

(actually: other weighting schemes are equally good or better...)

Why should this be good?

Adaboost proposed for classification by Freund & Schapire (1996)  
 data weights (rough original idea): large weights to previously heavily misclassified instances (sequential algorithm)  
 averaging weights  $a_m$ : large if in-sample performance in  $m$ th round was good

## 2.1. Boosting algorithms

classification of 2 lymph nodal status in breast cancer using gene expressions from  
 microarray data:  
 $n = 33$ ,  $p = 7129$  (for CART: gene-preselction, reducing to  $p = 50$ )

method	test set error	gain over CART
CART	22.5%	–
LogitBoost with trees	16.3%	28%
LogitBoost with bagged trees	12.2%	46%

Breiman (1998/99):  
 AdaBoost is **functional gradient descent (FGD)** procedure  
 a mix of statistical estimation and numerical optimization...

**2.2  $L_2$  Boosting**

(see also Friedman, 2001)

$L_2$  Boosting with base procedure  $\hat{\theta}(\cdot)$  is a "constrained minimization" of empirical risk  $n^{-1} \sum_{i=1}^n (Y_i - f(X_i))^2$  w.r.t.  $f(\cdot)$

$\rightsquigarrow$  useful for regression

$$\begin{aligned}
 m = 1 : (X_i, Y_i)_{i=1}^n &\rightsquigarrow \hat{\theta}_1(\cdot), f_1 = \nu \hat{\theta}_1 \\
 &\rightsquigarrow \text{resid. } U_i = Y_i - \hat{f}_1(X_i) \\
 m = 2 : (X_i, Y_i, U_i)_{i=1}^n &\rightsquigarrow \hat{\theta}_2(\cdot), \hat{f}_2 = \hat{f}_1 + \nu \hat{\theta}_2 \\
 &\rightsquigarrow \text{resid. } U_i = Y_i - \hat{f}_2(X_i)
 \end{aligned}$$

...

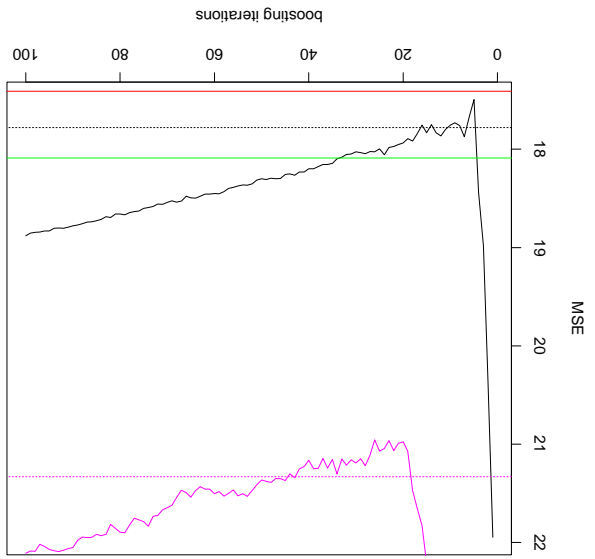
$$f^{m_{stop}}(\cdot) = \nu \sum_{m=1}^{m_{stop}} \hat{\theta}^m(\cdot), \quad m_{stop} \text{ a tuning parameter}$$

repeated greedy fitting (with shrinkage  $\nu$ ) of residuals

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

Any gain over classical methods? (for additive modeling)

Ozone data:  $n=300, p=8$



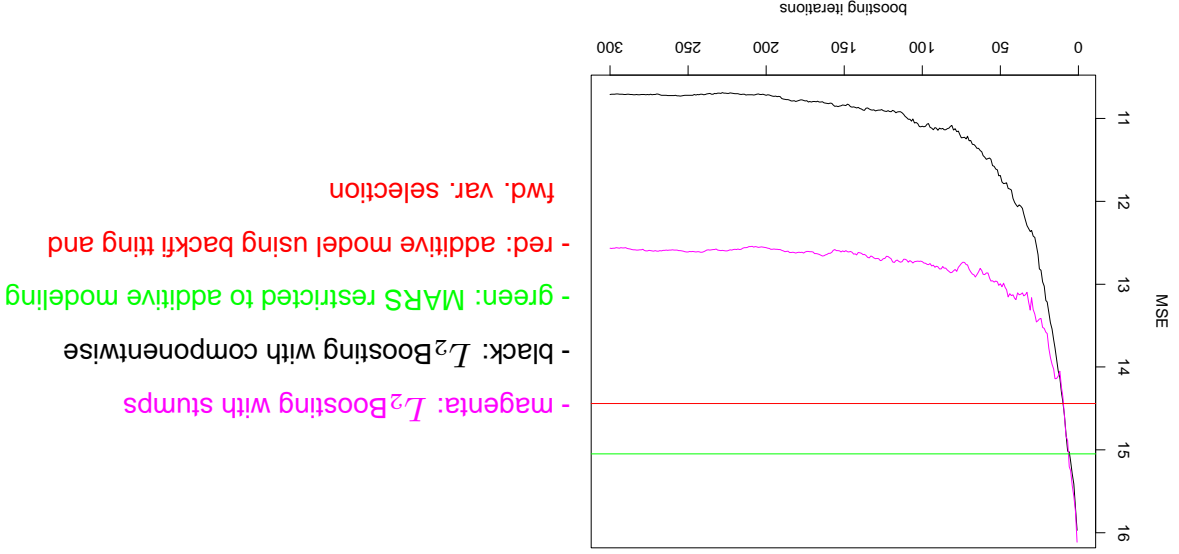
- magenta:  $L_2$  Boosting with stumps (horiz. line = cross-validated stopping)
  - black:  $L_2$  Boosting with componentwise smoothing spline (horiz. line = cross-validated stopping)
  - green: MARS restricted to additive modeling
  - red: additive model using backfitting
- selected predictor which reduces RSS most

$L_2$  Boosting with stumps or comp. smoothing splines also yields additive model:

$$\hat{\theta}_m(x) = \hat{g}_1(x) + \dots + \hat{g}_p(x)$$

Simulated data: non-additive regression function,  $n = 200, p = 100$

Regression:  $n=200, p=100$



- magenta:  $L_2$  Boosting with stumps
- black:  $L_2$  Boosting with componentwise smoothing spline
- green: MARS restricted to additive modeling
- red: additive model using backfitting and fwd. var. selection

similar for classification

very often: boosting performs comparatively well in high-dimensions  
(there is a lot of empirical evidence for this)

also SVM is often surprisingly accurate...

### 2.3. Choice of the base procedure

most popular in machine learning: tree algorithms (CART, C4.5)  
they do variable/feature selection

have seen: for componentwise smoothing splines or stumps  
→ boosting yields an additive model fit  
↔ we can use boosting for fitting in "quite many" structural models

Example: degree 2 nonparametric interaction modeling

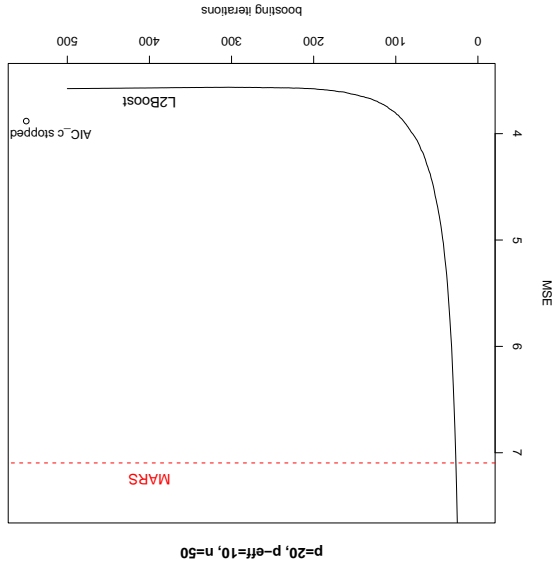
Friedman #1 model:

$$Y = 10 \sin(\pi X_1 X_2) + 20(X_3 - 0.5)^2 + 10X_4^2 + 5X_5 + \mathcal{N}(0, 1), \quad X = (X_1, \dots, X_{20}) \sim \text{unif}([0, 1]^{20})$$

$L_2$  Boosting with pairwise splines

sample size  $n = 50$

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



linear model

2.4.  $L_2$  Boosting for high-dimensional linear models

$$Y = f(X) + \varepsilon, \quad \sum_{j=1}^d g_j(x^{(j)}) = f(x), \quad u \ll d$$

or: a highly over-complete dictionary  $\{g_j(\cdot); j = 1, \dots, d\}$ ,  $d \gg u$

our approach:  $L_2$  Boosting with componentwise linear  $L_2$  regression

This **base procedure** fits a univariate linear regression model against the one predictor variable which reduces residual sum of squares most

R.V. Southwell in 1933  
 Professor in engineering  
 Oxford University



“Princeps Mathematicorum”

C.F. Gauss in 1803



Gauss-Southwell algorithm

For  $\nu = 1$ , this  $L_2$  Boosting is known as **Matching Pursuit** (Mallat and Zhang, 1993)  
 not full OLS on selected variables (even with  $\nu = 1$ )  
 assigns **variable amount of degrees of freedom for selected variables (shrinkage)**  
 this method does **variable selection** and

very different from forward variable selection  
 etc.

first round of estimation: selected predictor variable  $X^{(S_1)}$  (e.g.  $= X^{(3)}$ )  
 corresponding ordinary least squares  $\hat{\beta}_{S_1}$   
 use shrunken fit  $\hat{f}_1 = \nu \hat{\beta}_{S_1} X^{(S_1)}$  (e.g.  $\nu = 0.1$ )  
 second round of estimation: selected predictor variable  $X^{(S_2)}$  (e.g.  $= X^{(21)}$ )  
 corresponding OLS  $\hat{\beta}_{S_2}$   
 use shrunken fit  $\hat{f}_2 = \hat{f}_1 + \nu \hat{\beta}_{S_2} X^{(S_2)}$

because of

variable selection and

assigning variable amount of degrees of freedom (shrinkage) for selected variables

reminds to Lasso ( $\ell_1$ -penalized regression) (Tibshirani, 1996)

$$\beta_{Lasso} = \operatorname{argmin}_{\beta} \sum_{i=1}^n Y_i - \sum_{j=1}^d \beta_j X_{ij}^2 + \underbrace{\lambda}_{\geq 0} \sum_{j=1}^d |\beta_j|$$

and indeed: there is a relation (Efron, Hastie, Johnstone, Tibshirani, 2004)

but: the algorithms and estimates are not the same

Theorem for high dimensions (PB, 2004)

$L_2$  Boosting with comp. linear LS regression is **consistent** (for suitable number of

boosting iterations) if:

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

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  $m = m_n$ :

$$\mathbb{E}_X |f_{m_n,n}(X) - f^n(X)|^2 = o_P(1) \quad (n \rightarrow \infty)$$

"no" assumptions about the predictor variables/design matrix

in other words:

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

similar result has been given for the Lasso by Greenshtein and Ritov (2004)



### 3. $L_2$ Boosting, Lasso and LARS

Efron et al. (2004): intriguing relation between  $L_2$  Boosting and Lasso:  $\hat{\beta}_{Lasso} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - \sum_{j=1}^p \beta_j X_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j|$

for some special cases, roughly:

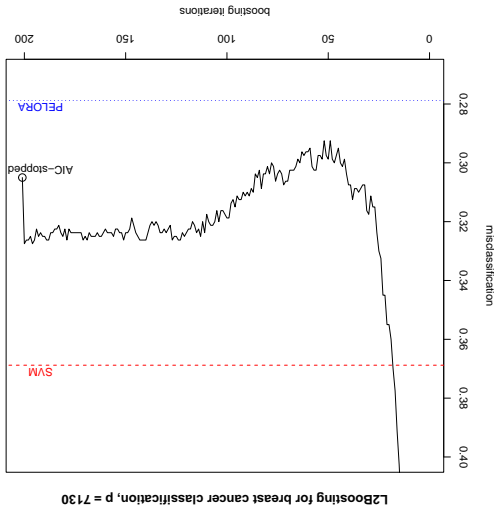
iterations of " $L_2$  Boosting with "infinitesimally" small  $\nu$  yield all Lasso solutions when varying  $\lambda$ "

↪ 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

42 out of  $p = 7130$  genes are selected (some of them biologically meaningful) good prediction and interesting gene selection



- black:  $L_2$  Boosting with componentwise  $n = 49, p = 7130$  gene expressions
- linear LS regression
- red: SVM with radial basis kernel
- blue: Pelora: a "biologically inspired" gene grouping method (Dettling & PB, 2004)

binary lymph node classification in breast cancer using gene expressions

and LARS is really fast

both: Lasso/LARS and  $L_2$  Boosting are very useful

for  $p \gg n$