

Approximate Inference (近似推論)

Chapter 19

楊 森 (Sen Yang) 2019/08/18

Inferenceは一部変数を与えられて、残る部分変数の確率分布を計算すること（つまり条件付き分布を計算すること）。この中に、潜在変数を含む確率モデルをトレーニングする場合、 $p(h|v)$ の計算は中心課題です。

Challenge of Inference (推論)

$p(h | v)$ の計算もしくは $p(h | v)$ 分布での期待値の計算

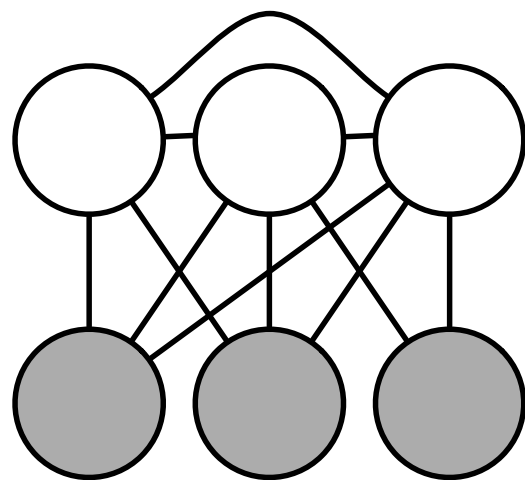
潜在変数
latent
variable

可視変数
visible
variable

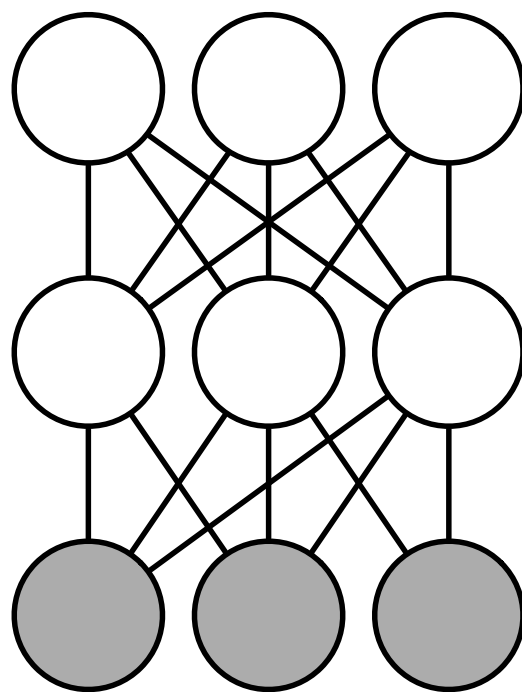
この計算はmaximum likelihood
learning (最尤学習) に必要

Inferenceが解けない場合の理由: 潜在変数間のインタラクションが存在する

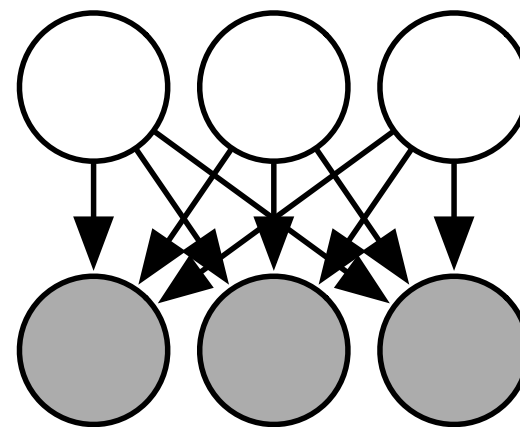
白丸: 潜在変数
灰色丸: 可視変数



semi-restricted
Boltzmann machine
(undirected)



deep Boltzmann
machine
(undirected)



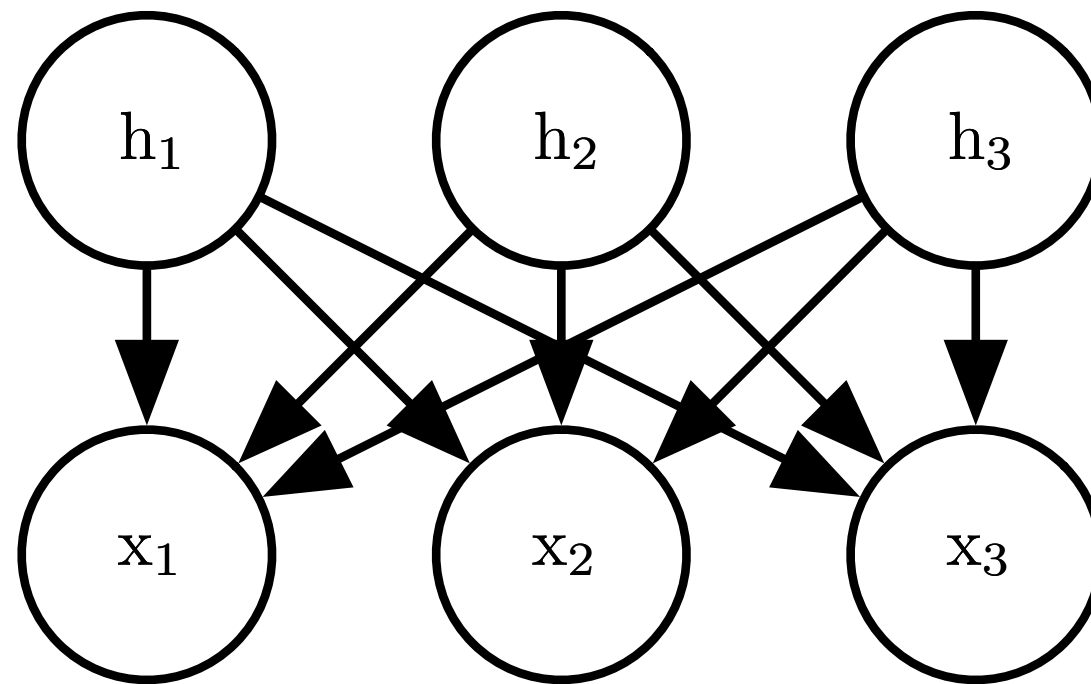
directed model

潜在層内にインタラクションがない
が、潜在層間にインタラクションが
存在する

子ノードが見えたら、
親ノード間インタラク
ションが発生する

Chapter 16 Structured
Probabilistic Modelの内容

参考： Linear Factor Models (Chapter 13)



$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \text{noise}$$

Linear factor modelは潜在変数の分布 $p(\mathbf{h})$ を階乗分布 (factorial distribution) にしている

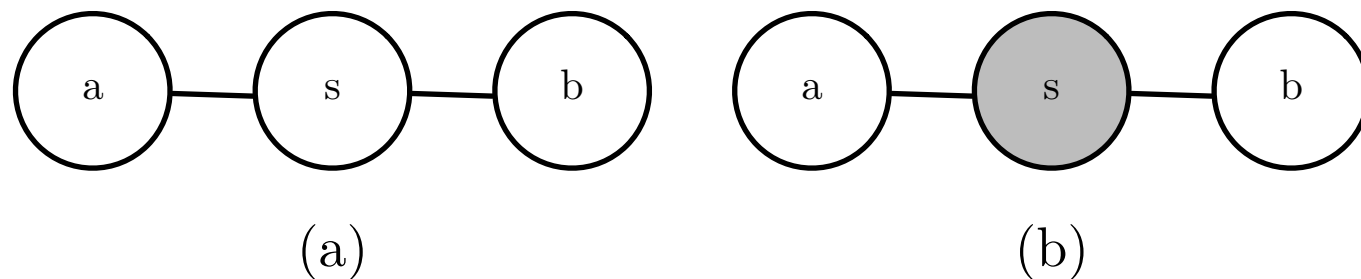
$$p(\mathbf{h}) = \prod_i p(h_i)$$

つまり潜在変数がお互いに独立しているのを仮定する

- ・ factorial分布だと、inferenceも簡単だし、samplingも簡単
- ・ factorial分布仮定は後ほど出る (Variational Inference and Learningでのmean field方法)

参考： Separation (undirected models) and D-Separation (directed models) (Chapter 16 16.2.5)

Separationの意味は、conditional independence (条件付独立)

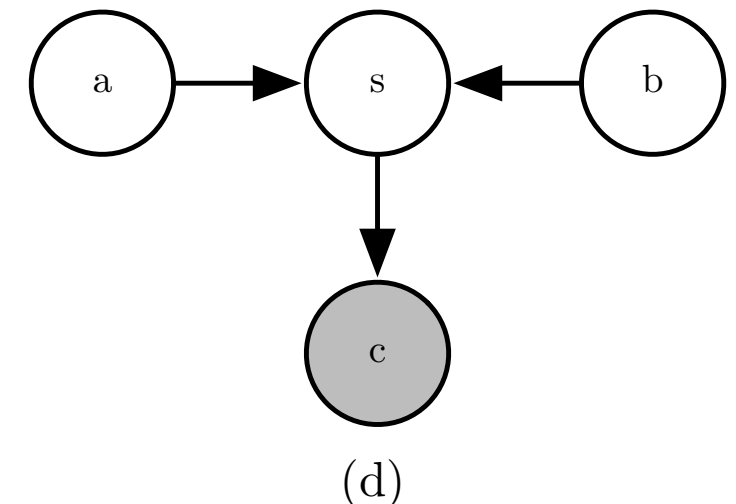
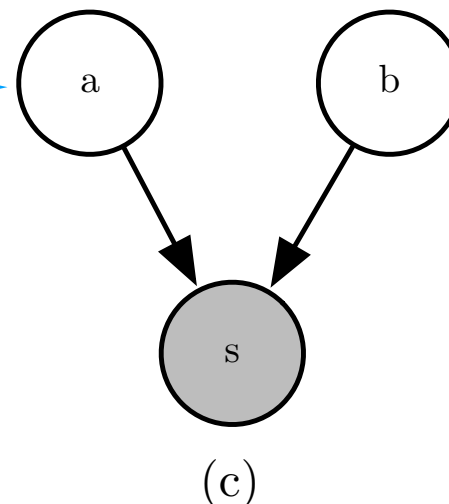


sが見えたら、aとbのパスは非アクティブ。 (aとbがseparate)
sが見えなかったら、aとbのパスはアクティブ。 (aとb not separate)

directed modelの特別な依存関係： 子ノードもしくは子ノードの後継が見えたら、親ノード間に依存がある (explain away effect)

例えば、s: 同僚が出社していない
a: 同僚が病気
b: 同僚が旅行中

もしsが見えたら、aとbは依存関係が出そう。例えば、sが見える場合、aが発生したら、bの可能性が非常に低くなる。



Inference as Optimization

最尤学習では $\log p(v; \theta)$ を計算すること。

そのため、 $p(v, h; \theta)$ で v の周辺分布を計算すれば良いが、コストが高い
代わりに、 $\log p(v; \theta)$ の下界 $\mathcal{L}(v, \theta, q)$ (ELBO, evidence lower bound) を計算

$$\mathcal{L}(v, \theta, q) = \log p(v; \theta) - D_{KL}(q(h | v) || p(h | v; \theta))$$

q は任意の h の分布。 D は常に非負数。 q と p が同じ分布の場合、ELBO イコール $\log p(v; \theta)$ 。 q が p に近づけば近づくほど ELBO がタイトになる。

(ELBO, evidence lower bound) の良く使う形: (変換は次のスライド)

$$\mathcal{L}(v, \theta, q) = \mathbb{E}_{h \sim q}[\log p(h, v)] + H(q)$$

ELBO の Optimization で良い q を出せる、つまり $p(h|v)$ の良い Inference を出せる

Inference as Optimization: ELBOを組み替える

$$\begin{aligned}\mathcal{L}(v, \theta, q) &= \log p(v; \theta) - \mathbb{E}_{h \sim q} \log \frac{q(h | v)}{p(h | v)} \\ &= \log p(v; \theta) - \mathbb{E}_{h \sim q} \log \frac{q(h | v)}{\frac{p(h, v; \theta)}{p(v; \theta)}} \\ &= \log p(v; \theta) - \mathbb{E}_{h \sim q} [\log q(h | v) - \log p(h, v; \theta) + \log p(v; \theta)] \\ &= - \mathbb{E}_{h \sim q} [\log q(h | v) - \log p(h, v; \theta)] \\ &= \mathbb{E}_{h \sim q} [\log p(h, v; \theta)] + H(q)\end{aligned}$$

VAE (Variational Autoencoder) の学習もELBOに基づく。次のスライドでまず1つ目ELBOを最適化する方法EM (Expectation Maximization) を紹介する。

Expectation Maximization (EM)

$$\begin{aligned}\mathcal{L}(v, \theta, q) &= -\mathbb{E}_{h \sim q}[\log q(h | v) - \log p(h, v; \theta)] \\ &= \mathbb{E}_{h \sim q}[\log p(h, v; \theta)] + H(q)\end{aligned}$$

収束までEステップとMステップを繰り返す。

E-step: q の更新

$$q(h^{(i)} | v) = p(h^{(i)} | v^{(i)}; \theta^{(0)})$$

$v^{(i)}$ は i 番目 training サンプル $\theta^{(0)}$ はステップ開始時のパラメータ

MステップでMonte Carloサンプリングで期待値を計算する為に q のサンプルを用意する。
(自分の理解)
Monte Carlo Sampling
は次のスライド

M-step: θ の更新

$$\sum_i \mathcal{L}(v^{(i)}, \theta, q) \text{ を } \theta \text{ に対して最大化にする}$$

M-stepで θ が変わると q も変わるはずだが、E-stepで取った q サンプルを使う。

E-stepで $q=p$ (exact inference) なので、 $ELBO=\log p(v;\theta)$ 。 M-stepで θ が変わって、 q 変わらないので、 $ELBO$ と $\log p(v;\theta)$ の差がどんどん出る。またE-stepに入ると、差が0になる。EMは近似推論 $q(h|v)$ だけじゃなく、近似事後 (posterior) 分布 $p(h|v)$ を出せる。

参考: Basics of Monte Carlo Sampling (Chapter 17 17.1.2)

厳密に計算できないsumもしくは積分 (integral, continued) を近似する方法
例えばある分布での期待値を計算する時 (分布がわからなくても、サンプリングできれば)

$$s = \sum_x p(x)f(x) = E_p[f(x)] \quad s = \int p(x)f(x) = E_p[f(x)]$$

n個サンプルを取って平均値を近似値として使う

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(x^{(i)})$$

このestimatorはunbiasedだ:

$$\mathbb{E}[\hat{s}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(x^{(i)})] = \frac{1}{n} \sum_{i=1}^n s = s$$

このestimatorのvarianceは:

$$\text{Var}[\hat{s}_n] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[f(x)] = \frac{\text{Var}[f(x)]}{n}$$

このbasic Monte CarloサンプリングでELBOにある期待値を全部計算できる

MAP (最大事後確率) Inference and Sparse Coding

最大事後確率推定と最尤推定の区別は次のスライドへ

考え方: 最尤 $q(h|v)$ ではなく、最尤 h (点推定) を推定して、ELBOを最適化する

$$\mathcal{L}(v, \theta, q) = \mathbb{E}_{h \sim q}[\log p(h, v)] + H(q)$$

①MAP推定問題になるように、 q をディラック分布に限定 $H(q) = -\mathbb{E}_{h \sim q} \log q(h|v)$

$$q(h|v) = \delta(h - \mu) = \begin{cases} +\infty, & h = \mu \\ 0, & h \neq \mu \end{cases} = -\int \delta(x) \log \delta(x) dx$$

$H(q)$ が定数なので、省略可

②ELBOに代入して、最尤 h 推定問題は、MAP推定問題になる

$$\mu^* = \underset{\mu}{\operatorname{argmax}} \log p(h = \mu, v) \quad \rightarrow \quad h^* = \underset{h}{\operatorname{argmax}} \log p(h, v)$$

③EMと同じように、E-stepにMAP推定で h^* を推定; M-stepに θ を更新して、 $\log p(h^*, v)$ を最適化する

潜在変数 h のMAP推定問題

$$\begin{aligned} &= \underset{h}{\operatorname{argmax}} \log \frac{p(h, v)}{p(v)} \\ &= \underset{h}{\operatorname{argmax}} \log p(h|v) \end{aligned}$$

参考：MAP（Maximum A Posteriori, 最大事後確率）推定と ML（Maximum Likelihood, 最尤）推定（Chapter 5 5.5, 5.6）

最尤推定：

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} p_{model}(\mathbb{X}; \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m p_{model}(x^{(i)}; \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log p_{model}(x^{(i)}; \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{x \sim \hat{p}_{data}} \log p_{model}(x; \theta)$$

最大事後確率推定：

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta | x) = \underset{\theta}{\operatorname{argmax}} \log p(x | \theta) + \log p(\theta)$$

最尤推定で、 θ は確率変数ではない。 θ は未知だが固定。（頻度論統計学の考え方、frequentist statistics）

最大事後確率推定で、 θ は見えないので、 θ は確率変数だ。（ベイズ統計学の考え方、Bayesian statistics）

log尤度

θ の事前確率分布を使う

MAP（最大事後確率） Inference and Sparse Coding (Chapter 13.4)

Deep Learningの中に、MAP推定は主にSparse Codingモデルに使われている

- ① Sparse CodingもLinear Factorモデルの一つ、 h の事前確率分布はラプラス分布:

$$p(h_i) = \frac{\lambda}{2} e^{-\lambda |h_i|}$$

- ② v は線形Decoder + ノイズで生成する: $p(v | h) = \mathcal{N}(v; Wh + b, \frac{1}{\beta} I)$

- ③ ページ2,4紹介した子ノード(v)が見えたら親ノード(h)間依存あり、独立ではない原因で $p(h|v)$ は解けないことで、MAP推定で、最尤 h だけを推定する:

$$h^* = \underset{h}{\operatorname{argmax}} p(h | v) = \underset{h}{\operatorname{argmax}} \log p(h | v) = \underset{h}{\operatorname{argmax}} \log p(v | h)$$

$$+ \underset{h}{\operatorname{argmax}} \log p(h) = \underset{h}{\operatorname{argmin}} \lambda \|h\|_1 + \beta \|v - Wh\|_2^2$$

- ④ training setの h を H に、 v を V に入れて、Sparse Codingモデルの最適化目標関数:

$$J(H, W) = \sum_{i,j} |H_{i,j}| + \sum_{i,j} (V - HW^T)_{i,j}^2$$

Sparse Codingモデルのトレーニングは J を最小化する

Variational (変分) Inference and Learning

考え方: $\mathbb{E}_q \log p(h, v)$ を簡単に計算できるように、 $q(h|v)$ を制限する。良く使われる制限は、factorial 分布:

$$q(h | v) = \prod_i q(h_i | v)$$

mean field (平均場、平均ば) 方法。変分方法の綺麗さは、 q を制限する時、何のパラメータも使わないこと。

連続潜在変数の場合: 変分法 (calculus of variations) で関数空間から一つの分布関数を決める。変分法紹介は後ほど。

変分法は Variational Inference 名前の由来が、離散の場合変分法はいらない

離散潜在変数の場合: 普通に ELBO を最適化できる

例えば、バイナリ Sparse Coding モデルの場合

$$q(h_i = 1 | v) = \hat{h}_i$$

$$\mathcal{L}(v, \theta, \hat{h}) = \mathbb{E}_{h \sim q} [\log p(h, v)] + H(q)$$

$$p(h_i = 1) = \sigma(b_i)$$

$$p(v | h) = \mathcal{N}(v; Wh, \beta^{-1})$$

収束まで更新。 \hat{h} を同時に更新できない。収束標準は、 \mathcal{L} が上がらない、もしくは \hat{h} が変わらない

$$\frac{\partial}{\partial \hat{h}_i} \mathcal{L}(v, \theta, \hat{h}) = 0$$

この不動点 (fixed-point) 更新推定ルールを適用後の式次スライドへ

Variational (変分) Inference and Learning

$\frac{\partial}{\partial \hat{h}_i} \mathcal{L}(v, \theta, \hat{h}) = 0$ を展開後、**RNN**になる:

$$\hat{h}_i = \sigma(b_i + \underbrace{v^T}_{\text{インプット}} \beta W_{:,i} - \frac{1}{2} W_{:,i}^T \beta W_{:,i} - \sum_{j \neq i} W_{:,j}^T \beta W_{:,i} \hat{h}_j)$$

この式はRNNなので、直接にトレーニングしても良い

組み替えて: $\hat{h}_i = \sigma(b_i + (v - \sum_{j \neq i} W_{:,j} \hat{h}_j)^T \beta W_{:,i} - \frac{1}{2} W_{:,i}^T \beta W_{:,i})$

隠れユニット*i*は、他の隠れユニットを使って、インプット(*v*)をEncodeしている

連続の場合、実践者は自分で変分法を解く必要がなく、mean fieldで良く使われる $q(h|v)$ 不動点更新式がある:

$$\tilde{q}(h_i | v) = \exp(\mathbb{E}_{h_{-i} \sim q(h_{-i}|v)} \log \tilde{p}(v, h))$$

この期待値を計算すると、 $q(h_i|v)$ の関数式が出る。

例 (p640-641) は略。

参考：変分法 (Calculus of Variations, Chapter 19 19.4.2)

Functional (関数の関数) : $J[f]$

Functional derivatives /
Variational derivatives :

$$\frac{\delta}{\delta f(x)} J$$

Functional derivativeの一つのプロパティ (式①) :

$$\frac{\delta}{\delta f(x)} \int g(f(x), x) dx = \frac{\partial}{\partial y} g(f(x), x) \rightarrow \frac{\partial}{\partial \theta_i} \sum_j g(\theta_j, j) = \frac{\partial}{\partial \theta_i} g(\theta_i, i)$$

分布関数のエントロピー (Functional) :

$$H[p] = - \mathbb{E}_x \log p(x) = - \int p(x) \log p(x) dx$$

左式の直感的な理解: θ_i 以外の微分は0なので、 θ_i での微分だけ残る

Functional微分で $H[p]$ が最大時の p を求める: しかし p が分布にならないかもしれないので、制限が必要。例えば、 $p(x)$ の積分が1になること、平均値、分散。ラグランジュの未定乗数法 (**Lagrange multiplier**) でこれらの制限を入れる。新しい p のFunctionalは:

$$\begin{aligned} \mathcal{L}[p] &= \lambda_1 \left(\int p(x) dx - 1 \right) + \lambda_2 (\mathbb{E}[x] - u) + \lambda_3 (\mathbb{E}[(x - \mu)^2] - \sigma^2) + H[p] \\ &= \int (\lambda_1 p(x) + \lambda_2 p(x)x + \lambda_3 p(x)(x - \mu)^2 - p(x) \log p(x)) dx - \lambda_1 - \mu \lambda_2 - \sigma^2 \lambda_3 \end{aligned}$$

参考：変分法 (Calculus of Variations, Chapter 19 19.4.2) 続

前ページのラグランジュを最小化(?)するため、Functional微分を0にする：

$$\frac{\delta}{\delta p(x)} \mathcal{L} = \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1 - \log p(x) = 0$$

前ページの変分
法の式①を使う

組み替えて： $p(x) = \exp(\lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1)$

制限を満たすために、 λ を設定して、 $p(x)$ が正規分布になる：

$$\lambda_1 = 1 - \log \sigma \sqrt{2\pi}, \lambda_2 = 0, \lambda_3 = -\frac{1}{2\sigma^2} \quad \rightarrow \quad p(x) = \mathcal{N}(x; \mu, \sigma^2)$$

だから分布がわからない時、正規分布を使う。正規分布の時のエントロピーが一番大きい。

Learned Approximate Inference (学んだ近似推定)

考え方: ELBOの最適化プロセスを下記のマッピングと見る

$$v \xrightarrow{f} q^* = \underset{q}{\operatorname{argmax}} \mathcal{L}(v, q)$$

Neural network (推定ネットワーク) をトレーニングして、 f に近似する

①Wake-Sleep: $p(v|h)$ で h から v までサンプリングして、 (v, h) をトレーニングデータとして推定ネットワークをトレーニングする。

②学んだ近似推定はVariational Autoencoderに使われて、生成モデルの主な手法の一つになっている。

推定ネットワークはターゲットがなし、ただELBOの一部として定義されて、ネットワークのパラメータはELBOを最大化するために調整される。