

Variational autoencoder

Prof. Ph.D. Woo Youn Kim
Chemistry, KAIST

Goals

10주	주제	Message Passing Neural Network (MPNN)
	목표	Understanding the most general expression of graph neural network
	내용	MPNN, molecular graph representation, GGNN, supervised learning of logP and TPSA
11주	주제	Molecular generative model 1
	목표	Understanding the principle of autoencoder and unsupervised learning
	내용	Molecular autoencoder, VAE, CVAE, de novo molecular design
12주	주제	Molecular generative model 2
	목표	Understanding difference between GAN and VAE
	내용	GAN, ARAE ARAE: conditional molecular design

Review

Bayes' Rule

$$P(\theta|x) = \frac{P(x|\theta)P(\theta)}{P(x)} = \frac{P(x|\theta)P(\theta)}{\sum_{\theta} P(x|\theta)P(\theta)}$$

Posterior probability = $\frac{(\text{likelihood}) \times (\text{prior probability})}{(\text{evidence})}$

Prior probability (사전확률) $P(\theta)$: probability of a parameter set θ .

Posterior probability (사후확률) $P(\theta|x)$: $P(\theta)$ given an observation x .

Evidence (증거) $P(x)$: probability of the observation.

Likelihood (가능도) $L(\theta|x) = P(x|\theta)$: $P(x)$ given the parameters θ

θ = disease

x = test result (T or F)

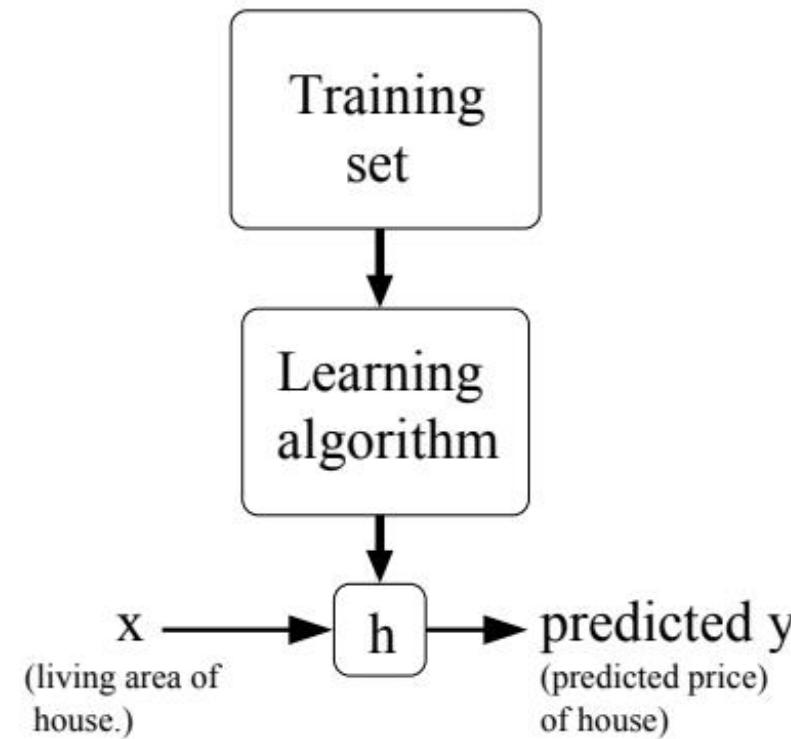
$P(x)$: probability of T and F

$P(\theta)$: probability of having a disease

$P(\theta|x)$: probability of having the disease given the test result

$P(x)$: probability of T and F.

Review



The principle of **maximum likelihood** says that we should choose θ so as to make the data as high probability as possible.

Review

$$Cost = - \sum_{i=1}^m [y^{(i)} \ln h(x^{(i)}) + (1 - y^{(i)}) \ln(1 - h(x^{(i)}))]$$



$$Cost = - \sum_{i=1}^m p_{true}(x^{(i)}) \ln p_{pred}(x^{(i)})$$

VS

Gibbs entropy formula

$$S = -k_B \sum_i p_i \ln p_i$$

Entropy in information theory

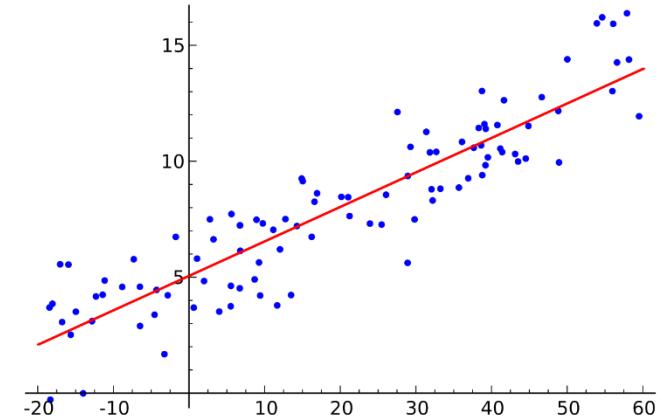
It is also called the **cross entropy**

$$S = - \sum_{i=1}^m p_i \ln p_i$$

Types of deep learning

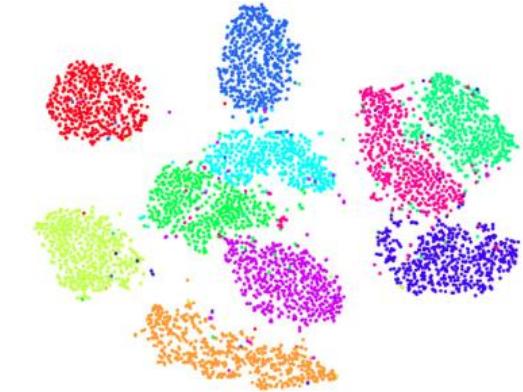
- **Supervised Learning: classification or regression**

The network makes its guesses, then compare its answers to the known “correct” ones and make adjustments according to its errors.



- **Unsupervised Learning: clustering**

Searching for a hidden pattern in a data set without known answers.



- **Reinforcement Learning: game, robotics, self-driving car, etc.**

A strategy built on observation.

no labeling

<https://www.youtube.com/watch?v=JFJkpVWTQVM>

Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification, regression,
object detection, semantic
segmentation, image captioning, etc.

Unsupervised Learning

Data: x

Just data, no labels!

Training data is cheap

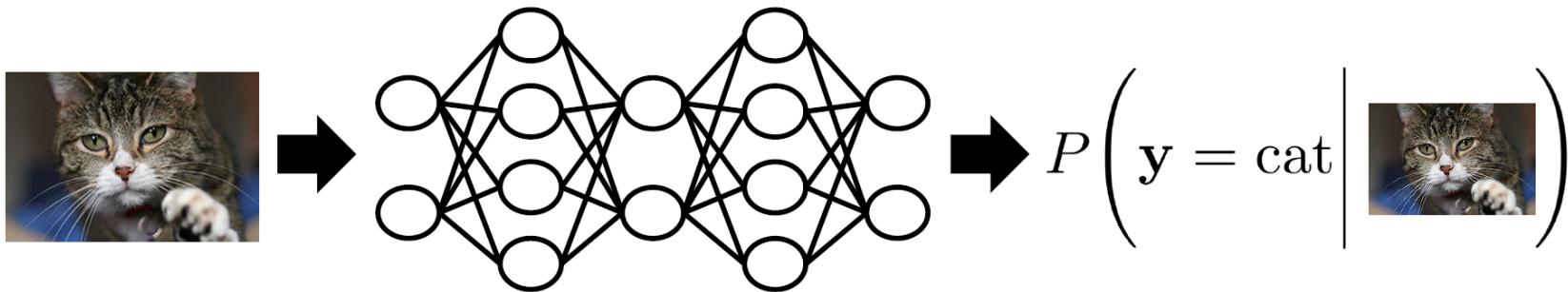


Goal: Learn some underlying hidden
structure of the data

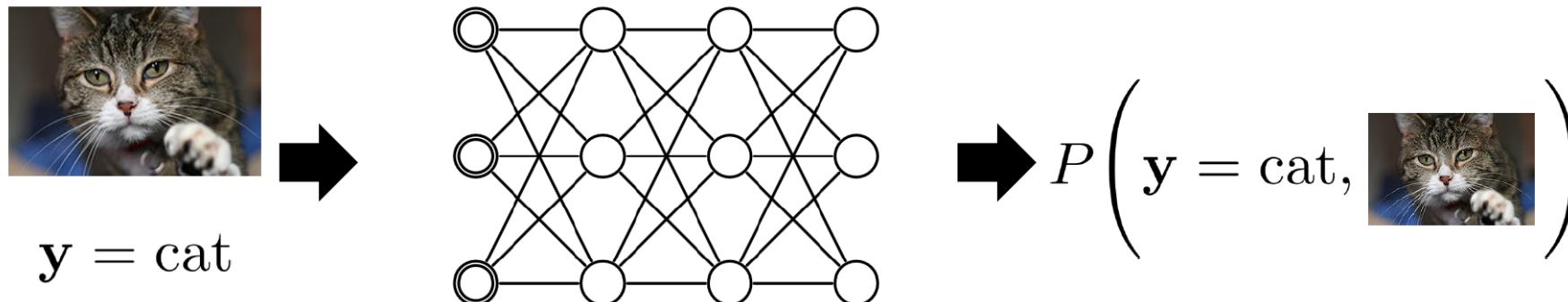
Examples: Clustering, dimensionality
reduction, feature learning, density
estimation, etc.

Generative model

- Given an observed variable x and a target variable y
- Discriminative model** is a model of a **conditional distribution** $P(y|x)$
e.g., neural networks (supervised)

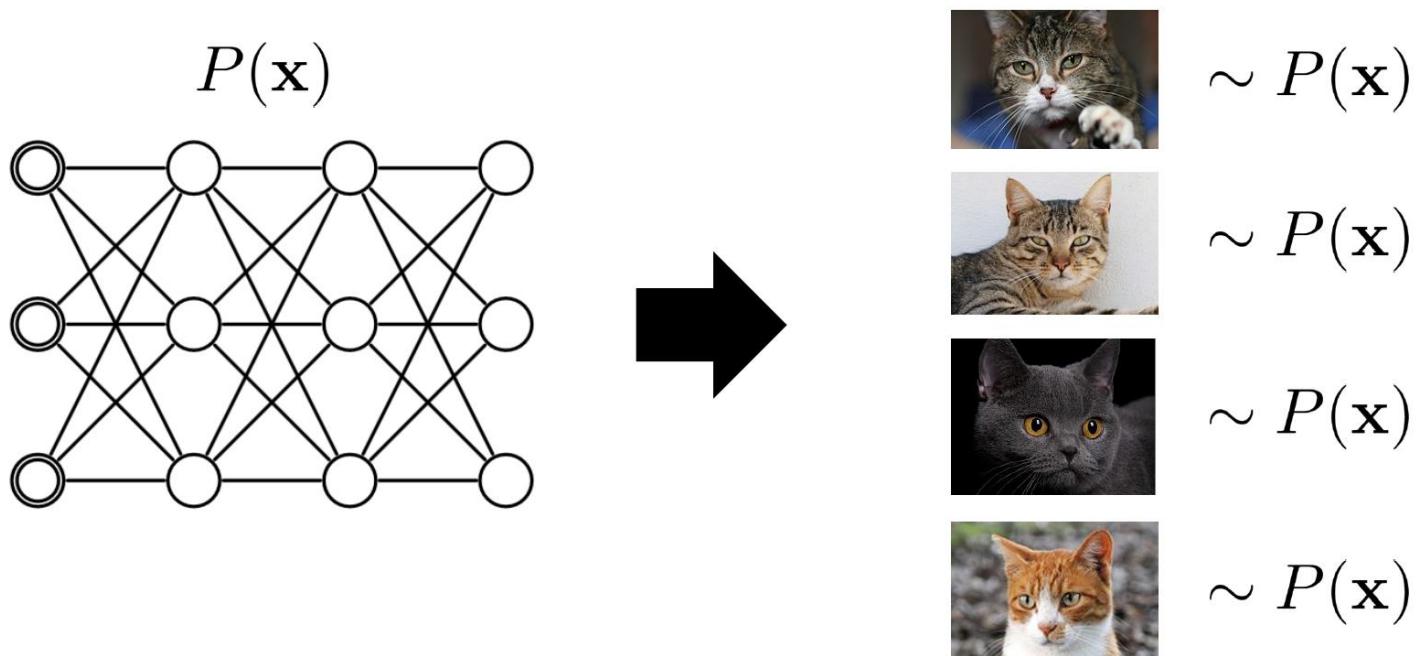


- Generative model** is a model of a **joint distribution** $P(x, y)$ (or $P(x)$)
e.g., RBM, VAE, GAN (unsupervised)



Generative model

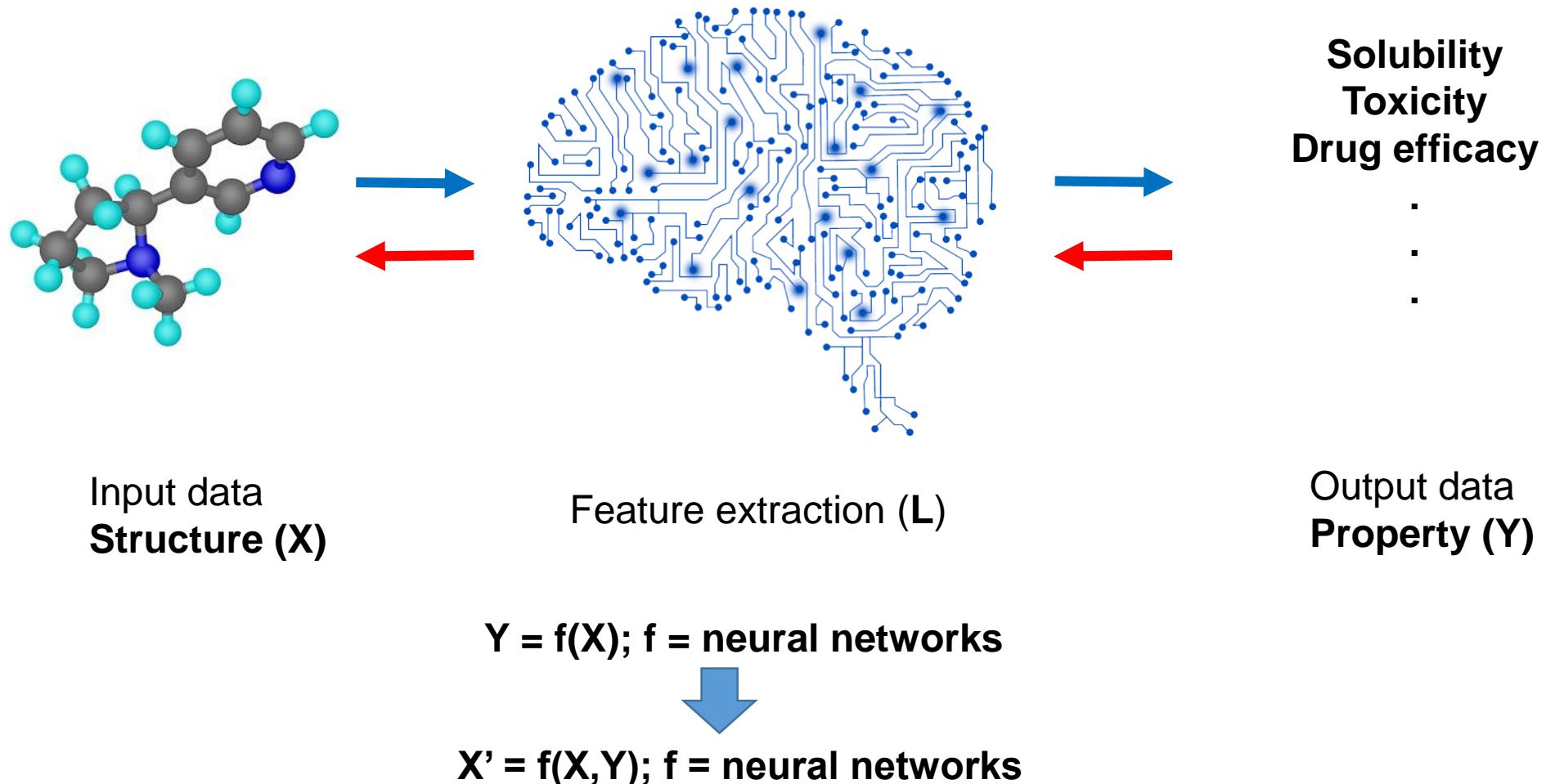
- Generative models model a full probability distribution given data
- $P(x,y)$ enables us to **generate new data** similar to existing (training) data



$P(\mathbf{x})$ = distribution of an observed variable \mathbf{x} (ex: pixel info [255, 255, 255])

Inverse design

- In chemistry or materials science, it is also called the inverse design



Taxonomy of Generative Models

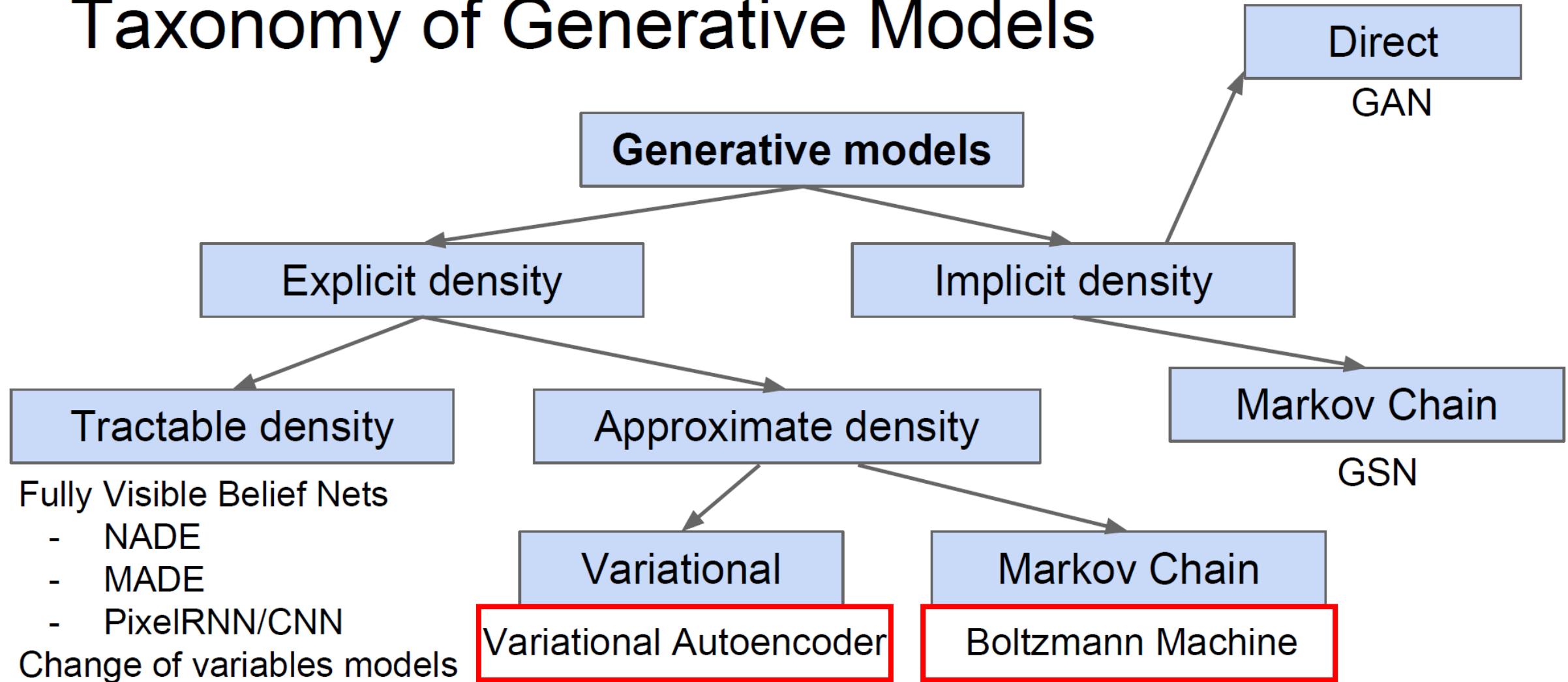


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Reference paper

RESEARCH ARTICLE

Open Access



Molecular generative model based on conditional variational autoencoder for de novo molecular design

Jaechang Lim¹, Seongok Ryu¹, Jin Woo Kim¹ and Woo Youn Kim^{1,2*} 

Abstract

We propose a molecular generative model based on the conditional variational autoencoder for de novo molecular design. It is specialized to control multiple molecular properties simultaneously by imposing them on a latent space. As a proof of concept, we demonstrate that it can be used to generate drug-like molecules with five target properties. We were also able to adjust a single property without changing the others and to manipulate it beyond the range of the dataset.

Keywords: Molecular design, Conditional variational autoencoder, Deep learning

Cite This: ACS Cent. Sci. 2018, 4, 268–276

Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules

Rafael Gómez-Bombarelli,^{†,‡,§,¶,#} Jennifer N. Wei,^{‡,§,#} David Duvenaud,^{¶,#} José Miguel Hernández-Lobato,^{§,#} Benjamín Sánchez-Lengeling,[‡] Dennis Sheberla,^{‡,§,#} Jorge Aguilera-Iparraguirre,[†] Timothy D. Hirzel,[†] Ryan P. Adams,^{▽,||} and Alán Aspuru-Guzik^{*,‡,§,⊥,#}

[†]Kyulux North America Inc., 10 Post Office Square, Suite 800, Boston, Massachusetts 02109, United States

[‡]Department of Chemistry and Chemical Biology, Harvard University, Cambridge, Massachusetts 02138, United States

[¶]Department of Computer Science, University of Toronto, 6 King's College Road, Toronto, Ontario M5S 3H5, Canada

[§]Department of Engineering, University of Cambridge, Trumpington Street, Cambridge CB2 1PZ, U.K.

[▽]Google Brain, Mountain View, California, United States

^{||}Princeton University, Princeton, New Jersey, United States

^{*}Biologically-Inspired Solar Energy Program, Canadian Institute for Advanced Research (CIFAR), Toronto, Ontario M5S 1M1, Canada

Contents

- Restricted Boltzmann machine
- Deep belief network & Auto-encoder
- Variational auto-encoder
- Molecular VAE
- Jointly trained VAE

Restricted Boltzmann machine

Example of generative model

1 Upload photo

The first picture defines the scene you would like to have painted.



2 Choose style

Choose among predefined styles or upload your own style image.



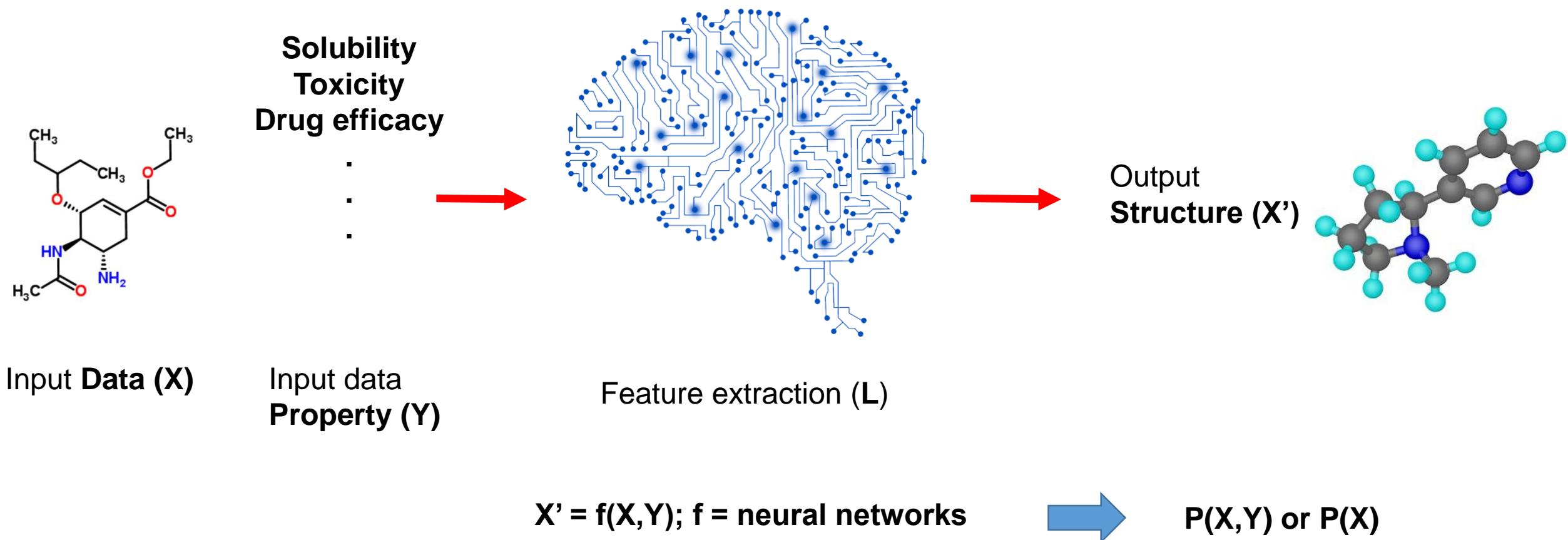
3 Submit

Our servers paint the image for you. You get an email when it's done.



- 손석희 아나운서의 목소리를 만들어내는 인공지능 모델
(박근혜 전 대통령, 문재인 대통령도 되었었는데, 현재 삭제된 상태)
: <https://carpedm20.github.io/tacotron/>

Molecular generative model



Boltzmann machine

- **Energy based model** (EBM) is an approximation of the distribution of observed variable x

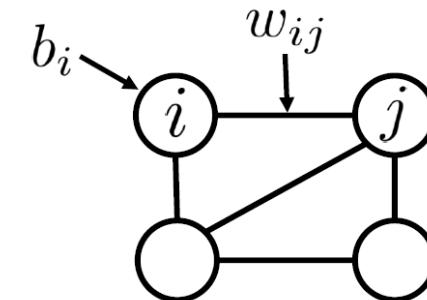
$$P(x) = \frac{\exp(-E(x))}{Z} \quad \text{Statistical probability (Boltzmann factor)}$$

where $E(x)$ is the energy function, and Z is the partition function

- Assignments with high energy appear less likely in EBM
- Given a graph $G(V,E)$, **Boltzmann machine** (BM) is a **joint distribution** on a binary vector

$$P(x) \propto \exp \left(\sum_{i \in V} b_i x_i + \sum_{(i,j) \in E} w_{ij} x_i x_j \right)$$

Onsite energy correlation energy



- To generate new data using BM, we need to learn parameters of BM

Boltzmann machine

- Given training data [$x=(x_1, \dots, x_n)$] ($P(x)$), learn the distribution $Q(x)$
 - $P(x)$: observed distribution, $Q(x)$: model distribution
- Find the optimal parameter (\mathbf{b}, \mathbf{w}) that minimizes the difference between $P(x)$ and $Q(x)$
- Use the Kullback–Leibler (KL) divergence as the **directed divergence** between two distributions (1951)

$$D_{KL}(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

KL divergence

$$D_{KL}(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

data entropy

- Physical interpretation

$$\begin{aligned} D_{KL}(P \parallel Q) &= \sum_x P(x) \log P(x) - \sum_x P(x) \log Q(x) & \beta = 1 \\ &= -S_{data} - \sum_x P(x) \log \frac{\exp(-\beta E(x))}{Z} \\ &= -S_{data} + \beta \sum_x E(x) P(x) + \sum_x P(x) \log Z \\ &= -S_P + \beta \langle E(x) \rangle_P + \log Z \\ &= \beta \left[\langle E(x) \rangle_P - k_B T S_P \right] + \beta [k_B T \log Z] \\ &= \beta (F_{data} - F_{model}) = -(S_{data} - S_{model} \text{ (cross entropy)}) \end{aligned}$$

Free energy difference between data and model distributions

KL divergence

$$D_{KL}(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

- Property

The KL divergence is always non-negative.

$$D_{KL}(P \parallel Q) = -(S_{\text{data}} - S_{\text{model}} \text{ (cross entropy)}) \geq 0$$

It is known as Gibbs's inequality, with zero if and only if $P = Q$ everywhere.

- Derivative

$$-\frac{\partial D_{KL}(P \parallel Q)}{\partial b_i} = \sum_x P(x)x_i - \sum_x Q(x)x_i = \langle x_i \rangle_{\text{data}} - \langle x_i \rangle_{\text{model}} = E_{\text{data}}[x_i] - E_{\text{model}}[x_i]$$

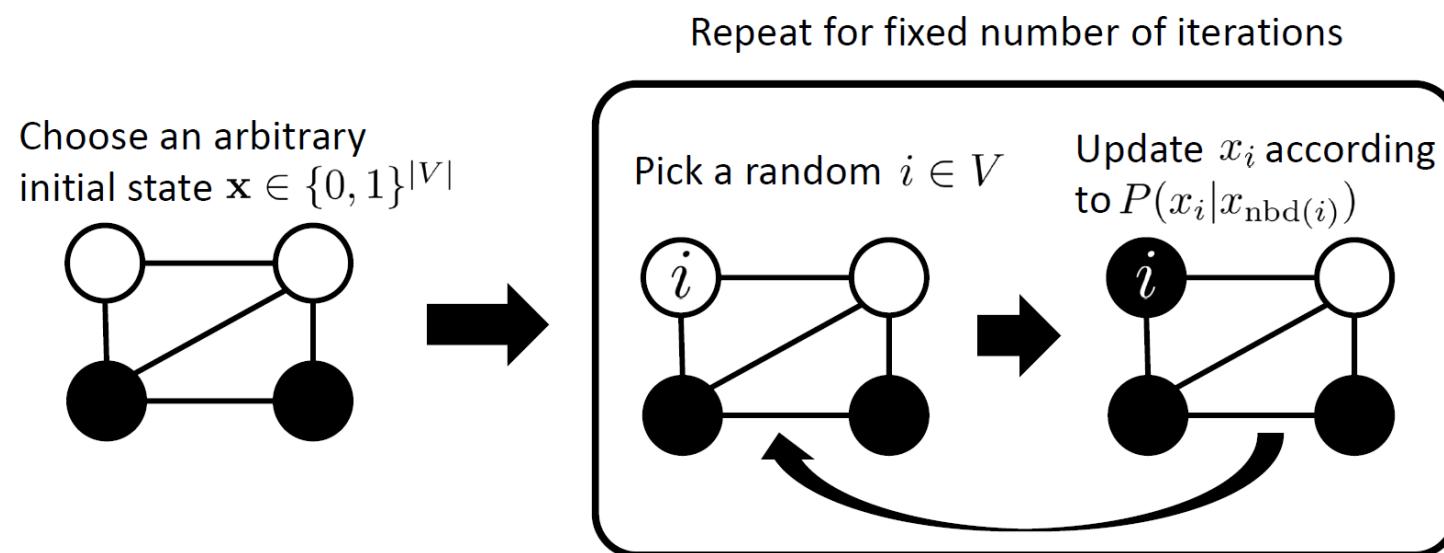
$$-\frac{\partial D_{KL}(P \parallel Q)}{\partial w_{ij}} = \langle x_i x_j \rangle_{\text{data}} - \langle x_i x_j \rangle_{\text{model}} = E_{\text{data}}[x_i x_j] - E_{\text{model}}[x_i x_j]$$

Boltzmann machine

- Update b , w using the gradient descent

$$b_i \leftarrow b_i - \alpha \frac{\partial D_{KL}(P \parallel Q)}{\partial b_i} \quad w_{ij} \leftarrow w_{ij} - \alpha \frac{\partial D_{KL}(P \parallel Q)}{\partial w_{ij}}$$

- **Gibbs sampler** is the most popular sampling algorithm in BM



Convergence problem

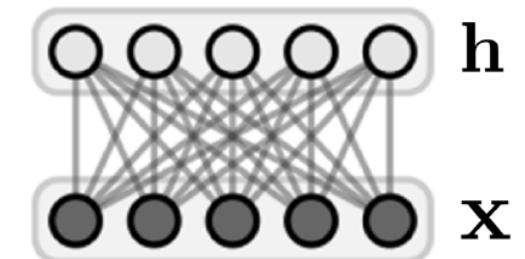
Restricted Boltzmann machine

- Restricted Boltzmann machine (RBM) is a bipartite Boltzmann machine with **visible nodes** and **hidden nodes**

$$P(\mathbf{x}) \propto \sum_{\mathbf{h}} \exp \left(\sum_i b_i x_i + \sum_j c_j h_j + \sum_{i,j} w_{ij} x_i h_j \right)$$



Higher order potential



- Hidden nodes can be described by **hidden features** of visible nodes
- In RBM structure, all hidden nodes are **conditionally independent** given visible nodes and vice versa
no connection between visible nodes \mathbf{x} and no connection between hidden nodes \mathbf{h}

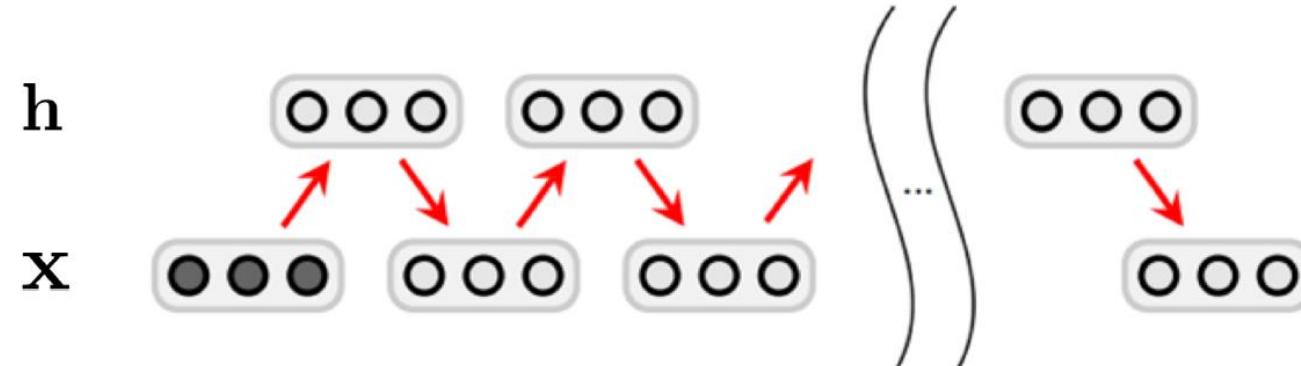
$$P(x_i = 1 | \mathbf{h}) = \sigma \left(\sum_j w_{ij} h_j + b_i \right) \quad P(h_j = 1 | \mathbf{x}) = \sigma \left(\sum_i w_{ij} x_i + c_j \right)$$

Restricted Boltzmann machine

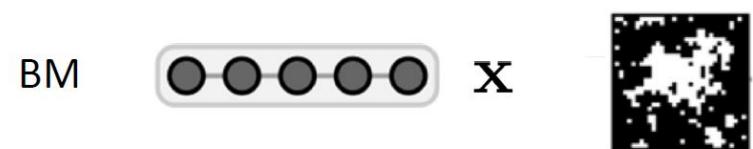
- Given training data $[x=(x_1, \dots, x_n)] (P(x))$, learn the distribution $Q(x)$
 - $P(x)$: observed distribution, $Q(x)$: model distribution
- Find the optimal parameter $(\mathbf{b}, \mathbf{w}, \mathbf{h})$ by minimizing the KL divergence
- Gradient descent converges to the **global optimum** with gradients

Restricted Boltzmann machine

- Due to conditional independence of RBM, **block Gibbs sampling** is possible



- Samples generated by BM and RBM



Deep belief network & Auto-encoder

Auto-encoder

Reducing the Dimensionality of Data with Neural Networks

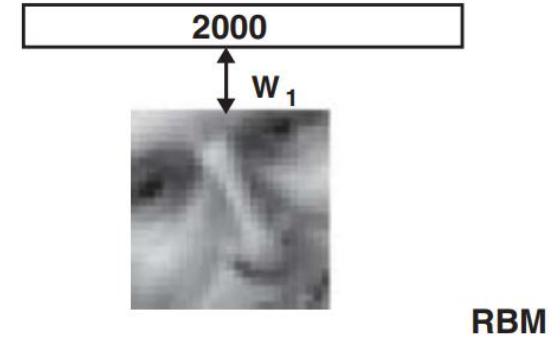
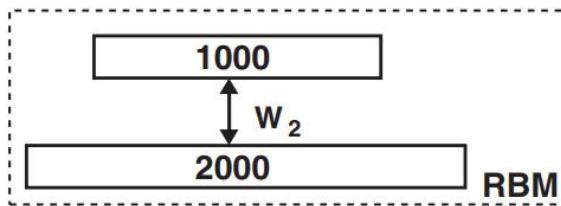
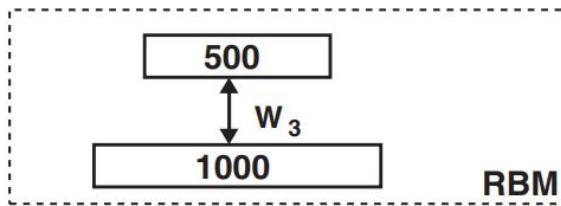
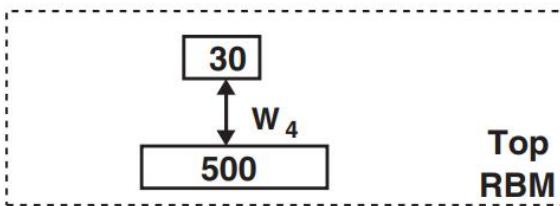
G. E. Hinton* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such “autoencoder” networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

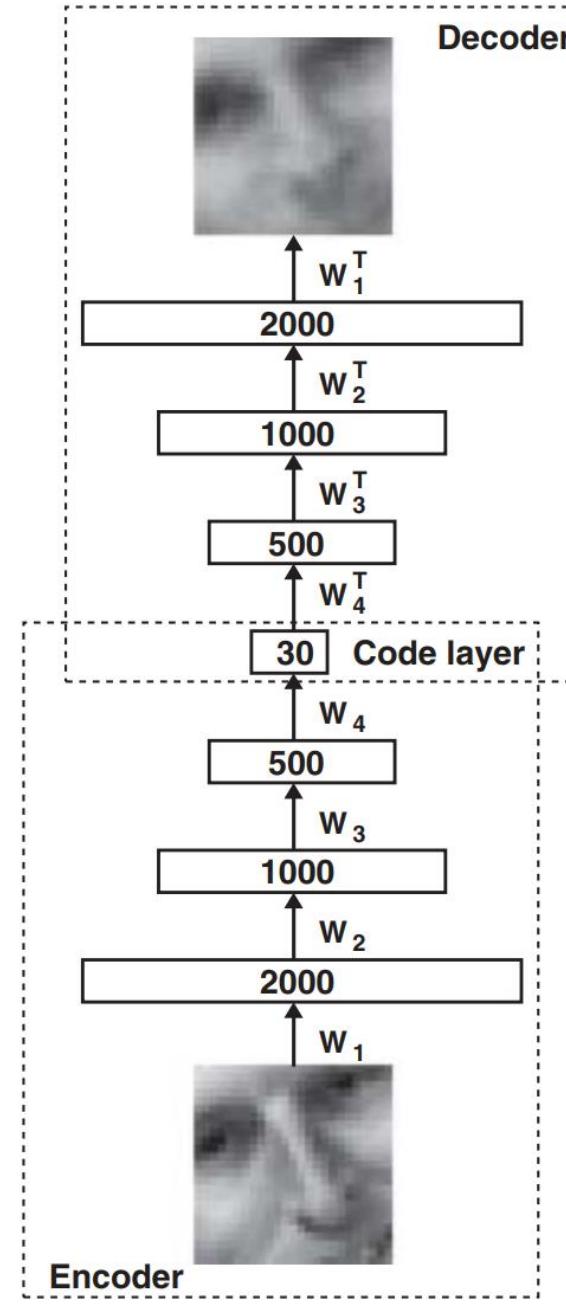
- ‘Principle Component Analysis (PCA)’ has been widely used in classification problems in an unsupervised manner.
- However, PCA has limitations in very complex problems.
- This paper proposed an “auto-encoder” method to reconstruct image from extracted feature without loss of accuracy.
- It is called the ‘Deep Belief Network’ which is stacked ‘Restricted Boltzmann Machine’.

Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *science* 313.5786 (2006): 504-507.

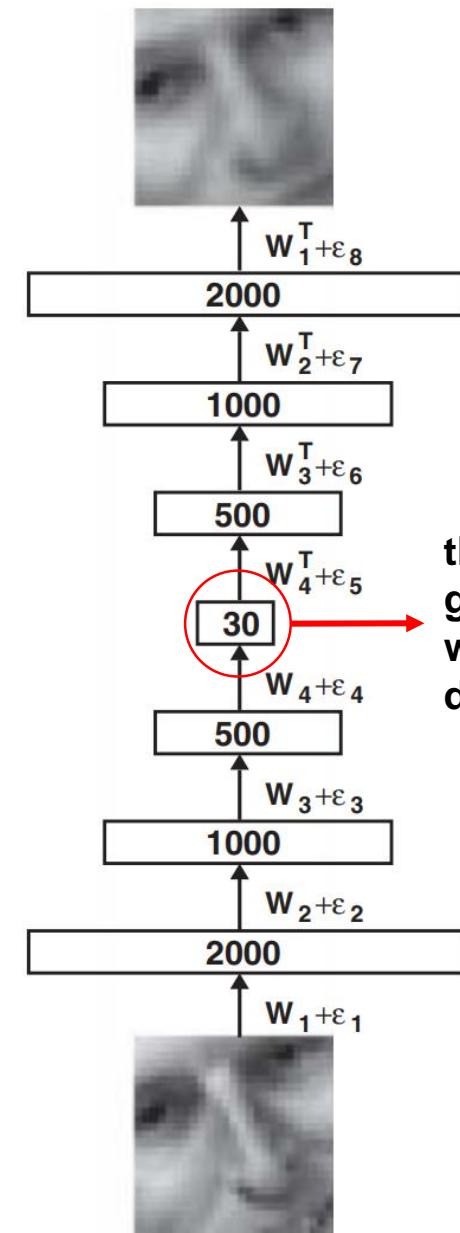
This is a stack
of RBMs.



Pretraining



Unrolling



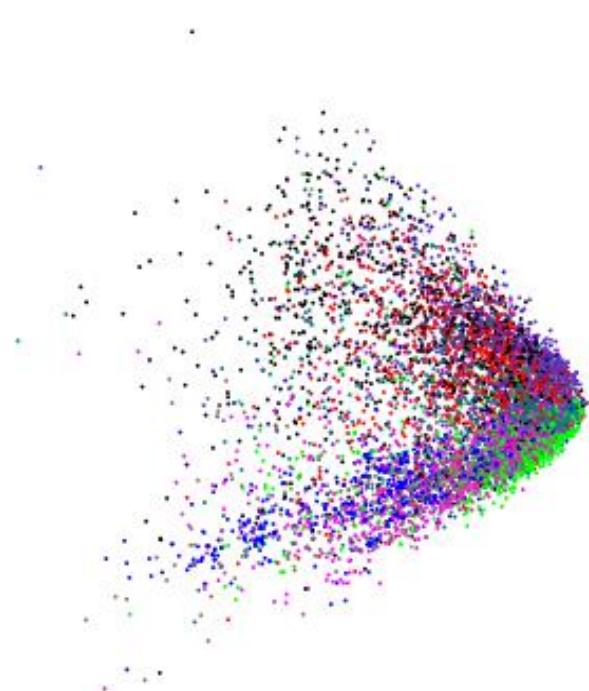
Fine-tuning

this hidden layer
gives a latent vector
with reduced
dimension.

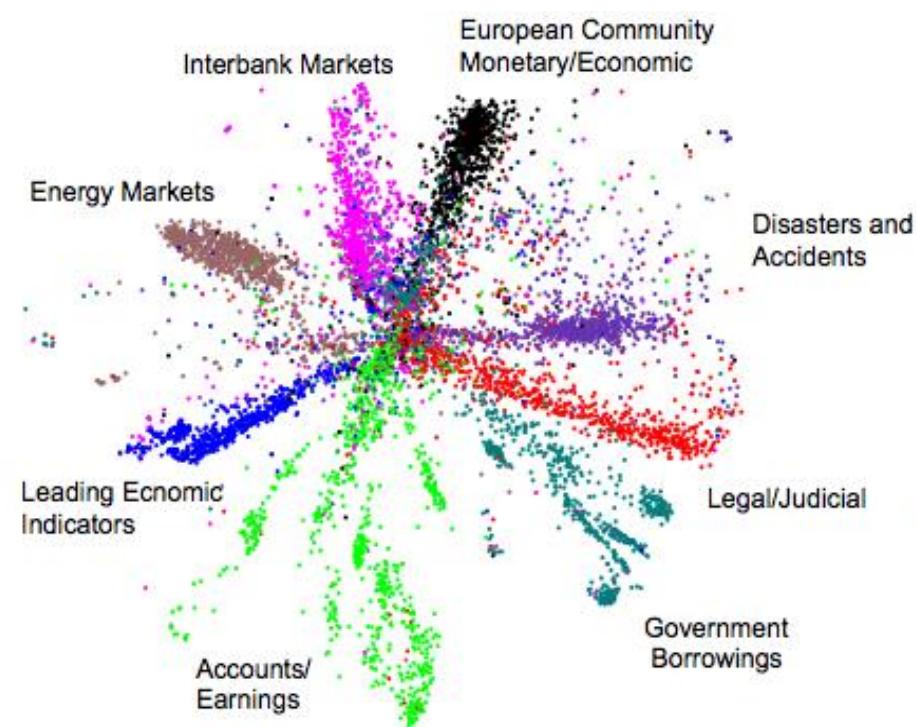
Unsupervised learning

From the raw data, **the deep auto-encoder** find the features through learning himself.

Principle Component Analysis(PCA)



Deep Auto-encoder



Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *science* 313.5786 (2006): 504-507.

Supervised vs Unsupervised

	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	classification or categorization	clustering
<i>Continuous</i>	regression	dimensionality reduction

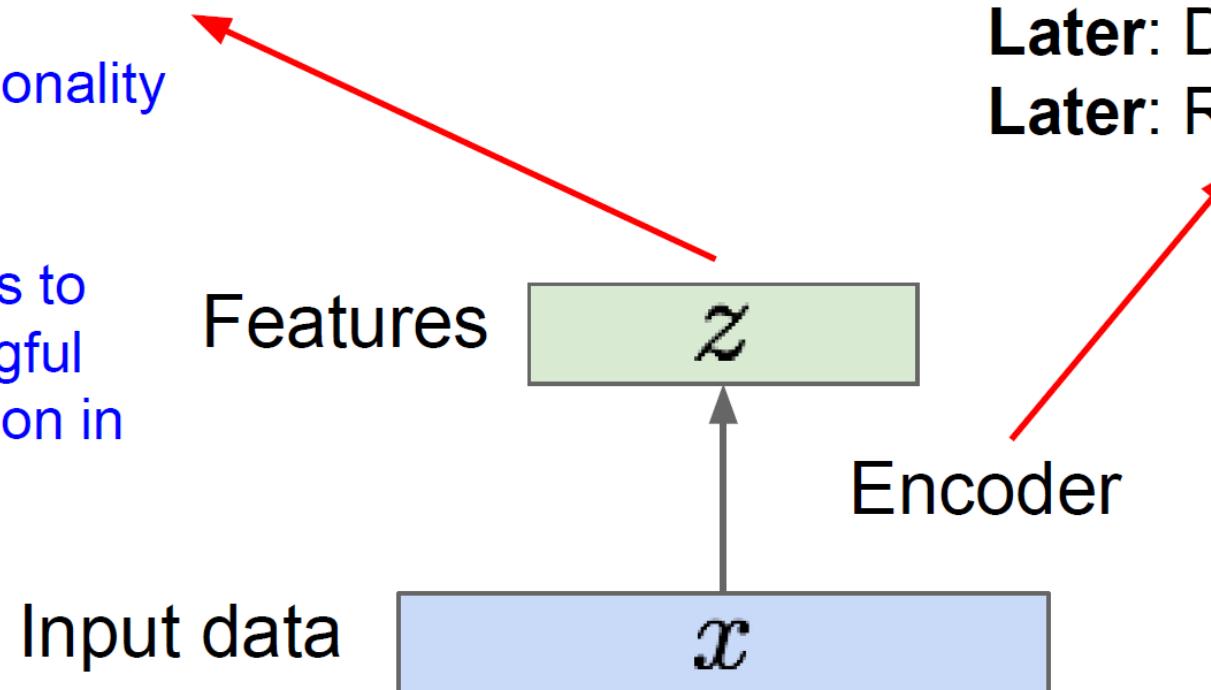
Auto-encoder

z usually smaller than x
(dimensionality reduction)

Q: Why dimensionality reduction?

A: Want features to capture meaningful factors of variation in data

Originally: Linear + nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN

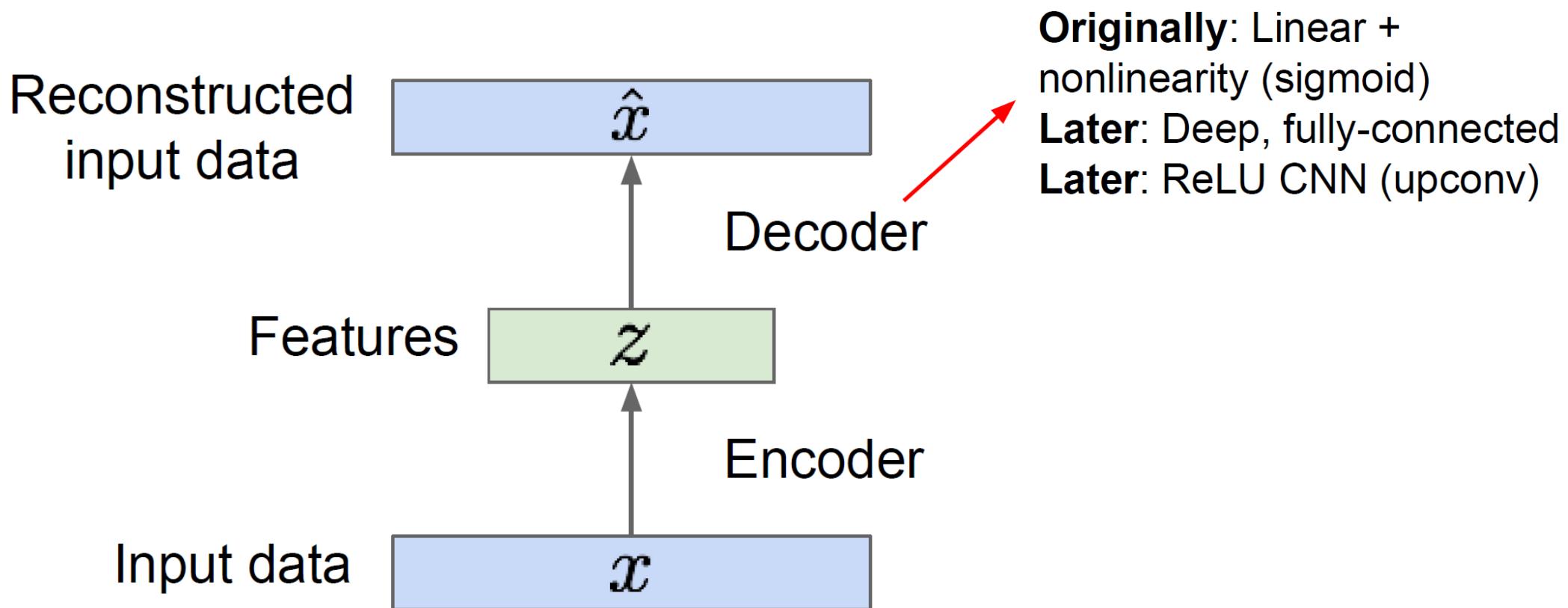


Auto-encoder

How to learn this feature representation?

Train such that features can be used to reconstruct original data

“Autoencoding” - encoding itself



Auto-encoder

Train such that features can be used to reconstruct original data

Reconstructed input data

Features

Input data

L2 Loss function:

Doesn't use labels

$$\|x - \hat{x}\|^2$$

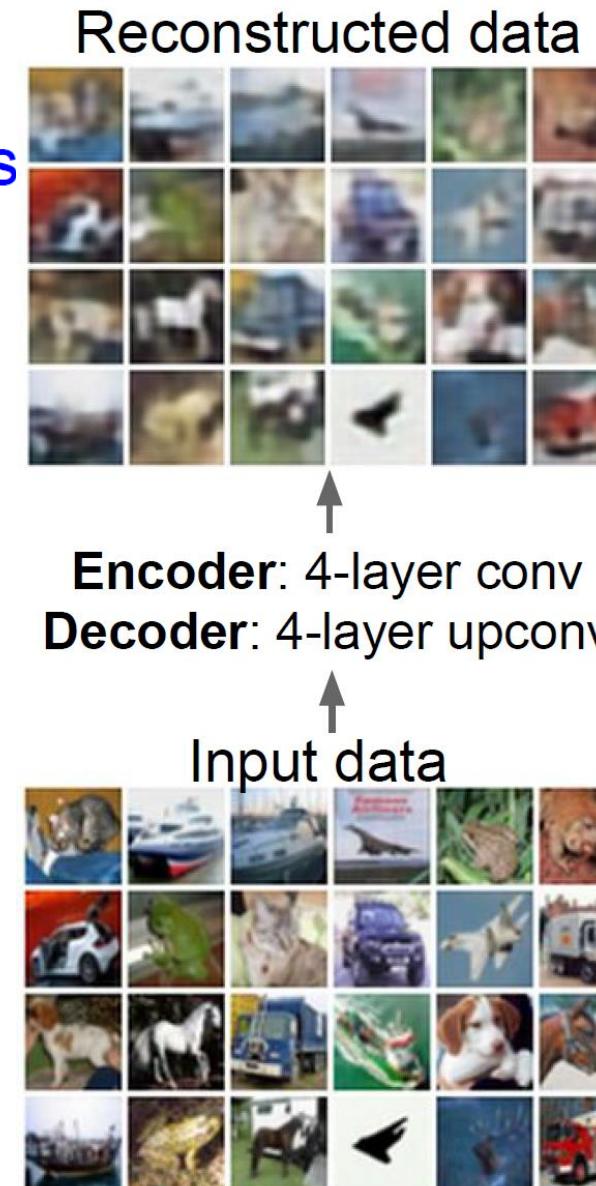
\hat{x}

Decoder

z

Encoder

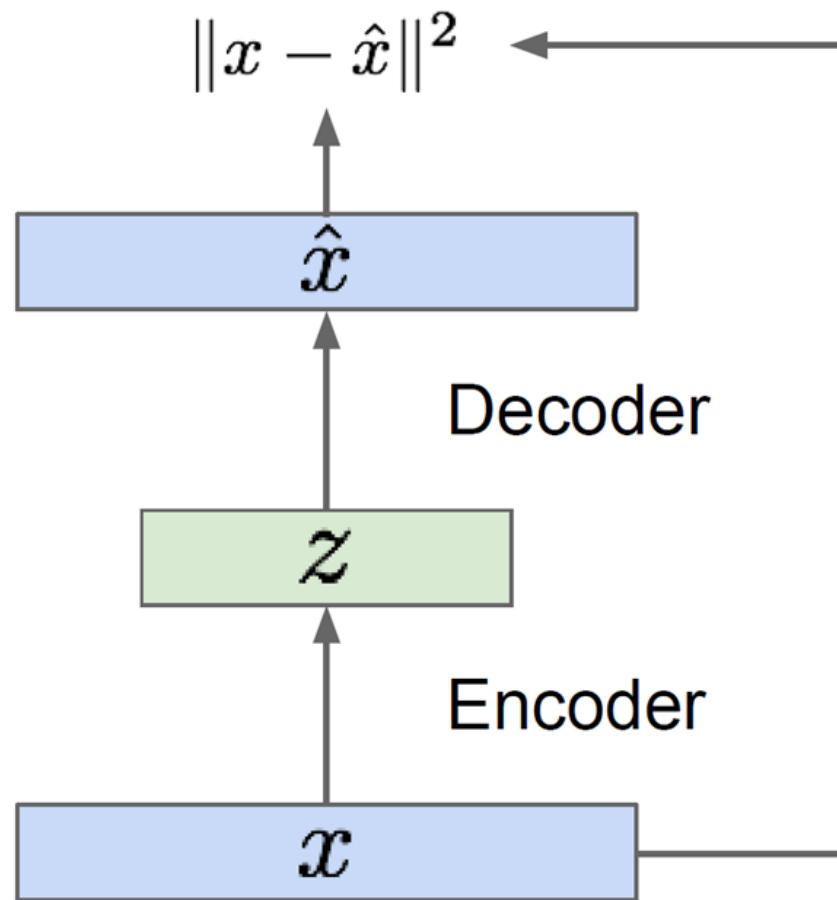
x



Variational auto-encoder

Auto-encoder vs VAE

L2 Loss function:



- Auto-encoder

$$q_\phi(z|x) : \text{NN encoder}$$

$$p_\theta(x|z) : \text{NN decoder}$$

- Variational auto-encoder

$$q_\phi(z|x) : \text{NN encoder}$$

$$p_\theta(z)$$

$$p_\theta(x|z) : \text{NN decoder}$$

→ Sampling new z'

$$p_\theta(x'|z') : \text{NN decoder}$$

New output x' which is different
from any x in the training set

VAE

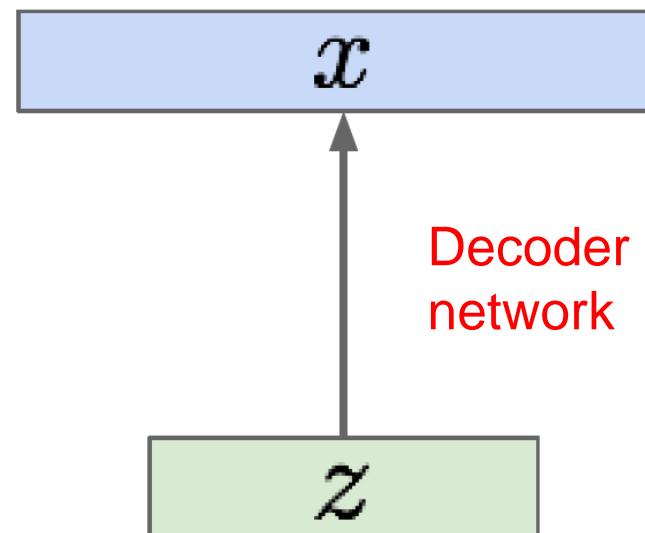
- Variation of auto-encoders enables us sample from the model to generate data!
- Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation \mathbf{z}

Sample from
true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample from
true prior

$$p_{\theta^*}(z)$$



- We want to estimate the true parameters θ^* of this generative model.
- No information on the true prior $p_{\theta^*}(z)$
→ choose prior $p(z)$ to be simple, e.g. Gaussian.
- Conditional $p(x|z)$ is complex (generates image) → represent with neural network

VAE

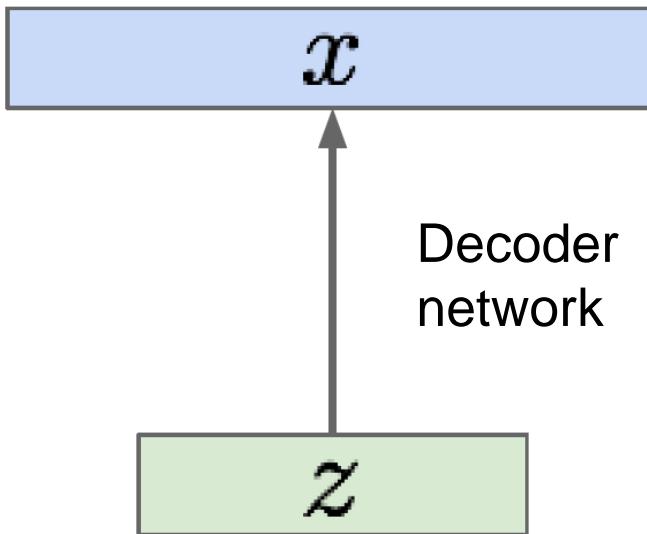
- We want to estimate the true parameters θ^* of this generative model.

Sample from
true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample from
true prior

$$p_{\theta^*}(z)$$



How to train the model?

Learn model parameters to maximize likelihood of training data

$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Now with latent z

$p_{\theta}(z)$: Gaussian prior

$p_{\theta}(x|z)$: NN decoder

Problem

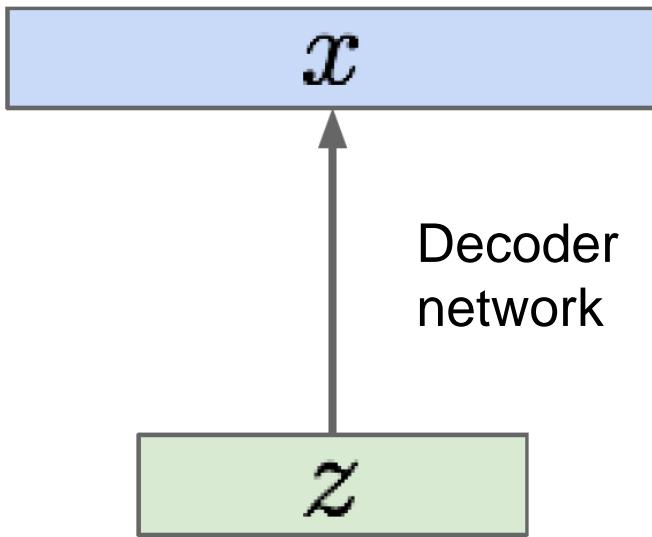
- We want to estimate the true parameters of this generative model.

Sample from
true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample from
true prior

$$p_{\theta^*}(z)$$



$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Now with latent z

$p_{\theta}(z)$: Gaussian prior

$p_{\theta}(x|z)$: NN decoder

Intractable integral over all possible z

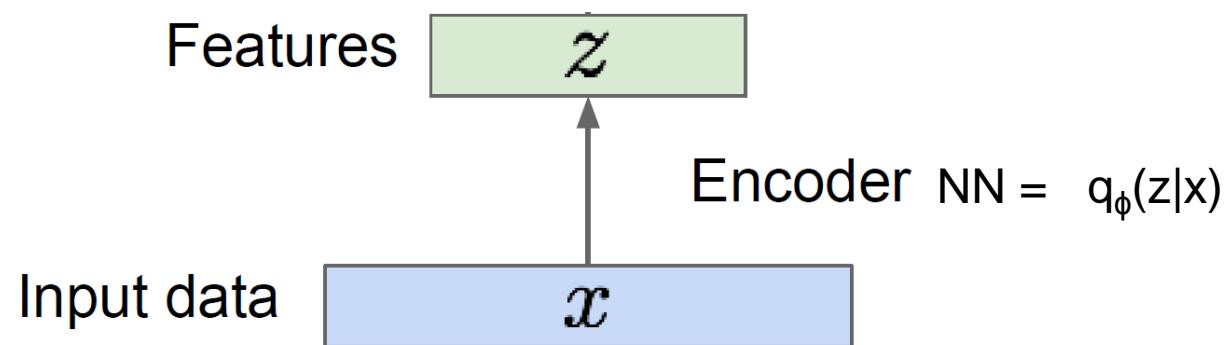
Posterior density is also intractable

$$p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$$

Intractable data likelihood

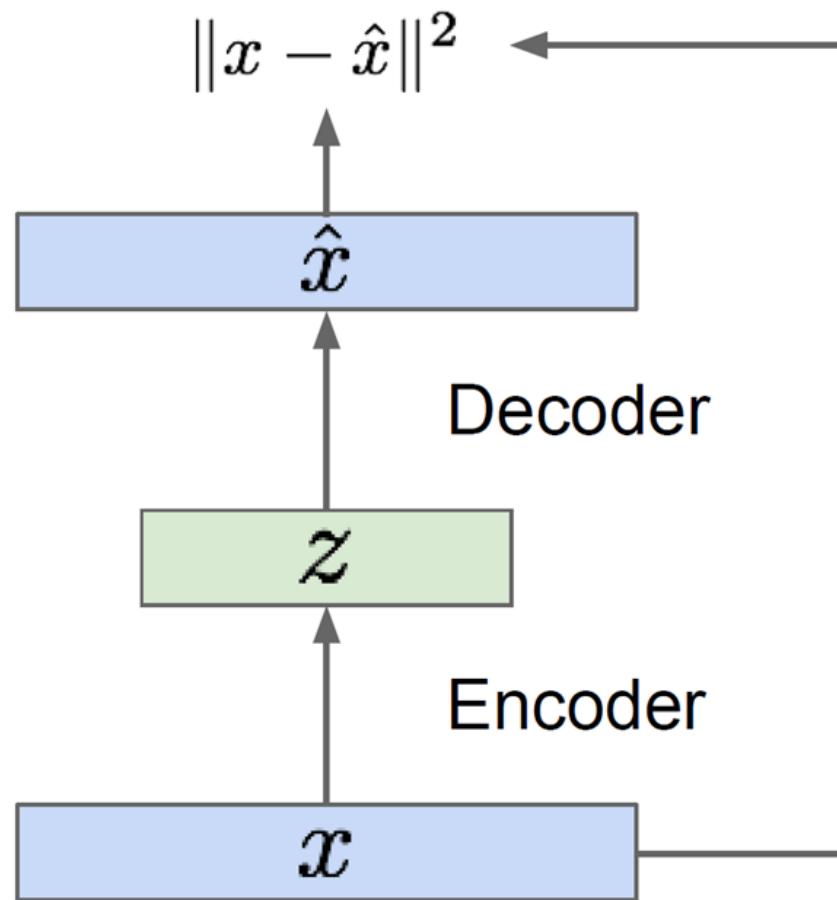
Solution

In addition to decoder network modeling $p_\theta(x|z)$, define additional encoder network $q_\phi(z|x)$ that approximates $p_\theta(z|x)$ to obtain z



Auto-encoder vs VAE

L2 Loss function:



- Auto-encoder

$$q_\phi(z|x) : \text{NN encoder}$$

$$p_\theta(x|z) : \text{NN decoder}$$

- Variational auto-encoder

$$q_\phi(z|x) : \text{NN encoder}$$

$$p_\theta(z) : \text{Gaussian}$$

$$p_\theta(x|z) : \text{NN decoder}$$

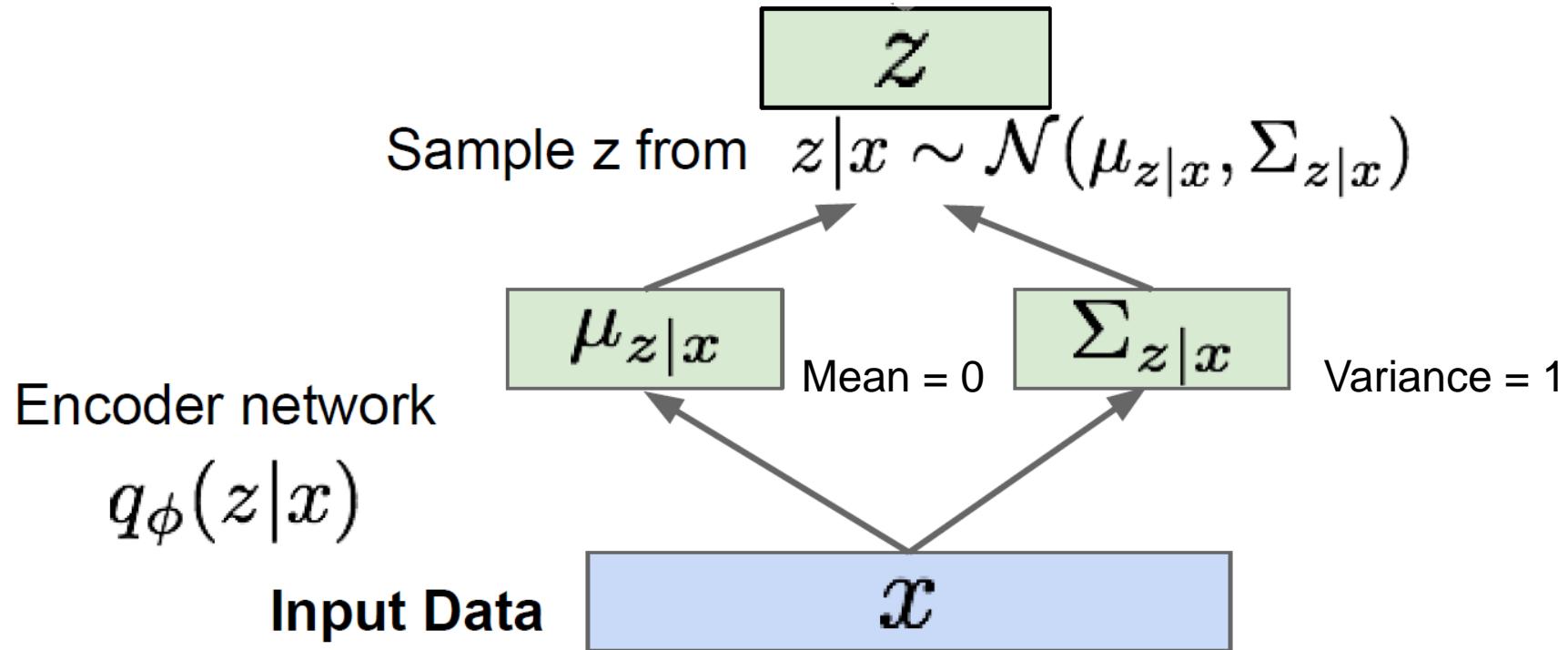
→ Sampling new z'

$$p_\theta(x'|z') : \text{NN decoder}$$

New output x' which is different
from any x in the training set

Latent vector with Gaussian prior

$p_\theta(z)$: Gaussian prior



Objective function of VAE

With our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)}))\end{aligned}$$

The expectation wrt. z (using encoder network) let us write nice KL terms

Objective function of VAE

With our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned}\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))\end{aligned}$$



Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick. see paper.)



This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!



$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :(But we know KL divergence always ≥ 0 .

Objective function of VAE

With our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned}\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0}\end{aligned}$$

Tractable lower bound which we can take gradient of and optimize! ($p_\theta(x|z)$ differentiable, KL term differentiable)

Objective function of VAE

With our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned}\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{> 0}\end{aligned}$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

Objective function of VAE

With our encoder and decoder networks, let's work out the (log) data likelihood:

Reconstruct the input data

$$\begin{aligned} \log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \end{aligned}$$

Make approximate posterior distribution close to prior

$$\begin{aligned} \text{L2 loss} \\ \text{Cross entropy} \\ \text{Softmax} &= \mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{> 0} \end{aligned}$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$

Variational lower bound (“ELBO”)

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

Molecular VAE

Reference paper

RESEARCH ARTICLE

Open Access



Molecular generative model based on conditional variational autoencoder for de novo molecular design

Jaechang Lim¹, Seongok Ryu¹, Jin Woo Kim¹ and Woo Youn Kim^{1,2*} 

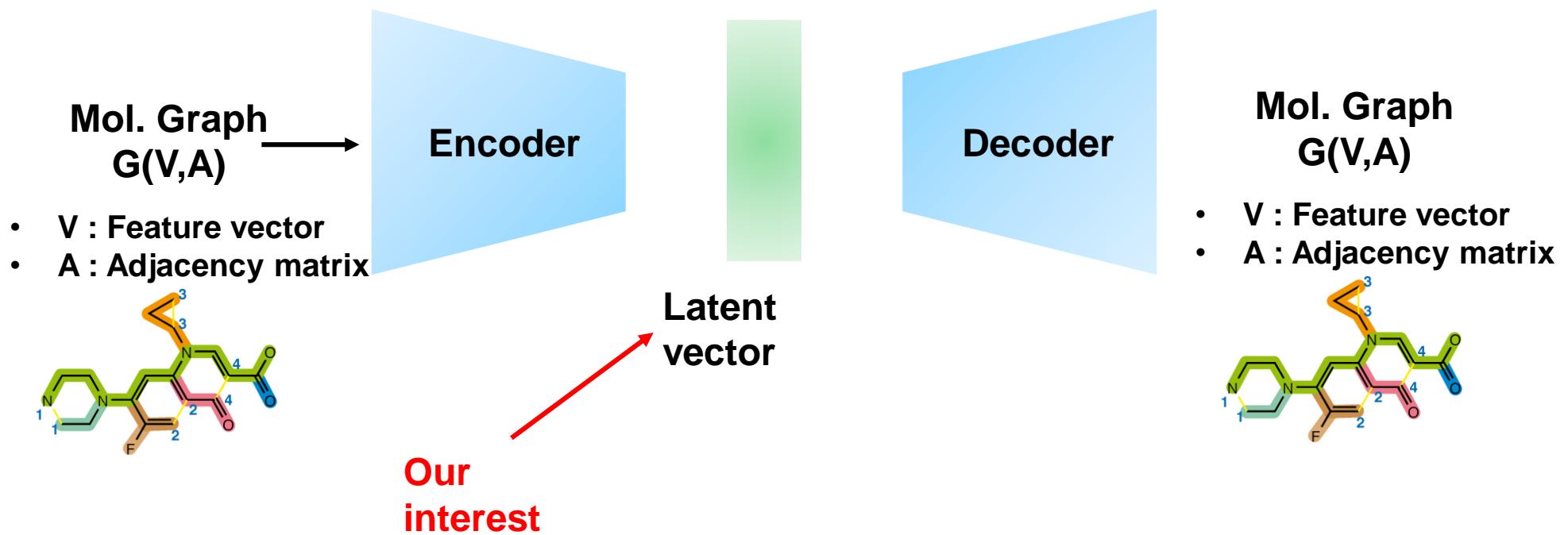
Abstract

We propose a molecular generative model based on the conditional variational autoencoder for de novo molecular design. It is specialized to control multiple molecular properties simultaneously by imposing them on a latent space. As a proof of concept, we demonstrate that it can be used to generate drug-like molecules with five target properties. We were also able to adjust a single property without changing the others and to manipulate it beyond the range of the dataset.

Keywords: Molecular design, Conditional variational autoencoder, Deep learning

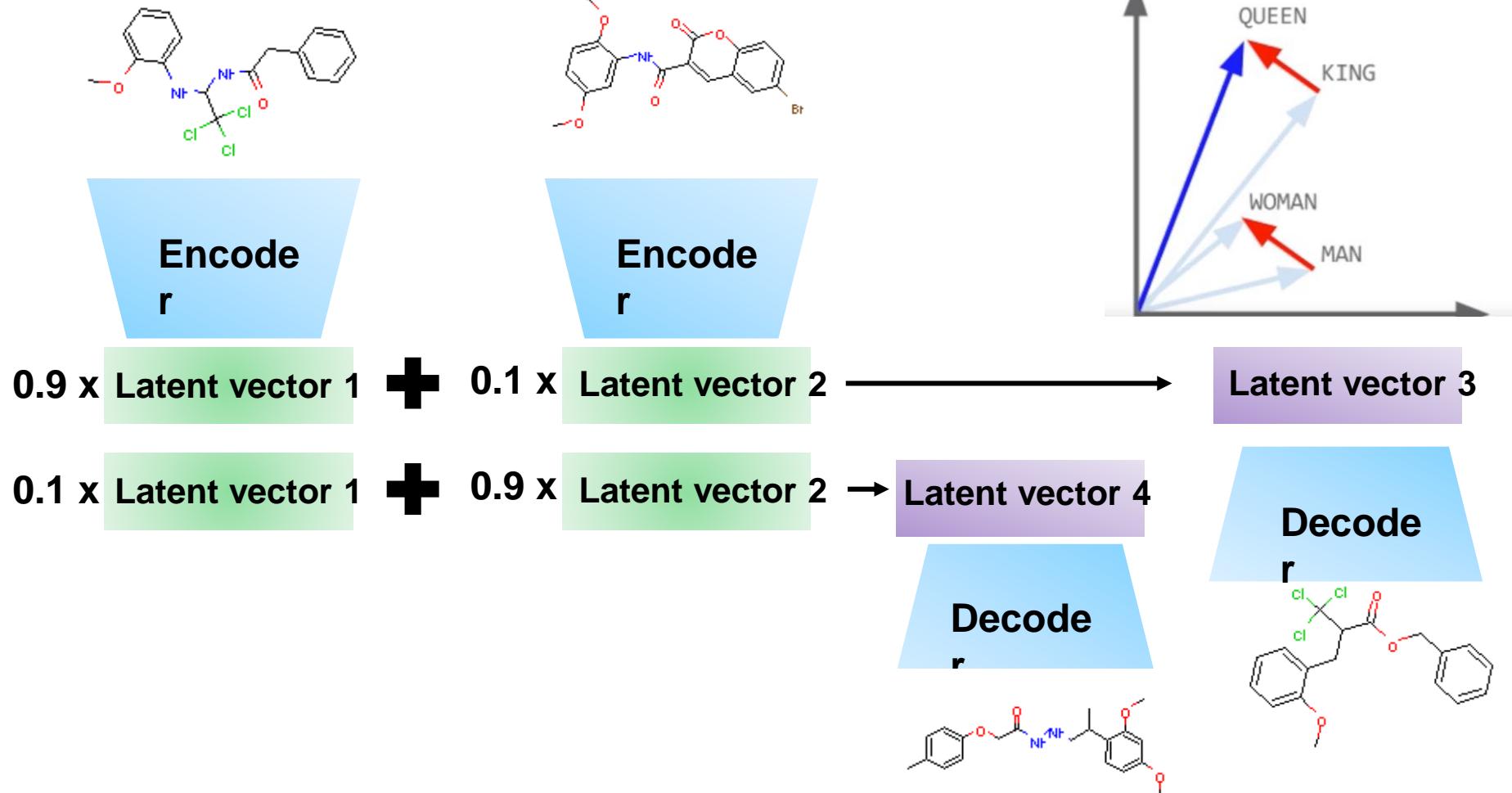
Molecular autoencoder

Molecular Structure Embedding via Autoencoder



Variational autoencoder (VAE)

✓ Embedding vector operation

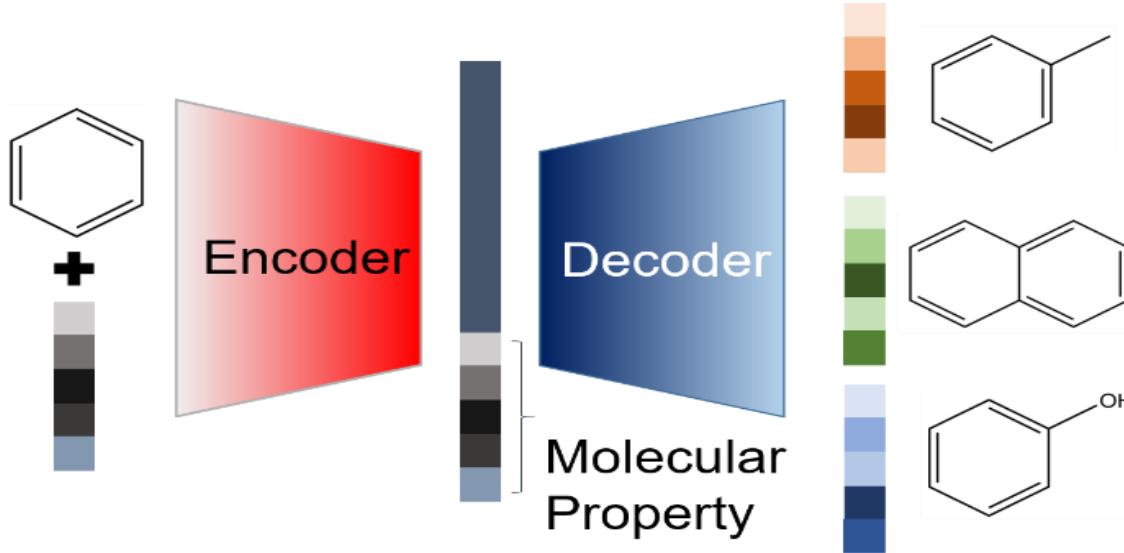


→ Represent molecules through vectors in a latent space.

→ Discontinuous graph/string representation to continuous vector representation

Conditional variational autoencoder (CVAE)

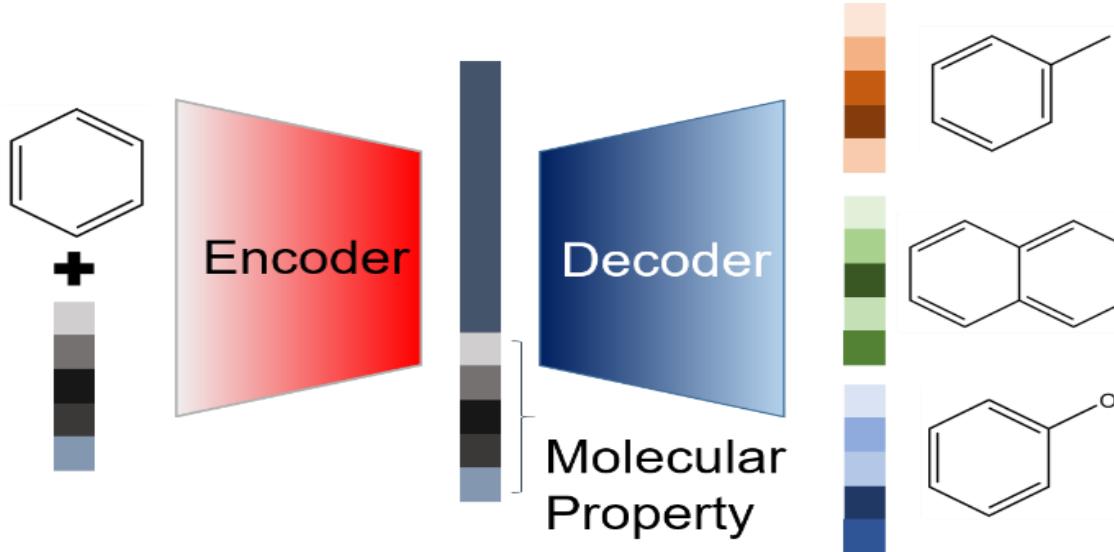
it “generates molecules with desired properties”.



- ✓ controlling several properties simultaneously
- ✓ embedding molecular properties into the condition vector
- ✓ Ex) Molecular weight, partition coefficient (LogP), HBD, HBA, TPSA
- ✓ 500,000 molecules randomly selected from the ZINC data set
- ✓ 400,000 for training and 100,000 for test

Conditional variational autoencoder (CVAE)

it “generates molecules with desired properties”.

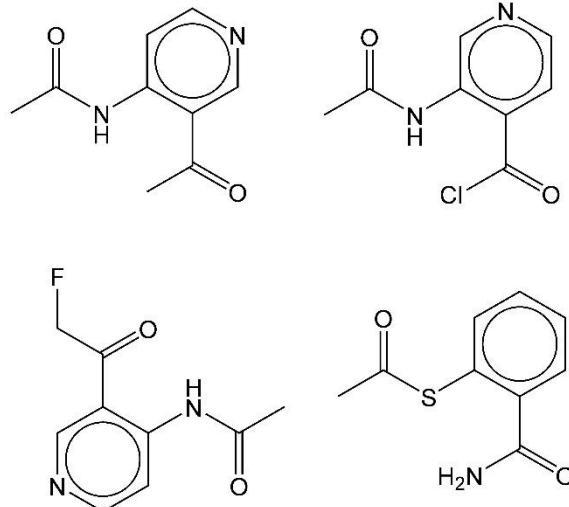


$$E[\log P(X|z)] - D_{KL}[Q(z|X) \parallel P(z)] \rightarrow E[\log P(X|z, c)] - D_{KL}[Q(z|X, c) \parallel P(z|c)]$$

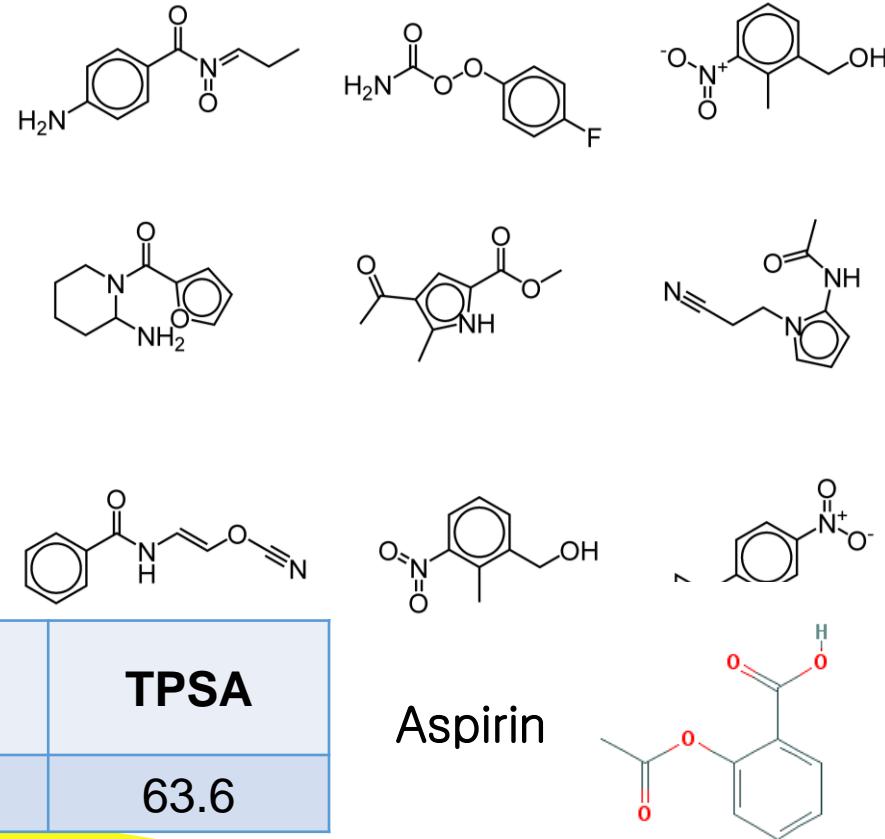
CVAE – Example 1

✓ 5 fixed target properties

A little change in molecular structure



A drastic change in molecular structure



Molecular weight	Log P	HBD	HBA	TPSA
180.04	1.31	1	3	63.6

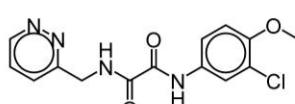
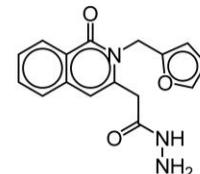
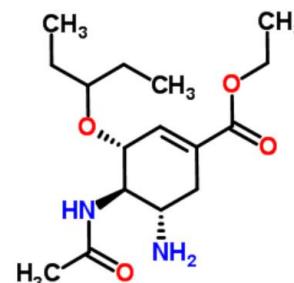
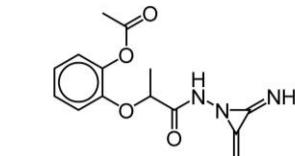
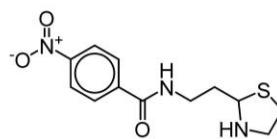
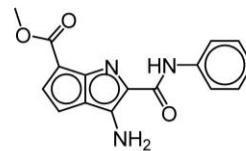
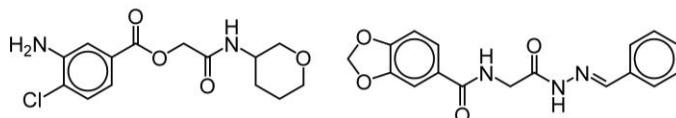


CVAE – Example 2

✓ 5 fixed target properties

Property of Tamiflu

Molecular weight	Log P	HBD	HBA	TPSA
312.2	1.285	2	5	90.64

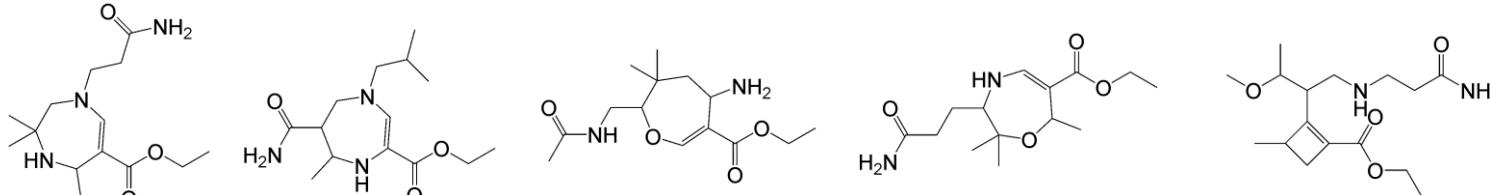


CVAE – Example 3

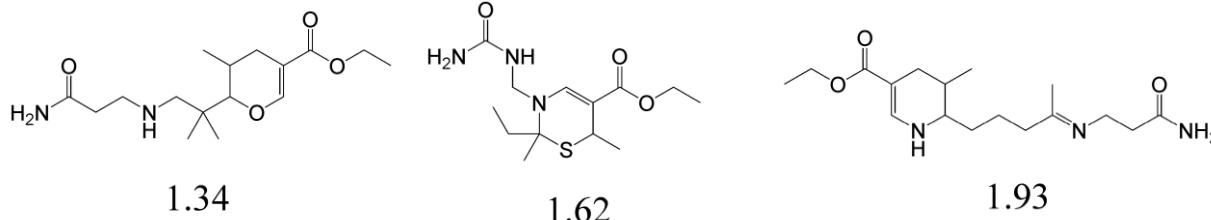
✓ Tuning one property while fixing the other four properties

Gradual change of LogP from 0.3 to 3.0 while keeping all other properties of Tamiflu fixed

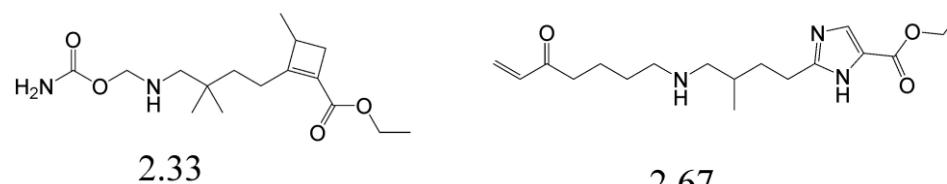
Molecular weight	Log P	HBD	HBA	TPSA
312.2	1.285	2	5	90.64



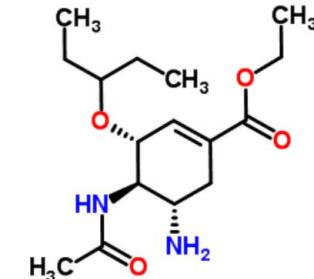
0.38 0.44 0.71



0.85 1.00
1.34 1.62 1.93



2.33 2.67



Jointly trained VAE

Cite This: ACS Cent. Sci. 2018, 4, 268–276

Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules

Rafael Gómez-Bombarelli,^{†,‡,§,¶,#} Jennifer N. Wei,^{‡,§,#} David Duvenaud,^{¶,#} José Miguel Hernández-Lobato,^{§,#} Benjamín Sánchez-Lengeling,[‡] Dennis Sheberla,^{‡,§,#} Jorge Aguilera-Iparraguirre,[†] Timothy D. Hirzel,[†] Ryan P. Adams,^{▽,||} and Alán Aspuru-Guzik^{*,‡,§,⊥,#}

[†]Kyulux North America Inc., 10 Post Office Square, Suite 800, Boston, Massachusetts 02109, United States

[‡]Department of Chemistry and Chemical Biology, Harvard University, Cambridge, Massachusetts 02138, United States

[¶]Department of Computer Science, University of Toronto, 6 King's College Road, Toronto, Ontario M5S 3H5, Canada

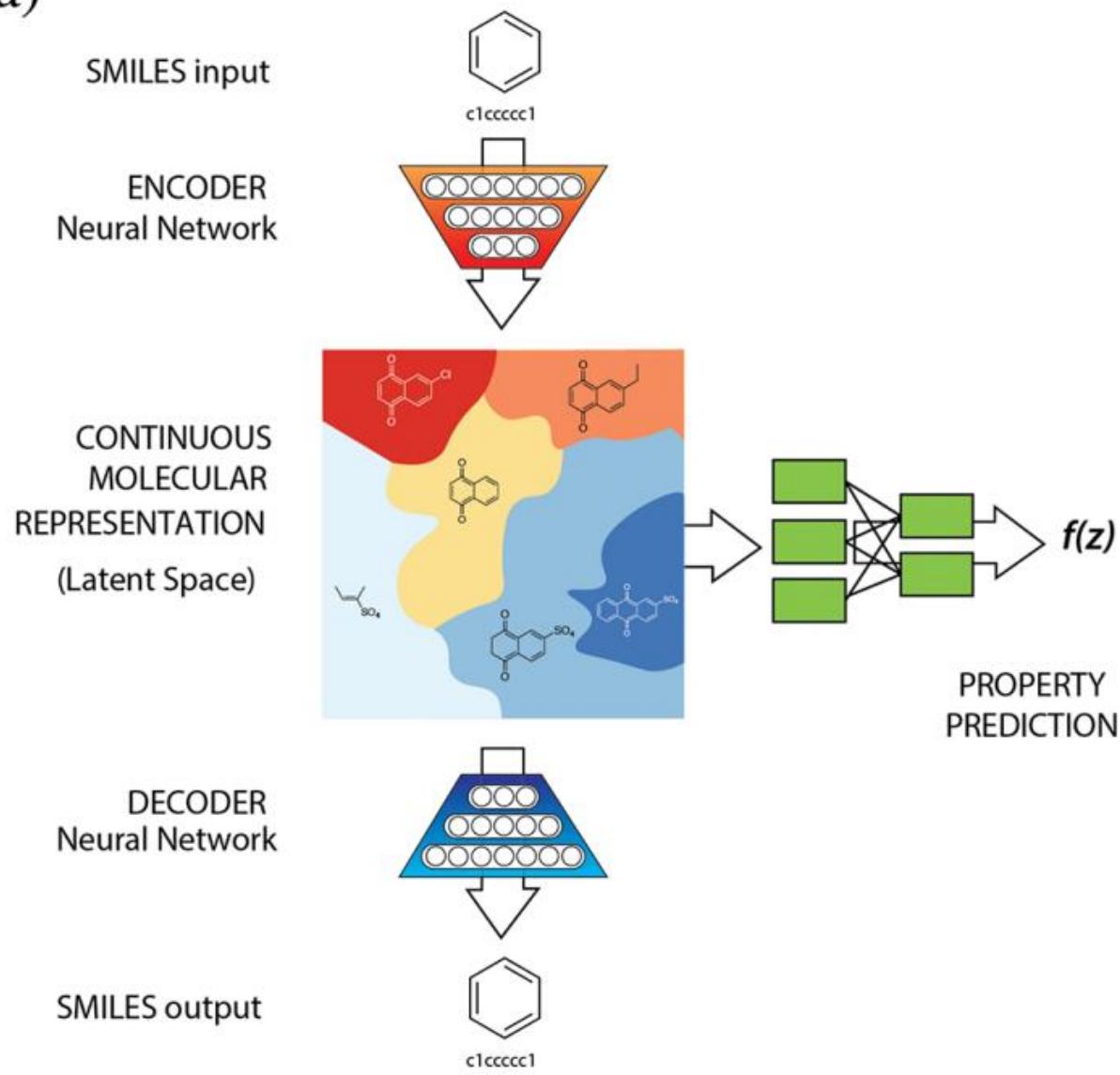
[§]Department of Engineering, University of Cambridge, Trumpington Street, Cambridge CB2 1PZ, U.K.

[▽]Google Brain, Mountain View, California, United States

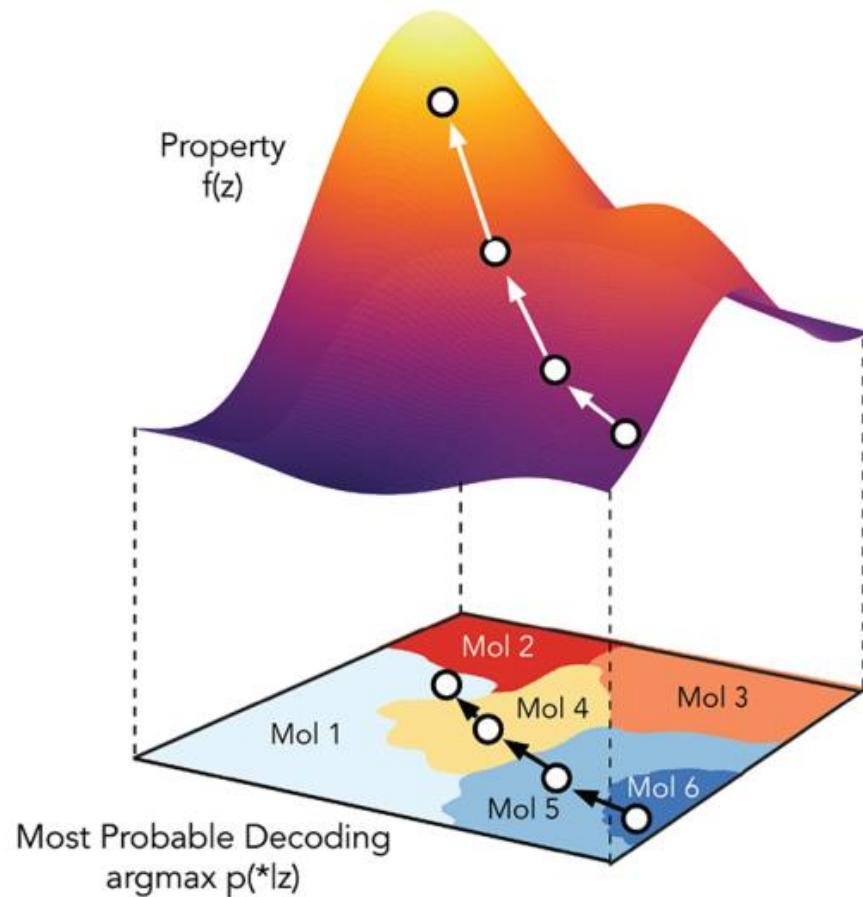
^{||}Princeton University, Princeton, New Jersey, United States

^{*}Biologically-Inspired Solar Energy Program, Canadian Institute for Advanced Research (CIFAR), Toronto, Ontario M5S 1M1, Canada

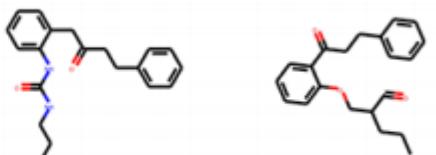
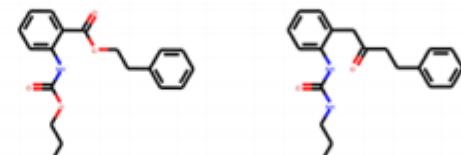
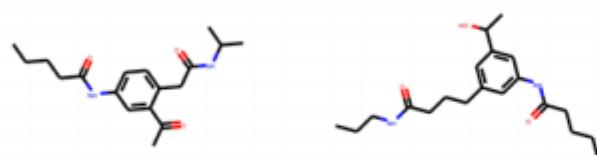
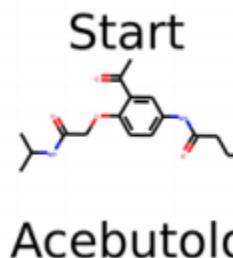
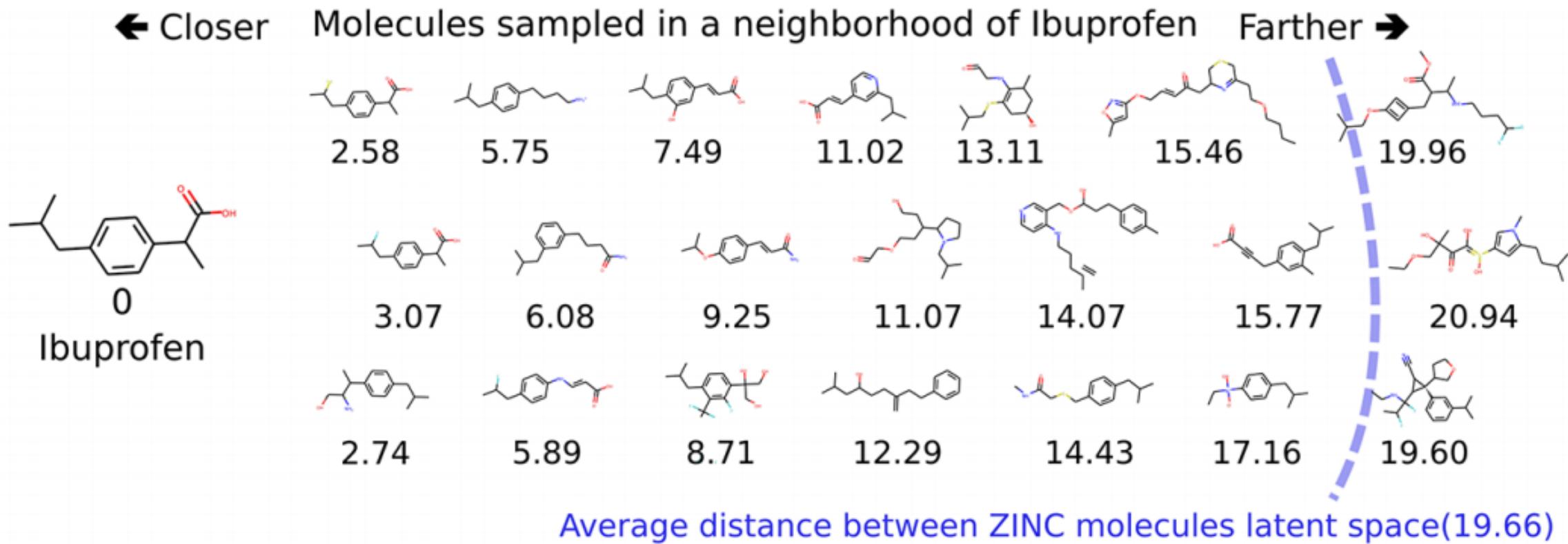
(a)



(b)

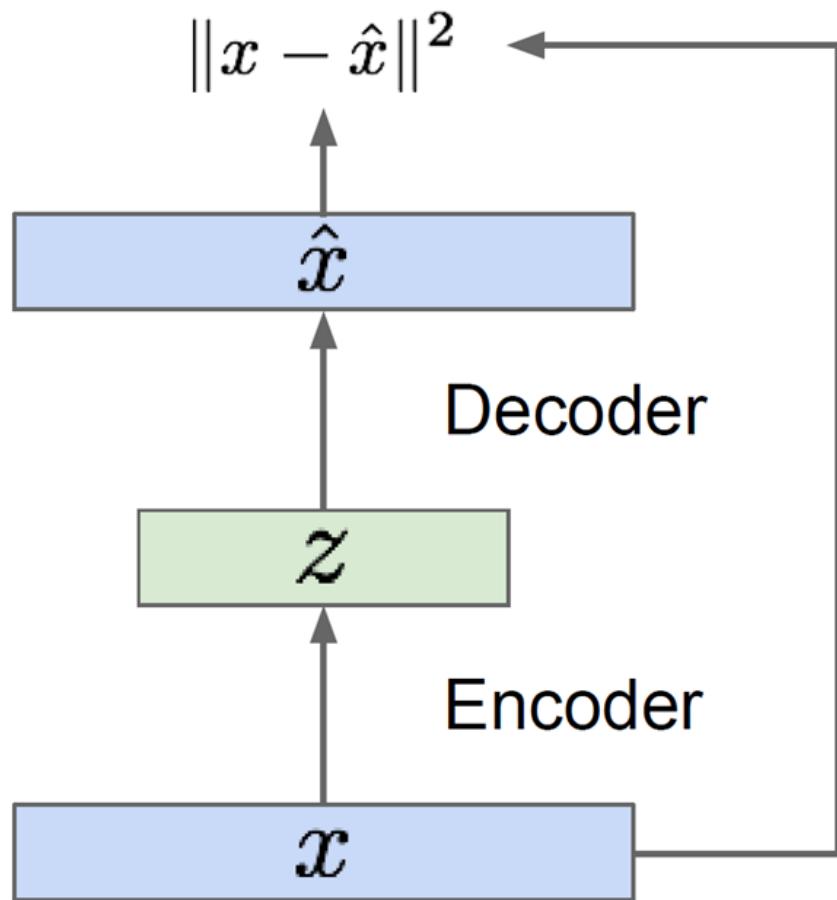


Sampling in latent space



Summary

L2 Loss function:



- Auto-encoder

$$q_\phi(z|x) : \text{NN encoder}$$

$$p_\theta(x|z) : \text{NN decoder}$$

- Variational auto-encoder

$$q_\phi(z|x) : \text{NN encoder}$$

$$p_\theta(z) : \text{Gaussian}$$

$$p_\theta(x|z) : \text{NN decoder}$$

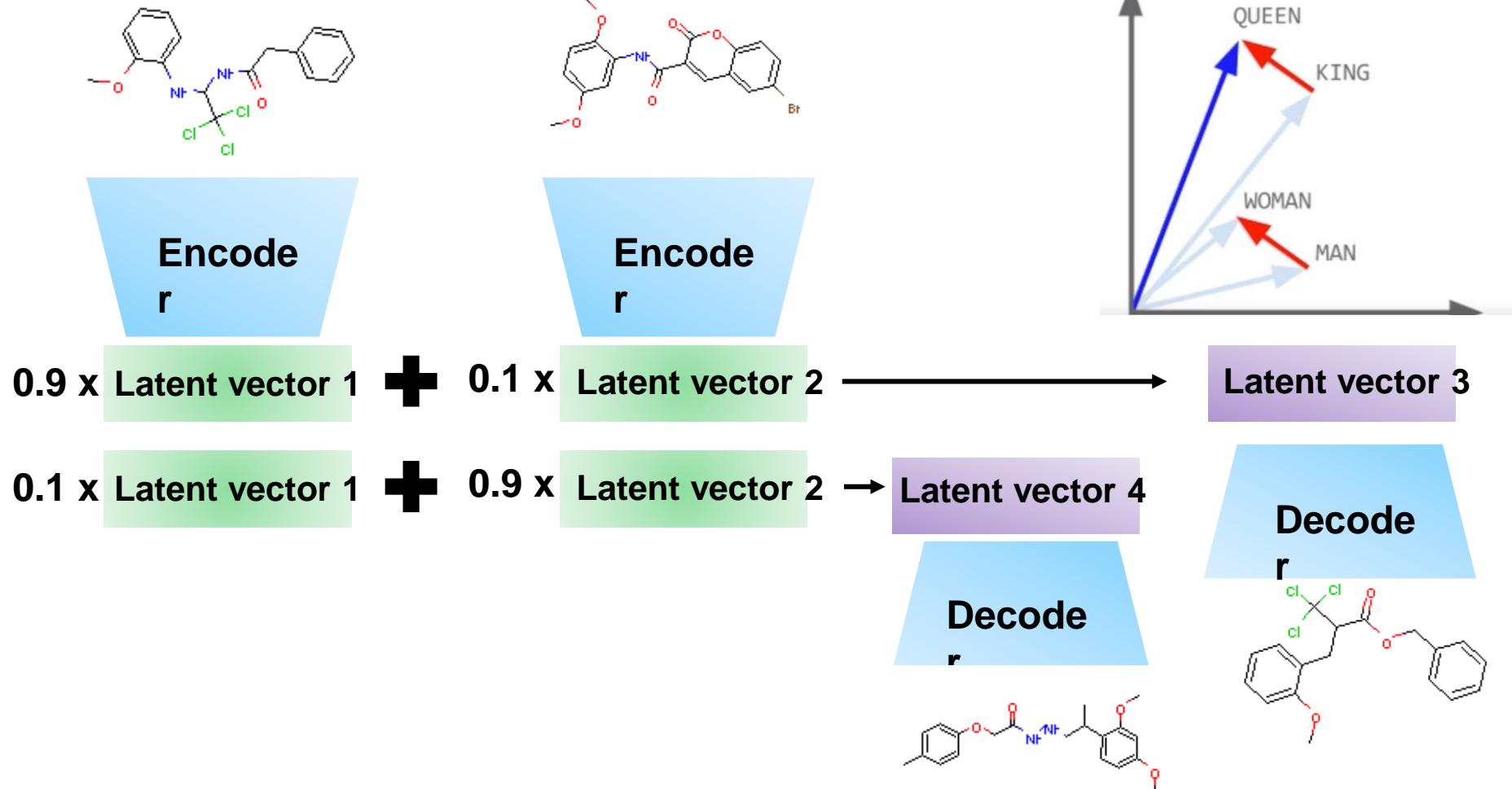
→ Sampling new z'

$$p_\theta(x'|z') : \text{NN decoder}$$

New output x' which is different
from any x in the training set

Summary

✓ Embedding vector operation



→ Represent molecules through vectors in a latent space.

→ Discontinuous graph/string representation to continuous vector representation

Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org
 Backfed Input Cell

 Input Cell

 Noisy Input Cell

 Hidden Cell

 Probabilistic Hidden Cell

 Spiking Hidden Cell

 Output Cell

 Match Input Output Cell

 Recurrent Cell

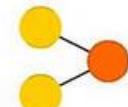
 Memory Cell

 Different Memory Cell

 Kernel

 Convolution or Pool

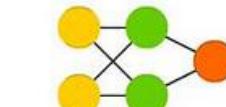
Perceptron (P)



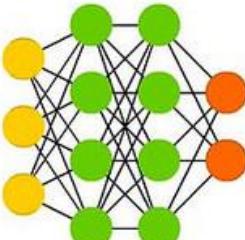
Feed Forward (FF)



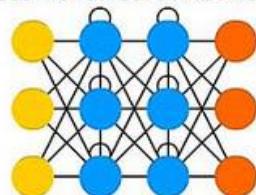
Radial Basis Network (RBF)



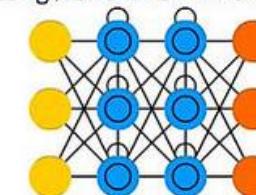
Deep Feed Forward (DFF)



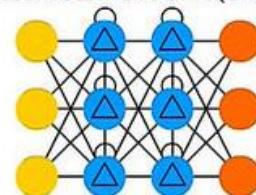
Recurrent Neural Network (RNN)



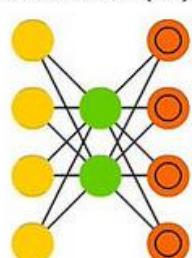
Long / Short Term Memory (LSTM)



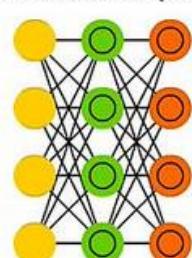
Gated Recurrent Unit (GRU)



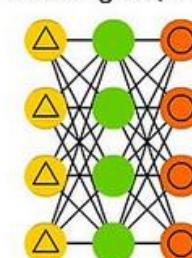
Auto Encoder (AE)



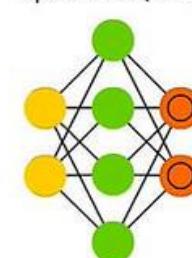
Variational AE (VAE)



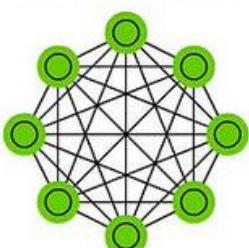
Denoising AE (DAE)



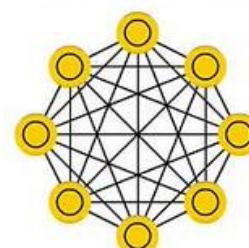
Sparse AE (SAE)



Markov Chain (MC)



Hopfield Network (HN)



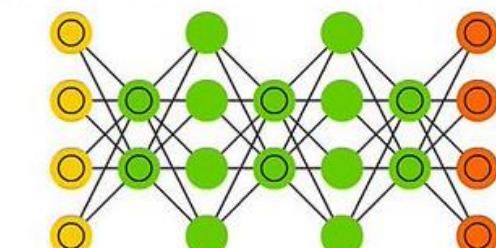
Boltzmann Machine (BM)



Restricted BM (RBM)



Deep Belief Network (DBN)



New terms

- Boltzmann machine
- Restricted Boltzmann machine
- Gibbs sampling
- Kullback–Leibler (KL) divergence
- Deep belief network
- Auto-encoder
- Variational auto-encoder (VAE)
- Conditional variational auto-encoder (CVAE)
- Jointly trained VAE