



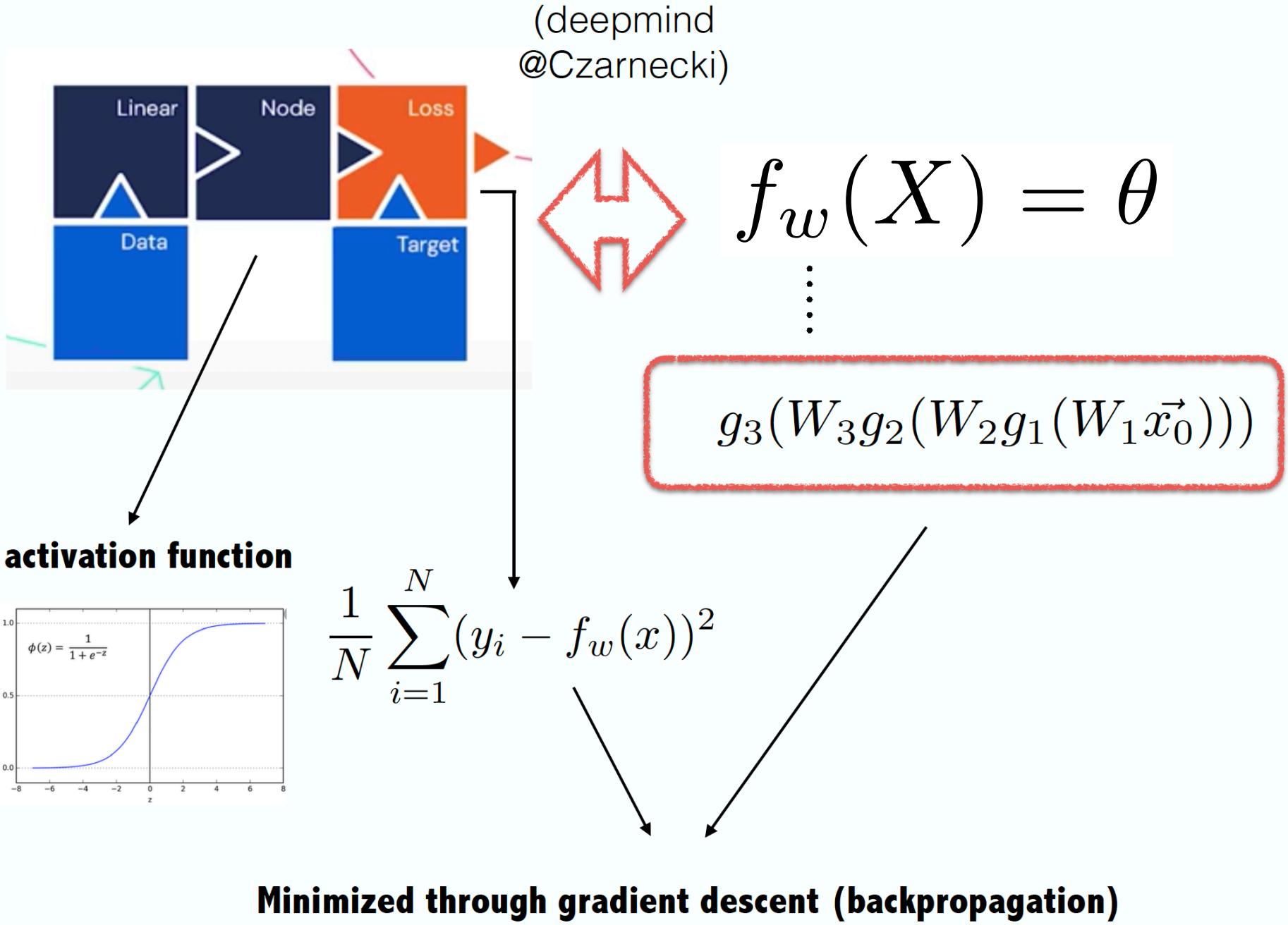
Bayesian deep learning and implicit likelihood inference

Marc Huertas-Company & Hubert Bretonnière

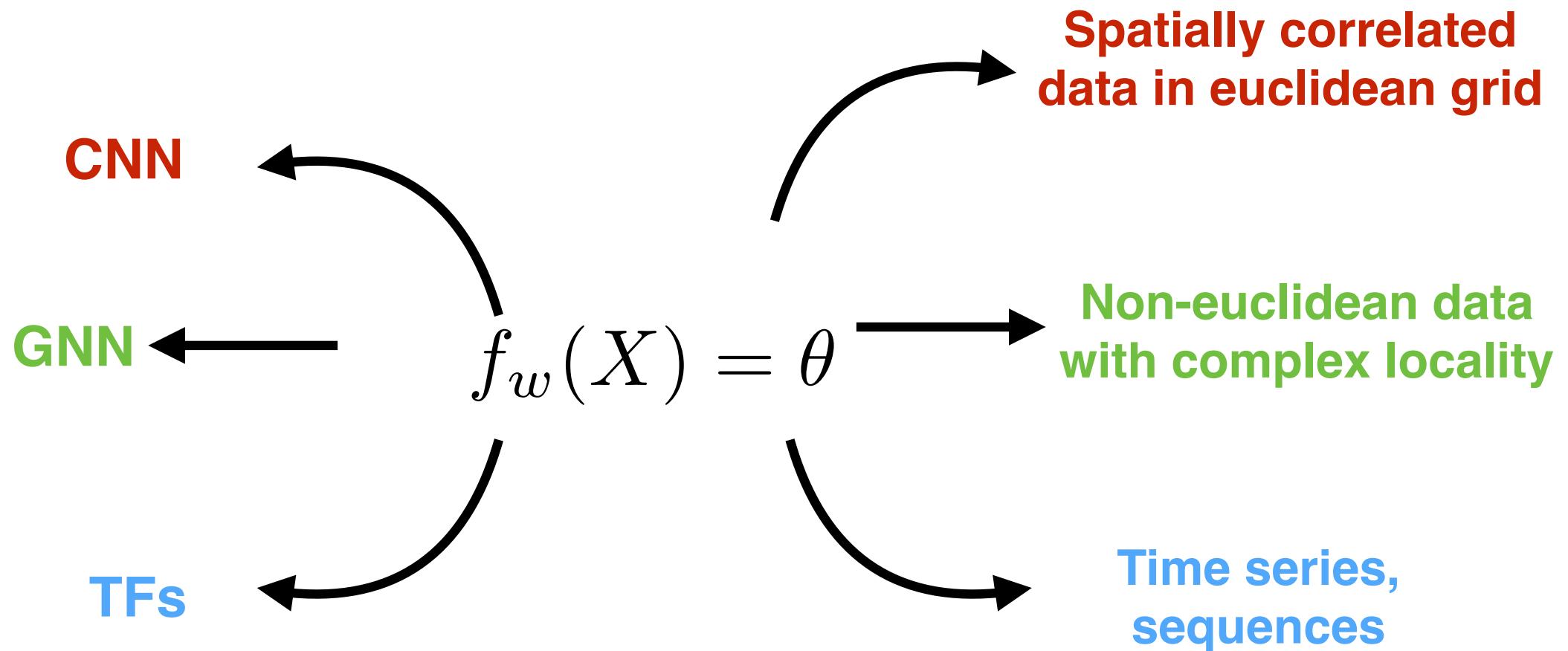
Program

- Basics of bayesian statistics (reminders) and implicit likelihood inference (simulation based inference)
- Neural density estimation techniques for SBI
- Applications to inference in cosmology / astrophysics

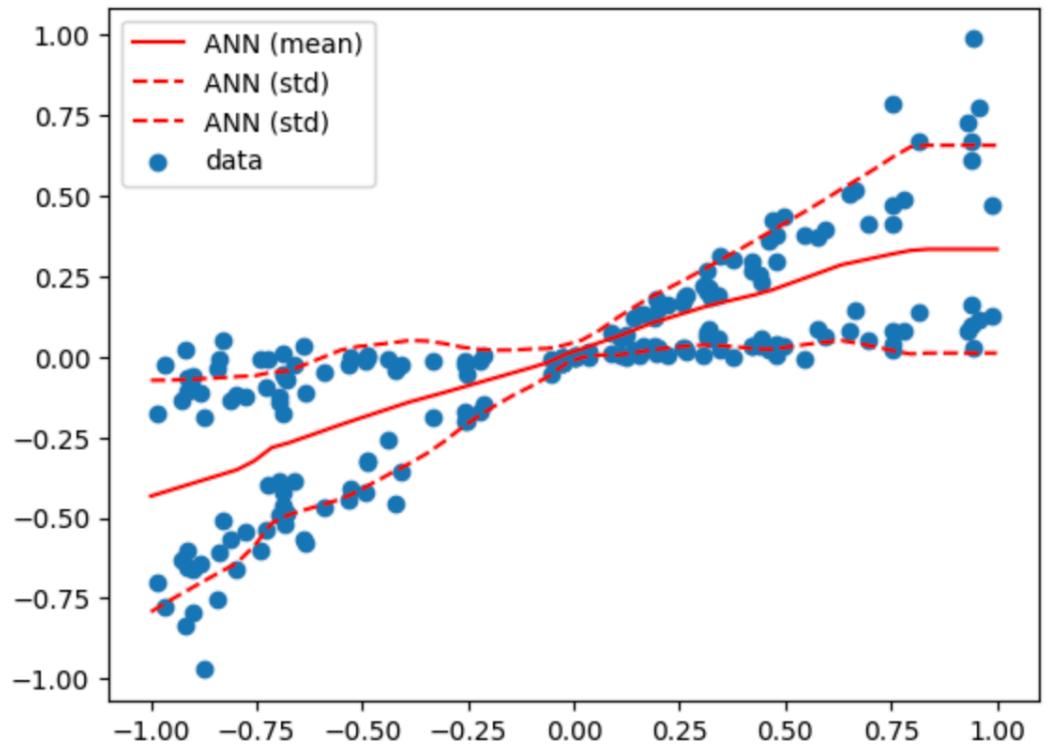
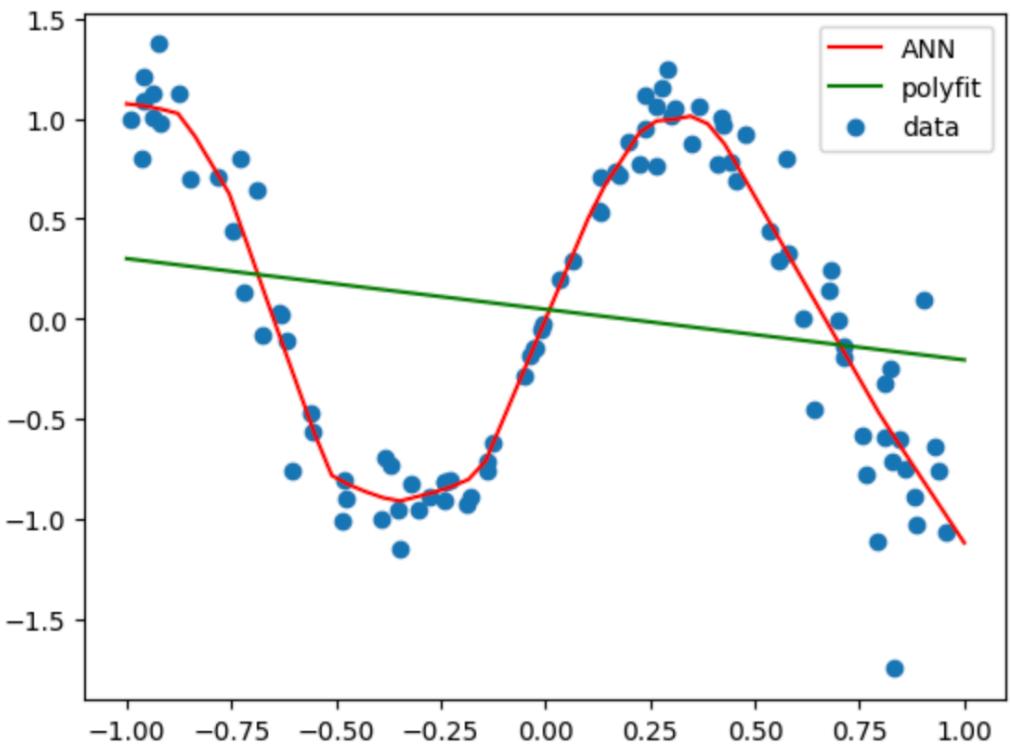
NNs as universal approximators (Y1)



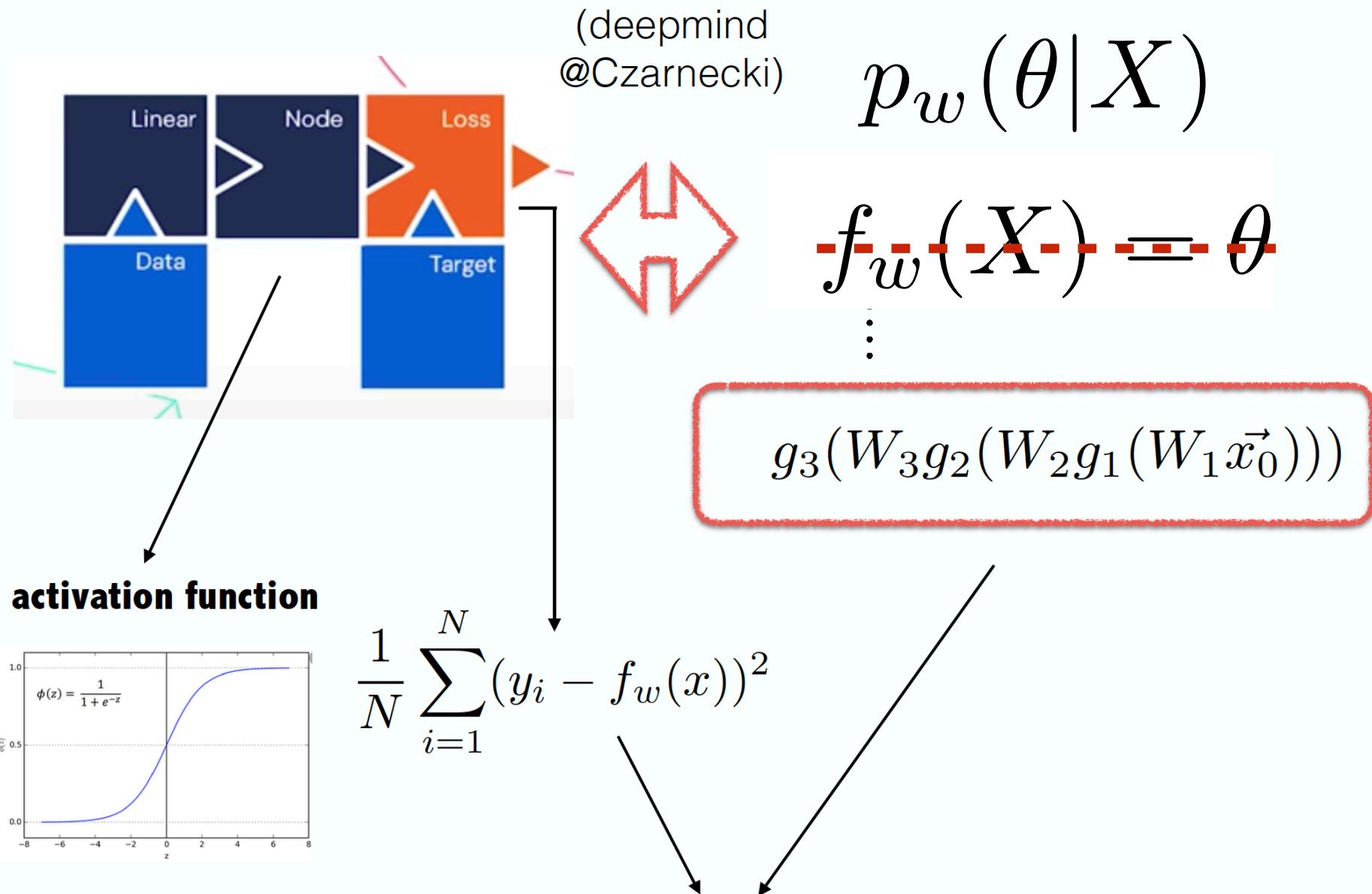
NNs as universal approximators (Y1)



relatively useful...



NNs as statistical models (Y2)



Minimized through gradient descent (backpropagation)

“standard” bayesian inference

goal of inference

$$p(\theta | X) = \frac{\textcolor{orange}{likelihoood} \quad \textcolor{red}{prior}}{p(X)} \quad \textcolor{blue}{posterior}$$

“standard” bayesian inference

goal of inference

$$p(\theta | X) = \frac{likelihood \ prior}{p(X)} \\ \textcolor{blue}{posterior}$$

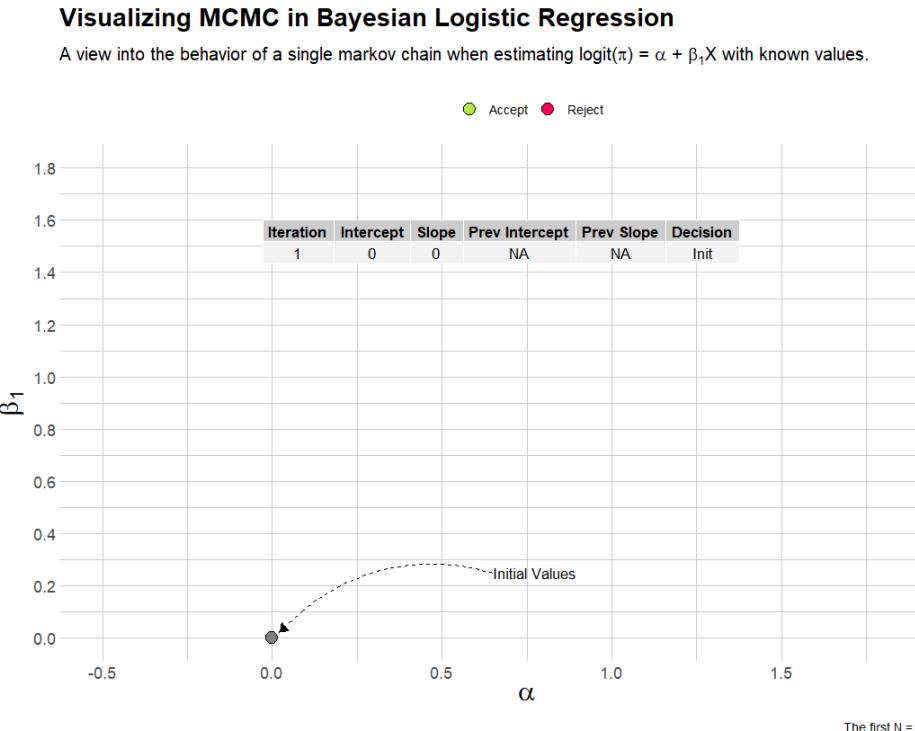
$$\ln p(X | \theta) = \ln \mathcal{L} = (X - m(\theta))^T C^{-1} (X - m(\theta))$$

↑
covariance matrix

model

“standard” bayesian inference

$$\ln p(\theta | X) = \ln p(X | \theta) + \ln p(\theta) + \text{const.}$$



sample $p(\theta | X)$ using montecarlo sampling to estimate posterior
(e.g. *MCMC, HMC, nested sampling*)

“standard” bayesian inference

Gaussian likelihood is often **an incorrect assumption (not enough data)**

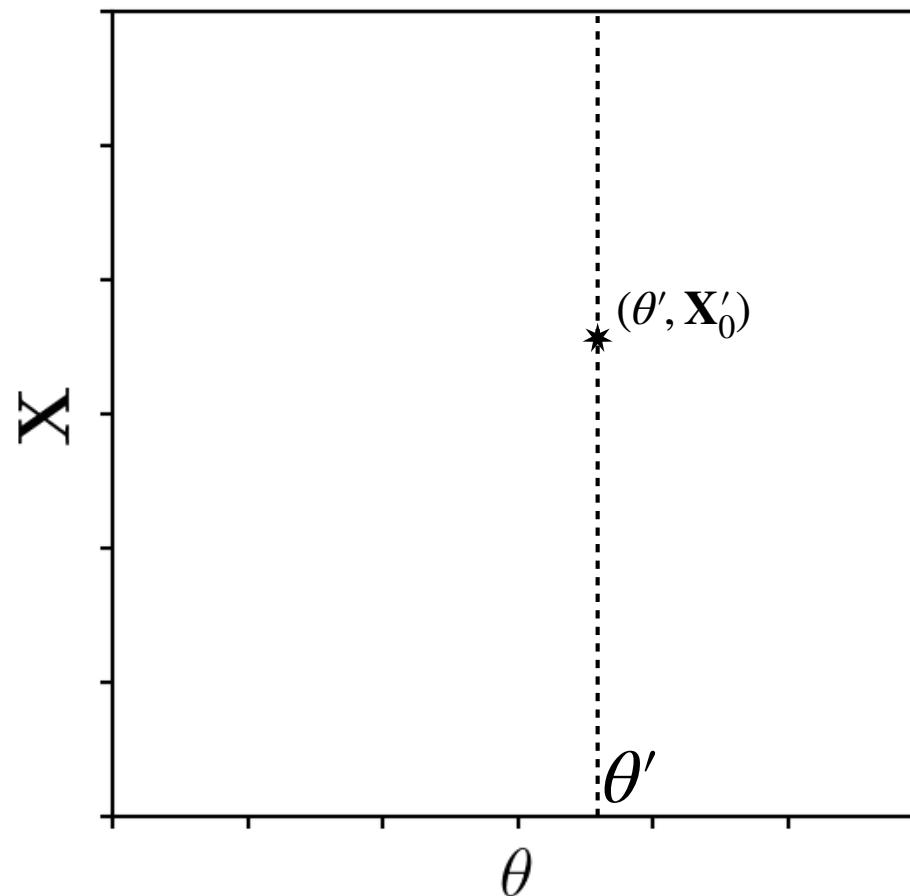
$$\ln p(X | \theta) = \ln \mathcal{L} \neq (X - m(\theta))^T \mathbf{C}^{-1} (X - m(\theta))$$

Sometimes the likelihood is intractable. e.g. Baryonic properties of DM haloes?

MCMC is slow! Impractical to infer *billions* of observations and/or on large parameter spaces.

forward model F that can simulate observations:

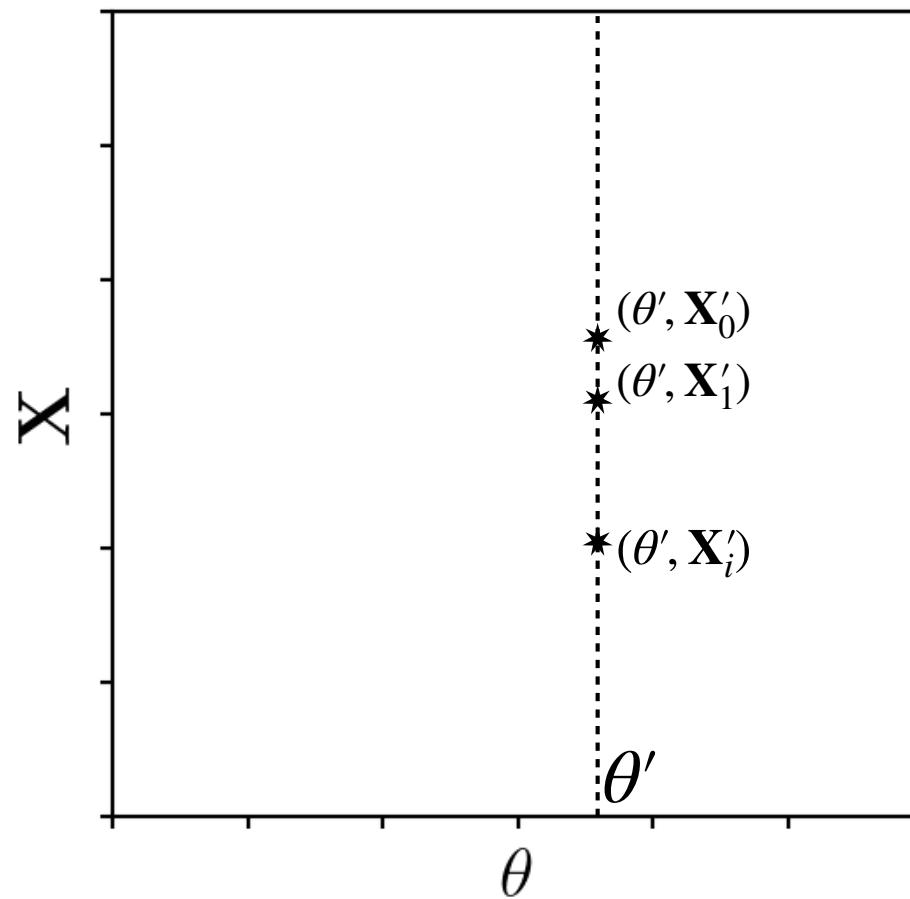
$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



F must **include noise model and observational systematics**

forward model F that can simulate observations:

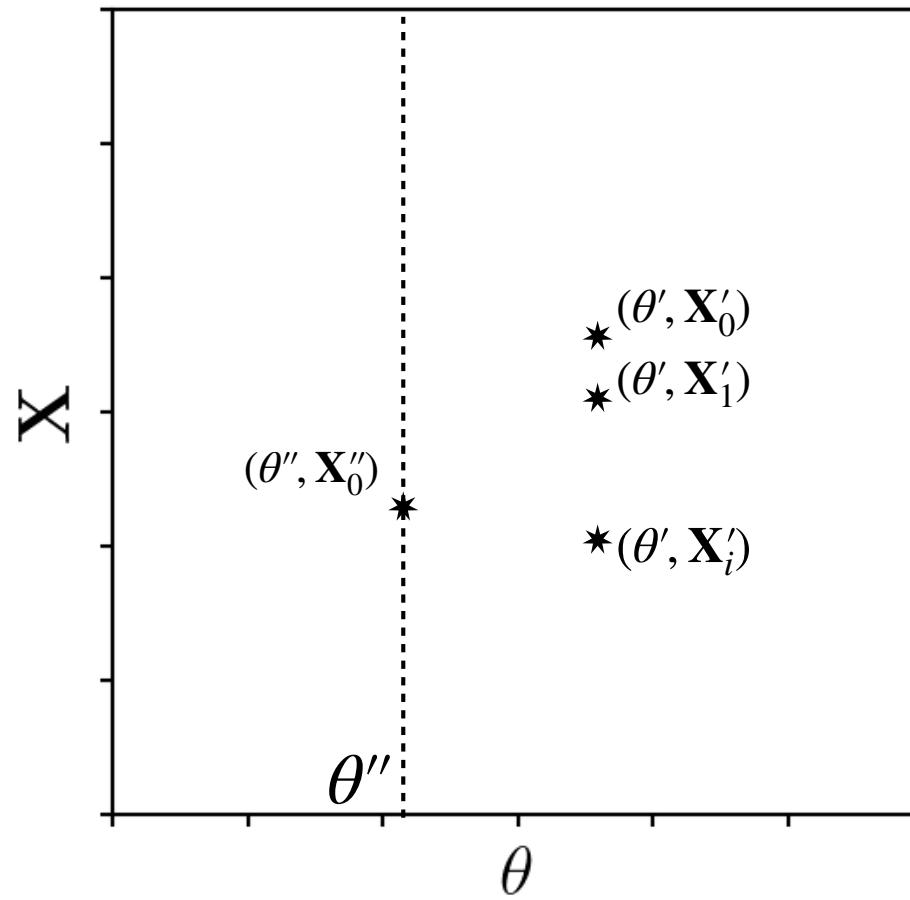
$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



F must include noise model

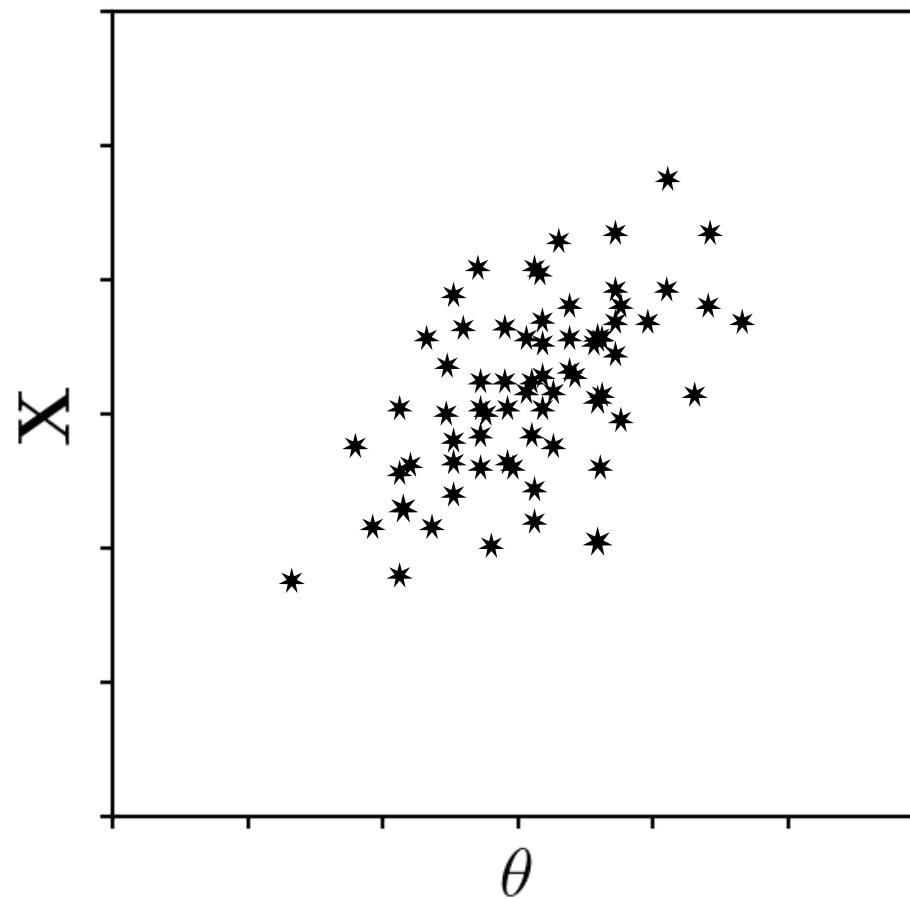
forward model F that can simulate observations:

$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



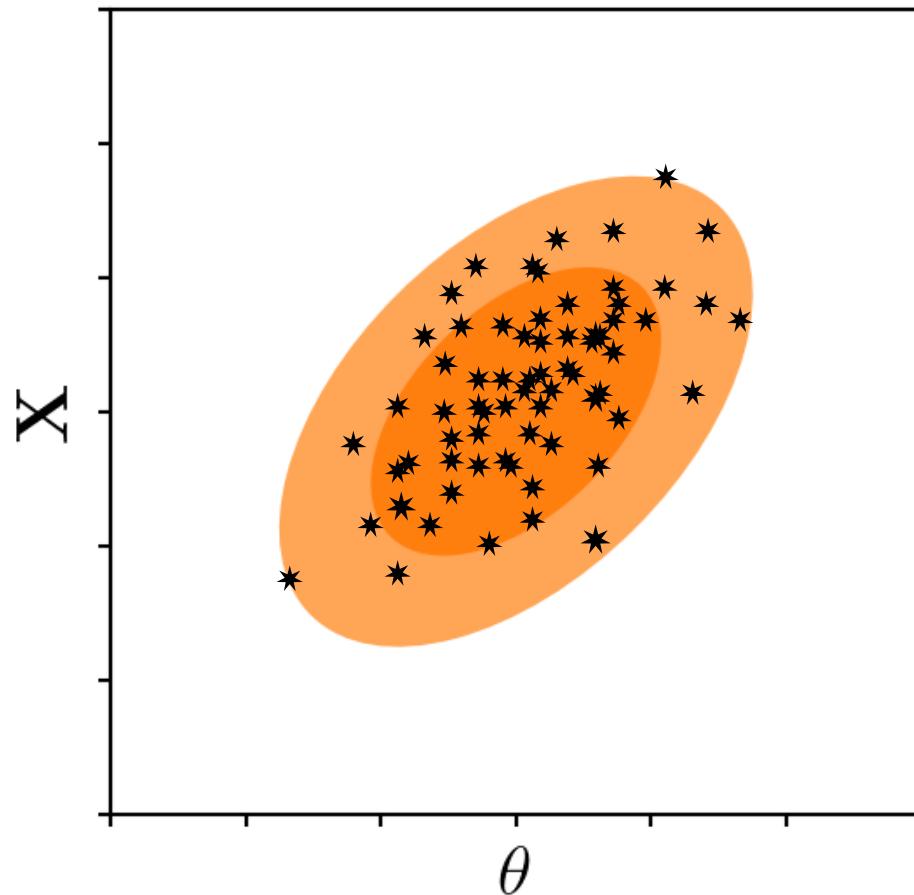
forward model F that can simulate observations:

$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



forward model F that can simulate observations:

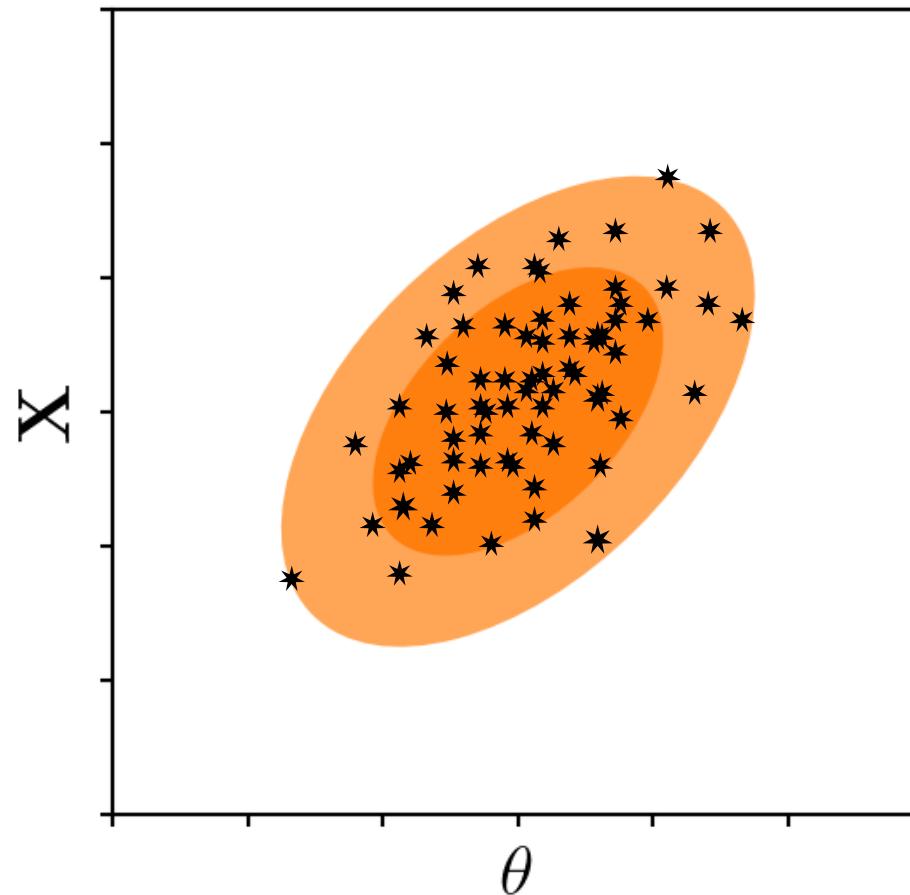
$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



F allows us to sample $(\theta', X') \sim p(\theta, X)$ ← Unknown

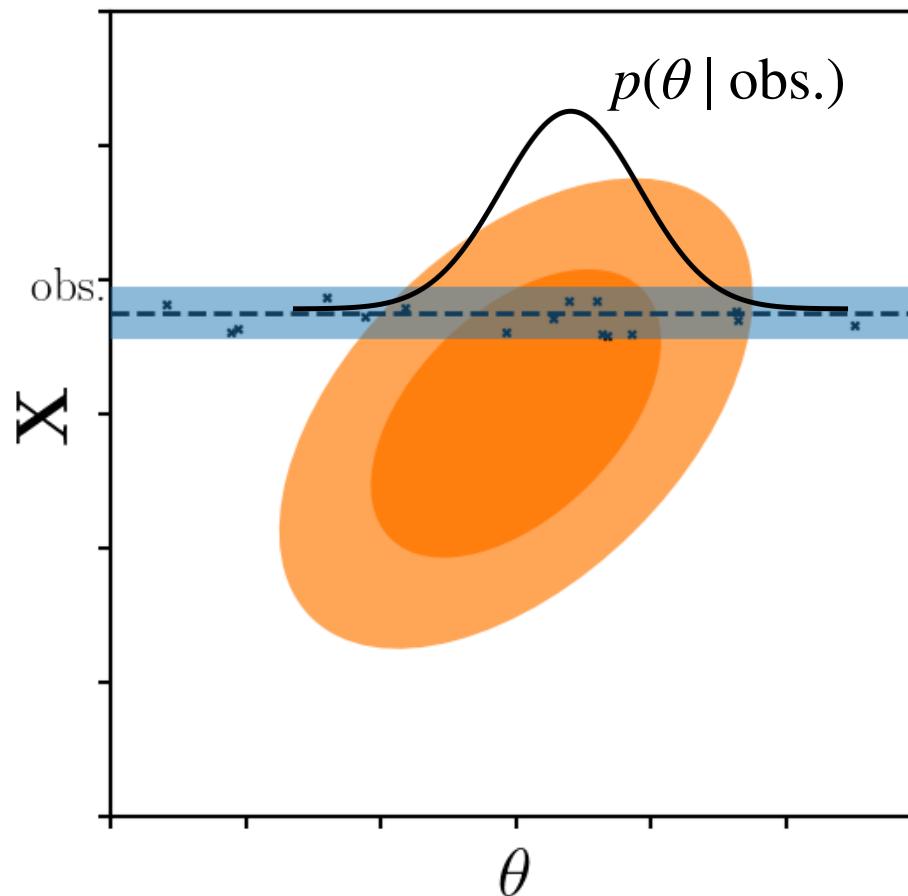
forward model F that can simulate observations:

$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



F allows us to sample $(\theta', X') \sim p(\theta, X)$ ← Unknown

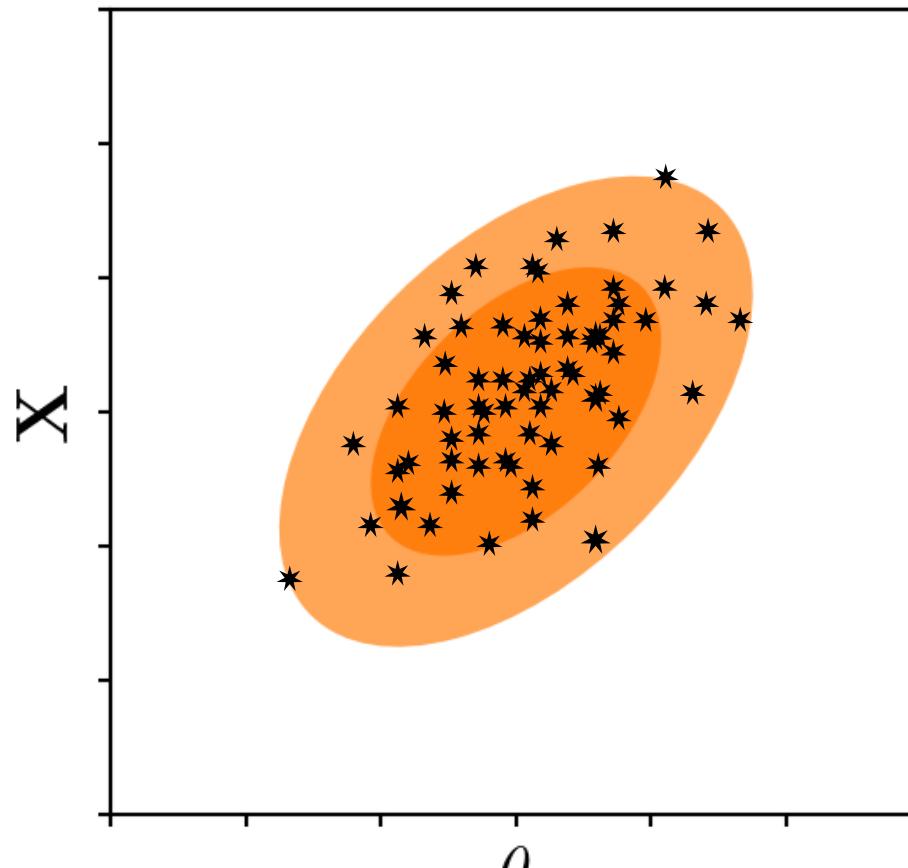
with F one way to infer the posterior is through *brute force* —
approximate bayesian computation (ABC)



only keep the simulations “close” to the observation

(This is called simulation-based inference, likelihood-free inference or implicit inference)

simulation-based inference is a **density estimation problem!**

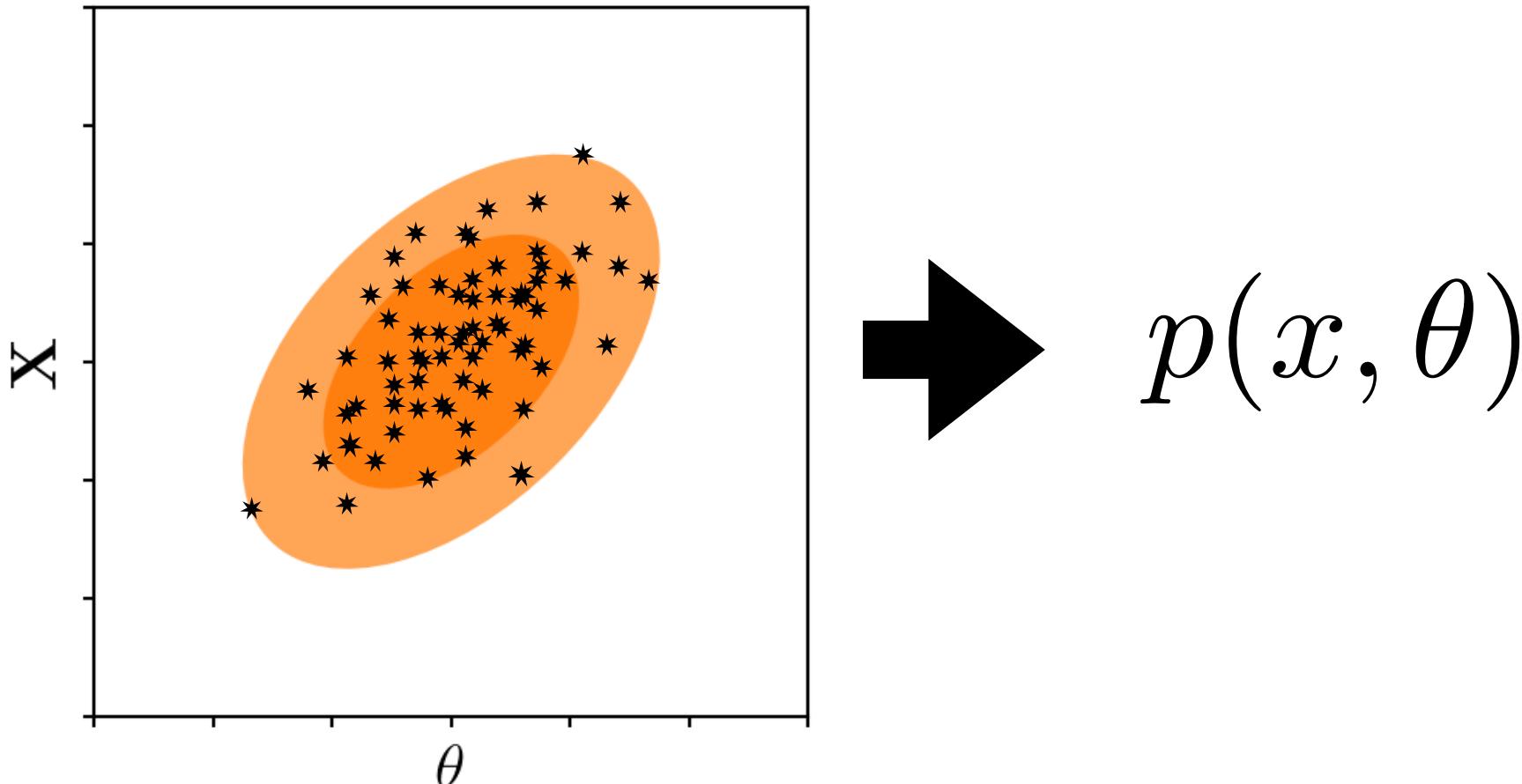


F allows us to sample $(\theta', X') \sim p(\theta, X)$

$$p(\theta, X)$$

Unknown

Density estimation



Given this dataset, the goal of density estimation is to fit a model to the data distribution such that we can synthesize new data points at will by sampling from the distribution.

- Assuming we have a perfect simulator, we have (data, parameter) pairs. How do we do inference?

$$P(\theta|x) \propto P(x|\theta) P(\theta)$$

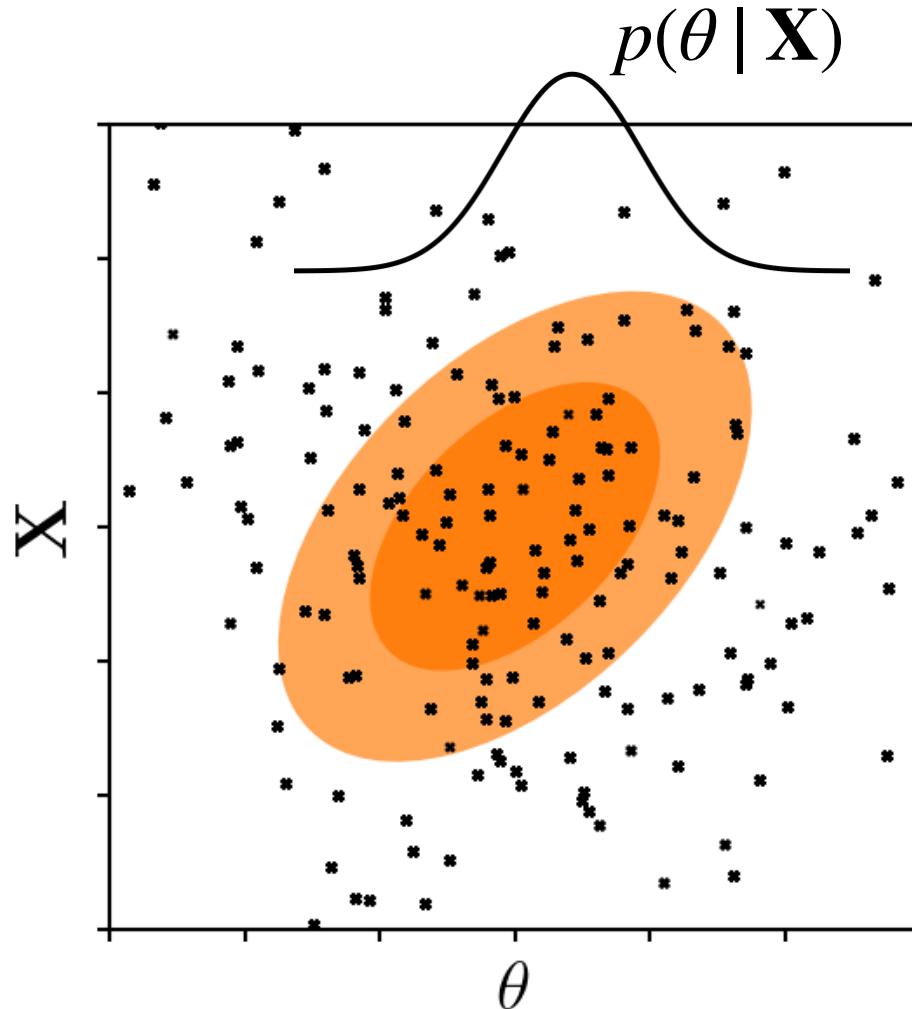
$\underbrace{P(\theta|x)}$ $\underbrace{P(x|\theta)}$ $\underbrace{P(\theta)}$

“Posterior” “Likelihood” “Prior”

Neural Posterior Estimation

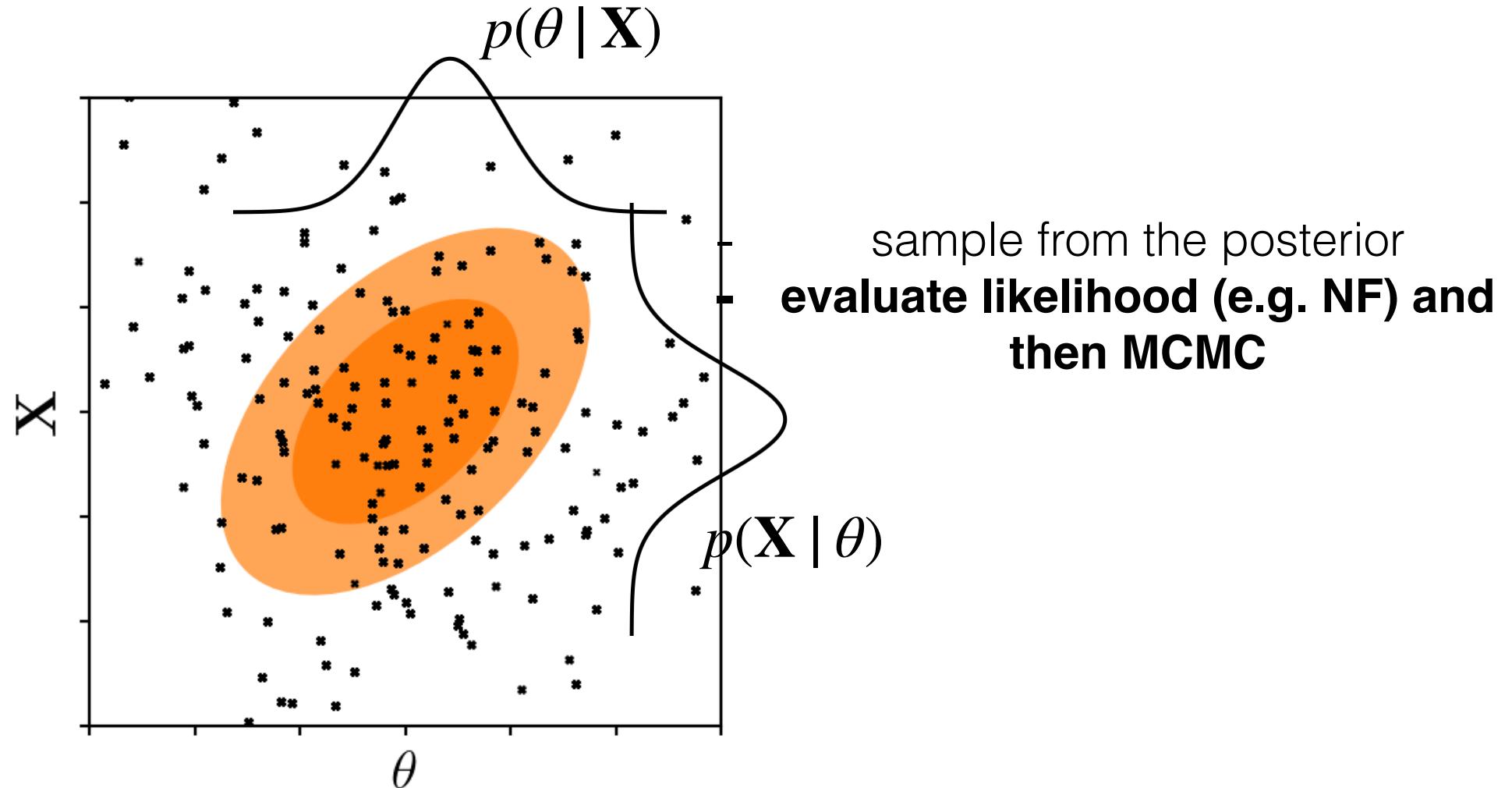
Neural Likelihood Estimation

simulation-based inference is a **density estimation problem!**



- sample from the posterior (amortized)

simulation-based inference is a **density estimation problem!**



- Assuming we have a perfect simulator, we have (data, parameter) pairs. How do we do inference?

$$P(\theta|x) \propto P(x|\theta) P(\theta)$$

↓ ↓ ↓
 “Posterior” “Likelihood” “Prior”

Neural Posterior Estimation

Neural Likelihood Estimation

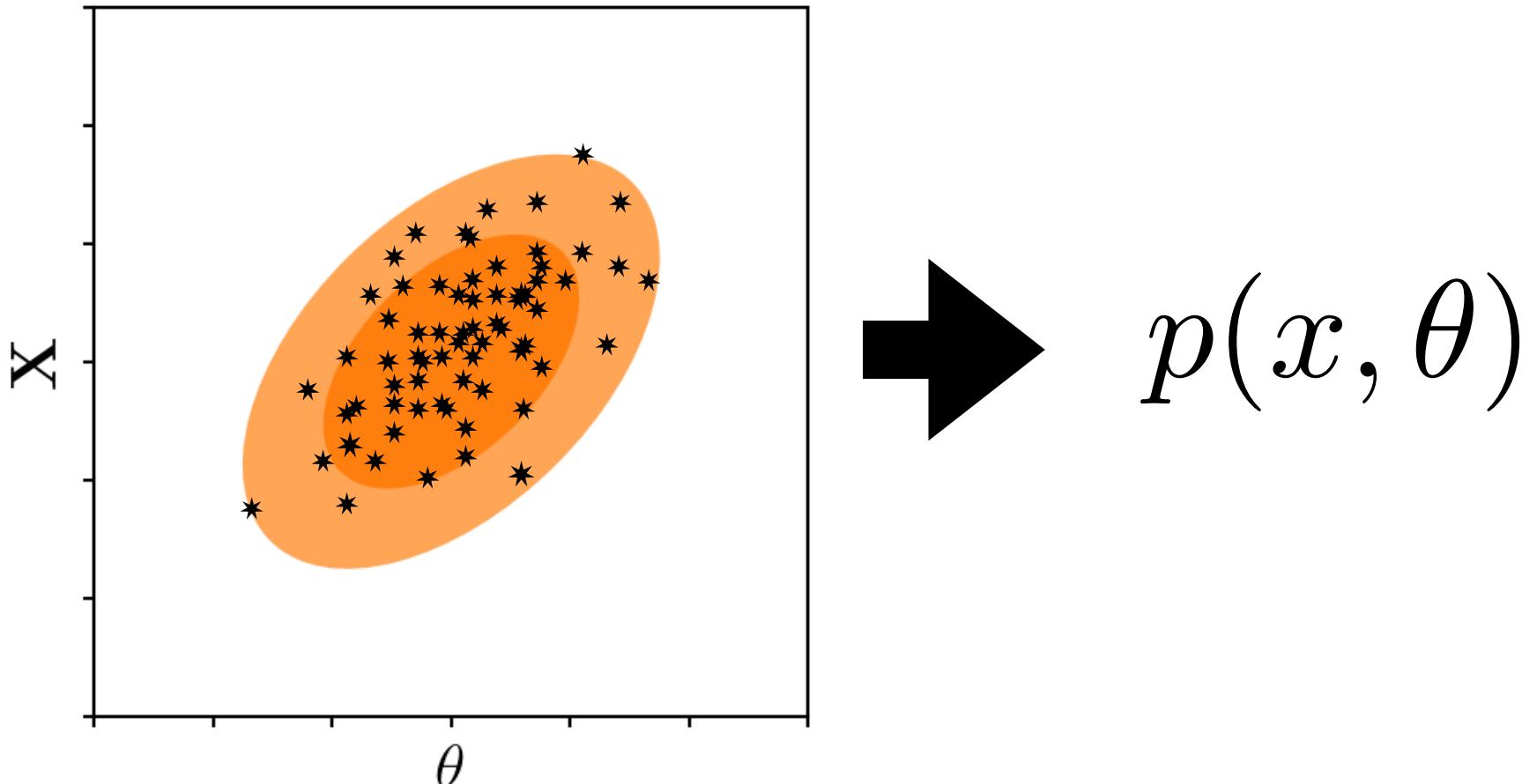
Amortized

(Need to sample many
posteriors for many X's)

Requires sampling (e.g. MCMC)

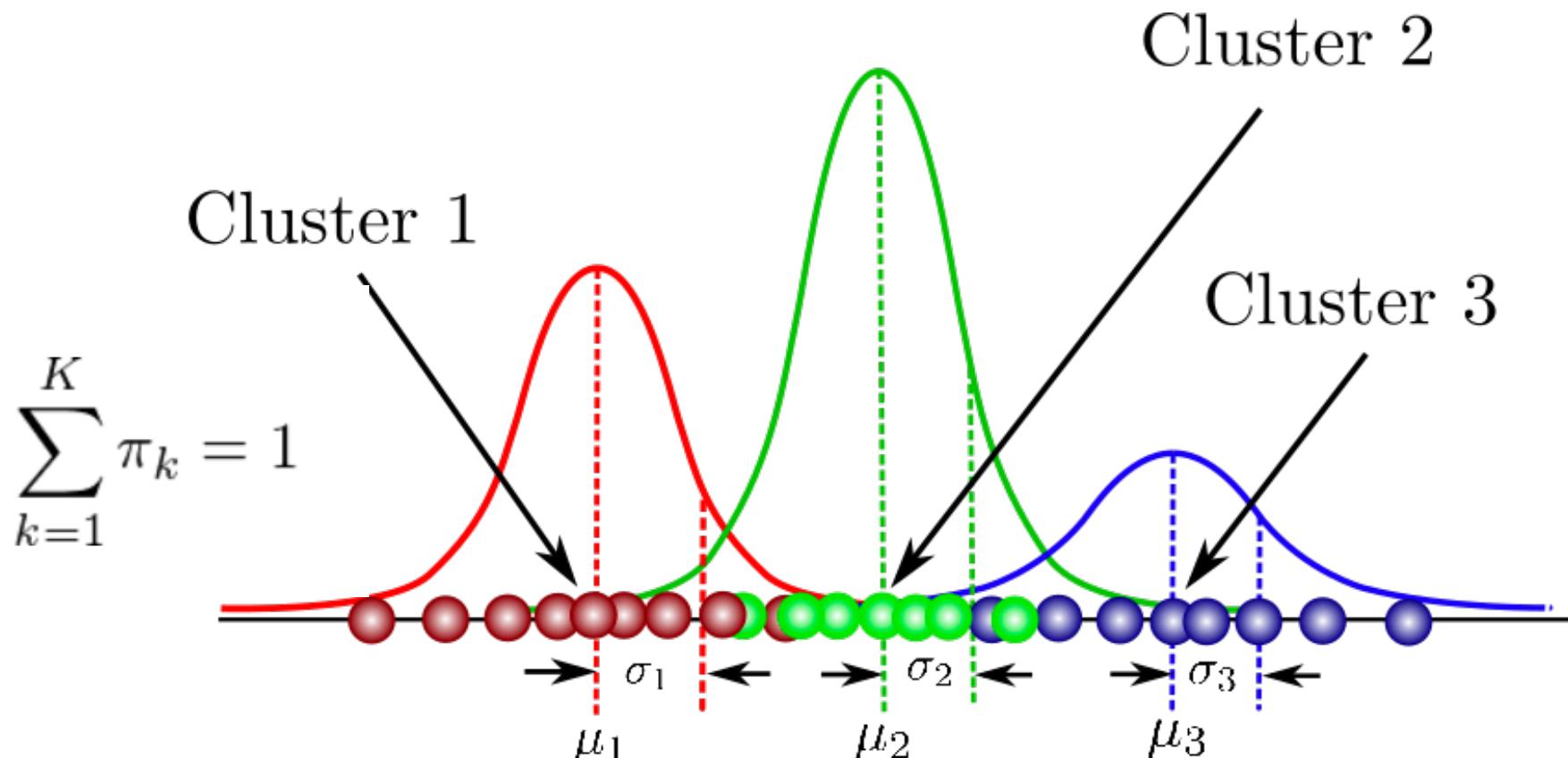
(Need to build model into hierarchical
sampling scheme)

Density estimation



Given this dataset, the goal of density estimation is to fit a model to the data distribution such that we can synthesize new data points at will by sampling from the distribution.

Standard kernel density can be used for that (GMMs)



means, sigmas and scale factors of each gaussian are free parameters

GMMs

- + Both sampling and evaluating are straightforward with GMMs
- – Scales poorly with increasing dimensionality
- – Depends on initial choices
- – Scales poorly for heavily multimodal distributions

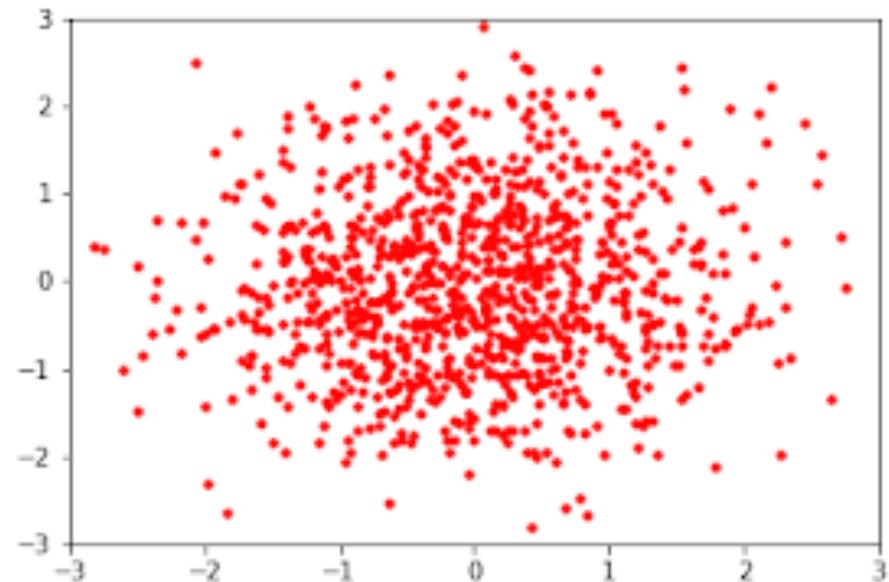
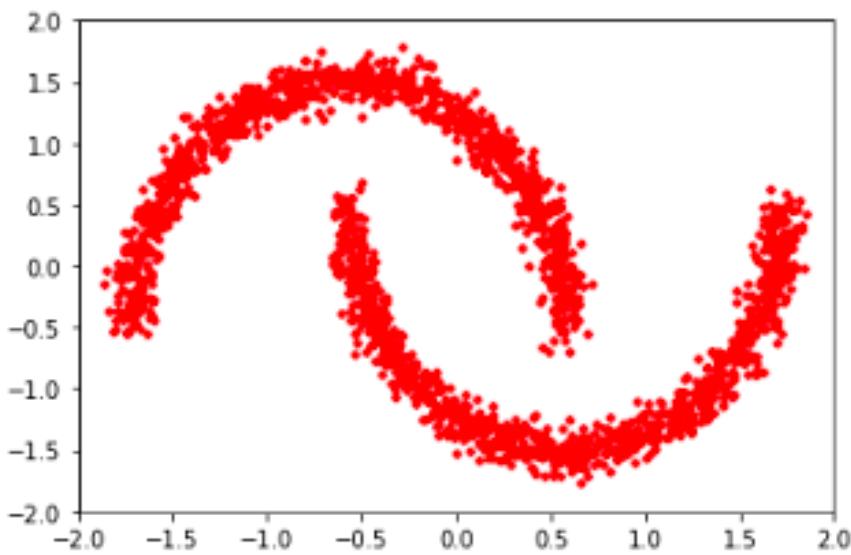
The deep learning revolution
has enabled fast density
estimation at very high
dimension

Pixel space!

Pixel space

$$p(\theta | x)$$

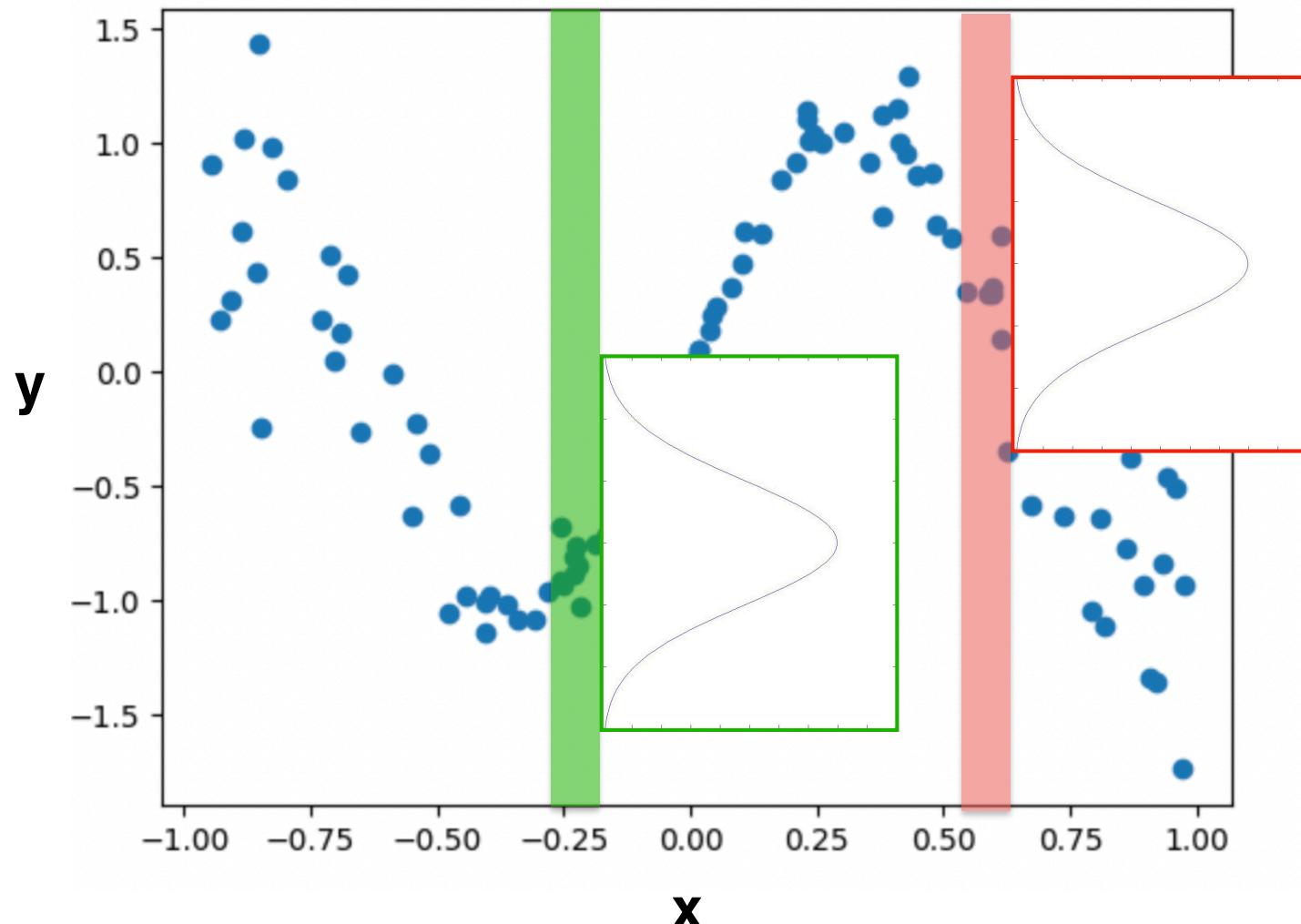
Not Gaussian



Neural Networks as statistical models

$$f_W(\vec{x}) = \vec{y} \quad \textcolor{red}{\leftrightarrow} \quad p(y|x; W)$$

$$p(y|x) = \mathcal{N}(\mu, 1)$$



Assume that the posterior distribution can be approximated by a Gaussian function with fixed standard deviation

we have a set of independent realizations:

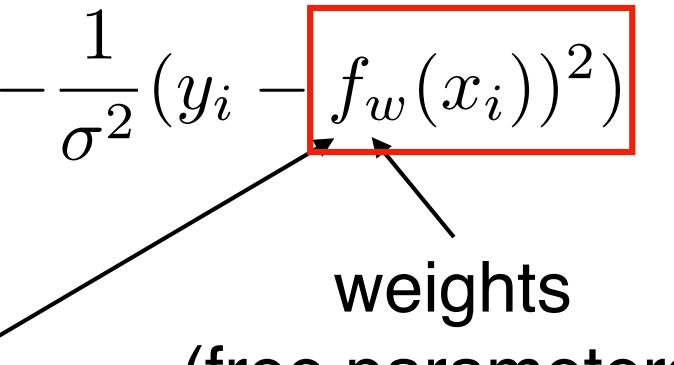
$$(x_i, y_i); i = 1, \dots, N$$

and want to find the underlying probability distribution of y given x -
p(y|x):

If we assume that $p(y|x)$ is a Gaussian distribution:

$$p_w(y|x) = L(w, \sigma^2) = \prod_{i=1}^N (2\pi\sigma^2)^{-N/2} \exp\left(-\frac{1}{\sigma^2}(y_i - f_w(x_i))^2\right)$$

weights
Network output (free parameters)



So the goal is to find the values of w (weights) that estimate the mean of a Gaussian distribution for y given a set of N independent observations

And then take the log:

$$\ln L = cst - \sum_{i=1}^N \frac{(x_i - f_w(x_i))^2}{2\sigma^2}$$

And now assume sigma = 1:

$$\ln L = - \sum_{i=1}^N (x_i - f_w(x_i))^2 = MSE$$

Regressions

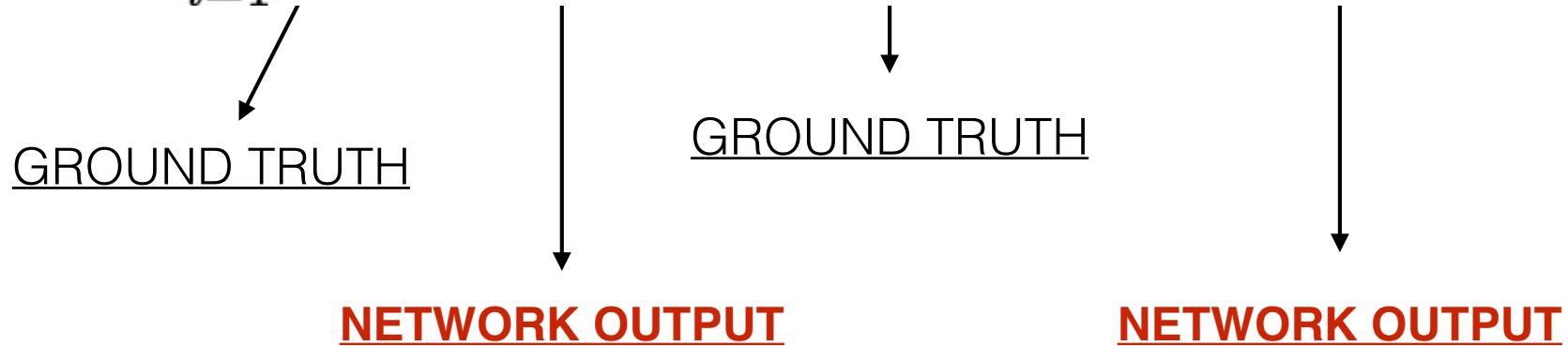
WHEN USING THE MEAN SQUARED ERROR, WE ARE ASSUMING THAT THE POSTERIOR $P(y|x)$ IS A NORMAL DISTRIBUTION WITH STANDARD DEVIATION EQUAL TO 1.

OUR PDF IS THEREFORE FULLY CHARACTERIZED BY ONE UNIQUE PARAMETER: THE MEAN OF THE NORMAL PDF.

THE MEAN IS ESTIMATED BY MAXIMIZING THE LIKELIHOOD OF THE OBSERVATIONS.

What about classification?

$$H = -\frac{1}{N} \sum_{i=1}^N y_i \log(f_w(x_i)) + (1 - y_i) \log(1 - f_w(x_i))$$



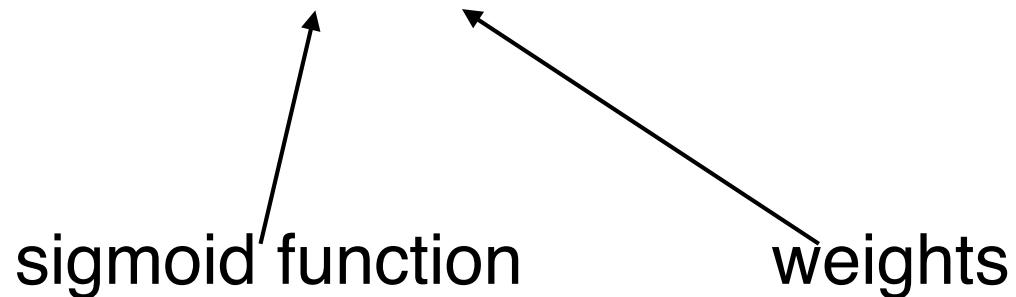
we have a set of independent realizations:

$$(x_i, y_i); i = 1, \dots, N$$

and want to find the underlying probability distribution of y given x:

Since y can only take 0 / 1 values, we assume it can be parametrized with a **Bernoulli distribution**:

$$P(y_i = 0|x_i) = 1 - \text{sigm}(f_w(x_i)) \quad P(y_i = 1|x_i) = \text{sigm}(f_w(x_i))$$



So the goal is to find the values of w (weights) that generate a Bernoulli distribution for y given a set of N independent observations

This can be achieved via Maximum Likelihood estimation:

The likelihood of a given observation under the Bernouilli assumption can be written as:

$$L(w; x_i, y_i) = [\text{sigm}(f_w(x_i))]^{y_i} [1 - \text{sigm}(f_w(x_i))]^{1-y_i}$$

which is equal to $[\text{sigm}(f_w(x_i))]$ if $y=1$ and $[1 - \text{sigm}(f_w(x_i))]$ if $y=0$

And the likelihood of the entire sample is the product of the likelihoods:

$$L(w; x_i, y_i) = \prod_{i=1}^N [\text{sigm}(f_w(x_i))]^{y_i} [1 - \text{sigm}(f_w(x_i))]^{1-y_i}$$

So, the log-likelihood:

$$l(w; x_i, y_i) = \sum_{i=1}^N y_i \times \log[\text{sigm}(f_w(x_i))] + (1 - y_i) \times \log[1 - \text{sigm}(f_w(x_i))]$$

THEREFORE EQUIVALENT TO THE CROSS-ENTROPY LOSS:

$$l(w; x_i, y_i) = \sum_{i=1}^N y_i \times \log[\text{sigm}(f_w(x_i))] + (1 - y_i) \times \log[1 - \text{sigm}(f_w(x_i))]$$

$$H = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(f_w(x_i)) + (1 - y_i) \cdot \log(1 - f_w(x_i))$$

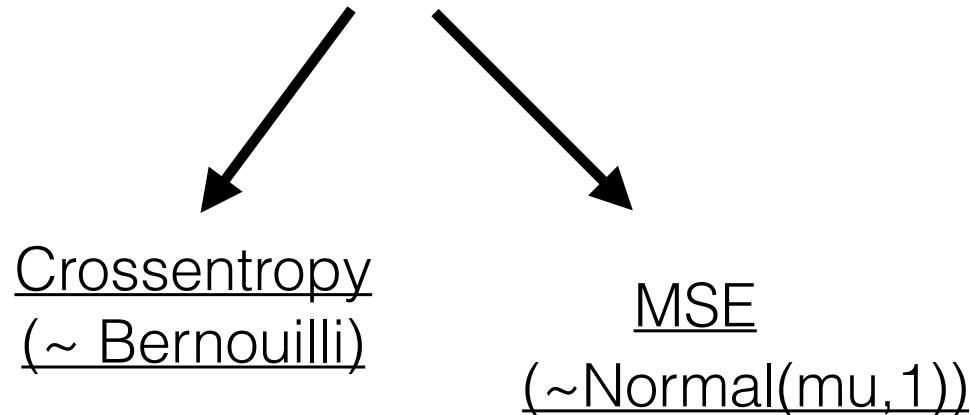
WHEN YOU DO A BINARY CLASSIFICATION WITH NEURAL NETWORKS YOU ARE SIMPLY FINDING THE WEIGHTS OF YOUR MODEL THAT MAXIMIZE THE LIKELIHOOD OF A BERNOULLI DISTRIBUTION

ASSUMES THAT THE POSTERIOR $p(y|x)$ CAN BE MODELLED BY A BERNOULLI DISTRIBUTION

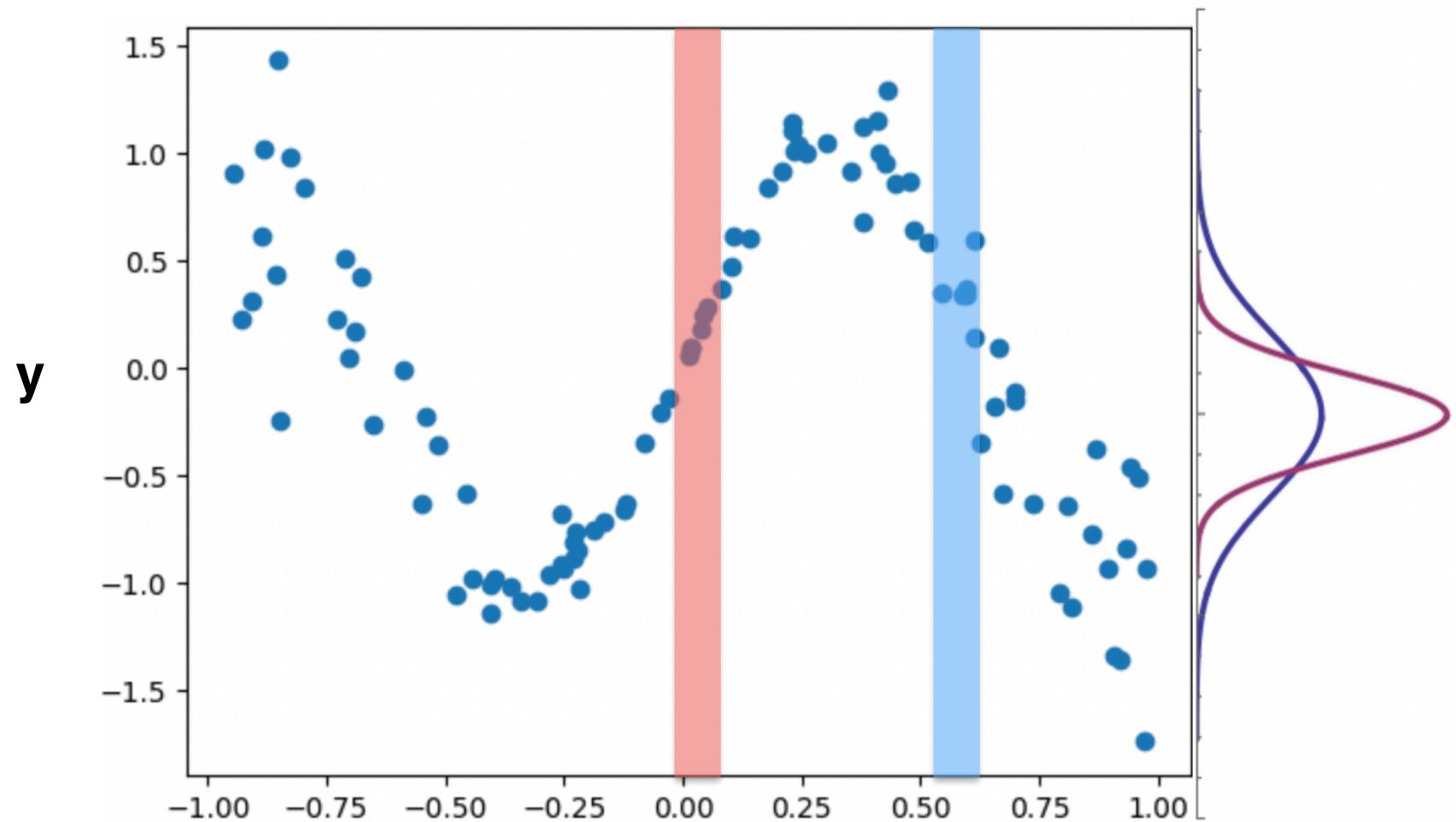
A standard NN performs Maximum Likelihood Estimation of the weights W

$$D = (x^{(i)}, y^{(i)})$$

$$p(D|w) = \prod_i p(y^{(i)}|x^{(i)}, W) \sim p_w(y|x)$$

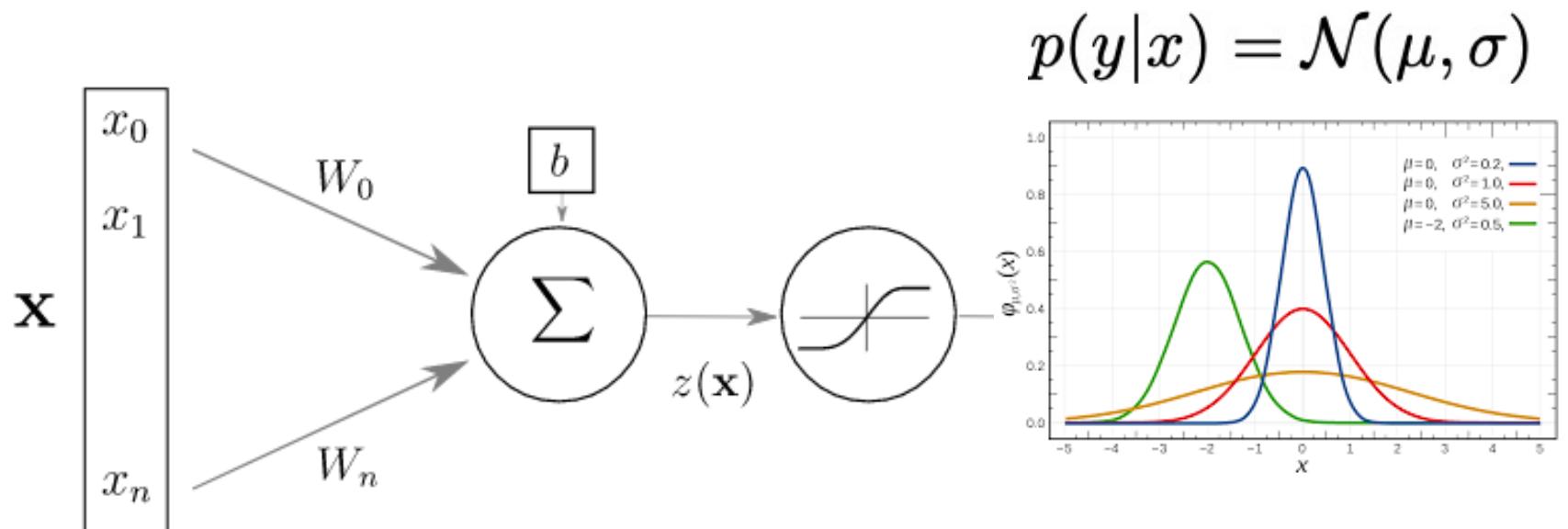
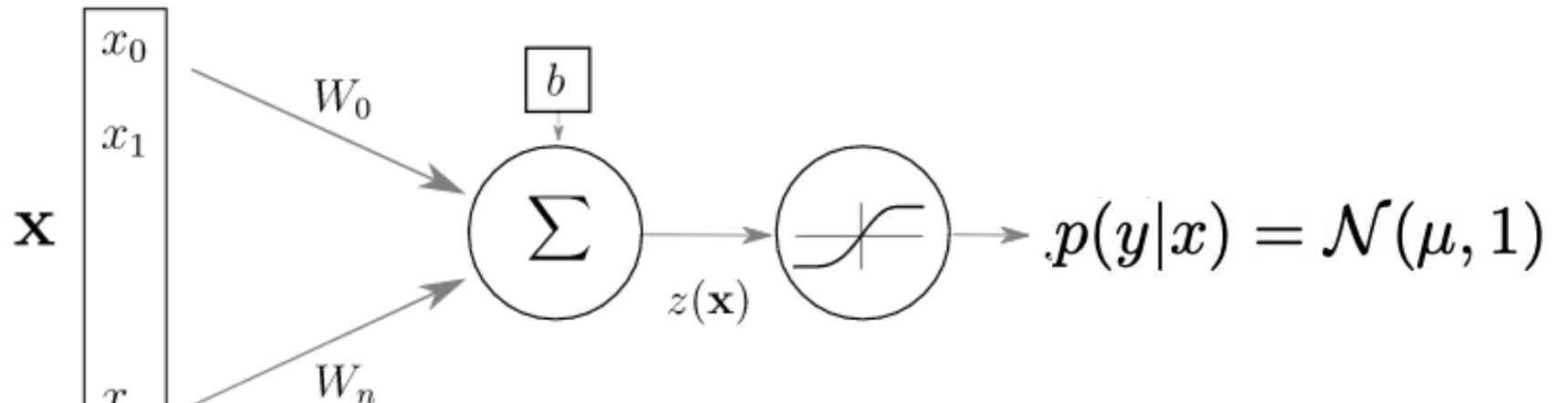


$$p(y|x) = \mathcal{N}(\mu, 1) ?$$

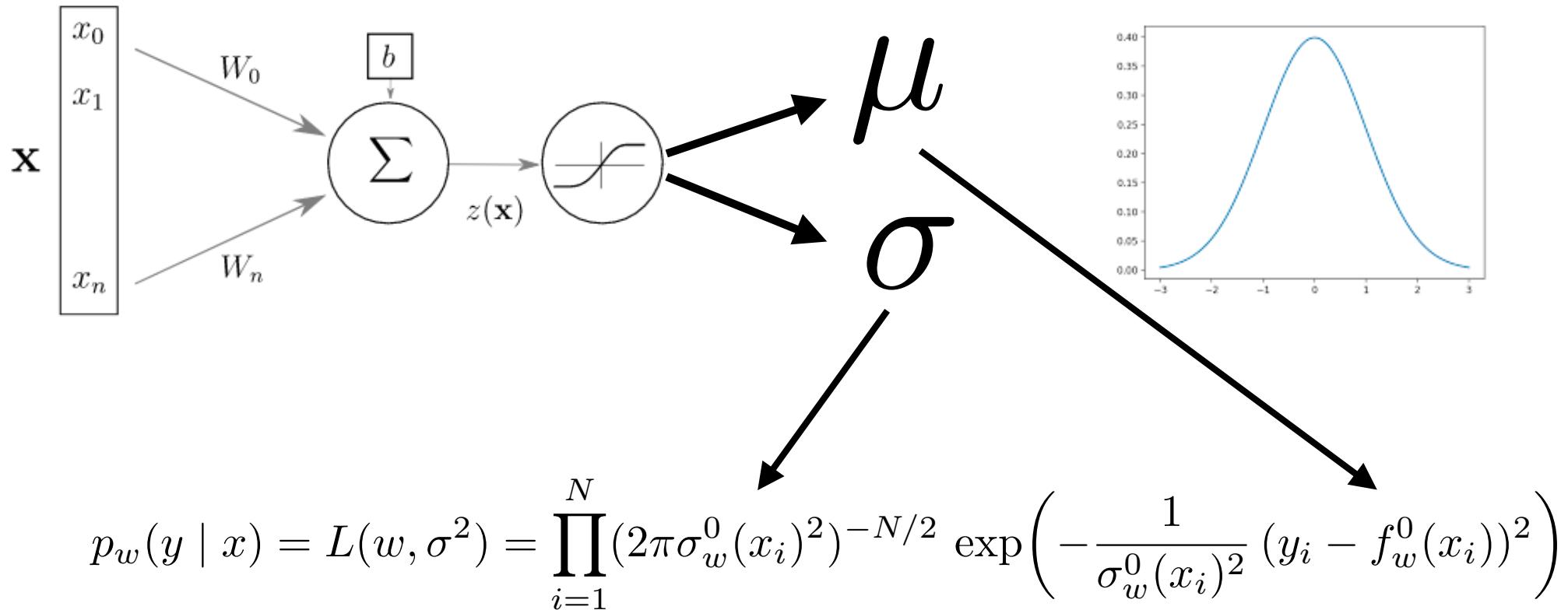


is one parameter enough to characterize the
posterior?

Mixture Density Network

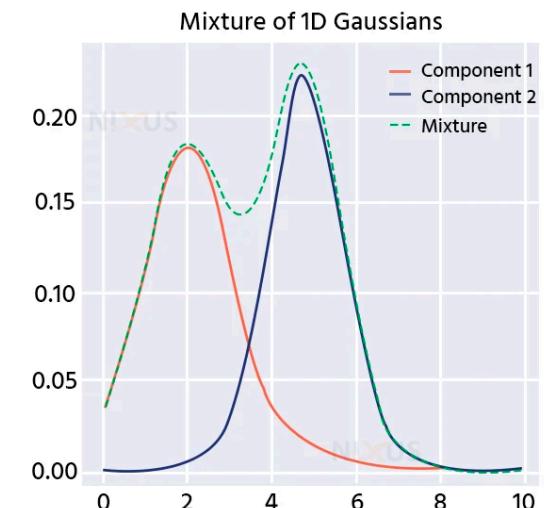
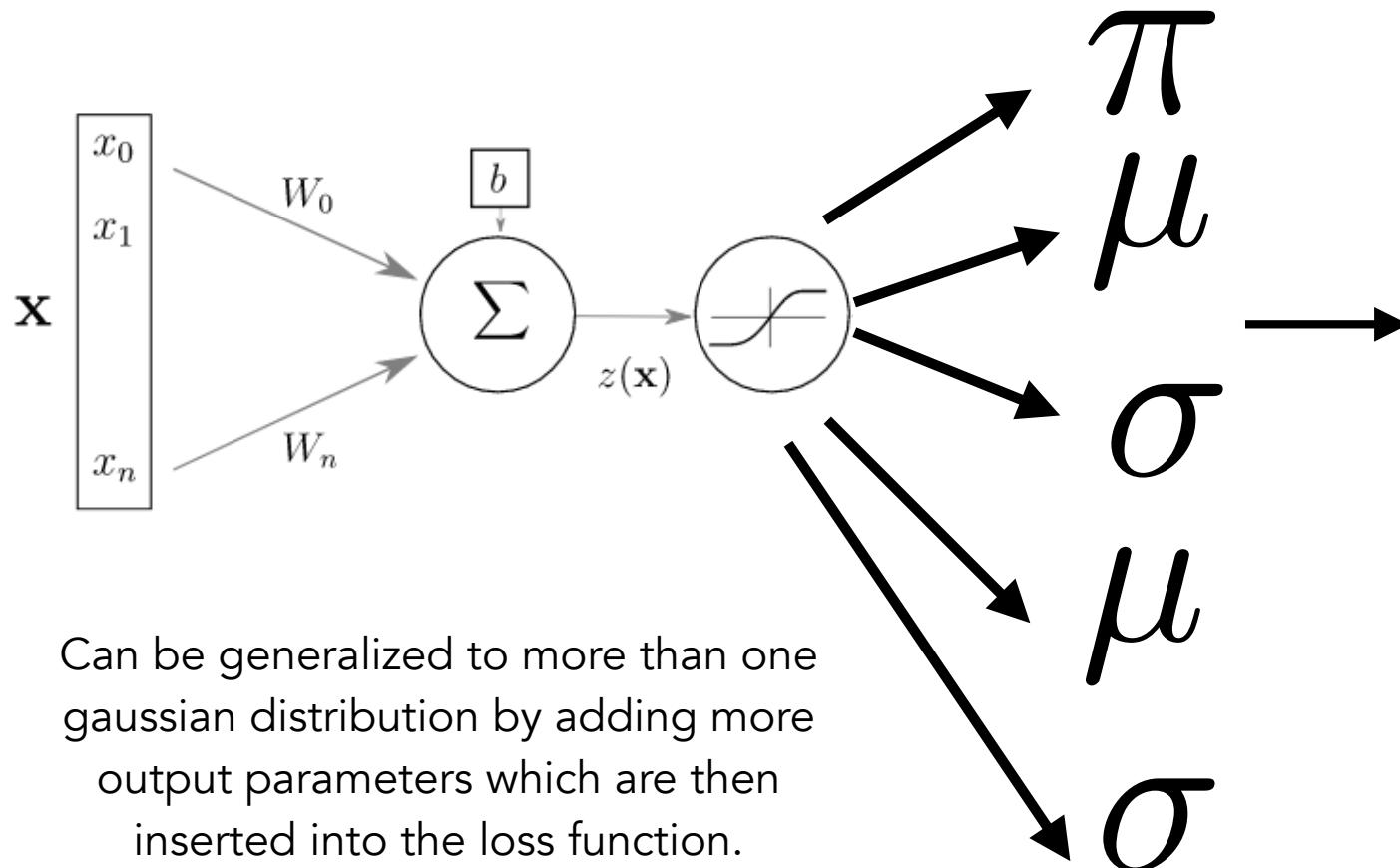


Mixture Density Network (MDN)



The neural network outputs the parameters of a distribution, on which one can back propagate (2 in the example above)

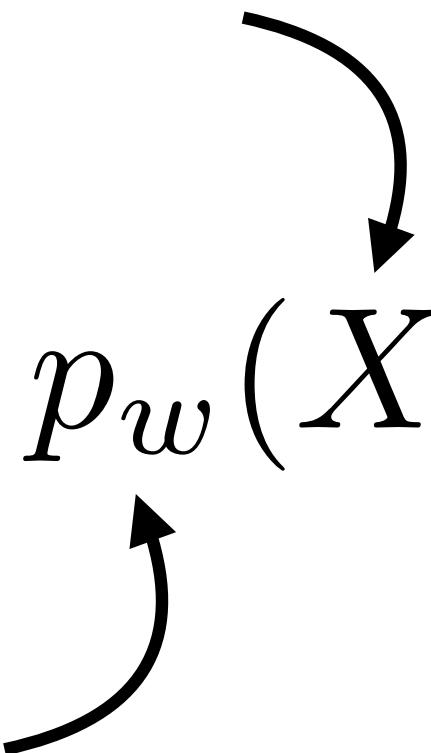
What if the posterior is not gaussian?



$$p_w(y | x) = L(w) = \prod_{i=1}^N \sum_{k=1}^2 \pi_w^{(k)}(x_i) \frac{1}{\sqrt{2\pi} \sigma_w^{(k)}(x_i)} \exp\left(-\frac{(y_i - f_w^{(k)}(x_i))^2}{2 [\sigma_w^{(k)}(x_i)]^2}\right)$$

The loss function remains the max of the (log) likelihood

High
dimensional data
(e.g. images)

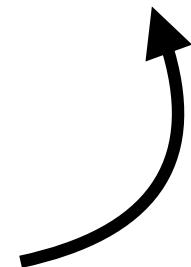
$$p_w(X)$$


Likely not
Gaussian

High
dimensional data
(e.g. images)

Evaluation

$$p_w(X)$$



Likely not
Gaussian

$$z \sim \mathcal{N}(0, 1)$$

(Easy to sample, easy to evaluate)

Sampling

Evaluation

High
dimensional data
(e.g. images)

$$p_w(X)$$

Likely not
Gaussian

NNs == universal
approximators

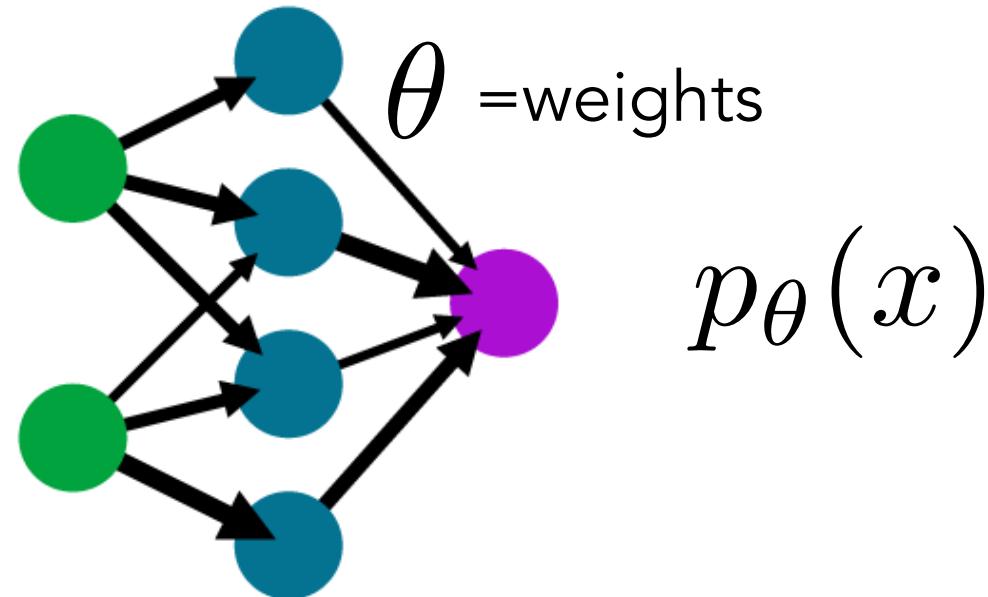
$$z \sim \mathcal{N}(0, 1)$$

(Easy to sample, easy to evaluate)

(Optimization problem)

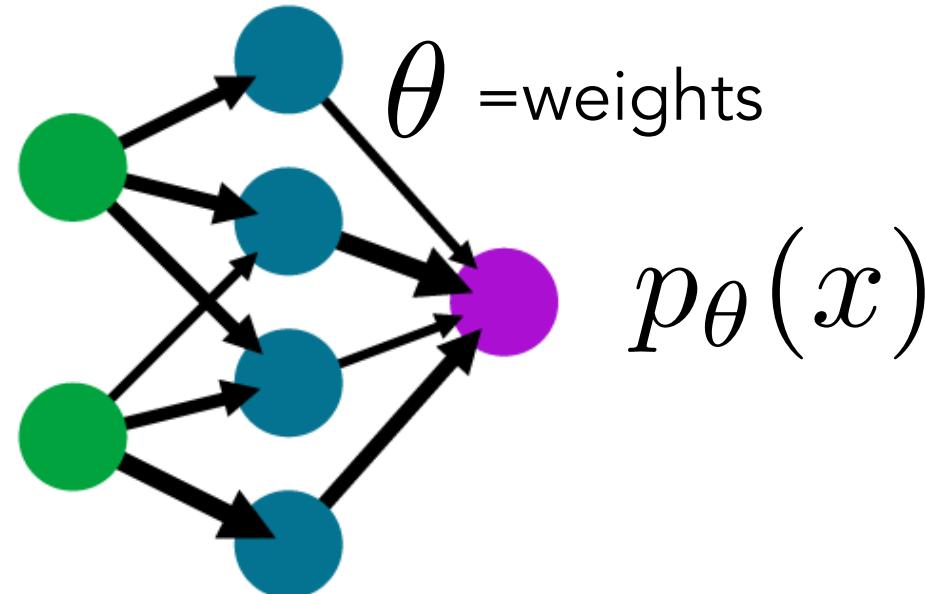
Sampling

$$z \sim \mathcal{N}(0, 1)$$
$$p(z)$$



$$z \sim \mathcal{N}(0, 1)$$

$$p(z)$$



Types of generative models

Likelihood based models	directly learn the distribution's probability density function via (approximate) maximum likelihood	VAEs, Normalizing Flows	$\max_{\theta} \sum_{i=1}^N \log p_{\theta}(x_i)$
Implicit generative models	the probability distribution is implicitly represented by a model of its sampling process	GANs	adversarial
Score based models	model the gradient of the log probability density function	Score based diffusion	$\nabla_x \log p_{\theta}(x)$

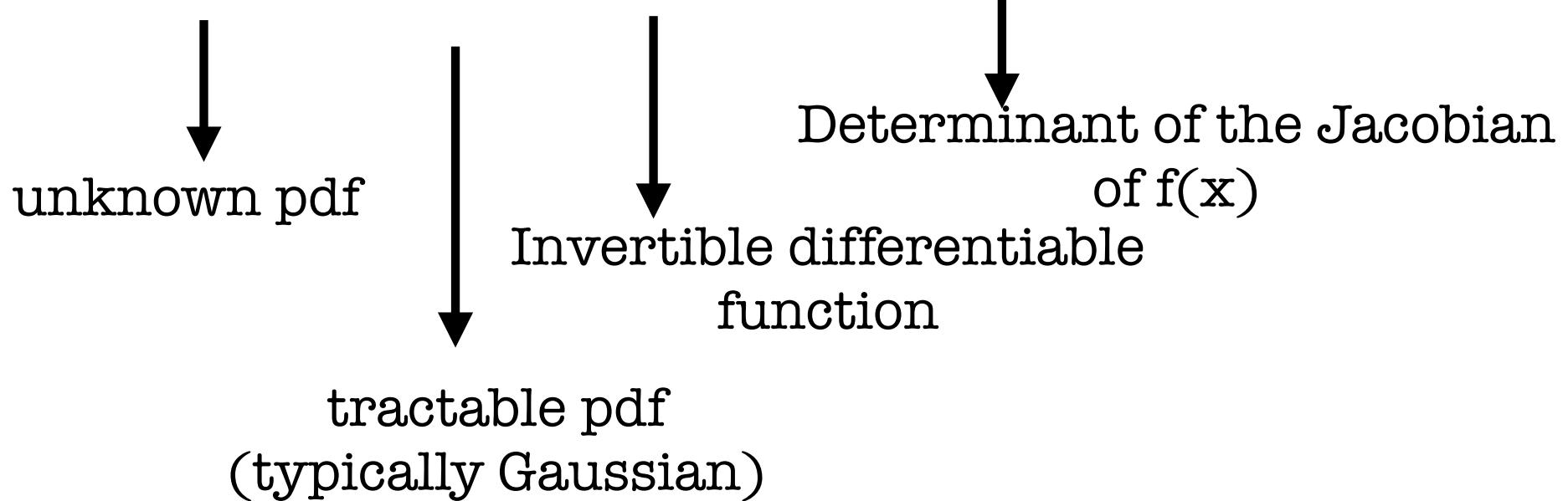
Types of generative models

Likelihood based models	directly learn the distribution's probability density function via (approximate) maximum likelihood	VAEs, Normalizing Flows
Implicit generative models	the probability distribution is implicitly represented by a model of its sampling process	GANs
Score based models	model the gradient of the log probability density function	Score based diffusion

Normalizing Flows

Based on change of variables:

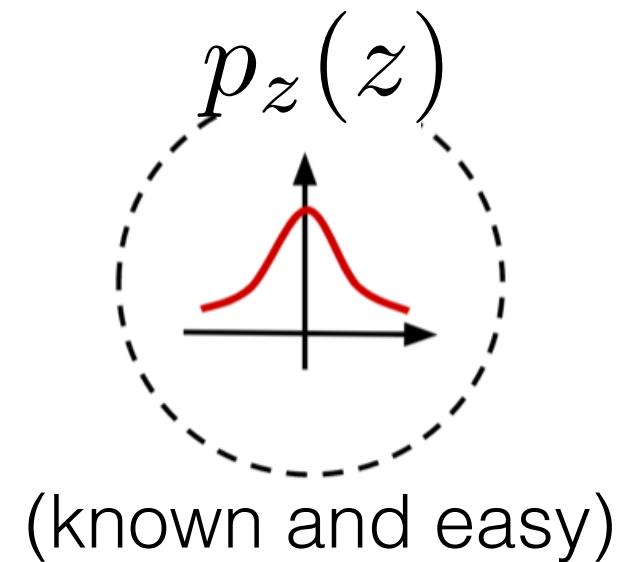
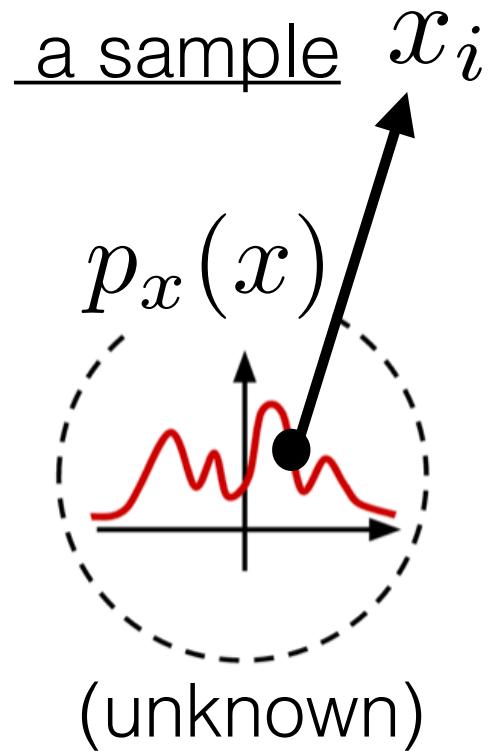
$$p_X(x) = p_Z(f(x)) |\det J_f(x)|, \quad J_f(x) = \frac{\partial f(x)}{\partial x^\top}$$



the neural network parametrizes the function f: $f_\theta(x)$

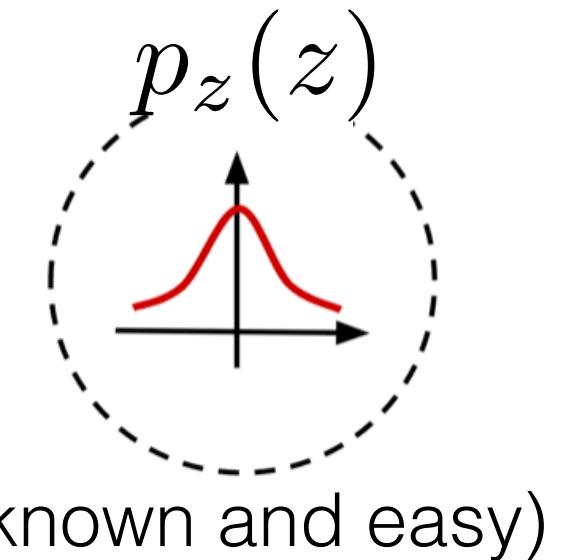
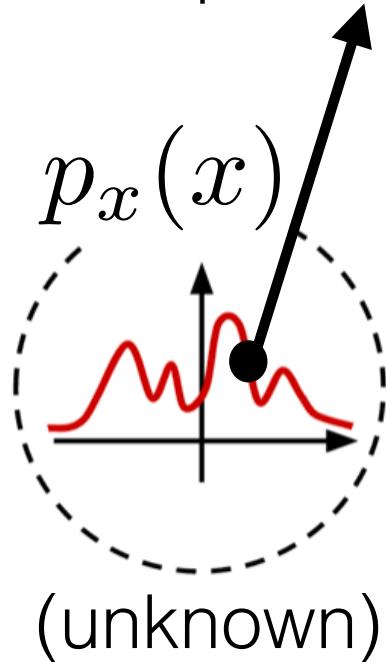
and optimize so that: $p_Z(f_\theta(x)) \cdot |\det J_f(x)|$ is maximum

training procedure

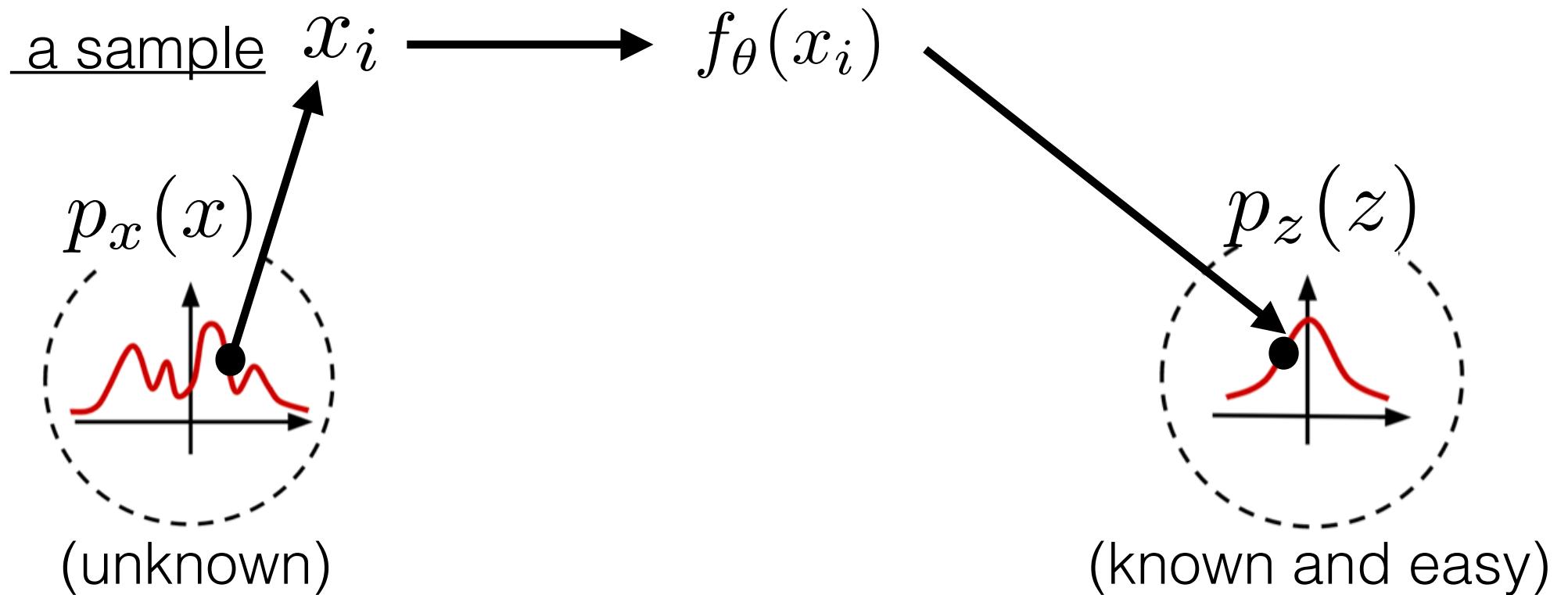


training procedure

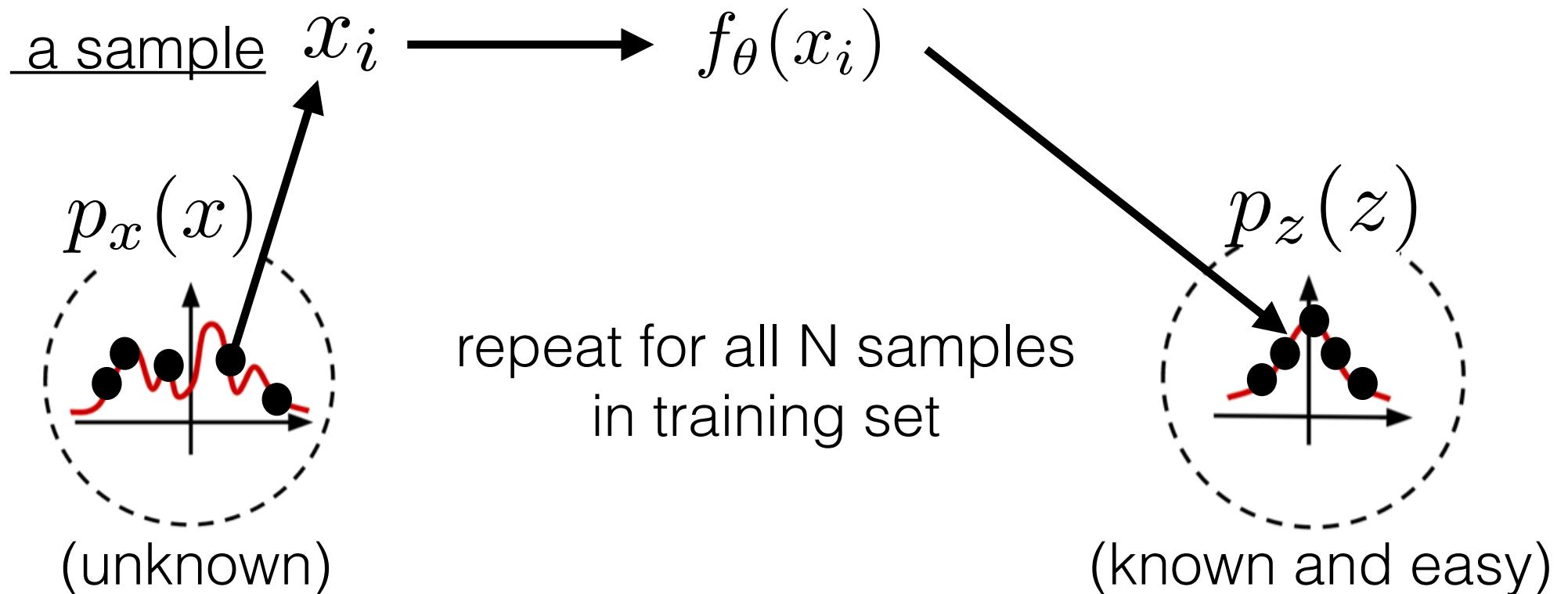
a sample $x_i \longrightarrow f_\theta(x_i)$



training procedure



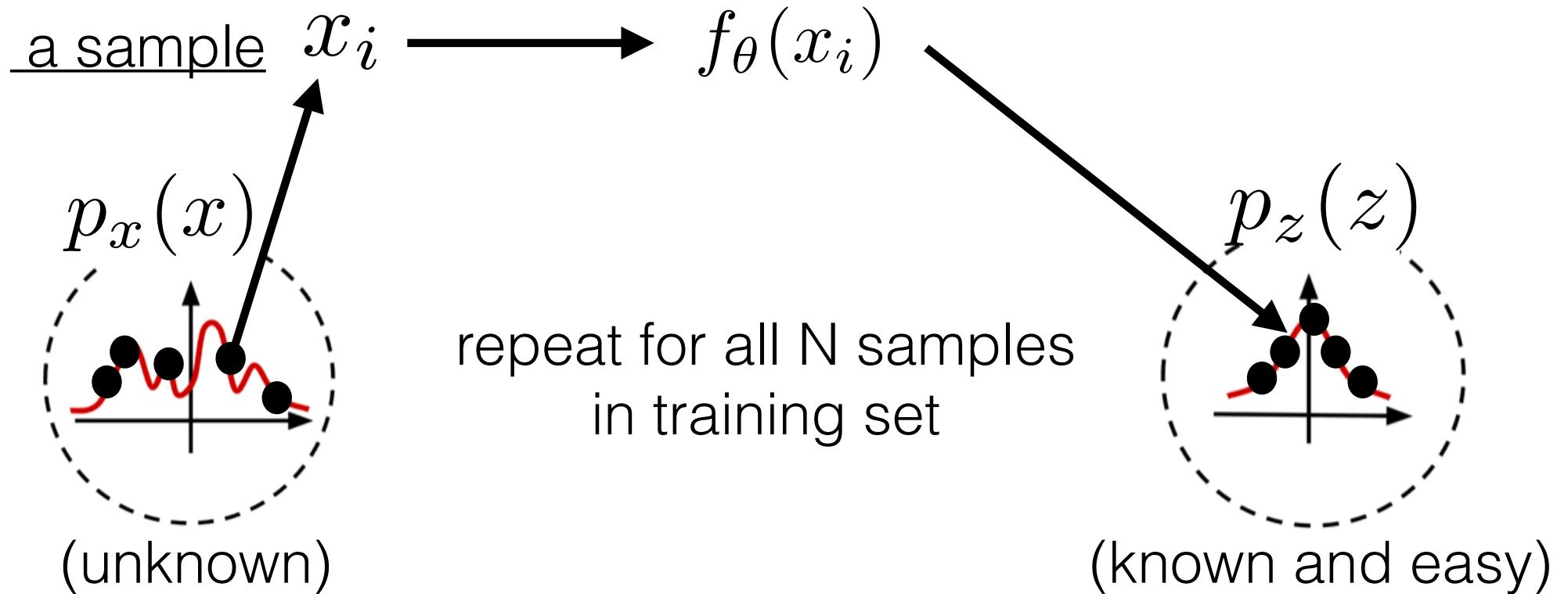
training procedure



And Maximize this

$$\sum_{i=1}^N \log p_z(x_i) + \log \det J_f(x_i)$$

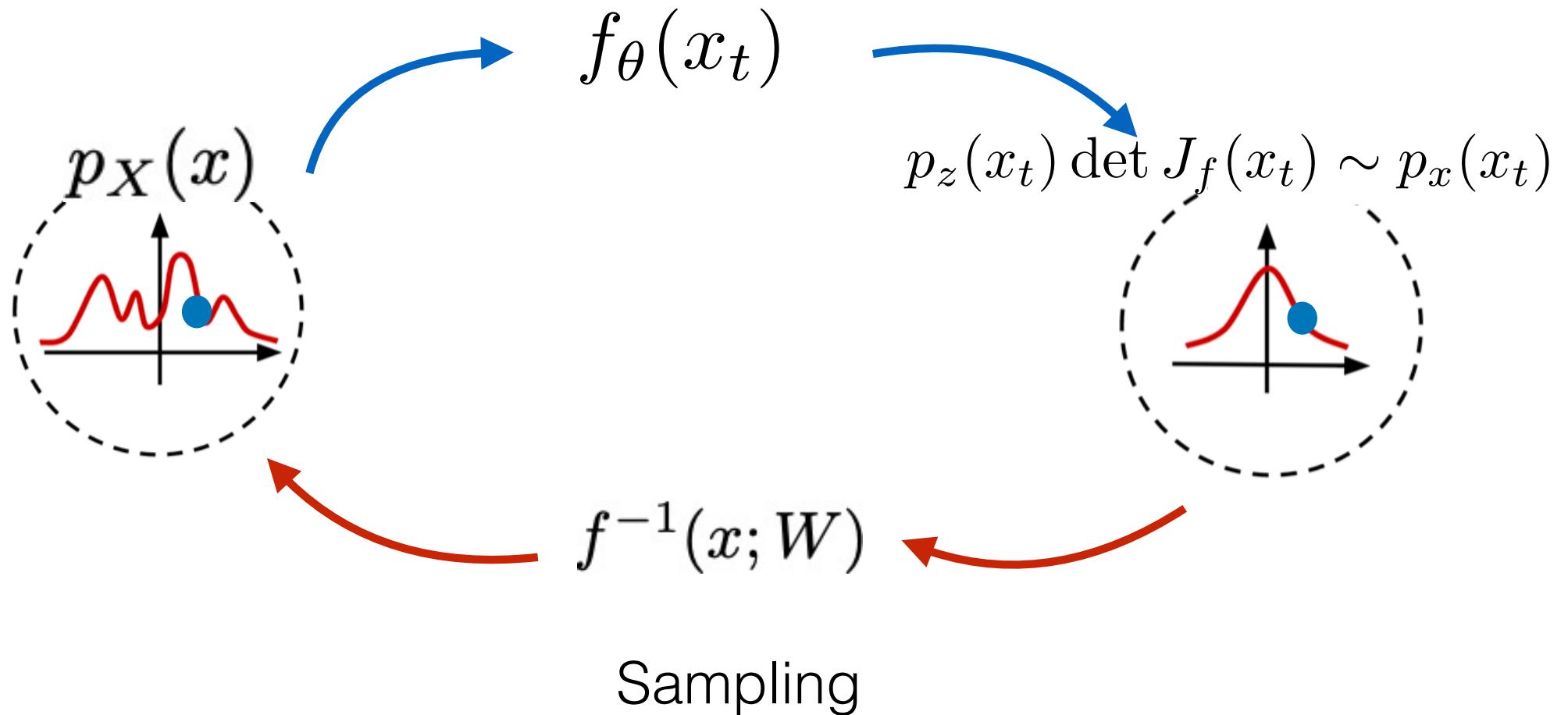
training procedure



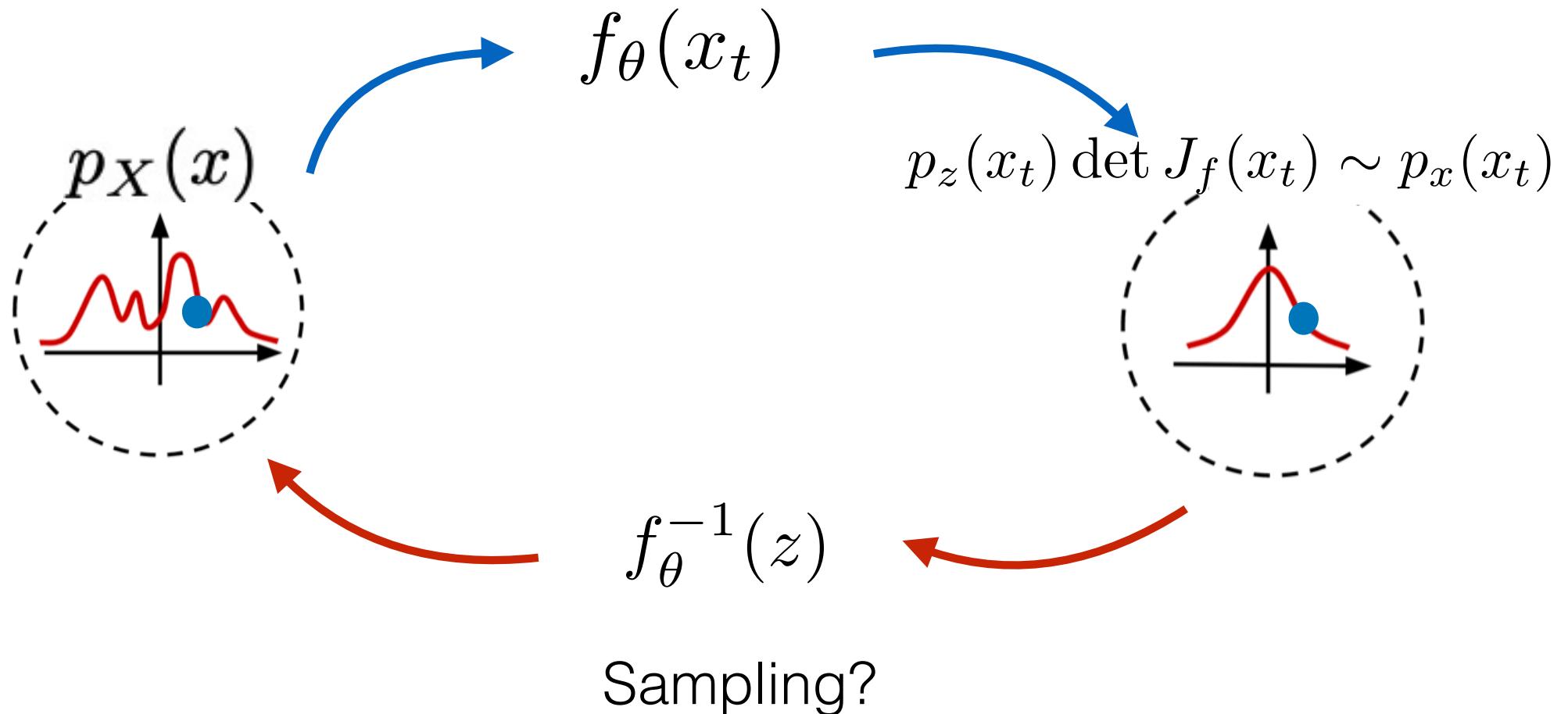
achieved by maximizing this loss function

$$\sum_{i=1}^N \log p_z(x_i) + \log \det J_f(x_i)$$

Likelihood evaluation



Likelihood evaluation



The price to pay is that we need to put some restrictions on the NN architecture

- 1. Easily invertible - to enable sampling
- 2. Jacobian tractable - to optimize the network

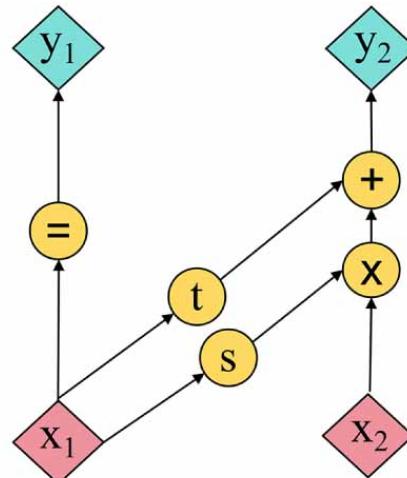
What functions satisfy these conditions?

An example are: RealNVPs (Real-valued Non-Volume Preserving)

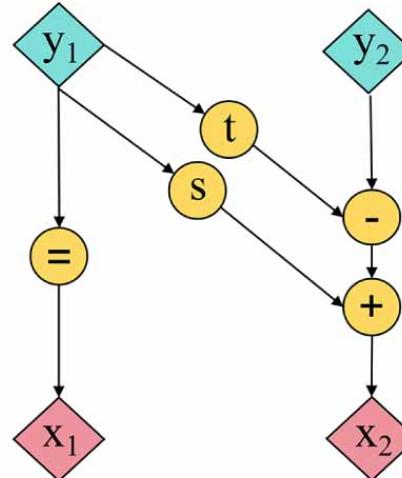
$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d})$$

NN



(a) Forward transformation



(b) Inverse transformation

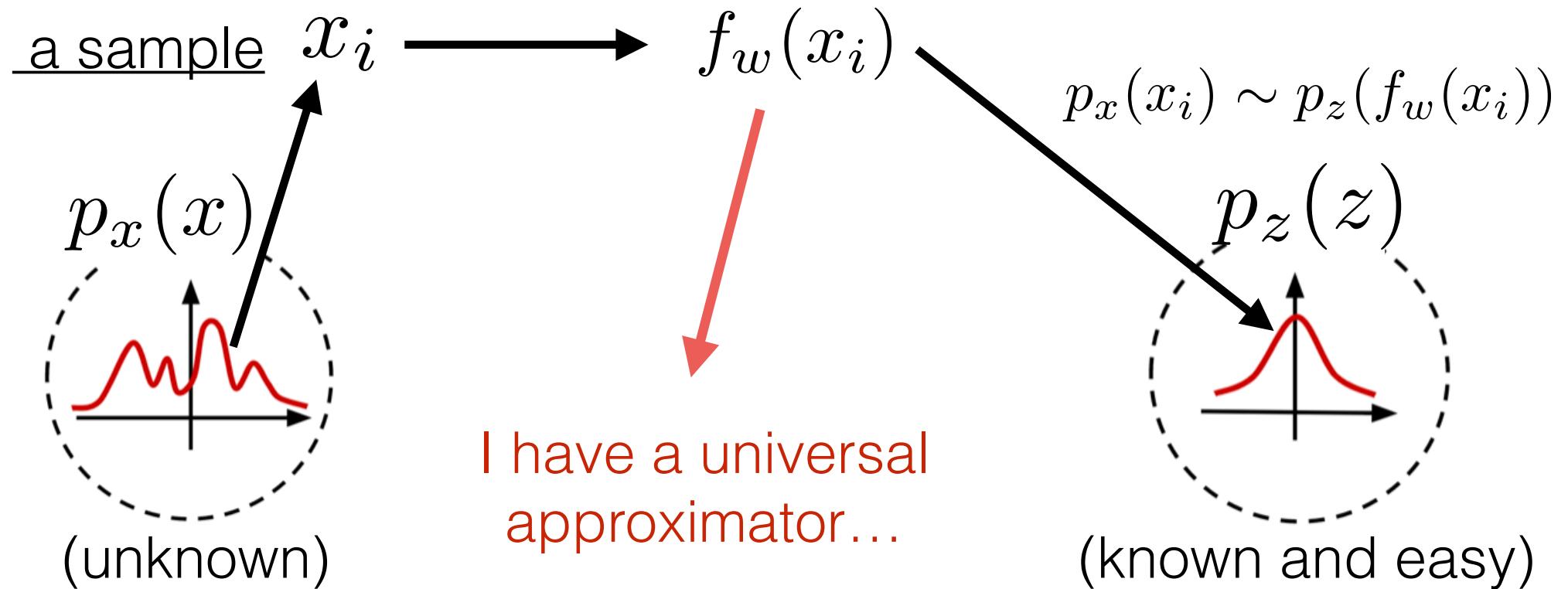
- Easily invertible:

$$\begin{cases} \mathbf{y}_{1:d} &= \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} &= \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d}) \end{cases} \Leftrightarrow \begin{cases} \mathbf{x}_{1:d} &= \mathbf{y}_{1:d} \\ \mathbf{x}_{d+1:D} &= (\mathbf{y}_{d+1:D} - t(\mathbf{y}_{1:d})) \odot \exp(-s(\mathbf{y}_{1:d})) \end{cases}$$

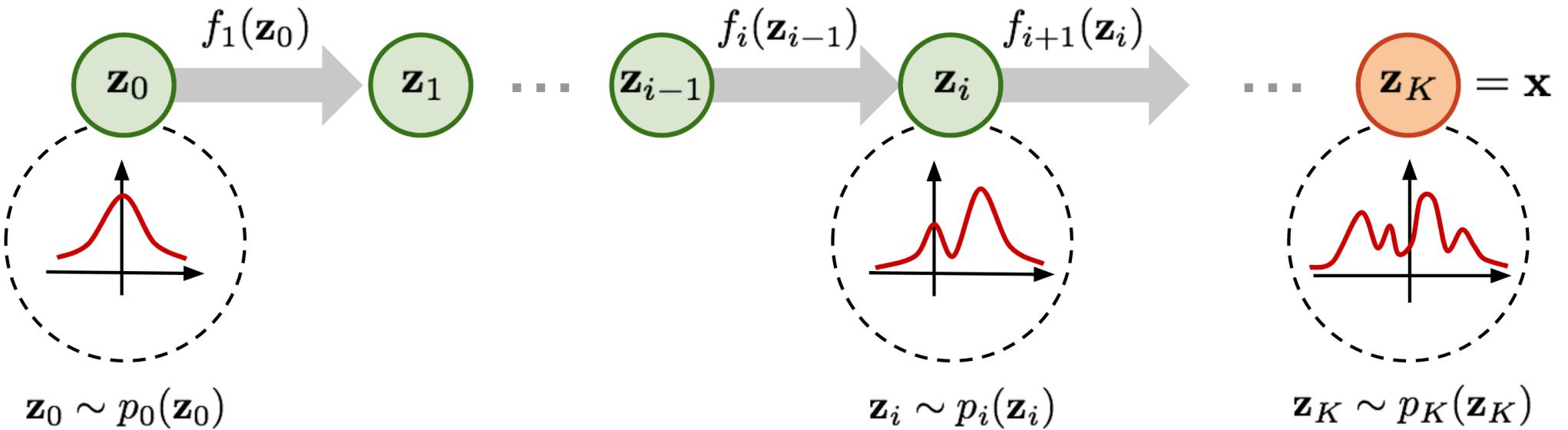
- Jacobian:

$$\mathbf{J} = \begin{bmatrix} \mathbb{I}_d & \mathbf{0}_{d \times (D-d)} \\ \frac{\partial \mathbf{y}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(\exp(s(\mathbf{x}_{1:d}))) \end{bmatrix}$$

$$\det(\mathbf{J}) = \prod_{j=1}^{D-d} \exp(s(\mathbf{x}_{1:d}))_j = \exp\left(\sum_{j=1}^{D-d} s(\mathbf{x}_{1:d})_j\right)$$



Normalizing flow: in practice we use a concatenation of neural networks
(called bijectors)



$z_i = f_i(z_{i-1})$ are **invertible** and **differentiable** transformations

$f = f_1 \circ f_2 \dots \circ f_{k-1} \circ f_k$ is also invertible and differentiable