

逻辑与推理2

主讲：郭春乐、刘夏雷
南开大学计算机学院

致谢：本课件主要内容来自浙江大学吴飞教授、
南开大学程明明教授

1. 和均是原子命题, “如果那么” 是由和组合得到的复合命题。下面对 “如果那么” 这一复合命题描述不正确的是()

- ☐ A “如果那么” 定义的是一种蕴涵关系(即充分条件)
- ☐ B “如果那么” 意味着命题包含着命题,即是子集
- ☒ C 无法用真值表来判断 “如果那么” 的真假
- ☐ D 当不成立时, “如果那么” 恒为真

提交

2. 下面哪个逻辑等价关系是不成立的 ()

A $\forall x \neg P(x) \equiv \neg \exists x P(x)$

B $\neg \forall x P(x) \equiv \exists x \neg P(x)$

C $\forall x P(x) \equiv \neg \exists x \neg P(x)$

D $\exists x P(x) \equiv \neg \forall x P(x)$

提交

3.下列句子中不是命题的是
()。

- ☒ A $x > 2$ 。
- ☒ B 今天天气真好啊！
- ☐ C 天津大学是中国近代第一所大学。
- ☐ D 所有实数的平方都大于或等于0。

提交

4.应用归结法证明以下命题集是不可满足的。

a) $\alpha \vee \beta$;

b) $\beta \rightarrow \gamma$;

c) $\neg \alpha \wedge \neg \gamma$;

解 应用归结法:

(1) $\alpha \vee \beta$ (已知);

(2) $\neg \beta \vee \gamma$ (b 进行蕴涵消除);

(3) $\alpha \vee \gamma$ (由 1 和 2);

(4) $\neg(\alpha \vee \gamma)$ (c 使用 De Morgan 定律)

(5) 3 和 4 矛盾, 因此原命题集是不可满足的。

5.证明苏格拉底三段论“所有人都是要死的，苏格拉底是人，所以苏格拉底是要死的”。

解：设 $F(x)$: x 是人； $G(x)$: x 是要死的人； a 是苏格拉底。

前提： $(\forall x)(F(x) \rightarrow G(x))$ ； $F(a)$ 。

结论： $G(a)$ 。

证明：(1) $(\forall x)(F(x) \rightarrow G(x))$ ；

(2) $F(a) \rightarrow G(a)$ (1的全称量词消去)；

(3) $F(a)$ ；

(4) $G(a)$ (2和3的假言推理)。

正常使用主观题需2.0以上版本雨课堂

作答

勘误

- P007

A是B的充分条件意味着A蕴含B，即B是A的子集。

- P028

公式（1）和（4）是单项箭头。

$$(1) \quad ((\forall x)A(x) \vee (\forall x)B(x) \Rightarrow (\forall x)(A(x) \vee B(x)))$$

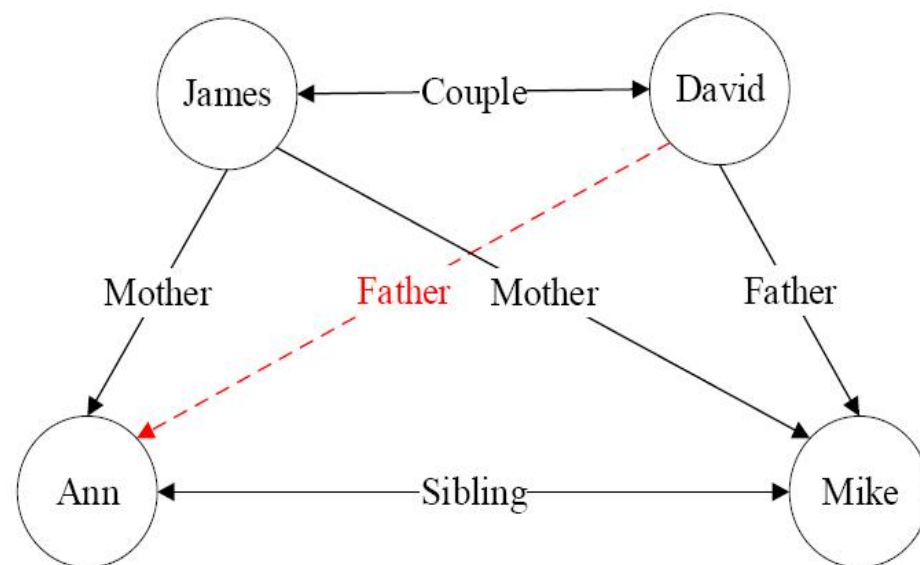
$$(4) \quad (\exists x)(A(x) \wedge B(x)) \Rightarrow (\exists x)A(x) \wedge (\exists x)B(x)$$

提纲

- 命题逻辑
- 谓词逻辑
- 知识图谱推理
- 因果推理

知识图谱：基本概念

- 知识图谱可视为包含多种关系的图。
 - 每个节点是一个实体（如人名、地名、事件和活动等），任意两个节点之间的边表示这两个节点之间存在的关系。
- 可将知识图谱中任意两个相连节点及其连接边表示成一个三元组 (*triplet*)
 - 即 $(left_{node}, relation, right_{node})$
 - 例 $(David, Father, Mike)$



知识图谱：知识工程

• Edward A. Feigenbaum于1977年提出了知识工程的概念

- 开发了基于知识的系统(knowledge-based systems)。该系统利用计算机程序包含大量知识、规则、和推理机制，为显示生活中的问题提供解决方案。
- 基于知识的系统主要表现形式是用于模仿专家决策过程的专家系统。
- Turing Award in 1994

Feigenbaum, E. A., [The Art of Artificial Intelligence: Themes and Case Studies of Knowledge Engineering](#),
Proceedings of the International Joint Conference on Artificial Intelligence(IJCAI), 1977

The art of artificial intelligence—Themes and case studies of knowledge engineering

by EDWARD A. FEIGENBAUM
Stanford University
Stanford, California

INTRODUCTION—AN EXAMPLE

This paper will examine emerging themes of knowledge engineering, illustrate them with case studies drawn from the work of the Stanford Heuristic Programming Project, and discuss general issues of knowledge engineering art and practice.

Let me begin with an example new to our workbench: a system called PUFF, the early fruit of a collaboration between our project and a group at the Pacific Medical Center (PMC) in San Francisco.*

A physician refers a patient to PMC's pulmonary function testing lab for diagnosis of possible pulmonary function disorder. For one of the tests, the patient inhales and exhales a few times in a tube connected to an instrument/computer combination. The instrument acquires data on flow rates and volumes, the so-called flow-volume loop of the patient's lungs and airways. The computer measures certain parameters of the curve and presents them to the diagnostician (physician or PUFF) for interpretation. The diagnosis is made along these lines: normal or diseased; restricted lung disease or obstructive airways disease or a combination of both; the severity; the likely disease type(s) (e.g., emphysema, bronchitis, etc.); and other factors important for diagnosis.

PUFF is given not only the measured data but also certain items of information from the patient record, e.g., sex, age, number of pack-years of cigarette smoking. The task of the PUFF system is to infer a diagnosis and print it out in English in the normal medical summary form of the interpretation expected by the referring physician.

Everything PUFF knows about pulmonary function diagnosis is contained in (currently) 55 rules of the IF... THEN... form. No textbook of medicine currently records these rules. They constitute the partly-public, partly-private knowledge of an expert pulmonary physiologist at PMC, and were extracted and polished by project engineers working intensively with the expert over a period of time. Here is an example of a PUFF rule (the unexplained acronyms refer to various data measurements):

* Dr. J. Osborn, Dr. R. Fallat, John Kunz, Diane McClung.

RULE 31

IF:

- 1) The severity of obstructive airways disease of the patient is greater than or equal to mild, and
- 2) The degree of diffusion defect of the patient is greater than or equal to mild, and
- 3) The tlc (body box) observed/predicted of the patient is greater than or equal to 110 and
- 4) The observed-predicted difference in rv/tlc of the patient is greater than or equal to 10

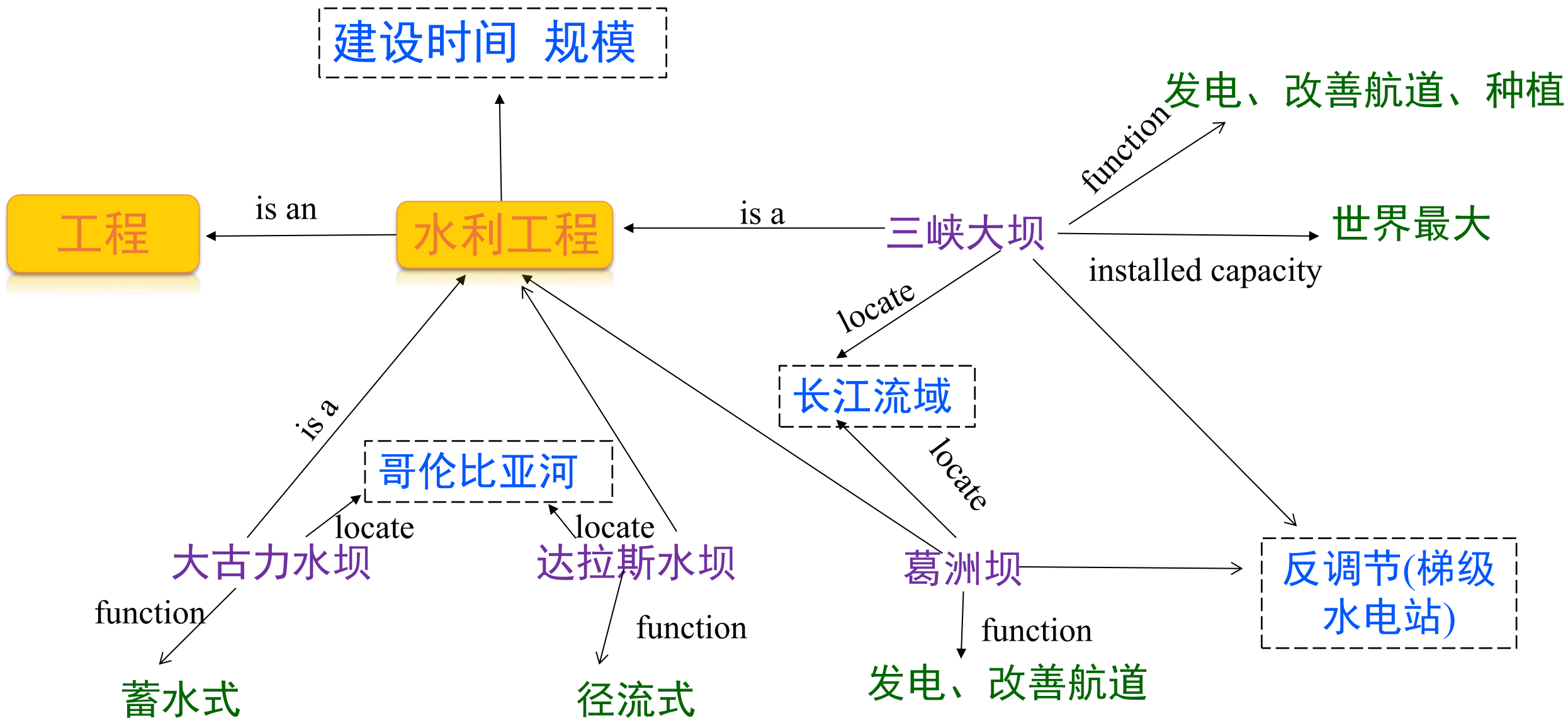
THEN:

- 1) There is strongly suggestive evidence (.9) that the subtype of obstructive airways disease is emphysema, and
- 2) It is definite (1.0) that "OAD, Diffusion Defect, elevated TLC, and elevated RV together indicate emphysema." is one of the findings.

One hundred cases, carefully chosen to span the variety of disease states with sufficient exemplary information for each, were used to extract the 55 rules. As the knowledge emerged, it was represented in rule form, added to the system and tested by running additional cases. The expert was sometimes surprised, sometimes frustrated, by the occasional gaps and inconsistencies in the knowledge, and the incorrect diagnoses that were logical consequences of the existing rule set. The interplay between knowledge engineer and expert gradually expanded the set of rules to remove most of these problems.

As cumulation of techniques in the art demands and allows, a new tool was not invented when an old one would do. The knowledge engineers pulled out of their toolkit a version of the MYCIN system (to be discussed later), with the rules about infectious diseases removed, and used it as the inference engine for the PUFF diagnoses. Thus PUFF, like MYCIN, is a relatively simple backward-chaining infer-

知识图谱：以水利工程为例



知识图谱：以水利工程为例

- **概念之间层次化关系(ontology):**

- 如：工程→水利工程
- 与Wordnet等早期本体知识构建不同，现有方法多在传统分类法(Taxonomy)中结合大众分类(Folksonomy)和机器学习来构建语义网络分类体系。

- **概念对应的例子或实体(instance/entity)**

- 如三峡大坝和葛洲坝等属于水利工程这一概念。
- 一般通过分类识别等手段实现。

知识图谱：以水利工程为例

- **概念或实体的属性：**

- 属性是对概念或实体内涵的描述，如水利工程具有建设时间和规模等属性、三峡大坝具有发电功能等属性。

- **概念或实体之间的关系：**

- 如三峡大坝和葛洲坝之间具有“梯级调节”关系。

- **概念或实体的属性描述和关系表达一般通过三元组来表示：**

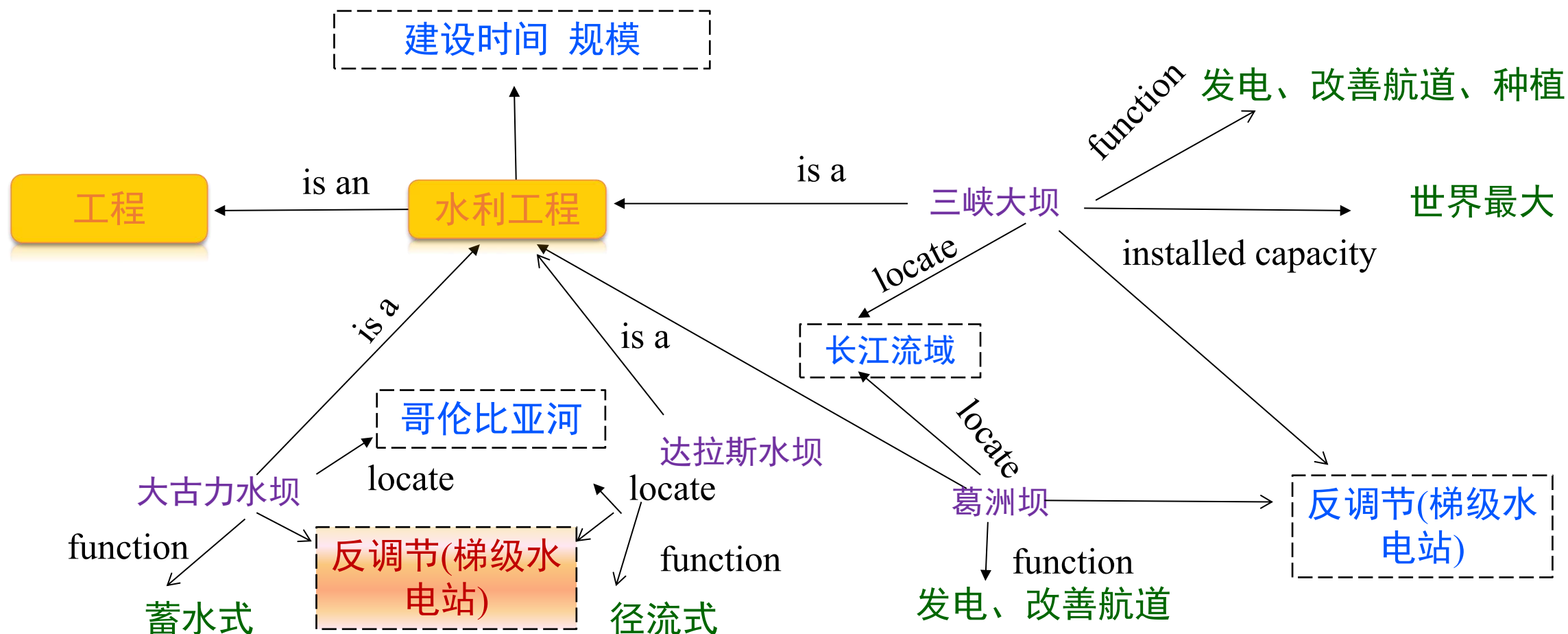
- (entity, relation, entity) 或 (subject, predicate, object)

- **学习概念或实体属性描述及其关联关系是丰富知识图谱的关键！**

知识图谱：以水利工程为例

• 知识图谱推理 (inference) :

- 通过机器学习等方法对知识图谱所蕴含关系进行挖掘。



知识图谱的构成

- 知识图谱一般可通过标注多关系图（labeled multi-relational graph）来表示。
 - 概念： 层次化组织
 - 实体： 概念的示例化描述
 - 属性： 对概念或实体的描述信息
 - 关系： 概念或实体之间的关联
 - 推理规则： 可产生语义网络中上述新的元素

知识图谱推理

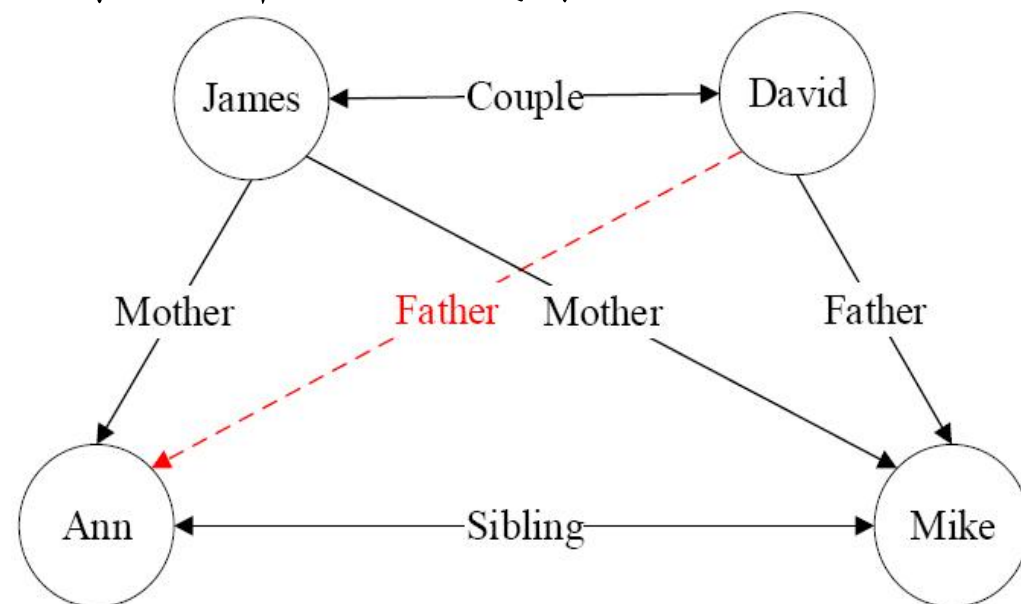
- 知识图谱中存在连线的两个实体可表达为三元组形式
 - $\langle left_node, relation, right_node \rangle$
 - 这种三元组也可以表示为一阶逻辑(first order logic, FOL)的形式
 - 为基于知识图谱的推理创造了条件
 - 例如从 $\langle \text{奥巴马}, \text{出生地}, \text{夏威夷} \rangle$ 和 $\langle \text{夏威夷}, \text{属于}, \text{美国} \rangle$ 两个三元组, 可推理得到 $\langle \text{奥巴马}, \text{国籍}, \text{美国} \rangle$ 。

知识图谱推理

- 可利用一阶谓词来表达刻画知识图谱中节点之间存在的关系
 - $\langle James, Couple, David \rangle$ 关系可用一阶逻辑的形式来描述

- $Couple(x, y)$ 是一阶谓词

- $Couple$ 是实体之间具有的关系
 - x 和 y 是谓词变量



- 图中可推知David和Ann具有父女关系

- 但这一关系在初始图(无红线)中并不存在，需要推理得到

知识图谱推理

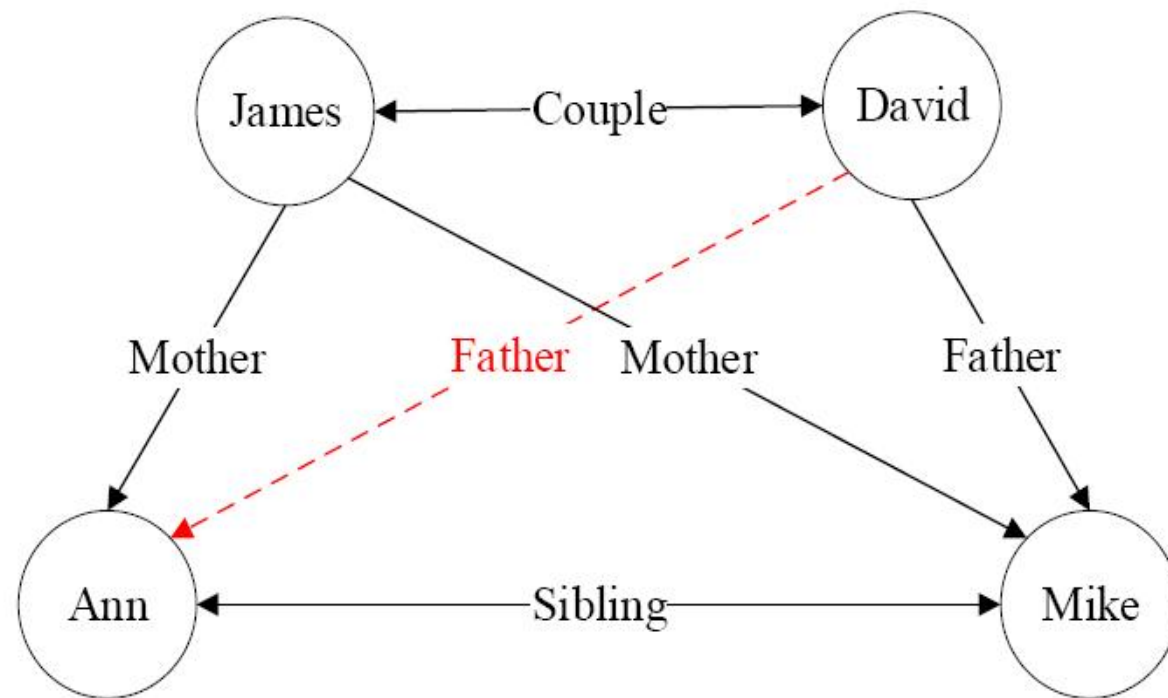
- 问题：如何从知识图谱中推理得到

father(David, Ann)



$(\forall x)(\forall y)(\forall z)(Mother(z, y) \wedge Couple(x, z) \rightarrow Father(x, y))$

如果能够学习得到这条规则，
该有多好？（从具体例子中学
习是归纳推理的范畴）



一个简单的家庭关系知识图谱

知识图谱推理: 归纳学习

- **归纳逻辑程序设计 (inductive logic programming, ILP) 算法**
 - ILP是机器学习和逻辑程序设计交叉领域的研究内容
 - ILP使用一阶谓词逻辑进行知识表示，通过修改和扩充逻辑表达式对现有知识归纳，完成推理任务。
 - 作为ILP的代表性方法，FOIL（First Order Inductive Learner）通过序贯覆盖实现规则推理。

知识图谱推理: FOIL (First Order Inductive Learner)

- 推理手段: positive examples + negative examples + background knowledge examples \Rightarrow hypothesis

$$\frac{(\forall x)(\forall y)(\forall z)(Mother(z, y) \wedge Couple(x, z) \rightarrow Father(x, y))}{\text{前提约束谓词 (学习得到)}} \quad \frac{}{\text{目标谓词 (已知)}}$$

前提约束谓词
(学习得到)

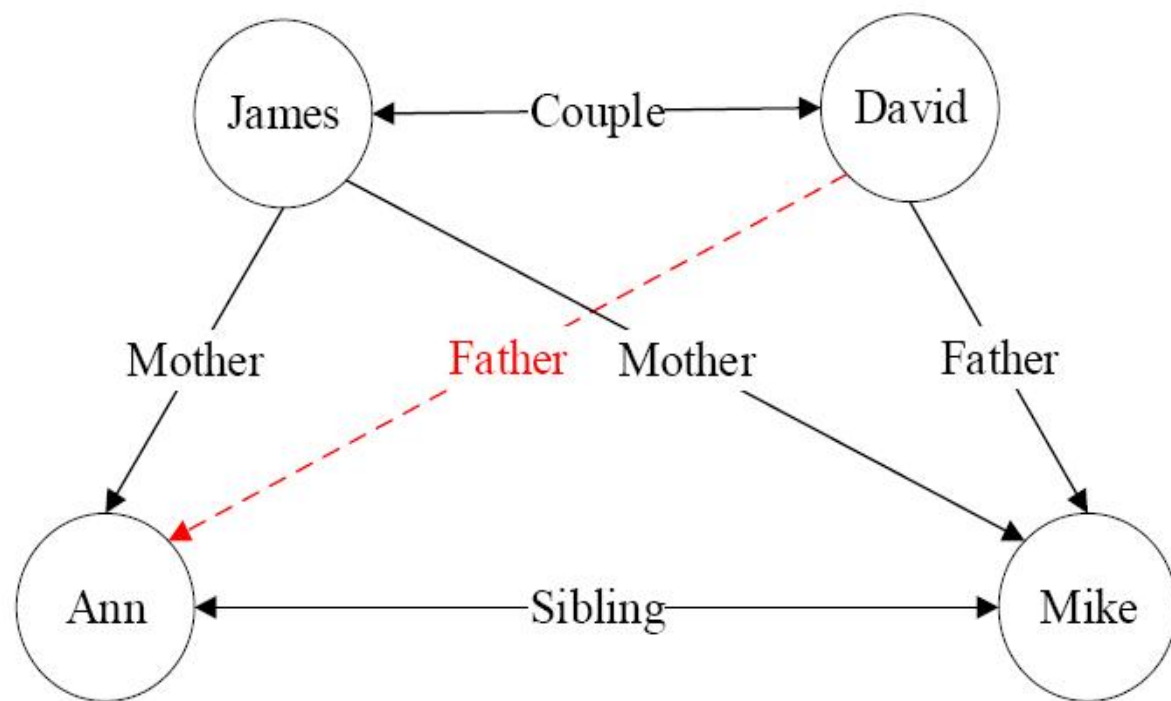
目标谓词
(已知)

知识图谱推理: FOIL (First Order Inductive Learner)

- 目标谓词: $Father(x, y)$
- 目标谓词只有一个正例 $Father(David, Mike)$
 - 反例在知识图谱中一般不会显式给出, 但可从知识图谱中构造出来。
 - 如从知识图谱中已经知道 $Couple(David, James)$ 成立
 - 则 $Father(David, James)$ 可作为目标谓词 $?$ 的一个反例
 - 记为 $\neg Father(David, James)$
- 只能在已知两个实体的关系且确定其关系与目标谓词相悖时, 才能将这两个实体用于构建目标谓词的反例
 - 而不能在不知两个实体是否满足目标谓词前提下将它们来构造目标谓词的反例

知识图谱推理: FOIL (First Order Inductive Learner)

- 目标谓词: $Father(x, y)$
- 背景知识: 知识图谱中目标谓词以外的其他谓词实例化结果
 - 如 $Sibling(Ann, Mike)$



一个简单的家庭关系知识图谱

知识图谱推理: FOIL (First Order Inductive Learner)

$$(\forall x)(\forall y)(\forall z)(Mother(z, y) \wedge Couple(x, z) \rightarrow \textit{Father}(x, y))$$

前提约束谓词 (学习得到)

目标谓词 (已知)

- 推理思路: 从一般到特殊, 逐步给目标谓词添加前提约束谓词, 直到所构成的推理规则不覆盖任何反例。
- 对目标谓词或前提约束谓词中的变量赋予具体值
 - 如将 $(\forall x)(\forall y)(\forall z)(Mother(z, y) \wedge Couple(x, z) \rightarrow \textit{Father}(x, y))$ 这一推理规则所包含的目标谓词 $\textit{Father}(x, y)$ 中 x 和 y 分别赋值为 $David$ 和 Ann

知识图谱推理: FOIL (First Order Inductive Learner)

- 哪些谓词好呢？可以作为目标谓词的前提约束谓词？

- FOIL中信息增益值(information gain)

- FOIL信息增益值计算方法如下：

$$FOIL_Gain = \widehat{m}_+ \cdot \left(\log_2 \frac{\widehat{m}_+}{\widehat{m}_+ + \widehat{m}_-} - \log_2 \frac{m_+}{m_+ + m_-} \right)$$

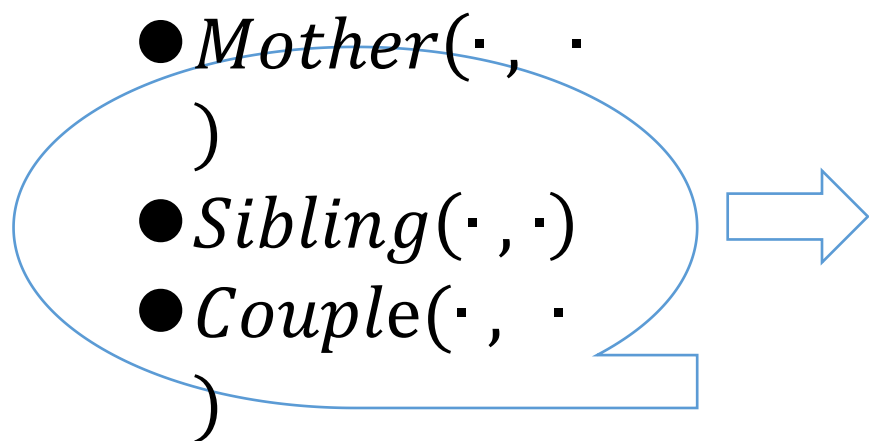
- 其中， \widehat{m}_+ 和 \widehat{m}_- 是增加前提约束谓词后所得新推理规则覆盖的正例和反例的数量， m_+ 和 m_- 是原推理规则所覆盖的正例和反例数量

知识图谱推理: FOIL (First Order Inductive Learner)

$$(\forall x)(\forall y)(\forall z)(Mother(z, y) \wedge Couple(x, z) \rightarrow \textit{Father}(x, y))$$

前提约束谓词 (学习得到)

目标谓词 (已知)



- 依次将谓词加入到推理规则中作为前提约束谓词，并计算所得新推理规则的FOIL增益值。
- 基于计算所得FOIL增益值来选择最佳前提约束谓词。

知识图谱推理: FOIL (First Order Inductive Learner)

$$(\forall x)(\forall y)(\forall z)(Mother(z, y) \wedge Couple(x, z) \rightarrow \textit{Father}(x, y))$$

前提约束谓词 (学习得到)

目标谓词 (已知)

| | | | |
|--------------|---|--------------------|---|
| 背景知识 样例集合 | Sibling(Ann, Mike) Couple(David, James) Mother(James, Ann) Mother(James, Mike) | 目标谓词 训练样例 集合 | Father(David, Mike) \neg Father(David, James) \neg Father(James, Ann) \neg Father(James, Mike) \neg Father(Ann, Mike) |
|--------------|---|--------------------|---|

知识图谱推理: FOIL (First Order Inductive Learner)

- 给定目标谓词，此时推理规则只有目标谓词
- 因此推理规则所覆盖的正例和反例的样本数分别是训练样本中正例和反例的数量
 - $m_+ = 1, m_- = 4$

| 推理规则 | | 推理规则涵盖的正例和反例数 | | FOIL信息增益值 |
|---------------------------|----------------------------------|---------------------------------------|---------------------------------------|--------------|
| 目标谓词 | 前提约束谓词 | 正例 | 反例 | 信息增益值 |
| $Father(x, y) \leftarrow$ | 空集 | $m_+ = 1$ | $m_- = 4$ | $FOIL_Gain$ |
| $Father(x, y) \leftarrow$ | $Mother(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 3$ | 0.32 |
| | $Sibling(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 2$ | 0.74 |
| | $Couple(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Couple(x, z)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 1$ | 1.32 |
| | $Couple(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Couple(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |

知识图谱推理: FOIL (First Order Inductive Learner)

- 将 $Mother(x, y)$ 作为前提约束谓词加入, 可得到推理规则
 - $Mother(x, y) \rightarrow Father(x, y)$
- 背景知识中 $Mother(x, y)$ 有两个实例
 - $Mother(James, Ann)$
 - $Mother(James, Mike)$
- 对于 $Mother(James, Ann)$ 这一实例
 - $x = James, y = Ann$, 将 $\textcircled{?}$ 和 $\textcircled{?}$ 代入 $Father(x, y)$, 可知在训练样本中 $Father(James, Ann)$ 是一个反例

| 推理规则 | | 推理规则涵盖的正例和反例数 | | FOIL信息增益值 |
|---------------------------|----------------------------------|---------------------------------------|---------------------------------------|--------------|
| 目标谓词 | 前提约束谓词 | 正例 | 反例 | 信息增益值 |
| $Father(x, y) \leftarrow$ | 空集 | $m_+ = 1$ | $m_- = 4$ | $FOIL_Gain$ |
| $Father(x, y) \leftarrow$ | $Mother(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 3$ | 0.32 |
| | $Sibling(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 2$ | 0.74 |
| | $Couple(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Couple(x, z)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 1$ | 1.32 |
| | $Couple(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Couple(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |

知识图谱推理: FOIL (First Order Inductive Learner)

- 将 $Mother(x, y)$ 作为前提约束谓词加入, 可得到推理规则
 - $Mother(x, y) \rightarrow Father(x, y)$
- 背景知识中 $Mother(x, y)$ 有两个实例
 - $Mother(James, Ann)$
 - $Mother(James, Mike)$
- 对于 $Mother(James, Mike)$ 这一实例
 - $x = James, y = Mike$, 将 $\textcircled{?}$ 和 $\textcircled{?}$ 代入 $Father(x, y)$, 可知在训练样本中 $Father(James, Mike)$ 是一个反例

| 推理规则 | | 推理规则涵盖的正例和反例数 | | FOIL信息增益值 |
|---------------------------|----------------------------------|---------------------------------------|---------------------------------------|--------------|
| 目标谓词 | 前提约束谓词 | 正例 | 反例 | 信息增益值 |
| $Father(x, y) \leftarrow$ | 空集 | $m_+ = 1$ | $m_- = 4$ | $FOIL_Gain$ |
| $Father(x, y) \leftarrow$ | $Mother(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 3$ | 0.32 |
| | $Sibling(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 2$ | 0.74 |
| | $Couple(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Couple(x, z)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 1$ | 1.32 |
| | $Couple(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Couple(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |

知识图谱推理: FOIL (First Order Inductive Learner)

$Mother(x, y) \rightarrow Father(x, y)$

- 覆盖正例和反例数量分别为 $\widehat{m}_+ = 0, \widehat{m}_- = 2$
- 由于 $\widehat{m}_+ = 0$, 代入 $\frac{\widehat{m}_+ \log \widehat{m}_+ + \widehat{m}_- \log \widehat{m}_-}{\widehat{m}_+ + \widehat{m}_-}$ 公式时会出现负无穷的情况, 此时 $\frac{\widehat{m}_+ \log \widehat{m}_+ + \widehat{m}_- \log \widehat{m}_-}{\widehat{m}_+ + \widehat{m}_-}$ 记为NA (Not Available)

| 推理规则 | | 推理规则涵盖的正例和反例数 | | FOIL信息增益值 |
|---------------------------|----------------------------------|---------------------------------------|---------------------------------------|--------------|
| 目标谓词 | 前提约束谓词 | 正例 | 反例 | 信息增益值 |
| $Father(x, y) \leftarrow$ | 空集 | $m_+ = 1$ | $m_- = 4$ | $FOIL_Gain$ |
| $Father(x, y) \leftarrow$ | $Mother(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 3$ | 0.32 |
| | $Sibling(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 2$ | 0.74 |
| | $Couple(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Couple(x, z)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 1$ | 1.32 |
| | $Couple(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Couple(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |

知识图谱推理: FOIL (First Order Inductive Learner)

- 如果将 $Couple(x, z)$ 作为前提约束谓词加入,

- 可得到如下推理规则

$$Couple(x, z) \rightarrow Father(x, y)$$

- 在背景知识中, $Couple(x, z)$ 只有一个实例

- 即 $x = David, z = James$

- 将其代入 $Father(x, y)$ 得到

$$Father(David, y)$$

| 推理规则 | | 推理规则涵盖的正例和反例数 | | FOIL信息增益值 |
|---------------------------|----------------------------------|---------------------------------------|---------------------------------------|--------------|
| 目标谓词 | 前提约束谓词 | 正例 | 反例 | 信息增益值 |
| $Father(x, y) \leftarrow$ | 空集 | $m_+ = 1$ | $m_- = 4$ | $FOIL_Gain$ |
| $Father(x, y) \leftarrow$ | $Mother(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 3$ | 0.32 |
| | $Sibling(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 2$ | 0.74 |
| | $Couple(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Couple(x, z)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 1$ | 1.32 |
| | $Couple(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Couple(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |

知识图谱推理: FOIL (First Order Inductive Learner)

- 训练样本中存在正例以及反例

- $Father(David, Mike)$
- $\neg Father(David, James)$,
- 即 $Couple(x, z) \rightarrow Father(x, y)$
覆盖正反例数量分别为1和1。
- 信息增益值为:

$$\begin{aligned} & \widehat{m}_+ \cdot \left(\log_2 \frac{\widehat{m}_+}{\widehat{m}_+ + \widehat{m}_-} - \log_2 \frac{m_+}{m_+ + m_-} \right) \\ &= 1 \cdot \left(\log_2 \frac{1}{1 + 1} - \log_2 \frac{1}{1 + 4} \right) = 1.32 \end{aligned}$$

| 推理规则 | | 推理规则涵盖的正例和反例数 | | FOIL信息增益值 |
|---------------------------|----------------------------------|---------------------------------------|---------------------------------------|--------------|
| 目标谓词 | 前提约束谓词 | 正例 | 反例 | 信息增益值 |
| $Father(x, y) \leftarrow$ | 空集 | $m_+ = 1$ | $m_- = 4$ | $FOIL_Gain$ |
| $Father(x, y) \leftarrow$ | $Mother(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Mother(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Mother(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 3$ | 0.32 |
| | $Sibling(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Sibling(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Sibling(z, y)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 2$ | 0.74 |
| | $Couple(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $Couple(x, z)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 1$ | 1.32 |
| | $Couple(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $Couple(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 2$ | NA |
| | $Couple(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |

知识图谱推理: FOIL (First Order Inductive Learner)

| | |
|-------------------------------|--|
| Back-ground knowledge | Sibling(Ann, Mike) Couple(David, James) Mother(James, Ann) Mother(James, Mike) |
| Positive and negative samples | Father(David, Mike) \neg Father(David, James) \neg Father(James, Ann) \neg Father(James, Mike) \neg Father(Ann, Mike) |

- $Couple(x, z)$ 加入后信息增益最大
- 将 $Couple(x, z)$ 加入推理规则, 得到 $Couple(x, z) \rightarrow Father(x, y)$ 新推理规则
- 将训练样例中与该推理规则不符的样例去掉。
这里不符指当 $Couple(x, z)$ 中 x 取值为 David 时, 与 $Father(David,)$ 或 $\neg Father(David,)$ 无法匹配的实例。
- 训练样本集中只有正例 $Father(David, Mike)$ 和负例 $\neg Father(David, James)$ 两个实例

知识图谱推理: FOIL (First Order Inductive Learner)

- 推理 $Mother(z, y)$ 加入信息增益最大

- 将 $Mother(z, y)$ 加入, 得到新规则

$Mother(z, y) \wedge Couple(x, z)$

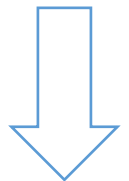
$\rightarrow Father(x, y)$

- 当 $x = David$ 、 $y = Mike$ 、 $z = James$ 时, 该推理规则覆盖训练样本集中正例 $Father(David, Mike)$ 且不覆盖任意反例, 因此算法学习结束。

| 推理规则 | | 推理规则涵盖的正例和反例数 | | FOIL 信息 增益 值 |
|--|--------------------------------|---------------------|---------------------|-----------------------|
| 现有规则 | 拟加入前提 约束谓词 | 正例 | 反例 | |
| $Father(x, y) \leftarrow Couple(x, z)$ | | $m_+ = 1$ | $m_- = 1$ | 1.32 |
| $Father(x, y) \leftarrow Couple(x, z)$ | $\wedge Mother(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Mother(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Mother(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Mother(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Mother(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge \mathbf{Mother(z, y)}$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 0$ | 1 |
| | $\wedge Sibling(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Couple(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $\wedge Couple(x, z)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 1$ | 0 |
| | $\wedge Couple(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Couple(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Couple(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Couple(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |

知识图谱推理: FOIL (First Order Inductive Learner)

$Mother(z, y) \wedge Couple(x, z)$
 $\rightarrow Father(x, y)$



已知:

$Mother(\text{James}, \text{Ann})$
 $Couple(\text{David}, \text{James})$

于是: $Father(\text{David}, \text{Ann})$

| 推理规则 | | 推理规则涵盖的正例和反例数 | | FOIL 信息增益值 |
|--|--------------------------------|------------------------------|------------------------------|------------|
| 现有规则 | 拟加入前提约束谓词 | 正例 | 反例 | |
| $Father(x, y) \leftarrow Couple(x, z)$ | | $m_+ = 1$ | $m_- = 1$ | 1.32 |
| $Father(x, y) \leftarrow Couple(x, z)$ | $\wedge Mother(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Mother(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Mother(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Mother(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Mother(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge \mathbf{Mother(z, y)}$ | $\widehat{m}_+ = \mathbf{1}$ | $\widehat{m}_- = \mathbf{0}$ | 1 |
| | $\wedge Sibling(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(x, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Sibling(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Couple(x, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 1$ | NA |
| | $\wedge Couple(x, z)$ | $\widehat{m}_+ = 1$ | $\widehat{m}_- = 1$ | 0 |
| | $\wedge Couple(y, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Couple(y, z)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Couple(z, x)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |
| | $\wedge Couple(z, y)$ | $\widehat{m}_+ = 0$ | $\widehat{m}_- = 0$ | NA |

知识图谱推理: FOIL (First Order Inductive Learner)

$$(\forall x)(\forall y)(\forall z)(Mother(z, y) \wedge Couple(x, z) \rightarrow \textit{Father}(x, y))$$

前提约束谓词 (学习得到)

目标谓词 (已知)

• 推理手段: *examples(positive + negative + background)* \Rightarrow **hypothesis**

| | | | |
|--------------|----------------------|----------------|-----------------------------|
| 背景知识 样例集合 | Sibling(Ann, Mike) | 目标谓词 训练样例集合 | Father(David, Mike) |
| | Couple(David, James) | | \neg Father(David, James) |
| | Mother(James, Ann) | | \neg Father(James, Ann) |
| | Mother(James, Mike) | | \neg Father(James, Mike) |
| | | | \neg Father(Ann, Mike) |

给定目标谓词，FOIL算法从实例（正例、反例、背景样例）出发，不断测试所得推理规则是否还包含反例，一旦不包含负例，则学习结束，展示了“**归纳学习**”能力。

知识图谱推理: FOIL (First Order Inductive Learner)

FOIL算法

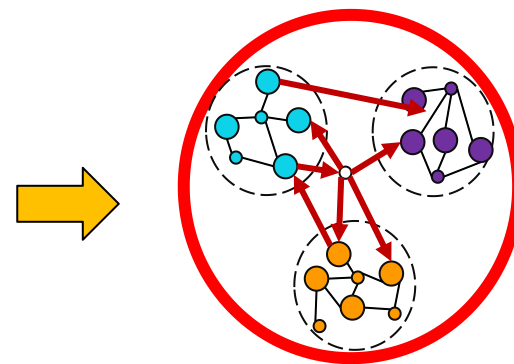
| | |
|-----|--|
| 输入: | 目标谓词 P , P 的训练样例 (正例集合 E^+ 和反例集合 E^-), 其他背景知识 |
| 输出: | 推导得到目标谓词 P 的推理规则 |
| 1 | 将目标谓词作为所学习推理规则的结论 |
| 2 | 将其他谓词逐一作为前提约束谓词加入推理规则, 计算所得到推理规则的FOIL信息增益值, 选取最优前提约束谓词以生成新推理规则, 并将训练样例集合中与该推理规则不符的样例去掉 |
| 3 | 重复2过程, 直到所得到的推理规则不覆盖任意反例 |

知识图谱推理



推理机

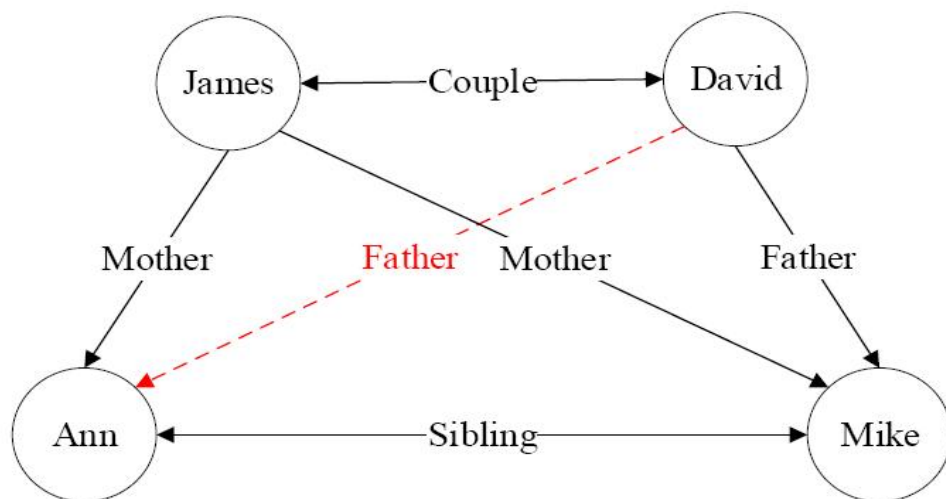
- 概念
- 实体
- 属性
- 关系



知识图谱的扩充

知识图谱推理：路径排序

$\text{Score}(\text{Father}(\text{David}, \text{Ann}))$



一个简单的家庭关系知识图谱

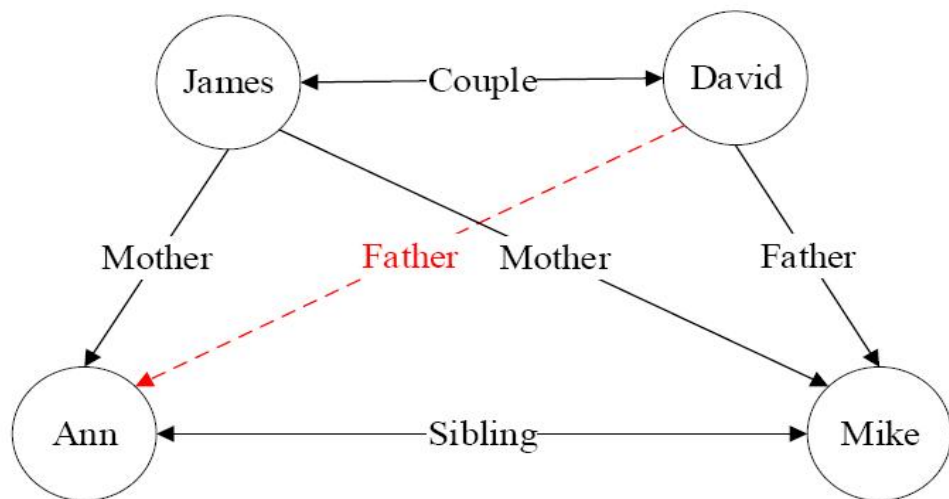
将实体之间的关联路径作为特征，
来学习目标关系的分类器

即：判断David和Ann之间的路径关联是否足够支持表述***Father***这一关系。

Ni Lao, William W. Cohen, Relational retrieval using a combination of path-constrained random walks,
Machine learning, 2010, 81(1): 53-67, 2010

知识图谱推理：路径排序

Score(Father(David, Ann))



一个简单的家庭关系知识图谱

$$\text{score}(s, t) = \sum_{\pi_j \in p_l} \theta_j \underbrace{P(s \rightarrow t; \pi_j)}$$

p_l 是链接节点 s 和节点 t 的所有路径集合

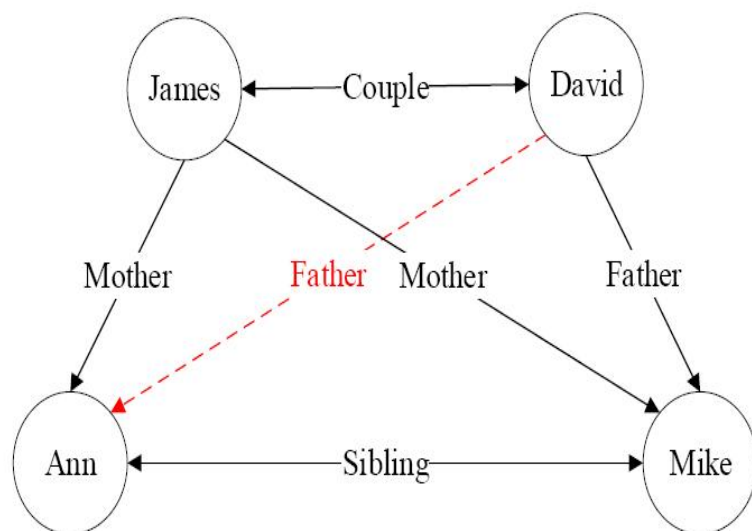
θ_j 是某一条路径 π_j 的权重

P 是路径 π_j 概率值大小

知识图谱推理：路径排序

$$\text{score}(s, t) = \sum_{\pi_j \in p_l} \theta_j P(s \rightarrow t; \pi_j)$$

Score(Father(David, Ann))

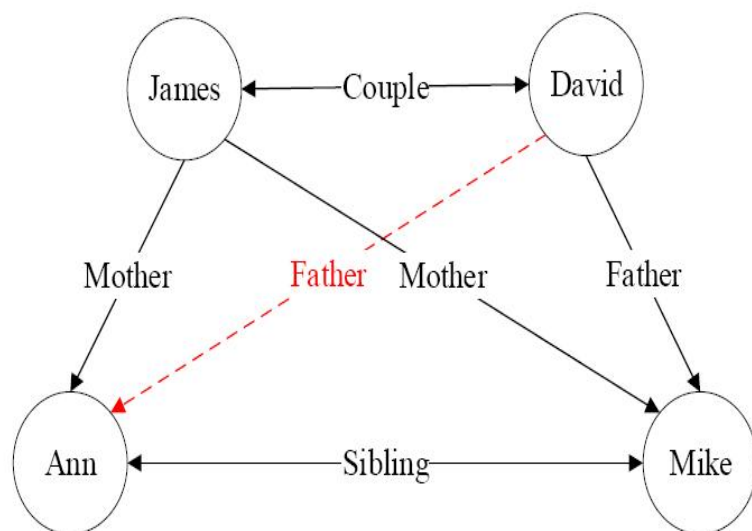


- **特征抽取**：生成并选择路径特征集合。生成路径的方式有随机游走（random walk）、广度优先搜索、深度优先搜索等。
- **特征计算**：计算每个训练样例的特征值 $P(s \rightarrow t; \pi_j)$ 。该特征值可以表示从实体节点 s 出发，通过关系路径 π_j 到达实体节点 t 的概率；也可以表示为布尔值，表示实体 s 到实体 t 之间是否存在路径 π_j ；还可以是实体 s 和实体 t 之间路径出现频次、频率等。
- **分类器训练**：根据训练样例的特征值，为目标关系训练分类器。当训练好分类器后，即可将该分类器用于推理两个实体之间是否存在目标关系。

知识图谱推理：路径排序

$$\text{score}(s, t) = \sum_{\pi_j \in p_l} \theta_j \text{father}(s \rightarrow t; \pi_j)$$

Score(Father(David, Ann))



给定目标关系：*Father(s, t)*

1. 对于目标关系*Father*，生成四组训练样例，一个为正例、三个为负例：

正例：(David, Mike)

负例：(David, James), (James, Ann), (James, Mike)

2. 从知识图谱采样得到路径，每一路径链接上述每个训练样例中两个实体：

(David, Mike)对应路径：*Couple* \rightarrow *Mother*

(David, James)对应路径：*Father* \rightarrow *Mother*⁻¹
(*Mother*⁻¹与*Mother*为相反关系)

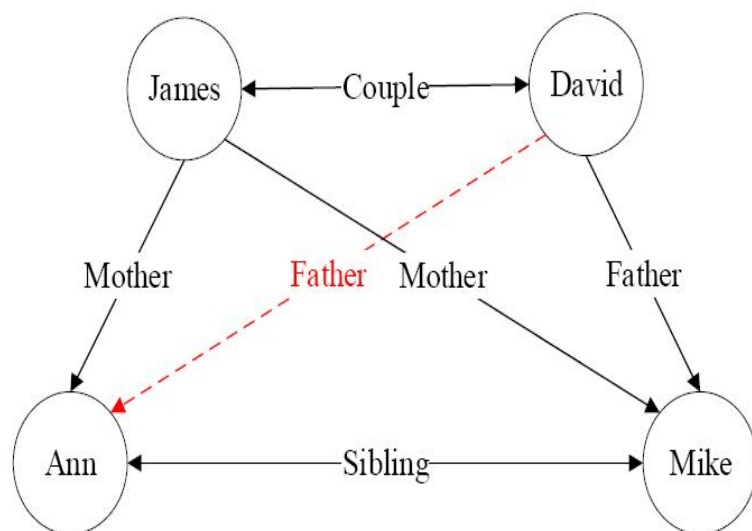
(James, Ann)对应路径：*Mother* \rightarrow *Sibling*

(James, Mike)对应路径：*Couple* \rightarrow *Father*

知识图谱推理： 路径排序

$$score(s, t) = \sum_{\pi_j \in p_l} \theta_j \text{father}(s \rightarrow t; \pi_j)$$

Score(Father(David, Ann))



3. 对于每一个正例/负例，判断上述四条路径可否链接其包含的两个实体，将可链接(记为1)和不可链接(记为0)作为特征，于是每一个正例/负例得到一个四维特征向量：

(David, Mike): $\{[1, 0, 0, 0], 1\}$

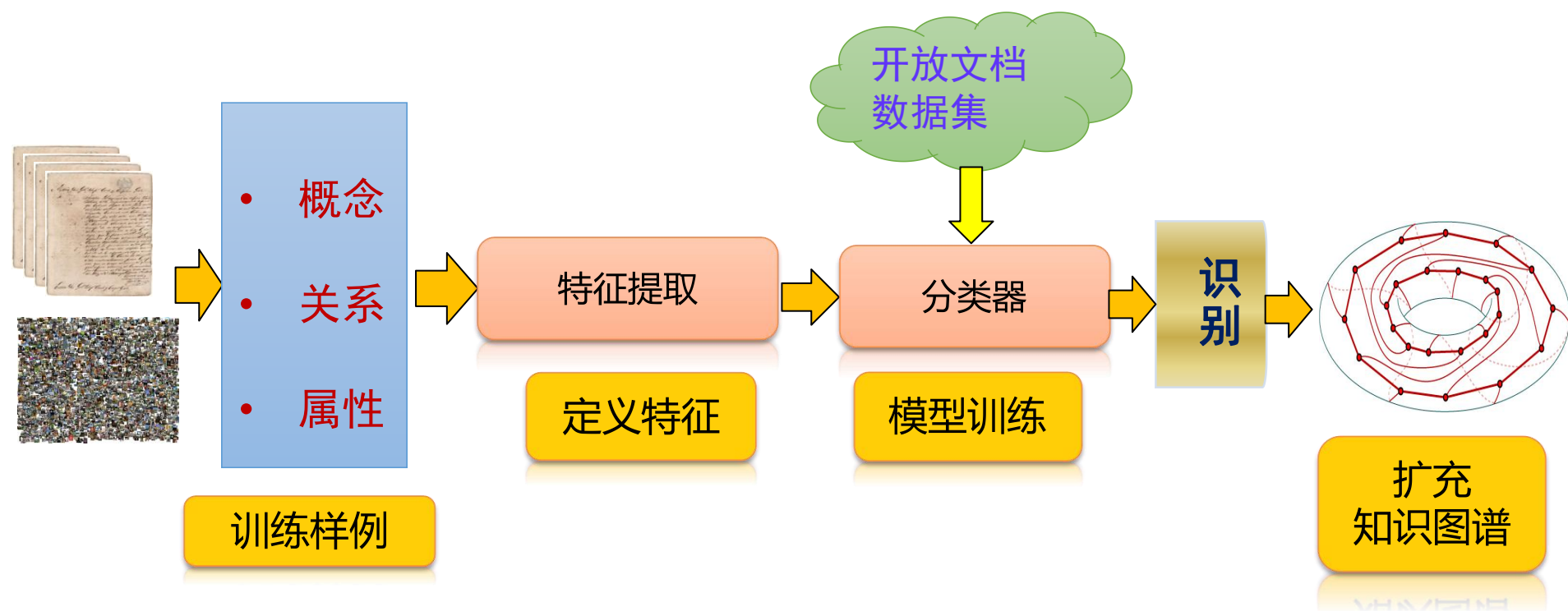
(David, James): $\{[0, 1, 0, 0], -1\}$

(James, Ann): $\{[0, 0, 1, 0], -1\}$

(James, Mike): $\{[0, 0, 1, 1], -1\}$

4. 依据训练样本，训练分类器 M

知识图谱推理：机器学习



知识图谱构造流程：以Wiki为例子

正文描述

William Henry "Bill" Gates III (born October 28, 1955) is an American business magnate, investor, programmer, inventor and philanthropist. [2][3][4] Gates is the former chief executive and chairman of Microsoft, the world's largest personal-computer software company, which he co-founded with Paul Allen.

Gates was born in Seattle, Washington, to William H. Gates, Sr. and Mary Maxwell Gates. His ancestry includes English, German, and Scots-Irish. [15][16] His father was a prominent lawyer, and his mother served on the board of directors for First Interstate BancSystem and the United Way. Gates's maternal grandfather was JW Maxwell, a national bank president. After being named one of *Good Housekeeping's* "50 Most Eligible Bachelors" in 1985, [71] Gates married Melinda French on January 1, 1994. They have three children: daughters Jennifer Katharine (b. 1996) and Phoebe Adele (b. 2002), and son Rory John (b. 1999). The family resides in

Bill Gates



Gates in 2013

| | |
|------------|---|
| Born | <u>William Henry Gates III</u> <u>October 28, 1955 (age 58)</u> Seattle, WA, US |
| Residence | Medina, WA, US |
| Alma mater | Harvard University (dropped out) |
| Children | <u>Jennifer, Rory, and Phoebe</u> |
| Parents | <u>William H. Gates, Sr.</u> <u>Mary Maxwell Gates</u> |
| Signature | <i>William H. Gates III</i> |

属性定义
及描述

44个
类别
标签

Categories: Bill Gates | 1955 births | American billionaires | American chairmen of corporations | American computer businesspeople | American computer programmers | American financiers | American humanitarians | American inventors | American investors | American nonprofit chief executives | American people of English descent | American people of German descent | American people of Scotch-Irish descent | American people of Scottish descent | American philanthropists | American Roman Catholics | American software engineers | American technology chief executives | American technology company founders | American technology writers | Big History | Bill & Melinda Gates Foundation people | Business people from Seattle | Businesspeople in software | Directors of Berkshire Hathaway | Directors of Microsoft | Fellows of the British Computer Society | Gates family | Giving Pledgers | Harvard University people | History of computing | History of Microsoft | Honorary Knights Commander of the Order of the British Empire | Lakeside School alumni | Living people | Members of the United States National Academy of Engineering | Microsoft employees | National Medal of Technology recipients | Personal computing | People from King County, Washington | Placards of the Order of the Aztec Eagle (Mexico) | Windows people | Writers from Seattle, Washington


Wiki中用户对“Bill Gates”这个实例的标注

知识图谱构造流程：以Wiki为例子

William Henry "Bill" Gates III (born October 28, 1955) is an American business magnate, investor, programmer, inventor and philanthropist. [2][3][4] Gates is the former chief executive and chairman of Microsoft, the world's largest personal-computer software company, which he co-founded with Paul Allen.

Gates was born in Seattle, Washington, to William H. Gates, Sr. and Mary Maxwell Gates. His ancestry includes English, German, and Scots-Irish. [15][16] His father was a prominent lawyer, and his mother served on the board of directors for First Interstate BancSystem and the United Way. Gates's maternal grandfather was JW Maxwell, a national bank president. After being named one of *Good Housekeeping's* "50 Most Eligible Bachelors" in 1985, [71] Gates married Melinda French on January 1, 1994. They have three children: daughters Jennifer Katharine (b. 1996) and Phoebe Adele (b. 2002), and son Rory John (b. 1999). The family resides in

Bill Gates



Gates in 2013

| | |
|-------------------|---|
| Born | William Henry Gates III October 28, 1955 (age 58) Seattle, WA, US |
| Residence | Medina, WA, US |
| Alma mater | Harvard University (dropped out) |
| Children | Jennifer, Rory, and Phoebe |
| Parents | William H. Gates, Sr. Mary Maxwell Gates |
| Signature | <i>William H. Gates III</i> |

- 从概念的专业分类 (**taxonomy**) 到大众分类 (**folksonomy**), 即用户趋向于用自我定义的标签对内容进行组织和分类。
- 44个类别标签: Bill Gates, 1955 births, American billionaires, American chairmen of corporations, American computer, business people, American computer programmers, American financiers, American humanitarians, American inventors, American investors, American nonprofit chief executives, American people of English descent, American people of German descent, American people of Scotch-Irish descent, American people of Scottish descent, American philanthropists, American Roman Catholics, American software engineers, American technology chief executives, American technology company founders, American technology writers, Big History, Bill & Melinda Gates Foundation people, Business people from Seattle, Businesspeople in software, Directors of Berkshire Hathaway, Directors of Microsoft, Fellows of the British Computer Society, Gates family, Giving Pledgers, Harvard University people, History of computing, History of Microsoft, Honorary Knights Commander of the Order of the British Empire, Lakeside School alumni, Living people, Members of the United States National Academy of Engineering, Microsoft employees, National Medal of Technology recipients, Personal computing, People from King County, Washington, Placards of the Order of the Aztec Eagle (Mexico), Windows people, Writers from Seattle, Washington


Categories: Bill Gates | 1955 births | American billionaires | American chairmen of corporations | American computer businesspeople | American computer programmers | American financiers | American humanitarians | American inventors | American investors | American nonprofit chief executives | American people of English descent | American people of German descent | American people of Scotch-Irish descent | American people of Scottish descent | American philanthropists | American Roman Catholics | American software engineers | American technology chief executives | American technology company founders | American technology writers | Big History | Bill & Melinda Gates Foundation people | Business people from Seattle | Businesspeople in software | Directors of Berkshire Hathaway | Directors of Microsoft | Fellows of the British Computer Society | Gates family | Giving Pledgers | Harvard University people | History of computing | History of Microsoft | Honorary Knights Commander of the Order of the British Empire | Lakeside School alumni | Living people | Members of the United States National Academy of Engineering | Microsoft employees | National Medal of Technology recipients | Personal computing | People from King County, Washington | Placards of the Order of the Aztec Eagle (Mexico) | Windows people | Writers from Seattle, Washington

知识图谱构造流程：以Wiki为例子

William Henry "Bill" Gates III (born October 28, 1955) is an American business magnate, investor, programmer, inventor and philanthropist. [2][3][4] Gates is the former chief executive and chairman of Microsoft, the world's largest personal-computer software company, which he co-founded with Paul Allen.

Gates was born in Seattle, Washington, to William H. Gates, Sr. and Mary Maxwell Gates. His ancestry includes English, German, and Scots-Irish. [15][16] His father was a prominent lawyer, and his mother served on the board of directors for First Interstate BancSystem and the United Way. Gates's maternal grandfather was JW Maxwell, a national bank president. After being named one of *Good Housekeeping's* "50 Most Eligible Bachelors" in 1985, [71] Gates married Melinda French on January 1, 1994. They have three children: daughters Jennifer Katharine (b. 1996) and Phoebe Adele (b. 2002), and son Rory John (b. 1999). The family resides in

Bill Gates



Gates in 2013

| | |
|-------------------|---|
| Born | William Henry Gates III October 28, 1955 (age 58) Seattle, WA, US |
| Residence | Medina, WA, US |
| Alma mater | Harvard University (dropped out) |
| Children | Jennifer, Rory, and Phoebe |
| Parents | William H. Gates, Sr. Mary Maxwell Gates |
| Signature | <i>William H. Gates III</i> |

- 手工构造的种子知识(Infobox)
- 比尔盖茨这个实体用了13个属性描述、奥巴马这个实体用了20多个属性描述(由于两者均属于persons这个类别，有些属性是共享的)
- Freebase中定义的cities这个概念时使用了将近200个属性，但是仍然远远不够。

Categories: Bill Gates | 1955 births | American billionaires | American chairmen of corporations | American computer businesspeople | American computer programmers | American financiers | American humanitarians | American inventors | American investors | American nonprofit chief executives | American people of English descent | American people of German descent | American people of Scotch-Irish descent | American people of Scottish descent | American philanthropists | American Roman Catholics | American software engineers | American technology chief executives | American technology company founders | American technology writers | Big History | Bill & Melinda Gates Foundation people | Business people from Seattle | Businesspeople in software | Directors of Berkshire Hathaway | Directors of Microsoft | Fellows of the British Computer Society | Gates family | Giving Pledgers | Harvard University people | History of computing | History of Microsoft | Honorary Knights Commander of the Order of the British Empire | Lakeside School alumni | Living people | Members of the United States National Academy of Engineering | Microsoft employees | National Medal of Technology recipients | Personal computing | People from King County, Washington | Placards of the Order of the Aztec Eagle (Mexico) | Windows people | Writers from Seattle, Washington

□如何发现和学习新的属性？


知识图谱构造流程：以Wiki为例子

- 基于机器学习算法进行概念、属性和关系学习，需要大量良好标注数据
- Wikipedia当中，英文文章占总文章数的38.95%，中文文章占4.75%，其他语言文章占56.30%。其中，38.60%的英文文章以Infobox方式被标注 (如实体属性或实体之间的关系)，21.43%的中文文章被标注，平均28.57%的其他语言文章被标注。

William Henry "Bill" Gates III (born October 28, 1955) is an American business magnate, investor, programmer, inventor and philanthropist. [2][3][4] Gates is the former chief executive and chairman of Microsoft, the world's largest personal-computer software company, which he co-founded with Paul Allen.

Gates was born in Seattle, Washington, to William H. Gates, Sr. and Mary Maxwell Gates. His ancestry includes English, German, and Scots-Irish. [15][16] His father was a prominent lawyer, and his mother served on the board of directors for First Interstate BancSystem and the United Way. Gates's maternal grandfather was JW Maxwell, a national bank president. After being named one of *Good Housekeeping's* "50 Most Eligible Bachelors" in 1985, [71] Gates married Melinda French on January 1, 1994. They have three children: daughters Jennifer Katharine (b. 1996) and Phoebe Adele (b. 2002), and son Rory John (b. 1999). The family resides in

Bill Gates



Gates in 2013

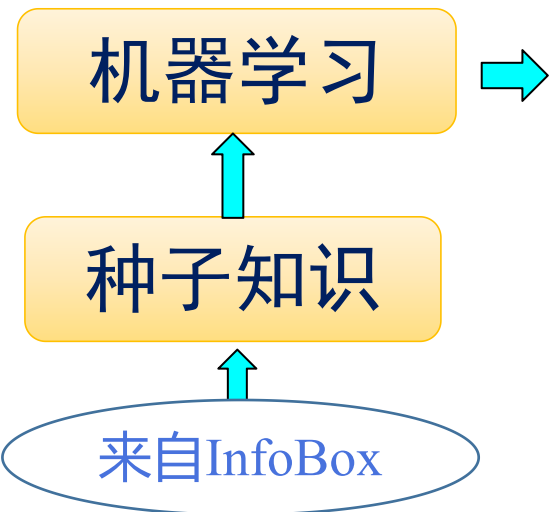
| | |
|-------------------|---|
| Born | William Henry Gates III October 28, 1955 (age 58) Seattle, WA, US |
| Residence | Medina, WA, US |
| Alma mater | Harvard University (dropped out) |
| Children | Jennifer, Rory, and Phoebe |
| Parents | William H. Gates, Sr. Mary Maxwell Gates |
| Signature | <i>William H. Gates III</i> |

用户对wiki中Bill Gates文章所标注的知识

知识图谱构造流程：以Wiki为例子

- 给定无结构化文档数据，通过机器学习方法对实体描述内容进行分类，同时提取描述实体的属性和对应属性值。

嫦娥三号类别: 2013 in China; Space probes launched in 2013; Chinese Lunar Exploration Program; Chinese space probes; Lunar rovers; Missions to the Moon; Spacecraft that orbited the Moon; Soft landings on the Moon; Space probes



| | |
|------|------------------------------|
| 名称 | 嫦娥三号月球探测器 |
| 所属国家 | 中华人民共和国 |
| 构成 | 着陆器、“玉兔号”月球车 |
| 重量 | 3,750千克 |
| 关键技术 | 7500牛变推力发动机、热控两相流体回路、可变热导热管等 |

类别分类、实体识别、属性填充

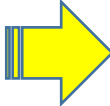
知识图谱构造流程：属性识别与填充

纳尔逊·罗利赫拉赫拉·曼德拉(Nelson Rolihlahla Mandela)
1918年7月18日出生于南非特兰斯凯一个酋长家庭，先后获南非大学文学士和威特沃特斯兰德大学律师资格，当过律师。曼德拉自幼性格刚强，崇敬民族英雄。他是家中长子而被指定为酋长继承人。但他表示：“决不愿以酋长身份统治一个受压迫的部族”，而要“以一个战士的名义投身于民族解放事业”。他毅然走上了追求民族解放的道路。

南非政府2013年12月6日（北京时间）宣布南非前总统曼德拉去世，享年95岁。

曼德拉，南非黑人领袖，因其在废除南非种族歧视政策方面作出了巨大贡献而于1993年荣获诺贝尔和平奖。

1918年7月18日，曼德拉出生于南非特兰斯凯，曼德拉自幼性格刚强，崇敬民族英雄。他是家中长子而被指定为酋长继承人。但他表示：“决不愿以酋长身份统治一个受压迫的部族”，而要“以一个战士的名义投身于民族解放事业”，他毅然……



| 人物 | |
|------|---------------------------|
| 中文名 | 纳尔逊·罗利赫拉赫拉·曼德拉 |
| 外文名 | Nelson Rolihlahla Mandela |
| 国籍 | 南非 |
| 出生地 | 南非特兰斯凯 |
| 出生日期 | 1918年7月18日 |
| 职业 | 政治 南非总统 |
| 主要成就 | 诺贝尔和平奖得主 |

无结构化文档

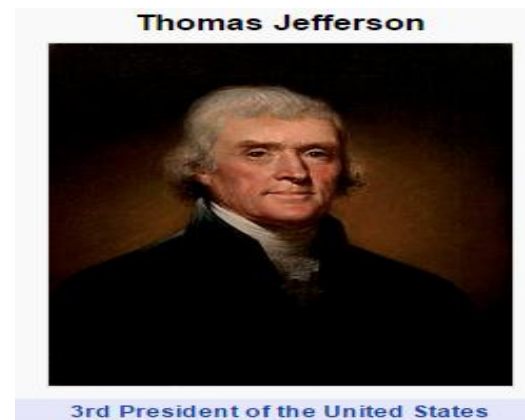
所提取的实体及其属性值形式的知识

知识图谱推理: 机器学习

infoboxes



从数据中学习
人名这个类别
所定义的各个
属性分类器



Born: February 22, 1732, Westmoreland County, Virginia, United States
Died: December 14, 1799, Mount Vernon, Virginia, United States
Spouse: Martha Washington (m. 1759–1799)
Presidential term: April 30, 1789 – March 4, 1797
Siblings: Lawrence Washington (1718–1752), more

Born:
Died:
Party:
Presidential term:
Children:

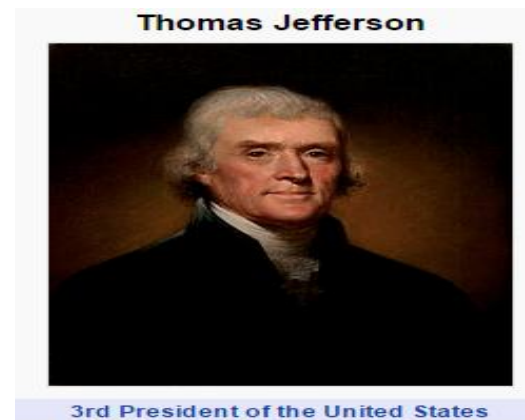


知识图谱推理: 机器学习

infoboxes



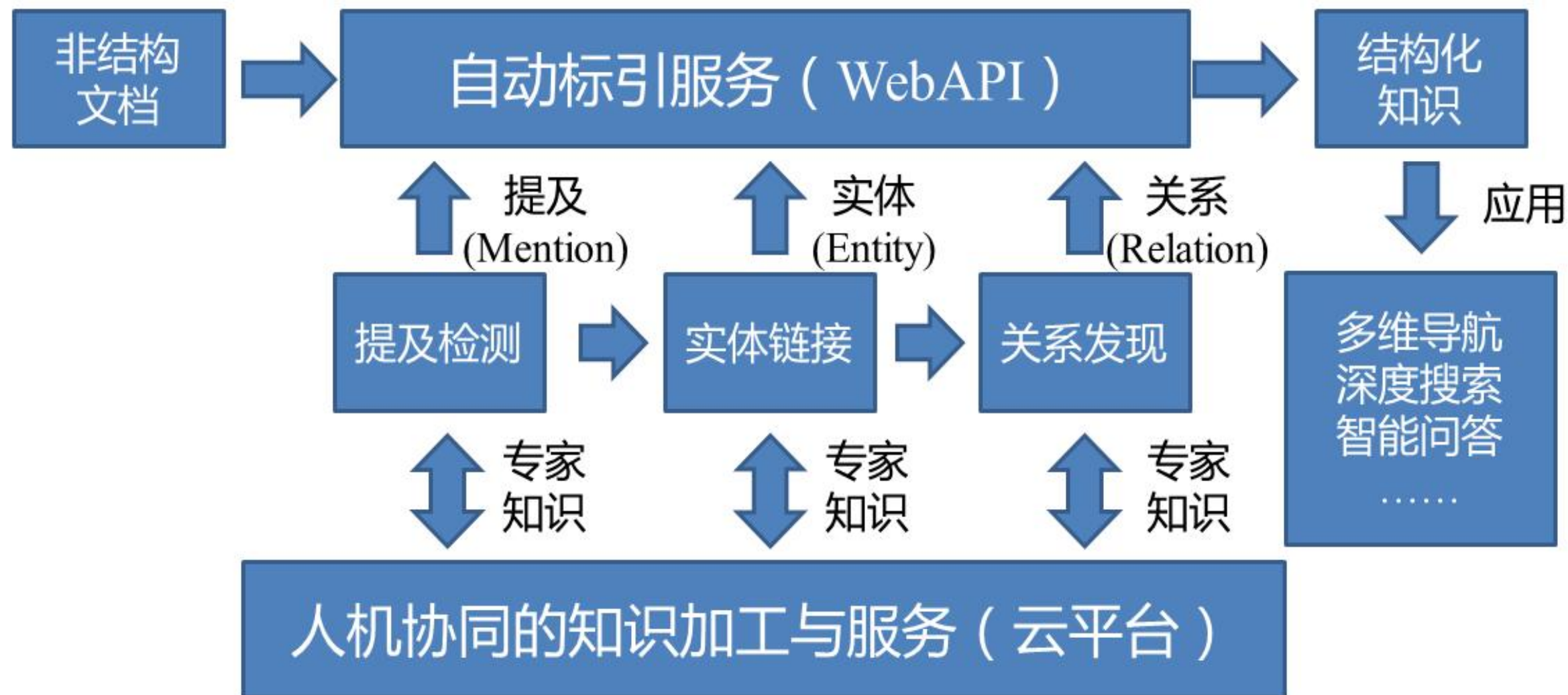
训练得到的
属性分类器



Born: February 22, 1732, Westmoreland County, Virginia, United States
Died: December 14, 1799, Mount Vernon, Virginia, United States
Spouse: Martha Washington (m. 1759–1799)
Presidential term: April 30, 1789 – March 4, 1797
Siblings: Lawrence Washington (1718–1752), more

Born: April 13, 1743, Shadwell, Virginia, United States
Died: July 4, 1826, Monticello, Virginia, United States
Party: Democratic-Republican Party
Presidential term: March 4, 1801 – March 4, 1809
Children: Martha Jefferson Randolph, Eston Hemings,

知识图谱：从数据到知识、从知识到决策



提纲

- 命题逻辑
- 谓词逻辑
- 知识图谱推理
- 因果推理

相关不意味着因果



致酒行（唐-李贺）

零落栖迟一杯酒，主人奉觞客长寿。
主父西游困不归，家人折断门前柳。
吾闻马周昔作新丰客，天荒地老无人识。
空将笺上两行书，直犯龙颜请恩泽。
我有迷魂招不得，**雄鸡一声天下白。**
少年心事当拿云，谁念幽寒坐呜呃。

公鸡打鸣与太阳升起

因果推理：Simpson's Paradox (辛普森悖论)

1973年伯克利本科生录取率

| | 男生 | | 女生 | |
|----|------|-----|------|-----|
| | 申请数 | 录取率 | 申请数 | 录取率 |
| 整体 | 8442 | 44% | 4321 | 35% |

男生录取率(44%)远高于女生(35%)

| 学院 | 男生 | | 女生 | |
|----|-----|------------|-----|-----|
| | 申请数 | 录取率 | 申请数 | 录取率 |
| A | 825 | 62% | 108 | 82% |
| B | 560 | 63% | 25 | 68% |
| C | 325 | 37% | 593 | 34% |
| D | 417 | 33% | 375 | 35% |
| E | 191 | 28% | 393 | 24% |
| F | 373 | 6% | 341 | 7% |

6个最大的院系中，4个院系女生录取率大于男生。如果按照这样的分类，女生实际上比男生的录取率还高一点点。

女生更愿意申请那些竞争压力很大的院系（比如英语系），但是男生却更愿意申请那些相对容易进的院系（比如工程学系）。

P.J. Bickel, E.A. Hammel, J.W. O’Connell, Sex bias in graduate admissions: Data from Berkeley, *Science*, 187(4175):398-404,1975

因果推理：Simpson's Paradox (辛普森悖论)

- 计算机学院学生的高个率高于文学院（左表）。
 - 分别比较两所学院男生和女生身高时，却发现计算机学院男生和女生的高个率均低于文学院
- 在总体样本上成立的某种关系却在分组样本里恰好相反。

| 身高(cm) | 计算机 | 文学院 |
|--------|------|------|
| 矮个人数 | 60 | 80 |
| 高个人数 | 290 | 270 |
| 高个率(%) | 82.9 | 77.1 |

左：计算机学院和文学院学生的身高情况

| | 计算机学院 | | 文学院 | |
|--------|-------|------|------|------|
| 身高(cm) | 男生 | 女生 | 男生 | 女生 |
| 矮个人数 | 35 | 25 | 10 | 70 |
| 高个人数 | 235 | 55 | 80 | 190 |
| 高个率(%) | 87 | 68.9 | 88.9 | 73.1 |

右：以性别分组后的计算机学院和文学院的学生身高情况

$$\frac{b}{a} < \frac{d}{c}, \frac{b'}{a'} < \frac{d'}{c'}$$



$$\frac{b + b'}{a + a'} > \frac{d + d'}{c + c'}$$

因果推理：Simpson's Paradox (辛普森悖论)

- 右表体现了男生比女生个子高这一现象

- 如计算机学院和文学院男生高个率都比女生高个率要大
- 性别会影响专业选择，计算机男生多，而文学院女生多。因此，当计算机学院的样本中包含更多的男生，就会看到左表所呈现的情况：计算机学院的高个率高于文学院。

| 身高(cm) | 计算机 | 文学院 |
|--------|------|------|
| 矮个人数 | 60 | 80 |
| 高个人数 | 290 | 270 |
| 高个率(%) | 82.9 | 77.1 |

左：计算机学院和文学院学生的身高情况

| | 计算机学院 | | 文学院 | |
|--------|-------|------|------|------|
| 身高(cm) | 男生 | 女生 | 男生 | 女生 |
| 矮个人数 | 35 | 25 | 10 | 70 |
| 高个人数 | 235 | 55 | 80 | 190 |
| 高个率(%) | 87 | 68.9 | 88.9 | 73.1 |

右：以性别分组后的计算机学院和文学院的学生身高情况

$$\frac{b}{a} < \frac{d}{c}, \frac{b'}{a'} < \frac{d'}{c'}$$



$$\frac{b + b'}{a + a'} > \frac{d + d'}{c + c'}$$

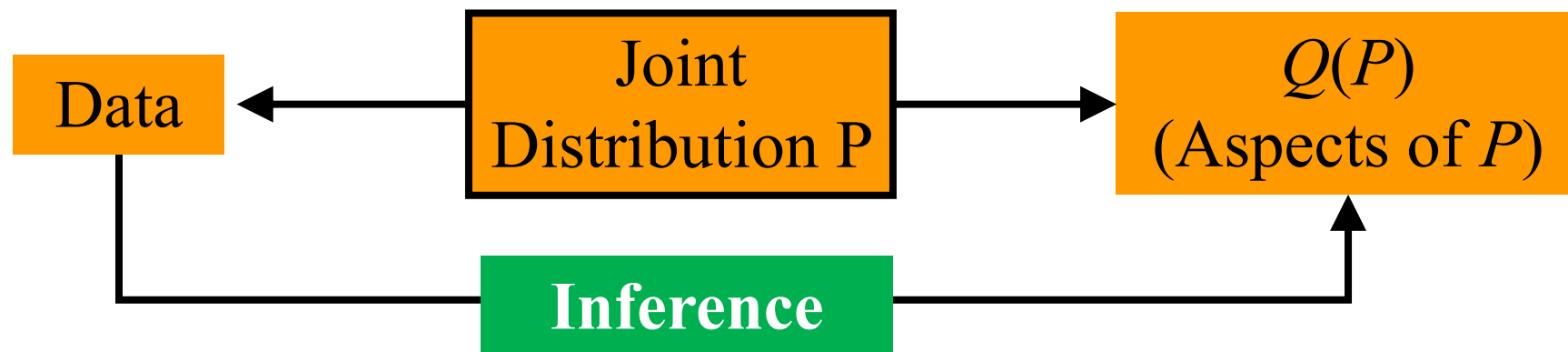
因果推理：Simpson's Paradox (辛普森悖论)

- 辛普森悖论表明，在某些情况下，忽略潜在的“第三个变量”（如性别就是专业和身高之外的第三个变量），可能会改变已有的结论，而我们常常却一无所知。从观测结果中寻找引发结果的原因，由果溯因，就是本节要介绍的因果推理

不能只满足于数字或图表，必须考虑数据生成
过程——因果模型

因果推理：Causal Inference

- 传统以统计建模为核心的推理手段



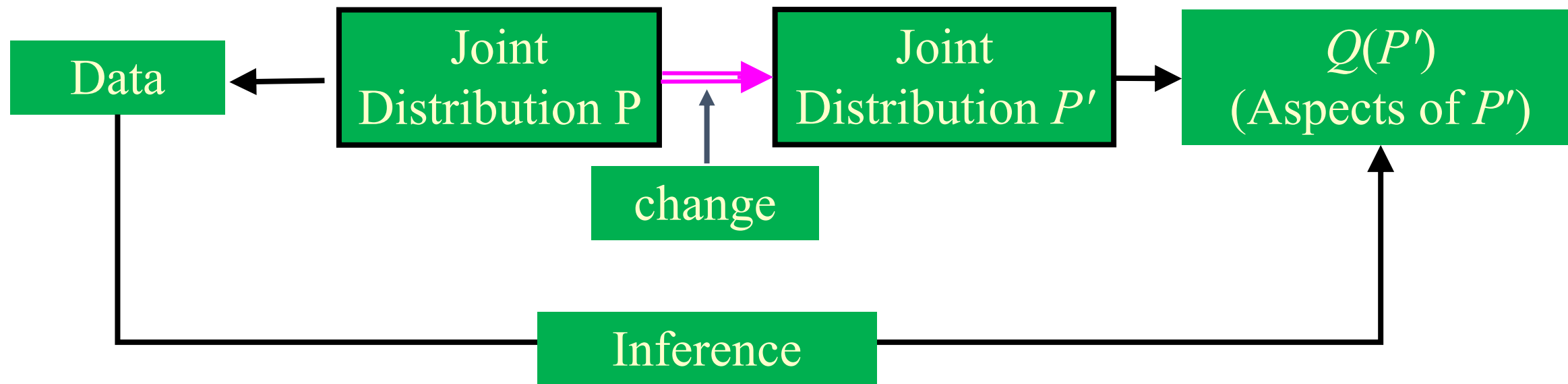
“The object of statistical methods is the reduction of data” (Fisher 1922).

购买了A商品的顾客是否会购买B商品(对A和B的联合分布建模)

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

因果推理：Causal Inference

- 如果商品价格涨价一倍，预测销售量 $P'(\text{sales})$ 的变化
- 如果放弃吸烟，预测癌症 $P'(\text{cancer})$ 的概率



数据分布从 P 变换到 P'

因果推理模型: 结构因果模型和因果图

- **结构因果模型(structural causal model, SCM)**
 - 也被称为因果模型 (causal model) 或Neyman–Rubin因果模型
 - Jerzy Neyman在1923年博士论文中(波兰语)提出的“潜在结果”(potential outcome) 的概念
 - 之后, Donald Rubin发展了“潜在结果”这一概念, 并将其和缺失数据的理论联系在一起。
- **因果图(causal diagram)**由Judea Pearl于 1995年提出。
- **每个结构因果模型 \leftrightarrow 都与一个因果图 \leftrightarrow 相对应**

因果推理的层级

- **Actions:** B will be true if we do A.
- **Counterfactuals:** B would be different if A were true

| | | |
|--------------------------------|--------------------------------------|--------------------------------|
| 可观测性问题 | What if we see A | $P(y A)$ |
| 决策行动问题 | What if we do A | $P(y do(A))$ (如果采取A行为, 则B真) |
| 反事实问题 (Counterfactual) | What if we did things differently | $P(y' A)$ (如果A为真, 则B将不同) |
| Options: with what probability | | |

因果推理：有向无环图

- **有向无环图(directed acyclic graphs, DAG)**

- 一个无回路的有向图，即从图中任意一个节点出发经过任意条边，均无法回到该节点。刻画了图中所有节点之间的依赖关系。

- **DAG 可用于描述数据的生成机制**

- 这样描述变量联合分布或者数据生成机制的模型，被称为 “贝叶斯网络” (Bayesian network)

因果推理：有向无环图（DAG）

- 对于任意的DAG，模型中 d 个变量的联合概率分布由每个节点与其父节点之间条件概率 $P(x_i | x_{pa(i)})$ 的乘积给出：

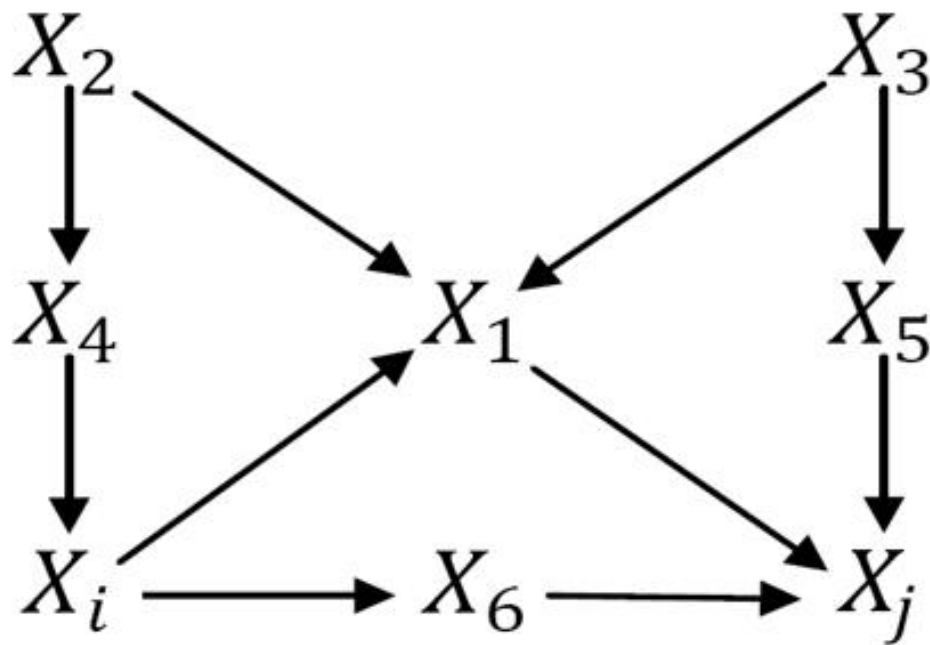
$$P(x_1, x_2, \dots, x_d) = \prod_{j=1}^d P(x_j | x_{pa(j)})$$

- 其中， $x_{pa(j)}$ 表示节点 x_j 的父节点集合（所有指向 x_j 的节点）

因果推理：有向无环图（DAG）

- 一个有向无环图唯一地决定了一个联合分布
- 一个联合分布不能唯一地决定有向无环图
 - 反过来的结论不成立
 - 如联合分布 $P(X_1, X_2) = P(x_1)P(x_2|x_1) = P(x_2)P(x_1|x_2)$

请写出下列有向无环图的联合概率形式，并区分内生变量和外生变量。



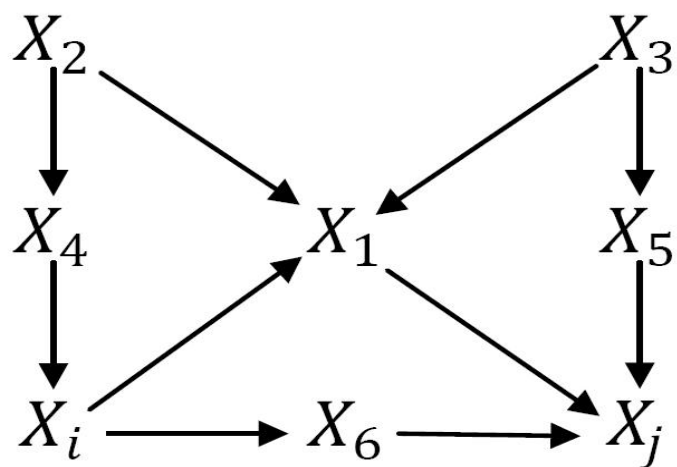
有向无环图DAG

正常使用主观题需2.0以上版本雨课堂

作答

因果推理：有向无环图（DAG）

- 一个有向无环图唯一地决定了一个联合分布
- 一个联合分布不能唯一地决定有向无环图
 - 反过来的结论不成立
 - 如联合分布 $P(X_1, X_2) = P(x_1)P(x_2|x_1) = P(x_2)P(x_1|x_2)$



有向无环图DAG

联合分布可表示为：

$$\begin{aligned} &P(X_1, X_2, X_3, X_4, X_5, X_6, X_i, X_j) \\ &= P(X_2) \times P(X_3) \times P(X_1|X_2, X_3, X_i) \\ &\times P(X_4|X_2) \times P(X_5|X_3) \times P(X_6|X_i) \times P(X_i|X_4) \\ &\times P(X_j|X_1, X_5, X_6) \end{aligned}$$

因果推理：有向无环图（DAG）

- 例：假设某个有向无环图中存在一条依赖路径 $X \rightarrow Y \rightarrow Z$ ，其中 X 表示气候好， Y 表示水果产量高， Z 表示水果价格低，给出 $P(\text{气候好}, \text{水果产量高}, \text{水果价格低})$ 的联合概率。
 - 使用乘积分解规则，将 $P(\text{气候好}, \text{产量高}, \text{价格低})$ 转换为：
$$P(\text{气候好}) \times P(\text{产量高} | \text{气候好}) \times P(\text{价格低} | \text{产量高})$$

假设：

$$P(\text{气候好}) = 0.5$$

$$P(\text{产量高} | \text{气候好}) = 0.8$$

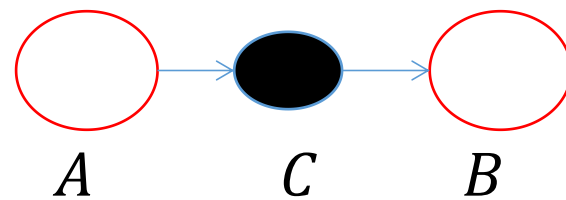
$$P(\text{价格低} | \text{产量高}) = 0.9$$

$$\begin{aligned} P(\text{气候好}, \text{产量高}, \text{价格低}) \\ = 0.5 \times 0.8 \times 0.9 = 0.36 \end{aligned}$$

因果推理：D-分离

- D-分离用于判断集合 A 中变量是否与集合 B 中变量相互独立（给定集合 C ），记为 $A \perp B \mid C$

D-分离的例子
(serial connection)



- 当 C 取值固定(可观测, observed), 有

$$P(A, B \mid C) = \frac{P(A, B, C)}{P(C)} = \frac{P(A)P(C \mid A) P(B \mid C)}{P(C)} = P(A \mid C)P(B \mid C)$$

- 可见 A 和 B 在 C 取值固定情况下, 是条件独立的
 - 注: 上式利用了 $P(A)P(C \mid A) = P(C)P(A \mid C)$

因果推理：D-分离

- D-分离用于判断集合 A 中变量是否与集合 B 中变量相互独立（给定集合 C ），记为 $A \perp B \mid C$

D-分离的例子
(diverging connection)



- 当 C 取值固定(observed), 有

$$P(A, B \mid C) = \frac{P(A, B, C)}{P(C)} = \frac{P(C)P(A \mid C)P(B \mid C)}{P(C)} = P(A \mid C)P(B \mid C)$$

- 可见 A 和 B 在 C 取值固定情况下，是条件独立的

- 如果 C 不固定，则有 $P(A, B) = \sum_C P(A \mid B)P(B \mid C)P(C)$
- 由于 $P(A, B) \neq P(A)P(B)$ ，因此 A 和 B 在条件 C 下不独立的

因果推理：D-分离

- D-分离用于判断集合 A 中变量是否与集合 B 中变量相互独立（给定集合 C ），记为 $A \perp B \mid C$

$$P(A, B, C) = P(A)P(B)P(C|A, B)$$

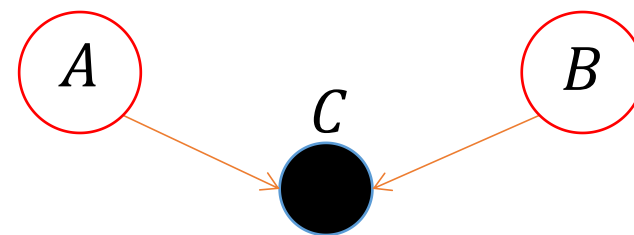
- 当 C 取值固定(observed)

$$P(A, B|C) = \frac{P(A, B, C)}{P(C)} = \frac{P(A)P(B)P(C|A, B)}{P(C)} \neq P(A)P(B)$$

- A 和 B 在条件 C 下是不独立的（是相关的）

D-分离的例子

(V-structure connection)



因果推理：D-分离

- D-分离用于判断集合 A 中变量是否与集合 B 中变量相互独立
(给定集合 C)，记为 $A \perp B \mid C$

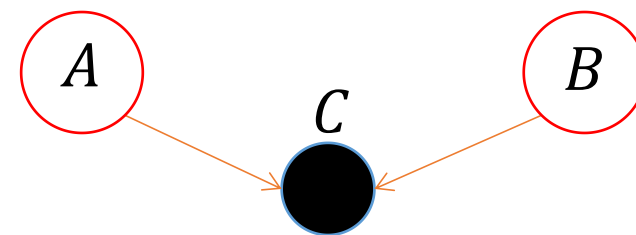
$$P(A, B, C) = P(A)P(B)P(C|A, B)$$

- 当 C 不作为观测点

$$P(A, B) = \sum_C (P(A)P(B)P(C|A, B)) = P(A)P(B) \sum_C P(C|A, B)$$

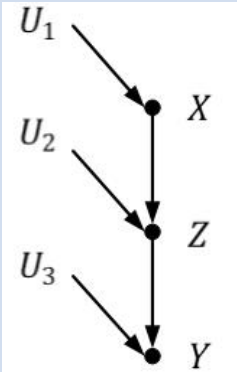
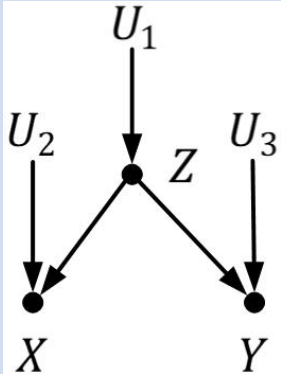
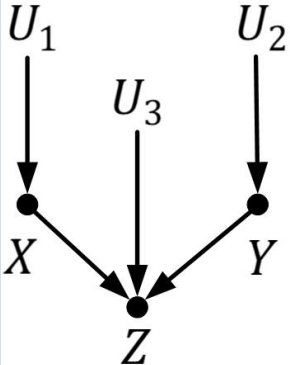
- A 和 B 在条件 C 下是独立的

D-分离的例子
(V-structure connection)



= 1

因果推理： ? -分离

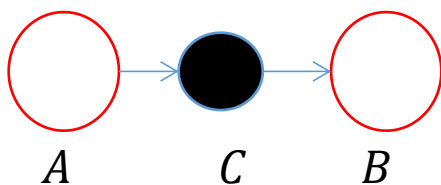
| 链结构(chain) | 分连结构(fork) | 汇连（或碰撞）结构(collider) |
|---|---|---|
|  |  |  |
| Z和X是相关的 | X和Z是相关的 | Z和X是相关的 |
| Y和Z是相关的 | Y和Z是相关的 | Z和Y是相关的 |
| Y和X很有可能是相关的 | Y和X很有可能是相关的 | Y和X是相互独立的 |
| 给定Z时，Y和X是条件独立的 | 给定Z时，Y和X是条件独立的 | 给定Z时，Y和X是相关的 |

因果推理：D-分离

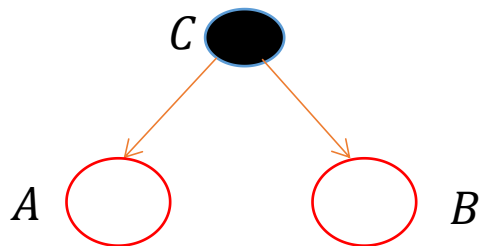
- **D-分离**：对于一个DAG图，如果 A 、 B 、 C 是三个集合（可以是单独的节点或者是节点的集合），为了判断 A 和 B 是否是 C 条件独立的，在DAG图中考虑所有 A 和 B 之间的路径(不管方向)。对于其中的一条路径，如果满足以下两个条件中的任意一条，则称这条路径是阻塞（block）的：路径中存在节点 X
 - X 是链结构或分连结构节点，且 $X \in C$
 - X 是汇连结构节点，并且 X 或 X 后代不包含在 C 中

因果推理：D-分离

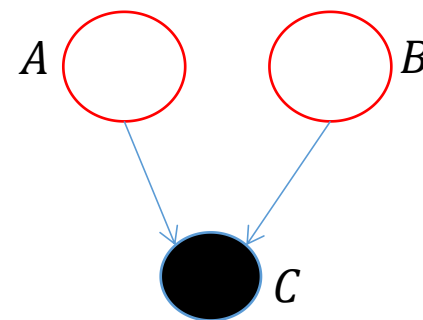
- **D-分离**：如果 X 和 Y 之间所有路径都是阻塞的，那么 X 和 Y 就是关于 C 条件独立的；否则 X 和 Y 不是关于 C 条件独立



链结构



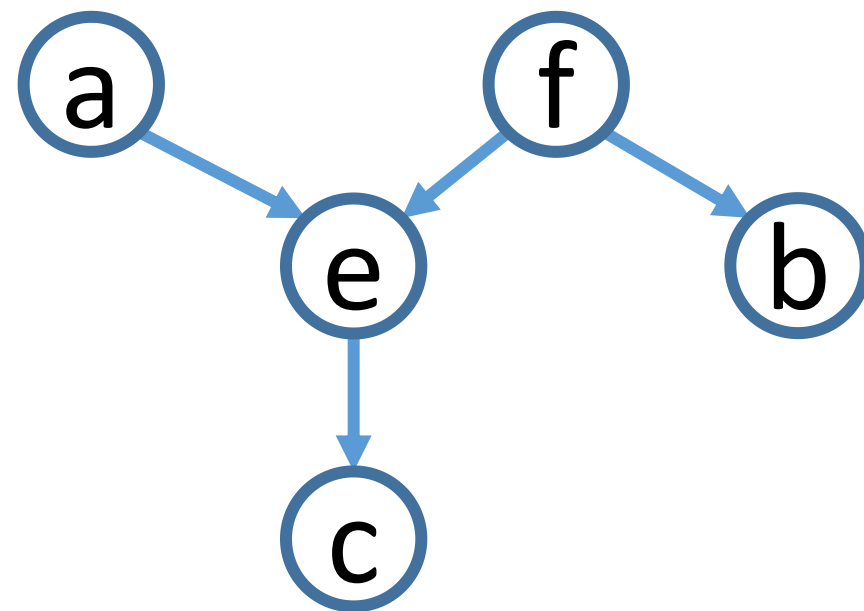
分连结构



汇连结构
(V结构或冲撞点, collider)

下面描述中正确的有

- ☐ A a和b是条件c下独立的
- ☐ B a和b是条件e下独立的
- ☒ C a和b是条件f下独立的



提交

因果推理：d-分离 (directional separation)

- d-分离方法可用于判断因果图上任意变量间相关性和独立性
- 在因果图上，若节点 X 和节点 Y 之间的每一条路径都是阻塞的
 - 称节点 X 和节点 Y 是有向分离的
 - 反之，称节点 X 和节点 Y 是有向连接的
- 当两个节点是有向分离时，意味着这两个节点相互独立

因果推理：干预(intervention)和do-算子(do-calculus)

- DAG中具有链接箭头的节点之间存在某种“因果关系”。
- 要在 DAG 上引入“因果”的概念，需要引进do 算子
 - do-calculus的意思可理解为“干预”(intervention)
- 在 DAG 中， $do(X_i) = x_i'$ ，表示将DAG中指向节点 X_i 的有向边全部切断，并且将 X_i 的值固定为常数 x_i'
- 在这样操作后，所得到新的DAG中变量联合分布为：

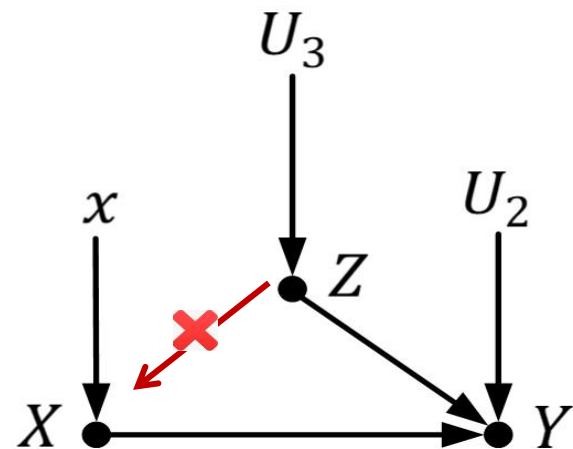
$$P(x_1, x_2, \dots, x_d | do(X_i) = x_i')$$

因果推理：干预(intervention)和do算子(do-calculus)

- **干预**指的是固定系统中某个变量，观察其他变量的变化
- 为了与自然取值进行区分，在对 X 进行干预时，引入**do算子**
 - 因此， $P(Y = y|X = x)$ 表示当 X 取值为 x 时， $Y = y$ 的概率；
 - 而 $P(Y = y|do(X = x))$ 表示对 X 取值进行了干预，固定其取值为 x 时， $Y = y$ 的概率。
- $P(Y = y|X = x)$ 反映了在取值为 x 的个体 X 上， Y 的总体分布；
- $P(Y = y|do(X = x))$ 反映的是如果将 X 每一个取值都固定为 x 时， Y 的总体分布。

因果推理：操纵图模型和操纵概率

- 对学院变量 X 进行干预并固定其取值为 x 时，可将所有指向 X 的边均移除。
- 因果效应 $P(Y = y \mid do(X = x))$ 等价于引入干预的**操纵图模型** (manipulated model)中条件概率 $P_m(Y = y \mid X = x)$



对 X 干预后操纵图模型

因果推理：操纵图模型和操纵概率

- 计算因果效应关键在于计算**操纵概率**(manipulated probability)

P_m 。 P_m 与无干预条件下概率 P 在如下两个方面的取值不变：

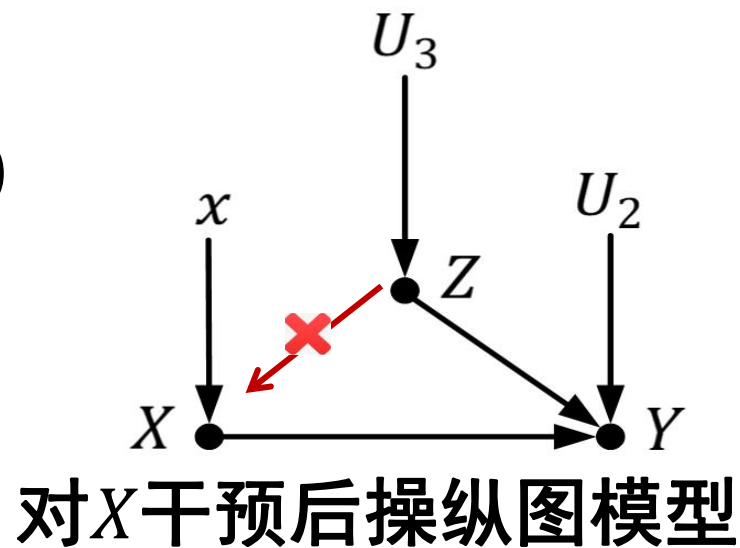
- 边缘概率 $P(Z = z)$ 不随干预而变化，因为 Z 的取值不会因为去掉从 Z 到 X 的箭头而变化，即：

$$P_m(Z = z) = P(Z = z)$$

- 条件概率 $P(Y = y \mid X = x, Z = z)$ 不变因为

Y 关于 X 和 Z 的函数 $f_Y = (X, Z)$ 并未改变，即：

$$P_m(Y = y \mid X = x, Z = z) = P(Y = y \mid X = x, Z = z)$$



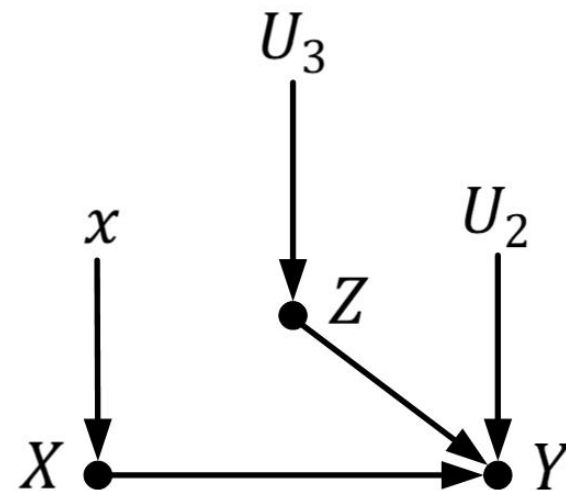
因果推理：调整公式

- 在干预图中， X 和 Z 是 D -分离的，因此是彼此独立的，即：

$$P_m(Z = z|X = x) = P_m(Z = z) = P(Z = z)$$

- 因果效应 $P(Y = y|do(X = x))$ 有：

$$\begin{aligned} P(Y = y|do(X = x)) &= P_m(Y = y|X = x) \\ &= \sum_z P_m(Y = y|X = x, Z = z)P_m(Z = z|X = x) \\ &= \sum_z P_m(Y = y|X = x, Z = z)P_m(Z = z) \\ &= \sum_z P(Y = y|X = x, Z = z)P(Z = z) \end{aligned}$$



对 X 干预后操纵图模型
 X 和 Z 相对于 Y 构成了V结构

因果推理：调整公式

- 因果效应 $P(Y = y|do(X = x))$ 有：

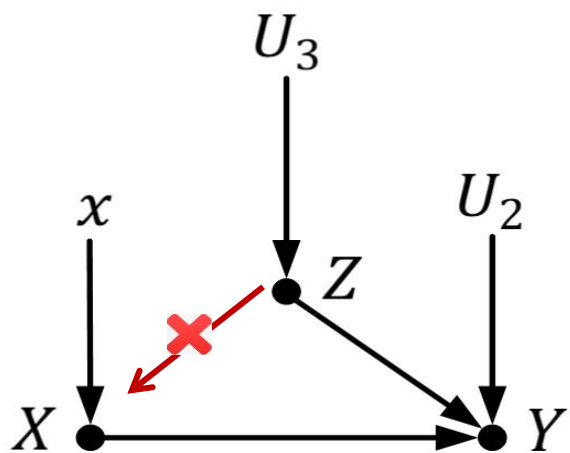
$$P(Y = y|do(X = x)) = \sum_Z P(Y = y|X = x, Z = z)P(Z = z)$$

- 该式被称为**调整公式** (adjustment formula)
 - 对于 Z 的每一个取值 z ，该式计算 X 和 Y 的条件概率并取均值，这个过程称为“ Z 调整” (adjusting for Z) 或“ Z 控制” (controlling for Z)
- 该式的右端只包含正常(无干预)条件下的概率 P
 - 即可用**正常(无干预)条件**下的条件概率来计算**干预后**的条件概率。

因果推理：调整公式

| | 计算机学院 | | 文学院 | |
|--------|-------|------|------|------|
| 身高(cm) | 男生 | 女生 | 男生 | 女生 |
| 矮个 | 35 | 25 | 10 | 70 |
| 高个人数 | 235 | 55 | 80 | 190 |
| 高个率 | 87 | 68.9 | 88.9 | 73.1 |

以性别分组后的计算机学院
和文学院的学生身高情况



对X干预后操纵图模型

$$P(Y = y|do(X = x)) \\ = \sum_z P(Y = y|X = x, Z = z)P(Z = z)$$

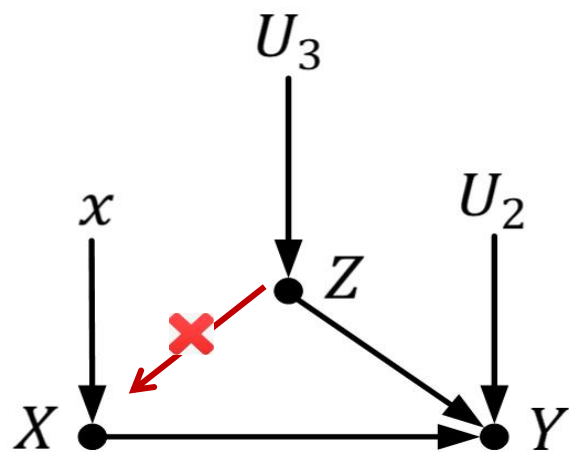
下面将调整公式用于计算对X取值进行
干预后计算机学院/文学院高个子率，
其中 $X = 1$ 表示计算机学院， $Y = 1$ 表示
高个， $Z = 1$ 表示表示男生，则有：

$$P(Y = 1|do(X = 1)) \\ = P(Y = 1|X = 1, Z = 1)P(Z = 1) \\ + P(Y = 1|X = 1, Z = 0)P(Z = 0)$$

因果推理：调整公式

| | 计算机学院 | | 文学院 | |
|--------|-------|------|------|------|
| 身高(cm) | 男生 | 女生 | 男生 | 女生 |
| 矮个 | 35 | 25 | 10 | 70 |
| 高个人数 | 235 | 55 | 80 | 190 |
| 高个率 | 87 | 68.9 | 88.9 | 73.1 |

以性别分组后的计算机学院
和文学院的学生身高情况



对X干预后操纵图模型

X干预取值为计算机学院的因果效应为：

$$\begin{aligned} P(Y = 1|do(X = 1)) \\ &= 0.87 \times \frac{(35 + 235 + 10 + 80)}{(350 + 350)} + 0.689 \times \frac{(25 + 55 + 70 + 190)}{(350 + 350)} \\ &= 0.782 \end{aligned}$$

X干预取值为文学院的因果效应为：



$$\begin{aligned} P(Y = 1|do(X = 0)) \\ &= 0.889 \times \frac{(35 + 235 + 10 + 80)}{(350 + 350)} + 0.731 \\ &\quad \times \frac{(25 + 55 + 70 + 190)}{(350 + 350)} = 0.812 \end{aligned}$$

其因果效应差：

$$\begin{aligned} ACE &= P(Y = 1|do(X = 1)) - P(Y = 1|do(X = 0)) \\ &= 0.782 - 0.812 = -0.03 \end{aligned}$$

可见计算机学院的高个率不如文学院的高个率。
即学院这一变量影响了高个率。

因果推理：因果图的不足

- 实际中难以得到一个完整的DAG
 - 用于阐述变量之间的因果关系或者数据生成机制，使得 DAG 的应用受到的巨大的阻碍
 - 从观测数据学习 DAG 的结构，是充满挑战的问题。
- Pearl 引入   **算子**是因果推断领域最主要贡献。
 - 从系统之外人为控制某些变量。但是，这依赖于一个假定：干预某些变量并不会引起 DAG 中其他结构的变化。
- DAG 作为一种简化的模型，在复杂系统中可能不完全适用
 - 需要将其拓展到动态系统（如时间序列），还有待研究。

因果推理：反事实推理 (counterfactual model)

- “反事实”框架是科学哲学家大卫·刘易斯 (David Lewis) 等人提出的推断因果关系的标准。
- 事实是指在某个特定变量(A)的影响下可观测到的某种状态或结果(B)。“反事实”是指在特定变量(A)取负向值时可观测到的状态或结果(B')
- 条件变量对于结果变量的因果性就是 A 成立时 B 的状态与 A 取负向值时“反事实”状态(B')之间的差异。
- 如果这种差异存在且在统计上是显著的，说明条件变量与结果变量存在因果关系。

推理总结

| 推理方法 | 推理方式 | 说明 |
|------|----------------------------|--|
| 归纳推理 | 如果 A_i (i 为若干取值), 那么B | 从若干事实出发推理出一般性规律 |
| 演绎推理 | 如果A, 那么B | A是B的前提、但不是唯一前提, 因此A是B的充分条件。当然, 在特殊情况下A也可作为B的充分必要条件 |
| 因果推理 | 因为A, 所以B | A是B的唯一前提, 因此“如果没有A, 那么没有B”也成立。 |

谢谢!