

贝叶斯统计

Bayesian Statistics

慕课网《行为程序员设计的统计学》
讲师：angelaychen
版权所有 侵权必究

- $H_0: \mu = 1w$
- $H_A: \mu > 1w$



H_0 为真



Data



$P(\text{data} \mid H_0 \text{为真})$



- $H_1: \mu = 1w$
- $H_2: \mu > 1w$
- $H_3: \mu < 1w$



Data



$P(H_1 \text{为真} \mid \text{data})$
 $P(H_2 \text{为真} \mid \text{data})$
 $P(H_3 \text{为真} \mid \text{data})$



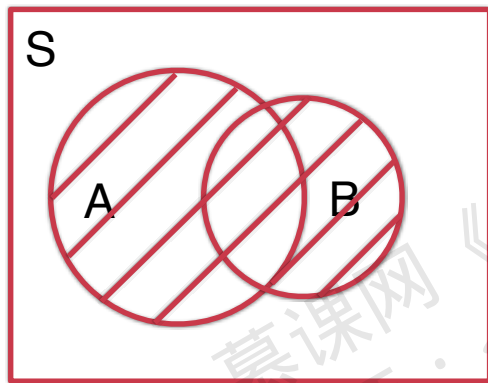
决策

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

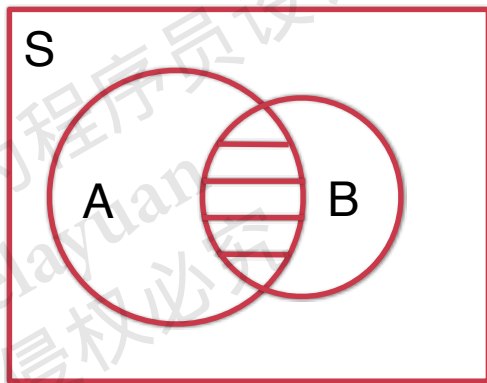
概率知识

事件间的关系

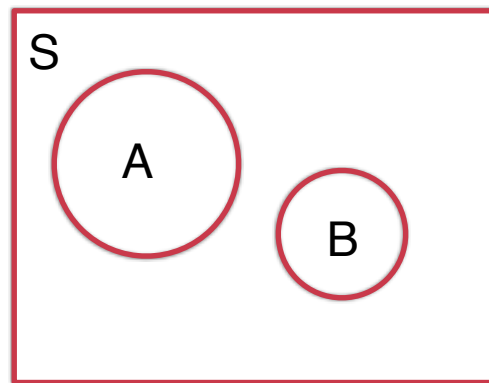
和事件 $A \cup B$



积事件 $A \cap B$



空集 $A \cap B = \emptyset$



A 与 B 是互不相容的/互斥的

条件概率和乘法定理

设A, B是两个事件, 且 $P(A) > 0$

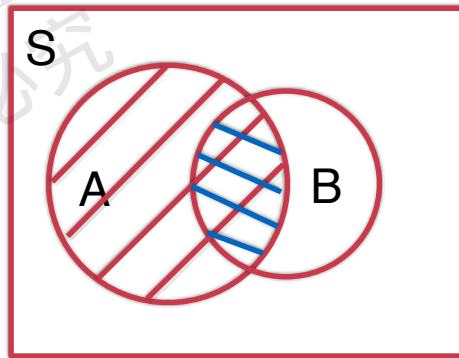
则称 $P(B|A) = \frac{P(AB)}{P(A)}$ 为在事件A发生的条件下

事件B发生的条件概率

$$P(AB) = P(B|A)P(A)$$

乘法定理

$$P(B) > 0, P(AB) = P(A|B)P(B)$$

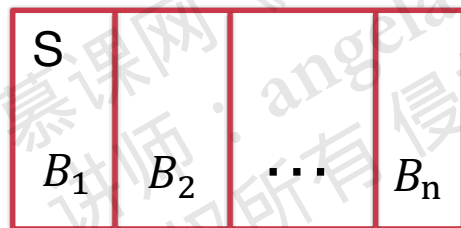


全概率公式

S 为试验 E 的样本空间, B_1, B_2, \dots, B_n 为 E 的一组事件

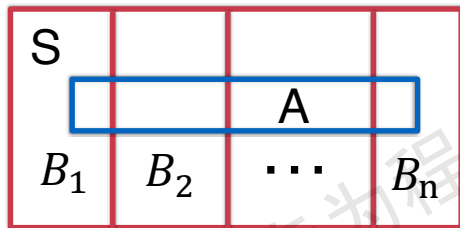
$B_i B_j = \emptyset \ (i \neq j)$ 且 $B_1 \cup B_2 \cup \dots \cup B_n = S$

称 B_1, B_2, \dots, B_n 为样本空间 S 的一个划分



对每次试验, 事件 B_1, B_2, \dots, B_n
必有一个且仅有一个发生

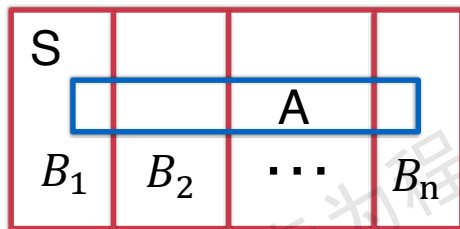
全概率公式



$$P(B_i) > 0 \quad (i = 1, 2, \dots, n)$$

$$P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots + P(A|B_n)P(B_n)$$

贝叶斯公式



$$P(A) > 0$$

$$P(B_i) > 0 \quad (i = 1, 2, \dots, n)$$

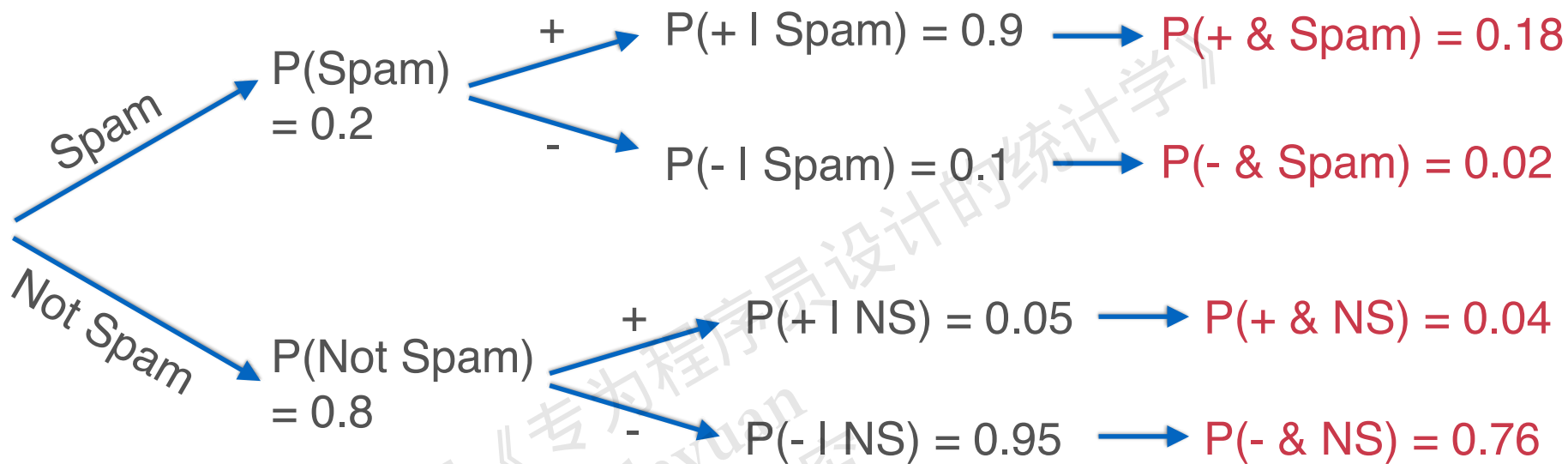
$$P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots + P(A|B_n)P(B_n)$$

$$P(B_i|A) = \frac{P(B_i A)}{P(A)} = \frac{P(A|B_i)P(B_i)}{\sum_{j=1}^n P(A|B_j)P(B_j)}$$

贝叶斯公式

概率树 (probability tree)

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

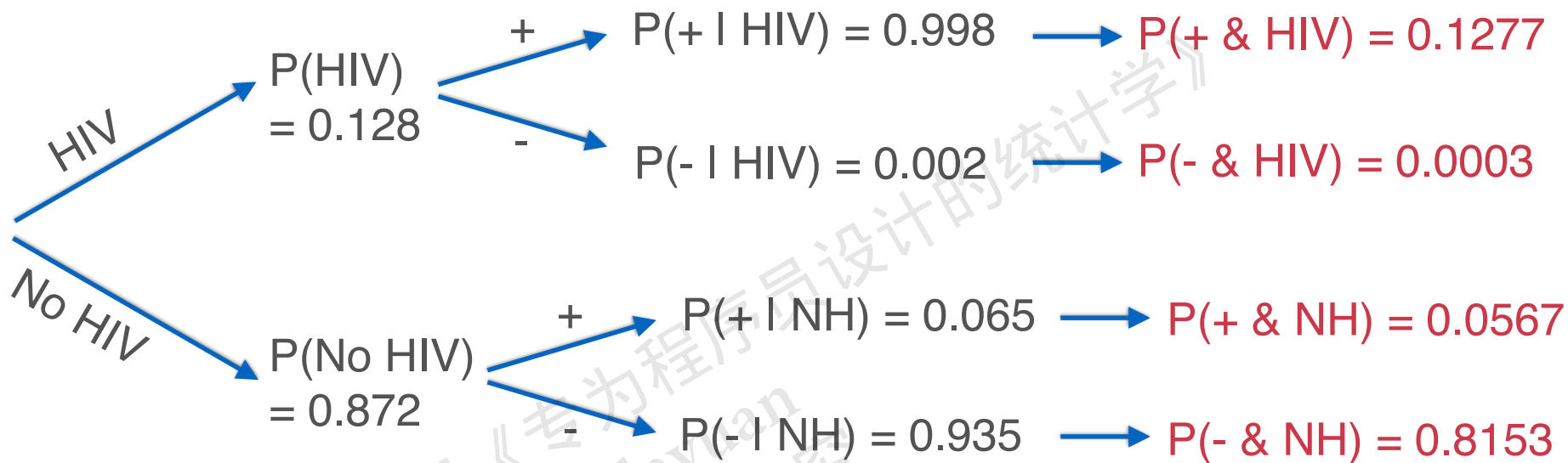


$$P(\text{Spam} | +) = 0.18 / (0.18 + 0.04) = 0.818$$

$$P(\text{NS} | +) = ?$$

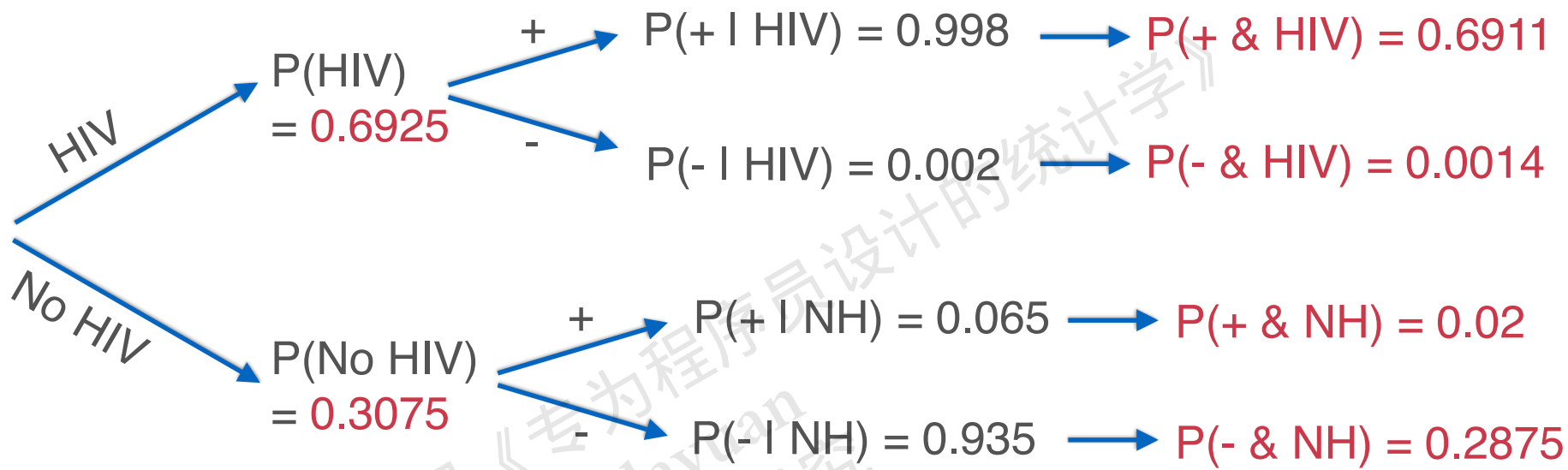
$$P(\text{Spam} | -) = 0.02 / (0.02 + 0.76) = 0.026$$

$$P(\text{NS} | -) = ?$$



$$P(\text{HIV} | +) = 0.1277 / (0.1277 + 0.0567) = 0.6925$$

$$P(\text{NH} | -) = 0.8153 / (0.8153 + 0.0003) = 0.9996$$



$$P(\text{HIV} | +) = 0.6911 / (0.6911 + 0.02) = 0.972$$

$$P(\text{NH} | -) = 0.2875 / (0.2875 + 0.0014) = 0.995$$

贝叶斯推断

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

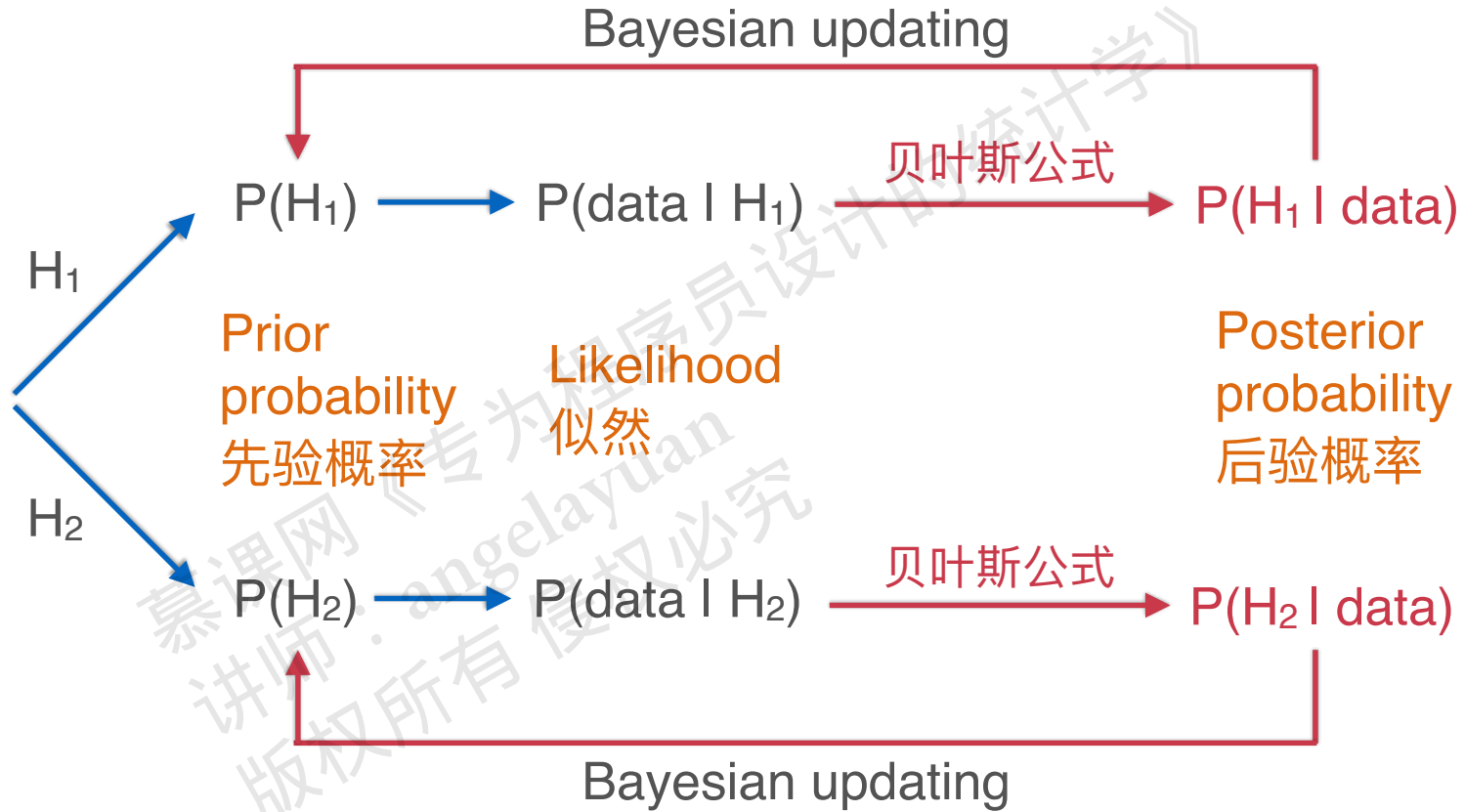
$$P(B_i|A) = \frac{P(A|B_i)P(B_i)}{\sum_{j=1}^n P(A|B_j)P(B_j)}$$

$A = \text{数据(data)}$

$B_i = \text{第}i\text{个假设}H_i$

$$P(H_i|\text{data}) = \frac{P(\text{data}|H_i)P(H_i)}{\sum_{j=1}^n P(\text{data}|H_j)P(H_j)}$$

$P(\text{data} | H_i) \longrightarrow P(H_i | \text{data})$



Posterior probability

Likelihood

Prior probability

$$P(H_i|\text{data}) = \frac{P(\text{data}|H_i)P(H_i)}{\sum_{j=1}^n P(\text{data}|H_j)P(H_j)}$$

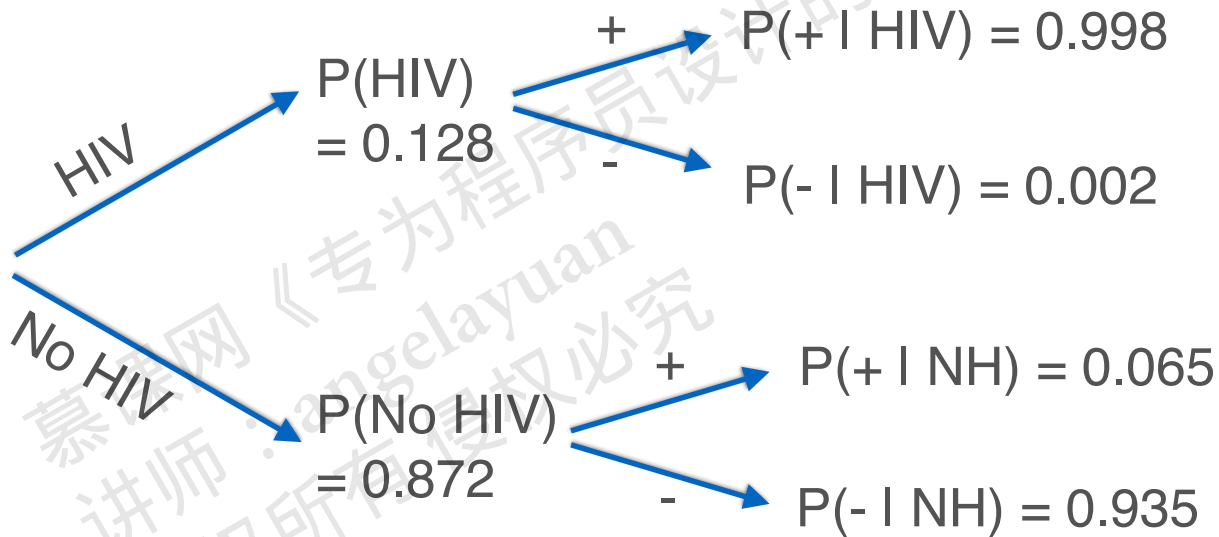
Marginal likelihood

The diagram illustrates the components of Bayes' theorem. It features the equation $P(H_i|\text{data}) = \frac{P(\text{data}|H_i)P(H_i)}{\sum_{j=1}^n P(\text{data}|H_j)P(H_j)}$. Four red arrows point from labels to specific parts of the equation: 'Posterior probability' points to $P(H_i|\text{data})$, 'Likelihood' points to $P(\text{data}|H_i)$, 'Prior probability' points to $P(H_i)$, and 'Marginal likelihood' points to the denominator $\sum_{j=1}^n P(\text{data}|H_j)P(H_j)$.

贝叶斯推断 - 例1

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

贝叶斯推断 - 例1



贝叶斯推断 - 例1

Hypothesis/Model	HIV	No HIV	Total
Prior: $P(\text{model})$	0.128	0.872	1
Data	+		
Likelihood $P(\text{data} \mid \text{model})$	0.998	0.065	
$P(\text{data} \mid \text{model}) \times P(\text{model})$	0.1277	0.0567	0.1844
Posterior: $P(\text{model} \mid \text{data})$	0.6925	0.3075	1

贝叶斯推断 - 例1

Hypothesis/Model	HIV	No HIV	Total
Prior: $P(\text{model})$	0.6925	0.3075	1
Data	+		
Likelihood $P(\text{data} \mid \text{model})$	0.998	0.065	
$P(\text{data} \mid \text{model}) \times P(\text{model})$	0.6911	0.02	0.7111
Posterior: $P(\text{model} \mid \text{data})$	0.972	0.028	1

贝叶斯因子(Bayes Factor)

$$\text{Bayes Factor} = \frac{P(\text{data} \mid H_1)}{P(\text{data} \mid H_2)}$$

数据更支持哪一个假设/模型

$$BF = \frac{P(+ \mid HIV)}{P(+ \mid NH)} = \frac{0.998}{0.065} = 15.35$$

$$BF = \frac{P(- \mid HIV)}{P(- \mid NH)} = \frac{0.002}{0.935} = 0.002$$

贝叶斯因子(Bayes Factor)

$$BF = \frac{P(+ | HIV)}{P(+ | NH)} = 15.35$$

$$BF = \frac{P(- | HIV)}{P(- | NH)} = 0.002$$

BF	解释
< 1	没有证据支持H ₁ (支持H ₂)
1 ~ 3	较弱的证据支持H ₁
3 ~ 10	中等程度的证据支持H ₁
10 ~ 30	较强的证据支持H ₁
30 ~ 100	非常强的证据支持H ₁
> 100	极强的证据支持H ₁

Odds和Odds ratio

Model	HIV	No HIV	odds of HIV = $\frac{P(\text{HIV})}{P(\text{No HIV})}$
Prior: P(model)	0.128	0.872	0.147
Posterior P(model data)	0.6925	0.3075	2.252

$$\text{odds ratio} = \frac{\text{Posterior odds}}{\text{Prior odds}} = 2.252 / 0.147 = 15.4$$

贝叶斯推断 - 例2

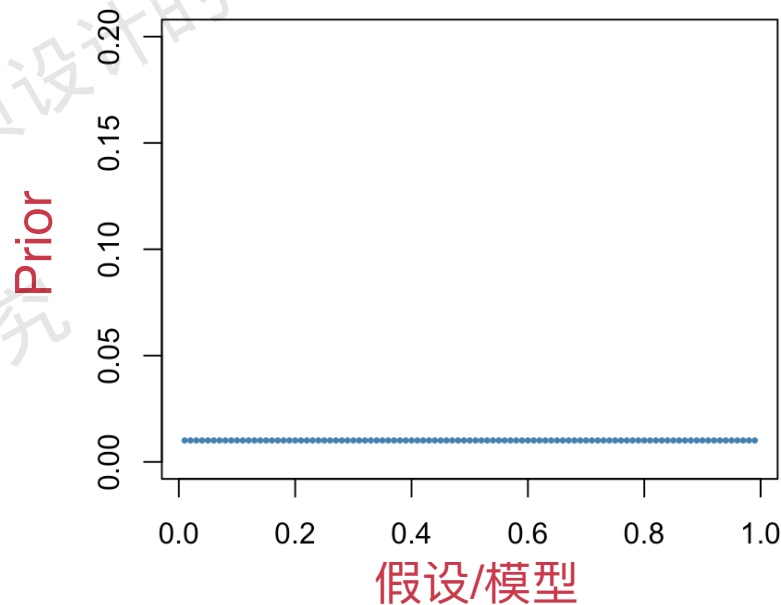
慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

贝叶斯推断 - 例2

- 研究问题: 治疗高血压的新药是否有效
- 方法: 让高血压患者服用新药一段时间, 然后检查高血压症状是否有所改善
- 使用症状得到改善的患者占总患者数的比例来近似代表新药的有效程度
- 数据: 100名高血压患者服用新药, 其中78名患者症状有所改善

贝叶斯推断 - 例2

- 明确假设(模型)和先验概率
 - 假设/模型: 新药的有效性为 0.01, 0.02, ..., 0.99 (以 0.01 为步长)
 - 先验概率: 各模型的概率相等, 均为 $1/99$



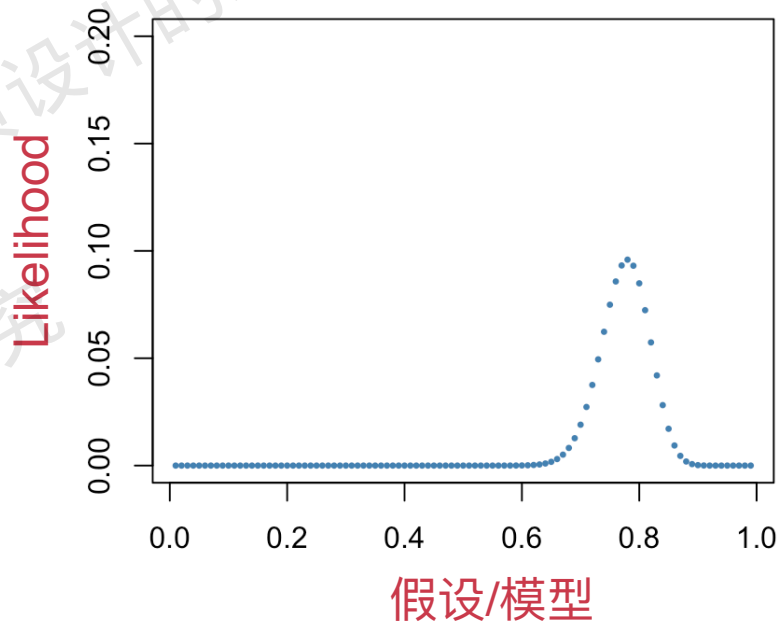
贝叶斯推断 - 例2

- 计算各假设下得到观测数据的可能性(likelihood)

$$P\{X = k\} = \binom{n}{k} p^k (1 - p)^{n-k}$$

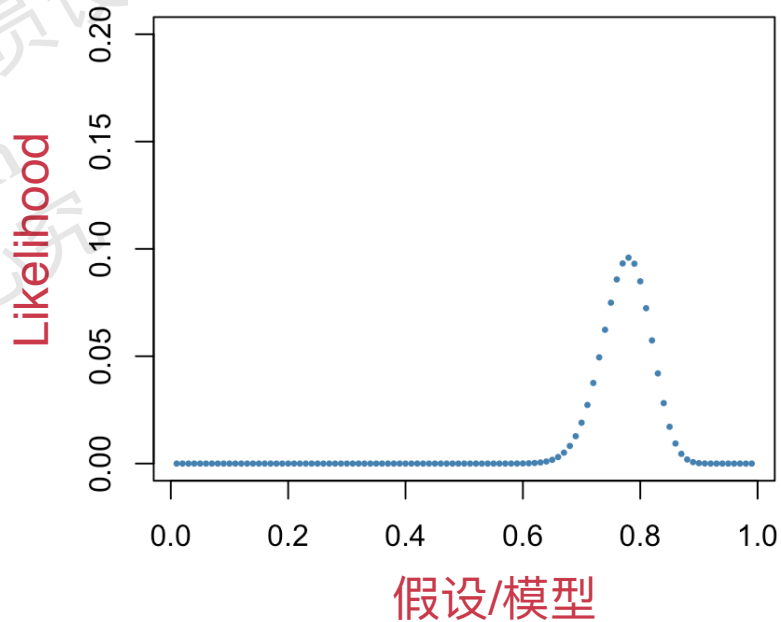
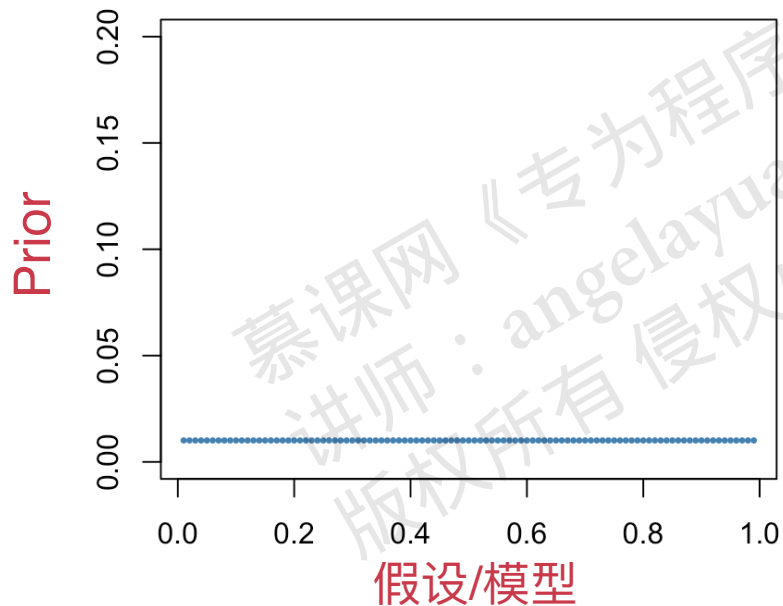
$$(k = 0, 1, 2, \dots, n; 0 < p < 1)$$

$k = 78, n = 100,$
 $p = 0.01, 0.02, \dots, 0.99$



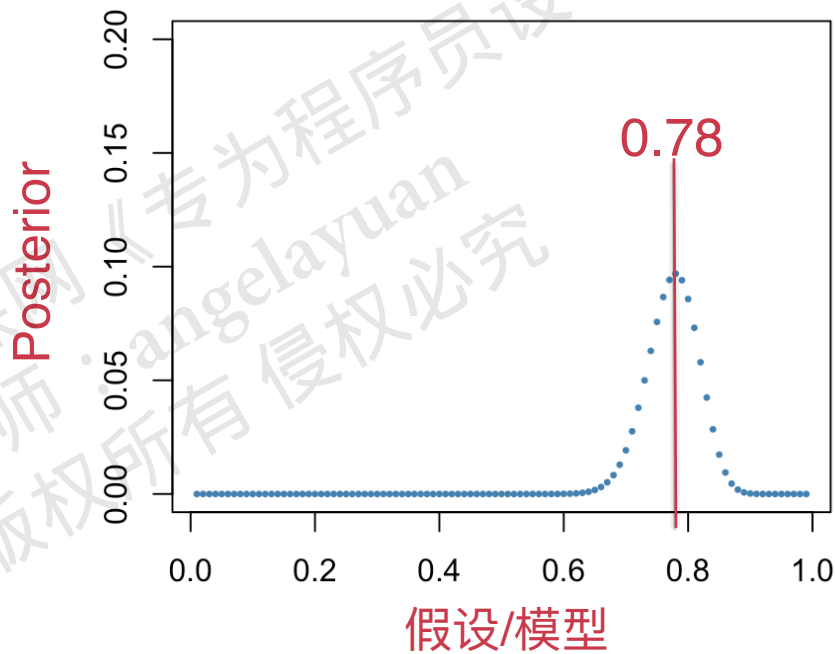
贝叶斯推断 - 例2

- 计算 marginal likelihood = $\text{sum}(\text{prior} \times \text{likelihood})$



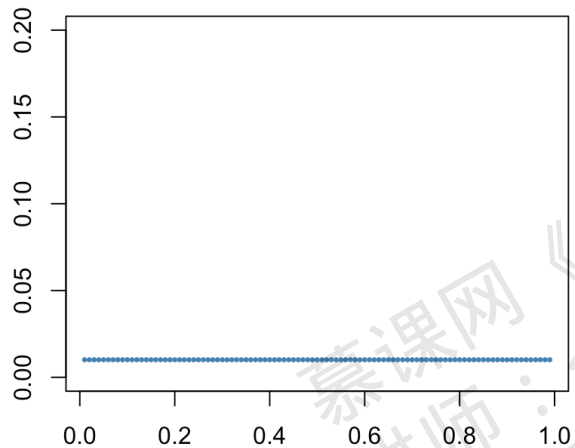
贝叶斯推断 - 例2

- 计算后验概率 = (likelihood x prior) / marginal likelihood



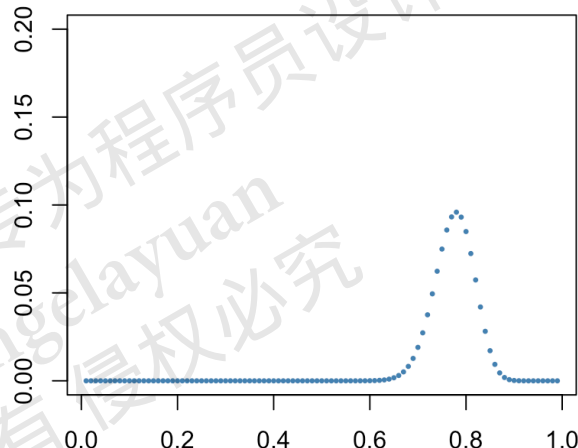
先验概率的选择对后验概率影响

Prior



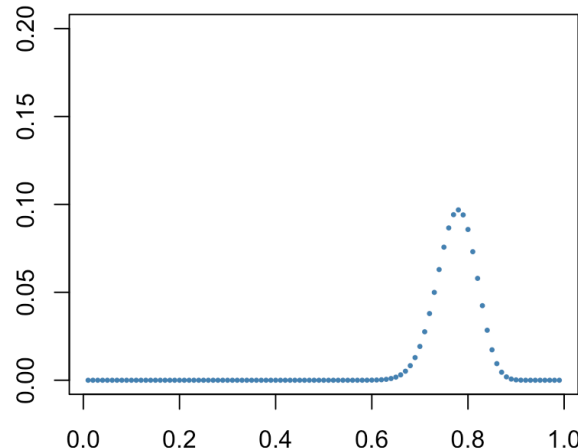
假设/模型

Likelihood



假设/模型

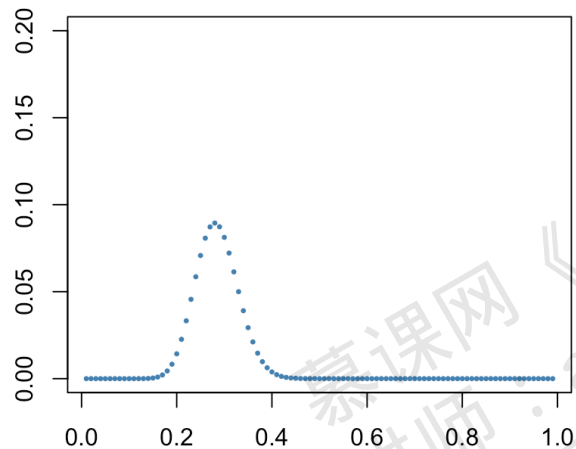
Posterior



假设/模型

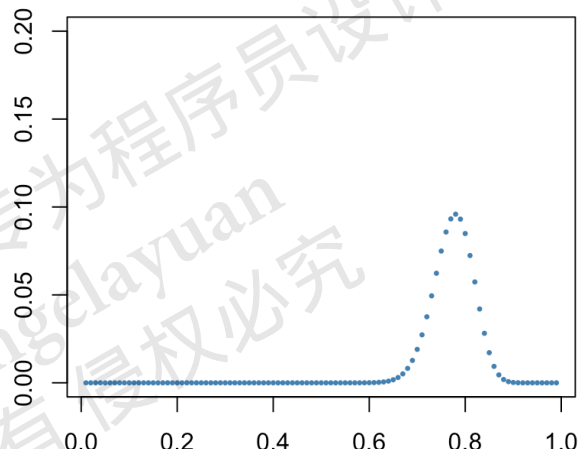
先验概率的选择对后验概率影响

Prior



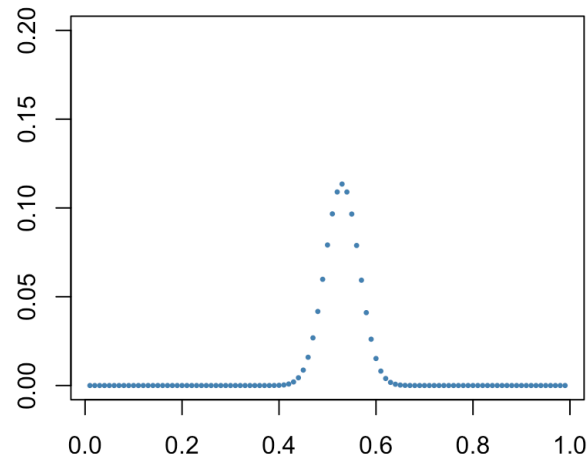
假设/模型

Likelihood



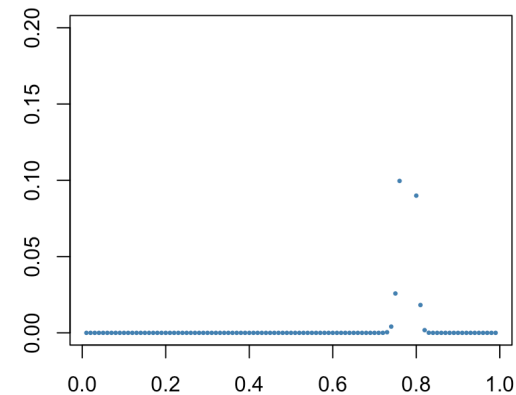
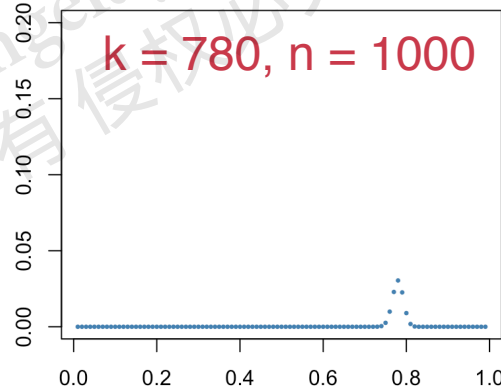
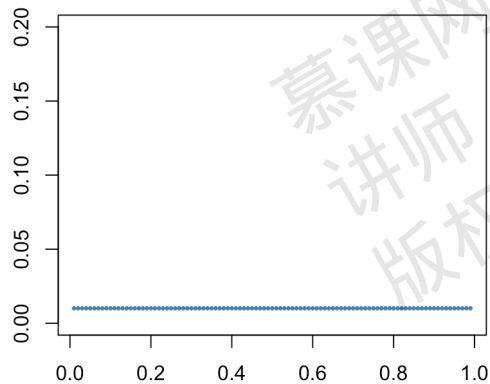
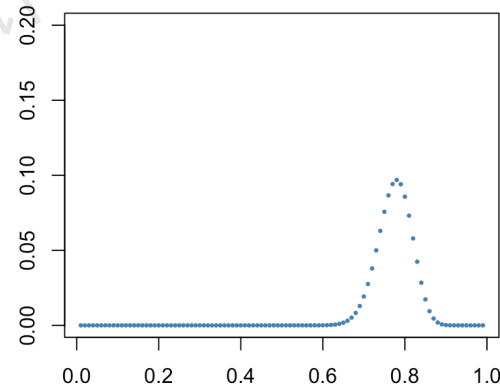
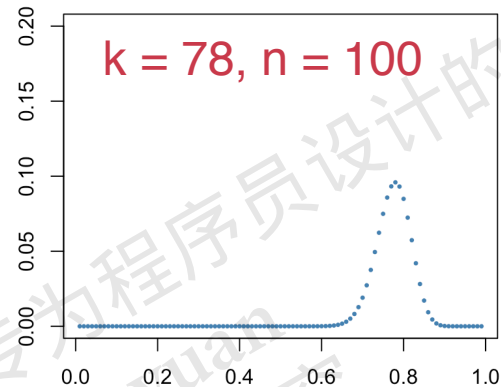
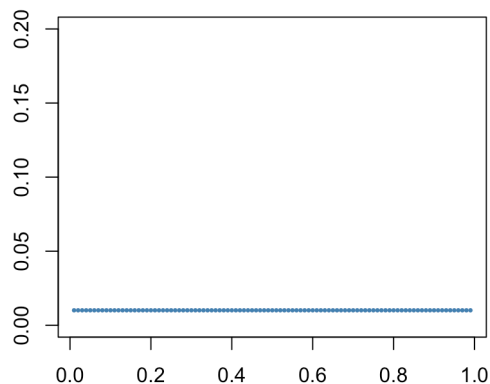
假设/模型

Posterior



假设/模型

样本容量对后验概率的影响



置信区间(credible interval)

慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

置信区间 (credible interval)

- 使用拒绝采样(rejection sampling)方法从后验概率分布中抽样, 然后计算分位数以得到置信区间
 - 抽样次数: $n = 100,000$
 - x 代表 p_{respond}
 - y 代表 p_{respond} 的后验概率

置信区间 (credible interval)

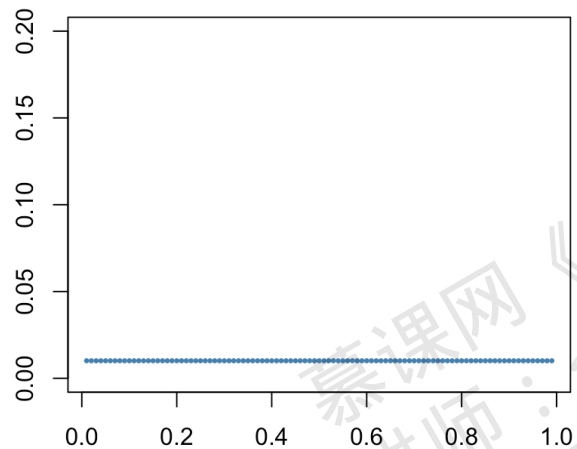
- 重复下列步骤 n 次
 - 从均匀分布 $U(0,1)$ 中抽取1个数, 赋值给 x
 - 从均匀分布 $U(0,1)$ 中抽取1个数, 赋值给 y
 - 从真实的后验概率分布中找到 x 对应的后验概率 $f(x)$
 - 如果 $y < f(x)$, 接受 x , 记录下该 x 的值

置信区间 (credible interval)

- 从所有被接受的 x 值中, 找到两个数值 L 和 H , 使落在这两个数值之间的 x 的值的个数占 x 值总数的95%, 则这两个数值构成的区间 (L,H) 就是一个95%的置信区间
- 置信区间的含义: p_{respond} 落在区间 (L,H) 内的概率是95%

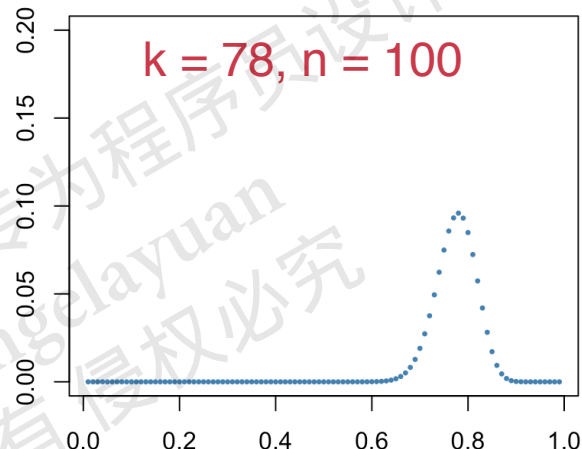
置信区间 (credible interval)

Prior



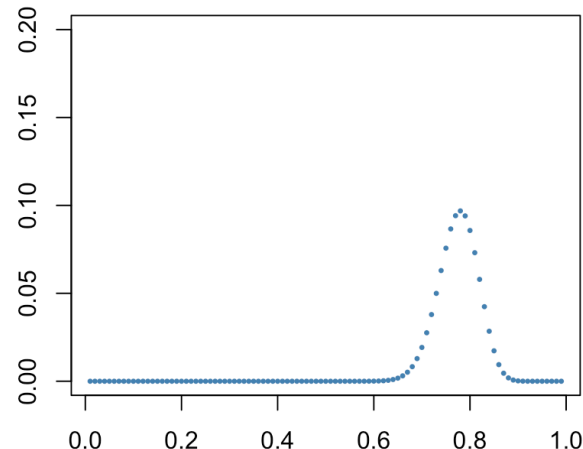
假设/模型

Likelihood



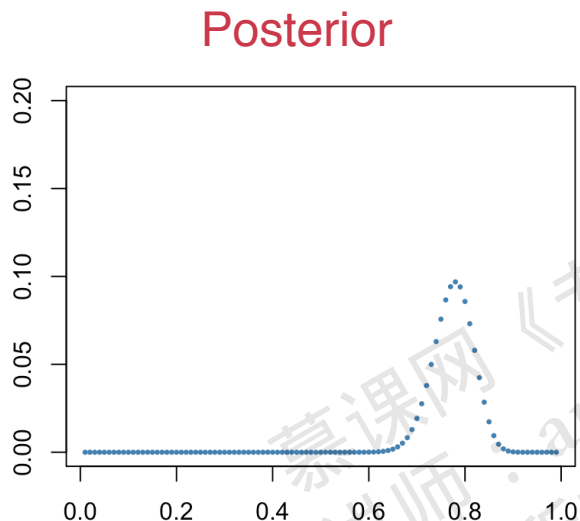
假设/模型

Posterior



假设/模型

置信区间 (credible interval)

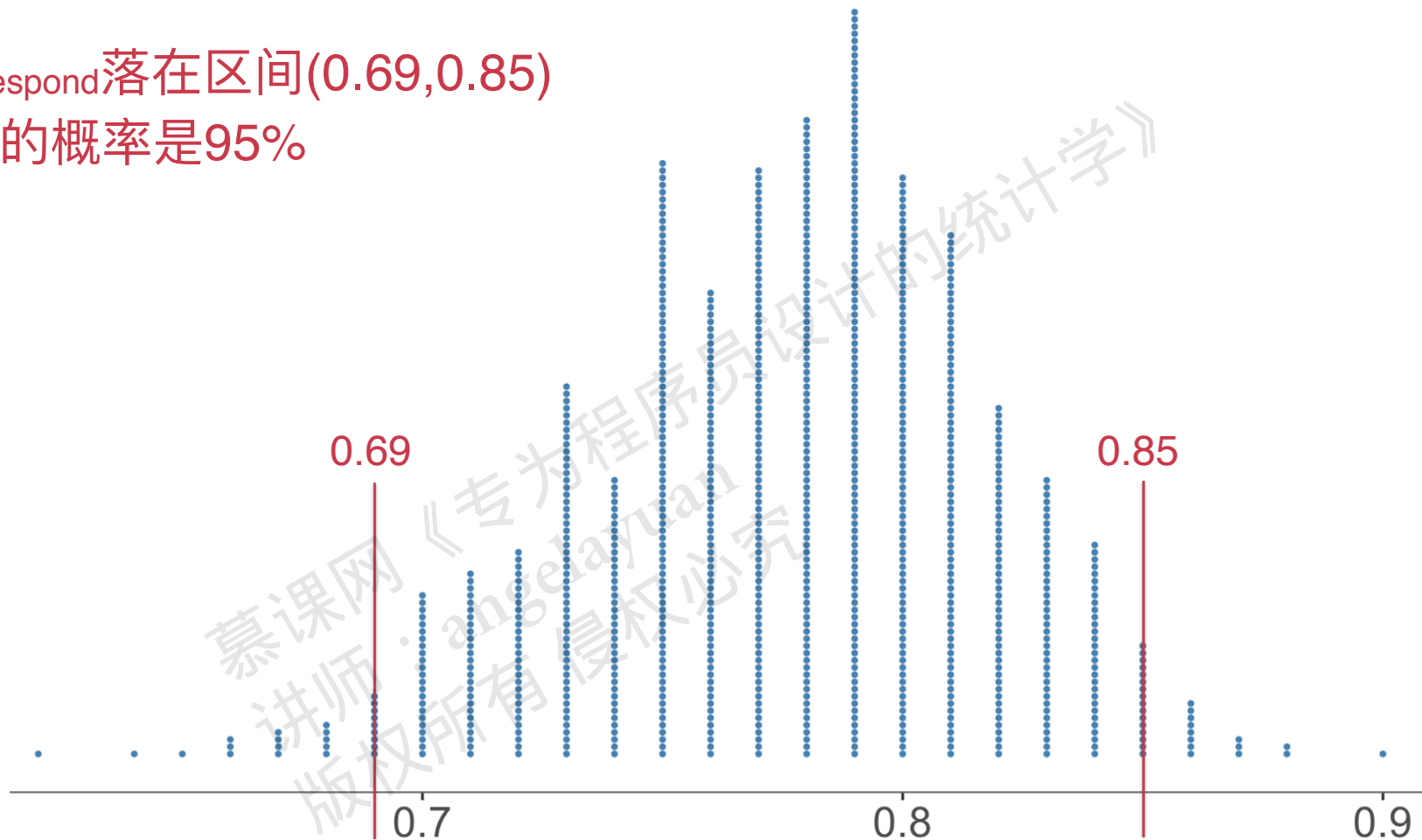


- $x = 0.8, y = 0.3$

- $\text{Posterior}(x = 0.8) = 0.086$

- $y > f(x)$, 拒绝 x

p_{respond} 落在区间(0.69,0.85)
内的概率是95%



慕课网《专为程序员设计的统计学》
讲师：angelayuan
版权所有 侵权必究

本章小结

贝叶斯统计

条件概率, 乘法定理, 全
概率公式, 贝叶斯公式

概率树

贝叶斯推断

- 疾病诊断; 参数估计
- 先验概率, 样本容量对后
验概率的影响

置信区间 (credible interval)

