

**A Consistent Model Selection Criterion
for L_2 -Boosting in
High-Dimensional Sparse Linear Models**

Tze Leung Lai, *Stanford University*

Ching-Kang Ing, *Academia Sinica, Taipei*

Zehao Chen, *Lehman Brothers*

High Dimensional Regression

- $$y_{in} = \sum_{j=1}^{p_n} \beta_{jn} x_{ij}(n) + \varepsilon_i, \quad i = 1, 2, \dots, n,$$

$p_n = O(\exp(n^\xi))$ for some $0 < \xi < 1$.

The n 's in β_{jn} , y_{in} and $x_{ij}(n)$ may be suppressed.

- Some typical examples for $n \ll p$:
 - ◇ Gene Expression Data: $n =$ sample size, $p = \#$ genes.
 - ◇ Sparse Signal Reconstruction: $n = \#$ measurements,
 $p = \#$ signals
 - ◇ Image Recovery: $n = \#$ grid points,
 $p = \#$ basis functions
 - ◇ Portfolio Selection: $p = \#$ assets, $n = \#$ time points

L_2 -Boosting: Pure Greedy Algorithm (PGA)

Step 1. Define $\mathbf{R}_0 = (y_1, \dots, y_n)'$ and $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})'$. Find a variable among $\{\mathbf{x}_1, \dots, \mathbf{x}_{p_n}\}$ that is **most correlated** to \mathbf{R}_0 . Call the variable $\mathbf{x}_{\hat{s}_1}$ and generate residual vector $\mathbf{R}_1 = \mathbf{R}_0 - \hat{\beta}_{\hat{s}_1} \mathbf{x}_{\hat{s}_1}$, where $\hat{\beta}_{\hat{s}_1}$ is the least squares estimate obtained by regressing \mathbf{R}_0 on $\mathbf{x}_{\hat{s}_1}$.

Step 2. Find a variable among $\{x_1, \dots, x_{p_n}\}$ that is **most correlated** to R_1 . Call the variable $x_{\hat{s}_2}$ and let $R_2 = R_1 - \hat{\beta}_{\hat{s}_2} x_{\hat{s}_2}$.

⋮

If iterations stop at step m , then the new outcome, y_{n+1} , is predicted by

$$\hat{y}_{\hat{s}_1, \dots, \hat{s}_m} = \hat{\beta}_{\hat{s}_1} x_{n+1, \hat{s}_1} + \dots + \hat{\beta}_{\hat{s}_m} x_{n+1, \hat{s}_m}.$$

L_2 -Boosting: Orthogonal Greedy Algorithm (OGA)

Step 1. Define $\mathbf{y}_n = (y_1, \dots, y_n)'$ and $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})'$. Find a variable among $\{\mathbf{x}_1, \dots, \mathbf{x}_{p_n}\}$ that is **most correlated** to \mathbf{y}_n . Call the variable $\mathbf{x}_{\hat{s}_1^o}$ and generate residual $\mathbf{R}_1 = \mathbf{y}_n - \mathbf{M}_{\hat{s}_1^o} \mathbf{y}_n$, where $\mathbf{M}_{\hat{s}_1^o}$ is the projection matrix into $\mathcal{L}(\mathbf{x}_{\hat{s}_1^o})$.

Step 2. Find a variable among $\{x_1, \dots, x_{p_n}\}$ that is **most correlated** to R_1 . Call the variable $x_{\hat{s}_2^o}$ and let $R_2 = y_n - M_{\hat{s}_1^o, \hat{s}_2^o} y_n$, where $M_{\hat{s}_1^o, \hat{s}_2^o}$ is the projection matrix into $\mathcal{L}(x_{\hat{s}_1^o}, x_{\hat{s}_2^o})$.

⋮

If iterations stop at step m , then the new outcome, y_{n+1} , is predicted by

$$\hat{y}_{\hat{s}_1^o, \dots, \hat{s}_m^o} = \hat{\beta}_{\hat{s}_1^o}^o x_{n+1, \hat{s}_1^o} + \dots + \hat{\beta}_{\hat{s}_m^o}^o x_{n+1, \hat{s}_m^o},$$

where $\hat{\beta}_{\hat{s}_1^o}^o, \dots, \hat{\beta}_{\hat{s}_m^o}^o$ are the least squares estimates.

Convergence Rates in L_2 -Boosting

Assumptions:

(K.1) $Ee^{t\varepsilon_1} + \sup_j E \exp(sx_{1j}) < \infty$ for $|t| \leq t_0$, $|s| \leq s_0$.

(K.2) $0 < \iota_1 < \lambda_{\min}(\mathbf{\Gamma}_n) \leq \lambda_{\max}(\mathbf{\Gamma}_n) < \iota_2 < \infty$ for all n ,
where $\mathbf{\Gamma}_n = E(\mathbf{x}_1 \mathbf{x}'_1)$

(K.2*) $0 < \eta_1 \leq \min_{1 \leq j \leq p_n} E(x_{1j}^2) \leq \max_{1 \leq j \leq p_n} E(x_{1j}^2) < \eta_2 < \infty$.

(K.3) $p_n = O(\exp(Cn^\xi))$ for some $0 < \xi < 1$ and $C > 0$.

(K.4) $\sup_{n \geq 1} \sum_{j=1}^{p_n} |\beta_j| < \infty$.

Theorem 1. Assume (K.1), (K.2), (K.3) and (K.4). Then, for any choice of $m = m_n$ satisfying $m_n = O(n^l)$ with $0 < l < (1 - \xi)/5$,

$$E \left\{ (y_{n+1} - \hat{y}_{\hat{s}_1^o, \dots, \hat{s}_m^o})^2 \mid \mathbf{y}_n, \mathbf{x}_1, \dots, \mathbf{x}_{p_n} \right\} - \sigma^2 = O_p(m^{-1}).$$

Theorem 2. Assume (K.1), (K.2*), (K.3) and (K.4). Then, for any choice of $m = m_n$ satisfying $m_n = o(\log n)$ and $0 < \theta < \frac{1}{3}$,

$$E \left\{ (y_{n+1} - \hat{y}_{\hat{s}_1^P, \dots, \hat{s}_m^P})^2 \mid \mathbf{y}_n, \mathbf{x}_1, \dots, \mathbf{x}_{p_n} \right\} - \sigma^2 = O_p(m^{-\theta}).$$

Theorem 3. Assume (K.1), (K.2), (K.3), (K.4) and
(K.5) For some $0 \leq \gamma < (1 - \xi)/5$,

$$\liminf_{n \rightarrow \infty} n^\gamma \min_{j \in O_n} \beta_j^2 > C^* > 0.$$

Then, for any choice of $m = m_n$ satisfying $m_n = O(n^l)$ with
 $0 < l < (1 - \xi)/5$ and $\lim_{n \rightarrow \infty} m_n/n^\gamma = \infty$,

$$\lim_{n \rightarrow \infty} P(D_n^o) = 1,$$

where

$$O_n = \{j : 1 \leq j \leq p_n, \beta_j \neq 0\}, \quad D_n^o = \left\{ O_n \subseteq \{\hat{s}_1^o, \dots, \hat{s}_{m_n}^o\} \right\}.$$

Theorem 3 indicates that with probability approaching 1, the variables chosen by the OGA after m_n steps contains all relevant variables. To prevent overfitting, the algorithm should be stopped at the first time when all relevant variables are included, i.e., choose the smallest correct model along the boosting path.

High-Dimensional Hannan–Quinn Criterion

$$\text{HQ}(k) = \log \hat{\sigma}^2(x_{\hat{s}_1^o}, \dots, x_{\hat{s}_k^o}) + \frac{(\log p_n)C_n k}{n}$$

$$\hat{k}_{\text{HQ}} = \arg \min_{1 \leq k \leq m_n} \text{HQ}(k),$$

where $C_n \rightarrow \infty$ and

$$\hat{\sigma}^2(x_{\hat{s}_1^o}, \dots, x_{\hat{s}_k^o}) = \frac{1}{n} \mathbf{y}'_n (\mathbf{I} - \mathbf{M}_{\hat{s}_1^o, \dots, \hat{s}_k^o}) \mathbf{y}_n.$$

The major difference between HQ and conventional Hannan–Quinn is the additional factor $\log p_n$, which is related to the maximum of i.i.d. $V_1, \dots, V_{p_n} \sim \chi^2(1)$:

$$\frac{\max_{1 \leq k \leq p_n} V_i}{2 \log p_n} \rightarrow 1 \text{ in probability.}$$

Define $A_j = \{\hat{s}_1^o, \dots, \hat{s}_j^o\}$, $j \geq 1$,

$$\tilde{k}_n^o = \begin{cases} m_n, & \text{if } O_n \not\subseteq A_{m_n} \\ \min\{j : 1 \leq j \leq m_n, O_n \subseteq A_j\}, & \text{if } O_n \subseteq A_{m_n}. \end{cases}$$

Theorem 4. *Assume (K.1), (K.2), (K.3), (K.4) and (K.5). Then, for any choice of $m = m_n$ satisfying $m_n = O(n^l)$ with $0 < l < (1 - \xi)/5$ and $\lim_{n \rightarrow \infty} m_n/n^\gamma = \infty$, and any choice of C_n satisfying $C_n \log p_n = o(n^{1-2\gamma})$ and $n^\gamma/C_n = o(1)$,*

$$\lim_{n \rightarrow \infty} P(A_{\hat{k}_{HQ}} = A_{\tilde{k}_n^o}) = 1.$$

- For $\gamma = 0$, one can take $C_n = \log n$, giving the high-dimensional version of BIC.

Sketch of Proof

1. Difference between OGA (or PGA) with its **population version**: Exponential bounds.
2. The convergence rates in Theorems 1 and 2 are used to prove Theorem 3.
3. Theorem 3, exponential bounds and extension of Hannan–Quinn-type arguments to prove Theorem 4.

Finite-Sample Performance

Example 1. $(\beta_1, \dots, \beta_5) = (3, -3.5, 4, -2.8, 3.2)$,

$\beta_j = 0$ for $j \geq 6$ in

$$y_i = \sum_{j=1}^5 \beta_j x_{ij} + \sum_{j=6}^p \beta_j x_{ij} + \varepsilon_i, \quad i = 1, \dots, n.$$

- ε_i i.i.d. $\mathcal{N}(0, 1)$, $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T$ independent of ε_i .
- $x_{ij} = z_{ij} + \eta w_i$, where $\eta > 0$, $(\mathbf{z}_i^T, w_i)^T \sim N(\mathbf{0}, \mathbf{I})$.
- 1000 simulation runs for each result.

Table 1: Frequency of all 5 relevant and i additional variables selected ($m_n = 30$) in 1000 simulations.

| η | n | p | method | i | | | | | | | | | Total |
|--------|-----|------|-----------|------|---|---|---|---|---|---|---|-------|-------|
| | | | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 9 | 25–30 | |
| 0 | 100 | 2000 | OGA+HDBIC | 1000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1000 |
| | | | OGA+BIC | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1000 | 1000 |
| | 200 | 4000 | OGA+HDBIC | 1000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1000 |
| | | | OGA+BIC | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1000 | 1000 |
| 2 | 100 | 2000 | OGA+HDBIC | 992 | 7 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1000 |
| | | | OGA+BIC | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1000 | 1000 |
| | 200 | 4000 | OGA+HDBIC | 1000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1000 |
| | | | OGA+BIC | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1000 | 1000 |

Example 2. $q = 10, n = 400, p = 4000$ in

$$y_i = \sum_{j=1}^q \beta_j x_{ij} + \sum_{j=q+1}^p \beta_j x_{ij} + \varepsilon_i, \quad i = 1, \dots, n,$$

- $(\beta_1, \dots, \beta_{10}) =$
 $(3.2, 3.2, 3.2, 3.2, 4.4, 4.4, 3.5, 3.5, 3.5, 3.5),$
 $\beta_j = 0$ for $j > 10$.
- ε_i i.i.d. $N(0, 2.25)$.
- $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T$, same as in Example 1.
- 100 simulations for each result.

Table 2: Frequency of all 9 relevant and i additional variables selected ($m_n = 40$) in 100 simulations.

| η | method | i | | | | | | | | | | | Total | |
|--------|-----------|-----|----|----|----|----|----|---|---|---|---|-------|-------|-----|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10–19 | | |
| 1 | OGA+HDBIC | 97 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 100 |
| | LASSO | 35 | 25 | 12 | 6 | 11 | 4 | 2 | 2 | 0 | 1 | 2 | 100 | |
| 3 | OGA+HDBIC | 64 | 29 | 6 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 100 | |
| | LASSO | 10 | 13 | 14 | 10 | 11 | 14 | 7 | 7 | 0 | 3 | 11 | 100 | |

Example 3. Whereas Example 2 satisfies the *neighborhood stability condition* of Meinshausen and Bühlmann (2006), which is an extension of the *irrepresentable condition* of Zhao and Yu (2006) to ensure model selection consistency of LASSO to random covariates, this condition is violated in taking $q = 10$ in Example 2 and

- $(\beta_1, \dots, \beta_q) = (3, 3.75, 4.5, 5.25, 6, 6.75, 7.5, 8.25, 9, 9.75)$.
- $\beta_j = 0$ for $j \geq 11$, $(x_{i1}, \dots, x_{iq})^T$ i.i.d. $N(\mathbf{0}, \mathbf{I}_q)$.
- ε_i i.i.d. $N(0, 1)$ and independent of $(x_{i1}, \dots, x_{ip})^T$.
- $x_{ij} = z_{ij} + (\sum_{l=1}^q b x_{il})$ for $q < j \leq p$, $\mathbf{z}_i \sim N(\mathbf{0}, \mathbf{I}_{p-q}/4)$.
- $qb^2 + 1/4 = 1$ (i.e, $b = \sqrt{3/4q}$).

Table 3: Frequency of all 10 relevant and i additional variables selected ($m_n = 40$) in 100 simulations.

| method | i | | | | | | | | | Total |
|-----------|-----|---|----|---|---|---|---|---|------------|-------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 30 or more | |
| OGA+HDBIC | 0 | 9 | 87 | 4 | 0 | 0 | 0 | 0 | 0 | 100 |
| LASSO | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 100 | 100 |

Conclusion

- By obtaining the convergence rates of orthogonal greedy L_2 -boosting (OGA), classical model selection criteria such as BIC can be modified to a large number p_n of regressors in sparse regression models. Consistency of model selection is formulated in terms of capturing all, and only, variables with nonzero coefficients under (K5).

- OGA with consistent model selection means that the estimate is asymptotically equivalent to OLS with the “correct” set of covariates.
- In the absence of (K5), consistency of model selection can be linked to the objective of the regression model, e.g., prediction, estimation of high-dimensional covariance matrix via modified Cholesky decomposition. For the latter, “consistent” model selection can be carried out via thresholding ([Bickel & Levina, 2008](#), and extensions).