

# 监督学习方法

司继春

上海对外经贸大学

2025年3月











## James-Stein估计量

- 实际上, 不仅仅可以向0进行收缩, 取任意的向量 $\vartheta$ , 估计量:

$$\hat{\theta}^{JS} = \left( 1 - \frac{K-3}{\|x-\vartheta\|_2^2} \right) (x-\vartheta) + \vartheta = x + \frac{K-3}{\|x-\vartheta\|_2^2} (\vartheta - x)$$

实现了向向量 $\vartheta$ 的收缩。

- 以上James-Stein估计量可以扩展到回归的情形。在线性回归中, 如果我们令OLS估计量向0收缩, 由于对 $Y$ 的预测为 $X\hat{\beta}$ , 于是我们就有了James-Stein估计量:

$$\hat{\beta}^{JS} = \hat{\beta}^{OLS} \left[ 1 - \frac{(K-2)\hat{\sigma}^2}{\hat{\beta}'X'X\hat{\beta}} \right]$$

与之前类似, 使用该估计量可以降低对 $y$ 进行估计的均方误差。



























# 稀疏一致性

- Lasso的理论性质依赖于稀疏性 (sparsity) :
  - 记 $\beta_0$ 为真实的 $\beta$ , 令 $S_0 = \{\beta_{0,k} \neq 0, k = 1, \dots, K\}$ 为活跃集 (active set), 记 $K_0 = |S_0|$ 为 $\beta_0$ 中不为0的分量的个数, 即稀疏指数 (sparsity index)。
  - 允许变量个数 $K \rightarrow \infty$ , 但是不能太快。
- 记 $\hat{S} = \{\hat{\beta}_k^s \neq 0, k = 1, \dots, K\}$ , 那么选择一致性 (selection consistency) 或稀疏一致性 (sparsistency) 定义即

$$P(\hat{S} = S_0) = 1$$

# oracle性质

- 而如果在选择一致性的基础上，估计量额外满足

$$\sqrt{N} \left( \hat{\beta}_{S_0}^s - \beta_{0,S_0} \right) \overset{a}{\sim} \mathcal{N} (0, \Sigma)$$

其中 $\Sigma$ 为估计量

$$\hat{\beta}^{oracle} = \arg \min_{\beta} \sum_{i=1}^N (y_i - x_i' \beta)^2$$

$$s.t. \beta_{S_0^c} = 0$$

的协方差矩阵，那么我们称估计量 $\hat{\beta}^s$ 具有**oracle性质**（oracle property, Fan和Li, 2001）。

- Oracle性质意味着估计量 $\hat{\beta}$ 的统计性质至少要与假设知道哪些系数为0的前提下所得到的估计量（ $\hat{\beta}^{oracle}$ ）性质一样好。



# Lasso的稀疏一致性

- 然而，Lasso估计量的选择一致性需要更强的条件，其中比较重要的是不可代表性 (irrepresentability) ，即存在 $\gamma > 0$ ，有

$$\max_{j \in S_0^c} \left\| (X'_{S_0} X_{S_0})^{-1} X'_{S_0} x_j \right\| \leq 1 - \gamma$$

以上条件实际上要求 $S_0$ 中的 $x$ 对于非 $S_0$ 中的 $x$ 的相关性足够的弱，或者解释能力足够的小。

- 最理想的情况下，两类 $x$ 应该正交，此时 $\gamma = 0$ 。
- 可以证明在一定的条件和适当的 $\lambda$ 的选取下，Lasso回归可以以一个比较高的概率达到选择一致性 (Hastie, Tibshirani和Wainwright, 2015) 。

# Lasso的oracle性质

- 然而具体到oracle性质，遗憾的是岭回归和Lasso回归通常不具有oracle性质：
  - 岭回归通常不具有选择一致性
  - 而Lasso回归虽然可能达到选择一致性，但是通常并不是一致估计量（渐进有偏）
- Oracle性质之所以吸引人主要在于其可以保证选择一致性的前提下，同时对非零参数达到了一致且渐进正态的估计，从而可以进一步对其进行统计推断，而简单的岭回归和Lasso回归却无法达到这样的效果。

# 选择后估计

- 基于此，一种方法是将Lasso回归作为变量选择的工具，即挑选出Lasso回归中系数不为0的变量，进而对 $\hat{S}$ 集合中所对应的参数用OLS进行估计，这种做法也被称为选择后估计（post-selection estimator），或者后Lasso估计（post Lasso estimator）。
- Leeb和Pötscher（2005，2008）指出，通常Lasso回归中的一致性或者收敛性都是点态收敛，而非一致收敛，这导致不同的 $\beta_0$ 会需要不同的样本量来达到需要的统计性质，而实践中往往 $\beta_0$ 是未知的而样本量是固定的，从而根据Lasso进行变量选择后的OLS估计的统计性质仍然是复杂的问题。
- Berk等人（2013）讨论了选择后估计的推断问题，提出可以通过拓宽置信区间的方法对参数进行合适的区间估计，而Lee等人（2016）则讨论了给定所选模型下的推断问题。

# 选择后估计

- 而另一方面，虽然有以上的缺点，选择后估计也的确有其优势。
- Belloni和Chernozhukov (2013) 比较了模型选择后估计与Lasso估计，发现在高维稀疏模型中，即使Lasso方法可能无法达到选择一致性，不过选择后估计的在收敛速度和偏差方面表现得至少与Lasso估计一样好
- 在一些情况下其收敛速度和偏差可以严格由于单纯的Lasso估计
- 甚至在极端情况下，如果Lasso完美的选择出了正确的模型，选择后估计就成为oracle估计量。



# 超参数

- 在以上的岭回归和Lasso回归中，拉格朗日乘子 $\lambda$ ，或者等价的 $t$ 控制了收缩的程度
  - 一个更大的 $\lambda$ 或者更小的 $t$ 意味着更“小”的模型或者更多的向0收缩，从而偏差大但是方差小；
  - 在这里同样有偏差-方差的权衡问题
  - 应该存在一个“最优”的 $\lambda$ 使得预测误差能够达到最小。
- 注意这里由于 $\hat{y} = x' \hat{\beta}^{lasso/r}$ ，从而参数 $\lambda$ 并不参与模型的预测， $\lambda$ 仅仅是为了“训练”在机器学习中通常将模型参数的过程称为“训练”（train）模型。
- 模型所需要的参数： $\hat{\beta}^{lasso/r}$ 是一个 $\lambda$ 的函数，给定一个 $\lambda$ 就会有一个对应的模型和参数，我们把这种训练模型之前需要设定的参数称为“超参数”（hyperparameter）或者“调整参数”（tuning parameter）。

# 超参数选择

- 最常用的方法是数据驱动的交叉验证法。
- 出于计算速度的要求，一般进行5-10折的交叉验证即可。
- 实际进行时会将样本随机分为 $S$ 折，然后将 $\lambda$ 从一个很大的数字到0进行格点（grid）化，针对每一个 $\lambda$ ，使用交叉验证的方法留出1折数据作为验证集，其他的 $S - 1$ 折作为训练集，再在留出的1折数据中进行预测，最终挑出能够使得均方误差最小的 $\lambda$ 。



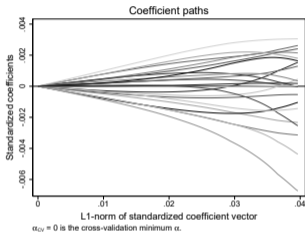
# 交叉验证

## 预测香港GDP增速

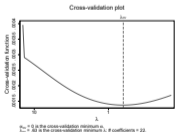
其中

- linear代表要在线性回归中加入惩罚项
- 选项sel(cv)代表使用交叉验证选取 $\lambda$
- alpha(0)代表进行岭回归估计
- rseed()是为了保证每次运行结果都相同的一个随机数种子。
- 估计完成后可以使用lassoinfo命令查看选取的最优 $\lambda$ ，也可以使用coefpath命令和cvplot命令可以分别画出系数路径图和交叉验证图。

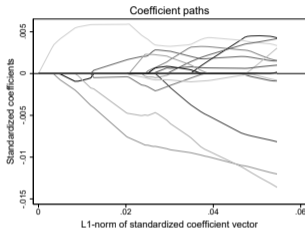
# Stata中的Lasso



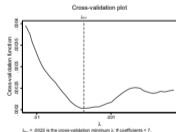
(a) 岭回归系数路径



(c) 岭回归交叉验证



(b) Lasso回归系数路径



(d) Lasso回归交叉验证

# Stata中的Lasso

## 预测香港GDP增速

- 此外，还可以在`sel(cv)`中加入`folds`选项以设定交叉验证的折数
  - 比如`sel(cv, folds(9))`即进行9折交叉验证。
- 如果需要进行Lasso回归，可以设定`alpha(1)`，或者直接使用`lasso`命令：

```
1 | lasso linear HongKong Australia-Thailand if _n<=18, rseed(5)  
   | sel(cv)
```

# Stata中的Lasso

## 预测香港GDP增速

- 除了以上的估计后命令外，还可以使用lassocoef命令查看Lasso回归保留了哪些变量
- 也可以使用preidct命令进行预测，在预测时可以使用经过收缩的Lasso回归系数，也可以使用选择后估计进行预测：

```
1 predict HongKong_lasso
2 label variable HongKong_lasso "Lasso"
3 predict HongKong_lasso_post, post
4 label variable HongKong_lasso_post "PostLasso"
```

其中加入post选项的是进行后Lasso估计进行的预测。

# Stata中的Lasso

## 预测香港GDP增速

- 如果需要手动进行后Lasso估计，可以使用Lasso估计之后保留的 `e(post_sel_vars)`，该宏记录了所有被Lasso命令保留的变量，其中第一个变量为被解释变量：

```
1 | lasso linear HongKong Australia-Thailand if _n<=18, rseed(5)  
   | sel(cv)  
2 | local post_vars=e(post_sel_vars)  
3 | reg `post_vars'
```



# 超参数选择

- 除了交叉验证外，一些信息准则，特别是BIC也会被用于 $\lambda$ 的选择中。
- 此外，文献中对于 $\lambda$ 的取值也有非常多的讨论（如Bickel、Ritov和Tsybakov, 2009; Belloni和Chernozhukov, 2011; Belloni等, 2012; Belloni、Chernozhukov和Hansen, 2014; Belloni、Chernozhukov和Wei, 2016）
- 其形式大概如：

$$\lambda = \frac{c}{\sqrt{N}} \hat{\sigma} \Phi^{-1} \left( 1 - \frac{\gamma}{2K} \right)$$

其中根据Belloni和Chernozhukov (2011) 建议,  $c = 1.1$ ,  $\gamma$ 为未能成功排除真值为0的参数的概率, Stata中取 $\gamma = \frac{0.1}{\ln(\max\{K, N\})}$ 。根据如上形式就可以使用插入 (plug-in) 法估计 $\lambda$ 。

# Stata中的Lasso

## 预测香港GDP增速

- 在Stata中，如果使用BIC选取超参数，可以使用sel(bic)选项：

```
1 lasso linear HongKong Australia-Thailand if _n<=18, rseed(5)  
  sel(bic)
```

而如果需要使用插入法，可以使用sel(plugin)：

```
1 lasso linear HongKong Australia-Thailand if _n<=18, rseed(5)  
  sel(plugin)
```

# 弹性网

作为一种收缩估计量，岭回归和Lasso回归分别使用了 $L^2$ 范数和 $L^1$ 范数，而类似的可以定义很多不同的收缩方法。

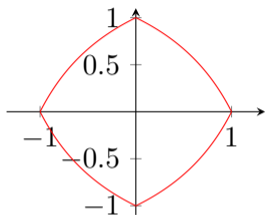
- 首先，我们可以将Lasso回归和岭回归结合起来，即同时使用 $L^2$ 范数和 $L^1$ 范数的线性组合作为惩罚项，即最小化

$$\hat{\beta}^{elastic}(\lambda, \alpha) = \arg \min_{\beta} \frac{1}{N} \sum_{i=1}^N (y_i - x'_i \beta)^2 + \lambda \left( \alpha \|\beta\|_1 + \frac{1-\alpha}{2} \|\beta\|_2^2 \right)$$

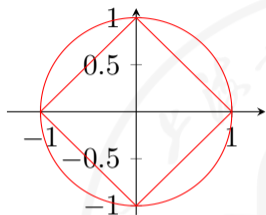
其中 $\alpha \in [0, 1]$ 为超参数。

- 易知当 $\alpha = 0$ 时以上即岭回归
- 而当 $\alpha = 1$ 时即Lasso回归
- 以上回归被称为弹性网（elastic net, Zou和Hastie, 2005）。

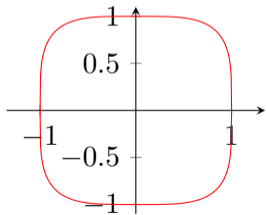
# 不同收缩方法的惩罚项



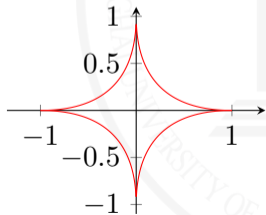
(a)  $\alpha = 0.5$ , 弹性网



(b)  $\gamma = 1, 2$ , Lasso回归和岭回归



(c)  $\gamma = 4$



(d)  $\gamma = 0.5$

# Bridge

- 如果我们将 $L^1$ 和 $L^2$ 范数改为更加一般的 $L^\gamma$ 范数，那么我们就得到了Bridge估计量，即最小化

$$\hat{\beta}^{bridge}(\lambda) = \arg \min_{\beta} \frac{1}{N} \sum_{i=1}^N (y_i - x'_i \beta)^2 + \lambda \sum_{k=1}^K |\beta_k|^\gamma$$

- 显然当 $\gamma = 1$ 时，即Lasso回归，而当 $\gamma = 2$ 时，就得到了岭回归。
- 值得注意的是，当 $\gamma < 1$ 时，可行集并不是一个凸集，当然以上优化问题也不是一个凸优化问题，在计算上存在困难，因而文献中经常假设 $\gamma \geq 1$ 。
- 然而， $\gamma < 1$ 时的Bridge回归在偏差和选择一致性上都更有优势（Bühlmann和van der Geer, 2011），相比较于Lasso具有更好的性质。

# 适应性Lasso

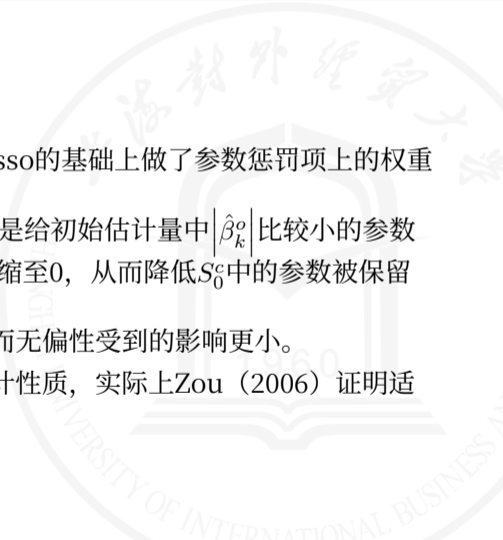
- 优化问题限制了 $\gamma < 1$ 时的模型应用，此时我们可以使用Zou (2006) 提出的适应性Lasso (adaptive Lasso)：

$$\hat{\beta}^{ada}(\lambda) = \arg \min_{\beta} \frac{1}{N} \sum_{i=1}^N (y_i - x_i' \beta)^2 + \lambda \sum_{k=1}^K w_k |\beta_k|$$

- 其中 $w_k = \frac{1}{|\hat{\beta}_k^o|^\theta}$ 为权重， $\theta > 0$ ，
  - $\hat{\beta}_j^o$ 为 $\beta$ 的一个一致估计量作为初始估计量，比如OLS估计量。
  - 如果 $K > N$ ，也可以使用岭回归或者Lasso回归等估计量。
  - 注意到如果 $|\hat{\beta}_k^o| = 0$ ，那么 $w_k = \infty$ ，则 $\hat{\beta}_k^{ada} = 0$ ，从而在估计适应性Lasso时，只需要对第一步Lasso回归中非0的参数进行估计即可。
- 以上步骤可以反复进行，即将 $\hat{\beta}^{ada}$ 作为初始估计量计算新的权重，再次计算适应性Lasso，当然一般实践中只进行一次迭代效果已经足够好。

# 适应性Lasso

- 适应性Lasso具有很多的优良性质。
  - 首先从计算上来看，适应性Lasso仅仅在Lasso的基础上做了参数惩罚项上的权重调整，其计算代价相比Lasso是完全相同的。
  - 从统计性质上，对于 $\theta > 0$ ，适应性Lasso总是给初始估计量中 $|\hat{\beta}_k^o|$ 比较小的参数以更大的惩罚，从而使得估计量更加容易收缩至0，从而降低 $S_0^c$ 中的参数被保留的可能性；
  - 而对于 $|\hat{\beta}_k^o|$ 比较大的参数以更小的惩罚，从而无偏性受到的影响更小。
  - 综合起来，适应性Lasso可以达到更好的统计性质，实际上Zou (2006) 证明适应性Lasso可以达到oracle性质。

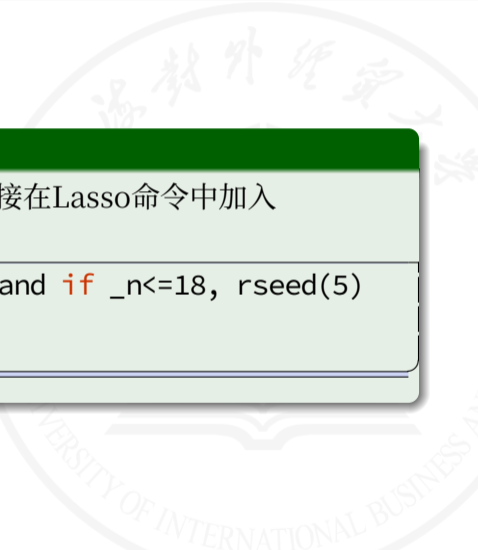


# Stata中的Lasso

## 预测香港GDP增速

- 在Stata中，如果需要做适应性Lasso，可以直接在Lasso命令中加入sel(adaptive)选项：

```
1 lasso linear HongKong Australia-Thailand if _n<=18, rseed(5)  
   sel(adaptive)  
2 lassocoef
```





# 平方根Lasso

- 以上Lasso的变种都集中在对惩罚项的设定上，而Belloni, Chernozhukov和Wang (2011) 则建议将目标函数中的残差平方和改为其平方根，即

$$\hat{\beta}^{ada}(\lambda) = \arg \min_{\beta} \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - x_i' \beta)^2} + \frac{\lambda}{N} \sum_{k=1}^K |\beta_k|$$

即平方根Lasso (square-root Lasso)。

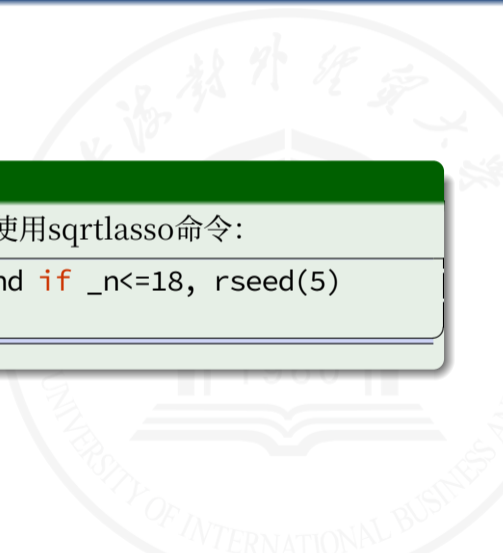
- 该方法的好处是在使用插入法计算 $\lambda$ 的最优选择时，无须估计 $\sigma$ 。
- 此外，该方法可以在误差项非正态分布时也达近乎oracle的收敛速度。

# Stata中的Lasso

## 预测香港GDP增速

- 在Stata中，如果需要做平方根Lasso，可以使用sqrtlasso命令：

```
1 sqrtlasso HongKong Australia-Thailand if _n<=18, rseed(5)
2 lassocoef
```





# 决策树与随机森林

- Lasso回归在处理高维问题上有其优势，然而仍然需要对函数形式进行假设。
- 我们之前已经学过一些非参数回归的方法，然而这些方法很容易面临维数的诅咒。
- 在机器学习中，针对不同类型被解释变量，通常有不同的方法可以在不假设函数形式的前提下进行预测，比如基于树的方法（tree-based methods）、神经网络（neural networks）等，都可以用于回归和分类。这一节我们主要介绍基于树的方法。















# 特征空间的划分

- 而对于分类树，如果可能的分类为 $l \in \{1, 2, \dots, L\}$ ，可以定义 $L$ 个阶梯函数：

$$f_l(x) = \sum_{m=1}^M p_{ml} \cdot \mathbb{1}\{x \in R_m\} p_m$$

为区域 $R_m$ 中 $y$ 取值为 $l$ 的概率，给定划分，我们可以使用 $R_m$ 内的 $y$ 取 $l$ 的比例对 $p_{ml}$ 进行估计：

$$\hat{p}_{ml} = \frac{\sum_{i=1}^N \mathbb{1}\{y_i = l\} \cdot \mathbb{1}\{x_i \in R_m\}}{\sum_{i=1}^N \mathbb{1}\{x_i \in R_m\}}$$

问题：如何得到划分？？







# 分类树的预测误差

- 现在再次考虑以上的二分类问题。如果我们将节点 $D$ 划分为两个子节点： $D_1, D_2$
- 为了评估划分之后的杂度，我们还需要进行加权
  - 比如对于基尼系数，需要计算

$$G(D_1, D_2) = \frac{|D_1|}{|D|} G_{D_1} + \frac{|D_2|}{|D|} G_{D_2}$$

其中 $|D|$ 代表 $D$ 节点的样本量。

- 而对于交叉熵，可以定义信息增益：

$$Gain(D, D_1, D_2) = CEnt_D - \left( \frac{|D_1|}{|D|} CEnt_{D_1} + \frac{|D_2|}{|D|} CEnt_{D_2} \right)$$

以及增益率：

$$Gain\_ratio(D, D_1, D_2) = \frac{Gain(D, D_1, D_2)}{- \left( \frac{|D_1|}{|D|} \log \frac{|D_1|}{|D|} + \frac{|D_2|}{|D|} \log \frac{|D_2|}{|D|} \right)}$$





# 过拟合与超参数选择

- 另一种方法是通过一些规则限定树的大小。比如，我们可以规定一个杂度下降的一个最小幅度，如果划分一个节点杂度下降的幅度小于这个最小幅度，就停止划分该节点，并将该节点作为叶子节点。
- 再比如，我们可以通过限制树的层数（深度）限制树的大小。
- 或者，我们可以限制叶子结点的数量。如果记 $|T|$ 为叶子结点的数量，我们可以要求 $|T|$ 小于等于某一个上限。此外，还可以将 $|T|$ 作为正则化项，通过最小化

$$\min \sum_{m=1}^{|T|} N_m Q_m + \alpha |T|$$

对于回归问题， $Q_m$ 即均方误差，而对于分类问题， $Q_m$ 则为基尼系数、交叉熵等的度量。

- 注意其中 $\alpha$ 是一个超参数，可以通过交叉验证等方法确定 $\alpha$ 。



# Boosting方法简介

- Boosting方法通常是一类可加模型 (additive model) :

$$f(x) = \sum_{m=1}^M \beta_m b(x, \gamma_m)$$

如果令  $f_0(x) = 0$ , boosting方法迭代地方法逐步解出  $\beta_m, \gamma_m$  从而得到最终的模型。

- 由于采用了迭代的方法, 因而不同的  $m$  之间并不是独立的。
- 常见的boosting方法如AdaBoot和gradient boosting的算法和性质可以参考 Hastie, Tibshirani和Friedman (2009), 在此不再赘述。





