# Very high-dimensional data: greedy boosting and convex Lasso-optimization

**Peter Bühlmann**

**ETH Zürich**

## 1. High-dimensional data

$(X_1, Y_1), \ldots, (X_n, Y_n)$ i.i.d. or stationary

$X_i \in \mathbb{R}^p$ predictor variable

$Y_i$ univariate response variable, e.g. $Y_i \in \mathbb{R}$ or $Y_i \in \{0, 1\}$

high-dimensional: $p \gg n$

areas of application: astronomy, biology, imaging, marketing research, text classification,...

# High-dimensional linear models

$$Y_i = \beta_0 + \sum_{j=1}^{p} \beta_j X_i^{(j)} + \varepsilon_i, \ i = 1, \ldots, n$$

$$p \gg n$$

How should we fit this model?

approaches include:

Ridge regression (Tikhonov regularization); variable selection via AIC, BIC, gMDL

(in a forward manner); Bayesian methods for regularization, ...

Boosting, Lasso, ...

our requirements:

- computationally feasible
- yields variable selection
- statistically accurate for prediction or selecting the correct variables

computational feasibility for high-dimensional problems

$$\rightsquigarrow$$

greedy methods

or

convex optimization

# 2. Greedy is good for $p \gg n$: $L_2$Boosting

specify a base procedure ("weak learner"):

$$\text{data} \quad \xrightarrow{\text{algorithm A}} \quad \hat{\theta}(\cdot) \quad \text{(a function estimate)}$$

e.g.: simple linear regression, tree (CART), ...

$L_2$Boosing with base procedure $\hat{\theta}(\cdot)$: repeated fitting of residuals

$$m = 1: \; (X_i, Y_i)_{i=1}^n \; \rightsquigarrow \hat{\theta}_1(\cdot), \; f_1 = \underbrace{\nu}_{\text{e.g.} = 0.1} \hat{\theta}_1 \; \rightsquigarrow \text{resid.} \; U_i = Y_i - f_1(X_i)$$

$$m = 2: \; (X_i, U_i)_{i=1}^n \; \rightsquigarrow \hat{\theta}_2(\cdot), \; f_2 = f_1 + \nu\hat{\theta}_2 \; \rightsquigarrow \text{resid.} \; U_i = Y_i - f_2(X_i)$$

$$\cdots \qquad\qquad \cdots$$

$$f_{m_{stop}}(\cdot) = \nu \sum_{m=1}^{m_{stop}} \hat{\theta}_m(\cdot) \quad \text{(greedy fitting of residuals)}$$

Tukey (1977): twicing for $m_{stop} = 2$ and $\nu = 1$

Componentwise linear least squares base procedure:
linear OLS regression against the one predictor variable which reduces residual
sum of squares most

$$\hat{\theta}(x) = \hat{\beta}_{\hat{S}} x^{(\hat{S})},$$

$$\hat{\beta}_j = \sum_{i=1}^n Y_i X_i^{(j)} \Big/ \sum_{i=1}^n (X_i^{(j)})^2, \quad \hat{S} = \arg\min_j \sum_{i=1}^n (Y_i - \hat{\beta}_j X_i^{(j)})^2$$

$L_2$ Boosting with componentwise linear LS:

first round of estimation: selected predictor variable $X^{(\hat{S}_1)}$ (e.g. $= X^{(3)}$)

corresponding $\hat{\beta}_{\hat{S}_1}$

use shrunken fit $\hat{f}_1(x) = \nu \hat{\beta}_{\hat{S}_1} x^{(\hat{S}_1)}$ (e.g. $\nu = 0.1$)

second round of estimation: selected predictor variable $X^{(\hat{S}_2)}$ (e.g.$= X^{(21)}$)

corresponding $\hat{\beta}_{\hat{S}_2}$

use shrunken fit $\hat{f}_2(x) = \hat{f}_1(x) + \nu \hat{\beta}_{\hat{S}_2} x^{(\hat{S}_2)}$

etc.

for $\nu = 1$, this is known as

Matching Pursuit (Mallat and Zhang, 1993)

Weak greedy algorithm (deVore & Temlyakov, 1997)

a version of Boosting (Schapire, 1992; Freund & Schapire, 1996)

Gauss-Southwell algorithm



C.F. Gauss in 1803

"Princeps Mathematicorum"



R.V. Southwell in 1933

Professor in engineering, Oxford

## Properties

variable selection

shrinkage towards zero for coefficients of selected variables

$\rightsquigarrow$ often much better performance than OLS on selected variables

("more stable" in Breiman's terminology)

"similar" to the Lasso (Efron, Hastie, Johnstone & Tibshirani, 2004)

but not the same

computational complexity:

$O(npm_{stop}) = O(p)$ if $p \gg n$, i.e. linear in dimension $p$

8

statistically consistent for very high-dimensional, sparse problems

Theorem (PB, 2004)

$L_2$Boosting with comp. linear LS regression is consistent (for suitable number of boosting iterations) if:

- $p_n = O(\exp(Cn^{1-\xi}))\ (0 < \xi < 1)$ (high-dimensional)

  essentially exponentially many variables relative to $n$

- $\sup_n \sum_{j=1}^{p_n} |\beta_{j,n}| < \infty$ $\ell^1$-sparseness of true function

i.e. for suitable, slowly growing $m = m_n$:

$$\mathbf{E}_X |f_{m_n,n}(X) - f_n(X)|^2 = o_P(1)\ (n \to \infty)$$

"no" assumptions about the predictor variables/design matrix

in other words:

consistency for de-noising sparse signal with highly over-complete dictionaries

binary lymph node classification in breast cancer using gene expressions:

a high noise problem

$n = 49$ samples, $p = 7130$ gene expressions

| CV-misclassif.err. | $L_2$Boosting | FPLR | Pelora | 1-NN | DLDA | SVM |
|---|---|---|---|---|---|---|
| | 24.8% | 35.25% | 27.8% | 43.25% | 36.12% | 36.88% |

$L_2$Boosting, Forward Penalized Logistic Regression (FPLR), Supervised Gene Grouping (Pelora)

no gene pre-selection ⤳ all these methods do multivariate gene selection

Nearest Neighbor (1-NN), Diagonal Linear Discriminant Analysis (DLDA), SVM with radial basis kernel

with gene pre-selection: the best 200 genes from 2-sample Wilcoxon score

⤳ no additional gene selection anymore

$L_2$Boosting selected 42 out of $p = 7129$ genes

for this data-set: not good prediction with all the different methods

(although we will improve to 16.3%)

but $L_2$Boosting may be a good(?) multivariate gene selection method

10

do variable selection such that predictive performance is good

(not necessarily optimal)

"classical": subset selection using BIC, AIC, gMDL, etc.

computationally infeasible for high-dimensional problems

remedies:

● forward selection

but often not competitive in terms of predictive performance

● $L_2$Boosting: seems quite interesting, but weak theoretical basis

● replace the computationally hard subset selection problem ($2^p$ sub-models)

by convex relaxation

## 3. Lasso-relaxation is good for $p \gg n$

consider again linear model (or highly overcomplete dictionary)

$$Y = f(X) + \varepsilon, \quad f(x) = \sum_{j=1}^{p} \beta_j x^{(j)}, \quad p \gg n$$

Lasso or $\ell^1$-penalized regression (Tibshirani, 1996):

$$\hat{\beta}_{Lasso} = \operatorname{argmin}_{\beta} \sum_{i=1}^{n} (Y_i - \sum_{j=1}^{p} \beta_j X_i^{(j)})^2 + \underbrace{\lambda}_{\geq 0; \text{ penalty par.}} \sum_{j=1}^{p} |\beta_j|$$

- does variable selection: some (many) $\beta_j$'s exactly equal to 0
- does shrinkage
- involves a convex optimization only

**this is convex relaxation:**

replace the computationally hard/infeasible subset selection ($\ell^0$-penalty)

by the convex $\ell^1$-penalized problem

"similar" properties of convex relaxation (Lasso) and greedy algorithm (Boosting):
variable selection; shrinkage;

consistency for prediction in high-dimensions (Greenshtein & Ritov (2004))

and indeed: there are relations

Efron, Hastie, Johnstone, Tibshirani (2004): for special design matrices,

iterations of $L_2$Boosting with "infinitesimally" small $\nu$

yield all Lasso solutions when varying $\lambda$

$\rightsquigarrow$ computationally interesting to produce all Lasso solutions in

one sweep of boosting

Least Angle Regression LARS (Efron et al., 2004) is computationally even more

clever and efficient than $L_2$Boosting

$O(np\min(n,p))$ essential operations to compute all Lasso solutions

$$= O(p) \text{ if } p \gg n$$

Zhao and Yu (2005): in general, when adding some backward step
the solutions from Lasso and Boosting coincide

greedy (plus backward steps) and convex relaxation are surprisingly similar

15

# 3.1. Variable selection and graphical modeling with the Lasso

<u>goal</u>: use the Lasso for determining presence/absence of associations between random variables (this includes regression)

assume that $X = X^{(1)}, \ldots, X^{(p)} \sim \mathcal{N}_p(\mu, \Sigma)$

| Gaussian conditional independence graph |

graph:

set of nodes $\Gamma = \{1, 2, \ldots, p\}$, corresponding to the $p$ random variables

set of edges $E \subseteq \Gamma \times \Gamma$ defined as:

there is an undirected edge between node $i$ and $j$

$\overset{\text{def}}{\Leftrightarrow}$ $X^{(i)}$ conditionally dependent of $X^{(j)}$ given all other $\{X^{(k)}; \ k \neq i, j\}$

$\Leftrightarrow$ $\Sigma^{-1}_{ij} \neq 0$

note: $\Sigma_{ij}^{-1}$ corresponds to $\beta_j^{(i)} = \Sigma_{ij}^{-1}/\Sigma_{ii}^{-1}$, where

$$X^{(i)} = \beta_j^{(i)} X^{(j)} + \sum_{k \neq i,j} \beta_k^{(i)} X^{(k)} + \text{error}^{(i)}$$

$\rightsquigarrow$ <span style="color:red">we can infer the graph from variable selection in regression</span>

$$\beta_j^{(i)} = 0 \Leftrightarrow \Sigma_{ij}^{-1} = 0$$

huge computational problem when using e.g. BIC: $p2^{p-1}$ least squares problems!

## Just relax!

replace the computationally **hard** problem by a **convex** problem:

compute the Lasso estimates $\hat{\beta}_i^{(j)}$

### Estimation of graph:

estimate an edge between node $i$ and $j$ if

$$\hat{\beta}_j^{(i)} \neq 0 \text{ and } \hat{\beta}_i^{(j)} \neq 0$$

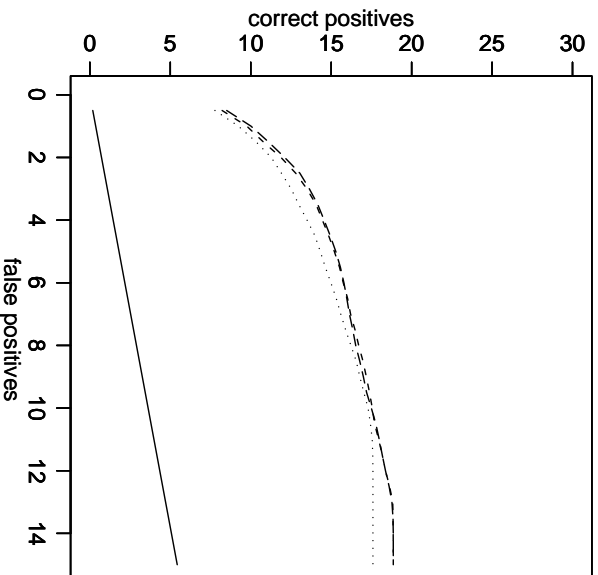(for finite samples: it could happen that only one of the $\hat{\beta}_j^{(i)}$, $\hat{\beta}_i^{(j)}$ is $\neq 0$)

note: depends on the tuning parameter $\lambda$ in Lasso

this involves only one convex optimization problem!

instead of checking exhaustively $2^{p-1} p$ least squares problems (e.g. using BIC)

# Comparison of Lasso and classical stepwise selection

$p = 10$

$p = 30$



dotted · · · ·    stepwise selection

dashed – – –    Lasso

ROC-curves for estimated graphs with $p = 10, 30$ nodes and $n = 40$ obs.

true graphs are sparse, having at most 4 edges out of every node

## Some theory for high dimensions

**Theorem** (Meinshausen & PB, 2004)

For $\lambda_n \sim C n^{-1/2+\delta/2}$,

$$\mathbb{P}[\text{estimated graph}(\lambda_n) = \text{true graph}] = 1 + O(\exp(-Cn^\delta)) \quad (n \to \infty)$$

$$(0 < \delta < 1)$$

if

- Gaussian data
- $p = p_n = O(n^r)$ for any $r > 0$ (high-dimensional)
- maximal number of edges out of a node $= O(n^k)$ $(0 < k < 1)$ (sparseness)
- plus some other technical conditions

justification for relaxation with a computationally simple convex problem!

## Choice of $\lambda$

Theorem doesn't say much about choosing $\lambda$...

first (not so good) idea: choose $\lambda$ to optimize prediction

e.g. via some cross-validation scheme

but: for prediction oracle solution

$$\lambda^* = \arg\min_{\lambda} \mathbb{E}[(X^{(i)} - \sum_{j \neq i} \hat{\beta}_j^{(i)}(\lambda) X^{(j)})^2]$$

asymptotically: the prediction optimal graph is too large

$$\mathbb{P}[\text{estimated graph}(\lambda^*) = \text{true graph}] \to 0 \ (p_n \to \infty, n \to \infty)$$

(Meinshausen & PB, 2004; related example by Meng et al., 2004)

we have a simple proposal for choosing the penalty parameter

which avoids connecting distinct connectivity components

## 4. Beyond Lasso (and Boosting)

linear model $Y = X\beta + \varepsilon$

for orthonormal design: $\mathbf{X}^T \mathbf{X} = I$: Lasso/LARS and $L_2$Boosting yield the

soft-threshold estimator:

$$\hat{\beta}_{soft}^{(j)} = \begin{cases} Z_j - \lambda, & \text{if } Z_j \geq \lambda, \\ 0, & \text{if } |Z_j| < \lambda, \\ Z_j + \lambda, & \text{if } Z_j \leq -\lambda. \end{cases} \qquad \text{where } Z_j = (\mathbf{X}^T \mathbf{Y})_j$$

**Is soft-thresholding or Lasso a good thing?**

- $\beta_1, \ldots \beta_p$ i.i.d. $\sim$ Double-Exponential, soft-thresholding and the Lasso yield the MAP (which often performs well)

- minimax results for soft-thresholding (Donoho & Johnstone, ...)

**but**: a different story in the very high-dimensional sparse case

assume:

- $p = p_n \sim C_1 \exp(C_2 n^{1-\xi}) \; (0 < \xi < 1)$
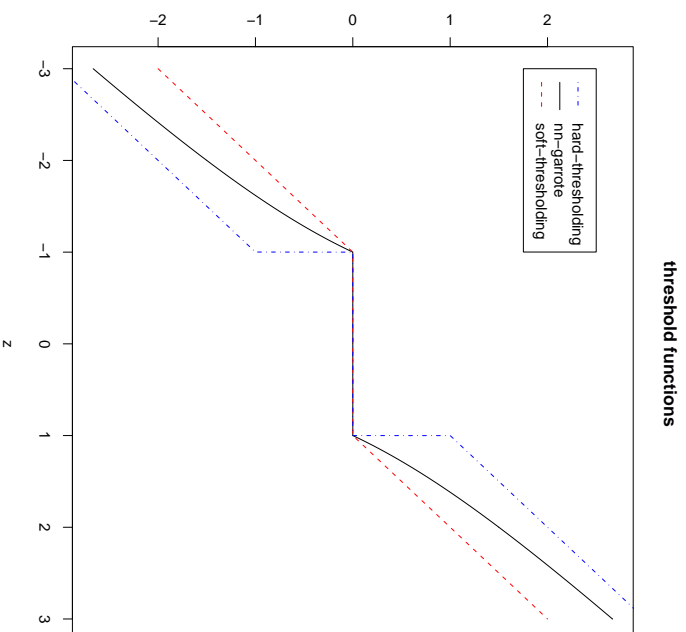- effective number of variables is finite (finite $\ell^0$-norm)

  non-effective variables are independent

Theorem (Meinshausen, 2005)

$$\mathbb{P}[\inf_{\lambda} \underbrace{L(\lambda)}_{\text{risk of Lasso}} > cn^{-r}] \to 1 \; (n \to \infty) \text{ for } r > \xi$$

while optimal rate is $n^{-1}$ (achieved e.g. by OLS with the true variables)

⤳ Lasso can have very poor convergence rate

reason: need large $\lambda$ for variable selection $\rightsquigarrow$ strong bias of soft-thresholding



threshold functions

Better:

- SCAD (Fan and Li, 2001)
- Nonnegative Garrote (Breiman, 1995)
- Bridge estimation
  (Frank and Friedman, 1993)

they all work for general $\mathbf{X}$

for non-orthogonal $\mathbf{X}$:

- non-convex optimization for SCAD or Bridge estimation
- NN-Garrote only for $p \leq n$

# 4.1. The relaxed Lasso (Meinshausen, 2005)

for $\lambda \geq 0$, $0 \leq \phi \leq 1$

$$\hat{\beta}_{\lambda,\phi} = \arg\min_{\beta} n^{-1} \sum_{i=1}^{n} (Y_i - \sum_{j \in \underbrace{\mathcal{M}_\lambda}_{\text{model from Lasso}(\lambda)}} \beta_j X_i^{(j)})^2 + \phi \lambda \|\beta\|_1$$

for $\phi = 0$: OLS on selected variables from Lasso($\lambda$)

for $\phi = 1$: Lasso($\lambda$)

amount of computation for finding all solutions over $\lambda$ and $\phi$:

often, the same computational complexity as for Lasso/LARS (surprising):

$$O(np \min(n, p)) = O(p) \text{ if } p \gg n$$

worst case: $O(np \min(n, p)^2) = O(p) \text{ if } p \gg n$  still linear in $p$

this is "quasi-convex" optimization: two levels of a convex problem

for orthonormal case:
$$\mathbf{X}^T \mathbf{X} = I$$

Theorem (Meinshausen, 2005)

with essentially the same assumptions as before

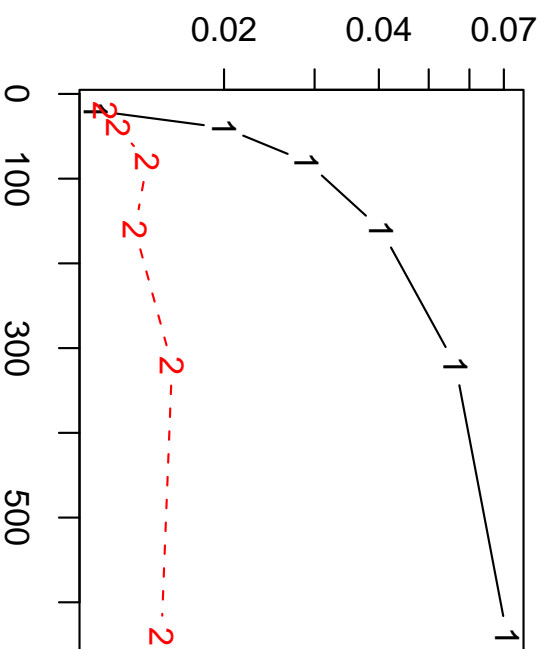$$\inf_{\lambda,\phi} L(\lambda,\phi) = O_P(n^{-1})(n \to \infty)$$

also: use the relaxed Lasso for variable selection and graphs/dependency networks

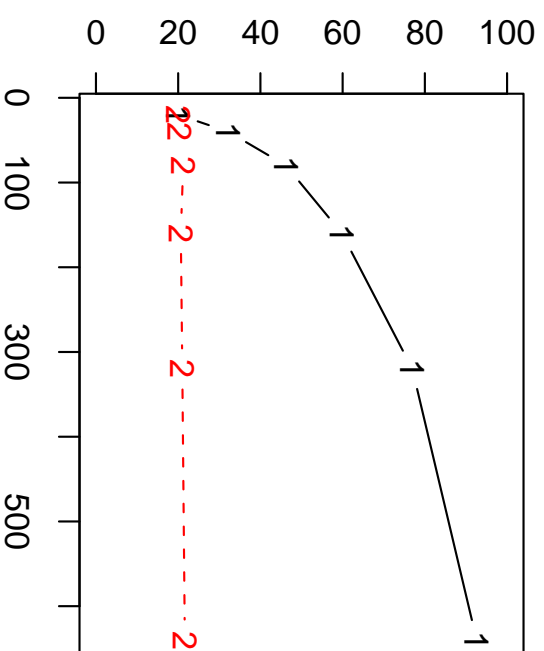$\rightsquigarrow$ prediction optimal (or cross-validated) tuning parameters yield (for some cases)
consistent variable selection and graph estimates



relaxed Lasso

$$n = 300, p = 20, \ldots 650, p_{eff} = 20$$
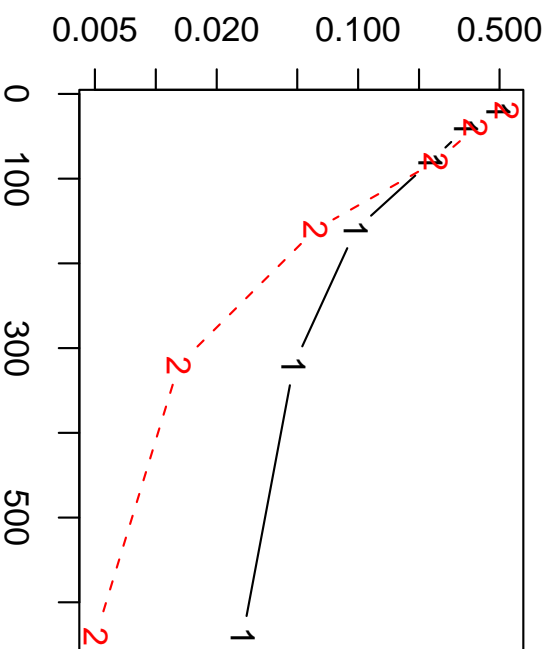
**L2–loss**

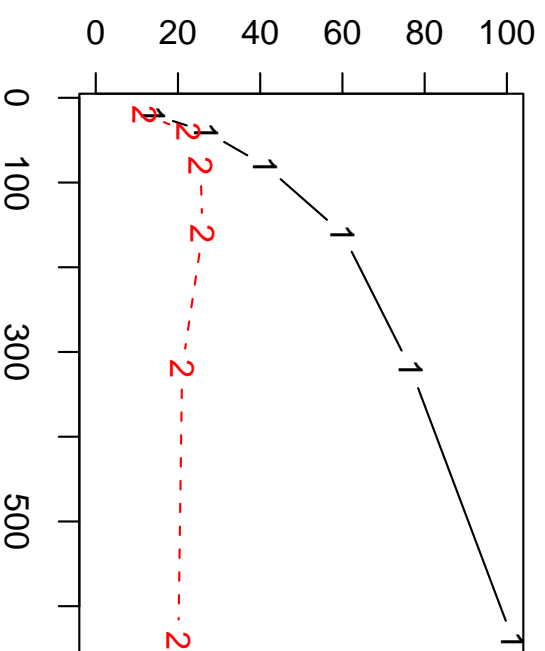**number of selected variables**

p

1: Lasso    2: relaxed Lasso

p

additional pure noise variables are much less damaging with the relaxed Lasso than
for Lasso and Boosting
and they are very disturbing for Ridge-type regularization (e.g. SVM)

$$n = p = 20, \ldots 650, p_{eff} = 20$$

L2-loss

number of selected variables

p

1: Lasso    2: relaxed Lasso

relaxed Lasso never substantially worse than the Lasso: the price for the flexibility of the relaxed Lasso is the larger search space $0 \le \phi \le 1$ (Lasso: $\phi = 1$)

**Results for high noise, binary lymph node classification**

cross-validated misclassification rate:

relaxed Lasso (tuned by 5-fold CV): 16.3%

Lasso (tuned by 5-fold CV): 21.0%

$L_2$Boosting (tuned by 5-fold CV): 24.8%

selected genes (on whole data set):

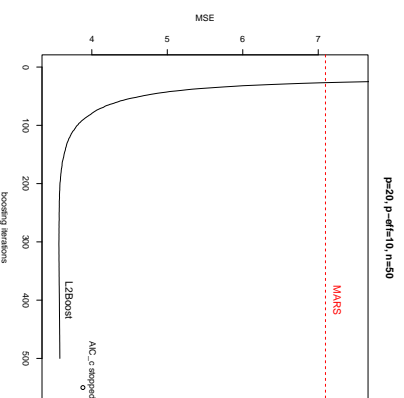relaxed Lasso: 2 genes (!)    Lasso: 23 genes    $L_2$Boosting: 42 genes

the 2 genes from relaxed Lasso are also selected by Lasso and $L_2$Boosting

note the identifiability problem among highly correlated predictor variables

# Conclusions

high-dimensional: blue greedy or convex?

the methods are similar and very useful

- Boosting is more generic: can be easily extended to e.g. the nonparametric setting

  nonparametric interaction modeling

  $L_2$Boosting with pairwise splines

  sample size $n = 50$

  $p = 20$, effective $p_{eff} = 5$



- Lasso is more explicit (and hence better understood)

  beyond Lasso (more sparse) is computationally feasible

  via relaxed Lasso doing "quasi-convex" optimization