



第1章 绪论

王博：自动化（人工智能）学院
wangbo@hdu.edu.cn



章节目录

- 基本术语
 - 假设空间
 - 归纳偏好
 - 发展历程
 - 应用现状
-



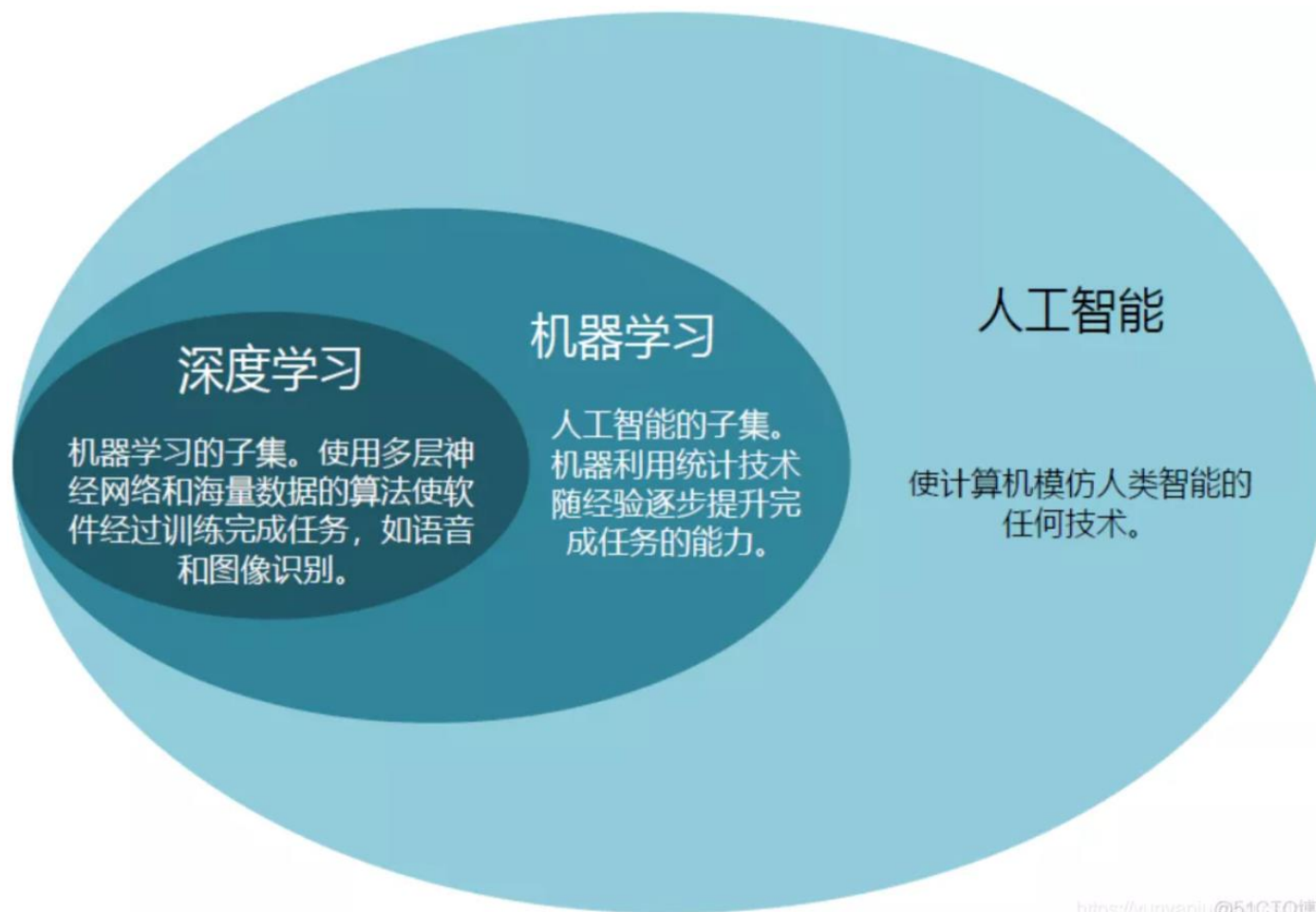
章节目录

- 基本术语
 - 假设空间
 - 归纳偏好
 - 发展历程
 - 应用现状
-



机器学习

机器学习致力于研究如何**通过计算**的手段，利用经验来改善系统自身的性能，从而在计算机上从数据中**产生“模型”**，用于对新的情况给出判断。





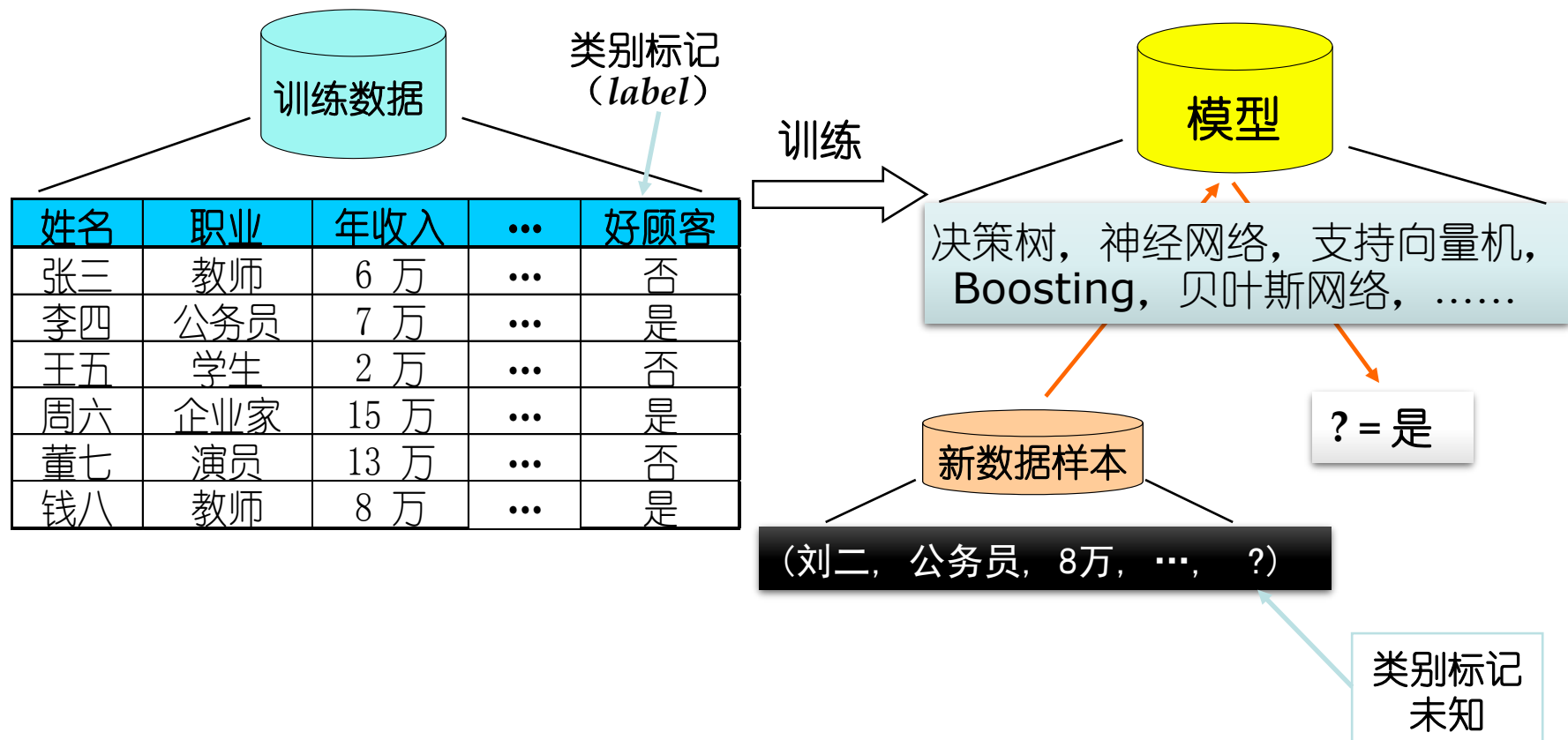
基本术语-数据

	特征				标记
	编号	色泽 ✓	根蒂	敲声	好瓜 ✓
训练集	1	青绿	蜷缩	浊响	是
	2	乌黑	蜷缩	沉闷	是
	3	青绿	硬挺	清脆	否
	4	乌黑	稍蜷	沉闷	否
测试集	1	青绿	蜷缩	沉闷	?



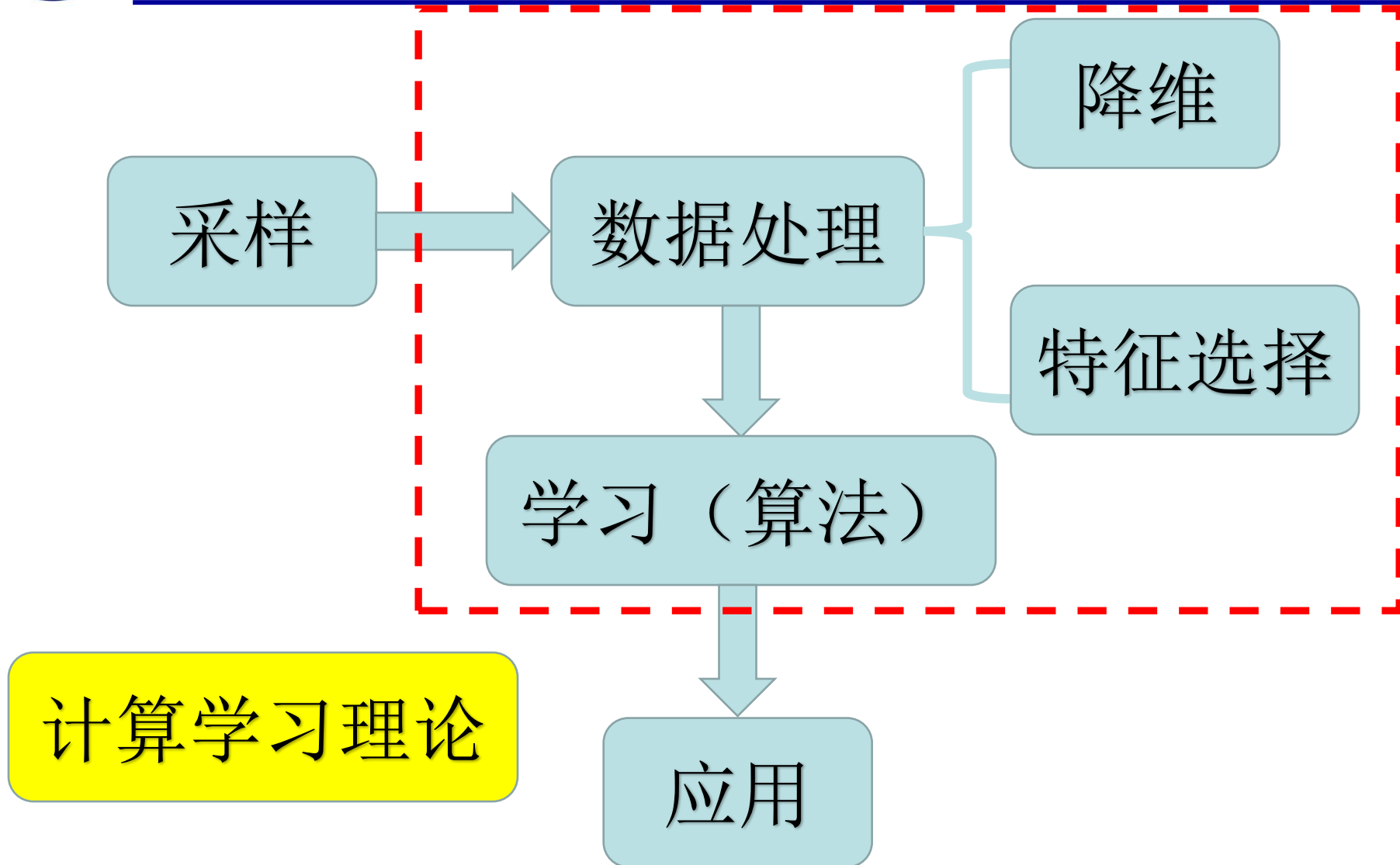
基本术语-数据-模型

使用学习算法 (*learning algorithm*)





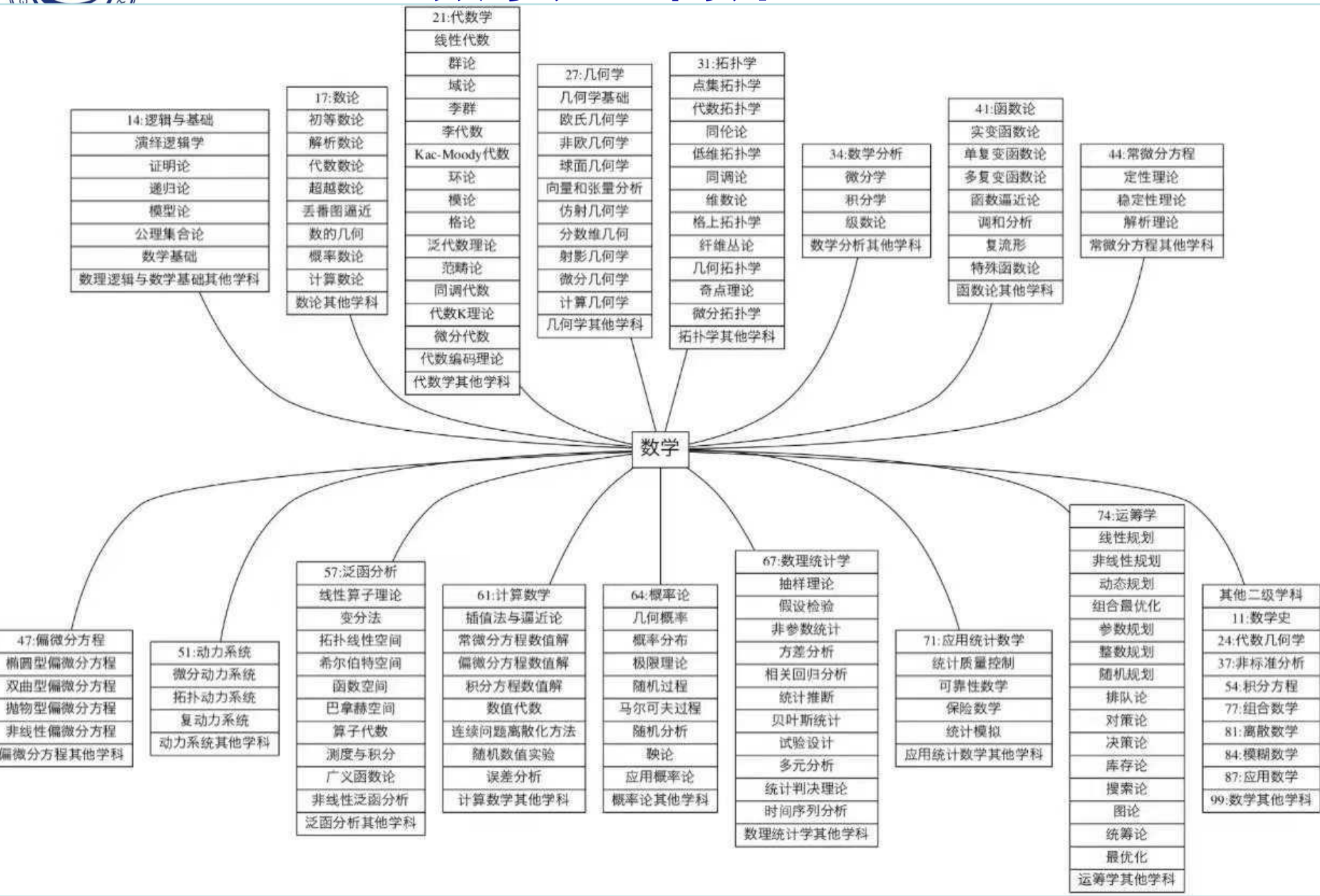
基本内容

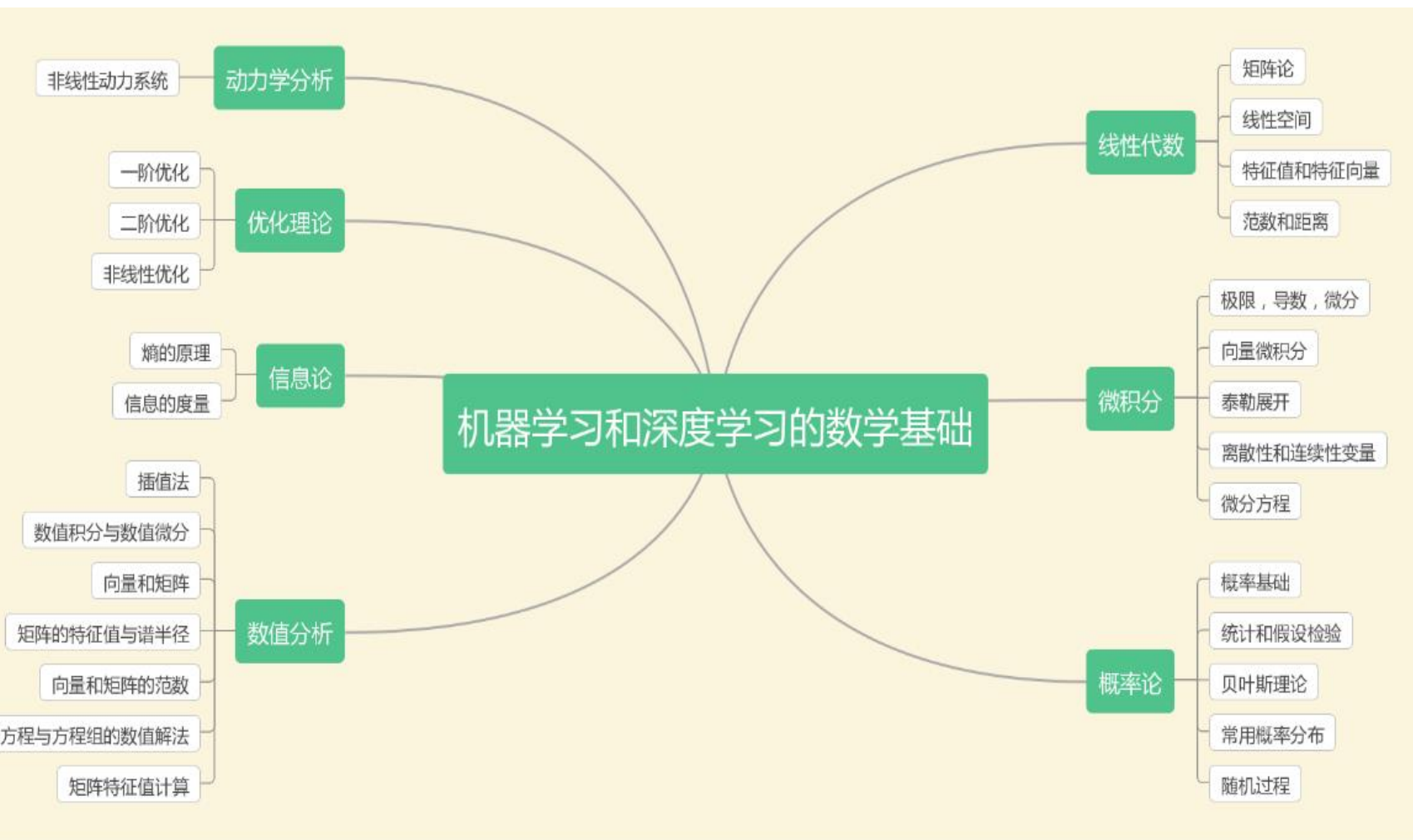






数学分支简介







基本术语-任务

预测目标:

- 分类:离散值
 - 二分类:好瓜;坏瓜
 - 多分类:冬瓜;南瓜;西瓜
- 回归:连续值
 - 瓜的成熟度
- 聚类:无标记信息



基本术语-任务

有无标记信息

- 监督学习：分类、回归
- 无监督学习：聚类
- 半监督学习：两者结合



基本术语-泛化能力

机器学习的目标是使得学到的模型能很好的适用于“新样本”，称模型适用于新样本的能力为**泛化 (generalization) 能力**。

一般而言训练样本越多越有可能通过学习获得强泛化能力的模型



章节目录

- 基本术语
 - 假设空间
 - 归纳偏好
 - 发展历程
 - 应用现状
-



假设空间

假设空间：由输入空间到输出空间的映射的集合。也就是由输入空间 X 到输出空间 Y 的映射 $f : X \rightarrow Y$ 所构成的集合，该空间是一个函数空间，即由函数所构成的集合（李航-机器学习方法）

集合元素个数数目： $|Y|^{|X|}$

学习过程：在所有假设空间组成的空间中进行搜索的过程。

编号	色泽	根蒂	敲声	好瓜
1	青绿	蜷缩	浊响	是
2	乌黑	蜷缩	沉闷	是
3	青绿	硬挺	清脆	否
4	乌黑	稍蜷	沉闷	否

$(\text{色泽}=?)\wedge(\text{根蒂}=?)\wedge(\text{敲声}=?)\leftrightarrow\text{好瓜}$

色泽2种、根蒂2种、敲声3种：假设空间： $2^{(2*2*3)}$



假设空间

可以自行选择（找子集）所考虑问题的假设空间的大小范围，并不一定要选择包含所有可能性的映射集作为假设空间。（周志华）

编号	色泽	根蒂	敲声	好瓜
1	青绿	蜷缩	浊响	是
2	乌黑	蜷缩	沉闷	是
3	青绿	硬挺	清脆	否
4	乌黑	稍蜷	沉闷	否

$$(\text{色泽}=\text{?}) \wedge (\text{根蒂}=\text{?}) \wedge (\text{敲声}=\text{?}) \leftrightarrow \text{好瓜}$$

通配符(*)考虑在内，在模型空间中搜索不违背训练集的假设空间大小：
 $3*3*4+1=37$ （考虑好瓜不存在）



版本空间

版本空间就是与训练集一致的所有假设所构成的集合，也就是假设空间的一个最大子集，该子集内的每一个元素都不与训练集相冲突。

由于每个假设都是一个完全确定的映射，可以将其与训练集中的每一个样例作比较，所得结果要么有冲突，要么没冲突。有冲突就不放进版本空间，都没冲突就放进版本空间



章节目录

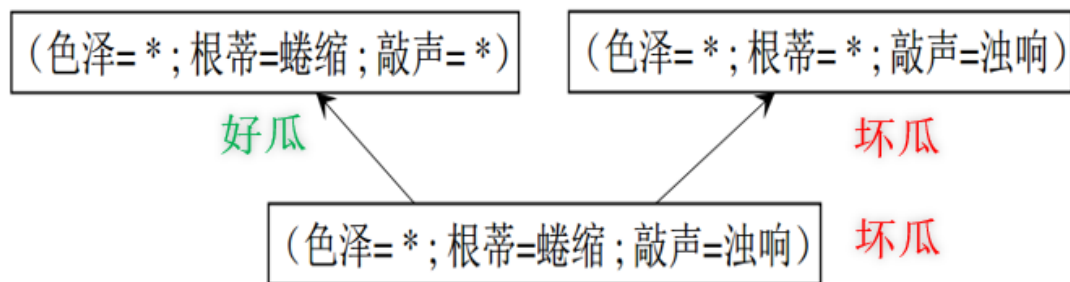
- 基本术语
 - 假设空间
 - 归纳偏好
 - 发展历程
 - 应用现状
-



归纳偏好

归纳偏好：学习过程中对某种类型假设的偏好称作归纳偏好

假设空间中有三个与训练集一致的假设，但他们对(色泽=青绿；"根蒂=" 蜷缩；敲声=沉闷)的瓜会预测出不同的结果：

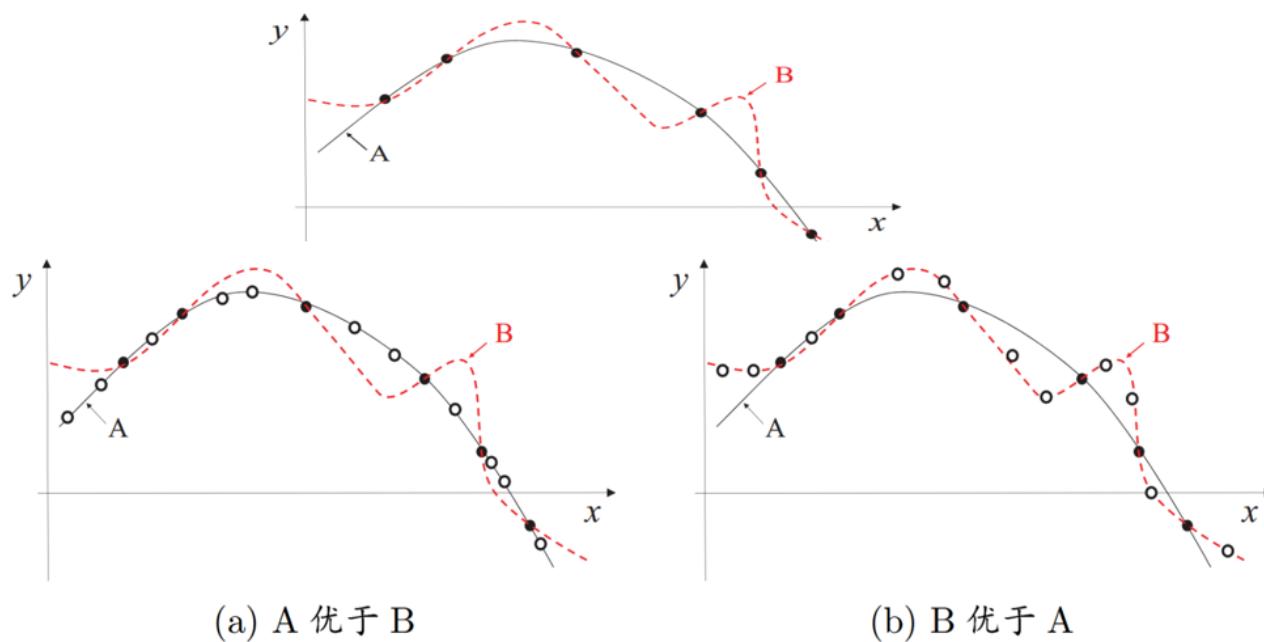


选取哪个假设作为学习模型？



归纳偏好

回归问题



没有免费的午餐. (黑点: 训练样本; 白点: 测试样本)

对于一个学习算法 L_a , 某些问题是比 L_b 好, 必然在一些问题上不如 L_b



归纳偏好

选取哪个假设作为学习模型？什么样的模型更好？

- “**奥卡姆剃刀**”是一种常用的、自然科学研究中最基本的原则，即“若有多个假设与观察一致，选最简单的那个”。
- KNN 中假设特征空间中相邻的样本倾向于属于同一类
- SVM 中假设好的分类器应该最大化类别边界距离
- 深度神经网络偏好性地认为，层次化处理信息有更好效果；
- 卷积神经网络 认为信息具有空间局部性，可用滑动卷积共享权重的方式降低参数空间；
- 循环神经网络 则将时序信息纳入考虑，强调顺序重要性；
- 图网络则认为中心节点与邻居节点的相似性会更好地引导信息流动。

具体的**现实问题**中，**学习算法本身所做的假设是否成立**，也即算法的归纳偏好是否与问题本身匹配，大多数时候直接决定了算法能否取得好的性能。



归纳偏好

最近邻居	假设在特征空间 (feature space) 中一小区域内大部分的样本是同属一类。给一个未知类别的样本, 猜测它与它最紧接的大部分邻居是同属一类。这是用于最近邻居法的偏置。这个假设是相近的样本应倾向同属于一类别。 KNN就是基于这种思想
最少特征数	除非有充分的证据显示一个特征是有效用的, 否则它应当被删除。这是特征选择算法背后所使用的假设。
最大边界	当要在两个类别间画一道分界线时, 试图去最大化边界的宽度。这是用于 支持向量机 的偏置
最小描述长度	当构成一个假设时, 试图去最小化其假设的描述长度。假设越简单, 越可能为真的。 奥卡姆剃刀 的理论基础
最大条件独立性	如果假说能转成贝叶斯模型架构, 则试着使用最大化条件独立性。这是用于 朴素贝叶斯 分类器的偏置。
最小交叉验证误差	当试图在假说中做选择时, 挑选那个具有最低交叉验证误差的假说。

脱离问题谈“什么学习算法好”毫无意义!!!



章节目录

- 基本术语
 - 假设空间
 - 归纳偏好
 - 发展历程
 - 应用现状
-



发展历程

推理期（20世纪50-70年代）：

- A. Newell和H. Simon的“逻辑理论家” (Logic Theorist) 程序以及此后的“通用问题求解” (General Problem Solving) 程序等在当时取得了令人振奋的结果；

知识期（70年代中期）：

- 大量专家系统问世，在很多应用领域取得大量成果，由人来总结知识再交给计算机相当困难。



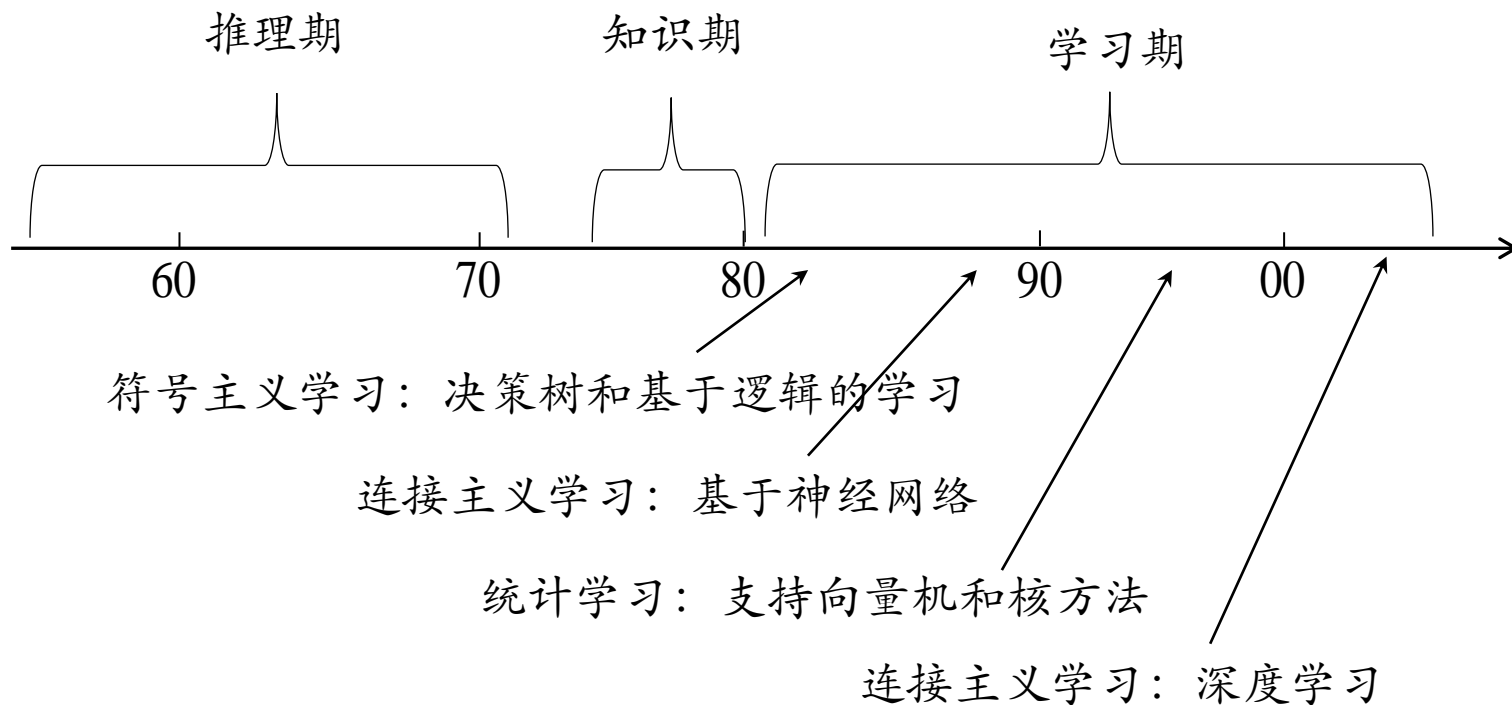
发展历程

学习期:

- 符号主义学习（20世纪80年代）
 - 决策树：以信息论为基础，最小化信息熵，模拟了人类对概念进行判定的树形流程
 - 基于逻辑的学习：使用一节逻辑进行知识表示，通过修改扩充逻辑表达式对数据进行归纳
- 连接主义学习
 - 感知机（20世纪50年代）
 - Hopfield神经网络（80年代）
 - BP神经网络（80年代）
- 统计学习（90年代）
 - 支持向量（60-70年代）开始研究
 - 支持向量机及核方法



发展历程





章节目录

- 基本术语
 - 假设空间
 - 归纳偏好
 - 发展历程
 - 应用现状
-



应用现状

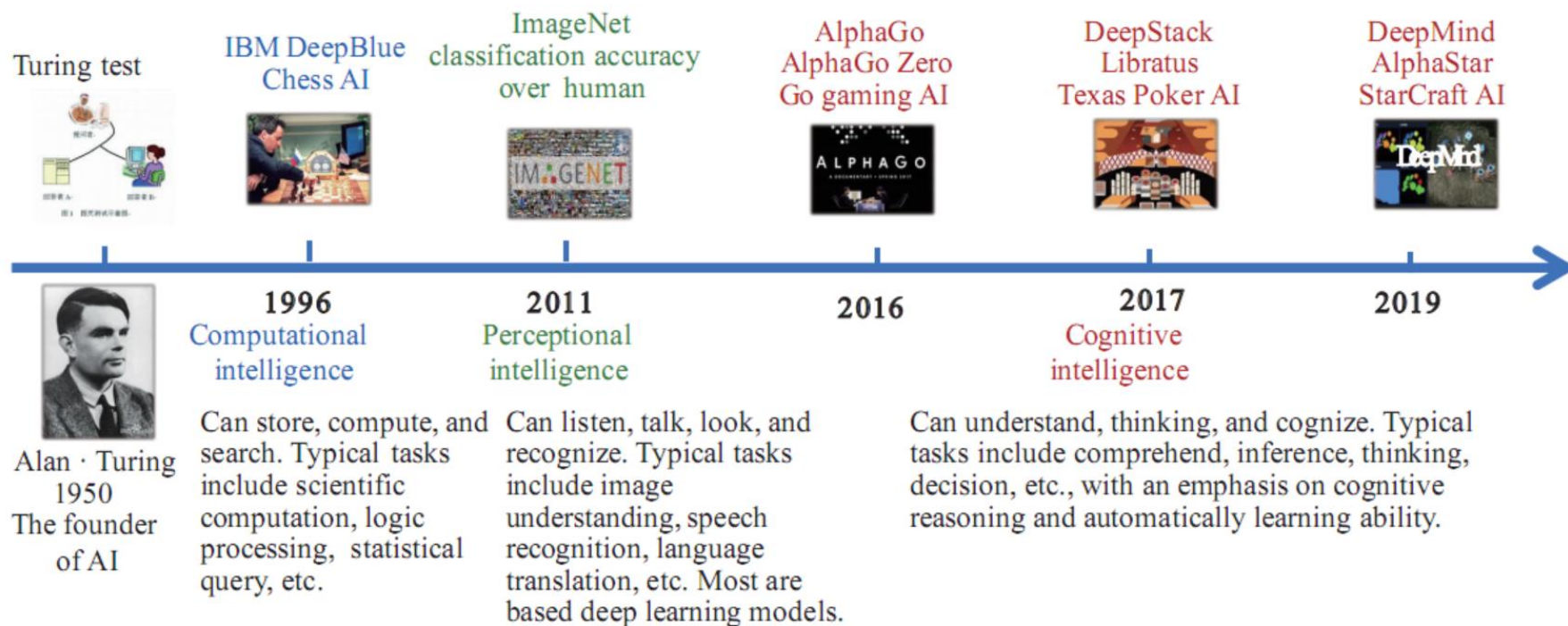
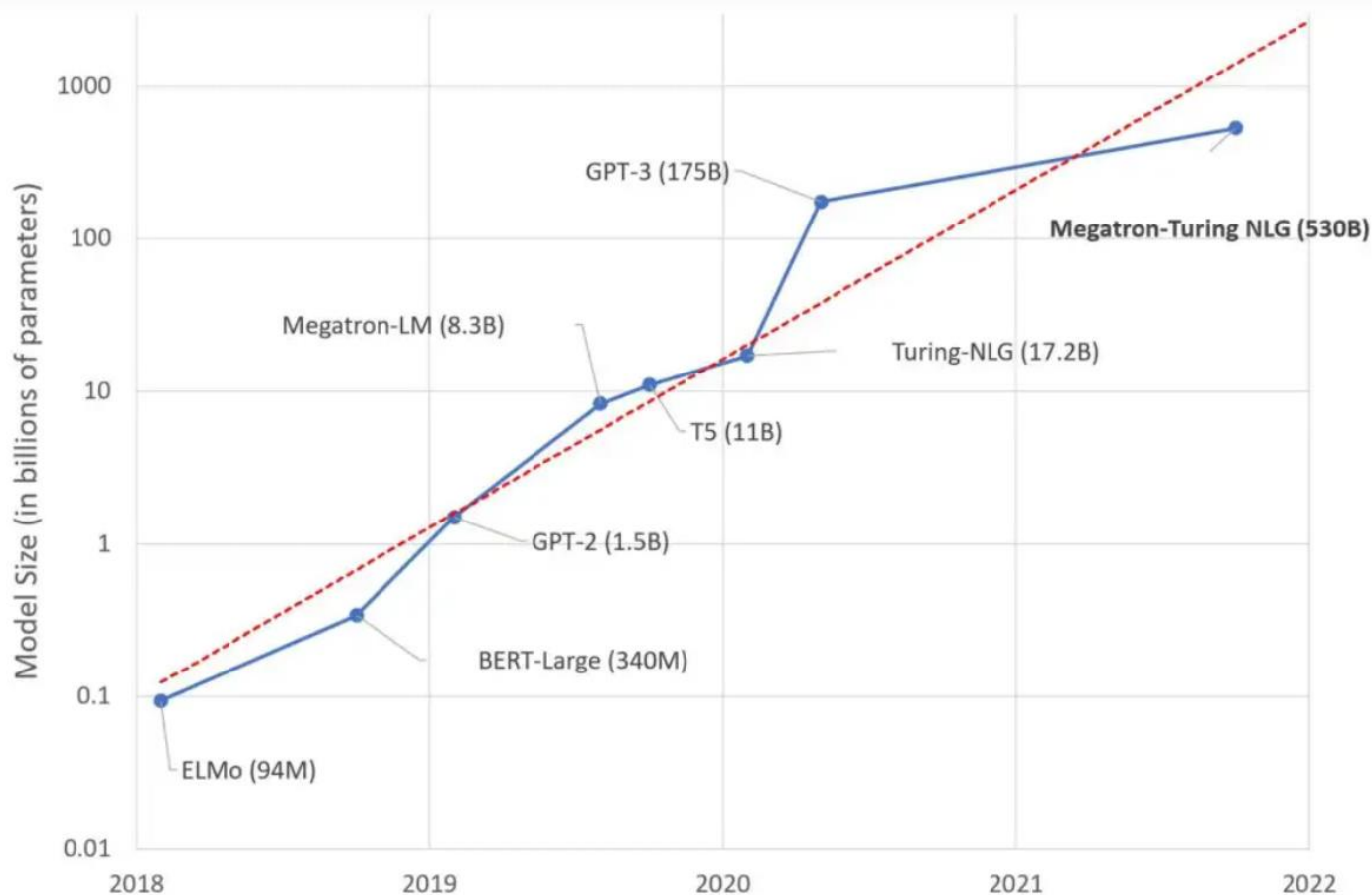


图 1 人机对抗的发展历史



应用现状



大脑的深度学习研究人员估计，人类大脑平均包含860亿个神经元和100万亿个突触。但不是所有的都用于语言。有趣的是，GPT-4预计将有大约100万亿个参数。



应用现状



千亿级参数通用大模型（文本编辑、编程、翻译、算术）