

LARS notes (Part1) (To page9,2 end)

Abstract:

Useful and less greedy version of traditional forward selection methods.

Main property:

- Implement the Lasso, lasso Modification: Calculates all possible Lasso estimates for a given problem.
- Different version: Another modification efficiently implements forward stagewise linear regression.
- A simple approximation for the degree of freedom of a LARS estimate is available, from which we derive a C_p estimate of prediction error. this allows a principled choice among the range of possible LARS estimates.

(Not quite understand the final part of the LARS goals.)

LARS relates: classic model-selection method known as "forward selection" or "forward stepwise regression."

- Forward Selection
 - Given a collection of possible predictors, select the one largest absolute correlation with the response y , say x_{j_1} , and perform simple linear regression of y on x_{j_1} , then leaves a residual vector which is orthogonal to x_{j_1} . Project the other predictors orthogonally to x_{j_1} and repeat the selection process. After k steps this results in a set of predictors x_{j_1}, \dots, x_{j_k} that are then used in the usual way to construct a k -parameter linear model.
 - Forward stagewise
 - More cautious version of forward selection -> take thousands tiny steps as it moves toward a final model.
 - Original motivation for the LARS algorithm.
 - LARS-Lasso-Stagewise connection is conceptually as well as computationally useful.
-

Model construction:

Predict response y from covariates x_1, \dots, x_n .

By location and scale transformations we always assume that the covariates have been standardized to have mean 0 and unit length, and that the response has mean 0.

$$\sum_{i=1}^n y_i = 0 \quad \sum_{i=1}^n x_{ij} = 0 \quad \sum_{i=1}^n x_{ij}^2 = 1 \quad \text{for } j = 1, 2, \dots, m \quad (1)$$

Regression coefficients : $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_m)'$ gives prediction vector $\hat{\mu}$

$$\hat{\mu} = \sum_{j=1}^m \mathbf{x}_j \hat{\beta}_j = X\hat{\beta} \quad [X_{n \times m} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m)] \quad (2)$$

Total squared error

$$S(\hat{\beta}) = \|\mathbf{y} - \hat{\mu}\|^2 = \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 \quad (3)$$

L_1 norm for lasso

$$T(\hat{\beta}) = \sum_{j=1}^m |\hat{\beta}_j| \quad (4)$$

$$\text{Lasso: minimize } S(\hat{\beta}) \text{ subject to } T(\hat{\beta}) \leq t \quad (5)$$

Quadratic programming techniques can be used to solve (5). though we will present an easier method here, closely related to the "homotopy method" of Osborne, Presnell and Turlach (2000a)."

Forward Stagewise Linear Regression.

- Begins with $\hat{\mu} = 0$, builds up the regression function in successive small steps.
- Let $\hat{\mu}$ is the current Stagewise estimate, let $\mathbf{c}(\hat{\mu})$ be the vector of *current correlations*
 $\hat{\mathbf{c}} = \mathbf{c}(\hat{\mu}) = X'(\mathbf{y} - \hat{\mu})$
- \hat{c}_j is proportional to the correlation between covariate x_j and current residual vector. Next step is taken in the direction of the greatest current correlation,

$$\hat{j} = \operatorname{argmax}_j |\hat{c}_j| \quad \text{and} \quad \hat{\mu} \rightarrow \hat{\mu} + \epsilon \cdot \operatorname{sign}(\hat{c}_{\hat{j}}) \cdot \mathbf{x}_{\hat{j}} \quad (6)$$

- Need to mentioned here: ϵ is a "small" constant, "small" is important, otherwise "big" choice like $\epsilon = |\hat{c}_j|$ leads to the standard forward selection technique. this could be over greedy.
-

The main point:

LARS is a stylized version of the stagewise procedure that uses a simple mathematical formula to accelerate the computations.

- Stagewise would choose $\hat{\gamma}_1$ equal to some value ϵ , than repeat many times, or make $\hat{\gamma}_1$ larger enough to make $\hat{\mu}_1$ equal \bar{y}_1 , the projection of y into $\mathcal{L}(\mathbf{x}_1)$.

LARS uses an intermediate value of $\hat{\gamma}_1$, the value that makes $\bar{y}_2 - \hat{\mu}$, equally correlated with x_1 and x_2 ; that is, $\bar{y}_2 - \hat{\mu}_1$ bisects the angle between x_1 and x_2 , so $c_1(\hat{\mu}_1) = c_2(\hat{\mu}_1)$.

- u_2 be the unit vector lying along the bisector. The next LARS estimate is

$$\hat{\mu}_2 = \hat{\mu}_1 + \hat{\gamma}_2 \mathbf{u}_2 \quad (8)$$

With $\hat{\gamma}_2$ chosen to make $\hat{\mu}_2 = \bar{y}_2$ in the case $m=2$. With $m > 2$ covariates, $\hat{\gamma}_2$ would be smaller, leading to another change of direction.

- LARS is motivated by the fact that it is easy to calculate the step size $\hat{\gamma}_1, \hat{\gamma}_2, \dots$, short-circuiting the small Stagewise steps.

We assume that the covariate vector x_1, x_2, \dots, x_m are linearly independent.

For \mathcal{A} a subset of the indices $\{1, 2, \dots, m\}$, define the matrix

$$X_{\mathcal{A}} = (\dots s_j \mathbf{x}_j \dots)_{j \in \mathcal{A}} \quad (9)$$

when the signs s_j equal ± 1 . Let

$$\mathcal{G}_{\mathcal{A}} = X'_{\mathcal{A}} X_{\mathcal{A}} \quad \text{and} \quad A_{\mathcal{A}} = (\mathbf{1}'_{\mathcal{A}} \mathcal{G}_{\mathcal{A}}^{-1} \mathbf{1}_{\mathcal{A}})^{-\frac{1}{2}} \quad (10)$$

$\mathbf{1}_{\mathcal{A}}$ being a vector of 1's of length equaling $|\mathcal{A}|$ the size of \mathcal{A} . The

$$\text{equiangular vector: } u_{\mathcal{A}} = X_{\mathcal{A}} w_{\mathcal{A}} \quad \text{where} \quad w_{\mathcal{A}} = A_{\mathcal{A}} \mathcal{G}_{\mathcal{A}}^{-1} \mathbf{1}_{\mathcal{A}} \quad (11)$$

is the unit vector making equal angles, less than 90° , with the column of $X_{\mathcal{A}}$,

$$X'_{\mathcal{A}} \mathbf{u}_{\mathcal{A}} = A_{\mathcal{A}} \mathbf{1}_{\mathcal{A}} \quad \text{and} \quad \|\mathbf{u}_{\mathcal{A}}\|^2 = 1 \quad (12)$$

We saw the previous part in a negative direction. First, we look the final part (12), which should be satisfied as the equal angular.

We need find a vector $u_{\mathcal{A}}$ satisfied $X' u = a \mathbf{1}$, now note that $\mathbf{1} = X' X (X' X)^{-1} \mathbf{1}$. That is, we have X' in the left hand side. so $u_{\mathcal{A}}$ have a candidate $X (X' X)^{-1} \mathbf{1}$, then what we need to do is just standardize it.

$$u = \frac{X(X'X)^{-1}\mathbf{1}}{\sqrt{\mathbf{1}'(X'X)^{-1}X'X(X'X)^{-1}\mathbf{1}}} \quad (13)$$

$$= \frac{X(X'X)^{-1}\mathbf{1}}{\sqrt{\mathbf{1}'(X'X)^{-1}\mathbf{1}}} \quad (14)$$

$$= \frac{XG^{-1}\mathbf{1}}{\sqrt{\mathbf{1}'G^{-1}\mathbf{1}}} \quad (15)$$

$$= XG^{-1}\mathbf{1} \times A \quad (16)$$

$$= XA^*G^{-1}\mathbf{1} \quad (17)$$

$$= Xw \quad (18)$$

This is how the equal-angular vector is constructed.

So the problem becomes to a nother one, why the form of vector making equal angles?

$$X'_{\mathcal{A}}\mathbf{u}_{\mathcal{A}} = A_{\mathcal{A}}\mathbf{1}_{\mathcal{A}} \quad (19)$$

That is, what is this formula mean?

Equal angular is equals to which angular????

That is, the $\cos(\theta)$ between any subset $x_{i_{\mathcal{A}}}$ and $\cos \langle u, x_{i_{\mathcal{A}}} \rangle = \frac{x'_{i_{\mathcal{A}}} u_{\mathcal{A}}}{\|x'_{i_{\mathcal{A}}} u_{\mathcal{A}}\|}$. In this case,

$X'_{\mathcal{A}}\mathbf{u}_{\mathcal{A}}$ is the vector about $x'_{i_{\mathcal{A}}} u_{\mathcal{A}}$, then if the angular is equal, \cos also should equal. However, $X'_{\mathcal{A}}$ projection to the direction in $u_{\mathcal{A}}$, then right hand side is the length of projection times 1. ($A_{\mathcal{A}}$ is a value rather a matrix.)

The describtion upon is so confusing!!!! need more clear idea about it!!!!

Then is the "fully" describe about the LARS.

- Start with $\hat{\mu}_0 = 0$ and build up $\hat{\mu}$ by steps.

Suppose current estimate is $\hat{\mu}_{\mathcal{A}}$, that $\hat{e} = X'(y - \hat{\mu}_{\mathcal{A}})$, is the current correlations. \mathcal{A} is active set which indices corresponding covariates with the greatest absolute current correlations

$$\hat{C} = \max_i \{|\hat{e}_j|\} \quad \text{and} \quad \mathcal{A} = \{j : |\hat{e}_j| = \hat{C}\} \quad (20)$$

Letting

$$s_j = \text{sign}\{\hat{e}_j\} \quad \text{for} \quad j \in \mathcal{A} \quad (21)$$

we compute $X_{\mathcal{A}}$, $A_{\mathcal{A}}$ and $u_{\mathcal{A}}$ we showed previous, the equal angle trisector and the inner product vector.

$$\mathbf{a} \equiv X'\mathbf{u}_{\mathcal{A}} \quad (22)$$

Then the next step of the LARS algorithm updates $\hat{\mu}_{\mathcal{A}}$, say to

$$\hat{\mu}_{\mathcal{A}+} = \hat{\mu}_{\mathcal{A}} + \hat{\gamma} \mathbf{u}_{\mathcal{A}} \quad (23)$$

where

$$\hat{\gamma} = \min_{j \in \mathcal{A}^c} \left\{ \frac{\hat{C} - \hat{c}_j}{A_{\mathcal{A}} - a_j}, \frac{\hat{C} + \hat{c}_j}{A_{\mathcal{A}} + a_j} \right\} \quad (24)$$

- 还是老样子从 $\hat{\mu}_0 = 0$ 开始，所以还是两件事，计算当前的correlation $\hat{c} = X'(y - \hat{\mu}_{\mathcal{A}})$ 。其中 \mathcal{A} 是一个indices set包括了每一步最大相关系数。这里有点奇怪。特别是 \hat{C} 的定义形式。

Emmmm，是不是可以这么解释，因为开始转向的时候， $c_1 = c_2$ 所以此时往 u_2 的方向走并不会影响 c_1 和 c_2 但是会影响 c_3 。所以等走走走， c_3 又会减小减小。好像不太对，residual的变化。。。如果按figure2的话，residual对 x_1 的角度一直在增大，角度增大导致correlation会减小，因为 x_1 的 γ 在增大， x_1 方向的解释越来越多，而同时，对 x_2 的角度一直在减小，也就是cos在增大，correlation在增大，一减一增两个过程直到这两个correlation相等

然后回到这个过程，也就是在找了一段时间之后，有记录: \mathcal{A} 。注意correlation的定义式：

$\hat{c} = X'(y - \hat{\mu}_{\mathcal{A}})$ X 是 $n \times p$ ， $p \times n \times n \times 1 = p \times 1$ ， \hat{c} 是此时的模型的已建模部分 $\hat{\mu}_{\mathcal{A}}$ 的residual关于 X 的correlation。

按角平分线的思路,这时候在第 i 维的变化是主要的， $i - 1$ 维上应该都是角平分线所以是一致的？

但是 \mathcal{A} 这样定义感觉很奇怪啊。 \mathcal{A} 里面只包含了correlation是最大的那几个，也就是说一直走一直走，走到有新的correlation能加进来，那再改变方向，重新计算 $u_{\mathcal{A}}$ ，否则就一直按这个方向走。比如说回到figure2的图例，在一开始，最大的correlation只有 x_1 这个方向，走一小步， x_2 方向的correlation还是小于 x_1 ，所以 \mathcal{A} 还是只有 $\{1\}$ 。直到residual的correlation到 x_1 和 x_2 一致，那么就有了两个同时达到max的correlation，那就得重新计算 $u_{\mathcal{A}}$ 。

那这么理解就没问题了，于是下一步是通过这些东西计算步长。

- Back to paper understanding mode

Define:

$$\mu(\gamma) = \hat{\mu}_{\mathcal{A}} + \gamma \mathbf{u}_{\mathcal{A}} \quad (25)$$

for $\gamma > 0$, so that the current correlation

$$c_j(\gamma) = x'_j(y - \mu(\gamma)) = \hat{c}_j - \gamma \alpha_j. \quad (26)$$

for $j \in \mathcal{A}$, around the equation (20), the definition of correlation, max correlation, direction, yield

$$|c_j(\gamma)| = \hat{C} - \gamma A_{\mathcal{A}} \quad (27)$$

也就是当前的 j 的correlation，是最大correlation减去步长乘以和投影向量的长度有关的一个玩意 ($A_{\mathcal{A}}$)。注意，这里 $j \in \mathcal{A}$ 。也就和之前描述的，因为走的是角平分线，所以这帮已经active的variable同进退。

然后下一步该考虑的就是不在 \mathcal{A} 里的 j 。

For $j \in \mathcal{A}^c$, the two formula upon shows that $c_j(\gamma)$ equals the maximal value at $\gamma = (\hat{C} - \hat{c}_j)/(A_{\mathcal{A}} - a_j)$. Likewise $-c_j(\gamma)$, the current correlation for the reversed covariate $-x_j$, achieves maximality at $(\hat{C} + \hat{c}_j)/(A_{\mathcal{A}} + a_j)$.

Therefore γ in (24), is the *smallest positive value of γ such that some new index \hat{j} joins the active set*; \hat{j} is the smallest positive value of γ such that some new index \hat{j} joins the active set; \hat{j} is the minimizing index in (24), the foot length of every j in \mathcal{A}^c . the new active set is $\mathcal{A} \cup \{\hat{j}\}$, the new maximum absolute correlation is $\hat{C}_+ = \hat{C} - \hat{\gamma}A_{\mathcal{A}}$.

The figure 10 shows the LARS in diabetes data. 10 iterations for procedure from start to end. The join order or LARS is same as Lasso. However, tracks of $\hat{\beta}_j$ are nearly but not exactly as either the LASSO or Stagewise tracks.

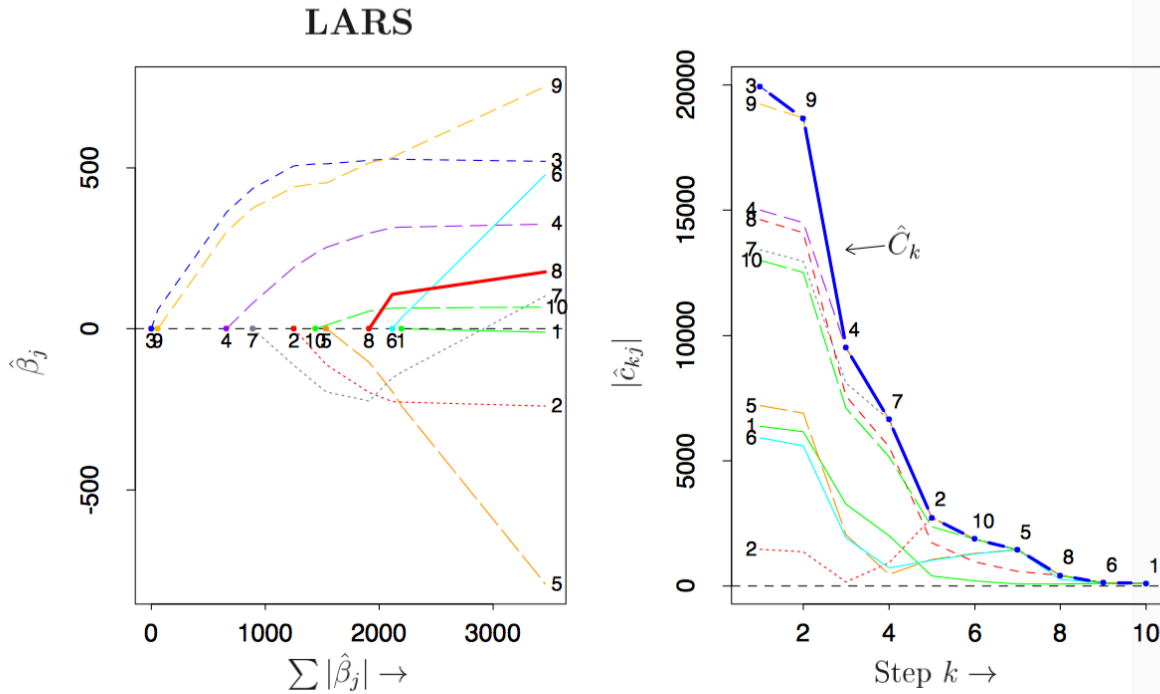


Figure 3. LARS analysis of the diabetes study. *Left:* estimates of regression coefficients $\hat{\beta}_j$, $j = 1, 2, \dots, 10$; plotted versus $\sum |\hat{\beta}_j|$; plot is slightly different than either Lasso or Stagewise, Figure 1. *Right:* Absolute current correlations as function of LARS step; variables enter active set (2.9) in order 3, 9, 4, 7, \dots , 1; heavy curve shows maximum current correlation \hat{C}_k declining with k .

The right panel shows the absolute current correlation goes down with the LARS step k .

$$\hat{C}_k = \max \{|\hat{c}_{kj}|\} = \hat{C}_{k-1} - \hat{\gamma}_{k-1}A_{k-1} \quad (28)$$

Declines with k .

Relation between LARS and OLS.

Suppose LARS has just completed step $k-1$, giving $\hat{\mu}_{k-1}$ and is embarking upon step k . The active set \mathcal{A}_k will have k members, giving X_k, \mathcal{G}_k, A_k and u_k . Similarly, let \bar{y}_k indicate the projection of y into $\mathcal{L}(X_k)$, which, since $\hat{\mu}_{k-1} \in \mathcal{L}(X_{k-1})$, is

$$\bar{y}_k = \hat{\mu}_{k-1} + X_k \mathcal{G}_k^{-1} X_k' (y - \hat{\mu}_{k-1}) = \hat{\mu}_{k-1} + \frac{\hat{C}_k}{A_k} u_k \quad (29)$$

因为等角性质和 A_k 里面的东西在 correlation 上同进退，所以有

$$X_k' (y - \hat{\mu}_{k-1}) = \hat{C}_k 1_{\mathcal{A}} \quad (30)$$

Since u_k is a unit vector, (29), \bar{y}_k 则有 $\bar{y}_k - \hat{\mu}_{k-1}$ 有长度

$$\bar{\gamma}_k \equiv \frac{\hat{C}_k}{A_k} \quad (31)$$

和 update 的那个公式进行比较，则 LARS 的估计 $\hat{\mu}_k$ 在 $\hat{\mu}_{k-1}$ 到 \bar{y}_k 的延长线上。

$$\hat{\mu}_k - \hat{\mu}_{k-1} = \frac{\hat{\gamma}_k}{\bar{\gamma}_k} (\bar{y}_k - \hat{\mu}_{k-1}) \quad (32)$$

可以发现一个问题, $\hat{\gamma}_k$ 总是比 $\bar{\gamma}_k$ 小，所以 LARS estimates always approaching but never reaching the OLS estimates \bar{y}_k .

有一个情况例外，如果 LARS 包含了所有的 covariates，然后。。。反正就和 OLS 等了。

因为一步到位的性质，LARS 算起来特别快。

以上的计算都没有好好看，但是大概意思结论和过程都不复杂所以先放着吧。如果有用再拿起来看。

因为 typero 打了那么多字好卡，所以暂时把文本分割。