# Distributed Lasso
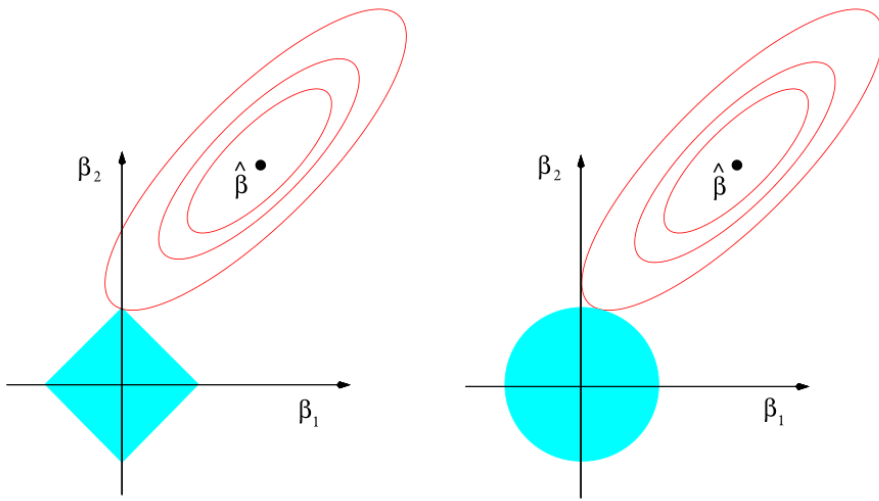
SEBASTIEN { DUBOIS, LEVY }
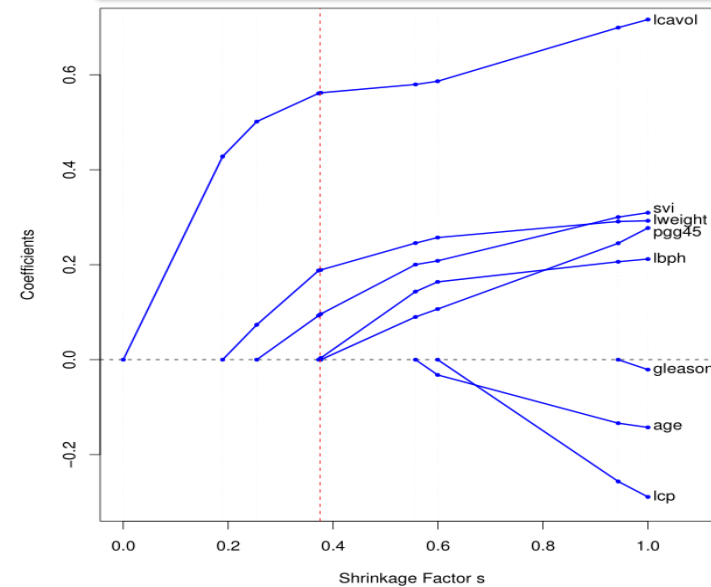
# The Lasso

$$\hat{\beta}^{\text{lasso}} = \operatorname*{argmin}_{\beta} \left\{ \frac{1}{2} \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j| \right\}$$

## Sparsity

## Coefficient Path

# Existing distributed methods

## By SGD (Spark)

- Sparsity not guaranteed
- Problem when $p \gg n$
- Coefficient path ?

## Shotgun (Distributed Coordinate Descent)

- Too much communication
- Convergence issues (need locks)
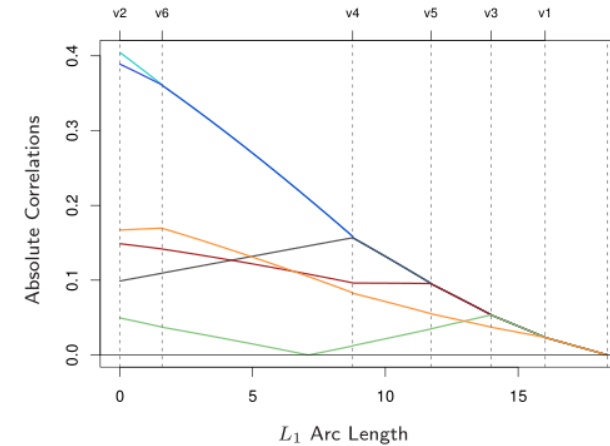- $n \gg p$ ? (Data stored by column)

# Least Angle Regression (LAR)

**Algorithm**

- Relax the penalty at each iteration

- Add features most correlated with the residual

- Increase coefficients to keep correlation tied



*Correlations with residual*

**Computations**

- Dot products between columns

- `k x k` matrix inversion

- `min(n,p)` iterations

$$\delta_k \leftarrow (X_{A_k}^T X_{A_k})^{-1} X_{A_k}^T r_k$$

*New coefficient direction*

# Distributed LAR

| Tall and Skinny $p \ll n$ | Small and Fat $p \gg n$ | Almost Square $p \sim n$ |
|---|---|---|
| • $p^2$ fits in memory, $n$ does not | • $n^2$ fits in memory, $p$ does not | • $p$ and n fit in memory, $n^2$ and $p^2$ do not |
| • Data stored by rows | • Data stored by columns | • If $I < \sqrt{n}$ both methods work |
| • $B$ machines, I iterations:<br>  • Communication:<br>    $O(p^2 B)$<br>  • Computations:<br>    $O\left(pI\left(p + \log(B) + \frac{n}{B}\right)\right)$ | • $B$ machines, I iterations:<br>  • Communication:<br>    $O(n^2 B)$<br>  • Computations:<br>    $O\left(nI\left(n + \log(B) + \frac{p}{B}\right)\right)$ | • More iterations:<br>  • Inversion using distributed block matrix<br>  • SGD to solve linear system<br>  • Forward Stagewise |