

PRML学习笔记——第三章

- PRML学习笔记——第三章
 - Linear Models for Regression
 - 3.1 Linear Basis Function Models
 - 3.1.1 Maximum likelihood and least squares
 - 3.1.2 Geometry of least squares
 - 3.1.3 Sequential learning
 - 3.1.4 Regularized least squares
 - 3.1.5 Multiple outputs
 - 3.2. The Bias-Variance Decomposition
 - 3.3. Bayesian Linear Regression
 - 3.3.1 Parameter distribution
 - 3.3.2 Predictive distribution
 - 3.3.3 Equivalent kernel
 - 3.4. Bayesian Model Comparison
 - 3.5. The Evidence Approximation

Linear Models for Regression

3.1 Linear Basis Function Models

最简单的用于回归的线性模型：

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$

被称为 *linear regression*. 这个 model 关于 parameter 是 linear 的，关于 input variable 也是 linear 的。一个一般的 linear model 只需要关于 parameter 线性：

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

其中的 $\phi_i(\mathbf{x})$ 可以是 non-linear 的，被称为 basis function. 这里的 w_0 是 'bias'，我们也可以改写成：

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

3.1.1 Maximum likelihood and least squares

假设target t 有确定的函数 $y(\mathbf{x}, \mathbf{w})$ 给出:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon$$

其中 ϵ 是一个服从mean是0,precision是 β 的gaussian random noise.可以写成:

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

对于一个dataset, $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, target t_1, \dots, t_N , 基于i.i.d Gaussian的假设, 有:

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

求使该表达式最大的 \mathbf{w} 等价于least square:

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

可以解出:

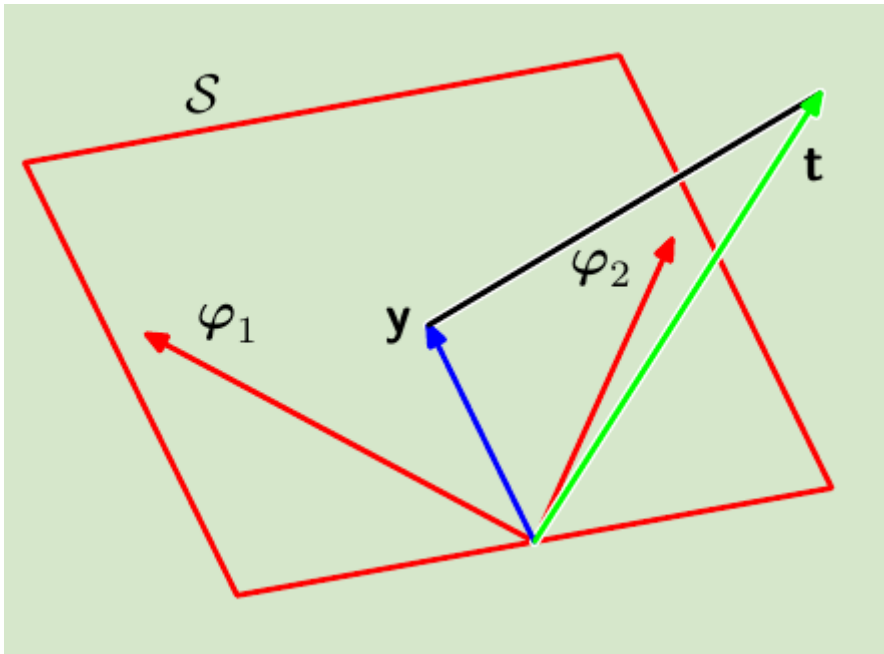
$$\mathbf{w}_{\text{ML}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t}$$

这被称为normal equations for least squares problem, $\boldsymbol{\Phi}$ 被称为design matrix.其中

$$\boldsymbol{\Phi}^\dagger \equiv (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T$$

被称为Moore-Penrose pseudo-inverse.

3.1.2 Geometry of least squares



在一个 N -dimensional space (axes是 t_1, \dots, t_N) , least-square regression function就是找一个orthogonal projection, 把data vector \mathbf{t} 投影到由basis function所span得到的subspace上。

3.1.3 Sequential learning

当data多的时候, 用Sequential method就变得值得。假设我们现在要minimize SSE: $E = \sum_n E_n$, 使用stochastic gradient descent(*sequential gradient descent*)求最优parameters \mathbf{w} 的一般形式:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

对于SSE最小的问题来说就是:

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

这也被称为*least-mean-squares*(LMS) algorithm.

3.1.4 Regularized least squares

在第一章的时候已经谈过在error function上增加regularization term防止over-fitting, 总的error function是 $E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$, 其中的 λ 用来控制两者的relative importance. 一个SSE问题带正则项的最简单形式是:

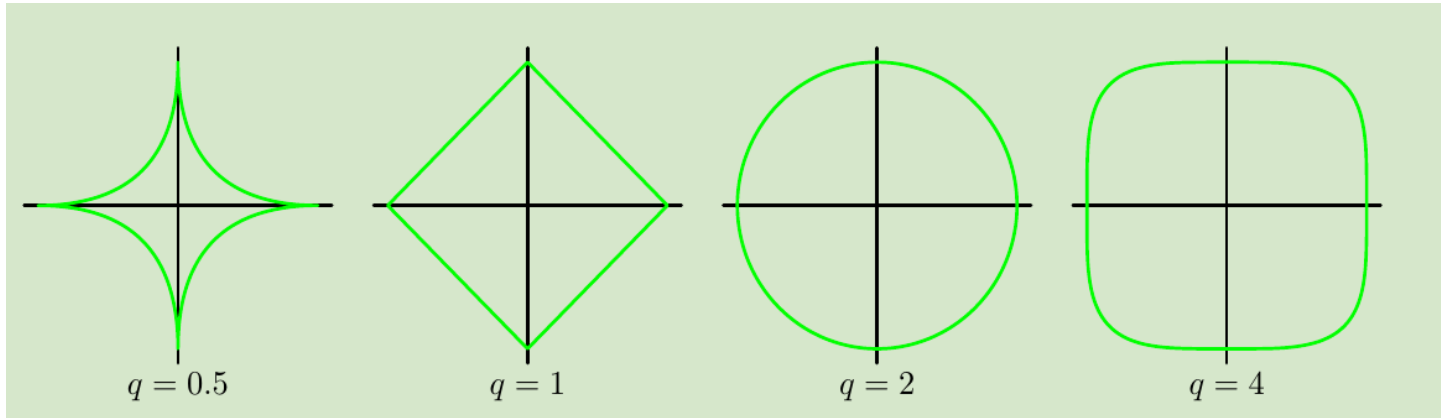
$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^\top \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

这个问题有closed-form solution:

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

更一般的正则化是：

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



对 q 取不同值时, *regular term*的contours

3.1.5 Multiple outputs

当我们需要predict $K > 1$ 的target的时候, 我们可以选择对每个 K 选择相同的basis function:

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

假设conditional distribution是isotropic Gaussian:

$$p(\mathbf{t} | \mathbf{x}, \mathbf{W}, \beta) = \mathcal{N}(\mathbf{t} | \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}), \beta^{-1} \mathbf{I})$$

同样maximum likelihood可以得到解:

$$\mathbf{W}_{\text{ML}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{T}$$

3.2. The Bias-Variance Decomposition

在已知 $p(t|\mathbf{x})$ 下, $h(\mathbf{x})$ 是最优Regression function:

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

但实际中并不知道 $p(t|\mathbf{x})$, 我们可以让expected Loss 最小来选择model $y(\mathbf{x})$.

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

这里面第二项与model无关，是data上的intrinsic noise.

假设我们在data \mathcal{D} 下得到一个prediction function $y(\mathbf{x}; \mathcal{D})$ ，那么expected loss中第一项的平方损失可以写成：

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

若关于data求expectation：

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

如此可以把开始的式子分解成三个部分：

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

where

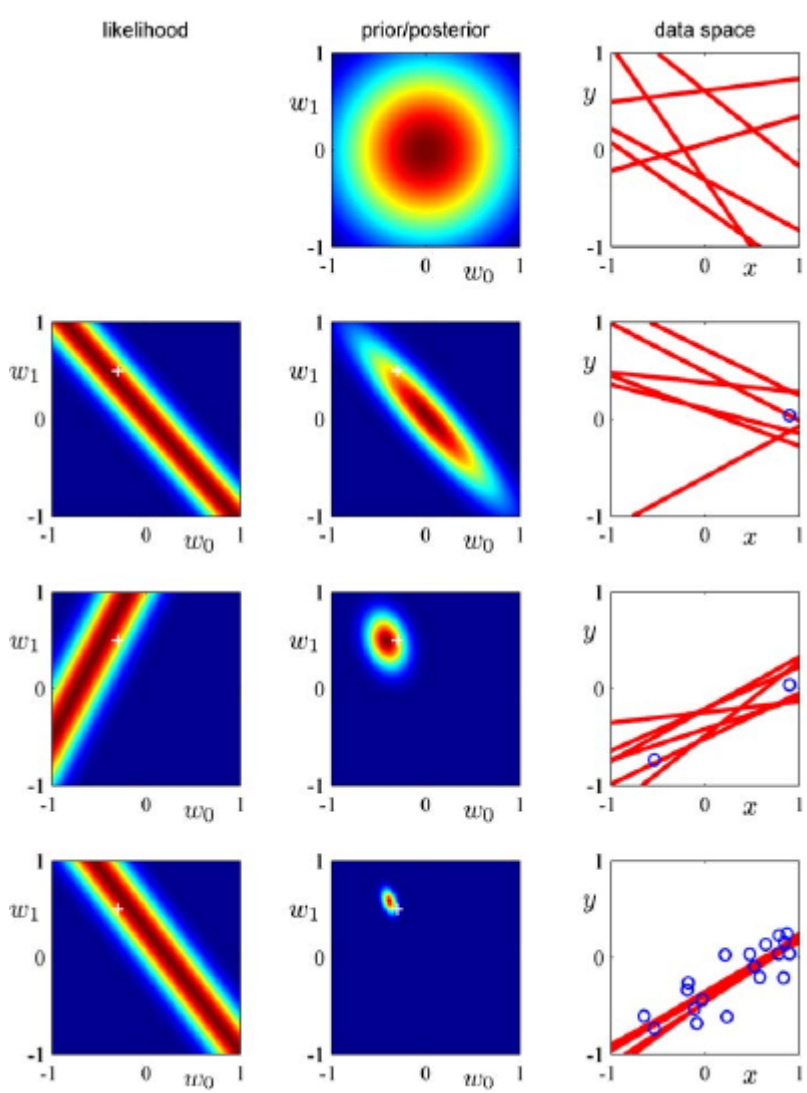
$$\begin{aligned} (\text{bias})^2 &= \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \\ \text{variance} &= \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x} \\ \text{noise} &= \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \end{aligned}$$

所以最终的问题是如何在**bias**和**variance**之间找一个balance，这是一个trade-off的问题。具体来说就是控制模型的complex和拟合效果。

3.3. Bayesian Linear Regression

3.3.1 Parameter distribution

考虑一个最简单的拟合线性函数 $y(x, \mathbf{w}) = w_0 + w_1 x$ 的例子.



仍然使用Gaussian和对应的共轭先验。使用sequential method估计posterior。图中可以看到随着data point增加, posterior越来越sharp.

3.3.2 Predictive distribution

实际中, 我们感兴趣的是 $p(t|\mathbf{t}, \alpha, \beta)$, 该predictive distribution的variance为

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

随着data point增加, 第二项会逐渐趋向0. 并且data point附近的variance会更小。

3.3.3 Equivalent kernel

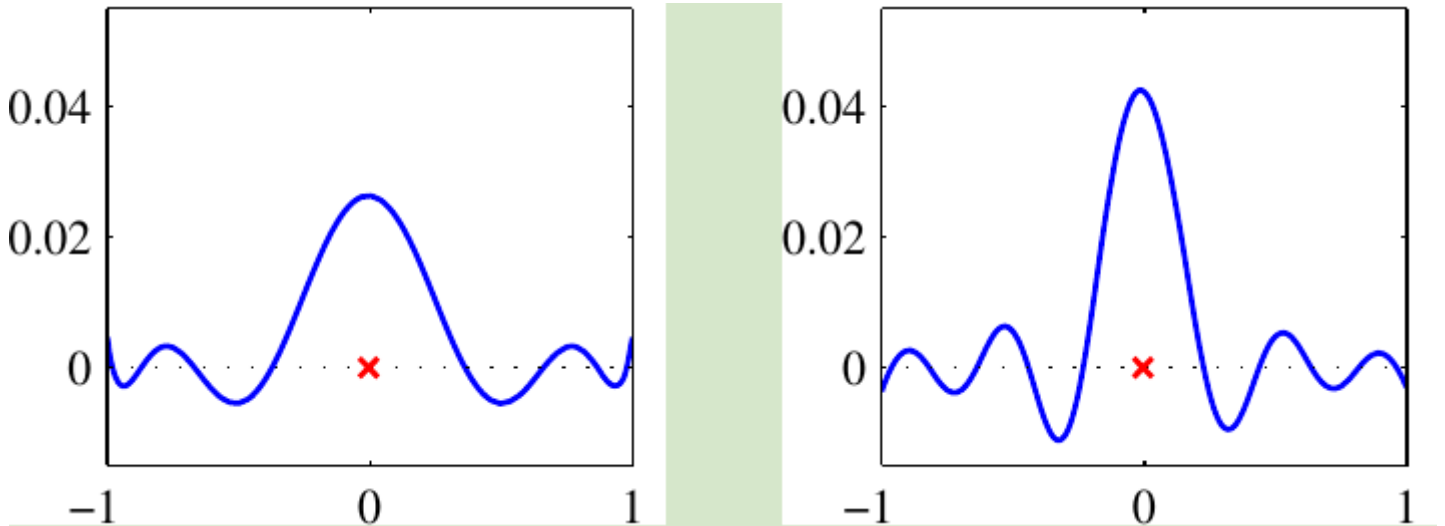
predictive mean(\mathbf{w} 取posterior的mean)能够写成:

$$y(\mathbf{x}, \mathbf{m}_N) = \mathbf{m}_N^T \phi(\mathbf{x}) = \beta \phi(\mathbf{x})^T \mathbf{S}_N \Phi^T \mathbf{t} = \sum_{n=1}^N \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n) t_n$$

$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

其中 $k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}')$ 被称为 *smoother matrix* or *equivalent kernel*.

note: K 做的就是对所有 train data 里的 target 做一个 weighted sum. 也可以称为 *linear smoothers*.



左边的是 basis function 取 Gaussian, 右边是取 sigmoid, 可以看到最终的 equivalent kernel 是相似的, 都是在 \mathbf{x} 附近的 data point 具有更大的 weight.

3.4. Bayesian Model Comparison

假设有不同的 models $\{\mathcal{M}_i\}$, evaluate 不同 model 的 posterior:

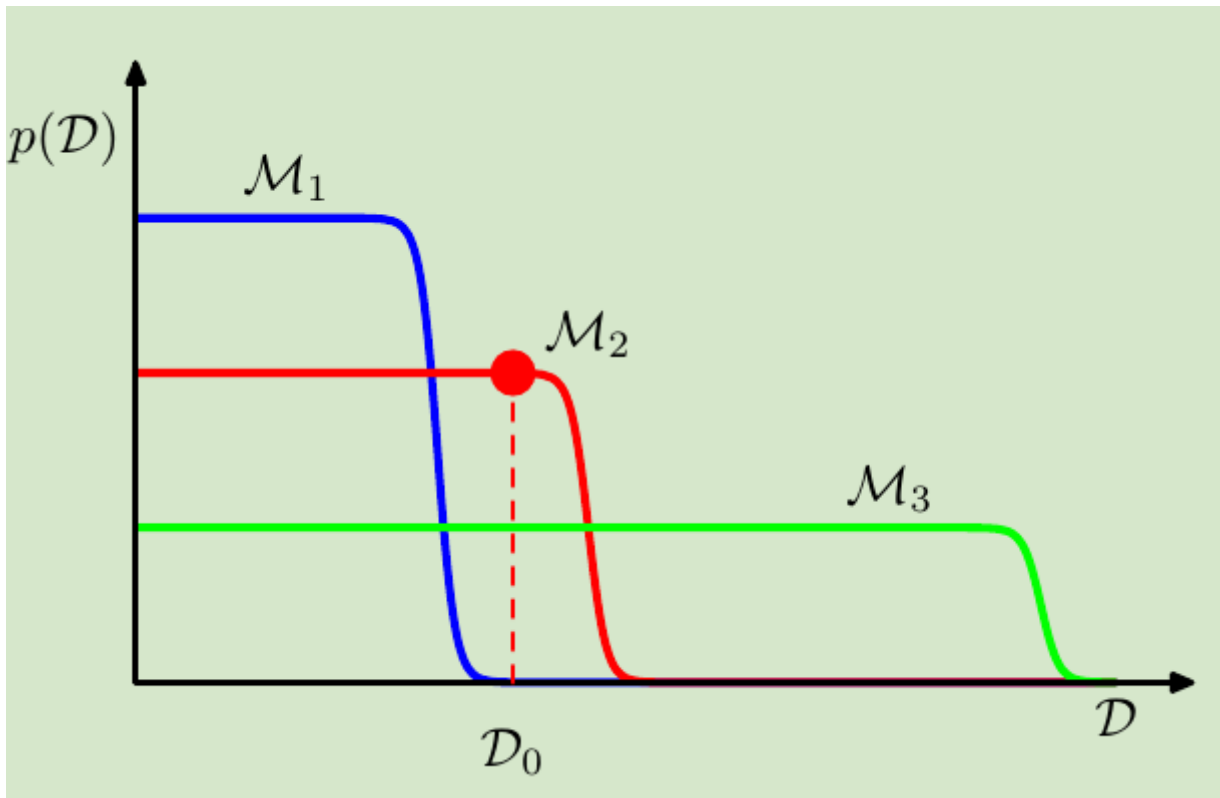
$$p(\mathcal{M}_i | \mathcal{D}) \propto p(\mathcal{M}_i) p(\mathcal{D} | \mathcal{M}_i)$$

帮助我们 model selection/model averaging. 其中的 $p(\mathcal{M}_i)$ 是 prior, $p(\mathcal{D} | \mathcal{M}_i)$ 是 model evidence. model evidence 有时也称 marginal likelihood.

note: marginal 体现在:

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

这其中 \mathcal{M}_i 是 hyper parameter, \mathbf{w} 是确定某个 model 类型后的 parameter. 我们要做的 model comparison 的核心就是 model evidence $p(\mathcal{D} | \mathcal{M}_i)$.



假设我们观测到的data是 \mathcal{D}_0 ，现在有三种model $\{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3\}$, complexity依次增加。图中可以看到，小的model complexity只在小部分类型的data里可能性大，而最complex的模型虽然表达能力强，在各种data下都有概率，但是由于 $p(\mathcal{D})$ 是normalize的，所以在某个具体观测data下，适当模型complex的是最可能的。（这也就说明了Bayes视角下对待over-fit的手段）

3.5. The Evidence Approximation

由于fully Bayesian treatment需要对parameter \mathbf{w} 和hyper-parameter α 、 β 积分，这是intractable.我们考虑一种近似手段，将hyper parameter的值由maximizing marginal likelihood先确定.predictive function为：

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

接着approximate：

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

其中 $\hat{\alpha}$ 、 $\hat{\beta}$ 是 $p(\alpha, \beta | \mathbf{t})$ 的尖峰.