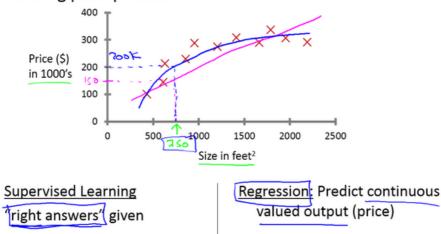
1.3 监督学习

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

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。

Housing price prediction.



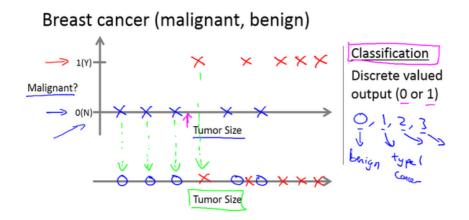
我们应用学习算法,可以在这组数据中画一条直线,或者换句话说,拟合一条直线,根据这条线我们可以推测出,这套房子可能卖\$150,000,当然这不是唯一的算法。可能还有更好的,比如我们不用直线拟合这些数据,用二次方程去拟合可能效果会更好。根据二次方程的曲线,我们可以从这个点推测出,这套房子能卖接近\$200,000。稍后我们将讨论如何选择学习算法,如何决定用直线还是二次方程来拟合。两个方案中有一个能让你朋友的房子出售得更合理。这些都是学习算法里面很好的例子。以上就是监督学习的例子。

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。



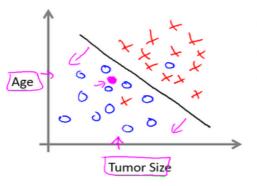
让我们来看一组数据:这个数据集中,横轴表示肿瘤的大小,纵轴上,我标出 1 和 0 表示是或者不是恶性肿瘤。我们之前见过的肿瘤,如果是恶性则记为 1,不是恶性,或者说良性记为 0。

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 \mathbf{x} ,良性的肿瘤改成用 \mathbf{o} 表示,恶性的继续用 \mathbf{x} 表示。来预测肿瘤的恶性与否。



- Clump Thickness
- Uniformity of Cell Size
- Uniformity of Cell Shape

•••

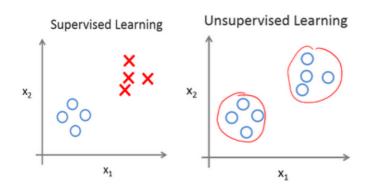
上图中,我列举了总共 5 种不同的特征,坐标轴上的两种和右边的 3 种,但是在一些学习问题中,你希望不只用 3 种或 5 种特征。相反,你想用无限多种特征,好让你的算法可以利用大量的特征,或者说线索来做推测。那你怎么处理无限多个特征,甚至怎么存储这些特征都存在问题,你电脑的内存肯定不够用。我们以后会讲一个算法,叫支持向量机,里面有一个巧妙的数学技巧,能让计算机处理无限多个特征。想象一下,我没有写下这两种和右边的三种特征,而是在一个无限长的列表里面,一直写一直写不停的写,写下无限多个特征,事实上,我们能用算法来处理它们。

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。

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。

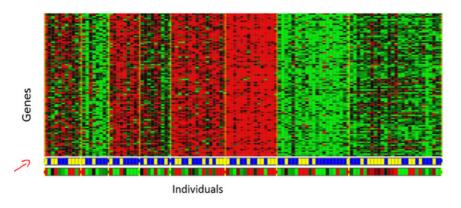


上个视频中,已经介绍了监督学习。回想当时的数据集,如图表所示,这个数据集中每条数据都已经标明是阴性或阳性,即是良性或恶性肿瘤。所以,对于监督学习里的每条数据,我们已经清楚地知道,训练集对应的正确答案,是良性或恶性了。

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,它能被用在很多地方。

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 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。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。

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)