

模式识别引论

An Introduction to Pattern Recognition

主讲: 李春光

www.pris.net.cn/teacher/lichunguang

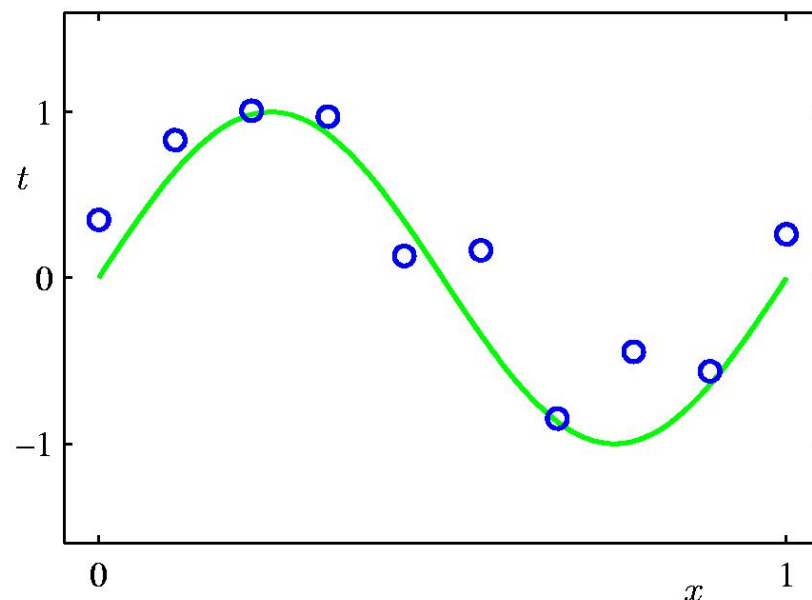
模式识别与智能系统实验室

网络搜索教研中心 信息与通信工程学院 北京邮电大学

引例：多项式曲线拟合

- 给定N个训练数据: (x, t)

- 数据其中绿色曲线为生成训练的真实曲线



- 多项式曲线拟合

- 使用多项式模型去构造
 - 定义为

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

引例：多项式曲线拟合

- 学习任务:

- 根据训练数据，估计模型中的参数

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

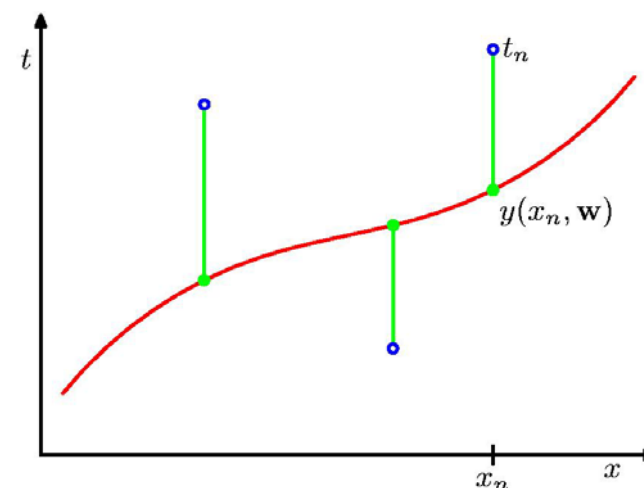
- 根据目标函数的不同，分为:

- 最小二乘法

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

- 正则化最小二乘法

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$



最大似然估计法(MLE)

- 似然函数为

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

- 对数似然函数为

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = - \underbrace{\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

- 通过最大化似然函数估计参数 \mathbf{w}_{ML} 和 β_{ML}

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

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2$$

最大后验概率(MAP)估计法

- 先验分布密度

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

- 似然函数

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

- 后验分布密度

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

- 等价于正则化的最小二乘法

$$\beta\tilde{E}(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2}\mathbf{w}^T\mathbf{w}$$

贝叶斯估计法(Bayesian Estimation)

- 先验分布密度

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

- 似然函数

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

- 后验分布密度

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

- 回归模型的预测性分布密度

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

其中

$$m(x) = \beta\phi(x)^T \mathbf{S} \sum_{n=1}^N \phi(x_n)t_n \quad s^2(x) = \beta^{-1} + \phi(x)^T \mathbf{S} \phi(x)$$

$$\mathbf{S}^{-1} = \alpha\mathbf{I} + \beta \sum_{n=1}^N \phi(x_n)\phi(x_n)^T \quad \phi(x_n) = (x_n^0, \dots, x_n^M)^T$$

贝叶斯线性回归 内容提要

- 引子: 曲线拟合问题
- 常用的几种分布
- 贝叶斯线性回归
 - 两个例子
 - 等价核
 - 先验分布中的超参数的处理

几种常用的分布

- 伯努利(Bernoulli)分布
- 二项 (Binomial)分布
- Beta分布
- 多项分布
- 狄利克雷(Dirichlet)分布
- 高斯分布

二值(Binary)变量的伯努利分布

- **伯努利(Bernoulli)分布**

- 比如，投一枚硬币观察结果为**Head or Tail**,
Head =1, Tail =0

$$p(x = 1|\mu) = \mu$$

$$\text{Bern}(x|\mu) = \mu^x(1 - \mu)^{1-x}$$

$$\mathbb{E}[x] = \mu$$

$$\text{var}[x] = \mu(1 - \mu)$$


伯努利分布参数的MLE

- Given

$$\mathcal{D} = \{x_1, \dots, x_N\}, \text{ } m \text{ heads (1), } N - m \text{ tails (0)}$$

$$p(\mathcal{D}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n}$$

$$\ln p(\mathcal{D}|\mu) = \sum_{n=1}^N \ln p(x_n|\mu) = \sum_{n=1}^N \{x_n \ln \mu + (1 - x_n) \ln(1 - \mu)\}$$


$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n = \frac{m}{N}$$

最大似然估计 (MLE) 的过拟合现象

- 投一枚硬币，观察其落地后是正面或反面
- 假如N次观测全部为正面：

$$\mathcal{D} = \{1, 1, 1\} \rightarrow \mu_{\text{ML}} = \frac{3}{3} = 1$$

- 则可以给出的预测是 后续所有观测全部为正面！
- 原因：**MLE**估计存在对数据集的**Overfitting**

二项(Binomial)分布

- 二项(**Binomial**)分布
 - 比如 投一枚硬币**N**次，出现**m**次**Head**朝上

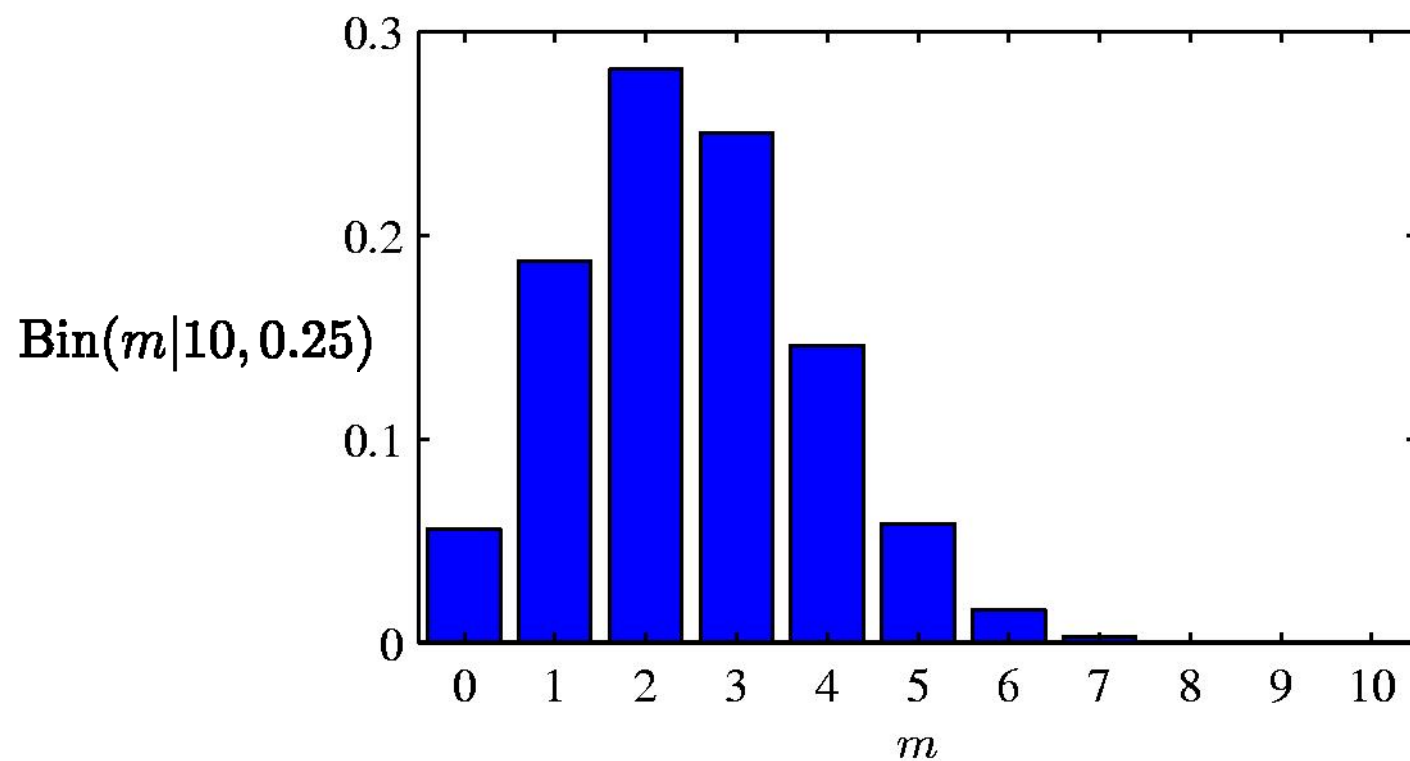
$$p(m \text{ heads} | N, \mu)$$

$$\text{Bin}(m | N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

$$\mathbb{E}[m] \equiv \sum_{m=0}^N m \text{Bin}(m | N, \mu) = N\mu$$

$$\text{var}[m] \equiv \sum_{m=0}^N (m - \mathbb{E}[m])^2 \text{Bin}(m | N, \mu) = N\mu(1 - \mu)$$

Binomial Distribution



Beta分布

- Beta分布

- 比如: 二项分布的参数 $\mu \in [0, 1]$ 可假设服从 **Beta** 分布

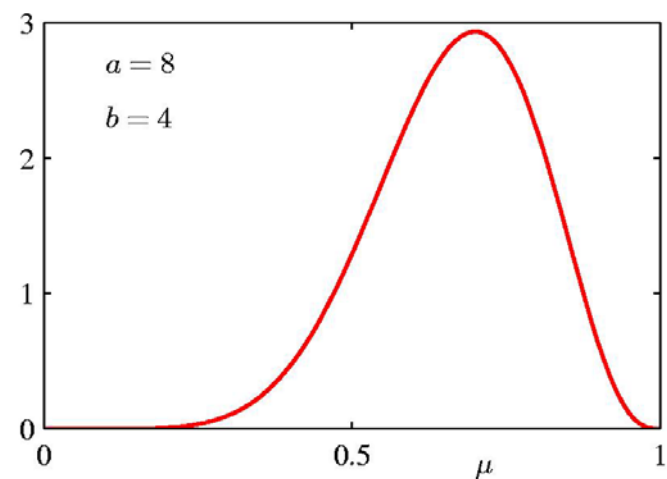
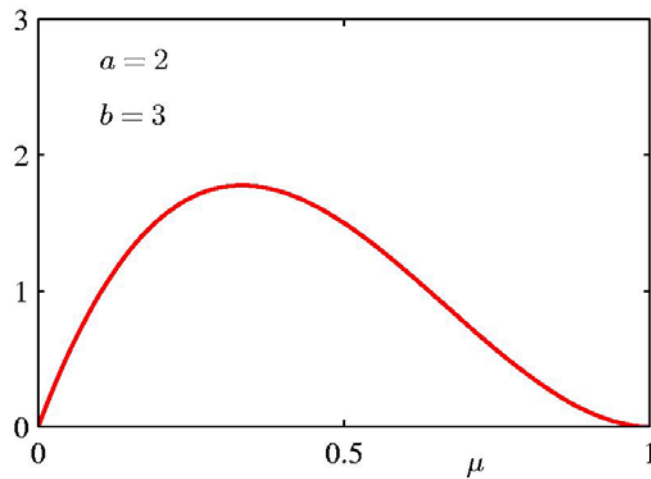
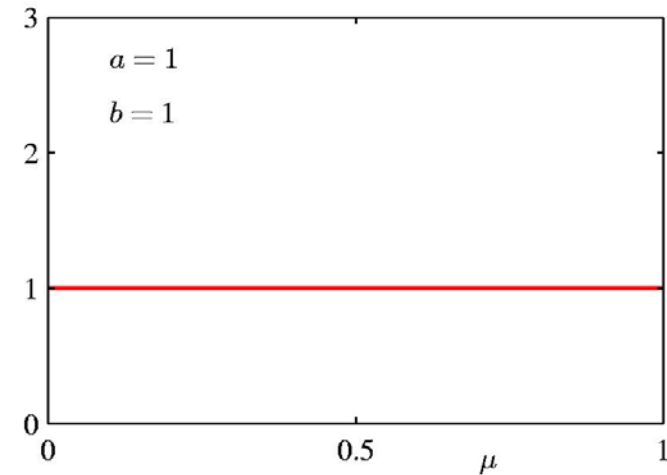
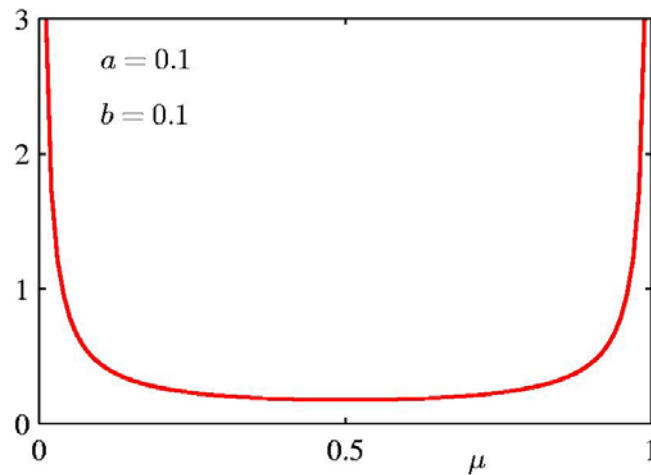
$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

$$\mathbb{E}[\mu] = \frac{a}{a+b}$$

$$\text{var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)}$$

$$\text{其中 } \Gamma(x) \equiv \int_0^\infty u^{x-1} e^{-u} \mathrm{d}u$$

Beta Distribution



Beta分布的应用举例

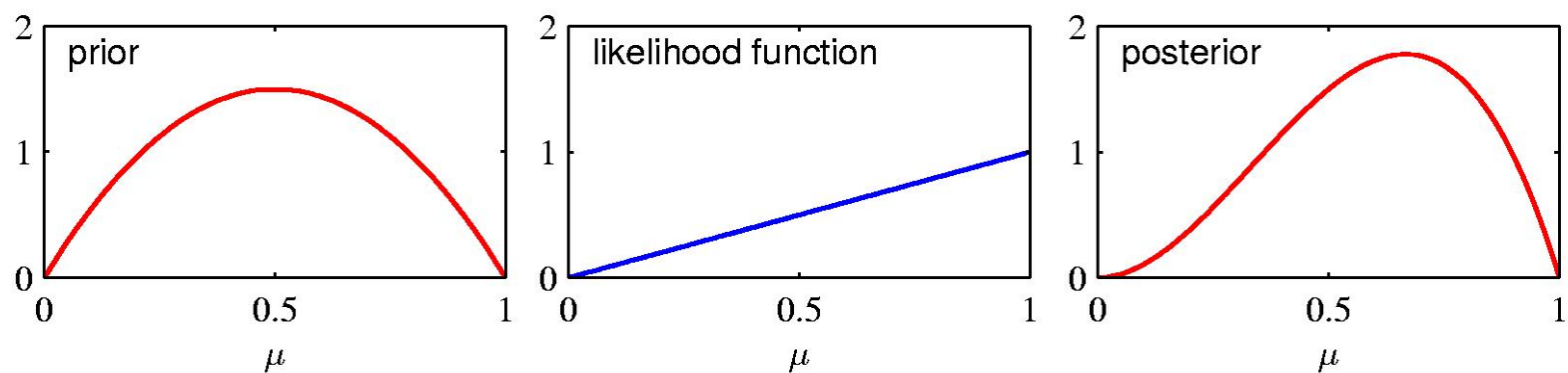
- Beta分布可以充当Binomial分布中参数的共轭先验分布 \rightarrow 参数的后验分布仍为Beta分布

$$\text{Posterior} = \text{Likelihood} \cdot \text{Prior}$$

$$\begin{aligned} \rightarrow p(\mu|a_0, b_0, \mathcal{D}) &\propto p(\mathcal{D}|\mu)p(\mu|a_0, b_0) \\ &= \left(\prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n} \right) \text{Beta}(\mu|a_0, b_0) \\ &\propto \mu^{m+a_0-1} (1 - \mu)^{(N-m)+b_0-1} \\ &\propto \text{Beta}(\mu|a_N, b_N) \end{aligned}$$

$$\text{其中 } a_N = a_0 + m \quad b_N = b_0 + (N - m)$$

Likelihood · Prior = Posterior



后验分布的应用举例

- 考虑投硬币试验. 假设给定基于观测数据集 \mathcal{D} 推理得到参数的后验分布, 请问: 在投硬币过程中, 下一次正面朝上的概率是多少?

$$\begin{aligned} p(x = 1 | a_0, b_0, \mathcal{D}) &= \int_0^1 p(x = 1 | \mu) p(\mu | a_0, b_0, \mathcal{D}) d\mu \\ &= \int_0^1 \mu p(\mu | a_0, b_0, \mathcal{D}) d\mu \\ &= \mathbb{E}[\mu | a_0, b_0, \mathcal{D}] = \frac{a_N}{b_N} \end{aligned}$$



$$p(x = 1 | \mathcal{D}) = \frac{m + a}{m + a + l + b}$$

其中, $a_N = m + a$, $b_N = l + b$

**Bayesian
Style...**

多状态伯努利分布

- 具有K个互斥状态的随机变量的分布
 - 比如：分类/聚类/指派问题中的K维指示向量

$$\mathbf{x} = (0, 0, 1, 0, 0, 0)^T$$

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \quad \forall k : \mu_k \geq 0 \quad \text{and} \quad \sum_{k=1}^K \mu_k = 1$$

$$\mathbb{E}[\mathbf{x}|\boldsymbol{\mu}] = \sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) \mathbf{x} = (\mu_1, \dots, \mu_K)^T = \boldsymbol{\mu}$$

$$\sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) = \sum_{k=1}^K \mu_k = 1$$

多状态伯努利分布中参数的MLE

- 给定数据集 $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})} = \prod_{k=1}^K \mu_k^{m_k}$$

$$\text{s.t. } \sum_k \mu_k = 1$$

使用Lagrange法



$$\sum_{k=1}^K m_k \ln \mu_k + \lambda \left(\sum_{k=1}^K \mu_k - 1 \right)$$



$$\mu_k = -m_k / \lambda \rightarrow \mu_k^{\text{ML}} = \frac{m_k}{N}$$

多项分布

- 多项(Multinomial)分布

- 比如 投一枚骰子N次，其中出现k点 m_k 次， $k=1,\dots,K$.

- $$\text{Mult}(m_1, m_2, \dots, m_K | \boldsymbol{\mu}, N) = \binom{N}{m_1 m_2 \dots m_K} \prod_{k=1}^K \mu_k^{m_k}$$

$$\mathbb{E}[m_k] = N\mu_k$$

$$\text{var}[m_k] = N\mu_k(1 - \mu_k)$$

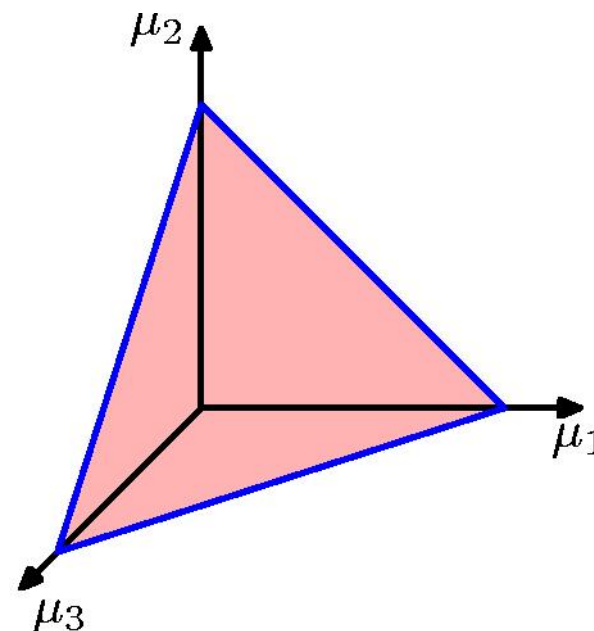
$$\text{cov}[m_j, m_k] = -N\mu_j\mu_k$$

狄利克雷(Dirichlet)分布

- 狄利克雷分布
 - 比如：多项分布中参数的先验分布可假设服从狄利克雷分布

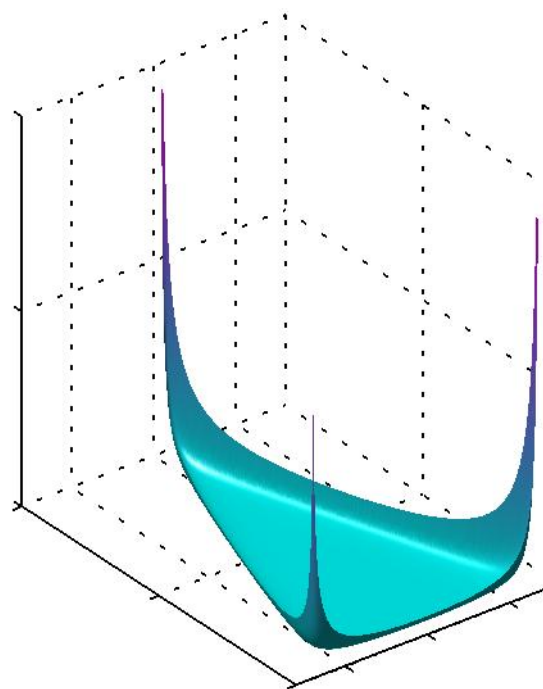
$$\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k-1}$$

$$\text{其中 } \alpha_0 = \sum_{k=1}^K \alpha_k$$

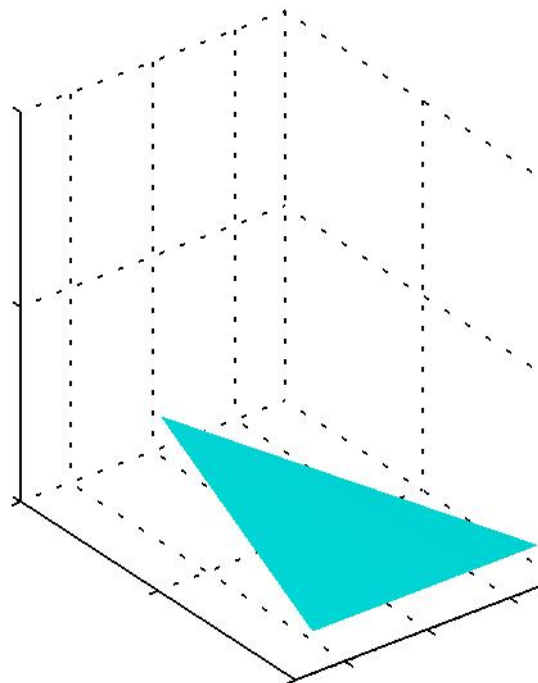


Dirichlet 分布

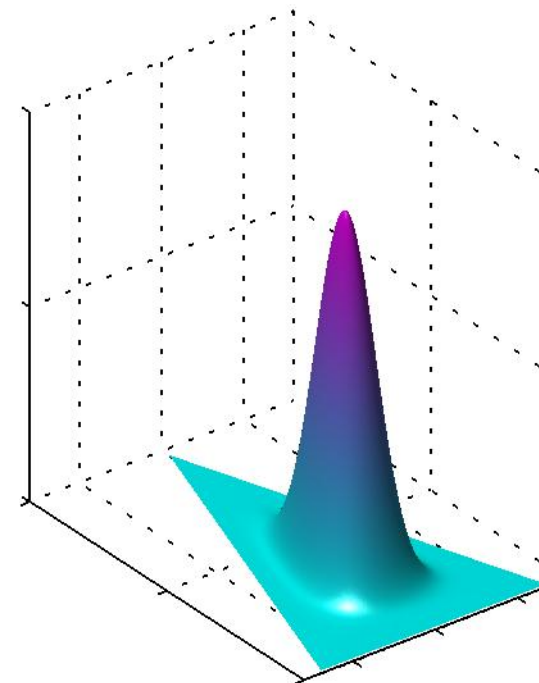
- 3个变量的Dirichlet分布



$$\alpha_k = 10^{-1}$$



$$\alpha_k = 10^0$$



$$\alpha_k = 10^1$$

Dirichlet分布的应用举例

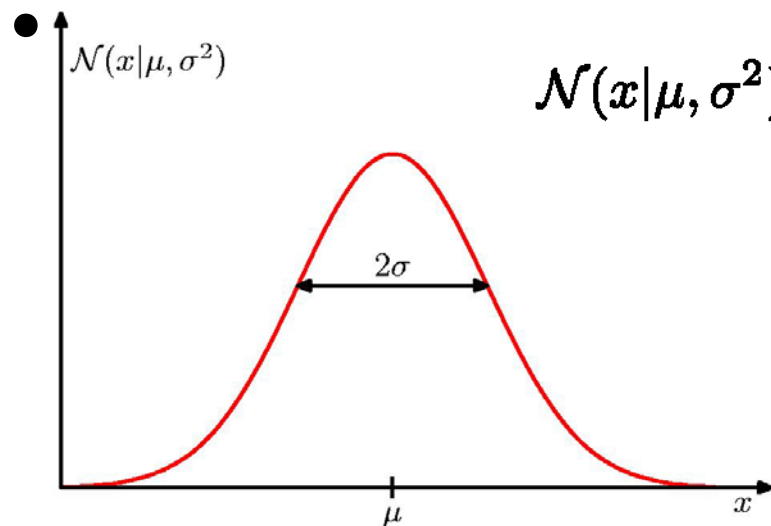
- **Dirichlet**分布可以充当多项分布中参数的共轭先验分布
 - 参数的后验分布仍为**Multinomial**分布

$$\text{Posterior} = \text{Likelihood} \cdot \text{Prior}$$

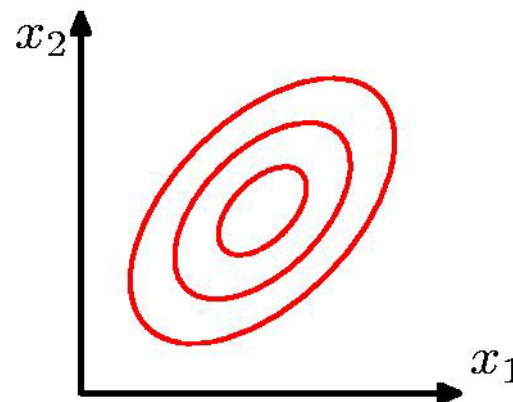
$$\longrightarrow p(\boldsymbol{\mu}|\mathcal{D}, \boldsymbol{\alpha}) \propto p(\mathcal{D}|\boldsymbol{\mu})p(\boldsymbol{\mu}|\boldsymbol{\alpha}) \propto \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1}$$

$$\begin{aligned} \longrightarrow p(\boldsymbol{\mu}|\mathcal{D}, \boldsymbol{\alpha}) &= \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha} + \mathbf{m}) \\ &= \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + m_1) \cdots \Gamma(\alpha_K + m_K)} \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1} \end{aligned}$$

高斯(Gaussian)分布



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



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

多变量/多元高斯分布的几何

- 多元高斯分布中的去相关

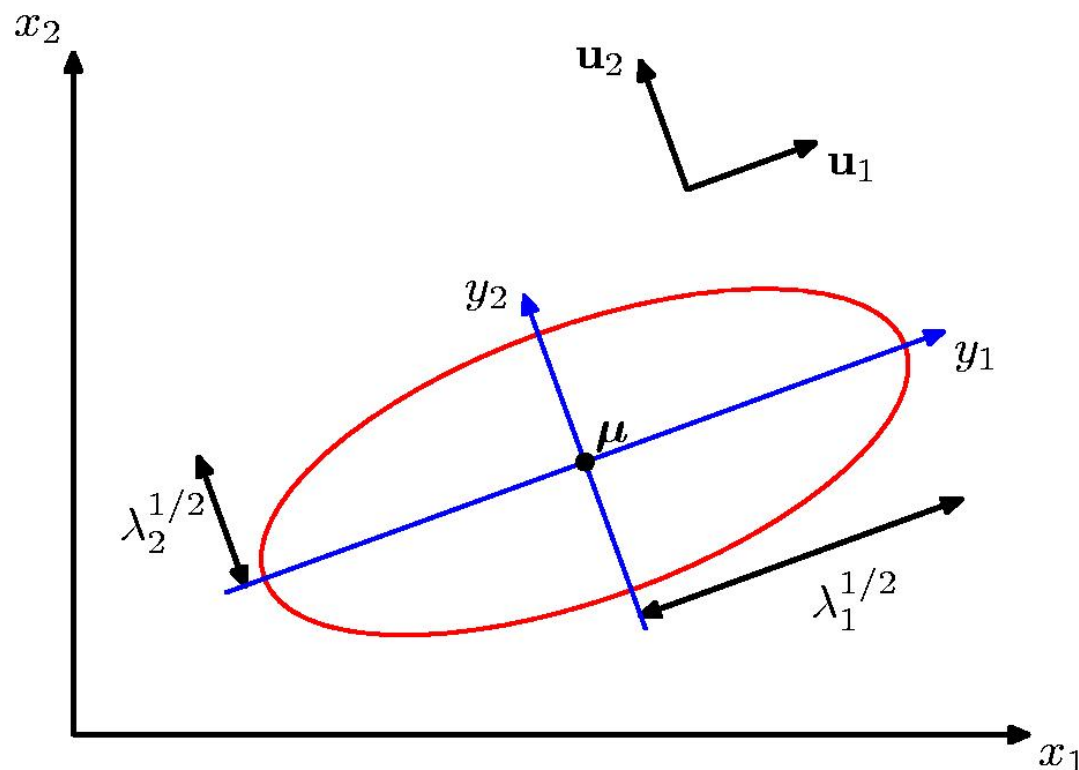
$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

- 马氏距离(Mahalanobis距离)

$$\boldsymbol{\Sigma}^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$$

$$\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$$

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



高斯分布中的参数估计

- 最大似然估计法(Maximal Likelihood Estimation)

- 给定**N**个数据点的数据集**X**，假设**i.i.d.**, **e.g.** 高斯分布


- independent and identically distributed: i.i.d.

- 数据集由特定参数下的高斯分布生成的概率为

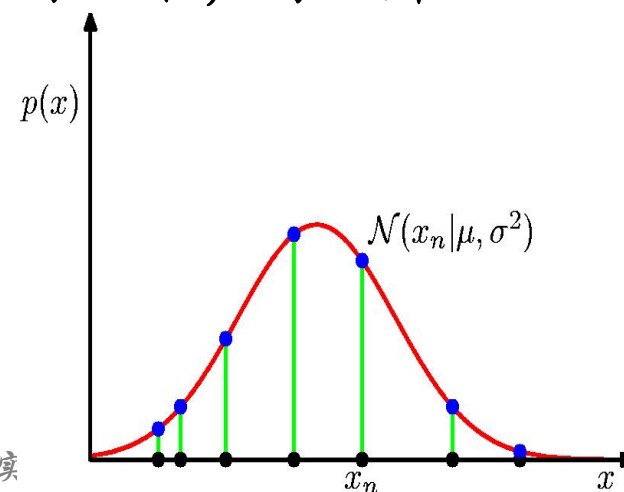
$$p(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2)$$

- 称为似然函数(Likelihood function)

- 最大似然估计原则: 在给定数据的情况下，寻找最大化似然函数的参数


$$\max_{\mu, \sigma^2} L(\mu, \sigma^2 | X) = p(X | \mu, \sigma^2)$$

$$\max_{\mu, \sigma^2} l(\mu, \sigma^2 | X) = \log p(X | \mu, \sigma^2)$$



高斯分布参数的最大似然估计

- 目标函数:

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

- 均值的估计

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$



- 方差的估计

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

高斯分布中均值的贝叶斯估计

- 假设方差已知, 给定 i. i. d. 数据 $\mathbf{x} = \{x_1, \dots, x_N\}$ 则, 似然函数为:

$$p(\mathbf{x}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \right\}.$$

– 假设均值服从高斯分布 $p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2)$.

– 后验分布为: Posterior = Likelihood · Prior

$$p(\mu|\mathbf{x}) \propto p(\mathbf{x}|\mu)p(\mu).$$

$$\longrightarrow p(\mu|\mathbf{x}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$$

高斯分布中均值的贝叶斯估计

- 后验分布

$$p(\mu|\mathbf{x}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$$

— 其中

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\mu_{\text{ML}}, \quad \mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$

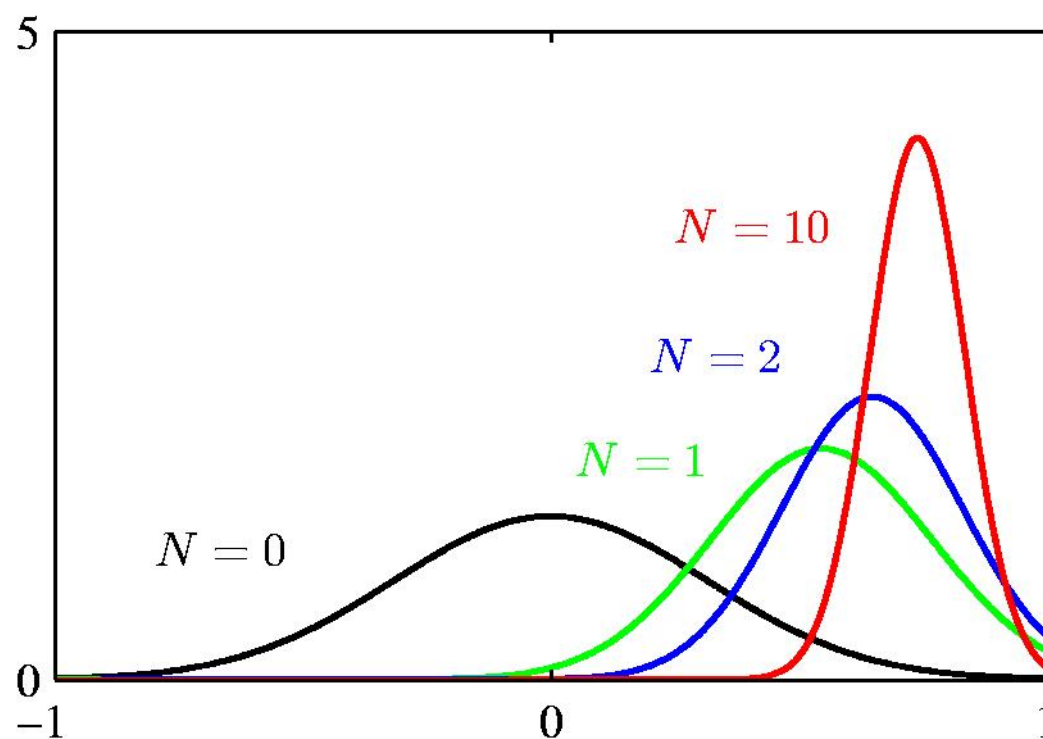
$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}.$$

	$N = 0$	$N \rightarrow \infty$
μ_N	μ_0	μ_{ML}
σ_N^2	σ_0^2	0

高斯分布中均值的贝叶斯估计

- 不同训练样本数目对均值估计的影响

$$p(\mu|\mathbf{x}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$$



高斯分布中均值的贝叶斯估计

- 序贯估计(Sequential Estimation)

$$\begin{aligned} p(\mu|\mathbf{x}) &\propto p(\mu)p(\mathbf{x}|\mu) \\ &= \left[p(\mu) \prod_{n=1}^{N-1} p(x_n|\mu) \right] p(x_N|\mu) \\ &\propto \mathcal{N}(\mu|\mu_{N-1}, \sigma_{N-1}^2) p(x_N|\mu) \end{aligned}$$

— 当第**N+1**个样本到来时，旧后验概率(基于**N**个样本的估计)变成“先验”

高斯分布中方差的贝叶斯估计

- 假设均值已知, 给定 i. i. d. 数据 $\mathbf{x} = \{x_1, \dots, x_N\}$ 则, 似然函数为:

$$p(\mathbf{x}|\lambda) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \lambda^{-1}) \propto \lambda^{N/2} \exp \left\{ -\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}.$$

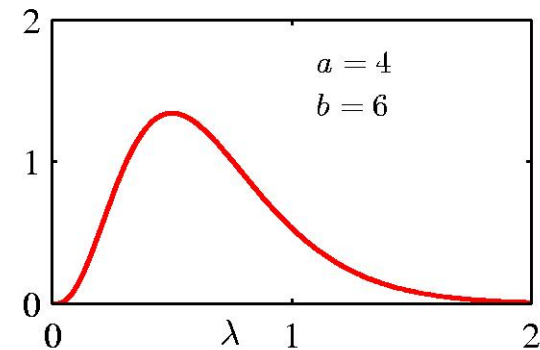
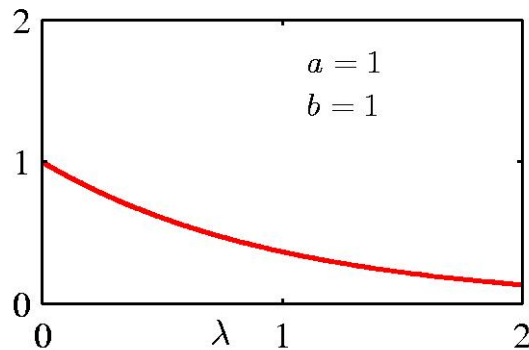
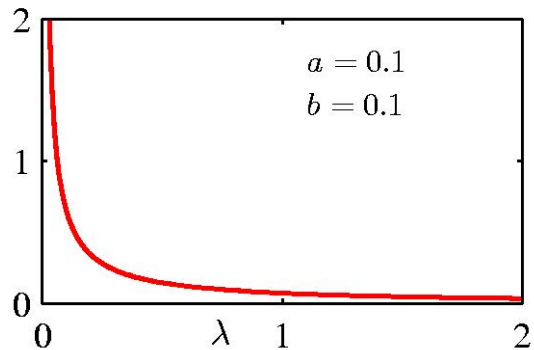
- 注意: 似然函数具有**Gamma**分布的形式
- 代替直接估计方差, 这里估计精度(**Precision**)参数

Gamma分布

- Gamma分布具有下述形式:

$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$

$$\mathbb{E}[\lambda] = \frac{a}{b} \qquad \text{var}[\lambda] = \frac{a}{b^2}$$




高斯分布中方差的贝叶斯估计

- 假设精度参数服从Gamma分布

$$\text{Gam}(\lambda|a_0, b_0)$$

— 则后验分布为

$$p(\lambda|\mathbf{x}) \propto \lambda^{a_0-1} \lambda^{N/2} \exp \left\{ -b_0 \lambda - \frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$$

 $\text{Gam}(\lambda|a_N, b_N)$

$$a_N = a_0 + \frac{N}{2}$$

$$b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{\text{ML}}^2.$$

高斯分布中均值-方差的贝叶斯估计

- 假设均值和方差均未知, 给定 i. i. d. 数据

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

则联合似然函数为:

$$p(\mathbf{x}|\mu, \lambda) = \prod_{n=1}^N \left(\frac{\lambda}{2\pi}\right)^{1/2} \exp\left\{-\frac{\lambda}{2}(x_n - \mu)^2\right\}$$
$$\propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^N \exp\left\{\lambda\mu \sum_{n=1}^N x_n - \frac{\lambda}{2} \sum_{n=1}^N x_n^2\right\}.$$

如何选择先验?

高斯分布中均值-方差的贝叶斯估计

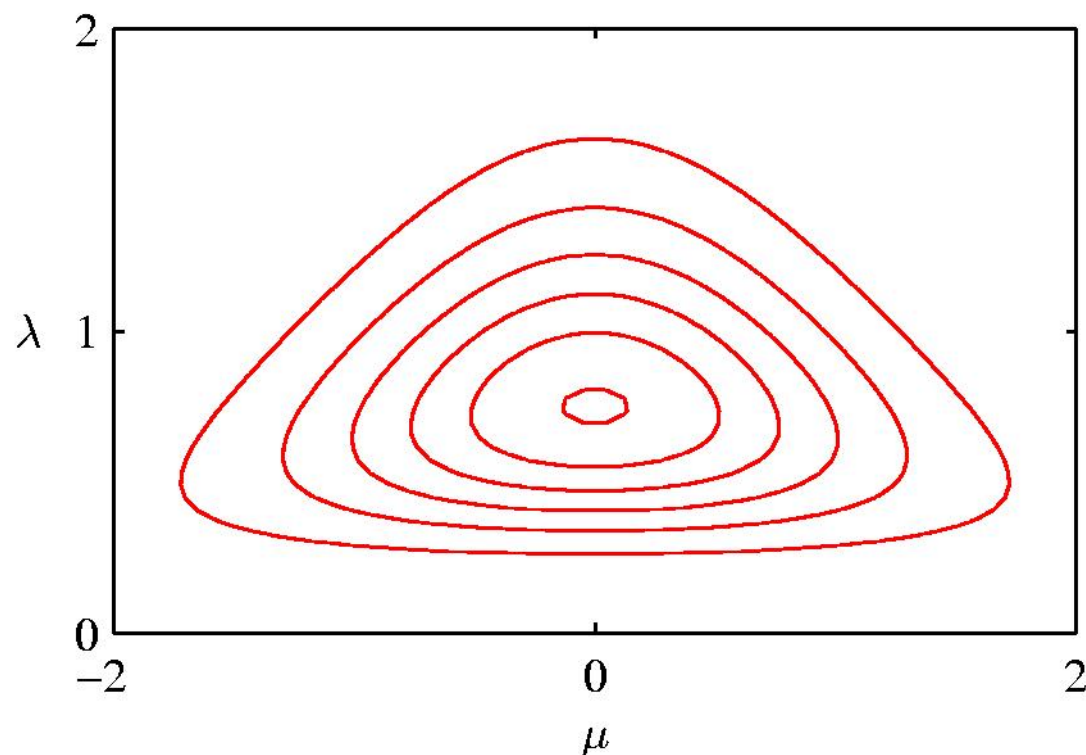
- 高斯-伽玛(Gaussian-Gamma)先验分布

$$p(\mu, \lambda) = \mathcal{N}(\mu | \mu_0, (\beta\lambda)^{-1}) \text{Gam}(\lambda | a, b)$$
$$\propto \underbrace{\exp \left\{ -\frac{\beta\lambda}{2} (\mu - \mu_0)^2 \right\}}_{\text{Quadratic in } \mu} \underbrace{\lambda^{a-1} \exp \{-b\lambda\}}_{\text{Gamma distribution over } \lambda}$$

- Quadratic in μ .
- Linear in λ .
- Gamma distribution over λ .
- Independent of μ .

Gaussian-Gamma分布

- $p(\mu, \lambda) = \mathcal{N}(\mu|\mu_0, (\beta\lambda)^{-1})\text{Gam}(\lambda|a, b)$
 $\propto \exp\left\{-\frac{\beta\lambda}{2}(\mu - \mu_0)^2\right\} \lambda^{a-1} \exp\{-b\lambda\}$



多元高斯分布的共轭先验

- 多变量高斯分布的均值/协方差矩阵的共轭先验(**Conjugate priors**)

- 假设协方差矩阵已知，则均值向量的共轭先验为高斯分布

- 假设均值向量已知，则协方差矩阵的共轭先验是威沙特(Wishart)分布

$$\mathcal{W}(\Lambda|\mathbf{W}, \nu) = B|\Lambda|^{(\nu-D-1)/2} \exp\left(-\frac{1}{2}\text{Tr}(\mathbf{W}^{-1}\Lambda)\right).$$

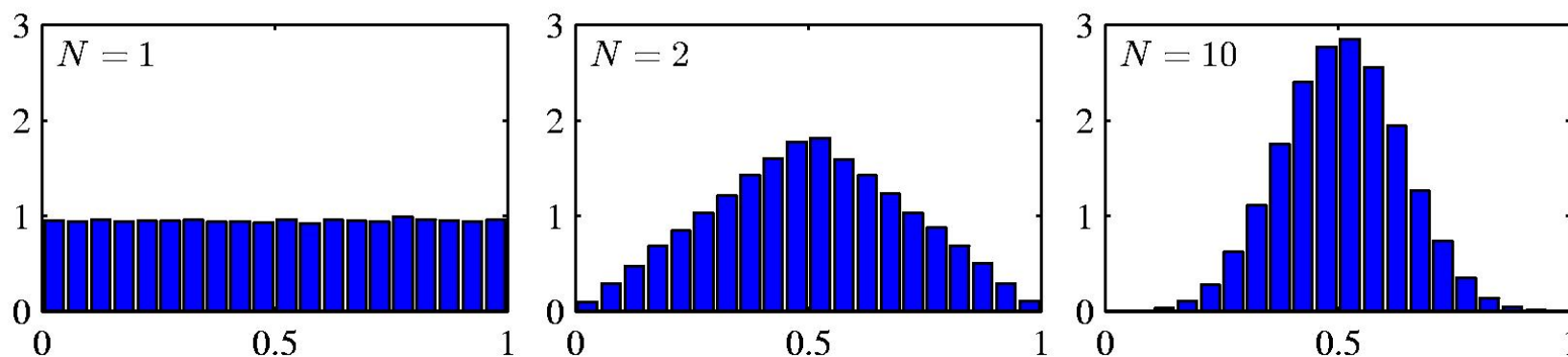
- 假设均值向量未知、协方差矩阵未知，则两者的联合共轭先验分布为高斯-威沙特分布

$$p(\mu, \Lambda|\mu_0, \beta, \mathbf{W}, \nu) = \mathcal{N}(\mu|\mu_0, (\beta\Lambda)^{-1}) \mathcal{W}(\Lambda|\mathbf{W}, \nu)$$

中心极限(Central Limit)定理

- N 个i.i.d. (独立同分布) 随机变量的和的分布，当 N 增大时，趋于高斯分布

— 举例: N 个在 $[0,1]$ 区间内均匀(uniform)分布的随机变量的和的分布

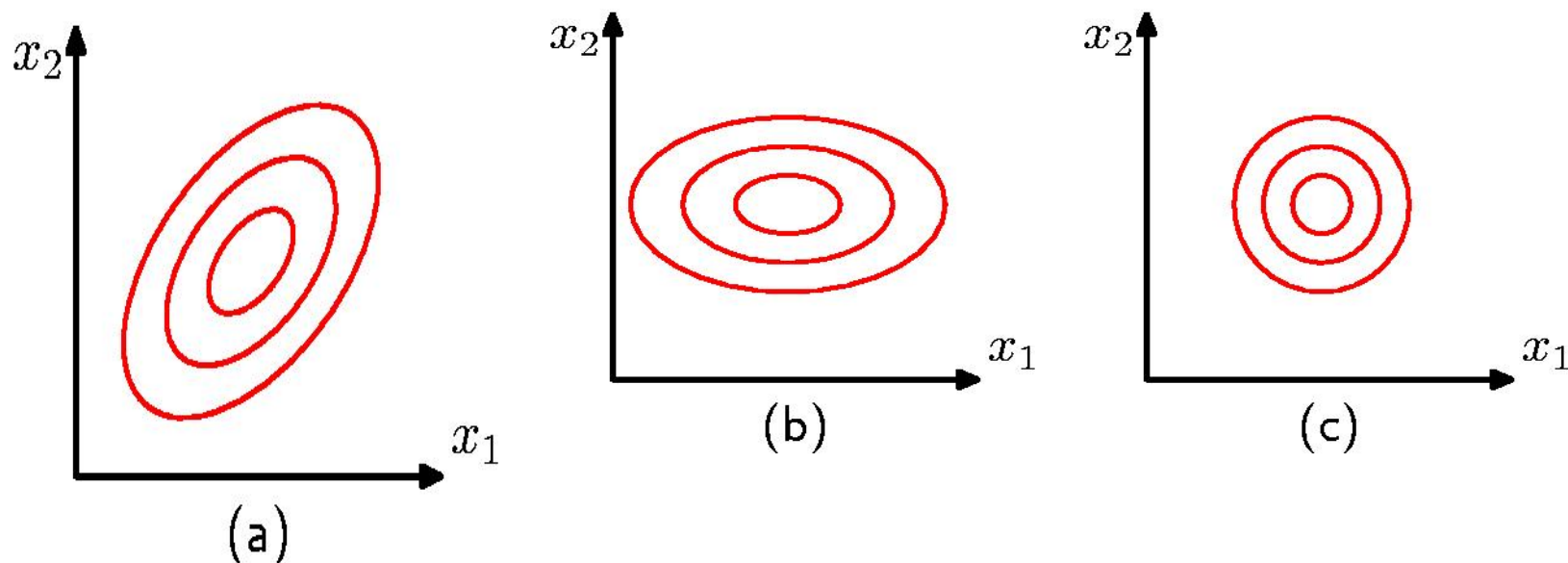


多元高斯分布的矩

- 二阶矩和协方差矩阵

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$$

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



多元高斯随机变量的分块

- 假设 \mathbf{x} 服从多元高斯分布

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

- 令 $\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$ $\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}$ $\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$
 $\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$ $\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$

- 条件分布和边缘分布仍然为高斯分布.

多元高斯随机变量的分块

- 条件分布仍为高斯分布，且

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\mu_{a|b}, \Sigma_{a|b})$$

– 凑全平方项: completing the square

$$-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu) = -\frac{1}{2}\mathbf{x}^T \Sigma^{-1}\mathbf{x} + \mathbf{x}^T \Sigma^{-1}\mu + \text{const}$$

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu) = & \\ & -\frac{1}{2}(\mathbf{x}_a - \mu_a)^T \Lambda_{aa}(\mathbf{x}_a - \mu_a) - \frac{1}{2}(\mathbf{x}_a - \mu_a)^T \Lambda_{ab}(\mathbf{x}_b - \mu_b) \\ & -\frac{1}{2}(\mathbf{x}_b - \mu_b)^T \Lambda_{ba}(\mathbf{x}_a - \mu_a) - \frac{1}{2}(\mathbf{x}_b - \mu_b)^T \Lambda_{bb}(\mathbf{x}_b - \mu_b). \end{aligned}$$

分块矩阵求逆运算的恒等式

- 使用到下列矩阵恒等式

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}$$

– 其中 $\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}$

– 推导: 基于分块矩阵方程的块消除

分块的多元高斯随机变量

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

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

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\Sigma}_{a|b} \{ \boldsymbol{\Lambda}_{aa}\boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \}$$

$$= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)$$

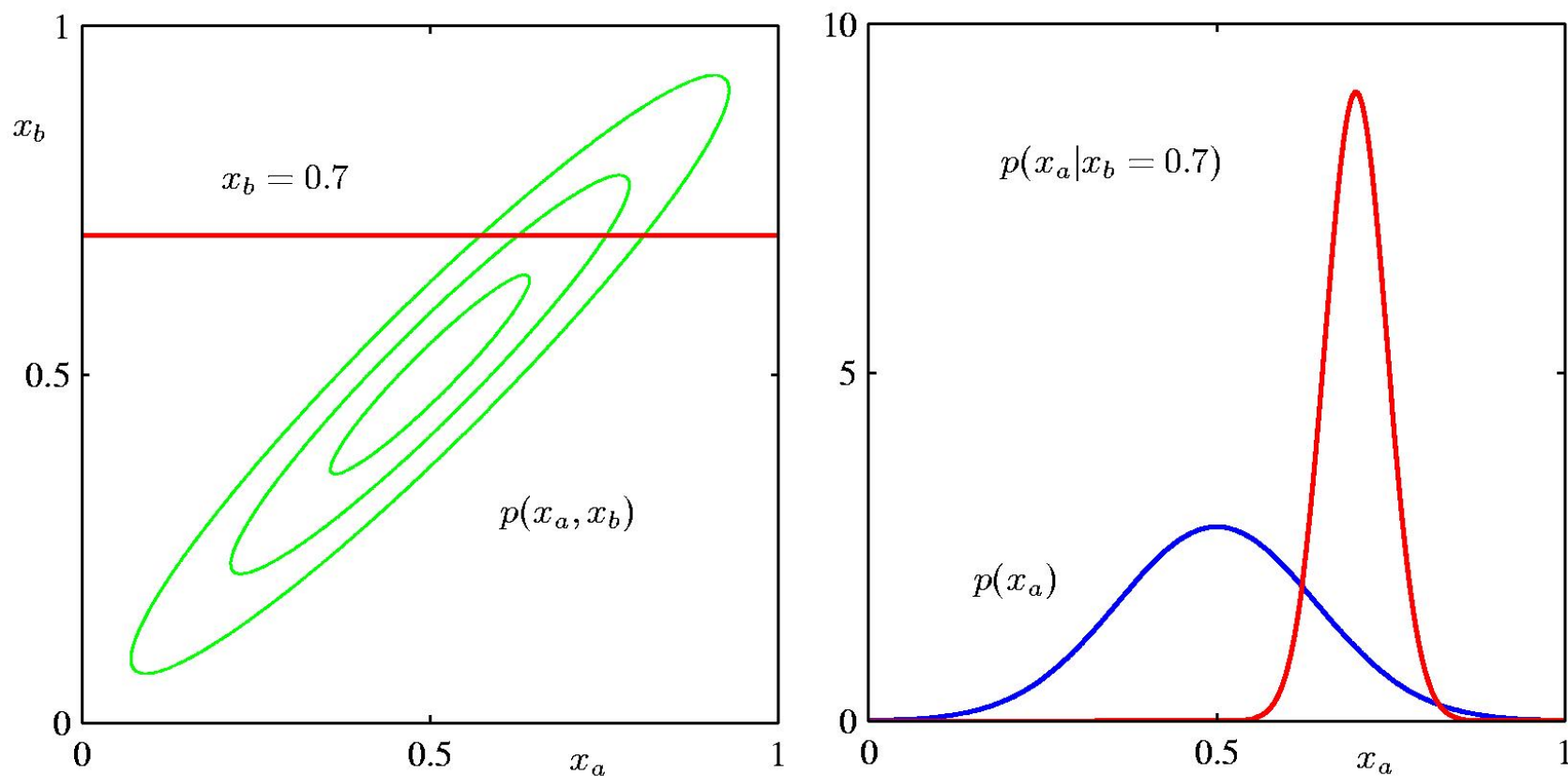
$$= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$$

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

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

分块的条件分布和边缘分布

•



Bayes定理用于高斯随机变量

- 边缘分布密度和条件概率密度

$$\begin{aligned}p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})\end{aligned}$$

$$\begin{aligned}p(\mathbf{y}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \\p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})\end{aligned}$$

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

贝叶斯线性回归 内容提要

- 引子: 曲线拟合问题
- 常用的几种分布
- 贝叶斯线性回归
 - 两个例子
 - 等价核
 - 先验分布中的超参数的处理

贝叶斯线性回归

- 假设观测数据来采样于一个确定性 (**deterministic**) 函数，训练数据中存在i.i.d. 加性高斯噪声(**additive Gaussian noise**)

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{where} \quad p(\epsilon|\beta) = \mathcal{N}(\epsilon|0, \beta^{-1})$$

- 由于噪声的存在， t 的取值具有不确定性，因此我们使用下述高斯分布：

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

THE NEXT:

- 确定先验分布
- 根据训练数据推理后验分布
- 计算回归模型的预测性分布

先验分布

- 共轭先验(Conjugate prior):
 - 一种特定形式的先验分布，它能使后验分布与先验分布具有相同的泛函形式，从而简化后续的贝叶斯分析

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

- 举例:
 - **Bernoulli**分布的共轭先验是**Beta**分布
 - 多项分布的共轭先验是狄利克雷(Dirichlet)分布
 - 高斯分布的共轭先验是高斯分布

贝叶斯线性回归

- 先验分布

- 假设观测数据中存在i.i.d.加性高斯噪声 (additive Gaussian noise)
- 选择共轭先验: $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)$.
 - 一种特定形式的先验分布, 使得后验分布与先验分布具有相同的泛函形式, 从而简化后续的贝叶斯分析

- 后验分布

$$\text{Posterior} = \text{Likelihood} \cdot \text{Prior}$$

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

其中 $\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t})$
 $\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi.$

贝叶斯线性回归

- 通常选择先验分布为

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$$

- 则后验分布为:

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

- 其中 $\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$
 $\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi.$

Bayesian Estimation:
得到的不是参数w本身,
而是w的后验分布

- 后验分布的对数为:

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

Maximum A Posterior: 使用w的后验分布的最大值所对应的w

贝叶斯线性回归 内容提要

- 引子: 曲线拟合问题
- 常用的几种分布
- 贝叶斯线性回归
 - 两个例子
 - 等价核
 - 先验分布中的超参数的处理

考虑一个例子:

- 数据生成:

$$t = f(x, \mathbf{a}) + \varepsilon$$

其中 $x \sim U[-1, 1]$, $f(x, \mathbf{a}) = a_0 + a_1 x = -0.3 + 0.5x$

$$\varepsilon \sim N(0, \sigma^2)$$

- 线性回归模型:

$$t = y(x, \mathbf{w}) + \varepsilon = \mathbf{w}^T x + \varepsilon$$

– 任务: 基于给定数据, 估计线性回归模型

- 如果使用MLE: 即确定参数 w
- 如果使用MAP估计: 即确定参数 w

考虑一个例子:

- 数据生成: $t = f(x, \mathbf{a}) + \varepsilon$

其中 $x \sim U[-1, 1]$, $f(x, \mathbf{a}) = a_0 + a_1 x = -0.3 + 0.5x$

$$\varepsilon \sim N(0, \sigma^2)$$

- 线性回归模型:

$$t = y(x, \mathbf{w}) + \varepsilon = \mathbf{w}^T x + \varepsilon \quad t \sim N(t | \mathbf{w}^T x, \sigma^2)$$

– 任务: 基于给定数据, 估计线性回归模型

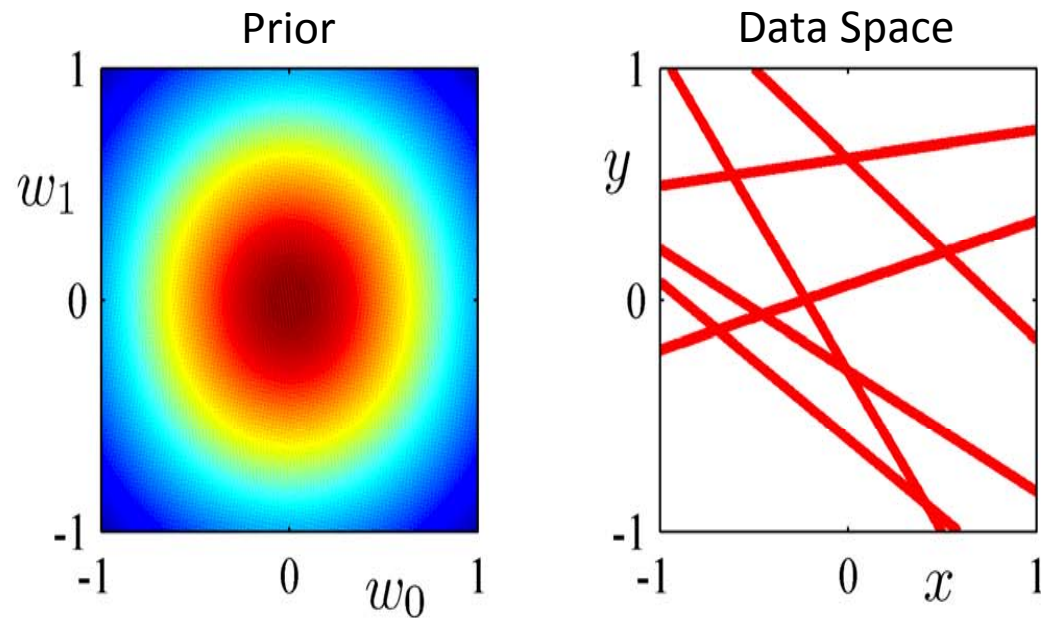
- 在贝叶斯估计中, 我们基于给定训练数据去推理参数 \mathbf{w} 的后验分布

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

– 从后验分布中可以采样很多 \mathbf{w}

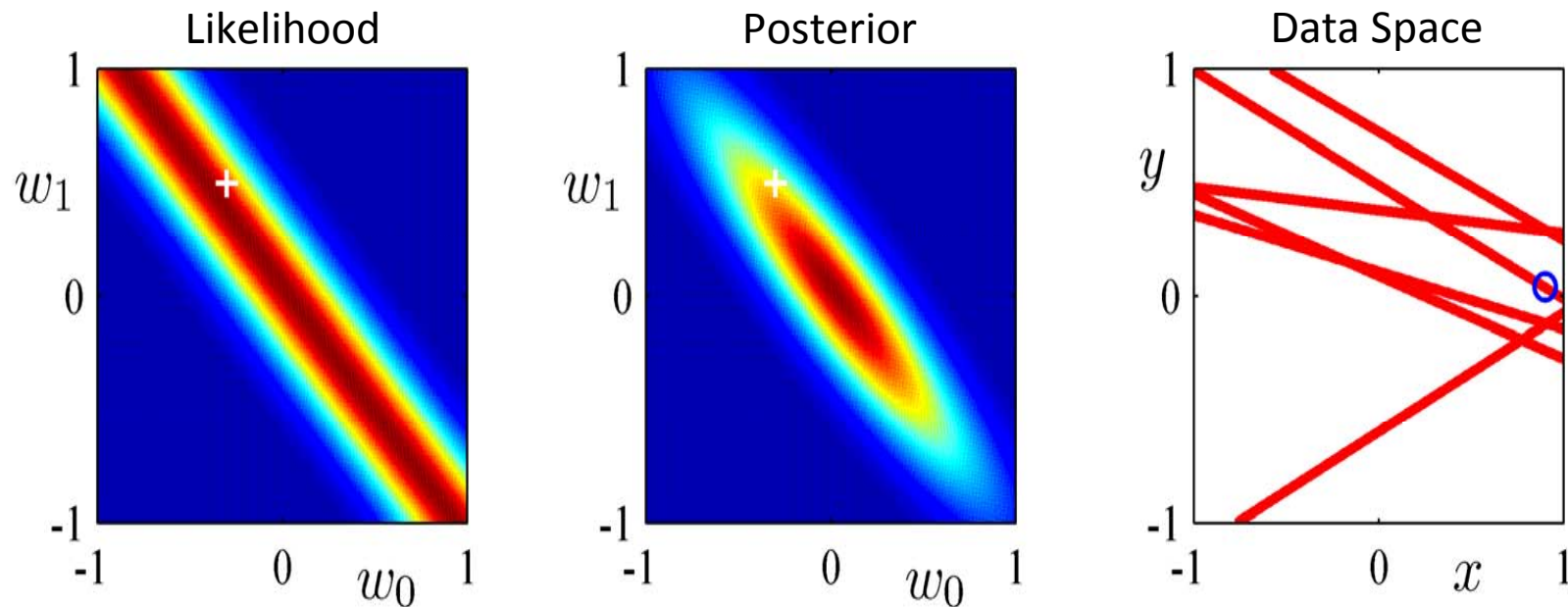
Bayesian Linear Regression (1)

0 data points observed



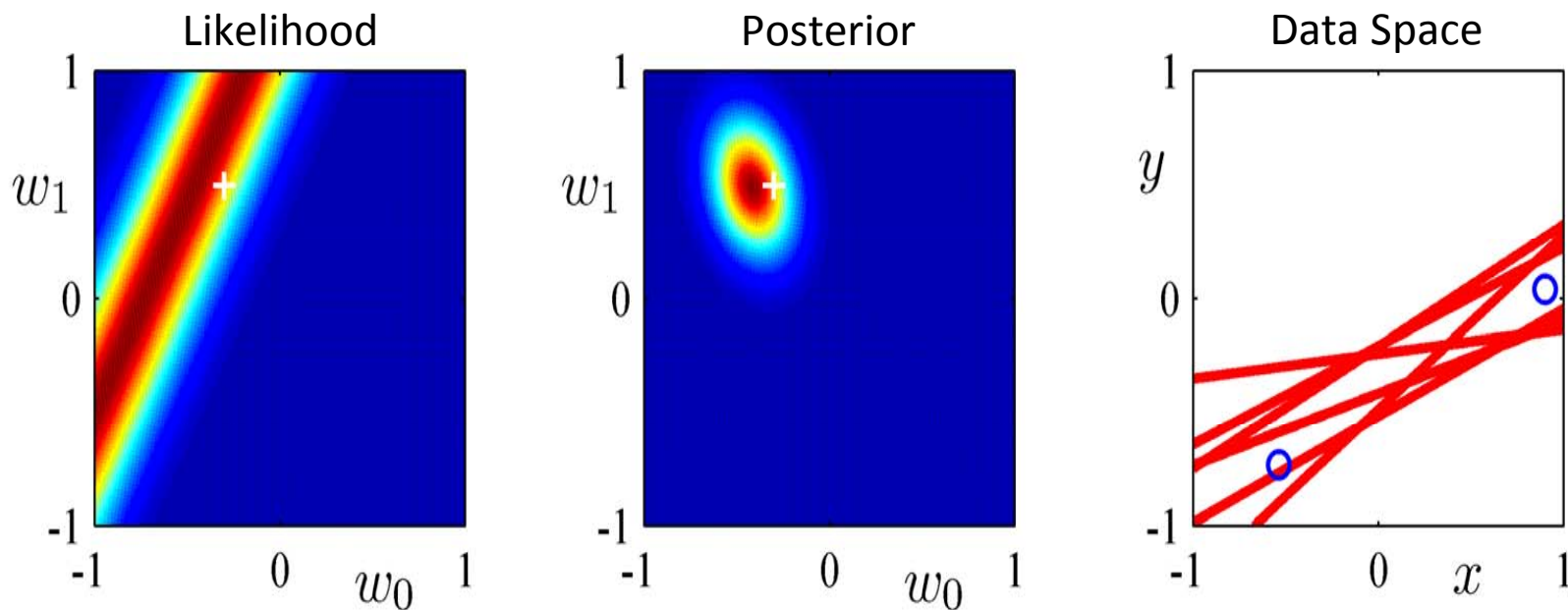
Bayesian Linear Regression (2)

1 data point observed



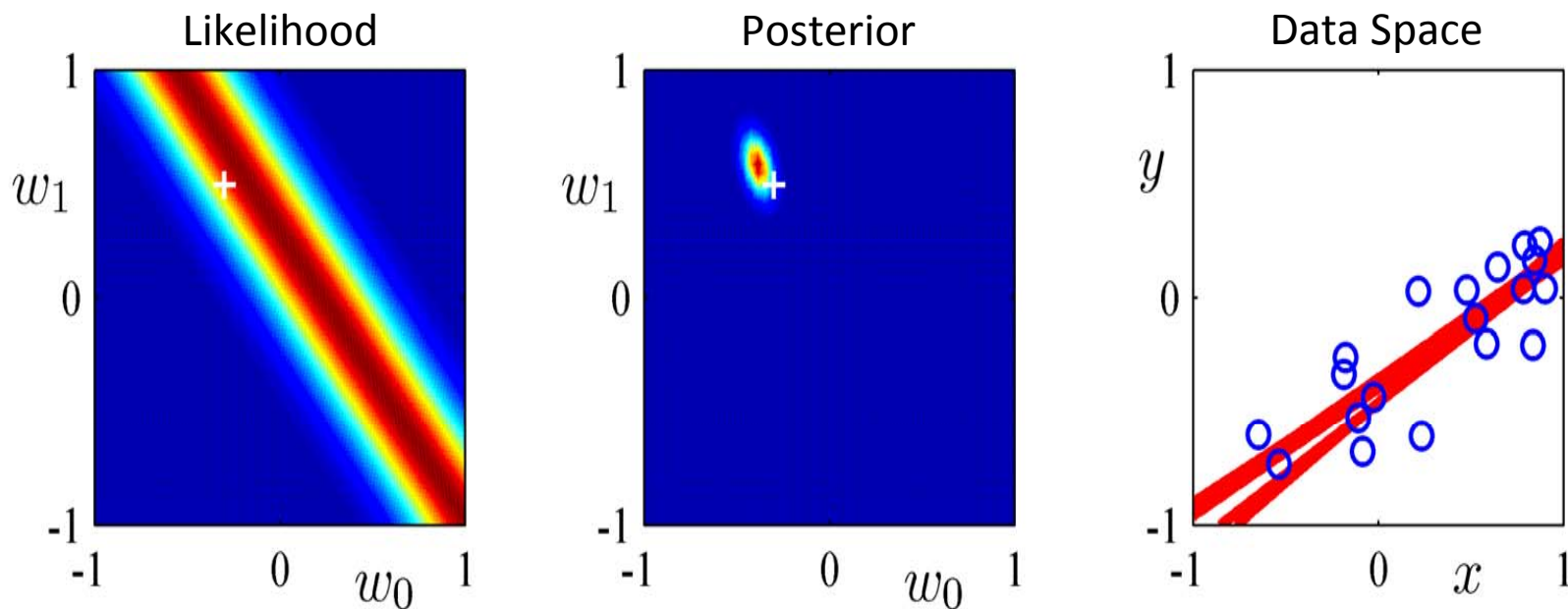
Bayesian Linear Regression (3)

2 data points observed



Bayesian Linear Regression (4)

20 data points observed



预测性分布(Predictive distribution)

- 给定一个新的 \mathbf{x} , 预测其输出 t :
 - “使用所有 \mathbf{w} 进行平均”

$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \end{aligned}$$

– 其中

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

例：使用高斯基函数的回归问题

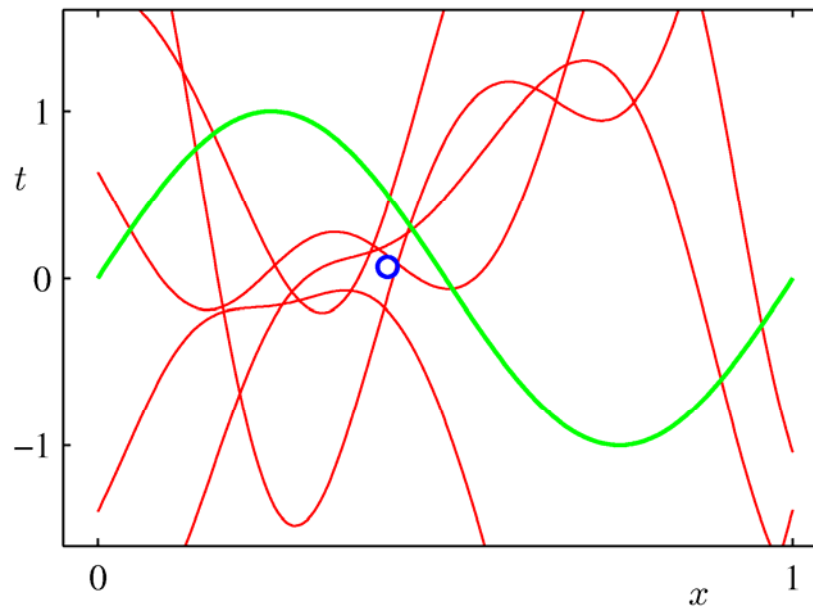
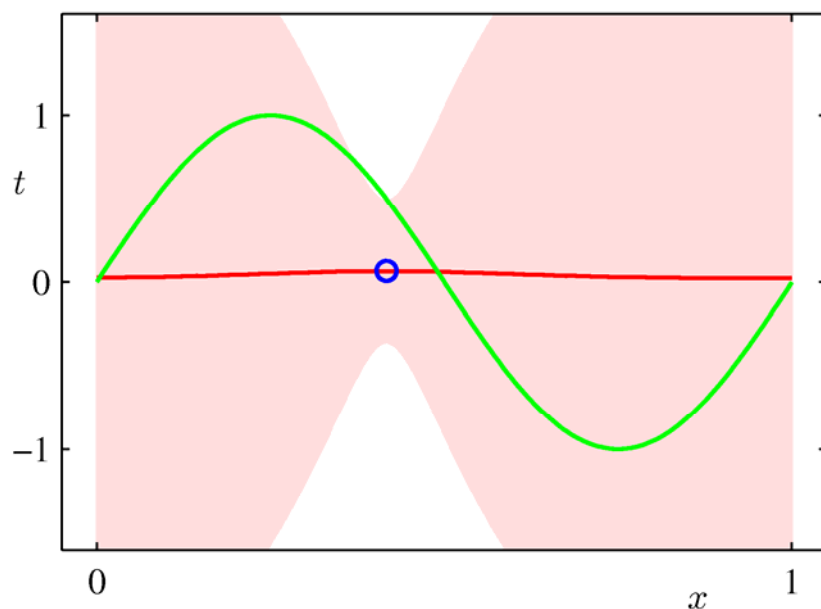
- 假设训练数据采样于一个正弦曲线
 - 我们使用**9**个高斯函数作为基函数
 - 使用不同数目的训练数据，可以获得不同的预测性分布

$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x})) \end{aligned}$$

例：使用高斯基函数的回归问题

- 使用1个训练样本获得的预测性分布

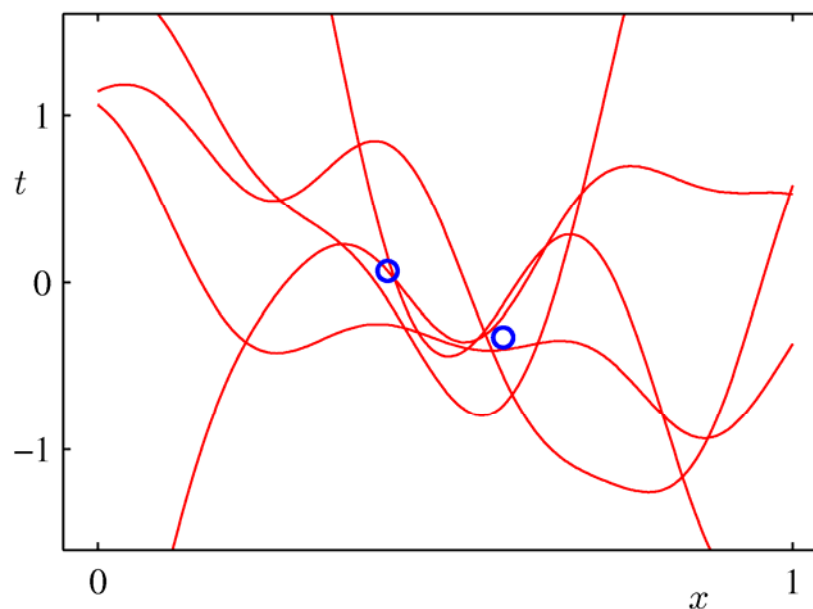
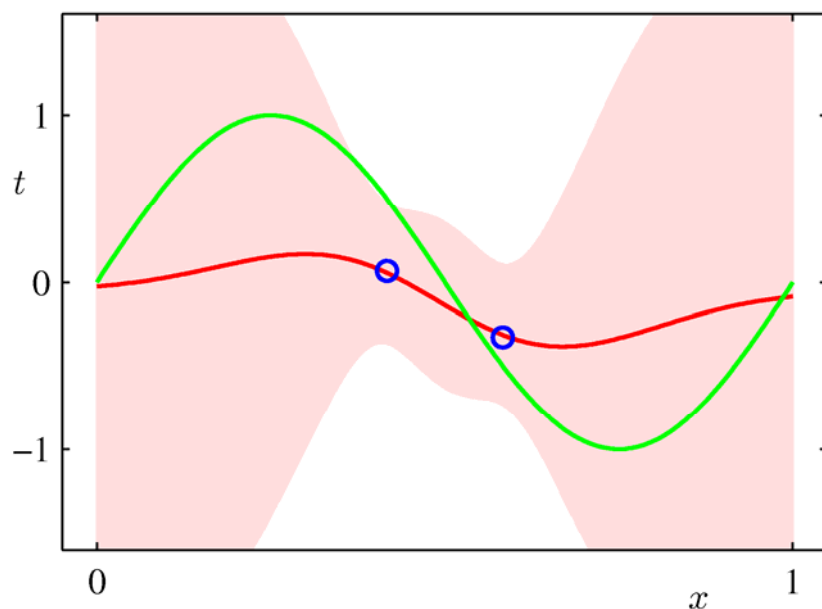
$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \end{aligned}$$



例：使用高斯基函数的回归问题

- 使用2个训练样本获得的预测性分布

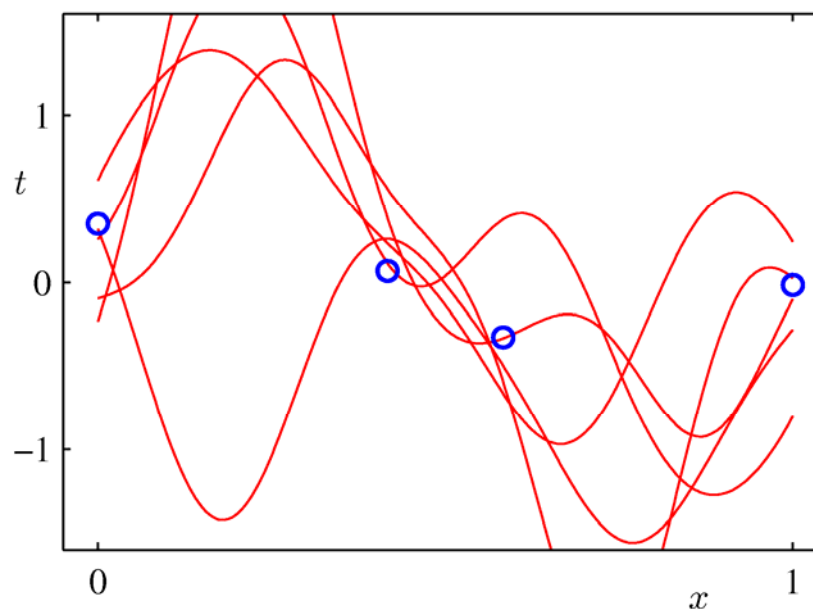
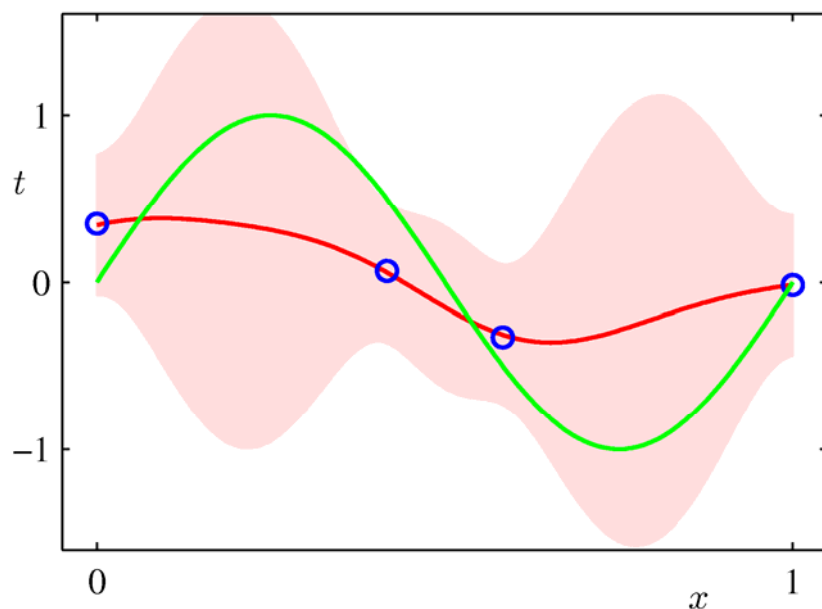
$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \end{aligned}$$



例：使用高斯基函数的回归问题

- 使用4个训练样本获得的预测性分布

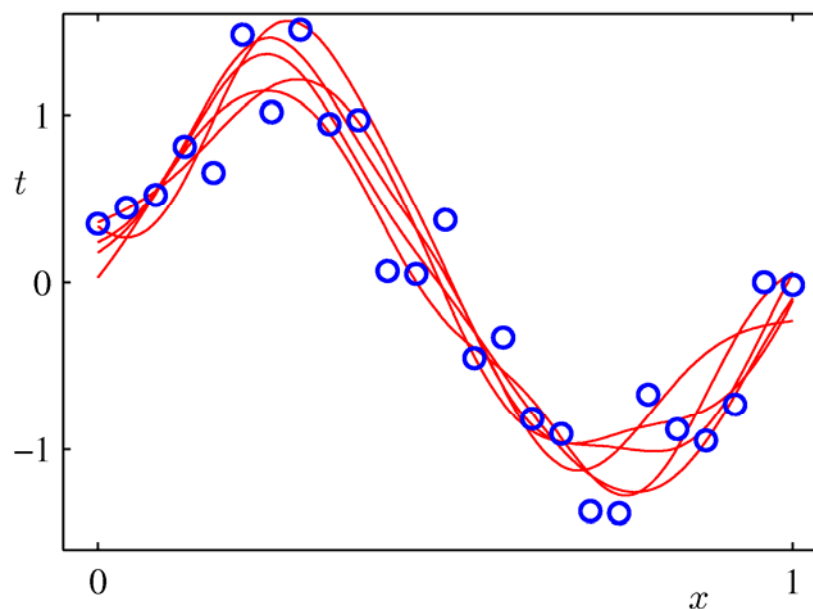
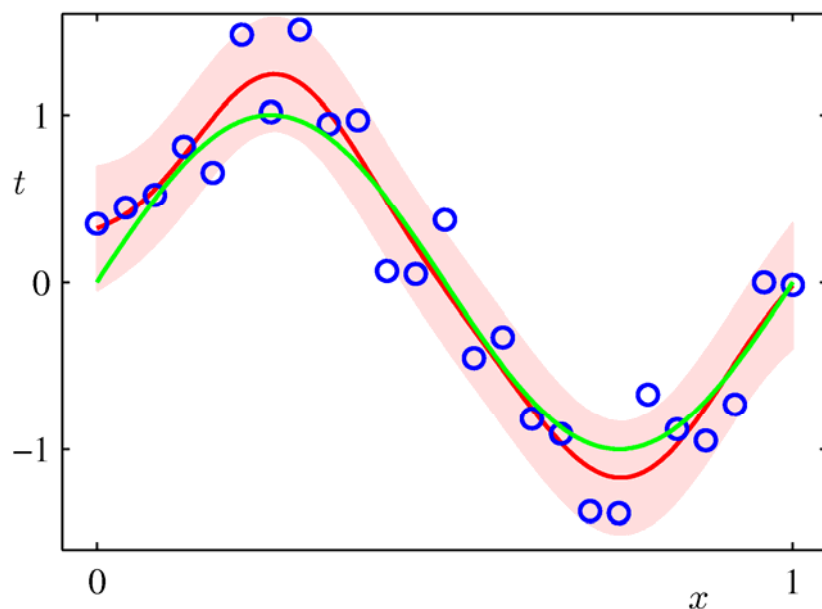
$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \end{aligned}$$



例：使用高斯基函数的回归问题

- 使用25个训练样本获得的预测性分布

$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \end{aligned}$$



贝叶斯线性回归 内容提要

- 引子: 曲线拟合问题
- 常用的几种分布
- 贝叶斯线性回归
 - 两个例子
 - 等价核
 - 先验分布中的超参数的处理

等价核 (Equivalent Kernel)

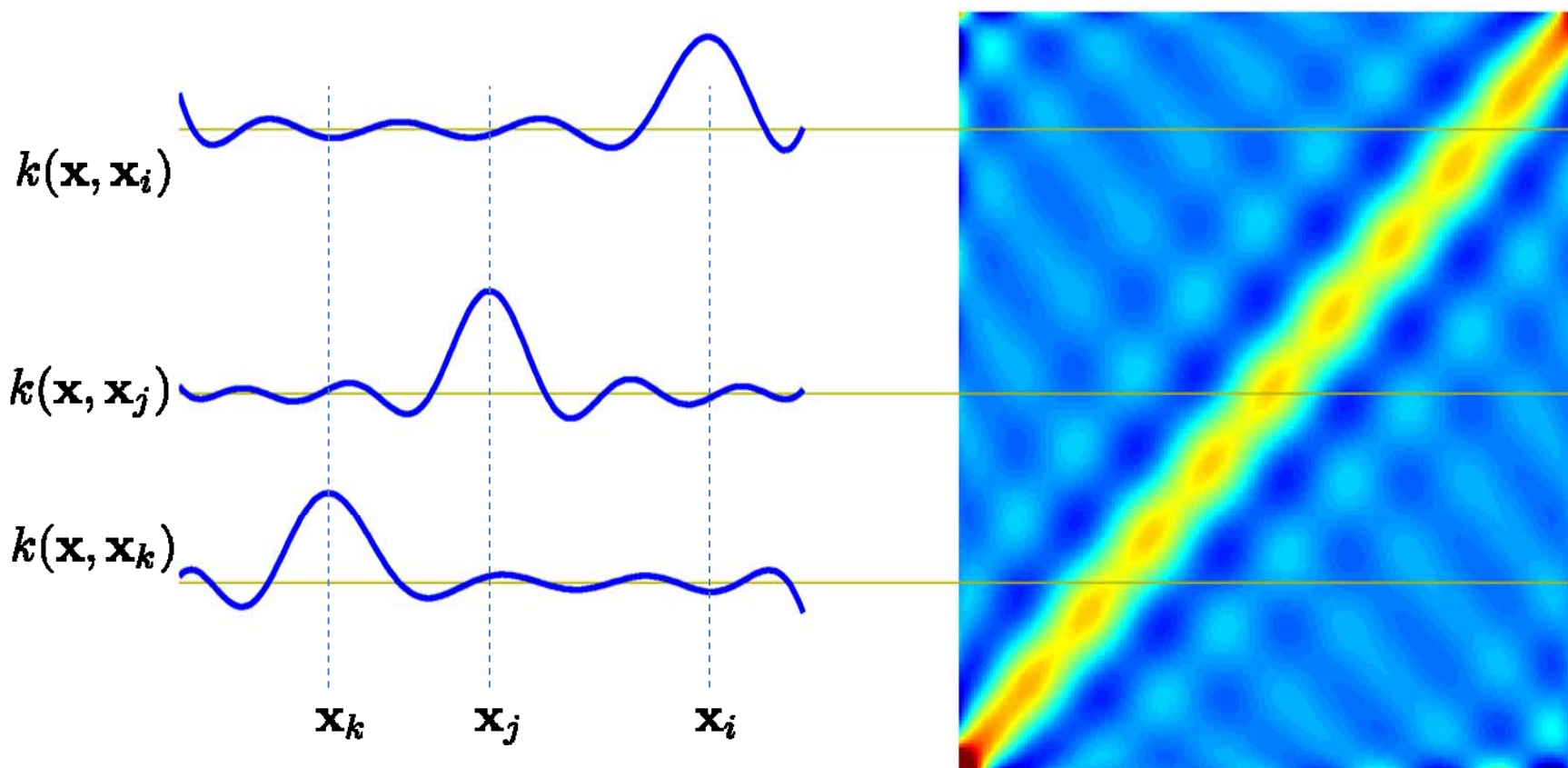
The predictive mean can be written

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

Equivalent kernel or smoother matrix.

This is a weighted sum of the training data target values, t_n .

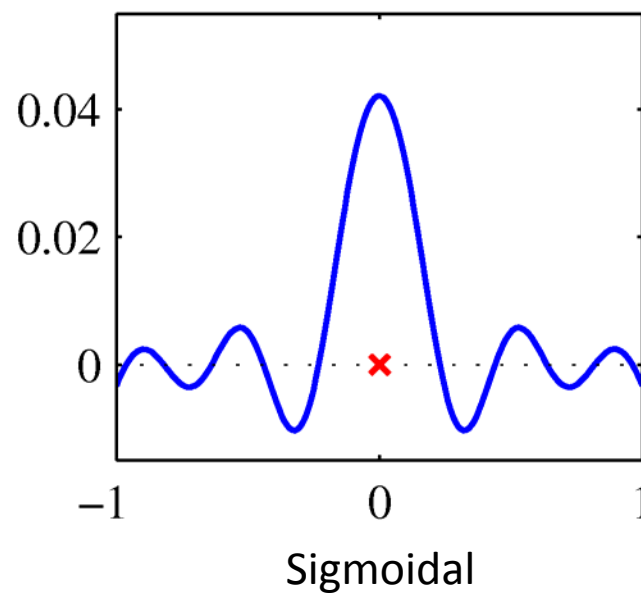
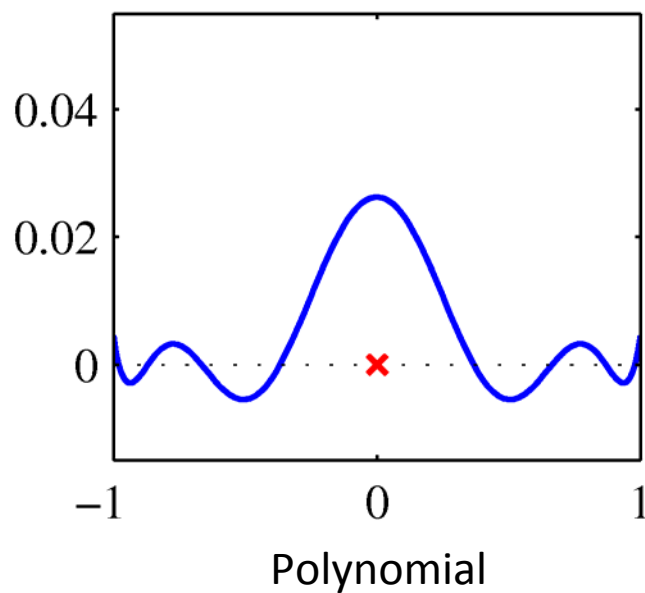
等价核 (Equivalent Kernel)



- t_n 对应的权值依赖于 \mathbf{x} 到 \mathbf{x}_n 的距离，距 \mathbf{x}_n 越近权值越大

等价核 (Equivalent Kernel)

- 非局部基(Nonlocal basis)也对应局部等价核



等价核 (Equivalent Kernel)

- 等价核作为协方差矩阵函数(covariance function)

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

- 若不引入基函数，而直接定义核函数，则得出高斯过程(Gaussian Processes)

等价核 (Equivalent Kernel)

- 基本性质：

- 归一化
$$\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1$$

- 注意到等价核的值可正可负

- 等价核也可以写成内积形式

$$k(\mathbf{x}, \mathbf{z}) = \psi(\mathbf{x})^T \psi(\mathbf{z})$$

- 其中

$$\psi(\mathbf{x}) = \beta^{1/2} \mathbf{S}_N^{1/2} \phi(\mathbf{x})$$

贝叶斯线性回归 内容提要

- 引子: 曲线拟合问题
- 常用的几种分布
- 贝叶斯线性回归
 - 两个例子
 - 等价核
 - 先验分布中的超参数的处理

先验分布中的超参数的处理方法

- 完整的贝叶斯预测性分布为:

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

- 但是，由于积分无法解析地进行，我们使用近似方法

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

- 其中 $(\hat{\alpha}, \hat{\beta})$ 是超参数联合分布 $p(\alpha, \beta|\mathbf{t})$ 的模式 (mode)，假设超参数的分布 **sharply peaked**
- 这种方法被称为:
 - **Empirical Bayes**
 - Type II or generalized maximum likelihood
 - **Evidence approximation**

如何估计超参数 (1)

- 由贝叶斯定理，得出

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

- 如果假设超参数的联合分布是均匀分布，则可以得出

$$\begin{aligned} p(\alpha, \beta | \mathbf{t}) &\propto p(\mathbf{t} | \alpha, \beta) \\ &= \int p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) d\mathbf{w}. \end{aligned}$$

– 若积分中的两项来自高斯分布，则有

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

如何估计超参数 (2)

- 寻找超参数，使得 $\ln p(\mathbf{t}|\alpha, \beta)$ 最大化
- 求微分，令其为0，则得到

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N}$$

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

$$- \text{其中 } (\beta \Phi^T \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i, \quad \gamma = \sum_i \frac{\lambda_i}{\alpha + \lambda_i}.$$

Q / A

- Any Questions...