

1 近似推論法

EM アルゴリズムが困難な場合

- EM アルゴリズムで行う計算

- **E ステップ**では、潜在変数の事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ を計算
- **M ステップ**では、完全データ対数尤度 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の期待値を計算

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \quad (1)$$

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \int_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) d\mathbf{Z} \quad (2)$$

そして、 $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$ を最大化するパラメータ $\boldsymbol{\theta}^{\text{new}}$ を求める

$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) \quad (3)$$

EM アルゴリズムが困難な場合

- EM アルゴリズムの困難さ

- 実際に扱うモデルでは、事後分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ の計算や、事後分布に従った期待値 $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$ の計算が、**不可能であることが多い**
- 隠れ変数の**次元が高すぎる**
- 事後分布が**複雑な形をしていて**、期待値を解析的に計算できない
- **連続変数**であれば、積分が閉形式の解を持たないかもしれない
- 空間の次元の問題や、被積分項の複雑さから、数値積分すら困難かもしれない
- **離散変数**であれば、期待値を計算するためには、**潜在変数の可能な全ての組み合わせについての和を取る**必要がある
- 隠れ変数の次元が高くなると、組み合わせ数が指数的に増大する
- 計算量が大きすぎて、期待値の厳密な計算がもはや不可能

EM アルゴリズムが困難な場合

- 近似法

- EM アルゴリズムが困難であるとき、何らかの方法で近似しなければならない
- 近似法は、確率的な近似と、決定的な近似の2つに分けられる

- 確率的な近似

- マルコフ連鎖モンテカルロ法などの手法がある
- 無限の計算資源があれば、厳密な結果が得られる
- 実際には計算量が有限であるため、得られる解は近似解となる

- 決定的な近似

- 事後分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ を解析的に近似する
- 事後分布に対して、何らかの仮定をおく
- 例えば、単純な項の積として分解できる、あるいは、(ガウス分布などの特別な) パラメトリックな分布であるといった仮定

- ここで扱う近似法
 - 変分推論法 (Variational inference) あるいは変分ベイズ法 (Variational Bayes) について扱う
- 変分推論 (Variational inference)
 - 18 世紀のオイラー、ラグランジュらによる変分法 (Calculus of variations) に起源をもつ
 - まずは、変分法について説明をしていく

- 関数と汎関数の違い

- 通常関数は、入力として値をとり、出力として関数の値を返す
- 通常関数は、値から値への写像である
- 関数の導関数は、入力値を微小に変えたときに、出力の関数値がどの程度変わるかを表す
- 汎関数 (Functional) とは、入力として関数を取り、出力として汎関数の値を返す
- 汎関数は、関数から値への写像である
- 汎関数微分 (Functional derivative) とは、入力関数が微小に変わったときに、出力の汎関数値がどの程度変わるかを表す
- 汎関数の微分を、変分という

- 汎関数の例

- エントロピー $H[p]$ は、確率分布 $p(x)$ を入力として、以下の量を返す汎関数

$$H[p] = - \int p(x) \ln p(x) dx \quad (4)$$

- 汎関数の最適化

- 多くの問題は、**汎関数の値を最適化する問題**として定式化できる
- 汎関数の最適化とは、**可能な全ての入力関数の中から**、汎関数の値を最大化、あるいは最小化するような**関数を選び出す**ことである
- 通常最適化では、可能な全てのパラメータ (入力値) の中から、関数を最大化、あるいは最小化するような 1 つのパラメータを選び出す
- 次は、いよいよ**変分**の計算について説明する

- 変分法

- 通常の微分を使えば、ある関数 $y(x)$ を最大化 (最小化) するような x の値が求められる
- **変分法**を使えば、汎関数 $F[y]$ を最大化 (最小化) するような、関数 $y(x)$ が求められる
- 従って、可能な全ての関数 $y(x)$ の中から、 $F[y]$ を最大 (最小) にするような関数が得られる

- 変分法によって解ける問題の例

- 2 点を結ぶ最短経路は? (答えは直線)
- 最速降下曲線は? (答えはサイクロイド)
- **エントロピーが最大**になるような確率分布は? (答えは**ガウス分布**)

- 通常の微分の表現

- 関数 $y(x + \epsilon)$ のテイラー展開は次のように記述できた

$$y(x + \epsilon) = \sum_{n=0}^{\infty} \frac{y^{(n)}(x)}{n!} \epsilon^n \quad (5)$$

$$= y(x) + \frac{dy}{dx} \epsilon + \frac{1}{2!} \frac{d^2 y}{dx^2} \epsilon^2 + \frac{1}{3!} \frac{d^3 y}{dx^3} \epsilon^3 + \dots \quad (6)$$

$$= y(x) + \frac{dy}{dx} \epsilon + O(\epsilon^2) \quad (7)$$

- これより微分 dy/dx は、次のように求められる
- 変数 x に微小な変化 ϵ を加え、このときの関数値 $y(x + \epsilon)$ を ϵ の累乗形として表現する
- 最後に $\epsilon \rightarrow 0$ の極限をとればよい

$$\frac{dy}{dx} = \lim_{\epsilon \rightarrow 0} \frac{y(x + \epsilon) - y(x)}{\epsilon} \quad (8)$$

- 多変数関数 $y(x_1, \dots, x_D)$ の偏微分の表現
 - 多変数関数 $y(x_1, \dots, x_D)$ のテイラー展開は次のように記述できた

$$D^n = \left(\epsilon_1 \frac{\partial y}{\partial x_1} + \dots + \epsilon_D \frac{\partial y}{\partial x_D} \right)^n \quad (9)$$

上記のような演算子 D を考えれば

$$\begin{aligned} & y(x_1 + \epsilon_1, \dots, x_D + \epsilon_D) \\ &= \sum_{n=0}^{\infty} \frac{1}{n!} (D^n y)(x_1, \dots, x_D) \end{aligned} \quad (10)$$

$$\begin{aligned} &= y(x_1, \dots, x_D) + \sum_{i=1}^D \frac{\partial y}{\partial x_i} \epsilon_i + \frac{1}{2!} \sum_{i=1}^D \sum_{j=1}^D \frac{\partial^2 y}{\partial x_i \partial x_j} \epsilon_i \epsilon_j + \\ & \quad \frac{1}{3!} \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D \frac{\partial^3 y}{\partial x_i \partial x_j \partial x_k} \epsilon_i \epsilon_j \epsilon_k + \dots \end{aligned} \quad (11)$$

であるから

$$\begin{aligned} & y(x_1 + \epsilon_1, \dots, x_D + \epsilon_D) \\ = & y(x_1, \dots, x_D) + \sum_{i=1}^D \frac{\partial y}{\partial x_i} \epsilon_i + O(\epsilon^2) \end{aligned} \quad (12)$$

- これより偏微分 $\partial y / \partial x_i$ は、次のように求められる

$$\begin{aligned} \frac{\partial y}{\partial x_i} = \lim_{\epsilon_i \rightarrow 0} \frac{1}{\epsilon_i} & (y(x_1, \dots, x_{i-1}, x_i + \epsilon_i, x_{i+1}, \dots, x_D) - \\ & y(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_D)) \end{aligned} \quad (13)$$

- 変分の表現

- 多少不正確だが、変分をどのように定義すればよいか考えてみる
- ここで、各 x_i に対する関数の値 $z_i = y(x_i)$ を個別の変数とみなして、次の関数 $F(z_1, \dots, z_D)$ について考えてみよう

$$\begin{aligned} & F(z_1 + \epsilon\eta(x_1), \dots, z_D + \epsilon\eta(x_D)) \\ &= F(z_1, \dots, z_D) + \sum_{i=1}^D \frac{\partial F}{\partial z_i} \epsilon\eta(x_i) + O(\epsilon^2) \end{aligned} \quad (14)$$

$z_i = y(x_i)$ を代入してみると

$$\begin{aligned} & F(y(x_1) + \epsilon\eta(x_1), \dots, y(x_D) + \epsilon\eta(x_D)) \\ &= F(y(x_1), \dots, y(x_D)) + \sum_{i=1}^D \frac{\partial F}{\partial y(x_i)} \epsilon\eta(x_i) + O(\epsilon^2) \end{aligned} \quad (15)$$

変分法

- ここで $D \rightarrow \infty$ の極限を取り、 x_1, \dots, x_D が、ある連続した区間 $[a, b]$ に含まれる、全ての実数を表すことにする
- このとき $y(x_1), \dots, y(x_D)$ は、実数の区間 $[a, b]$ で定義される連続関数 $y(x)$ として書けることが分かる
- 同様に $y(x_1) + \epsilon\eta(x_1), \dots, y(x_D) + \epsilon\eta(x_D)$ は、実数の区間 $[a, b]$ で定義される連続関数 $y(x) + \epsilon\eta(x)$ として、まとめることができる
- 関数 $\eta(x)$ も、実数の区間 $[a, b]$ で定義される連続関数
- $\epsilon\eta(x)$ は、 $y(x)$ に加わる摂動として、考えることができる

変分法

- 関数 F は、関数 $y(x)$ や $y(x) + \epsilon\eta(x)$ を入力として受け取る、汎関数 $F[y]$ として解釈できるから、次のように書ける

$$F(y(x_1) + \epsilon\eta(x_1), \dots, y(x_D) + \epsilon\eta(x_D)) = F[y(x) + \epsilon\eta(x)] \quad (16)$$

$$F(y(x_1), \dots, y(x_D)) = F[y(x)] \quad (17)$$

- 以下の項は、入力を $y(x)$ に摂動を加えて $y(x) + \epsilon\eta(x)$ へと微小に変化させたときの、汎関数の ($F[y(x)]$ から $F[y(x) + \epsilon\eta(x)]$ への) 変化量を表している

$$\sum_{i=1}^D \frac{\partial F}{\partial y(x_i)} \epsilon\eta(x_i) \quad (18)$$

- 点 x_i における汎関数 F の変化量を、 x_1, \dots, x_D の範囲について、即ち、実数の区間 $[a, b]$ について足し合わせていると解釈する

変分法

- $D \rightarrow \infty$ のとき、 x_1, \dots, x_D は区間 $[a, b]$ における全ての実数を表すから、総和を積分に置き換えられそうである
- 汎関数の微分 $\frac{\delta F}{\delta y(x)}$ を使えば、次のように書ける

$$\begin{aligned} & \sum_{i=1}^D \frac{\partial F}{\partial y(x_i)} \epsilon \eta(x_i) \\ \Rightarrow & \int_a^b \frac{\delta F}{\delta y(x)} \epsilon \eta(x) dx = \epsilon \int_a^b \frac{\delta F}{\delta y(x)} \eta(x) dx \end{aligned} \quad (19)$$

- 結局、変分 $\frac{\delta F}{\delta y(x)}$ は次のように定義できる

$$F[y(x) + \epsilon \eta(x)] = F[y(x)] + \epsilon \int_a^b \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2) \quad (20)$$

変分法

- $F[y]$ は、区間 $[a, b]$ で定義される関数 y を受け取るとする
- 変分 $\delta F/\delta y$ は、入力関数 $y(x)$ に、任意の微小な変化 $\epsilon\eta(x)$ を加えたときの、汎関数 $F[y]$ の変化量として定義できる
- $\eta(x)$ は x についての任意の関数

Figure D.1 A functional derivative can be defined by considering how the value of a functional $F[y]$ changes when the function $y(x)$ is changed to $y(x) + \epsilon\eta(x)$ where $\eta(x)$ is an arbitrary function of x .

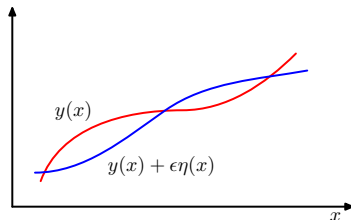


図 1: $y(x)$ と $y(x) + \epsilon\eta(x)$ の表現

- 変分法の例

- 次の図 2 を使って、実際に変分を求めてみよう

- 汎関数 $F[y]$ は、次のように定義されとする

$$F[y] = \int_a^b y(x) dx \quad (21)$$

- 汎関数の値 $F[y(x)], F[y(x) + \epsilon\eta(x)]$ は次のようになる

$$F[y(x)] = \int_a^b y(x) dx \quad (22)$$

$$F[y(x) + \epsilon\eta(x)] = \int_a^b (y(x) + \epsilon\eta(x)) dx \quad (23)$$

- $F[y(x) + \epsilon\eta(x)]$ は次のように分解できる

$$F[y(x) + \epsilon\eta(x)] = \int_a^b y(x)dx + \epsilon \int_a^b \eta(x)dx \quad (24)$$

$$= F[y(x)] + \epsilon \int_a^b \eta(x)dx \quad (25)$$

- ここで、変分の定義式は

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int_a^b \frac{\delta F}{\delta y(x)} \eta(x)dx + O(\epsilon^2) \quad (26)$$

であったので、上の2つの式を見比べれば、変分 $\delta F/\delta y$ は結局

$$\frac{\delta F}{\delta y(x)} = 1 \quad (27)$$

となることが分かる

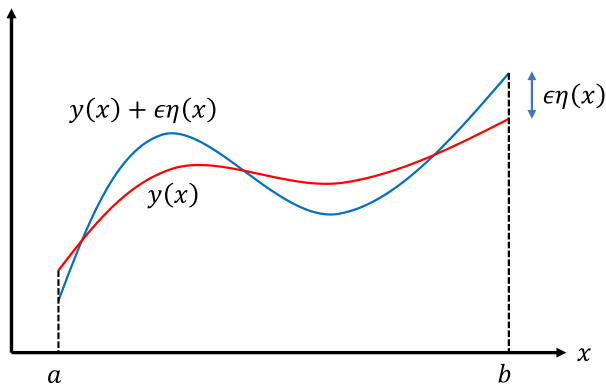


図 2: 区間 $[a, b]$ で定義された関数 $y(x)$ の表現

● 汎関数の最適化

- 汎関数 $F[y]$ が最大 (最小) となるとき、関数 $y(x)$ の微小な変化に対して、汎関数は変化しないはず
- 即ち、汎関数が最大 (最小) となるとき、 $F[y(x) + \epsilon\eta(x)] = F[y(x)]$ が成り立つ
- 従って、変分の定義式から、以下が成り立つ

$$\int_a^b \frac{\delta F}{\delta y(x)} \eta(x) dx = 0 \quad (28)$$

- 上式は任意の $\eta(x)$ について成立しなければならない
- 従って、変分 $\delta F/\delta y$ は、任意の x について 0 とならなければならない
- 汎関数 $F[y]$ が最大 (最小) となるとき、 $\delta F/\delta y = 0$ が成立することが分かった (通常の微分と同じ)

- 変分法の例

- 様々な汎関数について、変分を導出してみよう
- また、その汎関数が最大 (最小) となるときに成り立つ条件を、導出してみよう

- 汎関数の例 (1)

- $y(x)$ とその微分 $y'(x) = dy/dx$ 、そして x によって決まる関数 $G(y(x), y'(x), x)$ があるとする
- 汎関数 $F[y]$ を、 $G(y(x), y'(x), x)$ を区間 $[a, b]$ にわたって積分した結果を出力する関数として、次のように定める

$$F[y] = \int_a^b G(y(x), y'(x), x) dx \quad (29)$$

- 積分区間は無限であってもよいとする ($a = -\infty, b = \infty$ でもよい)

変分法

- $y(x)$ に摂動 $\epsilon\eta(x)$ を加えたときの、汎関数の値 $F[y(x) + \epsilon\eta(x)]$ を使って、変分 $\delta F/\delta y$ を調べてみる

$$F[y(x) + \epsilon\eta(x)] = \int_a^b G(y(x) + \epsilon\eta(x), y'(x) + \epsilon\eta'(x), x) dx \quad (30)$$

ここで、被積分項のテーラー展開を考えれば

$$\begin{aligned} & G(y(x) + \epsilon\eta(x), y'(x) + \epsilon\eta'(x), x) \\ = & G(y(x), y'(x), x) + \frac{\partial G}{\partial y} \epsilon\eta(x) + \\ & \frac{\partial G}{\partial y'} \epsilon\eta'(x) + \frac{\partial G}{\partial x} \cdot 0 + O(\epsilon^2) \end{aligned} \quad (31)$$

$$= G(y(x), y'(x), x) + \epsilon \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) + O(\epsilon^2) \quad (32)$$

であるから

$$\begin{aligned}
 & F[y(x) + \epsilon \eta(x)] \\
 = & \int_a^b G(y(x) + \epsilon \eta(x), y'(x) + \epsilon \eta'(x), x) dx \\
 = & \int_a^b \left(G(y(x), y'(x), x) + \right. \\
 & \left. \epsilon \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) + O(\epsilon^2) \right) dx \quad (33)
 \end{aligned}$$

$$\begin{aligned}
 = & \int_a^b G(y(x), y'(x), x) dx + \\
 & \epsilon \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx + O(\epsilon^2) \quad (34)
 \end{aligned}$$

$$= F[y(x)] + \epsilon \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx + O(\epsilon^2) \quad (35)$$

ここで

$$\begin{aligned} & \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx \\ = & \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) \right) dx + \int_a^b \left(\frac{\partial G}{\partial y'} \eta'(x) \right) dx \end{aligned} \quad (36)$$

$$\begin{aligned} = & \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) \right) dx + \\ & \left[\frac{\partial G}{\partial y'} \eta(x) \right]_a^b - \int_a^b \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \eta(x) dx \end{aligned} \quad (37)$$

$$= \left[\frac{\partial G}{\partial y'} \eta(x) \right]_a^b + \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (38)$$

である

- 途中の式変形では、部分積分を使っていることに注意
- いま、積分区間の両端において、 $y(x)$ の値は固定されているとする
- これを **固定端条件** という (図 3)
- このとき、 $\eta(a) = \eta(b) = 0$ であるから、上式の最初の項が消えて

$$\begin{aligned} & \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx \\ &= \left[\frac{\partial G}{\partial y'} \eta(x) \right]_a^b + \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \end{aligned} \quad (39)$$

$$= \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (40)$$

のようになる

- 従って、摂動を加えたときの汎関数の値 $F[y(x) + \epsilon\eta(x)]$ は

$$\begin{aligned} & F[y(x) + \epsilon\eta(x)] \\ = & F[y(x)] + \epsilon \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx + O(\epsilon^2) \quad (41) \end{aligned}$$

となる

- 上式を、変分の定義式と比べれば

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int_a^b \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2) \quad (42)$$

変分 $\delta F/\delta y$ は次のように書ける

$$\frac{\delta F}{\delta y(x)} = \frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \quad (43)$$

変分法

- 汎関数 $F[y]$ が最大 (最小) になるとき、変分 $\delta F/\delta y$ が 0 になる
- 従って、汎関数が最大 (最小) になるとき、以下の方程式が成り立つ

$$\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0 \quad (44)$$

- これをオイラー-ラグランジュ方程式という
- オイラー-ラグランジュ方程式は、次のような考え方で導出することもできる
- $F[y]$ が最大 (最小) であれば、摂動 $\epsilon\eta(x)$ によって $y(x)$ が少し変化しても、 $F[y]$ の値は変化しないはず
- 従って、 $F[y]$ が最大 (最小) であるとき、 $F[y]$ の ϵ による微分は 0 になるはず

- これを数式で表現すると、次のようになる

$$\left. \frac{\partial F[y]}{\partial \epsilon} \right|_{\epsilon=0} = 0 \quad (45)$$

左辺は通常の偏微分であり、これを計算すると

$$\begin{aligned} & \frac{\partial F[y]}{\partial \epsilon} \\ = & \frac{\partial}{\partial \epsilon} \int_a^b G(y, y', x) dx \end{aligned} \quad (46)$$

$$= \int_a^b \frac{\partial}{\partial \epsilon} G(y, y', x) dx \quad (47)$$

$$= \int_a^b \left(\frac{\partial G}{\partial y} \frac{\partial y}{\partial \epsilon} + \frac{\partial G}{\partial y'} \frac{\partial y'}{\partial \epsilon} + \frac{\partial G}{\partial x} \frac{\partial x}{\partial \epsilon} \right) dx \quad (48)$$

$$= \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx \quad (49)$$

$$= \left[\frac{\partial G}{\partial y'} \eta(x) \right]_a^b + \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (50)$$

$$(\because \eta(a) = \eta(b) = 0)$$

$$= \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (51)$$

$$= 0$$

- 上の式変形では、 $y = y(x) + \epsilon \eta(x)$ であるから

$$\frac{\partial y}{\partial \epsilon} = \eta(x) \quad (52)$$

$$\frac{\partial y'}{\partial \epsilon} = \frac{\partial}{\partial \epsilon} \left(\frac{\partial y}{\partial x} \right) = \frac{\partial}{\partial \epsilon} (y'(x) + \epsilon \eta'(x)) = \eta'(x) \quad (53)$$

$$\frac{\partial x}{\partial \epsilon} = 0 \quad (54)$$

が成立することを利用している

- 任意の $\eta(x)$ について、上式が恒等的に成り立つためには

$$\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0 \quad (55)$$

でなければならないことが分かり、先程と同様に、オイラー-ラグランジュ方程式を得る

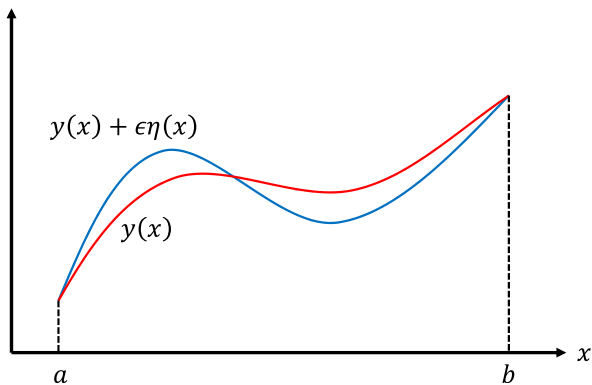


図 3: 制約条件を含んでいる場合の表現

- 汎関数の例 (2)

- 上では $G(y(x), y'(x), x)$ について考えて、変分を導出した
- $y(x)$ と x のみによって決まり、 $y'(x)$ には依存しない関数 $G(y(x), x)$ を考えよう
- 汎関数 $F[y]$ は、先程と同様に以下で表されとする

$$F[y] = \int_a^b G(y(x), x) dx \quad (56)$$

- このとき変分 $\delta F / \delta y$ を求めるのは、非常に簡単である
- 先程の式に、 $\partial G / \partial y' = 0$ を代入すれば直ちに得られる

$$\frac{\delta F}{\delta y(x)} = \frac{\partial G}{\partial y} \quad (57)$$

- あるいは以下のように書ける

$$\frac{\delta}{\delta y(x)} \int_a^b G(y(x), x) dx = \frac{\partial}{\partial y} G(y(x), x) \quad (58)$$

- $F[y]$ が最大 (最小) であるとき、以下のオイラー-ラグランジュ方程式が成り立つ

$$\frac{\partial G}{\partial y} = 0 \quad (59)$$

- 汎関数の例 (3)

- 今度は、 $y'(x)$ と x のみによって決まり、 $y(x)$ には依存しない関数 $G(y'(x), x)$ を考えよう
- この場合も変分 $\delta F / \delta y$ を求めるのは簡単である
- $G(y(x), y'(x), x)$ の変分の式に、 $\partial G / \partial y = 0$ を代入すればよい

$$\frac{\delta F}{\delta y(x)} = -\frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \quad (60)$$

- オイラー-ラグランジュ方程式は次のようになる

$$-\frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0 \quad (61)$$

● 汎関数の例 (4)

- $y(x)$ と $y'(x)$ によって決まる関数 $G(y(x), y'(x))$ を考えよう
- このときのオイラー-ラグランジュ方程式を導出してみよう
- $G(y(x), y'(x))$ を x で微分すれば

$$\begin{aligned} & \frac{d}{dx} G(y, y') \\ = & \frac{\partial}{\partial y} G(y, y') \frac{dy}{dx} + \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} \end{aligned} \quad (62)$$

$$= y' \frac{\partial}{\partial y} G(y, y') + \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} \quad (63)$$

となるから

$$y' \frac{\partial}{\partial y} G(y, y') = \frac{d}{dx} G(y, y') - \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} \quad (64)$$

また、オイラー-ラグランジュ方程式の両辺に y' を掛けたものは

$$y' \frac{\partial}{\partial y} G(y, y') - y' \frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y, y') \right) = 0 \quad (65)$$

これらを連立させて

$$y' \frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y, y') \right) = \frac{d}{dx} G(y, y') - \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} \quad (66)$$

$$y' \frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y, y') \right) + \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} = \frac{d}{dx} G(y, y') \quad (67)$$

$$\frac{d}{dx} \left(y' \cdot \frac{\partial}{\partial y'} G(y, y') \right) = \frac{d}{dx} G(y, y') \quad (68)$$

$$\int \left(\frac{d}{dx} \left(y' \cdot \frac{\partial}{\partial y'} G(y, y') \right) \right) dx = \int \left(\frac{d}{dx} G(y, y') \right) dx + C \quad (69)$$

$$G(y, y') = y' \cdot \frac{\partial}{\partial y'} G(y, y') + C \quad (70)$$

となるので、結局オイラー-ラグランジュ方程式は

$$G - y' \frac{\partial G}{\partial y'} = \text{Const.} \quad (71)$$

と書ける

変分法で解ける問題の例

- 変分法で解ける問題の例 (1)

- 2 点 $P(0, 0)$ 、 $Q(a, b)$ を結ぶ最短経路は?
- 2 点を結ぶ経路 $y = f(x) (0 \leq x \leq a)$ の長さ l は、次のようになる

$$l = \int_0^a \sqrt{1 + \left(\frac{dy}{dx}\right)^2} dx = \int_0^a \sqrt{1 + y'^2} dx \quad (72)$$

- 被積分項が y' のみの関数となっていることが分かる

変分法で解ける問題の例

- $G(y'(x), x)$ の場合の公式を使えば、 l の $y = f(x)$ による変分が求まる

$$\frac{\delta l}{\delta f(x)} = -\frac{d}{dx} \left(\frac{\partial}{\partial y'} \sqrt{1 + y'^2} \right) \quad (73)$$

$$= -\frac{d}{dx} \left(\frac{1}{2} \frac{1}{\sqrt{1 + y'^2}} \frac{\partial}{\partial y'} (1 + y'^2) \right) \quad (74)$$

$$= -\frac{d}{dx} \frac{y'}{\sqrt{1 + y'^2}} \quad (75)$$

- l を最小化するような $y = f(x)$ は、上式の変分を 0 と等置すれば

$$-\frac{d}{dx} \frac{y'}{\sqrt{1 + y'^2}} = 0 \quad (76)$$

$$\therefore \frac{y'}{\sqrt{1 + y'^2}} = \text{Const.} \quad (77)$$

変分法で解ける問題の例

- これは、 y' が定数であることを意味している
- 従って、 $y = C_0x + C_1$ と書ける
- 以上より、2 点間を結ぶ最短経路は直線である
- $y = f(x)$ の形について、具体的な仮定は特に置いていないことに注意
- 変分法では、関数そのものを最適化する
- 従って、関数の具体的な形については、特に仮定する必要がない

変分法のまとめ

- 変分のまとめ

- これまでの計算で、次の変分が明らかとなった

$$\frac{\delta}{\delta y(x)} \int G(y(x), y'(x), x) dx = \frac{\partial}{\partial y} G(y(x), y'(x), x) - \frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y(x), y'(x), x) \right) \quad (78)$$

$$\frac{\delta}{\delta y(x)} \int G(y(x), x) dx = \frac{\partial}{\partial y} G(y(x), x) \quad (79)$$

$$\frac{\delta}{\delta y(x)} \int G(y'(x), x) dx = -\frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y'(x), x) \right) \quad (80)$$

変分法のまとめ

- ここまでの話の流れ

- 1 変分の定義や、変分の計算法について調査した
- 2 **変分** (汎関数の微分) とは、入力関数が微小に変化したときの、出力値の変化量として定義される
- 3 汎関数が特定の形で表せるとき、変分がどのようなになるか計算した
- 4 汎関数が最大 (最小) になるとき、**オイラー-ラグランジュ方程式**が成立した
- 5 変分法を用いて、2 点間を結ぶ最短経路が**直線**になることを確認した

- これからの話の流れ

- 変分最適化を、どのように推論問題に適用するのかについて調べていく

- 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
- $p(Z|X, \theta)$ 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

- 変分推論の目的

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

- 注意点

- 観測変数と潜在変数をまとめて、 X, Z とおく
- $p(X)$ は、確率モデルからデータ 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、 $p(X)$ 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

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p) \quad (81)$$

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

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \quad (82)$$

$$\text{KL}(q||p) = - \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} \quad (83)$$

- 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、エビデンス下界ともいう
 - $\mathcal{L}(q)$ を q について最大化し、従って $\text{KL}(q||p)$ を最小化できれば、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いれば、 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ 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

- $q(\mathbf{Z})$ 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 $q(\mathbf{Z})$ 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

$$q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i) \quad (84)$$

- 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

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (85)$$

$$= \int \left(\prod_i q_i(\mathbf{Z}_i) \right) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (86)$$

$$= \int \prod_i q_i(\mathbf{Z}_i) \left(\ln p(\mathbf{X}, \mathbf{Z}) - \sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \quad (87)$$

$$\begin{aligned}
 &= \int \prod_i q_i(\mathbf{Z}_i) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} - \\
 &\quad \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \quad (88)
 \end{aligned}$$

ここで第 1 項は

$$\begin{aligned}
 &\int \prod_i q_i(\mathbf{Z}_i) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} \\
 &= \int q_j(\mathbf{Z}_j) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} \quad (89)
 \end{aligned}$$

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

$$= \int q_j(\mathbf{Z}_j) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (90)$$

$$= \int q_j(\mathbf{Z}_j) (\ln p(\mathbf{X}, \mathbf{Z})) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) \left(\prod_{i \neq j} d\mathbf{Z}_i \right) d\mathbf{Z}_j \quad (91)$$

$$= \int q_j(\mathbf{Z}_j) \left(\int (\ln p(\mathbf{X}, \mathbf{Z})) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \right) d\mathbf{Z}_j \quad (92)$$

$$= \int q_j(\mathbf{Z}_j) \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j \quad (93)$$

- 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)

$$\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (94)$$

- 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

$$\mathbb{E}_{i \neq j}[\ln p(\mathbf{X}, \mathbf{Z})] = \int (\ln p(\mathbf{X}, \mathbf{Z})) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z} \quad (95)$$

- 第2項は

$$\begin{aligned} & \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \\ &= \sum_i \int \prod_i q_i(\mathbf{Z}_i) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z} \end{aligned} \quad (96)$$

$$\begin{aligned} &= \sum_i \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \\ &= \int \prod_k q_k(\mathbf{Z}_k) (\ln q_j(\mathbf{Z}_j)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M + \end{aligned} \quad (97)$$

$$\sum_{i \neq j} \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (98)$$

但し

$$\int \prod_k q_k(\mathbf{Z}_k) (\ln q_j(\mathbf{Z}_j)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (99)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \left(\prod_{k \neq j} q_k(\mathbf{Z}_k) \right) \left(\prod_{k \neq j} d\mathbf{Z}_k \right) d\mathbf{Z}_j \quad (100)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \left(\int \prod_{k \neq j} q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right) d\mathbf{Z}_j \quad (101)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \prod_{k \neq j} \underbrace{\left(\int q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right)}_{=1} d\mathbf{Z}_j \quad (102)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \quad (103)$$

であるほか

$$\begin{aligned} & \sum_{i \neq j} \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \\ &= \sum_{i \neq j} \int q_i(\mathbf{Z}_i) q_j(\mathbf{Z}_j) \left(\prod_{k \neq i, j} \ln q_k(\mathbf{Z}_k) \right) \\ & \quad (\ln q_i(\mathbf{Z}_i)) \left(\prod_{k \neq i, j} d\mathbf{Z}_k \right) d\mathbf{Z}_i d\mathbf{Z}_j \quad (104) \\ &= \sum_{i \neq j} \underbrace{\left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j \right)}_{=1} \left(\int \prod_{k \neq i, j} \ln q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right) \end{aligned}$$

$$\int \ln q_i(\mathbf{Z}_i) q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (105)$$

$$= \sum_{i \neq j} \prod_{k \neq i, j} \underbrace{\left(\int \ln q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right)}_{=\text{Const.}} \underbrace{\int \ln q_i(\mathbf{Z}_i) q_i(\mathbf{Z}_i) d\mathbf{Z}_i}_{=\text{Const.}} \quad (106)$$

$$= \sum_{i \neq j} \text{Const.} = \text{Const.} \quad (107)$$

となるから結局

$$\begin{aligned} & \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \\ &= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{Const.} \end{aligned} \quad (108)$$

- 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

$$\begin{aligned}\mathcal{L}(q) = & \int q_j(\mathbf{Z}_j) \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \\ & \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{Const.}\end{aligned}\quad (109)$$

但し

$$\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (110)$$

- $\mathcal{L}(q)$ 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
- $q_j(\mathbf{Z}_j)$ 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

$$\mathcal{L}(q) = \int q_j(\mathbf{Z}_j) \ln \frac{\tilde{p}(\mathbf{X}, \mathbf{Z}_j)}{q_j(\mathbf{Z}_j)} d\mathbf{Z}_j + \text{Const.} \quad (111)$$

$$= -\text{KL}(q_j(\mathbf{Z}_j) \parallel \tilde{p}(\mathbf{X}, \mathbf{Z}_j)) + \text{Const.} \quad (112)$$

- 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

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j}[\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (113)$$

- 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

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (114)$$

- 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

$$q_j^*(\mathbf{Z}_j) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j} \quad (115)$$

- 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)

- 下界 $\mathcal{L}(q)$ の最適化

- 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

$$\ln q_j(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (116)$$

$$\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (117)$$

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
 - 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

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \boldsymbol{\Lambda} = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \quad (118)$$

- 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

$$\ln q_j(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (119)$$

- 従って、 $q_1^*(z_1)$ を計算する式は次のようになる

$$\ln q_1^*(z_1) = \mathbb{E}_{z_2} [\ln p(\mathbf{z})] + \text{Const.} \quad (120)$$

$$\mathbb{E}_{z_2} [\ln p(\mathbf{z})] = \int \ln p(\mathbf{z}) q_2(z_2) dz_2 \quad (121)$$

- 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

$$\ln q_1^*(z_1) = \mathbb{E}_{z_2} [\ln p(\mathbf{z})] + \text{Const.} \quad (122)$$

$$= \mathbb{E}_{z_2} [\ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})] + \text{Const.} \quad (123)$$

但し

$$\begin{aligned} & \ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ = & \ln \left(\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Lambda}^{-1}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T (\boldsymbol{\Lambda}^{-1})^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right) \right) \\ = & \ln \left(\frac{1}{2\pi} \frac{1}{|\boldsymbol{\Lambda}|^{-\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{z} - \boldsymbol{\mu}) \right) \right) \\ = & -\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{z} - \boldsymbol{\mu}) + \text{Const.} \end{aligned} \quad (124)$$

z_1 に依存する項だけを取り出せば

$$\begin{aligned} & -\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{z} - \boldsymbol{\mu}) \\ = & -\frac{1}{2} [z_1 - \mu_1, z_2 - \mu_2] \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \begin{bmatrix} z_1 - \mu_1 \\ z_2 - \mu_2 \end{bmatrix} \\ = & -\frac{1}{2} \left[\left\{ (\Lambda_{11}(z_1 - \mu_1) + \Lambda_{21}(z_2 - \mu_2)) \right\} (z_1 - \mu_1) + \right. \\ & \left. \left\{ \Lambda_{12}(z_1 - \mu_1) + \Lambda_{22}(z_2 - \mu_2) \right\} (z_2 - \mu_2) \right] \\ = & -\frac{1}{2} (\Lambda_{11}(z_1 - \mu_1)^2 + 2\Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \\ & \Lambda_{22}(z_2 - \mu_2)^2) \quad (\because \Lambda_{21} = \Lambda_{12}) \\ = & -\frac{1}{2} \Lambda_{11}(z_1 - \mu_1)^2 - \Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \text{Const.} \quad (125) \end{aligned}$$

これを代入して

$$\begin{aligned} & \ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ = & -\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 - \Lambda_{12} (z_1 - \mu_1) (z_2 - \mu_2) + \text{Const.} \quad (126) \end{aligned}$$

従って

$$\begin{aligned} & \ln q_1^*(z_1) \\ = & \mathbb{E}_{z_2} [\ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})] + \text{Const.} \\ = & \mathbb{E}_{z_2} \left[-\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 - \Lambda_{12} (z_1 - \mu_1) (z_2 - \mu_2) \right] + \text{Const.} \\ = & -\frac{1}{2} \Lambda_{11} z_1^2 + \Lambda_{11} \mu_1 z_1 - \Lambda_{12} z_1 (\mathbb{E}[z_2] - \mu_2) + \text{Const.} \quad (127) \end{aligned}$$

- これより、 $q_1^*(z_1)$ は次のように書ける

$$\begin{aligned} & q_1^*(z_1) \\ \propto & \exp\left(-\frac{1}{2}\Lambda_{11}z_1^2 + \Lambda_{11}\mu_1z_1 - \Lambda_{12}z_1(\mathbb{E}[z_2] - \mu_2)\right) \\ = & \exp\left(-\frac{1}{2}\Lambda_{11}\left(z_1 - (\mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(\mathbb{E}[z_2] - \mu_2))\right)^2 + \right. \\ & \left. \frac{1}{2}\Lambda_{11}\left(\mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(\mathbb{E}[z_2] - \mu_2)\right)^2\right) \\ \propto & \exp\left(-\frac{1}{2\Lambda_{11}^{-1}}\left(z_1 - (\mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(\mathbb{E}[z_2] - \mu_2))\right)^2\right) \\ = & \mathcal{N}(z_1|m_1, \Lambda_{11}^{-1}) \end{aligned} \tag{128}$$

但し m_1 は次のようにおいた

$$m_1 = \mu_1 - \Lambda_{11}^{-1}\Lambda_{12}(\mathbb{E}[z_2] - \mu_2) \tag{129}$$

- 対称性から、 $q_2^*(z_2)$ も次のように求められる

$$\ln q_2^*(z_2) = \mathbb{E}_{z_2}[\ln p(\mathbf{z})] + \text{Const.} \quad (130)$$

$$= \mathbb{E}_{z_2}[\ln \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})] + \text{Const.} \quad (131)$$

$$q_2^*(z_2) = \mathcal{N}(z_2|m_2, \Lambda_{22}^{-1}) \quad (132)$$

但し m_2 は次のようにおいた

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21}(\mathbb{E}[z_1] - \mu_1) \quad (133)$$

- これより $q_1^*(z_1), q_2^*(z_2)$ は

$$q_1^*(z_1) = \mathcal{N}(z_1|m_1, \Lambda_{11}^{-1}) \quad (134)$$

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12}(\mathbb{E}[z_2] - \mu_2) \quad (135)$$

$$q_2^*(z_2) = \mathcal{N}(z_2|m_2, \Lambda_{22}^{-1}) \quad (136)$$

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21}(\mathbb{E}[z_1] - \mu_1) \quad (137)$$

● 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

- $q_1(z_1), q_2(z_2)$ 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

$$KL(p||q) \equiv KL(p(\mathbf{Z}|\mathbf{X})||q(\mathbf{Z})) \quad (138)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})} d\mathbf{Z} \quad (139)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) (\ln q(\mathbf{Z}) - \ln p(\mathbf{Z}|\mathbf{X})) d\mathbf{Z} \quad (140)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln q(\mathbf{Z}) d\mathbf{Z} - \underbrace{\int p(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z}}_{q \text{ には依存しない定数項}} \quad (141)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln \prod_i q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (142)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \sum_i \ln q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (143)$$

$$= - \sum_i \int p(\mathbf{Z}|\mathbf{X}) \ln q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (144)$$

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

$$\begin{aligned} & \sum_i \int p(\mathbf{Z}|\mathbf{X}) \ln q_i(\mathbf{Z}_i) d\mathbf{Z} \\ &= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z} \end{aligned} \quad (145)$$

$$= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (146)$$

$$= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \left(\prod_{i \neq j} d\mathbf{Z}_i \right) \quad (147)$$

$$= \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \quad (148)$$

- $q_j(\mathbf{Z}_j)$ 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

$$\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j = 1 \quad (\text{規格化条件}) \quad (149)$$

$$q_j(\mathbf{Z}_j) \geq 0 \quad (150)$$

- 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

$$\begin{aligned} \mathcal{L}[q_j] = & - \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \\ & \lambda \left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j - 1 \right) \end{aligned} \quad (151)$$

- 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

$$\frac{\delta}{\delta y(x)} \int G(y(x), x) dx = \frac{\partial}{\partial y} G(y(x), x) \quad (152)$$

- 従って

$$\begin{aligned}
 & \frac{\delta}{\delta q_j(\mathbf{Z}_j)} \mathcal{L}[q_j] \\
 = & -\frac{\delta}{\delta q_j(\mathbf{Z}_j)} \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \\
 & \frac{\delta}{\delta q_j(\mathbf{Z}_j)} \lambda \left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j - 1 \right) \quad (153)
 \end{aligned}$$

$$\begin{aligned}
 = & -\frac{\partial}{\partial q_j} \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) + \\
 & \lambda \frac{\partial}{\partial q_j} q_j(\mathbf{Z}_j) \quad (154)
 \end{aligned}$$

$$= - \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} + \lambda = 0 \quad (155)$$

- これより、未定乗数 λ は

$$\lambda q_j(\mathbf{Z}_j) = \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} \quad (156)$$

$$\Rightarrow \int \lambda q_j(\mathbf{Z}_j) d\mathbf{Z}_j = \int \underbrace{\left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right)}_{=p(\mathbf{Z}_j|\mathbf{X})} d\mathbf{Z}_j \quad (157)$$

$$\Rightarrow \underbrace{\lambda \int q_j(\mathbf{Z}_j) d\mathbf{Z}_j}_{=1} = \underbrace{\int p(\mathbf{Z}_j|\mathbf{X}) d\mathbf{Z}_j}_{=1} \quad (158)$$

$$\Rightarrow \lambda = 1 \quad (159)$$

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

$$- \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} + \lambda = 0 \quad (160)$$

$$\Rightarrow q_j^*(\mathbf{Z}_j) = \underbrace{\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i}_{=p(\mathbf{Z}_j|\mathbf{X})} \quad (161)$$

$$\Rightarrow q_j^*(\mathbf{Z}_j) = p(\mathbf{Z}_j|\mathbf{X}) \quad (162)$$

- $q_j^*(\mathbf{Z}_j)$ 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

- 今回の場合は $p(z|x) = p(z)$ の場合を考えており、かつ $p(z) = \mathcal{N}(z|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ であった
- 従って、 $q_1^*(z_1)$ は

$$\begin{aligned} & q_1^*(z_1) \\ &= \int p(\mathbf{z}) d\mathbf{z}_2 \\ &= \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) d\mathbf{z}_2 \end{aligned} \tag{163}$$

$$= \frac{1}{2\pi} \frac{1}{|\boldsymbol{\Lambda}|^{-\frac{1}{2}}} \int \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{z} - \boldsymbol{\mu})\right) d\mathbf{z}_2 \tag{164}$$

ここで

$$-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{z} - \boldsymbol{\mu})$$

$$= -\frac{1}{2} (\Lambda_{11}(z_1 - \mu_1)^2 + 2\Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \Lambda_{22}(z_2 - \mu_2)^2) \quad (165)$$

$$= -\frac{1}{2} \Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 + \frac{1}{2} (2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}(z_2 - \mu_2)^2) \quad (166)$$

そして

$$= -\frac{1}{2} (2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}(z_2 - \mu_2)^2) \\ = -\frac{1}{2} (\Lambda_{22}z_2^2 - 2\Lambda_{22}\mu_2z_2 + 2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}\mu_2^2) \quad (167)$$

$$= -\frac{1}{2} (\Lambda_{22}z_2^2 - 2(\Lambda_{22}\mu_2 - \Lambda_{12}(z_1 - \mu_1))z_2 + \Lambda_{22}\mu_2^2) \quad (168) \\ = -\frac{1}{2} \left(\Lambda_{22} \left(z_2 - \Lambda_{22}^{-1} (\Lambda_{22}\mu_2 - \Lambda_{12}(z_1 - \mu_1)) \right)^2 - \right.$$

$$\Lambda_{22} \left(\Lambda_{22}^{-1} (\Lambda_{22} \mu_2 - \Lambda_{12}(z_1 - \mu_1)) \right)^2 + \Lambda_{22} \mu_2^2 \quad (169)$$

$$= -\frac{1}{2} \left(\Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 - \Lambda_{22}^{-1} m^2 + \Lambda_{22} \mu_2^2 \right) \quad (170)$$

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

$$\begin{aligned} & -\frac{1}{2} (z - \mu)^T \Lambda (z - \mu) \\ = & -\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 + \Lambda_{12} (z_1 - \mu_1) \mu_2 + \Lambda_{22} \mu_2^2 - \\ & \frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 + \frac{1}{2} \Lambda_{22}^{-1} m^2 \end{aligned} \quad (171)$$

これより

$$\int \exp \left(-\frac{1}{2} (z - \mu)^T \Lambda (z - \mu) \right) dz_2$$

$$\begin{aligned}
 &= \exp \left(-\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 + \Lambda_{12} (z_1 - \mu_1) \mu_2 + \right. \\
 &\quad \left. \Lambda_{22} \mu_2^2 + \frac{1}{2} \Lambda_{22}^{-1} m^2 \right) \\
 &\quad \int \exp \left(-\frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 \right) dz_2 \quad (172)
 \end{aligned}$$

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

$$\begin{aligned}
 &\int \exp \left(-\frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 \right) dz_2 \\
 &= (2\pi \Lambda_{22}^{-1}) \cdot \int \frac{1}{(2\pi \Lambda_{22}^{-1})^{\frac{1}{2}}} \exp \left(-\frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 \right) dz_2 \\
 &= 2\pi \Lambda_{22}^{-1} \quad (173)
 \end{aligned}$$

となって、 z_2 を積分により消去できるほか、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

$$\begin{aligned}
 & -\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 + \Lambda_{22}\mu_2^2 + \frac{1}{2}\Lambda_{22}^{-1}m^2 \\
 = & -\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 + \Lambda_{22}\mu_2^2 + \\
 & \frac{1}{2}\Lambda_{22}^{-1}(\Lambda_{22}\mu_2 - \Lambda_{12}(z_1 - \mu_1))^2 \\
 = & -\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 + \Lambda_{22}\mu_2^2 + \\
 & \frac{1}{2}\Lambda_{22}\mu_2^2 - \mu_2\Lambda_{12}(z_1 - \mu_1) + \\
 & \frac{1}{2}\Lambda_{22}^{-1}\Lambda_{12}^2(z_1 - \mu_1)^2 \tag{174} \\
 = & -\frac{1}{2}(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)(z_1 - \mu_1)^2 + \text{Const.}
 \end{aligned}$$

$$= -\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) (z_1 - \mu_1)^2 + \text{Const.} \quad (175)$$

$$= -\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) z_1^2 + \frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) \mu_1 z_1 + \text{Const.} \quad (176)$$