

Deep Generative Modeling with Variational Auto-Encoders

Jakub M. Tomczak

What is **intelligence**?

What is **intelligence**?

...

What is **intelligence**?

...

What is **artificial intelligence**?

INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?



INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?

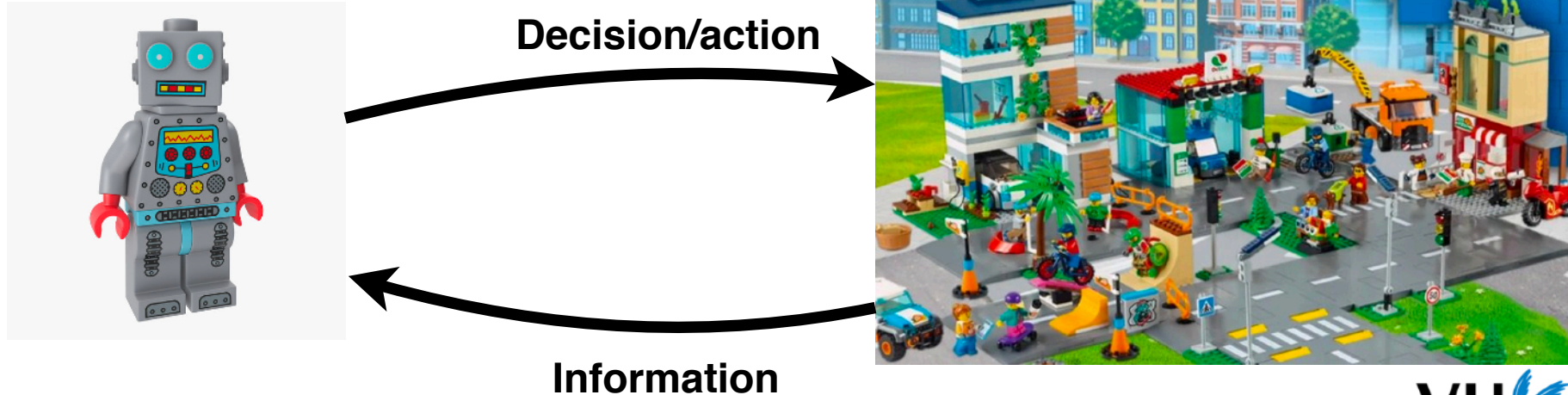


INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?



INFORMATION, INTELLIGENCE AND ARTIFICIAL INTELLIGENCE

What is **intelligence**?

...

What is **artificial intelligence**?



Decision/action



Information

0001101010011...

What is **artificial intelligence**?

- **Information** processing
- **Information** storing
- **Information** transmission



What is **artificial intelligence**?

- **Information** processing
- **Information** storing
- **Information** transmission
- **Decision** making



What is **artificial intelligence**?

- **Information** processing
- **Information** storing
- **Information** transmission
- **Decision** making



Learning
Knowledge representation
Models...



What is **artificial intelligence**?

- **Information** processing
- **Information** storing
- **Information** transmission
- **Decision** making



Learning
Knowledge representation
Models...



The question is how to formalize the problem of AI?

Information (a quick recap)



Claude Shannon

Information (a quick recap)

We have a random source of data x .



Claude Shannon

Information (a quick recap)

We have a random source of data x .

We can quantify the **uncertainty** of this source by calculating **the entropy**:

$$\mathbb{H}[x] = - \sum_x p(x) \log p(x)$$



Claude Shannon

Information (a quick recap)

We have a random source of data x .

We can quantify the **uncertainty** of this source by calculating **the entropy**:

$$\mathbb{H}[x] = - \sum_x p(x) \log p(x)$$

Entropy is max if all x 's are equiprobable.

Entropy is min if the probability of one value is 1.



Claude Shannon

Information (a quick recap)

We have a random source of data x .

We can quantify the **uncertainty** of this source by calculating **the entropy**:

$$\mathbb{H}[x] = - \sum_x p(x) \log p(x)$$

Optimal message length \approx the entropy.



Claude Shannon

Information (a quick recap)

We have two random sources: x and y .

We can quantify the **uncertainty** of them by calculating **the joint entropy**:

$$\mathbb{H}[x, y] = - \sum_{x, y} p(x, y) \log p(x, y)$$

or **the conditional entropy**:

$$\mathbb{H}[y | x] = - \sum_{x, y} p(x, y) \log p(y | x)$$



Claude Shannon

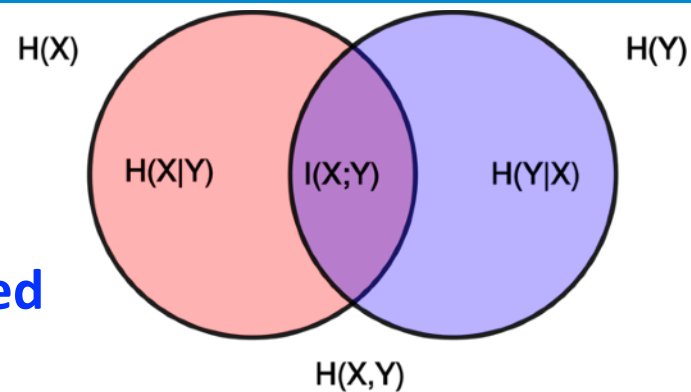
Mutual Information (a quick recap)

We have two random sources: x and y .

Mutual Information (a quick recap)

We have two random sources: x and y .

We can quantify how much **information is shared**
by the two sources:

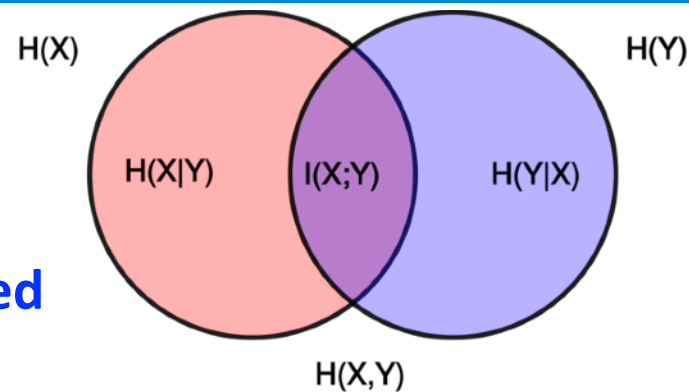


$$I[x; y] = H[y] - H[y | x]$$

Mutual Information (a quick recap)

We have two random sources: x and y .

We can quantify how much **information is shared**
by the two sources:



$$I[x; y] = H[y] - H[y | x]$$

or **how much knowing one source reduces uncertainty about the other.**

We have two random sources: x (e.g., images) and y (e.g., decisions).

We have two random sources: x (e.g., images) and y (e.g., decisions).

We have also a **model** m (a representation of a world).

We have two random sources: x (e.g., images) and y (e.g., decisions).

We have also a **model** m (a representation of a world).

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m :


$$\mathbb{I}[(x, y); m] = \mathbb{H}[x, y] - \mathbb{H}[x, y \mid m]$$

We have two random sources: x (e.g., images) and y (e.g., decisions).

We have also a **model** m (a representation of a world).

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m :

$$\mathbb{I}[(x, y); m] = \mathbb{H}[x, y] - \mathbb{H}[x, y | m]$$



Entropy of the world
(model has no influence on that)



That's the “real” goal!

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m

(or minimize $\mathbb{H}[x, y | m]$, i.e., minimize uncertainty of the world):

$$\mathbb{H}[x, y | m] = \sum_{x, y, m} p(x, y, m) [\log p(y | x, m) + \log p(x | m)]$$


The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m

(or minimize $\mathbb{H}[x, y | m]$, i.e., minimize uncertainty of the world):

$$\mathbb{H}[x, y | m] = \sum_{x, y, m} p(x, y, m) [\log p(y | x, m) + \log p(x | m)]$$



A model for
decision making



A model for
understanding
the world.

The **goal** of AI is to **maximize** the **mutual information** between (x, y) and m (or minimize $\mathbb{H}[x, y | m]$, i.e., minimize uncertainty of the world).

In order to achieve that, AI should focus on learning **two models**:

- **A model for decision making:** $p(y | x, m)$
- **A model for understanding the world:** $p(x | m)$

WHAT HAPPENS IF WE LEARN ONLY DECISION MAKING

The bulk of AI is focused on the decision making part **only!**

WHAT HAPPENS IF WE LEARN ONLY DECISION MAKING

The bulk of AI is focused on the decision making part **only**!

Example: Let's say we have a model that is well trained.



$$p(y = \text{cat}|\mathbf{x}) = 0.90$$

$$p(y = \text{dog}|\mathbf{x}) = 0.05$$

$$p(y = \text{horse}|\mathbf{x}) = 0.05$$

WHAT HAPPENS IF WE LEARN ONLY DECISION MAKING

The bulk of AI is focused on the decision making part **only**!

Example: Let's say we have a model that is well trained.



$p(y = \text{cat}|\mathbf{x}) = 0.90$
 $p(y = \text{dog}|\mathbf{x}) = 0.05$
 $p(y = \text{horse}|\mathbf{x}) = 0.05$

+



noise

=



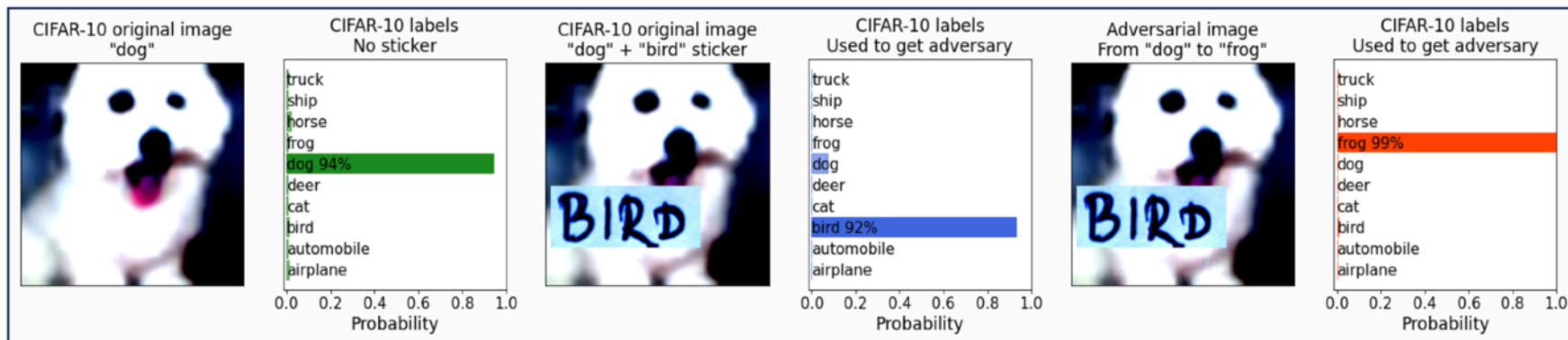
$p(y = \text{cat}|\mathbf{x}) = 0.05$
 $p(y = \text{dog}|\mathbf{x}) = 0.05$
 $p(y = \text{horse}|\mathbf{x}) = 0.90$

But after adding a little noise it could fail completely...

IS LEARNING CLASSIFIERS ENOUGH?

Let's assume we have a perfectly trained neural net.

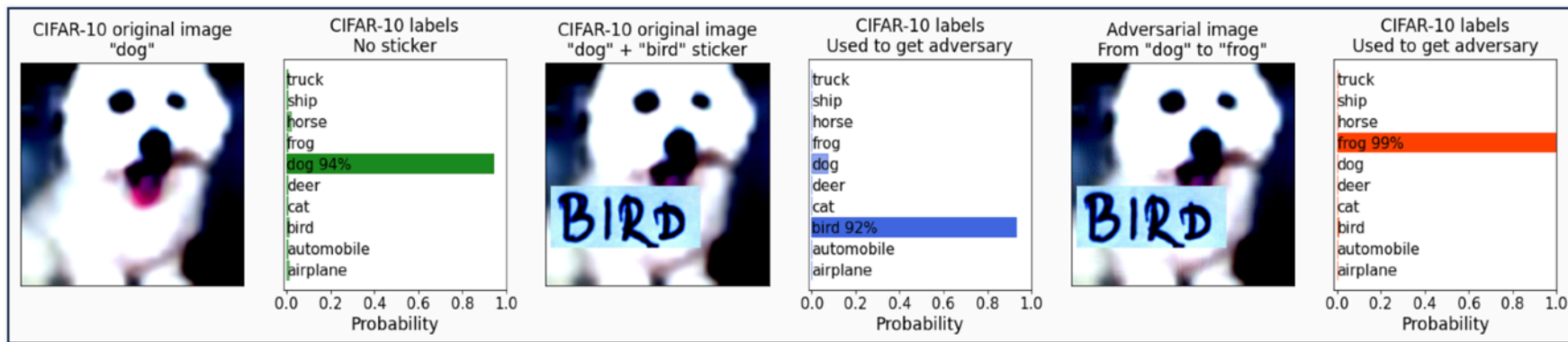
What happens if we add (adversarial) noise to an image?



IS LEARNING CLASSIFIERS ENOUGH?

Let's assume we have a perfectly trained neural net.

What happens if we add (adversarial) noise to an image?

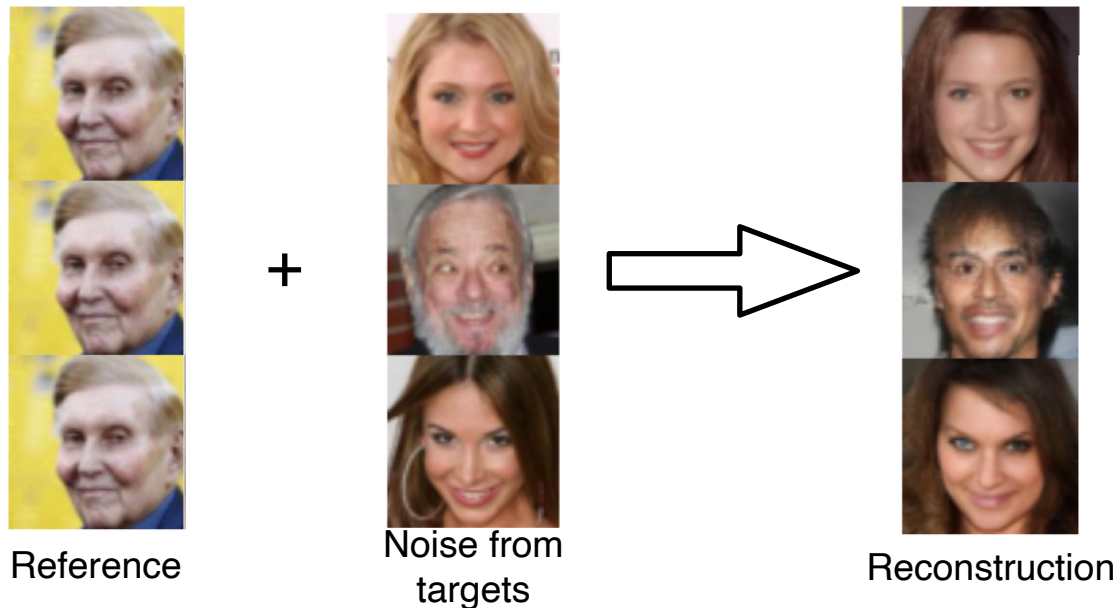


It fails completely...

IS LEARNING CLASSIFIERS ENOUGH?

Let's assume we have a perfectly trained neural net.

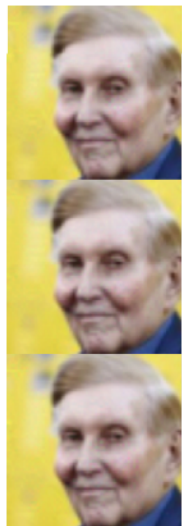
What happens if we add (adversarial) noise to an image?



IS LEARNING CLASSIFIERS ENOUGH?

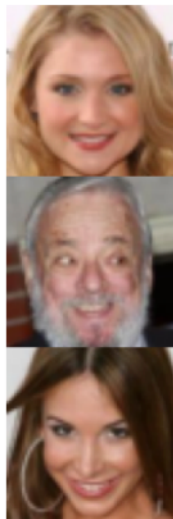
Let's assume we have a perfectly trained neural net.

What happens if we add (adversarial) noise to an image?

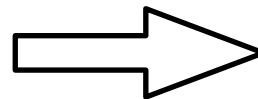


Reference

+



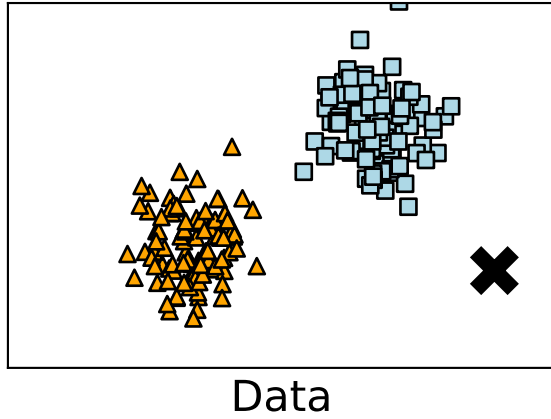
Noise from
targets



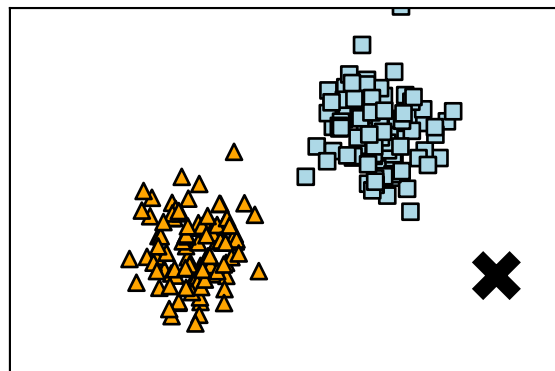
Reconstruction

It fails completely...

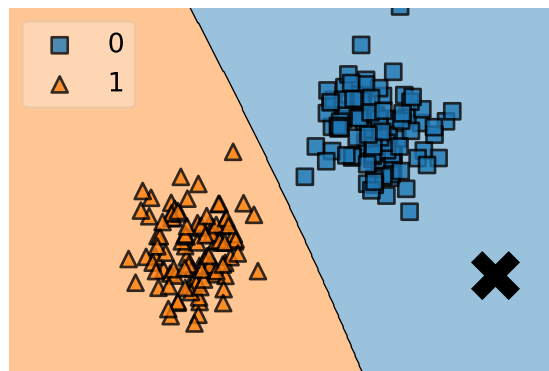
DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?



DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?



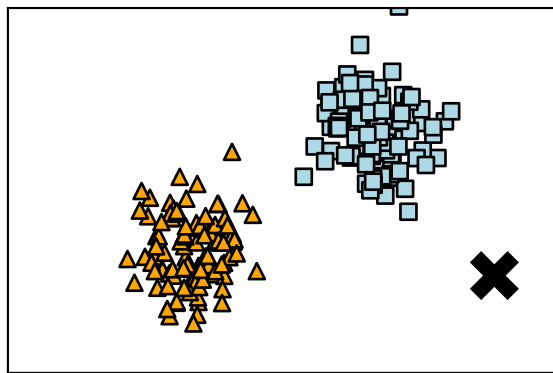
Data



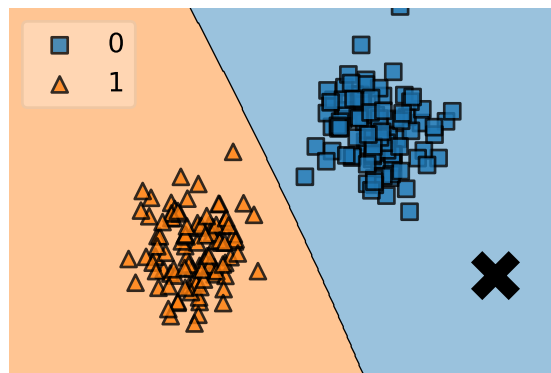
$p(y|\mathbf{x})$

$p(\text{blue}|\mathbf{x})$ is high
= certain decision!

DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?

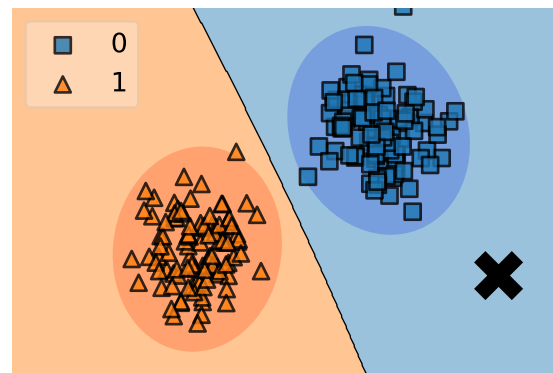


Data



$p(y|\mathbf{x})$

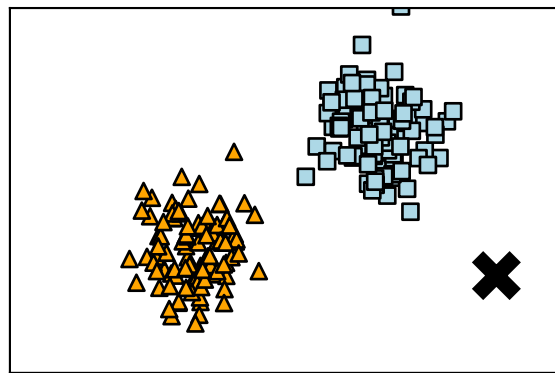
$p(\text{blue}|\mathbf{x})$ is high
= certain decision!



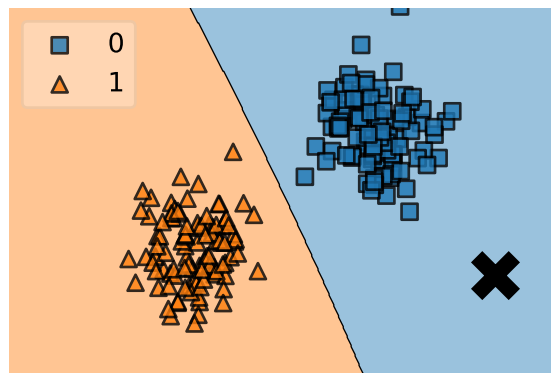
$p(\mathbf{x}, y) = p(y|\mathbf{x}) p(\mathbf{x})$

$p(\text{blue}|\mathbf{x})$ is high
and $p(\mathbf{x})$ is low
= uncertain decision!

DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?

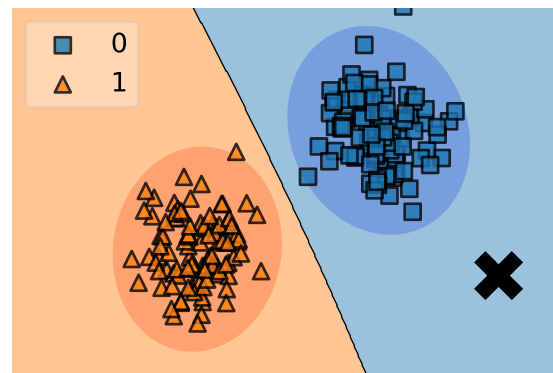


Data



$p(y|\mathbf{x})$

$p(\text{blue}|\mathbf{x})$ is high
= certain decision!



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

$p(\text{blue}|\mathbf{x})$ is high
and $p(\mathbf{x})$ is low
= uncertain decision!

Thus, learning the conditional is only a part of the story!
How can we learn $p(\mathbf{x})$?

DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?

We clearly see that training a neural network (i.e., a conditional distribution):

$$p(y | \mathbf{x}) = \text{softmax} (NN(\mathbf{x}))$$

is **not enough!**



Granny Smith	0.1%
iPod	99.7%
library	0.0%
pizza	0.0%
toaster	0.0%
dough	0.0%

DEEP GENERATIVE MODELING: WHY DO WE NEED THEM?

We clearly see that training a neural network (i.e., a conditional distribution):

$$p(y | \mathbf{x}) = \text{softmax} (NN(\mathbf{x}))$$

is **not enough!**

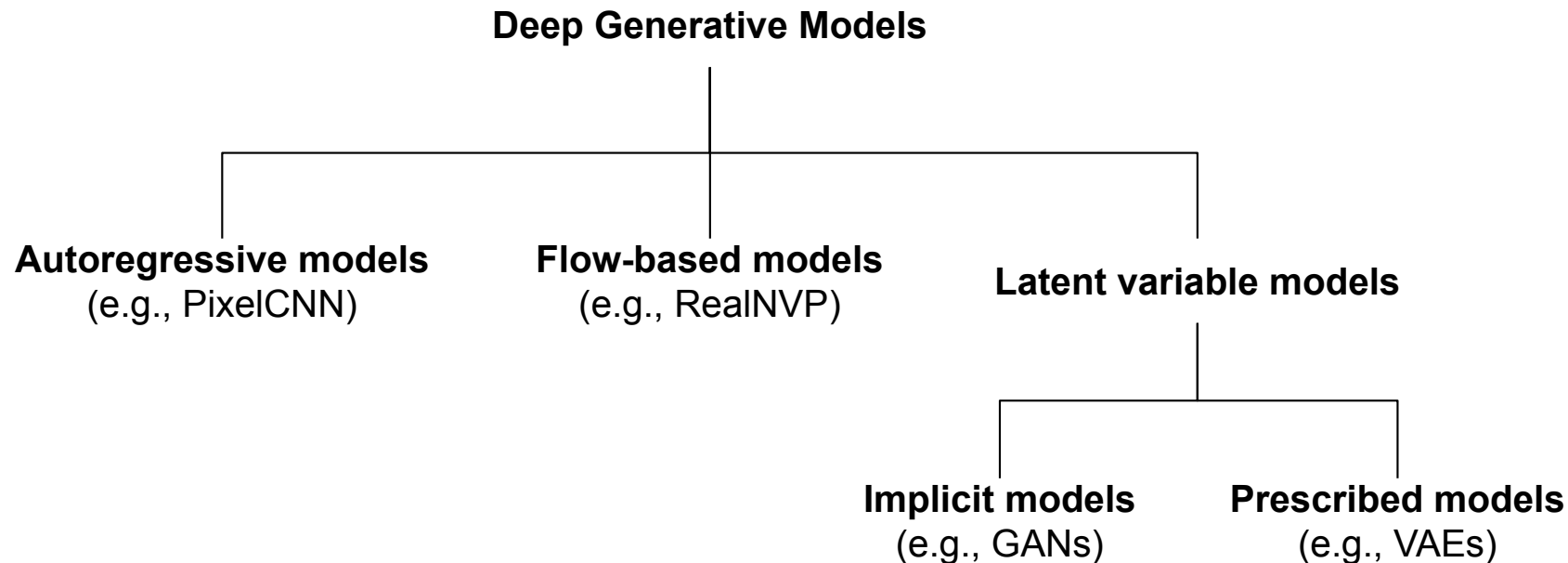
What can we do then?

Or, how to modify the **wrong certainty?**

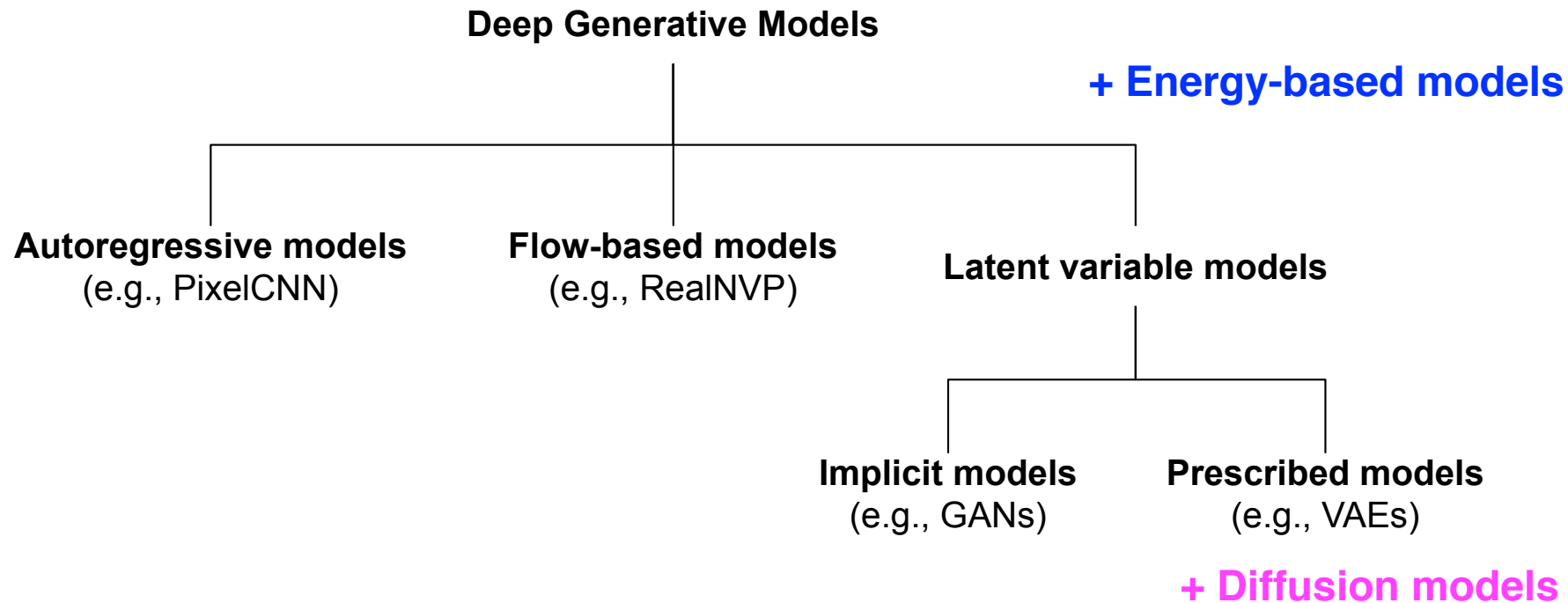


Granny Smith	0.1%
iPod	99.7%
library	0.0%
pizza	0.0%
toaster	0.0%
dough	0.0%

DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?



DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?



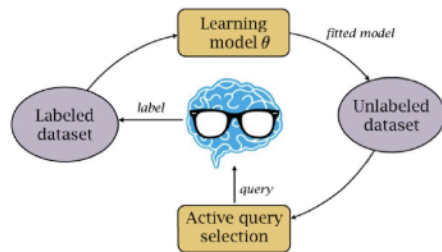
DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?

Generative models	Training	Likelihood	Sampling	Lossy compression	Lossless compression
Autoregressive models	stable	exact	slow	no	yes
Flow-based models	stable	exact	fast/slow	no	yes
Implicit models	unstable	no	fast	no	no
Prescribed model	stable	approximate	fast	yes	no

DEEP GENERATIVE MODELING: WHERE CAN WE USE IT?

" i want to talk to you . "
" i want to be with you . "
" i do n't want to be with you . "
i do n't want to be with you .
she did n't want to be with him .
he was silent for a long moment .
he was silent for a moment .
it was quiet for a moment .
it was dark and cold .
there was a pause .
it was my turn .

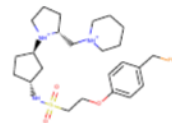
Text analysis



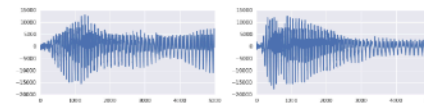
Active Learning



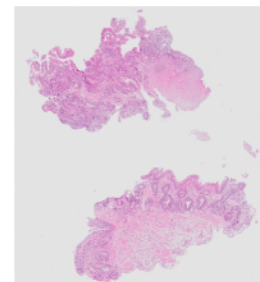
Image analysis



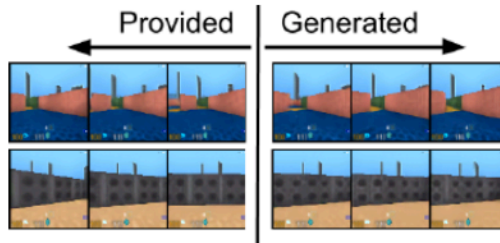
Graph analysis



Audio analysis



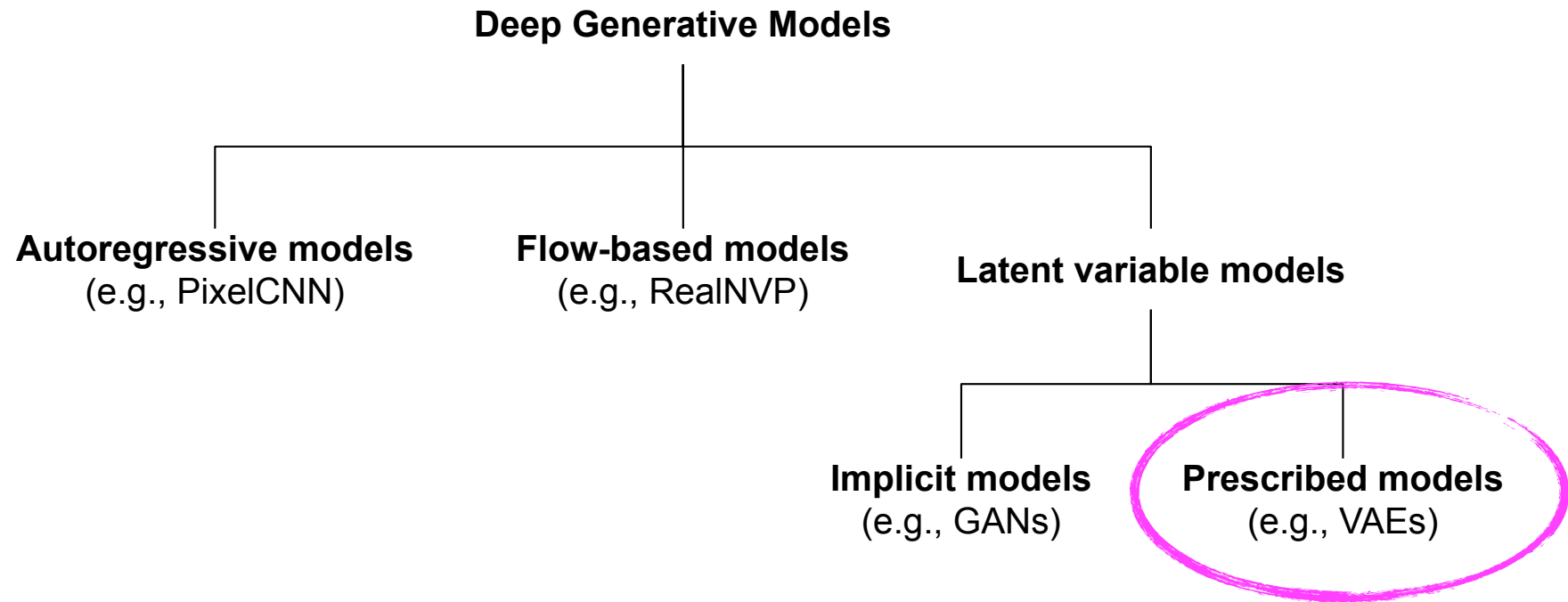
Medical data



Reinforcement Learning

and more...

DEEP GENERATIVE MODELING: HOW WE CAN FORMULATE IT?

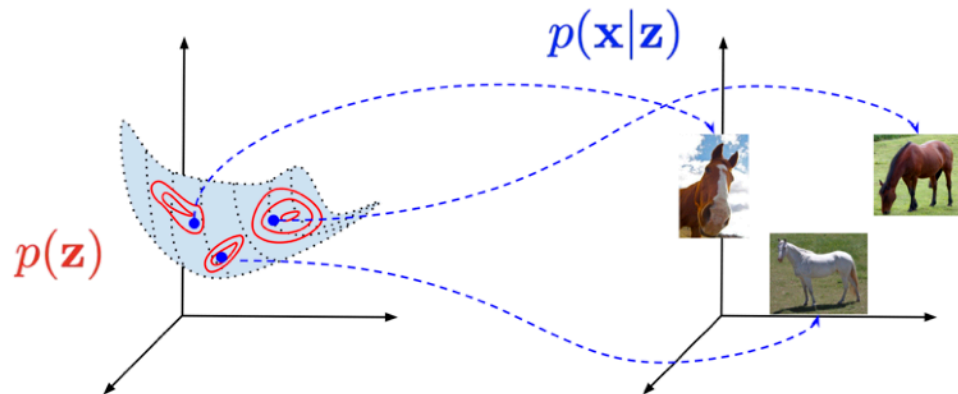


VARIATIONAL AUTO-ENCODERS

Let's consider a **latent variable model** where we distinguish:

- **latent variables** $\mathbf{z} \in \mathcal{Z}^M$
- **observable variables** $\mathbf{x} \in \mathcal{X}^D$

Latent variables lie on a **low-dimensional manifold**.



Generative process:

1. $\mathbf{z} \sim p(\mathbf{z})$
2. $\mathbf{x} \sim p(\mathbf{x} | \mathbf{z})$

VARIATIONAL AUTO-ENCODERS

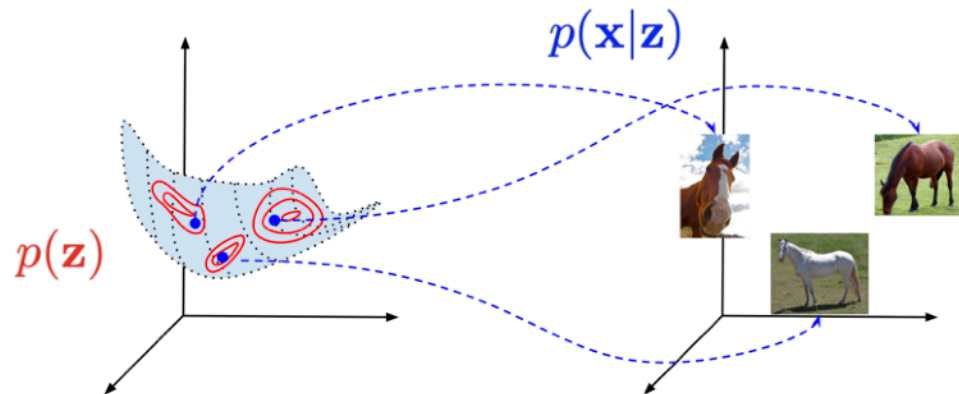
Let's consider a **latent variable model** where we distinguish:

- **latent variables** $\mathbf{z} \in \mathcal{Z}^M$
- **observable variables** $\mathbf{x} \in \mathcal{X}^D$

Latent variables lie on a **low-dimensional manifold**.

The objective function:

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



Generative process:

1. $\mathbf{z} \sim p(\mathbf{z})$
2. $\mathbf{x} \sim p(\mathbf{x} | \mathbf{z})$

VARIATIONAL AUTO-ENCODERS

Let's consider a **latent variable model** where we distinguish:

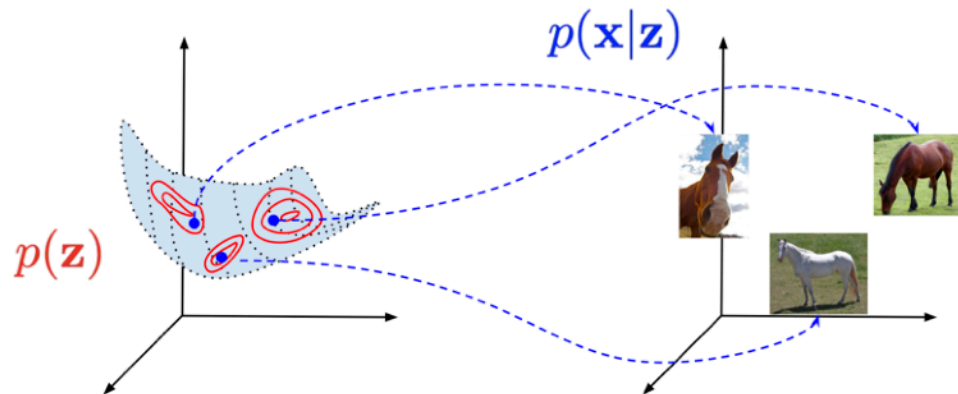
- **latent variables** $\mathbf{z} \in \mathcal{Z}^M$
- **observable variables** $\mathbf{x} \in \mathcal{X}^D$

Latent variables lie on a **low-dimensional manifold**.

The objective function:

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

The integral is intractable...



Generative process:

1. $\mathbf{z} \sim p(\mathbf{z})$
2. $\mathbf{x} \sim p(\mathbf{x} | \mathbf{z})$

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

$$= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z}$$

$$= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]$$

$$\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right]$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) \right] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z}) \right]$$

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

$$= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z}$$

Variational posteriors

$$= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]$$

$$\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right]$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) \right] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z}) \right]$$

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) \right] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z}) \right]\end{aligned}$$

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \quad \text{Jensen's inequality} \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) \right] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z}) \right]\end{aligned}$$

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

$$= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z}$$

$$= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]$$

$$\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]$$

Jensen's inequality

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right]$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) \right] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z}) \right]$$

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) \right] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z}) \right]\end{aligned}$$

$$\begin{aligned}\ln p(\mathbf{x}) &= \ln \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \\&= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) + \ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z}) \right] \\&= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln p(\mathbf{x} | \mathbf{z}) \right] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z}) \right]\end{aligned}$$

VARIATIONAL AUTO-ENCODERS

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

$$= \ln \int \frac{q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z}$$

$$= \ln \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]$$

$$\geq \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} \ln \left[\frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{q_\phi(\mathbf{z})} \right]$$

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} | \mathbf{z})]}_{\text{Reconstruction error}} + \underbrace{\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{z}) - \ln q_\phi(\mathbf{z})]}_{\text{“Regularization” term}}$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z})} [\ln q_\phi(\mathbf{z}) - \ln p(\mathbf{z})]$$

$$\ln p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln q_{\phi}(\mathbf{z} | \mathbf{x}) - \ln p(\mathbf{z})]$$

ELBO: Evidence Lower Bound

$$\ln p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln q_{\phi}(\mathbf{z} | \mathbf{x}) - \ln p(\mathbf{z})]$$

We consider **amortized inference**: $q_{\phi}(\mathbf{z} | \mathbf{x})$

In other words, a single parameterization for each new input \mathbf{x} .

$$\ln p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln q_{\phi}(\mathbf{z}|\mathbf{x}) - \ln p(\mathbf{z})]$$

We consider **amortized inference**: $q_{\phi}(\mathbf{z}|\mathbf{x})$

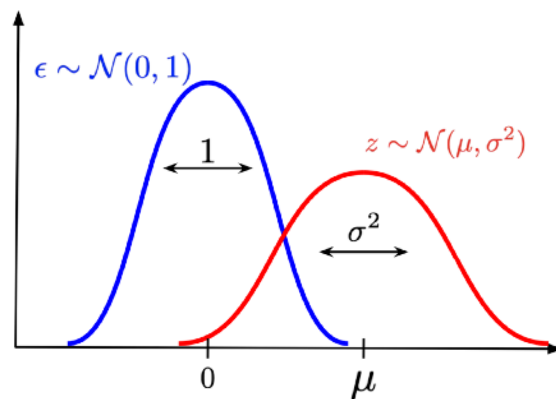
In other words, a single parameterization for each new input \mathbf{x} .

Moreover, we use **reparameterization trick**:

Every Gaussian variable could be defined as:

$$z = \mu + \sigma \cdot \epsilon$$

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



$$\ln p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln q_{\phi}(\mathbf{z} | \mathbf{x}) - \ln p(\mathbf{z})]$$

We consider **amortized inference**: $q_{\phi}(\mathbf{z} | \mathbf{x})$

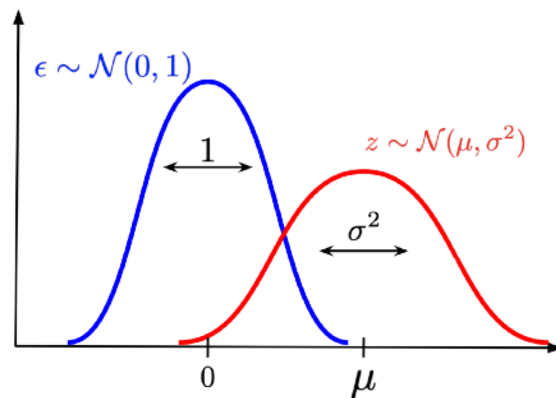
In other words, a single parameterization for each new input \mathbf{x} .

Moreover, we use **reparameterization trick**:

It reduces the variance of the gradients.

It allows to get randomness outside \mathbf{z} .

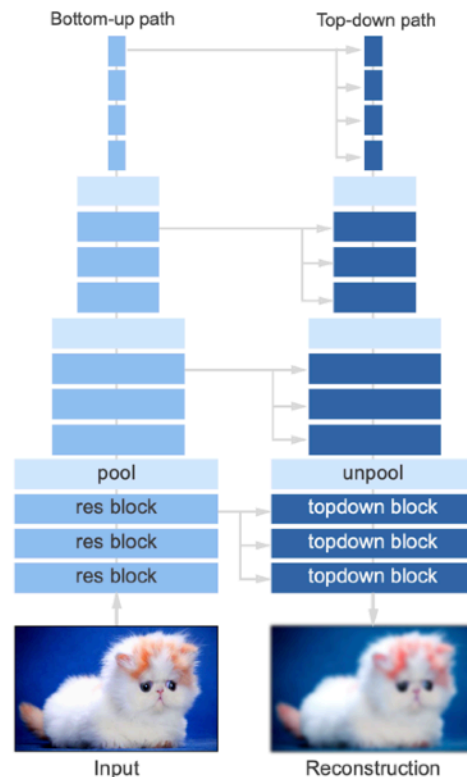
$$\mathbf{z} = \mu + \sigma \cdot \epsilon$$



VARIATIONAL AUTO-ENCODERS



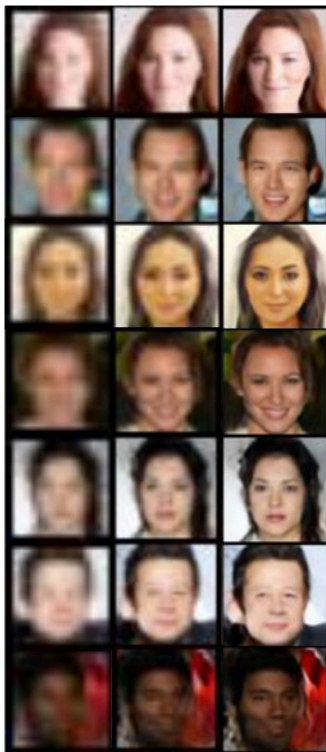
Generations



Very Deep VAE

VARIATIONAL AUTO-ENCODERS

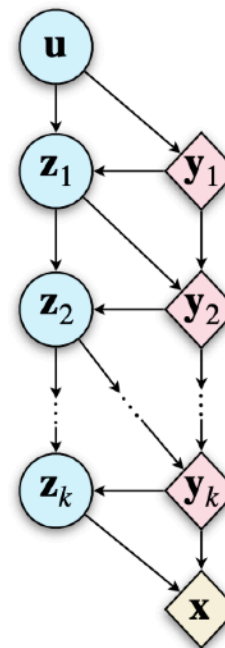
i) selfVAE - downscale - 3lvl



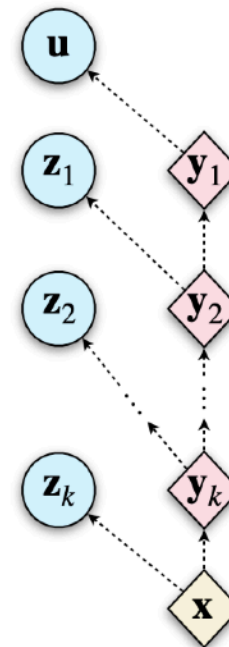
ii) selfVAE - sketch



Generations



i) Generative Model



ii) Inference Model

Hierarchical VAE

CONCLUSION

- Here: **the likelihood-based generative models.**



CONCLUSION

- Here: **the likelihood-based generative models.**
- We **skipped** Generative Adversarial Nets & others.



CONCLUSION

- Here: **the likelihood-based generative models.**
- We **skipped** Generative Adversarial Nets & others.
- Why generative modeling?

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



CONCLUSION

- Here: **the likelihood-based generative models**.
- We **skipped** Generative Adversarial Nets & others.
- Why generative modeling?

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

- Important directions:
 - ➡ Better uncertainty quantification
 - ➡ New parameterization (new neural networks)
 - ➡ Out-of-Distribution
 - ➡ Continual learning



If you are interested in going deeper into deep generative modeling, please take a look at my blog: [\[Blog\]](#)

- **Intro:** [\[Link\]](#)
- **ARMs:** [\[Link\]](#)
- **Flows:** [\[Link\]](#), [\[Link\]](#)
- **VAEs:** [\[Link\]](#), [\[Link\]](#)
- **Hybrid modeling:** [\[Link\]](#)

THANK YOU FOR YOUR ATTENTION

Jakub M. Tomczak
Computational Intelligence group
Vrije Universiteit Amsterdam

Webpage: <https://jmtomczak.github.io/>

Github: <https://github.com/jmtomczak>

Twitter: <https://twitter.com/jmtomczak>