

§3.5 線形回帰の例.

線形回帰モデルを、Gauss分布を使って構築.

3.5.1 モデルの構築.

$x_n \in \mathbb{R}^M$: 入力値, $w \in \mathbb{R}^M$: パラメータ,

$\varepsilon_n \sim \mathcal{N}(\varepsilon_n | 0, \lambda^{-1})$: サイズ. ($\lambda > 0$: given)

出力値 $y_n \in \mathbb{R}$ を次のようにモデル化:

$$y_n = w^T x_n + \varepsilon_n \quad \leftarrow w^T x_n \text{ は } \varepsilon_n \text{ たとえば } \mathcal{N}(0, \lambda^{-1}).$$

$\Rightarrow y_n$ の確率分布は,

$$p(y_n | x_n, w) = \mathcal{N}(y_n | w^T x_n, \lambda^{-1})$$

と表せる.

- やりたいこと: w の学習.

w について、次のような事前分布を仮定する:

$$p(w) = \mathcal{N}(w | m, \Lambda^{-1}),$$

$$m \in \mathbb{R}^M, \Lambda \in M_M(\mathbb{R}), \Lambda > 0 : \text{ハイパーパラメータ}.$$

\rightarrow 分布の組みあわせ、構成する分布の種類、ハイパーパラメータの値の設定により、一つの確率モデル（データの生成過程についての一つの仮説）を明記できる。

モデルが“明記できると”

サンプリングにより、仮説 ω カバーしていき実現例を視覚化できる。

→ モデルの特性・妥当性をチェックできる。

3.5.2 事後分布と予測分布の計算.

- データ $\mathcal{Y} := \{y_1, \dots, y_N\}$, $\mathcal{X} := \{x_1, \dots, x_N\}$ を観測後の事後分布を求める。

Bayes' Thm. 5'),

$$\begin{aligned} p(w|\mathcal{Y}, \mathcal{X}) &= \frac{p(\mathcal{Y}|\mathcal{X}, w)p(w)}{p(\mathcal{Y}|\mathcal{X})} \\ &\propto p(\mathcal{Y}|\mathcal{X}, w)p(w) \\ &= \left(\prod_{n=1}^N p(y_n|x_n, w) \right) p(w) \\ &= \left(\prod_{n=1}^N \mathcal{N}(y_n|w^T x_n, \lambda^{-1}) \right) \mathcal{N}(w|m, \Lambda^{-1}) \end{aligned}$$

$$\begin{aligned} \log p(w|\mathcal{Y}, \mathcal{X}) &= \sum_{n=1}^N \log \mathcal{N}(y_n|w^T x_n, \lambda^{-1}) + \log \mathcal{N}(w|m, \Lambda^{-1}) + \text{const.} \\ &= \sum_{n=1}^N \left(-\frac{1}{2} \left(\lambda (y_n - w^T x_n)^2 - \underline{\log \lambda} \right) \right) \\ &\quad - \frac{1}{2} \left((w - m)^T \Lambda (w - m) - \underline{\log (\det \Lambda)} \right) + \text{const.} \end{aligned}$$

$$\begin{aligned}
&= -\frac{\lambda}{2} \sum_{n=1}^N (y_n - w^\top x_n)^2 - \frac{1}{2} (w - m)^\top \Lambda (w - m) + \text{const.} \\
&= -\frac{1}{2} \left(\lambda \sum_{n=1}^N \underbrace{(y_n^2 - 2w^\top x_n y_n + w^\top x_n x_n^\top w)}_{\text{const.}} \right. \\
&\quad \left. + w^\top \Lambda w - 2w^\top \Lambda m + \underbrace{m^\top \Lambda m}_{\text{const.}} \right) + \text{const.} \\
&= -\frac{1}{2} \left(w^\top \left(\lambda \sum_{n=1}^N x_n x_n^\top + \Lambda \right) w - 2w^\top \left(\lambda \sum_{n=1}^N x_n y_n + \Lambda m \right) \right) + \text{const.}
\end{aligned}$$

↑ M次元Gauss分布のpdfのlogΣとTを

$$\begin{aligned}
&\therefore p(w | y, x) = \mathcal{N}(w | \hat{m}, \hat{\Lambda}^{-1}), \\
&\hat{\Lambda} := \lambda \sum_{n=1}^N x_n x_n^\top + \Lambda, \quad \hat{m} := \hat{\Lambda}^{-1} \left(\lambda \sum_{n=1}^N x_n y_n + \Lambda m \right).
\end{aligned}$$

- 新規入力値 x_* : given のときの出力値 y_* の予測分布を求める。

Bayes' Thm. より,

$$p(w | y_*, x_*) = \frac{p(y_* | x_*, w) p(w)}{p(y_* | x_*)}.$$

$$\therefore \log p(y_* | x_*)$$

$$= \log p(y_* | x_*, w) - \log p(w | y_*, x_*) + \text{const.}$$

上の事後分布の結果を用いると、

$$p(w | y_*, x_*) = \mathcal{N}(w | m(y_*), (\lambda x_* x_*^\top + \Lambda)^{-1}),$$

$$m(y_*) := (\lambda x_* x_*^\top + \Lambda)^{-1} (\lambda x_* y_* + \Lambda m).$$

$$\therefore \log p(w | y_*, x_*)$$

$$= \log \mathcal{N}(w | m(y_*), (\lambda x_* x_*^T + \Lambda)^{-1})$$

$$= -\frac{1}{2} ((w - m(y_*))^T (\lambda x_* x_*^T + \Lambda) (w - m(y_*)))$$

$$-\underbrace{\log(\det(\lambda x_* x_*^T + \Lambda))}_{\text{const.}} + \text{const.}$$

$$= -\frac{1}{2} (w - m(y_*))^T (\lambda x_* x_*^T + \Lambda) (w - m(y_*)) + \text{const.}$$

$$= -\frac{1}{2} \left(\underbrace{w^T (\lambda x_* x_*^T + \Lambda) w}_{\text{const.}} - w^T (\lambda x_* x_*^T + \Lambda) m(y_*) \right.$$

$$\left. - m(y_*)^T (\lambda x_* x_*^T + \Lambda) w + m(y_*)^T (\lambda x_* x_*^T + \Lambda) m(y_*) \right) + \text{const}$$

$$= -\frac{1}{2} \left(m(y_*)^T (\lambda x_* x_*^T + \Lambda) m(y_*) \xrightarrow{(\lambda x_* x_*^T + \Lambda)^T = \lambda x_* x_*^T + \Lambda} \right. \\ \left. - 2 w^T (\lambda x_* x_*^T + \Lambda) m(y_*) \right) + \text{const.}$$

$$= -\frac{1}{2} \left((\lambda x_* y_* + \Lambda m)^T (\lambda x_* x_*^T + \Lambda)^{-1} (\lambda x_* y_* + \Lambda m) \right. \\ \left. - 2 w^T (\lambda x_* y_* + \Lambda m) \right) + \text{const.}$$

$$= -\frac{1}{2} \left(\lambda^2 x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} x_* y_*^2 + 2 \lambda x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} \Lambda m y_* \right. \\ \left. + \underbrace{m^T \Lambda (\lambda x_* x_*^T + \Lambda)^{-1} \Lambda m}_{\text{const.}} - 2 \lambda w^T x_* y_* \right) + \text{const.}$$

$$= -\frac{1}{2} \left(\lambda^2 x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} x_* y_*^2 \right. \\ \left. + 2 \lambda \left(x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} \Lambda m - w^T x_* \right) y_* \right) + \text{const.}$$

また、

$$\begin{aligned}
& \log p(y_* | x_*, w) \\
&= \log \mathcal{N}(y_* | w^T x_*, \lambda^{-1}) \\
&= -\frac{1}{2} \left(\lambda (y_* - w^T x_*)^2 - \underbrace{\log \lambda}_{\text{const.}} \right) + \text{const.} \\
&= -\frac{1}{2} \lambda \left(y_*^2 - 2w^T x_* y_* + \underbrace{(w^T x_*)^2}_{\text{const.}} \right) + \text{const.} \\
&= -\frac{1}{2} \lambda (y_*^2 - 2w^T x_* y_*) + \text{const.} \\
&\therefore \log p(y_* | x_*) \\
&= -\frac{1}{2} \lambda (y_*^2 - 2w^T x_* y_*) \\
&\quad + \frac{1}{2} \left(\lambda^2 x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} x_* y_*^2 \right. \\
&\quad \left. + 2\lambda \left(x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} \Lambda m - w^T x_* \right) y_* \right) + \text{const.} \\
&= -\frac{1}{2} \left((\lambda - \lambda^2 x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} x_*) y_*^2 \right. \\
&\quad \left. - 2\lambda x_*^T (\cancel{w} + (\lambda x_* x_*^T + \Lambda)^{-1} \Lambda m - \cancel{w}) y_* \right) + \text{const.} \\
&= -\frac{1}{2} \left((\lambda - \lambda^2 x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} x_*) y_*^2 \right. \\
&\quad \left. - 2\lambda x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} \Lambda m y_* \right) + \text{const.} \\
&\quad \uparrow \text{Gauss分布のpdfをlogで取る} \\
&\therefore p(y_* | x_*) = \mathcal{N}(y_* | \mu_*, \lambda_*^{-1}), \\
&\lambda_* := \lambda - \lambda^2 x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} x_*, \\
&\mu_* := \lambda_*^{-1} \lambda x_*^T (\lambda x_* x_*^T + \Lambda)^{-1} \Lambda m.
\end{aligned}$$

もう少し計算をすらわる。

Sherman-Morrison(レミー) ,

$$(\Lambda + \lambda x_* x_*^\top)^{-1} = \Lambda^{-1} - (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} \lambda \Lambda^{-1} x_* x_*^\top \Lambda^{-1}.$$

$$\therefore \lambda_*$$

$$\begin{aligned} &= \lambda - \lambda^2 x_*^\top (\Lambda^{-1} - (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} \lambda \Lambda^{-1} x_* x_*^\top \Lambda^{-1}) x_* \\ &= \lambda - \lambda^2 x_*^\top \Lambda^{-1} x_* + \lambda^3 (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} (x_*^\top \Lambda^{-1} x_*)^2 \\ &= \frac{\lambda + \lambda^2 x_*^\top \Lambda^{-1} x_* - \lambda^2 x_*^\top \Lambda^{-1} x_* - \lambda^3 (x_*^\top \Lambda^{-1} x_*)^2 + \lambda^3 (x_*^\top \Lambda^{-1} x_*)^2}{1 + \lambda x_*^\top \Lambda^{-1} x_*} \\ &= \lambda (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} \end{aligned}$$

$$\lambda_*^{-1} = \lambda^{-1} (1 + \lambda x_*^\top \Lambda^{-1} x_*) = \lambda^{-1} + x_*^\top \Lambda^{-1} x_*.$$

$$\mu_*$$

$$\begin{aligned} &= \lambda_*^{-1} \lambda x_*^\top (\Lambda^{-1} - (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} \lambda \Lambda^{-1} x_* x_*^\top \Lambda^{-1}) \Lambda m \\ &= \lambda_*^{-1} \lambda (1 - \lambda (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} x_*^\top \Lambda^{-1} x_*) x_*^\top m \\ &= \lambda_*^{-1} \lambda (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} (1 + \lambda x_*^\top \Lambda^{-1} x_* - \lambda x_*^\top \Lambda^{-1} x_*) m^\top x_* \\ &= \lambda_*^{-1} \lambda (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} m^\top x_* \\ &= \lambda^{-1} (1 + \lambda x_*^\top \Lambda^{-1} x_*) \lambda (1 + \lambda x_*^\top \Lambda^{-1} x_*)^{-1} m^\top x_* \\ &= m^\top x_* . \end{aligned}$$

結局、

$$p(y_* | \mathbf{x}_*) = \mathcal{N}(y_* | \mu_*, \lambda_*^{-1}),$$

$$\mu_* = \mathbf{m}^\top \mathbf{x}_*, \quad \lambda_*^{-1} = \lambda (1 + \lambda \mathbf{x}_*^\top \Lambda^{-1} \mathbf{x}_*)^{-1}.$$

- N 個のデータを観測したあとの予測分布は、上の式で
 $\mathbf{m} \leftarrow \hat{\mathbf{m}}$ は、 $\Lambda \leftarrow \hat{\Lambda}$ に直すのがいい。

3.5.3 モデルの比較.

- **モデル選択**：データセット \mathcal{D} に対して、複数のモデルの良さを比較する。
 - 予測分布の可視化をするのも一案。可視化できない場合は…?
- Bayes推論では、**周辺尤度（モデルエビデンス）** $p(\mathcal{D})$ を
 比較するのが一般的。
↑ モデルが生成する尤もしくは。
- 線形回帰モデルの例では。

$$p(\mathcal{D}) = p(y | \mathcal{X})$$

$$= \frac{p(w) \prod_{n=1}^N p(y_n | \mathbf{x}_n, w)}{p(w | y, \mathcal{X})}.$$

← これは常に正しいです。

対数とく、

$$\log p(y | \mathcal{X})$$

$$= \log p(w) + \sum_{n=1}^N \log p(y_n | \mathbf{x}_n, w) - \log p(w | y, \mathcal{X})$$

$$\begin{aligned}
&= \log \mathcal{N}(\mathbf{w} | \mathbf{m}, \Lambda^{-1}) + \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{w}^T \mathbf{x}_n, \lambda^{-1}) \\
&\quad - \log \mathcal{N}(\mathbf{w} | \hat{\mathbf{m}}, \hat{\Lambda}^{-1}) \\
&= -\frac{1}{2} \left((\mathbf{w} - \mathbf{m})^T \Lambda (\mathbf{w} - \mathbf{m}) - \log(\det \Lambda) + M \log 2\pi \right) \\
&\quad - \frac{1}{2} \sum_{n=1}^N \left(\lambda (y_n - \mathbf{w}^T \mathbf{x}_n)^2 - \log \lambda + \log 2\pi \right) \\
&\quad + \frac{1}{2} \left((\mathbf{w} - \hat{\mathbf{m}})^T \hat{\Lambda} (\mathbf{w} - \hat{\mathbf{m}}) - \log(\det \hat{\Lambda}) + M \log 2\pi \right) \\
&= -\frac{1}{2} \left(\cancel{\mathbf{w}^T \Lambda \mathbf{w}} - \cancel{2\mathbf{w}^T \Lambda \mathbf{m}} + \mathbf{m}^T \Lambda \mathbf{m} - \log(\det \Lambda) \right. \\
&\quad \left. + \lambda \sum_{n=1}^N (y_n^2 - \cancel{2\mathbf{w}^T \Lambda \mathbf{x}_n y_n} + \cancel{(\mathbf{w}^T \mathbf{x}_n)^2}) - N \log \lambda + N \log 2\pi \right. \\
&\quad \left. - \cancel{\mathbf{w}^T \hat{\Lambda} \mathbf{w}} + \cancel{2\mathbf{w}^T \hat{\Lambda} \hat{\mathbf{m}}} - \hat{\mathbf{m}}^T \hat{\Lambda} \hat{\mathbf{m}} + \log(\det \hat{\Lambda}) \right) \\
&= -\frac{1}{2} \left(\lambda \sum_{n=1}^N y_n^2 - N \log \lambda + N \log 2\pi \right. \\
&\quad \left. + \mathbf{m}^T \Lambda \mathbf{m} - \log(\det \Lambda) - \hat{\mathbf{m}}^T \hat{\Lambda} \hat{\mathbf{m}} + \log(\det \hat{\Lambda}) \right).
\end{aligned}$$

$$\rightarrow M^* := \underset{M}{\operatorname{argmax}} \log p(y | x)$$

とすれば

- モデルエビデンスはデータ数 N に依存することに注意。