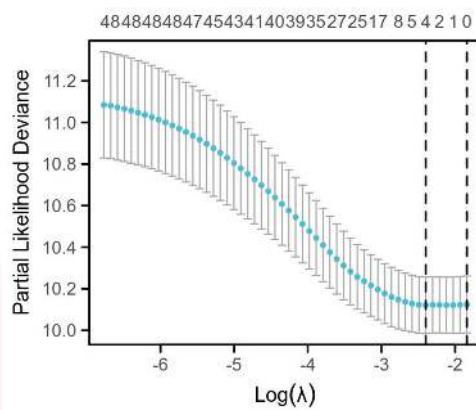


## 临床意义 - 预后 Lasso 系数筛选



网址: <https://www.xiantao.love>



更新时间: 2023.02.07

## 目录

基本概念 .....	3
应用场景 .....	3
分析流程 .....	4
结果解读 .....	6
数据格式 .....	8
参数说明 .....	9
方法 .....	9
点 .....	11
误差线 .....	12
标题文本 .....	13
风格 .....	14
图片 .....	14
结果说明 .....	15
主要结果 .....	15
补充结果 .....	16
方法学 .....	17
如何引用 .....	18
常见问题 .....	19

## 基本概念

- **Lasso 回归**：在线性回归的基础上，通过增加**惩罚项**（ $\lambda \times \text{斜率的绝对值}$ ），减少模型的过拟合，提高模型的泛化能力。另外一种也是通过增加惩罚项来减少模型的过拟合的方法是岭回归，对应的惩罚项是（ $\lambda \times \text{斜率的平方}$ ）。惩罚项在机器学习领域也叫做正则化，其中，Lasso 回归的惩罚项是**L1 正则化**（曼哈顿距离（参数绝对值求和）），而岭回归的惩罚项是**L2 正则化**（欧氏距离（参数平方值求和））
- Lasso 可用于 logistics、Cox 其中，此模块就是 Lasso 在预后中的应用。预后 Lasso 常常出现在构建预后模型或者筛选变量上，最常出现两种图，一种是 系数( $\lambda$ )筛选的图，另外一种为变量轨迹图。Lasso 的  $\lambda$  筛选一般会采用**交叉验证**的手段进行筛选，常见的会有五折和十折交叉验证。

## 应用场景

将预后 Lasso 系数筛选过程中各个  $\lambda$  值（惩罚项）对应的统计量(似然偏差值或 C 指数)进行可视化，以**构建预后模型或者筛选变量**。当样本较少或者变量较多（少于样本数一半的变量）时，可以用 Lasso 直接构建预后模型或者筛选变量。

## 分析流程

上传数据 → 数据处理(清洗) → lasso 预后分析 → lasso 系数筛选可视化

➤ 数据格式: xlsx / csv / txt 文件格式:

- 第 1 列数据作为结局变量(事件发生情况), 需要是数值类型数据, 用 (0 和 1, 0 表示未发生事件, 1 表示发生了事件) 或 (1 和 2, 1 表示未发生事件, 2 表示发生了事件) 表示, 注: 第 1 列(结局变量)不能都是删失

	A	B	C	D	E	F	G	H	I	J
1	event	time	Gene 1	Gene 2	Gene 3	Gene 4	Gene 5	Gene 6	Gene 7	Gene 8
2	1	306	-0.0222806	-2.7704356	-0.4670543	0.65971003	-0.5151171	0.02662303	0.96486822	-1.0649103
3	1	455	-1.1832171	-0.316118	0.60603766	-0.3035985	0.14872787	-1.2101381	1.2133766	-0.195116
4	0	1010	-0.6229744	1.8458862	1.51766889	-0.8392173	-0.2736995	-1.9147945	0.98693195	-0.1859204
5	1	210	-0.9614322	-0.136129	0.70702704	-2.2407777	-0.1158965	-1.6783932	0.58134326	1.24883956
6	1	883	-2.0090579	0.75447248	-1.3601112	0.74345661	1.24201298	0.37301567	0.65579689	-0.6545812
7	0	1022	0.79356585	-0.2366209	-0.5012338	0.93805595	-1.219659	-1.625082	0.32808147	1.04612915
8	1	310	-0.2919467	-0.1946606	0.20888995	-0.7444609	-1.5937526	-0.1180966	-1.2945132	1.31273942
9	1	361	0.70980194	-0.255891	1.34543741	-1.0402237	-0.0404575	0.70215222	-0.5927644	2.2450957
10	1	218	0.25708681	0.37272211	-0.0135116	-1.0462137	0.86385945	0.76747574	0.84760712	-1.5613683
11	1	166	2.50492511	1.37171585	-0.1141563	-1.7225671	-0.0523161	-1.2741515	1.36151475	-0.6789088
12	1	170	0.44265243	-0.4562389	-0.0150315	1.58488562	0.05882662	-1.2996843	-1.5015927	-1.0092582

- 第 2 列数据作为时间变量(具体时间/生存时间, 必须以天作为单位, 并且时间要长于 1 年以上), 需要是数值类型数据, 注: 第 2 列(时间变量)不能都是同一个时间, 并且不能出现小于 0(负数)和非数值的情况

- 第 3 列开始直至后面每 1 列都代表一个变量/样本/分子, 必须是数值类型数据

➤ 数据处理: 分别对第 1 列 (事件)、第 2 列 (时间)、第 3 列开始后的所有变量进行清洗 (去除掉数据中的非数值或者不符合条件的数据)

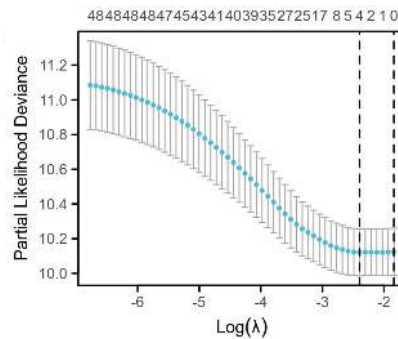
➤ Lasso 预后分析:

- 构建 lasso 预后模型

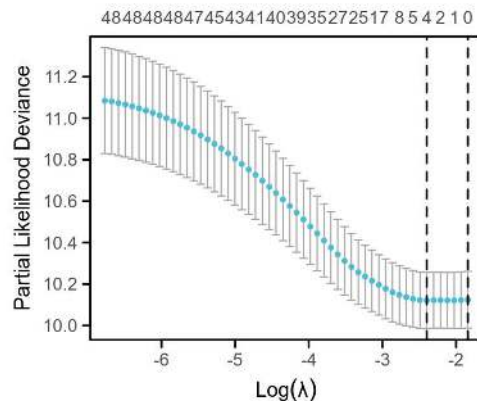
- 计算模型的 lambda 值
- 通过 lambda 值计算变量的系数值
- 筛选掉 lambda 值对应系数为 0 的变量(系数为 0 表示变量之间不存在相关关系，在预后模型中没有实质上的意义)

➤ Lasso 系数筛选可视化

- Lasso 预后分析得到的 lambda 值取对数，对应到 lasso 系数筛选可视化结果的横坐标值
- Lasso 预后分析得到的不同指标下的似然偏差值（deviance）（默认）或 C 指数（c-index）对应到 lasso 系数筛选可视化结果的纵坐标
- 进行可视化，结果如下：



## 结果解读



- 下方 x 轴：表示 Lasso 回归中惩罚项  $\lambda$  值取对数 ( $\log(\lambda)$ )
- 上方 x 轴的数字：表示每个  $\lambda$  值对应的非 0 系数的变量个数
  - 这些数字对应的值是说：不同  $\lambda$  值计算得到模型中所有变量系数不为 0 的变量的个数，而不是所有的变量（要是数值与变量个数对应不上，则是因为缺少的那些变量间不存在相关关系（系数为 0）被筛选掉了）
  - 由于可视化结果是 ggplot2 格式，故不能展示全部的数值
- y 轴：表示在不同指标下的似然偏差值 (deviance)（默认）或 C 指数 (c-index)
- 每个点：表示数据在进行交叉验证过程中，每个  $\lambda$  对应的似然偏差值（默认）或 C 指数的均值
- 每条竖线（误差线）：表示数据在进行交叉验证过程中，每个  $\lambda$  对应的似然偏差值的标准误
- 左边虚线：表示评价指标最佳的  $\lambda$  值 ( $\lambda_{\min}$ )
- 右边虚线：表示评价指标在最佳值 1 个标准误范围的模型的  $\lambda$  值 ( $\lambda_{1se}$ )
- 当选择的指标为 deviance 时，y 值（似然偏差值）越小对应的模型越好

- 当选择的指标为 c-index 时，y 值（C 指数）越大对应的模型越好
- 当  $\lambda_{\min}$  和  $\lambda_{1se}$  一样(图中只有 1 根虚线并且在最右侧)，说明模型没有筛选出来任何一个非 0 系数的变量； $\lambda_{\min}$  可能对模型过于严格， $\lambda_{1se}$  对应的变量越少，模型会更加简洁；两个都可以选，比较常用的是  $\lambda_{\min}$ ，如果  $\lambda$  对应的变量较多，也会用  $\lambda_{1se}$



## 数据格式

	A	B	C	D	E	F	G	H	I	J
1	event	time	Gene 1	Gene 2	Gene 3	Gene 4	Gene 5	Gene 6	Gene 7	Gene 8
2	1	306	-0.0222806	-2.7704356	-0.4670543	0.65971003	-0.5151171	0.02662303	0.96486822	-1.0649103
3	1	455	-1.1832171	-0.316118	0.60603766	-0.3035985	0.14872787	-1.2101381	1.2133766	-0.195116
4	0	1010	-0.6229744	1.8458862	1.51766889	-0.8392173	-0.2736995	-1.9147945	0.98693195	-0.1859204
5	1	210	-0.9614322	-0.136129	0.70702704	-2.2407777	-0.1158965	-1.6783932	0.58134326	1.24883956
6	1	883	-2.0090579	0.75447248	-1.3601112	0.74345661	1.24201298	0.37301567	0.65579689	-0.6545812
7	0	1022	0.79356585	-0.2366209	-0.5012338	0.93805595	-1.219659	-1.625082	0.32808147	1.04612915
8	1	310	-0.2919467	-0.1946606	0.20888995	-0.7444609	-1.5937526	-0.1180966	-1.2945132	1.31273942
9	1	361	0.70980194	-0.255891	1.34543741	-1.0402237	-0.0404575	0.70215222	-0.5927644	2.2450957
10	1	218	0.25708681	0.37272211	-0.0135116	-1.0462137	0.86385945	0.76747574	0.84760712	-1.5613683
11	1	166	2.50492511	1.37171585	-0.1141563	-1.7225671	-0.0523161	-1.2741515	1.36151475	-0.6789088
12	1	170	0.44265243	-0.4562389	-0.0150315	1.58488562	0.05882662	-1.2996843	-1.5015927	-1.0092582

数据要求:

- 列数: 至少需要 7 列以上的数据, 最多 202 列 (200 个变量) 的数据
- 行数: 至少需要 20 个以上的样本(20 行), 暂时支持最多 1500 个以上的样本
- 第 1 列是事件发生情况, 用 0 和 1 (默认, 或 1、2) 表示, 0 表示未发生事件, 1 表示发生了事件。例如, 事件可以定义为死亡, 当受试发生了死亡, 该受试的事件就定义为 1, 当受试未发生死亡 (删失), 该受试的事件就定义为 0
- 第 2 列是具体时间, 时间的具体数值不能都是一样的。另外, 如果样本对应的时间为 0 或小于 0, 这些样本在进行 lasso 分析前都会进行筛选 (lasso 回归要求)
- 第 3 列以及以后的列为准备进行筛选的变量, 变量必须是数值类型



## 参数说明

(说明：标注了颜色的为常用参数。)

## 方法

方法

lambda指标

deviance

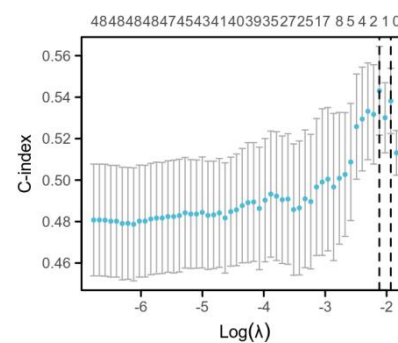
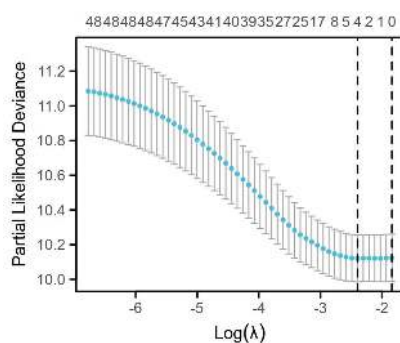
验证方法

十折交叉验证

种子号

2022

- lambda 指标：可以选择 lasso 系数筛选的指标：deviance（似然偏差值）（默认）或 c-index（C 指数），如下：左侧为 deviance，右侧为 c-index
  - 当选择的指标为 deviance 时，y 值（似然偏差值）越小对应的模型越好
  - 当选择的指标为 c-index 时，y 值（C 指数）越大对应的模型越好



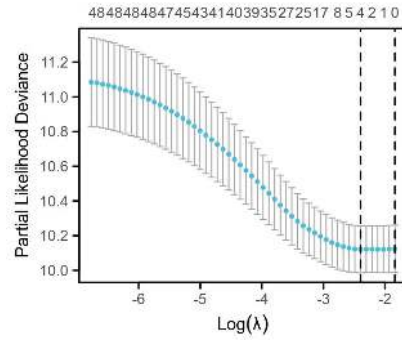
- 验证方法（交叉验证的倍数）：可选三折、五折、七折、十折交叉验证。例如十折交叉验证，就是把数据分成 10 份，轮流把 9 份数据作为训练集训练模型，另外 1 份作为验证集验证模型。默认是选择十折交叉验证，如下：

方法

筛选指标 deviance

交叉验证 十折交叉验证

种子号 2022

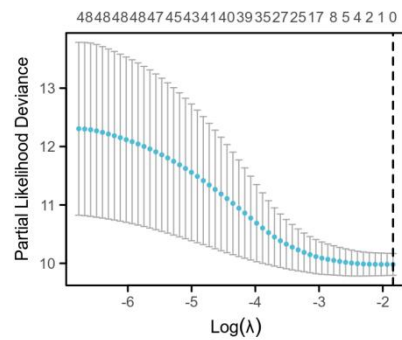


方法

筛选指标 deviance

交叉验证 三折交叉验证

种子号 2022



- 种子号：可填入其他的数字，默认为 2022。由于在进行交叉验证的过程中会涉及到对数据的抽样和分训练集和验证集，故不同的种子号对应的结果都会有不同，但是只要是同一份数据同一个种子号，对应的结果都是一样的

## 点



点

描边色

填充色

样式 圆形

大小 0.5

不透明度 1

- 描边色：可以修改图中点的描边色
- 填充色：可以修改图中点的填充色
- 样式：可以选择图中点样式类型，可选择圆形、正方形、菱形、三角形、倒三角
- 大小：可修改点的大小
- 不透明度：可修改点的透明度。0 为完全透明，1 为完全不透

## 误差线

误差线

颜色

粗细

0.50pt

不透明度

1

- 颜色：可以修改图中误差线(竖线)的颜色，如下：

误差线

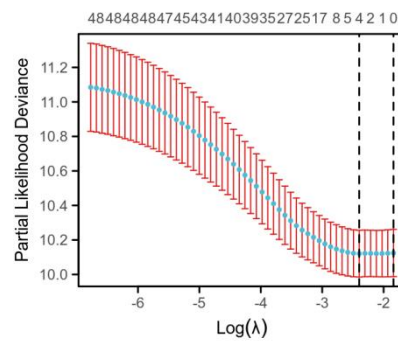
颜色

粗细

0.50pt

不透明度

1



- 粗细：可以修改图中误差线的线条粗细
- 不透明度：可以修改图中误差线的不透明度，1 表示完全不透明，0 表示完全透明

## 标题文本

标题	
大标题	大标题内容
x轴标题	x轴标题内容
y轴标题	y轴标题内容

- 大标题：大标题文本
- x 轴标题：x 轴标题文本
- y 轴标题：y 轴标题文本

补充：在要换行的中间插入\n。如果需要上标，可以用两个英文输入法下的大括号括住，比如  $\{2\}$ ；如果需要下标，可以用两个英文输入法下的中括号括住，比如  $[2]$

## 风格



风格

边框 ☒

网格 ☐

文字大小 7pt

- 外框：是否添加外框，默认添加
- 网格：是否添加网格
- 文字大小：控制整体文字大小，默认为 7pt



## 图片



图片

宽度 (cm) 6

高度 (cm) 5

字体 Arial

- 宽度：图片横向长度，单位为 cm
- 高度：图片纵向长度，单位为 cm
- 字体：可以选择图片中文字的字体

## 结果说明

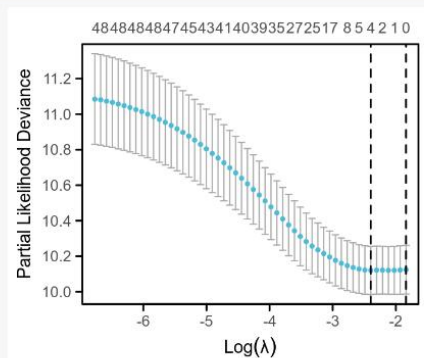
## 主要结果

### 预后lasso系数筛选

预后lasso系数筛选: 可视化lasso系数筛选过程中各个lambda值 (惩罚值) 对应的统计量 (似然偏差值或C指数)

作用: 预后lasso常用于构建预后模型或者筛选变量

交叉验证: 十折交叉验证 || 种子号: 2022



预后lasso系数筛选.pdf

预后lasso系数筛选.tiff

预后lasso系数筛选.pptx

变量-Lasso系数.xlsx

Lasso-RiskScore.xlsx

对于主要结果的解释:

· 下方x轴: 表示lambda对数值 ( $\log(\lambda)$ )

- 变量-系数的 excel 表, 文件内一共有两个 sheet, 其中一个是 lambda.min 对应的变量和系数情况, 另外一个为 lambda.1se 对应的变量和系数的情况
  - 一般 lasso 是看非 0 系数的, 系数为 0 的变量为剔除的变量。
- 包含有预后资料、系数非 0 的变量以及对应的 RiskScore 的 excel 表, 分别 2 个表, 分别对应 lambda.min 和 lambda.1se 对应的情况。可用于: 临床意义-预后-风险因子图的绘制 (只需要修改 RiskScore 的列名); 基础绘图-时间依赖 ROC 曲线 (提取时间资料和 RiskScore 列) 基础绘图-生存曲线 (提取时间资料和 RiskScore 列)

## 补充结果

### Lasso-交叉验证

交叉验证：十折交叉验证

种子号：2022

评价指标(统计量/图中y值)：deviance

	lambda值	Index	统计量	标准误(SE)	系数非0的个数
lambda.min	0.090953	7	10.121	0.1356	4
lambda.1se	0.15894	1	10.124	0.1368	0

说明：

· lambda.min表示：统计量(似然偏差值最小或C指数最大)对应的lambda

· lambda.1se表示：统计量(似然偏差值最小或C指数最大)且在1倍标准误以内对应的lambda

· Index表示：lambda.min与lambda.1se在所有的lambda值中的位置

补充：

· lambda.min与lambda.1se均可作为cutoff，但lambda.min相对严格，lambda.1se对应的变量个数更少，模型相对更简洁

· 模型的变量(表格中的系数非0的个数列)尽量控制在10个左右

· lasso可作为筛选变量的方法，如果筛选出来的变量仍很多，可以对这些筛选出来的变量进一步多因素Cox回归，构建Cox模型

这里提供 lasso-交叉验证表格：可以查看种子号为 2022、指标为 deviance 且使用十折交叉验证方法进行 lasso 系数筛选时的数据信息

- lambda.min 代表：统计量(似然偏差值最小或 C 指数最大)对应的 lambda
- lambda.1se 代表：统计量(似然偏差值最小或 C 指数最大)且在 1 倍标准误以内对应的 lambda
- Index 代表：lambda.min 与 lambda.1se 在所有的 lambda 值中的位置
- lambda.min 与 lambda.1se 均可作为 cutoff，但 lambda.min 相对严格
- lambda.1se 对应的变量个数更少，模型相对更简洁
- 模型的变量(表格中的系数非 0 的个数列)尽量控制在 10 个左右
- lasso 可作为筛选变量的方法，如果筛选出来的变量仍很多，可以对这些筛选出来的变量进一步多因素 Cox 回归，构建 Cox 模型



## 方法学

统计分析和可视化均在 R 4.2.1 版本中进行

涉及的 R 包：survival 包（用于构建模型）、glmnet（用于分析）

处理过程：

- (1) 使用 survival 包对清洗过后的数据进行模型构建，使用 glmnet 包进行分析得到 lambda 值、最大似然数或 C 指数等
- (2) 对数据进行可视化



## 如何引用

生信工具分析和可视化用的是 R 语言，可以直接写自己用 R 来进行分析和可视化即可，可以无需引用仙桃，如果想要引用仙桃，可以在致谢部分 (Acknowledge) 致谢仙桃学术 ([www.xiantao love](http://www.xiantao love))。

方法学部分可以参考对应说明文本中的内容以及一些文献中的描述。



## 常见问题

1. 图里面看不到两根竖（虚）线，只有 1 根？

答：

这个情况说明 lasso 筛选得到的 cutoff 重合了，如果两根竖线都在图的最右侧，非 0 系数的变量个数为 0 或者 1（上方 x 轴），说明 lasso 无法筛选出来变量或无法构建模型。

2. lambda.min 和 lambda.1se 对应的系数非 0 的变量个数为 0，如何才能让结果能好？

答：

由于 lasso 系数筛选过程中涉及到交叉验证和样本抽样的过程，所以不同的种子号是可能会对应不同的一个情况，如果是想要结果更加“好看”一些，可以手动修改种子号。

3. 为什么图上方非 0 系数变量的个数与数据中的变量个数对应不上？为什么看不到所有的数字，只是一小部分？

答：

①图上方的这些数字对应的值是说：不同 lambda 值计算得到模型中所有变量系数不为 0 的变量的个数，而不是所有的变量（要是数值与变量个数对应不上，则

是因为缺少的那些变量间不存在相关关系(系数为 0)被筛选掉了，或者变量在数据处理过程中就被筛选掉了)

②由于可视化结果是 ggplot2 格式，故不能展示全部的数值

