

# Why does least angle regression work?

Sven Laur  
swen@math.ut.ee

Helsinki University of Technology

## Minimisation goal of the LASSO algorithm

GIVEN: an output vector  $\mathbf{y}$  and a design matrix  $X$  with columns  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .

FIND: a coefficient vector  $\boldsymbol{\beta}$  that minimises

$$E_{\text{lasso}} = \frac{1}{2} \cdot \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 + \lambda \cdot \|\boldsymbol{\beta}\|_1 \quad (1)$$

EQUIVALENT FORMULATION: Find a coefficient vector  $\boldsymbol{\beta}$  that minimises

$$E_{\text{ols}} = \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\beta}\|_1 \leq t \quad (2)$$

CORRESPONDENCE: Task (1) is Lagrange functional of Task (2).

## Explicit gradient of the cost function

Divide the set of feasible solutions  $\mathbb{R}^n$  into octants  $\text{sign}(\beta_i) = \text{const.}$

Let  $\mathbf{s}$  be the sign vector, i.e.  $s_i = \text{sign}(\beta_i)$ . Then in each octant

$$E_{\text{lasso}} = \frac{1}{2} \cdot (\mathbf{y} - X\boldsymbol{\beta})^2 + \lambda \cdot \mathbf{s}^T \boldsymbol{\beta}$$
$$\nabla_{\boldsymbol{\beta}} E_{\text{lasso}} = X^T X \boldsymbol{\beta} - X^T \mathbf{y} + \lambda \cdot \mathbf{s}$$

If the minimum is an internal point, then the solution has a form

$$\boldsymbol{\beta}^* = (X^T X)^{-1} (X^T \mathbf{y} - \lambda \mathbf{s})$$

## What happens in the boundaries?

For the minimisation over a boundary, we explicitly require

$$\beta_i = 0 \quad \text{for all} \quad i \in \mathcal{N}$$

$$\beta_i \in \mathbb{R} \quad \text{for all} \quad i \in \mathcal{A}$$

Hence, the cost function simplifies

$$E_{\text{lasso}} = \frac{1}{2} \cdot (\mathbf{y} - X_{\mathcal{A}}\boldsymbol{\beta}_{\mathcal{A}})^2 + \lambda \cdot \mathbf{s}_{\mathcal{A}}^T \boldsymbol{\beta}_{\mathcal{A}}$$

$$\nabla_{\boldsymbol{\beta}_{\mathcal{A}}} E_{\text{lasso}} = X_{\mathcal{A}}^T X_{\mathcal{A}} \boldsymbol{\beta}_{\mathcal{A}} - X_{\mathcal{A}}^T \mathbf{y} + \lambda \cdot \mathbf{s}_{\mathcal{A}}$$

and thus

$$\boldsymbol{\beta}_{\mathcal{A}}^* = (X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} (X_{\mathcal{A}}^T \mathbf{y} - \lambda \mathbf{s}_{\mathcal{A}})$$

$$\boldsymbol{\beta}_{\mathcal{N}}^* = \mathbf{0}$$

## Geometrical interpretation of $\beta^*$

If  $\beta^*$  is an internal point, then the corresponding prediction vector

$$\boldsymbol{\mu} = X\boldsymbol{\beta}^* = X(X^T X)^{-1}(X^T \mathbf{y} - \lambda \mathbf{s}) = \boldsymbol{\mu}_{\text{ols}} - \lambda \cdot \underbrace{X(X^T X)^{-1} \mathbf{s}}_{\mathbf{u}}$$

where  $\mathbf{u}$  is an equiangular to the vectors  $s_1 \mathbf{x}_1, \dots, s_n \mathbf{x}_n$

$$X^T \mathbf{u} = X^T X (X^T X)^{-1} \mathbf{s} = \mathbf{s} = (\pm 1, \dots, \pm 1)^t$$

To summarise, a small change in  $\lambda$  moves  $\boldsymbol{\mu}$  in the direction of  $\mathbf{u}$ .

## What happens in the boundaries?

Let  $\beta^*$  be the internal point of a boundary with working set  $\mathcal{A}$ , i.e.

$$\beta_i = 0 \quad \text{for all } i \in \mathcal{N}$$

$$\beta_i \neq 0 \quad \text{for all } i \in \mathcal{A}$$

Then the corresponding prediction vector

$$\mu = X_{\mathcal{A}}\beta_{\mathcal{A}}^* = X_{\mathcal{A}}(X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1}(X_{\mathcal{A}}^T \mathbf{y} - \lambda \mathbf{s}_{\mathcal{A}}) = \mu_{\mathcal{A}} - \lambda \cdot \underbrace{X_{\mathcal{A}}(X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} \mathbf{s}_{\mathcal{A}}}_{\mathbf{u}_{\mathcal{A}}}$$

where  $\mathbf{u}_{\mathcal{A}}$  is equiangular to the vectors  $s_i \mathbf{x}_i$ ,  $i \in \mathcal{A}$

$$X_{\mathcal{A}}^T \mathbf{u}_{\mathcal{A}} = X^T X_{\mathcal{A}}(X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} \mathbf{s}_{\mathcal{A}} = \mathbf{s}_{\mathcal{A}} = (\pm 1, \dots, \pm 1)^t$$

To summarise, a small change in  $\lambda$  moves  $\mu$  in the direction of  $\mathbf{u}_{\mathcal{A}}$ .

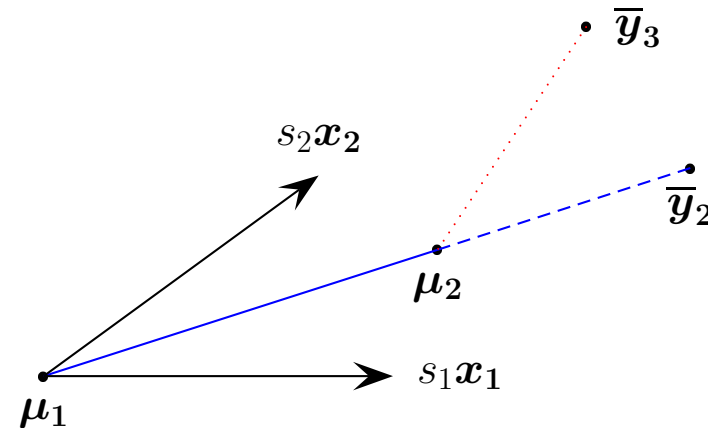
## Informal description of LARS

LARS is a greedy optimisation algorithm:

- Starts from the extreme boundary:  $\mathcal{A} = \emptyset$ ,  $\beta_0 = \mathbf{0}$  and  $\mu_0 = \mathbf{0}$ .
- Moves along the “optimal” path in space vector space  $\langle \mathbf{x}_i, i \in \mathcal{A} \rangle$ .
- Occasionally, extends to higher dimension.
- Always chooses the most profitable vector  $\mathbf{x}_i$  to add.
- Finally, reaches the ordinary least squares solution.

W.l.o.g. we can assume that the working set  $\mathcal{A}_k = \{1, \dots, k\}$ .

## The LARS path. Steady phase



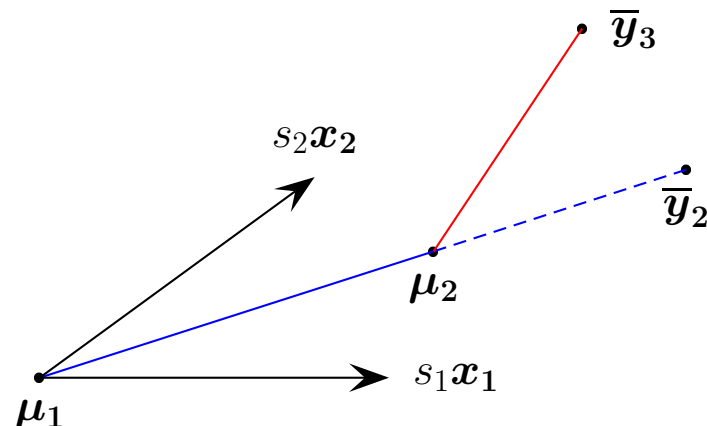
Equiangular vector  $\mathbf{u}_k = \bar{\mathbf{y}}_k - \boldsymbol{\mu}_{k-1}$  and  $\boldsymbol{\mu}(\gamma)$  moves along  $\boldsymbol{\mu}_k \rightarrow \boldsymbol{\mu}_{k+1}$

$$\boldsymbol{\mu}(\gamma) = \boldsymbol{\mu}_{k-1} + \gamma \mathbf{u}_k = \bar{\mathbf{y}}_k - \gamma \mathbf{u}_k$$

Parameter  $\lambda$  decreases and  $t$  increases in the path.



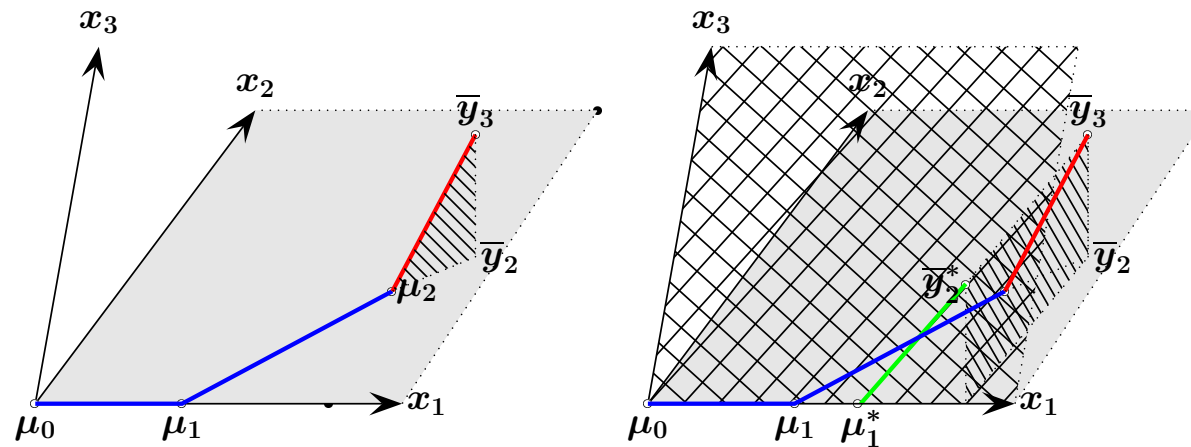
## The LARS path. Regime change



Two paths  $\boldsymbol{\mu}(\gamma) = \mathbf{y}_k - \gamma \mathbf{u}_k$  and  $\boldsymbol{\mu}(\gamma) = \bar{\mathbf{y}}_{k+1} - \gamma \mathbf{u}_{k+1}$  intersect at  $\boldsymbol{\mu}_k$ .

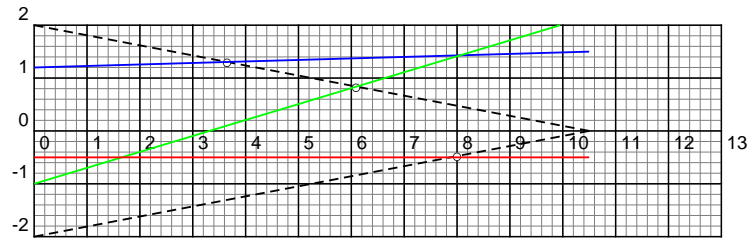
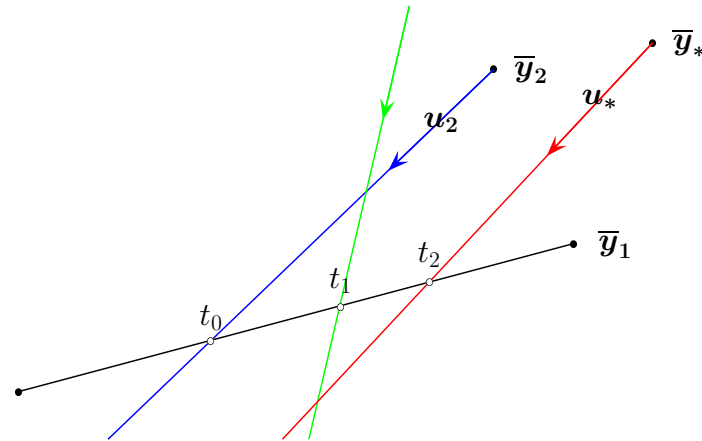
- Therefore,  $\mathbf{c}(\gamma) = X^T(\mathbf{y} - \boldsymbol{\mu}(\gamma))$  must have  $k + 1$  equal components.

## The LARS path. Greedy nature



The LARS algorithm chooses the most advantageous dimension to extend.

# Guessing the correct take-off point



The first  $c_j(\gamma)$  that intersects with the boundary reveals the next vector.

## Guessing the correct take-off point

Consider a the LARS in steady phase. The path point  $\boldsymbol{\mu}(\gamma) = \boldsymbol{\mu}_{k-1} + \gamma \boldsymbol{u}_k$  can belong to the optimal path for  $\langle \boldsymbol{x}_1, \dots, \boldsymbol{x}_{k+1} \rangle$  iff

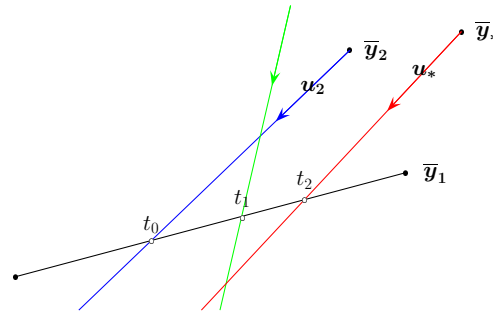
$$\begin{aligned}c_{k+1} &= \boldsymbol{x}_{k+1}^T (\boldsymbol{y} - \boldsymbol{\mu}(\gamma)) = \boldsymbol{x}_{k+1}^T (\bar{\boldsymbol{y}}_{k+1} - \boldsymbol{\mu}(\gamma)) = c_{k+1} - \gamma a_{k+1} \\c_k &= \boldsymbol{x}_k^T (\boldsymbol{y} - \boldsymbol{\mu}(\gamma)) = \boldsymbol{x}_k^T (\bar{\boldsymbol{y}}_{k+1} - \boldsymbol{\mu}(\gamma)) = c^* - \gamma a_k\end{aligned}$$

are equal.

The LARS algorithm chooses the next vector in a greedy way

$$\begin{aligned}\hat{\gamma} &= \min \{ \gamma > 0 : |c^* - \gamma a_k| = |c_j - \gamma a_j|, j \in \mathcal{N} \} \\ \hat{j} &= \operatorname{argmin} \{ \gamma > 0 : |c^* - \gamma a_k| = |c_j - \gamma a_j|, j \in \mathcal{N} \}\end{aligned}$$

## Guessing the correct take-off point



We have to minimise  $\|\mathbf{y} - X\boldsymbol{\beta}\|_2^2$  w.r.t.  $\|\boldsymbol{\beta}\|_1 \leq t$

- Moving along the line increases  $t$ .
- Choosing the first take-off point guarantees that, we remain on optimum line after the direction change.
- Formally, for any  $t$  the covariance vector  $\mathbf{c}(t) = X^T(\mathbf{y} - \boldsymbol{\mu}(t))$  satisfies

$$|c_j(t)| \leq |c_i(t)| \quad i \in \mathcal{A} \quad j \in \mathcal{N}$$

## What about signs of $\beta_{k+1}$ ?

If the LARS extends to next dimension it must correctly guess the signs

$$\boldsymbol{\mu}(\lambda) = \boldsymbol{\mu}_{\mathcal{A}} - \lambda \cdot \underbrace{X_{\mathcal{A}}(X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} \mathbf{s}_{\mathcal{A}}}_{\mathbf{u}_{\mathcal{A}}}$$

- Sign variables  $s_i$  for  $i = 1, \dots, k$  are known from previous step.
- In steady phase  $|\beta_{k+1}(\lambda)|$  grows monotonically—caused by equiangularity.

$$X^T \underbrace{X X^T (X^T X)^{-1} s_j \mathbf{e}_j}_{\mathbf{v}_j} = s_j \mathbf{e}_j \quad \Rightarrow \quad \mathbf{x}_i \perp \mathbf{v}_j \quad \text{and} \quad s_j \mathbf{x}_j \uparrow \uparrow \mathbf{v}_j$$

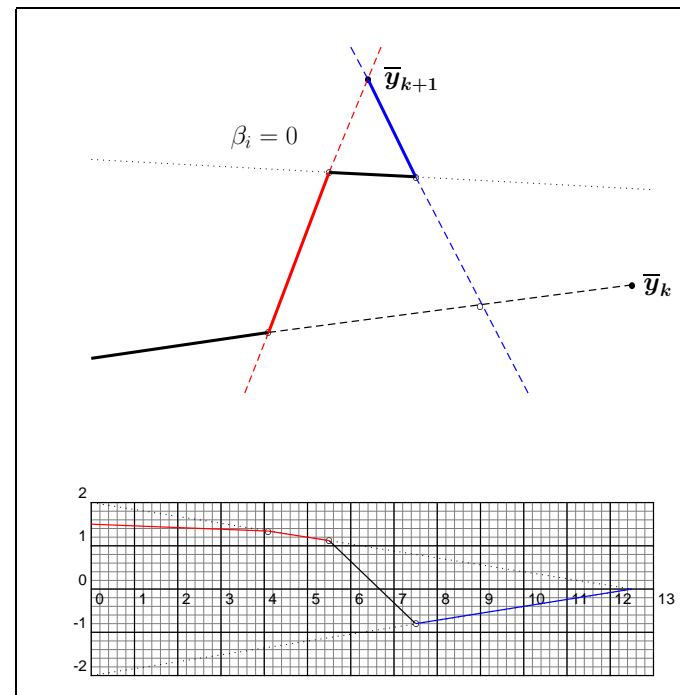
$$\mathbf{u}_{k+1} = \mathbf{v}_1 + \dots + \mathbf{v}_{k+1}$$

I.e.  $c_{k+1} = \mathbf{x}_{k+1}^T (\mathbf{y} - \boldsymbol{\mu}_k)$  has the same sign as  $s_{k+1}$ .

- Unfortunately, this is not true for other coordinates.

## When does the LARS and LASSO coincide

If the design matrix is orthogonal, we might be a sign problem.



The take-off point is correct, but there are more turns in the path.

## Quick fix to LARS algorithm

Check for sign changes:

- Compute  $\beta_k$  and  $\beta_{k+1}$
- If there is no sign change, i.e.  $\beta_{k,i}\beta_{k+1,i} \geq 0$ , proceed as usual.
- Otherwise, find largest intermediate vector  $\beta$  such that  $\beta_{k,i}\beta_{k+1,i} \geq 0$ .
  - Find corresponding  $\mu$  and store it as  $\mu_{k+1}$ .
  - Remove  $k + 1$  from the working list. Recompute direction  $u$ .
  - Proceed with the next LARS step.



## Quick recap to Stagewise algorithm

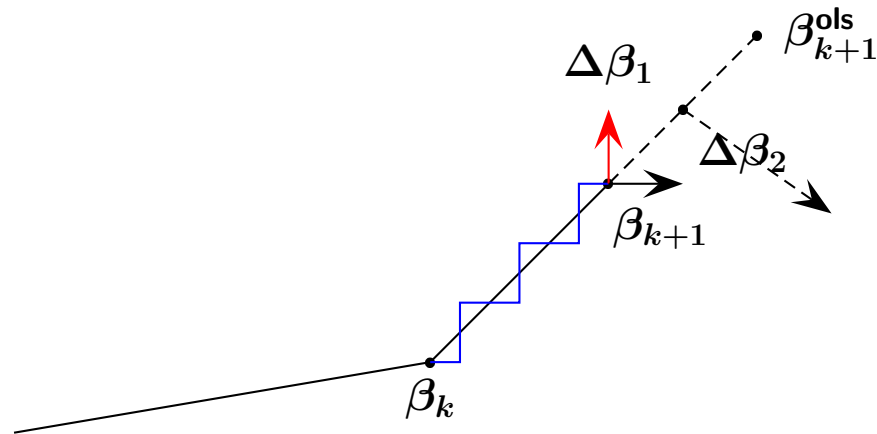
Let  $\varepsilon > 0$  be small enough (infinitely small).

- Choose the coordinate  $i$  that has the biggest impact on squared error. Make a  $\varepsilon$ -step in appropriate direction towards  $c_i$ .

$$\nabla_{\beta} E_{\text{sqe}} = X^{\text{T}}(X\beta - \mathbf{y}) = X^{\text{T}}(\boldsymbol{\mu} - \mathbf{y}) = -\mathbf{c}$$

- If we take infinitesimal steps, we follow the minimising path, except:
  - active correlations  $c_j$  and  $\Delta\beta_j$  have same sign.
- **Fix to the LARS algorithm.** When a  $j$ th coordinate of  $\mathbf{u}_{k+1}$  has different sign than  $c$ , determine active Stagewise coordinates with a projection. Update working list  $\mathcal{A}$  and the vector  $\mathbf{u}_{\mathcal{A}}$ .

## Why does the fix work?



The Stagewise algorithm with infinitesimal  $\epsilon > 0$  assures that

$$|c_j(t)| \leq |c_i(t)| \quad i \in \mathcal{A} \quad j \in \mathcal{N} \quad \text{and} \quad \Delta\beta_i c_i \geq 0$$

Projection is a clever way to determine active working set.

## Final remarks

- The LASSO algorithm minimises the true objective, but sometimes make more steps than the LARS algorithm.
- The LARS skips several LASSO steps. Hopefully, the LASSO and LARS paths are different for small regions.
- The Stagewise algorithm provides the most heuristic approach, but is more widely applicable.
- For large datasets (design matrices) they all perform relatively similarly.

The column vectors of the large design matrix are almost orthogonal *with high probability*.