

# PRML学习笔记——第二章

- PRML学习笔记——第二章
  - Probability Distributions
    - 2.1 Binary Variable
      - 2.1.1 The beta distribution
    - 2.2 Multinomial Variable
      - 2.2.1 The Dirichlet distribution
    - 2.3 The Gaussian Distribution
      - 2.3.1 Conditional Gaussian distributions
      - 2.3.2 Marginal Gaussian distribution
      - 2.3.3 Bayes' theorem for Gaussian variable
      - 2.3.4 Maximum likelihood for the Gaussian
      - 2.3.5 Sequential estimation
      - 2.3.6 Bayesian inference for the Gaussian
      - 2.3.7 Student's t-distribution
      - 2.3.8 Periodic variables
      - 2.3.9 Mixtures of Gaussians
    - 2.4 The Exponential Family
      - 2.4.1 Maximum likelihood and sufficient statistics
      - 2.4.2 Conjugate priors
      - 2.4.3 Noninformative priors
    - 2.5 Nonparametric Methods
      - 2.5.1 Kernel density estimators
      - 2.5.2 Nearest-neighbour methods

## Probability Distributions

在第一章中已经强调了概率在机器学习中的重要性，本章会对一些特别的概率分布exploration。

### 2.1 Binary Variable

对于一个binary random variable  $x \in \{0, 1\}$ ，我们denote该事件的probability:

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

因此我们可以得出*bernouli distribution*:

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

容易得出：

$$\begin{aligned}\mathbb{E}[x] &= \mu \\ \text{var}[x] &= \mu(1 - \mu)\end{aligned}$$

假设有数据集  $\mathcal{D} = \{x_1, \dots, x_N\}$ ，并且满足*i.i.d*于bernouli distribution，那么

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

就是likelihood，现在就可以通过maximize likelihood来得到 $\mu$ 的似然解：

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

也即样本均值。

当这个binary event重复 $N$ 次的时候，我们会得到*binomial distribution*：

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

### 2.1.1 The beta distribution

Beta distribution:

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

mean 和 variance:

$$\begin{aligned}\mathbb{E}[\mu] &= \frac{a}{a+b} \\ \text{var}[\mu] &= \frac{ab}{(a+b)^2(a+b+1)}.\end{aligned}$$

由于具有conjugacy性质，所以它的posterior distribution：

$$p(\mu | m, l, a, b) \propto \mu^{m+a-1} (1 - \mu)^{l+b-1}$$

其中的 $m$ 是发生 $x = 1$ 的事件次数,  $l = N - m$ . 可以看到prior和posterior具有相同的 $\mu$ 形式, 只是系数不同 (系数可以单独通过normalization得到)。从这个式子中我们也能知道每多观测一次事件, 只需要多乘一项 $\mu$ 或者 $1 - \mu$ , 再进行normalization。这也就能够进行sequential的预测, 避免每更新一次观测就要重新考虑所有data。

如果我们的目标是尽可能好的预测下一次 $x$ 的结果, 那么就是在预测:

$$p(x = 1 | \mathcal{D}) = \int_0^1 p(x = 1 | \mu) p(\mu | \mathcal{D}) d\mu = \int_0^1 \mu p(\mu | \mathcal{D}) d\mu = \mathbb{E}[\mu | \mathcal{D}]$$

得到:

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

该式表明了当实验次数足够多 ( $m, l \rightarrow \infty$ ) 的时候, 结果就是极大似然估计。对于有限的数据集,  $\mu$ 的估计就是在prior和极大似然估计之间。

## 2.2 Multinomial Variable

现在将binary variable推广到多个变量, 假设一共有 $K$ 个state, 事件 $\mathbf{x} = \{0, \dots, 1, \dots, 0\}$ , 就denote第 $k$ 个state( $x_k = 1$ )。

multinomial distribution定义:

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

其中的 $m_k$ 满足约束:

$$\sum_{k=1}^K m_k = N$$

### 2.2.1 The Dirichlet distribution

Multinomial distribution的conjugate prior distribution是Dirichlet distribution:

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

其中  $0 \leq \mu_k \leq 1$  and  $\sum_k \mu_k = 1$ ,  $\alpha_0 = \sum_{k=1}^K \alpha_k$ .

通过Bayes定理，可以得到同是Dirichlet distribution的posterior：

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

## 2.3 The Gaussian Distribution

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

其中唯一依赖 $\mathbf{x}$ 的项是

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

这个 $\Delta$ 被称为*Mahalanobis distance*(马氏距离)。

note: 当 $\boldsymbol{\Sigma}$ 是identity matrix时，这个距离就也就是Euclidean distance

不失一般性，将 $\boldsymbol{\Sigma}$ 假定为symmetric matrix，那么可以做特征值分解：

$$\boldsymbol{\Sigma} \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

我们选择orthonormal的 $\mathbf{u}$ ，有(写成分块矩阵再左右乘上 $U^T$ 、 $U$ 可证)：

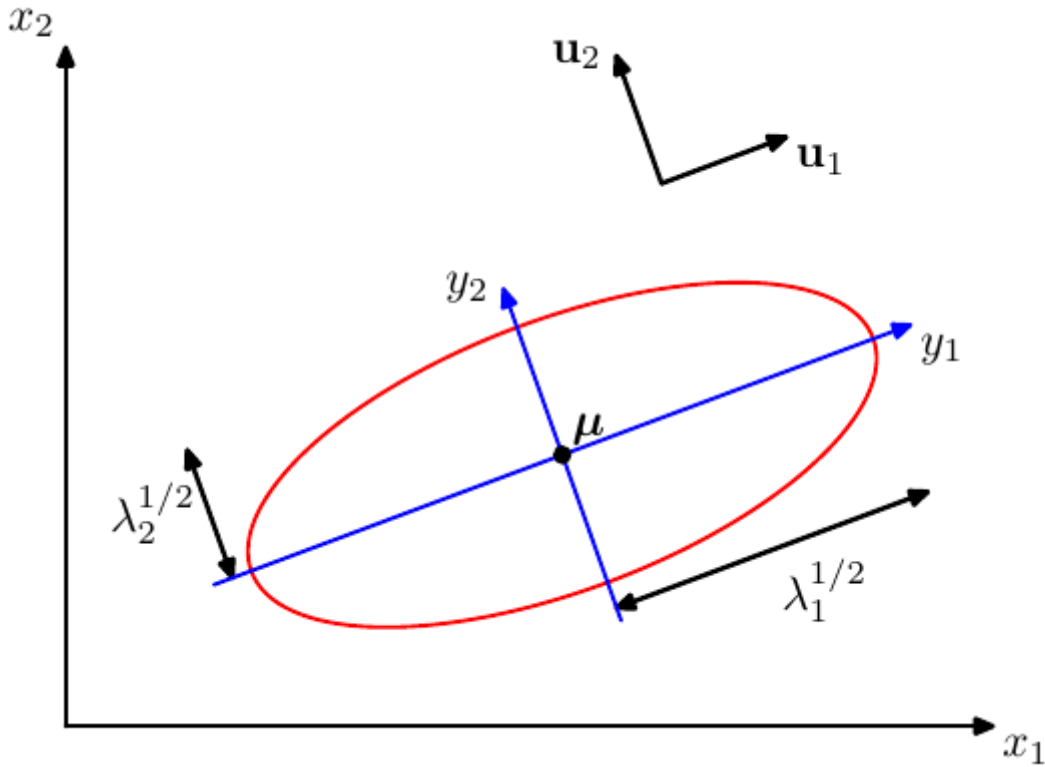
$$\boldsymbol{\Sigma} = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$

那么二次型变为：

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

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

显然，如果考虑空间中等概率密度的区域，二维情况就是一个椭圆线。



原来的 $\mathbf{x}$ 经过shift、rotate得到normalize后的变量 $\mathbf{y}$ , 有 $\mathbf{y} = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu})$

现在考虑从 $\mathbf{x}$ 到 $\mathbf{y}$ 的坐标变换, 得到一个Jacobian matrix  $\mathbf{J}$ :

$$J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ij}$$

因此Jacobian matrix的行列式:

$$|\mathbf{J}|^2 = |\mathbf{U}^T|^2 = |\mathbf{U}^T| |\mathbf{U}| = |\mathbf{U}^T \mathbf{U}| = |\mathbf{I}| = 1$$

$$|\mathbf{J}| = 1.$$

同时Covariance matrix的determinant能被写成:

$$|\boldsymbol{\Sigma}|^{1/2} = \prod_{j=1}^D \lambda_j^{1/2}$$

因此在 $\mathbf{y}$ 的坐标系下的Gaussian distribution为

$$p(\mathbf{y}) = p(\mathbf{x}) |\mathbf{J}| = \prod_{j=1}^D \frac{1}{(2\pi\lambda_j)^{1/2}} \exp\left\{-\frac{y_j^2}{2\lambda_j}\right\}$$

### 2.3.1 Conditional Gaussian distributions

首先我们把变量划分为 $a, b$ 两个子集，将多维高斯分布用分块矩阵来表示。

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

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

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

记 *precision matrix*:

$$\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$$

也就有：

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

note: 这里的 $\boldsymbol{\Lambda}_{aa}$ 并不是简单的等于 $\boldsymbol{\Sigma}_{aa}^{-1}$ 。

一个重要的性质是：joint Gaussian distribution对应的marginal和conditional distribution都是Gaussian distribution.

现在将Gaussian distribution的指数项用分块矩阵表示：

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = & \\ & -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ & - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

对于一个general的Gaussian distribution，它的指数项为：

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

note: 对于求conditional probability  $p(\mathbf{x}_a|\mathbf{x}_b)$ ， $\mathbf{x}_b$ 被看作一个constant.

如此就能通过比较系数来得出未知数：

$$\begin{aligned}\Sigma_{a|b} &= \Lambda_{aa}^{-1} \\ \mu_{a|b} &= \Sigma_{a|b} \{ \Lambda_{aa} \mu_a - \Lambda_{ab} (\mathbf{x}_b - \mu_b) \} \\ &= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \mu_b)\end{aligned}$$

现在利用一个恒等式结果：

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

其中  $\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}$  被成为 *Schur complement* (舒尔补) .可以得到：

$$\begin{aligned}\Lambda_{aa} &= (\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1} \\ \Lambda_{ab} &= -(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}\end{aligned}$$

所以conditional distribution的结果用mean和covariance表示就是：

$$\begin{aligned}\mu_{a|b} &= \mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(\mathbf{x}_b - \mu_b) \\ \Sigma_{a|b} &= \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}.\end{aligned}$$

### 2.3.2 Marginal Gaussian distribution

前面得到了joint distribution是gaussian distribution时的conditional distribution，现在：

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

同样的，关注指数项的二次型：

$$-\frac{1}{2}\mathbf{x}_b^T \Lambda_{bb} \mathbf{x}_b + \mathbf{x}_b^T \mathbf{m} = -\frac{1}{2}(\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m})^T \Lambda_{bb} (\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m}) + \frac{1}{2}\mathbf{m}^T \Lambda_{bb}^{-1} \mathbf{m}$$

其中  $\mathbf{m} = \Lambda_{bb} \mu_b - \Lambda_{ba} (\mathbf{x}_a - \mu_a)$  .

等式右边第一项是dependent  $\mathbf{x}_b$ ，第二项independent  $\mathbf{x}_b$ (但dependent  $\mathbf{x}_a$ )，由于积分在 $\mathbf{x}_b$ 上，所以只需考虑这个积分形式：

$$\int \exp \left\{ -\frac{1}{2}(\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m})^T \Lambda_{bb} (\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m}) \right\} d\mathbf{x}_b$$

而这个积分里面相当于是未normalize的Gaussian distribution，那么积分结果就是这个normalize系数

的倒数（会成为marginal distribution的系数项）。又由于确定Gaussian distribution只需要考虑指数项：

$$\begin{aligned} & \frac{1}{2} [\Lambda_{bb} \boldsymbol{\mu}_b - \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a)]^T \Lambda_{bb}^{-1} [\Lambda_{bb} \boldsymbol{\mu}_b - \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a)] \\ & - \frac{1}{2} \mathbf{x}_a^T \Lambda_{aa} \mathbf{x}_a + \mathbf{x}_a^T (\Lambda_{aa} \boldsymbol{\mu}_a + \Lambda_{ab} \boldsymbol{\mu}_b) + \text{const} \\ & = -\frac{1}{2} \mathbf{x}_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \mathbf{x}_a \\ & + \mathbf{x}_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1} \boldsymbol{\mu}_a + \text{const} \end{aligned}$$

对比系数仍然能够得到：

$$\begin{aligned} \mathbb{E}[\mathbf{x}_a] &= \boldsymbol{\mu}_a \\ \text{cov}[\mathbf{x}_a] &= \Sigma_{aa}. \end{aligned}$$

这个表示形式要比conditional distribution简单多了。

综上：partitioned Gaussians

给定Gaussian distribution  $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \Sigma)$  ,其中

$$\Lambda \equiv \Sigma^{-1}, \mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}.$$

Conditional distribution

$$\begin{aligned} p(\mathbf{x}_a | \mathbf{x}_b) &= \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \Lambda_{aa}^{-1}) \\ \boldsymbol{\mu}_{a|b} &= \boldsymbol{\mu}_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned}$$

Marginal distribution

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

### 2.3.3 Bayes' theorem for Gaussian variable

Marginal and Conditional Gaussians

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

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

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

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



where

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

### 2.3.4 Maximum likelihood for the Gaussian

假设有观测  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ , 并且是由独立同分布的Gaussian distribution中得到, 我们可以用 maximum likelihood来estimate分布中的parameter:

$$\ln p(\mathbf{X} | \boldsymbol{\mu}, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

可以解得似然解:

$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \text{ 和 } \Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T$$

### 2.3.5 Sequential estimation

Sequential methods能够一次只处理一个data point。在一些on-line application和涉及到large data的时候, 一次处理全部数据是不可行的。

考虑上一节中的Gaussian distribution对 $\mu$ 的似然估计。现在将 $\mathbf{x}_N$ 从式子中拆分出来:

$$\begin{aligned} \boldsymbol{\mu}_{\text{ML}}^{(N)} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ &= \frac{1}{N} \mathbf{x}_N + \frac{1}{N} \sum_{n=1}^{N-1} \mathbf{x}_n \\ &= \frac{1}{N} \mathbf{x}_N + \frac{N-1}{N} \boldsymbol{\mu}_{\text{ML}}^{(N-1)} \\ &= \boldsymbol{\mu}_{\text{ML}}^{(N-1)} + \frac{1}{N} (\mathbf{x}_N - \boldsymbol{\mu}_{\text{ML}}^{(N-1)}) \end{aligned}$$

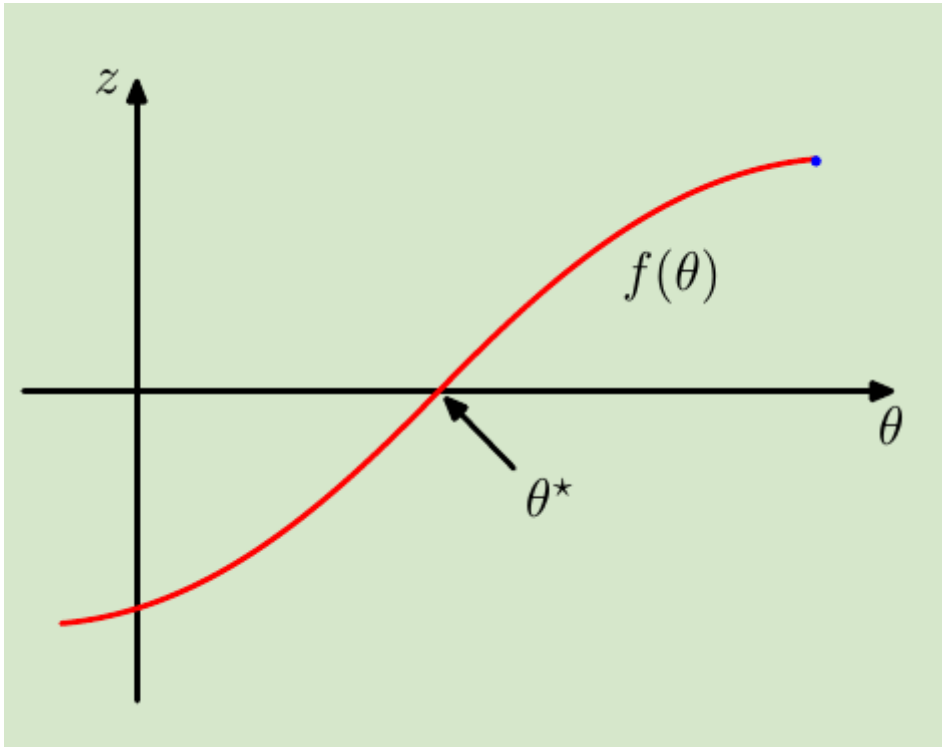
这样一看就像是在上一步估计得到的 $\boldsymbol{\mu}_{\text{ML}}^{(N-1)}$ 上再向'error signal'  $(\mathbf{x}_N - \boldsymbol{\mu}_{\text{ML}}^{(N-1)})$ 更新一步。并且随着 $N$ 的增加, 每次更新的data point会有更小的contribution.

下面介绍Robbins-Monro algorithm.考虑有一对随机变量 $\theta$ 和 $z$ , 定义函数 $f(\theta)$ :

$$f(\theta) \equiv \mathbb{E}[z | \theta] = \int z p(z | \theta) dz$$

这样定义的函数被称为regression functions.

我们的目标是找一个 $\theta^*$ 使得 $f(\theta^*) = 0$ . 假设 $\mathbb{E}[(z - f)^2 | \theta] < \infty$ , 不失一般性, 考虑 $f(\theta) > 0$  for  $\theta > \theta^*$ 与 $f(\theta) < 0$  for  $\theta < \theta^*$ .



该算法给了一种general算法去找 $f(\theta)$ 的根, 其中 $f$ 是由条件期望给出的 $\mathbb{E}[z|\theta]$ .

算法策略为 $\theta^{(N)} = \theta^{(N-1)} + a_{N-1} z(\theta^{(N-1)})$ , 其中 $z(\theta^{(N-1)})$ 是当 $\theta$ 取 $\theta^{(N-1)}$ 时对 $z$ 的观测。系数 $\{a_N\}$ 是一个正数序列, 满足:

$$\begin{array}{l} \lim_{N \rightarrow \infty} a_N = 0 \quad \sum_{N=1}^{\infty} a_N = \infty \\ \sum_{N=1}^{\infty} a_N^2 < \infty \end{array}$$

note:

1. 第一个约束是为了保证序列的估计是收敛的
2. 第二个约束是为了保证不会不收敛到根
3. 第三个约束是为了保证累计的noise是个有限方差, 不会收敛失败

现在考虑一个一般的maximum likelihood如何用该算法序列解决。我们知道MLE就是求一个驻点:

$$\left. \frac{\partial}{\partial \theta} \left( \frac{1}{N} \sum_{n=1}^N -\ln p(x_n | \theta) \right) \right|_{\theta = \theta_{ML}} = 0$$

交换求导和求和位置并让 $N \rightarrow \infty$ :

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial \theta} -\ln p(x_n | \theta) = \mathbb{E}_x \left[ \frac{\partial}{\partial \theta} -\ln p(x | \theta) \right]$$

这也就把问题转换为了求the root of regression function.使用序列算法:

$$\theta^{(N)} = \theta^{(N-1)} + a_{N-1} \frac{\partial}{\partial \theta^{(N-1)}} \ln p(x_{N-1} \mid \theta^{(N-1)})$$

特别地，如果是求mean of Gaussian,那么

$$z = \frac{\partial}{\partial \mu_{\mathrm{ML}}} \ln p(x \mid \mu_{\mathrm{ML}}, \sigma^2) = \frac{1}{\sigma^2} (x - \mu_{\mathrm{ML}})$$

### 2.3.6 Bayesian inference for the Gaussian

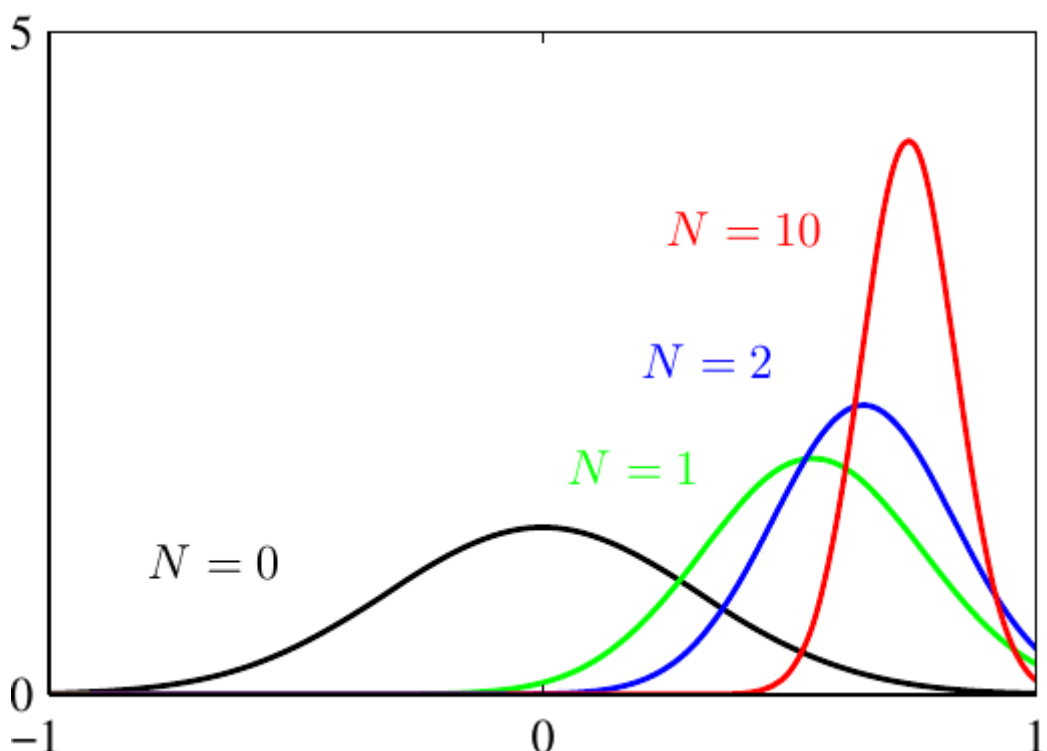
对于一个已知variance  $\sigma^2$ ,未知mean  $\mu$  的Gaussian distribution, 通过极大似然能得到一个  $\mu$  的具体估计值。但利用Bayes理论可以得到  $\mu$  的distribution, 为了简单我们取共轭先验Gaussian distribution,  $p(\mu) = \mathcal{N}(\mu \mid \mu_0, \sigma_0^2)$ , 我们求posterior  $p(\mu \mid \mathbf{X}) = \mathcal{N}(\mu \mid \mu_N, \sigma_N^2)$ .

其中

$$\begin{aligned} \mu_N &= \frac{\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N}{N\sigma_0^2 + \sigma^2} \mu_{\mathrm{ML}} \\ \sigma_N^2 &= \frac{1}{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}} \end{aligned}$$

从这个式子我们可以发现：

1. posterior 的  $\mu$  是介于prior和MLE之间的一个数，并且当  $N=0$  时，就是prior;当  $N \rightarrow \infty$  时，就是MLE.
2. 对posterior的  $\sigma$  来说,当  $N \rightarrow \infty$  时，  $\sigma_N$  会接近0，表示估计的  $\mu$  的precision 更好.
3. 类似于N,当  $\sigma_0 \rightarrow \infty$  时，  $\mu_N$  也会等于MLE的结果



图中的 $N=0$ 黑色curve就是prior, 当 $N$ 逐渐增大,  $\mu$ 的distribution就会接近MLE的结果 (mean是0.8, variance是0.1)

前面已经讨论了如何使用sequential method去估计Gaussian的mean, 现在我们会发现在bayes paradigm会很自然导出sequential的求解。

$$p(\boldsymbol{\mu} \mid D) \propto [p(\boldsymbol{\mu}) \prod_{n=1}^{N-1} p(\mathbf{x}_n \mid \boldsymbol{\mu})] p(\mathbf{x}_N \mid \boldsymbol{\mu})$$

由这个式子我们现在可以直接这么考虑: 前面 $N-1$ 个data point是这次estimate的**prior**, 而只有这次的第 $N$ 个data point才被考虑进**posterior**.

note: 这里有个前提假设: 这 $N$ 个样本点是i.i.d.的

前面推导了怎么对已知 $\sigma$ 未知 $\mu$ 的Gaussian作估计, 现在考虑已知 $\mu$ 未知 $\sigma$ 的情况. likelihood的形式是:

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

同样我们取共轭先验 (gamma distribution) 为了大大简化:

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

其中 $\lambda = 1/\sigma^2$ 表示precision.

我们计算posterior:

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

进而得到posterior distribution的parameter:

$$\begin{array}{l} a_N = a_0 + \frac{N}{2} \\ b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 \end{array} = b_0 + \frac{N}{2} \sigma_{\text{ML}}^2$$

从这个结果上分析:

1. 第一个式子可以理解prior  $a_0$ 就是已经有了 $2a_0$ 的先验观测.
2. 第二个式子可以理解有了 $2a_0$ 个variance是 $2b_0 / (2a_0) = b_0 / a_0$ 的有效的先验观测。 (\*\*why?\*\*)

然后我们再考虑当mean和variance都是未知的情况。为了找共轭先验, 我们先看下likelihood的形式:

$$\begin{array}{c} p(\mathbf{X} \mid \mu, \lambda) = \prod_{n=1}^N \left( \frac{\lambda}{2\pi} \right)^{1/2} \exp \left\{ -\frac{\lambda}{2} (x_n - \mu)^2 \right\} \propto \lambda^{N/2} \exp \left\{ -\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\} \end{array}$$

$$\frac{\lambda}{\sum_{n=1}^N x_n} \exp \left( -\frac{\lambda}{\sum_{n=1}^N x_n} \right)$$

现在我们希望prior具有和likelihood一样的依赖 $\mu$ 和 $\lambda$ 的形式:

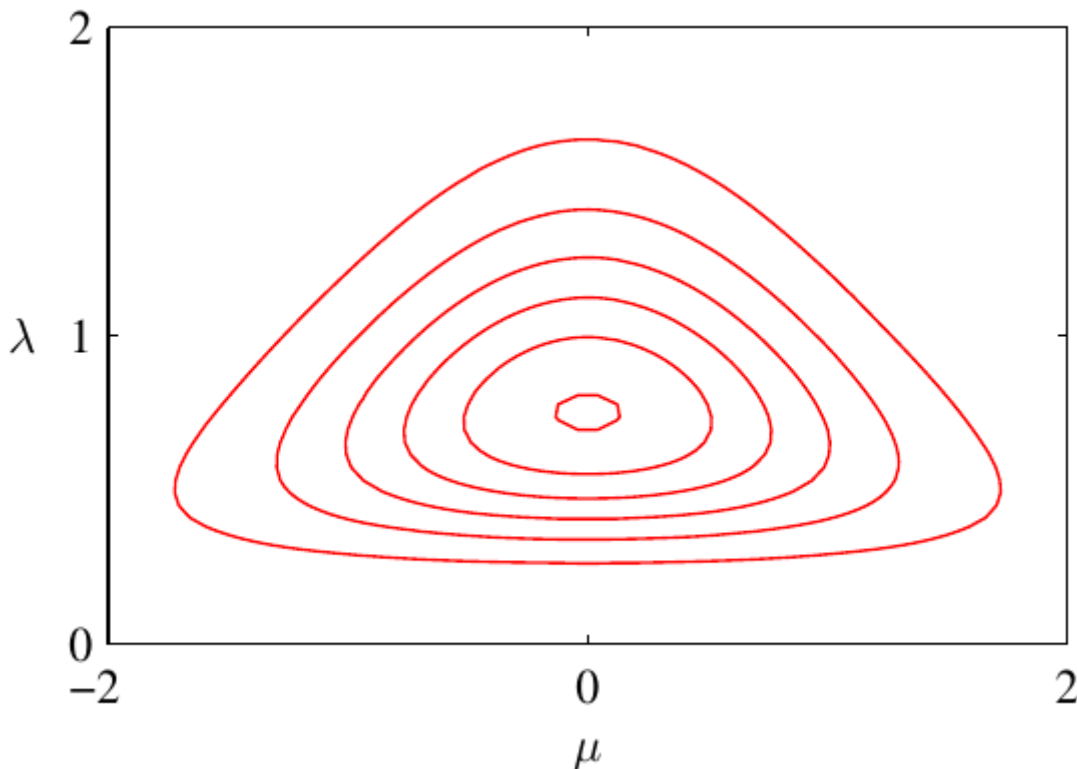
$$p(\mu, \lambda) \propto \lambda^{1/2} \exp \left( -\frac{\lambda}{2} \mu^2 \right) \exp \left( -\frac{c}{\lambda} \right) \propto \exp \left( -\frac{\beta \lambda}{2} \mu^2 - \frac{c}{\lambda} \right)$$

我们可以启发的写出prior的形式:

$$p(\mu, \lambda) = \mathcal{N}(\mu \mid \mu_0, (\beta \lambda)^{-1}) \operatorname{Gam}(\lambda \mid a, b)$$

其中 $\mu_0 = c / \beta$ ,  $a = 1 + \beta / 2$ ,  $b = d - c^2 / 2 \beta$ . 该分布被称为 *normal-gamma or Gaussian-gamma distribution*.

note: 值得注意的是这个distribution并不是简单的independent的Gaussian和Gamma distribution的乘积, 因为Gaussian的variance是 $\lambda$ 的线性关系。



该图就是 $\mu_0 = 0$ ,  $\beta = 2$ ,  $a = 5$ ,  $b = 6$ 的normal-gamma distribution

在多变量的Gaussian  $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{\Lambda}^{-1})$  里, 未知mean已知precision情况下的共轭prior仍然是Gaussian。但对于已知mean未知precision的情况下, 共轭prior是 *Wishart distribution*:

$$\mathcal{W}(\mathbf{\Lambda} \mid \mathbf{W}, \nu) \propto |\mathbf{\Lambda}|^{(\nu-D-1)/2} \exp \left( -\frac{1}{2} \operatorname{Tr}(\mathbf{W}^{-1} \mathbf{\Lambda}) \right)$$

其中的 $\nu$ 是该distribution的**自由度**,  $B$ 是normalize constant:

$$B(\mathbf{W}, \nu) = \frac{1}{|\mathbf{W}|^{\nu/2}} \frac{1}{(2\pi)^{D(D-1)/4}} \prod_{i=1}^D \frac{1}{\Gamma(\frac{\nu+1-i}{2})}$$

最后对于mean和precision都unknown的情况, 共轭prior类推出来就是 *normal-Wishart or Gaussian-Wishart* distribution:

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} \mid \boldsymbol{\mu}_0, \beta, \mathbf{W}, \nu) = \mathcal{N}(\boldsymbol{\mu} \mid \boldsymbol{\mu}_0, \beta \mathbf{W}) \boldsymbol{\Lambda}^{-1} \mathcal{W}(\boldsymbol{\Lambda} \mid \mathbf{W}, \nu)$$

### 2.3.7 Student's t-distribution

我们已经知道关于Gaussian precision的共轭prior是Gamma distribution。如果我们有一个Gaussian是  $\mathcal{N}(x \mid \mu, \tau^{-1})$  和一个Gamma prior是  $\text{Gam}(\tau \mid a, b)$ , 并且将 $\tau$ 积分积掉。就会得到关于 $x$ 的marginal distribution:

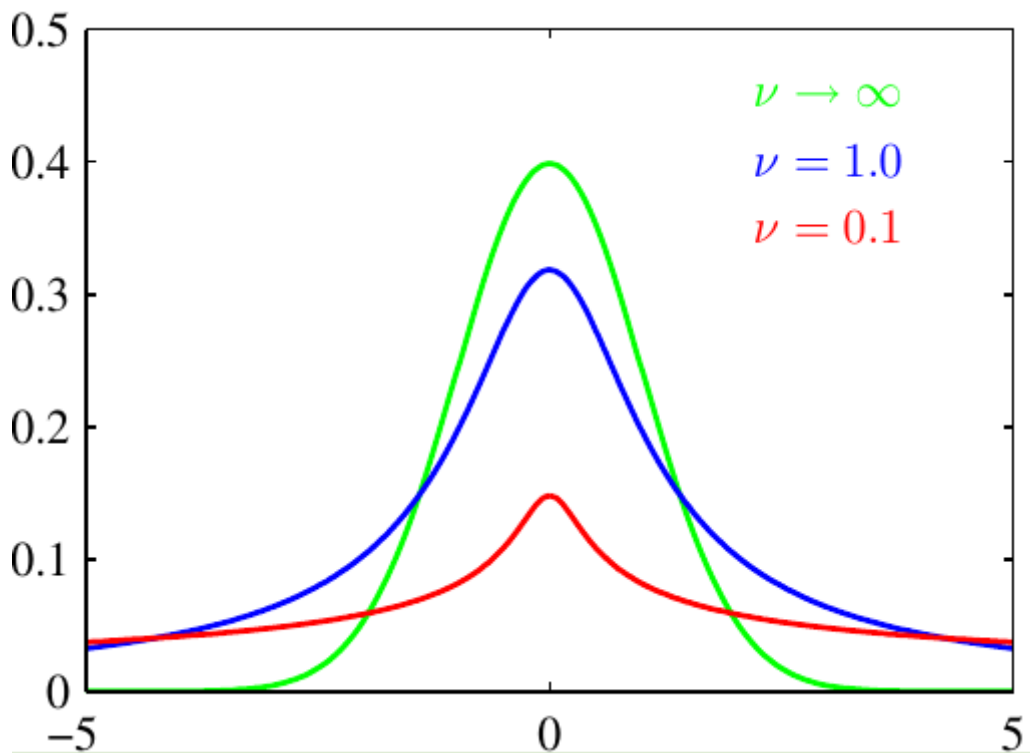
$$\begin{aligned} p(x \mid \mu, a, b) &= \int_0^\infty \mathcal{N}(x \mid \mu, \tau^{-1}) \text{Gam}(\tau \mid a, b) \mathrm{d}\tau \\ &= \int_0^\infty \frac{b^a}{\Gamma(a)} e^{-(b+\tau)} \tau^{a-1} \frac{1}{\Gamma(a)} \left(\frac{\tau}{2\pi}\right)^{1/2} \exp\left\{-\frac{\tau}{2}(x-\mu)^2\right\} \mathrm{d}\tau \\ &= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi}\right)^{1/2} \int_0^\infty \tau^{a-1} e^{-\tau\left[b+\frac{(x-\mu)^2}{2}\right]} \mathrm{d}\tau \end{aligned}$$

为了方便替换下变量  $v = 2a$ ,  $\lambda = a/b$ , 这样就得到了 *Student's t-distribution*:

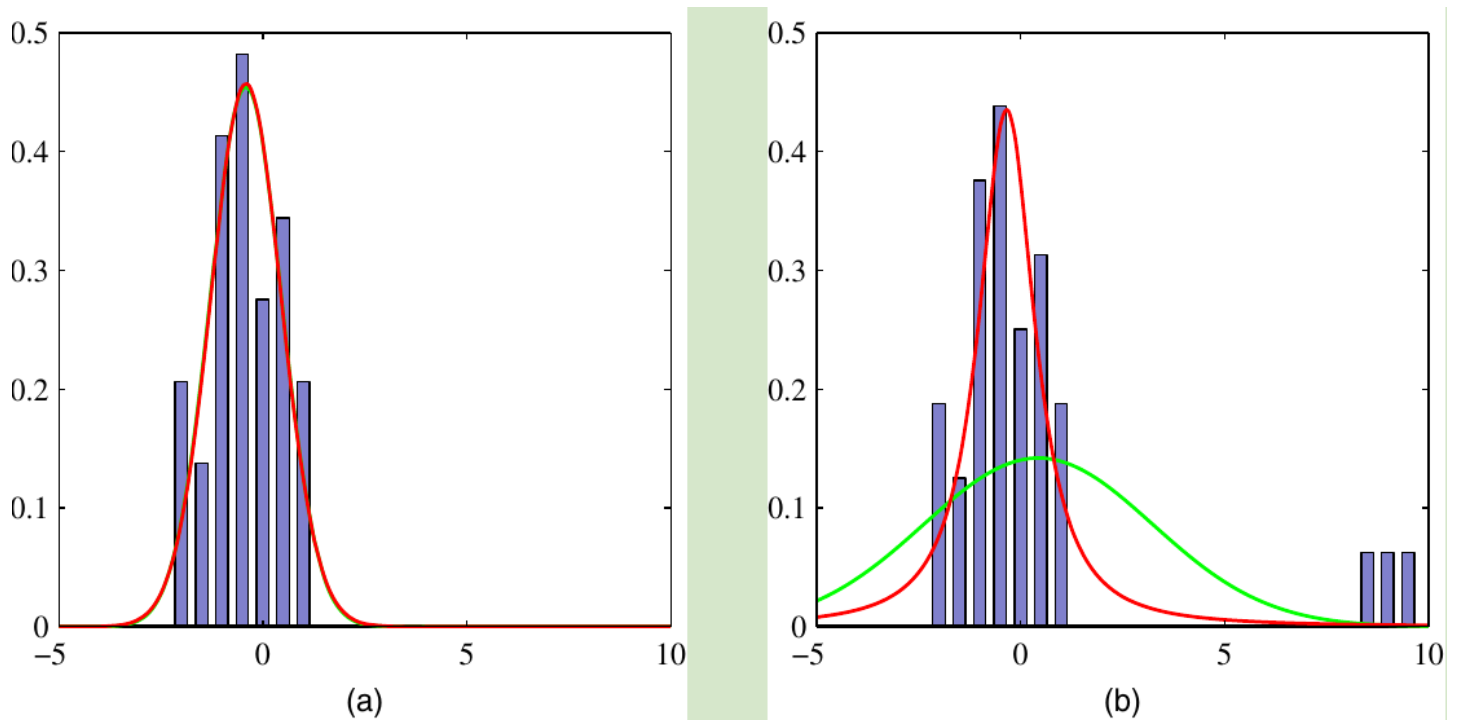
$$\text{St}(x \mid \mu, \lambda, \nu) = \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left(\frac{\lambda}{\pi}\right)^{1/2} \left[1 + \frac{\lambda(x-\mu)^2}{\nu}\right]^{-\nu/2 - 1/2}$$

$\lambda$ 控制precision,  $\nu$ 被称为 *degrees of freedom*, 当 $\nu=1$ 的时候, 就reduce成了 *Cauchy distribution*; 当 $\nu \rightarrow \infty$ 的时候, 就成了  $\mathcal{N}(x \mid \mu, \lambda^{-1})$  的 Gaussian distribution

相比于gaussian distribution, student's t-distribution具有一个重要的robustness性质



绿色curve表示Gaussian distribution, 可以看到general的student's t-distribution具有更长的'tails'



这张图中绿色curve表示Gaussian, 可以看到在没有outliers的时候, 两者几乎一样, 而加入一些outliers后绿色curve就明显affected

而第一章中已经说明了MLE就是基于noise服从Gaussian的解。也就是最小二乘并不具有robustness。

再做个变量替换 $\eta = \tau b/a$ , t-distribution就能写成:

$$\operatorname{St}(x \mid \mu, \lambda, \nu) = \int_0^\infty \mathcal{N}(x \mid \mu, \eta \lambda) \operatorname{Gam}(\eta \mid \nu/2, \nu/2) d\eta$$

generalize下得到多变量的distribution:

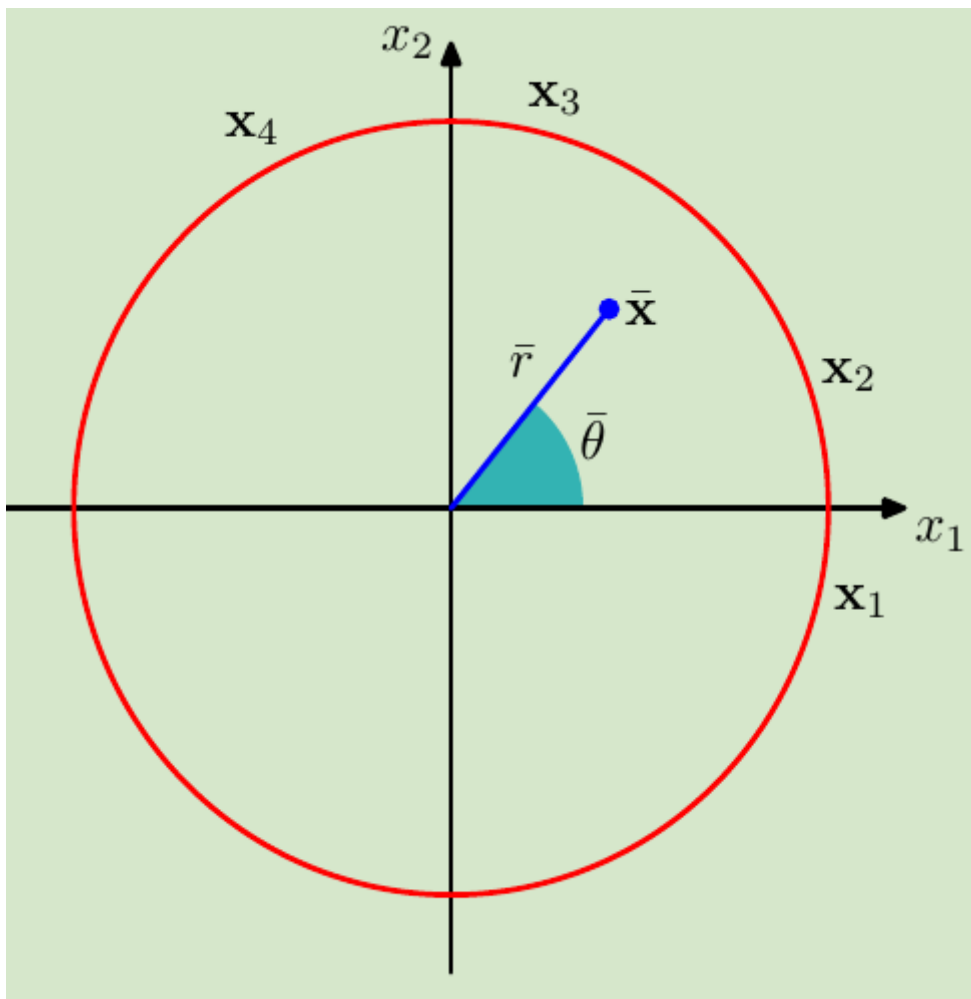
$$\int_0^\infty \left( \int \mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu \right) \left( \int \mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu \right)^{-1} d\boldsymbol{\mu} \\ \int \mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu = \frac{\Gamma(D/2 + \nu/2)}{\Gamma(\nu/2)} \frac{\boldsymbol{\Lambda}^{1/2}}{(\pi \nu)^{D/2}} \left[ 1 + \frac{\Delta^2}{\nu} \right]^{-D/2 - \nu/2}$$

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

### 2.3.8 Periodic variables

尽管Gaussian distribution具有非常重要的实际意义，对于一个periodic 变量，简单通过gaussian去measure并不理想，即使选择一些新的coordinate origin。

例如对于一个direction变量 $\theta$ ，有观测 $\mathcal{D} = \{\theta_1, \dots, \theta_N\}$ ，但 $(\theta_1 + \dots + \theta_N) / N$ 简单平均将会对coordinate强烈依赖。现在需要找一个invariant measure，考虑现在每个observations都是二维（更高维可类推）空间中单位圆上的points，也对应着二维的单位向量 $\mathbf{x}_1, \dots, \mathbf{x}_N$ 。有 $\mathbf{x}_n = (\cos \theta_n, \sin \theta_n)$





所有 $\{\mathbf{x}_N\}$ 的mean可以表示为 $\overline{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$

借助 $\mathbf{x}$ , 可以进而得到 $\theta$ 的mean.  $\overline{\mathbf{x}} = (\bar{r} \cos \bar{\theta}, \bar{r} \sin \bar{\theta})$ 可以得到

$$\bar{\theta} = \tan^{-1} \left( \frac{\sum_n \sin \theta_n}{\sum_n \cos \theta_n} \right)$$

定义在periodic variable上的概率density  $p(\theta)$ 应该满足约束:

$$\begin{aligned} p(\theta) &\geq 0 \quad \int_0^{2\pi} p(\theta) d\theta = 1 \quad p(\theta+2\pi) = p(\theta) \end{aligned}$$

现在我们给出一个在周期变量上的Gaussian-like distribution,

$$p(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp \left\{ -\frac{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2}{2\sigma^2} \right\}$$

现在作变量替换 $x_1 = r \cos \theta, \quad x_2 = r \sin \theta, \quad \mu_1 = r_0 \cos \theta_0, \quad \mu_2 = r_0 \sin \theta_0$ ,  $m = r_0 / \sigma^2$ , 得到:

$$p(\theta \mid \theta_0, m) = \frac{1}{2\pi I_0(m)} \exp \left\{ m \cos(\theta - \theta_0) \right\}$$

这被称为*von Mises distribution, or the circular normal.*, 其中 $\theta_0$ 就是distribution的mean,  $m$ 是*concentration* parameter, 类似于Gaussian的precision. 其中的 $I_0(m)$ 是第一类的零阶Bessel function:  $I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp \{ m \cos \theta \} d\theta$

当 $m$ 很大的时候就近似于Gaussian.

现在考虑von Mises distribution的MLE:

$$\ln p(\mathcal{D} \mid \theta_0, m) = -N \ln(2\pi) - N \ln I_0(m) + m \sum_{n=1}^N \cos(\theta_n - \theta_0) \\ \sum_{n=1}^N \sin(\theta_n - \theta_0) = 0$$

可解得:

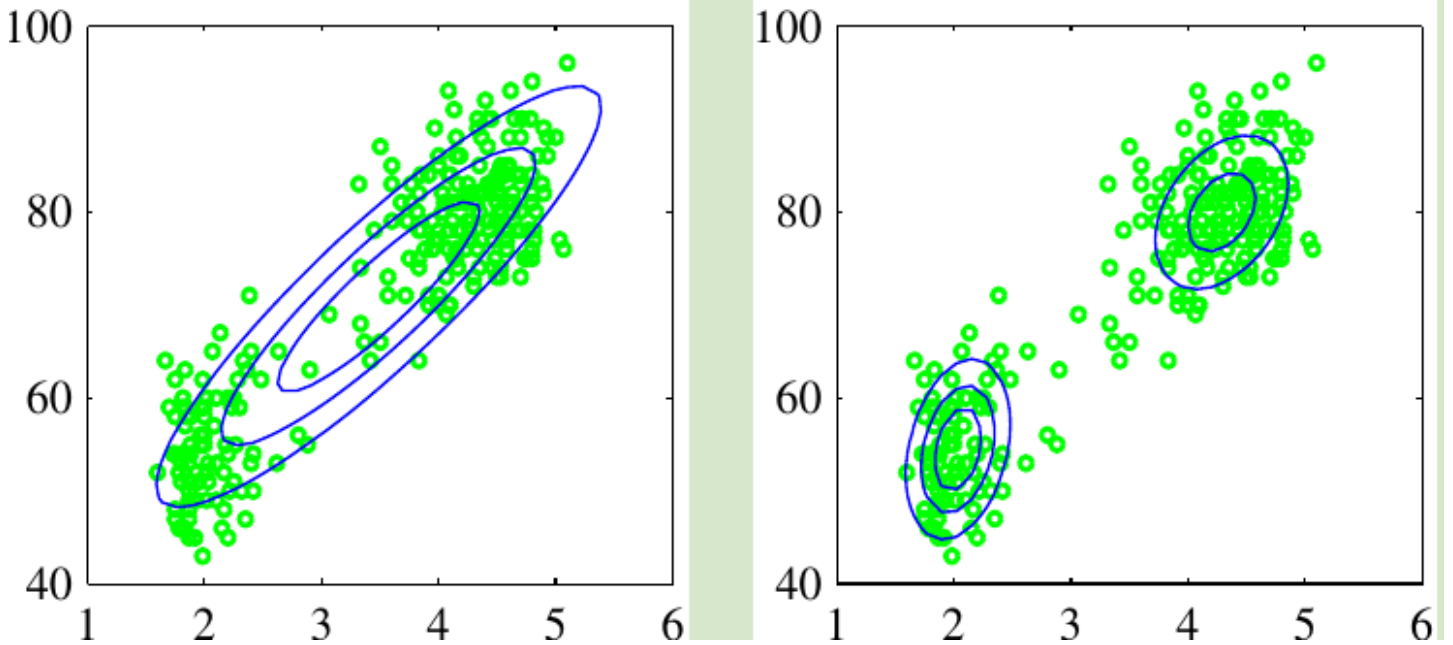
$$\theta_0^{\text{ML}} = \tan^{-1} \left( \frac{\sum_n \sin \theta_n}{\sum_n \cos \theta_n} \right)$$

$m$ 同理:

$$A(m) = \frac{1}{N} \sum_{n=1}^N \cos(\theta_n - \theta_0^{\text{ML}}) \\ A(m) = \frac{I_1(m)}{I_0(m)} \quad A(m_{\text{ML}}) = \left( \frac{1}{N} \sum_{n=1}^N \cos \theta_n \right) \cos \theta_0^{\text{ML}} - \left( \frac{1}{N} \sum_{n=1}^N \sin \theta_n \right) \sin \theta_0^{\text{ML}}$$

### 2.3.9 Mixtures of Gaussians

尽管Gaussian具有一些重要的analytical性质，但遇到real data sets时收到limitation。



绿色的point为真实的data point，左图试图使用一个Gaussian去拟合这个distribution，但是显然真实data point拥有两个clump，这样并不合理；右图使用两个Gaussian去拟合data，看上去就更加合理

我们因此考虑两个distribution的叠加：

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

这被称为*mixture of Gaussians*，其中

$$\sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1,$$

对于已知data去估计distribution的parameter，也就是posterior，往往使用MLE解：

$$\ln p(\mathbf{X} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left( \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

但是由于这个式子里面出现了logarithm中带summation，也就无法获得closed-form analytical solution. 一种方法是利用迭代的数值求解技术，后面会介绍*expected maximization*.

## 2.4 The Exponential Family

截止目前本章中的概率分布(除了Gaussian mixture)都属于*exponential family*。给定parameter  $\boldsymbol{\eta}$ ，在 $\mathbf{x}$ 上的exponential family的分布定义为：

$$p(\mathbf{x} \mid \boldsymbol{\eta}) = h(\mathbf{x}) g(\boldsymbol{\eta}) \exp \left( \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \right)$$

其中的 $\eta$ 被成为分布的 *natural parameters*, function  $g(\eta)$  是 normalize 系数。

首先看一个属于 exponential family 分布的例子, Bernoulli distribution:

$$p(x \mid \mu) = \text{Bern}(x \mid \mu) = \mu^x (1-\mu)^{1-x}.$$

通过对右边项取 logarithm 再 exponential, 可得:

$$p(x \mid \eta) = \frac{1}{\sigma(\eta)} \exp(\eta x)$$

其中

$$\eta = \ln \left( \frac{\mu}{1-\mu} \right) \quad \sigma(\eta) = \frac{1}{1 + \exp(-\eta)}$$

$\sigma$  被称为 *logistic sigmoid function*

接下来再看 multinomial distribution:

$$p(\mathbf{x} \mid \boldsymbol{\mu}) = \prod_{k=1}^M \mu_k^{x_k} = \exp \left( \sum_{k=1}^M x_k \ln \mu_k \right)$$

$$\text{换个写法就是: } p(\mathbf{x} \mid \boldsymbol{\eta}) = \exp \left( \boldsymbol{\eta}^T \mathbf{x} \right)$$

需要注意的是这里面只有  $M-1$  个自由度, 最后一个  $\mu_M$  由其他  $\mu_k$  确定, 可以把它拆分开:

$$\begin{aligned} \exp \left( \sum_{k=1}^M x_k \ln \mu_k \right) &= \exp \left( \sum_{k=1}^{M-1} x_k \ln \mu_k + \left( 1 - \sum_{k=1}^{M-1} x_k \right) \ln \left( 1 - \sum_{k=1}^{M-1} \mu_k \right) \right) \\ &= \exp \left( \sum_{k=1}^{M-1} x_k \ln \left( \frac{\mu_k}{1 - \sum_{j=1}^{M-1} \mu_j} \right) + \ln \left( 1 - \sum_{k=1}^{M-1} \mu_k \right) \right). \end{aligned}$$

$$\ln \left( \frac{\mu_k}{1 - \sum_{j=1}^{M-1} \mu_j} \right) = \eta_k$$

可以得到:

$$\mu_k = \frac{\exp(\eta_k)}{1 + \sum_{j=1}^{M-1} \exp(\eta_j)}$$

也被称为 *softmax function* 或 *normalized exponential*. 如此就有形式:

$$p(\mathbf{x} \mid \boldsymbol{\eta}) = \frac{\exp \left( \boldsymbol{\eta}^T \mathbf{x} \right)}{\sum_{k=1}^M \exp(\eta_k)}.$$

最后再看 Gaussian distribution:

$$\begin{aligned} p(x \mid \mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (x-\mu)^2 \right\} \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} x^2 + \frac{\mu}{\sigma^2} x - \frac{1}{2\sigma^2} \mu^2 \right\} \end{aligned}$$

经过整理：

$$\begin{aligned} \boldsymbol{\eta} &= \left( \begin{array}{c} \mu / \sigma^2 \\ -1 / 2 \sigma^2 \end{array} \right) \\ \mathbf{u}(x) &= \left( \begin{array}{c} x \\ x^2 \end{array} \right) h(\mathbf{x}) = (2\pi)^{-1/2} g(\boldsymbol{\eta}) = \left( -2 \boldsymbol{\eta}_2 \right)^{1/2} \exp \left( \frac{\boldsymbol{\eta}_1^2}{4 \boldsymbol{\eta}_2} \right) \end{aligned}$$

### 2.4.1 Maximum likelihood and sufficient statistics

由： $g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \right\} d\mathbf{x} = 1$ ，对两边关于 $\boldsymbol{\eta}$ 求导，得： $-\nabla \ln g(\boldsymbol{\eta}) = \mathbb{E}[\mathbf{u}(\mathbf{x})]$ 。而对于 $\mathbf{u}_i(\mathbf{x})$ 的covariance会是 $g$ 的二阶微分表达式，对于高阶矩类似性质。

现在考虑满足i.i.d的data  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  likelihood function：

$$p(\mathbf{X} \mid \boldsymbol{\eta}) = \left( \prod_{n=1}^N h(\mathbf{x}_n) \right) g(\boldsymbol{\eta})^N \exp \left\{ \boldsymbol{\eta}^T \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) \right\}$$

让其最大可得到： $-\nabla \ln g(\boldsymbol{\eta}_{\text{ML}}) = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$ ，原则上可以通过这个式子来解出 $\boldsymbol{\eta}_{\text{ML}}$ ，可以看到它只依赖于 $\sum_n \mathbf{u}(\mathbf{x}_n)$ （称为sufficient statistic）。意味着不需要存储整个data set。

note：当 $N \rightarrow \infty$ 时，右边项就是 $\mathbb{E}[\mathbf{u}(\mathbf{x})]$ ， $\boldsymbol{\eta}_{\text{ML}}$ 也会等于true value  $\boldsymbol{\eta}$ 。

### 2.4.2 Conjugate priors

Conjugate priors是指对给定的一个 $p(\mathbf{x} \mid \boldsymbol{\eta})$ ，我们寻找一个prior  $p(\boldsymbol{\eta})$ 与prior共轭，使得posterior和prior具有相同的形式：

$$p(\boldsymbol{\eta} \mid \boldsymbol{\chi}, \nu) = f(\boldsymbol{\chi}, \nu) g(\boldsymbol{\eta})^\nu \exp \left\{ \nu \boldsymbol{\eta}^T \boldsymbol{\chi} \right\}$$

其中 $f(\boldsymbol{\chi}, \nu)$ 是normalize项。由此得到的posterior是：

$$p(\boldsymbol{\eta} \mid \mathbf{X}, \boldsymbol{\chi}, \nu) \propto g(\boldsymbol{\eta})^{\nu+N} \exp \left\{ \boldsymbol{\eta}^T \left( \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) + \nu \boldsymbol{\chi} \right) \right\}$$

可以看到确实与prior是同样的形式。

### 2.4.3 Noninformative priors

在inference的时候往往需要用到prior, 有时候prior的选择会对最终的posterior有巨大影响。距离来说如果prior在某些值的概率为0,那么无论观测到怎样的data, 得到的posterior对应位置都是0.但是在很多情况下, 我们可能对分布的形式完全不知道, 这时我们就需要一种prior可以对posterior产生尽可能小的影响 (这就是*noninformative prior*) 。

假设现在有一个由 $\lambda$ 控制的distribution  $p(x|\lambda)$ , 现在就让 $\lambda$ 的prior是个constant。如果 $\lambda$ 是个有K个状态的离散变量, 我们就取每种状态的概率是 $1/K$ 。但在连续的情况下有两个问题。

1. 连续的prior可能无法被normalize, 因为对 $\lambda$ 的积分是发散的 (考虑取值无界情况) .这被称为*improper*. 实际中, 如果posterior是proper的, 那么可以采用这种improper的prior。例如, Gaussian的mean采用均匀分布, 但只要有一个观测, posterior就是正常的。
2. 变量的非线性变换对概率密度的影响。如果一个函数 $h(\lambda)$ 是constant, 那么做变量替换 $\lambda = \eta^2$ 仍然是个constant, 有 $\hat{h}(\eta) = h(\eta^2)$ 。但是如果 $p_{\lambda}(\lambda)$ 是constant, 那么
 
$$p_{\eta}(\eta) = p_{\lambda}(\lambda) \left| \frac{d\lambda}{d\eta} \right| = p_{\lambda}(\eta^2) 2\eta \propto \eta$$

从而关于 $\eta$ 的概率密度就不是constant了。但是当我们使用maximum likelihood的时候这并不会发生 (一般都是simple function) 。

## 2.5 Nonparametric Methods

前面的都是通过观测data来估计一个distribution中的一部分参数, 被称为*parametric method*。这样的 limitation就是如果使用一个poor distribution描述data, 就会产生糟糕的结果。

这里介绍nonparametric methods, 几乎不对distribution的形式做假设。一个最简单的方法就是直方图估计。将取值范围划分为N个小bins, 每个bins的宽度是 $\Delta_i$ , data落在第i个bin内的数量就是 $n_i$ , 由此估计data的概率密度为:

$$p_i = \frac{n_i}{N \Delta_i}$$

直方图估计有个limitation就是选取bins的宽度并不容易, 取太大太小都不行。还有就是受维度诅咒的影响。

### 2.5.1 Kernel density estimators

假设data是从未知distribution  $p(\mathbf{x})$ 中sample出来的, 对于包含 $\mathbf{x}$ 的区域 $\mathcal{R}$ 来说, probability mass为:

$$P = \int_{\mathcal{R}} p(\mathbf{x}) d\mathbf{x}$$

现在每个data point落在 $\mathcal{R}$ 上的事件看作服从概率为P的binomial distribution:

$$\text{Bin}(K \mid N, P) = \frac{N!}{K!(N-K)!} P^K (1-P)^{N-K}$$

因为 $\mathbb{E}[K]=NP$ , 对于一个大的 $N$ 有 $K \simeq NP$ , 对于一个小的 $\mathcal{R}$ 有 $N \simeq p(\mathbf{x}|V)$ , ( $V$ 是 $\mathcal{R}$ 的volume), 进而:  

$$p(\mathbf{x}) = \frac{K}{N V}$$

有这个式子就可以估计distribution.

note: 两种方法, 一种先fix  $K$ , 再确定 $V(K\text{-nearest-neighbour})$ ; 一种先fix  $V$ , 再确定 $K$ .

这里先考虑第二种方法。为了确定落在 $\mathcal{R}$ 中的 $K$ , 先定义一个 *kernel function*:

$$k(\mathbf{u}) = \begin{cases} 1, & \text{if } |\mathbf{u}| \leq 1/2 \\ 0, & \text{otherwise} \end{cases}$$

确定 $K$ :

$$K = \sum_{n=1}^N k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

其中 $h$ 是个超参, 这样就能估计distribution:

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

其实 *kernel function* 也可以取的smooth些(区别在cube的boundary上), 例如Gaussian kernel function, 会得到:

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{D/2}} \exp\left(-\frac{|\mathbf{x} - \mathbf{x}_n|^2}{2h^2}\right)$$

这个方法的limitation就是 $h$ 也不能取太大或太小。

## 2.5.2 Nearest-neighbour methods

这节介绍fix  $K$ , 然后再确定 $V$ 的方法。假设我们需要估计 $\mathbf{x}$ 的概率密度, 我们以point  $\mathbf{x}$ 为center作一个sphere, 使得这个sphere正好包入 $K$ 个data point, 如此就能得到 $V$ . 这样的方法就被称为 *K-nearest-neighbours*

note: 这样得到的 $p(\mathbf{x})$ 并不能保证积分后是1.

我们使用这个方法来解决一个分类问题。现在有 $N$ 个data point, 属于 $\mathcal{C}_k$ 类的有 $N_k$ 个data point. 那么conditional probability:

$$p(\mathbf{x} | \mathcal{C}_k) = \frac{K_k}{N_k V}$$

相似地, unconditional probability:

$$p(\mathbf{x}) = \frac{K}{N V}$$

prior为 $p(\text{mathcal{C}}_k) = \frac{N_k}{N}$ . 结合这几个式子可以得到:

$$p(\text{mathcal{C}}_k \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \text{mathcal{C}}_k) p(\text{mathcal{C}}_k)}{p(\mathbf{x})} = \frac{K_k}{K}.$$

通过最小化misclassify的概率, 可知解就是在这K个点中出现次数最多的class (maximize posterior). 特别地, 当K=1时, 就是*nearest-neighbour*.