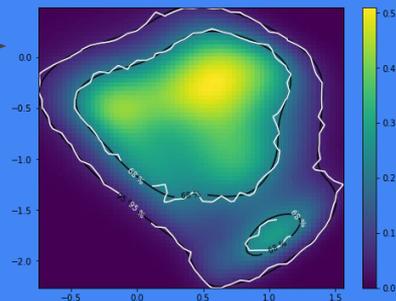
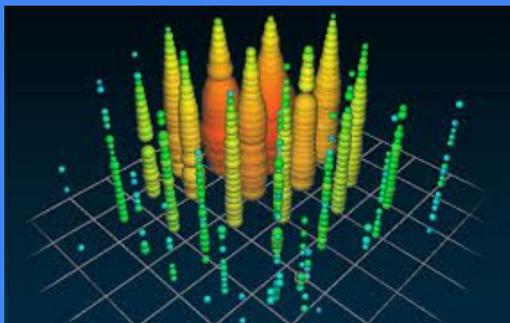


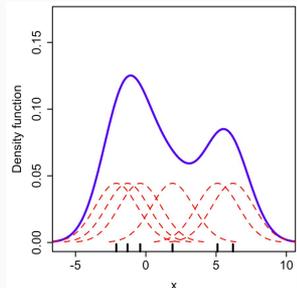


Normalizing flows

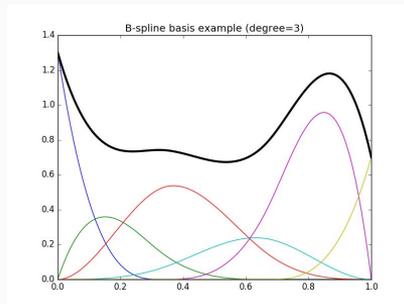
Differentiable Expectation values and more



PDF modelling:



Mixture models /
Kernel Density Functions

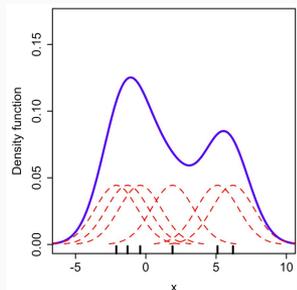


B-Splines

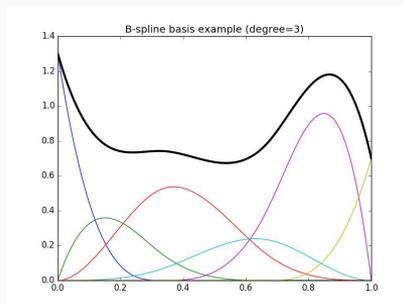
- (arbitrarily) Complex PDF shape
- Evaluate probability analytically
- Hard in high-D
- ...

Why normalizing flows?

PDF modelling:

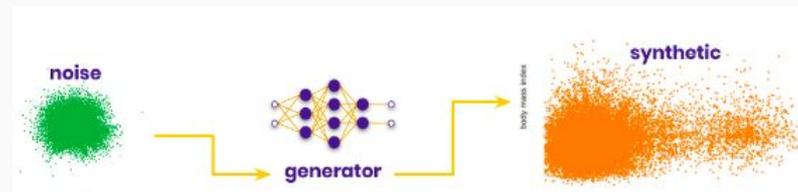


Mixture models /
Kernel Density Functions



B-Splines

Generative Models



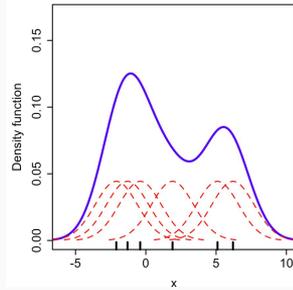
E.g. GANs

[statice.ai]

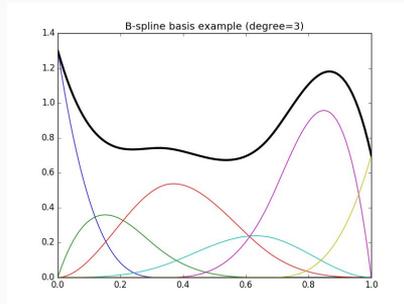
- (arbitrarily) Complex PDF shape
- Evaluate probability analytically
- Hard in high-D
- ...

- Complex data
- No PDF evaluation
- Easy in high-D
- ...

PDF modelling:

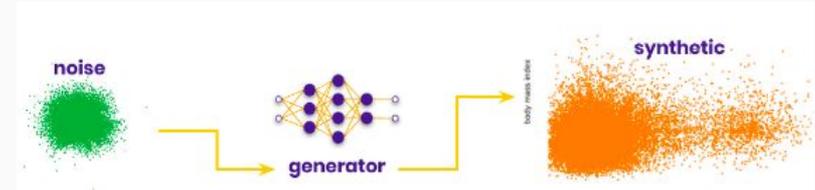


Mixture models /
Kernel Density Functions



B-Splines

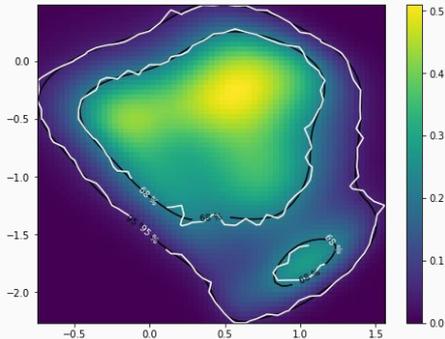
Generative Models



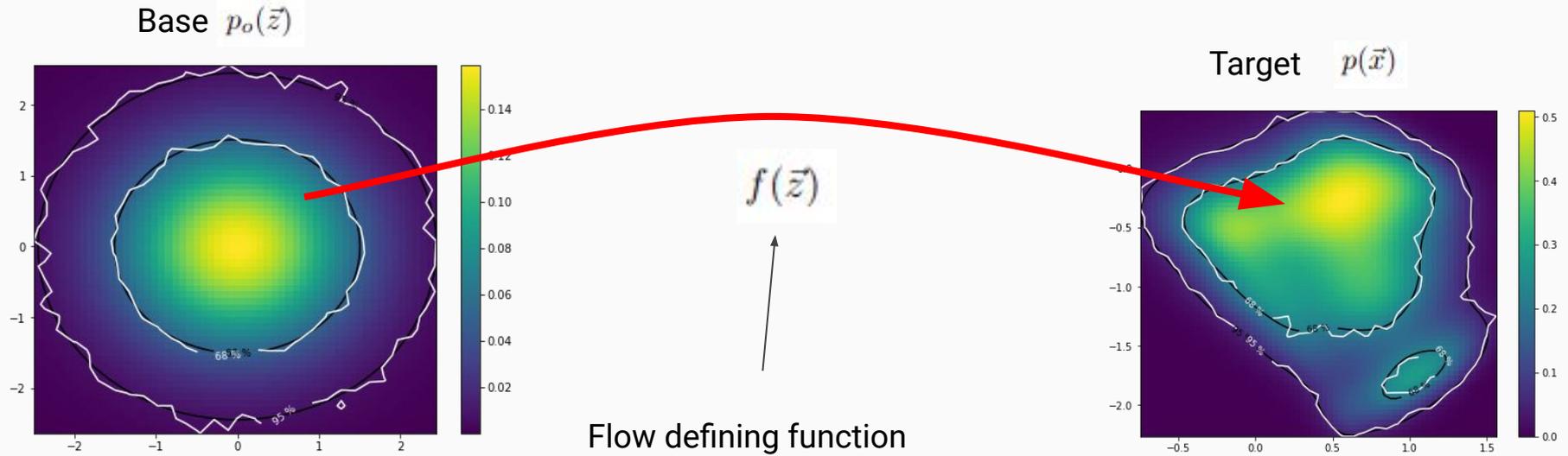
E.g. GANs

[statice.ai]

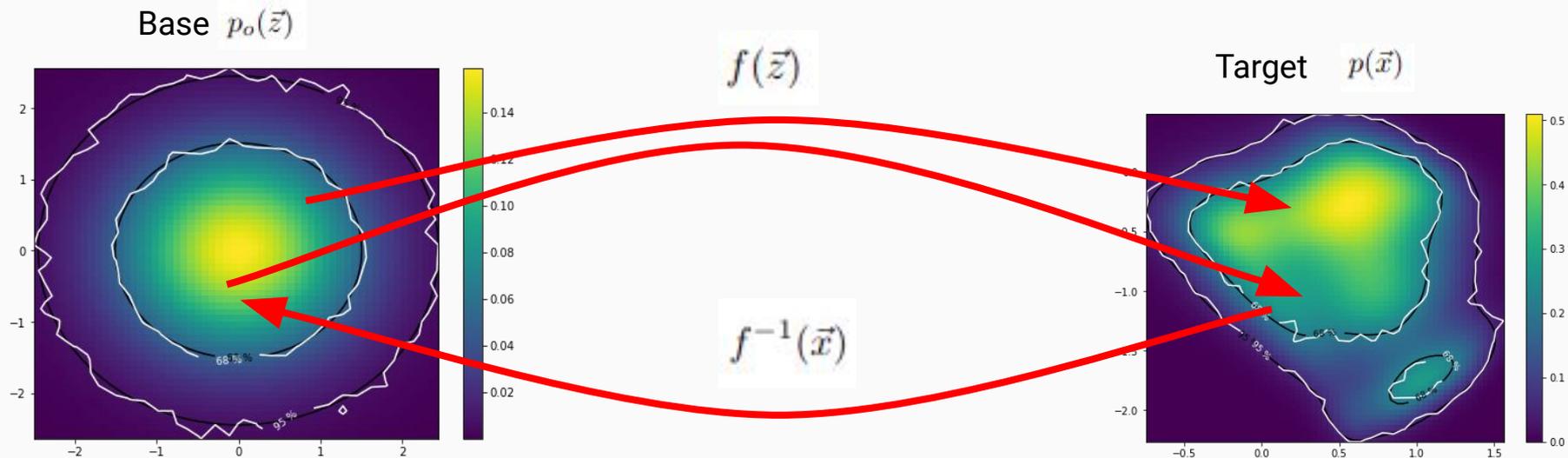
Normalizing flows: PDFs and generative models simultaneously



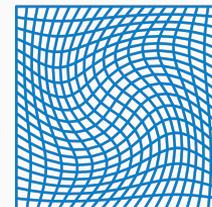
- (arbitrarily) Complex PDF shape
- Evaluate probability analytically
- Works in $D \gg 1$
- Can be interpreted as generalizations of the Gaussian distribution
- Works on manifolds
- Generate differentiable samples (-> differentiable expectation values)
- ...



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



- 1) The function has to be **bijective** for the inverse!
 - 2) The function has to be **differentiable**!
- } = Diffeomorphism

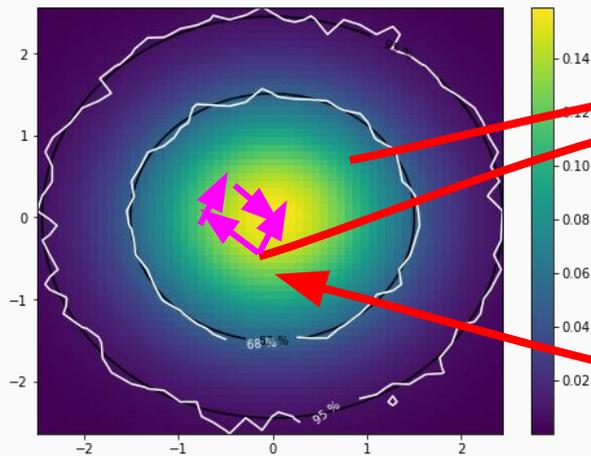




$$p(\vec{x}) = p_0(f^{-1}(\vec{x})) \cdot |\det J^{-1}(\vec{x})|$$

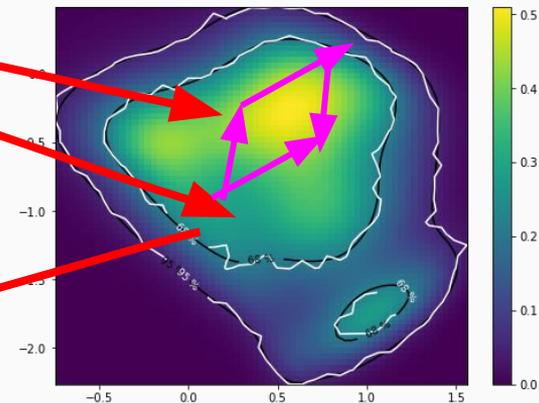
← Log probability includes
Local “stretching”/”squeezing”
Of volume element

Base $p_0(\vec{z})$



$f(\vec{z})$

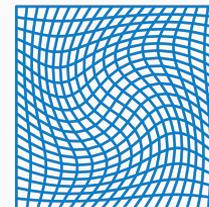
Target $p(\vec{x})$



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

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

= Diffeomorphism

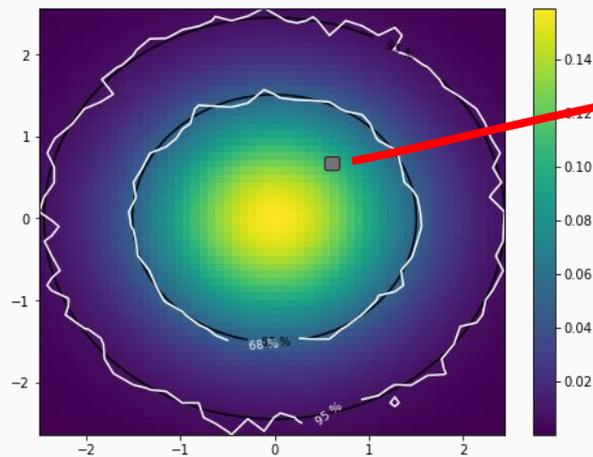


Sampling is also straight forward:

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

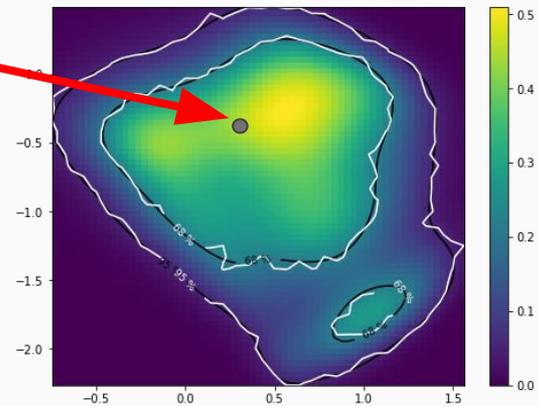
Must be able to sample from base distribution!

Base $p_o(\vec{z})$



$f(\vec{z})$

Target $p(\vec{x})$



Definition

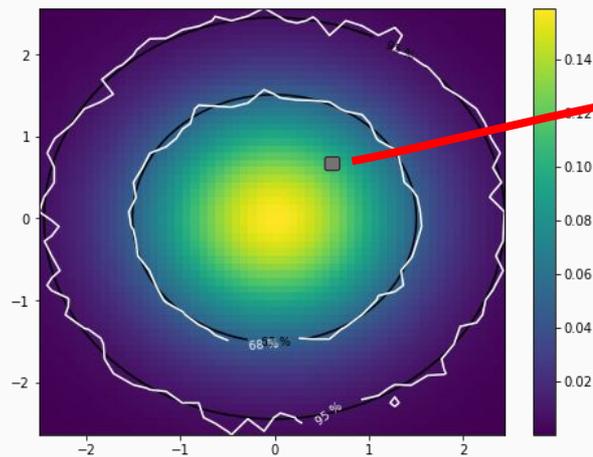
Sampling is also straight forward:

Must be able to sample from base distribution!

$$\vec{x}_\theta = f_\theta(\vec{z})$$

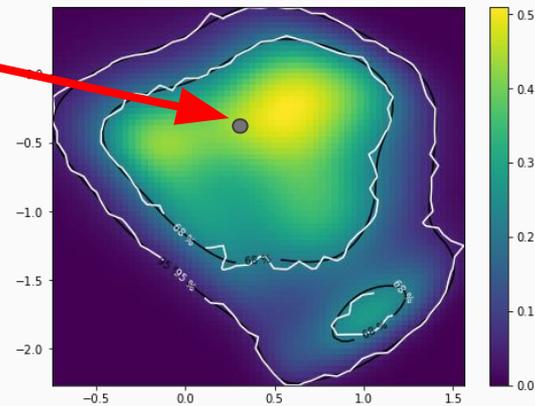
$$p_\theta(\vec{x}) = p_0(f_\theta^{-1}(\vec{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\vec{x})}{\partial \vec{x}} \right|$$

Base $p_0(\vec{z})$



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

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

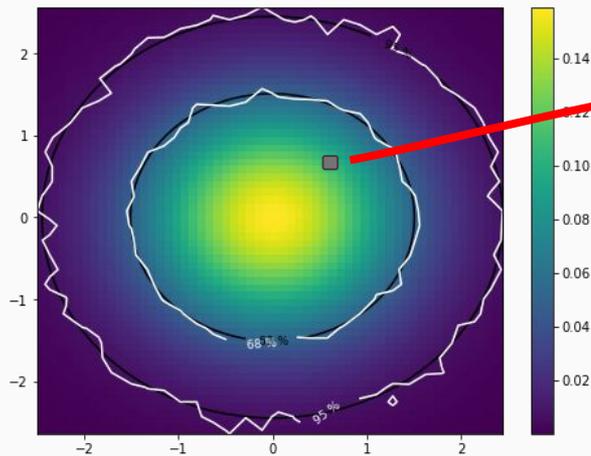




Sampling is also straight forward:

Must be able to sample from base distribution!

Base $p_0(\vec{z})$



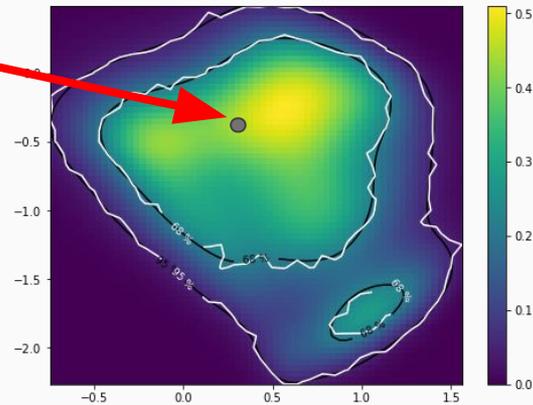
$$\vec{x}_\theta = f_\theta(\vec{z})$$

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

1) Differentiable prob.

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

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



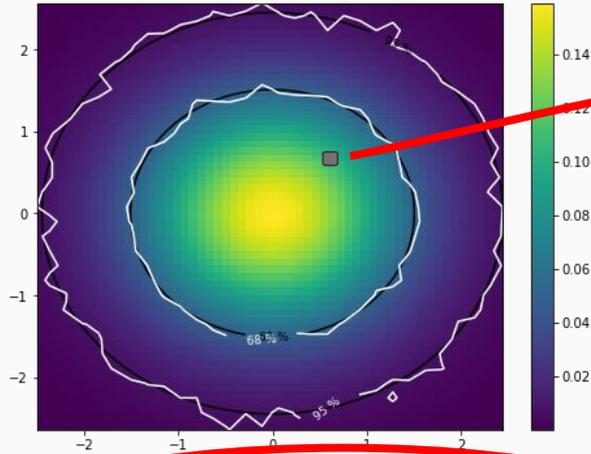
2) Differentiable samples!

Definition

Sampling is also straight forward:

Must be able to sample from base distribution!

Base $p_0(\vec{z})$



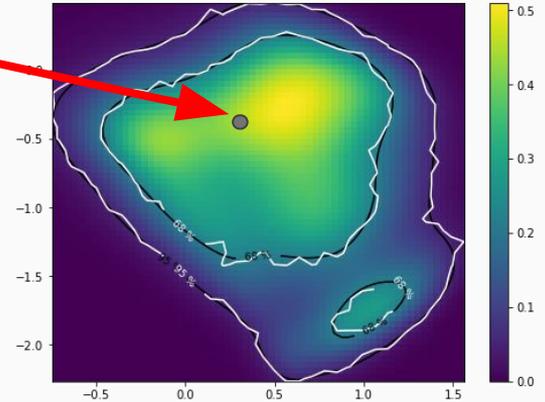
$$\vec{x}_\theta = f_\theta(\vec{z})$$

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

1) Differentiable prob.

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

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



2) Differentiable samples!

The simplest possible normalizing flow (1-d)



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

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

The simplest possible normalizing flow (1-d)



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

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

$$\theta = \{a, b\}$$

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



The simplest possible normalizing flow (1-d)

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

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

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



The simplest possible normalizing flow (1-d)

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

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

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



The simplest possible normalizing flow (1-d)

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

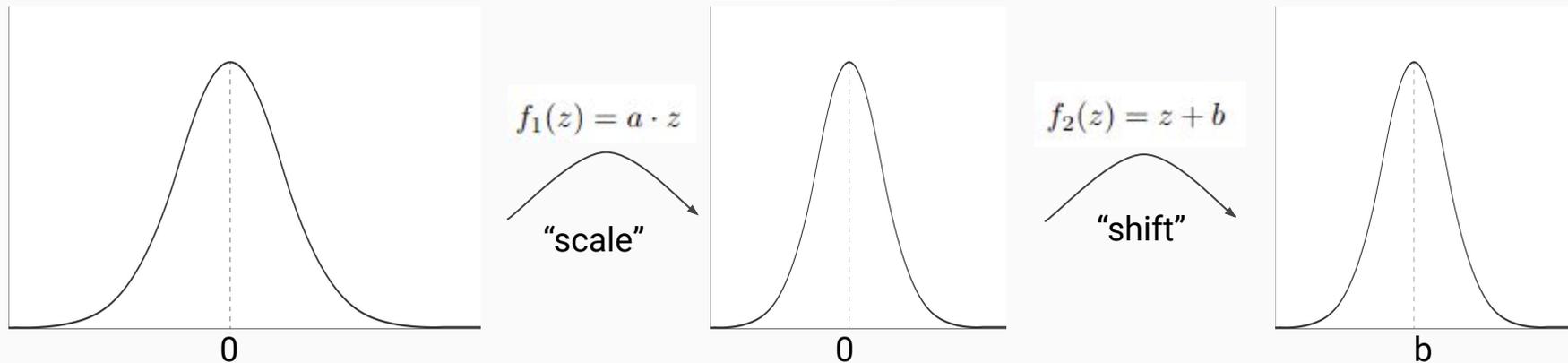
$$x = f(\vec{z}) = a \cdot z + b$$

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:

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



The simplest possible normalizing flow (n-d)

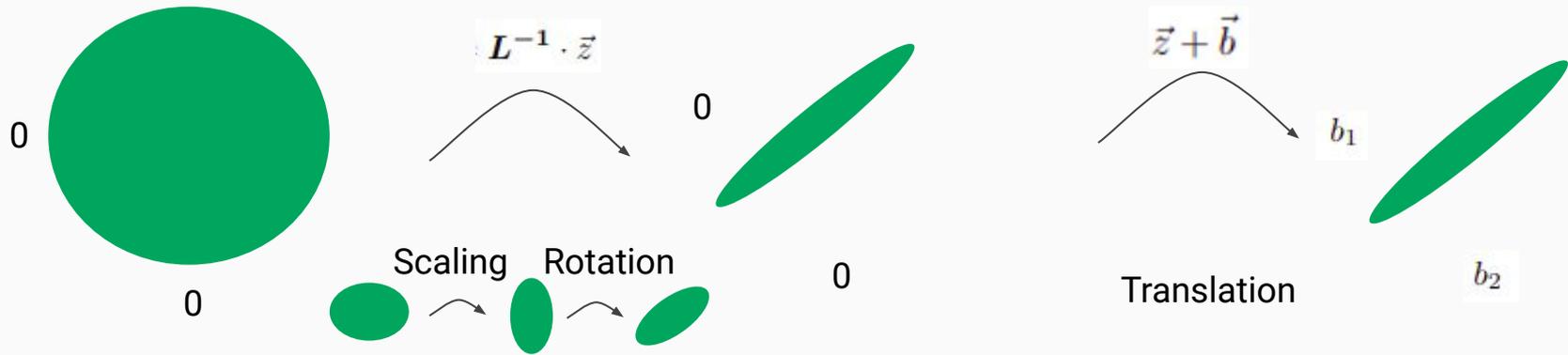
Let us try to see the distribution of the simplest possible normalizing flow in n dimensions!

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

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

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

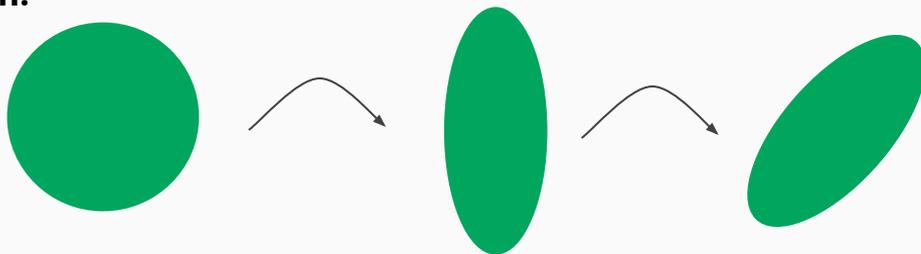
“Affine flow” = “general Gaussian distribution” in n-d





4) Gaussianization Flow(arXiv:2003.01941)

Gaussian:



Gaussianization flow generalized by non-linear scalings instead of “linear scalings”

$$T_{\theta}(\mathbf{x}) = \Psi_{\theta_L} \circ R_L \circ \Psi_{\theta_{L-1}} \circ \dots \circ \Psi_{\theta_1} \circ R_1 \mathbf{x}$$

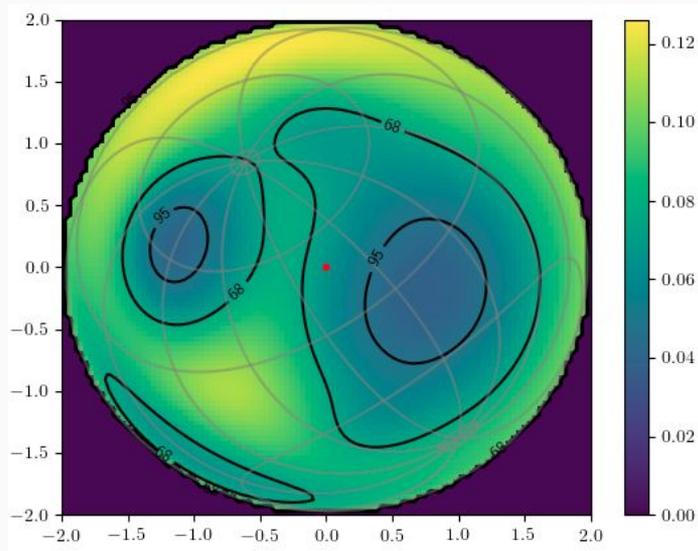


Step by step refinement of the PDF - provably approximates any distribution



“Manifold” normalizing flows

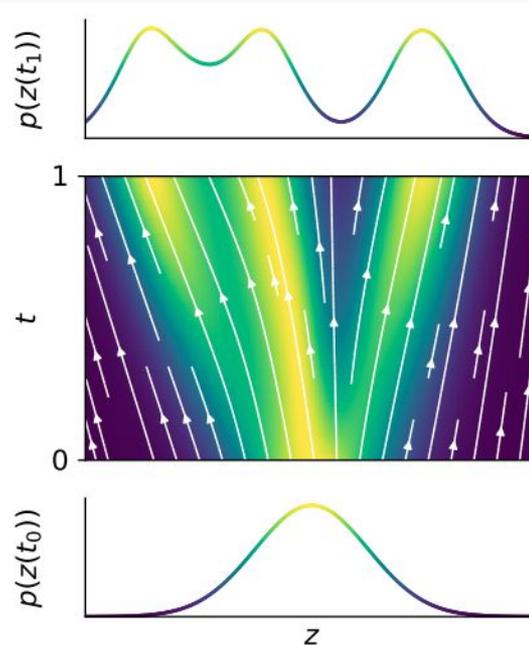
$$p(T(x)) = \frac{\pi(x)}{\sqrt{\det(E^T J^T J E)}}$$



(2002.02428)

“Continuous” normalizing flows

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



(FFJORD,
1810.01367)

Example: more expressive label predictions (e.g. For MODE applications)



Typically, label space in regression is low dimensional

MSE-loss
(-log-prob)



Std normal

(+ log_sigma predictions)
(-log-prob)



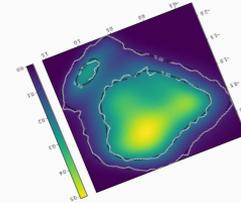
Scaled normal

-log-prob



MVN

-log-prob



Arbitrary shape



Example: more expressive label predictions (e.g. For MODE applications)

Typically, label space in regression is low dimensional

MSE-loss
(-log-prob)



Std normal

(+ log_sigma predictions)
(-log-prob)



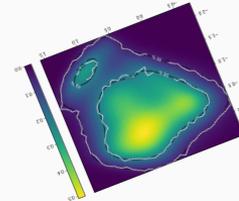
Scaled normal

-log-prob



MVN

-log-prob



Arbitrary shape

Regression always expressible as normalizing flow training

More expressive power



“Adding more labels does not help .. results even get worse. I just train separate models for each label.”

“With (autoregressive) normalizing flows the opposite should happen... adding more labels helps training and performance.”





$$D_{\text{KL},\text{joint}(x,z_o)}(\mathcal{P}_t; q_\phi) = \int_x \int_{z_o} \mathcal{P}_t(z_o, x) \cdot \ln \frac{\mathcal{P}_t(z_o, x)}{q_\phi(z_o, x)} dz_o dx$$



$$D_{\text{KL},\text{joint}(x,z_o)}(\mathcal{P}_t; q_\phi) = \int_x \int_{z_o} \mathcal{P}_t(z_o, x) \cdot \ln \frac{\mathcal{P}_t(z_o, x)}{q_\phi(z_o, x)} dz_o dx$$

$$\begin{aligned} \arg \min_{\phi} \hat{D}_{\text{KL},\text{joint}(x,z_o)}(\mathcal{P}_t; q_\phi) &= \arg \min_{\phi} \frac{1}{N} \sum_{x_i, z_{o,i}} \ln \left(\frac{\mathcal{P}_t(z_{o,i}; x_i)}{q_\phi(z_{o,i}; x_i)} \right) + \ln \left(\frac{\mathcal{P}_t(x_i)}{q(x_i)} \right) \\ &= \arg \min_{\phi} \frac{1}{N} \sum_{x_i, z_{o,i}} -\ln (q_\phi(z_{o,i}; x_i)) + \text{const} \\ &= \arg \min_{\phi} \frac{1}{N} \sum_{x_i, z_{o,i}} -\ln (q_\phi(z_{o,i}; x_i)) \\ &\equiv \arg \min_{\phi} \mathcal{L}_{\text{supervised}}(\phi) \end{aligned}$$



$$D_{\text{KL},\text{joint}(x,z_o)}(\mathcal{P}_t; q_\phi) = \int_x \int_{z_o} \mathcal{P}_t(z_o, x) \ln \frac{\mathcal{P}_t(z_o, x)}{q_\phi(z_o, x)} dz_o dx$$

Measure independent of parameters

$$\begin{aligned} \arg \min_{\phi} \hat{D}_{\text{KL},\text{joint}(x,z_o)}(\mathcal{P}_t; q_\phi) &= \arg \min_{\phi} \frac{1}{N} \sum_{x_i, z_{o,i}} \ln \left(\frac{\mathcal{P}_t(z_{o,i}; x_i)}{q_\phi(z_{o,i}; x_i)} \right) + \ln \left(\frac{\mathcal{P}_t(x_i)}{q(x_i)} \right) \\ &= \arg \min_{\phi} \frac{1}{N} \sum_{x_i, z_{o,i}} -\ln(q_\phi(z_{o,i}; x_i)) + \text{const} \\ &= \arg \min_{\phi} \frac{1}{N} \sum_{x_i, z_{o,i}} -\ln(q_\phi(z_{o,i}; x_i)) \\ &\equiv \arg \min_{\phi} \mathcal{L}_{\text{supervised}}(\phi) \end{aligned}$$

(Gradients straightforward -
We use it all the time)



$$I_{\theta} = \int p_{\theta}(x) F_{\theta}(x) dx \quad ?$$

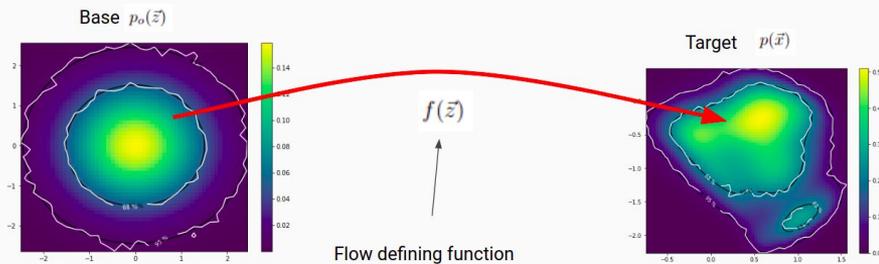
Differentiable expectation values

$$I_\theta = \int p_\theta(x) F_\theta(x) dx_\theta = \int p(z) F_\theta(f_\theta(z)) dz$$

$$\vec{x}_\theta = f_\theta(\vec{z})$$

Reparametrization trick -> calculate integral in base space (equivalent to target integral)

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



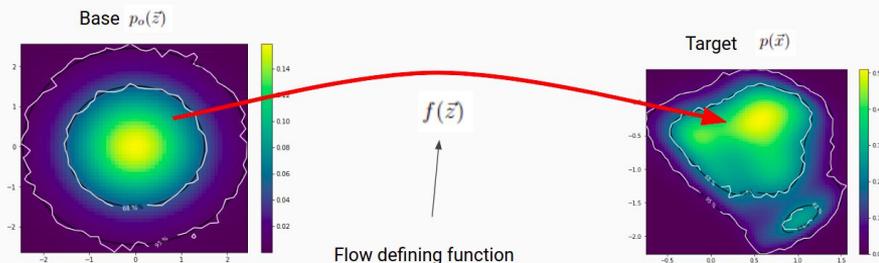
$$dx_\theta = dz \left| \det \frac{\partial f_\theta(\vec{z})}{\partial \vec{z}} \right|$$

Differentiable expectation values

$$I_\theta = \int p_\theta(x) F_\theta(x) dx_\theta = \int p(z) F_\theta(f_\theta(z)) dz$$

$$\vec{x}_\theta = f_\theta(\vec{z})$$

Reparametrization trick -> calculate integral in base space (equivalent to target integral)



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

$$dx_\theta = dz \left| \det \frac{\partial f_\theta(\vec{z})}{\partial \vec{z}} \right|$$

$$\approx \frac{1}{N} \cdot \sum_i^N F_\theta(x_\theta)$$



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

Examples:

n-th moment of p

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



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

Examples:

n-th moment of p

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

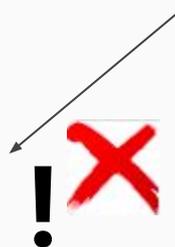
entropy

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



Wrong gradient!

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



.... Many other integrals in information theory
.... Required for VAE to work etc

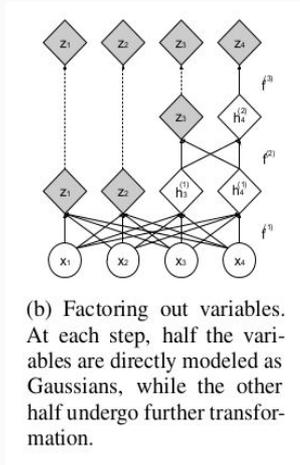
We need normalizing flows for our specific needs...



We need normalizing flows for our specific needs...

-> Need specific implementations

We don't really care for PDFs over 10k-dim images (what CS people are into)..



(real-NVP, 2017)

... But rather conditional PDFs over directions, energies, positions (maybe in total only a few 10s of dimension max) - and **correct joint correlation structure** over these PDFs!



https://github.com/thoglu/jammy_flows

“A generic pdf on a tensor product of manifolds”

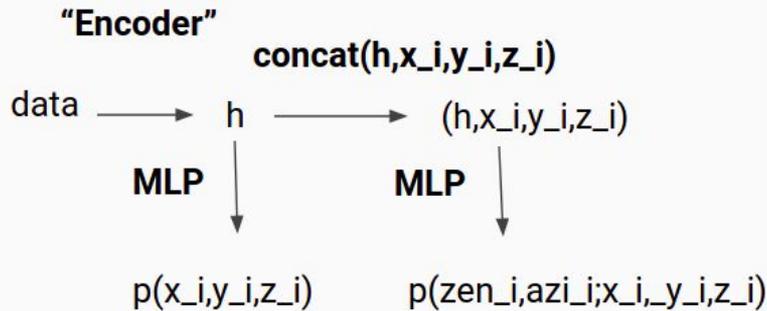
E.g. a 5-d conditional PDF: $p(x,y,z,zenith,azimuth ; data)$

```
import jammy_flows
```

```
pdf=jammy_flows.pdf("e3+s2", "gggg+n")
```

(generalizes notion of of “Inverse Autoregressive Flows” (Kingma et. al 2016))

(Inverse) autoregressive routing





Jammy_flows (conditional + non-conditional “autoregressive” PDFs)

https://github.com/thoglu/jammy_flows

“A generic pdf on a tensor product of manifolds”

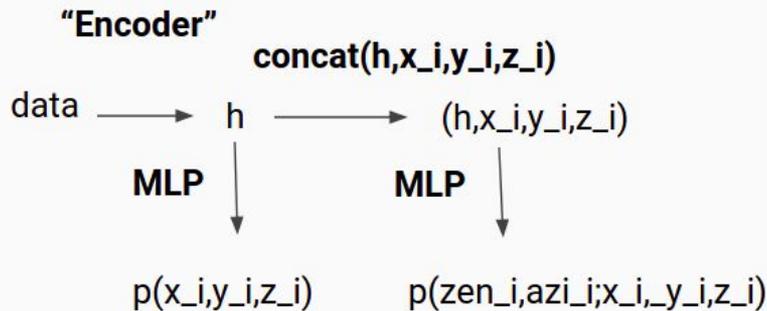
E.g. a 5-d conditional PDF: $p(x,y,z,\text{zenith},\text{azimuth} ; \text{data})$

```
import jammy_flows
```

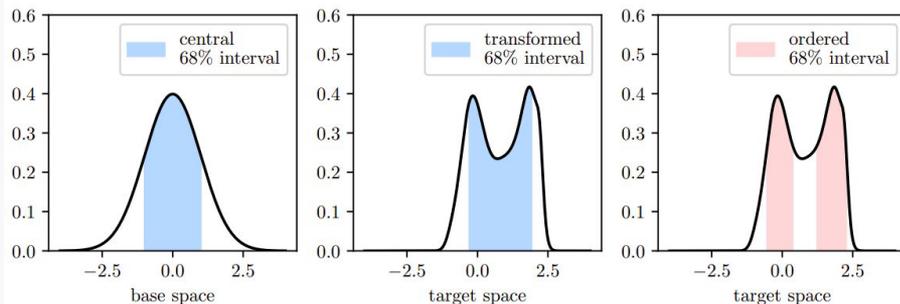
```
pdf=jammy_flows.pdf("e3+s2", "gggg+n")
```

(generalizes notion of of “Inverse Autoregressive Flows” (Kingma et. al 2016))

(Inverse) autoregressive routing



Gaussian base allows for coverage control:

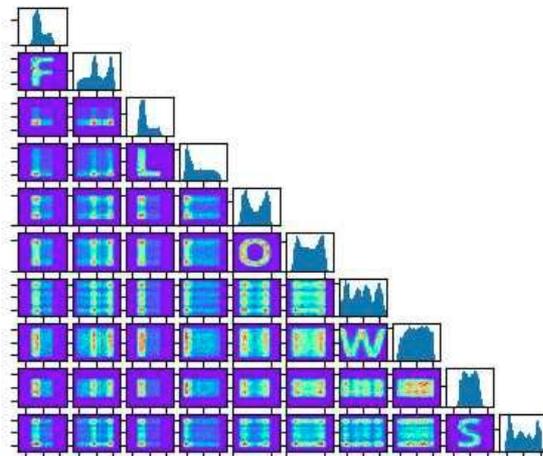
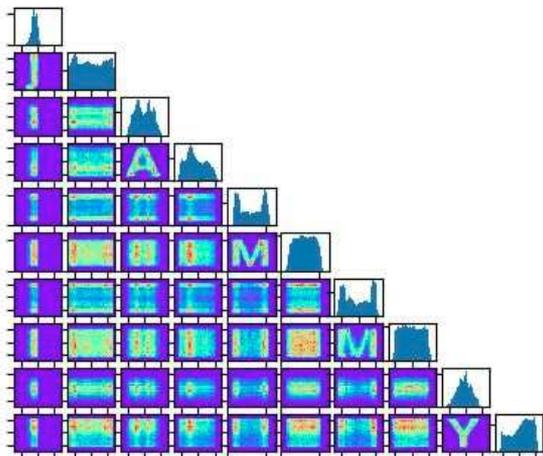


https://github.com/thoglu/jammy_flows

```
import jammy_flows
```

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

pdf structure: e4+s2+e4



“A generic pdf on a tensor product of manifolds”

Coverage for arbitrary distributions and other things ..

(<https://arxiv.org/abs/2008.05825>)

