

第九章

知识表征

肖承丽

1

知识图谱 vs 大语言模型

将需要的知识数据（结构化或非结构化数据）以图谱的形式进行展示。通常包含实体、关系和属性三个要素，例如人名、国家、语言等实体，以及人口、首都、官方语言等关系和属性。

- 通常需要人工标注和整理，随着知识的不断扩展和变化，知识图谱也需要不断更新。
- 它的作用仅限于图谱编码中的知识和信息，这使得它在处理模糊两可或不完整的信息时准确率较低。
- 但它的优势在于能够生成以图式，对于事实性的、专业性的知识有着非常高的准确度。

迷你翻转课堂-7
人类的知識表征
vs
机器的知識表征

在大量文本数据上进行训练，以学习模式、上下文以及单词和短语之间的关系。

- 无法将真实与想象、真实与虚构分开——“人工智能幻觉”，人工智能的自信反应，其训练数据似乎没有合理性。
- 语言大模型需要大量的计算能力和资源来进行训练和微调，其花费的时间和成本也不容小觑。
- 目前，语言大模型只在通用领域给出了较为惊艳的表现，至于在知识图谱广泛应用的垂直领域，语言大模型暂时未显示出领先的一面。

下次课上分享，5-10 分钟

2

如何让AI获得不同类型的知识？

□ 陈述性？

- Word2vec 将词汇转换成数字向量 (2013)
- 无监督学习，发现词语之间的关联，如北京-中国
- 中国 + 河流 → 长江
- 国王 - 男人 + 女人 → 女王
- 医生 - 男 + 女 → 家庭主妇
- 计算机程序员 - 男 + 女 → 家庭主妇

□ 程序性？

- 模仿学习 Mobile ALOHA (2024)

3

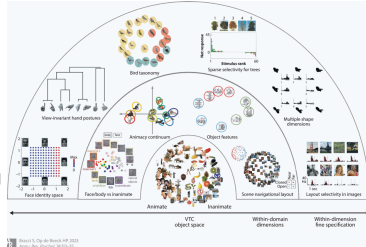
Understanding Human Object Vision: A Picture Is Worth a Thousand Representations

Annual Review of Psychology

人类物体视觉：与人类行为的全部行为目标联系起来，而非仅仅识别物体是什么

深度卷积神经网络：要有更大的表征多样性

- Hebart et al. (2020) 49种标签，“多形”、“圆形”、“动物相关”、“饮食相关”
- 对象域、行动域、社会域、空间/导航域
- 手和工具在脑中的表征重叠，二者对象域距离遥远，但行动域很近



4

知识

□ 信息在记忆中的存贮、整合和组织

- 信息——来源于感觉
- 知识 ≠ 信息
- 知识 = 经过组织（领悟）的信息

5

语言和知识

□ 语言作为种系发生学的划分标准

- 人类的语言发展水平远远超越了其他物种

□ 语言帮助我们了解我们的知识是怎样存储的

- 语义结构使我们能够辨认出哪些类别的事物被记忆存储，以及被存储的事物如何与头脑中的其他东西相联系
- 内容、结构、加工过程

6

试一试

你的推断与其他人非常类似！

“去年夏天，查理在卢浮宫见到了蒙娜丽莎。”

你有多大信心认为

1. 蒙娜丽莎和查理一起在塞纳河左岸喝咖啡？
2. 蒙娜丽莎向查理微笑？
3. 查理在巴黎吃了晚餐？
4. 查理带着欧元？
5. 查理的智商超过100？
6. 查理是个男性？

· 请你仔细听听其他人说出的简单句子，并且从单词在记忆中存储的方式这个角度进行一下分析。哪一种语义记忆模型更符合你的观察？

7

陈述性知识和程序性知识

□ 在任何一种知识表征的理论中，都有必要涵盖这两种类型的知识

- 陈述性知识
 - 知道是什么
 - 如，我把雨衣放在浴缸里，因为它是湿的。
- 程序性知识
 - 知道怎样做
 - 如，怎样洗澡

8

□ 知识在头脑中可能以多种不同的方式得到表征，包括：

关系	词汇表征	命题关系	图像	神经学成分
动作	珍妮=女人 打=动作 篮球=名词，运动	打{动作} (珍妮 {动作者}，篮球 {对象})		视觉皮层和联络皮层的一部分，也许还有部分运动皮层
属性	珍妮是个女人，她挺高的，她是个运动员。	珍妮是个女人。 珍妮挺高的。 珍妮是个运动员。		视觉皮层和联络皮层的一部分，也许还有右侧顶叶负责人脸识别的那部分区域
空间特性	珍妮拿着篮球。	珍妮的手上有一个篮球。		视觉皮层、联络皮层、运动和感觉皮层的一部分
类别成员	珍妮是以下类别的成员：女人、篮球运动员、高个子	有一种女人是珍妮。 有一种篮球运动员是珍妮。		联络皮层的一部分

9

“珍妮和她的朋友打篮球。”

关系	词汇表征	命题关系	图像	神经学成分
动作	珍妮=女人 打=动作 篮球=名词，运动	打{动作} (珍妮 {动作者}，篮球 {对象})		视觉皮层和联络皮层的一部分，也许还有部分运动皮层
属性	珍妮是个女人，她挺高的，她是个运动员。	珍妮是个女人。 珍妮挺高的。 珍妮是个运动员。		视觉皮层和联络皮层的一部分，也许还有右侧顶叶负责人脸识别的那部分区域
空间特性	珍妮拿着篮球。	珍妮的手上有一个篮球。		视觉皮层、联络皮层、运动和感觉皮层的一部分
类别成员	珍妮是以下类别的成员：女人、篮球运动员、高个子	有一种女人是珍妮。 有一种篮球运动员是珍妮。		联络皮层的一部分

10

第一节 早期的研究

1 语义组织

2 联想主义取向

11

语义组织

□ 在聚类模型中，语义组织通常被理解为

- 对意义相近的元素的分类和聚类
 - 如，里根，布什，尼克松（共和党总统）；克林顿，卡特尔，肯尼迪（民主党总统）

12

第一节 早期的研究

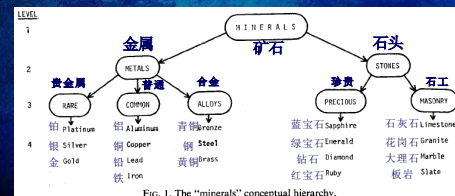
- 1 语义组织
- 2 联想主义取向

13

联想主义取向

□ 组织性变量 (Bower)

- 记忆中语义内容的组织化，对于记忆和回忆的实际影响远远大于人们此前的认定
- 概念层次对回忆产生的潜在影响(Bower et al., 1969)



14

□ 聚类模型 (Bousfield & Bower)

- 概念倾向于聚集成类
- 对“无关”单词的自由回忆表明，在类别上相似的单词会被一起回忆出来
 - 骆驼，驴子，马
 - John, Bob, Tom
 - 卷心菜，莴苣，菠菜

15

第二节 语义记忆：认知模型

□ 语义记忆的研究取向：联想主义观点⇒ 认知观点

- 假设：精细的**认知结构**表征着语义信息在记忆中的组织方式

16

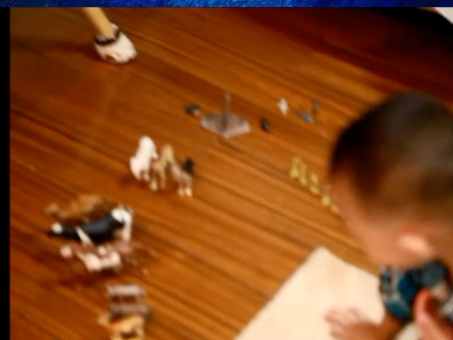
想一想

□ 请判断下列陈述的正误，并试着思考你是怎么得到这个答案的

- “所有金丝雀是鸟。”
- “一些动物是鸟。”
- “知更鸟是鸟。”
- “鸡是鸟。”
- “蝙蝠是鸟。”

17

4岁儿童的动物分类知识



18



19

集合-理论模型 (Meyer)

南京大学社会学院心理学系 肖承丽

□ 语义概念由元素（即信息的集合）来表征

- 样例集合
 - 鸟 = 金丝鸟、知更鸟、鹰...
- 属性集合
 - 鸟 = 翅膀、羽毛、鸣叫、飞行...

20

20

记忆提取是一个查证过程，即对两个或两个以上的信息集合进行搜索，找到其中重叠部分的样例

■ 如，“知更鸟是鸟”

知更鸟	鸟
物理对象	物理对象
活的	活的
会运动	会运动
有羽毛	有羽毛
红色胸羽	—
—	—
—	—

• 将两个集合的属性加以比较

• 集合属性的重叠程度构成了对命题正确性作出判断的基础

• 随着集合间距离增大，做出判断的反应时也会增加

21

21

两种逻辑关系

- 全称肯定(UA): 所有S都是P
 - “所有金丝雀都是鸟”
- 特称肯定(PA): 一些S是P
 - “一些动物是鸟”

• 判断特称肯定比判断全称肯定快

22

22



23

语义特征—比较模型 (Smith & Rosch)

与集合-理论模型在结构上相似，但是在一些重要假设上存在差别

- “单词的意义并非不可分析的单元，而是可以由一个语义特征集合来加以表征”

□ 如，“知更鸟”可以用下列特征来描述：有翅膀，两足动物，红色的胸羽，栖息在树上，喜欢吃虫子，难以驯服，是春天到来的预告

24

24

□ 单词的意义可以由一个语义特征集合来加以表征

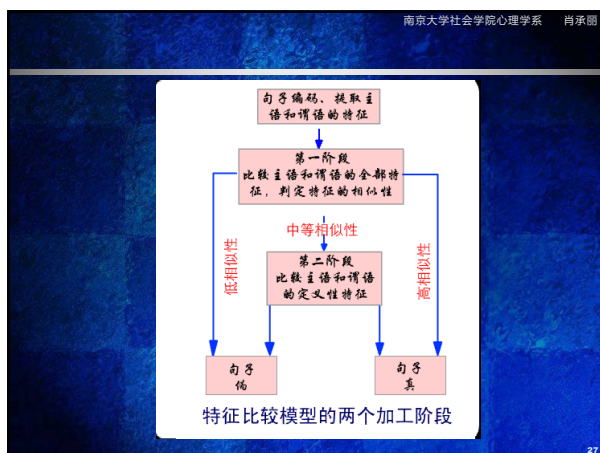
- 定义性特征 — 重要的，定义性的
 - 技术层面而言：“鸡是鸟。”
 - 鸡有喙，翅膀，羽毛
- 描述性特征 — 偶然的，描述性的
 - 宽泛而言：“蝙蝠是鸟。”
 - 蝙蝠会飞，有翅膀，看起来似乎是鸟

25

• 正确陈述的定义性特征和描述性特征均有很大重叠

语义模糊	陈述	被谓理性名词表征的特征	
		定义性	描述性
(正确陈述)	知更鸟是鸟	+	+
	麻雀是鸟	+	+
	长尾鸬鹚是鸟	+	+
技术层面而言	鸡是鸟	+	-
	鸭子是鸟	+	-
	鹅是鸟	+	-
宽泛而言	蝙蝠是鸟	-	+
	蝴蝶是鸟	-	+
	蛾子是鸟	-	+

26



27

□ 特征一比较模型可以解释这样的现象：一个类别中的某些成员要比另一些更加典型 (Rosch, 1977)

- 我看见一只鸟向南飞去。
- 鸟吃虫子。
- 树上有只鸟。
- 我听见有只鸟在我的窗台上唧唧叫。

■ 把鸟替换为：知更鸟，老鹰，鸵鸟，鸡

■ 评价每句句子的通顺程度

■ 类别的典型成员与类别的原型相似

28

□ 语义特征-比较模型的弱点 (Collins & Loftus, 1975)

- 没有任何一种单一属性是在定义某项事物中必不可少的
- 人们似乎很难判断一个特征到底是定义性的还是描述性的

29

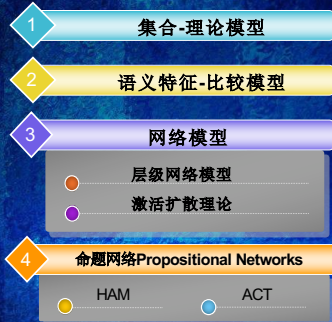
集合-理论模型&特征-比较模型

□ 两种模型都增进了人们对于语义记忆的理解：

1. 它们提供了关于语义记忆的多个维度的特定信息
2. 它们将语义分类的信息作为语义记忆的完整理论的起点，这种理论可以包容记忆技能的庞大网络
3. 由于这两个模型关系到复杂的记忆操作，因而他们触及了我们的知识表征这一更大的课题，其中最重要的部分涉及语义符号的存贮以及语义符号的回忆规律

30

语义记忆：认知模型



31

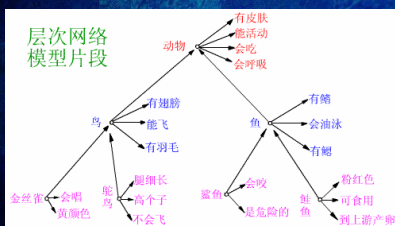
网络模型

- 在记忆中的知识，是相互独立的单元所连接而成的网络
- 单词的存储绑定到一个复杂的联系网络上
 - 如，“鸟”和“知更鸟”是用两者间的关系来储存的，那就是：知更鸟是鸟

32

层次网络模型 (Collins & Quillian)

- 对于任何一个单词，都用它与记忆中其他单词的相对关系来加以描述，其意义的是由它与其他单词的联系来表征的



- 信息存储所需空间的最小化
- 被认为是计算机存储设计中的经济选择

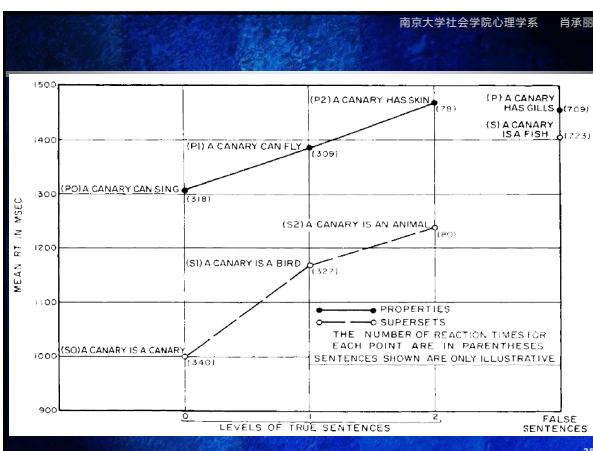
33

- 信息从语义记忆中的提取方式

- 如，“鲨鱼会运动。”
 - 鲨鱼——鱼——动物：会运动

- 所有这些概念结构中进行的搜索都需要花费时间

34



35

对层级网络模型的批评

- 典型性效应：网络中的联想强度是可变的
 - 如，“鸽子是鸟” & “企鹅是鸟”
- 某些联想关系破坏了系统的认知经济性
 - “鲨鱼会动”和“鱼会动”反应一样快
- 熟悉性效应 (Rips et al., 1973)
 - “狗是哺乳动物”比“狗是动物”反应更慢，因为我们更经常接触后者
- 不能解释否定判断：
 - 判断同一范畴的两个词比判断不同范畴的两个词需要花费更长的时间 (Glass & Holyoak, 1975)
 - “所有铁杉都是鸚鵡”比“所有铁杉都是雏菊”快

36

课堂演示

□小百合第六感测试

- 为什么大家的答案是类似的？
 - 这不是第六感
 - 典型性和原型：有些对象是其所属类别中“更好”的成员，他们更经常和更容易被选择。
 - 认知负荷增加了这种倾向性（这是为什么指导语要求你“快、快”）

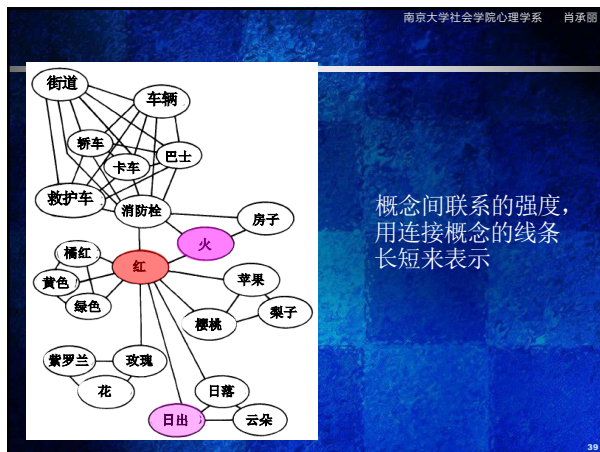
37

激活扩散理论

□Collins & Loftus (1975)

- 特定的记忆分布在概念空间里，并与其他与之有关的概念联结在一起

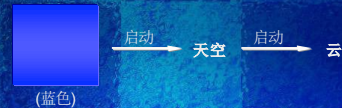
38



概念间联系的强度，
用连接概念的线条
长短来表示

39

- 激活过程在概念间扩散，可以用来解释启动实验



40

□批评

- 过于灵活、难以证伪
 - 没有决定连线长度的确定规则
 - 没有确定扩散后激活持续时间的规则
 - 没有说明多大程度的激活量才能激活一个节点的规则

41

语义记忆：认知模型



42

南京大学社会学院心理学系 肖承丽

命题网络(Propositional Networks)

- 复杂的观点可以通过简单的关系来表达
- 命题(Proposition): 可以作为独立断言的最小知识单位 (Anderson, 1985)

43

南京大学社会学院心理学系 肖承丽

人类联想记忆 (HAM)

- Human Associative Memory (Anderson & Bower, 1973)
- 基本表征单元是将概念连起来的**命题**，而不是单个的概念本身。
- 命题是**抽象的表征**，在形式上类似句子。
- 一个命题是由一小集**联想**构成的，每个联想则将两个概念结合在一起或联系起来。

HAM模型中的主要联想类型

联想	例子	联想	例子
上下文-事实	昨天在学校 约翰哭过	主语-谓语	他做 死了
	在家 我们吃		我的结婚 是教师
地点-时间	巴黎 1942	关系-宾语	高于 比尔
	学校 昨天		喜欢 雨

44

南京大学社会学院心理学系 肖承丽

命题树

- 结点(圆): 代表命题、上下文、事实、概念
- 指针: 代表联想

HAM模型的最大优点: 既可以表征语义记忆, 又可以表征情景记忆。

45

南京大学社会学院心理学系 肖承丽

思维的适应性控制 (ACT)

- Adaptive Control of Thought (Anderson, 1983)

- Working memory工作记忆: active memory
- Declarative memory 陈述性记忆: episodic & semantic
- Production memory²² 生式记忆: procedural knowledge

46

南京大学社会学院心理学系 肖承丽

产生系统: 人类的认知基础是一系列的**条件-行动对 (conditional-action pairs)**

IF-THEN clauses

- IF a is the father of b and b is the father of c
- THEN a is the grandfather of c
- IF the goal is to do an addition problem
- THEN the subgoal is to iterate through the columns of the problem
- IF the goal is to iterate through the columns of an addition problem and the rightmost column has not been processed
- THEN the subgoal is to iterate through the rows of the rightmost columns
-

67
39
72

47

南京大学社会学院心理学系 肖承丽

ACT里的知识表征

- A temporal string时间串
 - Encodes the order of a set of items, "one, two, three..."
- A spatial image空间图像
 - Encodes spatial representations, e.g., the coding of a square or triangle
- An abstract proposition抽象命题
 - Encodes meaning or semantic information, "Bill, John, hit"

48

□ A temporal string 时间串

- Records the sequential structure of events
- Recall the sequence of events in our daily experience.
 - E.g., the sequence of events in a movie or in a football game
- We are less able to fix, in absolute time, the occurrence of these events

49

□ A spatial image 空间图像

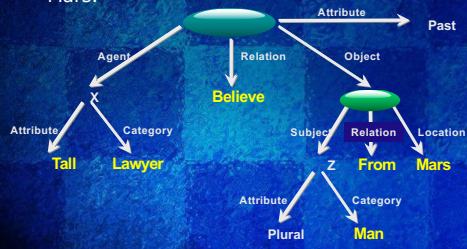
- Configural information: the type of information displayed in a figure, a form or even a letter
- But not the size

Z

50

□ An abstract proposition 抽象命题

- Similar to HAM
- E.g., "The tall lawyer believed the men were from Mars."



51

□ The Properties of the Three Representations

Process	Temporal String	Spatial Image	Abstract Proposition
Encoding process	Preserves temporal sequence	Preserves configural information	Preserves semantic relationships
Storage process	All or none of phrase units	All or none of image units	All or none of propositions
Retrieval process	All or none of phrase units	All or none of image units	All or none of propositions
Match process			
A. Degree of match	End-anchored at the beginning	Functions of distance and configurations	Function of set overlap
B. Salient properties	Ordering of any two elements, next element	Distance, direction, and overlap	Degree of connectivity
Execution: Construction of new structures	Combination of objects into linear strings, insertion	Synthesis of existing images, rotation	Insertion of objects into relational slots, filling in of missing slots

52

□ ACT (adaptive control of thought)

- Anderson has applied the system to a wide range of other conditions and cognitive tasks
 - Control of cognition
 - Memory for facts
 - Language acquisition
 - Spread of activation

53

第八章

知识表征(补充)

54

1、Connectionism (联结主义)

□ William James (*Psychology: The Briefer Course* [1892])

- “when two elementary brain-processes have been active together or in immediate succession, one of them, on recurring, tends to propagate its excitement into the other.”

□ Connectionism: A theory of mind that posits a large set of simple units connected in a parallel distributed network (PDP).

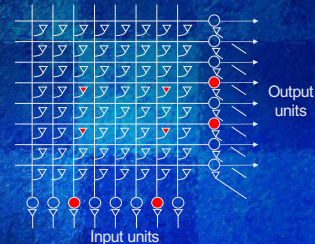
- Mental operations, such as memory, perception, thinking, and so on, are considered to be distributed throughout a highly complex neural network, which operates in a parallel manner.
- Based on the assumption: Units excite or inhibit each other throughout the system at the same time or in parallel.
- Not simply sequential

□ In connectionistic models the patterns themselves are not **stored**; what is stored is the connection strength between units, which allows these patterns to be recreated.

□ **Learning** consists of the acquisition of connection strengths that allow a net work of simple units to act as if they knew the rules. Rules are not learned; connections between simple units are.

□ PDP model is **neurally inspired**. The metaphor on which the model is based is the brain rather than the computer (see especially Collins and Quillian.)

□ Learning of associative relationships involves changing the strengths between the input units and the output units.



2图式(补充)

□ 请思考一下我们关于房子的知识:

- 房子是一种建筑物
- 房子里有一些房间
- 房子可以用木头、砖或石头来建造
- 房子可供人类居住
- 房子通常有直线和三角形
- 房子一般大于10平方米, 小于1000平方米

□ 图式依靠一些插槽 (slot) 结构来表征类别知识

■ 比如, 房子

插槽

- 上属: 建筑物
- 组成: 房间
- 材料: 木头, 砖, 石头
- 功能: 供人类居住
- 形状: 直线的, 三角形的
- 大小: 10~1000平方米

默认值

南京大学社会学院心理学系 肖承丽

□ 上属插槽

- 类似语义网络中的上属连接，指向上属概念
- 如果没有遇到冲突，那么一个概念就继承了其上属概念的特征
- 如，建筑物是房子的上属概念。建筑物的图式中存储了诸如房顶、墙壁以及建造在地面上这样的特征，这些信息不在房子的图式中表征，因为它可以由建筑物推论出来

□ 组成层次

- 比如墙和房间是房子的组成部分，有它们自己的图式定义，储存着它们自身的组成部分（比如窗户和天花板）。通过使用组成层次，我们能够推断出房子有窗户和天花板

61

南京大学社会学院心理学系 肖承丽

□ 图式是从具体实例中抽象出来的，它能够用来推断其所表征概念的实例的属性

- 比如我们知道某物是房子，那么我们就能够利用图式来推断它可能是由木头或砖制造的，并且有墙、窗户、天花板
- 基于图式的推论过程必须能够处理意外
 - 比如想象一座没有屋顶的房子
- 必须理解图式的插槽之间的约束关系
 - 比如听说一个房子建在地下，那么就可以推断出它没有窗户

62

南京大学社会学院心理学系 肖承丽

请参观这位心理学家的办公室



63

南京大学社会学院心理学系 肖承丽

□ 布鲁尔和特莱耶斯 (Brewer & Treyens, 1981)

- 30个被试被分别带入上图的房间，告知他们这是主试的办公室，他们要在这里等候，因为主试要去实验室看一下前一个被试是否已经完成了实验
- 35秒之后，主试返回并将等待的被试带到邻近的会议室，在这里，被试要写下他能记住的主试办公室中的所有东西
 - 包含在办公室图式中的项目回忆非常好
 - 29人回忆出一把椅子、一张书桌和墙壁
 - 不包含在办公室图式中的项目回忆较差
 - 8个人回忆出公告板、一个骷髅头
 - 错误记忆
 - 9个人回忆出有一些书（实际上并没有）

64

南京大学社会学院心理学系 肖承丽

决定性时刻包含以下一种或多种因素：

欣喜时刻

认知时刻

荣耀时刻

连接时刻

第四章 打破脚本

约翰点了一个汉堡包。端上的汉堡包是凉的。他留下了很少的一笔小费。

哈里特参加了杰克的生日派对。蛋糕很难吃。哈里特给杰克的妈妈留下了很少的一笔小费。

打造一流服务

赫恩的儿子在度假时把心爱的长颈鹿玩具乔希落在了酒店。没有乔希儿子不肯睡觉。

• 哄儿子说长颈鹿还在度假

• 打电话请酒店拍一张长颈鹿在泳池边的照片

几天后，收到乔希+一叠乔希度假的照片

65

南京大学社会学院心理学系 肖承丽



结果：赫恩心花怒放，儿子欢呼雀跃；该故事的博文被大家疯狂传播

66

课后思考

□遵循知识表征or打破图式，在日常生活中可以有什么应用or创新？

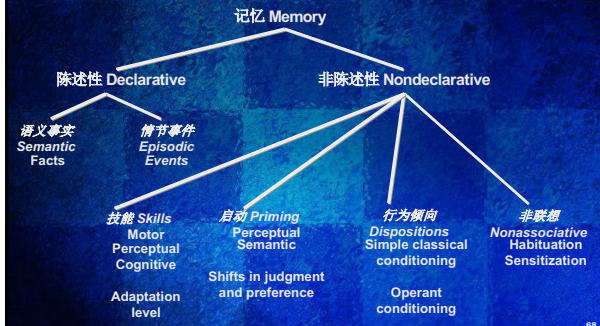
- 如，入职时如何提升归属感？
- 用户体验？
- 产品营销？

67

3、记忆结构分类

A taxonomy of memory structure

□Squire (1986) & Squire et al. (1990)



68

- The brain is organized around fundamentally different information storage systems
- This system accepts both conscious and unconscious memory as serious topics for research

69

Amnesia

- Studying people with amnesia lends support for the existence of STM and LTM
- Retrograde amnesia逆行性遗忘
 - Inability to recall information acquired prior to the onset of the disorder.
- Anterograde amnesia顺行性遗忘
 - The loss of information presented after the onset of the memory disorder.
- Both types of amnesia can be temporary or permanent
 - Temporary
 - Head trauma, electroconvulsive (shock) therapy (ECT)
 - Permanent
 - Prolonged use of ECT, severe traumas, stroke, cerebral vascular ruptures, massive consumption of alcohol (Korsakoff's syndrome)

70

Cognitive tasks and amnesia

□H. M. (Milner, Corkin, & Teuber, 1968)

- Some complex cognitive ability remained intact
- Episodic memory was seriously impaired
- (Chapter 6)

□Cohen & Squire (1980)

- Could acquire the skill involved in reading words from a mirror-reversed display
- Could neither remember the words nor the skill they had practiced
- Skill learning & episodic recall
- (See also K. C. in Chapter 7)

71

□Jacobson & Witherspoon (1982)

- Prime: "What is an example of a reed instrument (簧乐器)?"
- Spell test: spell [ri:d] (read/reed)
 - Control & Korsakoff groups: reed
- Recognition test
 - Korsakoff group: unable to recognize it they had heard
 - Control group: able
- Priming effect is unaffected by amnesia, while reorganization is affected
- An important distinction in the way knowledge is represented may be in terms of the **activation of neural pathways** and the **access** of episodic memory

72

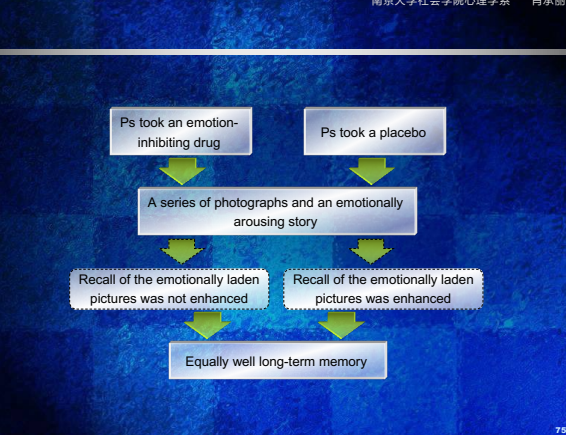
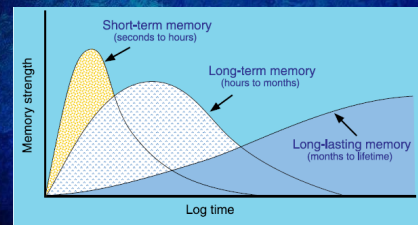
4、Memory Consolidation 记忆巩固

□ Perseveration-consolidation theory 重复-巩固理论 (Muller & Pilzecker)

- Memory for new material was disrupted by the learning of new information shortly after the original learning.
- It was as if memories, like photographic images, need a period of time to become "fixed".

□ McGaugh (2000)

- STM may not be 'bumped' into LTM as a sequential process
- Specific drugs can selectively block short-term memory or long-term memory, which suggests that the two memory processes are independent and operate in parallel.



5. Organization of concepts in the brain

How is that knowledge organized in the brain?

1. One obvious possibility is that all information we possess about any given object or concept is stored in **one location** in the brain.
2. Another possibility is that different kinds of information (features) about a given object are stored in **different locations** in the brain.
 - "Object concepts may be represented in the brain as distributed networks of activity in the areas involved in the processing of perceptual or functional knowledge"

知觉-功能理论 Perceptual-functional theories

□ Warrington and Shallice (1984) and Farah and McClelland (1991)

□ 重要区分:

- **知觉特征** (如, 那个物体看起来是什么样的?)
 - 对生物的语义知识
- **功能特征** (如, 那个物体有什么用?)
 - 对非生物的语义知识

□ 额外假设: 语义记忆中知觉属性的信息量远超功能属性 **semantic memory contains far more information about perceptual properties** of objects than of functional properties

■ 词典中对生物和非生物的解释:

□ 视觉: 功能描述

- 生物是 7.7:1
- 非生物是 1.4:1

□两个假设:

1. 脑损伤通常导致生物比非生物知识受损更严重.
 - 脑损伤可能损毁更多的知觉特征而非功能特征信息, 因为前者原本就存储得更多.
2. 脑神经成像应该能揭示知觉特征和功能特征分别激活不同的脑区

79

□许多脑损伤病人显示出**特定类型受损 (category-specific deficits)**, 即他们对特定类型的物体存在认知障碍

- 病人JBR: 识别生物比非生物图片更困难 (正确率6% vs 90%). (Warrington & Shallice, 1984)
- 特定类型受损病人的研究综述 (Martin & Caramazza, 2003):
 - 超过100名病人识别生物有障碍, 但识别非生物无碍
 - 仅有25位表现出相反的模式

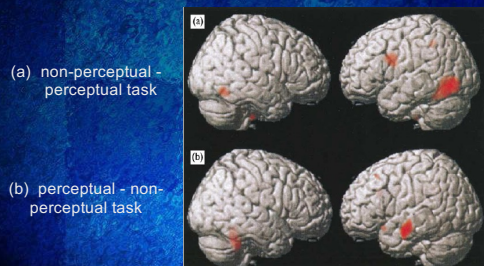
80

□对44位病人的综述(Gainotti, 2000)

- 38 位病人选择性地缺乏生物知识
 - 几乎所有人都有颞叶受损damage to the anterior, medial, and inferior parts of the **temporal lobes**
- 6 位病人选择性地缺乏人造物知识
 - 额-顶叶及往后延伸的区域受损

81

□请健康被试对生物或非生物提取知觉或非知觉信息 (Lee, Graham, Simons, Hodges, Owen, and Patterson, 2002)



82

多重属性取向 **Multiple-property approach**

□知觉-功能理论过于简单, 如

- 许多生物的属性既不感觉也不功能 (如, 肉食动物, 沙漠生物)
- 对功能的定义非常广泛, 通常包括一个物体的**用途 (uses)** 以及如何操作 (**how it is manipulated**)

83

功能知识应该进一步分为 "what for" 和 "how"

□Buxbaum and Saffran (2002)

- Apraxia失用症: 额顶叶受损的失用症病人保留了关于物体用途的知识, 但丢失了如何操作物体的知识.
- 非失用症病人: 颞叶受损, 表现出相反的模式

84

南京大学社会学院心理学系 肖承丽

□Canessa et al. (2008) fMRI; Healthy ps

action or "how" knowledge functional or "what for" knowledge

85

南京大学社会学院心理学系 肖承丽

□Cree and McRae (2003): 区分知觉和功能属性过于简单，脑损伤病人表现出各种类型的知识受损

TABLE 7.1: Cree and McRae's (2003) explanation of why brain-damaged patients show various patterns of deficit in their knowledge of different categories. From Smith and Kosslyn (2007). Copyright © Pearson Education, Inc. Reproduced with permission.

Deficit pattern	Shared properties
1. Multiple categories consisting of living creatures	Visual motion, visual parts, colour
2. Multiple categories of non-living things	Function, visual parts
3. Fruits and vegetables	Colour, function, taste, smell
4. Fruits and vegetables with living creatures	Colour
5. Fruits and vegetables with non-living things	Sound, colour
6. Inanimate foods with living things (especially fruits and vegetables)	Function, taste, smell
7. Musical instruments with living things	Function

86

南京大学社会学院心理学系 肖承丽

Grounded cognition 扎根认知

□Barsalou (2008)

■ “拒绝语义记忆是通道无关的符号的标准观点.....关注模拟在语义记忆中的作用.....模拟是重新激活在体验时获取的知觉、运动和内省状态”

“reject the standard view that amodal symbols represent knowledge in semantic memory . . . [they] focus on the roles of simulation in cognition. . . . Simulation is the re-enactment of perceptual, motor, and introspective states acquired during experience (p. 618).

87

南京大学社会学院心理学系 肖承丽

□Hauk, Johnsrude, and Pulvermüller (2004).

■ 当被试看到单词如“舔”，“捡”，“踢” (“lick”, “pick”, and “kick”)，这些动词也会激活运动皮层的部分区域，其位置与执行这些动作激活的区域重叠（或非常接近）

88

南京大学社会学院心理学系 肖承丽

□Pulvermüller, Hauk, Nikulin, and Ilmoniemi (2005).

TMS: 单词呈现150ms后对左侧大脑（语言中枢）施以阈下单次磁刺激

- 手臂区域的TMS导致对手臂词语(e.g. pick)比腿部词语(e.g. kick)的反应更快
- 腿部区域的TMS导致对腿部词语比手臂词语的反应更快

89

南京大学社会学院心理学系 肖承丽

90



91