

Wenmin Wang



Principles of Machine Learning

Principles of Machine Learning

The Three Perspectives

 Springer

Part III Paradigms

- 8 Supervised Learning Paradigm
- 9 Unsupervised Learning Paradigm
- 10 Reinforcement Learning Paradigm
- 11 Other Learning Quasi-Paradigm

9 Unsupervised Learning Paradigm

Contents

- 9.1 Definition
- 9.2 Working Principle
- 9.3 Classic Tasks
- 9.4 Beyond the Classic
- 9.5 Energy Models
- 9.6 Autoencoding Models
- 9.7 Generative Adversarial Models
- 9.8 Normalizing Flow Models
- 9.9 Autoregressive Models
- 9.10 Autoregressive Models

What is Unsupervised Learning

Unsupervised learning is second one of three major paradigms of machine learning. It is fundamentally different from the supervised learning paradigm.

Definition: Unsupervised learning

Unsupervised learning (UL) is a learning paradigm to machine learning with unlabeled data, which analyzes or processes the data based on a certain established criterion to obtain desired results.

Several unsupervised tasks are existed, but since their relatively independent natures, it is hard to abstract common principles, as in supervised learning.

However, human learning is mainly unsupervised, therefore unsupervised machine learning becomes even more important.

What is Unsupervised Learning

It is a way of “self-learning” algorithm based on a certain given criterion, without a specific “teacher”.

Supervised learning vs. unsupervised learning

Supervised Learning	Unsupervised Learning
The training data for supervised learning algorithm is labeled.	All data for unsupervised learning algorithm is unlabeled.
The algorithm is trained based on labeled data.	Analysis or processing is carried out based on unlabeled data.

Formal Description

Unsupervised learning (UL) can be represented as $UL = \langle \mathcal{X}, \mathcal{Y}, H, \mathcal{L} \rangle$, a 4-tuple. Where: \mathcal{X} : input space, \mathcal{Y} : output space, H : hypothesis set, \mathcal{L} : loss function.

Using $P(\mathbf{x})$ from \mathcal{X} to obtain n independent and identically distributed (i.i.d.) data:

$$D = \{\mathbf{x}_i \mid \mathbf{x}_i \in \mathcal{X} \text{ and } i = 1, \dots, n\}.$$

Let target function of unsupervised learning be:

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

Hypothesis set H is the algorithm designed for a task of unsupervised learning:

$$H : \mathcal{X} \rightarrow \mathcal{Y} .$$

Formal Description

By directly analyzing D , find the best hypothesis $h \in H$, $h(\mathbf{x}) = \hat{y}$, so that the loss between the hypothesis $h(\mathbf{x})$ and the target function $f(\mathbf{x})$ is minimized:

$$\mathcal{L}(h(\mathbf{x}), f(\mathbf{x})) = \mathcal{L}(\hat{y}, y) = \arg \min_{h \in H} \mathbb{E}[h(\mathbf{x}) - f(\mathbf{x})].$$

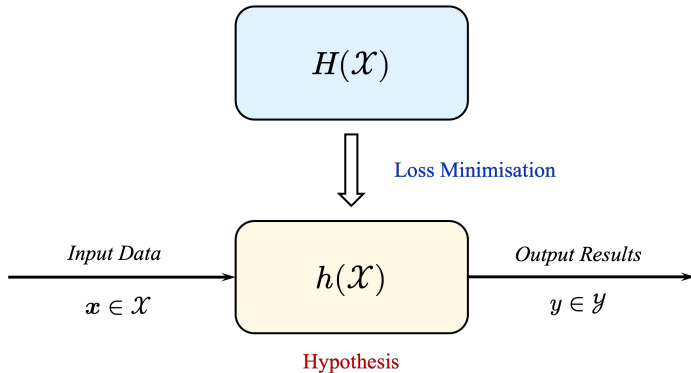
The loss function depends on specific task and its model, because it varies with the different tasks and models of unsupervised learning.

Here, $\mathcal{L}(\cdot, \cdot)$ is an abstract representation of the loss function, the purpose is to minimize the error between hypothesis and target function.

In other words, using the loss function to make the hypothesis as close as possible to its target function.

Illustrated Description

Unsupervised Learning Algorithm



A schematic diagram of unsupervised learning, it can be seen that it does not have a training process with labeled data.

Classic Tasks

The classic tasks of unsupervised learning mainly include following types.

Clustering:

Dividing data into several clusters based on the analysis of input data. It can be divided into classical and neo-classical clustering methods.

Dimensionality reduction:

Mapping high-dimensional data into low-dimensional space. It can be divided into linear, non-linear, and deep dimensionality reduction.

Association:

Discovering associative relationships between data objects in large databases and establishing corresponding association rules. The *shopping basket analysis* in supermarket is known as one of the examples.

Classic Tasks

Differences of the three tasks are reflected in their input and output spaces.

Clustering:

Output is k clusters $\{G_j \mid j = 1, \dots, k\} \subseteq \mathcal{Y}$, obtained by grouping input data. Where G_j is also a subset of input data, $G_j \subseteq \mathcal{X}$; and k is unknown in advance.

Dimensionality Reduction:

Output is new data after reducing the dimensionality of input data, $h : \mathcal{X} \rightarrow \mathcal{Y}$. Where, $\mathcal{X} \subseteq \mathbb{R}^h$ is high-dimensional input space, $\mathcal{Y} \subseteq \mathbb{R}^l$ is low-dimensional output space, $l \ll h$ and $1 \leq l \leq 3$.

Association:

Output is a set of association rules. Let $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$ be a set of items, if there are a subset of items $X \subseteq \mathcal{X}$ and a single item $Y \in \mathcal{X}$, then $X \Rightarrow Y \in \mathcal{Y}$.

Why Beyond

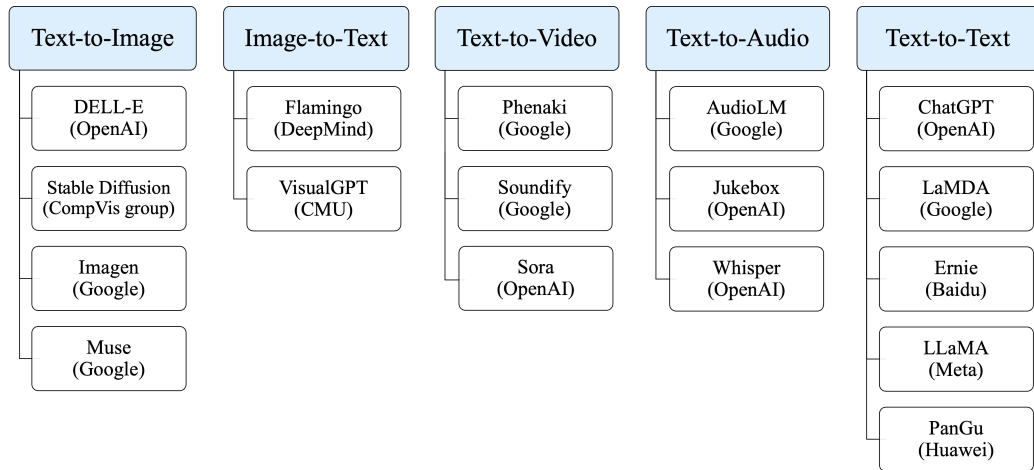
The paradigm of unsupervised learning should not be limited to the three classic tasks.

Because humans have a wide range of unsupervised learning abilities, and machines should also have relatively broad unsupervised learning functions.

With the continuous development of deep learning, unsupervised learning is increasingly attracting people's attention.

Some researchers combine deep neural networks with some other approaches in AI, ML, probability theory, statistics, physics, etc., and propose some emerging unsupervised models, namely *generative models*.

Representative Generative Models



About the Energy Models

Energy models, are also referred to as energy-based models (EBMs).

Be a type of generative model that introduces statistical mechanics into machine learning.

Two important properties in axioms of probability theory:

- Non-negativity: $p_{\theta}(\mathbf{x}) \geq 0$;
- Normalization: $\int p_{\theta}(\mathbf{x}) d\mathbf{x} = 1$.

Energy models can avoid the normalization problem, making the design of machine learning algorithms more flexible.

Energy models can provide a unified framework for many probabilistic and non-probabilistic learning.

Basic Energy Model

It refers to non-normalized probability scalar as energy, and characterizes its dependence by association of energy with input and latent variables.

Inference in the model includes finding the latent variables (values) that minimize energy given a set of input variables (values).

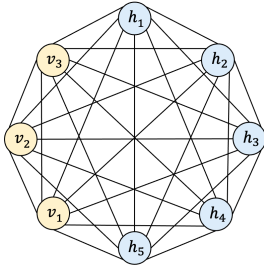
It learns a function that associates low energy with the correct values of latent variables and high energy with incorrect values.

Without loss of generality, here we only introduce the single-variable energy model:

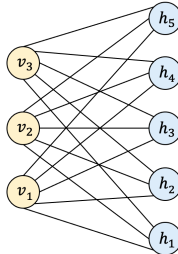
$$p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}(\mathbf{x})}, \quad Z_{\theta}(\mathbf{x}) = \int e^{-E_{\theta}(\mathbf{x})} d\mathbf{x}.$$

Where $E_{\theta}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}$ is an energy function, and $Z_{\theta}(\mathbf{x})$ is normalizing constant.

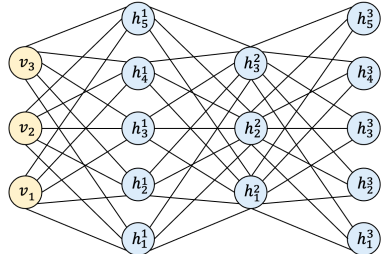
Boltzmann Models



(a)



(b)



(c)

The three types of Energy-based Boltzmann models.

(a) Boltzmann machine (BM), (b) restricted Boltzmann machine (RBM), and (c) deep Boltzmann machine (DBM).

Boltzmann Models

Boltzmann Machine:

A Boltzmann machine (BM) is a undirected network of symmetrically connected.

It includes a set of visible units $\mathbf{v} \in \{0, 1\}^D$, and a set of hidden units $\mathbf{h} \in \{0, 1\}^P$, where D and P denote the numbers of visible and hidden units, respectively.

The energy of the state $\{\mathbf{v}, \mathbf{h}\}$ is defined as follows:

$$E_{\theta}(\mathbf{v}, \mathbf{h}) = -\frac{1}{2}\mathbf{v}^T \mathbf{L} \mathbf{v} - \frac{1}{2}\mathbf{h}^T \mathbf{J} \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}.$$

Where $\theta = \{\mathbf{W}, \mathbf{L}, \mathbf{J}\}$ is the parameter, and $\{\mathbf{W}, \mathbf{L}, \mathbf{J}\}$ are three symmetric interaction terms: \mathbf{W} is visible-to-hidden, \mathbf{L} is visible-to-visible, and \mathbf{J} is hidden-to-hidden.

It is trained by contrastive divergence and finds its equilibrium state through Gibbs sampling. But, this requires an exponential time in terms of the number of units.

Boltzmann Models

Restricted Boltzmann Machine:

A restricted Boltzmann machine (RBM) is a special type of Boltzmann machine.

The restriction refers to limiting the connections between visible and hidden units to solve the exponential time problem in BM.

RBM becomes a two-layer undirected network: one is visible, and the other is hidden.

Its state $\{\mathbf{v}, \mathbf{h}\}$ is defined by the following energy:

$$E_{\theta}(\mathbf{v}, \mathbf{h}) = -\mathbf{v}^T \mathbf{W} \mathbf{h}.$$

Where $\theta = \{\mathbf{W}\}$, i.e., there is only visible to hidden, and no other two items in BM.

RBM can be trained using the maximum likelihood method. The sampling can be done by Gibbs sampling or Markov chain Monte Carlo (MCMC) methods.

Boltzmann Models

Deep Boltzmann Machine:

A deep Boltzmann machine (DBM) is an undirected deep network with one visible layer and multiple hidden layers.

The DBM is equivalent to adding several hidden layers to the RBM.

For the DBM that has one visible layer and three hidden layers, its state $\{\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3\}$ is defined by the following energy:

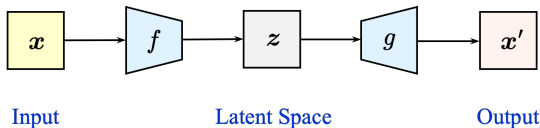
$$E_{\theta}(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3) = -\mathbf{v}^{\top} \mathbf{W}^1 \mathbf{h}^1 - \mathbf{h}^{1\top} \mathbf{W}^2 \mathbf{h}^2 - \mathbf{h}^{2\top} \mathbf{W}^3 \mathbf{h}^3.$$

Where the parameter $\theta = \{\mathbf{W}^1, \mathbf{W}^2, \mathbf{W}^3\}$ denotes the three symmetric interaction terms from visible to hidden, hidden to hidden, and hidden to hidden, respectively.

Autoencoders

Autoencoders:

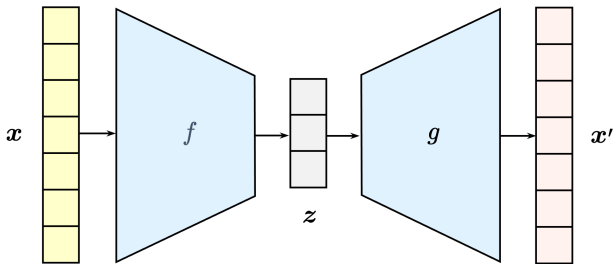
Let $\mathcal{X} \subseteq \mathbb{R}^u$ be a u -dimensional space, $\mathcal{Z} \subseteq \mathbb{R}^v$ be a v -dimensional latent space, $x, x' \in \mathcal{X}$ are input and output data respectively, and $z \in \mathcal{Z}$ is intermediate code in the latent space. Let f and g denote encoder and decoder. The f encodes x to z , i.e., $z = f(x)$; and the g decodes z into x' , i.e., $x' = g(z)$. Therefore, a neural network is known as an autoencoder (AE), if it satisfies $x' = g(z) = g(f(x))$ and $x' \neq x$.



Undercomplete Autoencoders

Undercomplete Autoencoders:

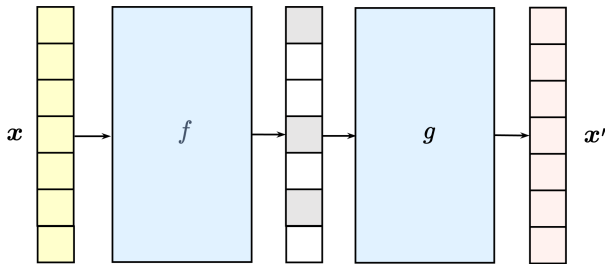
It is called an undercomplete autoencoder (UAE), iff the dimension of its latent space $\mathcal{Z} \subseteq \mathbb{R}^v$ is less than the dimension of the data space $\mathcal{X} \subseteq \mathbb{R}^u$, i.e., $v < u$. Its loss function can be expressed as $\mathcal{L}(x, x') = \mathcal{L}(x, g(f(x)))$.



Sparse Autoencoders

Sparse Autoencoders:

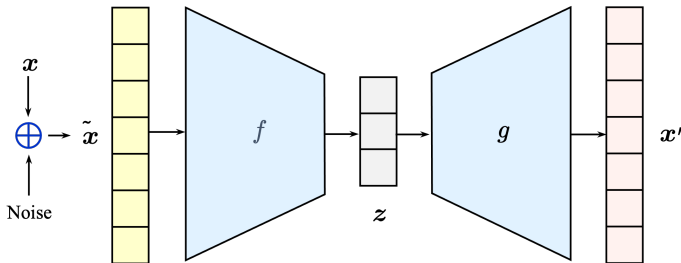
A sparse autoencoder (SAE) is an AE with a sparsity constraint, but the dimension of latent space can be equal to or even greater than input space. Its loss function is $\mathcal{L}(x, x') + \Omega(z) = \mathcal{L}(x, g(f(x))) + \Omega(z)$, where $\Omega(z)$ is a sparsity penalty term, can be adopted by KL divergence, L^1 or L^2 regularization.



Denoising Autoencoders

Denoising Autoencoders:

A denoising autoencoder (DAE) learns useful features in the input data by changing the reconstruction error term in the loss function: $\mathcal{L}(x, x') = \mathcal{L}(x, g(f(\tilde{x})))$, where $\tilde{x} = x \oplus \text{noise}$.



Variational Autoencoders

Variational Autoencoders:

A variational autoencoder (VAE) is an ANN that uses probabilistic graphical models and variational Bayes method.

Generally, we need a probabilistic encoder and a probabilistic decoder, i.e.,

$$\mathbf{z} \sim p_{\theta}(\mathbf{z}|\mathbf{x}) \quad \text{and} \quad \mathbf{x}' \sim p_{\theta}(\mathbf{x}|\mathbf{z}).$$

According to Bayes' theorem, the encoder should be solved by:

$$p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\theta}(\mathbf{z})}{p_{\theta}(\mathbf{x})}, \quad p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}.$$

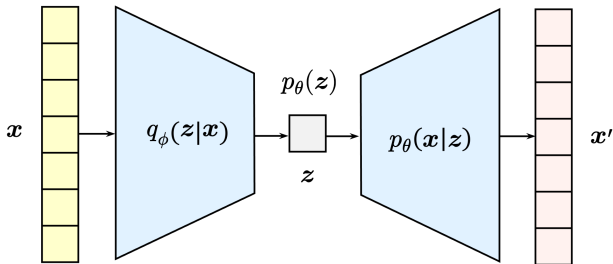
Since $p_{\theta}(\mathbf{x})$ is calculated by integrating \mathbf{z} in latent space, which requires exponential time. So, $p_{\theta}(\mathbf{z}|\mathbf{x})$ is a *computational intractability problem*.

Variational Autoencoders

VAE uses *variational inference* $q_{\phi}(z|x) \approx p_{\theta}(z|x)$ to obtain the approximate result.

Therefore, the $q_{\phi}(z|x)$ is treated as probabilistic encoder, i.e.,

$$z \sim q_{\phi}(z|x) \quad \text{and} \quad x' \sim p_{\theta}(x|z).$$

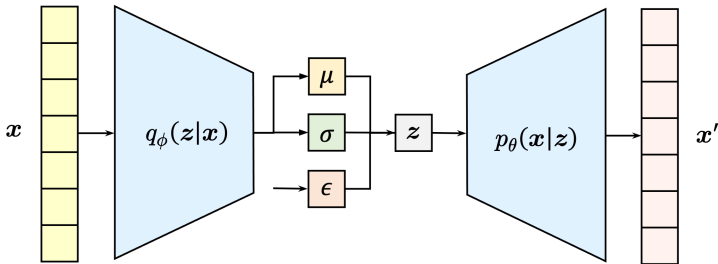


Variational Autoencoders

Re-Parameterization Trick:

It is a differentiable transformation of z into ϕ , \mathbf{x} , and ϵ , that is $z = f(\phi, \mathbf{x}, \epsilon)$. Suppose $z \sim q_\phi(z|\mathbf{x}) = \mathcal{N}(\mu, \sigma^2)$. Given $\epsilon \sim \mathcal{N}(0, 1)$, then by re-parameterization:

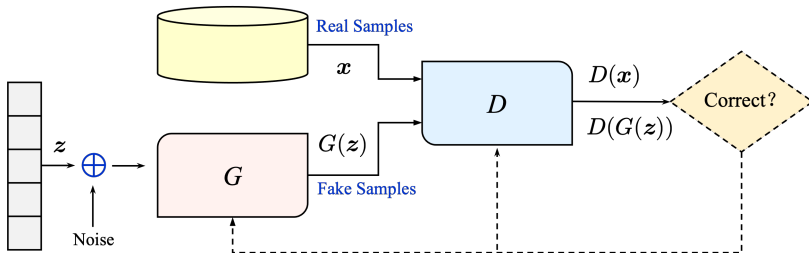
$$z = \mu + \sigma \odot \epsilon.$$



Generative Adversarial Models

Generative Adversarial Networks:

Generative adversarial network (GAN) is based on the minimax theorem in zero-sum game, a generator and a discriminator operate in a manner of two player game.



The structure of GAN, consists of a generator (G) and a discriminator (D), both of which use multi-layer neural networks.

Generative Adversarial Networks

The training process:

G first generates a fake sample $G(z)$, then D takes a real sample x and $G(z)$ as inputs, and obtains $D(x)$ and $D(G(z))$ respectively.

The G tries to make $G(z)$ true, while the D tries to make $D(x)$ true and $D(G(z))$ false.

Since GAN is based on the minimax theorem, that is:

$$\min_G \max_D V(G, D) = E_{x \sim P_{\text{data}}(x)} [\log D(x)] + E_{z \sim P_z(z)} [\log (1 - D(G(z)))] .$$

The training process of these two modules is like a game between two players. After multiple rounds of adversarial training, their performance is continuously optimized.

When the D finally cannot distinguish the difference between $D(x)$ and $D(G(z))$, they reach a zero-sum state, it is considered that the GAN has reached its optimum.

Normalizing Flow Models

Normalizing Flow Models:

A normalized flow model uses the change of variables in probability to transform complex distributions into simple (normal) distributions, and then performs inverse transformations to restore simple (normal) to complex distributions.

Change of Variables:

It refers to a basic technique in mathematics to simplify problems, in which the original variables are replaced with functions of other variables.

Let $\mathbf{z} \in \mathbb{R}^D$ be a latent variable, $p_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^D$ be its probability density function: $\mathbf{z} \sim p_\theta(\mathbf{z})$. Where $p_\theta(\mathbf{z})$ is a simple probability distribution, e.g. Gaussian distribution.

Using invertible function $g_\theta(\mathbf{z})$ to generate $\mathbf{x} \in \mathbb{R}^D$, and through its inverse function $g_\theta^{-1}(\mathbf{x})$ to transform \mathbf{x} back to \mathbf{z} , i.e., $\mathbf{x} = g_\theta(\mathbf{z})$ and $\mathbf{z} = g_\theta^{-1}(\mathbf{x})$.

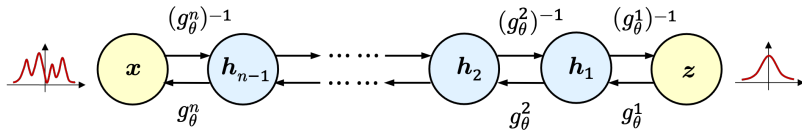
Normalizing Flow Models

Let probability density function $g_{\theta} = g_{\theta}^n \circ \dots \circ g_{\theta}^2 \circ g_{\theta}^1$, then $\mathbf{x} = g_{\theta}(\mathbf{z})$ can be written as:

$$\mathbf{x} = g_{\theta}(\mathbf{z}) = g_{\theta}^n \circ \dots \circ g_{\theta}^2 \circ g_{\theta}^1(\mathbf{z}).$$

Similarly, let the inverse probability density function $g_{\theta}^{-1} = (g_{\theta}^1)^{-1} \circ (g_{\theta}^2)^{-1} \circ \dots \circ (g_{\theta}^n)^{-1}$, then $\mathbf{z} = g_{\theta}^{-1}(\mathbf{x})$ will be:

$$\mathbf{z} = g_{\theta}^{-1}(\mathbf{x}) = (g_{\theta}^1)^{-1} \circ (g_{\theta}^2)^{-1} \circ \dots \circ (g_{\theta}^n)^{-1}(\mathbf{x}).$$



Autoregressive Models

Autoregressive Property:

Given a set of variables $\{x_1, x_2, \dots, x_t\}$ over time t , it is called autoregressive property if x_t depends linearly on some or all its previous variables $\{x_1, x_2, \dots, x_{t-1}\}$.

k -Order Autoregressive Models:

A k -order autoregressive model $AR(k)$ is a statistical model with the autoregressive property, which is defined as follows:

$$x_t = \sum_{i=1}^k \varphi_i x_{t-i} + \varepsilon_t.$$

Where, $\varphi_1, \dots, \varphi_k$ are the parameters of the model, and ε_t is white noise.

And, if $k = t$, then $x_t = (x_1 + x_2 + \dots + x_{t-1}) + \varepsilon_t$.

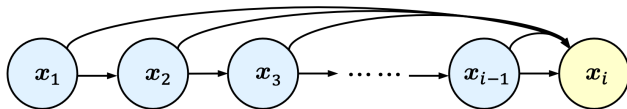
Autoregressive Models

Deep Autoregressive Models:

A deep autoregressive model (DAM) is a deep network model with autoregressive property, can be modeled based on the chain rule of joint probability.

Given a set of random variables $\{x_1, x_2, \dots, x_t\}$ over time t , and $p_\theta(\cdot)$ denotes a probability density function, then its joint probability can be expressed as:

$$p_\theta(x_1, x_2, \dots, x_t) = \prod_{i=1}^t p_\theta(x_i | x_1, x_2, \dots, x_{i-1}).$$



The directed graph of $p_\theta(x_i | x_1, x_2, \dots, x_{i-1})$.

Diffusion Models

Diffusion models are inspired by diffusion principle in nonequilibrium thermodynamics.

Diffusion:

It is the net process of movement of molecules or particles under a concentration gradient, generally from a region of higher concentration to a region of lower concentration.

E.g., a tea bag immersed in a cup of hot water will diffuse into the water and change its color.

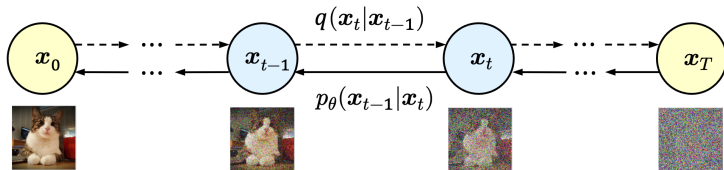
Diffusion Probabilistic Models:

A diffusion probabilistic model is a type of latent variable generative model, that consists of the forward diffusion process to slowly add random noise to data and the reverse diffusion process to reconstruct desired data from the noise.

Denoising Diffusion Models

Denoising Diffusion Probability Models:

A denoising diffusion probability model uses a Markov chain that gradually adds Gaussian noise to original image to obtain an approximate posterior probability $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$. The joint probability distribution $p_\theta(\mathbf{x}_{0:T})$ is its reverse process, learning a Gaussian transition, starting from $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$.



Schematic diagram of the denoising diffusion probability model.

Denoising Diffusion Processes

Forward Diffusion Process:

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \prod_{t=1}^T \mathcal{N}\left(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}\right).$$

Reverse Diffusion Process:

$$p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=T}^1 p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = p(\mathbf{x}_T) \prod_{t=T}^1 \mathcal{N}\left(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sum_{\theta}(\mathbf{x}_t, t)\right).$$

Thank You