



链滴



多元线性回归解决机器学习问题的一般方法

作者: [vcjmhg](#)

原文链接: <https://ld246.com/article/1591100963946>

来源网站: 链滴

许可协议: [署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

前言

线性回归作为机器学习中经典方法之一，以其形式简单、易于建模而被广泛应用。虽然说线性回归形式简单但却蕴含着机器学习中一些非常重要的思想。许多更能更为强大的非线性模型可在线性模型基础上通过引入层级结构或者高维映射而得到。因此线性模型在机器学习领域是极为重要的模型之一

而线性回归中多元线性回归较之一元线性回归应用范围更广，应用过程也更加复杂。因此有必要多元线性回归在机器学习相关领域一些问题的一般过程进行归纳和总结。当然这一系列步骤成立的前在于：通过一定分析，可以预见该问题可以通过多元线性回归来解决。

多元线性回归解决问题的一般方法

问题抽象

当遇到一个具体问题，首先要对具体的问题抽象成数学语言，并以恰当的数学符号来表示，从建立起数学模型，使得问题更加直观便于分析。而一般来说线性回归遇到的问题一般是这样的：给定 d 个属性描述的示例 (x_1, x_2, \dots, x_d) ，其中 x_i 是在第 i 个属性上的取值，线性模型 (linear model) 试图通过属性的线性组合来进行预测的函数，即

$$f(x) = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b$$

一般向量形式写成

$$f(x) = w^T x + b$$

其中 $w = (w_1, w_2, w_3, \dots, w_d)$ 和 b 学得之后，模型就得以确定。因此通过该步的数学抽象就可以明确，整个线性回归的整个过程主要就是确定合适的 w 和 b 使得训练的模型尽可能地与实际情况相吻合。

数据预处理

其实这一步对所有算法都是适用的。毕竟机器学习算法都是建立在数据的基础之上的，因此对于数据的预处理对所有算法来说都是极其重要的。该过程大体上分为以下几步：

-

- 数据收集：简单来说就是通过各种手段收集自己所需要的数据，例如爬虫、开源数据集等

- 数据清洗：包括数据格式的转化，数据的清洗（处理噪声数据缺失值数据），以及数据的采样

- 数据等价转换：包括统一数据的度量（这在距离计算时十分重要）、零均值化、属性分解以及合

属性的分解是，一个属性能够分解成多个属性，只有某些子属性对于输出有显著的影响。那我们可以只存储这些子属性对于输出具有显著影响，那我们就可以只存储这些子属性，而不用去存储原来属性。其中数据的合并和属性的分解是对立的。将一些子属性合并成一个新属性后，这个属性对于输出的影响会更加显著。那就将这个属性进行合并。

通过数据预处理阶段后，数据会变得更加规整，便于训练模型时直接使用。

确定假设函数

经过数据预处理之后，一般来说就可以确定出一些与输出相关性强的属性，因而可以开始尝试构建解决问题的假设模型。比如说如果确定出与输出相关的属性是 n 维的，那么就可以建立起如下的假设函数：

$$f(x) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

其中：

-

- $x = (x_0, x_1, x_2, \dots, x_n)^T$ 为输入， $w = (w_0, w_1, w_2, \dots, w_n)^T$ 为 $n+1$ 维的列向量，其中 $w_0 = 1, x_0 = 1$

$w = (w_1, w_2, \dots, w_n)^T$ 为参数， $n \times 1$ 维的列向量， b 是 bias

构建代价函数 LMS

在确定假设函数之后，需要确定模型中的参数，具体如何确定那？显然关键在于衡量 $f(x)$ 与 y 之间的差别因此需要建立起一个代价函数来进行衡量。均方误差是回归任务中最常见的性能度量，因此我们以均误差为例，建立起如下的代价函数：

$$E(w, b) = \sum_{i=1}^n (f(x_i) - y_i)^2$$

模型训练

对于不同的参数，成本函数有不同的值，而我们所需要的是让代价函数最小化即：

$$\begin{aligned} (w, b) &= \underset{(w, b)}{\arg \min} \sum_{i=1}^n (f(x_i) - y_i)^2 \\ &= \underset{(w, b)}{\arg \min} \sum_{i=1}^n (y_i - w x_i - b)^2 \end{aligned}$$

求解 w 和 b 使得 $E(w, b)$ 最小化的过程，称之为线性回归型的最小二乘参数估计 (parameter estimation)。通过最小二乘法我们可以推导出通用的线性回归型如下：

$$f(\hat{\mathbf{x}}_i) = \hat{\mathbf{x}}_i^T \mathbf{X}^{-1} \mathbf{X}^T \mathbf{y}$$

其中：

\mathbf{X} : 数据集 \mathbf{D} 所表示的

$m \times (d+1)$ 大小的矩阵

$\mathbf{X}^{-1} \mathbf{X}^T$ 为 $\mathbf{X}^T \mathbf{X}$ 的逆矩阵

如果训练集的矩阵维数不是太大，并且 $\mathbf{X}^T \mathbf{X}$ 刚好为满秩矩阵。我们可以直接套上边公式得出结果。

当然现实任务中 $\mathbf{X}^T \mathbf{X}$ 往往不是满秩矩阵。例如在许多任务中我们会遇到大量的量，其数目甚至超过样例数导致 $\mathbf{X}^T \mathbf{X}$ 的列数多于行数，此 $\mathbf{X}^T \mathbf{X}$ 显然不是满秩的。此时需要解出多个 \hat{w} ，它们都能使均方差最小化，具体选择哪一种需要通过学习算法的归纳偏好来决定，见的做法是引入正则化项。当然如果训练集的矩阵维数很大这种方法，显然不太合适，此时我们可以考虑通过梯度下降算法来求出参数。关于梯度下降算法的具体使用，可以参考周志华老师的《机器学习》，此处不再详述。

预测与结果优化

在参数求出来之后，我们便建立起了一个针对该问题的线性回归模型，然后便可以使用该模型进行预测。当然可能此时得出的模型与现实问题特别吻合，可能造成输入结果与实际结果有出入。此时，要对结果进行优化。

总结

本文尝试从实际问题入手，归纳出了多元线性回归解决实际问题的一般方法的步骤。一般来说，

元线性回归在解决实际问题时，需要经过问题抽象、数据预处理、假设函数的确定、构建代价函数、型训练和预测与结果优化六个部分。一个实际问题经过这几个步骤之后，一般来说便可以建立起一个较合适的多元线性模型。 </p>

<h2 id="参考文献">参考文献</h2>

 https://www.cnblogs.com/ordil/p/10262569.html

 https://zhuannlan.zhihu.com/p/49222906

 https://blog.csdn.net/tuqinag/article/details/54730360

 https://baike.baidu.com/item/%E6%95%B0%E6%8D%AE%E9%A2%84%E5%A4%84%E7%90%86/27112883](https://baike.baidu.com/item/数据预处理/2711288#3)

周志华 《机器学习》[M].北京：清华大学出版社,2016,74-8

