

ベイズ線形回帰

岡山大学 異分野基礎科学研究所

大槻純也



前回の復習

ベイズの定理

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w})p(\mathbf{w})$$

事後確率分布

尤度関数

事前確率分布

Training set

$$\mathbf{x} = \{x_1, x_2, \dots, x_N\}$$

$$\mathbf{t} = \{t_1, t_2, \dots, t_N\}$$

Maximum posterior (MAP) 推定

$p(\mathbf{w}) = \text{const}$ なら **最尤推定**

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{x}, \mathbf{t})$$

→ (正則化付き)

最小二乗法に一致

MAP推定は...

パラメータ \mathbf{w} の分布を無視する近似

物理の言葉を使うと

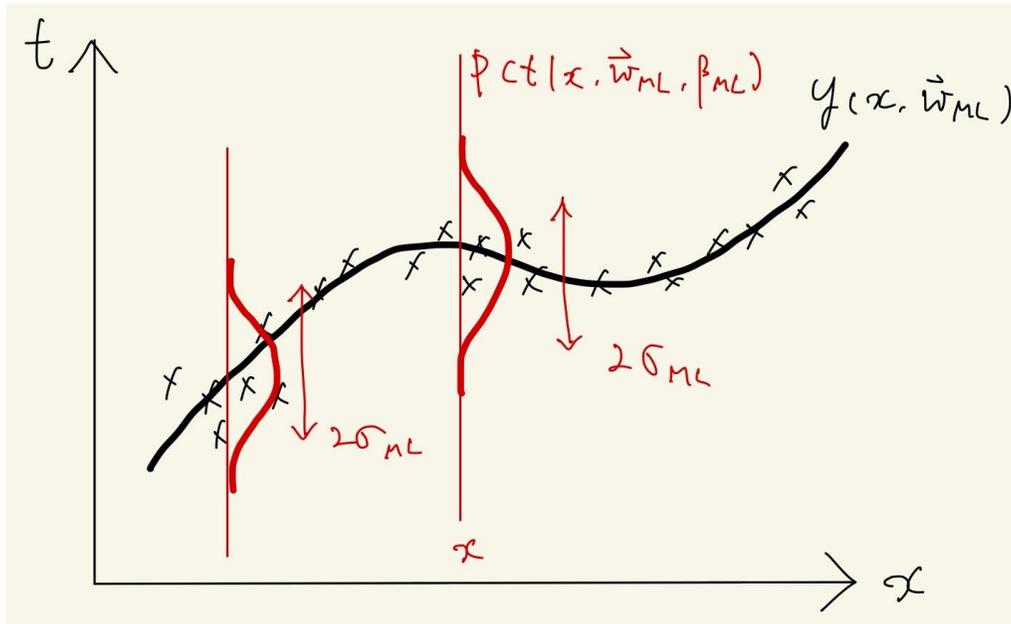
- パラメータ \mathbf{w} の **ゆらぎを無視**する近似
- 平均場近似、鞍点近似

予測分布 (Predictive distribution)

最尤推定の予測分布 (Predictive distribution)

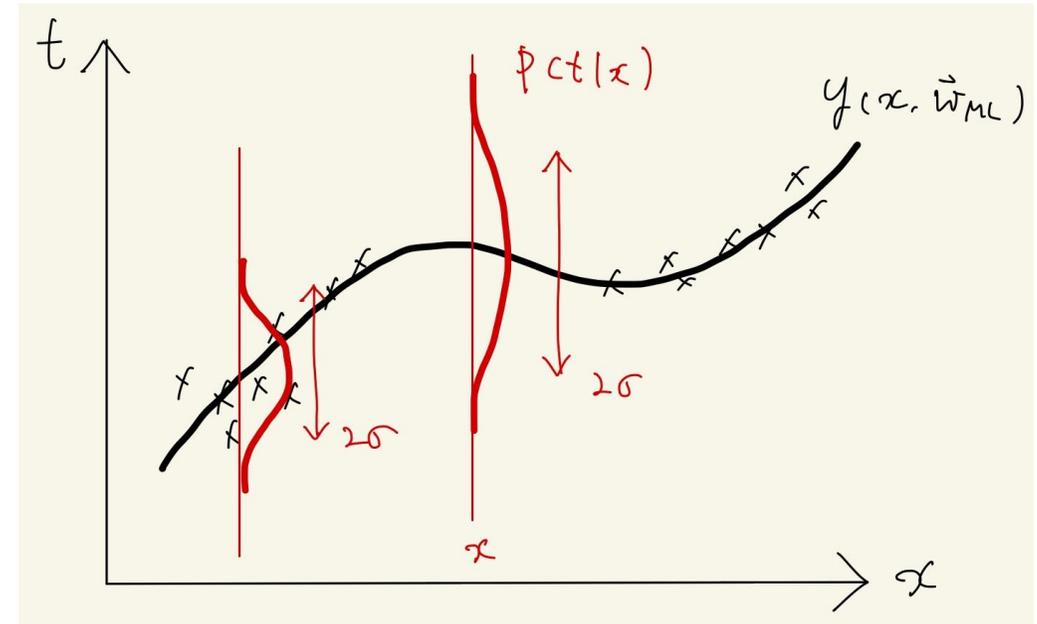
$$p(t|x, \mathbf{w}_{ML}, \beta_{ML}) = \mathcal{N}(t|y(x, \mathbf{w}_{ML}), \beta_{ML}^{-1})$$

しかしこのように、
データ点が一様に存在しない場合



σ は x に依らない

$$\sigma = \sigma_{ML} = \beta_{ML}^{-1/2}$$



σ は x に依存すべき

$$\sigma = \sigma(x)$$

やりたいこと：ベイズ線形回帰

パラメータ w の分布を考慮に入れる

= 最尤推定値 w_{ML} まわりのゆらぎ

w の事後確率分布

$$p(w|x, t) \propto p(t|x, w)p(w)$$

ベイズ統計に基づく予測分布 (Predictive distribution)

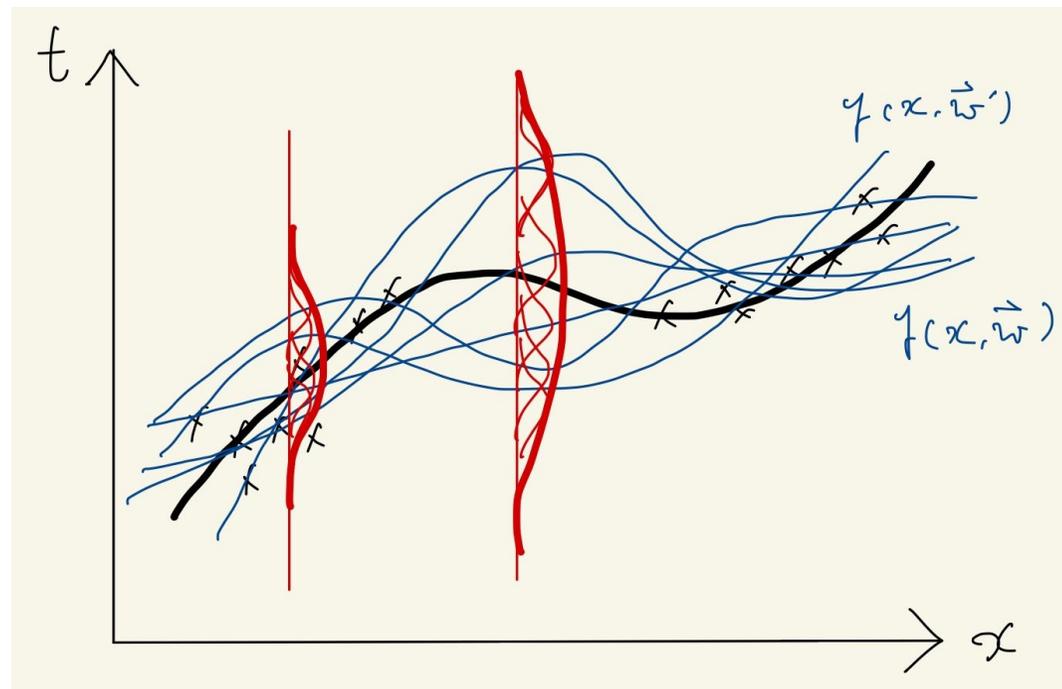
$$p(t|x, \mathbf{x}, t) = \int dw p(t|x, w)p(w|x, t)$$

sum rule

パラメータ w の事後確率分布

パラメータ w が与えられたときの予測分布

座標 x に依存した予測の不確定性を記述する



必要な計算

予測分布 (Predictive distribution)

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int d\mathbf{w} p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t})$$

$$p(t|x, \mathbf{w}) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) \propto \exp \left[-\frac{\beta}{2} (t - \mathbf{w}^T \phi(x))^2 \right]$$

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}) = \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1}) \prod_{j=1}^M \mathcal{N}(w_j|0, \alpha^{-1}) \\ \propto \exp \left[-\frac{\beta}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right]$$

Bayes' theorem

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}) p(\mathbf{w})$$

Gaussianの積の積分もまたGaussianになる

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x)) \quad \text{分散が } x \text{ に依存する}$$

これを導出するには
多変量ガウス分布 (multivariate Gaussian)
の関係式が必要

多変量ガウス分布 (multivariate Gaussian distribution)

1次元のガウス分布

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$

μ : mean (平均)
 σ^2 : variance (分散)

D次元のガウス分布

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$

$\boldsymbol{\mu}$: mean (平均)
 $\boldsymbol{\Sigma}$: covariance matrix (共分散行列)
 $|\boldsymbol{\Sigma}|$: determinant (行列式)

$$\mathbb{E}[\mathbf{x}] = \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathbf{x} d\mathbf{x} = \boldsymbol{\mu} \quad \text{mean}$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathbf{x}\mathbf{x}^T d\mathbf{x} = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$$

$$\text{cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] = \boldsymbol{\Sigma} \quad \text{covariance}$$

Covariance (共分散)

Covariance matrix Σ の性質

- ◆ Real symmetric (対称行列)

$$\Sigma_{ij} = \Sigma_{ji} \quad \text{よって固有値は実数}$$

- ◆ Positive definite (正定値)
全ての固有値が正

$$\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i, \quad \lambda_i > 0$$

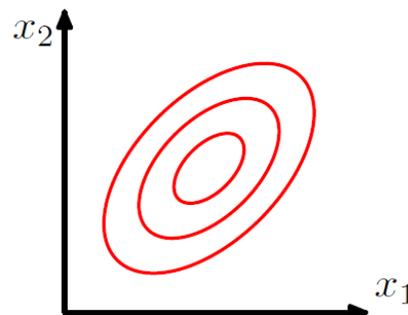
なぜなら

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \prod_{i=1}^D \frac{1}{(2\pi\lambda_i)^{D/2}} \exp\left(-\frac{y_i^2}{2\lambda_i}\right)$$

$$y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})$$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12} & \Sigma_{22} \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{pmatrix} \quad \Sigma \propto \mathbf{I}$$

$$\Sigma_{12} > 0$$

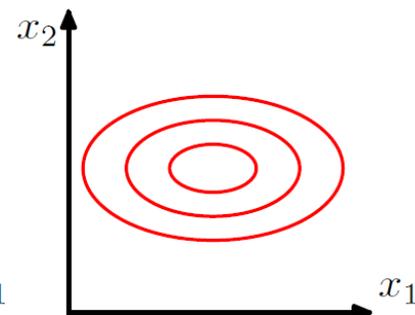


(a)

正の相関

x_1 と x_2 が近い値
を取りやすい

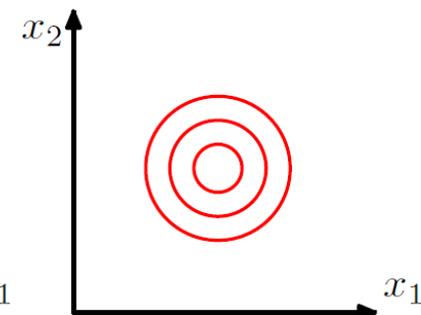
$$\Sigma_{11} > \Sigma_{22}$$



(b)

相関なし

x_1 と x_2 の分布の
幅が異なる



(c)

相関なし

x_1 と x_2 が同じ分
布し従う

PRML Fig. 2.8

準備

D 次元ガウス分布を考え、
 \mathbf{x} を M 成分と $D - M$ 成分に分割する

$$p(\mathbf{x}_a, \mathbf{x}_b) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

$\boldsymbol{\Sigma}$ は対称行列なので $\boldsymbol{\Sigma}_{ab} = \boldsymbol{\Sigma}_{ba}^T$

Precision matrix $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$ を定義

$$\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

$\boldsymbol{\Sigma}$ と $\boldsymbol{\Lambda}$ の関係

$$\boldsymbol{\Lambda}_{aa} = (\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba})^{-1}$$

$$\boldsymbol{\Lambda}_{ab} = -(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba})^{-1} \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1}$$

計算する量によって $\boldsymbol{\Sigma}$ と $\boldsymbol{\Lambda}$ のどちらで表記
がシンプルになる場合がある。

Marginal Gaussian distribution

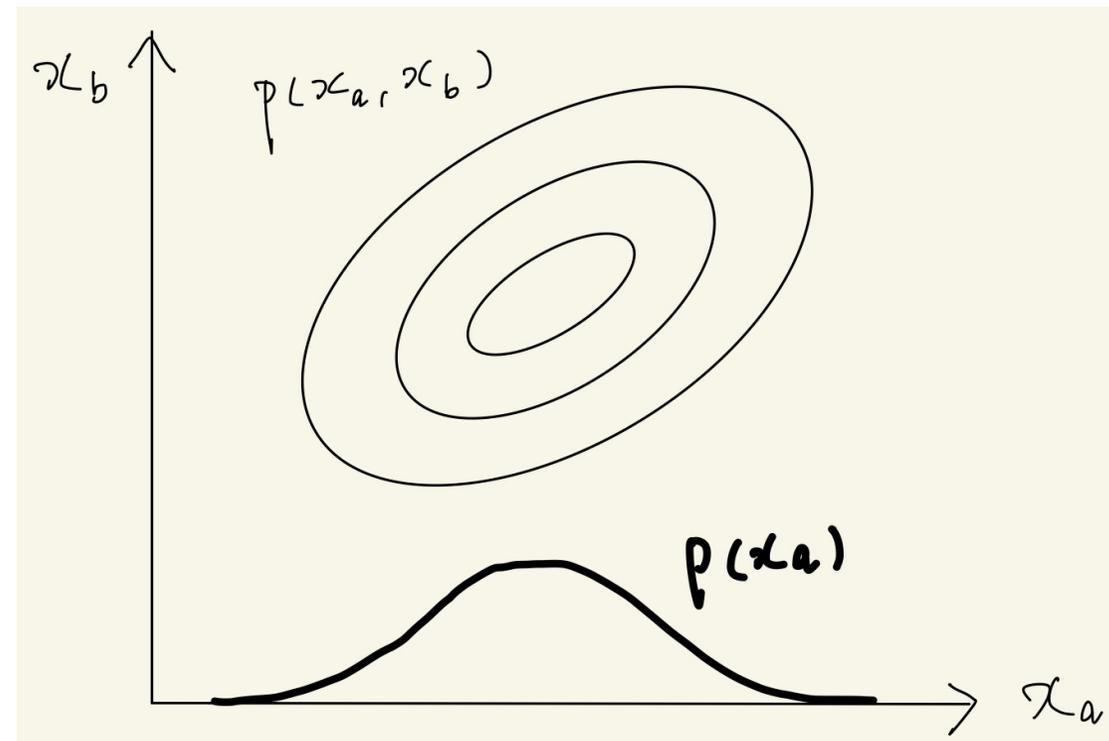
変数 \mathbf{x}_b について周辺化（積分）する

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$$

結果は、 $\boldsymbol{\mu}$ や $\boldsymbol{\Sigma}$ から aa 成分のみを取り出したものになる

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

$\boldsymbol{\Lambda}$ ではなく $\boldsymbol{\Sigma}$ で aa 成分を抜き出す点に注意。
 結果はシンプルだが、証明は意外と大変
 (PRML § 2.3.2)



周辺分布は積分に対応
 Gaussianを積分してもGaussian

Conditional Gaussian distribution

条件付き確率を得るには

$$p(\mathbf{x}_a, \mathbf{x}_b) = p(\mathbf{x}_a | \mathbf{x}_b) p(\mathbf{x}_b)$$

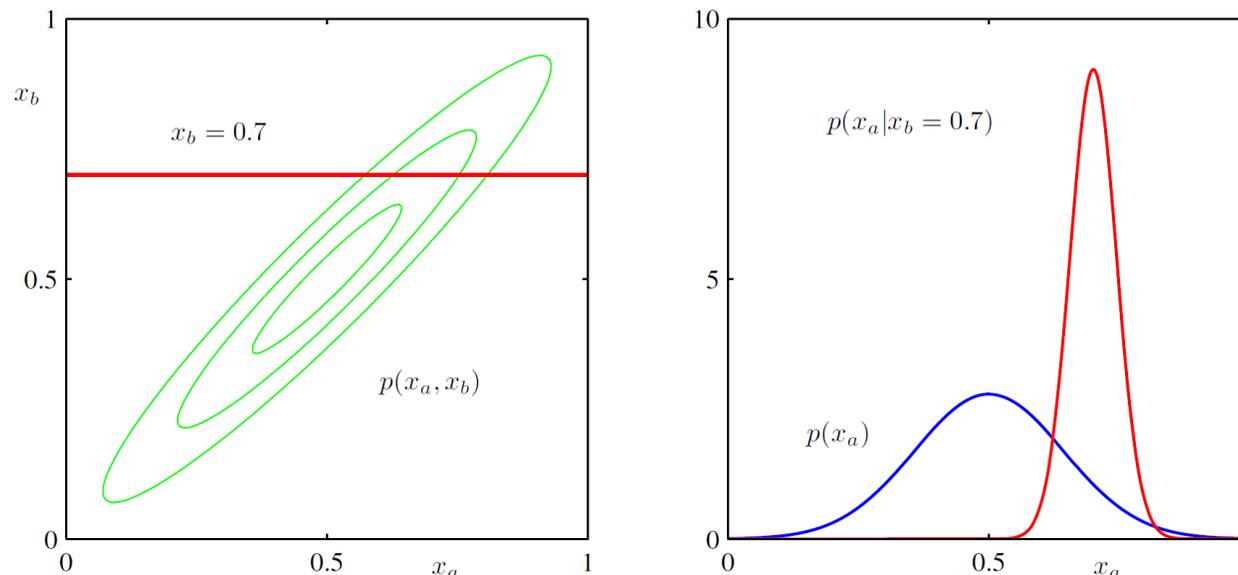
product rule

の両辺を比較する。結果は

$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$$

$$\begin{aligned} \boldsymbol{\mu}_{a|b} &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

$$\begin{aligned} \boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Lambda}_{aa}^{-1} \\ &= \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba} \end{aligned}$$



PRML Fig. 2.9

条件付き確率はカットに対応
GaussianをカットしてもGaussian

公式集

実際に必要になったら公式集を見るのがよい (Bishop, PRML)
どこに載っているかを覚えておくことが重要

2.3. The Gaussian Distribution

2.3.1 Conditional Gaussian distributions

2.3.2 Marginal Gaussian distributions

Partitioned Gaussians

Given a joint Gaussian distribution $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$ and

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad (2.94)$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \quad \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}. \quad (2.95)$$

Conditional distribution:

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1}) \quad (2.96)$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b). \quad (2.97)$$

Marginal distribution:

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}). \quad (2.98)$$

2.3.3 Bayes' theorem for Gaussian variables

Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for \mathbf{x} and a conditional Gaussian distribution for \mathbf{y} given \mathbf{x} in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (2.113)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \quad (2.114)$$

the marginal distribution of \mathbf{y} and the conditional distribution of \mathbf{x} given \mathbf{y} are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \quad (2.115)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \quad (2.116)$$

where

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}. \quad (2.117)$$

ここまでは準備。いよいよ計算

ベイズ統計に基づく予測分布 (Predictive distribution)

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int d\mathbf{w} p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) \quad \dots \textcircled{1}$$

② \mathbf{w} が与えられたときの予測分布

$$p(t|x, \mathbf{w}) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) \propto \exp \left[-\frac{\beta}{2} (t - \mathbf{w}^T \phi(x))^2 \right]$$

線形モデル

$$y(x, \mathbf{w}) = \mathbf{w}^T \phi(x)$$

③ \mathbf{w} の事後確率分布

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}) = \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1}) \prod_{j=1}^M \mathcal{N}(w_j|0, \alpha^{-1})$$

$$\propto \exp \left[-\frac{\beta}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right]$$

Bayes' theorem

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}) p(\mathbf{w})$$

尤度関数 事前分布

結果はガウス分布 (ガウス分布の畳み込み積分もまたガウス分布)

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x)) \quad \text{分散が } x \text{ に依存する}$$

今日やること

多変量ガウス分布 (multivariate Gaussian)
 の関係式を使って計算する

式③

$$\begin{aligned} p(\mathbf{w}|\mathbf{x}, \mathbf{t}) &= \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1}) \prod_{j=1}^M \mathcal{N}(w_j|0, \alpha^{-1}) && \text{1変数ガウス分布} \\ &= \mathcal{N}(\mathbf{t}|\Phi\mathbf{w}, \beta^{-1}\mathbf{I})\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) && \text{多変量ガウス分布} \\ &\propto \exp\left[-\frac{\beta}{2}(\mathbf{t} - \Phi\mathbf{w})^T(\mathbf{t} - \Phi\mathbf{w}) - \frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right] \end{aligned}$$

線形モデル

$$y(x, \mathbf{w}) = \phi(x)^T \mathbf{w}$$

計画行列

$$\Phi_{nj} = \phi_j(x_n)$$

$N \times M$ 行列

\mathbf{w} に関して平方完成すると

これは公式を使わなくても計算できる

しかし公式を使いこなすのもテクニックのひとつ

対応する公式は、PRML 式(2.116)

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{ML}}, S)$$

$$\mathbf{w}_{\text{ML}} = \beta S \Phi^T \mathbf{t} \quad \text{最尤推定におけるパラメータの最適値}$$

$$S = (\beta \Phi^T \Phi + \alpha \mathbf{I})^{-1}$$

最尤推定における共分散行列
(最適値周りの揺らぎ)

式①

$$\begin{aligned} p(t|x, \mathbf{x}, \mathbf{t}) &= \int d\mathbf{w} p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) \\ &= \int d\mathbf{w} \mathcal{N}(t|\phi(x)^T \mathbf{w}, \beta^{-1}) \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{ML}}, S) \\ &\propto \int d\mathbf{w} \exp \left[-\frac{\beta}{2} (t - \phi(x)^T \mathbf{w})^2 - \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{ML}})^T S^{-1} (\mathbf{w} - \mathbf{w}_{\text{ML}}) \right] \end{aligned}$$

\mathbf{w} に関して平方完成して積分すると

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$

$$m(x) = \beta \phi(x)^T S \Phi^T \mathbf{t}$$

$$s^2(x) = \beta^{-1} + \phi(x)^T S \phi(x)$$

この計算は結構大変

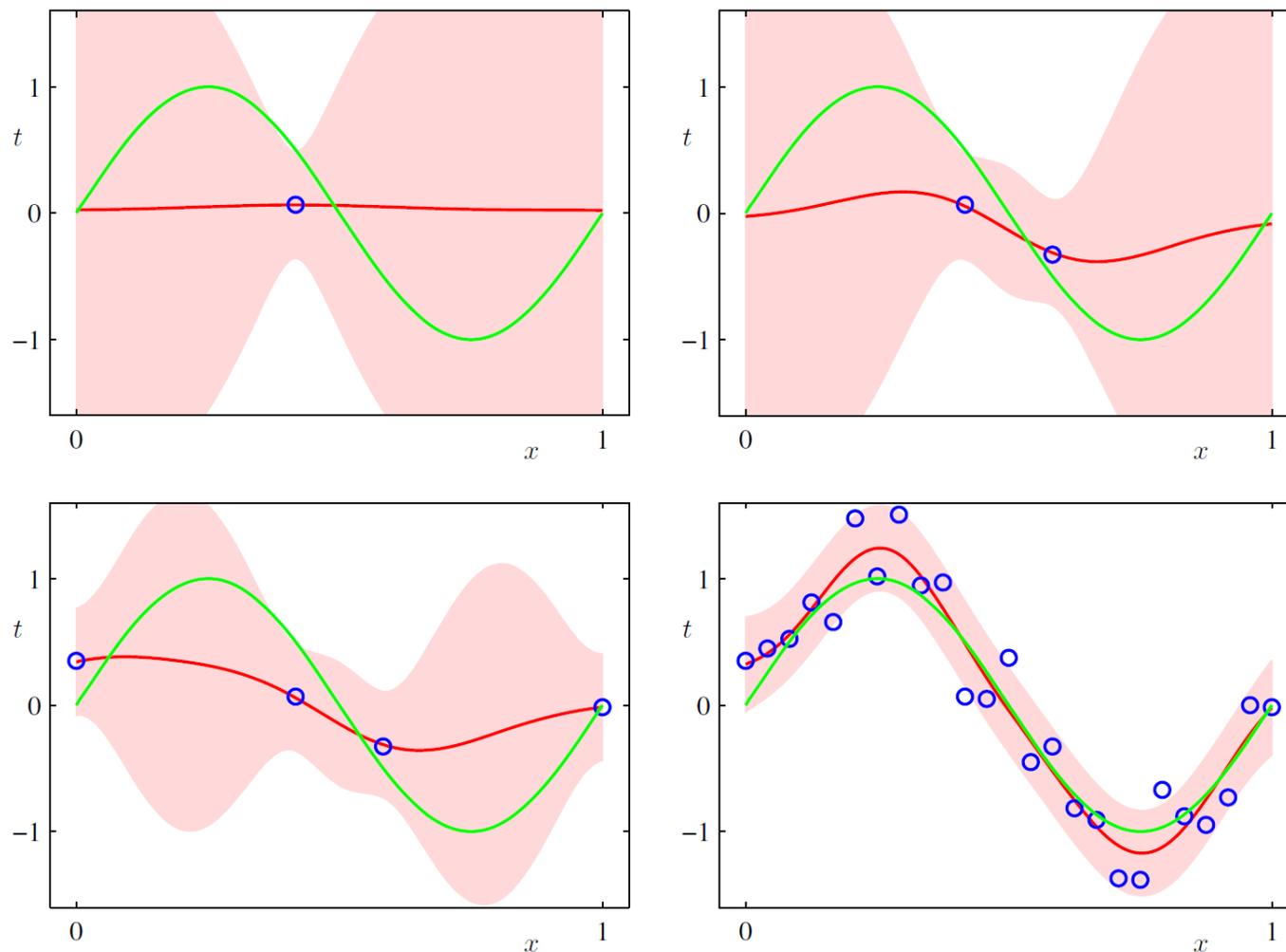
腕試しに計算してみるのもよい

公式を使う場合は、PRML 式(2.115)

\mathbf{w} の分布による予測分布の拡がり

これが最尤推定と正しいベイズ推定の違い

予測分布



$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$

緑：答え

青：測定データ

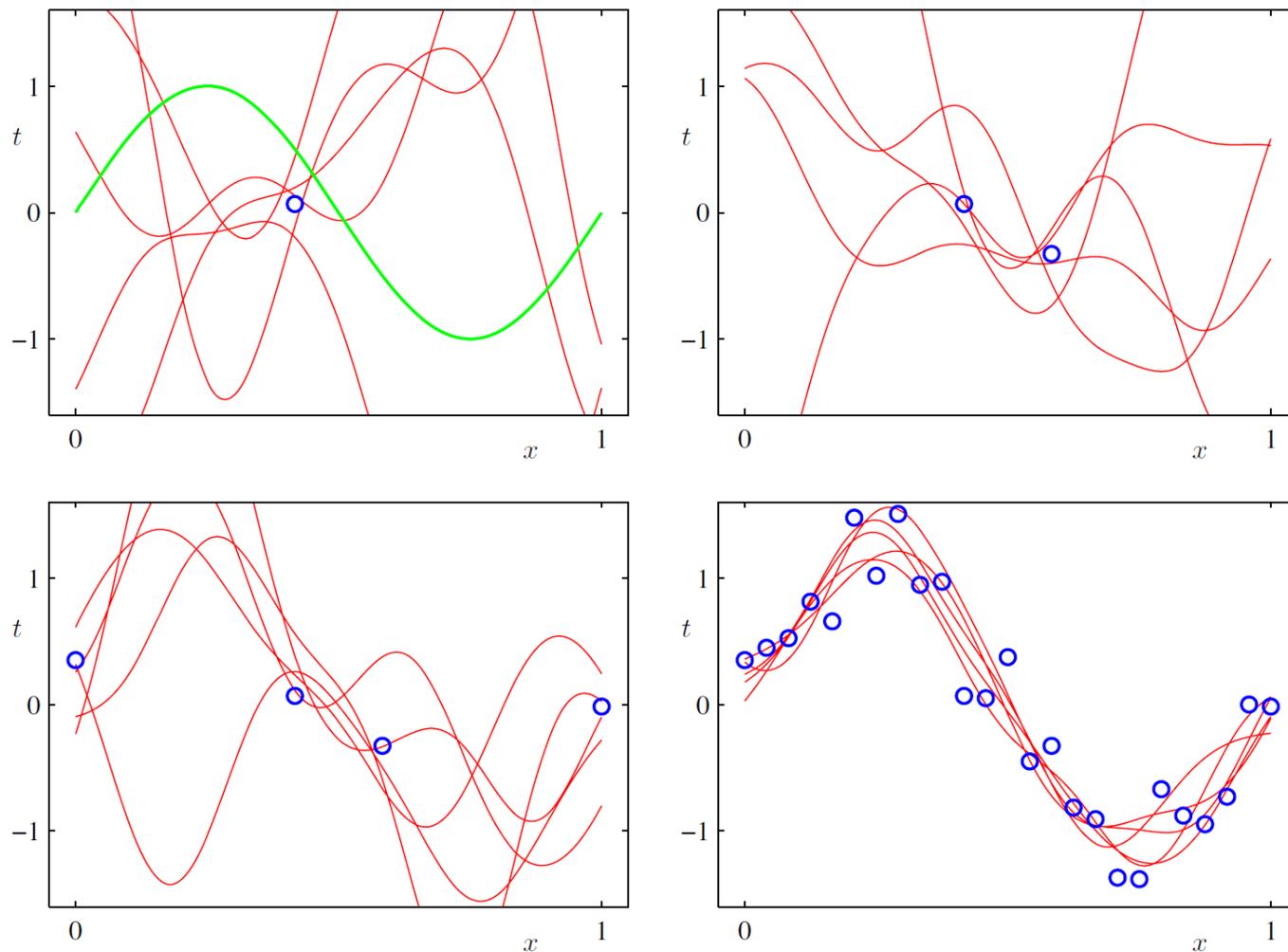
赤： $m(x)$

ピンク： $[m(x) - s(x), m(x) + s(x)]$

基底関数：
ガウス関数（あとで）

PRML **Figure 3.8** Examples of the predictive distribution (3.58) for a model consisting of 9 Gaussian basis functions of the form (3.4) using the synthetic sinusoidal data set of Section 1.1. See the text for a detailed discussion.

予測曲線



赤線：
予測曲線の例

$$y(x, \mathbf{w}) = \phi^T \mathbf{w}$$

パラメータ \mathbf{w} の事後確率分布に従ってサンプルを抽出

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t})$$

実際の予測分布は、赤線を中心とした分散 β^{-1} のガウス分布
それを積分したものが前頁の図

PRML **Figure 3.9** Plots of the function $y(x, \mathbf{w})$ using samples from the posterior distributions over \mathbf{w} corresponding to the plots in Figure 3.8.

まとめ

ベイズ線形回帰による予測分布

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$

$$m(x) = \beta \phi(x)^T S \Phi^T \mathbf{t}$$

$$s^2(x) = \beta^{-1} + \phi(x)^T S \phi(x)$$

$$S = (\beta \Phi^T \Phi + \alpha \mathbf{I})^{-1}$$

最尤推定における予測曲線

$$m(x) = \phi(x)^T \mathbf{w}_{\text{ML}} \quad \mathbf{w}_{\text{ML}} = \beta S \Phi^T \mathbf{t}$$

\mathbf{w} の分布による予測分布の広がり

これが最尤推定と正しいベイズ推定の違い

最尤推定における共分散行列
(最適値周りの揺らぎ)

補足：基底関数

1次元データ(x, t)

$$y(x, \mathbf{w}) = \sum_j w_j \phi_j(x)$$

基底関数

空間的に局在している関数が好ましい
 関数が広がりを持っていると、一つのデータ点が全ての w_j に影響する
 その意味で、多項式はあまりよくない

ちなみに、これまでは入力が1次元としてきたが、実際は何次元でも定式化は変わらない。

$$\mathbf{x} = (x_1, x_2, \dots, x_D)$$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=1}^M w_j \phi_j(\mathbf{x})$$

$\phi_j(x) \rightarrow \phi_j(\mathbf{x})$ と置き換えるだけ

基底関数の例

多項式 $\phi_j(x) = x^j$

スプライン (区分多項式)

ガウス関数

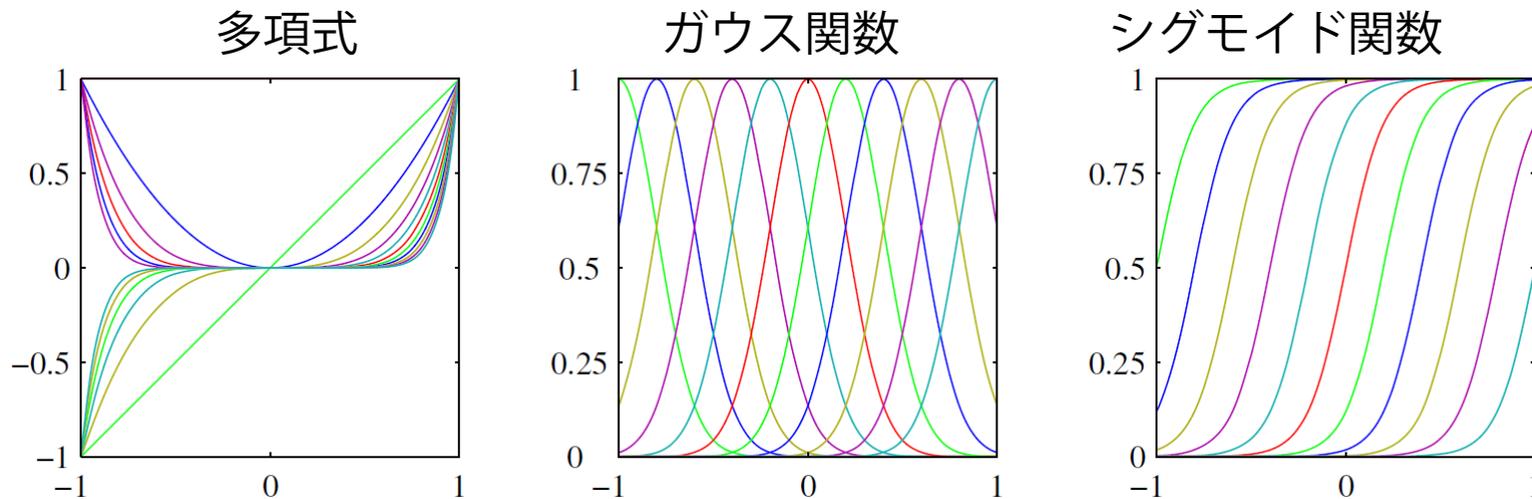
$$\phi_j(x) = \exp\left[-\frac{(x - \mu_j)^2}{2s^2}\right]$$

シグモイド (sigmoid) 関数

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

logistic関数



PRML Fig. 3.1