

重要知识点回顾

◆ 第三章：通过搜索进行问题求解

- 无信息搜索策略
- 有信息（启发式）搜索策略

无信息搜索算法不提供有关某个状态与目标状态的接近程度的任何线索：

- 广度优先搜索 Breadth-first search
- 一致代价搜索 Uniform-cost search
- 深度优先搜索 Depth-first search

广度优先搜索

- 优先扩展最浅的未被扩展的节点
- 边界 (frontier) 可以实现为一个FIFO队列, 即, 新节点 (总是比其父节点更深) 进入队列的队尾, 而旧节点, 即比新节点浅的节点, 首先被扩展。

一致代价搜索 (Dijkstra算法)

- 优先扩展代价最小的未被扩展的节点
- 边界 (frontier) 可以实现为一个按路径代价排序的队列, 最浅层的优先

深度优先搜索

- 优先扩展最深的未被扩展的节点
- 边界 (frontier) 可以实现为一个LIFO队列, 即后进先出。

- **有信息搜索 (informed search) 策略**使用关于目标位置的特定领域线索来比无信息搜索策略更有效地找到解。
- 线索以**启发式函数 (heuristic function)** 的形式出现, 记为 $h(n)$:

$h(n)$ = 从节点 n 的状态到目标状态的最小代价路径的代价估计值

例如, 在寻径问题中, 我们可以通过计算地图上两点之间的直线距离来估计从当前状态到目标的距离。

贪心最佳优先搜索 (greedy best-first search)

- 首先扩展 $h(n)$ 值最小的节点，即看起来最接近目标的节点，因为这样可能可以更快找到解。
- 评价函数 $f(n) = h(n)$

A*搜索

- 主要思想：避免扩展代价已经很高的路径
- 评价函数 $f(n) = g(n) + h(n)$

$g(n)$ = 从初始节点到节点n的路径代价

$h(n)$ = 从节点n的状态到目标状态的最小代价路径的代价估计值

$f(n)$ = 经过n到一个目标状态的最优路径的代价估计值

- 对于**可容许的启发式 (admissible heuristic)** 函数, A*搜索是代价最优的, 即

$$h(n) \leq h^*(n)$$

这里 $h^*(n)$ 经过节点n到目标状态的真实代价. ($h(n) \geq 0$, 对于任意目标G, $h(G) = 0$.)

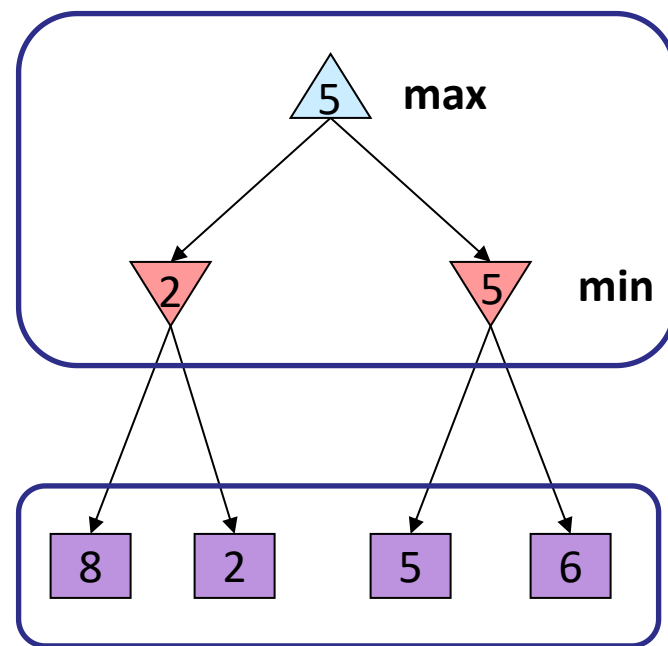
◆ 第五章：对抗搜索和博弈

- 极小化极大搜索
- 评价函数
- α - β 剪枝
- 期望最大搜索

极小化极大搜索

- 确定性的，零和游戏：
 - 井字棋, 象棋, 围棋
 - 一个玩家最大化结果
 - 另一个玩家最小化结果
- 极小化极大搜索：
 - 状态空间搜索树
 - 玩家交替进行操作
 - 计算每个节点的**极小化极大值**:
针对对手所能获得的最佳效用值

极小化极大值:
递归计算



终止状态值:
效用值

评价函数

- 在深度受限搜索中评价函数给出非终止状态的期望效用的估计值

$$UTILITY(loss, p) \leqslant EVAL(s, p) \leqslant UTILITY(win, p)$$

- 理想的函数：返回位置的实际的极小化极大值
- 在实际中：特征的线性加权，例如：

$$f_1(s) = (\text{白棋皇后数} - \text{黑棋皇后数})$$

$$EVAL(s) = w_1 f_1(s) + w_2 f_2(s) + \cdots + w_n f_n(s) = \sum_{i=1}^n w_i f_i(s)$$

- **一个好的评价函数：**

首先，计算时间不能太长！（加快搜索速度）

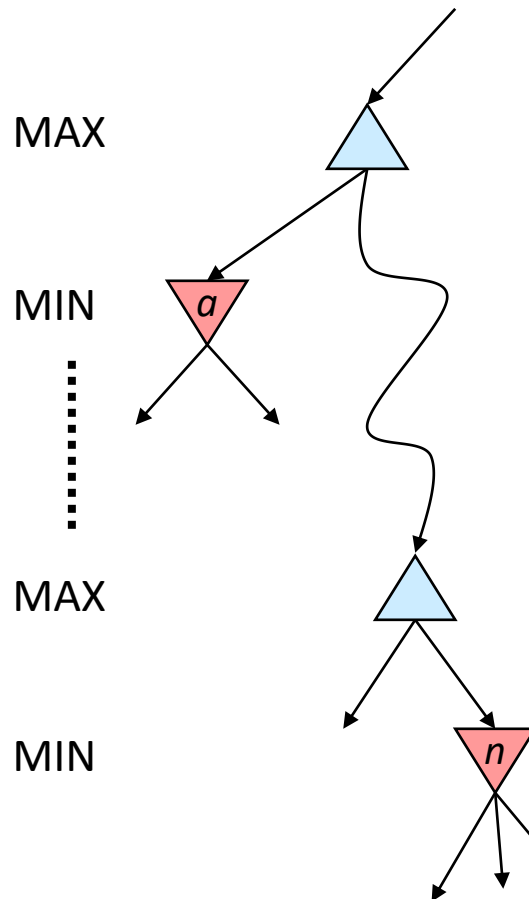
其次，评价函数应与实际的获胜机会密切相关。

博弈中的优化决策

α - β 剪枝 (alpha-beta pruning)

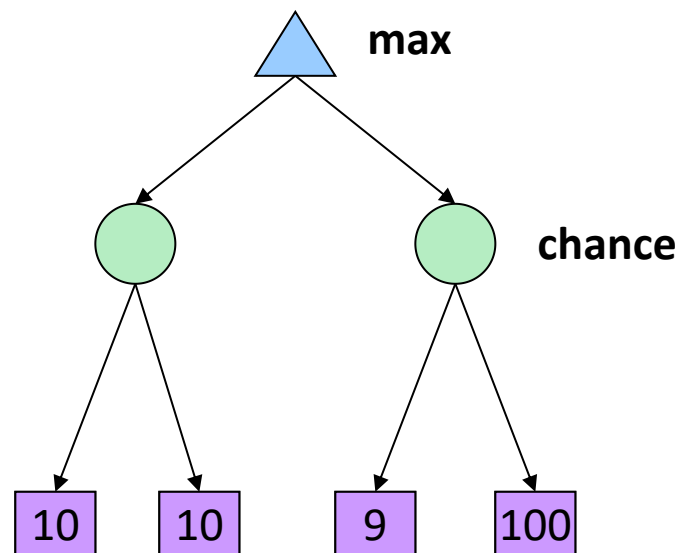
α - β 剪枝的一般情况:

- 计算某一节点n的MIN值;
- 开始遍历n的后继节点;
- n节点的值是随着后继节点的遍历而下降的;
- 令 α 为MAX当前能够获得的最大值;
- 如果n的值变得比a小, 它的值便不会被MAX采用, 所以我们停止考虑n的其他后继节点。



期望最大(expectimax)搜索

- 随机性产生的原因
 - 显式的随机性：游戏中的掷色子。
 - 不可预测的对手：吃豆人中幽灵的随机响应。
 - 动作可能存在一定的失败概率：当机器人移动时，轮胎可能会打滑。
- 这时的状态值应该反应平均情况，即期望最大(expectimax)，而不是最坏情况，即极小化极大(minimax)。
- **期望最大(expectimax)搜索**：计算最佳博弈策略下的平均分
 - Max 节点与极小化极大搜索中的一样。
 - 机会节点类似Min节点，但是结果状态不确定。
 - 计算它们的期望效用值，即，对后继状态值使用对应的概率进行加权平均。



◆ 第七章：逻辑智能体

- 逻辑推断
- 命题逻辑：语法
- 命题逻辑：语义
- 命题定理证明

前述的例子不仅阐明了什么是蕴含，还展示了如何用蕴含的定义来推导出结论，即进行**逻辑推断**。如果将 KB 的所有推论的集合比作干草堆而将 α 比做一根针，那么**蕴含正如草堆中的针一样，而推断就像找到这根针的过程**。如果一个推断算法 i 可以从 KB 中推导出 α ，则记为

$$KB \vdash_i \alpha$$

可靠性：推断算法 i 是可靠的，如果

对于任何 $KB \vdash_i \alpha$ ，那么 $KB \models \alpha$ 也为真

完备性：推断算法 i 是完备的，如果

对于任何 $KB \models \alpha$ ，那么 $KB \vdash_i \alpha$ 也为真

命题逻辑的**语法**定义合法的语句。**原子语句 (atomic sentence)** 由单个**命题符号 (proposition symbol)** 构成。使用括号和被称作**逻辑联结词 (logical connective)** 的运算符可以将简单语句构造成**复合语句 (complex sentence)**。常用的联结词有5个：

- \neg (非)。类似 $\neg W_{1,3}$ 这样的语句称为 $W_{1,3}$ 的**否定**。一个**文字**要么是原子语句，即**正文字**，要么是原子语句的否定，即**负文字**。

- \wedge (与)。主要联结词是 \wedge 的语句称为**合取式**，例如 $W_{1,3} \wedge P_{3,1}$ ，其各部分称为**合取子句**。（ \wedge 看起来像是“And”中的“A”。）

- \vee (或)。主要联结词是 \vee 的语句称为**析取式**，例如 $(W_{1,3} \wedge P_{3,1}) \vee W_{2,2}$ ，其各部分为**析取子句**，本例中分别为 $(W_{1,3} \wedge P_{3,1})$ 和 $W_{2,2}$ 。

- \Rightarrow (蕴涵)。如 $(W_{1,3} \wedge P_{3,1}) \Rightarrow W_{2,2}$ 这样的语句称为**蕴涵式 (implication)** 或条件式，其**前提 (premise)** 或**前件 (antecedent)** 是 $(W_{1,3} \wedge P_{3,1})$ ，其**结论 (conclusion)** 或**后件 (consequent)** 是 $W_{2,2}$ 。蕴涵式也被称为**规则 (rule)** 或**if-then** 声明。有时，蕴涵符号在一些书籍中写作 \supset 或 \rightarrow 。

- \Leftrightarrow (当且仅当)。语句 $W_{1,3} \Leftrightarrow \neg W_{2,2}$ 是**双向蕴涵式 (biconditional)**。

语义定义了用于判定特定模型中语句真值的规则。命题逻辑中，模型就是对每个命题符号设定真值，即真 (true) 或假 (false)。命题逻辑的语义必须指定在给定模型下如何计算任一语句的真值。**这是以递归的方式实现的。**所有语句都是由原子语句和5个联结词构建的。

- $\neg P$ 为真，当且仅当在 m 中 P 为假。
- $P \wedge Q$ 为真，当且仅当在 m 中 P 和 Q 都为真。
- $P \vee Q$ 为真，当且仅当在 m 中 P 或 Q 中至少一个为真。
- $P \Rightarrow Q$ 为真，除非在 m 中 P 为真而 Q 为假。
- $P \Leftrightarrow Q$ 为真，当且仅当在 m 中 P 和 Q 都为真或都为假。

逻辑等价 (logical equivalence)：如果两个语句在相同的模型集合中都为真，则这两个语句逻辑等价，可以写作： $\alpha \equiv \beta$

等价的另一种定义为任意两条语句是等价的，当且仅当它们互相蕴含：

$$\alpha \equiv \beta \text{ 当且仅当 } \alpha \models \beta \text{ 且 } \beta \models \alpha$$

标准的逻辑等价： $(\alpha \wedge \beta) \equiv (\beta \wedge \alpha)$ \wedge 的交换律

$$(\alpha \vee \beta) \equiv (\beta \vee \alpha) \quad \vee \text{ 的交换律}$$

$$((\alpha \wedge \beta) \wedge \gamma) \equiv (\alpha \wedge (\beta \wedge \gamma)) \quad \wedge \text{ 的结合律}$$

$$((\alpha \vee \beta) \vee \gamma) \equiv (\alpha \vee (\beta \vee \gamma)) \quad \vee \text{ 的结合律}$$

$$\neg(\neg\alpha) \equiv \alpha \quad \text{双重否定律}$$

$$(\alpha \Rightarrow \beta) \equiv (\neg\beta \Rightarrow \neg\alpha) \quad \text{假言易位}$$

$$(\alpha \Rightarrow \beta) \equiv (\neg\alpha \vee \beta) \quad \text{蕴涵消去}$$

$$(\alpha \Leftrightarrow \beta) \equiv ((\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)) \quad \text{等价消去}$$

$$\neg(\alpha \wedge \beta) \equiv (\neg\alpha \vee \neg\beta) \quad \text{德摩根律}$$

$$\neg(\alpha \vee \beta) \equiv (\neg\alpha \wedge \neg\beta) \quad \text{德摩根律}$$

$$(\alpha \wedge (\beta \vee \gamma)) \equiv ((\alpha \wedge \beta) \vee (\alpha \wedge \gamma)) \quad \wedge \text{ 对 } \vee \text{ 的分配律}$$

$$(\alpha \vee (\beta \wedge \gamma)) \equiv ((\alpha \vee \beta) \wedge (\alpha \vee \gamma)) \quad \vee \text{ 对 } \wedge \text{ 的分配律}$$

命题定理证明：有效性和可满足性

- **有效性 (validity)** : 如果一条语句在所有模型中都为真, 则这条语句是有效的。

例如 $\text{True}, A \vee \neg A, A \Rightarrow A, (A \wedge (A \Rightarrow B)) \Rightarrow B$

蕴含的定义可以推导出演绎定理:

$KB \models a$ 当且仅当 $(KB \Rightarrow a)$ 是有效的

命题定理证明：有效性和可满足性

- **有效性 (validity)** : 如果一条语句在所有模型中都为真, 则这条语句是有效的。

例如 $\text{True}, A \vee \neg A, A \Rightarrow A, (A \wedge (A \Rightarrow B)) \Rightarrow B$

蕴含的定义可以推导出演绎定理:

$KB \models a$ 当且仅当 $(KB \Rightarrow a)$ 是有效的

- **可满足性 (satisfiability)** : 如果一条语句在某些模型中为真或能够被满足, 则这条语句是可满足的。

如果一条语句没有模型使其为真, 则这条语句是不可满足的。

例如, $A \wedge \neg A$

可满足可以推导出归谬法:

$KB \models a$ 当且仅当 $(KB \wedge \neg a)$ 是不可满足的

- **单调性**: 它表明蕴含的语句集只能随着信息被加入知识库而增长。

如果 $KB \models a$, 则 $KB \wedge \beta \models a$

- **肯定前件** (Modus Ponens, mode that affirms的拉丁语) , 写作

$$\frac{\alpha \Rightarrow \beta, \alpha}{\beta}$$

- **合取消去** (and-elimination) , 即可以从一个合取式推导出任一合取子句:

$$\frac{\alpha \wedge \beta}{\alpha}$$

- 所有逻辑等价都可以用作推断规则。例如, 等价消去可以产生两条推断规则:

$$\frac{\alpha \Leftrightarrow \beta}{(\alpha \Rightarrow \beta)(\beta \Rightarrow \alpha)} \text{ 和 } \frac{(\alpha \Rightarrow \beta)(\beta \Rightarrow \alpha)}{\alpha \Leftrightarrow \beta}$$

命题定理证明：通过归结证明

单元归结 (unit resolution) 规则

$$\frac{\ell_1 \vee \cdots \vee \ell_k, \quad m}{\ell_1 \vee \cdots \vee \ell_{i-1} \vee \ell_{i+1} \vee \cdots \vee \ell_k}$$

单元归结规则可以推广为**全归结规则**

$$\frac{\ell_1 \vee \cdots \vee \ell_k, \quad m_1 \vee \cdots \vee m_n}{\ell_1 \vee \cdots \vee \ell_{i-1} \vee \ell_{i+1} \vee \cdots \vee \ell_k \vee m_1 \vee \cdots \vee m_{j-1} \vee m_{j+1} \vee \cdots \vee m_n}$$

归结使用两个子句并产生一个新的子句，该新子句包含除一对互补文字以外的原始子句的所有文字，例如

$$\frac{P_{1,1} \vee P_{3,1}, \quad \neg P_{1,1} \vee \neg P_{2,2}}{P_{3,1} \vee \neg P_{2,2}}$$

一次只能归结一对互补文字，例如

$$\frac{P \vee \neg Q \vee R, \quad \neg P \vee Q}{\neg Q \vee Q \vee R}$$

命题定理证明：通过归结证明

形式为子句合取式的语句被称为**合取范式 (conjunctive normal form) 或CNF**。把语句转换为CNF的过程：

$$B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1})$$

1. 消去 \Leftrightarrow , 替换 $a \Leftrightarrow b$ 为 $(a \Rightarrow b) \wedge (b \Rightarrow a)$.

$$(B_{1,1} \Rightarrow (P_{1,2} \vee P_{2,1})) \wedge ((P_{1,2} \vee P_{2,1}) \Rightarrow$$

2. 消去 \Rightarrow , 替换 $a \Rightarrow b$ 为 $\neg a \vee b$.

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge (\neg(P_{1,2} \vee P_{2,1}) \vee B_{1,1})$$

3. 将 \neg 内移，通过使用德摩根律和双重否定律：

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge ((\neg P_{1,2} \wedge \neg P_{2,1}) \vee B_{1,1})$$

4. 应用分配律：

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge (\neg P_{1,2} \vee B_{1,1}) \wedge (\neg P_{2,1} \vee B_{1,1})$$

◆ 第十二章：不确定性的量化

- 贝叶斯法则
- 朴素贝叶斯模型

- 乘积法则推可以写成两种形式:

$$P(a \wedge b) = P(a|b)P(b) \quad \text{和} \quad P(a \wedge b) = P(b|a)P(a) .$$

- 联立两式右侧, 除以 $P(a)$, 我们可以得到贝叶斯法则:

$$P(b|a) = \frac{P(a|b)P(b)}{P(a)} .$$

- 通常, 我们把一些未知原因 (*cause*) 的结果 (*effect*) 视为证据, 并想要确定这个原因。在这种情况下, 贝叶斯法则变为了:

$$P(\text{cause}|\text{effect}) = \frac{P(\text{effect}|\text{cause})P(\text{cause})}{P(\text{effect})}$$

- 条件概率 $P(\text{effect}|\text{cause})$ 量化因果方向上的关系, 而 $P(\text{cause}|\text{effect})$ 描述诊断方向上的关系。

- 单个原因直接影响许多结果，给定原因时，所有这些结果都是条件独立的。此时，完全联合分布可以写作：

$$P(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = P(\text{Cause}) \prod_i P(\text{Effect}_i | \text{Cause})$$

- 这样的概率分布叫作朴素贝叶斯 (*naive Bayes*) 模型。“朴素”是因为它经常（作为一种简化假设）用于在给定原因变量时，“结果”变量不是严格独立的情况。
- 考虑观测到的结果 $E=e$, 剩余结果变量 Y 是未观测的：

$$P(\text{Cause} | e) = \alpha \sum_y P(\text{Cause}, e, y)$$

- 根据朴素贝叶斯，我们有

$$\begin{aligned} P(\text{Cause} | e) &= \alpha \sum_y P(\text{Cause}) P(y | \text{Cause}) \left(\prod_j P(e_j | \text{Cause}) \right) \\ &= \alpha P(\text{Cause}) \left(\prod_j P(e_j | \text{Cause}) \right) \sum_y P(y | \text{Cause}) \\ &= \alpha P(\text{Cause}) \prod_j P(e_j | \text{Cause}) \end{aligned}$$

◆ 第十三章：概率推理

- 贝叶斯网络中的精确推断：枚举法
- 贝叶斯网络中的精确推断：变量消元法
- 贝叶斯网络中的近似推理：直接采样
- 贝叶斯网络中的近似推理：拒绝采样
- 贝叶斯网络中的近似推理：似然加权

通过枚举进行推断

- 一般形式:

- 证据变量: $E_1 \dots E_k = e_1 \dots e_k$
 - 查询变量: Q
 - 隐藏变量: $H_1 \dots H_r$
- $\left. \begin{array}{l} X_1, X_2, \dots, X_n \\ \text{所有变量} \end{array} \right\}$

- 目标概率: $P(Q|e_1 \dots e_k)$

- **步骤1:** 在联合概率分布表中选择和证据变量一致的条目。

- **步骤2:** 对隐藏变量H求和消元来获得查询和证据变量的联合概率。

$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} \underbrace{P(Q, h_1 \dots h_r, e_1 \dots e_k)}_{X_1, X_2, \dots, X_n}$$

- **步骤3:** 归一化。

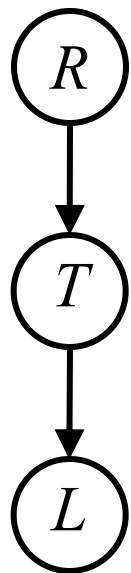
$$P(Q|e_1 \dots e_k) = \frac{1}{Z} P(Q, e_1 \dots e_k) \quad Z = \sum_q P(Q, e_1 \dots e_k)$$

变量消元法

- 查询语句: $P(Q|E_1 = e_1, \dots, E_k = e_k)$
- 从初始因子开始:
 - 局部的 CPTs (会由证据实例化)
- 如果还有隐藏变量未被处理, 则循环:
 - 选择一个隐藏变量 H
 - 连接所有含有隐藏变量 H 的因子
 - 求和消元 H
- 连接所有剩余的因子并归一化

贝叶斯网络中的精确推断

变量消元法



$$P(L) = ?$$

- 枚举法

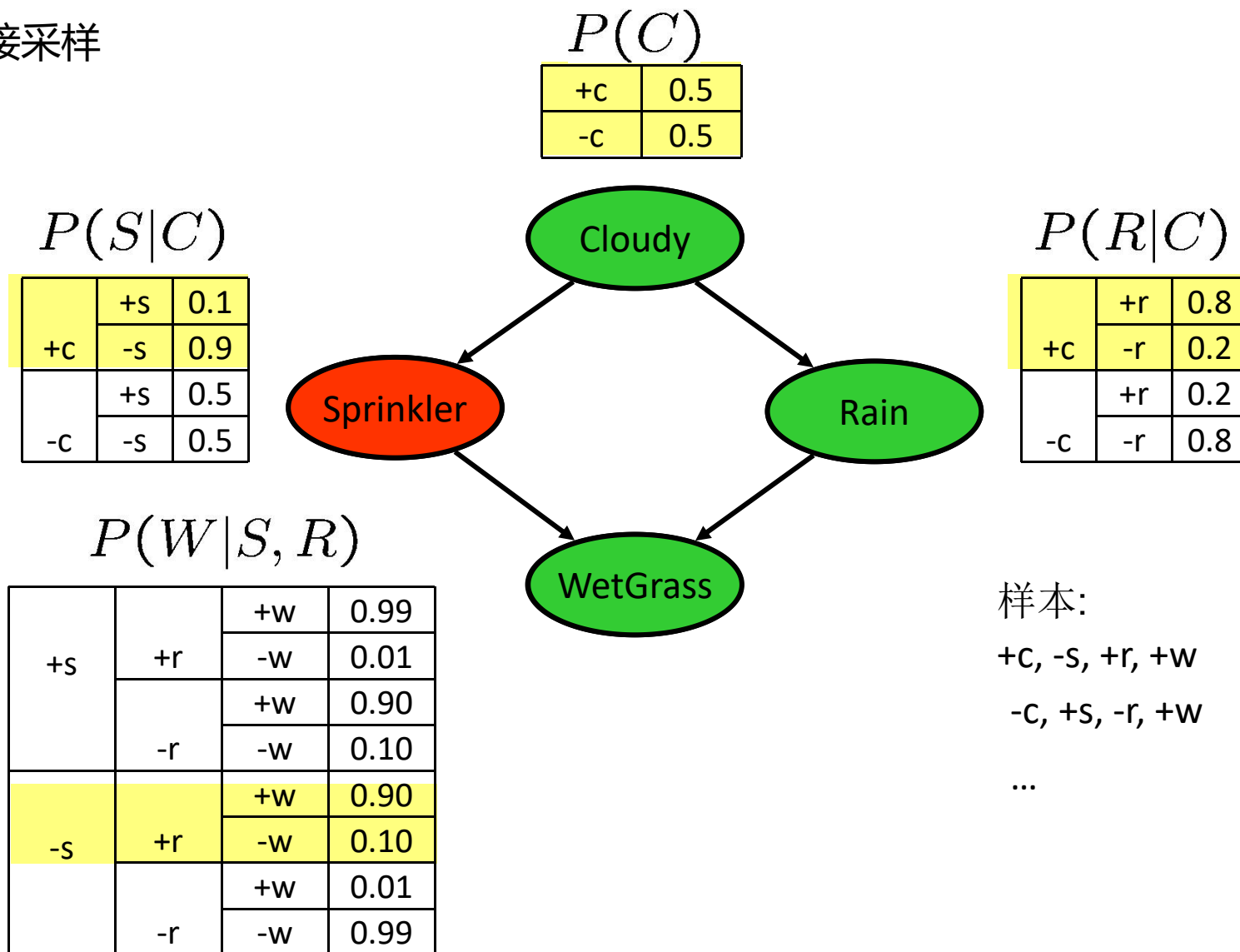
$$= \sum_t \sum_r \underbrace{P(L|t)P(r)P(t|r)}_{\text{连接 } r} \underbrace{\quad}_{\text{连接 } t} \underbrace{\quad}_{\text{消元 } r} \underbrace{\quad}_{\text{消元 } t}$$

- 变量消元法

$$= \sum_t P(L|t) \underbrace{\sum_r P(r)P(t|r)}_{\text{连接 } r} \underbrace{\quad}_{\text{消元 } r} \underbrace{\quad}_{\text{连接 } t} \underbrace{\quad}_{\text{消元 } t}$$

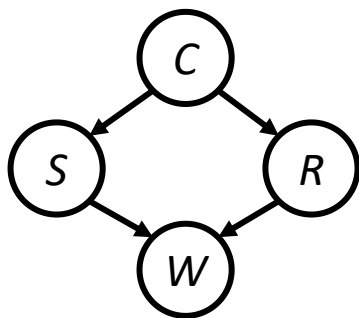
贝叶斯网络中的近似推断

直接采样



拒绝采样

- 假设我们想要计算 $P(C \mid +s)$
 - 使用直接采样的方法生成样本
 - 忽略 (拒绝) 所有与证据变量不一致的样本, 在这例子中 $S = +s$
 - 计算剩余样本中 $C = +c$ 或者 $-c$ 的个数
 - 与目标条件概率是一致的(即, 当样本足够多, 会逼近真实概率值)



+C, -S, +r, +W

+C, +S, +r, +W

-C, +S, +r, -W

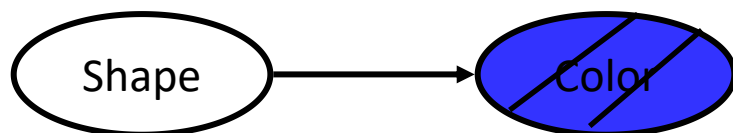
+C, -S, +r, +W

-C, -S, -r, +W

贝叶斯网络中的近似推断

似然加权

- 思路：固定证据变量的取值，对其他变量进行采样
- 问题：样本分布与真实分布不一致
- 解决方案：给定父节点，使用证据变量的概率进行加权



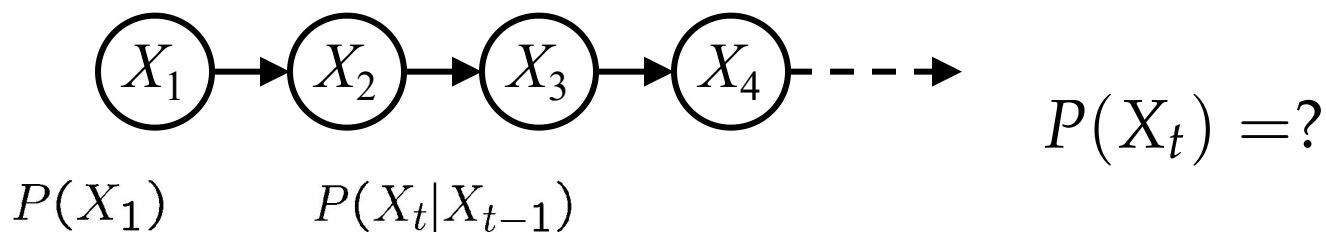
pyramid, blue
pyramid, blue
sphere, blue
cube, blue
sphere, blue

◆ 第十四章：时间上的概率推理

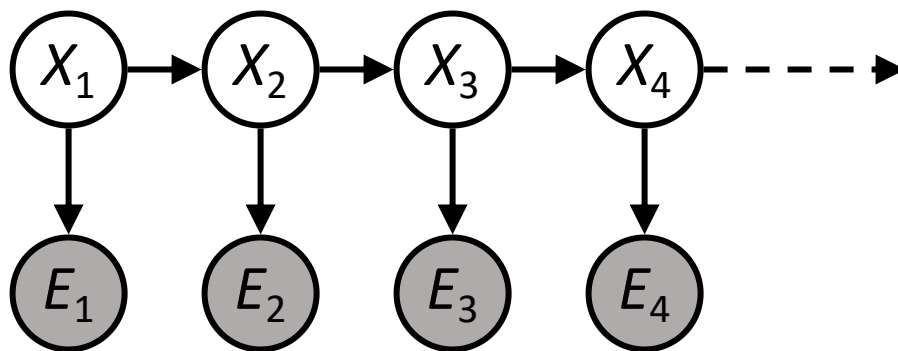
- 马尔可夫模型
- 隐马尔可夫模型

马尔可夫模型

- \mathbf{X}_t : 表示在时刻 t 的状态变量
- 初始状态概率分布: 随机变量 \mathbf{X} 的先验概率分布
- 转移模型: 给定前一状态的值时, 最新状态变量的概率分布: $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-1})$
- 稳态假设: 转移概率在任何时刻都相同。



- 隐马尔可夫模型的定义:
 - 初始概率: $P(X_1)$
 - 转移概率: $P(X_t | X_{t-1})$
 - 输出概率: $P(E_t | X_t)$



◆ 第十九章：样例学习

- 决策树学习
- 线性回归与分类
- 非参数模型：最近邻模型
- 集成学习：自适应提升法

- 决策树 (*decision tree*) 将属性值向量映射到单个输出值 (即 “决策”)。
 - 执行一系列测试来实现其决策, 它从根节点出发, 沿着适当的分支, 直到到达叶节点为止。
 - 树中的每个内部节点对应于一个输入属性的测试
 - 该节点的分支用该属性的所有可能值进行标记
 - 叶节点指定了函数要返回的值
- 布尔型的决策树等价于如下形式的逻辑语句:

$$Output \Leftrightarrow (Path_1 \vee Path_2 \vee \dots)$$

线性回归与分类

单变量线性回归

输入 x 和输出 y

$$y = w_1x + w_0$$

线性函数

$$h_w(x) = w_1x + w_0$$

线性回归：找到最匹配这些数据的线性函数 h_w

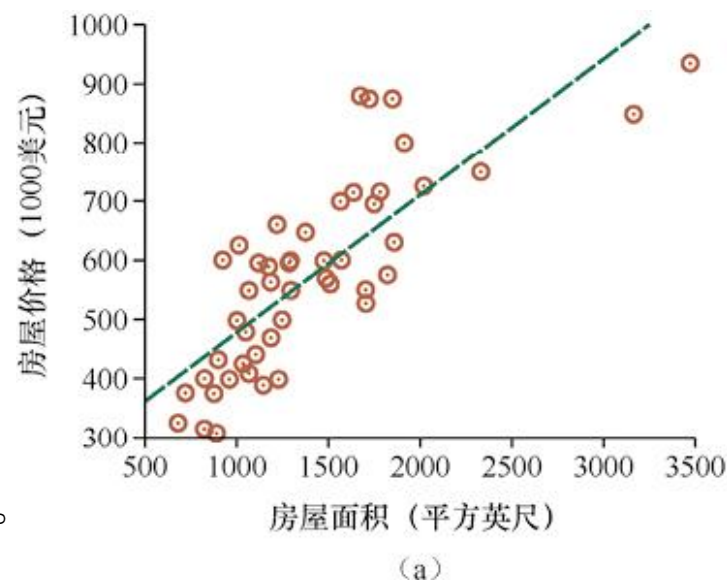
找到对应的权重值 (w_0, w_1) 使得其经验损失最小。

平方误差损失函数, L_2 , 对所有训练样本求和：

$$Loss(h_w) = \sum_{j=1}^N L_2(y_j, h_w(x_j)) = \sum_{j=1}^N (y_j - h_w(x_j))^2 = \sum_{j=1}^N (y_j - (w_1x_j + w_0))^2.$$

$$\frac{\partial}{\partial w_0} \sum_{j=1}^N (y_j - (w_1x_j + w_0))^2 = 0 \text{ and } \frac{\partial}{\partial w_1} \sum_{j=1}^N (y_j - (w_1x_j + w_0))^2 = 0.$$

$$w_1 = \frac{N(\sum x_j y_j) - (\sum x_j)(\sum y_j)}{N(\sum x_j^2) - (\sum x_j)^2}; \quad w_0 = (\sum y_j - w_1(\sum x_j))/N.$$



最近邻模型

k -近邻

- 给定待查询的 \mathbf{x}_q , 寻找最接近 \mathbf{x}_q 的 k 个样例。
- 为了实现分类, 我们寻找 \mathbf{x}_q 的一组邻居, 并以占比最大的输出值为分类结果。例如 $k = 3$ 并且输出值为 *Yes, No, Yes*, 则分类结果为 *Yes*。
- 为了实现回归, 我们可以取 k 个邻居的平均值或中位数, 也可以在最近邻邻居上求解一个线性回归问题。
- 闵可夫斯基距离 (Minkowski distance) 或 L^p 范数

$$L^p(\mathbf{x}_j, \mathbf{x}_q) = \left(\sum_i |x_{j,i} - x_{q,i}|^p \right)^{1/p}.$$

- $p=2$, 欧几里得距离
- $p=1$, 曼哈顿距离
- 对于布尔属性值, 汉明距离
- 马氏距离: 考虑维度之间的协方差

$$D_M(x, y) = \sqrt{(x - y)^T \Sigma^{-1} (x - y)}$$

自适应提升法

- **加权训练集：**给每个样例赋予一个权重 $w_j \geq 0$ ，该权重描述了样例在训练过程中应计数的次数。
- 从所有样例具有相等的权重开始，根据该训练集，训练第一个假设 h_1 .
- 我们希望下一个假设能在被分类错误的样例中表现得更好，因此我们将增加它们的权重，同时减小正确分类的样例的权重。
- 基于这个重新进行加权得到的训练集，我们训练得到假设 h_2 。这一过程将以这种方式不断进行，直到生成 K 个假设。
- 类似于贪心算法，即不会回退，一旦算法选择了某个假设 h_i ，它就永远不会抛弃该选择，而是会添加新的假设：

$$h(\mathbf{x}) = \sum_{i=1}^K z_i h_i(\mathbf{x})$$

◆ 第二十章：概率模型学习

- 贝叶斯学习
- 最大后验学习
- 最大似然学习
- 最大似然参数学习：离散模型
- 贝叶斯参数学习

贝叶斯学习

- 将学习看作假设空间中概率分布的贝叶斯更新:

H 是假设变量, 值为 h_1, h_2, \dots , 先验分布 $P(H)$

第 j 个观测 d_j 给出了随机变量的输出 D_j

训练数据 $d = d_1, \dots, d_n$

- 给定到目前为止的数据, 每一个假设有一个后验分布:

$$P(h_i|d) = \frac{1}{n} P(d|h_i)P(h_i)$$

这里 $P(d|h_i)$ 被称为似然

- 预测为在假设上的概率加权平均:

$$P(X|d) = \sum_i P(X|d, h_i)P(h_i|d) = \sum_i P(X|h_i)P(h_i|d)$$

最大后验学习

- 在假设空间上求和通常是非常困难的 $\sum_i P(X|h_i)P(h_i|d)$
- 最大后验 (MAP)学习: 选择 h_{MAP} 来最大化 $P(h_i|d)$

即, 最大化 $P(d|h_i)P(h_i)$ 或 $\log P(d|h_i) + \log P(h_i)$

$$P(X|d) \approx P(X|h_{\text{MAP}})$$

- $-\log_2 P(d|h_i) - \log_2 P(h_i)$ 负对数项可以被看作
 - 给定假设编码数据所需的比特数 + 编码假设所需的比特数
 - 这是最小描述长度(MDL)学习的基本思想

最大似然学习

- 当数据集很大时，假设的先验分布就不那么重要了，因为来自数据的证据足够强大，足以淹没假设的先验分布。
- 最大似然 (ML) 学习: 选择 h_{ML} 来最大化 $P(d|h_i)$
- 即, 简单的获取对数据的最佳拟合; 对于假设空间具有均匀先验分布, 等同于最大后验学习 (例如所有的假设都同样复杂)
- 最大似然学习是“标准”的（非贝叶斯）统计学习方法

最大似然参数学习：离散模型

- 可能含有樱桃味和酸橙味糖果的糖果袋，其中糖果口味的比例完全未知。
- 参数 θ 表示樱桃味糖果所占的比例，其对应的假设为 h_θ
- 如果我们假设所有的比例有相同的先验可能性，那么采用最大似然估计是合理的
- 现在假设我们已经打开了 N 颗糖果，其中有 c 颗为樱桃味，则该特定数据集的似然为：

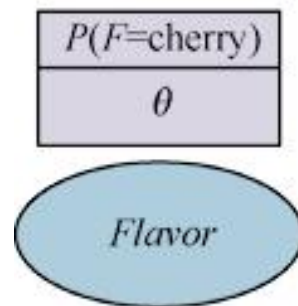
$$P(\mathbf{d} | h_\theta) = \prod_{j=1}^N P(d_j | h_\theta) = \theta^c \cdot (1-\theta)^\ell$$

- 最大似然假设所需的参数即为使得上式最大化的参数。由于 \log 函数是单调函数，我们可以最大化对数似然来简化计算：

$$L(\mathbf{d} | h_\theta) = \log P(\mathbf{d} | h_\theta) = \sum_{j=1}^N \log P(d_j | h_\theta) = c \log \theta + \ell \log (1-\theta)$$

- 对上式关于 θ 进行求导，并令导数为0可得：

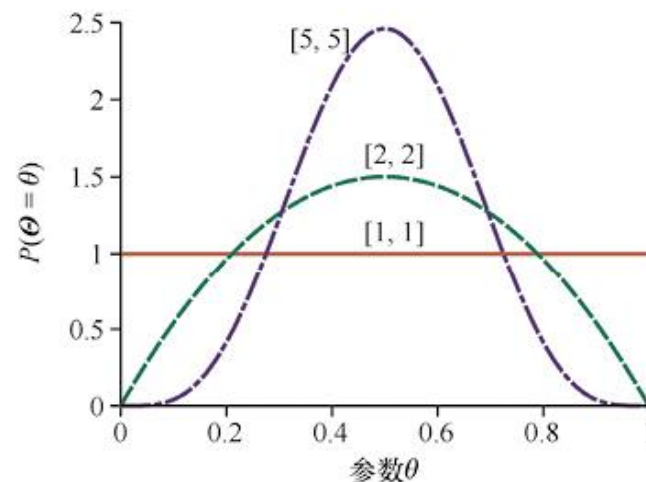
$$\frac{dL(\mathbf{d} | h_\theta)}{d\theta} = \frac{c}{\theta} - \frac{\ell}{1-\theta} = 0 \quad \Rightarrow \quad \theta = \frac{c}{c+\ell} = \frac{c}{N}$$



贝叶斯参数学习

- 基于贝叶斯方法的参数学习过程从一个关于假设的先验分布开始，随着新数据出现而不断更新该分布。
- 从贝叶斯角度看，随机变量 Θ 定义了假设空间， θ 是 Θ 的一个未知值。
- 假设先验是先验分布 $P(\Theta)$ 。因此, $P(\Theta = \theta)$ 是糖果袋中含有比例 θ 的樱桃味糖果的先验概率。
- $P(\theta) = Uniform[0,1](\theta)$, 均匀分布是 β 分布的一个特例。
- β 分布由两个超参数 a 和 b 定义：

$$\text{beta}[a, b](\theta) = \alpha \theta^{a-1} (1 - \theta)^{b-1},$$



贝叶斯参数学习

- 假设我们观测到了一颗樱桃味的糖果，那么我们有

$$\begin{aligned} P(\theta | D_1 = \text{cherry}) &= \alpha P(D_1 = \text{cherry} | \theta) P(\theta) \\ &= \alpha' \theta \cdot \text{beta}[a, b](\theta) = \alpha' \theta \cdot \theta^{a-1} (1 - \theta)^{b-1} \\ &= \alpha' \theta^a (1 - \theta)^{b-1} = \text{beta}[a + 1, b](\theta) . \end{aligned}$$

