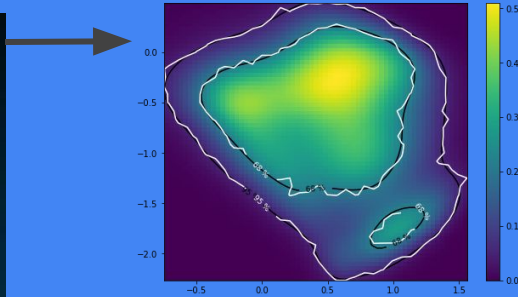# Normalizing flows
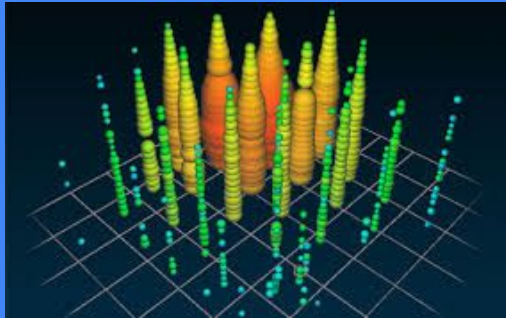
**An overview and its revolutionary potential for high-energy physics when combined with deep learning**



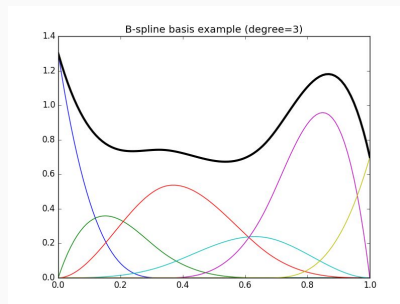Thorsten Glüsenkamp, 22.4. 2022, IoP Machine Learning workshop

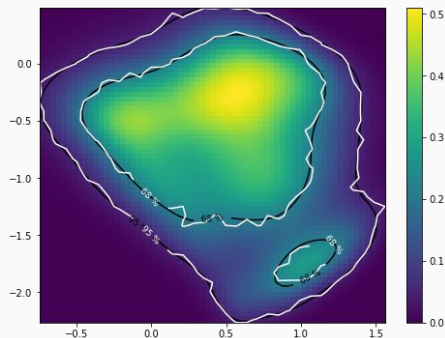Some well-known techniques to construct PDFs:



Mixture models /
Kernel Density Functions



B-Spline representation of a PDF

- (arbitrarily) Complex PDF shape
- Evaluate probability analytically
- Works in D > 1

**Normalizing flows are also PDFs, but richer functionality:**
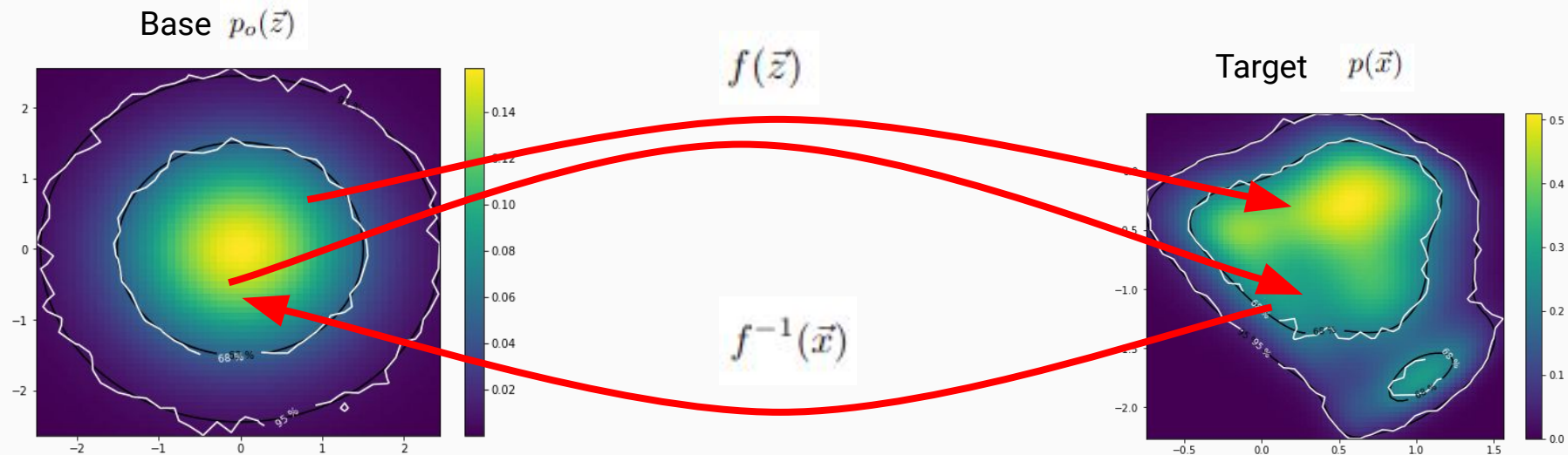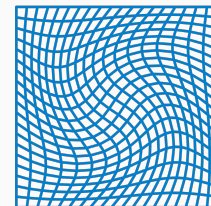


- (arbitrarily) Complex PDF shape
- Evaluate probability analytically
- Works in D > 1
- Generate differentiable samples (-> differentiable expectation values)
- Coverage of the PDF
- Can be interpreted as generalizations of the Gaussian distribution
- Works on manifolds
- …

2

Base $p_o(\vec{z})$

Target $p(\vec{x})$

$f(\vec{z})$

Flow defining function

$$p(\vec{x}) = p_0(f^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f^{-1}(\vec{x})}{\partial \vec{x}} \right|$$

Base $p_o(\vec{z})$

$f(\vec{z})$

Target $p(\vec{x})$

$f^{-1}(\vec{x})$

1) The function has to be **bijective** for the inverse!
2) The function has to be **differentiable!**

= Diffeomorphism

4

$$p(\vec{x}) = p_0(f^{-1}(\vec{x})) \cdot \left| \det J^{-1}(\vec{x}) \right|$$

Log probability includes
Local "stretching"/"squeezing"
Of volume element

Base $p_o(\vec{z})$

$f(\vec{z})$

Target $p(\vec{x})$

$f^{-1}(\vec{x})$
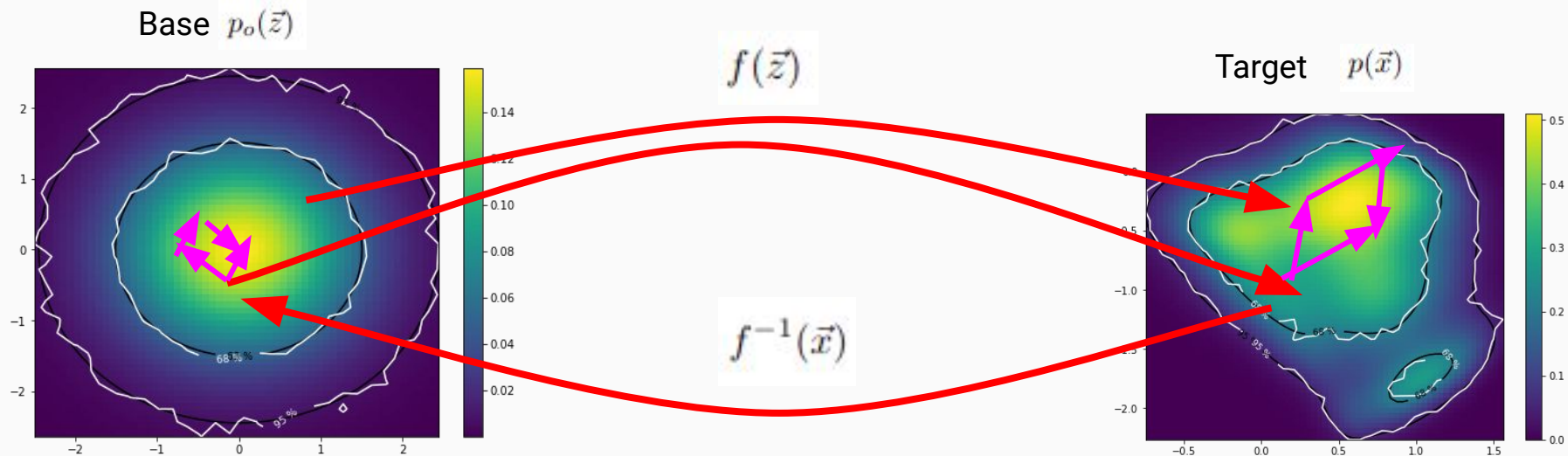


1) The function has to be **bijective** for the inverse!
2) The function has to be **differentiable!**

= Diffeomorphism

5

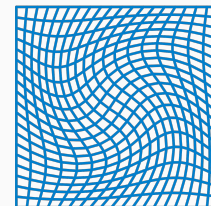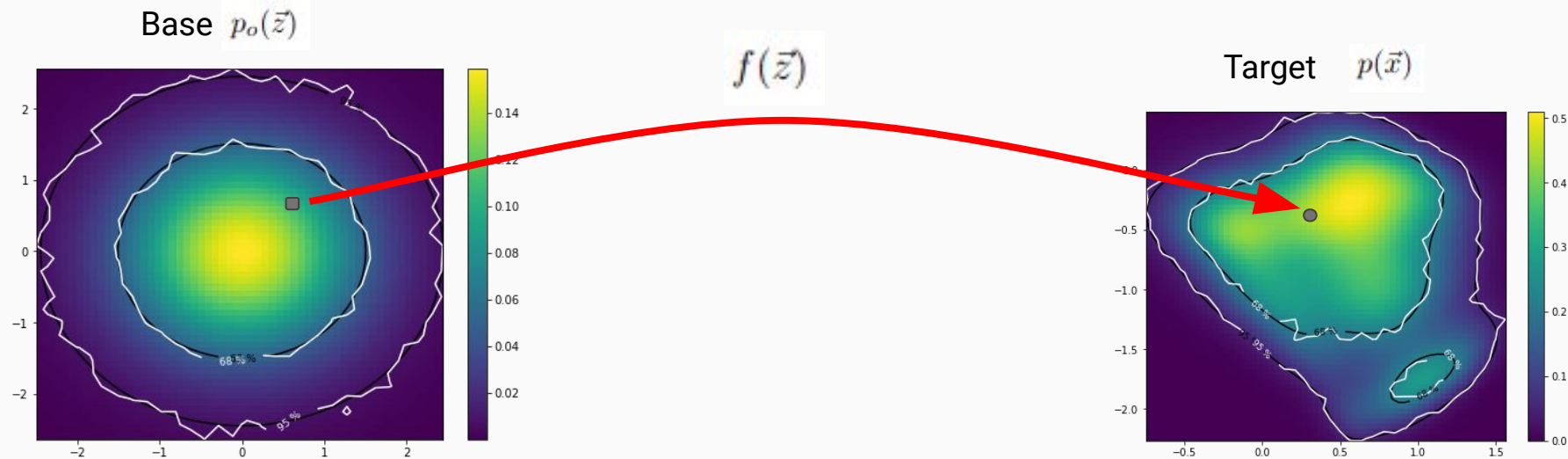Sampling is also straight forward: $$\vec{x} = f(\vec{z})$$

**Must be able to sample from base distribution!**

Base $p_o(\vec{z})$

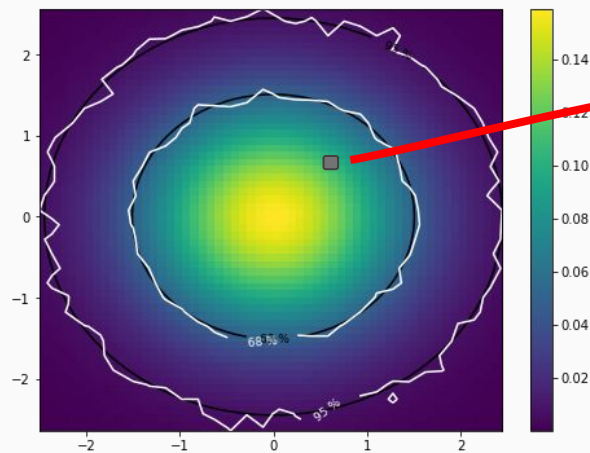$f(\vec{z})$

Target $p(\vec{x})$

Sampling is also straight forward:

**Must be able to sample from base distribution!**

$$\vec{x_\theta} = f_\theta(\vec{z})$$

$$p_\theta(\vec{x}) = p_0(f_\theta^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\vec{x})}{\partial \vec{x}} \right|$$

Base $p_o(\vec{z})$

$f_\theta(\vec{z})$

Target $p_\theta(\vec{x})$

Sampling is also straight forward:

**Must be able to sample from base distribution!**

$$\vec{x_\theta} = f_\theta(\vec{z})$$

**1) Differentiable prob.**

$$p_\theta(\vec{x}) = p_0(f_\theta^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\vec{x})}{\partial \vec{x}} \right|$$

Base $p_o(\vec{z})$

Target $p_\theta(\vec{x})$

$$f_\theta(\vec{z})$$



**2) Differentiable samples!**

Sampling is also straight forward:

**Must be able to sample from base distribution!**

$$\vec{x}_\theta = f_\theta(\vec{z})$$
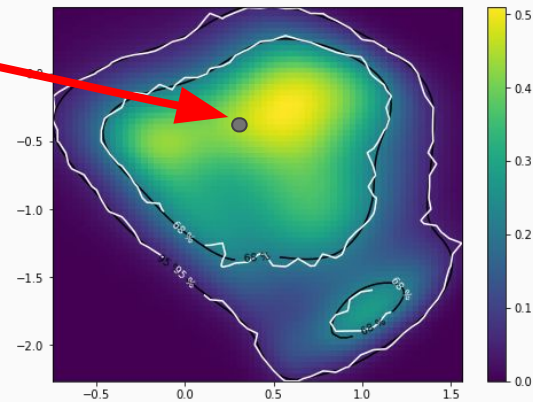
**1) Differentiable prob.**

$$p_\theta(\vec{x}) = p_0(f_\theta^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\vec{x})}{\partial \vec{x}} \right|$$
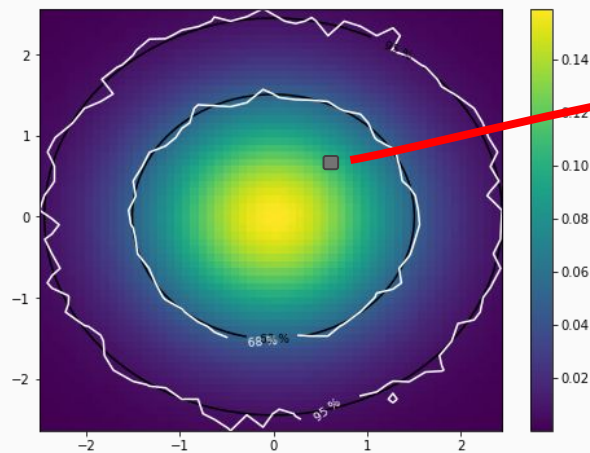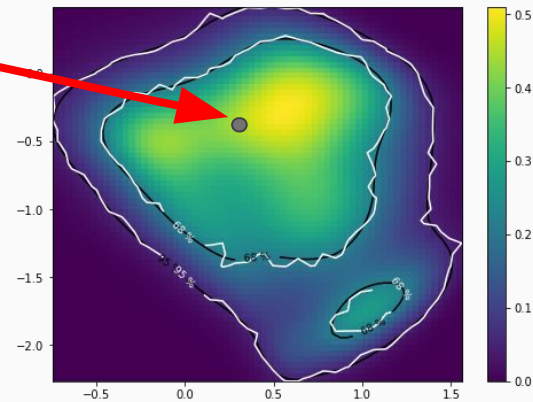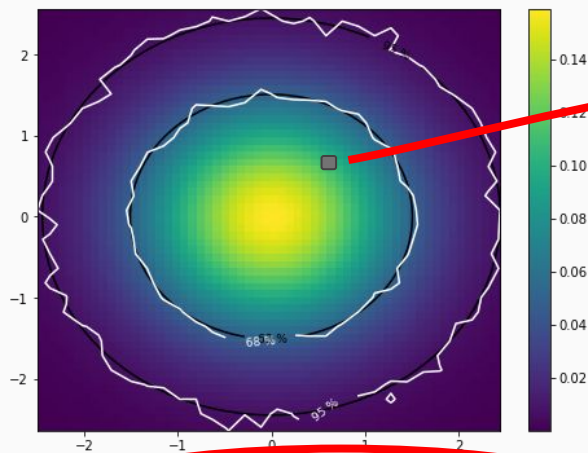
Base $p_o(\vec{z})$

$f_\theta(\vec{z})$

Target $p_\theta(\vec{x})$



**2) Differentiable samples!** **Unique and important!**

9

$$I_\theta = \int p_\theta(x) F_\theta(x) dx \approx \frac{1}{N} \cdot \sum_i^N F_\theta(x_\theta)$$

$$I_\theta = \int p_\theta(x) F_\theta(x) dx \approx \frac{1}{N} \cdot \sum_i^N F_\theta(x_\theta)$$

Examples:

n-th moment of p

$$\int p_\theta(x) x^n dx \approx \frac{1}{N} \cdot \sum_i^N x_\theta^n$$

entropy

$$-\int p_\theta(x) \log(p_\theta(x)) dx \approx \frac{1}{N} \cdot \sum_i^N -\log p_\theta(x_\theta)$$

…. Many other integrals in information theory

$$I_\theta = \int p_\theta(x) F_\theta(x) dx \approx \frac{1}{N} \cdot \sum_i^N F_\theta(x_\theta)$$

Examples:

n-th moment of p

$$\int p_\theta(x) x^n dx \approx \frac{1}{N} \cdot \sum_i^N x_\theta^n$$

entropy

$$-\int p_\theta(x) \log (p_\theta(x))\, dx \approx \frac{1}{N} \cdot \sum_i^N -\log p_\theta(x_\theta)$$

**Wrong gradient!**

.... Many other integrals in information theory

$$\frac{1}{N} \cdot \sum_i^N -\log p_\theta(x)$$

!

Let us try to see the distribution of the simplest possible normalizing flow in 1 dimension!

$$x = f(\vec{z}) =$$

Let us try to see the distribution of the simplest possible normalizing flow in 1 dimension!

$$x = f(\vec{z}) = a \cdot z + b \qquad \theta = \{a, b\}$$

$$p_\theta(\vec{x}) = p_0(f_\theta^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\vec{x})}{\partial \vec{x}} \right|$$

Let us try to see the distribution of the simplest possible normalizing flow in 1 dimension!

$$x = f(\vec{z}) = a \cdot z + b$$

Use Standard normal base distribution
(we will always use the standard normal
for convenience)

Base distribution

$$p_0(z) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot \exp(-0.5 \cdot z^2)$$

Change of variables formula

$$p_\theta(\vec{x}) = p_0(f_\theta^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\vec{x})}{\partial \vec{x}} \right|$$

Inv. Flow function

$$z = f^{-1}(\vec{x}) = \frac{x - b}{a}$$

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi \cdot a^2}} \cdot \exp\left(-0.5 \cdot \left(\frac{x - b}{a}\right)^2\right)$$

Let us try to see the distribution of the simplest possible normalizing flow in 1 dimension!

$$x = f(\vec{z}) = a \cdot z + b$$

Use Standard normal base distribution
(we will always use the standard normal
for convenience)

"Linear flow" = "general Gaussian distribution" in 1d

Base distribution

$$p_0(z) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot \exp(-0.5 \cdot z^2)$$

Change of variables formula

$$p_\theta(\vec{x}) = p_0(f_\theta^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\vec{x})}{\partial \vec{x}} \right|$$

Inv. Flow function

$$z = f^{-1}(\vec{x}) = \frac{x - b}{a}$$

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi \cdot a^2}} \cdot \exp\left(-0.5 \cdot \left(\frac{x-b}{a}\right)^2\right)$$

16

Let us try to see the distribution of the simplest possible normalizing flow in 1 dimension!
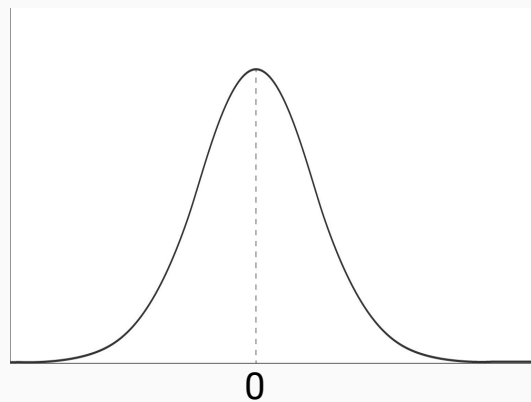
$$x = f(\vec{z}) = a \cdot z + b$$

Use Standard normal base distribution
(we will always use the standard normal
for convenience)

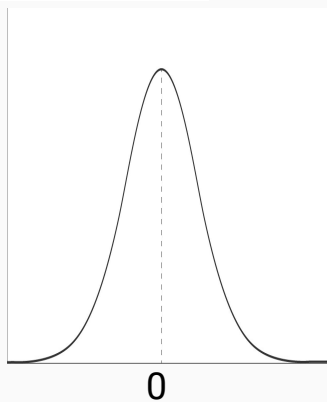"Linear flow" = "general Gaussian distribution" in 1d

Technically, this is 2-step flow:

$$f(z) = f_2(f_1(z))$$



$$f_1(z) = a \cdot z$$
"scale"

$$f_2(z) = z + b$$
"shift"

0                                    0                                    b
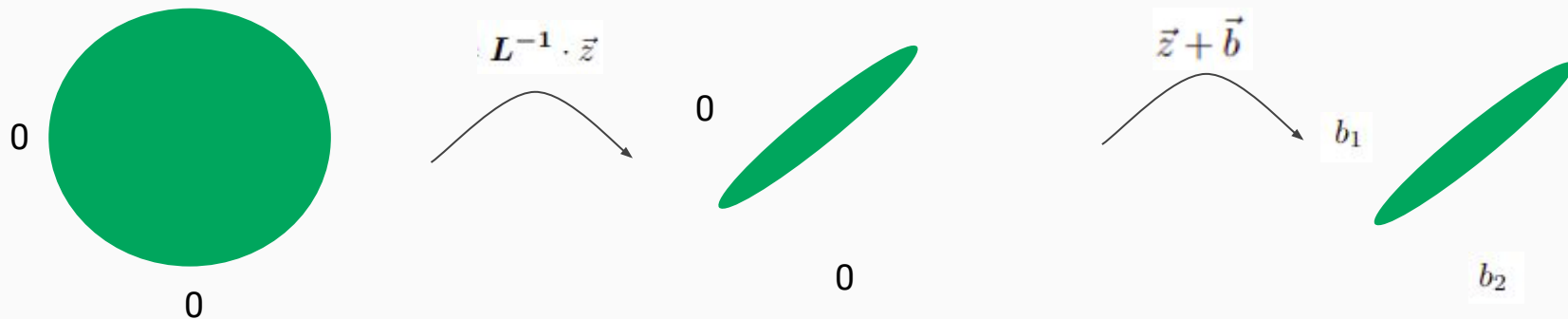
Let us try to see the distribution of the simplest possible normalizing flow in n dimensions!

Use Standard normal base distribution (we will always use the standard normal for convenience)

$$\vec{x} = f(\vec{z}) = L^{-1} \cdot \vec{z} + \vec{b}$$

$$C^{-1} = L^T \cdot L$$

"Affine flow" = "general Gaussian distribution" in n-d



$L^{-1} \cdot \vec{z}$

$\vec{z} + \vec{b}$

$b_1$

$b_2$

One possibility:
Neural Spline Flows (Durkan et al. 2019)

$$p(\vec{x}) = p_0(f^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f^{-1}(\vec{x})}{\partial \vec{x}} \right|$$
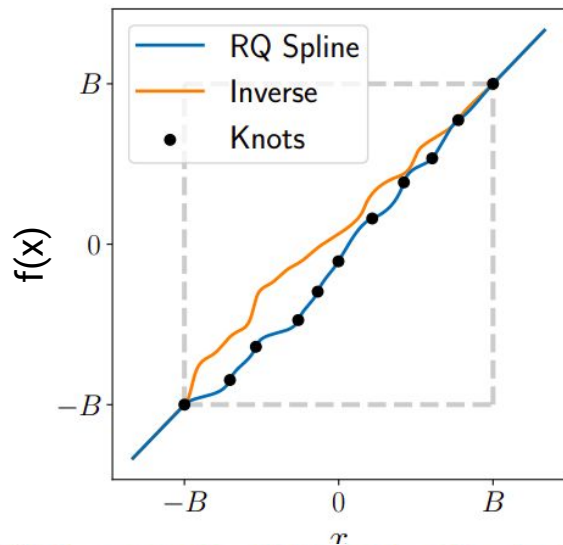


RQ Spline
Inverse
Knots

**Flow function:** $f^{-1}(\vec{x})$

$$\frac{\alpha^{(k)}(\xi)}{\beta^{(k)}(\xi)} = y^{(k)} + \frac{(y^{(k+1)} - y^{(k)})\left[s^{(k)}\xi^2 + \delta^{(k)}\xi(1-\xi)\right]}{s^{(k)} + \left[\delta^{(k+1)} + \delta^{(k)} - 2s^{(k)}\right]\xi(1-\xi)}$$
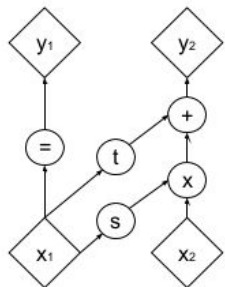
**Derivative:** $\dfrac{\partial f^{-1}(\vec{x})}{\partial \vec{x}}$

$$\frac{\mathrm{d}}{\mathrm{d}x}\left[\frac{\alpha^{(k)}(\xi)}{\beta^{(k)}(\xi)}\right] = \frac{(s^{(k)})^2\left[\delta^{(k+1)}\xi^2 + 2s^{(k)}\xi(1-\xi) + \delta^{(k)}(1-\xi)^2\right]}{\left[s^{(k)} + \left[\delta^{(k+1)} + \delta^{(k)} - 2s^{(k)}\right]\xi(1-\xi)\right]^2}$$
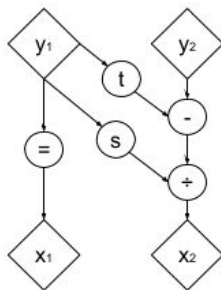
which passes through the knots, with the given derivatives at the knots. Defining $s_k = (y^{k+1} - y^k)/(x^{k+1} - x^k)$ and $\xi(x) = (x - x^k)/(x^{k+1} - x^k)$, the expression for the rational-

**1)** "Affine Coupling layer" (Dinh et al. 2014 (NICE) / Dinh et al. 2017 (real-NVP))



(a) Forward propagation



(b) Inverse propagation

$$y_{1:d} = x_{1:d}$$
$$y_{d+1:D} = x_{d+1:D} \odot \exp\left(s(x_{1:d})\right) + t(x_{1:d}),$$

$$\frac{\partial y}{\partial x^T} = \begin{bmatrix} \mathbb{I}_d & 0 \\ \frac{\partial y_{d+1:D}}{\partial x_{1:d}^T} & \mathrm{diag}\left(\exp\left[s\left(x_{1:d}\right)\right]\right) \end{bmatrix},$$
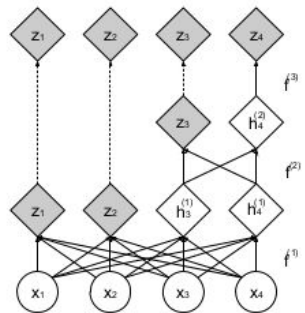
Induces coupling between different layers

Functions **s** and **t** can be arbitrarily complex (neural networks)
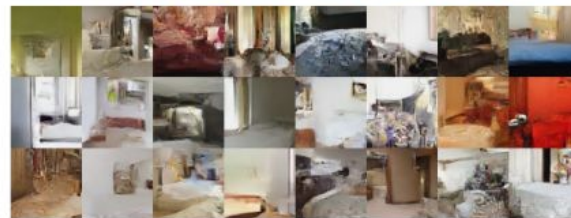
These functions are also called "conditioners"

**1)** "Affine Coupling layer" (Dinh et al. 2014 (NICE) / Dinh et al. 2017 (real-NVP))



(b) Factoring out variables. At each step, half the variables are directly modeled as Gaussians, while the other half undergo further transformation.
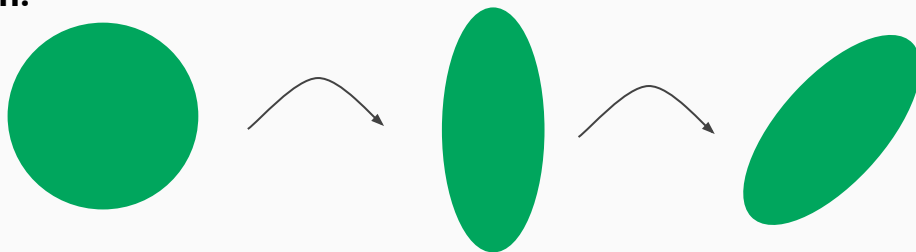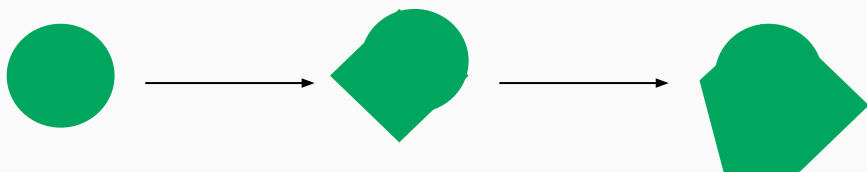
Can be scaled to very high dimensions (here images)

**4)** Gaussianization Flow(arXiv:2003.01941)

**Gaussian:**



**Gaussianization flow generalized by non-linear scalings instead of "linear scalings"**

$$T_{\boldsymbol{\theta}}(\mathbf{x}) = \Psi_{\boldsymbol{\theta}_L} \circ R_L \circ \Psi_{\boldsymbol{\theta}_{L-1}} \circ \cdots \circ \Psi_{\boldsymbol{\theta}_1} \circ R_1 \mathbf{x}$$



...
....

Step by step refinement of the PDF - provably approximates any distribution

## 4) Gaussianization Flow(arXiv:2003.01941)

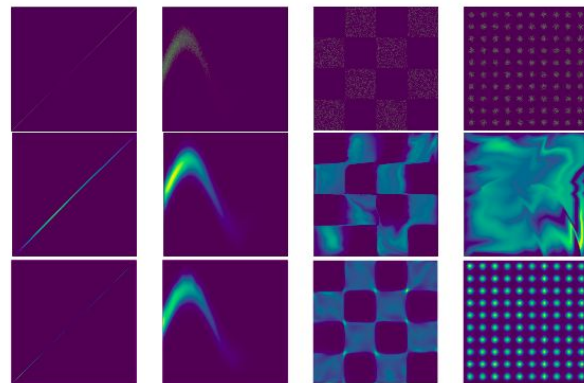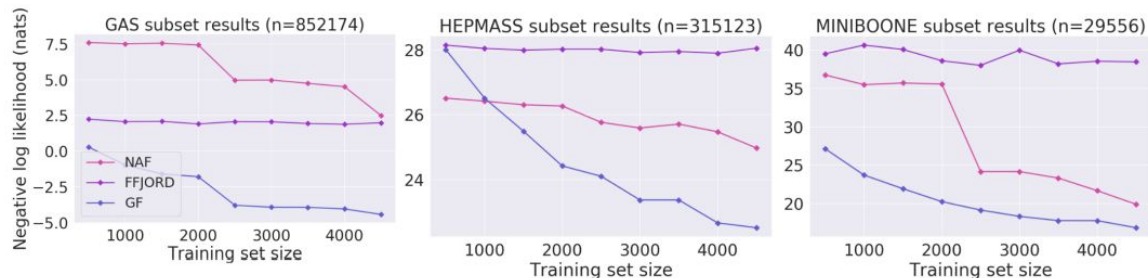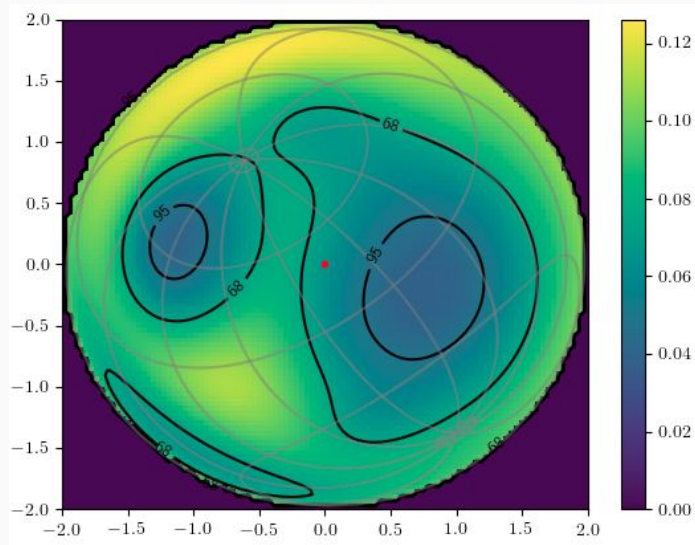**Can also be initialized to quite good values by data!**





Figure 2: 2D density estimation results. **Top:** Ground truth samples. **Middle:** Glow. **Bottom:** GF.

| Method | POWER | GAS | HEPMASS | MINIBOONE | BSDS300 | MNIST | FMNIST |
|---|---|---|---|---|---|---|---|
| Real NVP | -0.17 | -8.33 | 18.71 | 13.55 | -153.28 | 1.06 | **2.85** |
| Glow | -0.17 | -8.15 | 18.92 | 11.35 | -155.07 | 1.05 | 2.95 |
| FFJORD | -0.46 | -8.59 | **14.92** | 10.43 | **-157.40** | **0.99** | - |
| RBIG | 1.02 | 0.05 | 24.59 | 25.41 | -115.96 | 1.71 | 4.46 |
| GF(ours) | **-0.57** | **-10.13** | 17.59 | **10.32** | -152.82 | 1.29 | 3.35 |
| MADE | 3.08 | -3.56 | 20.98 | 15.59 | -148.85 | 2.04 | 4.18 |
| MAF | -0.24 | -10.08 | 17.70 | 11.75 | -155.69 | 1.89 | - |
| TAN | -0.48 | -11.19 | 15.12 | 11.01 | -157.03 | - | - |
| MAF-DDSF | -0.62 | -11.96 | 15.09 | 8.86 | -157.73 | - | - |

"Manifold" normalizing flows

"Continuous" normalizing flows

$$p(T(x)) = \frac{\pi(x)}{\sqrt{\det(E^\top J^\top J E)}}$$

$$\log p(\mathbf{z}(t_1)) = \log p(\mathbf{z}(t_0)) - \int_{t_0}^{t_1} \mathrm{Tr}\left(\frac{\partial f}{\partial \mathbf{z}(t)}\right) dt.$$
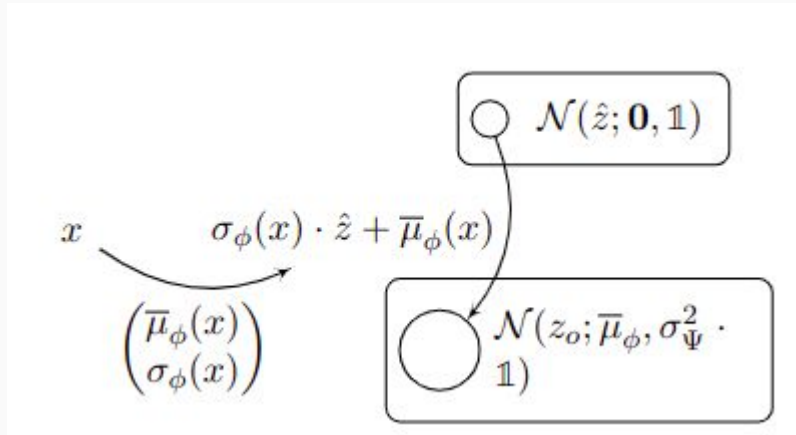


(FFJORD, 1810.01367)

(2002.02428)

**p(x) -> p(x;y) ?**

**Connection to neural networks...**

### Conditional affine flow



$$\mathcal{N}(\hat{z}; \mathbf{0}, \mathbb{1})$$

$$\sigma_\phi(x) \cdot \hat{z} + \overline{\mu}_\phi(x)$$

$$x$$

$$\begin{pmatrix} \overline{\mu}_\phi(x) \\ \sigma_\phi(x) \end{pmatrix}$$

$$\mathcal{N}(z_o; \overline{\mu}_\phi, \sigma_\Psi^2 \cdot \mathbb{1})$$

### Conditional general flow



$$\mathcal{N}(\hat{z}; \mathbf{0}, \mathbb{1})$$

$$x \qquad \rho_\phi(\hat{z}; x)$$

$$\overline{F}_\phi(x)$$

$$q_\phi(z_o; x)$$

Instead of flow parameters, one optimizes NN parameters
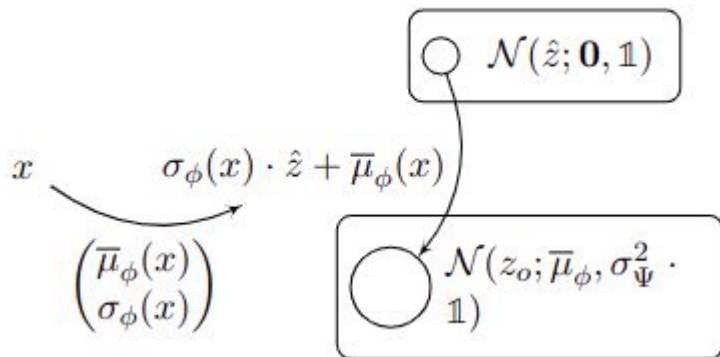
Conditional affine flow



Conditional general flow



→ Instead of flow parameters, one optimizes NN parameters

→ Conditional normalizing flow shows that **MSE loss comes from conditional Flow that only consists of a shift** (and unit scaling) $0.5 \cdot (x - \mu)^2 = \ln(p(x))$

Conditional affine flow

$$\mathcal{N}(\hat{z}; \mathbf{0}, \mathbb{1})$$

$$x \qquad \sigma_\phi(x) \cdot \hat{z} + \overline{\mu}_\phi(x)$$

$$\begin{pmatrix} \overline{\mu}_\phi(x) \\ \sigma_\phi(x) \end{pmatrix}$$

$$\mathcal{N}(z_o; \overline{\mu}_\phi, \sigma_\Psi^2 \cdot \mathbb{1})$$

Conditional general flow

$$\mathcal{N}(\hat{z}; \mathbf{0}, \mathbb{1})$$

$$x \qquad \rho_\phi(\hat{z}; x)$$
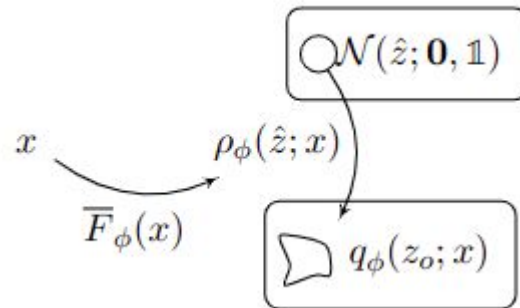
$$\overline{F}_\phi(x)$$

$$q_\phi(z_o; x)$$

→ Instead of flow parameters, one optimizes NN parameters

→ Conditional normalizing flow shows that **MSE loss comes from conditional Flow that only consists of a shift** (and unit scaling) $\quad 0.5 \cdot (x - \mu)^2 = \ln(p(x))$

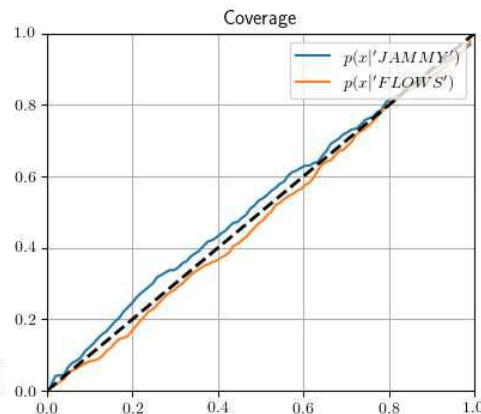→ The process of **predicting parameters by a neural network** is also called "**amortization**" 28

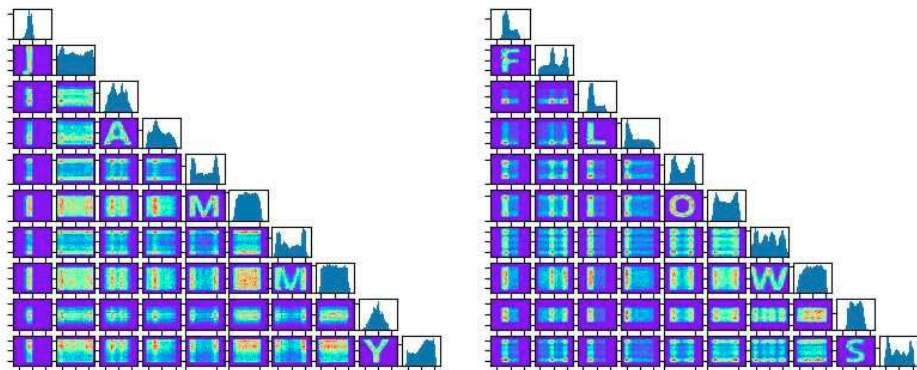**import jammy_flows**

**pdf=jammy_flows.pdf("e4+s2+e4", "gggg+n+gggg")**

**pdf.sample(nsamples=1000)**

A package to describe amortized (conditional) normalizing-flow PDFs defined jointly on tensor products of manifolds with coverage control. The connection between different manifolds is fixed via an autoregressive structure.
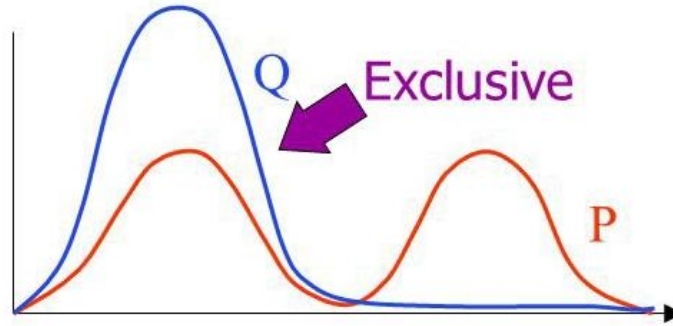


29

**They allow to do statistical analysis with probabilistic machine learning...**

"Inclusive" KL-diverengce vs "exclusive" KL-divergence
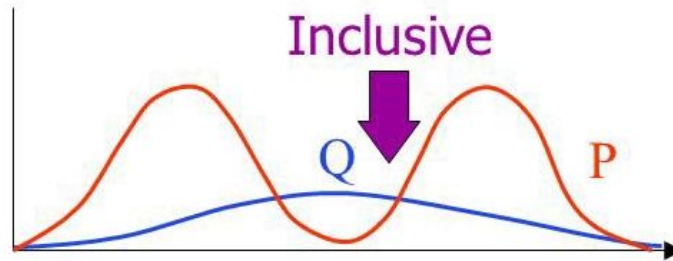
Minimising

$$\mathrm{KL}(Q||P)$$

$$= \sum_H Q(H) \ln \frac{Q(H)}{P(H|V)}$$

Minimising

$$\mathrm{KL}(P||Q)$$

$$= \sum_H P(H|V) \ln \frac{P(H|V)}{Q(H)}$$

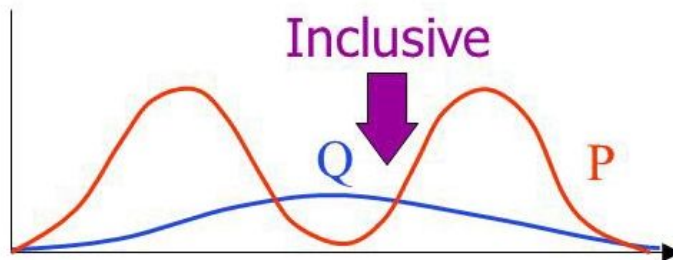"Inclusive" KL-diverengce vs "exclusive" KL-divergence



In most settings:

P = "True" PDF
(not accessible,
"Nature",
Samples from MC simulation
Draw from P)

Q = "Approximating PDF",
Parametrized by us,
"Surrogate model"

What is a Monte Carlo simulation?
Samples from some "true" distribution



$$\arg\min_{\phi} \hat{D}_{\mathrm{KL,joint(x,z_o)}}(\mathcal{P}_t; q_\phi) = \arg\min_{\phi} \frac{1}{N} \sum_{S \in x_i, z_{o,i}} \ln\left(\frac{\mathcal{P}_t(z_{o,i}; x_i)}{q_\phi(z_{o,i}; x_i)}\right) + \ln\left(\frac{\mathcal{P}_t(x_i)}{q(x_i)}\right)$$

$$= \arg\min_{\phi} \frac{1}{N} \sum_{S \in x_i, z_{o,i}} -\ln\left(q_\phi(z_{o,i}; x_i)\right)$$

Computer scientists call this
"**conditional ML objective**"
Of supervised learning

Physicists should call it
"**variational inference objective**"
For the **variational Posterior**

What is a Monte Carlo simulation?
Samples from some "true" distribution

$$\underset{\phi}{\arg\min} \, \hat{D}_{\mathrm{KL,joint(x,z_o)}}(\mathcal{P}_t; q_\phi) = \underset{\phi}{\arg\min} \, \frac{1}{N} \sum_{S \in x_i, z_{o,i}} \ln\left(\frac{\mathcal{P}_t(z_{o,i}; x_i)}{q_\phi(z_{o,i}; x_i)}\right) + \ln\left(\frac{\mathcal{P}_t(x_i)}{q(x_i)}\right)$$

$$= \underset{\phi}{\arg\min} \, \frac{1}{N} \sum_{S \in x_i, z_{o,i}} -\ln\left(q_\phi(z_{o,i}; x_i)\right)$$

Sample from "systematics distribution" during MC generation

Computer scientists call this
"**conditional ML objective**"
Of supervised learning

*Including systematics is trivial!*

Physicists should call it
"**variational inference objective**"
For the **variational Posterior**

"Bayesian"

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)}$$

MCMC

"Frequentist"

Maximum likelihood principle

KL-Div posterior                                    KL-Div likelihood

"Bayesian"                                          "Frequentist"

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)}$$

(90's)                                              (Aikake 70's)

MCMC            Variational inference               Maximum likelihood principle

KL-Div posterior             Joint KL              KL-Div likelihood
                             (deep learning)

"Bayesian"                                          "Frequentist"

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)}$$

              (90's)                                (Aikake 70's)

MCMC        Variational inference        Maximum likelihood principle

"All of deep learning is probability distribution matching"



| supervised learning | ext. supervised learning | VAE learning |
|---|---|---|
| **Input:** data $x_i$, labels $z_{o,i}$ | **Input:** data $x_i$, labels $z_{o,i}$ | **Input:** data $x_i$ |
| $D_{\mathrm{KL}}(\mathcal{P}_t \| q)[x, z_o]$ | $D_{\mathrm{KL}}(\mathcal{P}_t \| q)[x, z_o, z_u]$ | $D_{\mathrm{KL}}(\mathcal{P}_t \| q)[x, z_u]$ |
| VI posterior (via supervised loss): $q_\phi(z_o; x) \approx \mathcal{P}_t(z_o; x)$ | VI posterior (via mixed loss): $q_{\phi,\varphi}(z_o, z_u; x) \sim \mathcal{P}_t(z_o, z_u; x)$ | VI posterior (via ELBO): $q_\phi(z_u; x) \sim \mathcal{P}_t(z_u; x)$ |

"Inklusive KL divergence"

"exclusive KL divergence"

(2008.05825)

Minimising
KL*(P||Q)*
$= \sum_H P(H|V) \ln \frac{P(H|V)}{Q(H)}$

Inclusive

$$\rho_{\Psi,\text{intrinsic}}$$

$$\frac{K(\theta_1,\ldots,\theta_{d-2})}{S_d} \cdot \sin^{d-1}(\theta_{d-1})$$

$$q_\Psi(\theta_1,\ldots,\theta_{d-1},\phi_d)$$

$$\rho_{\text{tot}} = \rho_2 \circ \rho_1$$

$$\theta_{d-1}$$

$$\rho_2(r_f) = \theta_{d-1} = \arccos\left(\frac{r_f^2-1}{r_f^2+1}\right)$$

$$\frac{K(\theta_1,\ldots,\theta_{d-2})}{(2\pi)^{d/2}} \cdot r_g^{d-1} \cdot \exp\left(-\frac{r_g^2}{2}\right)$$

$$\frac{K(\theta_1,\ldots,\theta_{d-2})}{S_d} \cdot \left(\frac{2}{r_f^2+1}\right)^d \cdot r_f^{d-1}$$

$$\rho_1(r_g) = r_f = \text{CDF}_{r,f}^{-1}\left(\text{CDF}_{r,g}(r_g)\right)$$

$$-2 \cdot \left(\ln p_b(\hat{z}_i) - \ln p_b(0)\right) \sim \chi^2$$

$$p_b(\hat{z}) \qquad \hat{z}_i \cdot \qquad \rho^{-1}(z_{o,i}) \qquad z_{o,i} \cdot \qquad q_\phi(z_o)$$

$f_k(x)$     $\chi_k^2$

- $k=1$
- $k=2$
- $k=3$
- $k=4$
- $k=6$
- $k=9$

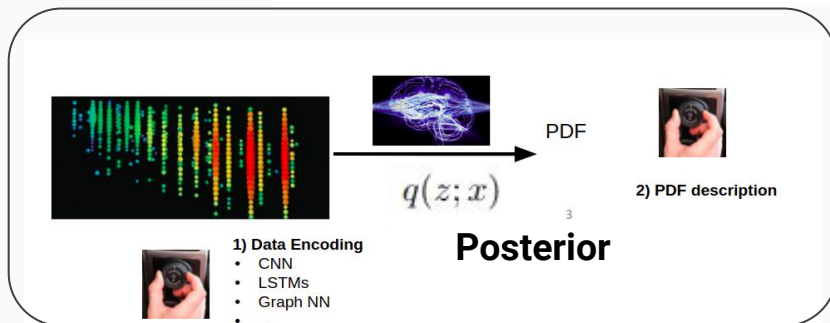(2008.05825)

(2008.05825)

ext. supervised learning

**Input:** data $x_i$, labels $z_{o,i}$

$$D_{\mathrm{KL}}(\mathcal{P}_t | q)[x, z_o, z_u]$$

VI posterior
(via mixed loss):
$$q_{\phi,\varphi}(z_o, z_u; x) \sim \mathcal{P}_t(z_o, z_u; x)$$

PDF

**2) PDF description**

$q(z; x)$

**Posterior**

**1) Data Encoding**
- CNN
- LSTMs
- Graph NN

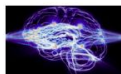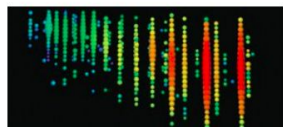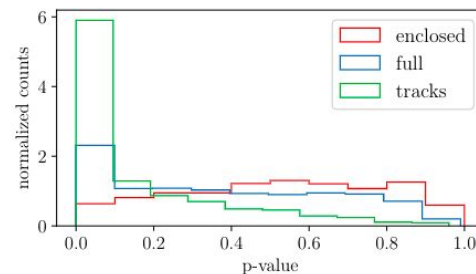(2008.05825)

**Likelihood** $p(x; z)$

ext. supervised learning

**Input:** data $x_i$, labels $z_{o,i}$

$$D_{KL}(\mathcal{P}_t|q)[x, z_o, z_u]$$

VI posterior
(via mixed loss):
$$q_{\phi,\varphi}(z_o, z_u; x) \sim \mathcal{P}_t(z_o, z_u; x)$$

PDF

$q(z; x)$

2) PDF description

3

**Posterior**

1) Data Encoding
- CNN
- LSTMs
- Graph NN
- ...

$$p_{val} = \int_{x,z} \boldsymbol{I}_{T(x,z)>T(x_{obs},z)} p_\theta(x; z) q_\phi(z; x_{obs}) dx dz$$

$$T(x, z) = \ln p_\theta(x; z)/N_d$$

(2008.05825)

**Likelihood** $p(x; z)$

- Normalizing flows + neural networks allow to model **complex (conditional) PDFs**

- Probabilistic interpretation of machine possible with normalizing flows (supervised / unsupervised learning etc. are just **PDF matching of Posterior/likelihood**)

- **Systematics** trivial (include in training)
- **Coverage** for arbitrary shaped distributions
- **Goodness-of-Fit** (potentially "single-cut" final level event selections that are "optimal")

(2008.05825 for more info)

**https://github.com/thoglu/jammy_flows**
**(first release in ~ few weeks)**