

W8PRML

3.4-3.6

東工大 山本航平

目次

- 3.4 ベイズモデル比較
- 3.5 エビデンス近似
 - 3.5.1 エビデンス関数の評価
 - 3.5.2 エビデンス関数の最大化
 - 3.5.3 有効パラメータ数
- 3.6 固定された基底関数の限界

3.4 ベイズモデル比較

- モデルの選択や正則化パラメータの決定
 - 1章ではクロスバリデーションを用いてモデルの選択を行った
 - ここではベイズの立場から考える
- ここでは一般論を述べ、3.5では線形回帰の正則化パラメータの具体例を述べる

L個のモデルを比較

- L個のモデル $\{\mathcal{M}_i\} (i = 1, \dots, L)$
 - このモデルはデータ \mathcal{D} 上の確率分布
 - 目的値 t 上の分布, 入力 \mathbf{X} は既知とする
- このモデルのどれかによってデータは生成されているがそのどれかはわからない
- このとき, データ集合が与えられた時の事後分布を評価する

$$p(\mathcal{M}_i | \mathcal{D}) \propto p(\mathcal{M}_i) p(\mathcal{D} | \mathcal{M}_i) \quad (3.66)$$

モデルに対する好み
(ここでは定数と仮定)

モデルエビデンス
(周辺尤度)

モデルエビデンス(周辺尤度)の評価

- ベイズ因子(=エビデンスの比) $\frac{p(\mathcal{D}|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_j)}$
- モデルの事後分布がわかると予測分布がわかる

$$p(t|\mathbf{x}, \mathcal{D}) = \sum_{i=1}^L p(t|\mathbf{x}, \mathcal{M}_i, \mathcal{D})p(\mathcal{M}_i|\mathcal{D}) \quad (3.67)$$

- この分布は**混合分布**である
 - 全体の予測分布が個々のモデルの予測分布の事後確率に関する重み付き平均で与えられる
=モデル平均(の近似)を求めたい

モデル選択

- モデル平均のもっとも単純な近似
→ 一番もっともらしいモデルを選ぶ (= **モデル選択**)
- エビデンスは以下で与えられる (\mathbf{w} はモデルのパラメータ)

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)d\mathbf{w} \quad (3.68)$$

- モデルエビデンスはパラメータの事後分布を計算するときの分母

$$p(\mathbf{w}|\mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D}|\mathbf{w}, \mathcal{M}_i)p(\mathbf{w}|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_i)} \quad (3.69)$$

モデルエビデンスの別の解釈

- パラメータがひとつの例 $p(\mathcal{D}) = \int p(\mathcal{D}|w)p(w)dw$ (3.70')

- 仮定:

- 事前確率は平坦で幅が Δw_{prior} ($p(w) = \frac{1}{\Delta w_{\text{prior}}}$)
- パラメータの事後分布が最頻値 (w_{MAP}) の近傍で尖っている
- その幅が $\Delta w_{\text{posterior}}$

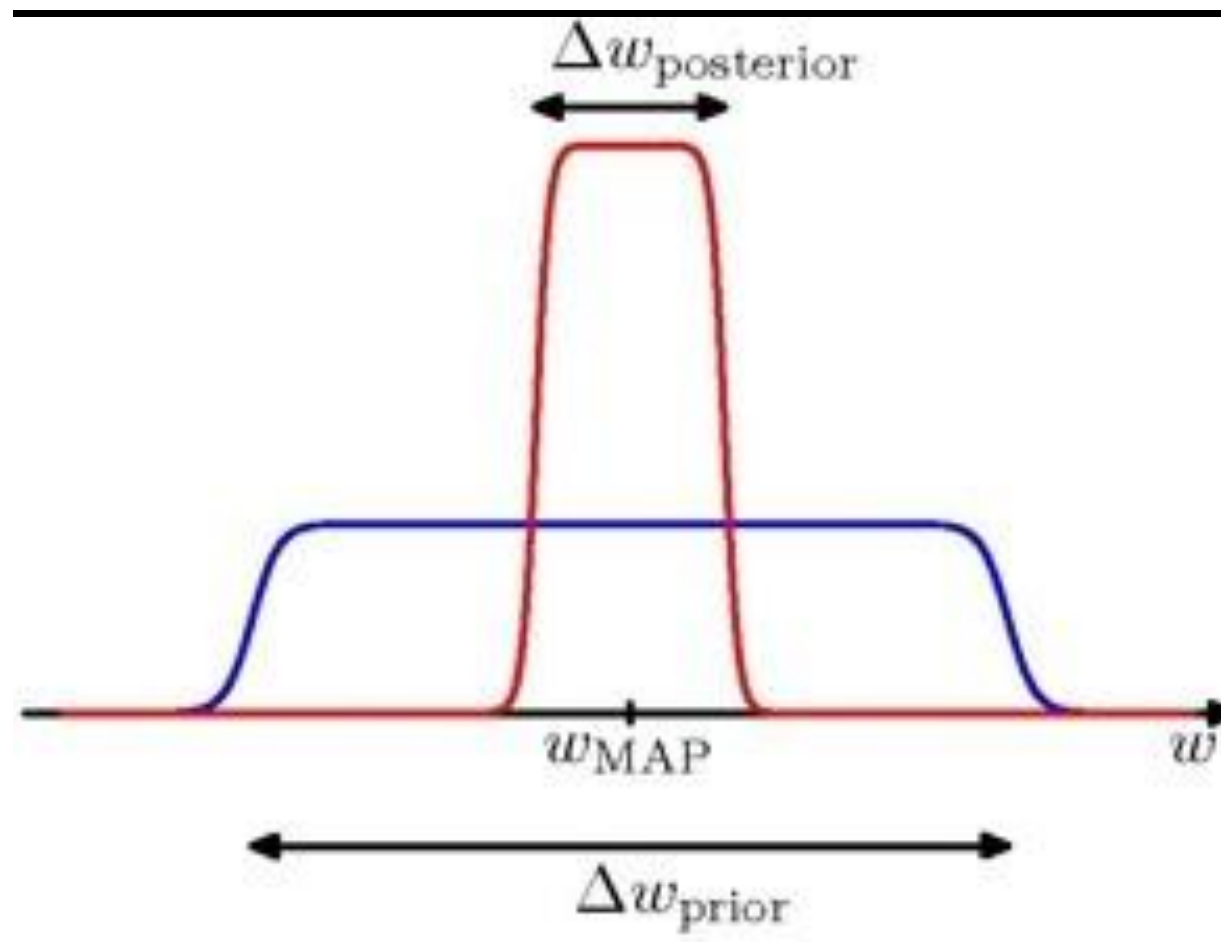
- 仮定を用いると, 掛け算で(3.70')を近似できる

$$p(\mathcal{D}) = \int p(\mathcal{D}|w)p(w)dw \simeq p(\mathcal{D}|w_{\text{MAP}}) \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \quad (3.70)$$

- 対数をとる

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|w_{\text{MAP}}) + \ln \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \quad (3.71)$$

近似の様子を図に示す



(3.71)の解釈

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|w_{\text{MAP}}) + \ln \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \quad (3.71)$$

- 第一項
 - 一番尤もらしいパラメータによるデータへのフィッティング度
- 第二項
 - ペナルティ項
 - 事後分布がデータに強くフィットするようにパラメータを調整するとペナルティは大きくなる

M個のパラメータを含む場合

- 全てのパラメータにおける $\frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}}$ の値が等しいとき、以下がなりたつ

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D} | \mathbf{w}_{\text{MAP}}) + M \ln \left(\frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right) \quad (3.72)$$

- モデルの複雑さを増すと...

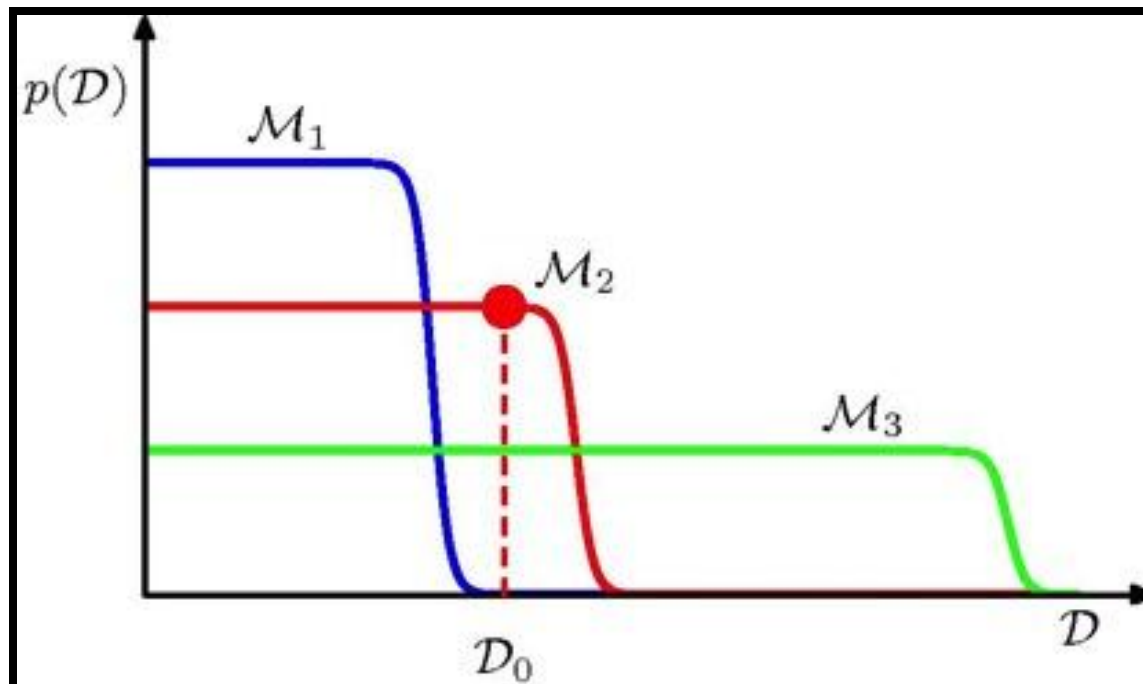
- 第一項は増加
- 第二項は減少



バランスが大事

さらなる解釈

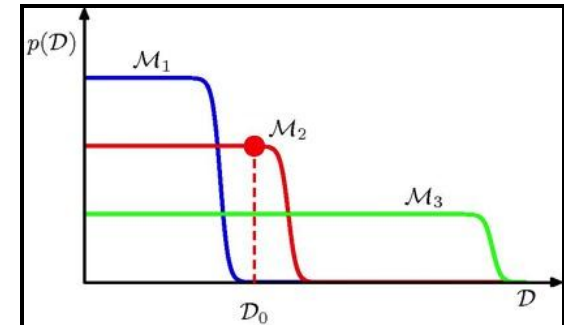
- 3つの、複雑さの異なるモデルを考えた時、なぜ周辺尤度最大化で中間の複雑さのモデルが選ばれるのか



横軸：
各データ集合空間

さらなる解釈

- データ \mathcal{D} の生成
 - パラメータのを事前分布 $p(\mathbf{w})$ に従って選択
 - それに基づいてデータを $p(\mathcal{D}|\mathbf{w})$ からサンプリング
 - 単純なモデル (M_1) は多様性に乏しい(狭い)
 - 複雑なモデル (M_3) はデータの範囲が広い
- 単純なモデル $\rightarrow \mathcal{D}_0$ を生成できない
- 複雑なモデル $\rightarrow \mathcal{D}_0$ が選ばれることが相対的に少ない



ベイズモデル比較の枠組み

- ベイズモデル比較においては、
モデル集合の中に真の分布が含まれていると仮定
- 2つのモデル $\mathcal{M}_1, \mathcal{M}_2$ を考える(前者が真のモデル)
- 期待ベイズ因子
$$\int p(\mathcal{D}|\mathcal{M}_1) \ln \frac{p(\mathcal{D}|\mathcal{M}_1)}{p(\mathcal{D}|\mathcal{M}_2)} d\mathcal{D} \quad (3.73)$$
 - KLダイバージェンスの例
 - 2つの分布が等しい時にだけ0, それ以外は常に正
 - 従って, 平均的には常に正しいモデルのベイズ因子のほうが大きい
- ベイズの枠組みでは
 - モデルの形に関して仮定をおいているのでテスト集合を用意して性能評価すべき

3.5 エビデンス近似

- 線形基底関数モデルを完全にベイズ的に取り扱う
 - 超パラメータ α , β の導入
 - 通常のパラメータだけでなく, 超パラメータも周辺化して予測
- 全ての変数の上で完全に積分するのは難しい
- 周辺尤度関数を最大にするように超パラメータを決める
 - 経験ベイズ, 第二種の最尤推定, 一般化最尤推定, エビデンス近似

予測分布(α と β と \mathbf{w} 全部を周辺化)

$$p(t|\mathbf{t}) = \int \int \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) p(\alpha, \beta|\mathbf{t}) d\mathbf{w} d\alpha d\beta \quad (3.74)$$

- ただし,

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \quad (3.8)$$

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \quad (3.49)$$

- 事後分布 $p(\alpha, \beta|\mathbf{t})$ が $\hat{\alpha}, \hat{\beta}$ の近傍で鋭く尖っている時

$$p(t|\mathbf{t}) \simeq p(t|\mathbf{t}, \hat{\alpha}, \hat{\beta}) = \int p(t|\mathbf{w}, \hat{\beta}) p(\mathbf{w}|\mathbf{t}, \hat{\alpha}, \hat{\beta}) d\mathbf{w} \quad (3.75)$$

尖ってる α, β の場所を求める

- $\hat{\alpha}, \hat{\beta}$ の事後分布は

$$p(\alpha, \beta | \mathbf{t}) \propto p(\mathbf{t} | \alpha, \beta) p(\alpha, \beta) \quad (3.76)$$

- 事前分布が平坦のとき
- エビデンス関数 $p(\mathbf{t} | \alpha, \beta)$ を最大化することで, $\hat{\alpha}, \hat{\beta}$ を求める
- 対数エビデンスを最大化する方法
 - EMアルゴリズム
 - 解析的に求める(3.5.2節)

3.5.1 エビデンス関数の評価

- エビデンス関数は以下の通り

$$p(\mathbf{t}|\alpha, \beta) = \int p(\mathbf{t}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)d\mathbf{w} \quad (3.77)$$

- これをごによごによします
- (3.11), (3.12), (3.52)よりエビデンス関数は以下の様に書ける

$$p(\mathbf{t}|\alpha, \beta) = \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \int \exp\{-E(\mathbf{w})\}d\mathbf{w} \quad (3.78)$$

- ただし, M はパラメータ \mathbf{w} の次元数であり,

$$\begin{aligned} E(\mathbf{w}) &= \beta E_D(\mathbf{w}) + \alpha E_W(\mathbf{w}) \\ &= \frac{\beta}{2} \|\mathbf{t} - \Phi\mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \end{aligned} \quad (3.79)$$

(3.78)の確認 (演習3.17)

- 3.11, 3.12から $p(\mathbf{t}|\mathbf{w}, \beta) = \exp \left\{ \frac{N}{2} \ln \frac{\beta}{2\pi} - \beta E_D(\mathbf{w}) \right\}$
 $= \left(\frac{\beta}{2\pi} \right)^{\frac{N}{2}} \exp \{ -\beta E_D(\mathbf{w}) \}$
- 3.52から $p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \exp \left\{ -\frac{1}{2} \mathbf{w}^T (\alpha^{-1}\mathbf{I})^{-1} \mathbf{w} \right\}$
- 従って

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

$$= \left(\frac{\beta}{2\pi} \right)^{\frac{N}{2}} \left(\frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \int \exp \left\{ -\beta E_D(\mathbf{w}) - \alpha \frac{1}{2} \mathbf{w}^T \mathbf{w} \right\}$$



$E_D(\mathbf{w})$



$-\mathbf{E}(\mathbf{w})$

(3.79)を平方完成(演習3.18)

$$E(\mathbf{w}) = E(\mathbf{m}_N) + \frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^T \mathbf{A}(\mathbf{w} - \mathbf{m}_N) \quad (3.80)$$

• ただし,

$$\mathbf{A} = \alpha \mathbf{I} + \beta \Phi^T \Phi \quad (3.81)$$

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N \quad (3.82)$$

$$\mathbf{m}_N = \beta \mathbf{A}^{-1} \Phi^T \mathbf{t} \quad (3.84)$$

• また, $\mathbf{A} = \nabla \nabla E(\mathbf{w})$ と対応している(ヘッセ行列)

• (3.54)から, $\mathbf{A} = \mathbf{S}_N^{-1}$ が成り立つ \rightarrow (3.84)は事後分布の平均

- 以上から(演習3.19help!)

$$\begin{aligned} & \int \exp \{-E(\mathbf{w})\} d\mathbf{w} \\ &= \exp \{-E(\mathbf{m}_N)\} \int \exp \left\{ \frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^T \mathbf{A} (\mathbf{w} - \mathbf{m}_N) \right\} \\ &= \exp \{-E(\mathbf{m}_N)\} (2\pi)^{\frac{M}{2}} |\mathbf{A}|^{-\frac{1}{2}} \end{aligned} \quad (3.85)$$

- 対数周辺尤度は

$$\ln p(\mathbf{t}|\alpha, \beta) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - E(\mathbf{m}_N) - \frac{1}{2} \ln |\mathbf{A}| - \frac{N}{2} \ln 2\pi \quad (3.86)$$

演習(3.18)

- まず展開

$$\frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

$$= \frac{\beta}{2} (\mathbf{t}^T \mathbf{t} - 2\mathbf{t}^T \Phi \mathbf{w} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}) + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

$$= \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - 2\beta \mathbf{t}^T \Phi \mathbf{w} + \mathbf{w}^T (\beta \Phi^T \Phi + \alpha \mathbf{I}) \mathbf{w})$$

A

- さらに式変形

$$\frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - 2\beta \mathbf{t}^T \Phi \mathbf{w} + \mathbf{w}^T \mathbf{A} \mathbf{w})$$

$$= \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - 2\beta \mathbf{t}^T \Phi \mathbf{A}^{-1} \mathbf{A} \mathbf{w} + \mathbf{w}^T \mathbf{A} \mathbf{w})$$

$$= \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - 2\mathbf{m}_N^T \mathbf{A} \mathbf{w} + \mathbf{w}^T \mathbf{A} \mathbf{w})$$

$$= \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - 2\mathbf{m}_N^T \mathbf{A} \mathbf{w} + \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{m}_N^T \mathbf{A} \mathbf{m}_N - \mathbf{m}_N^T \mathbf{A} \mathbf{m}_N)$$

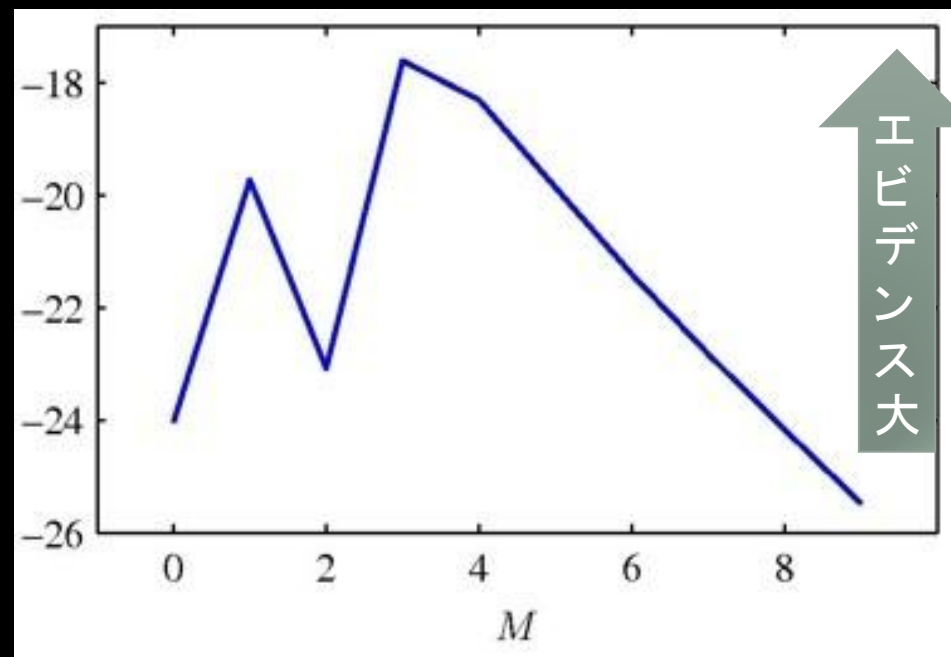
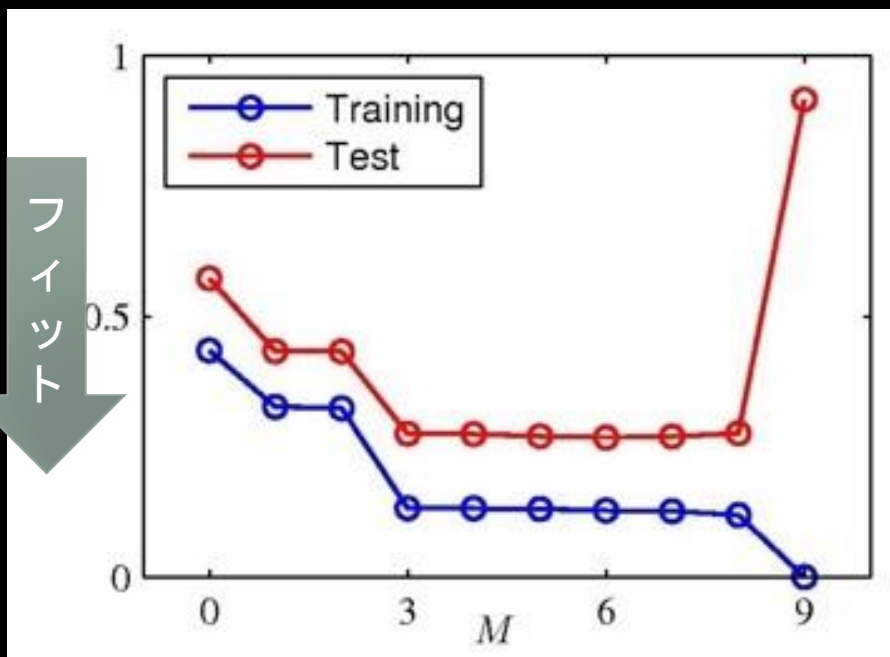
$$= \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - \mathbf{m}_N^T \mathbf{A} \mathbf{m}_N) + \frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^T \mathbf{A} (\mathbf{w} - \mathbf{m}_N)$$

 $(\mathbf{A}^{-1})^T = \mathbf{A}^{-1}$

- 第一項を更に変形

$$\begin{aligned} & \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - \mathbf{m}_N^T \mathbf{A} \mathbf{m}_N) \\ &= \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - 2\mathbf{m}_N^T \mathbf{A} \mathbf{m}_N + \mathbf{m}_N^T \mathbf{A} \mathbf{m}_N) \\ &= \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - 2\mathbf{m}_N^T \mathbf{A} \mathbf{A}^{-1} \Phi^T \mathbf{t} \beta + \mathbf{m}_N^T (\alpha \mathbf{I} + \beta \Phi^T \Phi) \mathbf{m}_N) \\ &= \frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - 2\mathbf{m}_N^T \Phi^T \mathbf{t} \beta + \beta \mathbf{m}_N^T \Phi^T \Phi \mathbf{m}_N + \alpha \mathbf{m}_N^T \mathbf{m}_N) \\ &= \frac{1}{2} (\beta (\mathbf{t} - \Phi \mathbf{m}_N)^T (\mathbf{t} - \Phi \mathbf{m}_N) + \alpha \mathbf{m}_N^T \mathbf{m}_N) \\ &= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N \end{aligned}$$

回帰の問題に戻ります



- 左図1.4, 右図3.14
- エビデンス最大は $M=3$ (データを説明できる最も簡単なモデル)
(よくデータにフィットしておりかつ複雑さが一番小さいから)

3.5.2 エビデンス関数の最大化

- $p(\mathbf{t}|\alpha, \beta)$ を α に関して最大化する問題を定義
- まず $\ln |\mathbf{A}|$ の α に関する導関数を求めるために以下の固有ベクトル方程式を考える

$$(\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad (3.87)$$

- (3.81)より, \mathbf{A} は固有値 $\alpha + \lambda_i$ を持つ. したがって

$$\frac{d}{d\alpha} \ln |\mathbf{A}| = \frac{d}{d\alpha} \ln \prod_i (\lambda_i + \alpha) = \frac{d}{d\alpha} \sum_i \ln(\lambda_i + \alpha) = \sum_i \frac{1}{\lambda_i + \alpha} \quad (3.88)$$

- (3.86)を α で微分して, 0とおく

$$0 = \frac{M}{2\alpha} - \frac{1}{2} \mathbf{m}_N^T \mathbf{m}_N - \frac{1}{2} \sum_i \frac{1}{\lambda_i + \alpha} \quad (3.89)$$

- 両辺に 2α をかけて整理

$$\alpha \mathbf{m}_N^T \mathbf{m}_N = M - \alpha \sum_i \frac{1}{\lambda_i + \alpha} = \gamma \quad (3.90)$$

- Σ には M 個の項が含まれるので

$$\gamma = \sum_i \frac{\lambda_i}{\lambda_i + \alpha} \quad (3.91)$$

- 周辺尤度を最大にする α は

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N} \quad (3.92)$$

α の求め方

1. (式(3.53, 3.91))から \mathbf{m}_N, γ を始めに求める
(これらは α に依存する)
2. (3.92)を使って α を最推定
3. この仮定を収束するまで続ける
 - $\Phi^T \Phi$ は変化しないので, 固有値を最初に選ぶだけで良い
 - β に関しても同じ手順で求めることができる

βの場合

- ※ λ_i が β に比例することに注意する $\frac{d\lambda_i}{d\beta} = \frac{\lambda_i}{\beta}$

$$\frac{d}{d\beta} \ln |\mathbf{A}| = \frac{d}{d\beta} \sum_i \ln(\lambda_i + \alpha) = \frac{1}{\beta} \sum_i \frac{1}{\lambda_i + \alpha} = \frac{\gamma}{\beta} \quad (3.93)$$

- よって

$$0 = \frac{N}{2\beta} - \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{m}_N^T \phi(\mathbf{x}_n)\}^2 \quad (3.94)$$

- 整理すると

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N \{t_n - \mathbf{m}_N^T \phi(\mathbf{x}_n)\}^2 \quad (3.95)$$

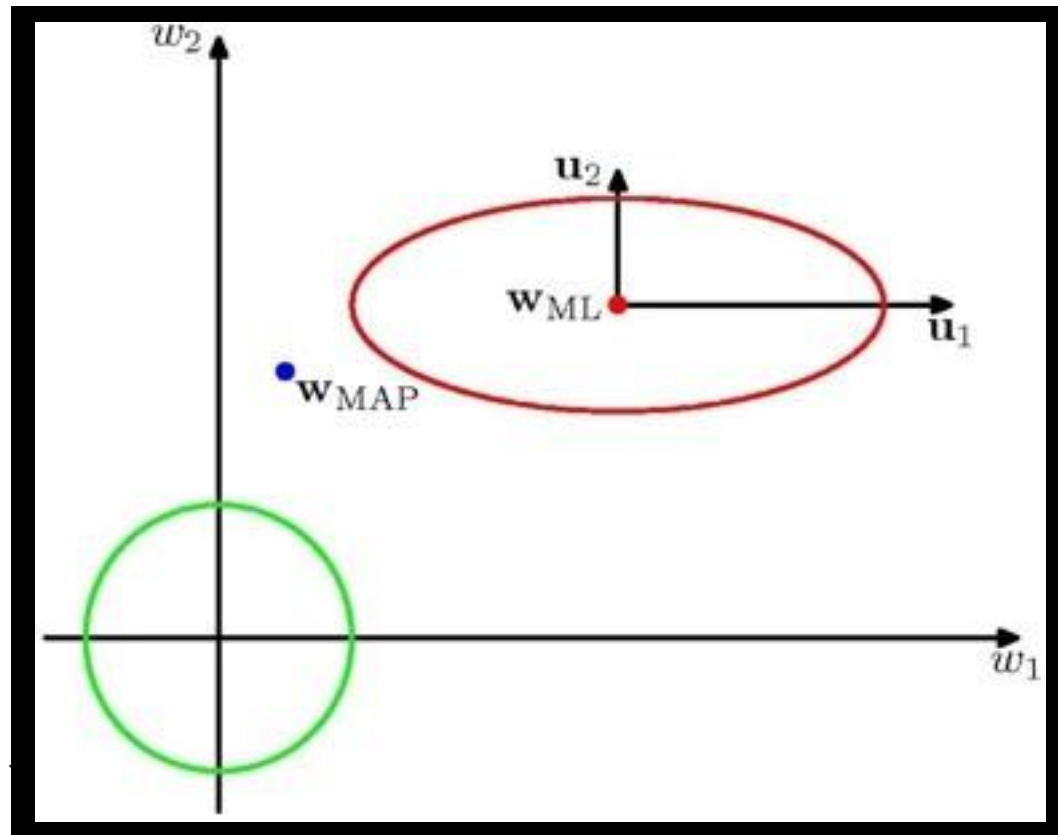
3.5.3 有効パラメータ数

- ベイズ解 α の解釈
$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N} \quad , \quad \gamma = \sum_i \frac{\lambda_i}{\lambda_i + \alpha}$$
- 尤度関数と事前分布の等高線を描く(図3.15, 次のスライド)
 - 固有値は全て正 ($\beta \Phi^T \Phi$ が正定値行列)なので

$$0 \leq \frac{\lambda_i}{\lambda_i + \alpha} \leq 1 \quad 0 \leq \gamma \leq M$$
- $\lambda_i \gg \alpha \rightarrow \frac{\lambda_i}{\lambda_i + \alpha} \simeq 1 \rightarrow \gamma$ によく影響する (well-determined)
- $\lambda_i \ll \alpha \rightarrow \frac{\lambda_i}{\lambda_i + \alpha} \simeq 0 \rightarrow \gamma$ にあまり影響しない
- したがって, γ は有効(well-determined)なパラメータの数

図3.15

- 緑: 事前分布
- 赤: 尤度関数
- ここでは $\lambda_1 < \alpha < \lambda_2$
- λ_2 は λ_1 よりも影響が大きい
→ λ_2 のほうが最尤に近い



- ぼくにはよくわからない

$\beta(3.95)$ の解釈

- 1変数 x のガウス分布の分散の最尤推定値(式1.56より)

$$\sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2 \quad (3.96)$$

- μ_{ML} はデータに含まれるノイズまでフィットしているので、
上式はバイアスをもつ
→平均の推定に自由度の1つを使ってしまっている
したがって、分散の不偏推定量は

$$\sigma_{\text{MAP}}^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2 \quad (3.97)$$

バイアスを取り
除くための"-1"

- 線形回帰の場合，有効なパラメータ数は γ
- 式(3.95)では，そのぶんの補正をおこなっている

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N \{t_n - \mathbf{m}_N^T \phi(X_n)\}^2 \quad (3.95)$$

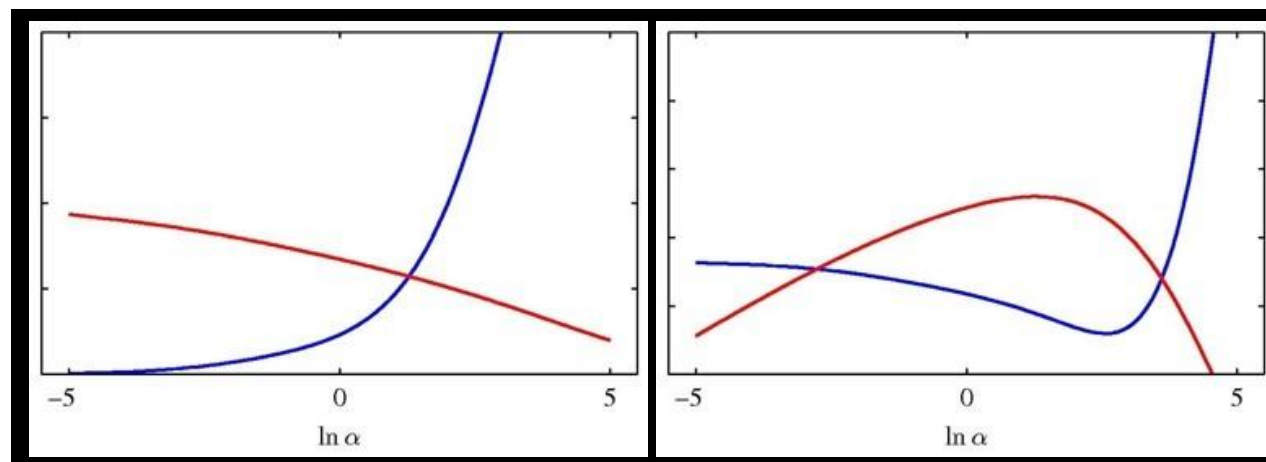


例:

三角関数をガウス基底関数モデルで近似

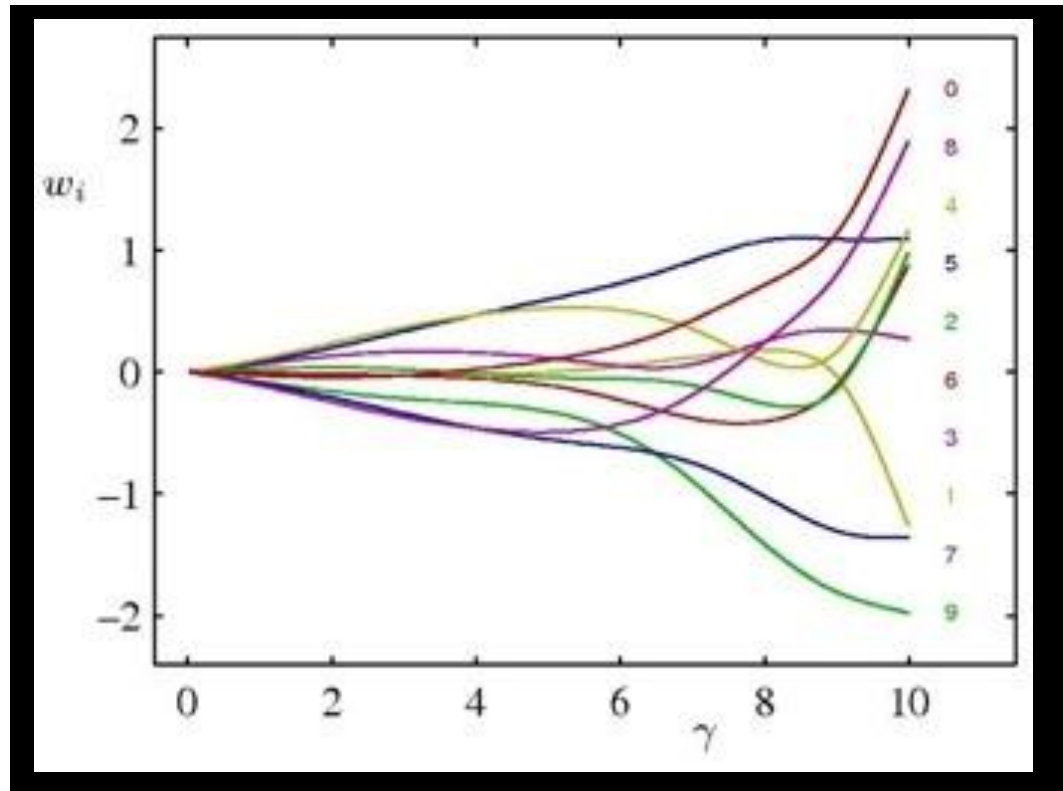
- 9個の基底関数を使う $M=10$ (9+バイアスパラメータ)
- β を真の値11.1に設定したあとに, エビデンスの枠組みで α を決定する
- 左: 青線 $\rightarrow \gamma$, 赤線 $\rightarrow 2E_W(\mathbf{m}_N)$, 交点が最適な α の値
- 右: 青線 \rightarrow テスト集合に対する誤差, 赤線 \rightarrow 対数エビデンス
- 右図は, (左図の)交点の位置で誤差最小, エビデンス最大になっていることがわかる

$$\ln p(\mathbf{t}|\alpha, \beta)$$



各パラメータの値と γ の関係

- 各グラフが各パラメータに対応
- α を $0 \leq \alpha \leq \infty$ の範囲で変化させると, γ が $0 \leq \gamma \leq M$ の範囲で変化する
- α の値がパラメータの大きさを制御している



データ点の数がパラメータの数と比べて十分大きい時

- $N \gg M$ の時, 全てのパラメータはwell-determined
 - 式(3.87)における $\Phi^T \Phi$ はデータ点に関する陰的な和を含んでいるため, データ集合のサイズとおもに固有値の数も増加するため
- このとき, $\gamma = M$
- すると, α と β の再推定方程式は

$$\alpha = \frac{M}{2E_W(\mathbf{m}_N)} \quad \beta = \frac{N}{2E_D(\mathbf{m}_N)} \quad (3.98), (3.99)$$

- ただし, $E_W(\mathbf{m}_N), E_D(\mathbf{m}_N)$ は式(3.25)(3.26)で定義した

$$E_W(\mathbf{m}_N) = \frac{1}{2} \mathbf{m}_N^T \mathbf{m}_N \quad E_D(\mathbf{m}_N) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$

3.6 固定された基底関数の限界

- メリット
 - パラメータの線形性を仮定している
→ 最小二乗問題の閉じた解が求まる
 - ベイズ推定の計算が簡単になる
 - 基底関数を適切に選べば、任意の非線形変換をモデル化できる
- デメリット
 - 訓練データの観測前に基底関数を固定する
 - 入力空間の次元数に対して、指数的に基底関数を増やす必要がある
- デメリットの解決法 → 実際のデータがもつ性質を利用する
 - 実データは限られた非線形多様体上に存在する
 - 入力空間中のデータがある場所にのみ基底関数を配置できる
 - 目標変数がデータ多様体中の少数の可能な方向にしか強く依存しない

3.86を β で微分

$$\ln p(\mathbf{t}|\alpha, \beta) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - E(\mathbf{m}_N) - \frac{1}{2} \ln |\mathbf{A}| - \frac{N}{2} \ln 2\pi \quad (3.86)$$

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N$$

$$\mathbf{m}_N = \beta \mathbf{A}^{-1} \Phi^T \mathbf{t}$$

- というかこれの微分→展開しただけだった