

基于theta ipod的新的超参数选择方法

李 昊¹ 张玉环²

1.Hankuk University of Foreign Studies 韩国首尔 100071

2.北京市第十中学 北京 100039

摘要: 线性回归模型是应用广泛的一类模型估计回归系数的方法是最小二乘法, 然而最小二乘法很容易受到异常点的影响不稳健的。真实数据中都会存在异常点很难避免, 这时候利用最小二乘法进行估计估计结果会比较不理想, 在一定程度上限制了线形模型在实际科研中的应用。面对有异常值的数据的处理方法, 大部分的处理方法为利用 cook' s distance 找出异常值后, 删除再进行回归, 这种方法在一些情况下是不合理的, 因为在实际科学领域中异常值也是包含有一定的信息, 不能完全删除, 特别是在医学药学等领域。在此方法中我们可以利用不同的 loss function, 来决定留下多少有关于异常值的信息, theta ipod (Outlier Detection Using Nonconvex Penalized Regression (Yiyuan She & Art B. Owen, 2011)) 是一个基于最小二乘法下的一个探测异常值并且估计出具有鲁棒性的参数。他的基本思想是在最小二乘法估计的基础上将一个转换参数, 并且为这个转换因子添加一个L1惩罚项函数, 当这个点的转换参数变为0则认为这个点不是异常值。并且为惩罚项函数会依赖于一个超参数 lambda, 最后使用最小 bic 得分的 lambda 来建造模型, 但是在论文中出现的 bic 方法是非常不严谨的, 因为 bic 在 lambda 端点值时会存在极端最小值。所以我提出了一种新的利用 bic 来选择 lambda 的方法, 我们考虑了方差对 bic 值的影响, 先讲总体方差利用样本估计出来。之后添加到新的模型里, 可以得到更加准确的超参数, 进而得到很准确的参数估计值并且准确判断出异常值和异常值对模型影响的大小。在最后一节数据模拟章节给出了具体例子方便参考。

关键词: theta ipod; BIC; 异常值

1 theta-ipod with Huber' s loss 算法介绍

theta ipod (Outlier Detection Using Nonconvex Penalized Regression (Yiyuan She & Art B. Owen, 2011)) 是一个探测异常值并且估计出具有鲁棒性的参数的方法, 他的模型是在原来的一般回归模型中引入一个新的参数 γ , 用于判断数据的距离当前离散程度并且利用阈值函数修正, 为了防止过度拟合我们还为新的参数 γ 添加了一个惩罚项 λ , 使用迭代算法就可以减小异常值影响并且估计出一个具有鲁棒性的回归参数, 并且有效的防止过度拟合, 并且式 1.1 为凸函数可以利用迭代算法就可以很容易找出使得式 1.1 最小的参数估计。

数学模型为用 1.1 来估计 1.2 中的 β 。其中 1.2 的 $\rho()$ 为 Huber' s 损失函数

$$g(\beta, \gamma, \sigma) = \frac{1}{2\sigma} \sum_{i=1}^n (y_i - x_i^T \beta - \lambda_i) + \frac{1}{2} c n \sigma + \lambda \sum_{i=1}^n |\gamma_i| \quad (1.1)$$

$$l(\beta, \sigma) = \sum_{i=1}^n \rho\left(\frac{y_i - x_i^T \beta}{\sigma}; \lambda\right) + \frac{1}{2} c n \sigma \quad (1.2)$$

并且可进一步证明在 $\Psi(t; \lambda) = t - \Theta(t; \lambda)$ 条件下两个模型的所估计的 β 参数是一致的

证明如下

$$\Psi(t; \lambda) = \rho'(t; \lambda)$$

$$\begin{aligned} X^T X \hat{\beta} &= X^T (y - \hat{\gamma}) \\ X^T X \hat{\beta} &= X^T (y - \Theta(\hat{r}; \lambda \hat{\sigma})) \\ X^T (y - \Theta(\hat{r}; \lambda \hat{\sigma})) &= X^T (y - \hat{r} + \Psi(\hat{r}; \lambda \hat{\sigma})) \\ X^T (y - \hat{r}) + X^T (\Psi(\hat{r}; \lambda \hat{\sigma})) & \\ \hat{r} &= y - X \hat{\beta} \\ X^T X \hat{\beta} + X^T (\Psi(\hat{r}; \lambda \hat{\sigma})) &= X^T X \hat{\beta} \\ X^T (\Psi(\hat{r}; \lambda \hat{\sigma})) &= 0 \end{aligned}$$

$$\Psi(\hat{r}; \lambda) = 0$$

Theta-ipod 迭代算法

将有异常值的数据进行最小二乘估计所求出的 β_{ols} 作为迭代起始点

$$\gamma^{(0)} = y - X^T (\hat{\beta}_{ols})$$

直到 γ 不再变化循环下列

$$y_{adj} = y - \gamma^{(j)}$$

$$\beta^{(j)} = (X^T X)^{-1} X^T y_{adj}$$

$$r^{(j)} = y - X^T \beta^{(j)}$$

$$\gamma^{(j+1)} = \Theta(r^{(j)}, \lambda)$$

最后返回 γ 和 β 的值

2 theta-ipod with Huber' s loss 算法新的超参数选择方法

在 Outlier Detection Using Nonconvex Penalized

Regression (Yiyuan She & Art B. Owen, 2011) 一文中提到的 λ 选择方法是创建一个从0开始到 y 的最大值的向量之后分别计算

$$BIC(\lambda) = m \log(RSS/m) - k(\log(m)+1) \quad (2.1)$$

选择最小的 BIC 所对应的 λ ，利用此 λ 求出 β

$$2.1 \text{ 式中 } RSS = \left\| \left(I - X(X^T X)^{-1} X^T \right) (y - \hat{\gamma}) \right\|^2$$

2.1 式中 m = 样本数 - 参数个数

2.1 式中 k = 当前 λ 下的 γ 不为 0 的个数

利用 bic 作为参数选择方法是很明智的他可以在大样本的情况下很快计算出结果

但是 Outlier Detection Using Nonconvex Penalized Regression (Yiyuan She & Art B. Owen, 2011) 论文中的 bic 计算方法有一个致命缺点，在他的论文中选择超参数 λ 的部分中也提到：“曲线有时在 λ 范围的末端附近具有狭窄的局部最小值。为了抵消这种影响，我们将平滑样条拟合到数据点集 ($DF(\lambda)$, $BIC^*(\lambda)$) 并选择具有最大邻域的局部最小值。可以使用平滑样条的局部最大值来确定邻域大小。当然，可能还有其他合理的方法来解决这个问题”。经过推导实验，发现他并没有包括原来样本的方差信息，所以这导致了他在 λ 在靠近 0 时有极端最小值 (图 2.1)。我们将总体的方差的预测值设置为是最小二乘法中的样本残差减去样本残差中位数的绝对值的中位数除以 0.6745 (0.6745 是一个校正因子，当误差服从正太分布时可以保证方差估计的一致性)。

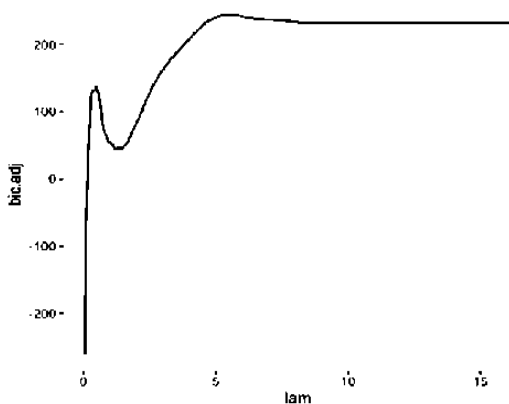


图 2.1

我们提出一种新的 bic 计算方法，计算公式如下：

$$NEWBIC(\lambda) = n \log(\sigma^2) + RSS/(\sigma^2) + k(\log(m)+1) \quad (2.2)$$

2.2 式中 σ 是最小二乘法中的样本残差减去样本残差中位数的绝对值的中位数除以 0.6745

2.2 式中 n 为样本个数

$$2.2 \text{ 式中 } RSS = \left\| \left(I - X(X^T X)^{-1} X^T \right) (y - \hat{\gamma}) \right\|^2$$

2.2 式中 m = 样本数 - 参数个数

2.2 式中 k = 当前 λ 下的 γ 不为 0 的个数

利用此方法就可以解决 Outlier Detection Using Nonconvex Penalized Regression (Yiyuan She & Art B. Owen, 2011) 论文中计算 bic 的漏洞，在 bic 中添加了数据中的 y 的离散信息，非常有效的解决了 λ 范围的末端附近具有狭窄的局部最小值，利用此 newbic 计算出的 bic 和 λ 的折线图 (图 2.2)，可以看出 bic 符合典型机器学习参数的选择曲线。

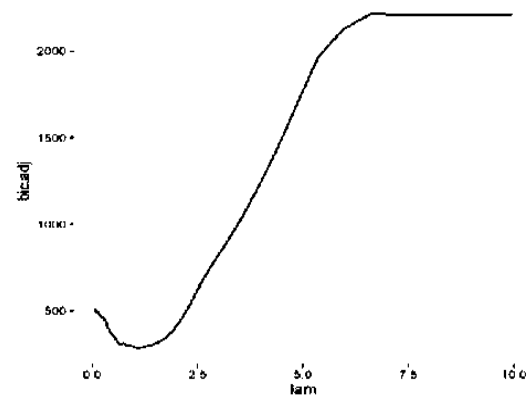


图 2.2

3 数据模拟

生成 100 个服从平均分布的自变量 x 将 β 的真值设为 (0, 2)，利用 $X^T \beta$ 计算出 y 值并将其中百分之 10 的点加上 6 作为异常值。运行 theta-ipod 算法并利用 newbic 筛选出最合适的 λ ，画出每一个样本根据不同 λ 的值的 γ 路径 (图 3.1) 和根据不同 λ 的值 β 值的变化路径，其中的虚线为最优 λ 时 β 和每个样本的 γ 值。

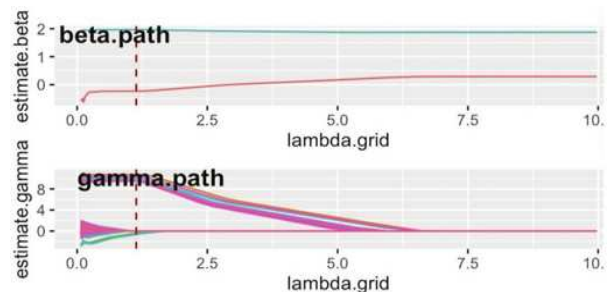


图 3.1

从图 3.1 可以看出原论文的所筛选出最小的 bic 时的最小值可以发现 lambda 无限接近于 0，这会导致将所有样本点变为异常值，所估计的参数也会变得不准确。在 r 语言中利用 lm 求出的 OLS 估计值为 (1.412, 2.096) 可以看出 OLS 估计收到了很大的影响，尤其是 $\beta_{(0)}$ 。通

过上述算法我们可以得出当 newbic 最小时所对应的 λ 为 1.134418, 在当前 λ 下 β 的预测值 (蓝色线条) 为 $(-0.2320448, 1.9688299)$, 可以发现与真值 (粉色虚线) 的 $(0, 2)$ 所差无几并且将异常值准确的判断了出来如图 (3.2) 所示。

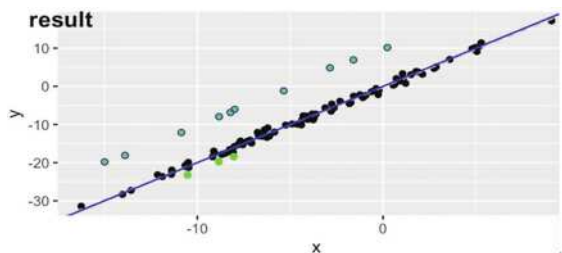


图 3.2

上图中的空心原点为异常值, 绿色点为 λ 值为

1.134418 时所探测出的异常值, 虽然有一些点被污染但是还是将异常值充分的找了出来。

当然我们也可以利用交叉验证的方法来选择 λ 值, 经过实验验证, 交叉验证的方法所需时间非常久, 而且抽样方差也会非常大 (当验证数据中有太多的异常值时和没有异常值时), 所以 mse 和 λ 的曲线会非常的混乱而且偏差非常严重, 所以并不推荐使用。当然应该也会有更加准确的 λ 选择方法来避免正常点被污染。

参考文献:

- [1][She and Owen, 2011] She, Y., Owen, A. B. (2011). Outlier detection using nonconvex penalized regression. Journal of the American Statistical Association, 106, 626–639.