# CORRECTED PROOF OF THE RESULT OF 'A PREDICTION ERROR PROPERTY OF THE LASSO ESTIMATOR AND ITS GENERALIZATION' BY HUANG (2003)

SAHARON ROSSET[1*] AND JI ZHU[2]

*IBM T.J. Watson Research Center and University of Michigan*

## Summary

The Lasso achieves variance reduction and variable selection by solving an $\ell_1$-regularized least squares problem. Huang (2003) claims that 'there always exists an interval of regularization parameter values such that the corresponding mean squared prediction error for the Lasso estimator is smaller than for the ordinary least square estimator'. This result is correct. However, its proof in Huang (2003) is not. This paper presents a corrected proof of the claim, which exposes and uses some interesting fundamental properties of the Lasso.

*Key words:* lasso; least squares estimator; piecewise linear; prediction error.

## 1. Introduction

The Lasso (Tibshirani, 1996) achieves variance reduction and variable selection by solving an $\ell_1$-penalized least squares problem:

$$\hat{\boldsymbol{\beta}}(\gamma) = \arg \min_{\boldsymbol{\beta}} \| \boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta} \|^2 + \gamma \sum_j |\beta_j| \,.$$

For 'large' values of $\gamma$ (defined in a relative sense), $\hat{\boldsymbol{\beta}}(\gamma)$ has many zero components. This was a major motivating factor for the Lasso, as it implies that the Lasso is appropriate for 'sparse' models, where ridge regression is unlikely to succeed, since it forces all coefficients to be non-zero (Friedman *et al.*, 2004). In some cases, the Lasso or equivalent procedures have provable optimality properties, such as in the case of wavelet shrinkage (Donoho *et al.*, 1995).

Recently, it has been shown (Osborne, Presnell & Turlach, 2000; Efron *et al.*, 2004) that the path of optimal solutions for the Lasso, $\{\hat{\boldsymbol{\beta}}(\gamma), 0 \le \gamma \le \infty\}$ is piecewise linear, and thus the Lasso can be solved efficiently for *all* values of $\gamma$ using an incremental algorithm. A simple example to illustrate the piecewise linear property can be seen in Figure 1, where we show the Lasso optimal solution paths for a four-variable synthetic dataset. The plot shows the optimal Lasso solutions $\hat{\boldsymbol{\beta}}(\gamma)$ as a function of the $\ell_1$ norm $\|\hat{\boldsymbol{\beta}}(\gamma)\|_1$. Each line represents one coefficient and gives its values at the optimal solution for the range of $\|\hat{\boldsymbol{\beta}}(\gamma)\|_1$ values. We observe that between points marked '+' the lines are straight, i.e. the coefficient paths are piecewise-linear, and the one-dimensional curve $\hat{\boldsymbol{\beta}}(\gamma)$ is piecewise linear in $\mathbb{R}^4$.
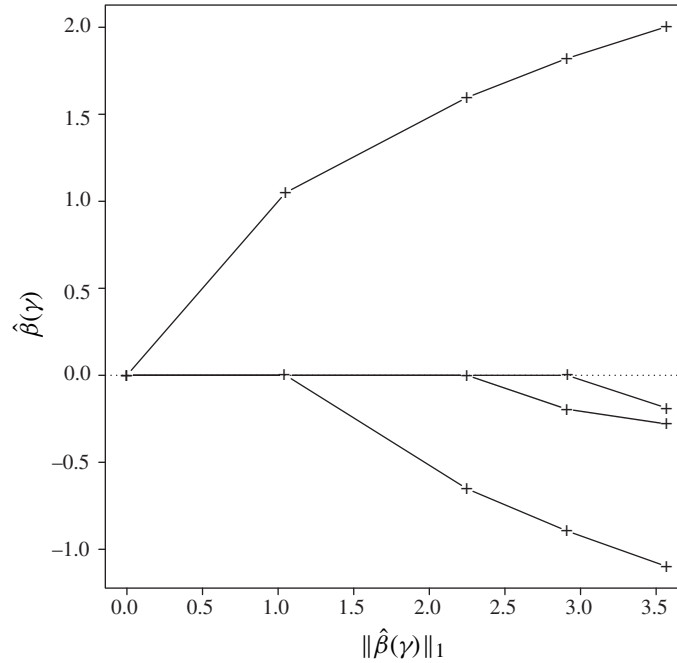
Figure 1.  Piecewise linear solution paths for the Lasso on a simple four-variable example

Using this property, Efron *et al.* (2004) suggest the LAR-Lasso algorithm, which allows generation of the whole regularized solution path, $\{\hat{\boldsymbol{\beta}}(\gamma), 0 \leq \gamma \leq \infty\}$, for 'approximately' the computational cost of one least-squares calculation on the full dataset (the exact cost depends on some rather complicated properties of the regularized path, but the assumptions required to attain the above property are quite mild).

The paper by Huang (2003) describes another interesting and desirable property of the Lasso (Huang, 2003 p.218, Theorem 2), shown here with modified notation:

*There exists a value $\gamma_0 > 0$ such that for $\gamma \in (0, \gamma_0]$:*

$$\mathrm{E}\left(\|\boldsymbol{y}_{\mathrm{new}} - \boldsymbol{X}\hat{\boldsymbol{\beta}}(\gamma)\|^2\right) < \mathrm{E}\left(\|\boldsymbol{y}_{\mathrm{new}} - \boldsymbol{X}\hat{\boldsymbol{\beta}}\|^2\right)$$

*That is, the mean squared prediction error of the Lasso estimator $\hat{\boldsymbol{\beta}}(\gamma)$ is smaller than that of the least square estimator $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}(0)$ when the tuning parameter $\gamma$ is small enough.*

This result is correct. However, its proof in that paper is not. Specifically, Theorem 3 therein is incorrect, and is fundamental in the proof of Theorem 2. In this note, we present a corrected version of Theorem 3, and a corrected proof of Theorem 2.

The notation throughout this paper is as in Huang (2003), except that we drop the $(\ell)$ and $(0)$ superscripts for the Lasso solutions. As in Huang (2003) we assume throughout that the predictor matrix $\boldsymbol{X}$ is fixed. We denote by $\boldsymbol{y}$ the stochastic response vector used for fitting the model and by $\boldsymbol{y}_{\mathrm{new}}$ a new, independent copy.

Unfortunately, our proof is not nearly as short and elegant as the proof using the incorrect Theorem 3. However, we believe it is of independent interest, as it exposes and uses some interesting fundamental properties of the Lasso — in particular Lemma 2 below and its proof which describes the 'piecewise linear' pieces of the Lasso path analytically.

## 2. Corrected results

### 2.1. Corrected Theorem 3 of Huang (2003 p.219)

The original theorem in the paper reads, with modified notation:

*There exists a value $\gamma_1 > 0$ such that for $\gamma \in [0, \gamma_1]$*

$$\hat{\boldsymbol{\beta}}(\gamma) = \hat{\boldsymbol{\beta}} - \tfrac{1}{2}\gamma(X^{\mathsf{T}}X)^{-1}s(\hat{\boldsymbol{\beta}}) \qquad almost\ surely.$$

Our corrected version is:

*With probability 1 there exists a sample dependent value $\gamma_1(\boldsymbol{y}) > 0$ such that for $\gamma \in [0, \gamma_1(\boldsymbol{y})]$*

$$\hat{\boldsymbol{\beta}}(\gamma) = \hat{\boldsymbol{\beta}} - \tfrac{1}{2}\gamma(X^{\mathsf{T}}X)^{-1}s(\hat{\boldsymbol{\beta}}). \tag{1}$$

This corrected version does not assume there is a $\gamma_1$ which applies to (almost surely) all possible samples, but rather that it is sample dependent. The proof of Theorem 3 in Huang (2003 p.226) actually proves this corrected version and requires no modification.

### 2.2. Corrected proof of main result

We consider only the right derivative of the expected error at $\gamma = 0$ and prove it is negative. This concludes an existence proof for $\gamma_0$ in Theorem 2.

Define $\tilde{\boldsymbol{\beta}}(\gamma)$ as the solution which just extends the last piece of the Lasso path backwards,

$$\tilde{\boldsymbol{\beta}}(\gamma) = \hat{\boldsymbol{\beta}} - \tfrac{1}{2}\gamma(X^{\mathsf{T}}X)^{-1}s(\hat{\boldsymbol{\beta}}),$$

and compare this to (1); $\tilde{\boldsymbol{\beta}}(\gamma) = \hat{\boldsymbol{\beta}}(\gamma)$ if and only if $\gamma \leq \gamma_1(\boldsymbol{y})$.

**Proof of Theorem 2.** Consider the right derivative of the error at $\gamma = 0$:

$$\frac{\partial}{\partial\gamma_+}\big(\mathrm{E}\big(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}(\gamma)\|^2\big) - \mathrm{E}\big(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}\|^2\big)\big)_{\gamma=0}$$

$$= \lim_{\gamma\searrow 0}\frac{\mathrm{E}(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}(\gamma)\|^2) - \mathrm{E}(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}\|^2)}{\gamma}.$$

We re-write the numerator as

$$\mathrm{E}\big(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}(\gamma)\|^2\big) - \mathrm{E}\big(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}\|^2\big) = \big(\mathrm{E}\big(\|\boldsymbol{y}_{\mathrm{new}} - X\tilde{\boldsymbol{\beta}}(\gamma)\|^2\big) - \mathrm{E}\big(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}\|^2\big)\big)$$
$$- \big(\mathrm{E}\big(\|\boldsymbol{y}_{\mathrm{new}} - X\tilde{\boldsymbol{\beta}}(\gamma)\|^2\big) - \mathrm{E}\big(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}(\gamma)\|^2\big)\big).$$

The proof given for Theorem 2 in Huang (2003 p.227) proves that

$$\lim_{\gamma\searrow 0}\frac{\mathrm{E}(\|\boldsymbol{y}_{\mathrm{new}} - X\tilde{\boldsymbol{\beta}}(\gamma)\|^2) - \mathrm{E}(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}\|^2)}{\gamma} < 0.$$

And so all we have left to prove is

$$\lim_{\gamma\searrow 0}\frac{\mathrm{E}(\|\boldsymbol{y}_{\mathrm{new}} - X\tilde{\boldsymbol{\beta}}(\gamma)\|^2) - \mathrm{E}(\|\boldsymbol{y}_{\mathrm{new}} - X\hat{\boldsymbol{\beta}}(\gamma)\|^2)}{\gamma} = 0. \tag{2}$$

We start by proving a couple of useful lemmas.

**Lemma 1.** $\Pr(\hat{\boldsymbol{\beta}}(\gamma) \neq \tilde{\boldsymbol{\beta}}(\gamma)) \to 0$ *as* $\gamma \to 0$.

**Proof.** This follows from the corrected Theorem 3, because if $\gamma < \gamma_1(\boldsymbol{y})$ then $\hat{\boldsymbol{\beta}}(\gamma) = \tilde{\boldsymbol{\beta}}(\gamma)$.

**Lemma 2.** *There exists $M > 0$ such that $\|\hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma)\| \leq M\gamma$ for all $\gamma > 0$. So we can say that $\hat{\boldsymbol{\beta}}(\gamma)$ is 'linearly' close to $\tilde{\boldsymbol{\beta}}(\gamma)$.*

**Proof.** We first observe that the Lasso optimal solution path $\hat{\boldsymbol{\beta}}(\gamma)$ is piecewise linear as a function of $\gamma$ (see Efron *et al.*, 2004 for details). One result of Efron *et al.* (2004) indicates that within each linear piece of the solution path, the set of predictor variables with non-zero coefficients is constant, i.e. $\mathcal{A} = \{j: \hat{\boldsymbol{\beta}}_j(\gamma) \neq 0, j = 1, \ldots, p\}$ does not change within each linear piece, and the derivative of $\hat{\boldsymbol{\beta}}(\gamma)$ with respect to $\gamma$ is equal to $-\frac{1}{2}(X_{\mathcal{A}}^{\top} X_{\mathcal{A}})^{-1} s_{\mathcal{A}}$, where $X_{\mathcal{A}}$ is the corresponding sub-matrix of $X$, and $s_{\mathcal{A}}$ is the vector containing the signs of $\hat{\boldsymbol{\beta}}_{\mathcal{A}}(\gamma)$. Hence it implies that there exist $0 = \gamma_0 < \gamma_1 < \gamma_2 < \cdots < \gamma_m = \infty$ and $\mathcal{A}_1, \ldots, \mathcal{A}_m \subset \{1, \ldots, p\}$ and $s_0 \in \{-1, +1\}^p$, $s_1 \in \{-1, +1\}^{|\mathcal{A}_1|}, \ldots, s_m \in \{-1, +1\}^{|\mathcal{A}_m|}$ such that if $\gamma_k \leq \gamma < \gamma_{k+1}$ then (with a slight abuse of notation)

$$\hat{\boldsymbol{\beta}}(\gamma) = \hat{\boldsymbol{\beta}} - \tfrac{1}{2}\gamma_1 (X^{\top}X)^{-1} s_0 - \tfrac{1}{2}(\gamma_2 - \gamma_1)(X_{\mathcal{A}_1}^{\top} X_{\mathcal{A}_1})^{-1} s_1 - \cdots - \tfrac{1}{2}(\gamma - \gamma_k)(X_{\mathcal{A}_k}^{\top} X_{\mathcal{A}_k})^{-1} s_k.$$

Here $(X_{\mathcal{A}}^{\top} X_{\mathcal{A}})^{-1} s$ is actually an $|\mathcal{A}| \times 1$ vector, rather than a $p \times 1$ vector. For notational simplicity, we have omitted the zero components, but this does not affect our claims below.

By the triangle inequality we then get

$$\|\hat{\boldsymbol{\beta}}(\gamma) - \hat{\boldsymbol{\beta}}\| \leq \tfrac{1}{2}\gamma \max_j \|(X_{\mathcal{A}_j}^{\top} X_{\mathcal{A}_j})^{-1} s_j\|,$$

which uses the specific data-dependent sequence of $\mathcal{A}_j$, $s_j$, but we can easily generalize it to a data-independent result by observing that

$$\left| \left\{ s \in \{-1, +1\}^{|\mathcal{A}|}: \mathcal{A} \subset \{1, \ldots p\} \right\} \right| = 3^p - 1,$$

and thus we can find

$$M = \max_{(\mathcal{A}, s)} \|(X_{\mathcal{A}}^{\top} X)^{-1} s\|$$

and get the data-independent bound

$$\|\hat{\boldsymbol{\beta}}(\gamma) - \hat{\boldsymbol{\beta}}\| \leq \tfrac{1}{2}\gamma M. \tag{3}$$

Next, we use the definition of $\tilde{\boldsymbol{\beta}}(\gamma)$ to bound

$$\|\tilde{\boldsymbol{\beta}}(\gamma) - \hat{\boldsymbol{\beta}}\| \leq \tfrac{1}{2}\gamma \max_s \|(X^{\top}X)^{-1} s\| \leq \tfrac{1}{2}\gamma M, \tag{4}$$

and combining (3) and (4) proves Lemma 2.

Now, consider the numerator of (2) again. We define

$$\Delta = \Delta(\gamma, \boldsymbol{y}, \boldsymbol{y}_{\text{new}}) = \|\boldsymbol{y}_{\text{new}} - X\tilde{\boldsymbol{\beta}}(\gamma)\|^2 - \|\boldsymbol{y}_{\text{new}} - X\hat{\boldsymbol{\beta}}(\gamma)\|^2,$$

and get

$$|\mathrm{E}(\Delta)| \leq \left| \mathrm{E}\left(\Delta \, \mathrm{I}\left(\|\hat{\boldsymbol{\beta}}\| < C\right)\right) \right| + \left| \mathrm{E}\left(\Delta \, \mathrm{I}\left(\|\hat{\boldsymbol{\beta}}\| \geq C\right)\right) \right|.$$

We now analyse the two components of the right-hand side separately, via two additional lemmas.

**Lemma 3.**
$$\lim_{\gamma \searrow 0} \left| \mathrm{E}\left( \frac{1}{\gamma} \Delta \, \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| < C \right) \right) \right| = 0 \qquad \text{for all } C.$$

**Proof.** We fix $C$. First we re-phrase the numerator:

$$\mathrm{E}\left( \Delta \, \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| < C \right) \right) = \mathrm{E}\left( \left( 2\boldsymbol{y}_{\mathrm{new}}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right) \right.\right.$$
$$\left.\left. - \left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right)^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) + \tilde{\boldsymbol{\beta}}(\gamma) \right) \right) \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| < C \right) \right). \quad (5)$$

The expectation is over the distribution of both $\boldsymbol{y}$ and $\boldsymbol{y}_{\mathrm{new}}$. All the $\boldsymbol{\beta}$ quantities we have depend on $\boldsymbol{y}$ only. So our next step is to remove dependence on $\boldsymbol{y}$, by bounding the expectation by its maximum

$$\left| \mathrm{E}\left( \Delta \, \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| < C \right) \right) \right| \le \Pr\left( \hat{\boldsymbol{\beta}}(\gamma) \ne \tilde{\boldsymbol{\beta}}(\gamma) \right) \max_{\boldsymbol{y}, \|\hat{\boldsymbol{\beta}}\| < C} \left| 2\mathrm{E}_{\boldsymbol{y}_{\mathrm{new}}}\left( \boldsymbol{y}_{\mathrm{new}}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right) \right.\right.$$
$$\left.\left. - \left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right)^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) + \tilde{\boldsymbol{\beta}}(\gamma) \right) \right) \right|. \quad (6)$$

The next step is to bound the two expressions inside the maximum. Denote by $\lambda$ the maximal absolute singular value of $\boldsymbol{X}$ (so $\lambda^2$ is the maximal eigenvalue of $\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}$). Then Lemma 2 gives us

$$\left| \mathrm{E}_{\boldsymbol{y}_{\mathrm{new}}}\left( \boldsymbol{y}_{\mathrm{new}}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right) \right) \right| \le \|\mathrm{E}(\boldsymbol{y}_{\mathrm{new}})\| \lambda M \gamma, \quad (7)$$

$$\left| \left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right)^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) + \tilde{\boldsymbol{\beta}}(\gamma) \right) \right| \le \lambda^2 M \gamma (2C + \gamma M), \quad (8)$$

and combining (6), (7) and (8) with Lemma 1 we get

$$\lim_{\gamma \searrow 0} \frac{|\mathrm{E}(\Delta \, \mathrm{I}(\|\hat{\boldsymbol{\beta}}\| < C))|}{\gamma} \le \lim_{\gamma \searrow 0} \Pr\left( \hat{\boldsymbol{\beta}}(\gamma) \ne \tilde{\boldsymbol{\beta}}(\gamma) \right) \left( 2\|\mathrm{E}(\boldsymbol{y}_{\mathrm{new}})\| \lambda M + \lambda^2 M(2C + \gamma M) \right) = 0.$$

This proves Lemma 3.

**Lemma 4.**
$$\lim_{C \to \infty} \lim_{\gamma \searrow 0} \left| \mathrm{E}\left( \frac{1}{\gamma} \Delta \, \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| \ge C \right) \right) \right| = 0.$$

**Proof.** Using the same algebra as in (5) we can write

$$\left| \mathrm{E}\left( \frac{1}{\gamma} \Delta \, \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| \ge C \right) \right) \right| \le \left| \mathrm{E}\left( 2\boldsymbol{y}_{\mathrm{new}}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right) \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| \ge C \right) \right) \right|$$
$$+ \left| \mathrm{E}\left( \left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right)^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) + \tilde{\boldsymbol{\beta}}(\gamma) \right) \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| \ge C \right) \right) \right|. \quad (9)$$

The first expression is bounded as follows:

$$\left| \mathrm{E}\left( \boldsymbol{y}_{\mathrm{new}}^{\mathsf{T}} \boldsymbol{X}\left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right) \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| \ge C \right) \right) \right| \le \|\mathrm{E}(\boldsymbol{y}_{\mathrm{new}})\| \lambda M \gamma \Pr\left( \|\hat{\boldsymbol{\beta}}\| \ge C \right). \quad (10)$$

Next, (i) $\|\hat{\boldsymbol{\beta}}(\gamma) + \tilde{\boldsymbol{\beta}}(\gamma)\| \le 2\|\hat{\boldsymbol{\beta}}\| + 2M\gamma$, by Lemma 2, and (ii) $\mathrm{E}(\|\hat{\boldsymbol{\beta}}\| \, \mathrm{I}(\|\hat{\boldsymbol{\beta}}\| \ge C)) \to 0$ as $C \to \infty$ for the second term, by the fact that $\mathrm{E}(\|\hat{\boldsymbol{\beta}}\|) \le \sum_j \mathrm{E}(|\hat{\beta}_j|) < \infty$ as $\mathrm{E}(\hat{\beta}_j) = \beta_j^0$, the true parameter. This also implies $\Pr(\|\hat{\boldsymbol{\beta}}\| \ge C) \to 0$.

Thus we can now write

$$\left| \mathrm{E}\left( \left( \hat{\boldsymbol{\beta}}(\gamma) - \tilde{\boldsymbol{\beta}}(\gamma) \right)^{\top} X^{\top} X \left( \hat{\boldsymbol{\beta}}(\gamma) + \tilde{\boldsymbol{\beta}}(\gamma) \right) \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| \geq C \right) \right) \right|$$
$$\leq 2\lambda^2 M \gamma \left( \mathrm{E}\left( \|\hat{\boldsymbol{\beta}}\| \, \mathrm{I}\left( \|\hat{\boldsymbol{\beta}}\| \geq C \right) \right) + 2\gamma M \right), \qquad (11)$$

and putting (9), (10) and (11) together we get

$$\lim_{\gamma \searrow 0} \frac{|\mathrm{E}(\Delta \, \mathrm{I}\, (\|\hat{\boldsymbol{\beta}}\| \geq C \,))|}{\gamma}$$
$$\leq 2\|\mathrm{E}(\boldsymbol{y}_{\mathrm{new}})\| \lambda M \, \mathrm{Pr}(\|\hat{\boldsymbol{\beta}}\| \geq C \,) + 2\lambda^2 M \, \mathrm{E}\left( \|\hat{\boldsymbol{\beta}}\| \, \mathrm{I}\, (\|\hat{\boldsymbol{\beta}}\| \geq C \,) \right) \to 0 \ \text{ as } \ C \to \infty \,,$$

which concludes the proof of Lemma 4.

Putting Lemma 3 and Lemma 4 together completes our proof, since for any $C$ it gives us an upper bound on (2), and this bound converges to 0 as $C \to \infty$.

## References

DONOHO, D., JOHNSTONE, I., KERKYACHAIRAN, G. & PICARD, D. (1995). Wavelet shrinkage: asymptopia? (with Discussion). *J. Roy. Statist. Soc. Ser. B* **57**, 301–369.

EFRON, B., HASTIE, T., JOHNSTONE, I.M. & TIBSHIRANI, R. (2004). Least angle regression (with Discussion). *Ann. Statist.* **32**, 407–499.

FRIEDMAN, J., HASTIE, T., ROSSET, S., TIBSHIRANI, R. & ZHU, J. (2004). Discussion of 'Consistency in boosting' by W. Jiang, G. Lugosi, N. Vayatis & T. Zhang. *Ann. Statist.* **32**, 102–107.

HUANG, F. (2003). A prediction error property of the lasso and its generalization. *Aust. N. Z. J. Stat.* **45**, 217–228.

OSBORNE, M.R., PRESNELL, B. & TURLACH, B. (2000). On the lasso and its dual. *J. Comput. Graph. Statist.* **9**, 319–337.

TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58**, 267–288.