

1、机器学习的定义：

Arthur Samuel 将机器学习定义为让计算机在没有明确编程的情况下学习的研究领域；
那么什么叫机器学习中的“模型”：

机器学习模型其实是一种映射，可以将其看作是在给定输入情况（ x ）下，输出一定结果的函数 $f(x)$ ；

机器学习模型本质上是一种接受数据作为输入并生成输出的函数；

从而可以引申出机器学习的定义：通过输入海量训练数据对模型进行训练，使模型掌握数据蕴含的潜在规律，进而对新输入的数据进行更准确的分类或预测；

2.算法和模型的区别：

机器学习中的算法是在数据上运行以创造机器学习模型的过程；

模型是在训练数据上运行机器学习算法后保存下来的，表示用于预测所需的规则等

3.机器学习模型分类：

A：有监督学习：

使用带标签的数据（包含输入数据 x ，人工标注输出数据 y ），学习从 x 到 y 的映射关系；

主要应用于分类，回归，序列标注；

B：无监督学习：

仅使用无标签数据，学习数据自身的内在结构；

主要应用于聚类，降维（如 PCA），密度估计，生成模型；

C：半监督学习：

使用少量标签数据和大量的无标签数据，核心逻辑是当标签数据稀缺时，利用数据分布的连续性，让无标签数据辅助标签数据“拓展规律”；

主要应用于半监督分类，半监督回归；

D：强化学习：

通过与环境的动态交互，学习最优决策策略，无固定标签，但是有“延迟奖励信号”；智能体对环境的操作称为“动作”，环境对动作的反馈称为“奖励”而强化学习的核心目标是最大化“累积奖励”；

主要应用于游戏 AI，机器人控制，资源调度；

E：自监督学习：

看作无监督学习的“进阶版”，通过从数据自身构造“伪标签”实现监督信号的自生成；

主要应用于预训练语言模型，视觉表征学习等；

4.有监督学习模型的工作原理：

（1）：初始化：模型的“初始状态”设定

任何可学习模型均由固定结构和可调整参数构成；

结构：由人类设计的函数形式，定义了从输入到输出的计算逻辑；

参数：模型中可学习的变量，初始值通常随机确定

（2）：前向计算，从输入到预测的“正向传导”

当输入数据 x 进入模型后，模型通过固定结构函数生成预测值 \hat{y}

例如对于线性函数， $\hat{y} = \omega^T x + b$

（3）：量化预测与真实值的偏差：

模型生成 \hat{y} 后，需要通过损失函数计算与真实标签 y 的误差，如回归任务均方误差

$$(MSE) L = \frac{1}{2}(y - \hat{y})^2$$

（4）：基于量化偏差对参数进行更新

模型核心“学习”内容体现在用损失指导参数调整，用于减少未来预测的损失，这一过程通过优化算法实现：

第一步：梯度计算：通过链式法则计算损失函数对每个参数的偏导数，表示“参数微小变化对损失的影响幅度与方向”；

对于神经网络需要通过反向传播算法逐层计算梯度

第二步：参数迭代：依据梯度更新参数；

（5）：迭代收敛

前向计算-偏差量化-参数更新这一过程会在训练集上重复迭代，参数逐步调整，损失函数值逐步下降

（6）：预测应用

训练完成后，模型参数固定，进入“预测模型”，对新输入的 x 仅进行前向计算输出结果；

5.无监督学习模型的工作原理：

无监督学习模型的核心是在无标签数据中自主挖掘内在结构或统计规律，将数据自身的统计规律（相似度，协方差，概率分布）转化为可优化的目标函数。让模型通过优化该目标，自发捕捉数据的内部规律；

下以 K-means 模型工作原理为例：

1.初始化：随机选取 k 个数据点作为初始的“簇中心”， k 为预设的簇数量；

2.分配样本：计算每个样本点 x_i 与所有簇中心的距离，将 x_i 分配给最近的簇；

3.更新中心：对每个簇，计算该簇内所有点的均值作为新的簇中心，本质上是让新中心“代表”簇内数据的平均特征；

4.迭代收敛：重复 2-3，直到簇中心的变化小于阈值；

无监督学习模型的共通原理：

都是利用数据自身的统计特性作为优化目标，替代监督学习中的标签误差；

都是从数据中提取可重复的模式，让模型理解什么是数据的典型特征；

无监督模型也有“计算-误差-参数调整”迭代实现学习；最终使目标函数收敛到最优；

6.模型“学习”的本质：

模型学习的本质是在参数化函数簇中，通过优化算法利用数据的统计信号，引导参数从初始随机状态向匹配数据规律的最优配置迭代收敛；

7.什么是 AI：

人工智能是一个跨学科的研究领域，同时指代该领域所构建的技术系统；

人工智能是研究如何通过技术手段设计和构建具有类智能行为能力的系统

（1）AI 研究对象与目标：

AI 的研究对象是“智能行为的计算实现”而非对生物智能的复制，其目标是构建能执行传统上需要人类智能完成的任务的系统；

（2）AI “智能”的体现：

AI 的智能体现在对信息的符号/数值计算能力，其类智能行为通过以下机制实现：

符号操作：基于逻辑规则对抽象符号进行推理；

数据驱动的模式学习：通过统计方法从数据中提取规律，实现对未观测信息的预测；

环境交互与反馈优化：通过感知-动作循环与环境动态交互，以目标为导向调整策略；

8.深度学习与机器学习的区别：

深度学习是机器学习的一个子领域，通过构建和训练多层神经网络来解决复杂的模式识别和预测问题。深度学习常包含多个隐藏层，这些隐藏层能自动从原始数据中提取出高层次的特征，从而实现对复杂任务的高效处理；

那么区别在于：

（1）模型结构：

传统机器学习通常使用简单的模型结构，如线性回归，决策树等，这些模型通常只有少数几个隐藏层；

深度学习使用多层神经网络，层数可以到达几十甚至上百层，这种深层次的模型结构有利于使模型学到更复杂的特征表示；

（2）特征工程：

传统机器学习需要人工进行特征选择，即手动选择和输入特征；

深度学习能够自动从原始数据中学习和提取特征，减少对人工特征工程的依赖；

（3）数据需求：

传统机器学习通常需要较少的数据量，可以通过少量的标注数据进行有效的训练；

深度学习需要大量的标注数据来训练复杂的模型

（4）计算资源：

传统机器学习通常在普通的计算资源上就能运行；

深度学习需要强大的计算资源；

（5）应用场景：

传统机器学习适用于各种任务，聚类，回归等；

深度学习特别适用于处理高维，非结构化数据等，在计算机视觉，自然语言处理，语音识别等领域取得了显著的成果；