

生成モデル

EM アルゴリズム, 変分オートエンコーダ

中田和秀

東京工業大学 工学院 経営工学系

機械学習入門

<http://www.iee.e.titech.ac.jp/~nakatalab/text/lecture.html>

概要

ここでは、データからデータの生成モデルを構築する方法について説明をする。一般的にこの構築は難易度が高いことが知られている。しかし、適切な生成モデルを作ることができれば、様々な予測タスクやクラスタリングやデータ生成など多様な形で応用することができる。

目次：

1. 潜在変数を含んだ確率分布
EM アルゴリズム
2. 変分オートエンコーダー (VAE)

記号の使い方：

- $A := B$ は、 B で A を定義する、 B を A に代入することを意味する
- $[n]$ は n までのインデックスの集合を表し $[n] := \{1, 2, \dots, n\}$

潜在変数を含んだ確率分布と EM アルゴリズム

潜在変数を含んだ確率変数

- z : 潜在変数
- x : 確率変数

$$p(x) = \sum_z p(z)p(x|z) \quad z \text{ が離散分布の場合}$$

$$= \int p(z)p(x|z)dz \quad z \text{ が連続分布の場合}$$

- 現実世界には、観測できない潜在的な何かに依存して確率が決まる状況は多い
- (潜在的な何かの存在はともかく) 複雑な分布 $p(x)$ を単純な2つの分布 $p(z)$ と $p(x|z)$ で表現できる

以下では、連続分布の場合で説明するが、離散分布の場合でも \int が \sum に変わるだけでまったく同じ議論ができる。

基本関係式

$p(\mathbf{x})$ はパラメタ θ を持つ確率分布 $p(\mathbf{x}; \theta)$ とする。

$q(z)$: 潜在変数に関する任意の確率分布¹ $\int q(z) dz = 1, q(z) \geq 0.$

基本関係式

$$\log p(\mathbf{x}; \theta) = LB(q(z), \theta) + KL(q(z), p(z|\mathbf{x}; \theta))$$

$$\text{ただし、 } LB(q(z), \theta) := \mathbb{E}_{q(z)} \left[\log \frac{p(\mathbf{x}, z; \theta)}{q(z)} \right]$$

$KL(q(z), p(z|\mathbf{x}; \theta)) \quad \dots \dots \quad \text{KL ダイバージェンス}$

※ $KL(q(z), p(z|\mathbf{x}; \theta)) \geq 0$ より、 $\log p(\mathbf{x}; \theta) \geq LB(q(z), \theta).$
つまり、 $LB(q(z), \theta)$ は、 $\log p(\mathbf{x}; \theta)$ の下界 (Lower Bound)

¹真の確率分布ではないことに注意

式変形

$$\begin{aligned} & LB(q(\mathbf{z}), \boldsymbol{\theta}) + KL(q(\mathbf{z}), p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})) \\ &= \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} \right] + \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} dz \\ &= \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} dz + \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} dz \\ &= \int q(\mathbf{z}) \left(\log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} + \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} \right) dz \\ &= \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} dz \\ &= \int q(\mathbf{z}) \log \frac{p(\mathbf{x}; \boldsymbol{\theta})p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} dz \\ &= \int q(\mathbf{z}) \log p(\mathbf{x}; \boldsymbol{\theta}) dz \\ &= \log p(\mathbf{x}; \boldsymbol{\theta}) \int q(\mathbf{z}) dz \\ &= \log p(\mathbf{x}; \boldsymbol{\theta}) \end{aligned}$$

LB の性質

$LB(q(z), \theta)$ に関しては次の 2 つの性質を使うので、先に説明しておく。

LB の性質 1

$$LB(q(z), \theta) = \mathbb{E}_{q(z)} [\log p(\mathbf{x}, z; \theta)] + H(q(z))$$

ただし、 $H(q(z)) := - \int q(z) \log q(z) dz$ エントロピー

$\theta = \{\theta_1, \theta_2\}$ で、 $p(\mathbf{x}, z; \theta) = p(z; \theta_1) p(\mathbf{x}|z; \theta_2)$ のとき、次の性質が成り立つ。

LB の性質 2

$$LB(q(z), \theta) = \mathbb{E}_{q(z)} [\log p(\mathbf{x}|z; \theta_2)] - KL(q(z), p(z; \theta_1))$$

式変形

$$\begin{aligned}LB(q(\mathbf{z}), \boldsymbol{\theta}) &= \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} \right] = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} \\ &= \int q(\mathbf{z}) \log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} - \int q(\mathbf{z}) \log q(\mathbf{z}) d\mathbf{z} \\ &= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})] + H(q(\mathbf{z}))\end{aligned}$$

$$\begin{aligned}LB(q(\mathbf{z}), \boldsymbol{\theta}) &= \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} \right] = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} \\ &= \int q(\mathbf{z}) \log \frac{p(\mathbf{z}; \boldsymbol{\theta}_1) p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}_2)}{q(\mathbf{z})} d\mathbf{z} \\ &= \int q(\mathbf{z}) \log p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}_2) d\mathbf{z} - \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z}; \boldsymbol{\theta}_1)} d\mathbf{z} \\ &= \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}_2)] - KL(q(\mathbf{z}), p(\mathbf{z}; \boldsymbol{\theta}_1))\end{aligned}$$

最尤推定

- 独立にサンプリングされたデータ $\mathcal{D} := \{\mathbf{x}_d\}_{d \in [D]}$
- 尤度: $p(\mathcal{D}) = \prod_{d \in [D]} p(\mathbf{x}_d; \boldsymbol{\theta})$

対数尤度の最大化

$$\max_{\boldsymbol{\theta}} \sum_{d \in [D]} \log p(\mathbf{x}_d; \boldsymbol{\theta})$$

$$\begin{aligned} \sum_{d \in [D]} \log p(\mathbf{x}_d; \boldsymbol{\theta}) &= \sum_{d \in [D]} LB_d(q_d(\mathbf{z}), \boldsymbol{\theta}) + KL(q_d(\mathbf{z}), p(\mathbf{z}|\mathbf{x}_d; \boldsymbol{\theta})) \\ &\geq \sum_{d \in [D]} LB_d(q_d(\mathbf{z}), \boldsymbol{\theta}) \quad (\because KL(q_d(\mathbf{z}), p(\mathbf{z}|\mathbf{x}_d; \boldsymbol{\theta})) \geq 0) \end{aligned}$$

EM アルゴリズムの導出

(1) $\sum_{d \in [D]} \log p(\mathbf{x}_d; \boldsymbol{\theta})$ の代わりに下界 $\sum_{d \in [D]} LB_d(q(\mathbf{z}); \boldsymbol{\theta})$ の最大化

$$\rightarrow \max_{\boldsymbol{\theta}} \sum_{d \in [D]} \log p(\mathbf{x}_d; \boldsymbol{\theta}) \implies \max_{\boldsymbol{\theta}} \sum_{d \in [D]} LB_d(q_d(\mathbf{z}); \boldsymbol{\theta})$$

(2) できるだけタイトな下界となるように $q_d(\mathbf{z})$ を決めたい。

$q_d(\mathbf{z}) := p(\mathbf{z} | \mathbf{x}_d; \boldsymbol{\theta})$ のとき、 $KL(q_d(\mathbf{z}), p(\mathbf{z} | \mathbf{x}_d; \boldsymbol{\theta})) = 0$ より、

$$\sum_{d \in [D]} \log p(\mathbf{x}_d; \boldsymbol{\theta}) = \sum_{d \in [D]} LB_d(q(\mathbf{z}); \boldsymbol{\theta})$$

一番よい下界となっている。

EM アルゴリズム 1

EM アルゴリズム 1

ステップ0 初期点 θ を定める。

ステップ1 $q_d(\mathbf{z}) := p(\mathbf{z}|\mathbf{x}_d; \theta)$ $d \in [D]$

ステップ2 $\theta := \operatorname{argmax}_{\theta} \sum_{d \in [D]} LB_d(q_d(\mathbf{z}), \theta)$

ステップ3 終了条件を満たしていなければ、ステップ1に戻る。

前の反復点を $\tilde{\theta}$ とすると、 $q_d(\mathbf{z}) := p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta})$ である。LB の性質 1 より、

$$LB_d(p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta}), \theta) = \mathbb{E}_{p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta})} [\log p(\mathbf{x}_d, \mathbf{z}; \theta)] + H(p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta}))$$

よって、

$$\max_{\theta} \sum_{d \in [D]} LB_d(p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta}), \theta) \iff \max_{\theta} \sum_{d \in [D]} \mathbb{E}_{p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta})} [\log p(\mathbf{x}_d, \mathbf{z}; \theta)]$$

※ $H(p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta}))$ は最適化に関係ない

EM アルゴリズム 2

EM アルゴリズムは次のように書くことができる
(この形で紹介されることが多い)

EM アルゴリズム 2

ステップ0 初期点 $\tilde{\theta}$ を定める。

ステップ1 $p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta})$ の計算

$d \in [D]$

ステップ2 $\theta := \operatorname{argmax}_{\theta} Q(\tilde{\theta}, \theta)$

ステップ3 終了条件を満たしていなければ、 $\tilde{\theta} := \theta$ としてステップ1に戻る。

ただし、 $Q(\tilde{\theta}, \theta) := \sum_{d \in [D]} \mathbb{E}_{p(\mathbf{z}|\mathbf{x}_d; \tilde{\theta})} [\log p(\mathbf{x}_d, \mathbf{z}; \theta)]$

混合モデルの場合

z が (有限) 離散分布

- 潜在変数 z の分布: k 番目の分布を取る確率 π_k
- k 番目の分布: $p_k(\mathbf{x}; \boldsymbol{\theta}_k)$

カテゴリ分布

- EM アルゴリズム 1 のステップ 1 :

$$q_d(z) := p(z|\mathbf{x}_d; \boldsymbol{\theta}) = \frac{p(z; \boldsymbol{\theta}) p(\mathbf{x}|z; \boldsymbol{\theta})}{\sum_{z'} p(z'; \boldsymbol{\theta}) p(\mathbf{x}_d|z'; \boldsymbol{\theta})}$$

$d \in [D]$ に対し、 \mathbf{x}_d が k 番目の分布に属する確率は

$$r_{dk} = \frac{\pi_k p_k(\mathbf{x}_d; \boldsymbol{\theta}_k)}{\sum_{k'} \pi_{k'} p_{k'}(\mathbf{x}_d; \boldsymbol{\theta}_{k'})} \quad (k \in [K])$$

- EM アルゴリズム 1 のステップ 2 :

$$\begin{aligned} \sum_{d \in [D]} LB_d(q_d(z), \boldsymbol{\theta}) &:= \sum_{d \in [D]} \mathbb{E}_{q_d(z)} \left[\log \frac{p(\mathbf{x}_d, z; \boldsymbol{\theta})}{q_d(z)} \right] \\ &= \sum_{d \in [D]} \sum_{k \in [K]} r_{dk} \log \frac{\pi_k p_k(\mathbf{x}_d; \boldsymbol{\theta}_k)}{r_{dk}} \end{aligned}$$

混合モデルに対する EM アルゴリズム

補助関数法と同様に $\sum_{d \in [D]} \sum_{k \in [K]} r_{dk} \log \frac{\pi_k p_k(\mathbf{x}_d; \boldsymbol{\theta}_k)}{r_{dk}}$ の最適化は分割できる。

混合モデルに対する EM アルゴリズム

ステップ0 初期点 $\pi, \boldsymbol{\theta}_k$ ($k \in [K]$) を決める。

ステップ1 $r_{dk} := \frac{\pi_k p(\mathbf{x}_d; \boldsymbol{\theta}_k)}{\sum_{k' \in [K]} \pi_{k'} p(\mathbf{x}_d; \boldsymbol{\theta}_{k'})}$ $d \in [D], k \in [K]$

ステップ2 $\pi_k := \frac{1}{D} \sum_{d \in [D]} r_{dk}$ ($k \in [K]$)

$k \in [K]$ に対して、重み付き最尤推定を行う。

$$\max_{\boldsymbol{\theta}_k} \sum_{d \in [D]} r_{dk} \log p_k(\mathbf{x}_d; \boldsymbol{\theta}_k)$$

ステップ3 終了条件を満たしていなければ、ステップ1に戻る

潜在変数が離散でカテゴリ分布のとき、EM アルゴリズムは補助関数法と等価

変分オートエンコーダー

潜在変数が離散分布でなく連続分布の場合を考える。

対数最尤の最大化

$$\max_{\boldsymbol{\theta}} \sum_{d \in [D]} \log p(\mathbf{x}_d; \boldsymbol{\theta})$$

$$p(\mathbf{x}_d; \boldsymbol{\theta}) = \int p(\mathbf{z}; \boldsymbol{\theta}) p(\mathbf{x}_d | \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z}$$

EM アルゴリズムを適用する場合は次の計算をする

EM アルゴリズム 1 (再掲)

ステップ0 適当に $\boldsymbol{\theta}$ を定める。

ステップ1 $q_d(\mathbf{z}) := p(\mathbf{z} | \mathbf{x}_d; \boldsymbol{\theta})$ $d \in [D]$

ステップ2 $\boldsymbol{\theta} := \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{d \in [D]} LB_d(q_d(\mathbf{z}), \boldsymbol{\theta})$

ステップ3 終了条件を満たさなければ、ステップ1に戻る

枠組み 1

ステップ1の計算が困難

- 分母の積分が計算できない
- z の取れる値が無限にある

$$q_d(z) = p(z|x_d; \theta) = \frac{p(x_d, z; \theta)}{p(x_d; \theta)} = \frac{p(z; \theta) p(x_d|z; \theta)}{\int p(z; \theta) p(x_d|z; \theta) dz}$$

⇒ $p(z|x_d; \theta) \simeq q(z|x_d; \phi)$ となる分布を $q_d(z)$ とする
 x_d から z の分布を作り、その関係性をパラメタ ϕ で表す。

- 2つの分布の近さ・遠さは KL ダイバージェンスで考える
- $KL(q(z|x_d; \phi), p(z|x_d; \theta))$ が最小となる分布 $q(z|x_d; \phi)$ を探す。

枠組み 2

$$\min_{\phi} KL(q(z|\mathbf{x}_d; \phi), p(z|\mathbf{x}_d; \theta))$$

基本関係式：

$$\log p(\mathbf{x}_d; \theta) = LB_d(q(z|\mathbf{x}_d; \phi), \theta) + KL(q(z|\mathbf{x}_d; \phi), p(z|\mathbf{x}_d; \theta))$$

左辺は ϕ に依存しない。

よって、KL ダイバージェンスの最小化は次の最大化と同等。

$$\max_{\phi} \sum_{d \in [D]} LB_d(q(z|\mathbf{x}_d; \phi), \theta)$$

最適化問題

EM アルゴリズム

ステップ1

$$q_d(\mathbf{z}) := p(\mathbf{z}|\mathbf{x}_d; \boldsymbol{\theta}) \implies \max_{\boldsymbol{\phi}} \sum_{d \in [D]} LB_d(q(\mathbf{z}|\mathbf{x}_d; \boldsymbol{\phi}), \boldsymbol{\theta})$$

ステップ2

$$\max_{\boldsymbol{\theta}} \sum_{d \in [D]} LB_d(q(\mathbf{z}|\mathbf{x}_d; \boldsymbol{\phi}), \boldsymbol{\theta})$$

これらをまとめて、次の最適化問題を解くことになる。

最適化問題

$$\max_{\boldsymbol{\theta}, \boldsymbol{\phi}} \sum_{d \in [D]} LB_d(q(\mathbf{z}|\mathbf{x}_d; \boldsymbol{\phi}), \boldsymbol{\theta})$$

変分オートエンコーダー

オートエンコーダー

$$\begin{array}{ccccc} \mathbf{x} & \xrightarrow{q(\mathbf{z}|\mathbf{x};\phi)} & \mathbf{z} & \xrightarrow{p(\mathbf{x}|\mathbf{z};\theta)} & \mathbf{x} \\ & \text{エンコーダ} & & \text{デコーダ} & \end{array}$$

\mathbf{x} を復元できるような、エンコーダとデコーダのセットを構築すると捉えることができる。

変分オートエンコーダー (Variational Auto-Encoder; VAE) と呼ぶ

必要な分布

- $p(\mathbf{z}; \theta)$ と $p(\mathbf{x}|\mathbf{z}; \theta)$

$$p(\mathbf{x}; \theta) = \int p(\mathbf{x}, \mathbf{z}; \theta) dz = \int p(\mathbf{z}; \theta) p(\mathbf{x}|\mathbf{z}; \theta) dz$$

- $p(\mathbf{z}|\mathbf{x}; \phi)$

デコーダーの設計

$p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = p(\mathbf{z}; \boldsymbol{\theta}) p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})$ を次のように定める。

(1)-A $p(\mathbf{z}; \boldsymbol{\theta})$

$\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$ とする。つまりパラメタ $\boldsymbol{\theta}$ に依存させない。

(1)-B $p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})$

- 深層学習によって、 \mathbf{z} から正規分布の平均ベクトルを生成

$$F_{\mu} : \mathbb{R}^k \rightarrow \mathbb{R}^n,$$
$$\boldsymbol{\mu}_D := F_{\mu}(\mathbf{z}; \boldsymbol{\theta})$$

※ F_{μ} は深層学習の回帰関数で $\boldsymbol{\theta}$ はパラメタ

- \mathbf{x} は正規分布 $N(\boldsymbol{\mu}_D, \mathbf{I})$ に従う。

エンコーダーの設計

$p(z|\mathbf{x}; \theta)$ を次のように定める。

(2) $p(z|\mathbf{x}; \phi)$

- 深層学習によって、 \mathbf{x} から正規分布の平均ベクトルを生成

$$G_{\mu} : \mathbb{R}^n \rightarrow \mathbb{R}^k,$$
$$\boldsymbol{\mu}_E := G_{\mu}(\mathbf{x}; \phi)$$

※ G_{μ} は深層学習の回帰関数で ϕ はパラメタ

- z は正規分布 $N(\boldsymbol{\mu}_E, \mathbf{I})$ に従う。

注意

エンコーダとデコーダの分散共分散行列として対角行列を考え、その対角成分も深層学習で学習することが多い。

本スライドでは、話を簡単にするため単位行列で固定する場合を考える。対角成分を考える場合でも、式が複雑になるだけで、以下の流れは同じ。

目的関数

正規分布の代入

$d \in [D]$ において

$$LB_d(q(\mathbf{z}|\mathbf{x}_d; \phi), \theta) = -\frac{n}{2} \log 2\pi + LB_1 + LB_2$$

$$LB_1 := -\frac{1}{2} \mathbb{E}_{\epsilon \sim N(\mathbf{0}, \mathbf{I})} [(\mathbf{x}_d - \boldsymbol{\mu}_D)^T (\mathbf{x}_d - \boldsymbol{\mu}_D)]$$

$$LB_2 := -\frac{1}{2} \boldsymbol{\mu}_E^T \boldsymbol{\mu}_E$$

ただし、 $\boldsymbol{\mu}_E := G_\mu(\mathbf{x}_d; \phi)$, $\mathbf{z} := \boldsymbol{\mu}_E + \epsilon$, $\boldsymbol{\mu}_D := F_\mu(\mathbf{z}; \theta)$

式変形 1

LB に関する性質 2 より、

$$LB_d(q(\mathbf{z}|\mathbf{x}_d; \phi), \theta) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x}_d; \phi)}[\log p(\mathbf{x}|\mathbf{z}; \theta_2)] - KL(q(\mathbf{z}|\mathbf{x}_d; \phi), p(\mathbf{z}; \theta_1))$$

第 1 項と第 2 項で分けて考える。

$$\log p(\mathbf{x}|\mathbf{z}; \theta) = -\frac{n}{2} \log 2\pi - \frac{1}{2}(\mathbf{x}_d - \boldsymbol{\mu}_D)^T (\mathbf{x}_d - \boldsymbol{\mu}_D) \text{ より、}$$

$$\begin{aligned} & \mathbb{E}_{q(\mathbf{z}|\mathbf{x}_d; \phi)}[\log p(\mathbf{x}|\mathbf{z}; \theta)] \\ &= -\frac{n}{2} \log 2\pi - \mathbb{E}_{\mathbf{z} \sim N(\boldsymbol{\mu}_E, \mathbf{I})} \left[\frac{1}{2}(\mathbf{x}_d - \boldsymbol{\mu}_D)^T (\mathbf{x}_d - \boldsymbol{\mu}_D) \right] \end{aligned}$$

$$\text{ただし、 } \boldsymbol{\mu}_D := F_{\mu}(\mathbf{z}; \theta)$$

$$= -\frac{n}{2} \log 2\pi - \mathbb{E}_{\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I})} \left[\frac{1}{2}(\mathbf{x}_d - \boldsymbol{\mu}_D)^T (\mathbf{x}_d - \boldsymbol{\mu}_D) \right]$$

$$\text{ただし、 } \mathbf{z} = \boldsymbol{\mu}_E + \boldsymbol{\epsilon}, \quad \boldsymbol{\mu}_D := F_{\mu}(\mathbf{z}; \theta)$$

式変形 2

$$\begin{aligned} & -KL(q(\mathbf{z}|\mathbf{x}_d; \phi), p(\mathbf{z}; \theta_1)) \\ &= -\mathbb{E}_{q(\mathbf{z}|\mathbf{x}_d; \phi)} \left[\log \frac{q(\mathbf{z}|\mathbf{x}_d; \phi)}{p(\mathbf{z})} \right] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{z}} \left[(\mathbf{z} - \boldsymbol{\mu}_E)^T (\mathbf{z} - \boldsymbol{\mu}_E) - \mathbf{z}^T \mathbf{z} \right] \end{aligned}$$

ただし $\mathbf{z} \sim N(\boldsymbol{\mu}_E, \mathbf{I})$

$$\begin{aligned} &= \frac{1}{2} \mathbb{E}_{\mathbf{z}} \left[-2\boldsymbol{\mu}_E^T \mathbf{z} + \boldsymbol{\mu}_E \boldsymbol{\mu}_E^T \right] \\ &= \frac{1}{2} \left(-2\boldsymbol{\mu}_E^T \mathbb{E}_{\mathbf{z}} [\mathbf{z}] + \boldsymbol{\mu}_E^T \boldsymbol{\mu}_E \right) \\ &= \frac{1}{2} \left(-2\boldsymbol{\mu}_E^T \boldsymbol{\mu}_E + \boldsymbol{\mu}_E^T \boldsymbol{\mu}_E \right) \\ &= -\frac{1}{2} \boldsymbol{\mu}_E^T \boldsymbol{\mu}_E \end{aligned}$$

確率的勾配降下法

確率的勾配降下法

ステップ0 初期点 θ , ϕ を決める。

ステップ1 バッチ $S \subset [D]$ を作る

ステップ2 探索方向の計算

$$\mathbf{d}_\theta := \frac{1}{|S|} \sum_{d \in S} \nabla_{\boldsymbol{\theta}} LB_d(q(\mathbf{z}|\mathbf{x}_d; \phi), \boldsymbol{\theta})$$
$$\mathbf{d}_\phi := \frac{1}{|S|} \sum_{d \in S} \nabla_{\boldsymbol{\phi}} LB_d(q(\mathbf{z}|\mathbf{x}_d; \phi), \boldsymbol{\theta})$$

ステップ3 ステップサイズ α を計算

ステップ4 パラメタの更新

$$\boldsymbol{\theta}_{k+1} := \boldsymbol{\theta}_k - \alpha \mathbf{d}_\theta$$

$$\boldsymbol{\phi}_{k+1} := \boldsymbol{\phi}_k - \alpha \mathbf{d}_\phi$$

ステップ5 終了条件を満たしていなければ、ステップ1に戻る

以下では、 $\nabla_{\boldsymbol{\theta}} LB_d(q(\mathbf{z}|\mathbf{x}_d; \phi), \boldsymbol{\theta})$ と $\nabla_{\boldsymbol{\phi}} LB_d(q(\mathbf{z}|\mathbf{x}_d; \phi), \boldsymbol{\theta})$ の計算法を考える

モンテカルロ法

$LB_1 := \mathbb{E}_{\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I})} \left[-\frac{1}{2} (\mathbf{x}_d - \boldsymbol{\mu}_D)^T (\mathbf{x}_d - \boldsymbol{\mu}_D) \right]$ は複雑な分布の期待値計算

→ モンテカルロ法で近似計算を行う

$$LB_1 = \mathbb{E}_{\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I})} \left[-\frac{1}{2} (\mathbf{x}_d - \boldsymbol{\mu}_D)^T (\mathbf{x}_d - \boldsymbol{\mu}_D) \right]$$

ただし、 $\boldsymbol{\mu}_E := G_\mu(\mathbf{x}_d; \phi)$, $\mathbf{z} = \boldsymbol{\mu}_E + \boldsymbol{\epsilon}$, $\boldsymbol{\mu}_D := F_\mu(\mathbf{z}; \theta)$

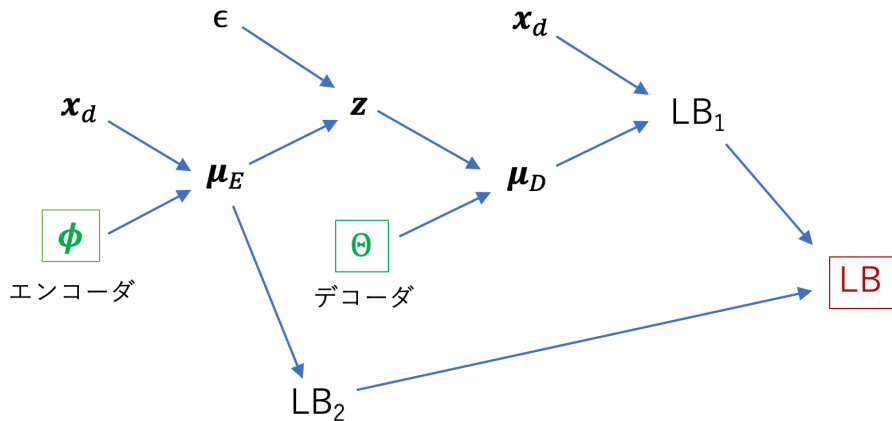
$$\simeq -\frac{1}{2C} \sum_{c \in [C]} (\mathbf{x}_d - \boldsymbol{\mu}_{Dc})^T (\mathbf{x}_d - \boldsymbol{\mu}_{Dc})$$

ただし、 $\boldsymbol{\mu}_E := G_\mu(\mathbf{x}_d; \phi)$, $\mathbf{z}_c = \boldsymbol{\mu}_E + \boldsymbol{\epsilon}_c$, $\boldsymbol{\mu}_{Dc} := F_\mu(\mathbf{z}_c; \theta)$

$\boldsymbol{\epsilon}_c$ は $N(\mathbf{0}, \mathbf{I})$ からのランダムサンプリング

以下では、 $LB_{1c} := -\frac{1}{2} (\mathbf{x}_d - \boldsymbol{\mu}_{Dc})^T (\mathbf{x}_d - \boldsymbol{\mu}_{Dc})$ について考える。

変数の依存関係



デコーダーの勾配

$$\begin{aligned}\frac{\partial LB_d(q(z|\mathbf{x}_d; \phi), \theta)}{\partial \theta} &= \frac{\partial LB_1}{\partial \theta} + \frac{\partial LB_2}{\partial \theta} \\ &\simeq \frac{1}{C} \sum_{c \in [C]} \frac{\partial LB_{1c}}{\partial \theta}\end{aligned}$$

$$c \in [C] \text{ に対して、 } \frac{\partial LB_{1c}}{\partial \theta} = \frac{\partial LB_{1c}}{\partial \mu_{Dc}} \frac{\partial \mu_{Dc}}{\partial \theta}$$

- $\frac{\partial LB_{1c}}{\partial \mu_{Dc}} = (\mathbf{x}_d - \mu_{Dc})^T$
- $\frac{\partial \mu_{Dc}}{\partial \theta} := \frac{\partial F_\mu}{\partial \theta}$ は深層学習の誤差逆伝播法で効率よく計算できる

エンコーダの勾配

$$\frac{\partial LB_d(q(z|x_d; \phi), \theta)}{\partial \phi} = \frac{\partial LB_1}{\partial \phi} + \frac{\partial LB_2}{\partial \phi} \simeq \frac{1}{C} \sum_{c \in [C]} \frac{\partial LB_{1c}}{\partial \phi} + \frac{\partial LB_2}{\partial \phi}$$

第1項: $\frac{\partial LB_{1c}}{\partial \phi} = \frac{\partial LB_{1c}}{\partial \mu_{Dc}} \frac{\partial \mu_{Dc}}{\partial z_c} \frac{\partial z_c}{\partial \mu_E} \frac{\partial \mu_E}{\partial \phi}$ について

$$\begin{aligned} \frac{\partial LB_{1c}}{\partial \mu_{Dc}} &= (\mathbf{x}_d - \mu_{Dc})^T & \frac{\partial \mu_{Dc}}{\partial z_c} &= \frac{\partial F_\mu}{\partial z_c} & \text{誤差逆伝播で計算} \\ \frac{\partial z_c}{\partial \mu_E} &= \mathbf{I} & \frac{\partial \mu_E}{\partial \phi} &= \frac{\partial G_\mu}{\partial \phi} & \text{誤差逆伝播で計算} \end{aligned}$$

第2項: $\frac{\partial LB_2}{\partial \phi} = \frac{\partial LB_2}{\partial \mu_E} \frac{\partial \mu_E}{\partial \phi}$ について

$$\frac{\partial LB_2}{\partial \mu_E} = -\mu_E^T \quad \frac{\partial \mu_E}{\partial \phi} = \frac{\partial G_\mu}{\partial \phi} \quad \text{誤差逆伝播で計算}$$