

1 近似推論法

1 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 様々な話題

- 変分推論が必要だった理由

- 潜在変数に関する事後分布 $p(Z|X, \theta)$ の計算は、困難であることが多い
- どのような場合に困難になるのか、次の図 1 に示す
- $p(Z|X, \theta)$ の厳密な計算は諦める代わりに、別の確率分布で近似したい
- 別の確率分布で近似するとき、単純な項の積として表現できるといった、何らかの仮定を置く

Figure 19.1: Intractable inference problems in deep learning are usually the result of interactions between latent variables in a structured graphical model. These can be due to edges directly connecting one latent variable to another, or due to longer paths that are activated when the child of a V-structure is observed. *(Left)* A **semi-restricted Boltzmann machine** (Osindero and Hinton, 2008) with connections between hidden units. These direct connections between latent variables make the posterior distribution intractable due to large cliques of latent variables. *(Center)* A deep Boltzmann machine, organized into layers of variables without intra-layer connections, still has an intractable posterior distribution due to the connections between layers. *(Right)* This directed model has interactions between latent variables when the visible variables are observed, because every two latent variables are co-parents. Some probabilistic models are able to provide tractable inference over the latent variables despite having one of the graph structures depicted above. This is possible if the conditional probability distributions are chosen to introduce additional independences beyond those described by the graph. For example, probabilistic PCA has the graph structure shown in the right, yet still has simple inference

図 2: 事後分布 $p(\mathbf{Z}|\mathbf{X})$ の計算が困難な場合

- 事後分布 $p(\mathbf{Z}|\mathbf{X})$ の計算が困難な場合 1
 - グラフィカルモデルにおいて、**潜在変数間の相互作用がある**場合、計算が困難になる
 - 左側の**半制限付きボルツマンマシン**では、全ての潜在変数の組み合わせ間で、接続がある
 - 従って、潜在変数間に**依存関係**が存在し、事後分布の計算が手に負えない
 - 白丸で描かれた潜在変数を $\mathbf{Z} = \{z_1, z_2, z_3\}$ 、灰色で描かれた観測データを \mathbf{X} とおくと、事後分布 $p(\mathbf{Z}|\mathbf{X})$ は、例えば次のようになる

$$\begin{aligned} & p(z_1, z_2, z_3 | \mathbf{X}) \\ = & p(z_1 | \mathbf{X}, z_2, z_3) p(z_2 | \mathbf{X}, z_3) p(z_3 | \mathbf{X}) \end{aligned} \quad (1)$$

- 事後分布 $p(\mathbf{Z}|\mathbf{X})$ の計算が困難な場合 2
 - 中央は、層間の結合がない潜在変数の層で構成される、**深層ボルツマンマシン**を表す
 - 潜在変数の層間の結合があるため、事後分布の計算が手に負えない
- 上の層の潜在変数を $\mathbf{Z}_1 = \{z_{11}, z_{12}, z_{13}\}$ 、中間層の潜在変数を $\mathbf{Z}_2 = \{z_{21}, z_{22}, z_{23}\}$ 、灰色で描かれた観測データを \mathbf{X} とおくと、事後分布は、例えば次のようになる

$$\begin{aligned} & p(\mathbf{Z}_1, \mathbf{Z}_2 | \mathbf{X}) \\ = & p(\mathbf{Z}_2 | \mathbf{X}) p(\mathbf{Z}_1 | \mathbf{Z}_2) \end{aligned} \tag{2}$$

$$\begin{aligned} = & p(z_{21}, z_{22}, z_{23} | \mathbf{X}) p(z_{11}, z_{12}, z_{13} | \mathbf{Z}_2) \\ = & p(z_{21} | \mathbf{X}) p(z_{22} | \mathbf{X}) p(z_{23} | \mathbf{X}) \\ & p(z_{11} | \mathbf{Z}_2) p(z_{12} | \mathbf{Z}_2) p(z_{13} | \mathbf{Z}_2) \end{aligned} \tag{3}$$

1 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 様々な話題

- 変分推論の目的

- 同時分布 $p(\mathbf{X}, \mathbf{Z})$ が求まっているときに、事後分布 $p(\mathbf{Z}|\mathbf{X})$ と、エビデンス $p(\mathbf{X})$ の近似を求める

- 注意点

- 観測変数と潜在変数をそれぞれ、 \mathbf{X}, \mathbf{Z} とおく
- $p(\mathbf{X})$ は、確率モデルからデータ \mathbf{X} が生起する確率である
- データからみたモデルの好みと解釈できるから、 $p(\mathbf{X})$ をモデルエビデンスという

- 周辺分布の対数 $\ln p(\mathbf{X})$ の分解

- EM アルゴリズムのときと同様であり、次のように分解できる

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p) \quad (4)$$

- 但し、 $\mathcal{L}(q)$ と $\text{KL}(q||p)$ は次のように定義した

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (5)$$

$$\text{KL}(q||p) = - \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} d\mathbf{Z} \quad (6)$$

- パラメータ θ の扱い

- EM アルゴリズムとは異なり、パラメータ θ がどこにも現れていない
- パラメータも潜在変数として扱っているので、パラメータベクトルは明示的には書かない

- ここでは、パラメータ θ については、あまり気にしない
- ここでは連続潜在変数について考えるが、離散潜在変数であれば、積分を Z に関する総和に置き換えればよい
- 下界 $\mathcal{L}(q)$ を最適化する動機
 - $\mathcal{L}(q)$ はエビデンスの対数 $\ln p(\mathbf{X})$ の下界であるから、エビデンス下界 (Evidence lower bound, ELBO) ともいう
 - または、負の変分自由エネルギー (Variational free energy) という
 - $\ln p(\mathbf{X})$ は q には依存しないため、定数項とみなせる
 - 従って、 $\mathcal{L}(q)$ を q について最大化することは、 $\text{KL}(q||p)$ の最小化に相当
 - このとき、分布 $q(\mathbf{Z})$ を真の事後分布 $p(\mathbf{Z}|\mathbf{X})$ に近づけられる
 - $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ が分かれば、データ \mathbf{X} から、潜在変数やパラメータ \mathbf{Z} が得られる

- 下界 $\mathcal{L}(q)$ の最適化

- EM アルゴリズムのときと同じように、下界 $\mathcal{L}(q)$ を、分布 $q(\mathbf{Z})$ について最大化する
- これは、KL ダイバージェンス $\text{KL}(q||p)$ を最小化することと等価である
- 従って、もし $q(\mathbf{Z})$ を任意の分布にしてよければ、 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ において、KL ダイバージェンスを 0 にすればよい
- しかしここでは、真の事後分布 $p(\mathbf{Z}|\mathbf{X})$ を求めることは不可能と仮定する

- 分布 $q(\mathbf{Z})$ の近似

- 計算コストを削減するために、 $q(\mathbf{Z})$ の形をある程度制限する
- 制限したクラスの $q(\mathbf{Z})$ の中で、KL ダイバージェンス $\text{KL}(q||p)$ を最小化するものを探す

● 変分推論の目的

- 分布のクラスを制限することで、 $q(\mathbf{Z})$ を計算可能にすること
- 表現力が豊かなクラスを使うことで、真の事後分布 $p(\mathbf{Z}|\mathbf{X})$ を良く近似する
- 計算可能な分布のクラスの中で、**可能な限り豊かな表現力を持つもの**を選びたい
- 表現力が豊かな分布を使うことは、真の事後分布を、精度良く近似することにつながるのであって、従って**過学習は発生しない**
- 分布 $q(\mathbf{Z})$ は、 $q(\mathbf{Z}|\mathbf{X})$ と書くこともある

- 分布 $q(\mathbf{Z})$ のクラスを制限する方法
 - 例えば分布 $q(\mathbf{Z})$ を、**パラメトリックな分布に限定**することができる
 - 即ち、パラメータベクトル ω によって $q(\mathbf{Z}|\omega)$ と記述されるような、分布に制限する
 - 分布 $q(\mathbf{Z})$ を、ガウス分布などの、何らかの特別なパラメトリックな分布と仮定することに相当

- クラスを制限する別の方法 (平均場近似)
 - 分布 $q(\mathbf{Z})$ のクラスを制限する別の方法として、平均場近似がある
 - 潜在変数 \mathbf{Z} を、 M 個の互いに排反なグループ $\{\mathbf{Z}_1, \dots, \mathbf{Z}_M\}$ に分割
 - 分布 $q(\mathbf{Z})$ が、これらのグループによって分解されると仮定

$$q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i) \quad (7)$$

- 分布 q について、これ以上の仮定はしない
- 従って、各因子 $q_i(\mathbf{Z}_i)$ の関数形については、何の制限も課さない
- 平均場近似とは、元々は物理学における用語である

- 下界 $\mathcal{L}(q)$ の最大化

- $q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i)$ と分解できるような分布 $q(\mathbf{Z})$ の中で、**下界 $\mathcal{L}(q)$ を最大にするもの**を探す
- $\mathcal{L}(q)$ を $q(\mathbf{Z})$ について最大化するために、 $\mathcal{L}(q)$ を各因子 $q_i(\mathbf{Z}_i)$ について**順番に最大化**していく
- $\mathcal{L}(q)$ に $q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i)$ を代入して、因子の一つ $q_j(\mathbf{Z}_j)$ に関する**依存項**を抜き出してみよう

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (8)$$

$$= \int \left(\prod_i q_i(\mathbf{Z}_i) \right) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (9)$$

$$= \int \prod_i q_i(\mathbf{Z}_i) \left(\ln p(\mathbf{X}, \mathbf{Z}) - \sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \quad (10)$$

$$= \int \prod_i q_i(\mathbf{Z}_i) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} - \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \quad (11)$$

ここで第 1 項は

$$\begin{aligned} & \int \prod_i q_i(\mathbf{Z}_i) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} \\ &= \int q_j(\mathbf{Z}_j) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} \end{aligned} \quad (12)$$

$d\mathbf{Z} = d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M$ であるから

$$= \int q_j(\mathbf{Z}_j) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (13)$$

$$= \int q_j(\mathbf{Z}_j) (\ln p(\mathbf{X}, \mathbf{Z})) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) \left(\prod_{i \neq j} d\mathbf{Z}_i \right) d\mathbf{Z}_j \quad (14)$$

$$= \int q_j(\mathbf{Z}_j) \left(\int (\ln p(\mathbf{X}, \mathbf{Z})) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \right) d\mathbf{Z}_j \quad (15)$$

$$= \int q_j(\mathbf{Z}_j) \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j \quad (16)$$

- 但し、新しい分布 $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ は以下の式で定義した (積分の結果であるため、定数項が出現する)

$$\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (17)$$

- 記法 $\mathbb{E}_{i \neq j}$ は、 $i \neq j$ をみたす全ての分布 $q_i(\mathbf{Z}_i)$ による、期待値を取ることを表す

$$\mathbb{E}_{i \neq j}[\ln p(\mathbf{X}, \mathbf{Z})] = \int (\ln p(\mathbf{X}, \mathbf{Z})) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (18)$$

- 第2項は

$$\begin{aligned} & \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \\ &= \sum_i \int \prod_i q_i(\mathbf{Z}_i) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z} \end{aligned} \quad (19)$$

$$\begin{aligned} &= \sum_i \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \\ &= \int \prod_k q_k(\mathbf{Z}_k) (\ln q_j(\mathbf{Z}_j)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M + \end{aligned} \quad (20)$$

$$\sum_{i \neq j} \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (21)$$

但し

$$\int \prod_k q_k(\mathbf{Z}_k) (\ln q_j(\mathbf{Z}_j)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (22)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \left(\prod_{k \neq j} q_k(\mathbf{Z}_k) \right) \left(\prod_{k \neq j} d\mathbf{Z}_k \right) d\mathbf{Z}_j \quad (23)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \left(\int \prod_{k \neq j} q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right) d\mathbf{Z}_j \quad (24)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \prod_{k \neq j} \underbrace{\left(\int q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right)}_{=1} d\mathbf{Z}_j \quad (25)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \quad (26)$$

であるほか

$$\begin{aligned} & \sum_{i \neq j} \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \\ &= \sum_{i \neq j} \int q_i(\mathbf{Z}_i) q_j(\mathbf{Z}_j) \left(\prod_{k \neq i, j} \ln q_k(\mathbf{Z}_k) \right) \\ & \quad (\ln q_i(\mathbf{Z}_i)) \left(\prod_{k \neq i, j} d\mathbf{Z}_k \right) d\mathbf{Z}_i d\mathbf{Z}_j \quad (27) \\ &= \sum_{i \neq j} \underbrace{\left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j \right)}_{=1} \left(\int \prod_{k \neq i, j} \ln q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right) \end{aligned}$$

$$\int \ln q_i(\mathbf{Z}_i) q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (28)$$

$$= \sum_{i \neq j} \prod_{k \neq i, j} \underbrace{\left(\int \ln q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right)}_{=\text{Const.}} \underbrace{\int \ln q_i(\mathbf{Z}_i) q_i(\mathbf{Z}_i) d\mathbf{Z}_i}_{=\text{Const.}} \quad (29)$$

$$= \sum_{i \neq j} \text{Const.} = \text{Const.} \quad (30)$$

となるから結局

$$\begin{aligned} & \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \\ &= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{Const.} \end{aligned} \quad (31)$$

- 従って、下界 $\mathcal{L}(q)$ から $q_j(\mathbf{Z}_j)$ に依存する項を取り出すと

$$\begin{aligned}\mathcal{L}(q) = & \int q_j(\mathbf{Z}_j) \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \\ & \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{Const.}\end{aligned}\quad (32)$$

但し

$$\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (33)$$

- $\mathcal{L}(q)$ を、 $i \neq j$ である全ての $q_i(\mathbf{Z}_i)$ について固定した上で、 $q_j(\mathbf{Z}_j)$ について最大化することになる
- $q_j(\mathbf{Z}_j)$ について可能な全ての分布の中で、 $\mathcal{L}(q)$ を最大にするようなものを選ぶ

- $\mathcal{L}(q)$ は次のように変形できる

$$\mathcal{L}(q) = \int q_j(\mathbf{Z}_j) \ln \frac{\tilde{p}(\mathbf{X}, \mathbf{Z}_j)}{q_j(\mathbf{Z}_j)} d\mathbf{Z}_j + \text{Const.} \quad (34)$$

$$= -\text{KL}(q_j(\mathbf{Z}_j) \parallel \tilde{p}(\mathbf{X}, \mathbf{Z}_j)) + \text{Const.} \quad (35)$$

- これより、 $\mathcal{L}(q)$ は $q_j(\mathbf{Z}_j)$ と $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ の間の、負の KL ダイバージェンスとなっている
- そして、 $\mathcal{L}(q)$ の $q_j(\mathbf{Z}_j)$ に関する最大化は、**KL ダイバージェンスの最小化**と等価
- KL ダイバージェンスを最小にするためには、 $q_j(\mathbf{Z}_j) = \tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ とすればよい
- 従って、 $q_j(\mathbf{Z}_j)$ の最適解は、次のように書ける

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j}[\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (36)$$

- 下界 $\mathcal{L}(q)$ を最大化する $\ln q_j(\mathbf{Z}_j)$ の解
 - $q_j(\mathbf{Z}_j)$ の最適解は次のように書けた

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (37)$$

- 上式は次のことを意味している
- 因子 $q_j(\mathbf{Z}_j)$ の最適解の対数 $\ln q_j^*(\mathbf{Z}_j)$ は、観測データ \mathbf{X} と潜在変数 \mathbf{Z} の同時分布の対数 $\ln p(\mathbf{X}, \mathbf{Z})$ を考え、 $i \neq j$ である他の因子 $q_i(\mathbf{Z}_i)$ について期待値を取ったものである
- 定数項は、正規化することで得られるので、結局次のようになる

$$q_j^*(\mathbf{Z}_j) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j} \quad (38)$$

- 正規化定数は必要に応じて計算すればよいので、取り敢えず無視できる

変分推論

- 最適解の式は、 $q(\mathbf{Z})$ の分解の数だけ得られるので、 $\{q_i(\mathbf{Z}_i)\}$ に関する M 本の連立方程式となる
- この方程式は、分布 $q(\mathbf{Z})$ が M 個の因子に分解されるという仮定の下で、下界 $\mathcal{L}(q)$ の最大値が満たすべき条件である
- $\ln q_j^*(\mathbf{Z}_j)$ の右辺は、 $i \neq j$ である $q_i(\mathbf{Z}_i)$ の期待値に依存するため、 $q_j^*(\mathbf{Z}_j)$ を陽に求めることができない
- そこで、下界 $\mathcal{L}(q)$ は次のように最適化される (重要)
- $i \neq j$ である全ての $q_i(\mathbf{Z}_i)$ を固定した状態で、 $q_j(\mathbf{Z}_j)$ を最適化することを、全ての $j = 1, \dots, M$ について繰り返す手続きを、座標降下法という

- 下界 $\mathcal{L}(q)$ の最適化 (座標降下法)

- 1 全ての因子 $q_j(\mathbf{Z}_j)$ を適当に初期化する

- 2 各因子を、以下の式を使って更新する

$$\ln q_j(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (39)$$

$$\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (40)$$

即ち、因子 $q_j(\mathbf{Z}_j)$ を、他の全ての因子の現在の値 $q_i(\mathbf{Z}_i)$ を使って改良する

- 3 (2) を、下界 $\mathcal{L}(q)$ が収束するまで繰り返す

- ここまでの話の流れ

- 1 $\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ であるから、エビデンス下界 $\mathcal{L}(q)$ を q について最大化すれば、 $\text{KL}(q||p) = 0$ とでき、従って $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ を得る
- 2 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ が分かれば、データ \mathbf{X} から、潜在変数やパラメータ \mathbf{Z} が得られる (パラメータは潜在変数 \mathbf{Z} に含まれている)
- 3 しかし、事後分布 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ は計算不可能なので、何らかの方法で近似するしかない
- 4 近似するといっても、計算可能でなければならないので、 $q(\mathbf{Z})$ の形には、通常何らかの制限を課す
- 5 $q(\mathbf{Z})$ を、パラメトリックな分布 $q(\mathbf{Z}|\omega)$ と仮定することがある

- 6 または、 $q(\mathbf{Z})$ を、 $\prod_i q_i(\mathbf{Z}_i)$ のように分解できるとする (平均場近似)
 - 7 平均場近似を行うとき、各因子 $q_j(\mathbf{Z}_j)$ の最適解は、 $\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.}$ であった
 - 8 全ての因子 $\{q_j(\mathbf{Z}_j)\}$ を同時に最適化することはできない
 - 9 下界 $\mathcal{L}(q)$ を、各因子 $q_j(\mathbf{Z}_j)$ について順番に最適化することはできる
- これからの話の流れ
 - MAP 推定と最尤推定は、変分推論の特殊な場合であることを確認する

1 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 様々な話題

- MAP 推定および最尤推定の変分推論からの導出
 - 変分推論で得たいのは、潜在変数に関する事後分布 $p(\mathbf{Z}|\mathbf{X})$ である
 - この変分推論の特殊ケースが、MAP 推定や最尤推定であることを導く
- 変分推論では、エビデンス下界 $\mathcal{L}(q)$ を q について最大化し、従って $\text{KL}(q||p)$ を最小化する
- いま、分布 $q(\mathbf{Z})$ が、次のデルタ関数であるとする
- 特定の値 $\mathbf{Z} = \hat{\mathbf{Z}}$ についてのみ確率が非零になる、無限に鋭い分布

$$q(\mathbf{Z}) = \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \quad (41)$$

- このとき、KL ダイバージェンス $KL(q||p)$ は次のようになる

$$\begin{aligned} & KL(q||p) \\ \equiv & KL(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X})) \\ = & - \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} d\mathbf{Z} \\ = & - \int q(\mathbf{Z}) \ln p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z} + \int q(\mathbf{Z}) \ln q(\mathbf{Z}) d\mathbf{Z} \quad (42) \end{aligned}$$

$$\begin{aligned} = & - \int \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \ln p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z} + \\ & \int \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \ln \delta(\mathbf{Z} - \hat{\mathbf{Z}}) d\mathbf{Z} \quad (43) \end{aligned}$$

$$= - \ln p(\hat{\mathbf{Z}}|\mathbf{X}) + \ln \delta(\hat{\mathbf{Z}} - \hat{\mathbf{Z}}) \quad (44)$$

$$= - \ln p(\hat{\mathbf{Z}}|\mathbf{X}) + \text{Const.} \quad (45)$$

- これを $q(\mathbf{Z})$ について最小化することは、 $\hat{\mathbf{Z}}$ について最小化することに対応する
- よって $\hat{\mathbf{Z}}$ の最適解 $\hat{\mathbf{Z}}^*$ は

$$\hat{\mathbf{Z}}^* = \arg \min_{\hat{\mathbf{Z}}} \left(-\ln p(\hat{\mathbf{Z}}|\mathbf{X}) \right) \quad (46)$$

$$= \arg \max_{\hat{\mathbf{Z}}} \ln p(\hat{\mathbf{Z}}|\mathbf{X}) \quad (47)$$

$$= \arg \max_{\hat{\mathbf{Z}}} p(\hat{\mathbf{Z}}|\mathbf{X}) \quad (48)$$

$$= \arg \max_{\hat{\mathbf{Z}}} p(\mathbf{X}|\hat{\mathbf{Z}})p(\hat{\mathbf{Z}}) \quad (49)$$

となるから、**MAP 推定** (**最大事後確率推定**) の式と一致する

- これより、MAP 推定は、KL ダイバージェンス $\text{KL}(q||p)$ の最小化、従って、**エビデンス下界 $\mathcal{L}(q)$ の最大化と等価**である

- 更に、 $p(\hat{\mathbf{Z}}) = \text{Const.}$ 、即ち $\hat{\mathbf{Z}}$ に関する事前の情報がないとすると

$$\begin{aligned}\hat{\mathbf{Z}}^* &= \arg \max_{\hat{\mathbf{Z}}} p(\mathbf{X}|\hat{\mathbf{Z}})p(\hat{\mathbf{Z}}) \\ &= \arg \max_{\hat{\mathbf{Z}}} p(\mathbf{X}|\hat{\mathbf{Z}})\end{aligned}\tag{50}$$

となるから、これは最尤推定の式に一致する

- 下界 $\mathcal{L}(q)$ の式から導くことも可能である

$$\begin{aligned}\mathcal{L}(q) &= \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \\ &= \int q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} - \int q(\mathbf{Z}) \ln q(\mathbf{Z}) d\mathbf{Z} \quad (51)\end{aligned}$$

$$\begin{aligned}&= \int \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \ln p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} - \\ &\quad \int \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \ln \delta(\mathbf{Z} - \hat{\mathbf{Z}}) d\mathbf{Z} \quad (52)\end{aligned}$$

$$= \ln p(\mathbf{X}, \hat{\mathbf{Z}}) - \ln \delta(\hat{\mathbf{Z}} - \hat{\mathbf{Z}}) \quad (53)$$

$$= \ln p(\mathbf{X}, \hat{\mathbf{Z}}) + \text{Const.} \quad (54)$$

$$= \ln p(\mathbf{X} | \hat{\mathbf{Z}}) p(\hat{\mathbf{Z}}) + \text{Const.} \quad (55)$$

- これより、下界 $\mathcal{L}(q)$ の最大化は、 $\hat{\mathbf{Z}}$ に関する $p(\hat{\mathbf{Z}} | \mathbf{X}) \propto p(\mathbf{X} | \hat{\mathbf{Z}}) p(\hat{\mathbf{Z}})$ の最大化、**即ち MAP 推定と等価**である

MAP 推定の例

- スパース符号化

- ここでは MAP 推定の例として、スパース符号化を扱う (変分推論の例ではない)
- データ x に対する潜在変数 z に、スパース性を持たせる
- そのために、スパース性を導く事前分布 (ラプラス事前分布) を、潜在変数に用いる
- データ x を D 次元、また潜在変数 z を K 次元とする
- データ x に対応する潜在変数を z_i と表し、 z_i の第 k 成分を z_{ik} とする

$$p(z_{ik}|\lambda) = \frac{\lambda}{2} \exp(-\lambda|z_{ik}|) \quad (56)$$

$$p(z_i|\lambda) = \prod_k p(z_{ik}) = \prod_k \frac{\lambda}{2} \exp(-\lambda|z_{ik}|) \quad (57)$$

MAP 推定の例

- データに関する事後分布 $p(\mathbf{x}_i | \mathbf{z}_i, \mathbf{W}, \mathbf{b}, \beta)$ を、次で定義する

$$\begin{aligned} & p(\mathbf{x}_i | \mathbf{z}_i, \mathbf{W}, \mathbf{b}, \beta) \\ &= \mathcal{N}(\mathbf{x}_i | \mathbf{W}\mathbf{z}_i + \mathbf{b}, \beta^{-1}\mathbf{I}) \end{aligned} \quad (58)$$

$$\begin{aligned} &= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\beta^{-1}\mathbf{I}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - (\mathbf{W}\mathbf{z}_i + \mathbf{b}))^T \right. \\ &\quad \left. (\beta^{-1}\mathbf{I})^{-1} (\mathbf{x}_i - (\mathbf{W}\mathbf{z}_i + \mathbf{b})) \right) \end{aligned} \quad (59)$$

$$\begin{aligned} &= \frac{1}{(2\pi\beta^{-1})^{\frac{D}{2}}} \exp \left(-\frac{\beta}{2} (\mathbf{x}_i - (\mathbf{W}\mathbf{z}_i + \mathbf{b}))^T \right. \\ &\quad \left. (\mathbf{x}_i - (\mathbf{W}\mathbf{z}_i + \mathbf{b})) \right) \end{aligned} \quad (60)$$

- データ \mathbf{x}_i は、対応する潜在変数 \mathbf{z}_i に、線形変換 $\mathbf{W}\mathbf{z}_i + \mathbf{b}$ を施し、更に分散 $\beta^{-1}\mathbf{I}$ のガウスノイズを足すことで、生成されると考える

MAP 推定の例

- パラメータについての補足

- 今回は、 $b = 0$ において**バイアス**を**無視**したものを考える
- λ と、精度 β は**ハイパーパラメータ**であり、予め決められているとする
- そこでこれ以降、次のように確率分布を表現する

$$p(z_i|\lambda) = p(z_i) \quad (61)$$

$$p(x_i|z_i, \mathbf{W}, \mathbf{b}, \beta) = p(x_i|z_i, \mathbf{W}) \quad (62)$$

- \mathbf{X} を、全てのデータ $\{x_i\}$ を集めた行列とする (第 i 行ベクトルが x_i^T)
- \mathbf{Z} についても同様に、全ての潜在変数 $\{z_i\}$ を集めた行列として定める (第 i 行ベクトルが z_i^T)

- スパース符号化の学習

- 事後分布 $p(z_i|x_i)$ は、**表現することすら困難**であるため、最尤推定による手法 (EM アルゴリズムなど) は利用できない

MAP 推定の例

- そこで、最尤推定の代わりに **MAP 推定** を利用することで、最適なパラメータが得られる
- 最大化するのは、以下の事後分布 $p(\mathbf{Z}|\mathbf{X}, \mathbf{W})$ である

$$p(\mathbf{Z}|\mathbf{X}, \mathbf{W}) = \frac{p(\mathbf{X}, \mathbf{Z}|\mathbf{W})}{p(\mathbf{X})} \quad (63)$$

$$\propto p(\mathbf{X}, \mathbf{Z}|\mathbf{W}) \quad (64)$$

$$= p(\mathbf{X}|\mathbf{Z}, \mathbf{W})p(\mathbf{Z}) \quad (65)$$

$$= \left(\prod_i p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) \right) \left(\prod_i p(\mathbf{z}_i) \right) \quad (66)$$

$$= \prod_i p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W})p(\mathbf{z}_i) \quad (67)$$

MAP 推定の例

- 対数を取って最大化してもよいので、最大化する量は結局

$$\ln p(\mathbf{Z}|\mathbf{X}, \mathbf{W}) = \ln \left(\prod_i p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) p(\mathbf{z}_i) \right) \quad (68)$$

$$= \sum_i (\ln p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) + \ln p(\mathbf{z}_i)) \quad (69)$$

$$= \sum_i \ln p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) + \sum_i \ln p(\mathbf{z}_i) \quad (70)$$

各項を求めると

$$\begin{aligned} & \sum_i \ln p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) \\ = & \sum_i \ln \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i, \beta^{-1}\mathbf{I}) \\ = & \sum_i \ln \left(\frac{1}{(2\pi\beta^{-1})^{\frac{D}{2}}} \exp \left(-\frac{\beta}{2} (\mathbf{x}_i - \mathbf{W}\mathbf{z}_i)^T (\mathbf{x}_i - \mathbf{W}\mathbf{z}_i) \right) \right) \end{aligned}$$

MAP 推定の例

$$\begin{aligned} &= \sum_i \left(-\frac{D}{2} \ln 2\pi - \frac{D}{2} \ln \beta^{-1} - \right. \\ &\quad \left. \frac{\beta}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) \right) \\ &= \sum_i \left(-\frac{\beta}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) \right) + \text{Const.} \\ &= -\frac{\beta}{2} \sum_i (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) + \text{Const.} \\ &= -\frac{\beta}{2} \sum_i (\mathbf{X} - \mathbf{Z} \mathbf{W}^T)_{i:} ((\mathbf{X} - \mathbf{Z} \mathbf{W}^T)_{i:})^T + \text{Const.} \quad (71) \end{aligned}$$

$$= -\frac{\beta}{2} \sum_{i,j} ((\mathbf{X} - \mathbf{Z} \mathbf{W}^T) \odot (\mathbf{X} - \mathbf{Z} \mathbf{W}^T))_{i,j} + \text{Const.} \quad (72)$$

$$= -\frac{\beta}{2} \sum_{i,j} (\mathbf{X} - \mathbf{Z} \mathbf{W}^T)_{i,j}^2 + \text{Const.} \quad (73)$$

MAP 推定の例

$$= -\frac{\beta}{2} \|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_2^2 + \text{Const.} \quad (74)$$

また

$$\begin{aligned} & \sum_i \ln p(\mathbf{z}_i) \\ &= \sum_i \ln \prod_j \left(\frac{\lambda}{2} \exp(-\lambda |z_{ij}|) \right) \\ &= \sum_i \sum_j \left(\ln \frac{\lambda}{2} - \lambda |z_{ij}| \right) \\ &= -\lambda \sum_{i,j} |z_{ij}| + \text{Const.} \end{aligned} \quad (75)$$

$$= -\lambda \sum_{i,j} |\mathbf{Z}_{i,j}| + \text{Const.} \quad (76)$$

$$= -\lambda \|\mathbf{Z}\|_1 + \text{Const.} \quad (77)$$

MAP 推定の例

のようになる

- \odot はアダマール積、 $A_{i:}$ は行列 A の第 i 行目のベクトルを表す
- また、行列のノルム $\|A\|_p$ は次で定義される

$$\|A\|_p = \left(\sum_{i,j} |A_{i,j}|^p \right)^{\frac{1}{p}} \quad (78)$$

$p = 1, p = 2$ のときは次のようになる

$$\|A\|_1 = \sum_{i,j} |A_{i,j}|, \quad \|A\|_2 = \sqrt{\sum_{i,j} A_{i,j}^2} \quad (79)$$

- これより、 $\|A\|_1$ は、**行列 A の各要素の絶対値の和**を表す
- また $\|A\|_2$ は、行列 A の各要素の二乗和の平方根を表し、**フロベニウスノルム $\|A\|_F$** ともよばれる

MAP 推定の例

- 従って、 $\|\mathbf{A}\|_2^2$ は、**行列 \mathbf{A} の各要素の二乗和**である
- 上記から

$$\begin{aligned} & \ln p(\mathbf{Z}|\mathbf{X}, \mathbf{W}) \\ &= \sum_i \ln p(\mathbf{x}_i|z_i, \mathbf{W}) + \sum_i \ln p(z_i) \\ &= -\frac{\beta}{2}\|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_2^2 - \lambda\|\mathbf{Z}\|_1 + \text{Const.} \end{aligned} \quad (80)$$

となるので、 $\ln p(\mathbf{Z}|\mathbf{X}, \mathbf{W})$ の最大化は、以下の**関数 $J(\mathbf{Z}, \mathbf{W})$ の最小化**である

$$J(\mathbf{Z}, \mathbf{W}) = \|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_2^2 + \|\mathbf{Z}\|_1 \quad (81)$$

- スパース符号化の学習方法のまとめ
 - 関数 $J(\mathbf{Z}, \mathbf{W})$ を、 \mathbf{Z} と \mathbf{W} について**交互に最小化**すればよい

MAP 推定の例

- $p(\mathbf{Z}|\mathbf{X}, \mathbf{W})$ の最大化は、 $p(\mathbf{X}|\mathbf{Z}, \mathbf{W})p(\mathbf{Z})$ の最大化、即ち $J(\mathbf{Z}, \mathbf{W})$ の最小化と等価であった
- 従って、 $J(\mathbf{Z}, \mathbf{W})$ の各パラメータ \mathbf{Z}, \mathbf{W} についての最小化は、事後確率を大きくする方向に働く
- 関数 $J(\mathbf{Z}, \mathbf{W})$ をもう一度みてみよう

$$J(\mathbf{Z}, \mathbf{W}) = \|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_2^2 + \|\mathbf{Z}\|_1 \quad (82)$$

- 第1項は、明らかに再構成誤差を表している
- 第2項は、データ \mathbf{X} の表現 \mathbf{Z} がスパースになるように付加した、正則化項である
- $\mathbf{Z}\mathbf{W}^T$ は、データ \mathbf{X} の内部表現 \mathbf{Z} に、重み \mathbf{W} を掛けることで、データ \mathbf{X} を復元したもの

1 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 様々な話題

- これまでの話の流れ

- 1 変分推論から、MAP 推定と最尤推定が導出できることを確認した
- 2 MAP 推定は、 $q(\mathbf{Z})$ を無限に鋭い確率分布 (デルタ関数) として、 \mathbf{Z} が特定の値しか取らないと仮定した場合であった
- 3 最尤推定は、MAP 推定において、 \mathbf{Z} の事前確率分布を設けない場合であった
- 4 MAP 推定の例として、スパース符号化を扱った

- これからの話の流れ

- 話題を変えて、分解 $q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i)$ によって $p(\mathbf{Z}|\mathbf{X})$ を近似するときの、弊害を調べる

- $q(\mathbf{Z})$ の分解による近似の性質
 - 変分推論では、真の事後分布 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ を、分解により近似する
 - 分解で近似することによって、**どのような不正確さが生じるのか?**
- ガウス分布の分解による近似
 - ガウス分布を、**分解されたガウス分布**で近似することを考えてみよう
 - 分解による近似で、どのような問題が起こるのか考えてみよう
 - 2つの変数 $\mathbf{z} = (z_1, z_2)$ 間には、**相関がある**とする
 - \mathbf{z} はガウス分布 $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ に従っているとする

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \boldsymbol{\Lambda} = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \quad (83)$$

- 精度行列 $\boldsymbol{\Lambda}$ は対称行列であるから、 $\Lambda_{12} = \Lambda_{21}$

- この分布 $p(\mathbf{z})$ を、分解されたガウス分布 $q(\mathbf{z}) = q_1(z_1)q_2(z_2)$ で近似する
- 各因子 $q_1(z_1), q_2(z_2)$ の関数形については何の仮定も置いていないことに注意

- 最適な因子 $q_1(z_1), q_2(z_2)$ の計算

- 最適な因子 $q_1^*(z_1)$ を、先程の結果を使って求める

$$\ln q_j(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (84)$$

- 従って、 $q_1^*(z_1)$ を計算する式は次のようになる

$$\ln q_1^*(z_1) = \mathbb{E}_{z_2} [\ln p(\mathbf{z})] + \text{Const.} \quad (85)$$

$$\mathbb{E}_{z_2} [\ln p(\mathbf{z})] = \int \ln p(\mathbf{z}) q_2(z_2) dz_2 \quad (86)$$

- 上式の右辺では、 z_1 に依存する項だけを考えればよい

- z_1 の関数を求めようとしているため
- z_1 に依存しない項は、全て定数項 (正規化定数) に含まれてしまうため
- 従って $q_1^*(z_1)$ は

$$\ln q_1^*(z_1) = \mathbb{E}_{z_2} [\ln p(\mathbf{z})] + \text{Const.} \quad (87)$$

$$= \mathbb{E}_{z_2} [\ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})] + \text{Const.} \quad (88)$$

但し

$$\begin{aligned} & \ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ = & \ln \left(\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Lambda}^{-1}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T (\boldsymbol{\Lambda}^{-1})^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right) \right) \\ = & \ln \left(\frac{1}{2\pi} \frac{1}{|\boldsymbol{\Lambda}|^{-\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{z} - \boldsymbol{\mu}) \right) \right) \\ = & -\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{z} - \boldsymbol{\mu}) + \text{Const.} \end{aligned} \quad (89)$$

z_1 に依存する項だけを取り出せば

$$\begin{aligned}
 & -\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{z} - \boldsymbol{\mu}) \\
 = & -\frac{1}{2} [z_1 - \mu_1, z_2 - \mu_2] \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \begin{bmatrix} z_1 - \mu_1 \\ z_2 - \mu_2 \end{bmatrix} \\
 = & -\frac{1}{2} \left[\left\{ (\Lambda_{11}(z_1 - \mu_1) + \Lambda_{21}(z_2 - \mu_2)) \right\} (z_1 - \mu_1) + \right. \\
 & \left. \left\{ \Lambda_{12}(z_1 - \mu_1) + \Lambda_{22}(z_2 - \mu_2) \right\} (z_2 - \mu_2) \right] \\
 = & -\frac{1}{2} (\Lambda_{11}(z_1 - \mu_1)^2 + 2\Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \\
 & \Lambda_{22}(z_2 - \mu_2)^2) \quad (\because \Lambda_{21} = \Lambda_{12}) \\
 = & -\frac{1}{2} \Lambda_{11}(z_1 - \mu_1)^2 - \Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \text{Const.} \quad (90)
 \end{aligned}$$

これを代入して

$$\begin{aligned} & \ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ = & -\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 - \Lambda_{12} (z_1 - \mu_1) (z_2 - \mu_2) + \text{Const.} \quad (91) \end{aligned}$$

従って

$$\begin{aligned} & \ln q_1^*(z_1) \\ = & \mathbb{E}_{z_2} [\ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})] + \text{Const.} \\ = & \mathbb{E}_{z_2} \left[-\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 - \Lambda_{12} (z_1 - \mu_1) (z_2 - \mu_2) \right] + \text{Const.} \\ = & -\frac{1}{2} \Lambda_{11} z_1^2 + \Lambda_{11} \mu_1 z_1 - \Lambda_{12} z_1 (\mathbb{E}[z_2] - \mu_2) + \text{Const.} \quad (92) \end{aligned}$$

- これより、 $q_1^*(z_1)$ は次のように書ける

$$\begin{aligned} & q_1^*(z_1) \\ \propto & \exp\left(-\frac{1}{2}\Lambda_{11}z_1^2 + \Lambda_{11}\mu_1z_1 - \Lambda_{12}z_1(\mathbb{E}[z_2] - \mu_2)\right) \\ = & \exp\left(-\frac{1}{2}\Lambda_{11}\left(z_1 - (\mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(\mathbb{E}[z_2] - \mu_2))\right)^2 + \right. \\ & \left. \frac{1}{2}\Lambda_{11}\left(\mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(\mathbb{E}[z_2] - \mu_2)\right)^2\right) \\ \propto & \exp\left(-\frac{1}{2\Lambda_{11}^{-1}}\left(z_1 - (\mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(\mathbb{E}[z_2] - \mu_2))\right)^2\right) \\ = & \mathcal{N}(z_1|m_1, \Lambda_{11}^{-1}) \end{aligned} \tag{93}$$

但し m_1 は次のようにおいた

$$m_1 = \mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(\mathbb{E}[z_2] - \mu_2) \tag{94}$$

- 対称性から、 $q_2^*(z_2)$ も次のように求められる

$$\ln q_2^*(z_2) = \mathbb{E}_{z_2}[\ln p(z)] + \text{Const.} \quad (95)$$

$$= \mathbb{E}_{z_2}[\ln \mathcal{N}(z|\mu, \Lambda^{-1})] + \text{Const.} \quad (96)$$

$$q_2^*(z_2) = \mathcal{N}(z_2|m_2, \Lambda_{22}^{-1}) \quad (97)$$

但し m_2 は次のようにおいた

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21}(\mathbb{E}[z_1] - \mu_1) \quad (98)$$

- これより $q_1^*(z_1), q_2^*(z_2)$ は

$$q_1^*(z_1) = \mathcal{N}(z_1|m_1, \Lambda_{11}^{-1}) \quad (99)$$

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12}(\mathbb{E}[z_2] - \mu_2) \quad (100)$$

$$q_2^*(z_2) = \mathcal{N}(z_2|m_2, \Lambda_{22}^{-1}) \quad (101)$$

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21}(\mathbb{E}[z_1] - \mu_1) \quad (102)$$

- $\mathbb{E}[z_2]$ は、 z_2 の $q_2(z_2) = \mathcal{N}(z_2|m_2, \Lambda_{22}^{-1})$ による平均であるから、
 $\mathbb{E}[z_2] = m_2$ である (同様に、 $\mathbb{E}[z_1] = m_1$)
- これより m_1, m_2 を連立させれば

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (m_2 - \mu_2) \quad (103)$$

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21} (m_1 - \mu_1) \quad (104)$$

であるから、 m_1 について解けば

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (m_2 - \mu_2) \quad (105)$$

$$= \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (\mu_2 - \Lambda_{22}^{-1} \Lambda_{21} (m_1 - \mu_1) - \mu_2) \quad (106)$$

$$= \mu_1 + \Lambda_{11}^{-1} \Lambda_{12} \Lambda_{22}^{-1} \Lambda_{21} (m_1 - \mu_1) \quad (107)$$

$$= \mu_1 + \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2 (m_1 - \mu_1) \quad (108)$$

$$= (1 - \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2) \mu_1 + \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2 m_1 \quad (109)$$

従って

$$(1 - \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2) \mu_1 = (1 - \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2) m_1 \quad (110)$$

精度行列 Λ が正則であれば

$$\Lambda_{11} \Lambda_{12} - \Lambda_{12} \Lambda_{21} \neq 0 \quad (111)$$

$$\Rightarrow \Lambda_{11} \Lambda_{12} - \Lambda_{12}^2 \neq 0 \quad (112)$$

$$\Rightarrow (\Lambda_{11} \Lambda_{12}) (1 - \Lambda_{11}^{-1} \Lambda_{12}^{-1} \Lambda_{12}^2) \neq 0 \quad (113)$$

$$\Rightarrow (1 - \Lambda_{11}^{-1} \Lambda_{12}^{-1} \Lambda_{12}^2) \neq 0 \quad (114)$$

が成立するから、分布 $p(z)$ が非特異 (精度行列が正則) ならば

$$m_1 = \mu_1 \quad (115)$$

が唯一の解であるほか、対称性から、 m_2 についても以下を得る

$$m_2 = \mu_2 \quad (116)$$

- ゆえに、 $q_1^*(z_1), q_2^*(z_2)$ は

$$q_1^*(z_1) = \mathcal{N}(z_1 | \mu_1, \Lambda_{11}^{-1}) \quad (117)$$

$$q_2^*(z_2) = \mathcal{N}(z_2 | \mu_2, \Lambda_{22}^{-1}) \quad (118)$$

- 注意点 1

- $q_1^*(z_1)$ は、 $q_2^*(z_2)$ を使って計算される $p(z)$ の期待値 $\mathbb{E}[z_2]$ に依存 (逆も成り立つ)
- $q_1^*(z_1)$ と、 $q_2^*(z_2)$ は相互に依存しているため、2つを同時に求めることはできない
- その代わりに、次のように最適化すればよい
- $q_1(z_1), q_2(z_2)$ を適当に初期化したあと、 $q_1^*(z_1), q_2^*(z_2)$ の式を使って、交互に $q_1(z_1)$ と $q_2(z_2)$ を更新していく (収束するまでこれを繰り返す)

- 注意点 2

変分推論

- $q_1(z_1), q_2(z_2)$ の具体的な関数形については、何の仮定も置かなかった
- $q_i^*(z_i)$ がガウス分布だという仮定は置いていないが、 $\text{KL}(q||p)$ を最適化する変分推論によって、結果的にガウス分布が得られた

- $KL(q||p)$ の最適化と $KL(p||q)$ の最適化の比較

- 上記の結果は、 $KL(q||p)$ の最適化 (エビデンス下界 $\mathcal{L}(q)$ の最適化) によって得た
- $KL(q||p)$ ではなく、 $KL(p||q)$ を最適化したらどうなるか?
- 変分推論ではない、もう一つの近似推論の方法である、EP 法で使われる考え方
- $q(Z)$ を $p(Z|X)$ に近づけたいのであれば、 $KL(q||p)$ と $KL(p||q)$ のどちらを最小化してもよいはず
- なぜなら、KL ダイバージェンスは、確率分布間の (擬似的な) 距離を表すため

- $KL(p||q)$ の最適化

- $q(Z)$ が平均場近似によって分解できるとき、 $KL(p||q)$ を最適化したい

- KL ダイバージェンス $KL(p||q)$ は、次のように書ける

$$KL(p||q) \equiv KL(p(\mathbf{Z}|\mathbf{X})||q(\mathbf{Z})) \quad (119)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})} d\mathbf{Z} \quad (120)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) (\ln q(\mathbf{Z}) - \ln p(\mathbf{Z}|\mathbf{X})) d\mathbf{Z} \quad (121)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln q(\mathbf{Z}) d\mathbf{Z} - \underbrace{\int p(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z}}_{q \text{ には依存しない定数項}} \quad (122)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln \prod_i q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (123)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \sum_i \ln q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (124)$$

$$= - \sum_i \int p(\mathbf{Z}|\mathbf{X}) \ln q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (125)$$

定数項は、 $p(\mathbf{Z}|\mathbf{X})$ のエントロピーであり、 q には依存しない

- 各因子 $q_j(\mathbf{Z}_j)$ について $\text{KL}(p||q)$ を最適化したい
- このとき、 $i \neq j$ となる、全ての $q_i(\mathbf{Z}_i)$ は**固定する**
- $q_j(\mathbf{Z}_j)$ に依存する項を取り出せば、次のようになる

$$\begin{aligned} & \sum_i \int p(\mathbf{Z}|\mathbf{X}) \ln q_i(\mathbf{Z}_i) d\mathbf{Z} \\ &= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z} \end{aligned} \quad (126)$$

$$= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (127)$$

$$= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \left(\prod_{i \neq j} d\mathbf{Z}_i \right) \quad (128)$$

$$= \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \quad (129)$$

- $q_j(\mathbf{Z}_j)$ は確率分布であるから、以下の条件を満たさなければならない

$$\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j = 1 \quad (\text{規格化条件}) \quad (130)$$

$$q_j(\mathbf{Z}_j) \geq 0 \quad (131)$$

- 従って $\text{KL}(p||q)$ を $q_j(\mathbf{Z}_j)$ について最適化するとき、**ラグランジュの未定乗数法**を使って、規格化条件を組み込む必要がある

- $q_j(\mathbf{Z}_j) \geq 0$ という条件は、 $\ln q_i(\mathbf{Z}_i)$ という項が既にあるから、何もしなくても常に満たされる (ラグランジュ関数に、制約条件を改めて取り入れる必要がない)
- 結局、ラグランジュ汎関数 $\mathcal{L}[q_j]$ は、次のようになる

$$\begin{aligned} \mathcal{L}[q_j] = & - \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \\ & \lambda \left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j - 1 \right) \end{aligned} \quad (132)$$

- 上記は $q_j(\mathbf{Z}_j)$ についての汎関数となっていることに注意
- 次の公式を使って、 $\mathcal{L}[q_j]$ を変分最適化する

$$\frac{\delta}{\delta y(x)} \int G(y(x), x) dx = \frac{\partial}{\partial y} G(y(x), x) \quad (133)$$

- 従って

$$\begin{aligned}
 & \frac{\delta}{\delta q_j(\mathbf{Z}_j)} \mathcal{L}[q_j] \\
 = & -\frac{\delta}{\delta q_j(\mathbf{Z}_j)} \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \\
 & \frac{\delta}{\delta q_j(\mathbf{Z}_j)} \lambda \left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j - 1 \right) \quad (134)
 \end{aligned}$$

$$\begin{aligned}
 = & -\frac{\partial}{\partial q_j} \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) + \\
 & \lambda \frac{\partial}{\partial q_j} q_j(\mathbf{Z}_j) \quad (135)
 \end{aligned}$$

$$= - \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} + \lambda = 0 \quad (136)$$

- これより、未定乗数 λ は

$$\lambda q_j(\mathbf{Z}_j) = \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} \quad (137)$$

$$\Rightarrow \int \lambda q_j(\mathbf{Z}_j) d\mathbf{Z}_j = \int \underbrace{\left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right)}_{=p(\mathbf{Z}_j|\mathbf{X})} d\mathbf{Z}_j \quad (138)$$

$$\Rightarrow \lambda \underbrace{\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j}_{=1} = \underbrace{\int p(\mathbf{Z}_j|\mathbf{X}) d\mathbf{Z}_j}_{=1} \quad (139)$$

$$\Rightarrow \lambda = 1 \quad (140)$$

- 結局、最適解 $q_j^*(\mathbf{Z}_j)$ は次のようになる

$$- \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} + \lambda = 0 \quad (141)$$

$$\Rightarrow q_j^*(\mathbf{Z}_j) = \underbrace{\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i}_{=p(\mathbf{Z}_j|\mathbf{X})} \quad (142)$$

$$\Rightarrow q_j^*(\mathbf{Z}_j) = p(\mathbf{Z}_j|\mathbf{X}) \quad (143)$$

- $q_j^*(\mathbf{Z}_j)$ の最適解は、 $p(\mathbf{Z}|\mathbf{X})$ を、 $i \neq j$ である全ての \mathbf{Z}_i について周辺化した分布
- これは閉じた解であり、繰り返しを必要としない

- 最適な因子 $q_1(z_1), q_2(z_2)$ の計算

- 今回は $p(\mathbf{z}|\mathbf{x}) = p(\mathbf{z})$ の場合を考えており、かつ $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ であった
- 従って、 $q_1^*(z_1)$ は、 $p(\mathbf{z})$ を z_2 について周辺化すればよいから

$$\begin{aligned} & q_1^*(z_1) \\ = & \int p(\mathbf{z}) dz_2 \\ = & \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) dz_2 \end{aligned} \tag{144}$$

$$= \frac{1}{2\pi} \frac{1}{|\boldsymbol{\Lambda}|^{-\frac{1}{2}}} \int \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{z} - \boldsymbol{\mu})\right) dz_2 \tag{145}$$

- ここで、指数の内側を、積分変数 z_2 に依存する項と、そうでない項に分ける

$$\begin{aligned} & -\frac{1}{2}(z - \mu)^T \Lambda (z - \mu) \\ = & -\frac{1}{2} (\Lambda_{11}(z_1 - \mu_1)^2 + \\ & 2\Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \Lambda_{22}(z_2 - \mu_2)^2) \quad (146) \end{aligned}$$

$$\begin{aligned} = & -\frac{1}{2} \Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 + \\ & -\frac{1}{2} (2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}(z_2 - \mu_2)^2) \quad (147) \end{aligned}$$

そして

$$\begin{aligned} & -\frac{1}{2} (2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}(z_2 - \mu_2)^2) \\ = & -\frac{1}{2} (\Lambda_{22}z_2^2 - 2\Lambda_{22}\mu_2z_2 + 2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}\mu_2^2) \quad (148) \end{aligned}$$

$$= -\frac{1}{2} (\Lambda_{22} z_2^2 - 2 (\Lambda_{22} \mu_2 - \Lambda_{12} (z_1 - \mu_1)) z_2 + \Lambda_{22} \mu_2^2) \quad (149)$$

$$= -\frac{1}{2} \left(\Lambda_{22} (z_2 - \Lambda_{22}^{-1} (\Lambda_{22} \mu_2 - \Lambda_{12} (z_1 - \mu_1)))^2 - \Lambda_{22} (\Lambda_{22}^{-1} (\Lambda_{22} \mu_2 - \Lambda_{12} (z_1 - \mu_1)))^2 + \Lambda_{22} \mu_2^2 \right) \quad (150)$$

$$= -\frac{1}{2} \left(\Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 - \Lambda_{22}^{-1} m^2 + \Lambda_{22} \mu_2^2 \right) \quad (151)$$

ゆえ ($m = \Lambda_{22} \mu_2 - \Lambda_{12} (z_1 - \mu_1)$ とおいた)

$$\begin{aligned} & -\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{z} - \boldsymbol{\mu}) \\ = & -\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 + \Lambda_{12} (z_1 - \mu_1) \mu_2 - \frac{1}{2} \Lambda_{22} \mu_2^2 - \\ & \frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 + \frac{1}{2} \Lambda_{22}^{-1} m^2 \end{aligned} \quad (152)$$

- これより、積分変数 z_2 の依存項だけを取り出せたので

$$\begin{aligned} & \int \exp \left(-\frac{1}{2} (z - \mu)^T \Lambda (z - \mu) \right) dz_2 \\ = & \exp \left(-\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 + \Lambda_{12} (z_1 - \mu_1) \mu_2 - \right. \\ & \quad \left. \frac{1}{2} \Lambda_{22} \mu_2^2 + \frac{1}{2} \Lambda_{22}^{-1} m^2 \right) \\ & \int \exp \left(-\frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 \right) dz_2 \end{aligned} \quad (153)$$

であって、右側の積分は、中身が (正規化されていない) ガウス分布であるから

$$\int \exp \left(-\frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 \right) dz_2$$

$$\begin{aligned}
 &= (2\pi\Lambda_{22}^{-1})^{\frac{1}{2}} \cdot \int \frac{1}{(2\pi\Lambda_{22}^{-1})^{\frac{1}{2}}} \exp\left(-\frac{1}{2}\Lambda_{22}(z_2 - \Lambda_{22}^{-1}m)^2\right) dz_2 \\
 &= (2\pi\Lambda_{22}^{-1})^{\frac{1}{2}} \tag{154}
 \end{aligned}$$

となって、 z_2 を積分により消去できる

- また指数の残りの部分から、 z_1 に依存する項だけを取り出して

$$\begin{aligned}
 &-\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 - \frac{1}{2}\Lambda_{22}\mu_2^2 + \frac{1}{2}\Lambda_{22}^{-1}m^2 \\
 = &-\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 - \frac{1}{2}\Lambda_{22}\mu_2^2 + \\
 &\quad \frac{1}{2}\Lambda_{22}^{-1}(\Lambda_{22}\mu_2 - \Lambda_{12}(z_1 - \mu_1))^2 \\
 = &-\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 - \frac{1}{2}\Lambda_{22}\mu_2^2 + \\
 &\quad \frac{1}{2}\Lambda_{22}\mu_2^2 - \mu_2\Lambda_{12}(z_1 - \mu_1) +
 \end{aligned}$$

$$\frac{1}{2} \Lambda_{22}^{-1} \Lambda_{12}^2 (z_1 - \mu_1)^2 \quad (155)$$

$$= -\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) (z_1 - \mu_1)^2 \quad (156)$$

$$= -\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) z_1^2 + \frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) \mu_1 z_1 + \text{Const.} \quad (157)$$

- これより結局、 z_2 による積分は次のようになる

$$\begin{aligned} & \int \exp \left(-\frac{1}{2} (z - \mu)^T \Lambda (z - \mu) \right) dz_2 \\ &= (2\pi \Lambda_{22}^{-1})^{\frac{1}{2}} \exp \left(-\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) (z_1 - \mu_1)^2 \right) \end{aligned} \quad (158)$$

- 従って、 $q_1^*(z_1)$ は次のようになる

$$= \frac{q_1^*(z_1)}{2\pi |\mathbf{\Lambda}|^{-\frac{1}{2}}} \int \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \mathbf{\Lambda}(\mathbf{z} - \boldsymbol{\mu})\right) dz_2 \quad (159)$$

$$= \frac{1}{2\pi} \frac{1}{(\Lambda_{11}\Lambda_{22} - \Lambda_{12}^2)^{-\frac{1}{2}}} (2\pi\Lambda_{22}^{-1})^{\frac{1}{2}} \exp\left(-\frac{1}{2}(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)(z_1 - \mu_1)^2\right) \quad (160)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)^{-\frac{1}{2}}} \exp\left(-\frac{1}{2}(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)(z_1 - \mu_1)^2\right) \quad (161)$$

$$= \mathcal{N}(z_1 | \mu_1, (\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)^{-1}) \quad (162)$$

- 共分散行列 Σ を、精度行列 Λ を使って次のように定めれば

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix} \quad (\Sigma_{12} = \Sigma_{21}) \quad (163)$$

次が成り立つから

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix} = \mathbf{I} \quad (164)$$

各成分に注目すれば

$$\begin{cases} \Sigma_{11}\Lambda_{11} + \Sigma_{12}\Lambda_{21} = 1 \\ \Sigma_{11}\Lambda_{12} + \Sigma_{12}\Lambda_{22} = 0 \end{cases} \quad (165)$$

これを Σ_{11} について解けば

$$\Sigma_{11}\Lambda_{11} + (-\Lambda_{22}^{-1}\Lambda_{12}\Sigma_{11})\Lambda_{21} = 1 \quad (166)$$

$$\Rightarrow \Sigma_{11}(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}\Lambda_{21}) = 1 \quad (167)$$

$$\Rightarrow \Sigma_{11} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) = 1 \quad (168)$$

$$\Rightarrow \Sigma_{11} = (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2)^{-1} \quad (169)$$

- これから、 $q_1^*(z_1)$ は次のようにも書ける

$$\begin{aligned} & q_1^*(z_1) \\ = & \frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{(\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2)^{-\frac{1}{2}}} \\ & \exp\left(-\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) (z_1 - \mu_1)^2\right) \end{aligned} \quad (170)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{\Sigma_{11}^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \Sigma_{11}^{-1} (z_1 - \mu_1)^2\right) \quad (171)$$

$$= \mathcal{N}(z_1 | \mu_1, \Sigma_{11}) \quad (172)$$

- $q_2^*(z_2)$ は、対称性から次のようになる

$$q_2^*(z_2) = \mathcal{N}(z_2 | \mu_2, (\Lambda_{22} - \Lambda_{11}^{-1} \Lambda_{12}^2)^{-1}) = \mathcal{N}(z_2 | \mu_2, \Sigma_{22}) \quad (173)$$

- 2つの解の比較

- $\text{KL}(q||p)$ の最小化によって次の解を得た

$$q_1^*(z_1) = \mathcal{N}(z_1 | \mu_1, \Lambda_{11}^{-1}) \quad (174)$$

$$q_2^*(z_2) = \mathcal{N}(z_2 | \mu_2, \Lambda_{22}^{-1}) \quad (175)$$

- $\text{KL}(p||q)$ の最小化では、次の解を得た

$$q_1^*(z_1) = \mathcal{N}(z_1|\mu_1, \Sigma_{11}) \quad (176)$$

$$\Sigma_{11} = (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2)^{-1} \quad (177)$$

$$q_2^*(z_2) = \mathcal{N}(z_2|\mu_2, \Sigma_{22}) \quad (178)$$

$$\Sigma_{22} = (\Lambda_{22} - \Lambda_{11}^{-1} \Lambda_{12}^2)^{-1} \quad (179)$$

- $p(z) = \mathcal{N}(z|\mu, \Sigma)$ の平均は $\mu = [\mu_1, \mu_2]^T$ であったので、いずれの場合も、平均は正しく捉えている
- しかし、両者の間では、分散が異なっている
- また、変数 z_1 と z_2 の間の相関は消えてなくなっている
- 両者の違いを、次の図 3 に示す
- 緑色の線が真の分布 $p(z)$ を表す
- 左側の赤線は、 $\text{KL}(q||p)$ の最小化によって得られた分布 $q(z)$
- 右側の赤線は、 $\text{KL}(p||q)$ の最小化によって得られた分布 $q(z)$

Figure 10.2 Comparison of the two alternative forms for the Kullback-Leibler divergence. The green contours corresponding to 1, 2, and 3 standard deviations for a correlated Gaussian distribution $p(\mathbf{z})$ over two variables z_1 and z_2 , and the red contours represent the corresponding levels for an approximating distribution $q(\mathbf{z})$ over the same variables given by the product of two independent univariate Gaussian distributions whose parameters are obtained by minimization of (a) the Kullback-Leibler divergence $\text{KL}(q\|p)$, and (b) the reverse Kullback-Leibler divergence $\text{KL}(p\|q)$.

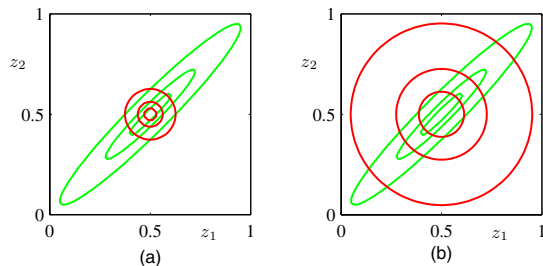


図 3: KL ダイバージェンスの 2 つの形の比較

● 2つの解の比較

- $KL(q||p)$ の最小化で得られる $q(z)$ は、分散が小さくなる方向に制御されていて、それと直交する方向の分散は、大きく過小評価されている
- 分解による近似では、一般に事後分布 $p(Z|X)$ をコンパクトに近似しすぎる
- $KL(p||q)$ の最小化で得られる $q(z)$ は、分散が大きくなる方向に制御されていて、それと直交する方向の分散は、過大評価されている
- 非常に低い確率しか持たないはずの領域にも、多くの確率質量が割り当てられている
- 変分推論では、計算コストの観点から $KL(q||p)$ の方を用いる
- $KL(q||p)$ の計算には、 q に関する期待値の評価が含まれるので、 q を単純な形に制限することで、必要な期待値を単純化できる
- $KL(p||q)$ の計算には、真の事後分布 p に関する期待値の計算が必要である

- 違いが生じる理由

- KL ダイバージェンス $KL(q||p)$ は次のようであった

$$KL(q||p) = - \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} d\mathbf{Z} \quad (180)$$

- $KL(q||p)$ が大きくなる主要因は、 $p(\mathbf{Z})$ がほとんど 0 で、 $q(\mathbf{Z})$ はそうではない領域
- 従って、 $KL(q||p)$ を最小化すると、 $q(\mathbf{Z})$ は、 $p(\mathbf{Z})$ が小さい領域を避けるようになる
- また、 $KL(p||q)$ は次のようであった

$$KL(p||q) = - \int p(\mathbf{Z}|\mathbf{X}) \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})} d\mathbf{Z} \quad (181)$$

変分推論

- $KL(p||q)$ が大きくなる主要因は、 $q(\mathbf{Z})$ がほとんど 0 で、 $p(\mathbf{Z})$ はそうではない領域
- 従って、 $KL(p||q)$ を最小化すると、 $q(\mathbf{Z})$ は、 $p(\mathbf{Z})$ が 0 でない領域にも、必ず確率を持たせるようになる
- 別の分布について、この両者の振舞いの違いを観察してみよう

- 多峰性のある分布を、単峰の分布で近似する場合
 - $KL(q||p)$ を最小化する変分近似では、多数ある峰のうちの 1 つを再現
 - $KL(p||q)$ を最小化する変分近似では、全ての峰を平均したような分布が得られる
- 多峰性のある分布を平均してしまうと、予測性能の悪化をもたらす
- これらの比較を次の図 4 に示す

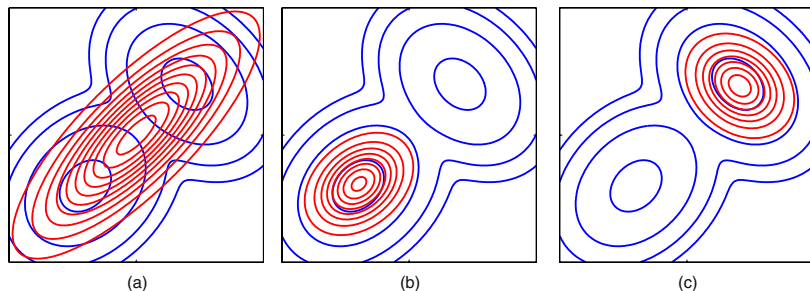


Figure 10.3 Another comparison of the two alternative forms for the Kullback-Leibler divergence. (a) The blue contours show a bimodal distribution $p(\mathbf{Z})$ given by a mixture of two Gaussians, and the red contours correspond to the single Gaussian distribution $q(\mathbf{Z})$ that best approximates $p(\mathbf{Z})$ in the sense of minimizing the Kullback-Leibler divergence $\text{KL}(p\|q)$. (b) As in (a) but now the red contours correspond to a Gaussian distribution $q(\mathbf{Z})$ found by numerical minimization of the Kullback-Leibler divergence $\text{KL}(q\|p)$. (c) As in (b) but showing a different local minimum of the Kullback-Leibler divergence.

図 4: KL ダイバージェンスの 2 つの形の別の比較

- ここまでの話の流れ

- 1 分解 $q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i)$ を使ったエビデンス下界 $\mathcal{L}(q)$ の最適化は、 $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ の最小化と等価である
- 2 $\text{KL}(q||p)$ と、 $\text{KL}(p||q)$ を最小化する変分近似を、2 変数のガウス分布を例として試した
- 3 $\text{KL}(q||p)$ の最小化を使って求めた $q(\mathbf{Z})$ は、事後分布 $p(\mathbf{Z}|\mathbf{X})$ を **コンパクトに近似**する傾向にあった
- 4 $\text{KL}(p||q)$ の最小化によって求めた $q(\mathbf{Z})$ は、事後分布 $p(\mathbf{Z}|\mathbf{X})$ を **大きく捉えて近似**する傾向にあった

- これからの話の流れ

- 変分推論の具体的な例について更に見ていく

変分推論

- 離散潜在変数のモデル (二値スパース符号化モデル) に変分推論を適用する
- 連続潜在変数の場合は、簡単な確率モデルを使って、変分推論を試してみる

1 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 様々な話題

離散潜在変数の変分推論

- 離散潜在変数をもつ変分推論の概要

- ここでは単純な場合を扱う
- データ \mathbf{x}_i に対する潜在変数 z_i は、各要素が二値であるとする
- データ \mathbf{x} は D 次元、潜在変数 z は K 次元とする
- 分布 $q(z_i|\mathbf{x}_i)$ は、平均場近似によって、次のように分解できるとする

$$q(z_i|\mathbf{x}_i) = \prod_k q(z_{ik}|\mathbf{x}_i) \quad (182)$$

- 潜在変数は二値であるから、 $q(z_{ik}=1|\mathbf{x}_i) = \widehat{z}_{ik}$ と書くことにする
- このとき、 $q(z_{ik}|\mathbf{x}_i)$ は次のようになる

$$q(z_{ik}|\mathbf{x}_i) = \widehat{z}_{ik}^{z_{ik}} (1 - \widehat{z}_{ik})^{(1-z_{ik})} \quad (183)$$

- 各パラメータ \widehat{z}_{ik} について、下界 $\mathcal{L}(q)$ を順番に最適化することを、 $\mathcal{L}(q)$ が収束するまで繰り返し行う

離散潜在変数の変分推論

- 即ち、以下の不動点方程式を、各 $\widehat{z_{ik}}$ について繰り返し解く

$$\frac{\partial}{\partial \widehat{z_{ik}}} \mathcal{L}(q) = 0 \quad (184)$$

- 離散潜在変数の場合は、単なる標準的な最適化問題を解くことになる
- 二値スパース符号化モデル
 - 二値スパース符号化モデルでのデータ生成過程は、次のようになる
 - 各データ x_i には、潜在変数 $z_i \in \{0, 1\}^K$ が対応する
 - z_i に対し、重み W を用いて線形変換を施し、更にガウスノイズを足し合わせることで、データ x_i が生成される
 - 潜在変数 z_i は、 K 次元ベクトルであり、その各要素は 0 または 1 である

離散潜在変数の変分推論

- 従って、確率分布は次のようになる

$$p(z_{ik} = 1) = \sigma(b_{ik}) \quad (185)$$

$$p(\mathbf{x}_i | \mathbf{z}_i) = \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \boldsymbol{\beta}^{-1}) \quad (186)$$

- $\sigma(\cdot)$ はシグモイド関数、 $\mathbf{b}_i = [b_{i1}, \dots, b_{iK}]^T$ は学習可能な**バイアス**、 \mathbf{W} は学習可能な**重み行列**、 $\boldsymbol{\beta}$ は学習可能な**対角精度行列**である
- $\boldsymbol{\beta} = \text{diag}(\beta_1, \beta_2, \dots, \beta_D)$ と書ける
- $p(\mathbf{z}_i)$ は次のように計算できる

$$p(\mathbf{z}_i) = \prod_k p(z_{ik}) \quad (187)$$

$$= \prod_k (\sigma(b_{ik}))^{z_{ik}} (1 - \sigma(b_{ik}))^{1-z_{ik}} \quad (188)$$

$$= \prod_k (\sigma(b_{ik}))^{z_{ik}} (\sigma(-b_{ik}))^{1-z_{ik}} \quad (189)$$

離散潜在変数の変分推論

- 事後分布 $p(z_i|x_i)$ は複雑であり、表現することも、計算することもできない
- 従って、最尤推定の手法 (EM アルゴリズム) により学習することはできない
- 例えばバイアス b_{ik} に関する微分を考えると

$$\begin{aligned} & \frac{\partial}{\partial b_{ik}} \ln p(z_i|x_i) \\ = & \frac{\partial}{\partial b_{ik}} \ln \frac{p(x_i, z_i)}{p(x_i)} \end{aligned} \quad (190)$$

$$= \frac{\partial}{\partial b_{ik}} (\ln p(x_i, z_i) - \ln p(x_i)) \quad (191)$$

であって、第 1 項の $p(x_i, z_i)$ は容易に計算できるが、第 2 項について考えると

$$\frac{\partial}{\partial b_{ik}} \ln p(x_i)$$

$$= \frac{1}{p(\mathbf{x}_i)} \frac{\partial}{\partial b_{ik}} p(\mathbf{x}_i) \quad (192)$$

$$= \frac{1}{p(\mathbf{x}_i)} \frac{\partial}{\partial b_{ik}} \sum_{\mathbf{z}_i} p(\mathbf{x}_i, \mathbf{z}_i) \quad (193)$$

$$= \frac{1}{p(\mathbf{x}_i)} \frac{\partial}{\partial b_{ik}} \sum_{\mathbf{z}_i} p(\mathbf{z}_i) p(\mathbf{x}_i | \mathbf{z}_i) \quad (194)$$

$$= \frac{1}{p(\mathbf{x}_i)} \sum_{\mathbf{z}_i} p(\mathbf{x}_i | \mathbf{z}_i) \frac{\partial}{\partial b_{ik}} p(\mathbf{z}_i) \quad (195)$$

$$= \sum_{\mathbf{z}_i} \frac{1}{p(\mathbf{x}_i)} \frac{p(\mathbf{x}_i, \mathbf{z}_i)}{p(\mathbf{z}_i)} \frac{\partial}{\partial b_{ik}} p(\mathbf{z}_i) \quad (196)$$

$$= \sum_{\mathbf{z}_i} p(\mathbf{z}_i | \mathbf{x}_i) \frac{1}{p(\mathbf{z}_i)} \frac{\partial}{\partial b_{ik}} p(\mathbf{z}_i) \quad (197)$$

$$= \sum_{\mathbf{z}_i} p(\mathbf{z}_i | \mathbf{x}_i) \frac{\partial}{\partial b_{ik}} \ln p(\mathbf{z}_i) \quad (198)$$

$$= \mathbb{E}_{\mathbf{z}_i \sim p(\mathbf{z}_i | \mathbf{x}_i)} \left[\frac{\partial}{\partial b_{ik}} \ln p(\mathbf{z}_i) \right] \quad (199)$$

何とか $p(\mathbf{z}_i)$ と $p(\mathbf{x}_i | \mathbf{z}_i)$ を使って計算できないかと試行錯誤したが、結局、 $p(\mathbf{z}_i | \mathbf{x}_i)$ に関する期待値の計算が必要になった

- $p(\mathbf{x}_i, \mathbf{z}_i)$ と、 $p(\mathbf{z}_i | \mathbf{x}_i)$ のグラフ構造を次の図 5 に示す
- 図では、 $\mathbf{h} = \mathbf{z}$ 、 $\mathbf{v} = \mathbf{x}$ のように読み替える必要がある

離散潜在変数の変分推論

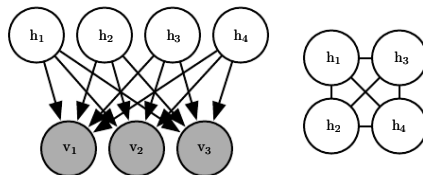


Figure 19.2: The graph structure of a binary sparse coding model with four hidden units. *(Left)* The graph structure of $p(\mathbf{h}, \mathbf{v})$. Note that the edges are directed, and that every two hidden units are co-parents of every visible unit. *(Right)* The graph structure of $p(\mathbf{h} \mid \mathbf{v})$. In order to account for the active paths between co-parents, the posterior distribution needs an edge between all of the hidden units.

図 5: 二値スパース符号化モデルのグラフ構造 (4 つの潜在変数をもつ場合)

- 変分推論による二値スパース符号化モデルの学習

- 最尤推定の手法を諦める代わりに、変分推論を使うことで、困難を解決できる
- 分布 $p(z_i|x_i)$ を $q(z_i)$ で表現し、更に平均場近似を行う
- 但し $z_i = [z_{i1}, z_{i2}, \dots, z_{iK}]^T$ とする

$$q(z_i|x_i) = \prod_k q(z_{ik}|x_i) \quad (200)$$

- 潜在変数の各要素は二値であるから、各 $q(z_{ik}|x_i)$ はベルヌーイ分布とすればよい
- 即ち、 $q(z_{ik} = 1|x_i) = \widehat{z_{ik}}$ とする
- $\widehat{z_{ik}} \neq 0, 1$ という制約を課すことで、 $\ln \widehat{z_{ik}}$ を計算できる
- このようにすれば、潜在変数 $z_{ik}, z_{il} (k \neq l)$ 間の相関を断ち切ることができる
- 先程の図 5 の右側から、潜在変数間の辺を、全て消し去ることになる

離散潜在変数の変分推論

- これで、平均場近似により、因数分解可能な q を表現することができる

$$q(\mathbf{z}_i | \mathbf{x}_i) = \prod_k q(z_{ik} | \mathbf{x}_i) \quad (201)$$

$$= \prod_k \widehat{z}_{ik}^{z_{ik}} (1 - \widehat{z}_{ik})^{1-z_{ik}} \quad (202)$$

$$q(\mathbf{Z} | \mathbf{X}) = \prod_i q(\mathbf{z}_i | \mathbf{x}_i) \quad (203)$$

$$= \prod_i \prod_k \widehat{z}_{ik}^{z_{ik}} (1 - \widehat{z}_{ik})^{1-z_{ik}} \quad (204)$$

- ソフトウェア上では、丸め誤差などによって \widehat{z}_{ik} が 0 や 1 になり、計算を続行できなくなるかもしれない
- これを回避するためには、パラメータ \tilde{z}_i を使って二値スパース符号化モデルを学習させる

離散潜在変数の変分推論

- そして、 $\hat{z}_i = \sigma(\tilde{z}_i)$ の関係によって、 $\hat{z}_i = [\hat{z}_{i1}, \dots, \hat{z}_{iK}]^T$ を得るようにする
- $\ln \hat{z}_{ik} = \ln \sigma(\tilde{z}_{ik}) = -\zeta(-\tilde{z}_{ik})$ によって、コンピュータ上で安全に $\ln \hat{z}_{ik}$ を計算できる ($\zeta(\cdot)$ はソフトプラス関数)
- 変分推論のために、まずはエビデンス下界 $\mathcal{L}(q)$ を計算する

$$\begin{aligned} & \mathcal{L}(q) \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z}|\mathbf{X})} \end{aligned} \quad (205)$$

$$= \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln \frac{p(\mathbf{Z})p(\mathbf{X}|\mathbf{Z})}{q(\mathbf{Z}|\mathbf{X})} \quad (206)$$

$$\begin{aligned} = & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) (\ln p(\mathbf{Z}) + \ln p(\mathbf{X}|\mathbf{Z}) - \ln q(\mathbf{Z}|\mathbf{X})) \quad (207) \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) (\ln p(\mathbf{Z}) - \ln q(\mathbf{Z}|\mathbf{X})) + \end{aligned}$$

$$\sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \quad (208)$$

- $q(\mathbf{Z}|\mathbf{X})$ についての期待値を取っていることに注意する
- 全ての \mathbf{Z} についての和を取ればよいので、第 1 項は次のようになる

$$\begin{aligned} & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) (\ln p(\mathbf{Z}) - \ln q(\mathbf{Z}|\mathbf{X})) \\ = & \sum_i \sum_k \sum_{z_{ik}} q(z_{ik}|\mathbf{x}_i) (\ln p(z_{ik}) - \ln q(z_{ik}|\mathbf{x}_i)) \end{aligned} \quad (209)$$

$$\begin{aligned} = & \sum_i \sum_k q(z_{ik} = 1|\mathbf{x}_i) (\ln p(z_{ik} = 1) - \ln q(z_{ik} = 1|\mathbf{x}_i)) + \\ & q(z_{ik} = 0|\mathbf{x}_i) (\ln p(z_{ik} = 0) - \ln q(z_{ik} = 0|\mathbf{x}_i)) \end{aligned} \quad (210)$$

$$\begin{aligned} = & \sum_i \sum_k \{ \widehat{z_{ik}} (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) + \\ & (1 - \widehat{z_{ik}}) (\ln (1 - \sigma(b_{ik})) - \ln (1 - \widehat{z_{ik}})) \} \end{aligned} \quad (211)$$

$$= \sum_i \sum_k \{ \widehat{z}_{ik} (\ln \sigma(b_{ik}) - \ln \widehat{z}_{ik}) + (1 - \widehat{z}_{ik}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z}_{ik})) \} \quad (212)$$

- また第 2 項は、次のようになる

$$\begin{aligned} & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \\ = & \sum_i \sum_{\mathbf{z}_i} q(\mathbf{z}_i|\mathbf{x}_i) \ln p(\mathbf{x}_i|\mathbf{z}_i) \end{aligned} \quad (213)$$

$$= \sum_i \sum_{\mathbf{z}_i} q(\mathbf{z}_i|\mathbf{x}_i) \ln \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i, \boldsymbol{\beta}^{-1}) \quad (214)$$

- 但し、 $\ln \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \boldsymbol{\beta}^{-1})$ は次のように分解できる

$$\begin{aligned} & \ln \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \boldsymbol{\beta}^{-1}) \\ = & \ln \left(\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\beta}^{-1}|^{\frac{1}{2}}} \right. \\ & \left. \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T \boldsymbol{\beta} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) \right) \right) \end{aligned} \quad (215)$$

$$\begin{aligned} = & -\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\beta}^{-1}| - \\ & \frac{1}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T \boldsymbol{\beta} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) \end{aligned} \quad (216)$$

- ここで、 $\beta = \text{diag}(\beta_1, \beta_2, \dots, \beta_D)$ であるので

$$\begin{aligned}\ln |\beta^{-1}| &= \ln |\beta|^{-1} = -\ln |\beta| \\ &= -\ln \prod_{j=1}^D \beta_j = -\sum_{j=1}^D \ln \beta_j\end{aligned}\quad (217)$$

となるほか

$$(\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T \beta (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) = \sum_{j=1}^D \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \quad (218)$$

であるから

$$\begin{aligned}&\ln \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \beta^{-1}) \\ &= -\frac{1}{2} \sum_j \ln 2\pi + \frac{1}{2} \sum_j \ln \beta_j - \frac{1}{2} \sum_j \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2\end{aligned}\quad (219)$$

$$= \frac{1}{2} \sum_j \left(-\ln 2\pi + \ln \beta_j - \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \right) \quad (220)$$

$$= \frac{1}{2} \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \right) \quad (221)$$

$$= \frac{1}{2} \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right) \quad (222)$$

- 上式の対数を外すと、確率分布の積になっていることが分かる

$$\begin{aligned} & \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \boldsymbol{\beta}^{-1}) \\ &= \exp \left(\sum_j \left(\frac{1}{2} \ln \frac{\beta_j}{2\pi} - \frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \right) \right) \end{aligned} \quad (223)$$

$$= \prod_j \exp \left(\frac{1}{2} \ln \frac{\beta_j}{2\pi} - \frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \right) \quad (224)$$

$$= \prod_j \exp\left(\frac{1}{2} \ln \frac{\beta_j}{2\pi}\right) \exp\left(-\frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2\right) \quad (225)$$

$$= \prod_j \left(\frac{\beta_j}{2\pi}\right)^{\frac{1}{2}} \exp\left(-\frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2\right) \quad (226)$$

$$= \prod_j \left(\frac{1}{(2\pi\beta_j^{-1})^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2\right)\right) \quad (227)$$

$$= \prod_j \mathcal{N}(x_{ij} | \mathbf{W}_{j:} \mathbf{z}_i, \beta_j^{-1}) \quad (228)$$

- 精度行列 β は対角行列であるから、 x_i の各成分は無相関である
- 従って、 $\mathcal{N}(x_i | \mathbf{W} \mathbf{z}_i, \beta^{-1})$ は、各成分 x_{ij} についての確率 $\mathcal{N}(x_{ij} | \mathbf{W}_{j:} \mathbf{z}_i, \beta_j^{-1})$ の積として、記述できる

- さて、第2項は

$$\begin{aligned} & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \\ &= \sum_i \sum_{\mathbf{z}_i} q(\mathbf{z}_i|\mathbf{x}_i) \ln \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i, \beta^{-1}) \end{aligned} \quad (229)$$

$$= \sum_i \mathbb{E}_{\mathbf{z}_i \sim q(\mathbf{z}_i|\mathbf{x}_i)} [\ln \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i, \beta^{-1})] \quad (230)$$

$$= \sum_i \mathbb{E}_{\mathbf{z}_i} \left[\frac{1}{2} \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right) \right] \quad (231)$$

$$= \frac{1}{2} \sum_i \sum_j \mathbb{E}_{\mathbf{z}_i} \left[\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \quad (232)$$

$$= \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \mathbb{E}_{\mathbf{z}_i} \left[\left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \right) \quad (233)$$

ここで

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} \left[\left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \\ = & \mathbb{E}_{\mathbf{z}_i} \left[x_{ij}^2 - 2x_{ij} \sum_k W_{jk} z_{ik} + \left(\sum_k W_{jk} z_{ik} \right)^2 \right] \end{aligned} \quad (234)$$

$$= x_{ij}^2 - 2x_{ij} \mathbb{E}_{\mathbf{z}_i} \left[\sum_k W_{jk} z_{ik} \right] + \mathbb{E}_{\mathbf{z}_i} \left[\left(\sum_k W_{jk} z_{ik} \right)^2 \right] \quad (235)$$

各項を順番に計算すると

$$\mathbb{E}_{\mathbf{z}_i} \left[\sum_k W_{jk} z_{ik} \right]$$

$$= \sum_k W_{jk} \mathbb{E}_{\mathbf{z}_i} [z_{ik}] \quad (236)$$

$$= \sum_k W_{jk} \mathbb{E}_{z_{ik}} [z_{ik}] \quad (237)$$

$$= \sum_k W_{jk} (q(z_{ik} = 1 | \mathbf{x}_i) \cdot 1 + q(z_{ik} = 0 | \mathbf{x}_i) \cdot 0) \quad (238)$$

$$= \sum_k W_{jk} \widehat{z_{ik}} \quad (239)$$

また

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} \left[\left(\sum_k W_{jk} z_{ik} \right)^2 \right] \\ = & \mathbb{E}_{\mathbf{z}_i} \left[\sum_k \sum_l W_{jk} W_{jl} z_{ik} z_{il} \right] \end{aligned} \quad (240)$$

$$= \mathbb{E}_{\mathbf{z}_i} \left[\sum_k \left(W_{jk}^2 z_{ik}^2 + \sum_{l \neq k} (W_{jk} W_{jl} z_{ik} z_{il}) \right) \right] \quad (241)$$

$$= \sum_k \left(W_{jk}^2 \mathbb{E}_{\mathbf{z}_i} [z_{ik}^2] + \sum_{l \neq k} W_{jk} W_{jl} \mathbb{E}_{\mathbf{z}_i} [z_{ik} z_{il}] \right) \quad (242)$$

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} [z_{ik}^2] \\ &= \mathbb{E}_{z_{ik}} [z_{ik}^2] \end{aligned} \quad (243)$$

$$= q(z_{ik} = 1 | \mathbf{x}_i) \cdot 1^2 + q(z_{ik} = 0 | \mathbf{x}_i) \cdot 0^2 \quad (244)$$

$$= \widehat{z_{ik}} \quad (245)$$

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} [z_{ik} z_{il}] \\ &= \mathbb{E}_{\mathbf{z}_i} [z_{ik}] \mathbb{E}_{\mathbf{z}_i} [z_{il}] \quad (\because z_{ik} \text{ と } z_{il} \text{ は独立であるため}) \end{aligned} \quad (246)$$

$$= \mathbb{E}_{z_{ik}} [z_{ik}] \mathbb{E}_{z_{il}} [z_{il}] \quad (247)$$

$$= \widehat{z_{ik}} \widehat{z_{il}} \quad (248)$$

従って

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} \left[\left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \\ &= x_{ij}^2 - 2x_{ij} \mathbb{E}_{\mathbf{z}_i} \left[\sum_k W_{jk} z_{ik} \right] + \mathbb{E}_{\mathbf{z}_i} \left[\left(\sum_k W_{jk} z_{ik} \right)^2 \right] \end{aligned} \quad (249)$$

$$\begin{aligned} &= x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} + \\ & \quad \sum_k \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) \end{aligned} \quad (250)$$

これより

$$\begin{aligned} & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \\ &= \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \mathbb{E}_{\mathbf{z}_i} \left[\left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \right) \end{aligned} \quad (251)$$

$$\begin{aligned} &= \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z}_{ik} + \right. \right. \\ &\quad \left. \left. \sum_k \left(W_{jk}^2 \widehat{z}_{ik} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z}_{ik} \widehat{z}_{il} \right) \right) \right) \end{aligned} \quad (252)$$

- よって、エビデンス下界 $\mathcal{L}(q)$ は

$$\begin{aligned} & \mathcal{L}(q) \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) (\ln p(\mathbf{Z}) - \ln q(\mathbf{Z}|\mathbf{X})) + \\ & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \end{aligned} \quad (253)$$

$$\begin{aligned} = & \sum_i \sum_k \{ \widehat{z}_{ik} (\ln \sigma(b_{ik}) - \ln \widehat{z}_{ik}) + \\ & (1 - \widehat{z}_{ik}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z}_{ik})) \} + \\ & \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z}_{ik} + \right. \right. \\ & \left. \left. \sum_k \left(W_{jk}^2 \widehat{z}_{ik} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z}_{ik} \widehat{z}_{il} \right) \right) \right) \end{aligned} \quad (254)$$

離散潜在変数の変分推論

- この式にはあまり美的魅力がない (Somewhat unappealing aesthetically) が、事後分布よりは計算しやすい
- 事後分布 $\ln p(\mathbf{Z}|\mathbf{X})$ を最大化する代わりに、上記の下界 $\mathcal{L}(q)$ を q について最大化することができる
- データ x_i に対するパラメータ $\{\widehat{z}_{i1}, \dots, \widehat{z}_{iK}\}$ をまとめて、ベクトル $\widehat{\mathbf{z}}_i$ として記述する
- パラメータ $\widehat{\mathbf{z}}_i = [\widehat{z}_{i1}, \dots, \widehat{z}_{iK}]^T$ は、データ x_i に対応する潜在変数 z の各要素が、1 となる確率を集めたベクトルである
- $\widehat{z}_{ik} = q(z_{ik} = 1|x_i)$ と定義されていることに注意
- よってベクトル $\widehat{\mathbf{z}}_i$ は、データ x_i に対応する二値のスパース符号である
- 勾配上昇法を利用しない理由
 - データ \mathbf{X} と、潜在変数 \mathbf{Z} についての勾配上昇法を用いれば、学習することが可能

離散潜在変数の変分推論

- 但し、その方法では、各 x_i について平均場パラメータ \hat{z}_i を保管する必要がある
 - 各事例について、動的に更新されるベクトルが必要であるため、そのアルゴリズムを、数十億もの事例に対して適用することは困難である
 - また、収束するまで繰り返し計算を行うため、データ x から、パラメータ \hat{z}_i を素早く抽出することができない
 - 現実にデプロイされるときは、 \hat{z}_i をリアルタイムで計算できなければならない
-
- 不動点方程式による平均場パラメータ \hat{z}_i の推定
 - 勾配上昇法の代わりに、不動点方程式を使ってパラメータ \hat{z}_{ik} を素早く推定できる
 - $\nabla_{\hat{z}_i} \mathcal{L}(q) = 0$ をみtas、 \hat{z}_i の極大値を見つけ出す

離散潜在変数の変分推論

- \widehat{z}_i の全ての成分について同時に解くことはできないので、各成分 \widehat{z}_{ik} について繰り返し解く
- 即ち、各パラメータ \widehat{z}_{ik} について、下界 $\mathcal{L}(q)$ を順番に最適化する手続きを、 $\mathcal{L}(q)$ の収束基準を満たすまで繰り返す

$$\frac{\partial}{\partial \widehat{z}_{ik}} \mathcal{L}(q) = 0 \quad (255)$$

- 平均場不動点方程式を導くためには、 $\mathcal{L}(q)$ を \widehat{z}_{ik} で微分する必要がある

$$\begin{aligned} & \frac{\partial}{\partial \widehat{z}_{ik}} \mathcal{L}(q) \\ = & \frac{\partial}{\partial \widehat{z}_{ik}} \sum_i \sum_k \{ \widehat{z}_{ik} (\ln \sigma(b_{ik}) - \ln \widehat{z}_{ik}) + \\ & (1 - \widehat{z}_{ik}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z}_{ik})) \} + \end{aligned}$$

$$\frac{\partial}{\partial \widehat{z_{ik}}} \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} + \sum_k \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) \right) \right) \quad (256)$$

- 前半部分は

$$\frac{\partial}{\partial \widehat{z_{ik}}} \sum_i \sum_k \{ \widehat{z_{ik}} (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) + (1 - \widehat{z_{ik}}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z_{ik}})) \} \quad (257)$$

$$= \frac{\partial}{\partial \widehat{z_{ik}}} \{ \widehat{z_{ik}} (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) + (1 - \widehat{z_{ik}}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z_{ik}})) \} \quad (258)$$

$$= \frac{\partial}{\partial \widehat{z_{ik}}} \widehat{z_{ik}} (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) +$$

$$\frac{\partial}{\partial \widehat{z_{ik}}} (1 - \widehat{z_{ik}}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z_{ik}})) \quad (259)$$

$$\begin{aligned} = & (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) + \widehat{z_{ik}} \left(-\frac{1}{\widehat{z_{ik}}} \right) + \\ & (- (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z_{ik}}))) + \\ & (1 - \widehat{z_{ik}}) \frac{1}{1 - \widehat{z_{ik}}} \end{aligned} \quad (260)$$

$$= \ln \sigma(b_{ik}) - \ln \widehat{z_{ik}} - 1 - \ln \sigma(-b_{ik}) + \ln(1 - \widehat{z_{ik}}) + 1 \quad (261)$$

$$= \ln \sigma(b_{ik}) - \ln \widehat{z_{ik}} - \ln \sigma(-b_{ik}) + \ln(1 - \widehat{z_{ik}}) \quad (262)$$

$$= -\zeta(-b_{ik}) - \ln \widehat{z_{ik}} + \zeta(b_{ik}) + \ln(1 - \widehat{z_{ik}}) \quad (263)$$

$$= (\zeta(b_{ik}) - \zeta(-b_{ik})) - \ln \widehat{z_{ik}} + \ln(1 - \widehat{z_{ik}}) \quad (264)$$

$$= b_{ik} - \ln \widehat{z_{ik}} + \ln(1 - \widehat{z_{ik}}) \quad (265)$$

離散潜在変数の変分推論

- ここで、シグモイド関数 $\sigma(\cdot)$ と、ソフトプラス関数 $\zeta(\cdot)$ に関する、以下の公式を用いた

$$\ln \sigma(x) = -\zeta(-x) \quad (266)$$

$$\zeta(x) - \zeta(-x) = x \quad (267)$$

- 後半部分は

$$\begin{aligned} & \frac{\partial}{\partial \widehat{z_{ik}}} \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} + \right. \right. \\ & \quad \left. \left. \sum_k \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) \right) \right) \quad (268) \\ &= \frac{1}{2} \sum_j \beta_j \frac{\partial}{\partial \widehat{z_{ik}}} \left(2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} - \right. \end{aligned}$$

$$\sum_k \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) \quad (269)$$

$$= \frac{1}{2} \sum_j \beta_j \frac{\partial}{\partial \widehat{z_{ik}}} \left(2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} - \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) - \sum_{m \neq k} \left(W_{jm}^2 \widehat{z_{im}} + \sum_{l \neq m} W_{jm} W_{jl} \widehat{z_{im}} \widehat{z_{il}} \right) \right) \quad (270)$$

$$= \frac{1}{2} \sum_j \beta_j \left(2x_{ij} W_{jk} - W_{jk}^2 - \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{il}} - \sum_{m \neq k} W_{jm} W_{jk} \widehat{z_{im}} \right) \quad (271)$$

$$= \frac{1}{2} \sum_j \beta_j \left(2x_{ij}W_{jk} - W_{jk}^2 - 2 \sum_{l \neq k} W_{jk}W_{jl}\widehat{z_{il}} \right) \quad (272)$$

$$= \sum_j \beta_j \left(x_{ij}W_{jk} - \frac{1}{2}W_{jk}^2 - \sum_{l \neq k} W_{jk}W_{jl}\widehat{z_{il}} \right) \quad (273)$$

$$= \sum_j x_{ij}\beta_j W_{jk} - \frac{1}{2} \sum_j W_{jk}\beta_j W_{jk} - \sum_j \beta_j \sum_{l \neq k} W_{jk}W_{jl}\widehat{z_{il}} \quad (274)$$

$$= \sum_j x_{ij}\beta_j W_{jk} - \frac{1}{2} \sum_j W_{jk}\beta_j W_{jk} - \sum_{l \neq k} \left(\sum_j W_{jl}\beta_j W_{jk} \right) \widehat{z_{il}} \quad (275)$$

$$= \mathbf{x}_i^T \beta \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \beta \mathbf{W}_{:k} \widehat{z}_{il} \quad (276)$$

- これより、 $\mathcal{L}(q)$ の \widehat{z}_{ik} による微分は次のようになる

$$\begin{aligned} & \frac{\partial}{\partial \widehat{z}_{ik}} \mathcal{L}(q) \\ = & b_{ik} - \ln \widehat{z}_{ik} + \ln(1 - \widehat{z}_{ik}) + \\ & \mathbf{x}_i^T \beta \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \beta \mathbf{W}_{:k} \widehat{z}_{il} \quad (277) \end{aligned}$$

離散潜在変数の変分推論

- これを 0 と等置して、 \widehat{z}_{ik} について解くと次のようになる

$$\ln \widehat{z}_{ik} - \ln(1 - \widehat{z}_{ik}) = b_{ik} + \mathbf{x}_i^T \boldsymbol{\beta} \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \boldsymbol{\beta} \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \boldsymbol{\beta} \mathbf{W}_{:k} \widehat{z}_{il} \quad (278)$$

右辺を A とおけば

$$\begin{aligned} \ln \widehat{z}_{ik} - \ln(1 - \widehat{z}_{ik}) &= A & (279) \\ \Rightarrow \ln \frac{\widehat{z}_{ik}}{1 - \widehat{z}_{ik}} &= A \\ \Rightarrow \frac{\widehat{z}_{ik}}{1 - \widehat{z}_{ik}} &= \exp A \\ \Rightarrow \frac{1}{1 - \widehat{z}_{ik}} - 1 &= \exp A \\ \Rightarrow 1 - \widehat{z}_{ik} &= \frac{1}{1 + \exp A} \end{aligned}$$

$$\begin{aligned}\Rightarrow \widehat{z_{ik}} &= 1 - \frac{1}{1 + \exp A} \\ \Rightarrow \widehat{z_{ik}} &= \frac{\exp A}{1 + \exp A} \\ \Rightarrow \widehat{z_{ik}} &= \frac{1}{1 + \exp(-A)}\end{aligned}\tag{280}$$

$$\Rightarrow \widehat{z_{ik}} = \sigma(A)\tag{281}$$

従って、不動点方程式は

$$\widehat{z_{ik}} = \sigma \left(b_{ik} + \mathbf{x}_i^T \boldsymbol{\beta} \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \boldsymbol{\beta} \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \boldsymbol{\beta} \mathbf{W}_{:k} \widehat{z_{il}} \right)\tag{282}$$

- 不動点方程式の観察

離散潜在変数の変分推論

- 不動点方程式は次で表された

$$\widehat{z}_{ik} = \sigma \left(b_{ik} + \mathbf{x}_i^T \beta \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \beta \mathbf{W}_{:k} \widehat{z}_{il} \right) \quad (283)$$

- 第2項 $\mathbf{x}_i^T \beta \mathbf{W}_{:k}$ は、潜在変数のユニット k に対する入力
- 第3項 $-\frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k}$ は、隠れユニット k から自身への入力
- 第4項 $-\sum_{l \neq k} \mathbf{W}_{:l}^T \beta \mathbf{W}_{:k} \widehat{z}_{il}$ は、他の隠れユニット $l \neq k$ から隠れユニット k への入力
- これより、平均場不動点方程式は、回帰結合型ニューラルネットワーク (RNN) との関係があることが分かる
- 隠れユニット k と l は、それらの重みベクトル $\mathbf{W}_{:l}$ と $\mathbf{W}_{:k}$ が互いに同調するとき (似たような重みを持っているとき) に、互いに抑制し合う

離散潜在変数の変分推論

- 即ち、2つの隠れユニット k, l が、共に入力を説明するとき (入力から同じような表現を抽出するとき)、**入力を最もよく説明するユニットのみがアクティブ**になる (強く活性化される)
- これはユニット間の競合の一形態である
- 従って、実際には多峰性の事後分布かもしれないが、そのうちの1つのみを選択される (図4の(b)と(c)参照)
- 不動点方程式を更に以下のように変形する

$$\widehat{z}_{ik} = \sigma \left(b_{ik} + \left(x_i - \sum_{l \neq k} \mathbf{W}_{:l} \widehat{z}_{il} \right)^T \beta \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k} \right) \quad (284)$$

- これより、ユニット k への入力は、 x_i ではなく $x_i - \sum_{l \neq k} \mathbf{W}_{:l} \widehat{z}_{il}$ であるとみなせる

離散潜在変数の変分推論

- ユニット k への入力は、他の全てのユニットによる x_i の再構成と、実際の入力 x_i との誤差である
 - ユニット k は、この残差誤差を符号化していると分かるので、スパース符号化は、**反復自己符号化器**とみなせる
 - スパース符号化では、入力 x_i の符号化 ($\widehat{z_{ik}}$ の計算) と、復号 ($\sum_{l \neq k} \mathbf{W}_{:l} \widehat{z_{il}}$ の計算) を繰り返す
 - この反復のたびに、再構成の誤差を修正していく
-
- **ダンピング**
 - 1つのユニットの更新則を (不動点方程式として) 導出した
 - 複数のユニットを同時に更新することは、二値スパース符号化モデルでは通常できない
 - 但し、**ダンピング**という発見的手法を使えば可能になる
 - 各要素 $\widehat{z_{ik}}$ についての最適値を計算し、その値の変化の方向に、他の要素 $\widehat{z_{il}}$ を小さいステップで動かす

離散潜在変数の変分推論

- 下界 $\mathcal{L}(q)$ が増加することはもはや保証されないが、多くの場合はうまくいく

離散潜在変数の変分推論のまとめ

● ここまでの話の流れ

- 1 離散潜在変数における変分推論の具体例として、二値スパース符号化モデルをみた
- 2 事後分布 $p(Z|X)$ が複雑になるので、最尤推定 (EM アルゴリズム) が使えない
- 3 代わりに、別の分布 $q(Z)$ を使って事後分布を近似することにした
- 4 エビデンス下界 $\mathcal{L}(q)$ を苦勞して求めた
- 5 更に、下界 $\mathcal{L}(q)$ を (各パラメータについて) 最大化するための、不動点方程式を導出した
- 6 不動点方程式を観察し、回帰結合型ニューラルネットワークや、自己符号化器との関係を考えて

● これからの話の流れ

離散潜在変数の変分推論のまとめ

- 連続潜在変数に対する変分推論を、簡単な確率モデルを使って、試してみよう

1 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 様々な話題

- 連続潜在変数をもつ変分推論の概要

- 平均場近似を行う場合は、以下の式によって最適な因子 $q_j^*(\mathbf{Z}_j|\mathbf{X})$ が得られる

$$\ln q_j^*(\mathbf{Z}_j|\mathbf{X}) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (285)$$

$$\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i(\mathbf{Z}_i|\mathbf{X}) d\mathbf{Z}_i \quad (286)$$

$$q_j^*(\mathbf{Z}_j|\mathbf{X}) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j} \quad (287)$$

- 上式は、下界 $\mathcal{L}(q)$ を最大化する q であり、従って連続潜在変数の場合の不動点方程式とみなせる
- 各因子 $q_j(\mathbf{Z}_j|\mathbf{X})$ を、上式を用いて順番に更新していく ($i \neq j$ である全ての q_i を固定した状態で、各 q_j について下界を最適化する)

連続潜在変数の変分推論

- このステップを、下界 $\mathcal{L}(q)$ が収束するまで繰り返し行う (座標降下法)
- 不動点方程式は、下界が最適な値に収束するかどうかには関係なく、 q_j の最適解が取る関数形を提供してくれる

- 扱う確率モデルの表現

- ここでは次のような確率モデルを対象として、変分推論を扱う

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i | 0, \mathbf{I}) \quad (288)$$

$$p(x_i | \mathbf{z}_i) = \mathcal{N}(x_i | \mathbf{w}^T \mathbf{z}_i, 1) \quad (289)$$

- 1次元のデータ $x \in \mathbb{R}$ に対して、2次元の潜在変数 $\mathbf{z} \in \mathbb{R}^2$ が存在する
- 同時分布 $p(x_i, \mathbf{z}_i) = p(x_i | \mathbf{z}_i)p(\mathbf{z}_i)$ を \mathbf{z}_i で積分消去すれば、 x_i についての単なるガウス分布となる

- 真の事後分布 $p(\mathbf{z}_i | x_i)$ の計算

連続潜在変数の変分推論

- 正規化定数を見捨てて次のように計算できる

$$\begin{aligned} & p(\mathbf{z}_i | x_i) \\ \propto & p(\mathbf{z}_i | x_i) p(x_i) \\ = & p(x_i, \mathbf{z}_i) \\ = & p(\mathbf{z}_i) p(x_i | \mathbf{z}_i) \\ = & \mathcal{N}(\mathbf{z}_i | 0, \mathbf{I}) \mathcal{N}(x_i | \mathbf{w}^T \mathbf{z}_i, 1) \\ \propto & \exp\left(-\frac{1}{2} \mathbf{z}_i^T \mathbf{z}_i\right) \exp\left(-\frac{1}{2} (x_i - \mathbf{w}^T \mathbf{z}_i)^T (x_i - \mathbf{w}^T \mathbf{z}_i)\right) \\ = & \exp\left(-\frac{1}{2} (z_{i1}^2 + z_{i2}^2)\right) \\ & \exp\left(-\frac{1}{2} (x_i - w_1 z_{i1} - w_2 z_{i2})^T (x_i - w_1 z_{i1} - w_2 z_{i2})\right) \\ = & \exp\left(-\frac{1}{2} (z_{i1}^2 + z_{i2}^2 + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 z_{i2}^2 - \right. \end{aligned}$$

$$2x_i w_1 z_{i1} - 2x_i w_2 z_{i2} + 2w_1 w_2 z_{i1} z_{i2} \Big) \quad (290)$$

- 正規化項を C とすれば次のように書ける

$$\begin{aligned} p(z_i | x_i) \\ = C \exp \left(-\frac{1}{2} (z_{i1}^2 + z_{i2}^2 + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 z_{i2}^2 - \right. \\ \left. 2x_i w_1 z_{i1} - 2x_i w_2 z_{i2} + 2w_1 w_2 z_{i1} z_{i2}) \right) \end{aligned} \quad (291)$$

- z_{i1} と z_{i2} を乗算する項が存在する
- 従って、真の事後分布は、 z_{i1} と z_{i2} のみの因子には分解できないことが分かる

- 平均場近似の計算

- 平均場近似を次のように表現する

$$q(\mathbf{z}_i | x_i) = q_1(z_{i1} | x_i) q_2(z_{i2} | x_i) \quad (292)$$

- q_2 を固定した状態で、最適な $q_1^*(z_{i1} | x_i)$ を求める

$$\begin{aligned} & \ln q_1^*(z_{i1} | x_i) \\ = & \mathbb{E}_{z_{i2} \sim q_2(z_{i2} | x_i)} [\ln p(\mathbf{z}_i, x_i)] + \text{Const.} \end{aligned} \quad (293)$$

$$\begin{aligned} = & \mathbb{E}_{z_{i2}} \left[\ln \left(C \exp \left(-\frac{1}{2} (z_{i1}^2 + z_{i2}^2 + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 z_{i2}^2 - \right. \right. \right. \\ & \left. \left. \left. 2x_i w_1 z_{i1} - 2x_i w_2 z_{i2} + 2w_1 w_2 z_{i1} z_{i2} \right) \right) \right] + \text{Const.} \\ = & \mathbb{E}_{z_{i2}} \left[-\frac{1}{2} (z_{i1}^2 + z_{i2}^2 + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 z_{i2}^2 - \right. \\ & \left. 2x_i w_1 z_{i1} - 2x_i w_2 z_{i2} + 2w_1 w_2 z_{i1} z_{i2}) \right] + \text{Const.} \end{aligned} \quad (294)$$

連続潜在変数の変分推論

- $z_{i2} \sim q_2(z_{i2}|x_i)$ による期待値を取っている
- $q_2(z_{i2}|x_i)$ から得る必要があるのは、結局 $\mathbb{E}_{z_{i2}}[z_{i2}]$ と、 $\mathbb{E}_{z_{i2}}[z_{i2}^2]$ の2つだけである
- $\langle z_{i2} \rangle = \mathbb{E}_{z_{i2}}[z_{i2}]$ 、 $\langle z_{i2}^2 \rangle = \mathbb{E}_{z_{i2}}[z_{i2}^2]$ と書くことにする
- このとき次式が得られる

$$\begin{aligned} & \ln q_1^*(z_{i1}|x_i) \\ = & -\frac{1}{2} \left(z_{i1}^2 + \langle z_{i2}^2 \rangle + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 \langle z_{i2}^2 \rangle - \right. \\ & \quad \left. 2x_i w_1 z_{i1} - 2x_i w_2 \langle z_{i2} \rangle + 2w_1 w_2 z_{i1} \langle z_{i2} \rangle \right) + \\ & \quad \text{Const.} \end{aligned} \tag{295}$$

連続潜在変数の変分推論

- これより、最適な $q_1^*(z_{i1}|x_i)$ は**ガウス分布**の形であると分かる

$$\begin{aligned} & q_1^*(z_{i1}|x_i) \\ = & C \exp \left(-\frac{1}{2} (z_{i1}^2 + \langle z_{i2}^2 \rangle + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 \langle z_{i2}^2 \rangle - \right. \\ & \quad \left. 2x_i w_1 z_{i1} - 2x_i w_2 \langle z_{i2} \rangle + 2w_1 w_2 z_{i1} \langle z_{i2} \rangle) \right) \\ \propto & \exp \left(-\frac{1}{2} (z_{i1}^2 + \langle z_{i2}^2 \rangle + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 \langle z_{i2}^2 \rangle - \right. \\ & \quad \left. 2x_i w_1 z_{i1} - 2x_i w_2 \langle z_{i2} \rangle + 2w_1 w_2 z_{i1} \langle z_{i2} \rangle) \right) \quad (296) \end{aligned}$$

- 対称性から、最適な $q_2^*(z_{i2}|x_i)$ も**ガウス分布**であることが分かる
- ガウス分布同士の積もガウス分布になるので、結局 $q(z_i|x_i) = q_1(z_{i1}|x_i)q_2(z_{i2}|x_i)$ はガウス分布である

連続潜在変数の変分推論

- $q(z_i|x_i)$ が 2 つの因子に分解できるとは仮定したが、**各因子の関数形については全く仮定していない**ことに注意
- ガウス分布は、下界 $\mathcal{L}(q)$ を q について変分最適化する過程で、**自然に出現した**

連続潜在変数の変分推論のまとめ

- ここまでの話の流れ

- 1 連続潜在変数における変分推論の例として、簡単な確率モデルを扱った
- 2 分布 $q(\mathbf{Z}|\mathbf{X})$ を、 $\prod_i q_i(\mathbf{Z}_i|\mathbf{X})$ のように因数分解できるという仮定 (平均場近似) のみを置いた
- 3 各因子 $q_i(\mathbf{Z}_i|\mathbf{X})$ の関数形については全く仮定を置かなかった
- 4 変分推論によって、最適解となる q_i の関数形が自然に導出できた

- これからの話の流れ

- 教科書で触れられている雑多な話題を扱う
- 近似推論が、推論アルゴリズムの精度に影響することについて考える
- 学習による近似推論の手法を考える

1 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 様々な話題

- 近似推論が精度に与える影響

- 近似推論では、以下の下界 $\mathcal{L}(q)$ を q について最適化する

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p) \quad (297)$$

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (298)$$

- 下界 $\mathcal{L}(q)$ を次のように分解する

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} - \int q(\mathbf{Z}) \ln q(\mathbf{Z}) d\mathbf{Z} \quad (299)$$

$$= \int q(\mathbf{Z}) (\ln p(\mathbf{Z}|\mathbf{X}) + \ln p(\mathbf{X})) d\mathbf{Z} + H[q] \quad (300)$$

$$= \int q(\mathbf{Z}) \ln p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z} + \ln p(\mathbf{X}) + H[q] \quad (301)$$

$$= \mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} [\ln p(\mathbf{Z}|\mathbf{X})] + \ln p(\mathbf{X}) + H[q] \quad (302)$$

学習と推論の相互作用

- 分布に含まれるパラメータを明示的に記述する
- $p(\mathbf{Z}|\mathbf{X}, \theta)$ と、 $p(\mathbf{Z}|\theta)$ をモデリングする
- このとき、 $p(\mathbf{X}, \mathbf{Z}|\theta) = p(\mathbf{Z}|\mathbf{X}, \theta)p(\mathbf{X}|\theta)$ と書ける
- 従って、下界 $\mathcal{L}(q)$ は次のようになる

$$\mathcal{L}(q, \theta) = \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{Z}|\mathbf{X}, \theta)] + \ln p(\mathbf{X}|\theta) + H[q] \quad (303)$$

- 下界 $\mathcal{L}(q, \theta)$ をパラメータ θ について増加させるとする
- これは、 $\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{Z}|\mathbf{X}, \theta)]$ を、 q を固定しつつ、 θ について増加させることに繋がる
- このとき、 $q(\mathbf{Z})$ が高い確率となる \mathbf{Z} について $p(\mathbf{Z}|\mathbf{X}, \theta)$ が増大し、 $q(\mathbf{Z})$ が低い確率となる \mathbf{Z} について $p(\mathbf{Z}|\mathbf{X}, \theta)$ が減少する
- これより、近似仮定が自己充足的予言となることが分かる
- 自己充足的予言とは、ある事象や状況に関する判断や思い込みが原因となり、その結果として、その判断や思い込みが現実化することである

学習と推論の相互作用

- 即ち、予め近似事後分布 $q(\mathbf{Z})$ に何らかの仮定を置いて訓練することで、**結果としてその仮定に沿うような分布が得られてしまう**
- 単峰性の近似事後分布 $q(\mathbf{Z})$ を使って訓練することを考える
- 得られる事後分布 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ は、厳密な推論によって得られた $p(\mathbf{Z}|\mathbf{X})$ と比べて、単峰性がはるかに強くなっている
- 事後分布が多峰性であるという仮定が、**推定結果にも現れてしまう**
- 変分推論によってモデルに課される損害
 - モデルを訓練した後に、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を推定し、下界 $\mathcal{L}(q, \boldsymbol{\theta})$ との**差が十分に小さい** ($\mathcal{L}(q, \boldsymbol{\theta}) \simeq \ln p(\mathbf{X}|\boldsymbol{\theta})$) ことを確認する
 - このとき、**ある特定のパラメータ $\boldsymbol{\theta}$ について**、変分近似が正確であるといえる
 - 但し、**一般的に変分近似が正確であるとはいえない**

学習と推論の相互作用

- また、変分推論が、学習過程に殆ど害を及ぼさないとも結論付けては**いけない**
- 変分近似による**真の損害**は、 $\ln p(\mathbf{X}|\theta)$ を最大にする θ^* が分からなければ、測定することができない
- ある θ について $\mathcal{L}(q, \theta) \simeq \ln p(\mathbf{X}|\theta)$ が成立しても、 $\mathcal{L}(q, \theta) \simeq \ln p(\mathbf{X}|\theta) \ll \ln p(\mathbf{X}|\theta^*)$ であるかもしれない
- あるパラメータ θ については変分近似は正確だが、実際には近似が全く上手く行っていないことになる
- また $\max_q \mathcal{L}(q, \theta^*) \ll \ln p(\mathbf{X}|\theta^*)$ であれば、 θ^* のときの事後分布は、制限された q にとって複雑すぎるため、訓練では θ^* を見つけることができない
- この問題は、変分推論とは別の方法で θ^* が求まる場合にしか、検知することができない

- 近似推論の実行の学習

- 不動点方程式や勾配に基づく最適化を、繰り返し行う近似推論は、計算コストが非常に高い
- 近似推論による最適化処理は、入力 X から、近似事後分布 $q^*(Z) = \arg \max_q \mathcal{L}(X, q)$ へと写像する関数 f とみなせる
- これにより、反復的な最適化処理を単なる関数とみなして、関数 $\hat{f}(X, \theta)$ を、ニューラルネットワークで近似することができる

- Wake-Sleep アルゴリズム

- データ X から潜在変数 Z を推論するときの困難は、正しい潜在変数 Z が分からず、従って教師あり訓練集合が得られないことである
- データ X から潜在変数 Z への写像は、モデル $(p(X|Z|\theta), p(Z|\theta))$ の選択に依存し、更に θ の変化に伴って変わり続ける

学習による近似推論

- Wake-Sleep アルゴリズムでは、 X と Z の両方のサンプルを、モデル分布から抽出する
- モデルの分布が高い確率を示すような X に対してしか、推論ネットワークを学習できない
- モデルの分布が、データの分布に似ていないとき、推論ネットワークはデータに似たサンプルを学習できない
- 推論学習の他の形式
 - 変分自己符号化器は、生成モデリングで主要なアプローチとなった
 - 推論ネットワークは、単に下界 $\mathcal{L}(q)$ を定義するために使用される
 - ネットワークのパラメータは、下界 $\mathcal{L}(q)$ が増加するように適用される