

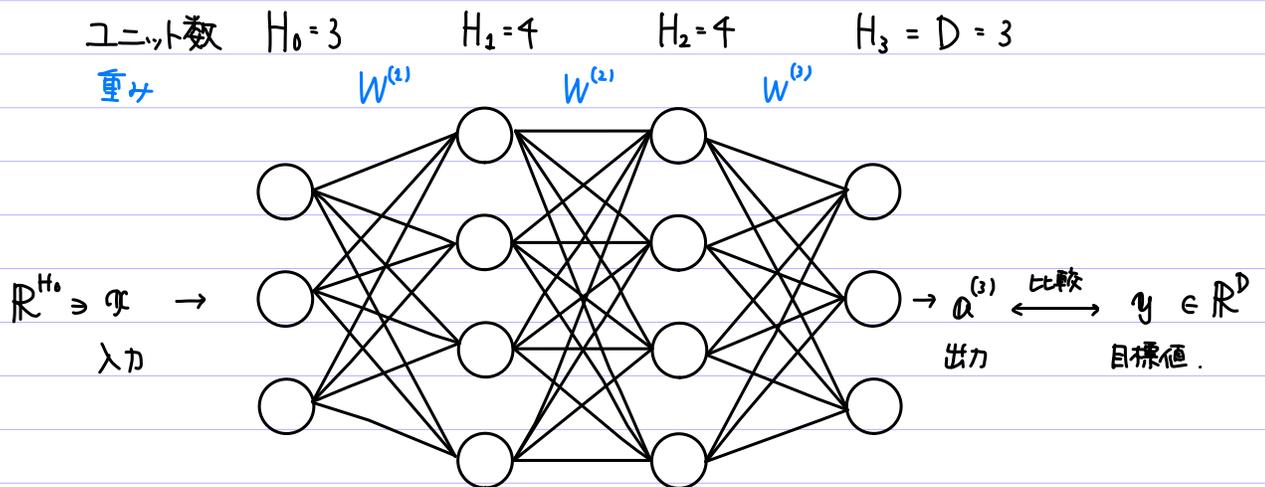
Ch.5 ニュラルネットワークのBayes推論.

Feedforward NN等の深層学習モデルをBayes的に扱って,

Ch.4でやった近似推論手法をNNモデルに適用していく.

用語・記法の復習 (Ch.2より)

順伝播型ニューラルネットワーク (feedforward NN)



入力層	隠れ層	出力層	層数
layer 0	layer 1	layer 3	$L = 3$

各層の出力: $z^{(0)} = x$, $z^{(l)} = \phi(a^{(l)})$ ($l=1, \dots, L-2$)

活性化: $a^{(l)} = W^{(l)} z^{(l-1)}$ ($l=1, \dots, L$)

重みの各要素を全てまとめて \mathcal{W} と書く. 1"クトルとして扱う.

重み \mathcal{W} をパラメータとする NN を $f(x; \mathcal{W})$ と表す. $f(x; \mathcal{W}) = a^{(L)}$.

§5.1 Bayes ニュラルネットワークモデルの近似推論法.

- NNのバッチ学習(全データ一度に学習)について考える.
- 回帰問題を扱う.(分類でも本筋は同じ)

5.1.1 Bayes ニュラルネットワークモデル

- ネットワークの挙動を支配するパラメータ \mathcal{W} に事前分布 $p(\mathcal{W})$ を設定することで Bayes 的に扱えるようにする.

入力 $\mathcal{X} = \{x_1, \dots, x_N\}$ に対する観測データ $\mathcal{Y} = \{y_1, \dots, y_N\}$: given.

- モデルを次のように設定:

$$p(\mathcal{Y}, \mathcal{W} | \mathcal{X}) = p(\mathcal{W}) p(\mathcal{Y} | \mathcal{X}, \mathcal{W}) = p(\mathcal{W}) \prod_{n=1}^N p(y_n | x_n, \mathcal{W})$$

$$\text{観測モデル: } p(y_n | x_n, \mathcal{W}) = \mathcal{N}(y_n | f(x_n; \mathcal{W}), \sigma_y^2 \mathbf{I}_D)$$

($\sigma_y^2 > 0$ は 固定のノイズパラメータ)

$$\text{事前分布: } p(\mathcal{W}) = \prod_{w \in \mathcal{W}} p(w) = \prod_{w \in \mathcal{W}} \mathcal{N}(w | 0, \sigma_w^2)$$

($\sigma_w^2 > 0$ は 固定のノイズパラメータ)

以下では, 簡単のため 出力次元を $D=1$ で考える.

5.1.2 Laplace 近似による学習.

5.1.2.1 事後分布の近似

1. 事後分布 $p(\mathcal{W} | \mathcal{Y}, \mathcal{X})$ の最大化問題の局所最適解 $\tilde{\mathcal{W}}$ を求める.

$\tilde{\mathcal{W}}$ が $p(\mathcal{W} | \mathcal{Y}, \mathcal{X})$ の局所最適解であるから,

$\tilde{\mathcal{W}}$ は $\log p(\mathcal{W} | \mathcal{Y}, \mathcal{X})$ の局所最適解でもある。

- $\log p(\mathcal{W} | \mathcal{Y}, \mathcal{X})$ に対して勾配降下法で“最小化する”は,

- $\log p(\mathcal{W} | \mathcal{Y}, \mathcal{X})$ の極小解 $\Leftrightarrow \log p(\mathcal{W} | \mathcal{Y}, \mathcal{X})$ の極大解 が得られる。

反復ステップ: 学習率 $\alpha > 0$ とし,

$$\mathcal{W}_{\text{new}} \leftarrow \mathcal{W}_{\text{old}} + \alpha \nabla_{\mathcal{W}} \log p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) \Big|_{\mathcal{W}=\mathcal{W}_{\text{old}}}$$

• $\nabla_{\mathcal{W}} \log p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) = \left(\frac{\partial}{\partial w} \log p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) \right)_{w \in \mathcal{W}}$ について,

$$p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) \propto p(\mathcal{W}) p(\mathcal{Y} | \mathcal{X}, \mathcal{W}) \quad \leftarrow \mathcal{W} \text{ に依存する項のみを考える} \\ \text{“なの”}$$

$$\log p(\mathcal{W} | \mathcal{Y}, \mathcal{X})$$

$$= \log p(\mathcal{Y} | \mathcal{X}, \mathcal{W}) + \log p(\mathcal{W}) + \text{const.}$$

$$= \sum_{n=1}^N \log \mathcal{N}(y_n | f(\mathcal{x}_n; \mathcal{W}), \sigma_y^2) + \sum_{w \in \mathcal{W}} \log \mathcal{N}(w | 0, \sigma_w^2) + \text{const.}$$

各項は,

$$\log \mathcal{N}(y_n | f(\mathcal{x}_n; \mathcal{W}), \sigma_y^2) = -\frac{1}{2\sigma_y^2} (y_n - f(\mathcal{x}_n; \mathcal{W}))^2 + \text{const.}$$

$$\log \mathcal{N}(w | 0, \sigma_w^2) = -\frac{w^2}{2\sigma_w^2} + \text{const.}$$

“なの”. 10-9 $w \in \mathcal{W}$ について微分すると,

$$\begin{aligned} & \frac{\partial}{\partial w} \log p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) \\ &= -\frac{1}{\sigma_y^2} \frac{\partial}{\partial w} \left(\underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - f(\mathcal{x}_n; \mathcal{W}))^2}_{=: E(\mathcal{W})} \right) - \frac{1}{\sigma_w^2} \frac{\partial}{\partial w} \left(\underbrace{\frac{1}{2} \sum_{w \in \mathcal{W}} w^2}_{=: \Omega_L(\mathcal{W})} \right) \\ &= -\left(\frac{1}{\sigma_y^2} \frac{\partial}{\partial w} E(\mathcal{W}) + \frac{1}{\sigma_w^2} \frac{\partial}{\partial w} \Omega_L(\mathcal{W}) \right). \end{aligned}$$

$= w$

$$\nabla_{\mathcal{W}} \log p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) = -\left(\frac{1}{\sigma_y^2} \nabla_{\mathcal{W}} E(\mathcal{W}) + \frac{1}{\sigma_w^2} \mathcal{W} \right).$$

第1項は backpropagation で計算できる.

→ ステップ4 $\mathcal{W}_{\text{new}} \leftarrow \mathcal{W}_{\text{old}} + \alpha \nabla_{\mathcal{W}} \log p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) \Big|_{\mathcal{W}=\mathcal{W}_{\text{old}}}$ はきろくと計算できる.

2. $p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) \approx q(\mathcal{W}) = \mathcal{N}(\mathcal{W} | \tilde{\mathcal{W}}, (\Lambda(\tilde{\mathcal{W}}))^{-1})$ と近似.

精度行列は,

$$\Lambda(\tilde{\mathcal{W}}) = -\nabla_{\mathcal{W}}^2 \log p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) = \frac{1}{\sigma_y^2} \underbrace{\nabla_{\mathcal{W}}^2 E(\mathcal{W})}_{\text{Hesse 行列}} + \frac{1}{\sigma_w^2} \mathbf{I}$$

5.1.2.2 予測分布の近似

次は予測分布 $p(y_* | \alpha_*, \mathcal{Y}, \mathcal{X}) = \int p(y_* | \alpha_*, \mathcal{W}) p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) d\mathcal{W}$

の近似をする. 上の事後分布の近似 $p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) \approx q(\mathcal{W})$ を使う.

$p(y_* | \alpha_*, \mathcal{W}) = \mathcal{N}(y_* | f(\alpha_*; \mathcal{W}), \sigma_y^2)$ と, NN $f(\alpha_*; \mathcal{W})$ が含まれるので

事後分布の近似としても依然として解析的に計算できない.

→ パラメータの事後分布の密度が $\tilde{\mathcal{W}}$ まわりに集中している, → サンプル数が増えればある程度は妥当か?

その近傍で $f(\alpha_*; \mathcal{W})$ が \mathcal{W} に関して線形近似できると仮定.

Taylor 展開をして 1次近似:

$$f(\alpha_*; \mathcal{W}) \approx f(\alpha_*; \tilde{\mathcal{W}}) + \underbrace{\left(\nabla_{\mathcal{W}} f(\alpha_*; \mathcal{W}) \Big|_{\mathcal{W}=\tilde{\mathcal{W}}} \right)^T}_{=: \mathbf{g}} (\mathcal{W} - \tilde{\mathcal{W}}).$$

結局, 予測分布は次で近似する.

$$\begin{aligned} & p(y_* | \alpha_*, \mathcal{Y}, \mathcal{X}) \\ & \approx \int \mathcal{N}(y_* | f(\alpha_*; \tilde{\mathcal{W}}) + \mathbf{g}^T (\mathcal{W} - \tilde{\mathcal{W}}), \sigma_y^2) \mathcal{N}(\mathcal{W} | \tilde{\mathcal{W}}, (\Lambda(\tilde{\mathcal{W}}))^{-1}) d\mathcal{W}. \end{aligned}$$

これに計算しよ。 $f = f(x_*; \tilde{\eta})$ と書く。

$$\begin{aligned} & \log \mathcal{N}(y_* | f + g^T(\eta - \tilde{\eta}), \sigma_y^2) \\ &= -\frac{1}{2\sigma_y^2} (y_* - (f + g^T(\eta - \tilde{\eta})))^2 + \text{const.} \\ &= -\frac{1}{2\sigma_y^2} (y_*^2 - 2(f + g^T(\eta - \tilde{\eta}))y_* + (f + g^T(\eta - \tilde{\eta}))^2) + \text{const.} \\ &= -\frac{1}{2\sigma_y^2} (y_*^2 - 2(f - g^T\tilde{\eta})y_* + (f - g^T\tilde{\eta})^2 - 2g^T\eta y_* + 2g^T\eta(f - g^T\tilde{\eta}) + (g^T\eta)^2) + \text{const.} \\ &= -\frac{1}{2\sigma_y^2} (y_* - (f - g^T\tilde{\eta}))^2 - \frac{1}{2\sigma_y^2} (\eta^T g g^T \eta - 2\eta^T g (y_* - (f - g^T\tilde{\eta}))) + \text{const.} \end{aligned}$$

また、

$$\begin{aligned} & \log \mathcal{N}(\eta | \tilde{\eta}, (\Lambda(\tilde{\eta}))^{-1}) \\ &= -\frac{1}{2} (\eta - \tilde{\eta})^T \Lambda(\tilde{\eta}) (\eta - \tilde{\eta}) + \text{const.} \\ &= -\frac{1}{2} (\eta^T \Lambda(\tilde{\eta}) \eta - 2\eta^T \Lambda(\tilde{\eta}) \tilde{\eta}) + \text{const.} \end{aligned}$$

これより、

$$\begin{aligned} & \log \mathcal{N}(y_* | f + g^T(\eta - \tilde{\eta}), \sigma_y^2) \mathcal{N}(\eta | \tilde{\eta}, (\Lambda(\tilde{\eta}))^{-1}) \\ &= -\frac{1}{2\sigma_y^2} (y_* - (f - g^T\tilde{\eta}))^2 \\ & \quad - \frac{1}{2} \left(\underbrace{\eta^T (\Lambda(\tilde{\eta}) + \frac{1}{\sigma_y^2} g g^T) \eta}_{=: \Lambda_*} - 2\eta^T (\underbrace{\Lambda(\tilde{\eta})\tilde{\eta} + \frac{1}{\sigma_y^2} g (y_* - (f - g^T\tilde{\eta}))}_{=: \Lambda_* \mathcal{M}(y_*)}) \right) + \text{const.} \\ &= \log \mathcal{N}(y_* | f - g^T\tilde{\eta}, \sigma_y^2) - \frac{1}{2} (\eta^T \Lambda_* \eta - 2\eta^T \Lambda_* \mathcal{M}(y_*)) + \text{const.} \\ &= \log \mathcal{N}(y_* | f - g^T\tilde{\eta}, \sigma_y^2) - \frac{1}{2} (\eta - \mathcal{M}(y_*))^T \Lambda_* (\eta - \mathcal{M}(y_*)) + \frac{1}{2} \mathcal{M}(y_*)^T \Lambda_* \mathcal{M}(y_*) + \text{const.} \\ &= \log \mathcal{N}(y_* | f - g^T\tilde{\eta}, \sigma_y^2) \mathcal{N}(\eta | \mathcal{M}(y_*), \Lambda_*^{-1}) + \frac{1}{2} \mathcal{M}(y_*)^T \Lambda_* \mathcal{M}(y_*) + \text{const.} \end{aligned}$$

と書くと、

$$\begin{aligned}
& \int \mathcal{N}(y_* | f + g^T(\alpha - \tilde{\alpha}), \sigma_y^2) \mathcal{N}(\alpha | \tilde{\alpha}, (\Lambda(\tilde{\alpha}))^{-1}) d\alpha \\
& \propto \mathcal{N}(y_* | f - g^T \tilde{\alpha}, \sigma_y^2) \exp\left(\frac{1}{2} \mathcal{M}(y_*)^T \Lambda_* \mathcal{M}(y_*)\right) \underbrace{\int \mathcal{N}(\alpha | \mathcal{M}(y_*), \Lambda_*^{-1}) d\alpha}_{=1} \\
& = \mathcal{N}(y_* | f - g^T \tilde{\alpha}, \sigma_y^2) \exp\left(\frac{1}{2} \mathcal{M}(y_*)^T \Lambda_* \mathcal{M}(y_*)\right).
\end{aligned}$$

∴ z",

$$\begin{aligned}
& \mathcal{M}(y_*)^T \Lambda_* \mathcal{M}(y_*) \\
& = (\Lambda(\tilde{\alpha}) \tilde{\alpha} + \frac{1}{\sigma_y^2} g (y_* - (f - g^T \tilde{\alpha})))^T (\Lambda(\tilde{\alpha}) + \frac{1}{\sigma_y^2} g g^T)^{-1} (\Lambda(\tilde{\alpha}) \tilde{\alpha} + \frac{1}{\sigma_y^2} g (y_* - (f - g^T \tilde{\alpha}))).
\end{aligned}$$

z"ある時, Sherman-Morrison-Woodbury (逆の公式)

$$\begin{aligned}
& (\Lambda(\tilde{\alpha}) + \frac{1}{\sigma_y^2} g g^T)^{-1} \quad \text{スカラー-ベクトル積の逆の逆数} \\
& = \Lambda(\tilde{\alpha})^{-1} - \frac{1}{\sigma_y^2} \Lambda(\tilde{\alpha})^{-1} g \left(1 + \frac{1}{\sigma_y^2} g^T \Lambda(\tilde{\alpha})^{-1} g\right)^{-1} g^T \Lambda(\tilde{\alpha})^{-1} \\
& = \Lambda(\tilde{\alpha})^{-1} - \underbrace{(\sigma_y^2 + g^T \Lambda(\tilde{\alpha})^{-1} g)^{-1}}_{=: \sigma(\alpha_*)^2} \Lambda(\tilde{\alpha})^{-1} g g^T \Lambda(\tilde{\alpha})^{-1} \\
& = \Lambda(\tilde{\alpha})^{-1} (\Lambda(\tilde{\alpha}) - \sigma(\alpha_*)^{-2} g g^T) \Lambda(\tilde{\alpha})^{-1}
\end{aligned}$$

z"のz",

$$\begin{aligned}
& \mathcal{M}(y_*)^T \Lambda_* \mathcal{M}(y_*) \\
& = (\tilde{\alpha} + \sigma_y^{-2} \Lambda(\tilde{\alpha})^{-1} g (y_* - (f - g^T \tilde{\alpha})))^T (\Lambda(\tilde{\alpha}) - \sigma(\alpha_*)^{-2} g g^T) (\tilde{\alpha} + \sigma_y^{-2} \Lambda(\tilde{\alpha})^{-1} g (y_* - (f - g^T \tilde{\alpha}))) \\
& = \underbrace{\tilde{\alpha}^T (\Lambda(\tilde{\alpha}) - \sigma(\alpha_*)^{-2} g g^T) \tilde{\alpha}}_{\text{const.}} + 2 \sigma_y^{-2} \tilde{\alpha}^T (\Lambda(\tilde{\alpha}) - \sigma(\alpha_*)^{-2} g g^T) \Lambda(\tilde{\alpha})^{-1} g (y_* - (f - g^T \tilde{\alpha})) \\
& \quad + \sigma_y^{-4} (y_* - (f - g^T \tilde{\alpha}))^T g^T \Lambda(\tilde{\alpha})^{-1} (\Lambda(\tilde{\alpha}) - \sigma(\alpha_*)^{-2} g g^T) \Lambda(\tilde{\alpha})^{-1} g \\
& = 2 \sigma_y^{-2} g^T \tilde{\alpha} \underbrace{(1 - \sigma(\alpha_*)^{-2} g^T \Lambda(\tilde{\alpha})^{-1} g)}_{=: \sigma_y^2 \sigma(\alpha_*)^{-2}} (y_* - (f - g^T \tilde{\alpha})) \\
& \quad + \sigma_y^{-4} (y_* - (f - g^T \tilde{\alpha}))^T \underbrace{g^T \Lambda(\tilde{\alpha})^{-1} g}_{\sigma(\alpha_*)^{-2} - \sigma_y^2} \underbrace{(1 - \sigma(\alpha_*)^{-2} g^T \Lambda(\tilde{\alpha})^{-1} g)}_{=: \sigma_y^2 \sigma(\alpha_*)^{-2}} + \text{const.}
\end{aligned}$$

$$= 2\sigma(\alpha_*)^{-2} \mathbf{g}^T \tilde{\boldsymbol{\eta}} (y_* - (\mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}})) + (\sigma_y^{-2} - \sigma(\alpha_*)^{-2}) (y_* - (\mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}}))^2 + \text{const.}$$

結局,

$$\begin{aligned} & \log \mathcal{N}(y_* | \mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}}, \sigma_y^2) \exp\left(\frac{1}{2} \mathcal{M}(y_*)^T \Lambda_* \mathcal{M}(y_*)\right) \\ &= -\frac{1}{2} \sigma_y^{-2} (y_* - (\mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}}))^2 \\ & \quad + \frac{1}{2} (2\sigma(\alpha_*)^{-2} \mathbf{g}^T \tilde{\boldsymbol{\eta}} (y_* - (\mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}})) + (\cancel{\sigma_y^{-2}} - \sigma(\alpha_*)^{-2}) (y_* - (\mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}}))^2) + \text{const.} \\ &= -\frac{1}{2} \sigma(\alpha_*)^{-2} ((y_* - (\mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}}))^2 - 2\mathbf{g}^T \tilde{\boldsymbol{\eta}} (y_* - (\mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}}))) + \text{const.} \\ &= -\frac{1}{2} \sigma(\alpha_*)^{-2} (y_* - (\mathbf{f} - \mathbf{g}^T \tilde{\boldsymbol{\eta}}))(y_* - (\mathbf{f} + \mathbf{g}^T \tilde{\boldsymbol{\eta}})) + \text{const.} \\ &= -\frac{1}{2} \sigma(\alpha_*)^{-2} (y_*^2 - 2\mathbf{f} y_*) + \text{const.} \\ &= \log \mathcal{N}(y_* | f(\alpha_*; \tilde{\boldsymbol{\eta}}), \sigma(\alpha_*)) + \text{const.} \end{aligned}$$

以上より,

$$p(y_* | \alpha_*, y. \setminus *) \approx \mathcal{N}(y_* | f(\alpha_*; \tilde{\boldsymbol{\eta}}), \sigma_y^2 + \mathbf{g}^T \Lambda(\tilde{\boldsymbol{\eta}})^{-1} \mathbf{g})$$

5.1.3 ハミルトニアンモンテカルロ法による学習.

HMCで $p(\mathcal{M} | \mathcal{Y}, \mathcal{F})$ からサンプリングをする.

→ サンプリングした \mathcal{M} を使って $y_* = f(x_*; \mathcal{M})$ と予測できる.

複数回サンプリングすれば, 予測の不確実性も分かる.

・ 離散変数をパラメータに含めないなら, $\nabla_{\mathcal{M}} \log p(\mathcal{M} | \mathcal{Y}, \mathcal{F})$ の計算は backpropagation でできる.

→ ポテンシャルの微分が出来るから HMC が使える

・ 計算時間は比較的長いものの, 理論的には真の事後分布からのサンプルが得られる.

5.1.3.1 重みパラメータの推論.

$$\log p(\mathcal{M} | \mathcal{Y}, \mathcal{F}) = \log p(\mathcal{Y} | \mathcal{F}, \mathcal{M}) + \log p(\mathcal{M}) + \text{const. term}$$

ポテンシャルエネルギー とし $\mathcal{U}(\mathcal{M}) = \log p(\mathcal{Y} | \mathcal{F}, \mathcal{M}) + \log p(\mathcal{M})$

と利用すればよい.

リ-プロダクト法で, この微分が必要になる (backpropagation でできる)

・ NN に対する HMC の適用に関する問題点.

× 複雑な分布からのサンプリングになるので, 少ないサンプリング数では分布の特徴を

とらえきかているかどうかわからない. 適切なサンプリング数も不明.

× ステップサイズ ϵ やステップ数 L の調整が難しい.

× 学習が低速.

5.1.3.2 ハイパーパラメータの推論.

- 上では σ_w^2, σ_y^2 をハイパーパラメータとして扱ってきた.

これにも事前分布を与えることで、推論ができるようにする.

- 以下, 精度パラメータ $\gamma_w = \sigma_w^{-2}, \gamma_y = \sigma_y^{-2}$ を利用する.

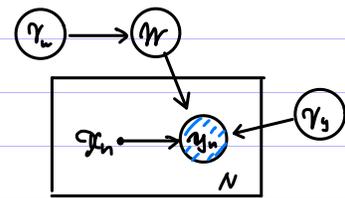
事前分布: $p(\gamma_w) = \text{Gam}(\gamma_w | a_w, b_w)$

← ガンマ分布は $\gamma_w > 0$ の確率分布. 正の実数をとる値の分布として使う.

$$p(\gamma_y) = \text{Gam}(\gamma_y | a_y, b_y)$$

($a_w, b_w, a_y, b_y > 0$: 固定値)

- モデル.



$$p(\mathcal{Y}, \mathcal{N}, \gamma_w, \gamma_y | \mathcal{X}) = p(\gamma_w) p(\gamma_y) p(\mathcal{N} | \gamma_w) \prod_{n=1}^N p(y_n | \mathcal{X}_n, \mathcal{N}, \gamma_y)$$

- 事後分布は, $p(\mathcal{N}, \gamma_w, \gamma_y | \mathcal{Y}, \mathcal{X})$.

- Gibbsサンプリングをする.

- γ_w, γ_y : given として, \mathcal{N} の条件付き分布は $p(\mathcal{N} | \mathcal{Y}, \mathcal{X}, \gamma_w, \gamma_y)$.

これは \mathcal{N} の事後分布なので, HMC で \mathcal{N} をサンプリングできる

- \mathcal{N}, γ_y : given として, γ_w の条件付き分布は,

$$p(\gamma_w | \mathcal{Y}, \mathcal{X}, \mathcal{N}, \gamma_y)$$

$$\propto p(\mathcal{Y}, \mathcal{N}, \gamma_y, \gamma_w | \mathcal{X})$$

$$\propto p(\mathcal{N} | \gamma_w) p(\gamma_w)$$

$p(\mathcal{N} | \gamma_w)$ は Gauss 分布, $p(\gamma_w)$ はその共役事前分布のガンマ分布なので,

条件付き分布 $p(\gamma_w | y, \mathcal{X}, W, \gamma_y)$ はガンマ分布. $K_w := (W \text{ のパラメータ数})$ とし,

Ch.3 の結果より

$$\gamma_w \sim \text{Gam}(\hat{a}_w, \hat{b}_w), \quad \hat{a}_w = a_w + \frac{1}{2} K_w, \quad \hat{b}_w = b_w + \frac{1}{2} \sum_{w \in W} w^2$$

とサンプリングすればよい.

• W, γ_w : given とし. γ_y の条件付き分布は,

$$p(\gamma_y | y, \mathcal{X}, W, \gamma_w)$$

$$\propto p(y, W, \gamma_y, \gamma_w | \mathcal{X})$$

$$\propto p(\gamma_y) p(y | \mathcal{X}, W, \gamma_y)$$

$p(y | \mathcal{X}, W, \gamma_y)$ は Gauss 分布で, $p(\gamma_y)$ はガンマ分布なので,

$p(\gamma_y | y, \mathcal{X}, W, \gamma_w)$ もガンマ分布.

$$\therefore \gamma_y \sim \text{Gam}(\hat{a}_y, \hat{b}_y),$$

$$\hat{a}_y = a_y + \frac{1}{2} N,$$

現在の W で表現できない NN の誤差の情報と
1111-パラメータに加味している.

$$\hat{b}_y = b_y + \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n; W))^2 = b_y + E(W).$$

$$\cdot \text{ガンマ分布の平均: } E[\gamma_y] = \frac{\hat{a}_y}{\hat{b}_y}$$

→ \hat{b}_y が大きい \Leftrightarrow 誤差 $E(W)$ が大きい.

γ_y の逆数

このとき $E[\gamma_y]$ は小さくなるので, y_n に対する分散は大きくなる

• 重みの 1111-パラメータ γ_w は共通としていたが, 層ごとパラメータを渡せる, などの

変更もできる.

§5.2 近似 Bayes 推論の効率化

- Bayes NN:

✗ パラメータの周辺化に伴う計算量が膨大.

学習に必要なデータ数が多い

→ ✗ バッチ学習は効率が良い.

- 巨大なネットワークや膨大なデータ量に対しても、近似的な Bayes 推論ができるようにする.

○ ミニバッチを利用した効率的な学習.

5.2.1 確率的勾配 Langevin 動力学法による学習.

- SGD

✓ 大規模な NN モデルの効率的な学習.

✗ 正則化項を追加した最適化や MAP 推定などでは、

パラメータの不確実性を取り扱えない.

→ 過適合をおこす.

不確実性に基づき予測やモデルの評価ができない.

- HMC ← 勾配を利用したサンプリング方法

- よく利用されている

✗ (MCMC 全般的に) 大規模データに対し、計算効率が悪い.

確率的勾配 Langevin 動力学法

(Stochastic gradient Langevin dynamics method)

・ SGD + Langevin dynamics method

✓ 計算効率が高い.

✓ 不確実性の推定が可能.

確率的 MCMC の一種.

・ SGD のアルゴリズムを用いて, NN の正則化項付きのコスト関数を最適化.

ミニバッチを指定する添字集合を \mathcal{S} と書く. $|\mathcal{S}| = M$.

ミニバッチ: $\mathcal{D}_{\mathcal{S}} = \{(x_n, y_n)\}_{n \in \mathcal{S}}$.

$\mathcal{X}_{\mathcal{S}} := \{x_n\}_{n \in \mathcal{S}}$, $\mathcal{Y}_{\mathcal{S}} := \{y_n\}_{n \in \mathcal{S}}$

MAP 推定

$$\mathcal{W}_{\text{MAP}} = \underset{\mathcal{W}}{\operatorname{argmax}} \log p(\mathcal{W} | \mathcal{Y}_{\mathcal{S}}, \mathcal{X}_{\mathcal{S}})$$

これは正則化項付きのコスト関数の最小化と等価だから:

$$\begin{aligned} & \mathbb{V}_{\mathcal{W}} \log p(\mathcal{W} | \mathcal{Y}_{\mathcal{S}}, \mathcal{X}_{\mathcal{S}}) \\ &= \mathbb{V}_{\mathcal{W}} \log p(\mathcal{Y}_{\mathcal{S}} | \mathcal{W}, \mathcal{X}_{\mathcal{S}}) + \mathbb{V}_{\mathcal{W}} \log p(\mathcal{W}) - \mathbb{V}_{\mathcal{W}} \log p(\mathcal{Y}_{\mathcal{S}} | \mathcal{X}_{\mathcal{S}}) \\ &= \sum_{n \in \mathcal{S}} \mathbb{V}_{\mathcal{W}} \log p(y_n | x_n, \mathcal{W}) + \mathbb{V}_{\mathcal{W}} \log p(\mathcal{W}). \end{aligned}$$

であり, ミニバッチに対するコスト関数が, 期待値としては全データに対する

コスト関数と等価になるよう係数を調整すると,

t 回目のパラメータの更新は

$$W \leftarrow W + \frac{\alpha_t}{2} \left(\frac{N}{M} \sum_{n \in \mathcal{S}} \nabla_{\theta} \log p(y_n | x_n, W) + \nabla_{\theta} \log p(W) \right)$$

係数は便宜上この設定が

とすればよい。学習率 α_t は最適解への収束のため

$$\sum_{i=1}^{\infty} \alpha_i = \infty, \quad \sum_{i=1}^{\infty} \alpha_i^2 < \infty$$

とみたらよいようにとる。

- Langevin 動力学法で事後分布 $p(W | y, \mathcal{X})$ からサンプル候補 W^* を得る。

ポテンシャルエネルギー $\mathcal{U} = -\log p(W | y, \mathcal{X})$,

運動量ベクトル $p \sim \mathcal{N}(p | 0, I)$ とし,

ステップサイズ $\varepsilon = \sqrt{\alpha_t}$ とすると,

$$\begin{aligned} W^* &= W - \frac{\varepsilon^2}{2} \nabla_{\theta} \mathcal{U} + \varepsilon p \\ &= W + \frac{\alpha_t}{2} \left(\sum_{n=1}^N \nabla_{\theta} \log p(y_n | x_n, W) + \nabla_{\theta} \log p(W) \right) + \sqrt{\alpha_t} p. \end{aligned}$$

- W^* の受容率は、 ε が小さい (i.e., α_t が小さい) ほど 1 に近づく。 → ミニバッチ版は誤差が小さくなるが。

- 上の二つのアプローチで、Langevin 動力学法のミニバッチ版 ε

$$W \leftarrow W + \frac{\alpha_t}{2} \left(\frac{N}{M} \sum_{n \in \mathcal{S}} \nabla_{\theta} \log p(y_n | x_n, W) + \nabla_{\theta} \log p(W) \right) + \sqrt{\alpha_t} p$$

$$p \sim \mathcal{N}(p | 0, I)$$

とする。

✓ $\nabla_{\theta} \log p(W | y, \mathcal{X})$ の不偏推定量が得られる。

✓ $t \rightarrow \infty$ でサンプル W の受容率が漸近的に 1 。

⇒ はじめは SGD の利点を活かして \mathcal{W} を効率的に探索.

t が大きくなるにつれて Langevin 動力学法による真の事後分布からの

近似的なサンプルを得ることができるようになる.

5.2.2 確率的変分推論法による学習

・ **確率的変分推論法** (stochastic variational inference method)

SGD + 変分推論 → 効率的に学習

・ $p(\mathcal{W} | \mathcal{Y}, \mathcal{X}) \approx q(\mathcal{W}; \xi)$ と近似.

ここで $\xi = \{ \mu_{ij}^{(q)}, \sigma_{ij}^{(q)2} \}_{i,j=1}^D$ とし $q(\mathcal{W}; \xi) = \prod_{i,j=1}^D \mathcal{N}(w_{ij}^{(q)} | \mu_{ij}^{(q)}, \sigma_{ij}^{(q)2})$ とする.

ELBO は,

$$\begin{aligned} \mathcal{L}(\xi) &= \int q(\mathcal{W}; \xi) \log \frac{p(\mathcal{W}, \mathcal{Y} | \mathcal{X})}{q(\mathcal{W}; \xi)} d\mathcal{W} \\ &= \int q(\mathcal{W}; \xi) \left(\log \frac{p(\mathcal{Y} | \mathcal{W}, \mathcal{X})}{q(\mathcal{W}; \xi)} + \log p(\mathcal{W}) \right) d\mathcal{W} \\ &= \int q(\mathcal{W}; \xi) \log p(\mathcal{Y} | \mathcal{W}, \mathcal{X}) d\mathcal{W} - D_{\text{KL}}[q(\mathcal{W}; \xi) \| p(\mathcal{W})] \\ &= \sum_{n=1}^N \int q(\mathcal{W}; \xi) \log p(y_n | f(\mathbf{x}_n; \mathcal{W})) d\mathcal{W} - D_{\text{KL}}[q(\mathcal{W}; \xi) \| p(\mathcal{W})]. \end{aligned}$$

・ $\mathcal{L}(\xi)$ を勾配降下法で ξ について最大化しようとすると、各ステップでの

$\nabla_{\xi} \mathcal{L}(\xi)$ の評価のために全データを読みこむ必要がある.

→ 効率的でない.

・ ミニバッチ \mathcal{D}_s に対する部分的な ELBO を以下のように設定:

$$\mathcal{L}_s(\xi) := \frac{N}{M} \sum_{n \in \mathcal{D}_s} \int q(\mathcal{W}; \xi) \log p(y_n | f(\mathbf{x}_n; \mathcal{W})) d\mathcal{W} - D_{\text{KL}}[q(\mathcal{W}; \xi) \| p(\mathcal{W})].$$

→ $E_{\xi}[\mathcal{L}_{\xi}(\xi)] = \mathcal{L}(\xi)$ (不偏推定量) となる.

(pf. は Ch.2 でやったのと同様.)

$\mathcal{L}(\xi)$ を最大化する代わりに, $\mathcal{L}_{\xi}(\xi)$ を最大化するようにすればよい.

5.2.3 勾配のモンテカルロ近似.

- $\mathcal{L}_{\xi}(\xi)$ に対して勾配降下法を行うには, $\nabla_{\xi} \mathcal{L}_{\xi}(\xi)$ を計算する必要がある
 - $\nabla_{\xi} D_{KL}[q_{\xi}(w; \xi) \| p(w)]$ は, と55も Gauss 分布なので解析的に計算できる
 - $\nabla_{\xi} \int q_{\xi}(w; \xi) \log p(y_n | f(w; \xi)) dW$ は解析的に計算できない.
- モンテカルロ法で積分を近似し, 勾配の推定を得る.

- 簡易的に表記する パラメータ $w \in \mathbb{R}$ に対して, ある関数 $f(w)$ と

分布 $q(w; \xi)$ を考えて, 勾配

$$I(\xi) = \nabla_{\xi} \int f(w) q(w; \xi) dw$$

を評価したい.

5.2.3.1 スコア関数推定

- 分布 $q(w; \xi)$ に対し, $\nabla_{\xi} \log q(w; \xi)$ を **スコア関数** という.

$$\begin{aligned} I(\xi) &= \int f(w) \nabla_{\xi} q(w; \xi) dw \\ &= \int f(w) q(w; \xi) \nabla_{\xi} \log q(w; \xi) dw \\ &= E_{q(w; \xi)} [f(w) \nabla_{\xi} \log q(w; \xi)]. \end{aligned}$$

微分と積分の交換

スコア関数

→ 分布 $q(w; \xi)$ から $w \in$ 複素数 サンプルリング, $f(w) \nabla_{\xi} \log q(w; \xi)$ の サンプル平均をとることによって $I(\xi)$ を推定できる.

✓ $\nabla_{\xi} \log q(w; \xi)$ が計算できれば利用できる.

✗ 分散が大きいの. → 制御変量法 (control variates method) と併用.

5.2.3.2 再パラメータ化勾配.

・ 再パラメータ化勾配 (reparametrization gradient)

$w \sim q(w; \xi)$ とする代わりに, 変分パラメータ ξ の新しい分布 $p(\varepsilon)$ から

$\varepsilon \sim p(\varepsilon)$ とサンプルリング, $w = g(\xi; \varepsilon)$ と変換することによって w を得る.

→ $I(\xi)$

$$\int q(w; \xi) dw = p(\varepsilon) d\varepsilon$$

↑ w が q に従うように p と g を設定する.

$$= \nabla_{\xi} \int f(g(\xi; \varepsilon)) p(\varepsilon) d\varepsilon$$

$$= \int f'(g(\xi; \varepsilon)) (\nabla_{\xi} g(\xi; \varepsilon)) p(\varepsilon) d\varepsilon$$

$$= \mathbb{E}_{p(\varepsilon)} [f'(g(\xi; \varepsilon)) \nabla_{\xi} g(\xi; \varepsilon)].$$

分布 $p(\varepsilon)$ から $\varepsilon \in$ 複素数 サンプルリング, $f'(g(\xi; \varepsilon)) \nabla_{\xi} g(\xi; \varepsilon)$ の

サンプル平均をとることによって $I(\xi)$ を推定できる.

ex) $\xi = \{\hat{\mu}, \hat{\sigma}\}$ とする. $q(w; \xi) = \mathcal{N}(w | \hat{\mu}, \hat{\sigma}^2)$ の場合を考える.

$\tilde{\varepsilon} \sim \mathcal{N}(0, 1)$ とし, $\tilde{w} = g(\xi; \varepsilon) = \hat{\mu} + \hat{\sigma} \tilde{\varepsilon}$ とすれば

$\tilde{w} \sim \mathcal{N}(w | \hat{\mu}, \hat{\sigma}^2)$ となる.

このとき, $\xi = \{\hat{\mu}, \hat{\sigma}\}$ に関する勾配は,

$$\frac{\partial}{\partial \mu} \int f(w) q(w; \xi) dw = \int f'(w) q(w; \xi) dw = \mathbb{E}_{q(w; \xi)} [f'(w)].$$

$$\begin{aligned} \frac{\partial}{\partial \hat{\mu}} \int f(w) q(w; \xi) dw &= \int f'(w) \tilde{\varepsilon} q(w; \xi) dw = \int f'(w) \frac{w - \hat{\mu}}{\sigma} q(w; \xi) dw \\ &= \mathbb{E}_{q(w; \xi)} \left[f'(w) \frac{w - \hat{\mu}}{\sigma} \right]. \end{aligned}$$

5.2.3.3 再パラメータ化勾配の一般化.

再パラメータ化勾配

✓ 勾配の分散を小さく抑えられる傾向がある.

✗ ξ に依存しない $p(\varepsilon)$ をつくれないことが多い.

・ ガンマ分布, ベータ分布などでは, 再パラメータ化勾配による勾配推定ができない.

一般化再パラメータ化勾配 (generalized reparametrization gradient):

g に関する制約を緩和して, ε の分布に対して ξ の依存性が残るのを許容する ($\varepsilon \sim p(\varepsilon; \xi)$).

陰関数微分 (implicit differentiation) を用いる方法.

・ g を求めるのは困難だが, g^{-1} は容易に得られるケース.

ex) ガンマ分布, Dirichlet分布, von Mises分布, etc.

cf.) von Mises分布: 角度 $\theta \in [0, 2\pi)$ の分布

$$\text{pdf: } p(\theta) = \frac{1}{2\pi I_0(\beta)} \exp(\beta \cos(\theta - \mu))$$

$$\mu \in [0, 2\pi), \beta \geq 0 : \text{パラメータ}$$

$$I_j(\beta) = \left(\frac{\beta}{2}\right)^j \sum_{i=0}^{\infty} \frac{1}{i! \Gamma(j+i+1)} \left(\frac{\beta}{2}\right)^{2i} \quad (j\text{-次第一種変形 Bessel 関数})$$

・ 逆変換 $\epsilon = g^{-1}(\xi, w)$ を ξ に関して微分

→ g を介さずに期待値の勾配を得られる.

・ w のサンプルを得るのには (g が計算できないので) 棄却サンプリングなどと使う必要がある.

・ 離散の確率変数に対する再パラメータ化勾配を用いる方法もある.

・ Gumbel ソフトマックス分布 (連続分布) で,

分布の温度パラメータ $\epsilon > 0$ に対することでカテゴリ分布に一致させることができる.

→ カテゴリ分布の Gumbel ソフトマックス分布による連続緩和ができる.

5.2.4 勾配近似による変分推論法.

再パラメータ化勾配を利用して, ミニバッチ $\mathcal{D}_s = \{(x_n, y_n)\}_{n \in \mathcal{S}}$ に対する

ELBO

$$\mathcal{L}_s(\xi)$$

$$:= \frac{N}{M} \sum_{n \in \mathcal{S}} \int q(\mathcal{W}; \xi) \log p(y_n | f(x_n; \mathcal{W})) d\mathcal{W} - D_{KL}[q(\mathcal{W}; \xi) \| p(\mathcal{W})]$$

を最大化する

→ 事後分布の近似精度を上げる.

$\epsilon \sim \mathcal{N}(0, I)$, $\mathcal{W} = g(\xi; \epsilon)$ と変数変換.

$q(\mathcal{W}; \xi) d\mathcal{W} = p(\epsilon) d\epsilon$ に注意して.

$$\mathcal{L}_s(\xi) = \frac{N}{M} \sum_{n \in \mathcal{S}} \int p(\epsilon) \log p(y_n | f(x_n; g(\xi; \epsilon))) d\epsilon - D_{KL}[q(\mathcal{W}; \xi) \| p(\mathcal{W})].$$

す="い雑下"けど計算効率的に
で 大抵なので仕方ない。複数サンプルで
してもちろんよい。

積分 \mathcal{E} , 1回のサンプル値 $\tilde{\mathcal{E}} \sim \mathcal{N}(\mathcal{E} | 0, I)$ で近似する.

↑ \tilde{y}_n, q_n, y_n などは $\mathcal{T} \rightarrow \mathcal{Y}$ がする

$$\mathcal{L}_d(\xi) \approx \mathcal{L}_{d, \tilde{\mathcal{E}}}(\xi)$$

$$= \frac{N}{M} \sum_{n \in \mathcal{D}} \log p(y_n | f(x_n; g(\xi; \tilde{\mathcal{E}}))) - D_{KL}[q(\mathcal{M}; \xi) \| p(\mathcal{M})]$$

$\mathcal{L}_{d, \tilde{\mathcal{E}}}(\xi)$ は $\mathcal{L}(\xi)$ の不偏推定量: $\mathbb{E}_{d, \tilde{\mathcal{E}}}[\mathcal{L}_{d, \tilde{\mathcal{E}}}(\xi)] = \mathcal{L}(\xi)$.

・計算すれば分かる.

・勾配を近似すると,

$$\nabla_{\xi} \mathcal{L}_d(\xi) \approx \nabla_{\xi} \mathcal{L}_{d, \tilde{\mathcal{E}}}(\xi)$$

$$= \frac{N}{M} \sum_{n \in \mathcal{D}} \nabla_{\xi} \log p(y_n | f(x_n; g(\xi; \tilde{\mathcal{E}}))) - \nabla_{\xi} D_{KL}[q(\mathcal{M}; \xi) \| p(\mathcal{M})]$$

↑ backpropagation で計算できる.

解析的に求まる.

よると, 以下のよう ξ を最適化する.

1. ミニバッチ \mathcal{D}_d を \mathcal{D} からとる.
2. M 個のノイズ $\tilde{\mathcal{E}}_i \sim \mathcal{N}(0, I)$ を取得.
3. $\nabla_{\xi} \mathcal{L}_{d, \tilde{\mathcal{E}}}(\xi)$ を計算.
4. $\xi \leftarrow \xi + \alpha \nabla_{\xi} \mathcal{L}_{d, \tilde{\mathcal{E}}}(\xi)$ と更新. 1へ.

5.2.5 期待値伝播法による学習.

・確率的逆伝播法 (probabilistic backpropagation method)

・順伝播: ネットワークを通して確率を伝播, 周辺尤度の評価.

・逆伝播: パラメータ学習のための周辺尤度の勾配計算.

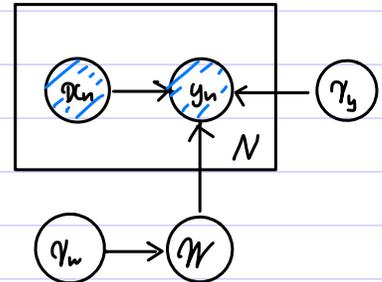
✓ データを逐次的に処理でき、大量データでの学習ができる。

✓ ハイパーパラメータの近似推論。

5.2.5.1 モデル

$y_n \in \mathbb{R}$: 1次元のラベルで考える。

モデル :



$$p(y, W, \gamma_y, \gamma_w | \mathcal{X}) = p(y | \mathcal{X}, W, \gamma_y) p(W | \gamma_w) p(\gamma_y) p(\gamma_w).$$

$$p(y | \mathcal{X}, W, \gamma_y) = \prod_{n=1}^N \mathcal{N}(y_n | f(x_n; W), \gamma_y^{-1})$$

$$p(\gamma_y) = \text{Gam}(\gamma_y | \alpha_{\gamma_y}, \beta_{\gamma_y}) \quad (\alpha_{\gamma_y}, \beta_{\gamma_y} : \text{ハイパーパラメータ})$$

$$p(W | \gamma_w) = \prod_{l=1}^L \prod_{i=1}^{H_l} \prod_{j=1}^{H_{l-1}} \mathcal{N}(w_{ij}^{(l)} | 0, \gamma_w^{-1}) \quad \leftarrow \text{全独立}$$

$$p(\gamma_w) = \text{Gam}(\gamma_w | \alpha_{\gamma_w}, \beta_{\gamma_w}) \quad (\alpha_{\gamma_w}, \beta_{\gamma_w} : \text{ハイパーパラメータ})$$

ここでは、NN $f(x_n; W)$ を次のように設定している :

$$z_n^{(l)} = x_n, \quad z_n^{(l)} = \text{ReLU}(a_n^{(l)}) \quad (l = 1, \dots, L-1)$$

$$a_n^{(l)} = \frac{1}{\sqrt{H_{l-1}}} W^{(l)} z_n^{(l-1)} \quad (l = 1, \dots, L)$$

← 係数調整

$$f(x_n; W) = z_n^{(L)} = a_n^{(L)}$$

事後分布

$$p(W, \gamma_y, \gamma_w | y, \mathcal{X})$$

$$\propto \frac{1}{p(y | \mathcal{X})} p(y, W, \gamma_y, \gamma_w | \mathcal{X})$$

$$\propto p(y | \mathcal{X}, W, \gamma_y) p(W | \gamma_w) p(\gamma_y) p(\gamma_w)$$

5.2.5.2 近似分布.

- 近似分布を次のとおり設定 (事前分布と同じ種類の分布を選択):

$$p(\mathcal{W}, \gamma_y, \gamma_w | \mathcal{Y}, \mathcal{X}) \approx q(\mathcal{W}, \gamma_y, \gamma_w),$$

$$q(\mathcal{W}, \gamma_y, \gamma_w) = \frac{\text{Gam}(\gamma_y | \alpha_{\gamma_y}, \beta_{\gamma_y})}{q(\gamma_y)} \frac{\text{Gam}(\gamma_w | \alpha_{\gamma_w}, \beta_{\gamma_w})}{q(\gamma_w)} \prod_{l=1}^L \prod_{i=1}^{H_l} \prod_{j=1}^{H_{l-1}} \mathcal{N}(w_{i,j}^{(l)} | m_{i,j}^{(l)}, v_{i,j}^{(l)}).$$

$q(w_{i,j}^{(l)})$
" $q(\gamma_w) = \prod_l \prod_i \prod_j q(w_{i,j}^{(l)})$ "

$$= q(\mathcal{W}) q(\gamma_y) q(\gamma_w).$$

5.2.5.3 初期化と事前分布因子の導入.

"一様"の意味.

$q(\mathcal{W}, \gamma_y, \gamma_w)$ が "無情報" になるように, パラメータを設定:

$$m_{i,j}^{(l)} = 0, v_{i,j}^{(l)} = \infty$$

$$\alpha_{\gamma_y} = 1, \beta_{\gamma_y} = 0$$

$$\alpha_{\gamma_w} = 1, \beta_{\gamma_w} = 0$$

この値

すくは更新されるので

お利便にしてください!!!

cf.) 無情報事前分布 (noninformative prior), 非正則事前分布 (improper prior)

- やること: 逐次学習

$$q_{i+1}(\theta) \approx \frac{1}{Z_{i+1}} f_{i+1}(\theta) q_i(\theta), \quad \theta = \{\mathcal{W}, \gamma_y, \gamma_w\}.$$

$f_{i+1}(\theta)$ として

$$p(\mathcal{W}, \gamma_y, \gamma_w | \mathcal{Y}, \mathcal{X})$$

$$\propto p(\mathcal{Y} | \mathcal{X}, \mathcal{W}, \gamma_y) p(\mathcal{W} | \gamma_w) p(\gamma_y) p(\gamma_w)$$

$$= \left(\prod_{n=1}^N \mathcal{N}(y_n | f(x_n; \mathcal{W}), \gamma_y^{-1}) \right) \left(\prod_{l=1}^L \prod_{i=1}^{H_l} \prod_{j=1}^{H_{l-1}} \mathcal{N}(w_{i,j}^{(l)} | 0, \gamma_w^{-1}) \right) p(\gamma_y) p(\gamma_w)$$

の各因子を次々と設定することによって $q(\theta)$ を更新していく.

事前分布因子

- ・ 残りの事前分布因子の情報 Σ を追加していく:

f とし $p(\gamma_w), p(\gamma_y), \mathcal{N}(w_{ij}^{(R)} | 0, \gamma_w^{-1})$ を設定して,

$q(\gamma_w), q(\gamma_y), q(\mathcal{W})$ のパラメータを更新.

- ・ 続いて, 尤度因子の情報 Σ を追加していく:

f とし $\mathcal{N}(y_n | f(x_n; \mathcal{W}), \gamma_y^{-1})$ を設定して

$q(\gamma_y), q(\mathcal{W})$ のパラメータを更新.

1. $p(\gamma_y)p(\gamma_w)$ の導入.

$q(\gamma_y)q(\gamma_w)$ を同じ形式にしたいので, q のパラメータを単純に

$$\alpha_{\gamma_y} \leftarrow \alpha_{\gamma_{y0}}, \quad \beta_{\gamma_y} \leftarrow \beta_{\gamma_{y0}}, \quad \alpha_{\gamma_w} \leftarrow \alpha_{\gamma_{w0}}, \quad \beta_{\gamma_w} \leftarrow \beta_{\gamma_{w0}}$$

と更新すればよい.

2. $\mathcal{N}(w_{ij}^{(R)} | 0, \gamma_w^{-1})$ の導入.

$$q_{\text{new}}(w_{ij}^{(R)}) q_{\text{new}}(\gamma_w)$$

$$\approx \frac{1}{Z_0} \mathcal{N}(w_{ij}^{(R)} | 0, \gamma_w^{-1}) q(w_{ij}^{(R)}) q(\gamma_w)$$

$$= \frac{1}{Z_0} \mathcal{N}(w_{ij}^{(R)} | 0, \gamma_w^{-1}) \mathcal{N}(w_{ij}^{(R)} | m_{ij}^{(R)}, v_{ij}^{(R)}) \text{Gam}(\gamma_w | \alpha_{\gamma_w}, \beta_{\gamma_w}),$$

$$Z_0 = Z(\alpha_{\gamma_w}, \beta_{\gamma_w})$$

$$= \int \mathcal{N}(w_{ij}^{(R)} | 0, \gamma_w^{-1}) \mathcal{N}(w_{ij}^{(R)} | m_{ij}^{(R)}, v_{ij}^{(R)}) \text{Gam}(\gamma_w | \alpha_{\gamma_w}, \beta_{\gamma_w}) dw_{ij}^{(R)} d\gamma_w$$

◦ $m_{ij}^{(l)}, v_{ij}^{(l)}$ の更新. §4.2.4.3 の結果より,

$$m_{ij}^{(l)} \leftarrow m_{ij}^{(l-1)} + v_{ij}^{(l-1)} \frac{\partial}{\partial m_{ij}^{(l-1)}} \log Z_0$$

$$v_{ij}^{(l)} \leftarrow v_{ij}^{(l-1)} - v_{ij}^{(l-1)2} \left(\left(\frac{\partial}{\partial m_{ij}^{(l-1)}} \log Z_0 \right)^2 - 2 \frac{\partial}{\partial v_{ij}^{(l-1)}} \log Z_0 \right).$$

◦ $\alpha_{\gamma_w}, \beta_{\gamma_w}$ の更新. §4.2.4.4 の結果より,

$$\alpha_{\gamma_w} \leftarrow \left(Z_0 Z_2 Z_1^{-1} \frac{\alpha_{\gamma_w} + 1}{\alpha_{\gamma_w}} - 1 \right)^{-1} \quad Z_1 = Z(\alpha_{\gamma_w} + 1, \beta_{\gamma_w}),$$

$$\beta_{\gamma_w} \leftarrow \left(Z_2 Z_1^{-1} \frac{\alpha_{\gamma_w} + 1}{\beta_{\gamma_w}} - Z_1 Z_0^{-1} \frac{\alpha_{\gamma_w}}{\beta_{\gamma_w}} \right)^{-1} \quad Z_2 = Z(\alpha_{\gamma_w} + 2, \beta_{\gamma_w}).$$

$Z(\alpha_{\gamma_w}, \beta_{\gamma_w})$ を求める必要はないから,

厳密解は得られないので"近似"できる.

$$Z(\alpha_{\gamma_w}, \beta_{\gamma_w})$$

$$= \int \mathcal{N}(w_{ij}^{(l)} | 0, \gamma_w^{-1}) \mathcal{N}(w_{ij}^{(l)} | m_{ij}^{(l-1)}, v_{ij}^{(l-1)}) \text{Gam}(\gamma_w | \alpha_{\gamma_w}, \beta_{\gamma_w}) dw_{ij}^{(l)} d\gamma_w$$

$$= \int \left(\int \mathcal{N}(w_{ij}^{(l)} | 0, \gamma_w^{-1}) \text{Gam}(\gamma_w | \alpha_{\gamma_w}, \beta_{\gamma_w}) d\gamma_w \right) \mathcal{N}(w_{ij}^{(l)} | m_{ij}^{(l-1)}, v_{ij}^{(l-1)}) dw_{ij}^{(l)}$$

$$= \int \text{St}(w_{ij}^{(l)} | 0, \frac{\alpha_{\gamma_w}}{\beta_{\gamma_w}}, 2\alpha_{\gamma_w}) \mathcal{N}(w_{ij}^{(l)} | m_{ij}^{(l-1)}, v_{ij}^{(l-1)}) dw_{ij}^{(l)}$$

↓ §3.2.4.4 の Student 分布

↓ Student の t 分布を
平均と分散の等しい Gauss 分布で"近似"

$$\approx \int \mathcal{N}(w_{ij}^{(l)} | 0, \frac{\beta_{\gamma_w}}{\alpha_{\gamma_w} - 1}) \mathcal{N}(w_{ij}^{(l)} | m_{ij}^{(l-1)}, v_{ij}^{(l-1)}) dw_{ij}^{(l)}.$$

" = z ",

$$\int \mathcal{N}(w | 0, \lambda_1^{-1}) \mathcal{N}(w | m, \lambda_2^{-1}) dw$$

$$= \int \frac{\sqrt{\lambda_2}}{\sqrt{2\pi}} \exp\left(-\frac{\lambda_2}{2} w^2\right) \frac{\sqrt{\lambda_1}}{\sqrt{2\pi}} \exp\left(-\frac{\lambda_1}{2} (w-m)^2\right) dw$$

$$= \frac{\sqrt{\lambda_2}}{\sqrt{2\pi}} \frac{\sqrt{\lambda_1}}{\sqrt{2\pi}} \int \exp\left(-\frac{1}{2} \left((\lambda_1 + \lambda_2) w^2 - 2\lambda_1 m w + \lambda_1 m^2 \right)\right) dw \quad \downarrow \text{平方完成}$$

$$= \frac{\sqrt{\lambda_2}}{\sqrt{2\pi}} \frac{\sqrt{\lambda_1}}{\sqrt{2\pi}} \int \exp\left(-\frac{\lambda_1 + \lambda_2}{2} \left(w - \frac{\lambda_1}{\lambda_1 + \lambda_2} m \right)^2\right) \exp\left(-\frac{1}{2} \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2} m^2\right) dw$$

$$\begin{aligned}
&= \frac{\sqrt{\lambda_2}}{\sqrt{2\pi}} \frac{\sqrt{\lambda_2}}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{\lambda_2 \lambda_2}{\lambda_2 + \lambda_2} m^2\right) \int \exp\left(-\frac{\lambda_2 + \lambda_2}{2} \left(w - \frac{\lambda_2}{\lambda_2 + \lambda_2} m\right)^2\right) dw \\
&= \frac{1}{\sqrt{2\pi}} \sqrt{\frac{\lambda_2 \lambda_2}{\lambda_2 + \lambda_2}} \exp\left(-\frac{1}{2} \frac{\lambda_2 \lambda_2}{\lambda_2 + \lambda_2} m^2\right) \underbrace{= \frac{\sqrt{2\pi}}{\lambda_2 + \lambda_2}} \\
&= \mathcal{N}\left(m \mid 0, \left(\frac{\lambda_2 \lambda_2}{\lambda_2 + \lambda_2}\right)^{-1}\right)
\end{aligned}$$

7.2.2".

$$Z(\alpha_{\gamma_w}, \beta_{\gamma_w}) \approx \mathcal{N}\left(m_{ij}^{(k)} \mid 0, \frac{\beta_{\gamma_w}}{\alpha_{\gamma_w} - 1} + V_{ij}^{(k)}\right).$$

$\frac{\partial}{\partial m_{ij}^{(k)}} \log Z_0$, $\frac{\partial}{\partial v_{ij}^{(k)}} \log Z_0$, Z_1 , Z_2 はこれ"計算"する.

5.2.5.4 尤度因子の導入.

3. $\mathcal{N}(y_n \mid f(\alpha_n; \mathcal{W}), \gamma_y^{-1})$ の導入.

$$\begin{aligned}
& q_{\text{new}}(\mathcal{W}) q_{\text{new}}(\gamma_y) \\
& \approx \frac{1}{Z} \mathcal{N}(y_n \mid f(\alpha_n; \mathcal{W}), \gamma_y^{-1}) q(\mathcal{W}) q(\gamma_y) \\
& = \frac{1}{Z} \mathcal{N}(y_n \mid f(\alpha_n; \mathcal{W}), \gamma_y^{-1}) q(\mathcal{W}) \text{Gam}(\gamma_y \mid \alpha_{\gamma_y}, \beta_{\gamma_y}). \\
& = \frac{1}{Z} \mathcal{N}(y_n \mid f(\alpha_n; \mathcal{W}), \gamma_y^{-1}) \left(\prod_{l=1}^L \prod_{i=1}^{H_l} \prod_{j=1}^{H_{l-1}} q(w_{ij}^{(l)}) \right) \text{Gam}(\gamma_y \mid \alpha_{\gamma_y}, \beta_{\gamma_y}).
\end{aligned}$$

$$Z = Z(\alpha_{\gamma_y}, \beta_{\gamma_y})$$

$$= \int \mathcal{N}(y_n \mid f(\alpha_n; \mathcal{W}), \gamma_y^{-1}) q(\mathcal{W}) \text{Gam}(\gamma_y \mid \alpha_{\gamma_y}, \beta_{\gamma_y}) d\mathcal{W} d\gamma_y$$

• $m_{ij}^{(k)}$, $v_{ij}^{(k)}$ の更新. §4.2.4.3 の結果より.

$$m_{ij}^{(k)} \leftarrow m_{ij}^{(k)} + v_{ij}^{(k)} \frac{\partial}{\partial m_{ij}^{(k)}} \log Z$$

$$v_{ij}^{(k)} \leftarrow v_{ij}^{(k)} - v_{ij}^{(k)2} \left(\left(\frac{\partial}{\partial m_{ij}^{(k)}} \log Z \right)^2 - 2 \frac{\partial}{\partial v_{ij}^{(k)}} \log Z \right).$$

• $\alpha_{\gamma_y}, \beta_{\gamma_y}$ の更新. §4.2.4.4 の結果から,

$$\alpha_{\gamma_y} \leftarrow \left(\sum \sum_2 \sum_1^{-1} \frac{\alpha_{\gamma_y} + 1}{\alpha_{\gamma_y}} - 1 \right)^{-1} \quad \sum_1 = \sum(\alpha_{\gamma_y} + 1, \beta_{\gamma_y}),$$

$$\beta_{\gamma_y} \leftarrow \left(\sum_2 \sum_1^{-1} \frac{\alpha_{\gamma_y} + 1}{\beta_{\gamma_y}} - \sum_1 \sum^{-1} \frac{\alpha_{\gamma_y}}{\beta_{\gamma_y}} \right)^{-1} \quad \sum_2 = \sum(\alpha_{\gamma_y} + 2, \beta_{\gamma_y}).$$

→ $\sum = \sum(\alpha_{\gamma_y}, \beta_{\gamma_y})$ が求まればパラメータの更新が可能.

$\mathcal{N} \sim q(\mathcal{N})$ のとき, $z^{(L)} = f(x_n; \mathcal{N})$ の分布は複雑になる.

$z^{(L)} \sim \mathcal{N}(z^{(L)} | m_{z^{(L)}}, \sigma_{z^{(L)}}^2)$ で近似したい.
← 中心極限定理.

つまり, $q(\mathcal{N}) d\mathcal{N} = p(z^{(L)}) dz^{(L)} \approx \mathcal{N}(z^{(L)} | m_{z^{(L)}}, \sigma_{z^{(L)}}^2) dz^{(L)}$ とする.

Remark 式 $q(\mathcal{N}) d\mathcal{N} = p(z^{(L)}) dz^{(L)}$ は, 変数変換の公式そのものではない.

(\because 一般には \mathcal{N} と $z^{(L)}$ の次元が違うので, $z^{(L)} = f(x_n; \mathcal{N})$ は全単射にならない)