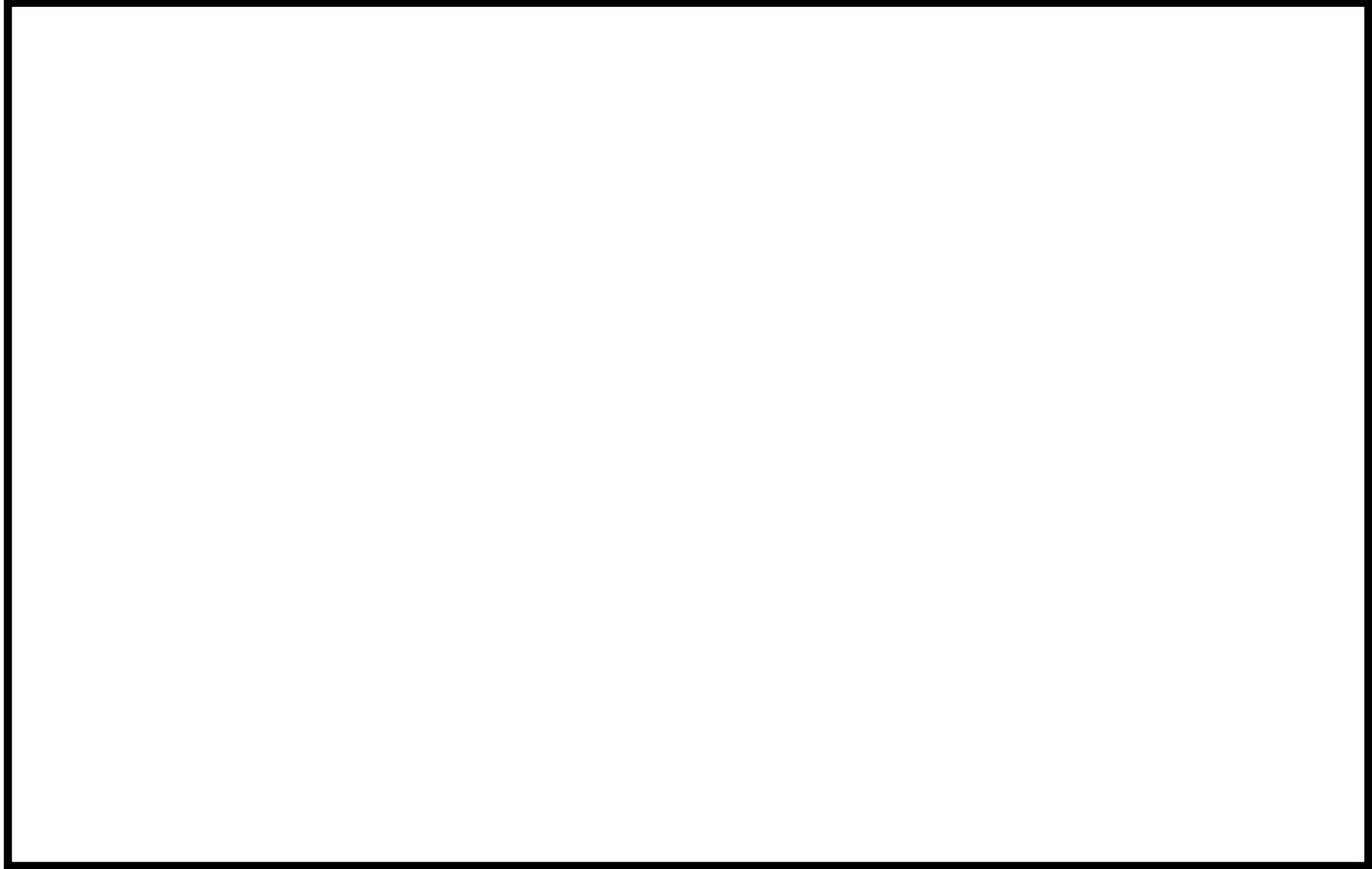


DEEP LEARNING

PART THREE - *DEEP GENERATIVE MODELS*

GENERATIVE MODELS



DATA

The diagram shows a large square divided into a 4x4 grid of smaller squares. The top-left 4x4 sub-grid is further divided into a 4x4 grid of even smaller squares. Three black dots are positioned in the top-left cell of this 4x4 sub-grid, arranged diagonally from the top-left corner to the bottom-right corner.

DATA

example 1

DATA

[illegible]

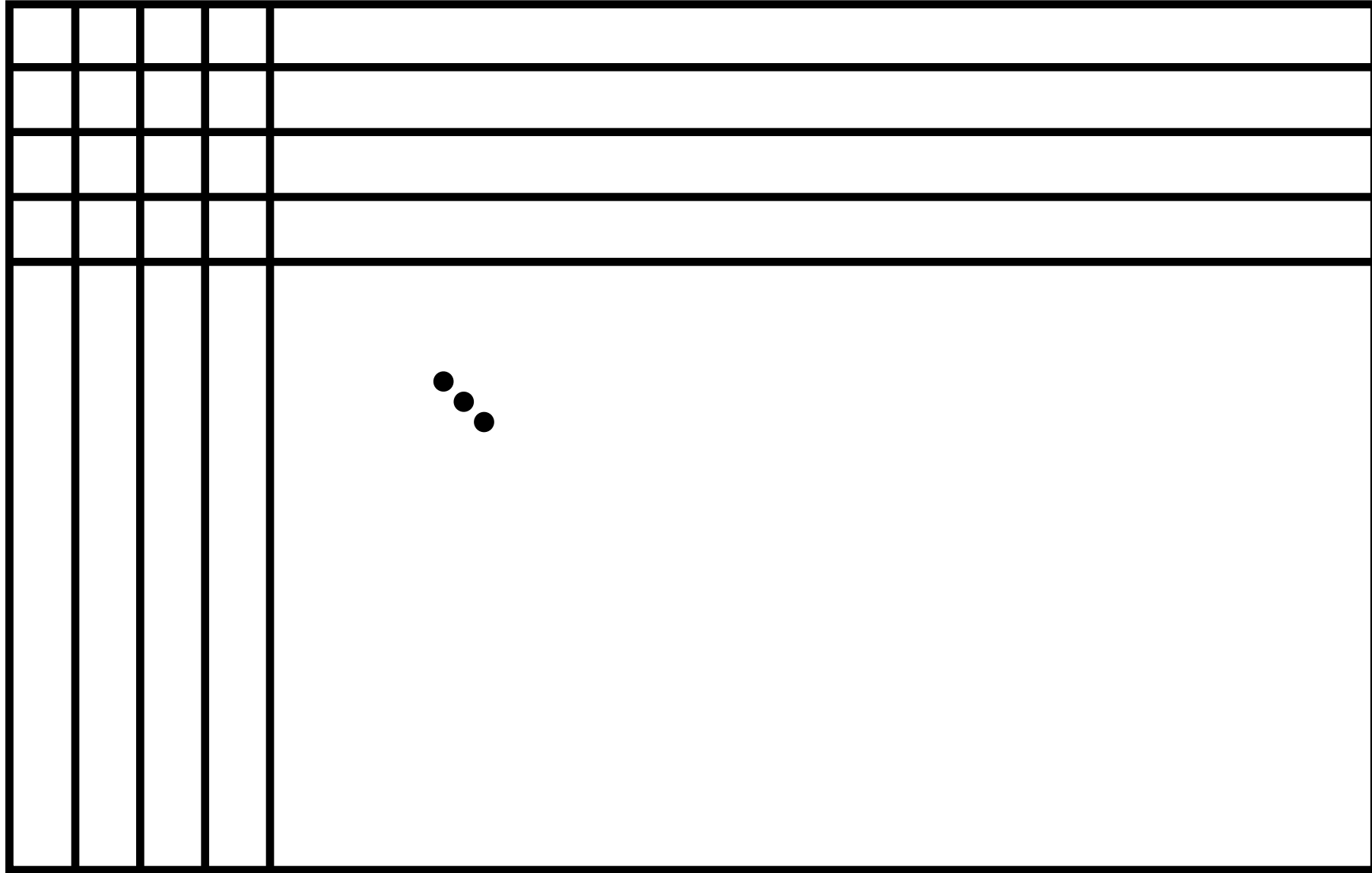
DATA

example 3

...

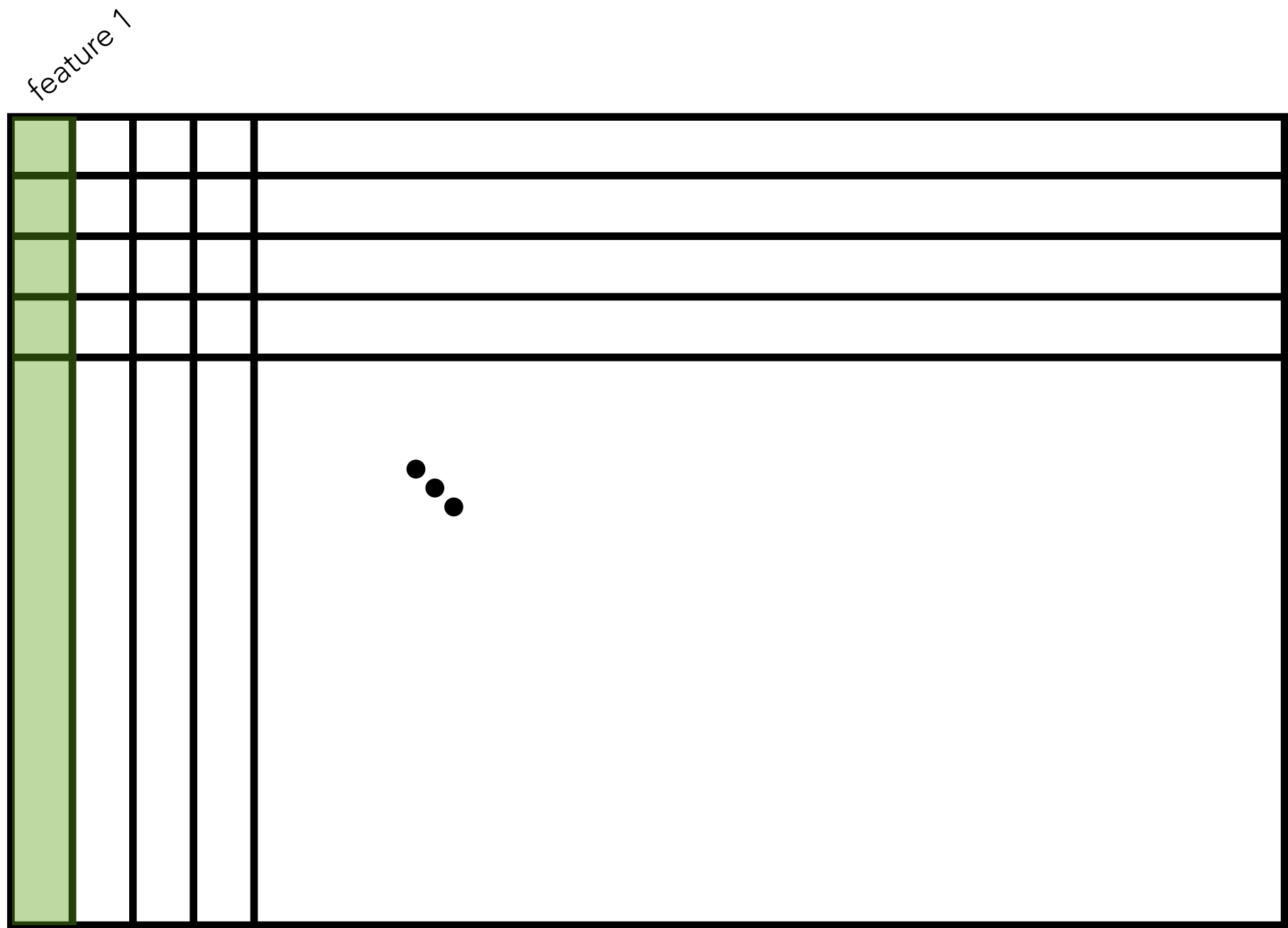
DATA

number of data examples



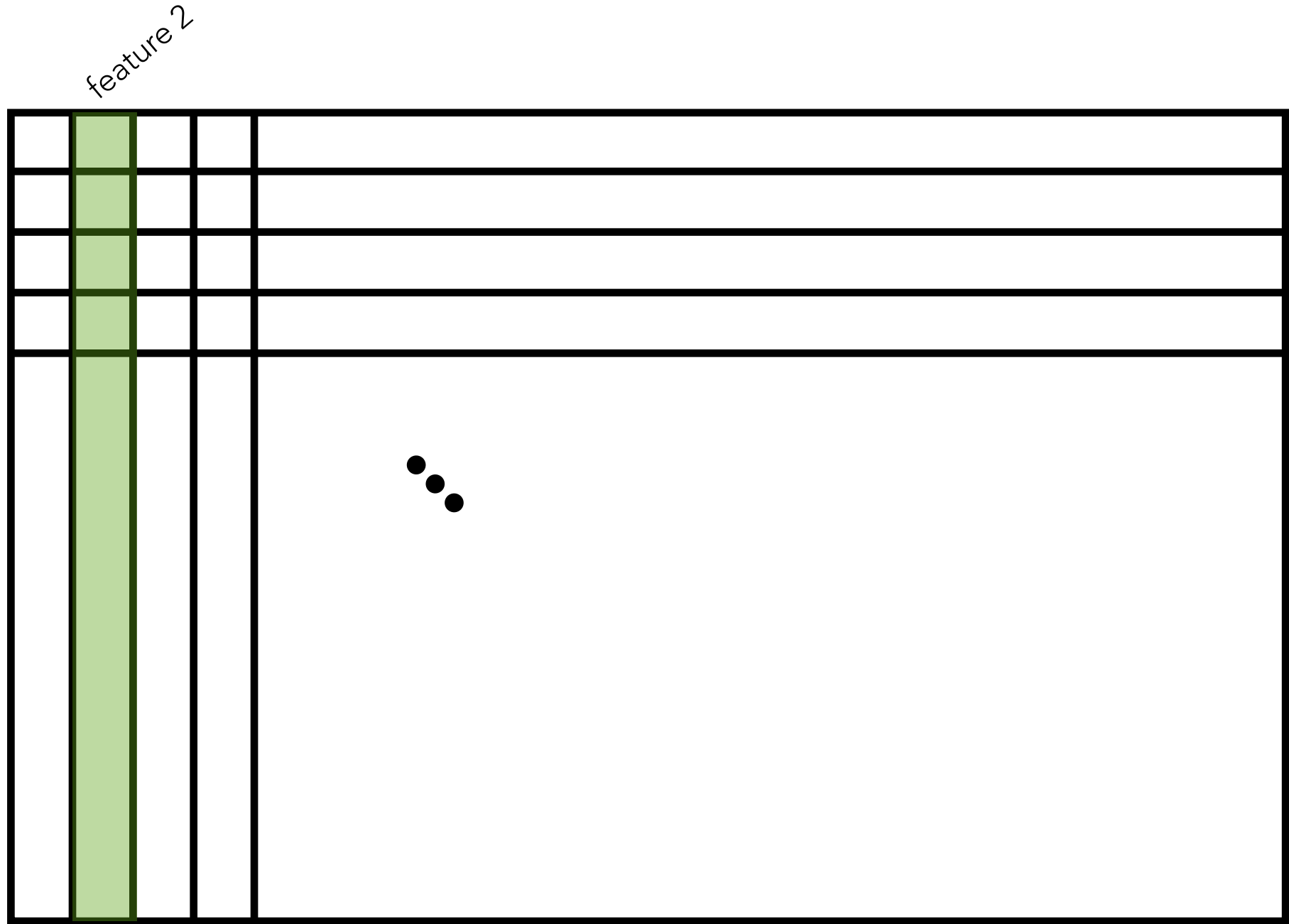
DATA

number of data examples



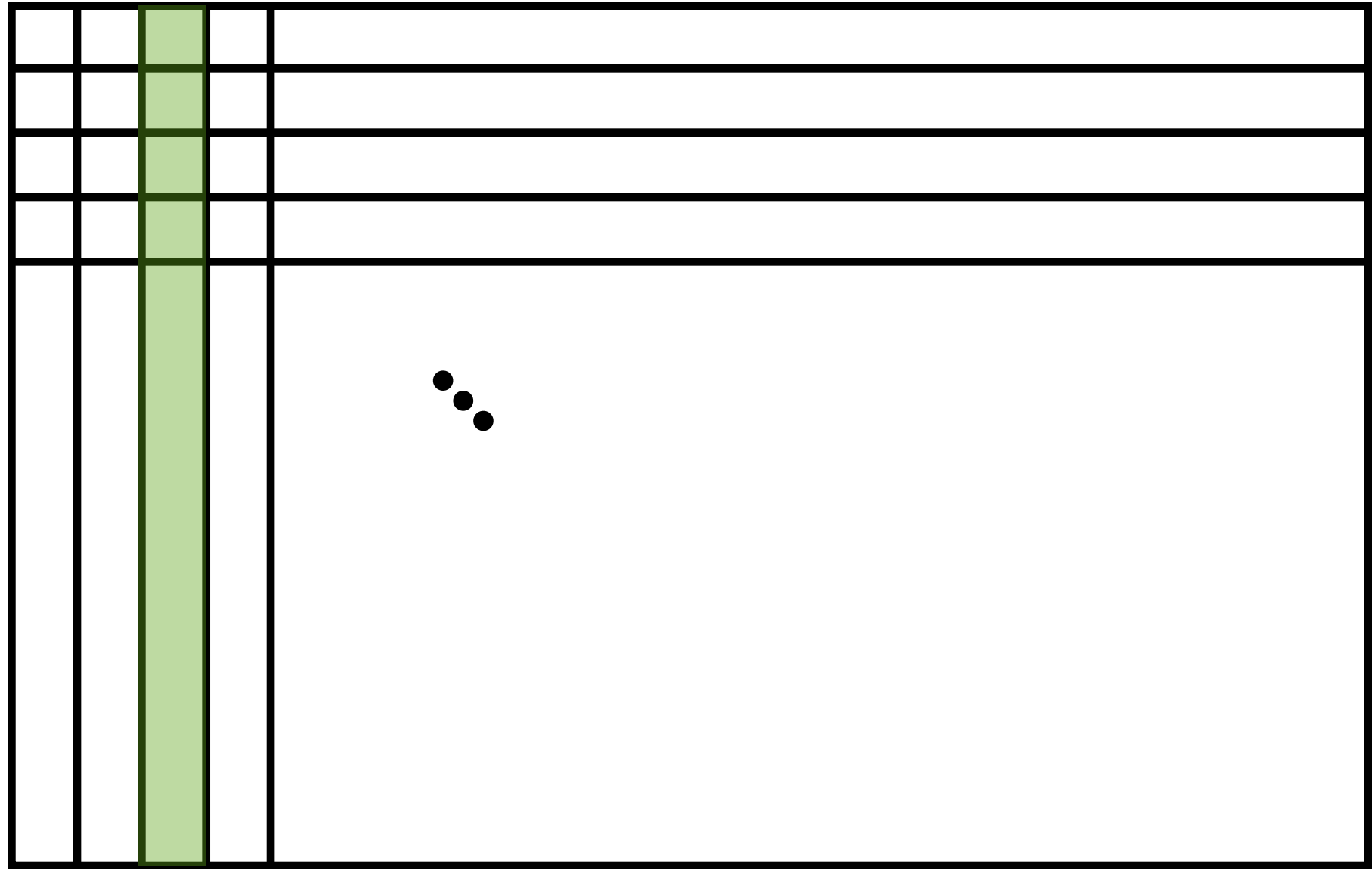
DATA

number of data examples



DATA

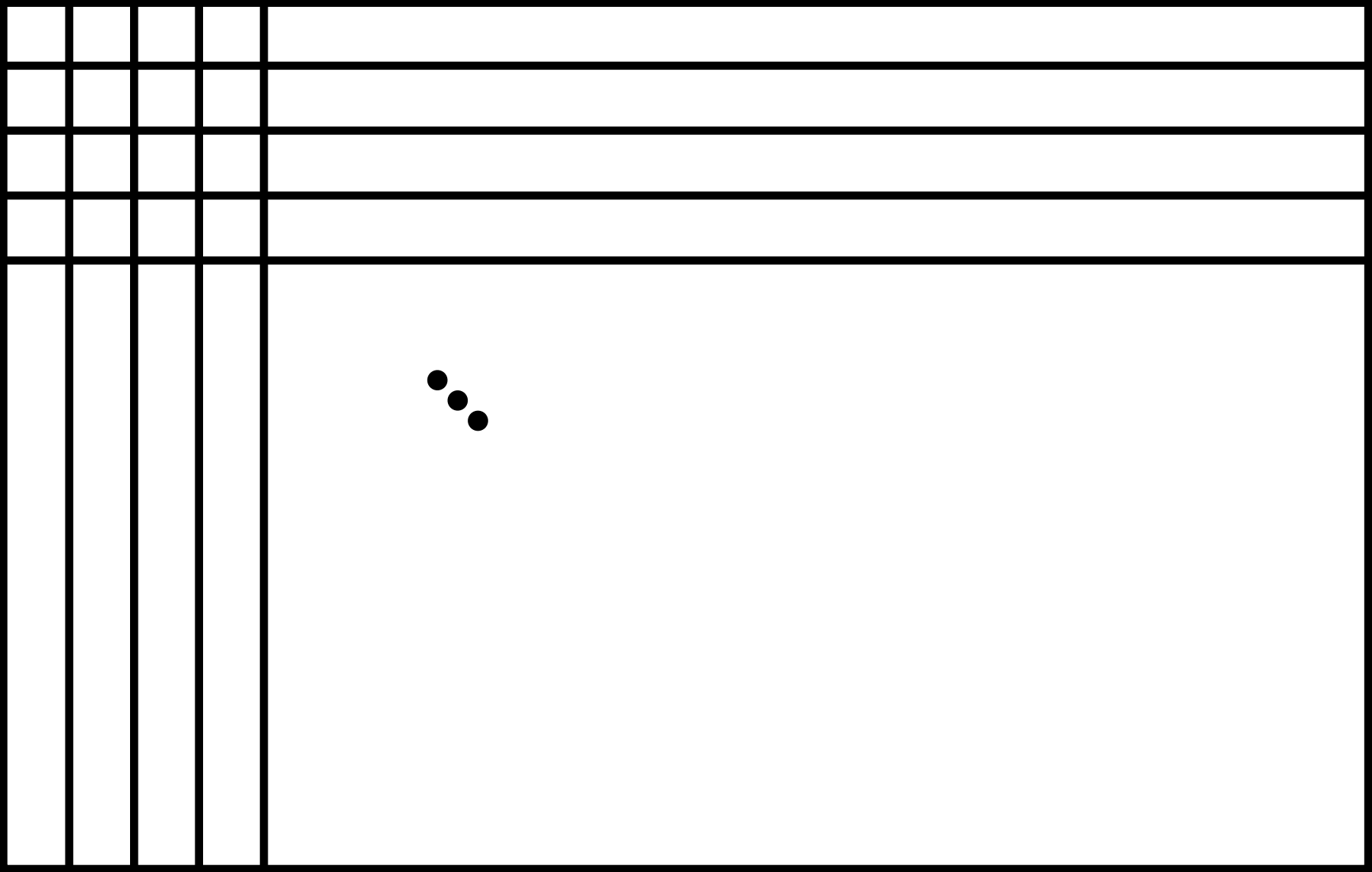
number of data examples



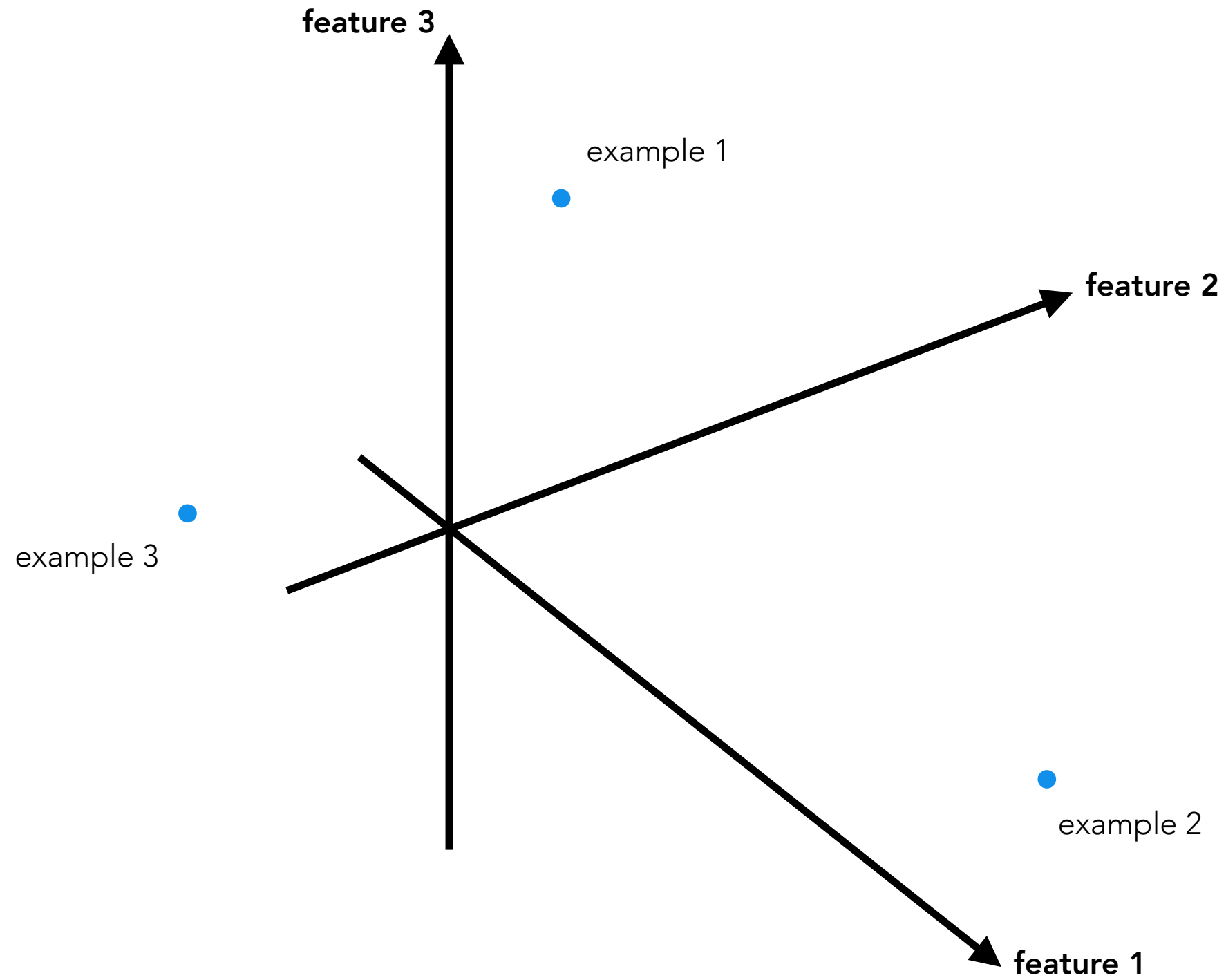
DATA

number of features

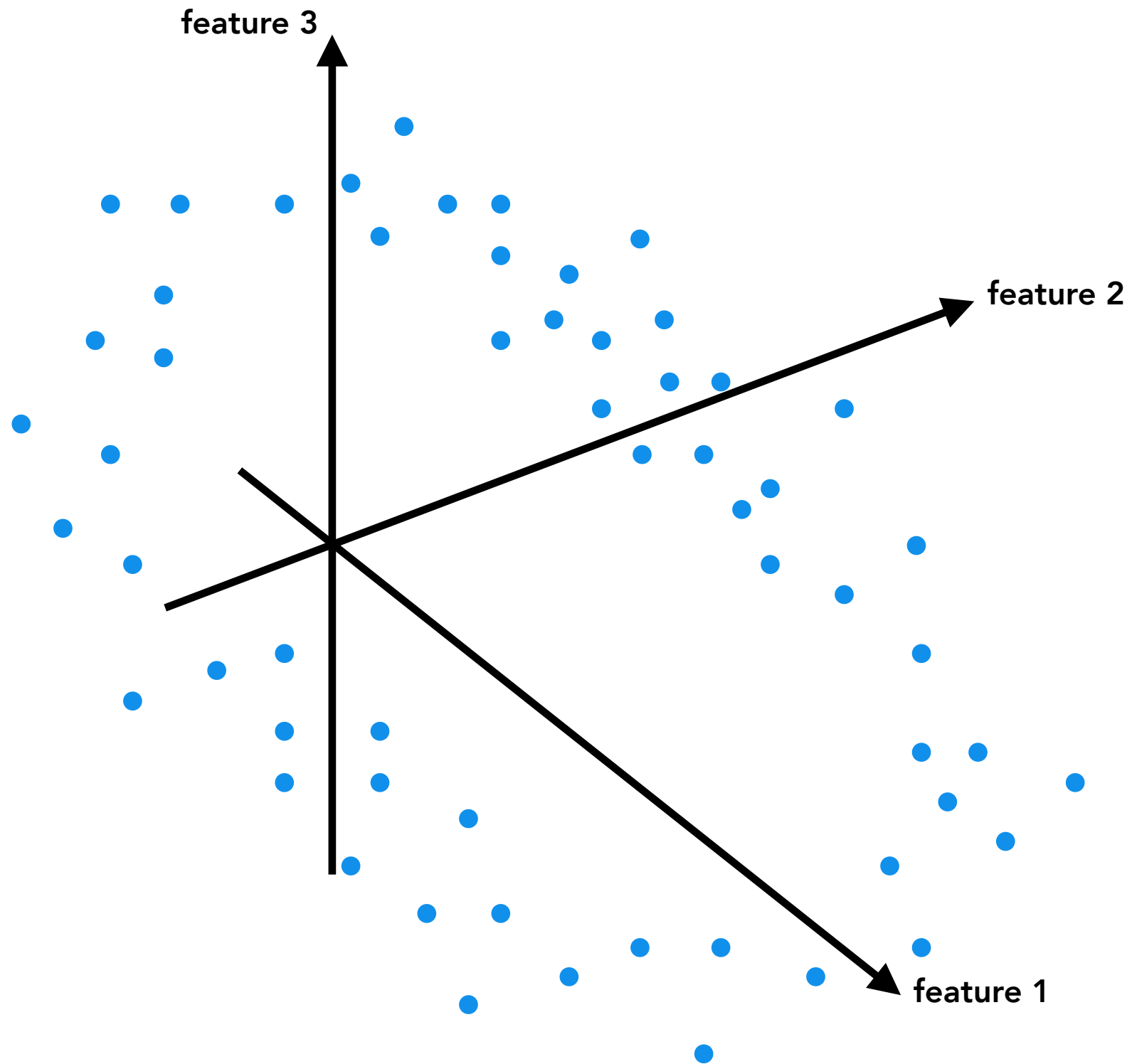
number of data examples



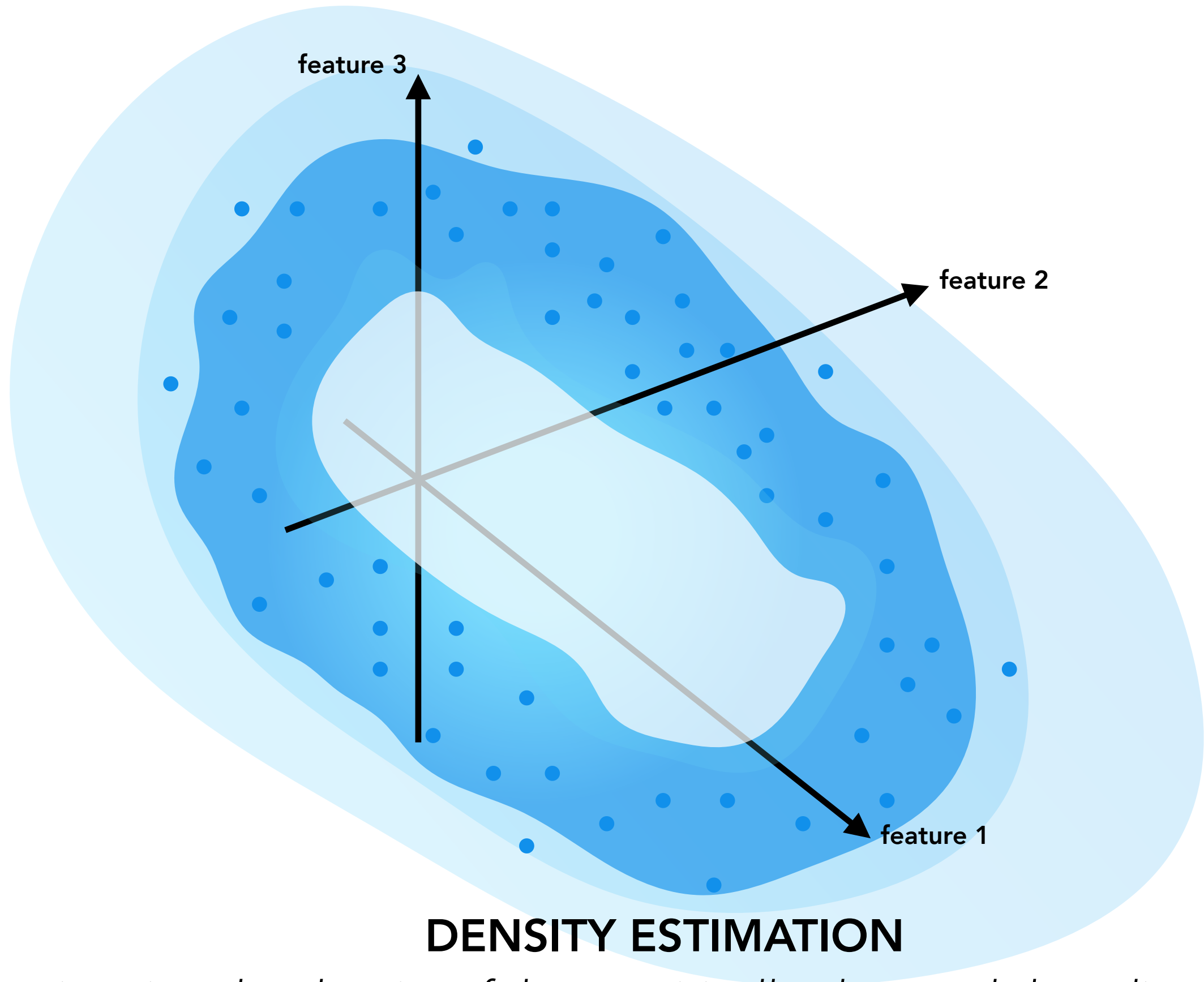
DATA



DATA DISTRIBUTION



DATA DISTRIBUTION



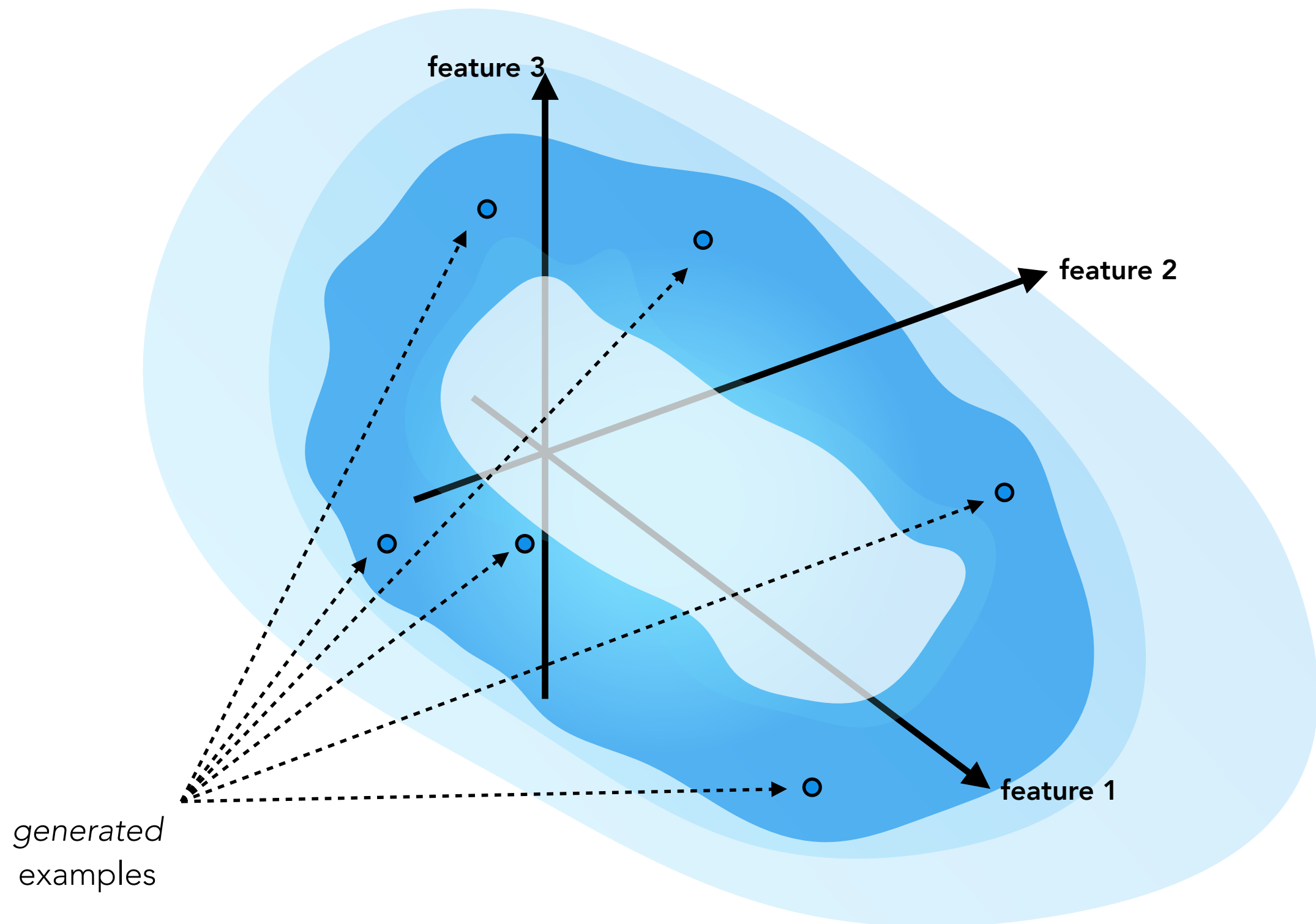
DENSITY ESTIMATION

estimating the density of the empirically observed data distribution

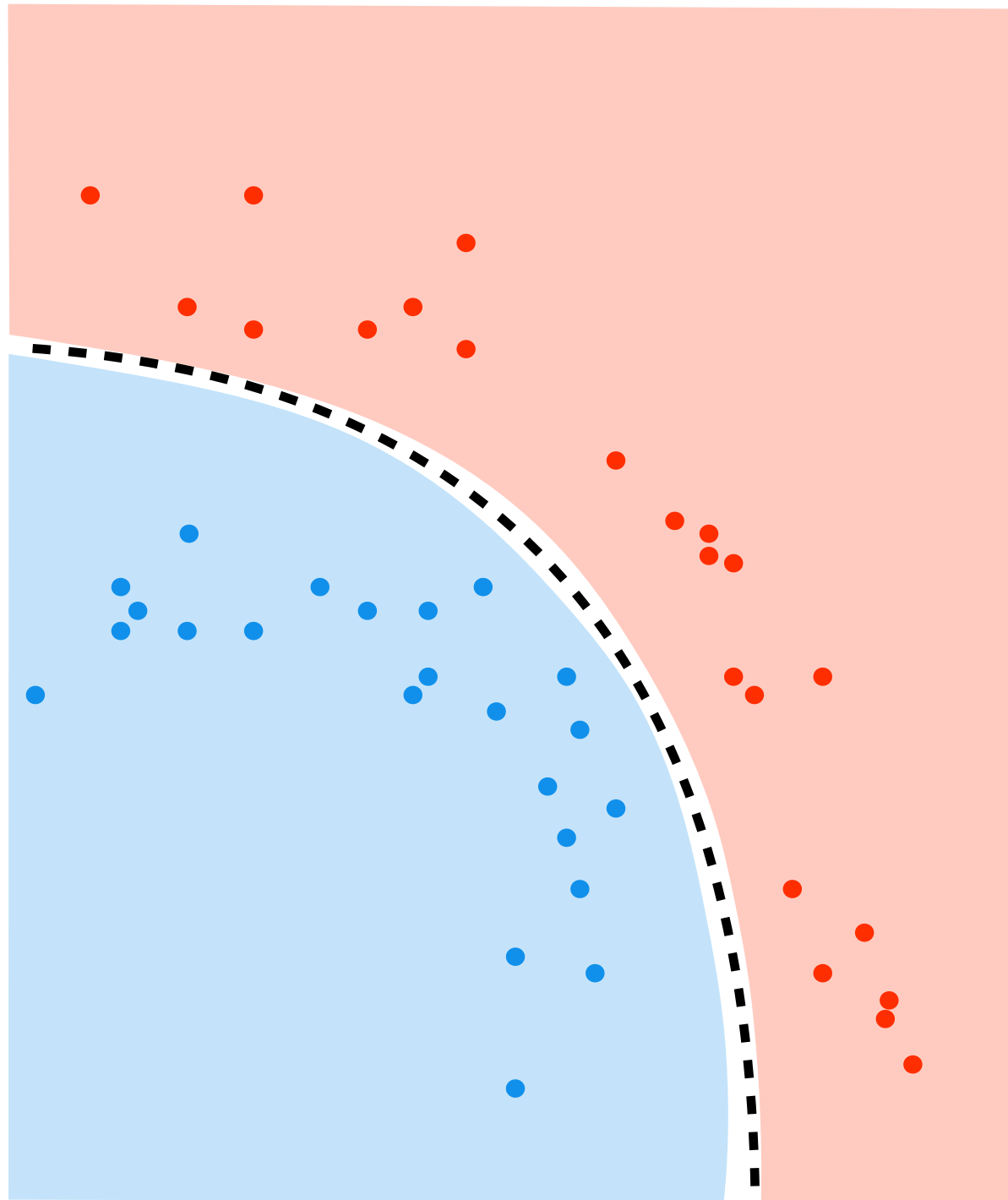
GENERATIVE MODEL

a model of the density of the data distribution

by modeling the data distribution,
generative models are able to **generate** new data examples



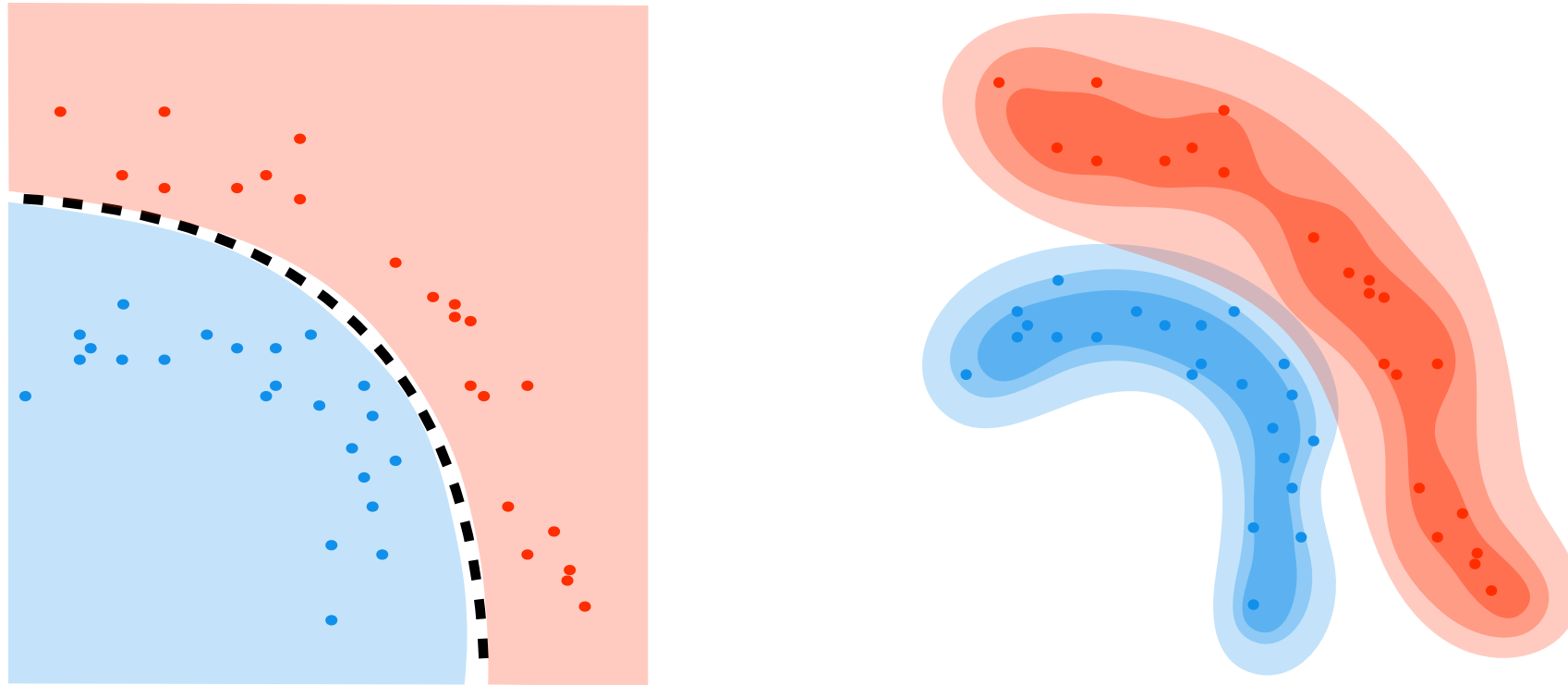
discriminative model



generative model



discriminative models vs. generative models



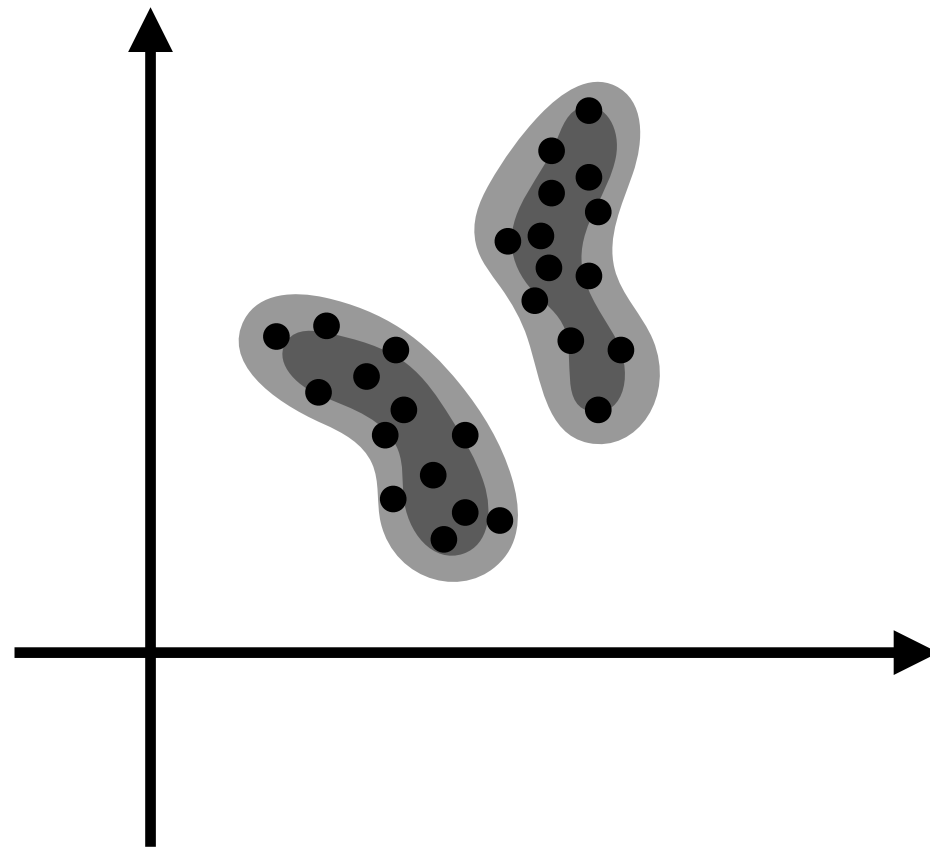
can both be trained using supervised learning

generative models are often easier to train with unsupervised methods

generative models typically require more modeling assumptions

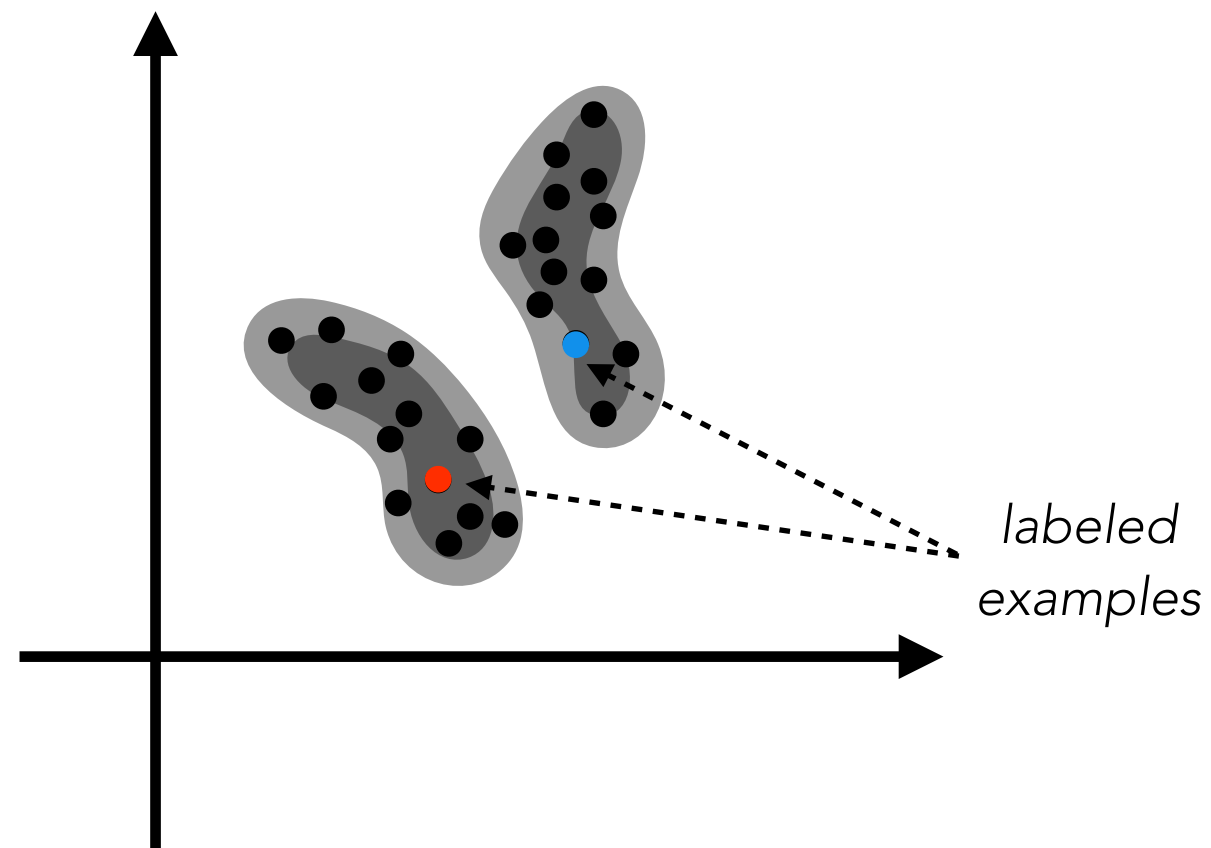
straightforward to quantify uncertainty with generative models

one of the main benefits of generative modeling is the ability to automatically extract structure from data



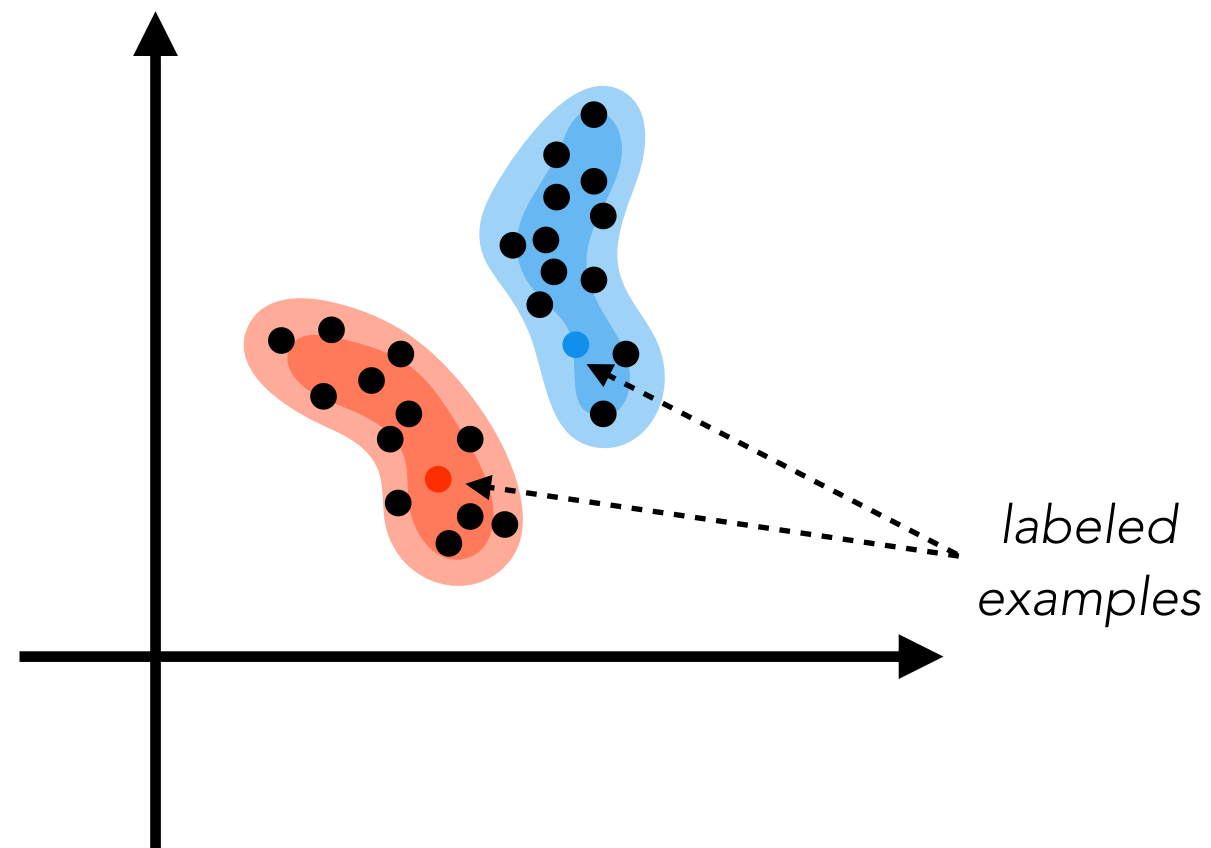
reducing the effective dimensionality of the data
can make it easier to learn and generalize on new tasks

one of the main benefits of generative modeling is the ability to automatically extract structure from data



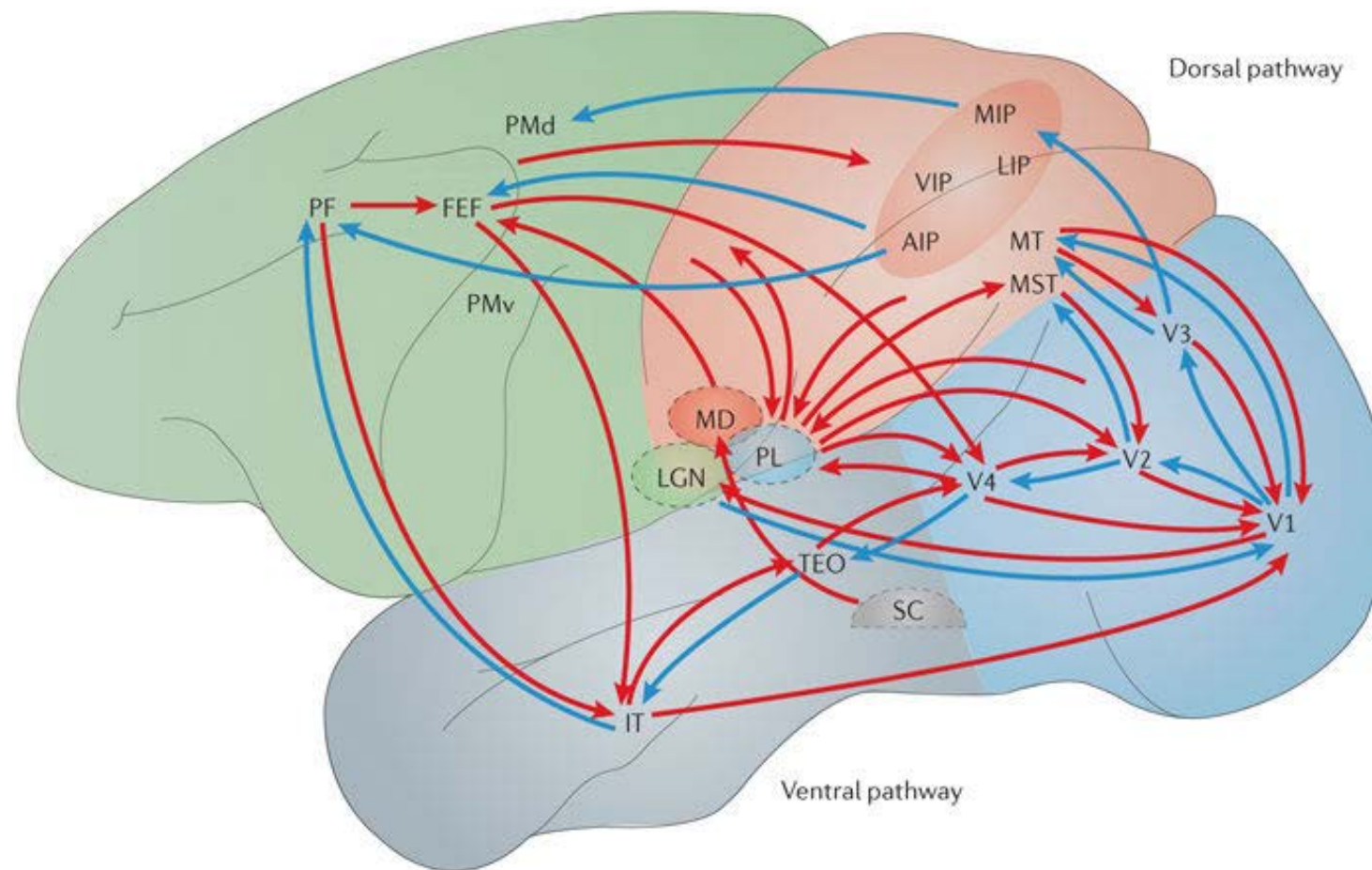
reducing the effective dimensionality of the data
can make it easier to learn and generalize on new tasks

one of the main benefits of generative modeling is the ability to automatically extract structure from data



reducing the effective dimensionality of the data
can make it easier to learn and generalize on new tasks

any model that has an output in the data space
can be considered a generative model



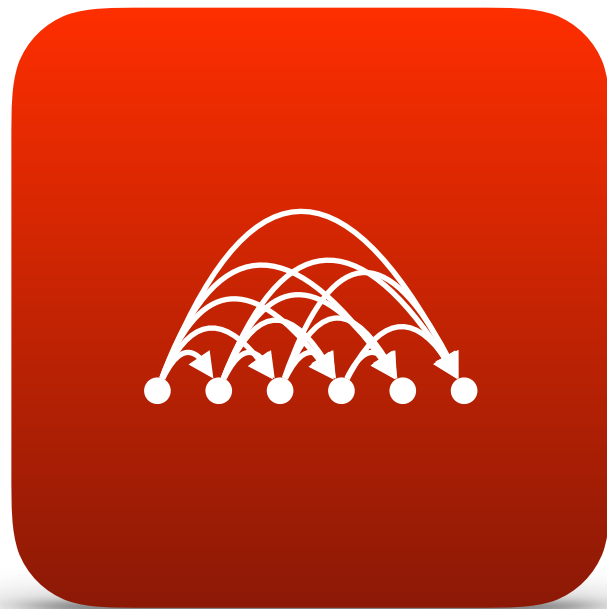
Nature Reviews | Neuroscience

nervous systems appear to use this mechanism in part
prediction of sensory input using "top-down" pathways

deep generative model

a generative model that uses deep neural networks
to model the data distribution

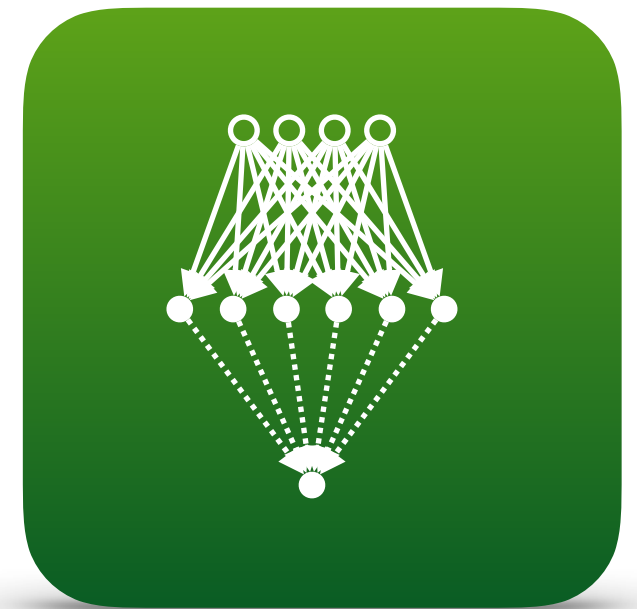
FAMILIES OF (DEEP) GENERATIVE MODELS



auto-regressive
models



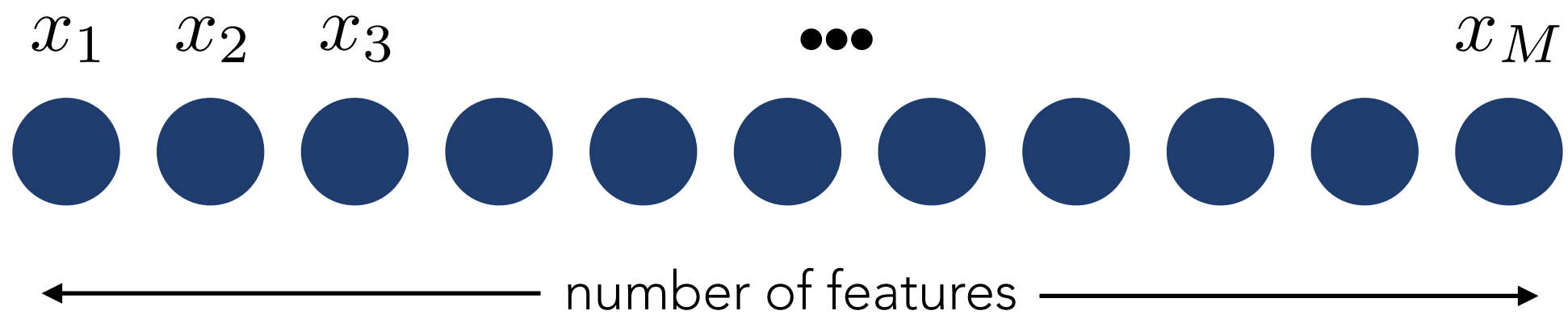
latent variable
models



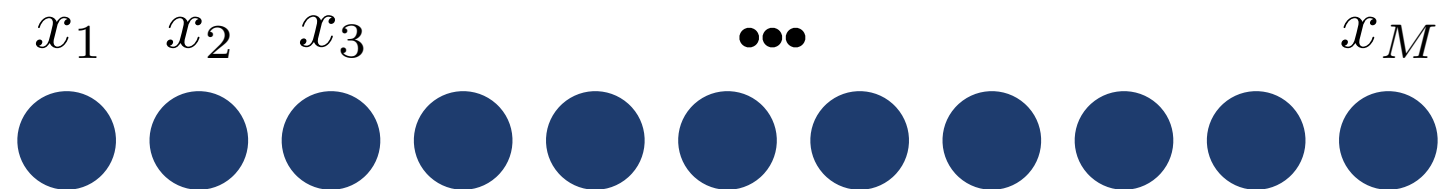
implicit
models

AUTO-REGRESSIVE MODELS

a data example



$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_M)$$



$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_M)$$

use chain rule of probability to split the joint distribution into a product of conditional distributions

definition of
conditional probability

$$p(a|b) = \frac{p(a, b)}{p(b)} \longrightarrow p(a, b) = p(a|b)p(b)$$

recursively apply to $p(x_1, x_2, \dots, x_M)$

$$p(x_1, x_2, \dots, x_M) = p(x_1|x_2, \dots, x_M)p(x_2, \dots, x_M)$$

⋮

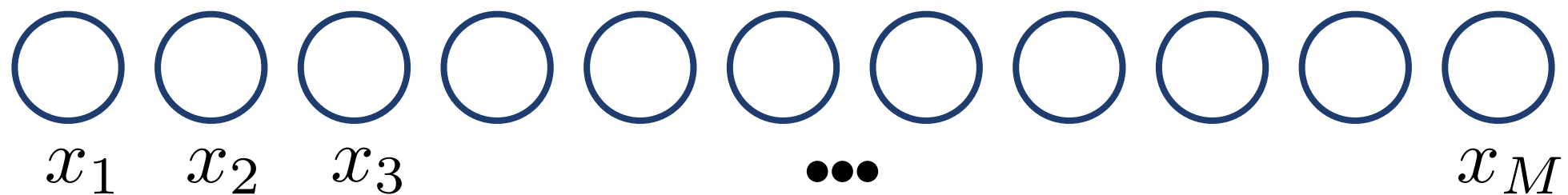
$$p(x_1, x_2, \dots, x_M) = p(x_1|x_2, \dots, x_M)p(x_2|x_3, \dots, x_M) \dots p(x_{M-1}|x_M)p(x_M)$$

note: conditioning order is arbitrary

$$p(x_1, \dots, x_M) = \prod_{j=1}^M p(x_j|x_1, \dots, x_{j-1})$$

model the conditional distributions of the data

learn to ***auto-regress*** to the missing values



model the conditional distributions of the data

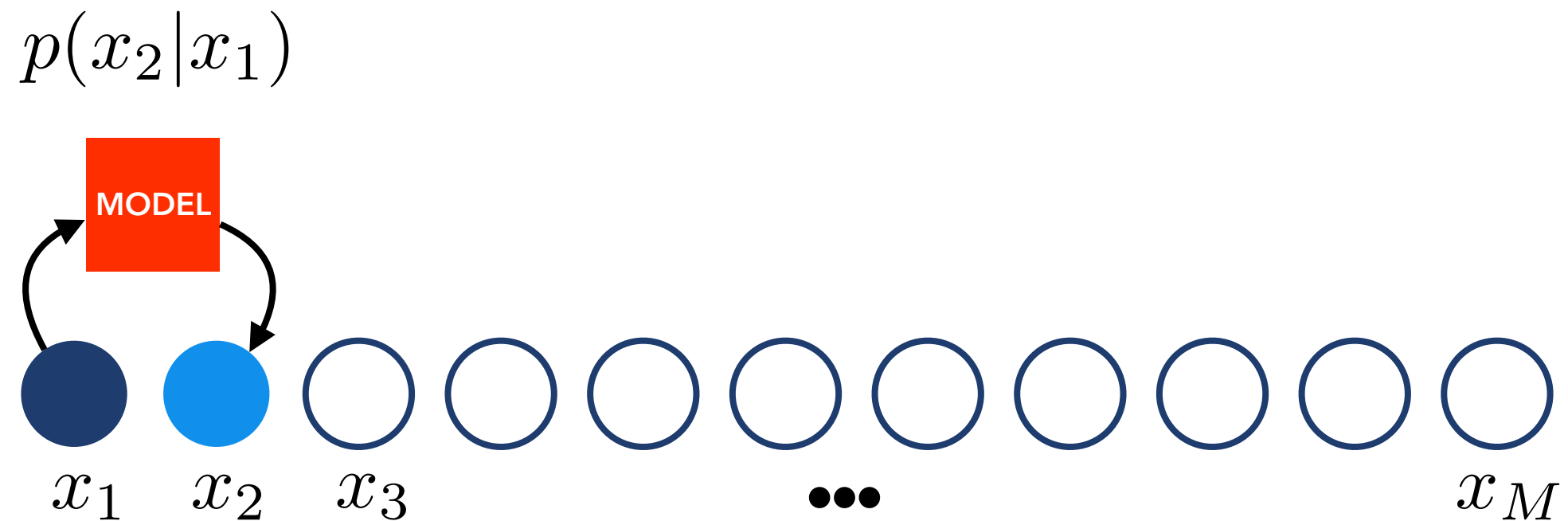
learn to **auto-regress** to the missing values

$$p(x_1)$$



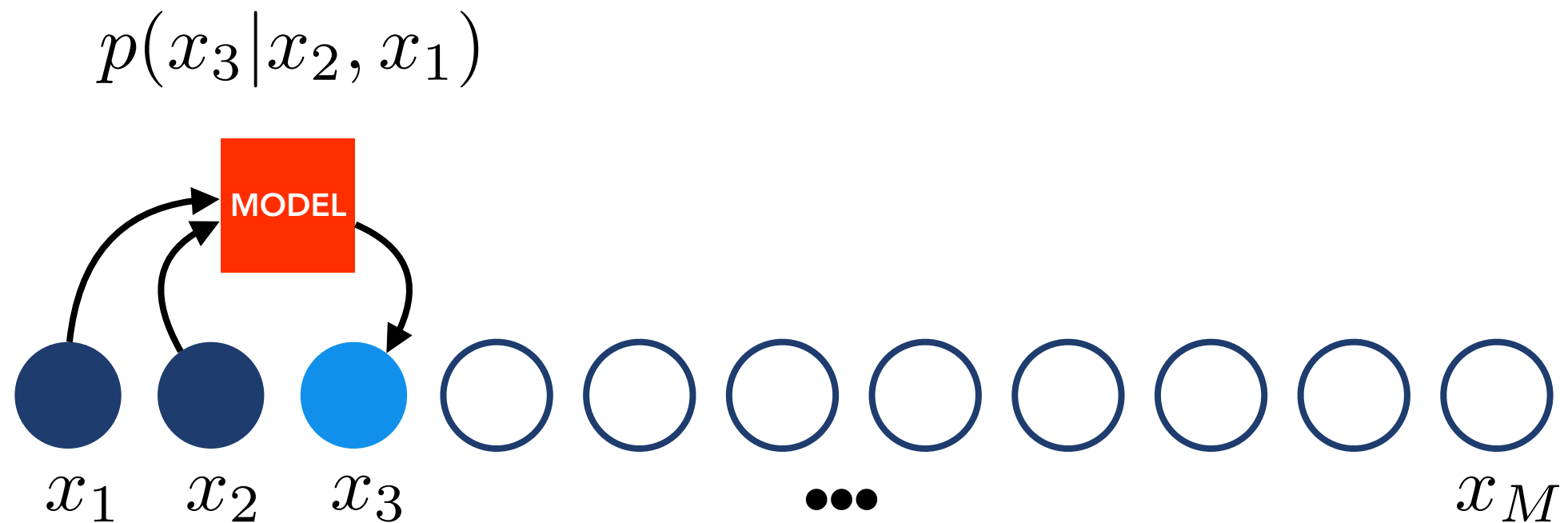
model the conditional distributions of the data

learn to **auto-regress** to the missing values



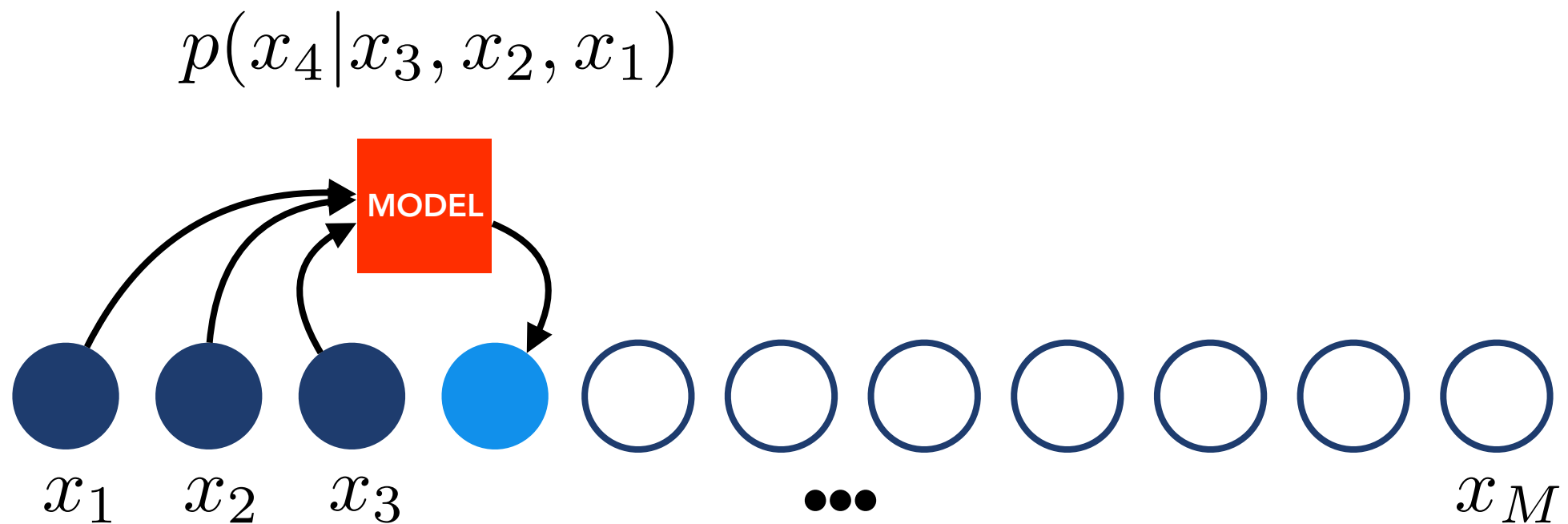
model the conditional distributions of the data

learn to **auto-regress** to the missing values



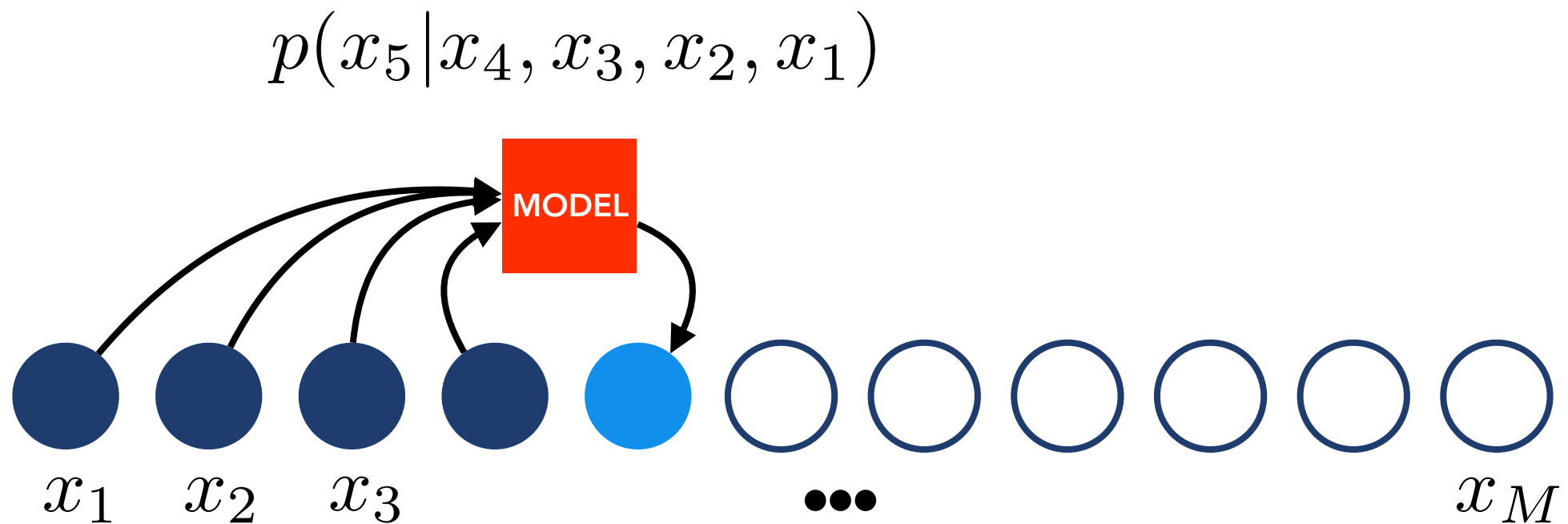
model the conditional distributions of the data

learn to **auto-regress** to the missing values



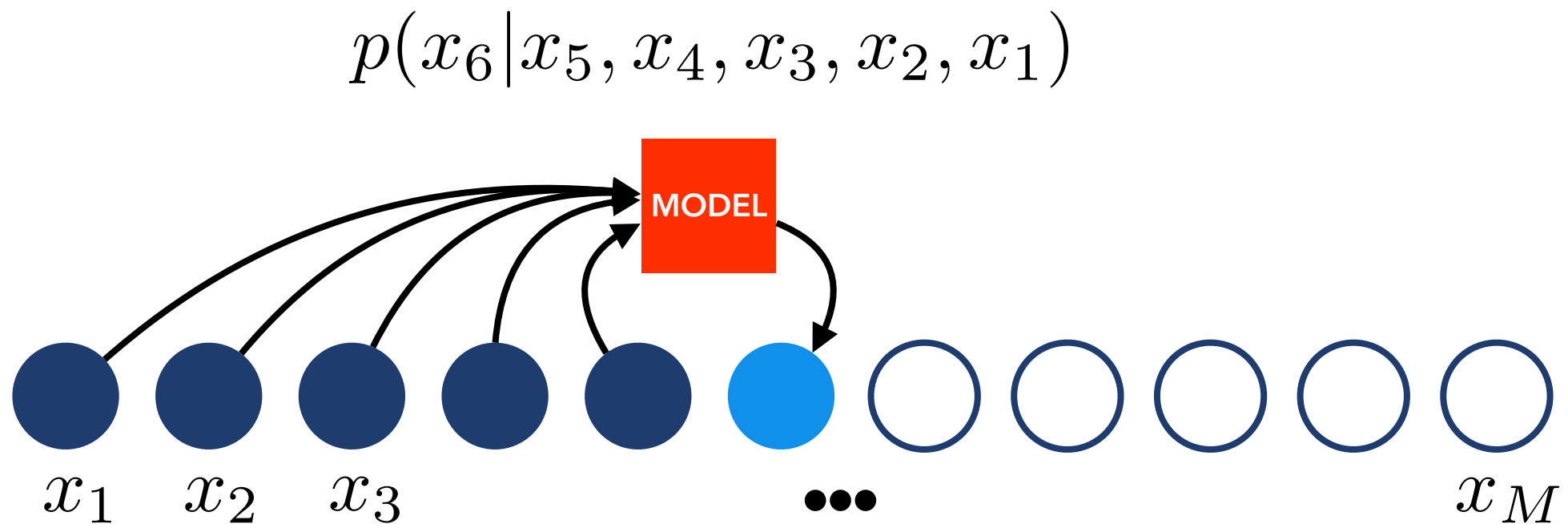
model the conditional distributions of the data

learn to **auto-regress** to the missing values



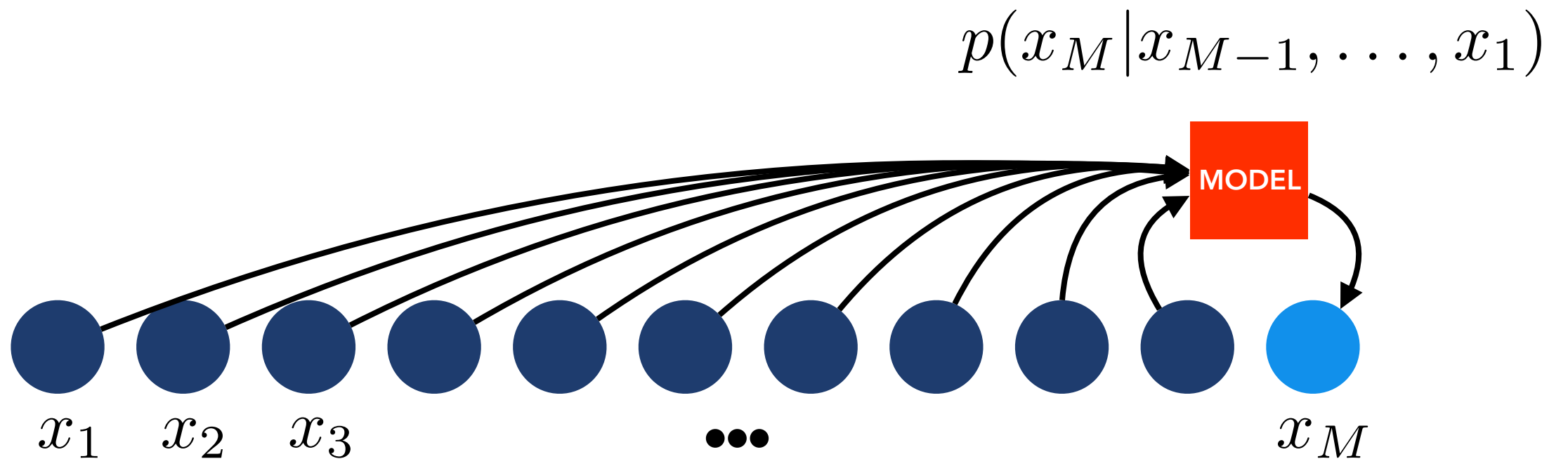
model the conditional distributions of the data

learn to **auto-regress** to the missing values



model the conditional distributions of the data


learn to **auto-regress** to the missing values



maximum likelihood

to fit the model to the empirical data distribution,
maximize the *likelihood* of the true data examples


likelihood:
$$p(\mathbf{x}) = \prod_{j=1}^M p(x_j | \mathbf{x}_{<j})$$



auto-regressive
conditionals

*optimize the parameters to assign high (log) probability
to the true data examples*

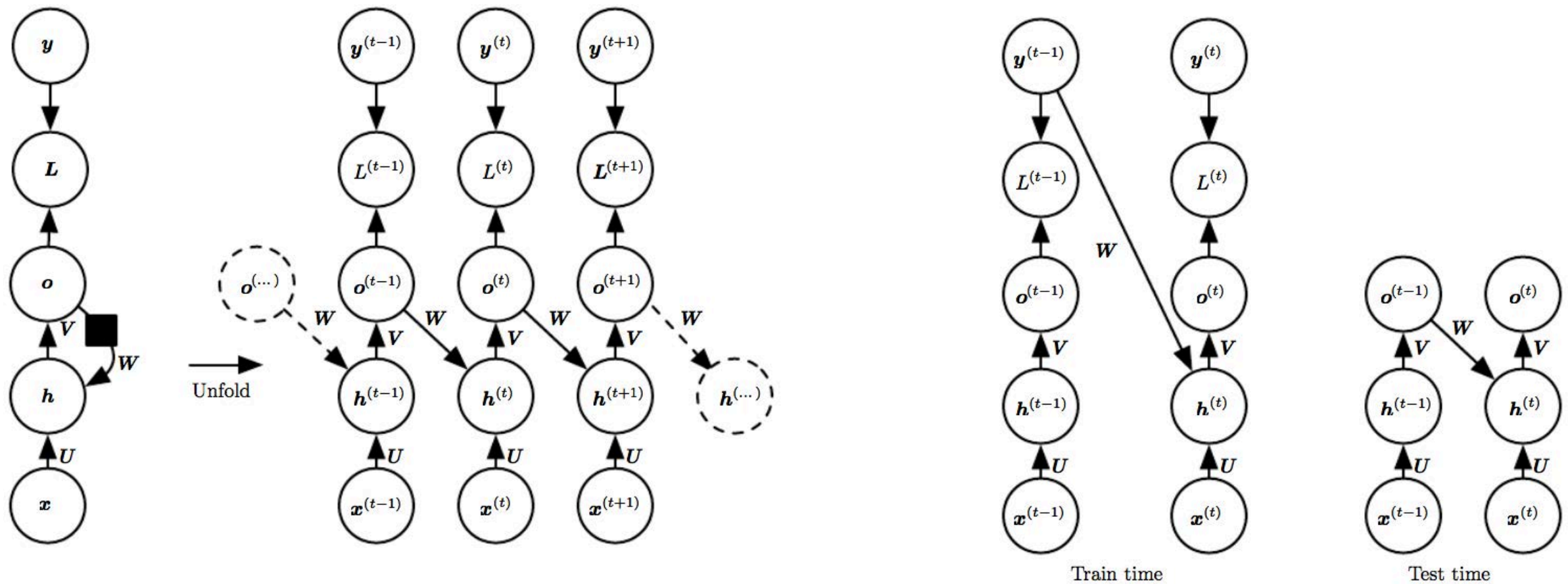
learning:
$$\theta^* = \operatorname{argmax}_{\theta} \log p(\mathbf{x})$$



logarithm for
numerical stability

models

can parameterize conditional distributions using a **recurrent neural network**

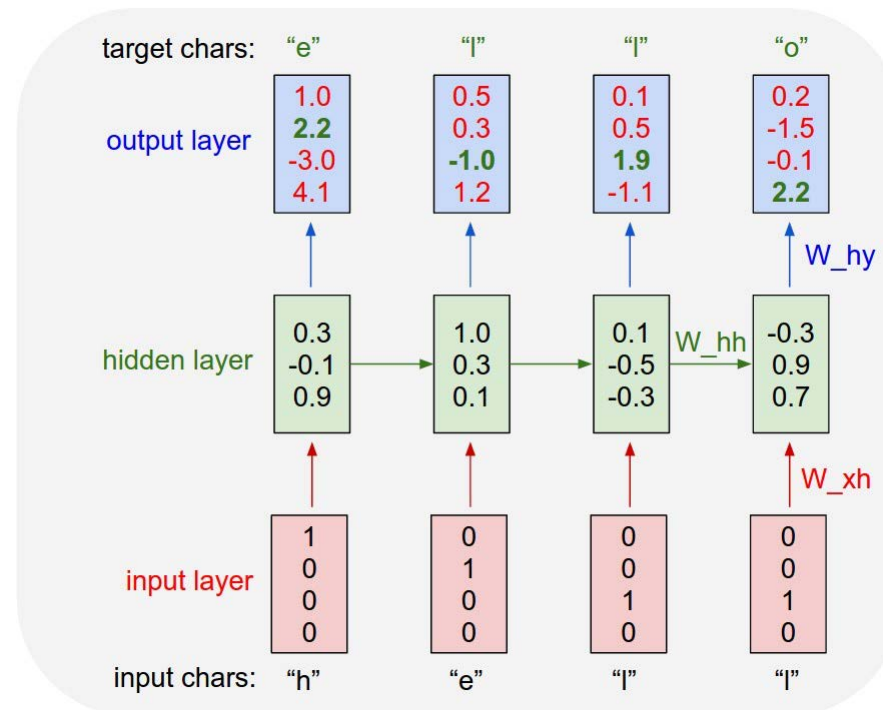


unrolling auto-regressive generation from an RNN

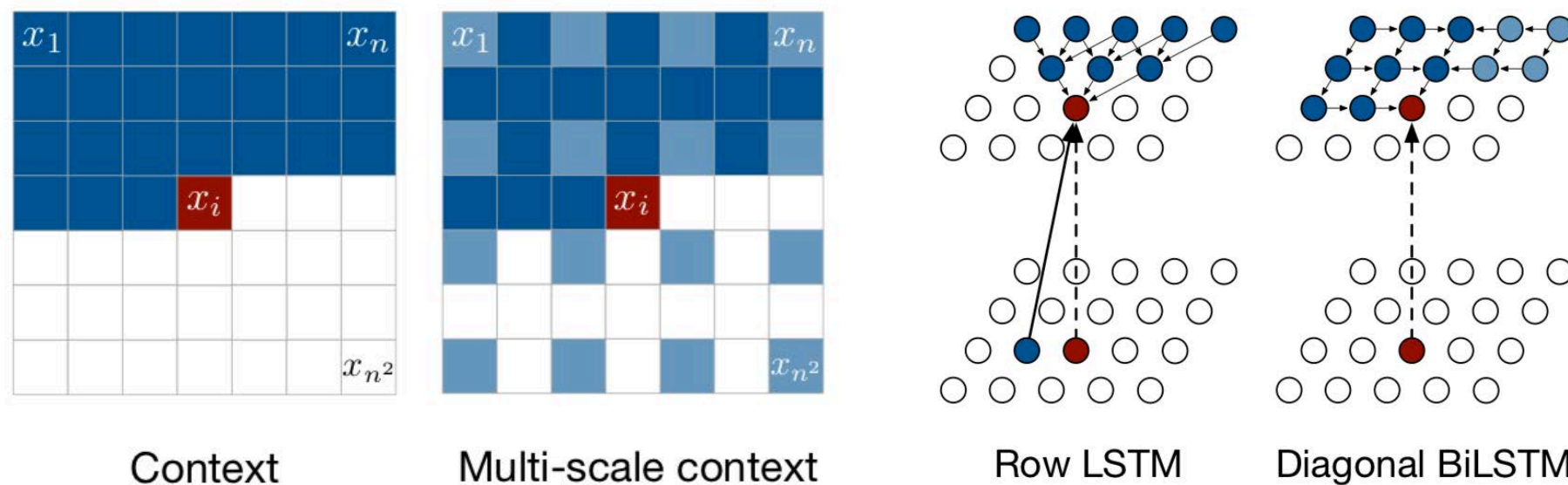
"teacher forcing"

models

can parameterize conditional distributions using a **recurrent neural network**



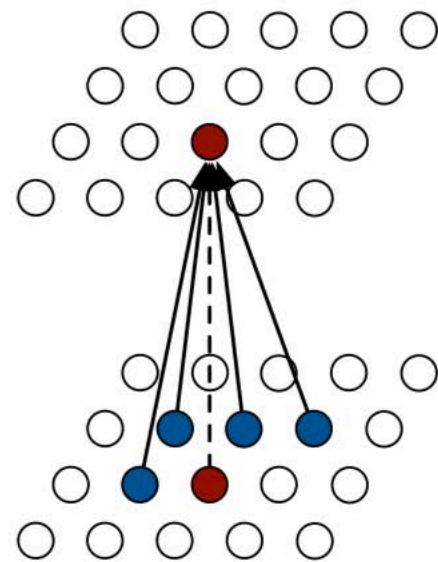
The Unreasonable Effectiveness of Recurrent Neural Networks, Karpathy, 2015



Pixel Recurrent Neural Networks, van den Oord et al., 2016

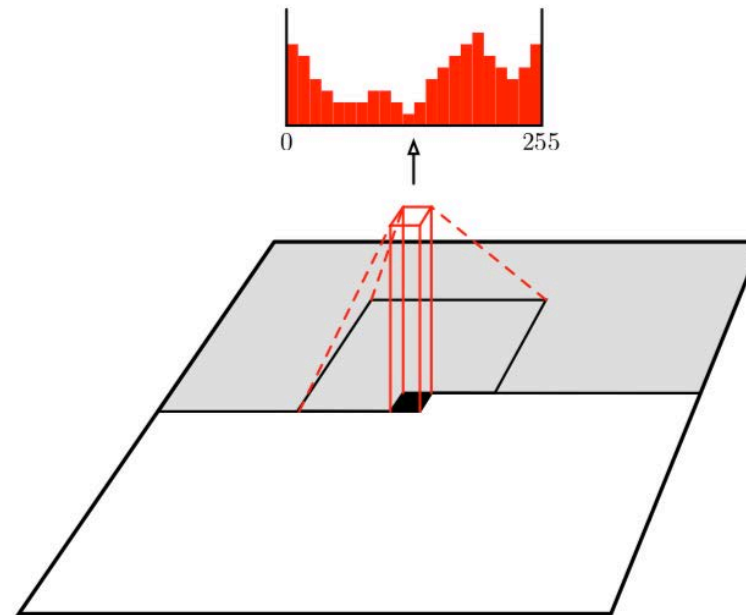
models

can also condition on a local window using **convolutional neural networks**

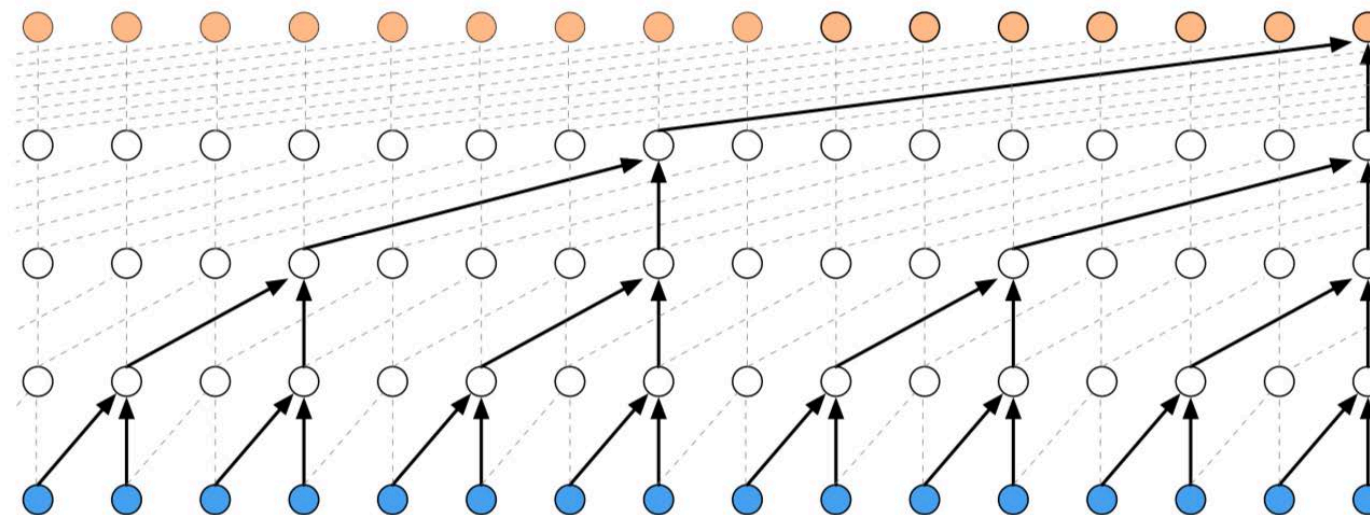
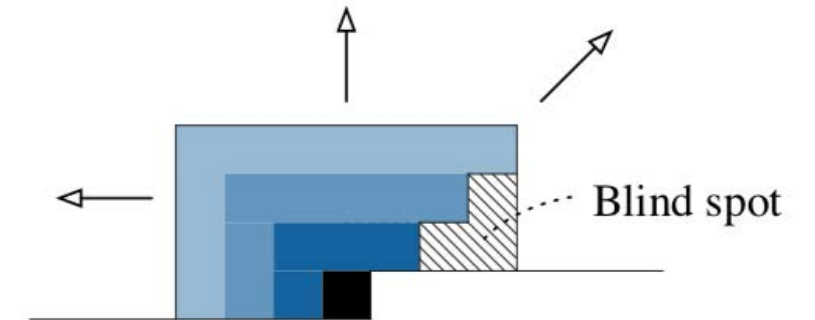


PixelCNN

Pixel Recurrent Neural Networks,
van den Oord et al., 2016



Conditional Image Generation with PixelCNN Decoders,
van den Oord et al., 2016



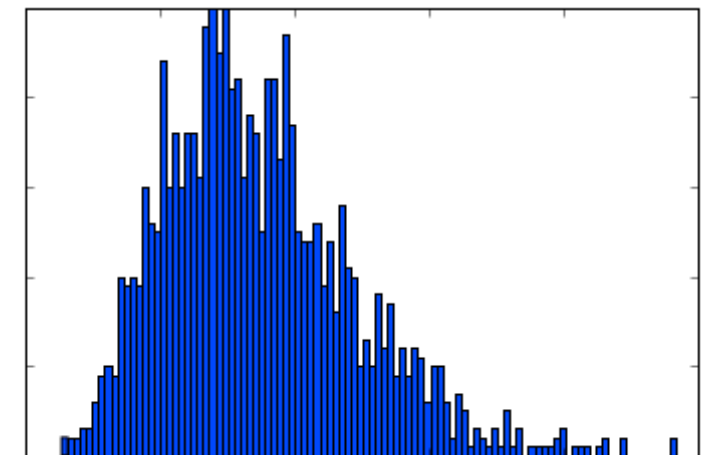
WaveNet: A Generative Model for Raw Audio, *van den Oord et al., 2016*

output distributions

need to choose a form for the conditional **output distribution**,
i.e. how do we express $p(x_j|x_1, \dots, x_{j-1})$?

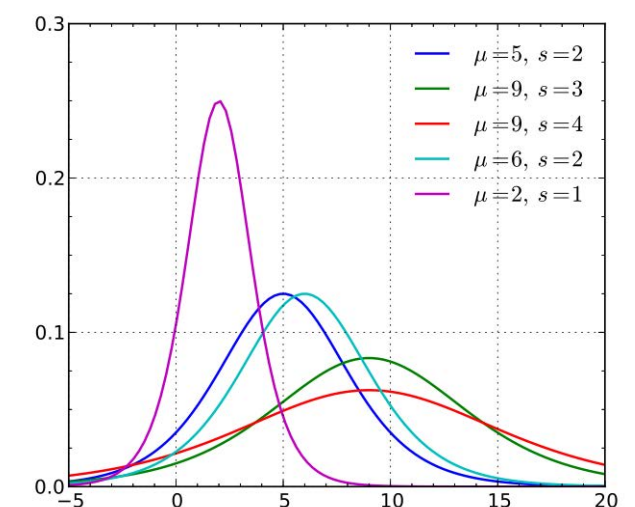
model the data as **categorical** variables

→ multinomial output



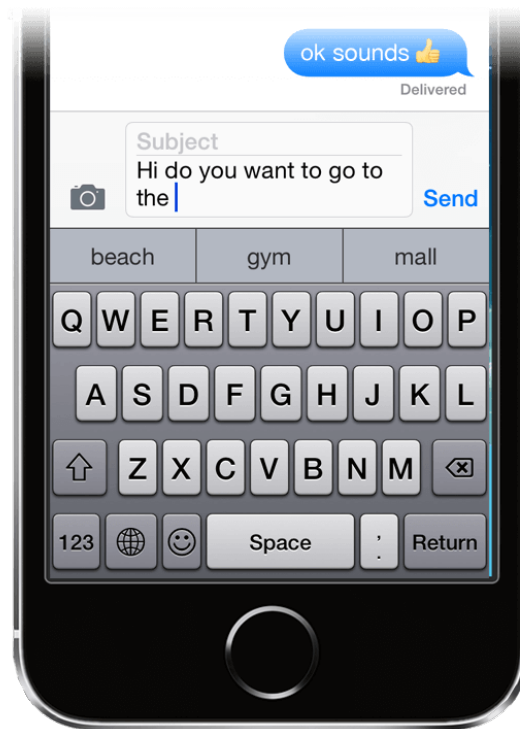
model the data as **continuous** variables

→ Gaussian, logistic, etc. output



example applications

text

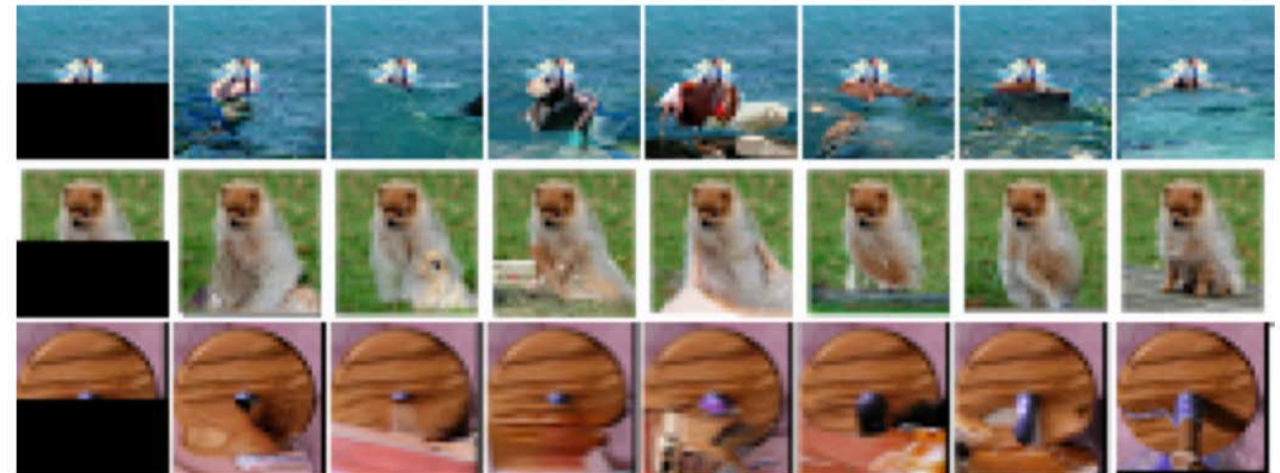


images

occluded

completions

original



Pixel Recurrent Neural Networks, *van den Oord et al., 2016*

speech



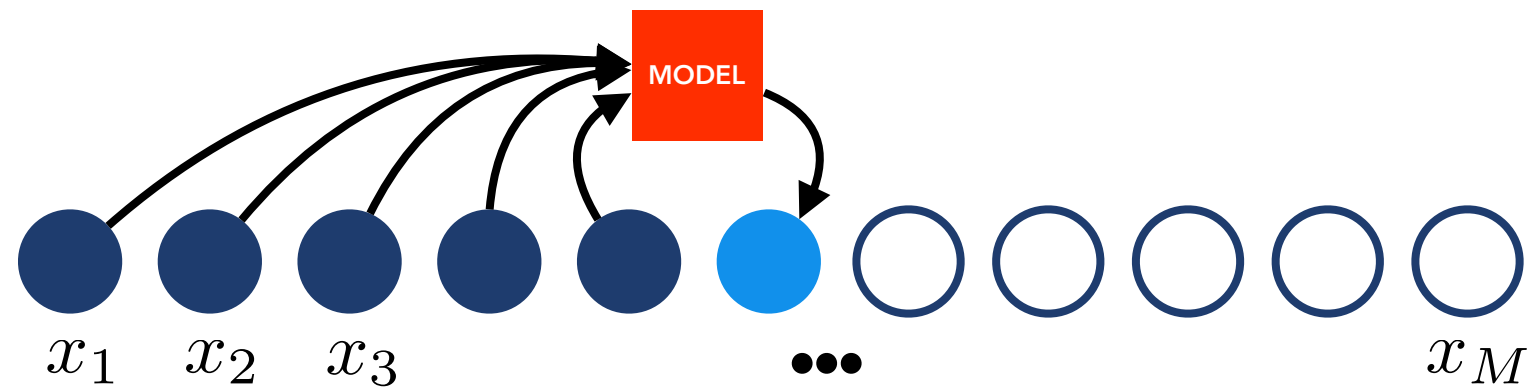
1 Second



WaveNet: A Generative Model for Raw Audio, *van den Oord et al., 2016*

recap: auto-regressive models

$$p(x_6 | x_5, x_4, x_3, x_2, x_1)$$



model conditional distributions to *auto-regress* to missing values

Pros

tractable and straightforward to
evaluate the (log) likelihood

great at capturing details

superior quantitative performance

Cons

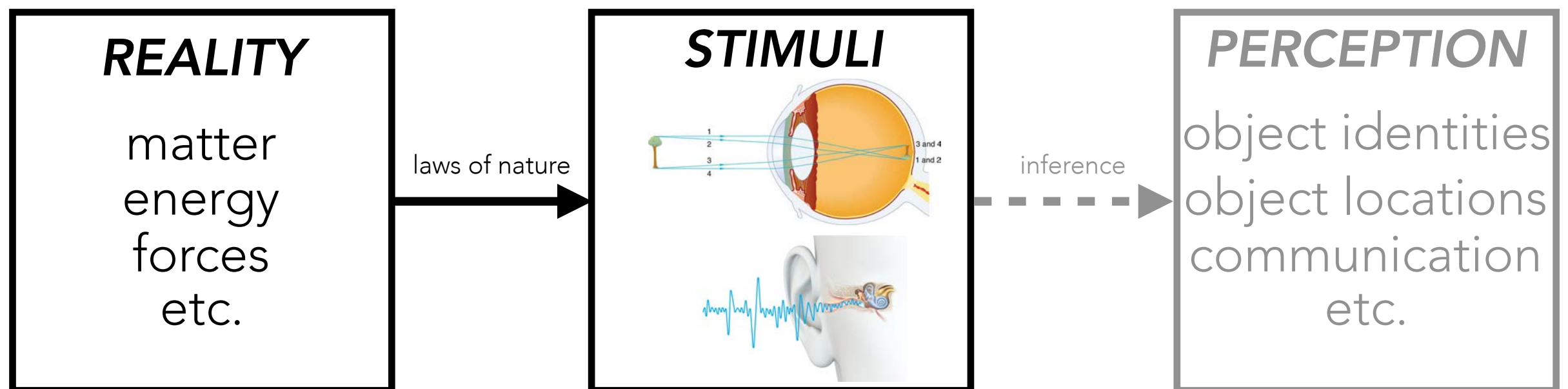
difficult to capture
“high-level” global structure

need to impose
conditioning order

sequential sampling is
computationally expensive

EXPLICIT LATENT VARIABLE MODELS

reality **generates** sensory stimuli from underlying latent phenomena

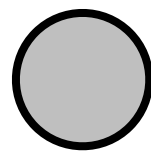


can use *latent variables* to help model these phenomena

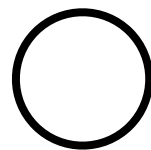
probabilistic graphical models provide a framework for modeling relationships between random variables

PLATE NOTATION

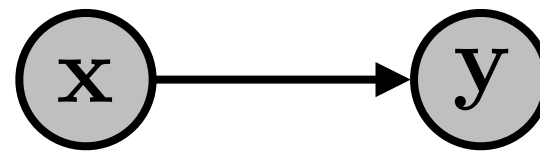
observed variable



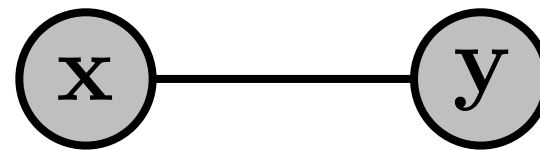
unobserved (latent)
variable



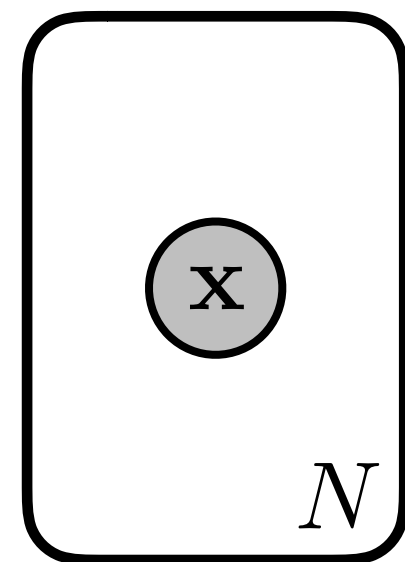
directed



undirected

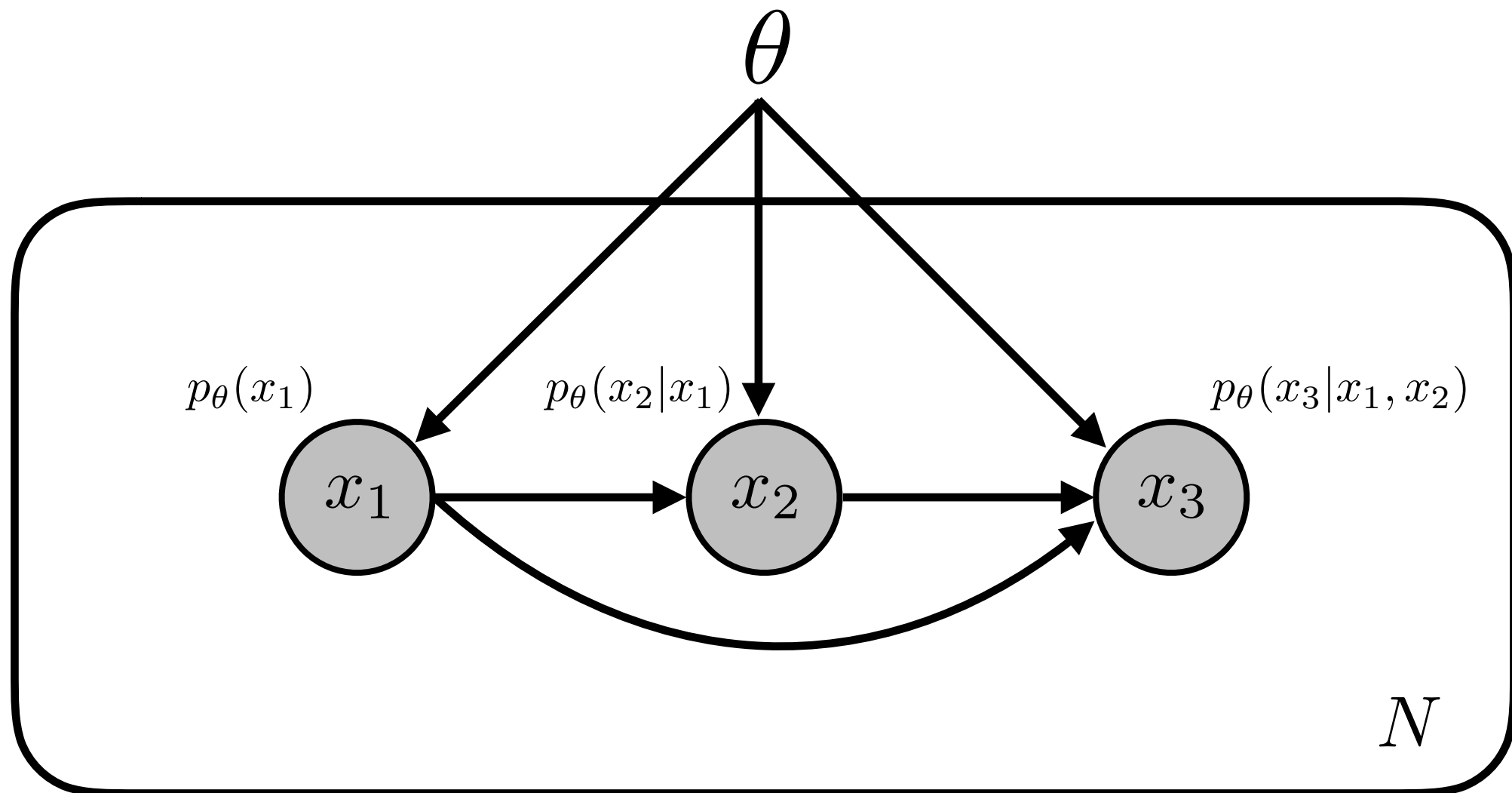


set of variables

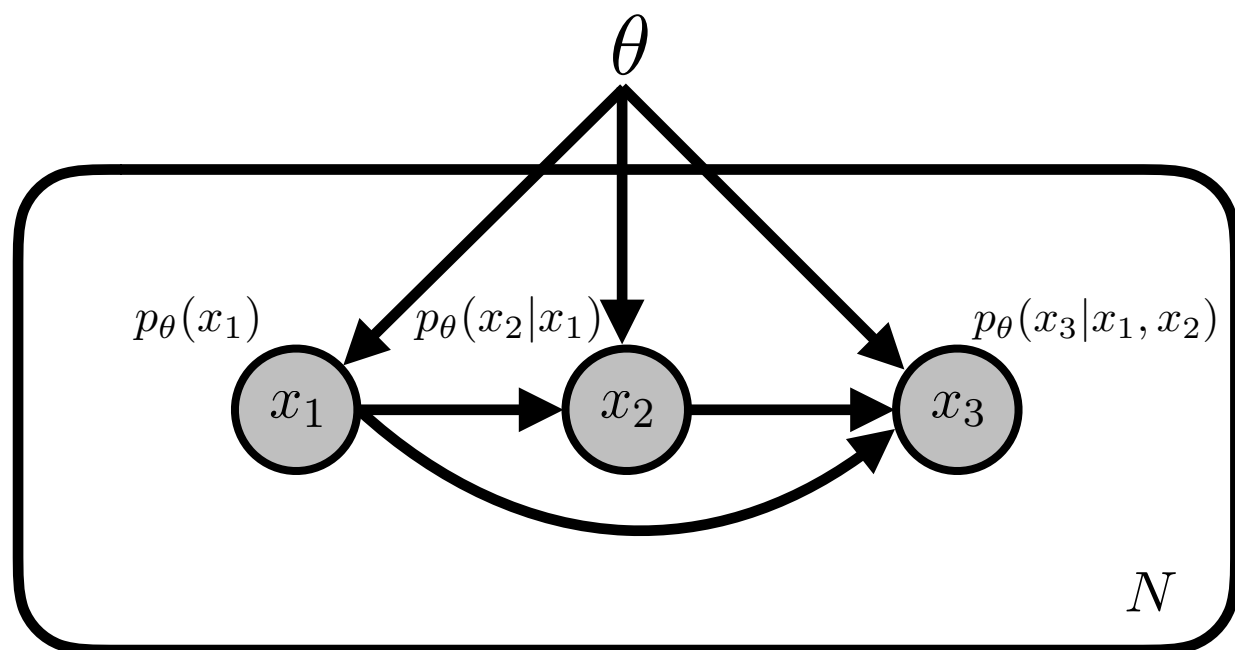


review exercise:

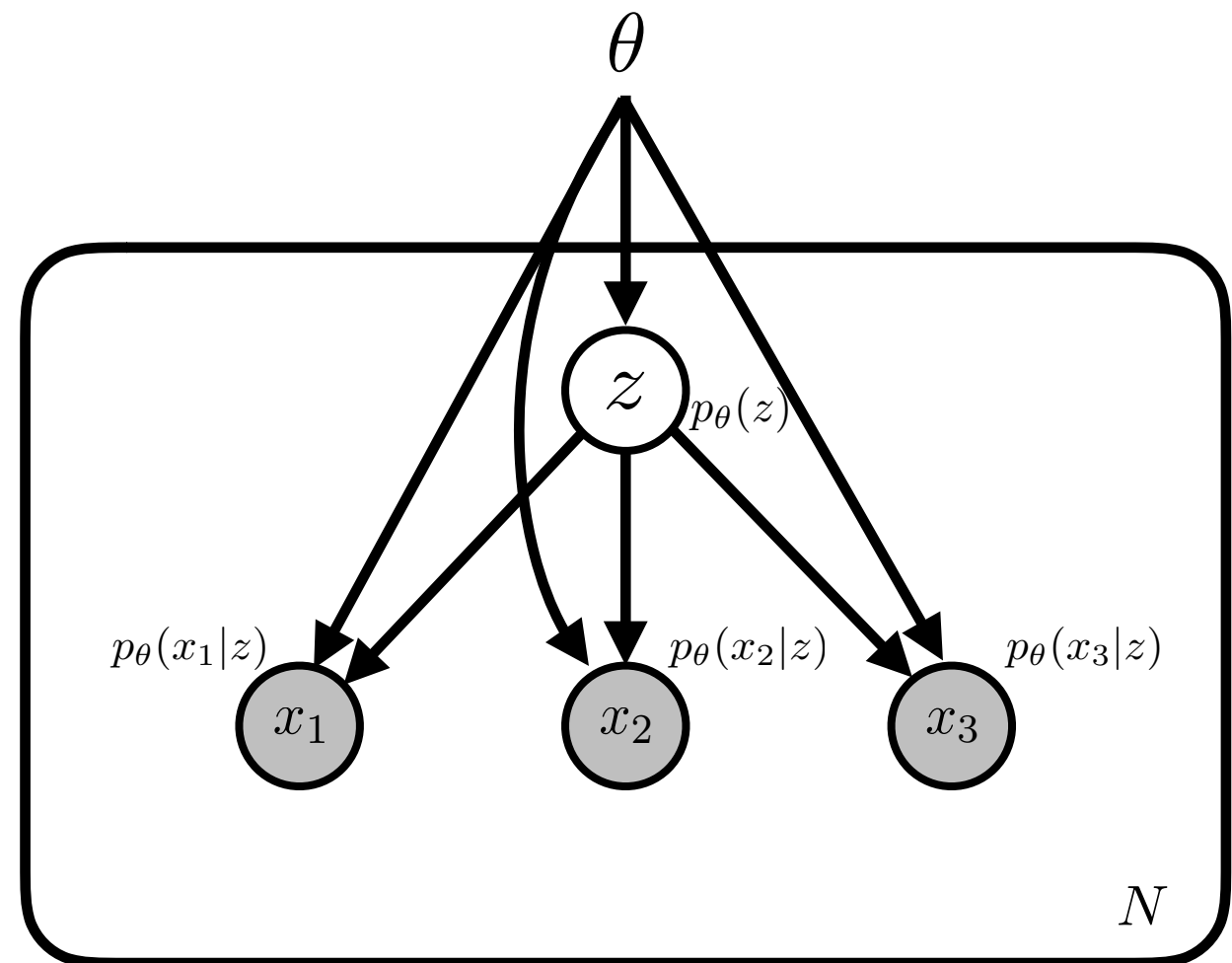
represent an auto-regressive model of 3 random variables
with plate notation



comparing *auto-regressive models* and *latent variable models*



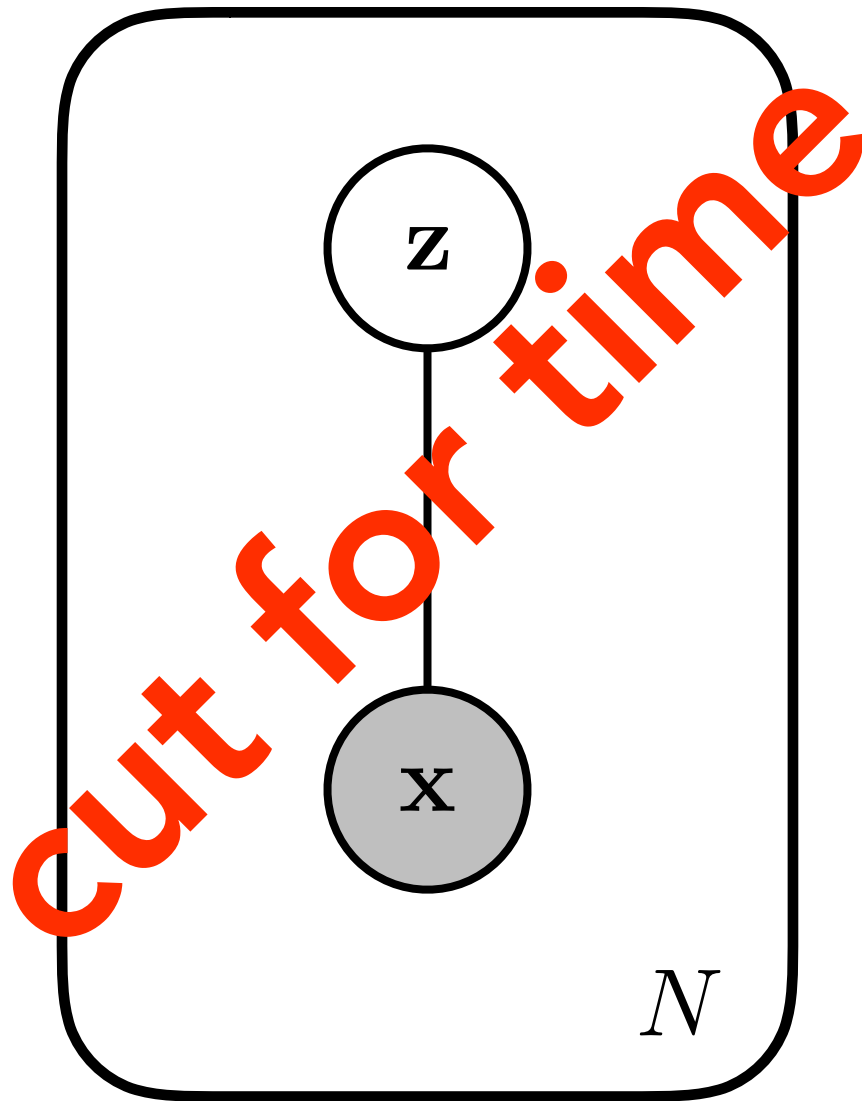
auto-regressive model



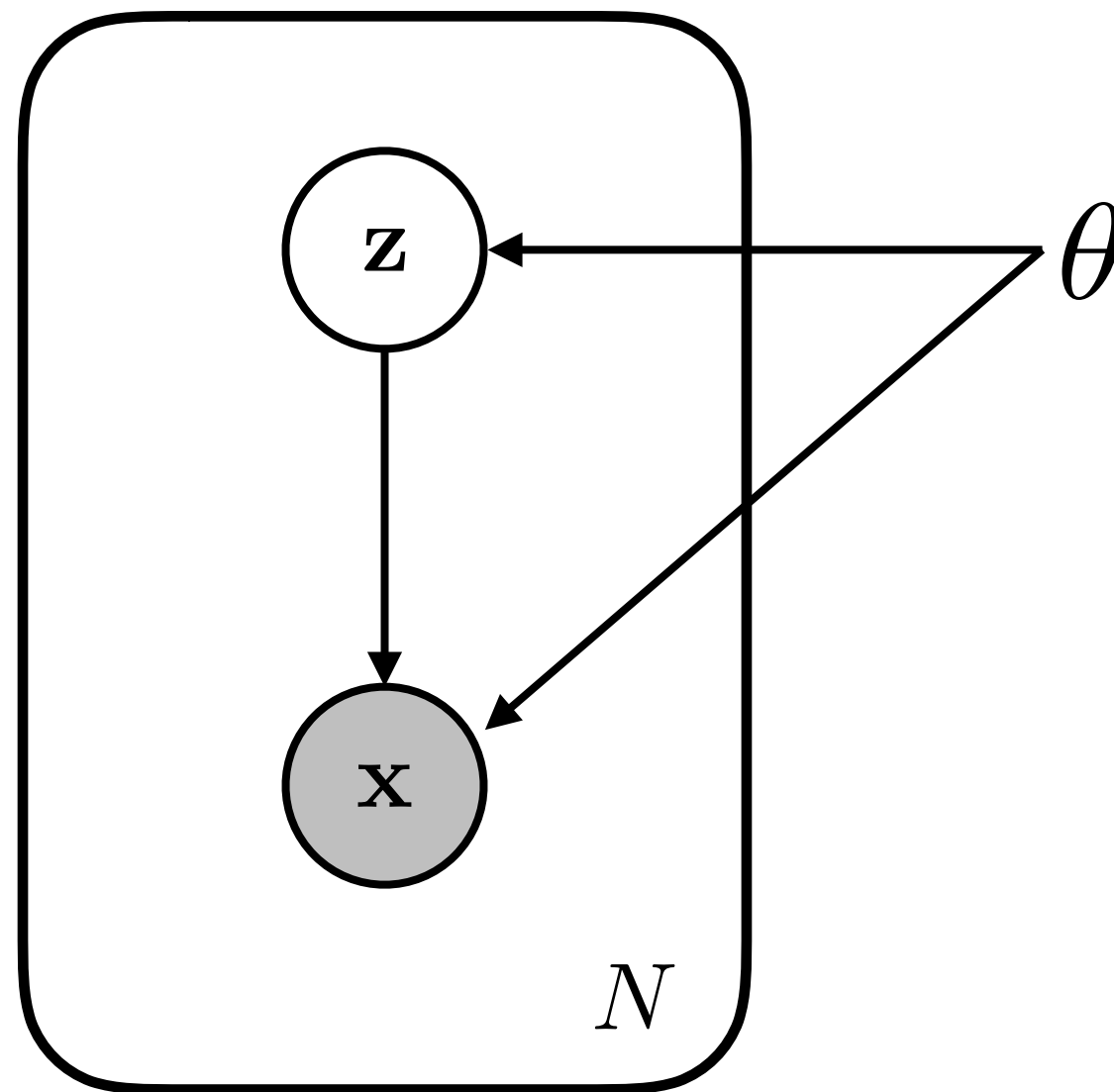
latent variable model

example: undirected latent variable model

→ restricted Boltzmann machine

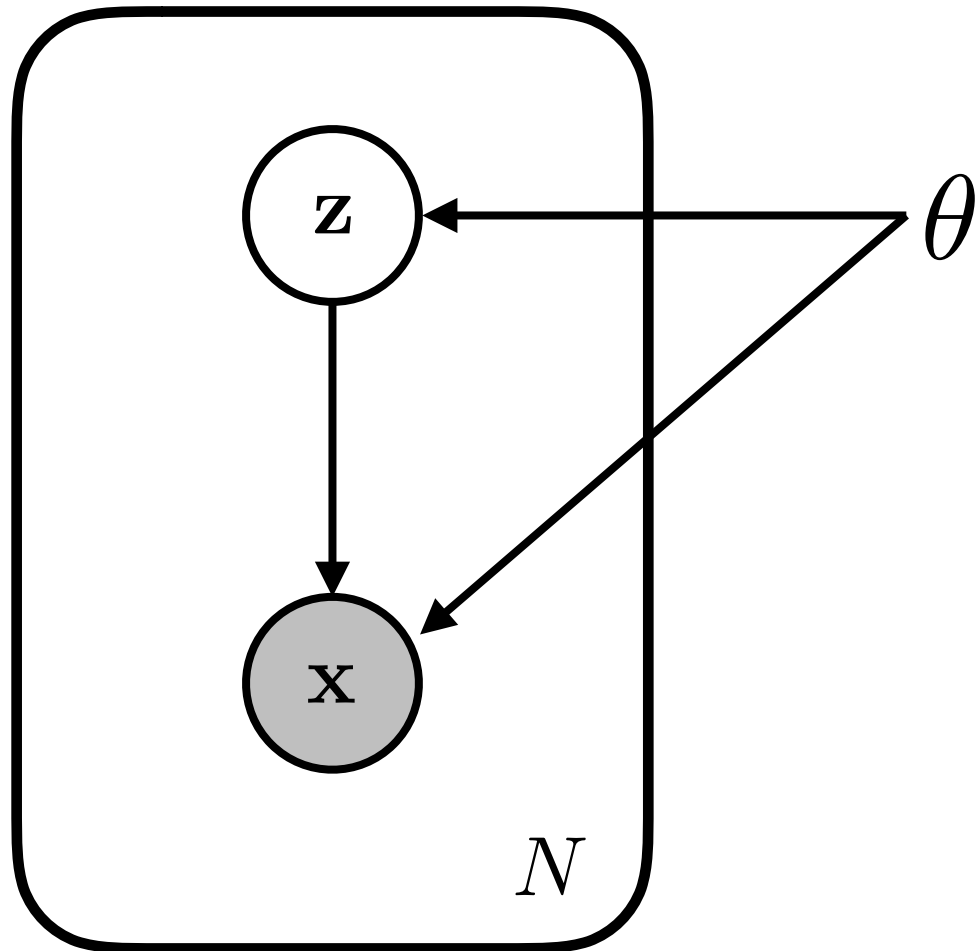


example: directed latent variable model



example: **directed** latent variable model

Generation



GENERATIVE MODEL

$$\overset{\text{joint}}{\curvearrowright} p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \overset{\text{prior}}{\curvearrowleft}$$

\uparrow
conditional likelihood

1. sample \mathbf{z} from $p(\mathbf{z})$
2. use \mathbf{z} samples to sample \mathbf{x} from $p(\mathbf{x}|\mathbf{z})$

intuitive example

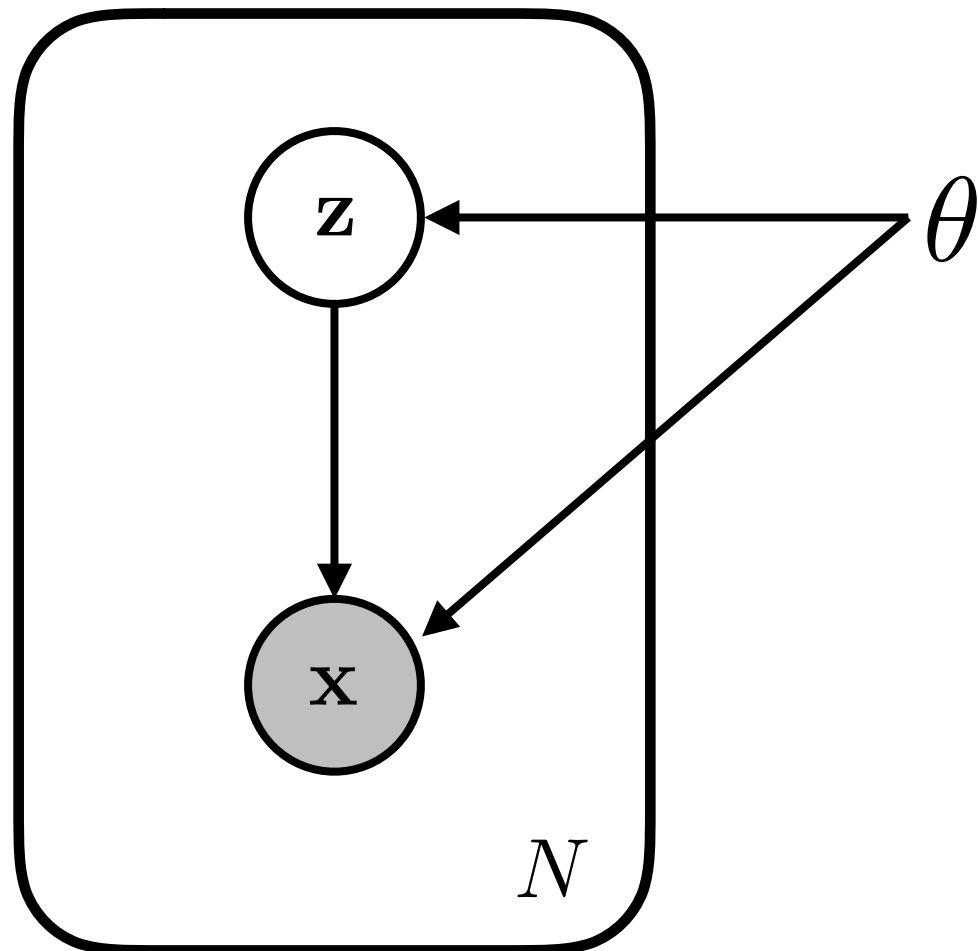
object $\sim p(\text{objects})$
lighting $\sim p(\text{lighting})$
background $\sim p(\text{bg})$

RENDER



example: **directed latent variable model**

Posterior Inference



INFERENCE

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})}$$

posterior

joint

marginal likelihood

use Bayes' rule (def. of cond. prob.)

provides conditional distribution
over latent variables

intuitive example

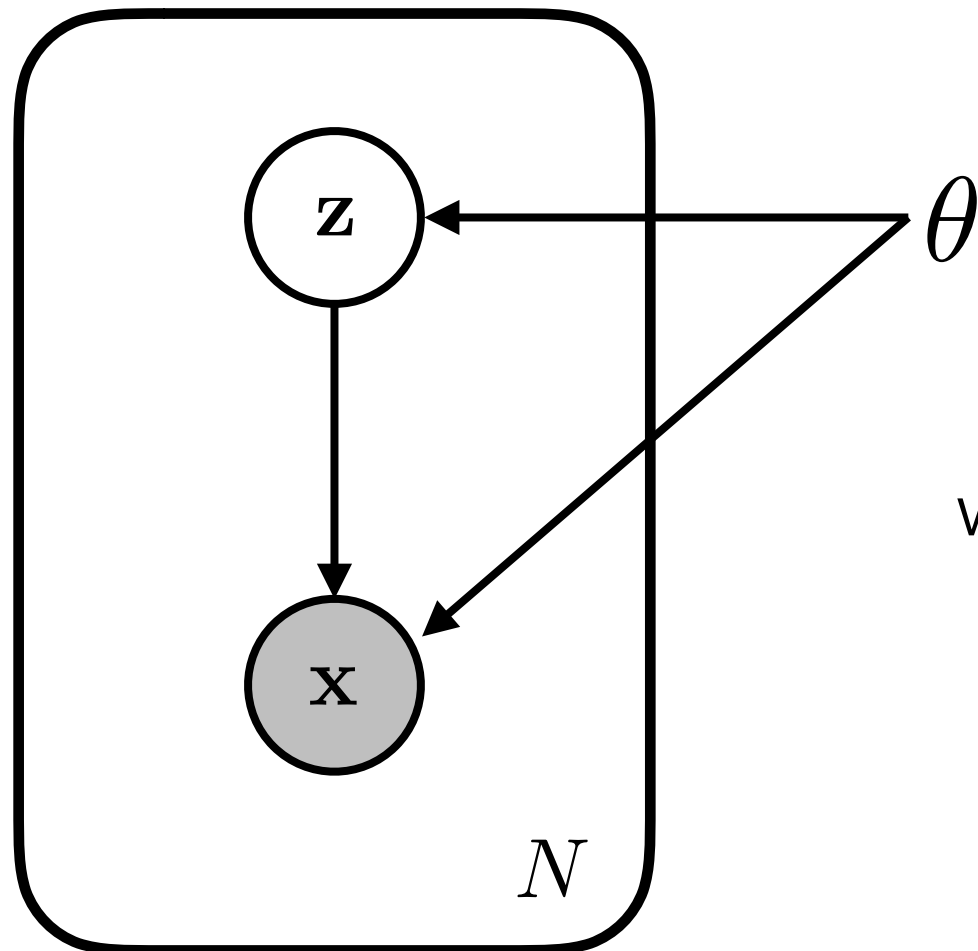


what is the probability that I am observing a cat
given these pixel observations?

$$p(\text{cat} | \text{image}) = \frac{p(\text{image} | \text{cat}) p(\text{cat})}{p(\text{image})}$$

example: **directed latent variable model**

Model Evaluation



MARGINALIZATION

$$\text{marginal likelihood} \rightarrow p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \leftarrow \text{joint}$$

to evaluate the likelihood of an observation,
we need to *marginalize* over all latent variables

i.e. consider all possible underlying states

intuitive example

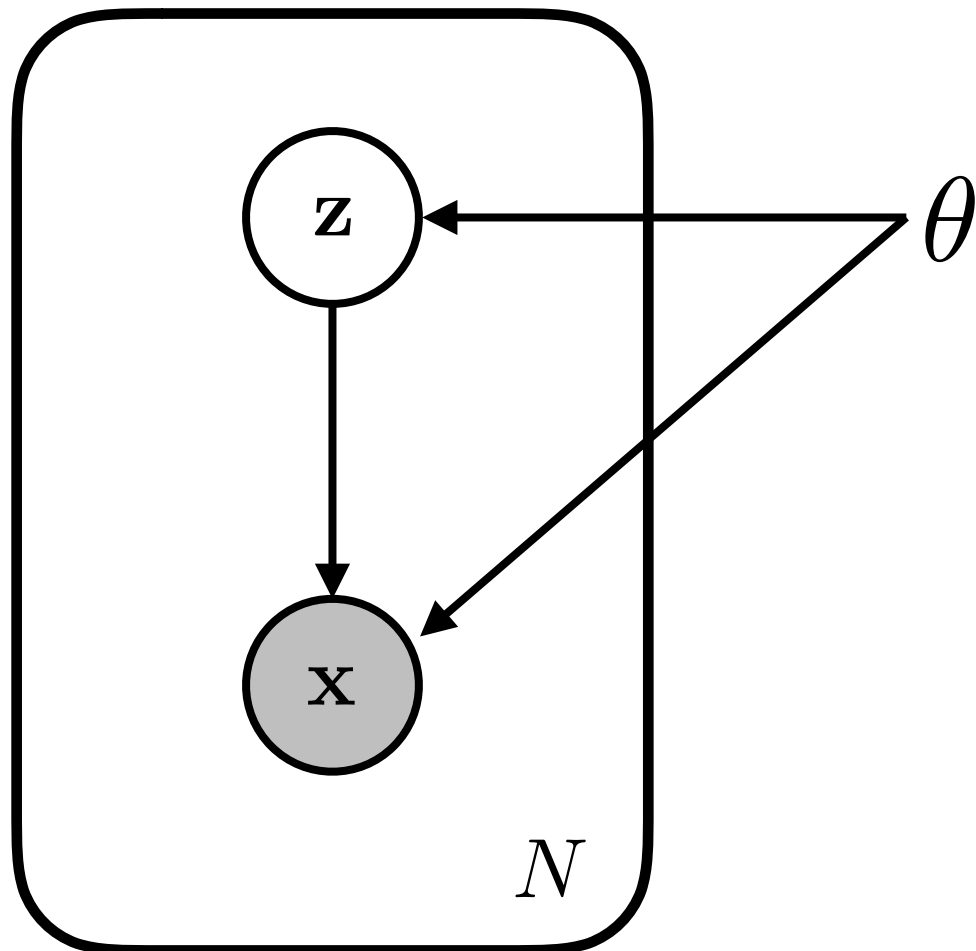


observation

how likely is this observation under my model?
(what is the probability of observing this?)

for all objects, lighting, backgrounds, etc.:
how plausible is this example?

example: **directed latent variable model**



to fit the model, we want to evaluate the *marginal* (log) likelihood of the data

$$\theta^* = \operatorname{argmax}_{\theta} \log p(\mathbf{x})$$

however, this is generally *intractable*, due to the integration over latent variables

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

integration in
high-dimensions

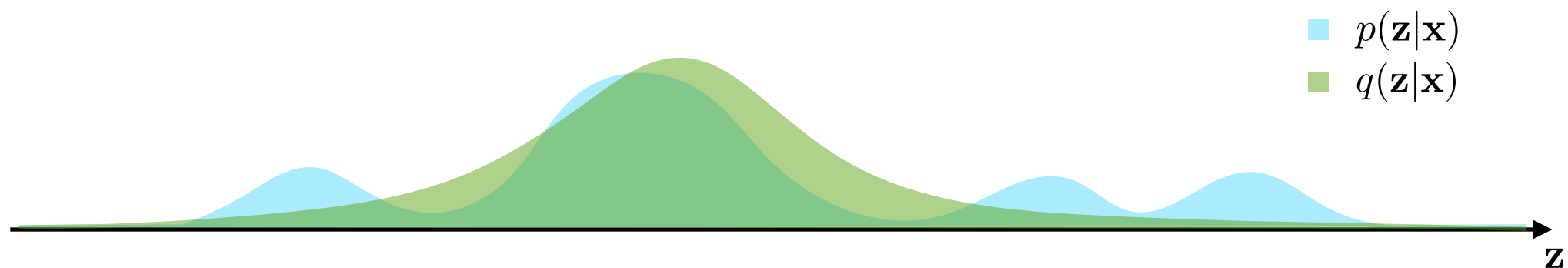
variational inference

main idea

instead of optimizing the (log) likelihood, optimize a **lower bound** on it

introduce an **approximate posterior**, then minimize KL-divergence to the true posterior

$$q^*(\mathbf{z}|\mathbf{x}) = \operatorname{argmin}_q D_{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))$$



evaluating KL-divergence involves evaluating $p(\mathbf{z}|\mathbf{x})$, instead maximize \mathcal{L} :

$$q^*(\mathbf{z}|\mathbf{x}) = \operatorname{argmax}_q \mathcal{L}$$

where \mathcal{L} is the *evidence lower bound (ELBO)*, defined as $\mathcal{L} \equiv \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})]$

\mathcal{L} provides a lower bound on $\log p(\mathbf{x})$, so we can use \mathcal{L} to (approximately) fit the model

$$\tilde{\theta}^* = \operatorname{argmax}_{\theta} \mathcal{L}$$

interpreting the ELBO

we can write the ELBO as

$$\begin{aligned}\mathcal{L} &\equiv \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \\ &= \underbrace{\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction}} - \underbrace{D_{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{regularization}}\end{aligned}$$

$q(\mathbf{z}|\mathbf{x})$ is optimized to represent the data while staying close to the prior

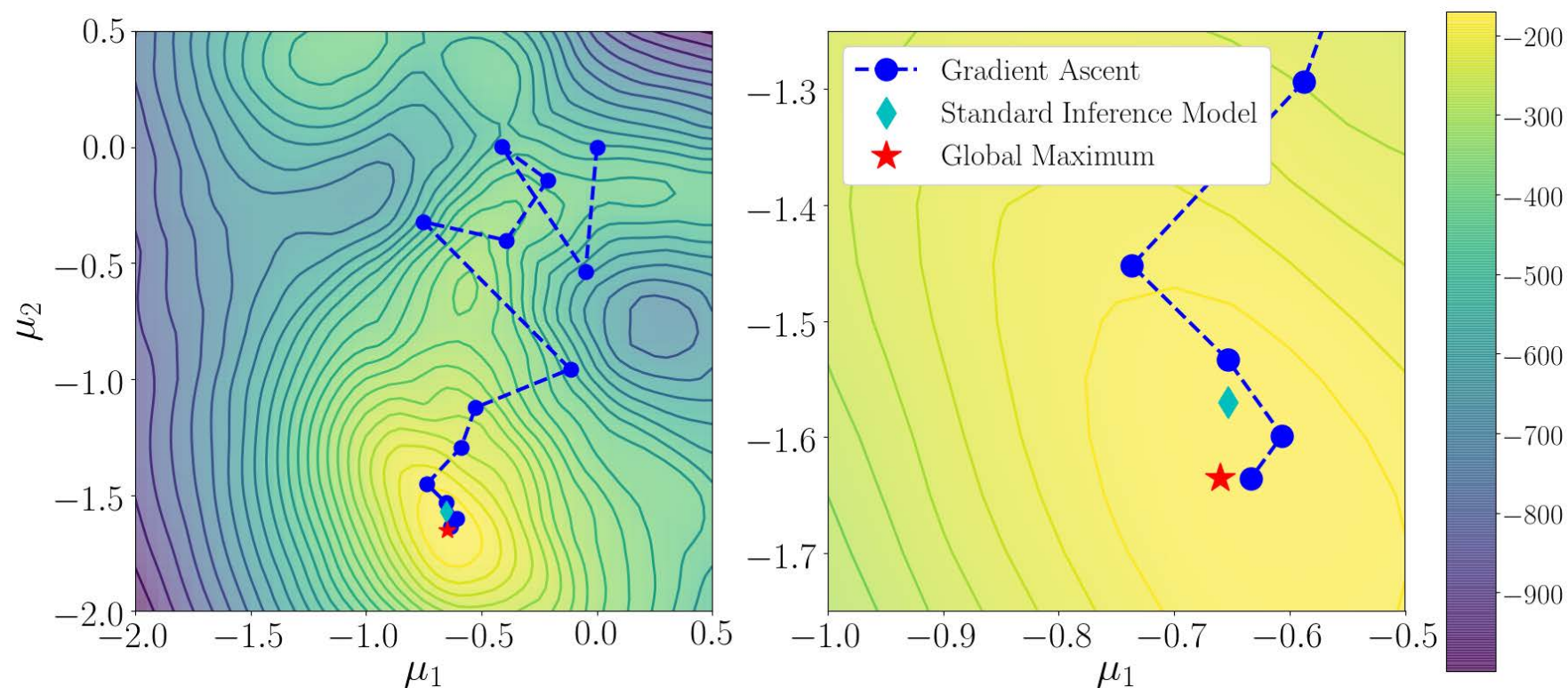
many connections to *compression, information theory*

resembles the “auto-encoding” framework

variational **inference** involves optimizing
the approximate posterior for each data example

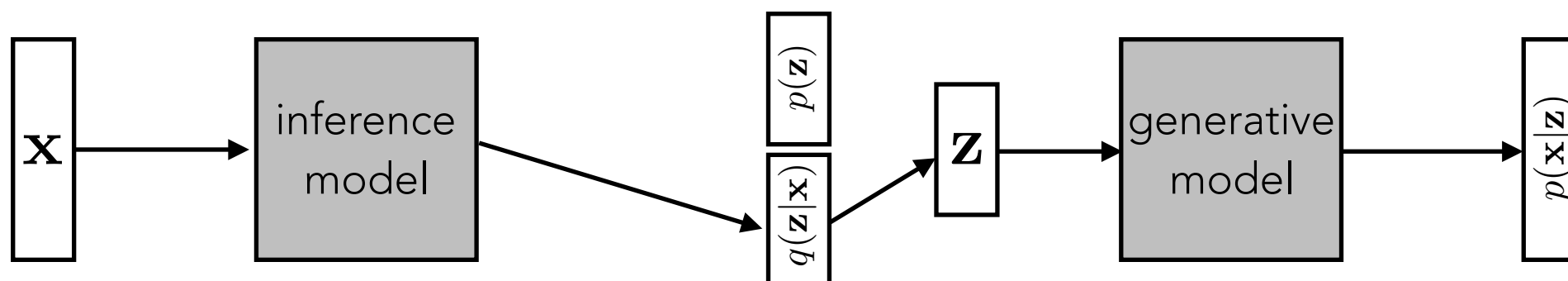
$$q^*(\mathbf{z}|\mathbf{x}) = \operatorname{argmax}_q \mathcal{L}$$

can be solved using gradient ascent and (stochastic) backpropagation / REINFORCE,
but can be computationally expensive

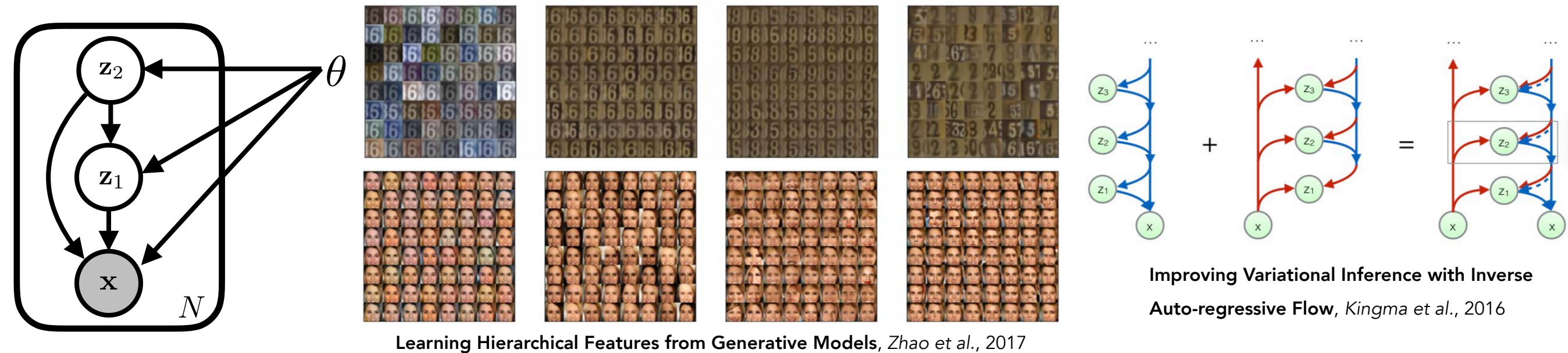


can instead *amortize* inference over data examples by learning
a separate **inference model** to output approximate posterior estimates

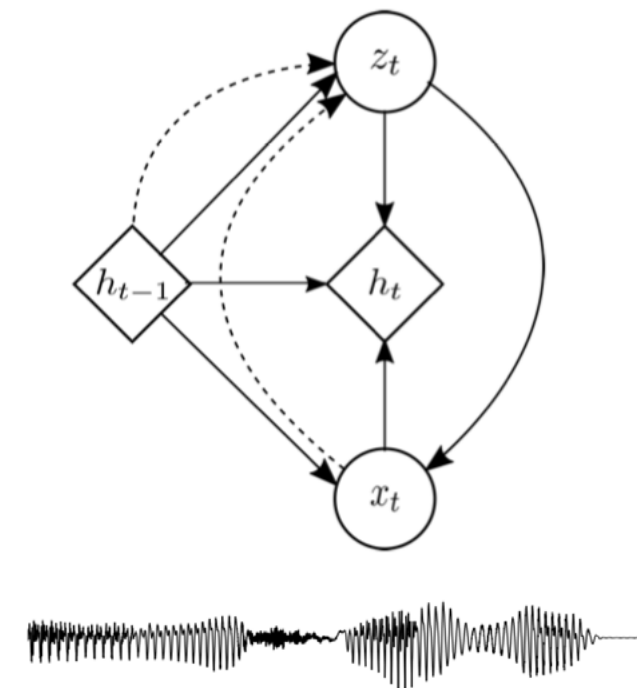
“variational auto-encoder”



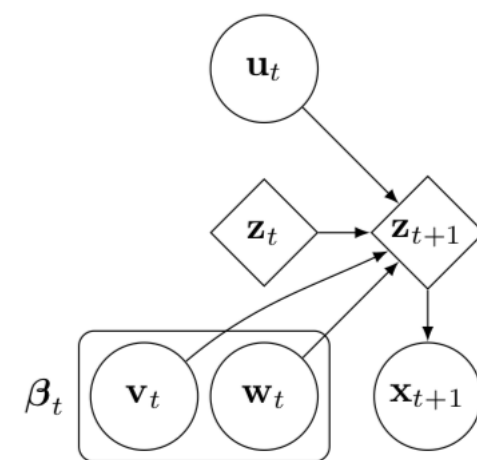
hierarchical latent variable models



sequential latent variable models

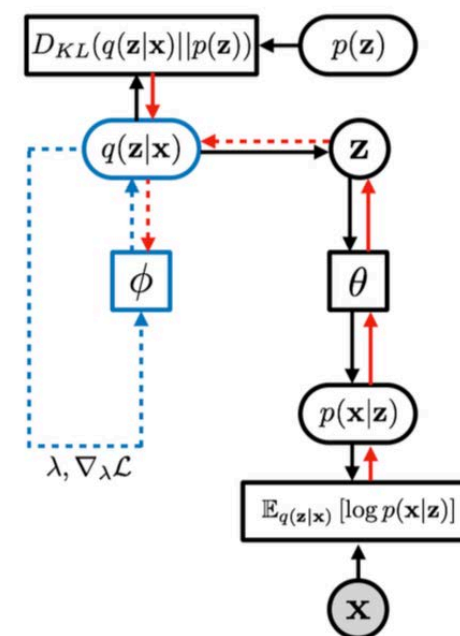


A Recurrent Latent Variable Model for Sequential Data, Chung et al., 2015

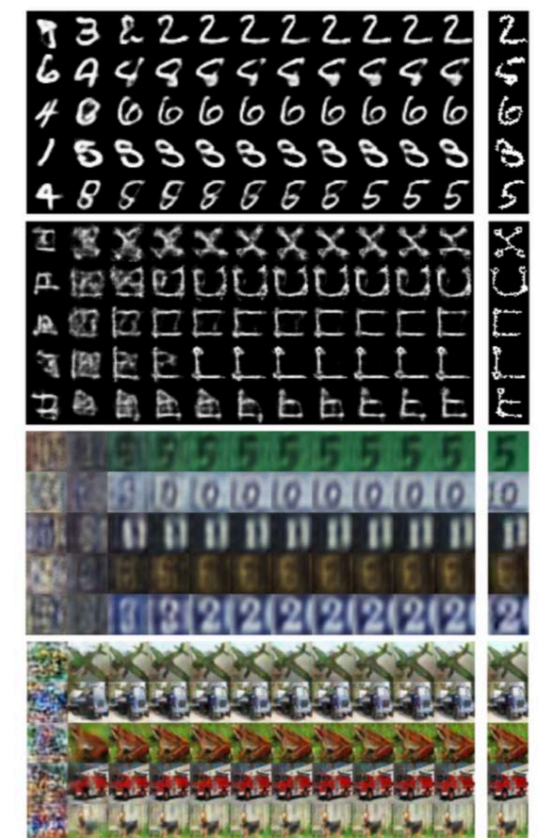


Deep Variational Bayes Filters: Unsupervised Learning of State Space Models from Raw Data, Karl et al., 2016

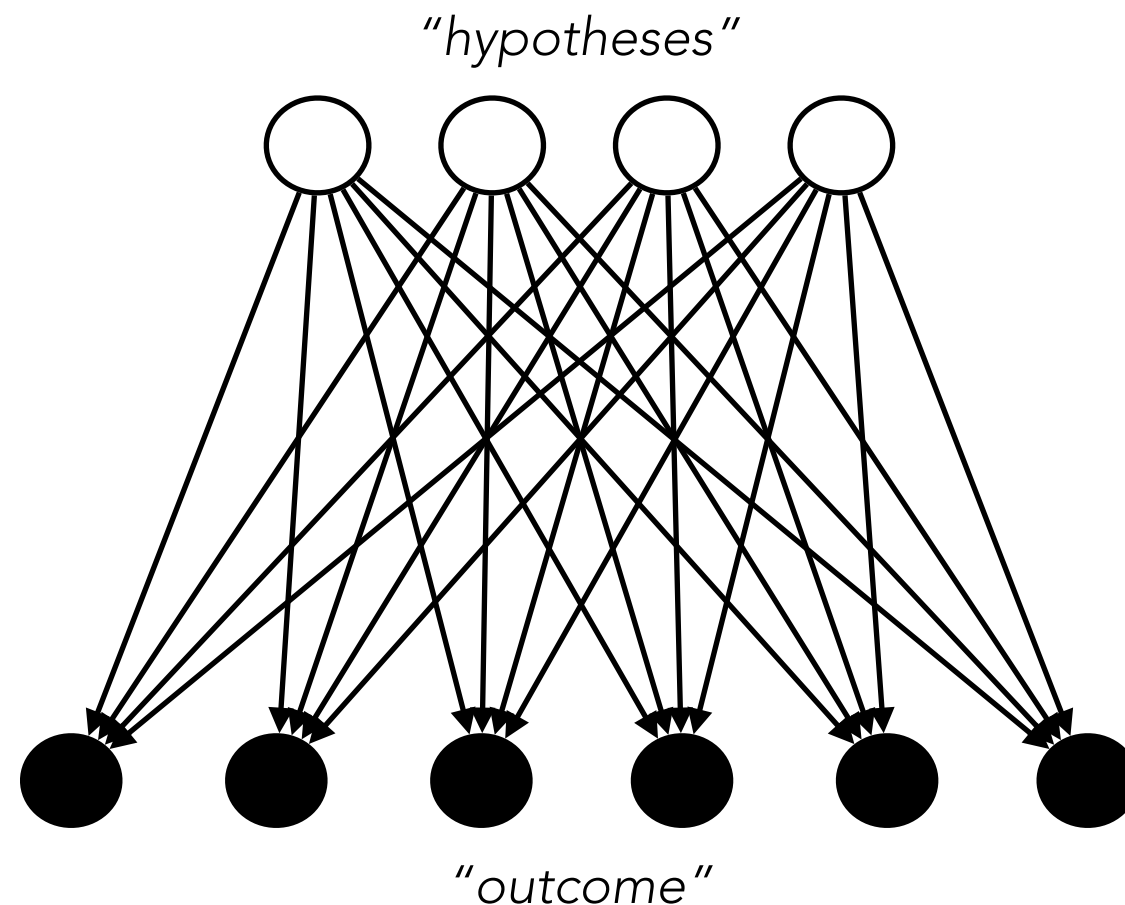
iterative inference models



Learning to Infer, Marino et al., 2017



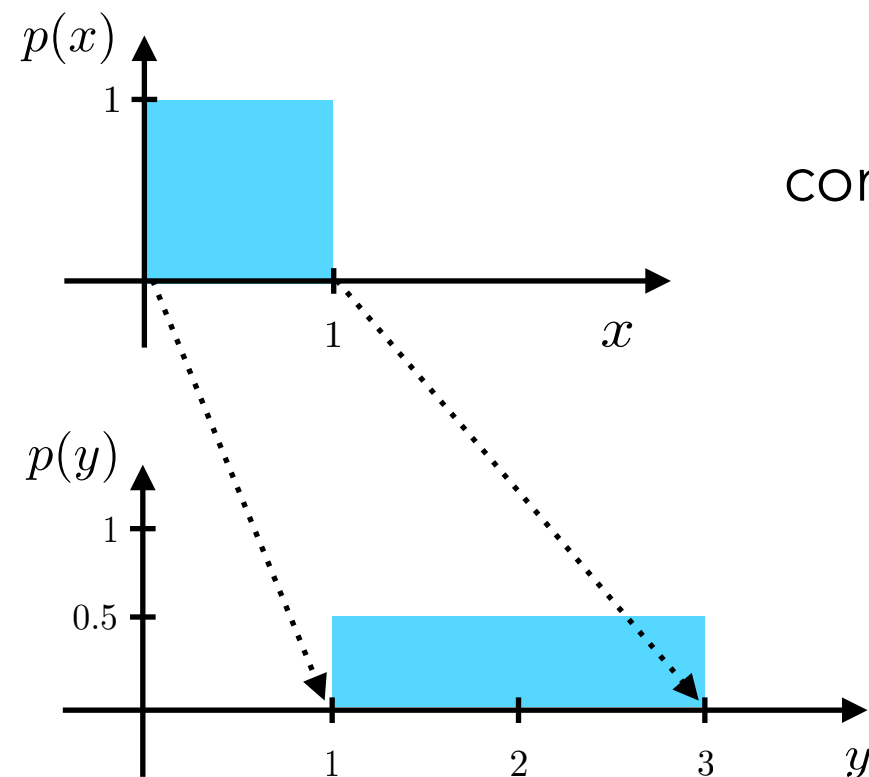
introducing latent variables to a generative model
generally makes evaluating the (log) likelihood intractable



*need to consider all possible "hypotheses"
to evaluate (marginal) likelihood of the "outcome"*

change of variables

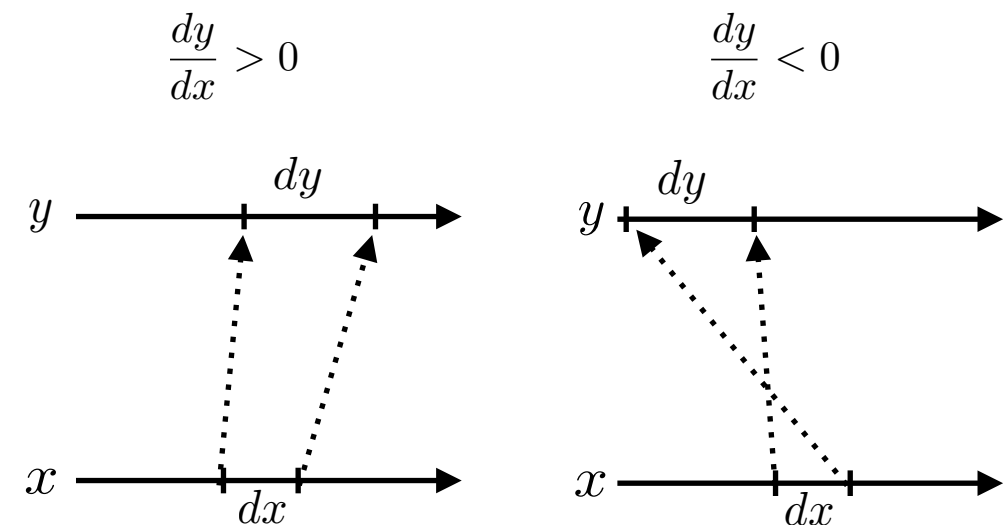
under certain conditions, we can use the change of variables formula to exactly evaluate the log likelihood



consider a variable in one dimension: $x \sim \text{Uniform}(0, 1)$

then let y be an *affine transformation* of x ,
e.g. $y = 2x + 1$.

to conserve probability mass, $p(y) = p(x) \left| \frac{dx}{dy} \right|$



change of variables

in higher dimensions, conservation of probability mass generalizes to

CHANGE OF VARIABLES FORMULA

$$p(\mathbf{y}) = p(\mathbf{x}) \left| \det \frac{d\mathbf{x}}{d\mathbf{y}} \right| = p(\mathbf{x}) |\det \mathbf{J}^{-1}|$$

where \mathbf{J} is the Jacobian matrix of the transformation, $\mathbf{J} = \frac{d\mathbf{y}}{d\mathbf{x}}$

$|\det \mathbf{J}^{-1}|$ expresses the local distortion in volume from the linear transformation

“law of the unconscious statistician” (LOTUS)

can evaluate the probability from one variable's distribution by evaluating the probability of a transformed variable and the volume transformation

for certain classes of transformations, this is tractable to evaluate

change of variables

to use the change of variables formula, we need to evaluate $|\det \mathbf{J}^{-1}|$

for an arbitrary $N \times N$ Jacobian matrix, this is worst case $O(N^3)$

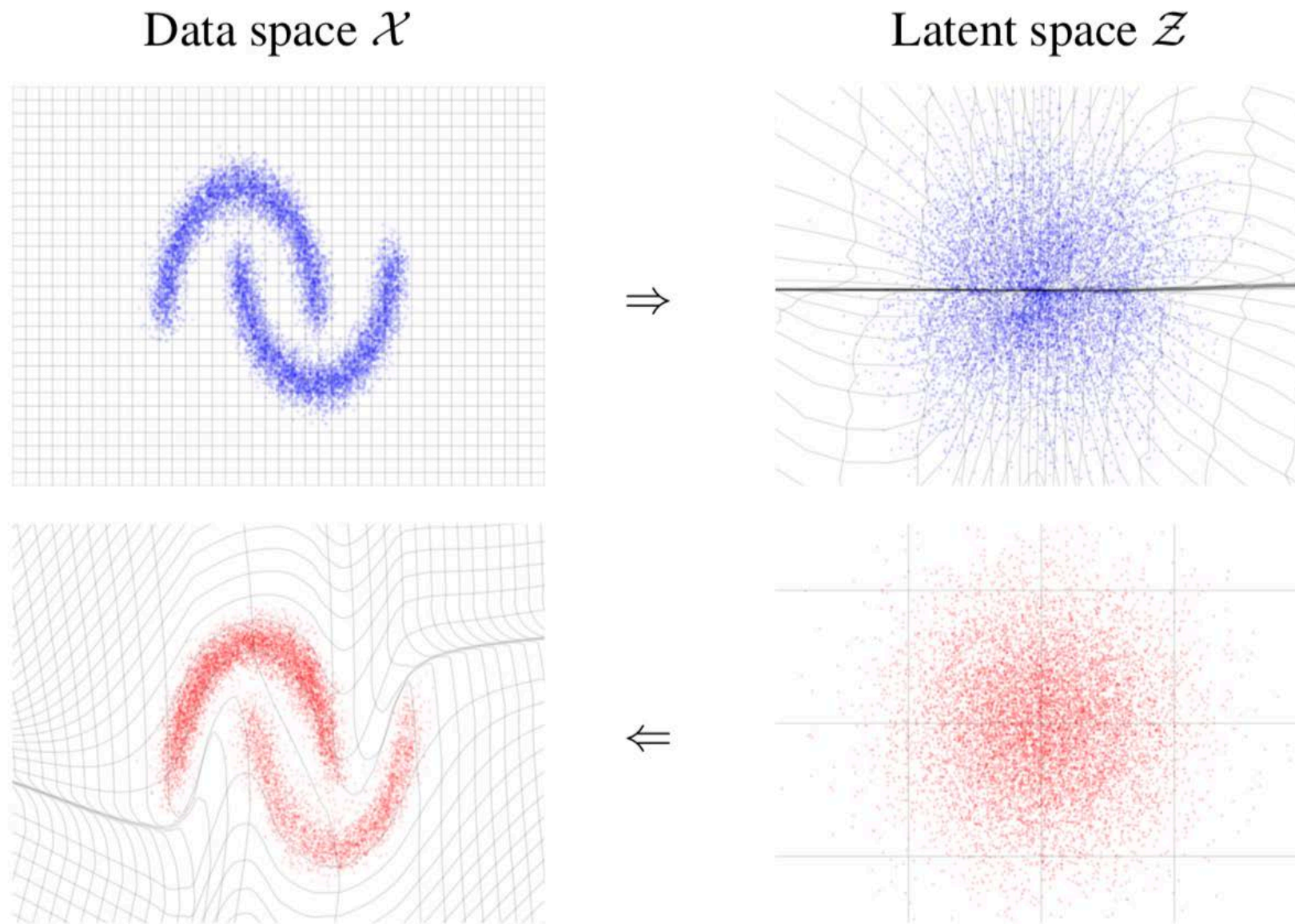


restricting the transformations to those with diagonal or triangular inverse Jacobians allows us to compute $|\det \mathbf{J}^{-1}|$ in $O(N)$.

→ *product of diagonal entries*

change of variables

can transform the data into a space that is easier to model

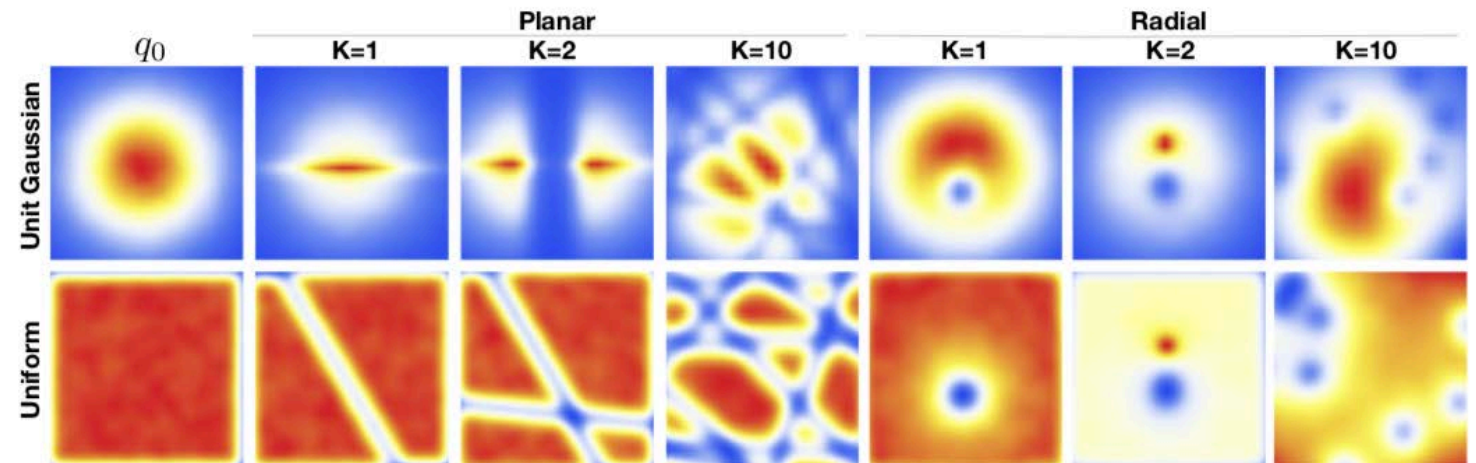
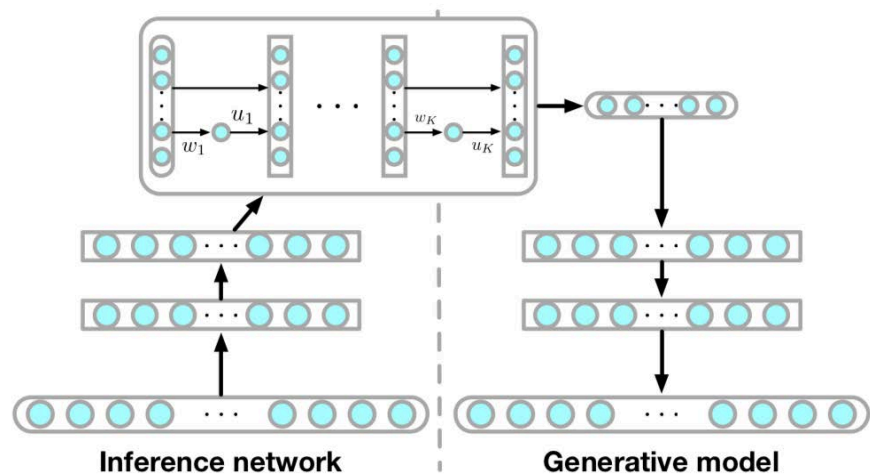


Density Estimation Using Real NVP, *Dinh et al.*, 2016

change of variables for variational inference: **normalizing flows**

use more complex approximate posterior, but evaluate a simpler distribution

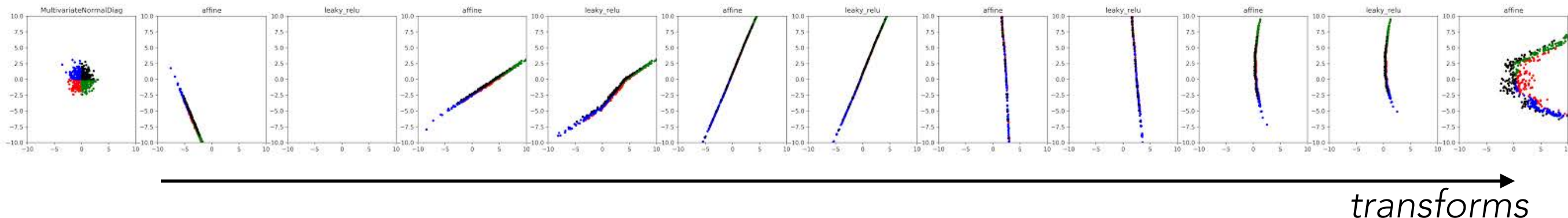
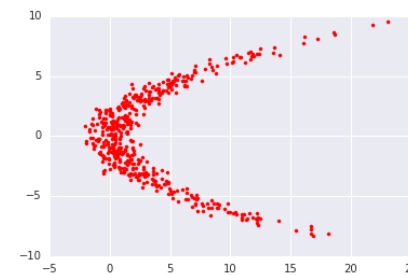
transform $q(\mathbf{z}|\mathbf{x})$



Variational Inference with Normalizing Flows, Rezende & Mohamed, 2015

chain together multiple transforms to get more expressive model

target distribution:



transforms

additive coupling layer Dinh et al., 2014

$$\begin{aligned}\mathbf{y}_{1:d} &= \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} &= \mathbf{x}_{d+1:D} + f(\mathbf{x}_{1:d})\end{aligned}$$

planar flow Rezende & Mohamed, 2015

$$\mathbf{y} = \mathbf{x} + f(\mathbf{x}) \odot g(h(\mathbf{x})^\top \mathbf{x} + b(\mathbf{x}))$$

affine coupling layer Dinh et al., 2016

$$\begin{aligned}\mathbf{y}_{1:d} &= \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} &= \mathbf{x}_{d+1:D} \odot \exp(f(\mathbf{x}_{1:d})) + g(\mathbf{x}_{1:d})\end{aligned}$$

inverse auto-regressive flow (IAF) Kingma et al., 2016

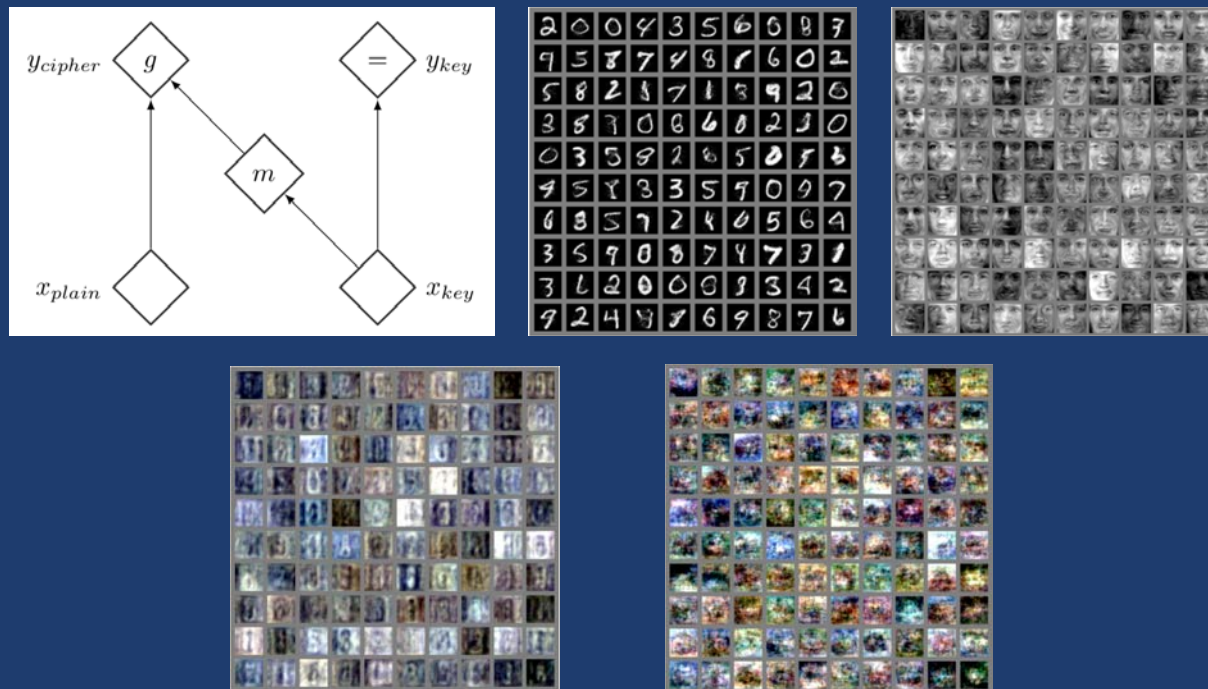
$$\mathbf{y} = \frac{\mathbf{x} - f(\mathbf{x})}{\exp(g(\mathbf{x}))}$$

masked auto-regressive flow (MAF) Papamakarios et al., 2017

$$\mathbf{y} = \mathbf{x} \odot \exp(g(\mathbf{x})) + f(\mathbf{x})$$

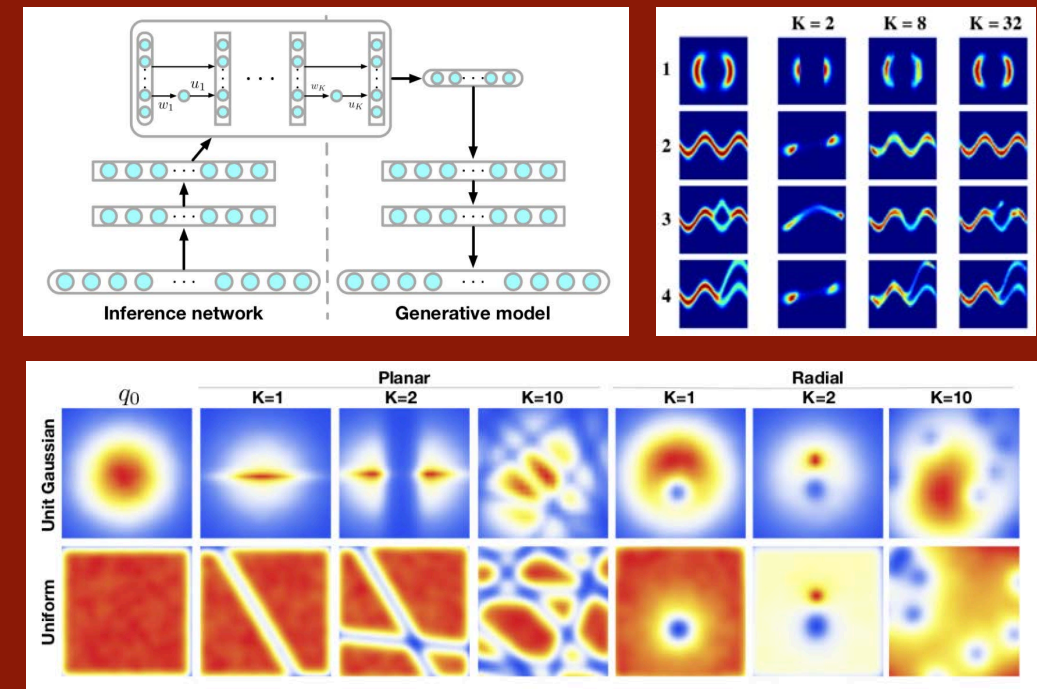
NICE: Non-linear Independent Components Estimation

Dinh et al., 2014



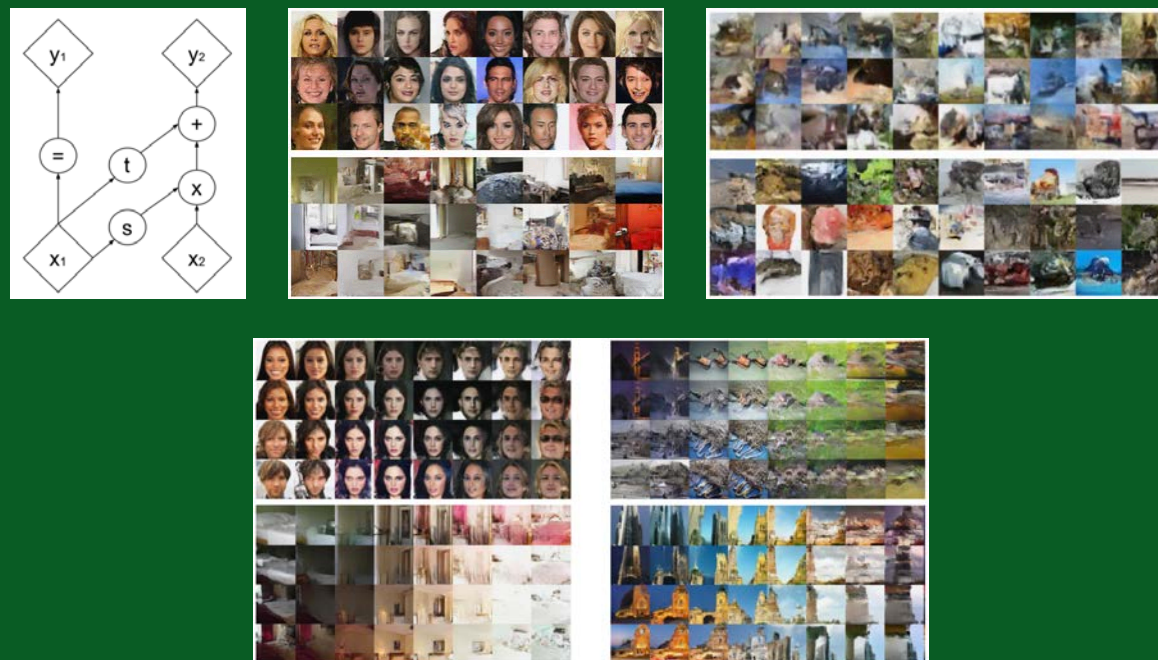
Variational Inference with Normalizing Flows

Rezende & Mohamed, 2015



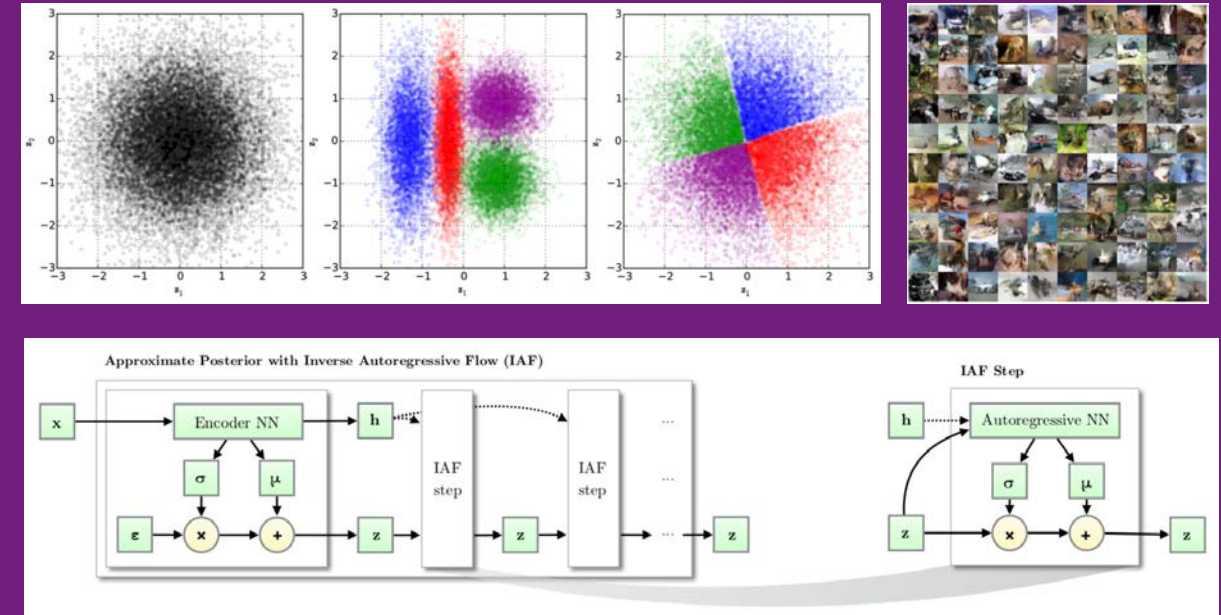
Density Estimation Using Real NVP

Dinh et al., 2016

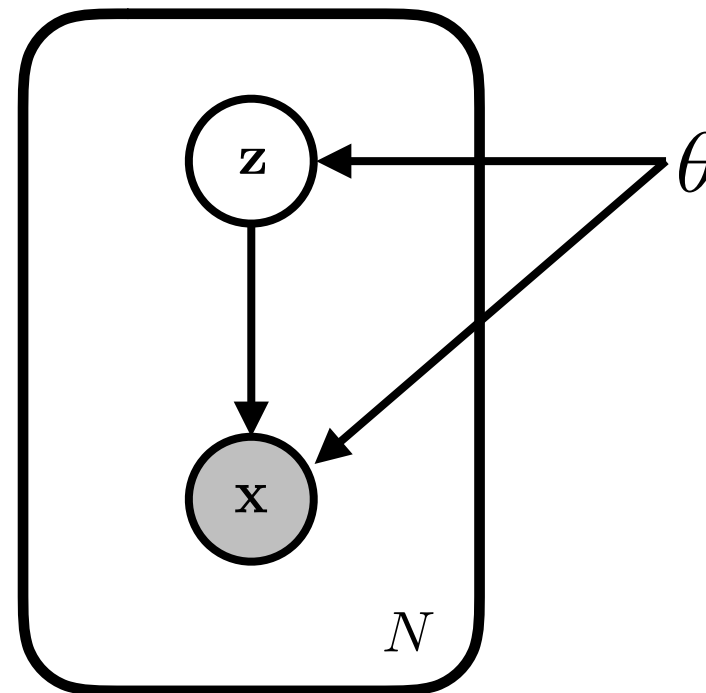


Improving Variational Inference with Inverse Autoregressive Flow

Kingma et al., 2016



recap: explicit latent variable models



model the data through *latent variables*

Pros

can capture abstract variables,
good for semi supervised learning

relatively fast sampling / training

theoretical foundations
from info. theory

Cons

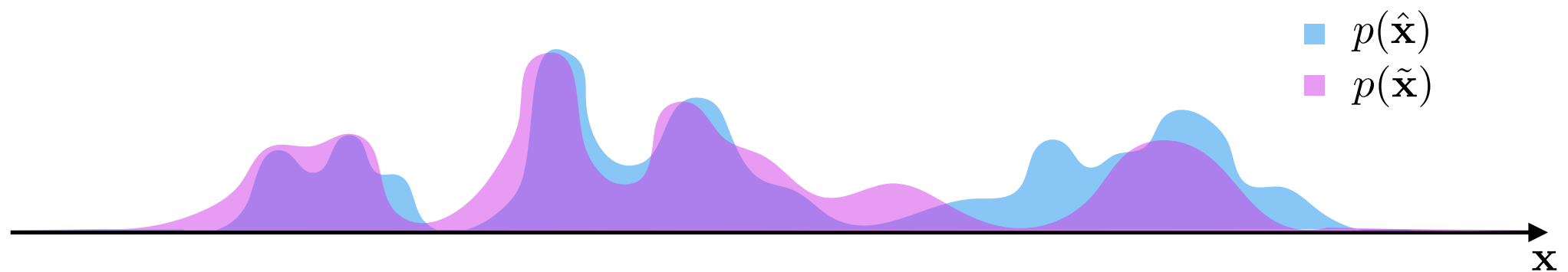
likelihood evaluation / inference
often intractable

requires additional assumptions
on latent variables

difficult to capture details

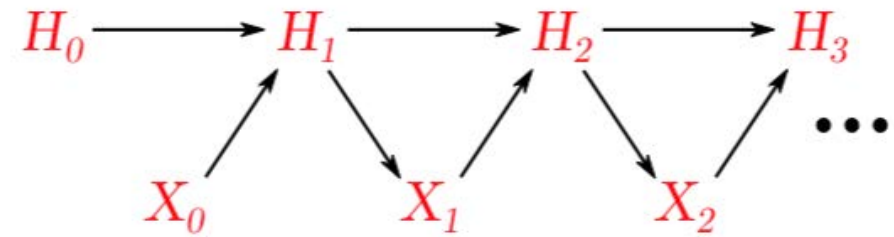
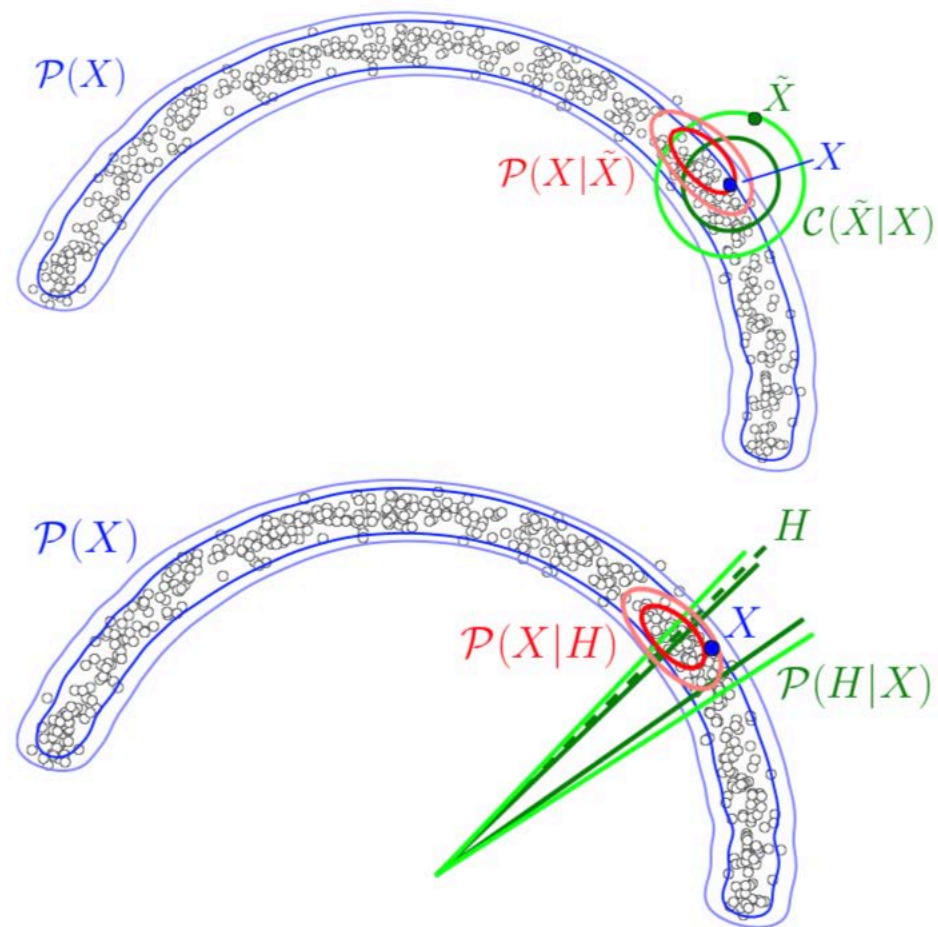
IMPLICIT LATENT VARIABLE MODELS

instead of using an *explicit* probability density,
learn a model that defines an *implicit density*



specify a stochastic procedure for generating the data
that does not require an explicit likelihood evaluation

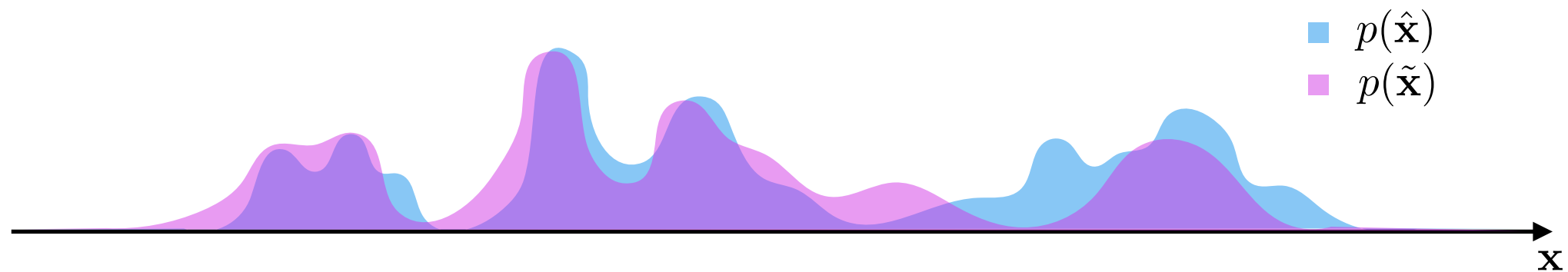
Generative Stochastic Networks (GSNs)



Deep Generative Stochastic Networks Trainable by Backprop, Bengio et al., 2013

train an auto-encoder to learn Monte Carlo sampling transitions

the generative distribution is *implicitly* defined by this transition



estimate density ratio through *Bayesian two-sample test*

data distribution $p(\hat{\mathbf{x}})$

generated distribution $p(\tilde{\mathbf{x}})$

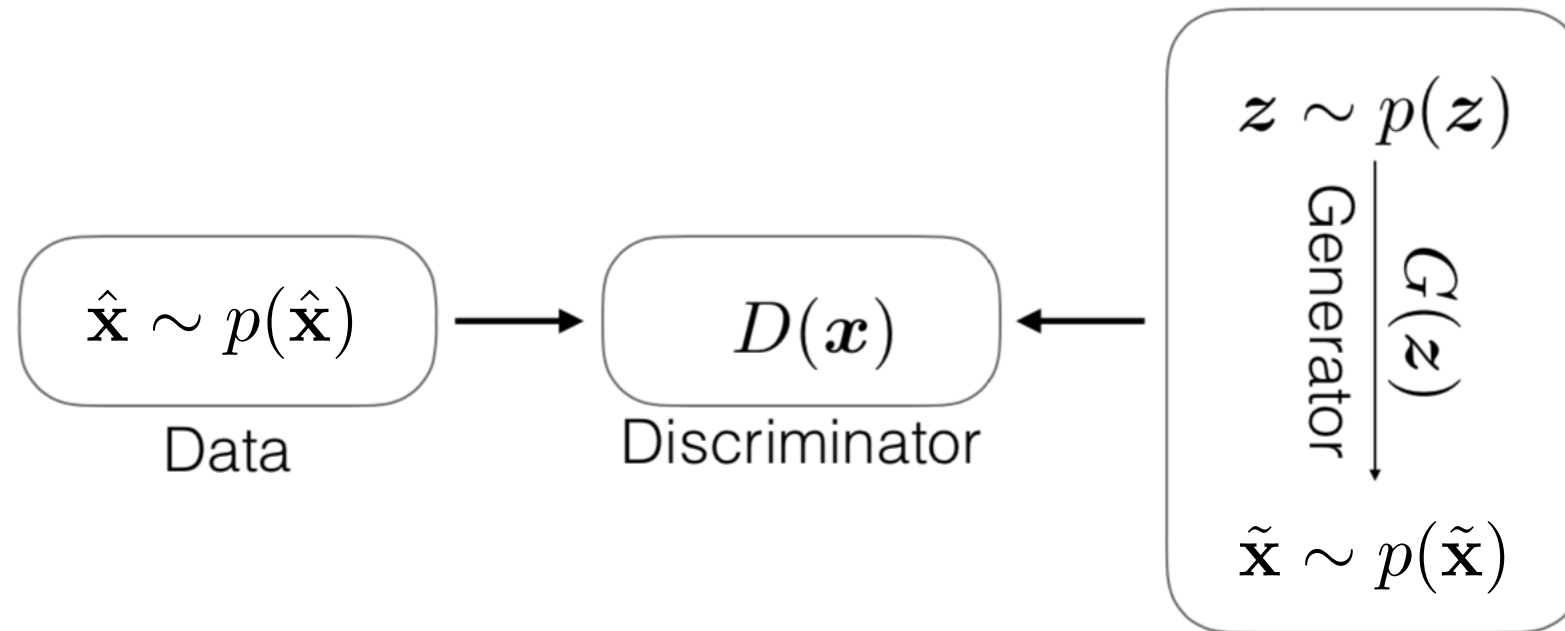
$$\frac{p(\hat{\mathbf{x}})}{p(\tilde{\mathbf{x}})} = \frac{p(\mathbf{x}|\text{data})}{p(\mathbf{x}|\text{gen.})}$$

$$\frac{p(\hat{\mathbf{x}})}{p(\tilde{\mathbf{x}})} = \frac{p(\text{data}|\mathbf{x})p(\mathbf{x})/p(\text{data})}{p(\text{gen.}|\mathbf{x})p(\mathbf{x})/p(\text{gen.})} \quad (\text{Bayes' rule})$$

$$\frac{p(\hat{\mathbf{x}})}{p(\tilde{\mathbf{x}})} = \frac{p(\text{data}|\mathbf{x})}{p(\text{gen.}|\mathbf{x})} \quad (\text{assuming equal dist. prob.})$$

density estimation becomes a sample discrimination task

Generative Adversarial Networks (GANs)



learn the discriminator:

$$p(\text{data}|\mathbf{x}) = D(\mathbf{x}) \qquad p(\text{gen.}|\mathbf{x}) = 1 - D(\mathbf{x})$$

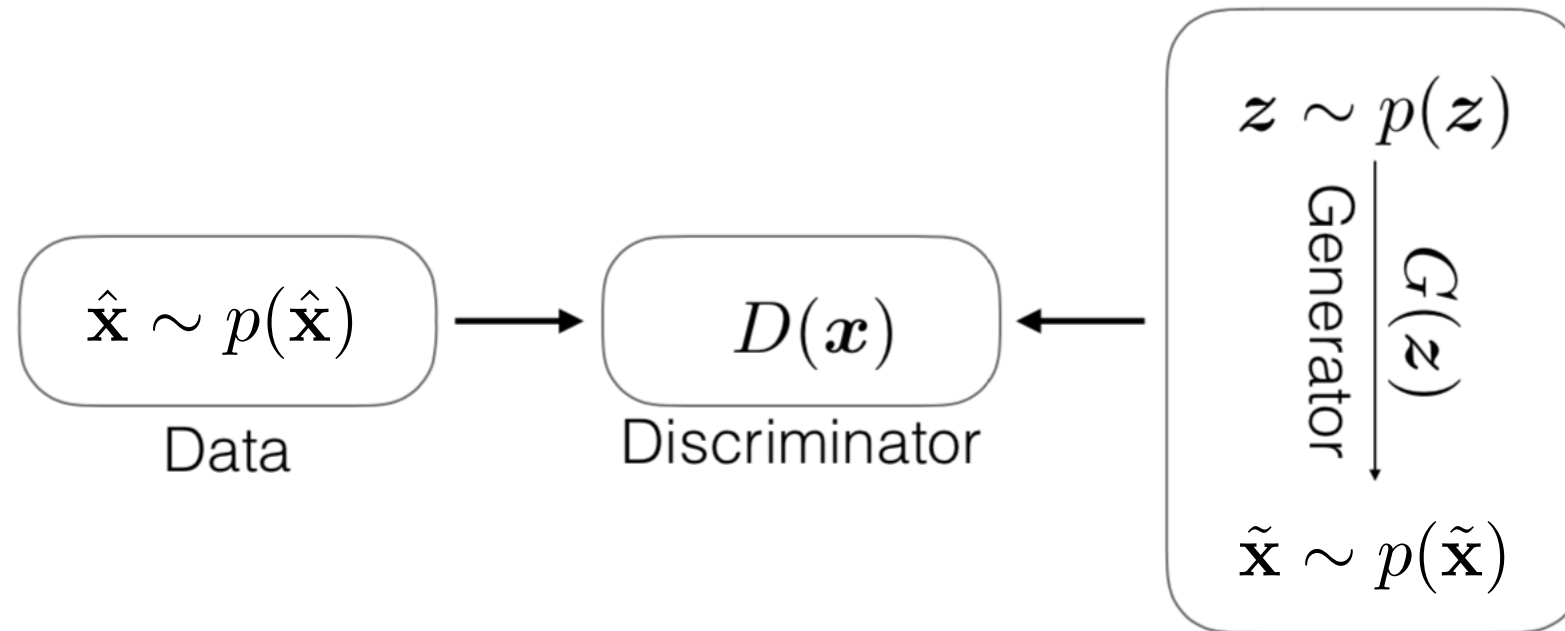
Bernoulli outcome: $y \in \{\text{data}, \text{gen.}\}$

$$\log p(y|\mathbf{x}) = \log D(\hat{\mathbf{x}}) + \log(1 - D(\tilde{\mathbf{x}}))$$

two-sample criterion:

$$\min_G \max_D \mathbb{E}_{p(\hat{\mathbf{x}})} [\log D(\hat{\mathbf{x}})] + \mathbb{E}_{p(\tilde{\mathbf{x}})} [\log(1 - D(\tilde{\mathbf{x}}))]$$

Generative Adversarial Networks (GANs)



two-sample criterion:

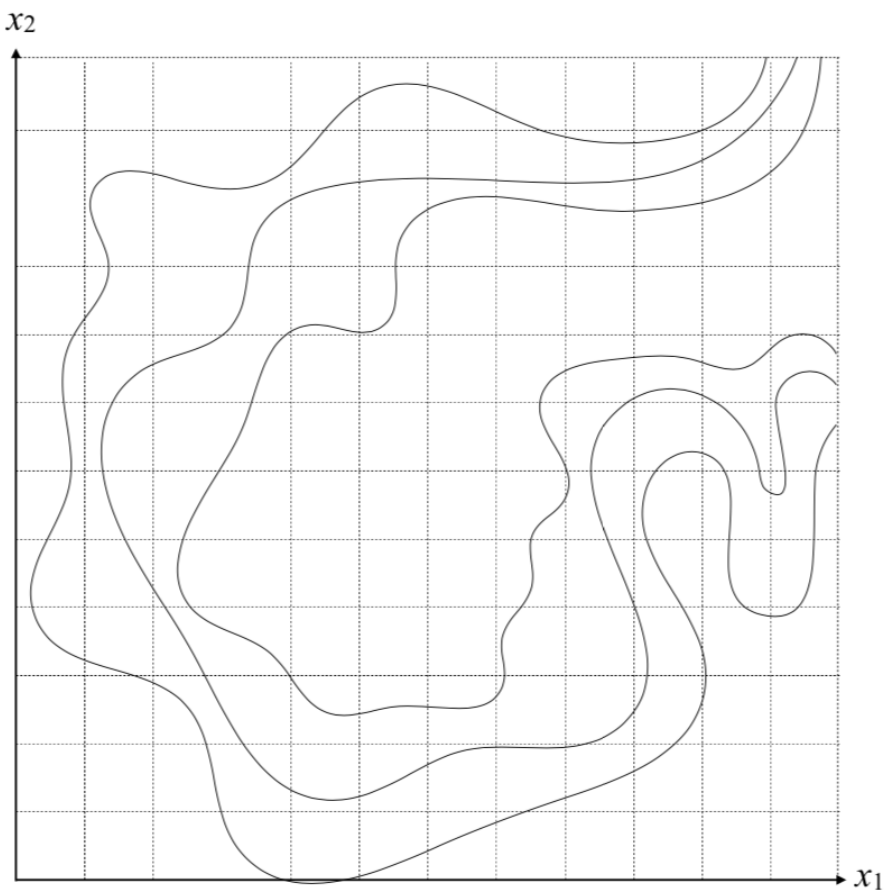
$$\min_G \max_D \mathbb{E}_{p(\hat{\mathbf{x}})} [\log D(\hat{\mathbf{x}})] + \mathbb{E}_{p(\tilde{\mathbf{x}})} [\log(1 - D(\tilde{\mathbf{x}}))]$$

in practice:

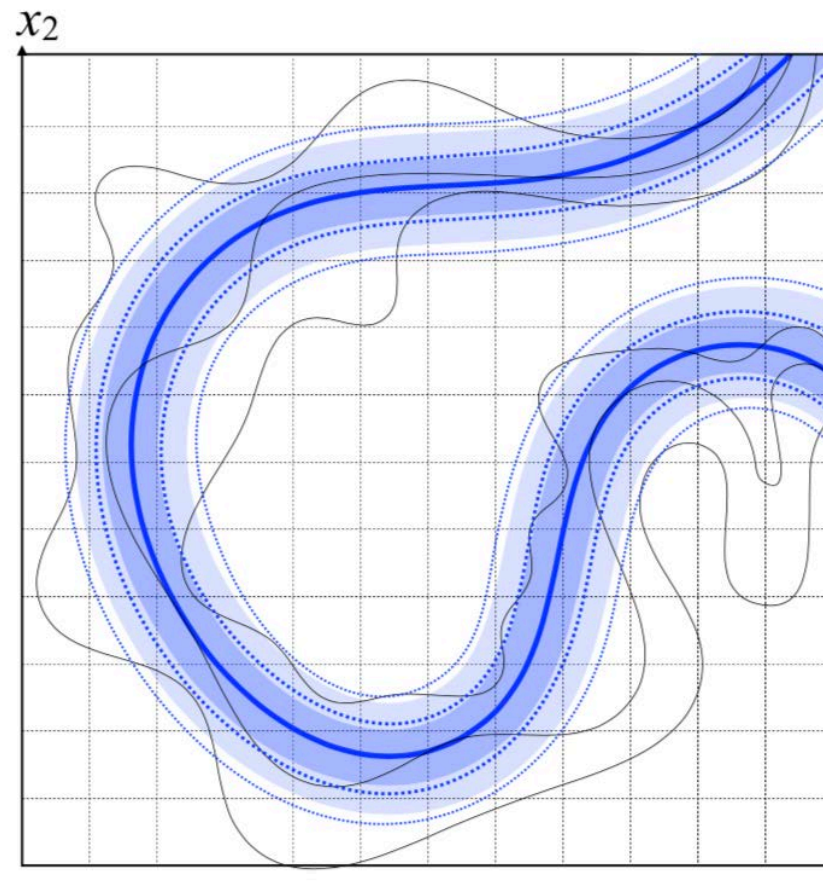
$$\max_D \mathbb{E}_{p(\hat{\mathbf{x}})} [\log D(\hat{\mathbf{x}})] + \mathbb{E}_{p(\tilde{\mathbf{x}})} [\log(1 - D(\tilde{\mathbf{x}}))]$$

$$\max_G \mathbb{E}_{p(\tilde{\mathbf{x}})} [\log D(\tilde{\mathbf{x}})]$$

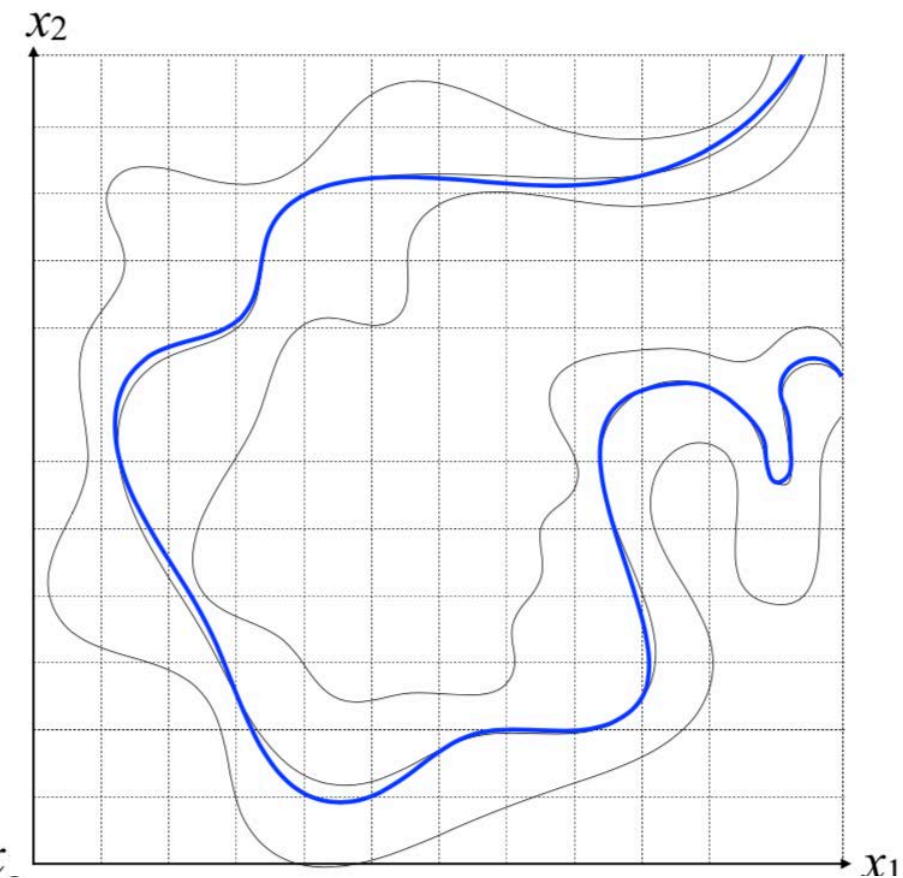
interpretation



data manifold



explicit model



implicit model

explicit models tend to cover the entire data manifold, but are constrained

implicit models tend to capture part of the data manifold, but can neglect other parts

→ "mode collapse"

Generative Adversarial Networks (GANs)

GANs can be difficult to optimize

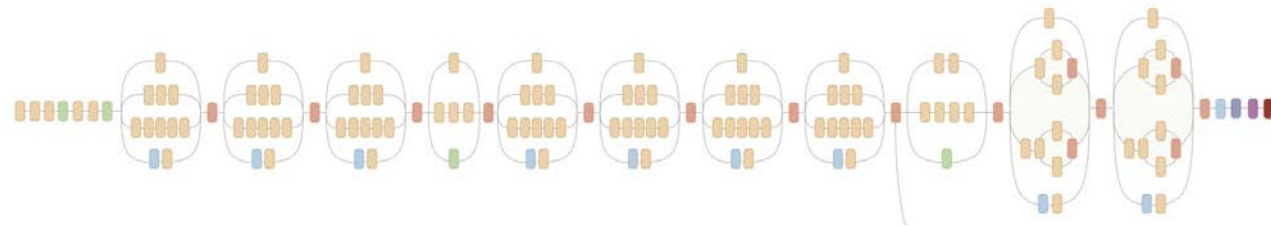
DCGAN	LSGAN	WGAN (clipping)	WGAN-GP (ours)
Baseline (G : DCGAN, D : DCGAN)			
			
G : No BN and a constant number of filters, D : DCGAN			
			
G : 4-layer 512-dim ReLU MLP, D : DCGAN			
			
No normalization in either G or D			
			
Gated multiplicative nonlinearities everywhere in G and D			
			
tanh nonlinearities everywhere in G and D			
			
101-layer ResNet G and D			
			

Improved Training of Wasserstein GANs, Gulrajani et al., 2017

evaluation

without an explicit likelihood, it is difficult to quantify the performance

inception score



use a pre-trained Inception v3 model to quantify class and distribution entropy

$$\text{IS}(G) = \exp \left(\mathbb{E}_{p(\tilde{\mathbf{x}})} D_{KL}(p(y|\tilde{\mathbf{x}}) || p(y)) \right)$$

$p(y|\tilde{\mathbf{x}})$ is the class distribution for a given image

→ should be highly peaked (low entropy)

$p(y) = \int p(y|\tilde{\mathbf{x}}) d\tilde{\mathbf{x}}$ is the marginal class distribution

→ want this to be uniform (high entropy)

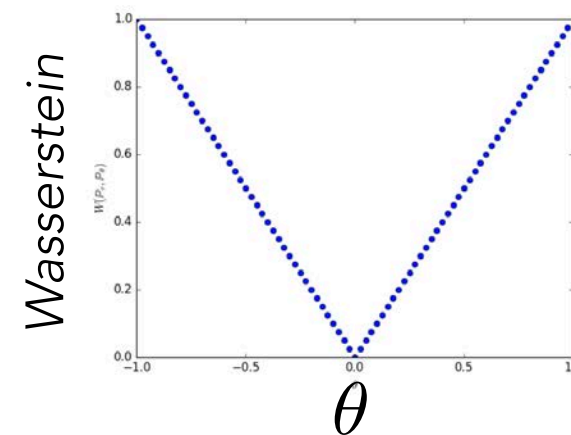
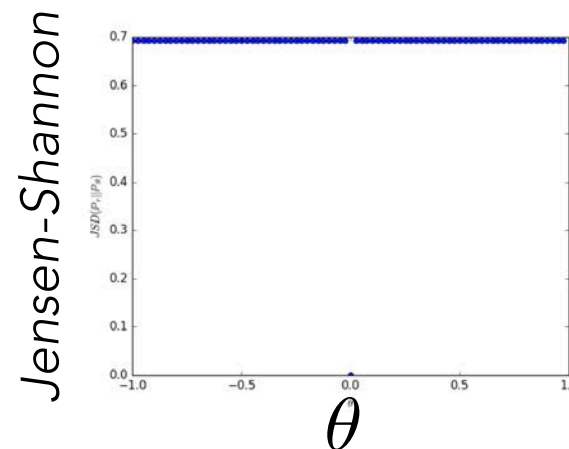
extensions: Wasserstein GAN

under an “ideal” discriminator, the generator minimizes the Jensen-Shannon divergence

$$D_{JS}(p(\hat{\mathbf{x}})||p(\tilde{\mathbf{x}})) = \frac{1}{2}D_{KL}(p(\hat{\mathbf{x}})||\frac{1}{2}(p(\hat{\mathbf{x}}) + p(\tilde{\mathbf{x}}))) + \frac{1}{2}D_{KL}(p(\tilde{\mathbf{x}})||\frac{1}{2}(p(\hat{\mathbf{x}}) + p(\tilde{\mathbf{x}})))$$

however, this metric can be **discontinuous**, making it difficult to train

θ is a gen. model parameter



can instead use the Wasserstein (Earth Mover's) distance (which is continuous and diff. almost everywhere):

$$W(p(\hat{\mathbf{x}}), p(\tilde{\mathbf{x}})) = \inf_{\gamma \in \Pi(p(\hat{\mathbf{x}}), p(\tilde{\mathbf{x}}))} \mathbb{E}_{(\hat{\mathbf{x}}, \tilde{\mathbf{x}}) \sim \gamma} [||\hat{\mathbf{x}} - \tilde{\mathbf{x}}||]$$

think of it as the “minimum cost of transporting points between two distributions”

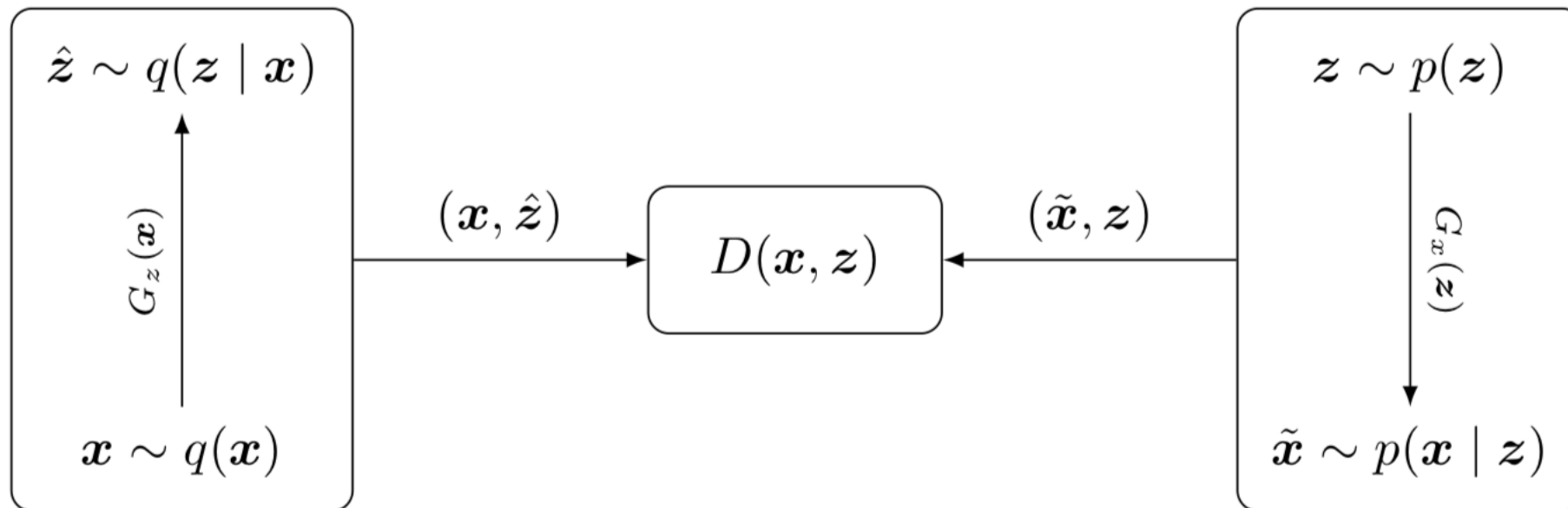
intractable to actually evaluate Wasserstein distance, but by constraining the discriminator, can evaluate

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{p(\hat{\mathbf{x}})} [D(\hat{\mathbf{x}})] - \mathbb{E}_{p(\tilde{\mathbf{x}})} [D(\tilde{\mathbf{x}})]$$

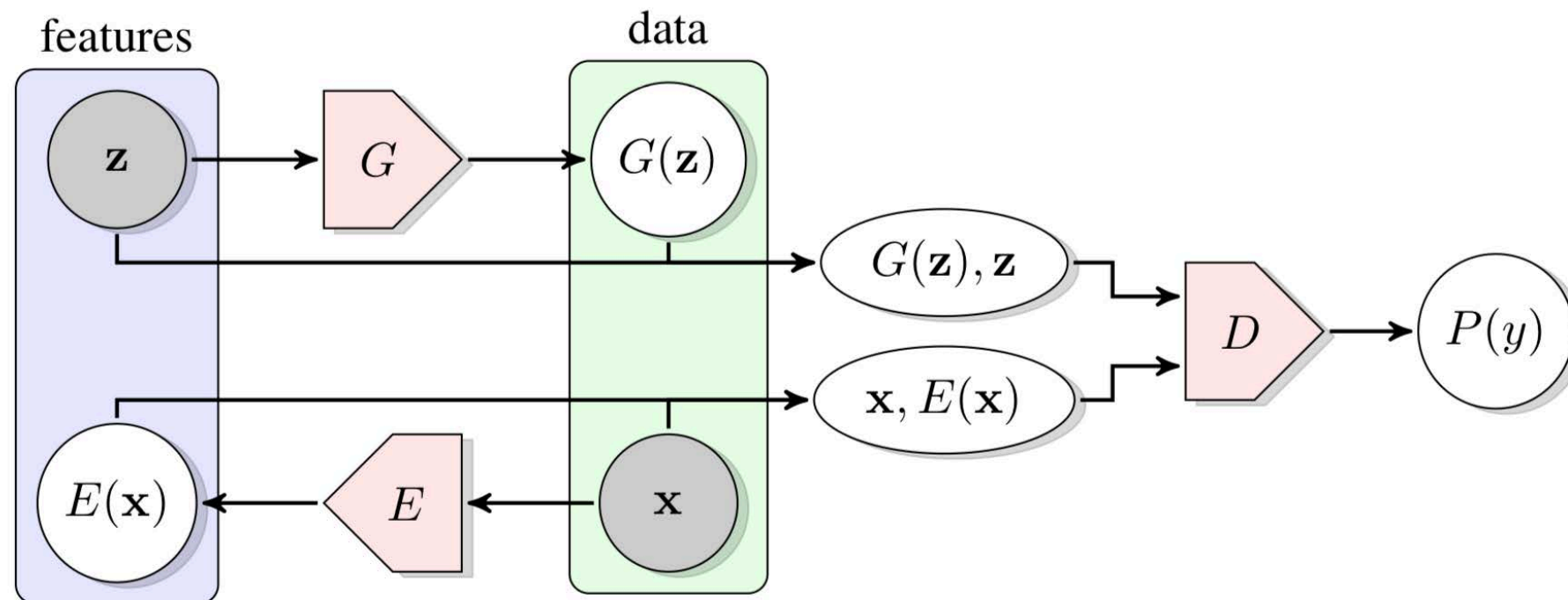
\mathcal{D} is the set of Lipschitz functions, which can be enforced through weight clipping or gradient penalty

extensions: inference

can we also learn to **infer a latent representation**?



Adversarially Learned Inference, Dumoulin et al., 2017



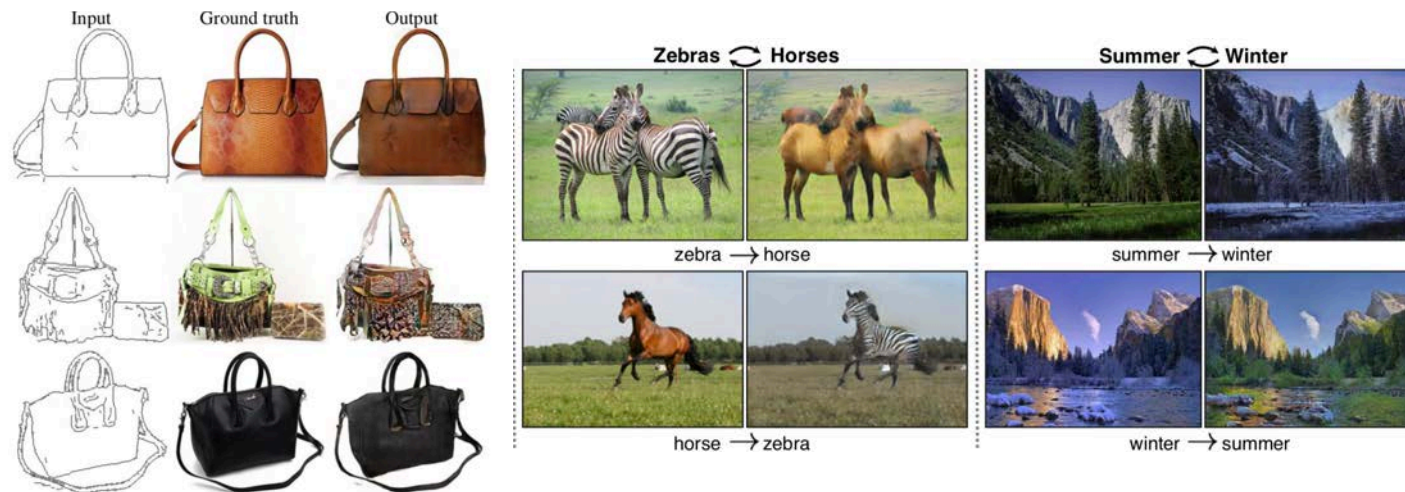
Adversarial Feature Learning, Donahue et al., 2017

applications

image to image translation

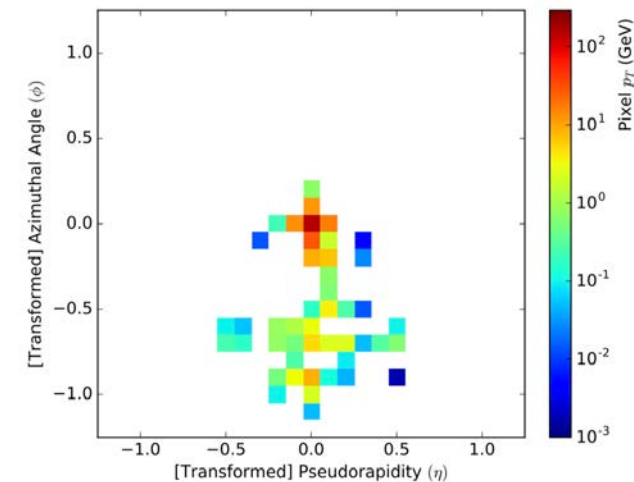


Image-to-Image Translation with Conditional Adversarial Networks, *Isola et al.*, 2016



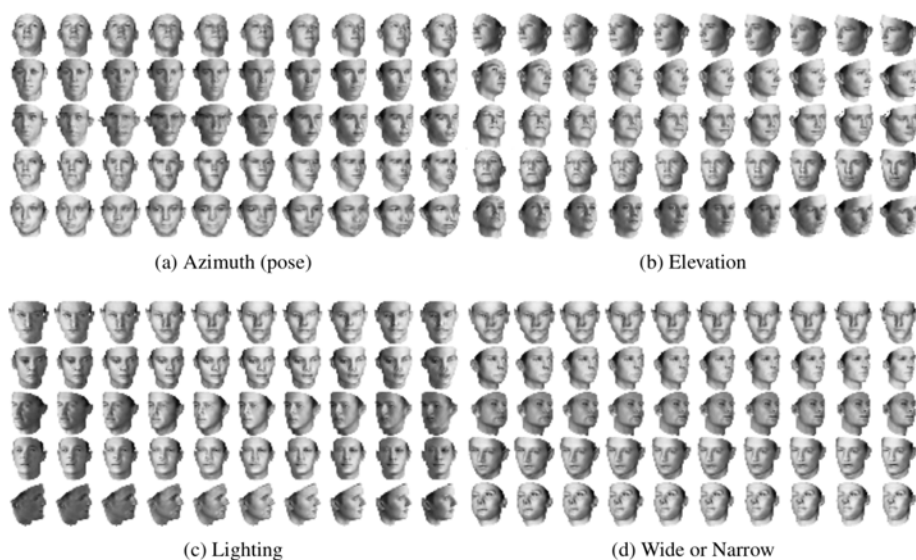
Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, *Zhu et al.*, 2017

experimental simulation



Learning Particle Physics by Example, *de Oliveira et al.*, 2017

interpretable representations



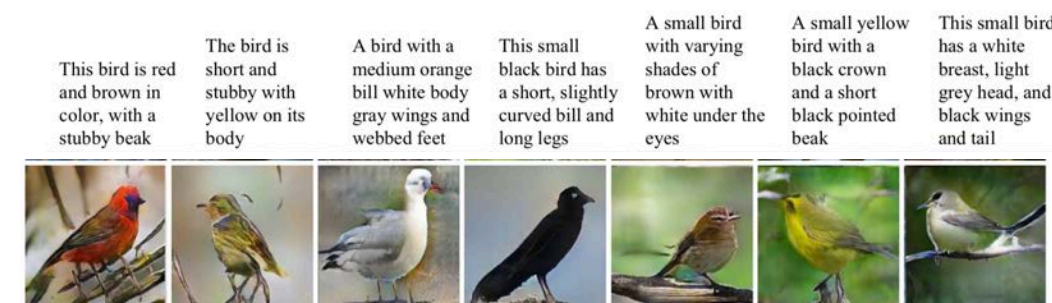
InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets, *Chen et al.*, 2016

music synthesis



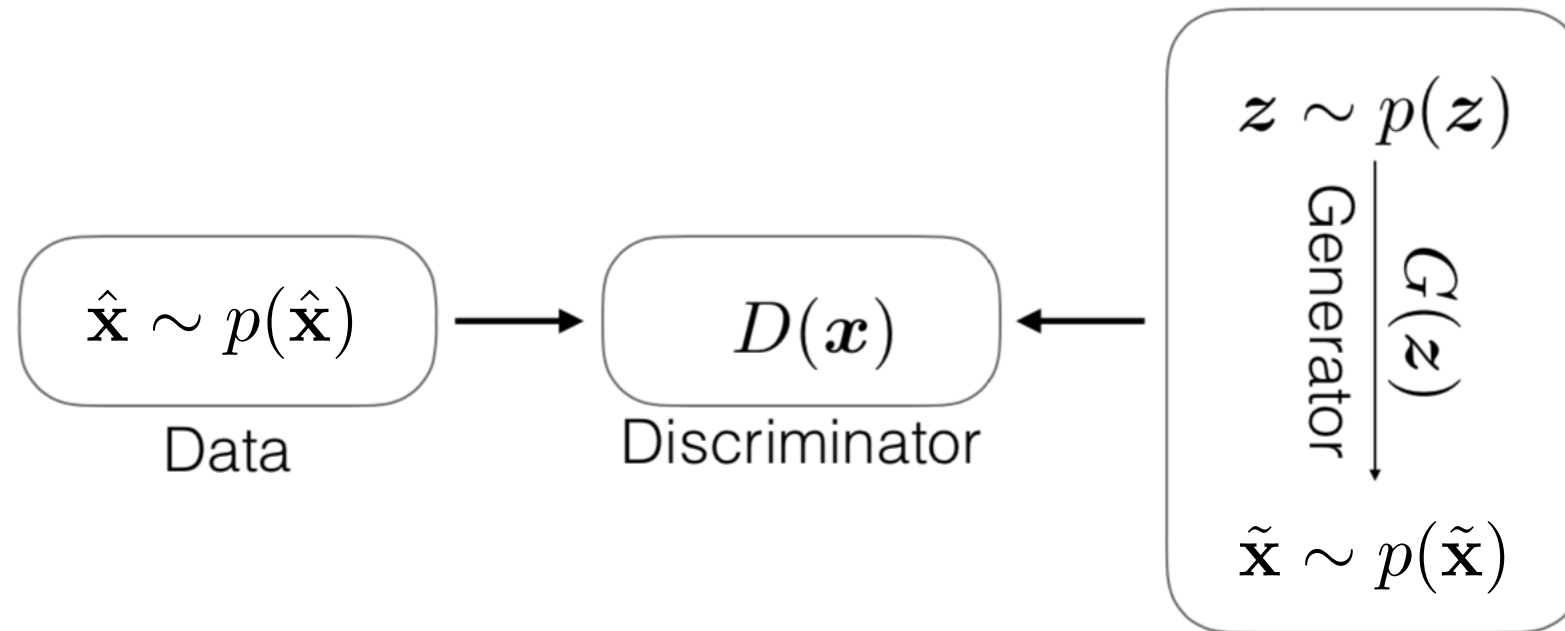
MIDINET: A CONVOLUTIONAL GENERATIVE ADVERSARIAL NETWORK FOR SYMBOLIC-DOMAIN MUSIC GENERATION, *Yang et al.*, 2017

text to image synthesis



StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks, *Zhang et al.*, 2016

recap: implicit latent variable models



Pros

able to learn flexible models

requires fewer modeling
assumptions

capable of learning latent
representation

Cons

difficult to evaluate

sensitive, difficult to optimize

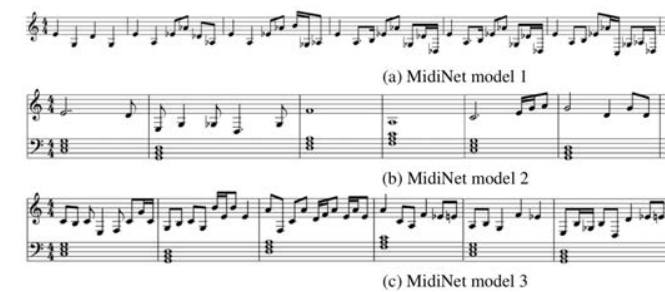
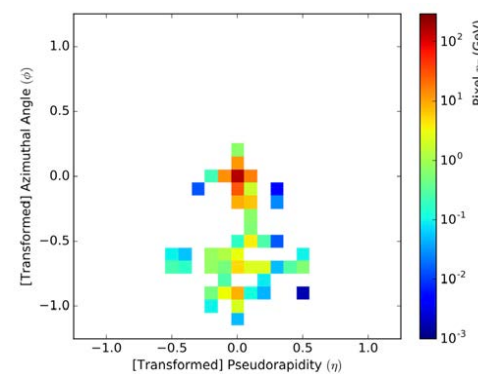
can be difficult to
incorporate model assumptions

DISCUSSION

generative models: *what are they good for?*

generative models model the data distribution

1. can generate and simulate data



2. can extract structure from data



generative models: *what's next?*

applying generative models to new forms of data



incorporating generative models into *complementary* learning systems



