

# 斯坦福大学机器学习教程笔记

# 斯坦福大学机器学习教程中文笔记

## 课程概述

课程地址: <https://www.coursera.org/course/ml>

**Machine Learning**(机器学习)是研究计算机怎样模拟或实现人类的学习行为,以获取新的知识或技能,重新组织已有的知识结构使之不断改善自身的性能。它是人工智能的核心,是使计算机具有智能的根本途径,其应用遍及人工智能的各个领域,它主要使用归纳、综合而不是演绎。在过去的十年中,机器学习帮助我们自动驾驶汽车,有效的语音识别,有效的网络搜索,并极大地提高了人类基因组的认识。机器学习是当今非常普遍,你可能会使用这一天几十倍而不自知。很多研究者也认为这是最好的人工智能的取得方式。在本课中,您将学习最有效的机器学习技术,并获得实践,让它们为自己的工作。更重要的是,你会不仅得到理论基础的学习,而且获得那些需要快速和强大的应用技术解决问题的实用技术。最后,你会学到一些硅谷利用机器学习和人工智能的最佳实践创新。

本课程提供了一个广泛的介绍机器学习、数据挖掘、统计模式识别的课程。主题包括:

(一) 监督学习(参数/非参数算法,支持向量机,核函数,神经网络)。

(二) 无监督学习(聚类,降维,推荐系统,深度学习推荐)。

(三) 在机器学习的最佳实践(偏差/方差理论;在机器学习和人工智能创新过程)。本课程还将使用大量的案例研究,您还将学习如何运用学习算法构建智能机器人(感知,控制),文本的理解(Web 搜索,反垃圾邮件),计算机视觉,医疗信息,音频,数据挖掘,和其他领域。

## 目录

第 1 周	7
1、 引言(Introduction)	7
1.1 欢迎	7
1.2 机器学习是什么?	10
1.3 监督学习	12
1.4 无监督学习	16
2、 单变量线性回归(Linear Regression with One Variable)	错误!未定义书签。
2.1 模型表示	错误!未定义书签。
2.2 代价函数	错误!未定义书签。
2.3 代价函数的直观理解 I	错误!未定义书签。
2.4 代价函数的直观理解 II	错误!未定义书签。
2.5 梯度下降	错误!未定义书签。
2.6 梯度下降的直观理解	错误!未定义书签。
2.7 梯度下降的线性回归	错误!未定义书签。
2.8 接下来的内容	错误!未定义书签。
3、 线性代数回顾(Linear Algebra Review)	错误!未定义书签。
3.1 矩阵和向量	错误!未定义书签。
3.2 加法和标量乘法	错误!未定义书签。
3.3 矩阵向量乘法	错误!未定义书签。
3.4 矩阵乘法	错误!未定义书签。
3.5 矩阵乘法的性质	错误!未定义书签。
3.6 逆、转置	错误!未定义书签。
第 2 周	错误!未定义书签。
4、 多变量线性回归(Linear Regression with Multiple Variables)	错误!未定义书签。
4.1 多维特征	错误!未定义书签。
4.2 多变量梯度下降	错误!未定义书签。
4.3 梯度下降法实践 1-特征缩放	错误!未定义书签。
4.4 梯度下降法实践 2-学习率	错误!未定义书签。
4.5 特征和多项式回归	错误!未定义书签。
4.6 正规方程	错误!未定义书签。
4.7 正规方程及不可逆性(选修)	错误!未定义书签。
5、 Octave 教程(Octave Tutorial)	错误!未定义书签。
5.1 基本操作	错误!未定义书签。
5.2 移动数据	错误!未定义书签。
5.3 计算数据	错误!未定义书签。
5.4 绘图数据	错误!未定义书签。
5.5 控制语句: for, while, if 语句	错误!未定义书签。
5.6 向量化	错误!未定义书签。
5.7 工作和提交的编程练习	错误!未定义书签。
第 3 周	错误!未定义书签。
6、 逻辑回归(Logistic Regression)	错误!未定义书签。
6.1 分类问题	错误!未定义书签。
6.2 假说表示	错误!未定义书签。

6.3 判定边界 .....	错误!未定义书签。
6.4 代价函数 .....	错误!未定义书签。
6.5 简化的成本函数和梯度下降 .....	错误!未定义书签。
6.6 高级优化 .....	错误!未定义书签。
6.7 多类别分类：一对多 .....	错误!未定义书签。
7、 正则化(Regularization) .....	错误!未定义书签。
7.1 过拟合的问题 .....	错误!未定义书签。
7.2 代价函数 .....	错误!未定义书签。
7.3 正则化线性回归 .....	错误!未定义书签。
7.4 正则化的逻辑回归模型 .....	错误!未定义书签。
第 4 周 .....	错误!未定义书签。
8、 神经网络：表述(Neural Networks: Representation) .....	错误!未定义书签。
8.1 非线性假设 .....	错误!未定义书签。
8.2 神经元和大脑 .....	错误!未定义书签。
8.3 模型表示 1.....	错误!未定义书签。
8.4 模型表示 2.....	错误!未定义书签。
8.5 特征和直观理解 1.....	错误!未定义书签。
8.6 样本和直观理解 II.....	错误!未定义书签。
8.7 多类分类 .....	错误!未定义书签。
第 5 周 .....	错误!未定义书签。
9、 神经网络的学习(Neural Networks: Learning) .....	错误!未定义书签。
9.1 代价函数 .....	错误!未定义书签。
9.2 反向传播算法 .....	错误!未定义书签。
9.3 反向传播算法的直观理解 .....	错误!未定义书签。
9.4 实现注意：展开参数 .....	错误!未定义书签。
9.5 梯度检验 .....	错误!未定义书签。
9.6 随机初始化 .....	错误!未定义书签。
9.7 综合起来 .....	错误!未定义书签。
9.8 自动驾驶 .....	错误!未定义书签。
第 6 周 .....	错误!未定义书签。
10、 应用机器学习的建议(Advice for Applying Machine Learning) .....	错误!未定义书签。
10.1 决定下一步做什么 .....	错误!未定义书签。
10.2 评估一个假设 .....	错误!未定义书签。
10.3 模型选择和交叉验证集 .....	错误!未定义书签。
10.4 诊断偏差和方差 .....	错误!未定义书签。
10.5 正则化和偏差/方差 .....	错误!未定义书签。
10.6 学习曲线 .....	错误!未定义书签。
10.7 决定下一步做什么 .....	错误!未定义书签。
11、 机器学习系统的设计(Machine Learning System Design) .....	错误!未定义书签。
11.1 首先要做什么 .....	错误!未定义书签。
11.2 误差分析 .....	错误!未定义书签。
11.3 类偏斜的误差度量 .....	错误!未定义书签。
11.4 查准率和查全率之间的权衡 .....	错误!未定义书签。
11.5 机器学习的数据 .....	错误!未定义书签。

第 7 周 .....	错误!未定义书签。
12、 支持向量机(Support Vector Machines) .....	错误!未定义书签。
12.1 优化目标 .....	错误!未定义书签。
12.2 大边界的直观理解 .....	错误!未定义书签。
12.3 大边界分类背后的数学（选修） .....	错误!未定义书签。
12.4 核函数 1.....	错误!未定义书签。
12.5 核函数 2.....	错误!未定义书签。
12.6 使用支持向量机 .....	错误!未定义书签。
第 8 周 .....	错误!未定义书签。
13、 聚类(Clustering).....	错误!未定义书签。
13.1 无监督学习：简介 .....	错误!未定义书签。
13.2 K-均值算法 .....	错误!未定义书签。
13.3 优化目标 .....	错误!未定义书签。
13.4 随机初始化 .....	错误!未定义书签。
13.5 选择聚类数 .....	错误!未定义书签。
14、 降维(Dimensionality Reduction).....	错误!未定义书签。
14.1 动机一：数据压缩 .....	错误!未定义书签。
14.2 动机二：数据可视化 .....	错误!未定义书签。
14.3 主成分分析问题 .....	错误!未定义书签。
14.4 主成分分析算法 .....	错误!未定义书签。
14.5 选择主成分的数量 .....	错误!未定义书签。
14.6 重建的压缩表示 .....	错误!未定义书签。
14.7 主成分分析法的应用建议 .....	错误!未定义书签。
第 9 周 .....	错误!未定义书签。
15、 异常检测(Anomaly Detection) .....	错误!未定义书签。
15.1 问题的动机 .....	错误!未定义书签。
15.2 高斯分布 .....	错误!未定义书签。
15.3 算法 .....	错误!未定义书签。
15.4 开发和评价一个异常检测系统 .....	错误!未定义书签。
15.5 异常检测与监督学习对比 .....	错误!未定义书签。
15.6 选择特征 .....	错误!未定义书签。
15.7 多元高斯分布（选修） .....	错误!未定义书签。
15.8 使用多元高斯分布进行异常检测（选修） .....	错误!未定义书签。
16、 推荐系统(Recommender Systems) .....	错误!未定义书签。
16.1 问题形式化 .....	错误!未定义书签。
16.2 基于内容的推荐系统 .....	错误!未定义书签。
16.3 协同过滤 .....	错误!未定义书签。
16.4 协同过滤算法 .....	错误!未定义书签。
16.5 向量化：低秩矩阵分解 .....	错误!未定义书签。
16.6 推行工作上的细节：均值归一化 .....	错误!未定义书签。
第 10 周 .....	错误!未定义书签。
17、 大规模机器学习(Large Scale Machine Learning).....	错误!未定义书签。
17.1 大型数据集的学习 .....	错误!未定义书签。
17.2 随机梯度下降法 .....	错误!未定义书签。

17.3	小批量梯度下降 .....	错误!未定义书签。
17.4	随机梯度下降收敛 .....	错误!未定义书签。
17.5	在线学习 .....	错误!未定义书签。
17.6	映射化简和数据并行 .....	错误!未定义书签。
18、	应用实例：图片文字识别(Application Example: Photo OCR) .....	错误!未定义书签。
18.1	问题描述和流程图 .....	错误!未定义书签。
18.2	滑动窗口 .....	错误!未定义书签。
18.3	获取大量数据和人工数据 .....	错误!未定义书签。
18.4	上限分析：哪部分管道的接下去做 .....	错误!未定义书签。
19、	总结(Conclusion) .....	错误!未定义书签。
19.1	总结和致谢 .....	错误!未定义书签。
附件	.....	错误!未定义书签。
	机器学习的数学基础 .....	错误!未定义书签。
	高等数学 .....	错误!未定义书签。
	线性代数 .....	错误!未定义书签。
	概率论和数理统计 .....	错误!未定义书签。

# 第 1 周

## 1、 引言(Introduction)

### 1.1 欢迎

参考视频: 1 - 1 - Welcome (7 min).mkv

第一个视频主要讲了什么是机器学习，机器学习能做些什么事情。

机器学习是目前信息技术中最激动人心的方向之一。在这门课中，你将学习到这门技术的前沿，并可以自己实现学习机器学习的算法。

你或许每天都在不知不觉中使用了机器学习的算法每次，你打开谷歌、必应搜索到你需要的内容，正是因为他们有良好的学习算法。谷歌和微软实现了学习算法来排行网页每次，你用 **Facebook** 或苹果的图片分类程序他能认出你朋友的照片，这也是机器学习。每次您阅读您的电子邮件垃圾邮件筛选器，可以帮你过滤大量的垃圾邮件这也是一种学习算法。对我来说，我感到激动的原因之一是有一天做出一个和人类一样聪明的机器。实现这个想法任重而道远，许多 **AI** 研究者认为，实现这个目标最好的方法是通过让机器试着模仿人的大脑学习我会在这门课中介绍一点这方面的内容。

在这门课中，你还讲学习到关于机器学习的前沿状况。但事实上只了解算法、数学并不能解决你关心的实际的问题。所以，我们将花大量的时间做练习，从而你自己能实现每个这些算法，从而了解内部机理。

那么，为什么机器学习如此受欢迎呢？原因是，机器学习不只是用于人工智能领域。

我们创造智能的机器，有很多基础的知识。比如，我们可以让机器找到 **A** 与 **B** 之间的最短路径，但我们仍然不知道如何让机器做更有趣的事情，如 **web** 搜索、照片标记、反垃圾邮件。我们发现，唯一方法是让机器自己学习怎么来解决问题。所以，机器学习已经成为计算机的一个能力。

现在它涉及到各个行业和基础科学中。我从事于机器学习，但我每个星期都跟直升机飞行员、生物学家、很多计算机系统程序员交流（我在斯坦福大学的同事同时也是这样）和平均每个星期会从硅谷收到两、三个电子邮件，这些联系我的人都对将学习算法应用于他们自

己的问题感兴趣。这表明机器学习涉及的问题非常广泛。有机器人、计算生物学、硅谷中大量的问题都收到机器学习的影响。

这里有一些机器学习的案例。比如说，数据库挖掘。机器学习被用于数据挖掘的原因之一是网络和自动化技术的增长，这意味着，我们有史上最大的数据集比如说，大量的硅谷公司正在收集 **web** 上的单击数据，也称为点击流数据，并尝试使用机器学习算法来分析数据，更好的了解用户，并为用户提供更好的服务。这在硅谷有巨大的市场。再比如，医疗记录。随着自动化的出现，我们现在有了电子医疗记录。如果我们可以把医疗记录变成医学知识，我们就可以更好地理解疾病。再如，计算生物学。还是因为自动化技术，生物学家们收集的大量基因数据序列、**DNA** 序列和等等，机器运行算法让我们更好地了解人类基因组，大家都知道这对人类意味着什么。再比如，工程方面，在工程的所有领域，我们有越来越大、越来越大的数据集，我们试图使用学习算法，来理解这些数据。另外，在机械应用中，有些人不能直接操作。例如，我已经在无人直升机领域工作了许多年。我们不知道如何写一段程序让直升机自己飞。我们唯一能做的就是让计算机自己学习如何驾驶直升机。

手写识别：现在我们能够非常便宜地把信寄到这个美国甚至全世界的原因之一就是当你写一个像这样的信封，一种学习算法已经学会如何读你信封，它可以自动选择路径，所以我们只需要花几个美分把这封信寄到数千英里外。

事实上，如果你看过自然语言处理或计算机视觉，这些语言理解或图像理解都是属于 **AI** 领域。大部分的自然语言处理和大部分的计算机视觉，都应用了机器学习。学习算法还广泛用于自定制程序。每次你去亚马逊或 **Netflix** 或 **iTunes Genius**，它都会给出其他电影或产品或音乐的建议，这是一种学习算法。仔细想一想，他们有百万的用户；但他们没有办法为百万用户，编写百万个不同程序。软件能给这些自定制的唯一的建议的唯一方法是通过学习你的行为，来为你定制服务。

最后学习算法被用来理解人类的学习和了解大脑。

我们将谈论如何用这些推进我们的 **AI** 梦想。几个月前，一名学生给我一篇文章关于最顶尖的 12 个 **IT** 技能。拥有了这些技能 **HR** 绝对不会拒绝你。这是稍显陈旧的文章，但在这个列表最顶部就是机器学习的技能。

在斯坦福大学，招聘人员联系我，让我推荐机器学习学生毕业的人远远多于机器学习的毕业生。所以我认为需求远远没有被满足现在学习“机器学习”非常好，在这门课中，我希望能告诉你们很多机器学习的知识。

在接下来的视频中，我们将开始给更正式的定义，什么是机器学习。然后我们会开始学



习机器学习的主要问题和算法你会了解一些主要的机器学习的术语,并开始了解不同的算法,用哪种算法更合适。

## 1.2 机器学习是什么？

参考视频: 1 - 2 - What is Machine Learning\_ (7 min).mkv

机器学习是什么？在本视频中，我们会尝试着进行定义，同时让你懂得何时会使用机器学习。实际上，即使是在机器学习的专业人士中，也不存在一个被广泛认可的定义来准确定义机器学习是什么或不是什么，现在我将告诉你一些人们尝试定义的示例。第一个机器学习的定义来自于 **Arthur Samuel**。他定义机器学习为，在进行特定编程的情况下，给予计算机学习能力的领域。**Samuel** 的定义可以回溯到 50 年代，他编写了一个西洋棋程序。这程序神奇之处在于，编程者自己并不是个下棋高手。但因为他太菜了，于是就通过编程，让西洋棋程序自己跟自己下了上万盘棋。通过观察哪种布局（棋盘位置）会赢，哪种布局会输，久而久之，这西洋棋程序明白了什么是好的布局，什么样是坏的布局。然后就牛逼大发了，程序通过学习后，玩西洋棋的水平超过了 **Samuel**。这绝对是令人瞩目的成果。

尽管编写者自己是个菜鸟，但因为计算机有着足够的耐心，去下上万盘的棋，没有人有这耐心去下这么多盘棋。通过这些练习，计算机获得无比丰富的经验，于是渐渐成为了比 **Samuel** 更厉害的西洋棋手。上述是个有点不正式的定义，也比较古老。另一个年代近一点的定义，由 **Tom Mitchell** 提出，来自卡内基梅隆大学，**Tom** 定义的机器学习是，一个好的学习问题定义如下，他说，一个程序被认为能从经验 **E** 中学习，解决任务 **T**，达到性能度量值 **P**，当且仅当，有了经验 **E** 后，经过 **P** 评判，程序在处理 **T** 时的性能有所提升。我认为经验 **E** 就是程序上万次的自我练习的经验而任务 **T** 就是下棋。性能度量值 **P** 呢，就是它在与一些新的对手比赛时，赢得比赛的概率。

在这些视频中，除了我教你的内容以外，我偶尔会问你一个问题，确保你对内容有所理解。说曹操，曹操到，顶部是 **Tom Mitchell** 的机器学习的定义，我们假设您的电子邮件程序会观察收到的邮件是否被你标记为垃圾邮件。在这种 **Email** 客户端中，你点击“垃圾邮件”按钮，报告某些 **Email** 为垃圾邮件，不会影响别的邮件。基于被标记为垃圾的邮件，您的电子邮件程序能更好地学习如何过滤垃圾邮件。请问，在这个设定中，任务 **T** 是什么？几秒钟后，该视频将暂停。当它暂停时，您可以使用鼠标，选择这四个单选按钮中的一个，让我知道这四个，你所认为正确的选项。它可能是性能度量值 **P**。所以，以性能度量值 **P** 为标准，这个任务的性能，也就是这个任务 **T** 的系统性能，将在学习经验 **E** 后得到提高。

本课中，我希望教你有关各种不同类型的学习算法。目前存在几种不同类型的学习算法。

主要的两种类型被我们称之为监督学习和无监督学习。在接下来的几个视频中，我会给出这些术语的定义。这里简单说两句，监督学习这个想法是指，我们将教计算机如何去完成任务，而在无监督学习中，我们打算让它自己进行学习。如果对这两个术语仍一头雾水，请不要担心，在后面的两个视频中，我会具体介绍这两种学习算法。此外你将听到诸如，强化学习和推荐系统等各种术语。这些都是机器学习算法的一员，以后我们都将介绍到，但学习算法最常用两个类型就是监督学习、无监督学习。我会在接下来的两个视频中给出它们的定义。本课中，我们将花费最多的精力来讨论这两种学习算法。而另一个会花费大量时间的任务是了解应用学习算法的实用建议。

我非常注重这部分内容，实际上，就这些内容而言我不知道还有哪所大学会介绍到。给你讲授学习算法就好像给你一套工具，相比于提供工具，可能更重要的，是教你如何使用这些工具。我喜欢把这比喻成学习当木匠。想象一下，某人教你如何成为一名木匠，说这是锤子，这是螺丝刀，锯子，祝你好运，再见。这种教法不好，不是吗？你拥有这些工具，但更重要的是，你要学会如何恰当地使用这些工具。会用与不会用的人之间，存在着鸿沟。尤其是知道如何使用这些机器学习算法的，与那些不知道如何使用的人。在硅谷我住的地方，当我走访不同的公司，即使是最顶尖的公司，很多时候我都看到人们试图将机器学习算法应用于某些问题。有时他们甚至已经为此花了六个月之久。但当我看着他们所忙碌的事情时，我想说，哎呀，我本来可以在六个月前就告诉他们，他们应该采取一种学习算法，稍加修改进行使用，然后成功的机会绝对会高得多所以在本课中，我们要花很多时间来探讨，如果你真的试图开发机器学习系统，探讨如何做出最好的实践类型决策，才能决定你的方式来构建你的系统，这样做的话，当你运用学习算法时，就不太容易变成那些为寻找一个解决方案花费6个月之久的人们中的一员。他们可能已经有了大体的框架，只是没法正确的工作于是这就浪费了六个月的时间。所以我会花很多时间来教你这些机器学习、人工智能的最佳实践以及如何让它们工作，我们该如何去做，硅谷和世界各地最优秀的人是怎样做的。我希望能帮你成为最优秀的人才，通过了解如何设计和构建机器学习和人工智能系统。

这就是机器学习，这些都是我希望讲授的主题。在下一个视频里，我会定义什么是监督学习，什么是无监督学习。此外，探讨何时使用二者。

## 1.3 监督学习

参考视频: 1 - 3 - Supervised Learning (12 min).mkv

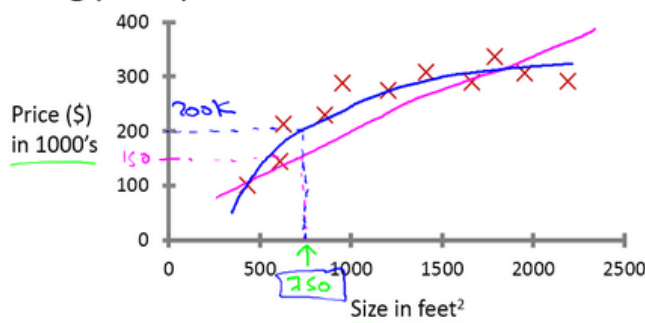
在这段视频中，我要定义可能是最常见一种机器学习问题：那就是监督学习。我将在后面正式定义监督学习。

我们用一个例子介绍什么是监督学习把正式的定义放在后面介绍。假如说你想预测房价。

前阵子，一个学生从波特兰俄勒冈州的研究所收集了一些房价的数据。你把这些数据画出来，看起来是这个样子：横轴表示房子的面积，单位是平方英尺，纵轴表示房价，单位是千美元。那基于这组数据，假如你有一个朋友，他有一套 750 平方英尺房子，现在他希望把房子卖掉，他想知道这房子能卖多少钱。

那么关于这个问题，机器学习算法将会怎么帮助你呢？

Housing price prediction.



Supervised Learning  
"right answers" given

Regression: Predict continuous  
valued output (price)

我们应用学习算法，可以在这组数据中画一条直线，或者换句话说，拟合一条直线，根据这条线我们可以推测出，这套房子可能卖\$150,000，当然这不是唯一的算法。可能还有更好的，比如我们不用直线拟合这些数据，用二次方程去拟合可能效果会更好。根据二次方程的曲线，我们可以从这个点推测出，这套房子能卖接近\$200,000。稍后我们将讨论如何选择学习算法，如何决定用直线还是二次方程来拟合。两个方案中有一个能让你朋友的房子出售得更合理。这些都是学习算法里面很好的例子。以上就是监督学习的例子。

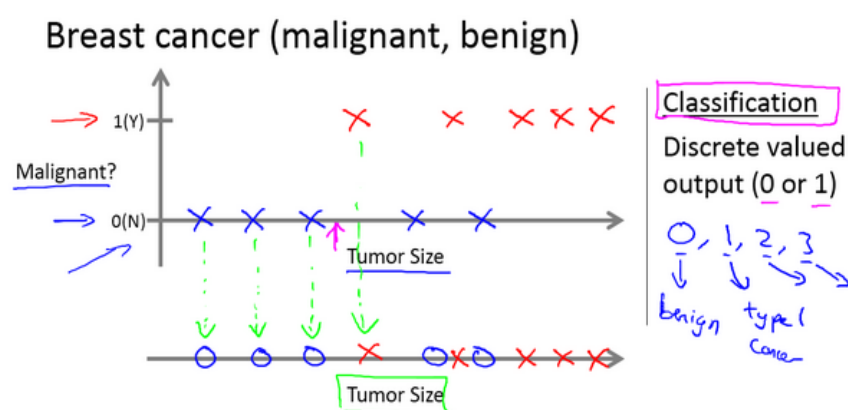
可以看出，监督学习指的就是我们给学习算法一个数据集。这个数据集由“正确答案”组成。在房价的例子中，我们给了一系列房子的数据，我们给定数据集中每个样本的正确价格，即它们实际的售价然后运用学习算法，算出更多的正确答案。比如你朋友那个新房子的

价格。用术语来讲，这叫做回归问题。我们试着推测出一个连续值的结果，即房子的价格。

一般房子的价格会记到美分，所以房价实际上是一系列离散的值，但是我们通常又把房价看成实数，看成是标量，所以又把它看成一个连续的数值。

回归这个词的意思是，我们在试着推测出这一系列连续值属性。

我再举另外一个监督学习的例子。我和一些朋友之前研究过这个。假设说你想通过查看病历来推测乳腺癌良性与否，假如有人检测出乳腺肿瘤，恶性肿瘤有害并且十分危险，而良性的肿瘤危害就没那么大，所以人们显然会很在意这个问题。



让我们来看一组数据：这个数据集中，横轴表示肿瘤的大小，纵轴上，我标出 1 和 0 表示是或者不是恶性肿瘤。我们之前见过的肿瘤，如果是恶性则记为 1，不是恶性，或者说良性记为 0。

我有 5 个良性肿瘤样本，在 1 的位置有 5 个恶性肿瘤样本。现在我们有一个朋友很不幸检查出乳腺肿瘤。假设说她的肿瘤大概这么大，那么机器学习的问题就在于，你能否估算出肿瘤是恶性的或是良性的概率。用术语来讲，这是一个分类问题。

分类指的是，我们试着推测出离散的输出值：0 或 1 良性或恶性，而事实上在分类问题中，输出可能不止两个值。比如说可能有三种乳腺癌，所以你希望预测离散输出 0、1、2、3。0 代表良性，1 表示第 1 类乳腺癌，2 表示第 2 类癌症，3 表示第 3 类，但这也是分类问题。

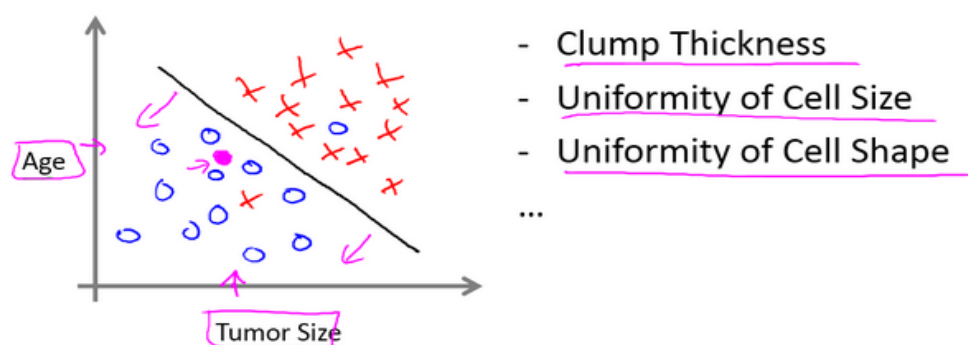
因为这几个离散的输出分别对应良性，第一类第二类或者第三类癌症，在分类问题中我们可以用另一种方式绘制这些数据点。

现在我用不同的符号来表示这些数据。既然我们把肿瘤的尺寸看做区分恶性或良性的特征，那么我可以这么画，我用不同的符号来表示良性和恶性肿瘤。或者说是负样本和正样本现在我们不全部画 x，良性的肿瘤改成用 o 表示，恶性的继续用 x 表示。来预测肿瘤的恶

性与否。

在其它一些机器学习问题中，可能会遇到不止一种特征。举个例子，我们不仅知道肿瘤的尺寸，还知道对应患者的年龄。在其他机器学习问题中，我们通常有更多的特征，我朋友研究这个问题时，通常采用这些特征，比如肿块密度，肿瘤细胞尺寸的一致性和形状的一致性等等，还有一些其他的特征。这就是我们即将学到最有趣的学习算法之一。

那种算法不仅能处理 2 种 3 种或 5 种特征，即使有无限多种特征都可以处理。



上图中，我列举了总共 5 种不同的特征，坐标轴上的两种和右边的 3 种，但是在一些学习问题中，你希望不只用 3 种或 5 种特征。相反，你想用无限多种特征，好让你的算法可以利用大量的特征，或者说线索来做推测。那你怎么处理无限多个特征，甚至怎么存储这些特征都存在问题，你电脑的内存肯定不够用。我们以后会讲一个算法，叫支持向量机，里面有一个巧妙的数学技巧，能让计算机处理无限多个特征。想象一下，我没有写下这两种和右边的三种特征，而是在一个无限长的列表里面，一直写一直写不停的写，写下无限多个特征，事实上，我们能用算法来处理它们。

现在来回顾一下，这节课我们介绍了监督学习。其基本思想是，我们数据集中的每个样本都有相应的“正确答案”。再根据这些样本作出预测，就像房子和肿瘤的例子中做的那样。我们还介绍了回归问题，即通过回归来推出一个连续的输出，之后我们介绍了分类问题，其目标是推出一组离散的结果。

现在来个小测验：假设你经营着一家公司，你想开发学习算法来处理这两个问题：

1. 你有一大批同样的货物，想象一下，你有上千件一模一样的货物等待出售，这时你想预测接下来的三个月能卖多少件？
2. 你有许多客户，这时你想写一个软件来检验每一个用户的账户。对于每一个账户，你要判断它们是否曾经被盗过？

那这两个问题，它们属于分类问题、还是回归问题？

问题一是一个回归问题，因为你知道，如果我有数千件货物，我会把它看成一个实数，

一个连续的值。因此卖出的物品数，也是一个连续的值。

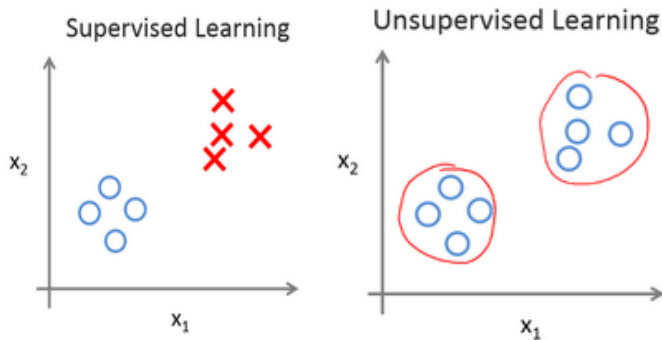
问题二是一个分类问题，因为我会把预测的值，用 0 来表示账户未被盗，用 1 表示账户曾经被盗过。所以我们根据账号是否被盗过，把它们定为 0 或 1，然后用算法推测一个账号是 0 还是 1，因为只有少数的离散值，所以我把它归为分类问题。

以上就是监督学习的内容。

## 1.4 无监督学习

参考视频: 1 - 4 - Unsupervised Learning (14 min).mkv

本次视频中，我们将介绍第二种主要的机器学习问题。叫做无监督学习。



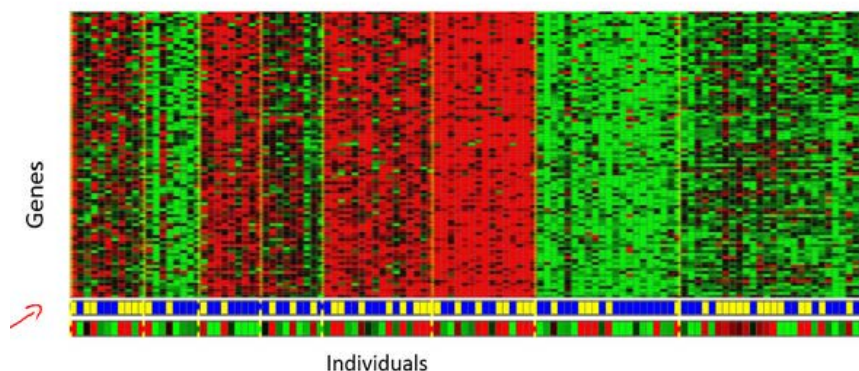
上个视频中，已经介绍了监督学习。回想当时的数据集，如图表所示，这个数据集中每条数据都已经标明是阴性或阳性，即是良性或恶性肿瘤。所以，对于监督学习里的每条数据，我们已经清楚地知道，训练集对应的正确答案，是良性或恶性了。

在无监督学习中，我们已知的数据。看上去有点不一样，不同于监督学习的数据的样子，即无监督学习中没有任何的标签或者是有相同的标签或者就是没标签。所以我们已知数据集，却不知如何处理，也未告知每个数据点是什么。别的都不知道，就是一个数据集。你能从数据中找到某种结构吗？针对数据集，无监督学习就能判断出数据有两个不同的聚集簇。这是一个，那是另一个，二者不同。是的，无监督学习算法可能会把这些数据分成两个不同的簇。所以叫做聚类算法。事实证明，它被用在很多地方。

聚类应用的一个例子就是在谷歌新闻中。如果你以前从来没见过它，你可以到这个 URL 网址 [news.google.com](https://news.google.com) 去看看。谷歌新闻每天都在，收集非常多，非常多的网络的新闻内容。它再将这些新闻分组，组成有关联的新闻。所以谷歌新闻做的就是搜索非常多的新闻事件，自动地把它们聚类到一起。所以，这些新闻事件全是同一主题的，所以显示到一起。

事实证明，聚类算法和无监督学习算法同样还用在很多其它的问题上。





其中就有基因学的理解应用。一个 **DNA** 微观数据的例子。基本思想是输入一组不同个体，对其中的每个个体，你要分析出它们是否有一个特定的基因。技术上，你要分析多少特定基因已经表达。所以这些颜色，红，绿，灰等等颜色，这些颜色展示了相应的程度，即不同的个体是否有着一个特定的基因。你能做的就是运行一个聚类算法，把个体聚类到不同的类或不同类型的组（人）.....

所以这个就是无监督学习，因为我们没有提前告知算法一些信息，比如，这是第一类的人，那些是第二类的人，还有第三类，等等。我们只是说，是的，这是有一堆数据。我不知道数据里面有什么。我不知道谁是什么类型。我甚至不知道人们有哪些不同的类型，这些类型又是什么。但你能自动地找到数据中的结构吗？就是说你要自动地聚类那些个体到各个类，我没法提前知道哪些是哪些。因为我们没有给算法正确答案来回应数据集中的数据，所以这就是无监督学习。

无监督学习或聚集有着大量的应用。它用于组织大型计算机集群。我有些朋友在大数据中心工作，那里有大型的计算机集群，他们想解决什么样的机器易于协同地工作，如果你能够让那些机器协同工作，你就能让你的数据中心工作得更高效。第二种应用就是社交网络的分析。所以已知你朋友的信息，比如你经常发 **email** 的，或是你 **Facebook** 的朋友、**谷歌+** 圈子的朋友，我们能否自动地给出朋友的分组呢？即每组里的人们彼此都熟识，认识组里的所有人？还有市场分割。许多公司有大型的数据库，存储消费者信息。所以，你能检索这些顾客数据集，自动地发现市场分类，并自动地把顾客划分到不同的细分市场中，你才能自动并更有效地销售或不同的细分市场一起进行销售。这也是无监督学习，因为我们拥有所有的顾客数据，但我们没有提前知道是什么的细分市场，以及分别有哪些我们数据集中的顾客。我们不知道谁是在一号细分市场，谁在二号市场，等等。那我们就必须让算法从数据中发现这一切。最后，无监督学习也可用于天文数据分析，这些聚类算法给出了令人惊讶、有趣、有用的理论，解释了星系是如何诞生的。这些都是聚类的例子，聚类只是无监督学习中的一种。

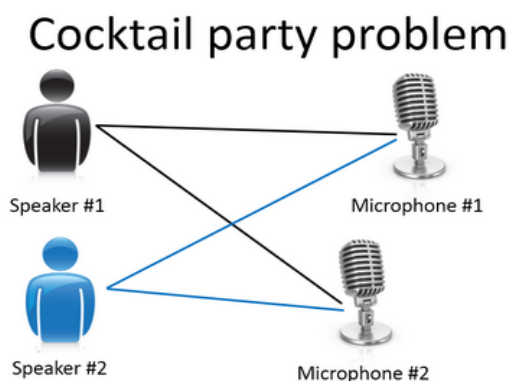
我现在告诉你们另一种。我先来介绍鸡尾酒宴问题。嗯，你参加过鸡尾酒宴吧？你可以想像下，有个宴会房间里满是人，全部坐着，都在聊天，这么多人同时在聊天，声音彼此重叠，因为每个人都在说话，同一时间都在说话，你几乎听不到你面前那人的声音。所以，可能在一个这样的鸡尾酒宴中的两个人，他俩同时都在说话，假设现在是在个有些小的鸡尾酒宴中。我们放两个麦克风在房间中，因为这些麦克风在两个地方，离说话人的距离不同每个麦克风记录下不同的声音，虽然是同样的两个说话人。听起来像是两份录音被叠加到一起，或是被归结到一起，产生了我们现在的这些录音。另外，这个算法还会区分出两个音频资源，这两个可以合成或合并成之前的录音，实际上，鸡尾酒算法的第一个输出结果是：

1, 2, 3, 4, 5, 6, 7, 8, 9, 10,

所以，已经把英语的声音从录音中分离出来了。

第二个输出是这样：

1, 2, 3, 4, 5, 6, 7, 8, 9, 10。



看看这个无监督学习算法，实现这个得要多么的复杂，是吧？它似乎是这样，为了构建这个应用，完成这个音频处理似乎需要你去写大量的代码或链接到一堆的合成器 **JAVA** 库，处理音频的库，看上去绝对是个复杂的程序，去完成这个从音频中分离出音频。事实上，这个算法对应你刚才知道的那个问题的算法可以就用一行代码来完成。

就是这里展示的代码：

```
[W,s,v] = svd((repmat(sum(x.*x,1),size(x,1),1).*x)*x');
```

研究人员花费了大量时间才最终实现这行代码。我不是说这个简单的问题，但它证明了，当你使用正确的编程环境，许多学习算法是相当短的程序。所以，这也是为什么在本课中，我们打算使用 **Octave** 编程环境。**Octave** 是免费的开源软件，使用一个像 **Octave** 或 **Matlab** 的工具，许多学习算法变得只有几行代码就可实现。

后面，我会教你们一点关于如何使用 **Octave** 的知识，你就可以用 **Octave** 来实现一些算法了。或者，如果你有 **Matlab**（盗版？），你也可以用 **Matlab**。事实上，在硅谷里，对大量机器学习算法，我们第一步就是建原型，在 **Octave** 建软件原型，因为软件在 **Octave** 中可以令人难以置信地、快速地实现这些学习算法。这里的这些函数比如 **SVM**（支持向量机）函数，**奇异值分解**，**Octave** 里已经建好了。如果你试图完成这个工作，但借助 **C++** 或 **JAVA** 的话，你会需要很多很多行的代码，并链接复杂的 **C++** 或 **Java** 库。所以，你可以实现这些算法，借助 **C++** 或 **Java** 或 **Python**，它只是用这些语言来实现会更加复杂。（编者注：这个是当时的情况，现在 **Python** 变主流了）

我已经见到，在我教机器学习将近十年后的现在，发现，学习可以更加高速，如果使用 **Octave** 作为编程环境，如果使用 **Octave** 作为学习工具，以及作为原型工具，它会让你对学习算法的学习和建原型快上许多。

事实上，许多人在大硅谷的公司里做的其实就是，使用一种工具像 **Octave** 来做第一步的学习算法的原型搭建，只有在你已经让它工作后，你才移植它到 **C++** 或 **Java** 或别的语言。事实证明，这样做通常可以让你的算法运行得比直接用 **C++** 实现更快，所以，我知道，作为一名指导者，我必须说“相信我”，但对你们中从未使用过 **Octave** 这种编程环境的人，我还是要告诉你们这一点一定要相信我，我想，对你们而言，我认为你们的时间，你们的开发时间是最有价值的资源。我已经见过很多人这样做了，我把你看作是机器学习研究员，或机器学习开发人员，想更加高产的话，你要学会使用这个原型工具，开始使用 **Octave**。

最后，总结下本视频内容，我有个简短的复习题给你们。

我们介绍了无监督学习，它是学习策略，交给算法大量的数据，并让算法为我们从数据中找出某种结构。

好的，希望你们还记得**垃圾邮件问题**。如果你有标记好的数据，区别好是垃圾还是非垃圾邮件，我们把这个当作**监督学习问题**。

**新闻事件分类**的例子，就是那个谷歌新闻的例子，我们在本视频中有见到了，我们看到，可以用一个聚类算法来聚类这些文章到一起，所以是**无监督学习**。

**细分市场**的例子，我在更早一点的时间讲过，你可以当作**无监督学习问题**，因为我只是拿到算法数据，再让算法去自动地发现细分市场。

最后一个例子，**糖尿病**，这个其实就像是我们的乳腺癌，上个视频里的。只是替换了好、坏肿瘤，良性、恶性肿瘤，我们改用糖尿病或没病。所以我们把这个当作**监督学习**，我们能够解决它，作为一个监督学习问题，就像我们在乳腺癌数据中做的一样。

好了，以上就是无监督学习的视频内容，在下一个视频中，我们将深入探究特定的学习算法，开始介绍这些算法是如何工作的，和我们还有你如何实现它们。