

3 Linear Models for Regression

3.1 线性基函数模型

3.1.1 最大似然与最小平方

0 平方误差

1 w_{ML} 求解

代码实现

2 $w_{0 ML}$ 求解

3 β_{ML}^{-1} 求解

3.1.2 最小平方的几何描述

3.1.3 Sequential learning

1 stochastic gradient descent(sequential gradient descent)

3.1.4 正则化最小平方

3.1.5 多个输出 (K 个)

3.1 回顾

3.2 偏置 - 方差分解

1 最优的预测

2 分解关于 $y(\mathbf{x})$ 的第一项

3 trade-off

3.2 回顾

3.3 贝叶斯线性回归

3.3.1 参数分布

+ 为简化, 考虑零均值高斯先验

代码实现

+ 例子: 线性基函数的贝叶斯学习过程

3.3.2 预测分布

3.3.3 等价核

+ 通过线性基函数模型输出引出等价核

+ 相关性质

3.3 回顾

3.4 贝叶斯模型比较

+ 从贝叶斯角度考虑模型选择问题

+ 贝叶斯模型比较会倾向于选择出正确的模型

3.4 回顾

3.5 证据近似

3.5.1 计算证据函数

代码实现

3.5.2 最大化证据函数

3.5 回顾

3 Linear Models for Regression

回归问题的目标: 在给定 D 维输入变量 \mathbf{x} (input)的情况下, 预测一个或多个连续目标变量 t (target)的值.

3.1 线性基函数模型

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

其中令 $\phi_0(\mathbf{x}) = 1$, w_0 是bias.

基函数(basis function) $\phi_j(\mathbf{x})$: 输入变量 \mathbf{x} 的非线性函数. 请注意这里模型是参数 \mathbf{w} 的线性函数 (这极大简化了对模型的分析), 但同时关于输入变量是非线性的.

基函数的选择，还可以有傅里叶基函数(?wavelets):

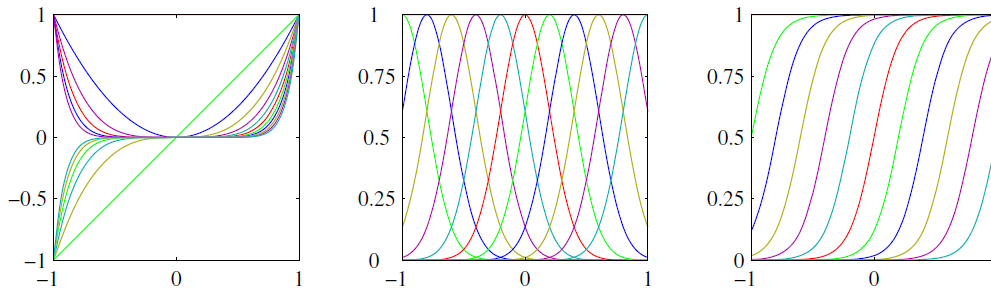


Figure 3.1 Examples of basis functions, showing polynomials on the left, Gaussians of the form (3.4) in the centre, and sigmoidal of the form (3.5) on the right.

3.1.1 最大似然与最小平方

维度表达: $\mathbf{x} \in \mathbb{R}^M$, $|\mathcal{D}| = N$, $\mathbf{x}_n \in \mathcal{D}$ 表示第 n 个数据, $\phi_j(\mathbf{x}_n)$ 表示对 \mathbf{x}_n 的第 j 个分量经过基函数.

0 平方误差

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \right\}^2 = \frac{1}{2} \sum_{n=1}^N \{ t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) \}^2$$

1 \mathbf{w}_{ML} 求解

如上, 给输出 $y(\mathbf{x}, \mathbf{w})$ 附加高斯噪声 ϵ (零均值, 精度(方差倒数)为 β), 此时后验为:

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

极大似然估计, 关于 \mathbf{w} 求偏导后为0, 得到:

$$0 = \sum_{n=1}^N t_n \boldsymbol{\phi}(\mathbf{x}_n)^T - \mathbf{w}^T \left(\sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^T \right)$$

normal equation: $\mathbf{w}_{ML} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t}$

其中design matrix $\boldsymbol{\Phi} \in \mathbb{R}^{N \times M}$:

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

代码实现

在实现时, $\boldsymbol{\Phi}$ 中的每一个高斯基函数均值是 $a \sim b$ 相同间隔的数, 就是直接输入数据的 \mathbf{x} 对应的高斯.

```
1 def fit(self, x:np.ndarray, t:np.ndarray):
2     # 3.17 MLE的 w = 伪逆(\Phi) @ t
3     self.w = np.linalg.pinv(X) @ t
4     self.var = np.mean(np.square(X @ self.w - t))
```

带有正则项的实现:

```
1 def fit(self, x:np.ndarray, t:np.ndarray):
2     self.w = np.linalg.inv(self.alpha * eye + X.T @ X) @ np.transpose(X) @ t
```

2 w_0 ML 求解

显示地写出偏置参数(bias) w_0 ，再极大似然估计后对 w_0 求偏导为0，得到：

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j$$

$$\text{where: } \bar{t} = \frac{1}{N} \sum_{n=1}^N t_n, \bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n)$$

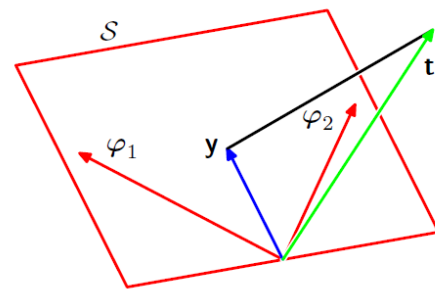
偏置 w_0 直接加到参数输出上，补偿了训练集上目标值的平均值，与基函数加权求和，之间的差。

3 β_{ML}^{-1} 求解

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{t_n - \mathbf{w}_{ML}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

3.1.2 最小平方的几何描述

Figure 3.2 Geometrical interpretation of the least-squares solution, in an N -dimensional space whose axes are the values of t_1, \dots, t_N . The least-squares regression function is obtained by finding the orthogonal projection of the data vector \mathbf{t} onto the subspace spanned by the basis functions $\phi_j(\mathbf{x})$ in which each basis function is viewed as a vector φ_j of length N with elements $\phi_j(\mathbf{x}_n)$.



$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

理论分析 平方误差形式中 \sum_n^N 可以被放到内层：

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \sum_{j=1}^M w_j \phi_j(\mathbf{x}_n) \right\}^2$$

$$= \frac{1}{2} \sum_{j=1}^M \{ \}$$

子空间 \mathcal{S} ： Φ 的第 j 列所张成的空间，即基函数作用于所有 \mathbf{x}_n 的某一个分量，在 n 个数据点估计的基函数值。

看不懂，为什么 M 维度？

3.1.3 Sequential learning

顺序算法中，模型在每观测到一个数据点后即更新。

1 stochastic gradient descent(sequential gradient descent)

在如上平方误差问题中：

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta \left(t_n - \mathbf{w}^{(\tau)T} \phi_n \right) \phi_n$$

3.1.4 正则化最小平方

误差函数添加正则项控制过拟合，使确定模型复杂度问题 从确定合适基函数数量 转换到 确定合适的正则化系数 λ ：

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

解出 \mathbf{w} 为：

$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

更一般的正则化项 (q 次)：

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

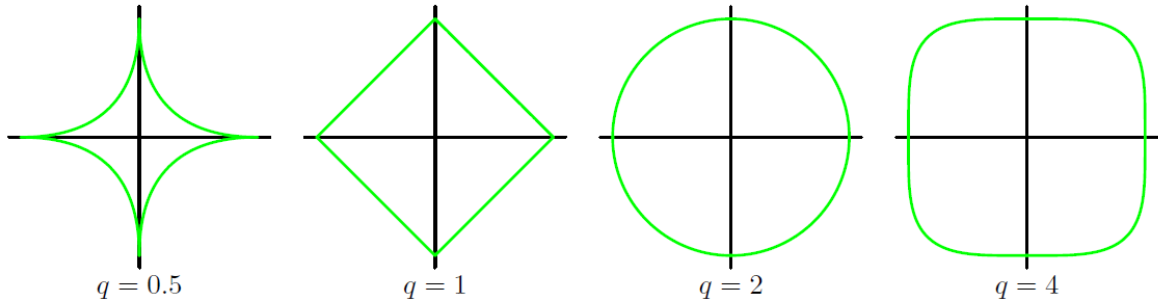


Figure 3.3 Contours of the regularization term in (3.29) for various values of the parameter q .

- $q = 1$: lasso，此时如果 λ 充分大，某些 w_j 将为 0，从而产生了一个稀疏模型，即输入中某些维度(基函数作用后)不产生作用。

定理：最小化上述 q 的一般正则化目标函数形式 \Leftrightarrow 在 $\sum_{j=1}^M |w_j|^q \leq \eta$ 条件下 最小化未带正则项的平方误差函数。

3.1.5 多个输出 (K 个)

多个输出对应的基函数都不一样 即对于某一个 \mathbf{x} ，有 $\phi(\mathbf{x}) \in \mathbb{R}^M$ ， $\mathbf{W} \in \mathbb{R}^{M \times K}$ ：

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

对数似然函数：

$$\begin{aligned} \ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{W}^T \phi(\mathbf{x}_n), \beta^{-1} \mathbf{I}) \\ &= \frac{NK}{2} \ln \left(\frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \|\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n)\|^2 \end{aligned}$$

可否分解成 K 个独立的回归问题？可以，因为 \mathbf{W} 仅与高斯噪声分布的均值有关，从 2.3.4 节中知道多元高斯分布均值的最大似然解与协方差无关。

3.1 回顾

从线性基函数模型出发：

- 3.1.1：该模型平方误差下最大似然的求解(\mathbf{w} 与附加高斯噪声的 β^{-1})。
- 3.1.2：最小平方误差的几何描述(但是没看懂)。
- 3.1.3：非常简要的梯度下降。
- 3.1.4：带有正则化项的线性基函数模型(介绍了一般形式的正则项及其作用)。
- 3.1.5： \mathbf{x} 是多个维度的，在后面代码实现中都是这样。

3.2 偏置 - 方差分解

考虑正则化系数 λ 的大小？

- 过大：导致一些 \mathbf{w} 变得很小，限制了模型描述数据中有趣且重要规律的灵活性。
- 过小：导致过拟合。

所以从频率角度研究模型复杂度：偏置 - 方差分解。

1 最优的预测

由条件概率分布的期望给出：

$$h(\mathbf{x}) = \mathbb{E}[t | \mathbf{x}] = \int t p(t | \mathbf{x}) dt$$

平方损失函数的期望：

$$\mathbb{E}[\mathcal{L}] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

找到一个 $y(\mathbf{x})$ 使第一项最小，而与 $y(\mathbf{x})$ 无关的第二项是由数据本身的噪声造成的，表示期望损失能够到达的最小值。

2 分解关于 $y(\mathbf{x})$ 的第一项

$h(\mathbf{x})$ 是由参数向量 \mathbf{w} 控制的函数 $y(\mathbf{x}, \mathbf{w})$ 进行建模的。每个数据集都独立同分布地从 $p(t, \mathbf{x})$ 中抽取，我们的学习算法在此基础上得到一个 $y(\mathbf{x}, \mathcal{D})$ ，对第一项分解：

$$\begin{aligned} \{y(\mathbf{x}) - h(\mathbf{x})\}^2 &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &\quad + 2 \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\} \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\} \end{aligned}$$

对上式关于 \mathcal{D} 求期望，最后交叉项为零(展开即可)。所以：

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ = \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{偏置})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{方差}} \end{aligned}$$

所以期望损失 $\mathbb{E}[\mathcal{L}]$ 被分解成了 偏置的平方 + 方差 + 噪声(注意，上述每一项都是期望/积分形式)：

$$\begin{aligned}\text{偏置}^2 &= \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \\ \text{方差} &= \int \mathbb{E}_{\mathcal{D}} \left[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 \right] p(\mathbf{x}) d\mathbf{x} \\ \text{噪声} &= \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt\end{aligned}$$

3 trade-off

150 3. LINEAR MODELS FOR REGRESSION

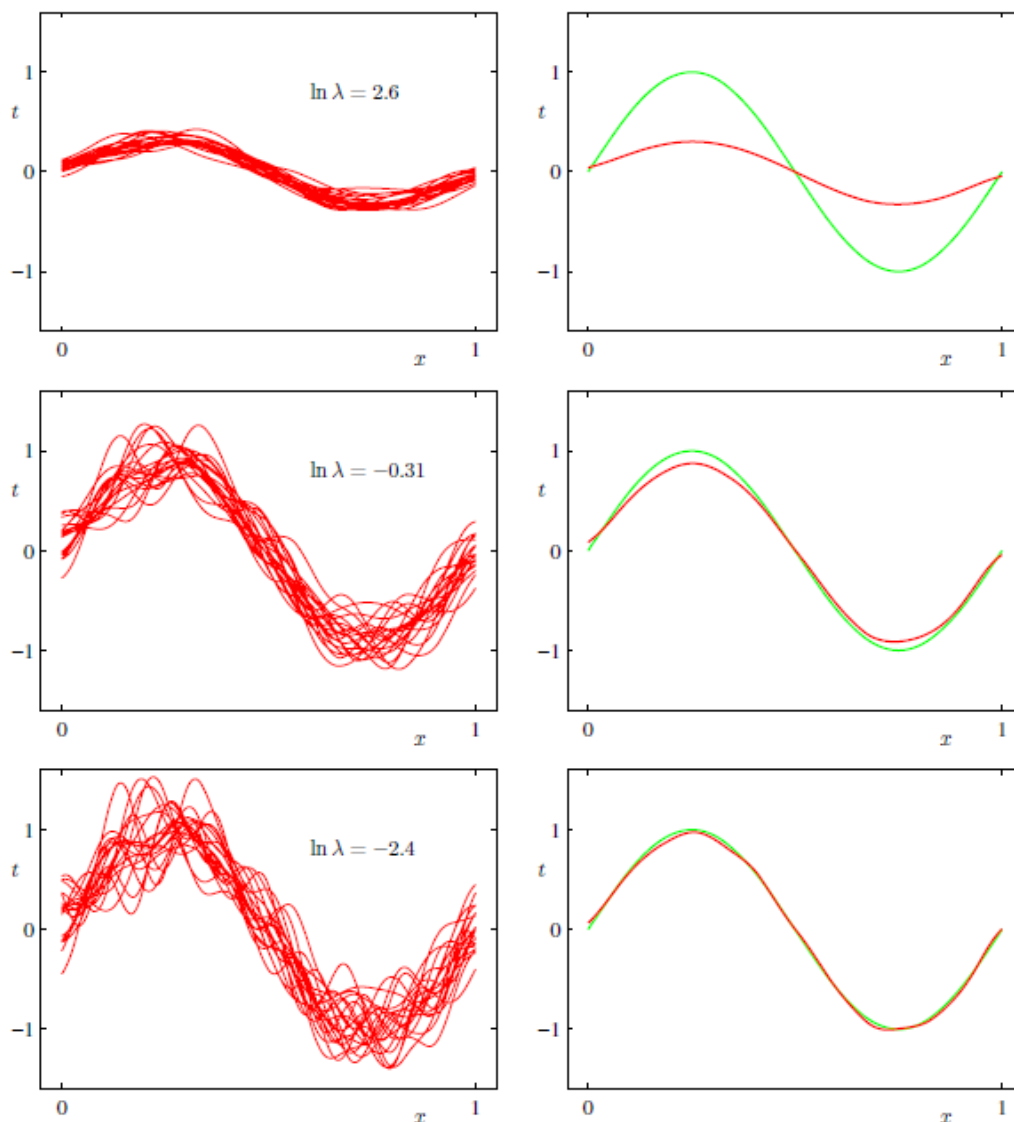


Figure 3.5 Illustration of the dependence of bias and variance on model complexity, governed by a regularization parameter λ , using the sinusoidal data set from Chapter 1. There are $L = 100$ data sets, each having $N = 25$ data points, and there are 24 Gaussian basis functions in the model so that the total number of parameters is $M = 25$ including the bias parameter. The left column shows the result of fitting the model to the data sets for various values of $\ln \lambda$ (for clarity, only 20 of the 100 fits are shown). The right column shows the corresponding average of the 100 fits (red) along with the sinusoidal function from which the data sets were generated (green).

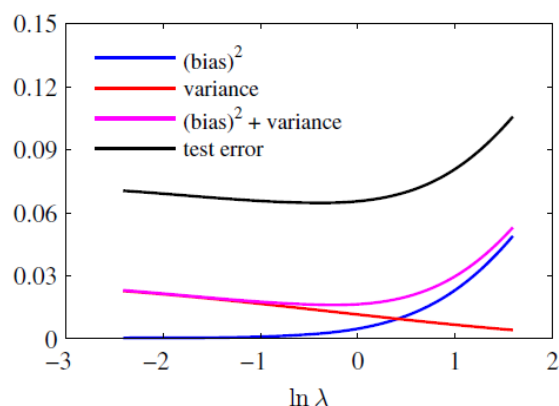
- 偏置 - 方差trade-off, 上图: 100个数据集, 每个数据集包含独立从sin函数中取样的25个数据点; 这100个数据集, 每一个拟合了带有多个高斯基函数的回归模型。
 - 上图: λ 较大时, (左图)这100个数据集拟合的函数较接近(偏置小), (右图)但是他们的平均离真实的sin函数(右图绿色)较远(方差大)。
- 定量计算:

$$\begin{aligned}\text{偏置}^2 &= \frac{1}{N} \sum_{n=1}^N \{\bar{y}(x_n) - h(x_n)\}^2 \\ \text{方差} &= \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \{y^{(l)}(x_n) - \bar{y}(x_n)\}^2\end{aligned}$$

其中 $\bar{y}(x) = \frac{1}{L} \sum y^{(l)}(x)$, 即100个数据集上拟合的函数之平均。这里所有数据以相同概率被采样。

- 局限性：偏置 - 方差分解需要对所有"相似"训练集求平均。

Figure 3.6 Plot of squared bias and variance, together with their sum, corresponding to the results shown in Figure 3.5. Also shown is the average test set error for a test data set size of 1000 points. The minimum value of $(\text{bias})^2 + \text{variance}$ occurs around $\ln \lambda = -0.31$, which is close to the value that gives the minimum error on the test data.



3.2 回顾

出发点：从频率角度研究模型复杂度， λ 对模型的影响量化分析。

- 1：平方损失函数分解为两部分，考虑与模型有关的第一项。
- 2：期望损失 $\mathbb{E}[\mathcal{L}]$ 被分解成了 偏置² + 方差 + 噪声。
- 3：多个模型(不同 λ)对数据集拟合，输出(所有函数的比较)和(所有的函数平均与真实采样函数的比较)，偏置方差都可以定量计算，考虑偏置方差的trade-off。

3.3 贝叶斯线性回归

MLE会产生过于复杂的模型并过拟合，而且需要对模型基函数的数量和形式进行选择。

3.3.1 参数分布

线性基函数拟合 共轭先验：

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_0, \mathbf{S}_0)$$

后验分布也是高斯分布：

$$\begin{aligned} p(\mathbf{w} \mid \mathbf{t}) &= \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N) \\ \mathbf{m}_N &= \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta \Phi^T \Phi \end{aligned}$$

最大后验估计(MAP)，最大化后验分布(众数)，所以此时 $\mathbf{w}_{MAP} = \mathbf{m}_N$,

+ 为简化, 考虑零均值高斯先验

那么如上后验分布参数为：

$$\begin{aligned} \mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi \end{aligned}$$

后验分布的对数似然：

$$\ln p(\mathbf{w} | \mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{常数}$$

后验分布关于 \mathbf{w} 的最大化等价于对平方和误差函数加上一个二次正则项进行最小化。形如(3.27)，其中 $\lambda = \frac{\alpha}{\beta}$ 。

代码实现

请结合下面一节。

```

1  # 获得先验:
2  # 先验受到上一层迭代的 self.w_mean, self.w_precision 影响.
3  mean_prev, precision_prev = self._get_prior(np.size(X, 1))
4  # 主要就是下面 由 先验(上一层的w_mean, w_precision), X 以及模型基本参数, 对w_mean,
   w_precision更新.
5  # 上面那句话是关键.
6  # (3.51):
7  w_precision = precision_prev + self.beta * X.T @ X
8  # (3.50) 移项:
9  w_mean = np.linalg.solve(
10     w_precision,
11     precision_prev @ mean_prev + self.beta * X.T @ t
12 )
13 self.w_mean = w_mean
14 self.w_precision = w_precision
15 # w_mean 就是一轮迭代下来的 超参数结果.
16 self.w_cov = np.linalg.inv(self.w_precision)

```

+ 例子: 线性基函数的贝叶斯学习过程

线性模型: $y(x, \mathbf{w}) = w_0 + w_1 x$ 。

生成数据: 从 $U(-1, 1)$ 中取 x_n , 代入 $f(x_n, a) = -0.3 + 0.5x_n$, 并附加 $\delta = 0.2$ 的高斯噪声得到 t_n

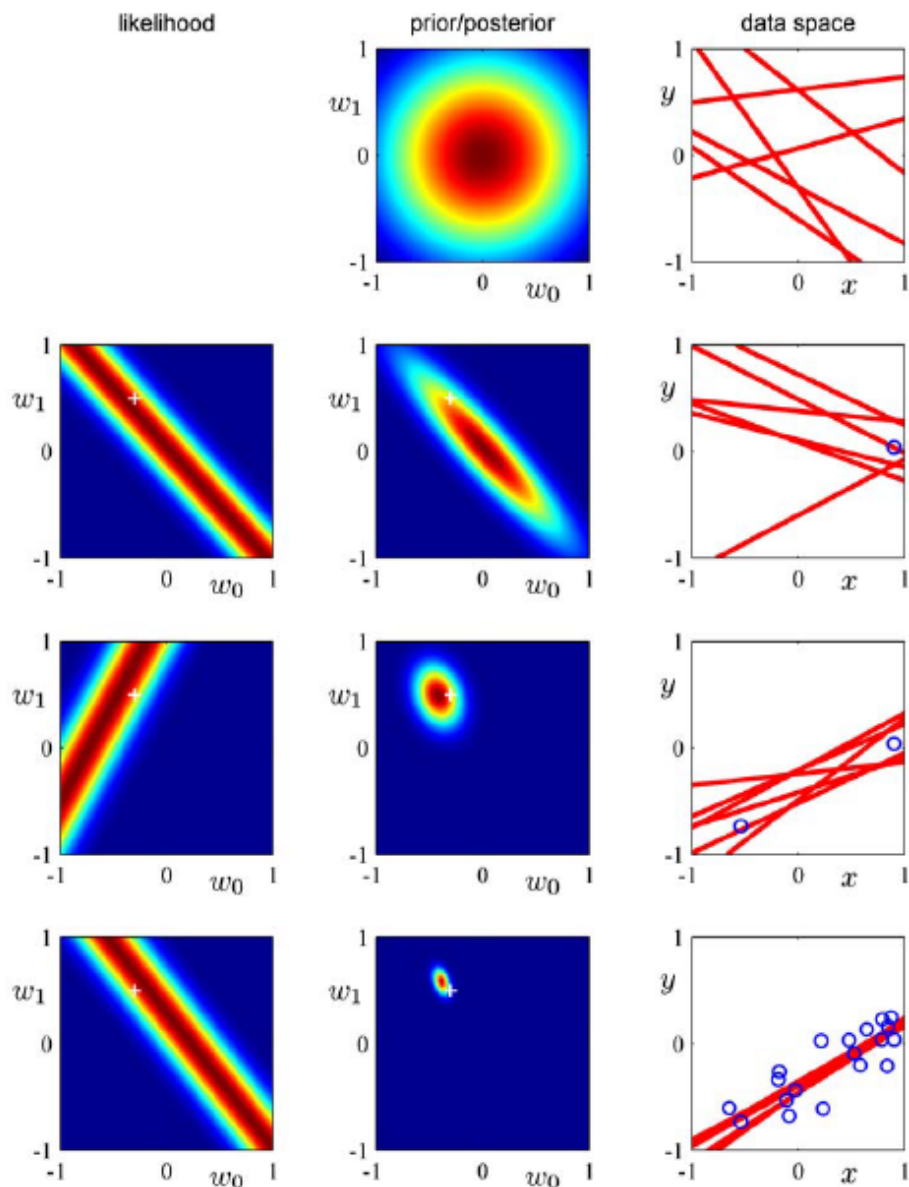


Figure 3.7 Illustration of sequential Bayesian learning for a simple linear model of the form $y(x, \mathbf{w}) = w_0 + w_1 x$. A detailed description of this figure is given in the text.

方差已知，所以 $\beta = (1/0.2)^2 = 25$ ， $\alpha = 2$ 。

- p.s. 上图左列和中间列都是 二维的概率分布图像。
- 上图第一行：上图第一幅（第一行中间列）代表 w_0 的先验，第一行右列表示从先验（中间列）中抽取的 w 后（就是本节最开始讲的生成数据）。
- 上图第二行：此时观测到了数据点（蓝色圆圈），左列是似然 $p(t, \mathbf{w})$ 关于 \mathbf{w} 的图像。
似然的限制为直线必须穿过数据点附近的位置，其中附近位置的范围由噪声精度确定。

◦ 迭代：

第二行左列的似然函数与第一行先验相乘，归一化后即得到第二行的后验，也是接下来第三行的先验。

- 上图第三行：因为两个点可以确定一条直线，所以得到后验概率分布也相对好。
- 如果继续往下，有无穷多个数据点，那么后验会变成一个 **delta** 函数。

以上考虑的是高斯先验，但是还有其他形式的先验：

$$p(\mathbf{w} | \alpha) = \left[\frac{q}{2} \left(\frac{\alpha}{2} \right)^{\frac{1}{q}} \frac{1}{\Gamma\left(\frac{1}{q}\right)} \right]^M \exp\left(-\frac{\alpha}{2} \sum_{j=0}^{M-1} |w_j|^q\right)$$

$q = 2$ 即为高斯。此时 \mathbf{w} 最大后验估计(如上) \Leftrightarrow 找到正则化误差函数的最小值。

3.3.2 预测分布

关注对于新的 \mathbf{x} 预测出的 t ，即 "预测分布"：

$$p(t | \mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t | \mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x})$$

其中方差 $\sigma_N^2(\mathbf{x})$ 中第一项表述数据中的噪声，第二项反映了与参数 \mathbf{w} 关联的不确定性。

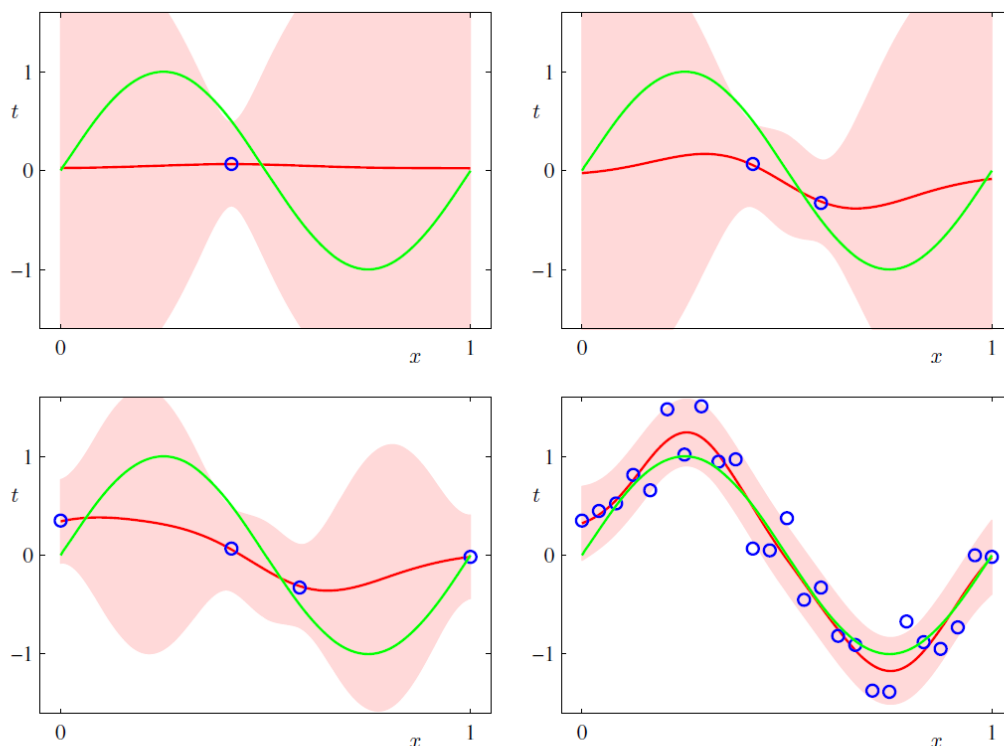


Figure 3.8 Examples of the predictive distribution (3.58) for a model consisting of 9 Gaussian basis functions of the form (3.4) using the synthetic sinusoidal data set of Section 1.1. See the text for a detailed discussion.

这里就相当于把上一节中说到的 $w_0 + w_1 x$ 换成了 $\sin(2\pi x)$ ：

- 其中绿色线是上一节采样的函数 $\sin(2\pi x)$ ，蓝色是生成的数据点(有高斯噪声)，红色阴影区域是均值两侧一个标准差范围。

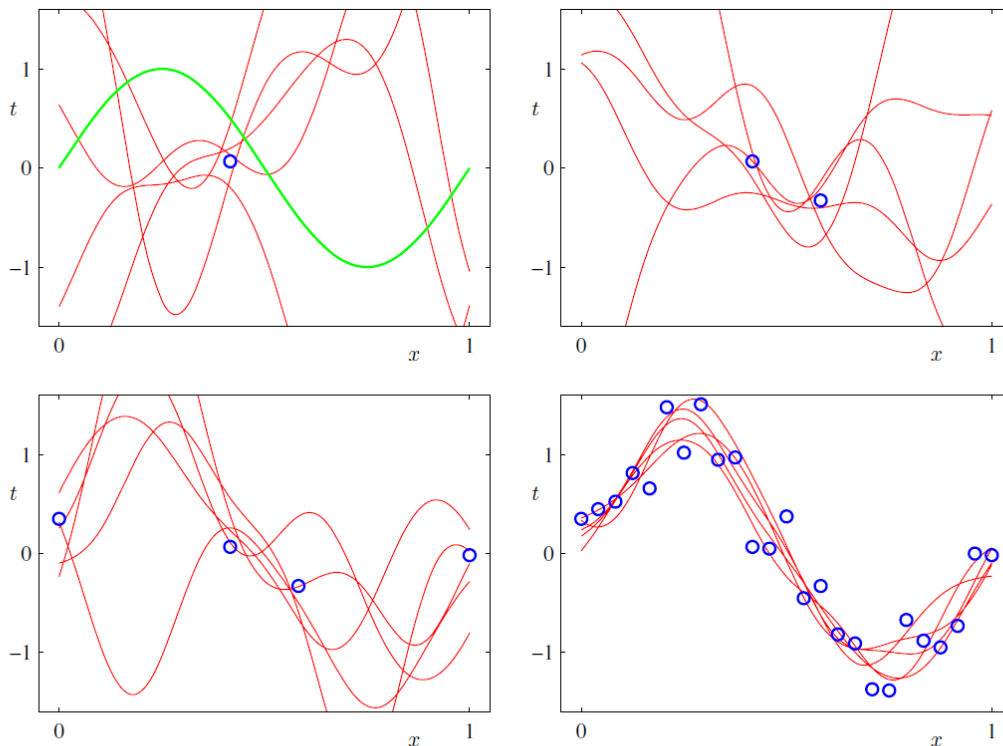


Figure 3.9 Plots of the function $y(x, \mathbf{w})$ using samples from the posterior distributions over \mathbf{w} corresponding to the plots in Figure 3.8.

上图是 $y(x, \mathbf{w})$ 图像，与上一节最右列一致，使用了服从 \mathbf{w} 上的后验概率的样本。

请注意上述过程的共轭先验是建立在 \mathbf{w}, β 都已知，如果未知则共轭为Gamma分布，此时预测分布为 \mathbf{t} 分布。

3.3.3 等价核

+ 通过线性基函数模型输出引出等价核

之前MAP，解取的是众数(对于高斯即均值)，(3.53) 零均值共轭先验的后验概率分布 的MAP结果代入线性基函数模型：

$$\begin{aligned} y(\mathbf{x}, \mathbf{w}) &= y(\mathbf{x}, \mathbf{m}_N) = \mathbf{m}_N^T \phi(\mathbf{x}) = \beta \phi(\mathbf{x})^T \mathbf{S}_N \Phi^T \mathbf{t} = \sum_{n=1}^N \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n) t_n \\ &= \sum_n^N k(\mathbf{x}, \mathbf{x}_n) t_n \end{aligned}$$

其中，以下被称为等价核(equivalent kernel)或平滑矩阵(smooth matrix)：

$$k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}')$$

+ 相关性质

- 可以表示为非线性函数的内积形式，即：

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= \psi(\mathbf{x})^T \psi(\mathbf{z}) \\ \text{where: } \psi(\mathbf{x}) &= \beta^{\frac{1}{2}} \mathbf{S}_N^{\frac{1}{2}} \phi(\mathbf{x}) \end{aligned}$$

- 考虑 $y(\mathbf{x})$ 和 $y(\mathbf{x}')$ 的协方差：

$$\begin{aligned} \text{cov}[y(\mathbf{x}), y(\mathbf{x}')] &= \text{cov}[\phi(\mathbf{x})^T \mathbf{w}, \mathbf{w}^T \phi(\mathbf{x}')] \\ &= \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}') = \beta^{-1} k(\mathbf{x}, \mathbf{x}') \end{aligned}$$

- 一个等价核定义了模型的权值，训练集里面的和新的 \mathbf{x} 组合，可以证明这些权值为1：

$$\sum_n^N k(\mathbf{x}, \mathbf{x}_n) = 1$$

3.3 回顾

比较有趣的一节：

- 3.3.1：共轭先验导出后验概率分布的形式，为简化可以采用零均值先验。之后讲解了贝叶斯线性回归的迭代学习过程，上一轮得到的后验作为下一轮的先验。
- 3.3.2： $w_0 + w_1 x$ 换成了 $\sin(2\pi x)$ ，具体的方差等也有写变化，但是没看懂。
- 3.3.3：等价核以及相关性质。

3.4 贝叶斯模型比较

+ 从贝叶斯角度考虑模型选择问题

模型比较的贝叶斯观点仅仅涉及到使用概率来表示模型选择的不确定性，数据由第 i 个模型生成的先验概率为： $p(\mathcal{M}_i)$ 。我们想估计后验分布：

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

由此有 多个模型的预测 混合分布，即对各个模型的预测分布加权平均：

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

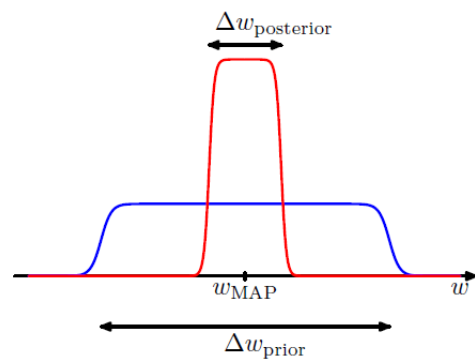
注意估计后验分布时的模型证据： $p(\mathcal{D} | \mathcal{M}_i)$ ，由全概率公式，贝叶斯公式(模型证据在其中作为归一化项)：

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

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

现在作如下假设：

Figure 3.12 We can obtain a rough approximation to the model evidence if we assume that the posterior distribution over parameters is sharply peaked around its mode w_{MAP} .

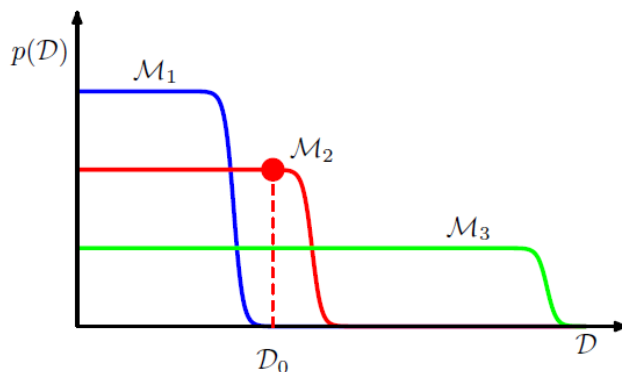


假设后验分布在 w_{MAP} 附近是尖峰(类似均匀分布的)，先验也是，宽度分别为 Δw_{sth} ：

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

接下来有个图看不懂：

Figure 3.13 Schematic illustration of the distribution of data sets for three models of different complexity, in which \mathcal{M}_1 is the simplest and \mathcal{M}_3 is the most complex. Note that the distributions are normalized. In this example, for the particular observed data set \mathcal{D}_0 , the model \mathcal{M}_2 with intermediate complexity has the largest evidence.



+ 贝叶斯模型比较会倾向于选择出正确的模型

考虑模型 $\mathcal{M}_1, \mathcal{M}_2$ ，将贝叶斯因子 $\frac{p(\mathcal{D}|\mathcal{M}_1)}{p(\mathcal{D}|\mathcal{M}_2)}$ 在数据集分布上进行平均

$$\int p(\mathcal{D} | \mathcal{M}_1) \ln \frac{p(\mathcal{D} | \mathcal{M}_1)}{p(\mathcal{D} | \mathcal{M}_2)} d\mathcal{D}$$

即 KL-散度。

3.4 回顾

这一节不是看的很懂。

- 不太明白那个 $p(\mathcal{D})$ 想表达什么。用近似均匀分布的后验和先验近似计算了。

3.5 证据近似

经验贝叶斯（第二类最大似然）：对 w 或者超参数的积分可能没有解析解，这里讨论一种近似方法：首先对参数 w 求积分，得到边缘似然函数，然后通过最大化边缘似然函数，确定超参数的值。

引入 α 和 β 上的超先验分布，预测分布可以通过对 w, α, β 求积分的方法得到：

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

上式如果定义了 α, β 上的共轭Gamma先验分布，那么对上式中的 w 就可以解析地得到 w 上的 \mathbf{t} 分布。虽然此时在 w 上的积分将不再有解析解（但可以用拉普拉斯近似的方法）。

如果后验在 $\hat{\alpha}, \hat{\beta}$ 附近有尖峰（近似均匀分布），那么上述积分可近似为：

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

3.5.1 计算证据函数

在证据框架中， $\hat{\alpha}, \hat{\beta}$ 可以通过最大化边缘似然函数 $p(\mathbf{t} | \alpha, \beta)$ 来获得，边缘似然函数 $p(\mathbf{t} | \alpha, \beta)$ 是全概率公式对 w 积分得到的：

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

接下来使用 对指数项配平方，然后使用高斯分布归一化形式：

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

$$E(\mathbf{w}) = \beta E_D(\mathbf{w}) + \alpha E_W(\mathbf{w})$$

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

如果忽略一些比例常数，对 \mathbf{w} 配平方：

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

其中令

$$\mathbf{A} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

有之前的 $\mathbf{S}_N^{-1} = \mathbf{A}$ ，以及

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

其中：

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

上式与之前定义的 \mathbf{m}_N 是完全一致的，表示后验概率分布均值。

由此 $E(\mathbf{w})$ 被整理成了较好的高斯形式，从而证据函数内的积分计算如下：

$$\int \exp\{-E(\mathbf{w})\} d\mathbf{w}$$

$$= \exp\{-E(\mathbf{m}_N)\} \int \exp\left\{-\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^T \mathbf{A} (\mathbf{w} - \mathbf{m}_N)\right\} d\mathbf{w}$$

$$= \exp\{-E(\mathbf{m}_N)\} (2\pi)^{\frac{M}{2}} |\mathbf{A}|^{-\frac{1}{2}}$$

所以证据函数(边缘似然函数)的对数：

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

代码实现

```

1 def _log_prior(self, w):
2     return -0.5 * self.alpha * np.sum(w ** 2)
3 def _log_likelihood(self, X, t, w):
4     return -0.5 * self.beta * np.square(t - X @ w).sum()
5
6 def _log_posterior(self, X, t, w):
7     # 式(3.82)
8     return self._log_likelihood(X, t, w) + self._log_prior(w)
9 def log_evidence(self, X: np.ndarray, t: np.ndarray):
10    # 式(3.86) 计算证据函数的对数。
11    N = len(t)
12    D = np.size(X, 1)
13    return 0.5 * (
14        D * np.log(self.alpha) + N * np.log(self.beta)
15        - np.linalg.slogdet(self.w_precision)[1] - D * np.log(2 * np.pi)
16    ) + self._log_posterior(X, t, self.w_mean)

```

3.5.2 最大化证据函数

- 证据函数的对数 对 α 求导：

上一节 证据函数的对数中出现了 $\ln |\mathbf{A}|$ ，对 $|\mathbf{A}|$ 展开之后才比较好求导，所以要用到其特征值：

- 由(3.81)，定义 \mathbf{A} 的特征值为 $\alpha + \lambda_i$ ：

$$\mathbf{A} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

- 证据函数的对数求导(关于 α)，整理：

$$0 = \frac{M}{2\alpha} - \frac{1}{2} m_N^T m_N - \frac{1}{2} \sum_i \frac{1}{\lambda_i + \alpha}$$

$$\text{Let: } \gamma = \alpha m_N^T m_N = M - \alpha \sum_i \frac{1}{\lambda_i + \alpha}$$

$$\alpha = \frac{\gamma}{m_N^T m_N}$$

- 注意这里 α 是一个隐式解：

γ 与 α 有关，最大后验概率的 m_N 也与 α 有关，所以 α 的估计是迭代的？

3.5 回顾

出发点：经验贝叶斯（第二类最大似然）：对 w 或者超参数的积分可能没有解析解，讨论一种近似方法：首先对参数 w 求积分，得到边缘似然函数，然后通过最大化边缘似然函数(证据函数)，确定超参数的值。

- 3.5.1：用高斯分布归一化系数来计算证据函数。
- 3.5.2：最大化上一节计算的证据函数(似然)，隐式解迭代计算直到收敛。