

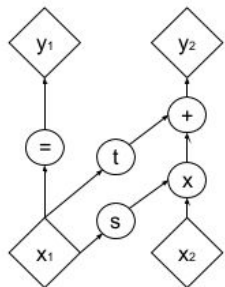
# Normalizing flows - part 2

## **Wednesday:**

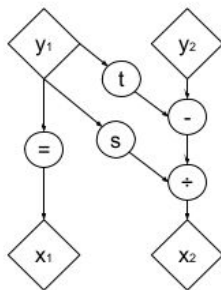
- 1) Various N-d non-linear generalizations**
- 2) Conditional PDFs: Amortization and the connection to deep learning**
- 3) A NF tool: jammy\_flows**
- 4) Hands-on 1**
- 5) Probabilistic deep learning**
- 6) Coverage / Systematics / Goodness-of-Fit**
- 7) Hands-on 2**

# Certain N-d non-linear normalizing-flow functions

## 1) “Affine Coupling layer” (Dinh et al. 2014 (NICE) / Dinh et al. 2017 (real-NVP))



(a) Forward propagation



(b) Inverse propagation

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

$$\frac{\partial y}{\partial x^T} = \begin{bmatrix} \mathbb{I}_d & 0 \\ \frac{\partial y_{d+1:D}}{\partial x_{1:d}^T} & \text{diag}(\exp[s(x_{1:d})]) \end{bmatrix},$$

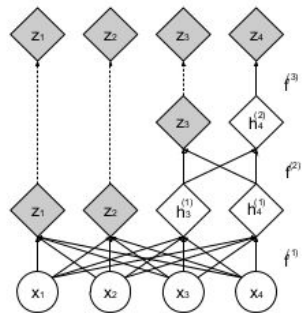
Induces coupling between different layers

Functions **s** and **t** can be arbitrarily complex (neural networks)

These functions are also called “conditioners”

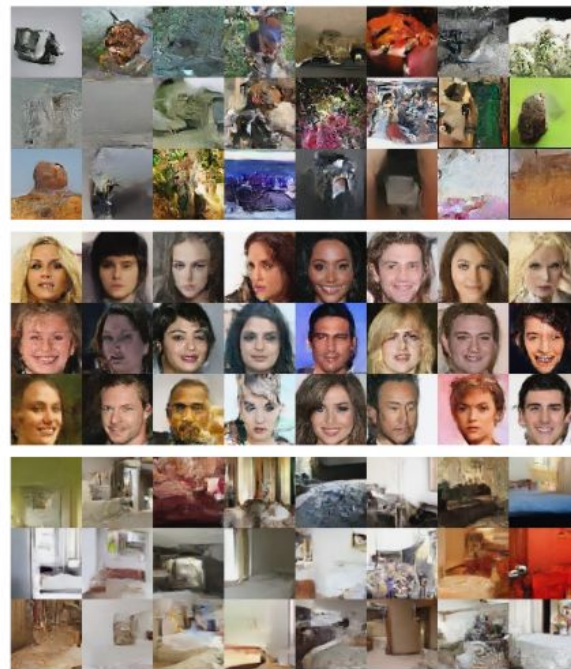
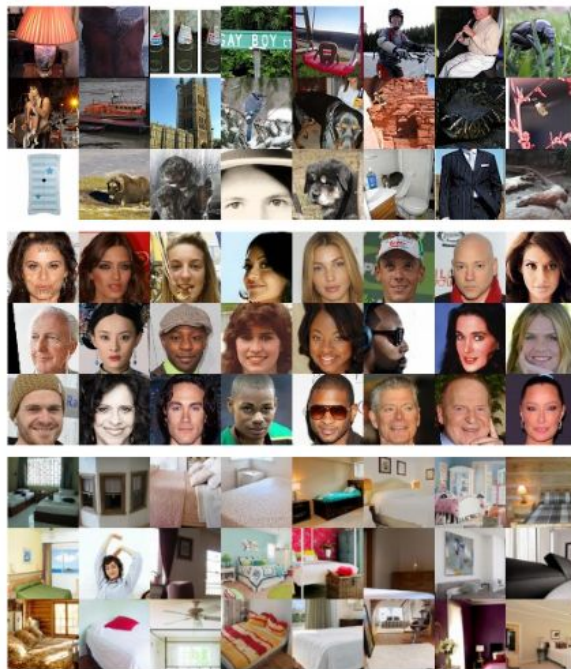
# Certain N-d non-linear normalizing-flow functions

## 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.

$D \odot \mu$



Can be scaled to very high dimensions (here images)

## Certain N-d non-linear normalizing-flow functions

### 2) “Autoregressive” structure

$$p(x_0) \cdot p(x_1; x_0) \cdot \dots \cdot p(x_i; x_{i-1}, x_{i-2}, \dots, x_0)$$

Flow function for “affine” example:

$$y_0 = \mu_0 + \sigma_0 \odot \epsilon_0, \quad y_i = \mu_i(\mathbf{y}_{1:i-1}) + \sigma_i(\mathbf{y}_{1:i-1}) \cdot \epsilon_i$$

Unlike Coupling layers, no dimensions remain untransformed, and  
Conditioning on previous transformed variables

# Certain N-d non-linear normalizing-flow functions

## 2) “Autoregressive” structure

$$p(x_0) \cdot p(x_1; x_0) \cdot \dots \cdot p(x_i; x_{i-1}, x_{i-2}, \dots, x_0)$$

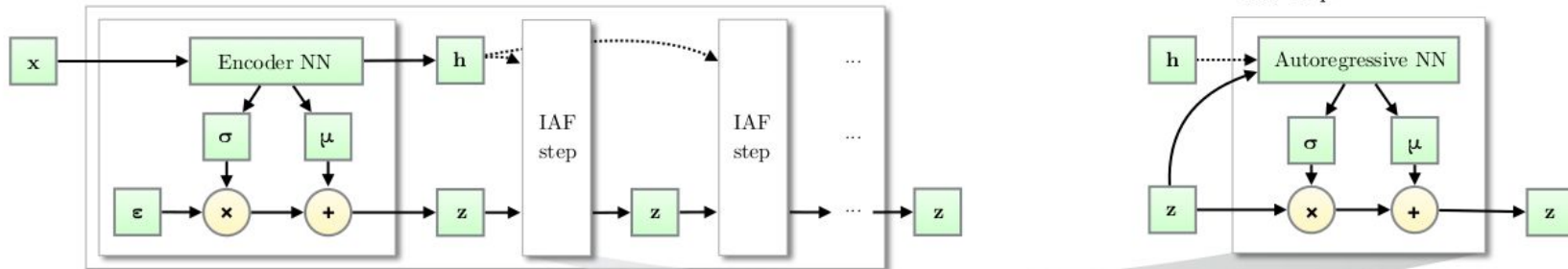
Flow function for “affine” example:

$$y_0 = \mu_0 + \sigma_0 \odot \epsilon_0, \quad y_i = \mu_i(\mathbf{y}_{1:i-1}) + \sigma_i(\mathbf{y}_{1:i-1}) \cdot \epsilon_i$$

Unlike Coupling layers, no dimensions remain untransformed, and  
Conditioning on previous transformed variables

One can also interpret the autoregressive transformations as acting in time  
Inverse Autoregressive flows (IAFs, Kingma et al. 2016)

Approximate Posterior with Inverse Autoregressive Flow (IAF)



# Certain N-d non-linear normalizing-flow functions

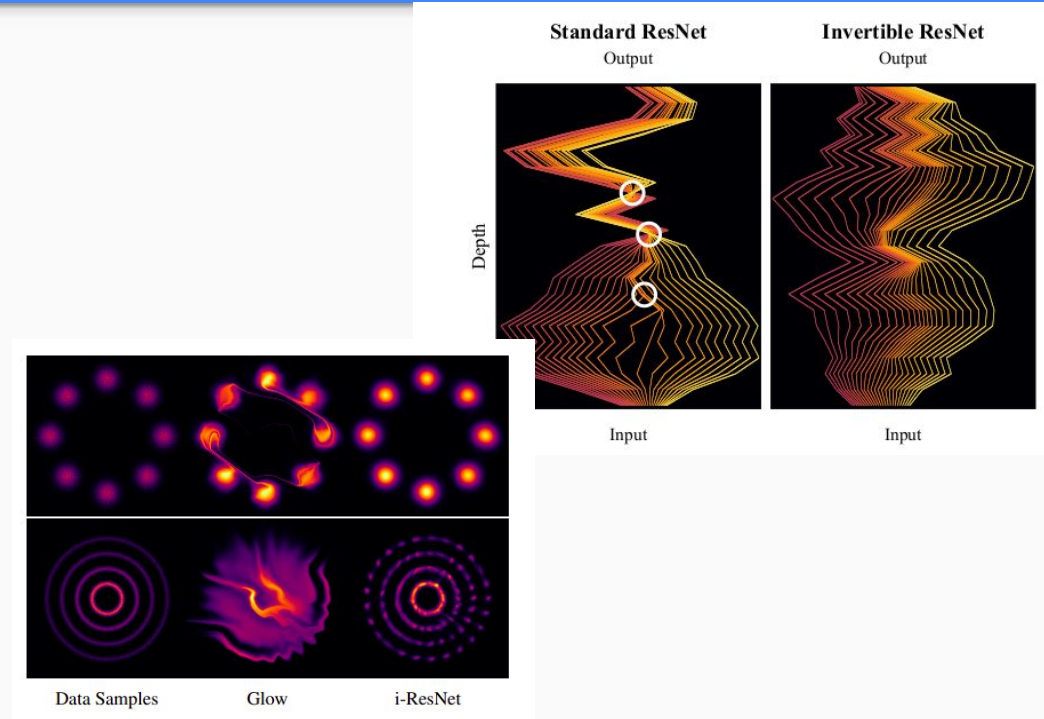
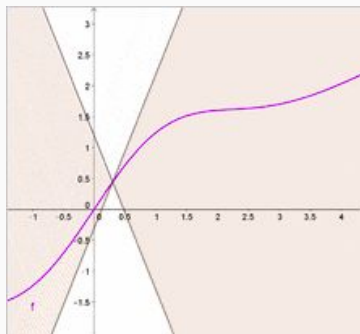
## 3) Use Invertible residual networks: (Behrman et al. 2019)

$$x_{t+1} \leftarrow x_t + g_{\theta_t}(x_t)$$

*invertible if*

$$\text{Lip}(g_{\theta_t}) < 1, \text{ for all } t = 1, \dots, T,$$

*where  $\text{Lip}(g_{\theta_t})$  is the Lipschitz-constant of  $g_{\theta_t}$ .*



# Certain N-d non-linear normalizing-flow functions

## 3) Use Invertible residual networks: (Behrman et al. 2019)

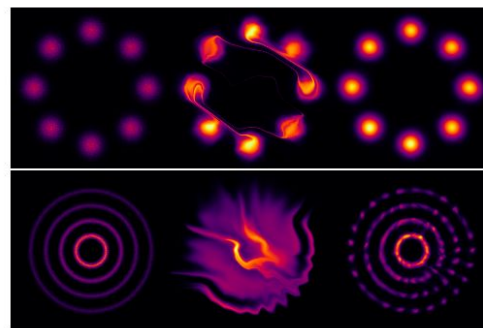
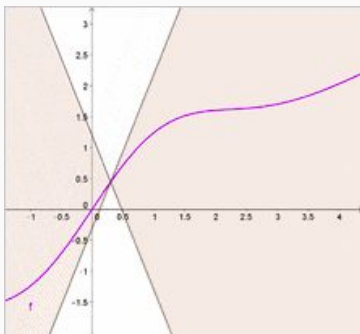
$$x_{t+1} \leftarrow x_t + g_{\theta_t}(x_t)$$

Flow parameters  
are NN network parameters

*invertible if*

$$\text{Lip}(g_{\theta_t}) < 1, \text{ for all } t = 1, \dots, T,$$

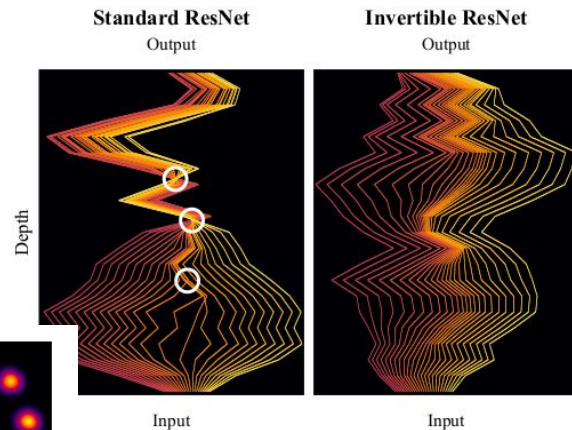
where  $\text{Lip}(g_{\theta_t})$  is the Lipschitz-constant of  $g_{\theta_t}$ .



Data Samples

Glow

i-ResNet



Lipschitz constraint similar to critic constraint in WGAN training



# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

$$D_{\text{KL}}(p(\mathbf{x}) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I})) \triangleq J(\mathbf{x}) = I(\mathbf{x}) + J_m(\mathbf{x})$$

“Procedure of **Gaussianization**”

-> Transform a random variable into Gaussian

By: 1) reducing correlations

2) making marginals = standard normal

“Multi information / total correlation”

-> correlation between variables

$$I(\mathbf{x}) = D_{\text{KL}} \left( p(\mathbf{x}) \parallel \prod_i^D p_i(x^{(i)}) \right)$$

$$J_m(\mathbf{x}) = \sum_{i=1}^D D_{\text{KL}} \left( p_i(x^{(i)}) \parallel \mathcal{N}(0, 1) \right).$$

Iterative Gaussianization (Chen and Gopinath, 2001)

Achieve 1) + 2) by iterated application of a Rotation (e.g. defined by PCA matrix) and nonlinear transformation of individual dimensions

New idea: **Learn** Rotations + nonlinear transformations in a normalizing flow  $T_{\theta}(\mathbf{x}) = \Psi_{\theta_L} \circ R_L \circ \Psi_{\theta_{L-1}} \circ \dots \circ \Psi_{\theta_1} \circ R_1 \mathbf{x}$

# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

“Multi information / total correlation”  
-> correlation between variables

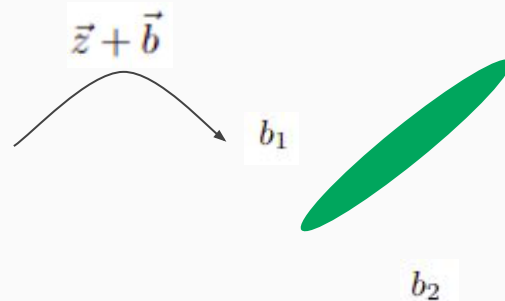
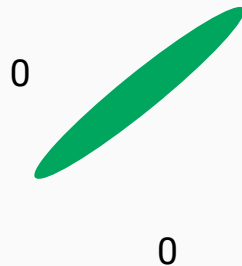
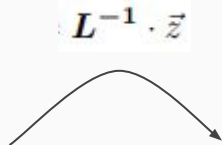
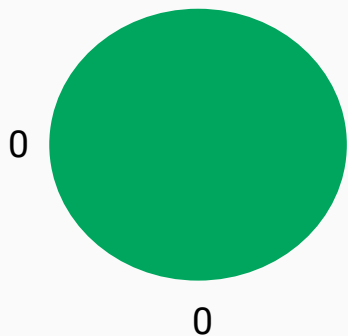
### A better intuitive understanding:

Start with Affine flow (Gaussian):

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot |\det(\mathbf{L})| \cdot \exp\left(-0.5 \cdot (\vec{x} - \vec{b})^T \cdot \mathbf{L}^T \cdot \mathbf{L} \cdot (\vec{x} - \vec{b})\right)$$

Inv. Flow function

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



# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

“Multi information / total correlation”  
-> correlation between variables

### A better intuitive understanding:

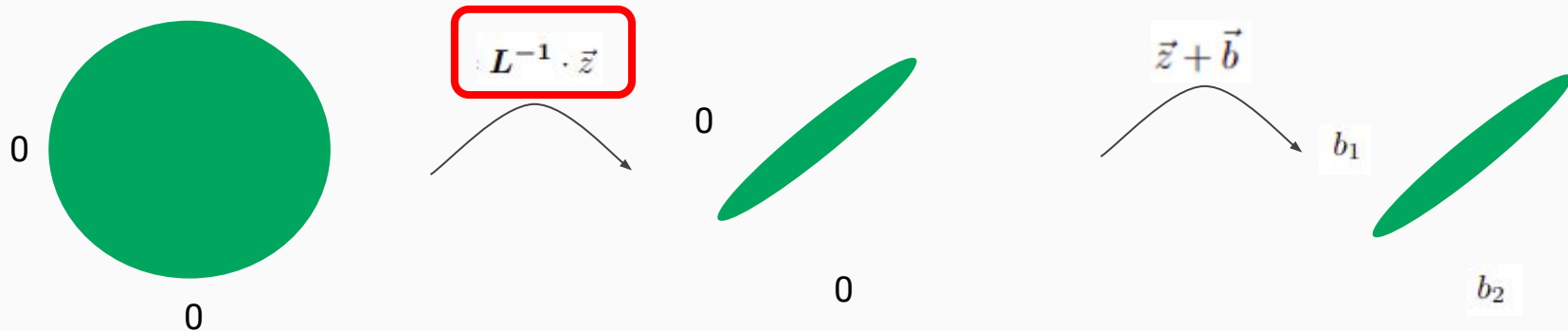
Start with Affine flow (Gaussian):

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot |\det(\mathbf{L})| \cdot \exp\left(-0.5 \cdot (\vec{x} - \vec{b})^T \cdot \mathbf{L}^T \cdot \mathbf{L} \cdot (\vec{x} - \vec{b})\right)$$

Inv. Flow function

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

Let us decompose this transformation further



# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

“Multi information / total correlation”  
-> correlation between variables

### A better intuitive understanding:

Start with Affine flow (Gaussian):

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot |\det(\mathbf{L})| \cdot \exp \left( -0.5 \cdot (\vec{x} - \vec{b})^T \cdot \mathbf{L}^T \cdot \mathbf{L} \cdot (\vec{x} - \vec{b}) \right)$$

Inv. Flow function

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

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot |\det(\mathbf{D})| \cdot \exp \left( -0.5 \cdot (\vec{x} - \vec{b})^T \cdot \mathbf{V} \cdot \mathbf{D} \cdot \mathbf{U}^T \cdot \mathbf{U} \cdot \mathbf{D} \cdot \mathbf{V}^T \cdot (\vec{x} - \vec{b}) \right)$$

# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

“Multi information / total correlation”  
-> correlation between variables

### A better intuitive understanding:

Start with Affine flow (Gaussian):

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot |\det(\mathbf{L})| \cdot \exp \left( -0.5 \cdot (\vec{x} - \vec{b})^T \cdot \mathbf{L}^T \cdot \mathbf{L} \cdot (\vec{x} - \vec{b}) \right)$$

Inv. Flow function

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

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot |\det(\mathbf{D})| \cdot \exp \left( -0.5 \cdot (\vec{x} - \vec{b})^T \cdot \mathbf{V} \cdot \mathbf{D} \cdot \mathbf{U}^T \cdot \mathbf{U} \cdot \mathbf{D} \cdot \mathbf{V}^T \cdot (\vec{x} - \vec{b}) \right)$$

$$p(x) = \frac{1}{\sqrt{2 \cdot \pi}} \cdot |\det(\mathbf{D})| \cdot \exp \left( -0.5 \cdot (\mathbf{D} \cdot \mathbf{O} \cdot (\vec{x} - \vec{b}))^T \cdot \mathbf{D} \cdot \mathbf{O} \cdot (\vec{x} - \vec{b}) \right)$$

$$f^{-1}(\vec{x}) = \mathbf{D} \cdot \mathbf{O}(\vec{x} - \vec{b})$$

Flow function

Is actually

**Scaling + Rotation**

$$\longrightarrow f(\vec{z}) = \mathbf{O} \cdot \mathbf{D}(\vec{x}) + \vec{b}$$

# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

Flow function

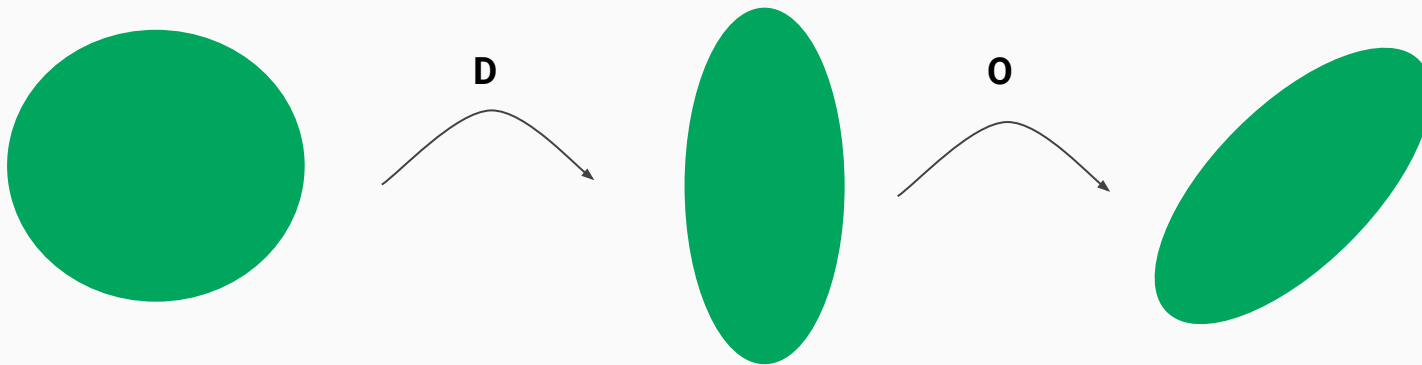
Is actually

**Scaling + Rotation**

$$\longrightarrow f(\vec{z}) = \mathbf{O} \cdot \mathbf{D}(\vec{x}) + \vec{b}$$

“Multi information / total correlation”

-> correlation between variables



# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

Flow function

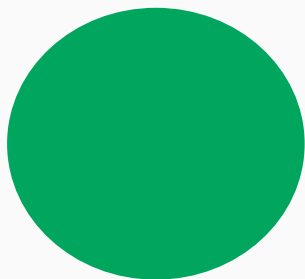
Is actually

**Scaling + Rotation**

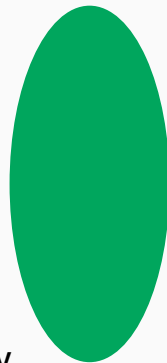
$$\longrightarrow f(\vec{z}) = \mathbf{O} \cdot \mathbf{D}(\vec{x}) + \vec{b}$$

“Multi information / total correlation”

-> correlation between variables

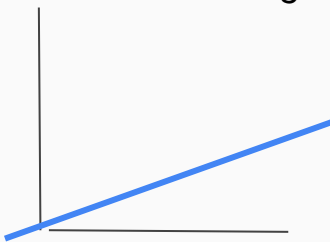


**D**

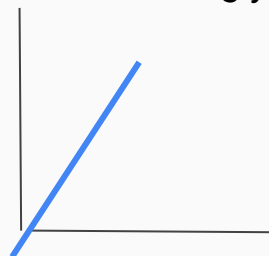


$$\mathbf{D} = \begin{pmatrix} a_x & 0 \\ 0 & a_y \end{pmatrix}$$

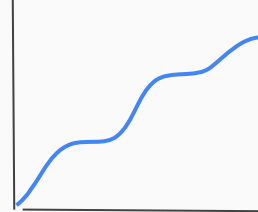
Linear Scaling x



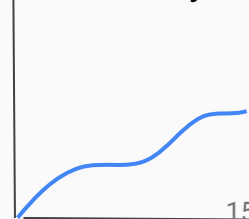
Linear Scaling y



Non-linear x



Non-linear y



# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

Flow function

Is actually

**Scaling + Rotation**

$$\longrightarrow f(\vec{z}) = \mathbf{O} \cdot \mathbf{D}(\vec{x}) + \vec{b}$$

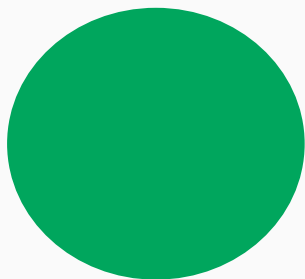
“Multi information / total correlation”

-> correlation between variables

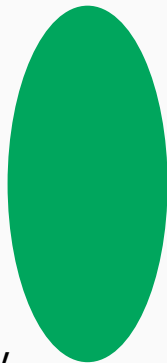
Apply many iterations of “rotation” and non-linear  
Dimension-wise scaling

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

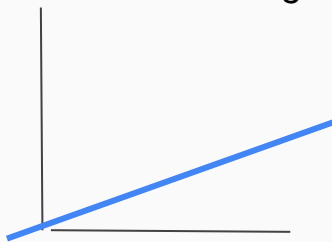
$$D = \begin{pmatrix} a_x & 0 \\ 0 & a_y \end{pmatrix}$$



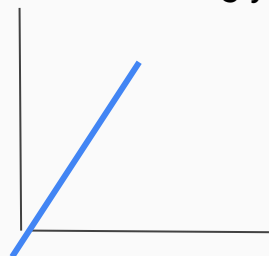
**D**



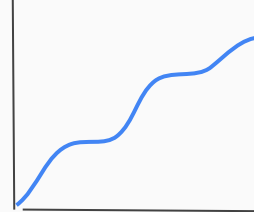
Linear Scaling x



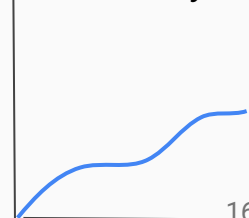
Linear Scaling y



Non-linear x



Non-linear y

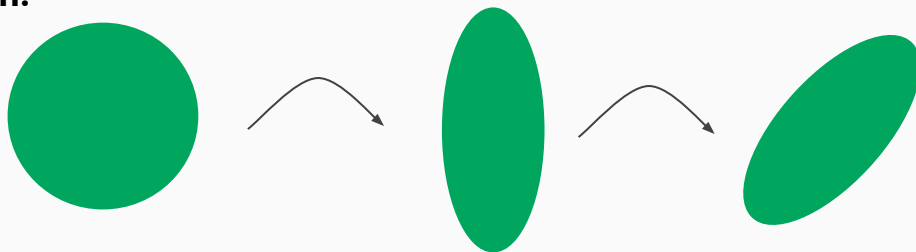




## Certain N-d non-linear normalizing-flow functions

### 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

# Certain N-d non-linear normalizing-flow functions

## 4) Gaussianization Flow(arXiv:2003.01941)

Can also be initialized to quite good values by data!

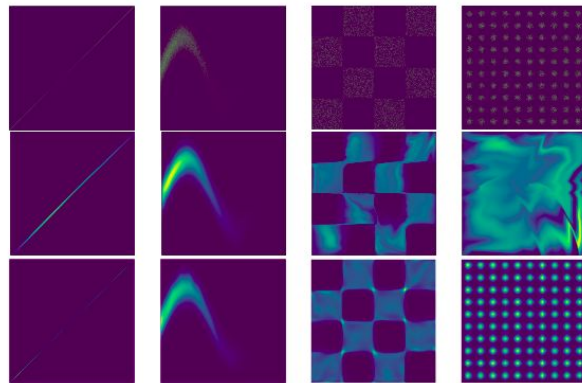
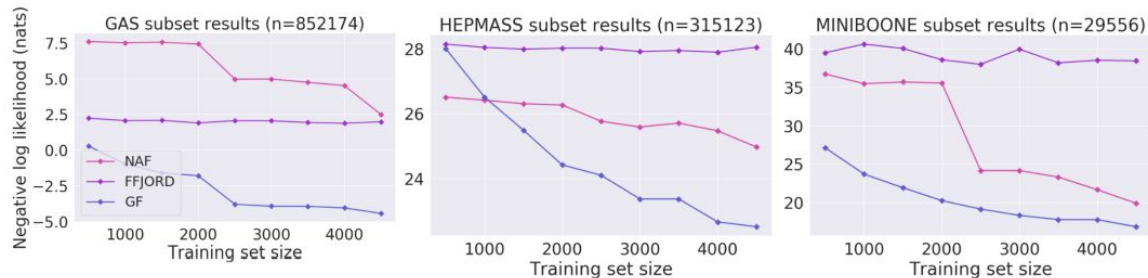


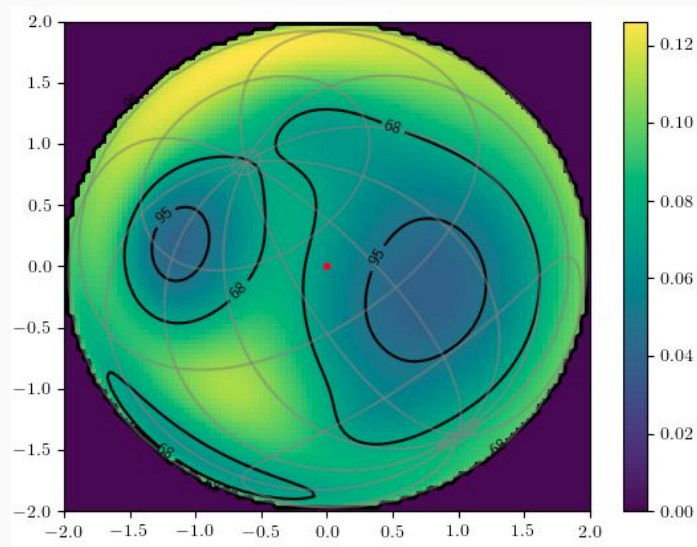
Figure 2: 2D density estimation results. **Top:** Ground truth samples. **Middle:** Glow. **Bottom:** GF.

Method	POWER	GAS	HEPMASS	MINIBOONE	BSDS300	MNIST	FMNIST
Real NVP	-0.17	-8.33	18.71	13.55	-153.28	1.06	<b>2.85</b>
Glow	-0.17	-8.15	18.92	11.35	-155.07	1.05	2.95
FFJORD	-0.46	-8.59	<b>14.92</b>	10.43	<b>-157.40</b>	<b>0.99</b>	-
RBIG	1.02	0.05	24.59	25.41	-115.96	1.71	4.46
GF(ours)	<b>-0.57</b>	<b>-10.13</b>	17.59	<b>10.32</b>	-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	-	-

# Other interesting current research

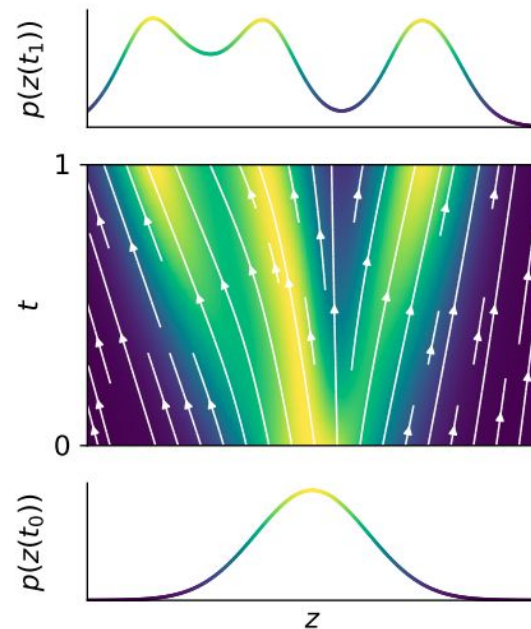
## “Manifold” normalizing flows

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



(2002.02428)

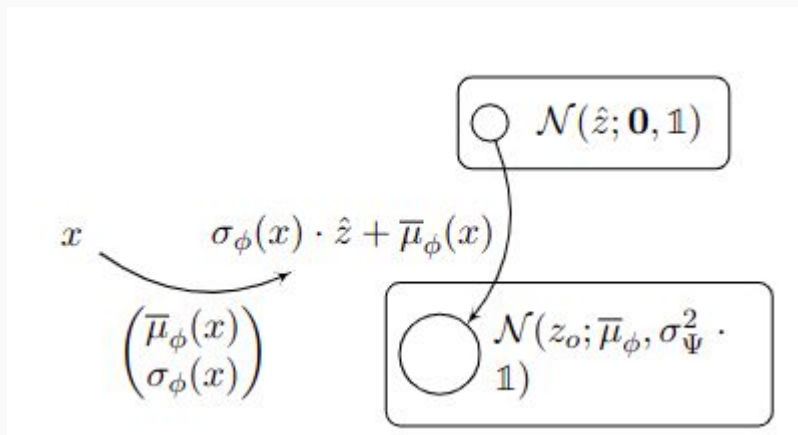
## “Continuous” normalizing flows



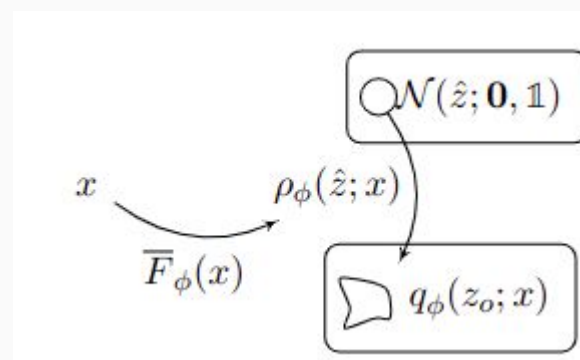
(FFJORD,  
1810.01367)

# Conditional PDFs .. parameters of flow are output of a neural network

## Conditional affine flow



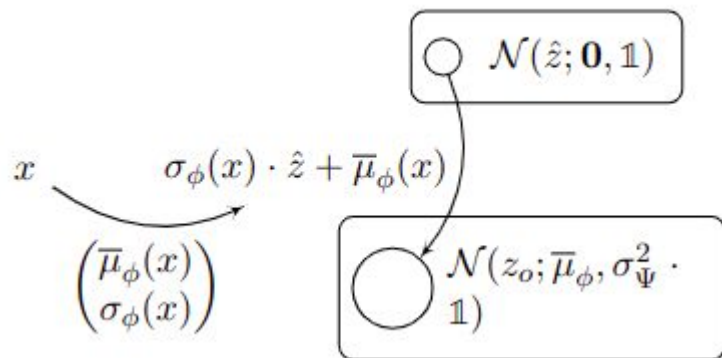
## Conditional general flow



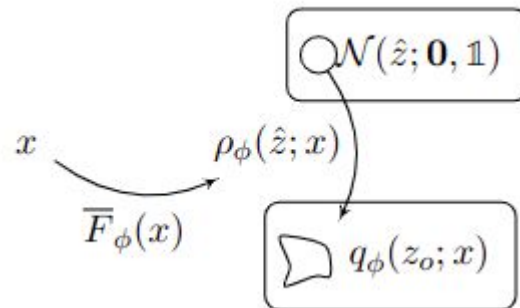
Instead of flow parameters, one optimizes NN parameters

# Conditional PDFs .. parameters of flow are output of a neural network

## Conditional affine flow



## Conditional general flow

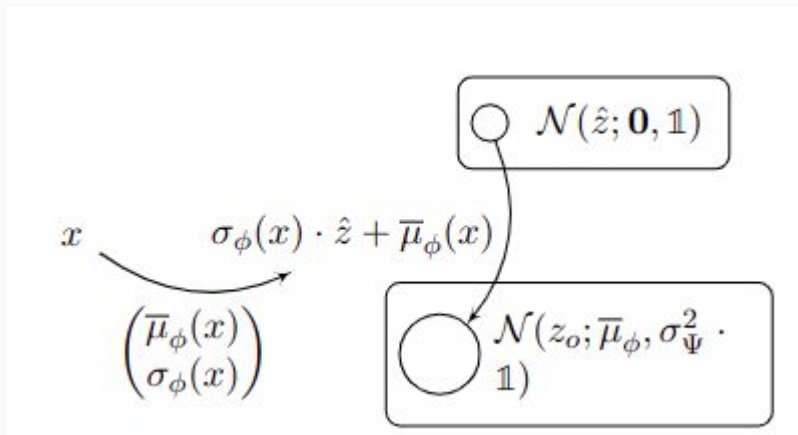


→ Instead of flow parameters, one optimizes NN parameters

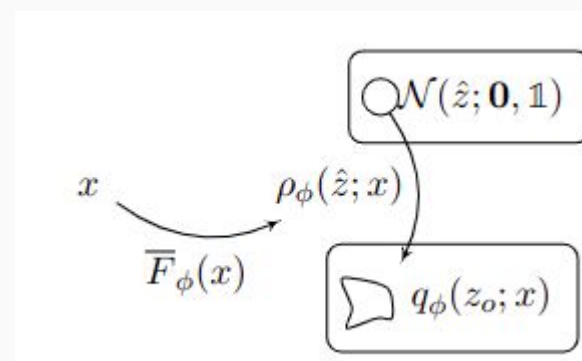
→ Conditional normalizing flow shows that **MSE loss comes from conditional Flow that only consists of a shift** (and unit scaling)  $0.5 \cdot (x - \mu)^2 = \ln(p(x))$

# Conditional PDFs .. parameters of flow are output of a neural network

## 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)  $0.5 \cdot (x - \mu)^2 = \ln(p(x))$
- The process of **predicting parameters by a neural network** is also called “amortization”

A package for flows on manifolds: jammy\_flows ([https://github.com/thoglu/jammy\\_flows](https://github.com/thoglu/jammy_flows))

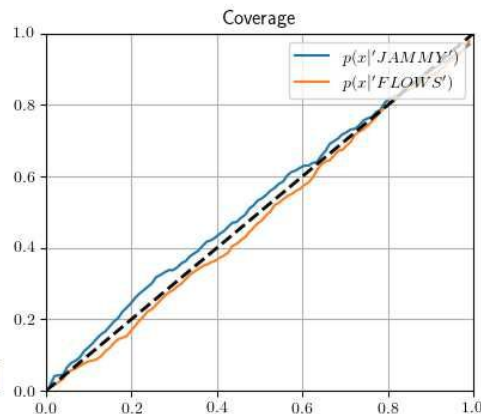
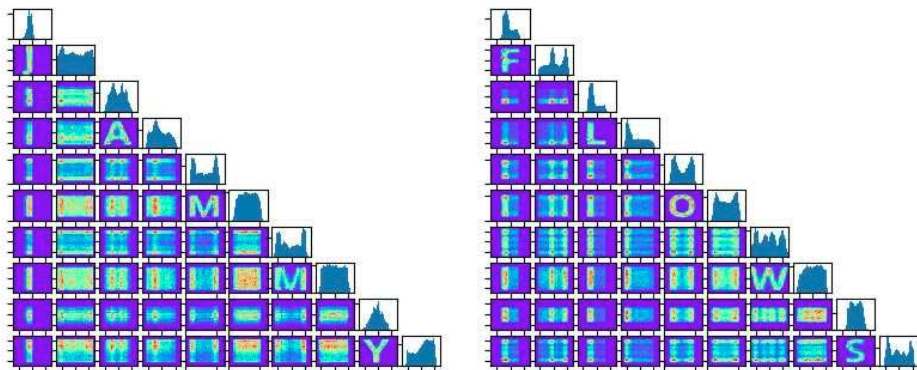
```
import jammy_flows
```

```
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.

pdf structure: e4+s2+e4



# Hands on 1

[https://colab.research.google.com/drive/1ySGCjEeVdKhCodwPxe\\_sw0bBoAymRe4o?usp=sharing](https://colab.research.google.com/drive/1ySGCjEeVdKhCodwPxe_sw0bBoAymRe4o?usp=sharing)



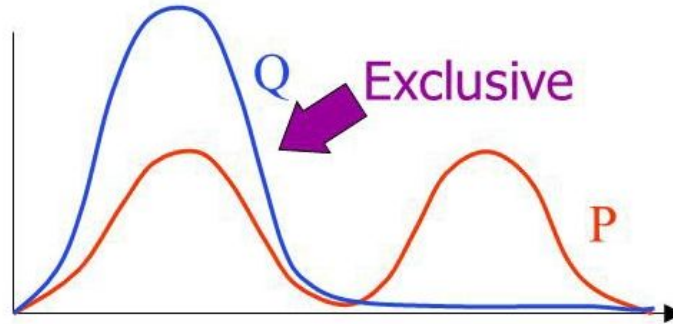
# KL-divergence introduction

“Inclusive” KL-divergence vs “exclusive” KL-divergence

Minimising

$$KL(Q||P)$$

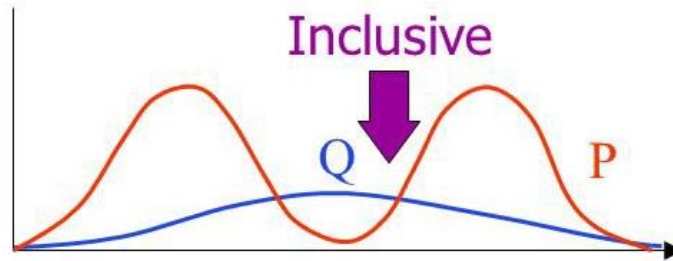
$$= \sum_H Q(H) \ln \frac{Q(H)}{P(H|V)}$$



Minimising

$$KL(P||Q)$$

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

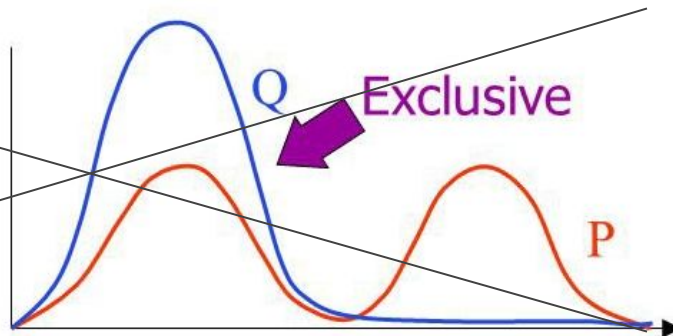


# KL-divergence introduction

“Inclusive” KL-divergence vs “exclusive” KL-divergence

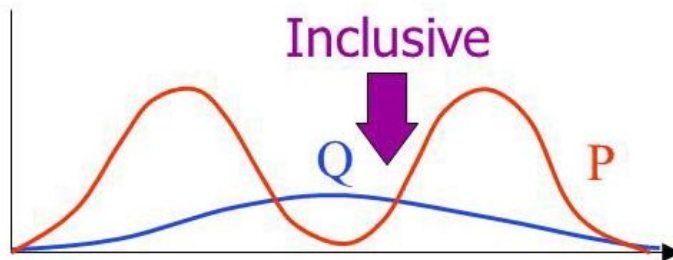
~~Minimising  
 $KL(Q||P)$~~

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



Minimising  
 $KL(P||Q)$

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



In most settings:

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

Q = “Approximating PDF”,  
Parametrized by us,  
“Surrogate model”

# Supervised learning as (inclusive) KL-divergence minimization

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



Computer scientists call this  
“conditional ML objective”  
Of supervised learning

Physicists should call it  
“variational inference objective”  
For the **variational Posterior**

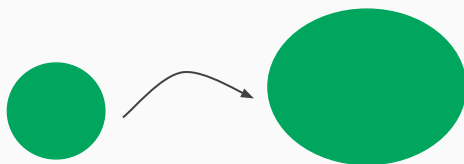
# Supervised learning as (inclusive) KL-divergence minimization

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



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**

# Deep learning generalizes classical statistical approaches

“Bayesian”

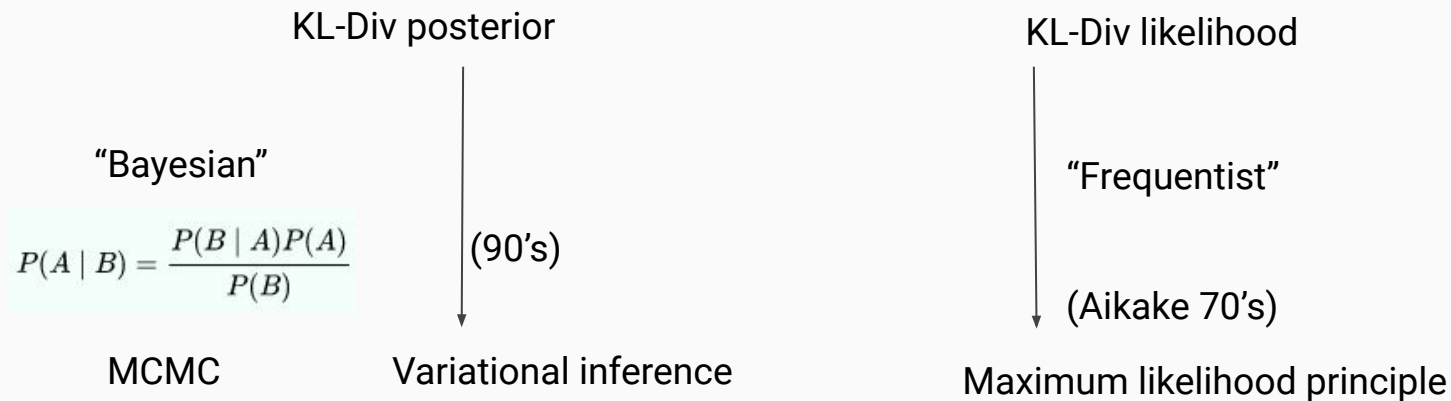
$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

MCMC

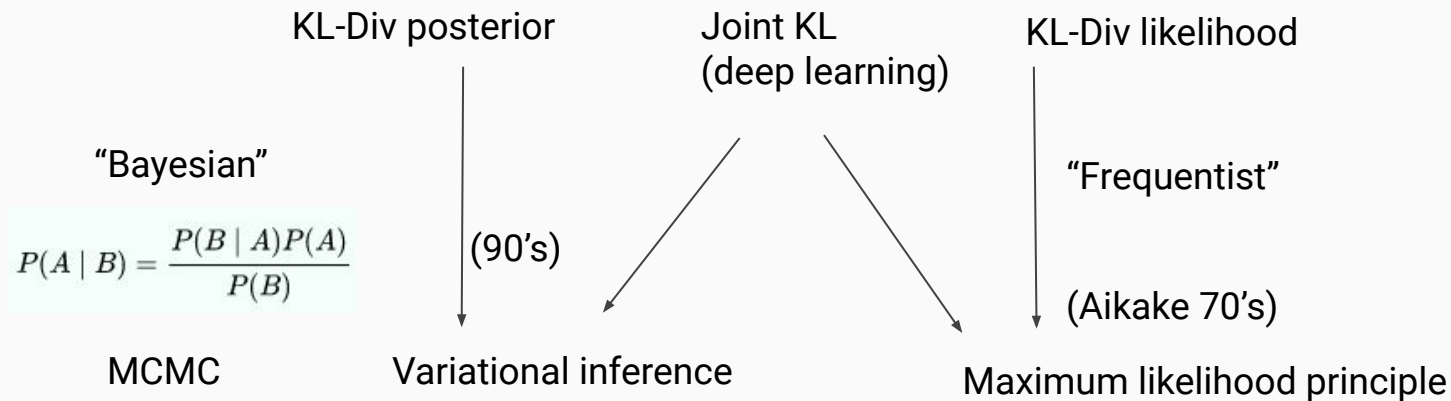
“Frequentist”

Maximum likelihood principle

# Deep learning generalizes classical statistical approaches

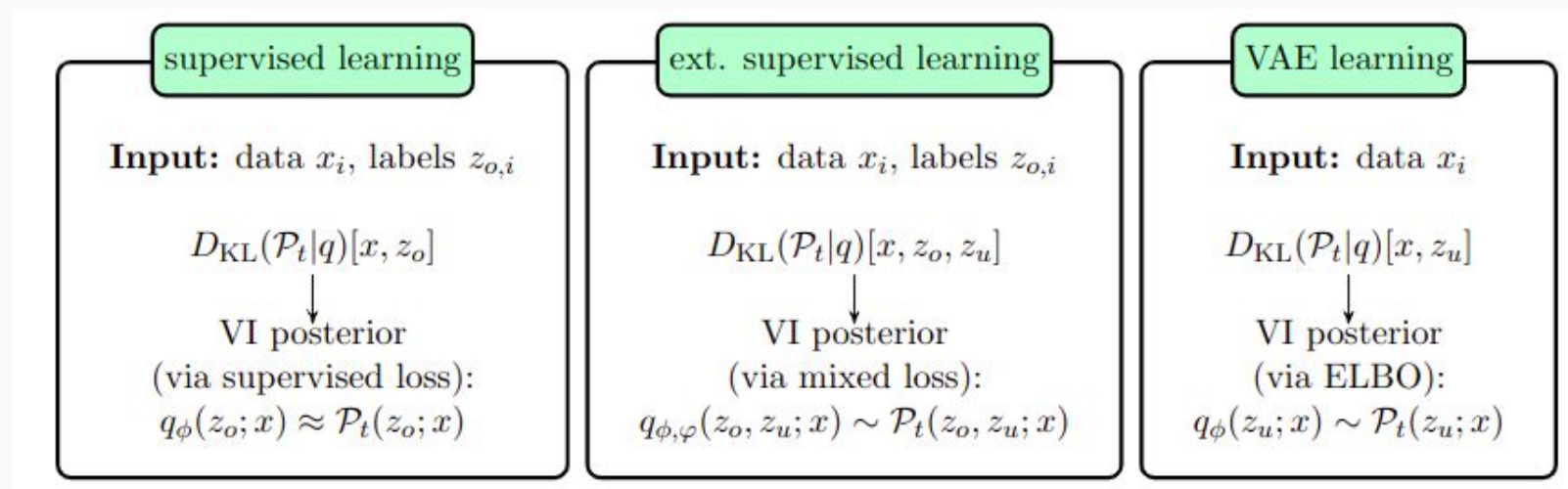


# Deep learning generalizes classical statistical approaches



# Probabilistic deep learning

“All of deep learning is probability distribution matching”



“Inklusive KL divergence”

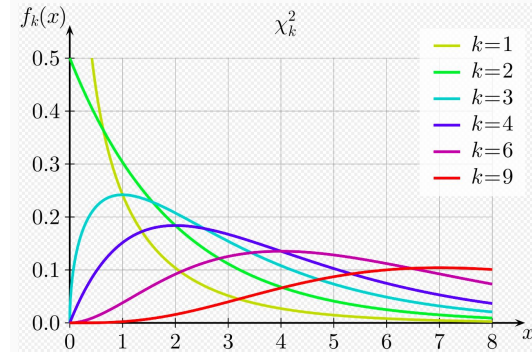
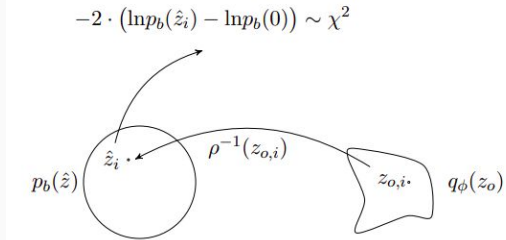
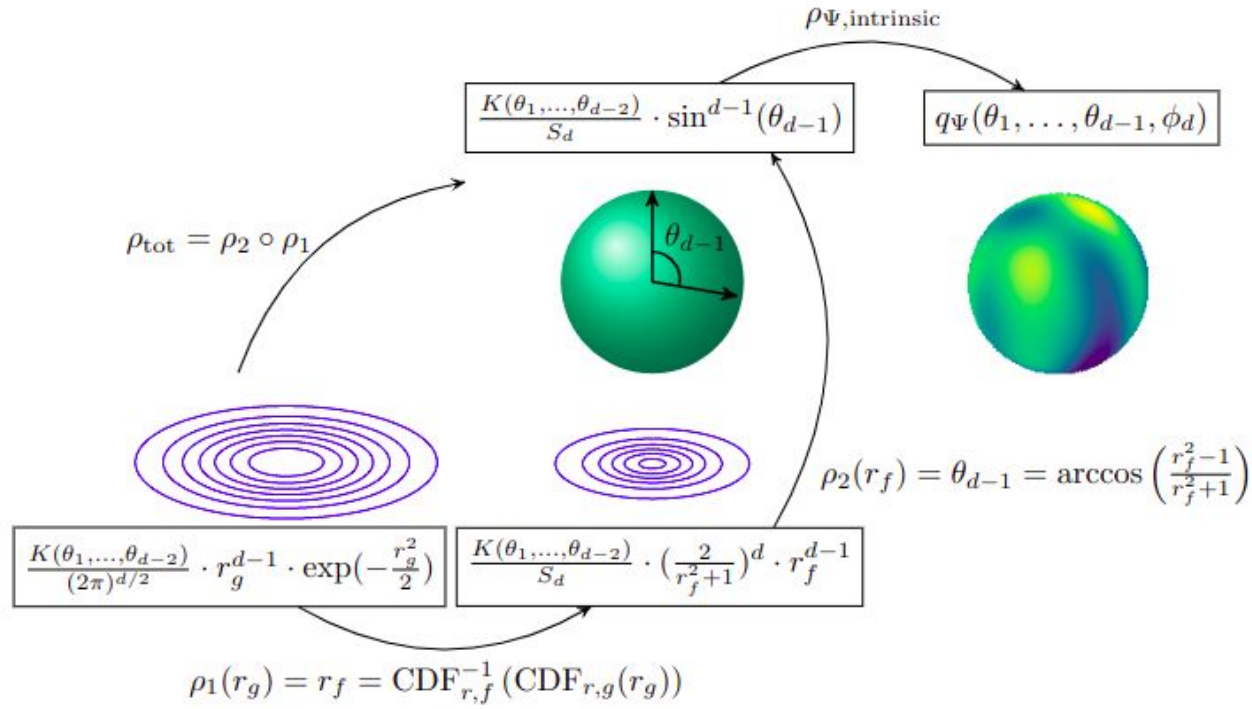
“exclusive KL divergence”

(2008.05825)

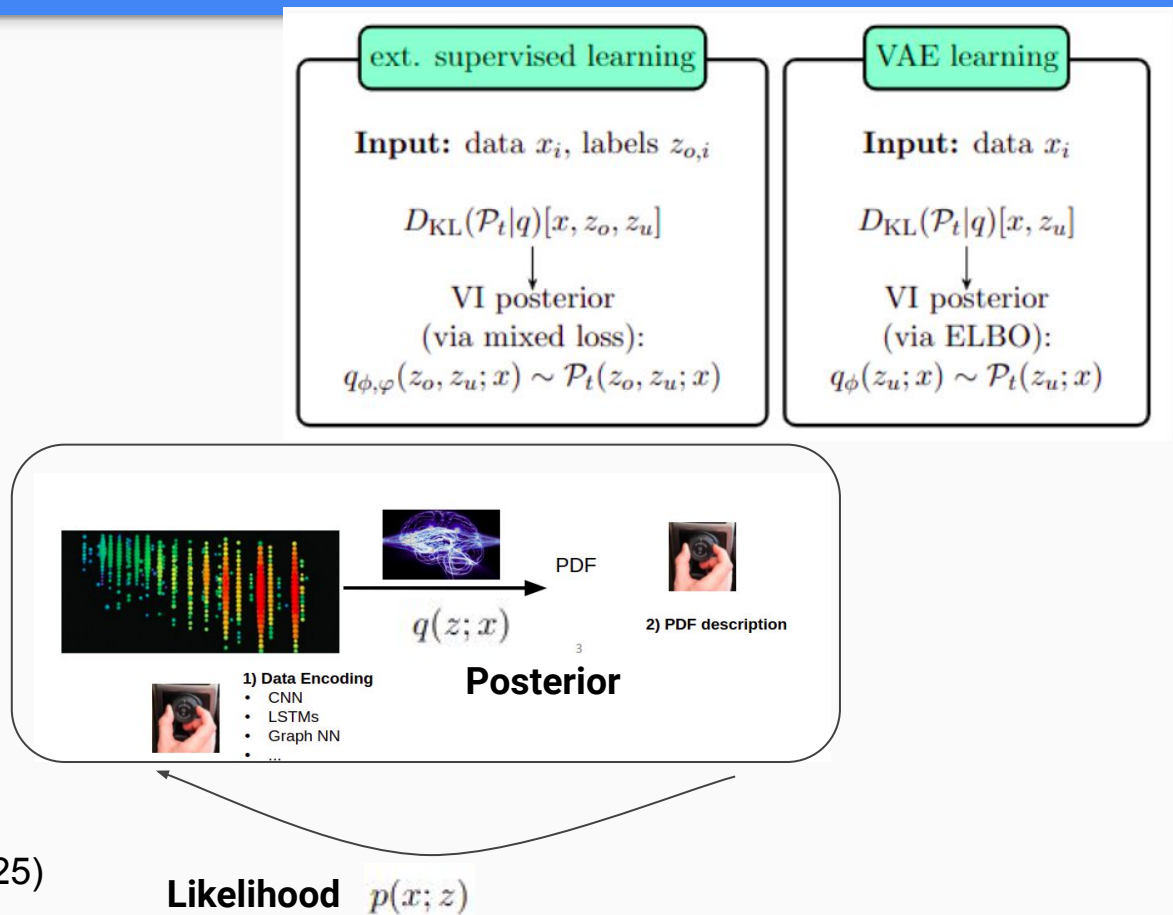




# Coverage for NFs, including NFs on manifolds

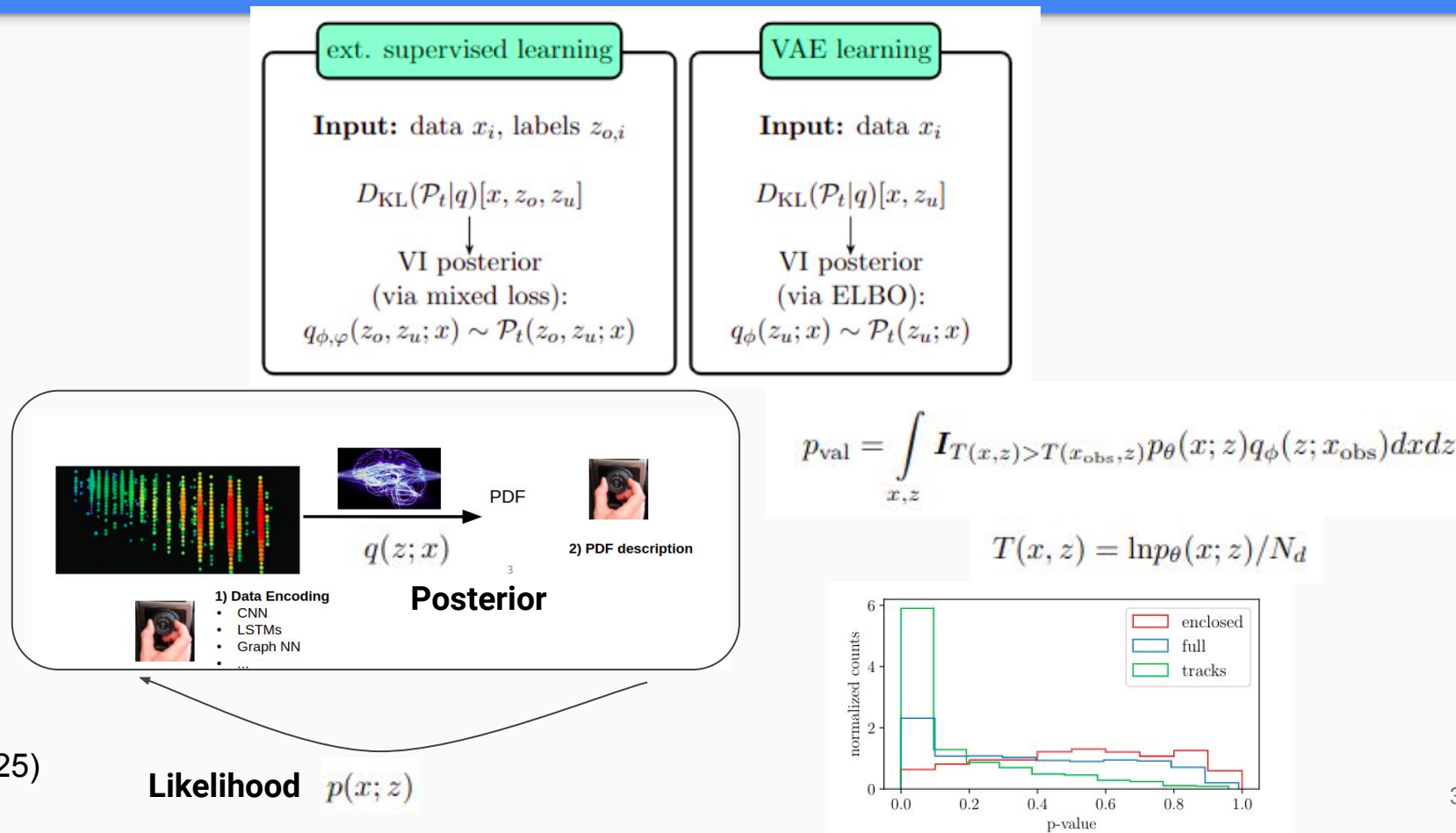


# Goodness of Fit



(2008.05825)

# Goodness of Fit



## Hands on 2

- 1) Check out  
[https://github.com/thoglu/jammy\\_flows.git](https://github.com/thoglu/jammy_flows.git)
- 2) Go to /examples/ subdirectory
- 3) Check out the notebook with examples
- 4) Call script jammy\_flows.py with suitable parameters