011*
*Uria, et.al. Neural Autoregressive Distribution Estimation, JMLR, 2016*

# AR with Gaussian conditionals

$$p(x) = p_1(x_1) \cdot p_2(x_2 \,|\, x_1) \cdot \ldots \cdot p_d(x_d \,|\, x_{<d})$$



$$\mathcal{N}(\mu_1, \sigma_1^2) \quad \mathcal{N}(\mu_2, \sigma_2^2) \qquad \mathcal{N}(\mu_d, \sigma_d^2)$$



$$x_1 = \sigma_1 \cdot z_1 + \mu_1 =: T_1(z_1)$$

$$x_2 = \sigma_2(z_1) \cdot z_2 + \mu_2(z_1) := T_2(z_2 \,; z_1)$$

$$x_d = \sigma_d(z_{<d}) \cdot z_d + \mu_d(z_{<d}) := T_d(z_d \,; z_{<d})$$

*Kingma, et.al. Improved Variational Inference with Inverse Autoregressive Flow, NeurIPS, 2016*

# masked autoregressive flows (MAFs)

deep autoregressive flows with Gaussian conditionals*



$$(T_\# q)(x) = q(z) \cdot \left| \nabla \mathbf{T}^{(1)} \right|^{-1} \cdot \left| \nabla \mathbf{T}^{(2)} \right|^{-1} \cdot \left| \nabla \mathbf{T}^{(3)} \right|^{-1} \cdot \left| \nabla \mathbf{T}^{(4)} \right|^{-1}$$

$$x_j = z_j \cdot \exp\left(\alpha_j(z_{<j})\right) + \mu_j(z_{<j}) =: T_j(z_j \, ; z_{<j})$$

triangular maps are fundamental blocks for complex models

*Papamakarios, et.al. Masked Autoregressive Flow for Density Estimation, NeurIPS, 2017*

*\* can consider mixture of gaussians*

# real-NVP



$$T_j(z_j \; ; \; z_{<l}) = \exp\Big( \alpha_j(z_{<l}) \cdot \mathbf{1}_{j \notin [l-1]} \Big) \cdot z_j + \mu_j(z_{<l}) \cdot \mathbf{1}_{j \notin [l-1]}$$

*Dinh, et.al. Density estimation using Real NVP, ICLR, 2017*

# neural autoregressive flows (NAFs)



$$x_j = \textbf{DNN}\left( z_j \; ; \; \mathbf{w}_j(z_{<j}) \right) =: T_j(z_j \; ; z_{<j})$$

**Strictly positive weights & strictly monotonic
activation function ensure that the map is increasing**

**Universal**

*Huang, et.al. Neural Autoregressive Flows, ICML, 2018*

# sum-of-squares polynomial flows

goal : learn a **universal** increasing triangular function



$$\mathbb{P}_r(z; \mathbf{a}) = \sum_{l=0}^{r} a_{l,k} z^l$$

*[JSY]. Sum-of-squares Polynomial Flow. ICML, 2019*

# Summarize

| Model | conditioner $C_j$ output | $T_j(z_j ; C_j(z_1, \ldots, z_{j-1}))$ | ◉ | 👻 | 🏛 | Δ |
|---|---|---|---|---|---|---|
| Mixture (e.g. McLachlan & Peel, 2004) | $\boldsymbol{\theta}_j$ | $S_j(z_j; \boldsymbol{\theta}_j)$ | ✗ | ✗ | ✓ | I |
| (Bengio & Bengio, 1999) | $\boldsymbol{\theta}_j(z_{<j})$ | $S_j(z_j; \boldsymbol{\theta}_j)$ | ✗ | ✗ | ? | I |
| MADE (Germain et al., 2015) | $\boldsymbol{\theta}_j(z_{<j})$ | $S_j(z_j; \boldsymbol{\theta}_j)$ | ✓ | ✓ | ? | I |
| NICE (Dinh et al., 2015) | $\mu_j(z_{<l})$ | $z_j + \mu_j \cdot \mathbf{1}_{j \notin [l]}$ | ✗ | ✗ | ? | E |
| NADE (Uria et al., 2016) | $\boldsymbol{\theta}_j(z_{<j})$ | $S_j(z_j; \boldsymbol{\theta}_j)$ | ✓ | ✗ | ? | I |
| IAF (Kingma et al., 2016) | $\sigma_j(z_{<j}), \mu_j(z_{<j})$ | $\sigma_j z_j + (1 - \sigma_j)\mu_j$ | ✓ | ✓ | ? | E |
| MAF (Papamakarios et al., 2017) | $\alpha_j(z_{<j}), \mu_j(z_{<j})$ | $z_j \exp(\alpha_j) + \mu_j$ | ✓ | ✓ | ? | E |
| Real-NVP (Dinh et al., 2017) | $\alpha_j(z_{<l}), \mu_j(z_{<l})$ | $\exp(\alpha_j \cdot \mathbf{1}_{j \notin [l]}) \cdot z_j + \mu_j \cdot \mathbf{1}_{j \notin [l]}$ | ✗ | ✗ | ? | E |
| NAF (Huang et al., 2018) | $\mathbf{w}_j(z_{<j})$ | $\text{DNN}(z_j ; \mathbf{w}_j)$ | ✓ | ✓ | ✓ | E |
| SOS | $\mathbf{a}_j(z_{<j})$ | $\mathfrak{P}_{2r+1}(z_j; \mathbf{a}_j)$ | ✓ | ✓ | ✓ | E |

*[JSY]. Sum-of-squares Polynomial Flow. ICML, 2019*

# Toy examples

Germain, et.al. MADE: Masked Autoencoder for Density Estimation, ICML, 2015
Papamakarios, et.al. Masked Autoregressive Flow for Density Estimation, NeurIPS, 2017
Oliva, et.al. Transformation Autoregressive Networks, ICML, 2018
Huang, et.al. Neural Autoregressive Flows, ICML, 2018
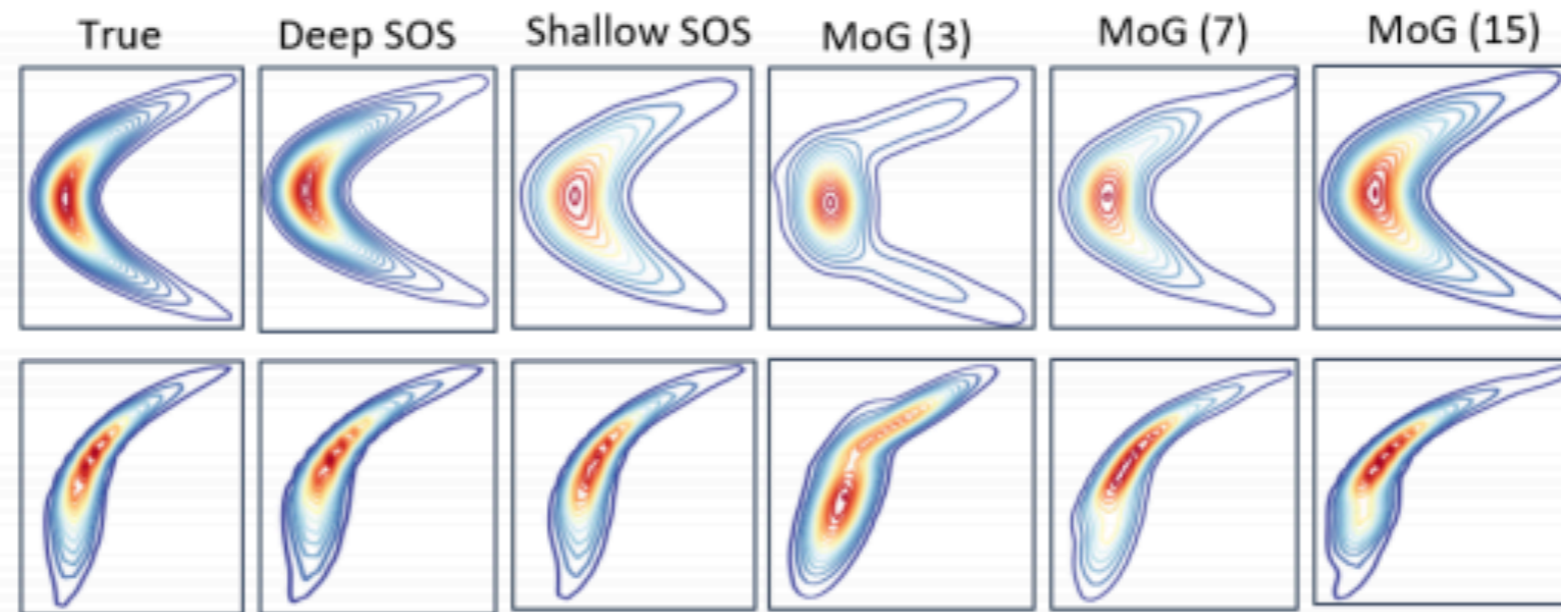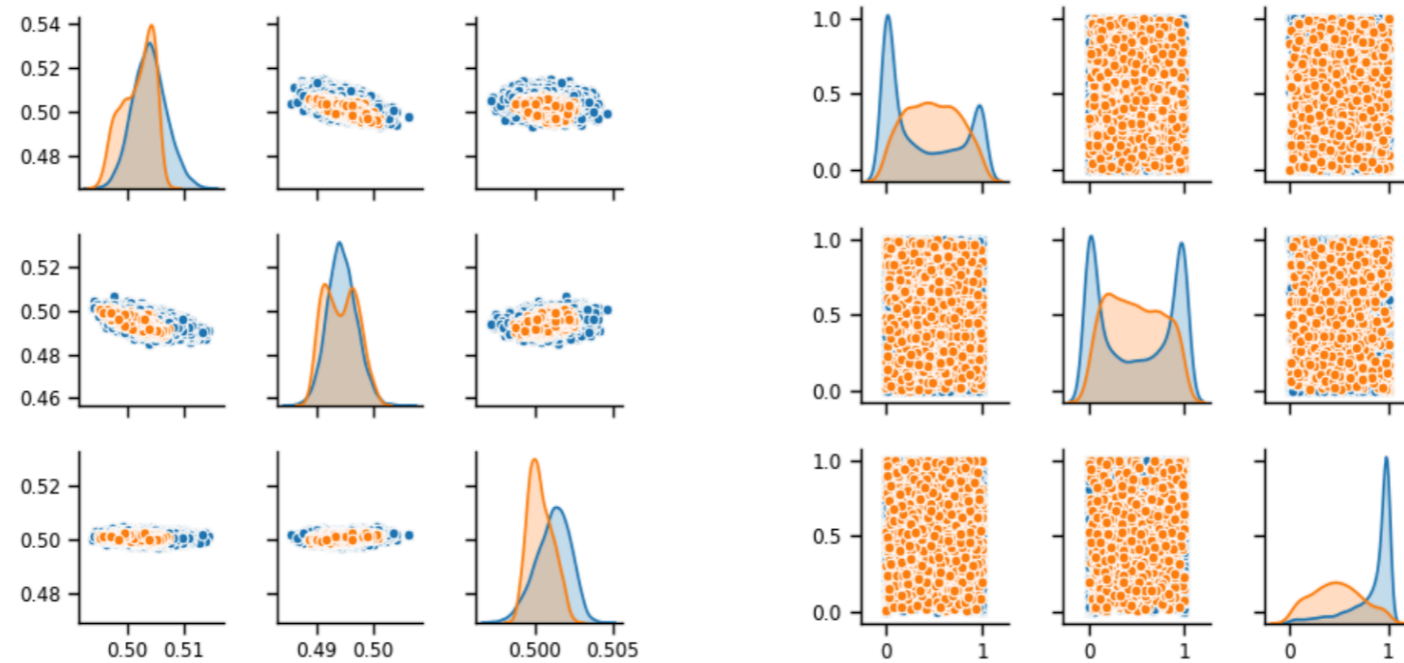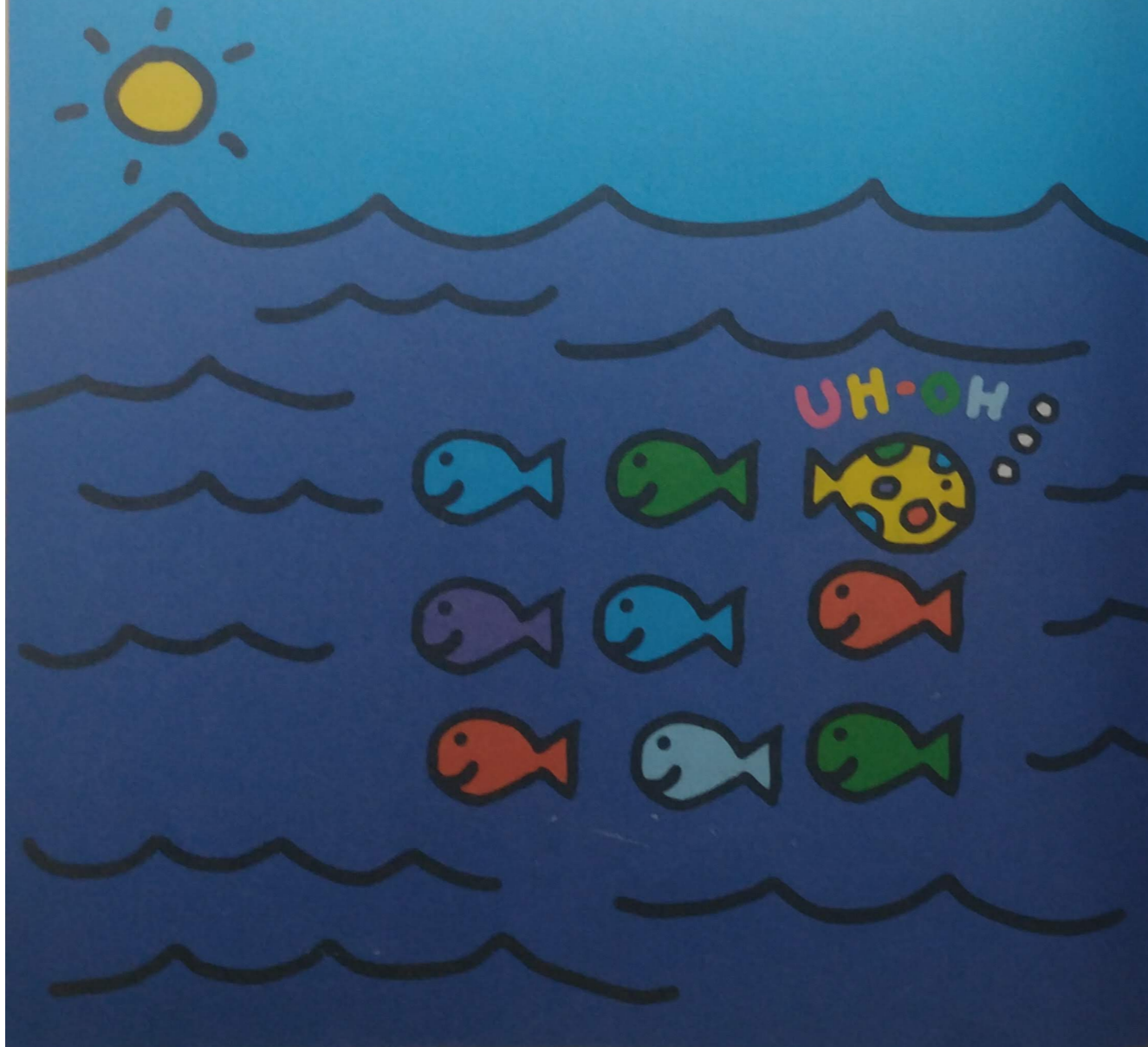
# Effect of ordering



*Figure 4.* **Top:** Left plot shows the target density given by $p(x_1, x_2) = \mathcal{N}(x_2 ; 0, 4)\mathcal{N}(x_1 ; 0.25x_2^2, 1)$. The second plot shows the density learnt by SOS flows with 3 blocks and a sum of 2 polynomials with degree 3 with ordering $(x_1, x_2)$. Third plot shows the density learnt by SOS flows with 1 block and a sum of 2 polynomials with degree 4 and ordering $(x_1, x_2)$. The last three plots estimate this density using a Mixture of Gaussian conditionals with varying components given in parenthesis and ordering $(x_1, x_2)$. **Bottom:** Same as Top but with target density given by $p(x_1, x_2) = \mathcal{N}(x_2 ; 2, 2)\mathcal{N}(x_1 ; 0.33x_1^3, 1.5)$.

# Application to novelty detection

## Multivariate triangular quantile maps

It's okay to try a different direction.

UH-OH

You might discover something new.

during training only nominal data is available.

# Two Approaches, One Idea

Novelty ≈ Low density region



Novelty $= [\![ -\hat{p} > -\alpha ]\!]$

Novelty $= [\![ \, | \hat{Q}^{-1} - \frac{1}{2} | > \alpha ]\!]$

*Ben-David and Lindenbaum. Learning Distributions by Their Density Levels: A Paradigm for Learning without a Teacher. JCSS 1997.*

*Steinwart, Hush and Scovel. A classification framework for anomaly detection. JMLR, 2005.*

*Schölkopf, Platt, Shawe-Taylor, Smola and Williamson. Estimating the Support of a High-Dimensional Distribution. Neural Computation, 2001.*

*Takeda and Sugiyama. v-Support Vector Machine as Conditional Value-at-Risk Minimization. ICML, 2008.*

# Triangular Quantile Map

Let $\mathbf{U} \sim \mathrm{Uniform}[0,1]^d$ and $\mathbf{X} \in \mathbb{R}^d$ any random vector. We call the increasing triangular map $\mathbf{Q} = \mathbf{Q_X} : [0,1]^d \to \mathbb{R}^d$ the triangular quantile map of $\mathbf{X}$ if $\mathbf{Q(U)} \sim \mathbf{X}$.

**Composable!**

Let $\mathbf{Y} = \mathbf{T(X)}$ for some increasing triangular map $\mathbf{T}$. Then, $\mathbf{Q_Y} = \mathbf{T} \circ \mathbf{Q_X}$.

- d=1: usual definition of quantile (inverse of cdf), advocated in (Parzen 1979)
- Precursors in Rosenblatt, Knothe, Ruschendorf, Decurninge …
- Other multivariate quantiles exist (e.g. Chernozhukov et al)

# One Stone, Two Birds

Novelty $\approx$ Low density region

Regularization

$$\min_{\mathbf{f}, \mathbf{Q}} \quad \gamma \mathrm{KL}(\mathbf{f}_\# p \| \mathbf{Q}_\# q) + \lambda \ell(\mathbf{f}) + \zeta g(\mathbf{Q})$$

density/quantile

dim reduction

[WSY]. Multivariate Triangular Quantiles for Novelty Detection. NeurIPS, 2019

# Implementation

$$\min_{\mathbf{f}, \mathbf{Q}} \quad \gamma \mathrm{KL}(\mathbf{f}_{\#}p \| \mathbf{Q}_{\#}q) + \lambda \ell(\mathbf{f}) + \zeta g(\mathbf{Q})$$

Parameterize Q using SOS flow

Solve by multiple gradient descent (no parameter tuning)

## Dimensionality reduction $\mathbf{Z} = \mathbf{f}(\mathbf{X})$

- **Density:** $\log |\mathbf{Q}'(\mathbf{Q}^{-1}(\mathbf{Z}))| + \|\mathbf{Q}^{-1}(\mathbf{Z})\|^2/2 > \alpha$

- **Quantile:** $\|\mathbf{Q}^{-1}(\mathbf{Z})) - \frac{1}{2}\|_{\infty} > \alpha$

  $\hookrightarrow$ **can be tuned**

# Donut

# Summary



$z \sim q(z)$

$T_1$

$z_1 \rightarrow x_1$

$T_2$

$z_2 \rightarrow x_2$

$T_d$

$z_d \rightarrow x_d$
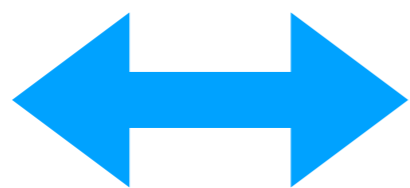
$x \sim p(x)$

**A probabilistic object**     deterministic map     **A probabilistic object**

**Complexity**         **Properties**

*[JKYB]. Tails of Lipschitz Triangular Flows. ICML, 2020*       *Spantini, Bigoni and Marzouk. Inference via low-dimensional couplings. JMLR, 2018*
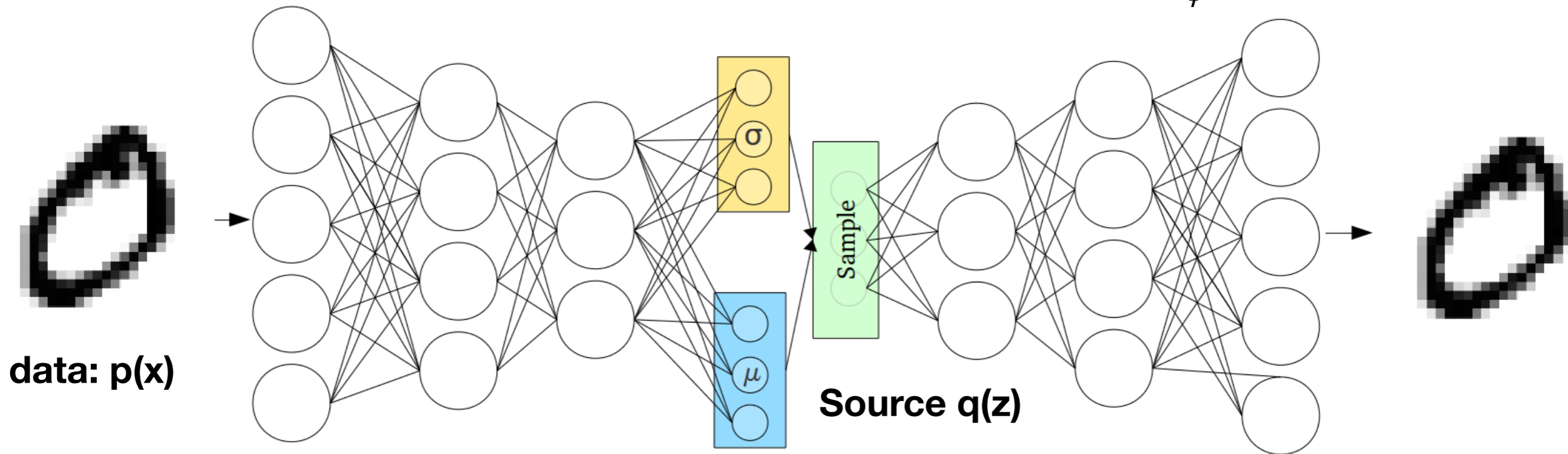
# Variational Auto-Encoder

Given data $\mathbf{x}_1, \ldots, \mathbf{x}_n$, estimate p(x)

$$\min_{\theta} \min_{\phi} \; \text{KL}\left[p(x)p_\theta(z|x) \middle\| q(z)q_\phi(x|z)\right]$$

**Encoder:** $p_\theta(z|x)$

**Decoder:** $q_\phi(x|z)$

σ

Sample

μ

**data: p(x)**

**Source q(z)**

*Kingma, D and Welling, M. Auto-Encoding Variational Bayes, ICLR, 2014*

*Rezende, D. et.al. Stochastic backpropagation and approximate inference in deep generative models, ICML, 2014*
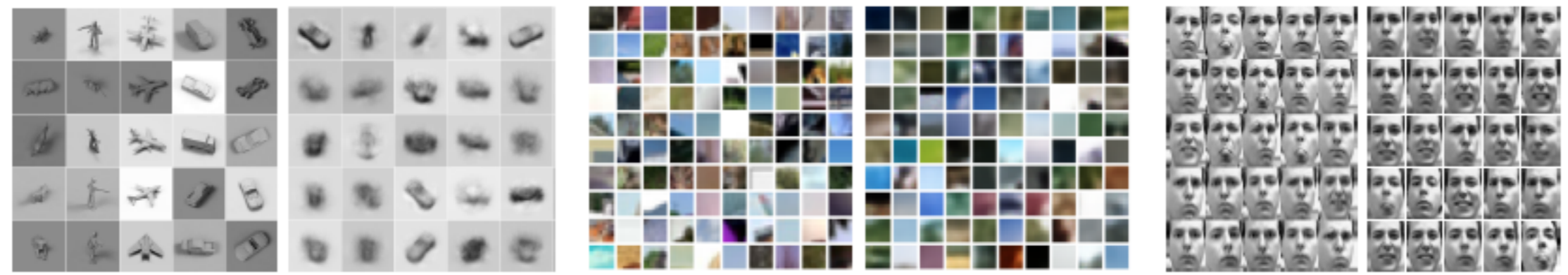
# VAE examples



(a) Learned Frey Face manifold

(b) Learned MNIST manifold

(a) NORB

(b) CIFAR

(c) Frey

*Kingma, D and Welling, M. Auto-Encoding Variational Bayes, ICLR, 2014*

*Rezende, D. et.al. Stochastic backpropagation and approximate inference in deep generative models, ICML, 2014*