

# Supplementary material. Tractable Bayesian variable selection: beyond normality

## 1. IMOM PRIOR

The product iMOM prior density on  $\theta_\gamma$  (Johnson and Rossell 2012) is given by

$$p_I(\theta_\gamma | \vartheta, \gamma) = \prod_{\gamma_j=1} \frac{(g_\theta \vartheta)^{\frac{1}{2}}}{\sqrt{\pi} \theta_j^2} \exp \left\{ -\frac{g_\theta \vartheta}{\theta_j^2} \right\}, \quad (1)$$

where by default  $g_\theta = 0.133$  assigns  $p(|\theta/\vartheta^{1/2}| > 0.2) = 0.99$ . Regarding the asymmetry parameter  $\tilde{\alpha} = \text{atanh}(\alpha)$ , the prior is  $p_I(\tilde{\alpha} | \gamma_{p+1} = 1) = \tilde{\alpha}^{-2} \sqrt{g_\alpha/\pi} e^{-g_\alpha/\tilde{\alpha}^2}$ , and the default prior dispersions are  $g_\alpha = 0.033$  to obtain  $P(|\tilde{\alpha}| \geq 0.1) = 0.99$  and  $g_\alpha = 0.136$  for  $P(|\alpha| \geq 0.2) = 0.99$ .

## 2. PROOFS

For simplicity, we drop  $\gamma$  from the notation in the proof of Propositions 1-4 and Corollary 1, given that all arguments are conditional on a given model  $\gamma$ .

### 2.1 Proof of Proposition 1

We start by stating a useful lemma stating that positive definite hessian plus continuous gradient guarantees concavity.

**Lemma 1.** *Let  $f(\theta)$  be a function with continuous gradient  $g(\theta)$ , for all  $\theta$ , and negative definite hessian  $H(\theta)$  almost everywhere with respect to the Lebesgue measure. Then,  $f(\theta)$  is strictly concave. If  $H(\theta)$  is negative semidefinite, then  $f(\theta)$  is concave.*

*Proof.* Let  $\theta_1$  and  $\theta_2$  be two arbitrary values and denote  $\theta_w = (1-w)\theta_1 + w\theta_2$  where  $w \in [0, 1]$ . Define  $h(w) = -f(\theta_w)$ , to show that  $f(\theta)$  is concave it suffices to see that  $h(w)$  is convex for

arbitrary  $(w, \theta_1, \theta_2)$ . Straightforward algebra shows that  $\frac{\partial}{\partial w} h(w) = -g(\theta_w)(\theta_2 - \theta_1)$  and further derivation shows that

$$\frac{\partial^2}{\partial w^2} h(w) = -(\theta_2 - \theta_1)^T H(\theta_w)(\theta_2 - \theta_1) > 0,$$

since  $H(\theta)$  is negative definite ( $\geq 0$  for negative semidefinite).

The second derivative  $\frac{\partial^2}{\partial w^2} h(w) > 0$  almost everywhere and the first derivative  $\frac{\partial}{\partial w} h(w)$  is continuous, which implies that  $\frac{\partial}{\partial w} h(w)$  is strictly increasing in  $w$  and hence  $h(w)$  is strictly convex (non-strictly convex when  $H(\theta)$  is negative semidefinite).  $\square$

### Proof of Proposition 1, Part (i)

The gradient  $g_1(\theta, \vartheta, \alpha)$  follows from straightforward algebra, which is obviously continuous with respect to  $\vartheta \in \mathbb{R}^+$  and  $\alpha \in [-1, 1]$ . To see continuity of  $g_1(\theta, \vartheta, \alpha)$  with respect to  $\theta$ , consider increasing a single  $\theta_j$  for some  $j \in \{1, \dots, p\}$  and fix the remaining elements in  $\theta$ , which we denote  $\theta_{(-j)}$ . Also denote  $x_{i(-j)}$  the subvector of  $x_i$  obtained by removing  $x_{ij}$ . Clearly,  $\log L_1(\theta, \vartheta, \alpha)$  is quadratic in  $\theta_j$  with coefficients that stay constant until  $\theta_j$  increases beyond a value  $t$  such that an observation  $i^*$  is added to or removed from  $A(\theta)$ , *i.e.*  $y_{i^*} < x_{i^*(-j)}^T \theta_{(-j)} + x_{i^*j} \theta_j$  for  $\theta_j \leq t$  and  $y_{i^*} > x_{i^*(-j)}^T \theta_{(-j)} + x_{i^*j} \theta_j$  for  $\theta_j > t$ . Taking the limit of the contribution of  $i^*$  to  $\log(L_1(\theta, \vartheta, \alpha))$  as either  $\theta_j \rightarrow t^-$  or  $\theta_j \rightarrow t^+$  we obtain

$$\lim_{\theta_j \rightarrow t^+} \frac{(y_{i^*} - x_{i^*} \theta_i)^2}{(1 + \alpha)^2} = \lim_{\theta_j \rightarrow t^-} \frac{(y_{i^*} - x_{i^*} \theta_i)^2}{(1 - \alpha)^2} = 0,$$

*i.e.*  $\log(L_1(\theta, \vartheta, \alpha))$  is continuous. Similarly, taking the limits for the contribution to the first partial derivative with respect to  $\theta_j$  gives

$$\lim_{\theta_j \rightarrow t^+} \frac{2(y_{i^*} - x_{i^*} \theta_i)}{(1 + \alpha)^2} = \lim_{\theta_j \rightarrow t^-} \frac{2(y_{i^*} - x_{i^*} \theta_i)}{(1 - \alpha)^2} = 0,$$

which proves that  $g_1(\theta, \vartheta, \alpha)$  is continuous.

### Proof of Proposition 1, Part (ii)

The form of  $H_1(\theta, \vartheta, \alpha)$  follows from easy algebra.

Proof of Proposition 1, Part (iii)

We start by noting that the maximum of the asymmetric-normal log-likelihood with respect to  $(\theta, \alpha)$  does not depend on  $\vartheta$ , hence we simply need to see that

$$H = \begin{pmatrix} X^T W^2 X & 2X^T \bar{W}^3 (y - X\theta) \\ 2(y - X\theta)^T \bar{W}^3 X & 3(y - X\theta)^T W^2 (y - X\theta) \end{pmatrix}, \quad (2)$$

is positive definite for almost all  $(\theta, \alpha)$ . Once we show this, by Part (i) and Lemma 1 we have that there is a unique maximum.

To see that  $H$  is positive definite, we shall show that all its leading principal minors are positive. Note that  $X^T W^2 X$  is the gram matrix corresponding to  $WX$  and is hence positive definite when  $\text{rank}(WX) = p$ , or equivalently when  $\text{rank}(X) = p$  given that the effect of  $W$  is to simply re-scale the rows of  $X$ . If  $\text{rank}(WX) < p$  then  $X^T W^2 X$  is positive semidefinite. Therefore, we just need to check that  $\det(H) > 0$ . Now, the usual formula for determinant based on submatrices gives that  $\det(H) = \det(X^T W^2 X) \det(B)$ , where  $B =$

$$\begin{aligned} & 3(y - X\theta)^T W^4 (y - X\theta) - 4(y - X\theta)^T \bar{W}^3 X (X^T W^2 X)^{-1} X^T \bar{W}^3 (y - X\theta) \\ & = 3(y - X\theta)^T W^2 \left( I - \frac{4}{3} \bar{W} X (X^T W^2 X)^{-1} X^T \bar{W} \right) W^2 (y - X\theta), \end{aligned} \quad (3)$$

is a scalar,  $I$  is the  $n \times n$  identity matrix, as usual  $W$  is an  $n \times n$  diagonal matrix with entries  $1/(1 \pm \alpha)^2$  where the  $\pm$  depends on whether  $i \in A(\theta)$  or  $i \notin A(\theta)$ , and similarly  $\bar{W}$  is diagonal with entries  $\pm(1 \pm \alpha)$ . All that is left is to see that  $B > 0$ . For ease of notation let us define  $Z = \bar{W}X$ , given that  $\bar{W}\bar{W} = \text{diag}(1/(1 \pm \alpha)^2) = W^2$  we can write

$$\begin{aligned} B & = 3(y - X\theta)^T W^2 \left( I - \frac{4}{3} Z (Z^T Z)^{-1} Z^T \right) W^2 (y - X\theta) = \\ & 4(y - X\theta)^T W^2 (I - Z (Z^T Z)^{-1} Z^T) W^2 (y - X\theta) - (y - X\theta)^T W^2 W^2 (y - X\theta) > 0 \\ & \Leftrightarrow 4 \frac{(y - X\theta)^T W^2 (I - Z (Z^T Z)^{-1} Z^T) W^2 (y - X\theta)}{(y - X\theta)^T W^2 W^2 (y - X\theta)} - 1 > 0. \end{aligned} \quad (4)$$

To complete the proof, note that  $a = W^2(y - X\theta) \in \mathbb{R}^n$  is simply a vector and that the hat matrix  $Z(Z^T Z)^{-1} Z^T$  is symmetric and idempotent, which implies that it has  $\text{rank}(Z)$  eigenvalues equal to 1 and  $n - \text{rank}(Z)$  eigenvalues equal to 0. Thus  $I - Z(Z^T Z)^{-1} Z^T$  has  $n - \text{rank}(Z)$  eigenvalues equal to 1 and the remaining  $\text{rank}(Z)$  eigenvalues equal to 0. Given that  $n > \text{rank}(Z)$  by assumption,  $I - Z(Z^T Z)^{-1} Z^T$  has at least one non-zero eigenvalue, which allows us to bound

$$\min_{a \in \mathbb{R}^n} \frac{a(I - Z(Z^T Z)^{-1} Z^T)a}{a^T a} \geq 1,$$

which from (4) gives that  $B \geq 3$  and hence that  $H$  is positive definite.

## 2.2 Proof of Proposition 2

Parts (i) and (ii) follow from straightforward algebra. For Part (iii) we first show that  $\log L_2(\theta, \vartheta, \alpha)$  is (non-strictly) concave in  $(\theta, \alpha)$  and then that when  $\text{rank}(X) = p$  it is strictly concave. To see non-strict concavity note that  $-|y_i - x_i^T \theta| / (\sqrt{\vartheta}(1 + \alpha)) = -\max\{y_i - x_i^T \theta, x_i^T \theta - y_i\} / (\sqrt{\vartheta}(1 + \alpha))$  is the maximum of two (non-strictly) concave functions in  $(\theta, \alpha)$  and hence also concave, from which it follows that  $L_3(\theta, \vartheta, \alpha)$  is a sum of concave functions and thus concave.

For ease of notation let  $\eta = (\theta, \vartheta, \alpha)$ , we now show that  $\log L_2(\eta)$  is strictly concave at any arbitrary  $\eta_1 = (\theta_1, \vartheta, \alpha_1)$  as long as  $\text{rank}(X) = p$ . It is useful to note that  $H_2(\theta, \vartheta, \alpha)$  is strictly negative definite in  $\alpha$ , as the corresponding minor  $-2|W^3(y - X\theta)| / \sqrt{\vartheta} < 0$ . From the definition of concavity and continuity of the log-likelihood, if  $\log L_2(\eta)$  were concave but non-strictly concave at  $\eta = \eta_1$  then for some  $\eta_2 = (\theta_2, \vartheta, \alpha_2) \neq \eta_1$  we would have that  $\log L_2(a\eta_1 + (1 - a)\eta_2) = a \log L_2(\eta_1) + (1 - a) \log L_2(\eta_2)$  for all  $a \in [0, 1]$ , *i.e.*  $\log L_2(\eta)$  would be locally linear (in fact, constant) along the direction defined by  $\eta_2 - \eta_1$ , and in particular  $\log L_2(\eta_1) = \log L_2(\eta_2)$ . From its form

$$\log L_2(\eta) = -\frac{n}{2} \log(\vartheta) - \frac{1}{\vartheta} \left( \frac{\sum_{i \in A(\theta)} |y_i - x_i^T \theta|}{1 + \alpha} + \frac{\sum_{i \notin A(\theta)} |y_i - x_i^T \theta|}{1 - \alpha} \right),$$

is locally linear in  $\theta$  but clearly non-linear in  $\alpha$ , implying that  $\alpha_2 = \alpha_1$ . More formally, it is easy to see that for fixed  $\theta_1 \neq \theta_2$  the roots of  $\log L_2(\eta_1) = \log L_2(\eta_2)$  in terms of  $\alpha_2$  are given by the roots of a quadratic polynomial that are not linear in  $\theta_2$ , thus the only possible linear solution is

$\alpha_2 = \alpha_1$ . The problem is hence reduced to showing that there is no  $\theta_2$  sufficiently close to  $\theta_1$  such that

$$|W(y - X\theta_1)| = |W(y - X\theta_2)|, \quad (5)$$

where  $|\cdot|$  denotes the  $L_1$  norm and as usual  $W$  is a diagonal matrix with  $(i, i)$  element  $(1 + \alpha)^{-1}$  if  $i \in A(\theta_1)$  and  $(1 - \alpha)^{-1}$  if  $i \notin A(\theta_1)$ , where we note that  $A(\theta_2) = A(\theta_1)$  for  $\theta_2$  sufficiently close to  $\theta_1$  and thus the same weighting matrix  $W$  can be used in left and right hand sides of (5). Expression (5) is the  $L_1$  error function featuring in median regression with re-scaled  $\tilde{y} = Wy$  and  $\tilde{X} = WX$ , which is concave as long as  $p = \text{rank}(WX) = \text{rank}(X)$ , as we wished to prove.

### 2.3 Proof of Proposition 3

#### Two-piece normal errors ( $k = 1$ )

The proof strategy is as follows: we first show that the average log-likelihood  $M_n(\theta_\gamma, \vartheta, \alpha) = \frac{1}{n} \log L_1(\theta_\gamma, \vartheta, \alpha)$  converges to its expected value  $M(\theta_\gamma, \vartheta, \alpha)$  uniformly across  $(\theta_\gamma, \vartheta, \alpha) \in \Gamma$ , and later show that  $M(\theta_\gamma, \vartheta, \alpha)$  has a unique maximum  $(\theta_\gamma^*, \vartheta_\gamma^*, \alpha_\gamma^*)$ , which jointly satisfy the conditions in Theorem 5.7 from van der Vaart (1998) for consistency of  $(\hat{\theta}_\gamma, \hat{\vartheta}_\gamma, \hat{\alpha}_\gamma) \xrightarrow{P} (\theta_\gamma^*, \vartheta_\gamma^*, \alpha_\gamma^*)$ .

We remark that Condition A3 is met for instance by deterministic sequences  $\{x_i\}$  satisfying the stated positive-definiteness condition and also by  $x_i \stackrel{i.i.d.}{\sim} \Psi$  as long as  $E(x_1 x_1^T) = \Sigma$  for some positive definite  $\Sigma$ , since then  $n^{-1} X^T X \xrightarrow{a.s.} \Sigma$  by the strong law of large numbers, and given that eigenvalues are continuous functions of  $X^T X$  by the continuous mapping theorem  $X^T X$  is positive definite almost surely as  $n \rightarrow \infty$ . Finally,  $\Gamma$  is assumed to contain the maximizer  $(\theta_\gamma^*, \vartheta_\gamma^*, \alpha_\gamma^*)$ .

By the law of large numbers and the *i.i.d.* assumption, we have that  $M_n(\theta_\gamma, \vartheta, \alpha) \xrightarrow{P} M(\theta_\gamma, \vartheta, \alpha)$ ,

for each  $(\theta_\gamma, \vartheta, \alpha) \in \Gamma$ . Next, we prove that the limit  $M$  is finite for all  $(\theta_\gamma, \vartheta, \alpha) \in \Gamma$ .

$$\begin{aligned}
|M(\theta_\gamma, \vartheta, \alpha)| &= \left| \mathbb{E} \left[ \log s_1(y_1 | x_1^T \theta_\gamma, \vartheta, \alpha) \right] \right| \leq \mathbb{E} \left[ \left| \log s_1(y_1 | x_1^T \theta_\gamma, \vartheta, \alpha) \right| \right] \\
&= \int \int \left| \log s_1(y_1 | x_1^T \theta_\gamma, \vartheta, \alpha) \right| dS_0(y_1 | x_1) d\Psi(x_1) \\
&= \int \int_{y_1 < x_1^T \theta_\gamma} \left| \log \frac{1}{\sqrt{\vartheta}} \phi \left( \frac{y_1 - x_1^T \theta_\gamma}{\sqrt{\vartheta}(1 + \alpha)} \right) \right| dS_0(y_1 | x_1) d\Psi(x_1) \\
&+ \int \int_{y_1 \geq x_1^T \theta_\gamma} \left| \log \frac{1}{\sqrt{\vartheta}} \phi \left( \frac{y_1 - x_1^T \theta_\gamma}{\sqrt{\vartheta}(1 - \alpha)} \right) \right| dS_0(y_1 | x_1) d\Psi(x_1).
\end{aligned}$$

For the first term in the last inequality we obtain, by integrating over the whole space, assumption A4 with  $j = 2$ , and the triangle inequality, the following upper bound

$$\begin{aligned}
&\int \int \left| \log \frac{1}{\sqrt{\vartheta}} \phi \left( \frac{y_1 - x_1^T \theta_\gamma}{\sqrt{\vartheta}(1 + \alpha)} \right) \right| dS_0(y_1 | x_1) d\Psi(x_1) \\
&\leq |\log \sqrt{2\pi\vartheta}| + \int \int \frac{(y_1 - x_1^T \theta_\gamma)^2}{2\vartheta(1 + \alpha)^2} dS_0(y_1 | x_1) d\Psi(x_1) < \infty.
\end{aligned}$$

Analogously for the second term. Now, let  $\vartheta = \vartheta^*$  be an arbitrary fixed value for the (squared) scale parameter. The aim now is to first show that the average log-likelihood  $M_n(\theta_\gamma, \vartheta^*, \alpha) = n^{-1} \log L_1(\theta_\gamma, \vartheta^*, \alpha)$  converges to its expected value  $M(\theta_\gamma, \vartheta^*, \alpha)$  uniformly in  $(\theta_\gamma, \alpha)$ , which implies that  $(\hat{\theta}_\gamma, \hat{\alpha}_\gamma) \xrightarrow{P} (\theta_\gamma^*, \alpha_\gamma^*)$ , and to then exploit that  $\hat{\vartheta}_\gamma$  and  $\vartheta_\gamma^*$  have simple expressions to show that  $\hat{\vartheta}_\gamma \xrightarrow{P} \vartheta_\gamma^*$ . To see that  $M_n(\theta_\gamma, \vartheta^*, \alpha)$  converges to  $M(\theta_\gamma, \vartheta^*, \alpha)$  uniformly in  $(\theta_\gamma, \alpha)$  we use the result in Proposition 1 that for positive-definite  $X^T X$  (which holds for  $n > n_0$ ) we have that  $M_n(\theta_\gamma, \vartheta^*, \alpha)$  is a sequence of concave functions in  $(\theta_\gamma, \alpha)$ , which by the convexity lemma in Pollard (1991) (see also Theorem 10.8 from Rockafellar (2015)) implies that

$$\sup_{(\theta_\gamma, \alpha) \in K} |M_n(\theta_\gamma, \vartheta^*, \alpha) - M(\theta_\gamma, \vartheta^*, \alpha)| \xrightarrow{P} 0, \tag{6}$$

for each compact set  $K \subseteq \Gamma$ , and also that  $M(\theta_\gamma, \vartheta^*, \alpha)$  is finite and concave in  $(\theta_\gamma, \alpha)$  and thus has a unique maximum  $(\theta_\gamma^*, \alpha_\gamma^*)$ . That is, for a distance measure  $d(\cdot)$  and every  $\varepsilon > 0$  we have

$$\sup_{d((\theta_\gamma^*, \vartheta^*, \alpha_\gamma^*), (\theta, \vartheta^*, \alpha)) \geq \varepsilon} M(\theta_\gamma, \vartheta^*, \alpha) < M(\theta_\gamma^*, \vartheta^*, \alpha_\gamma^*). \tag{7}$$

The consistency of  $(\widehat{\theta}_\gamma, \widehat{\alpha}_\gamma) \xrightarrow{P} (\theta_\gamma^*, \alpha_\gamma^*)$  follows directly from (6) and (7) together with Theorem 5.7 from van der Vaart (1998). To see that  $\widehat{\vartheta}_\gamma \xrightarrow{P} \vartheta_\gamma^*$ , note first that from

$$M(\theta_\gamma, \vartheta^*, \alpha) = -\log(\sqrt{2\pi\vartheta^*}) - \frac{1}{2\vartheta^*} \int \left[ \frac{(y_1 - x_1^T \theta_\gamma)^2}{(1 + \alpha)^2} I(y_1 < x_1^T \theta_\gamma) + \frac{(y_1 - x_1^T \theta_\gamma)^2}{(1 - \alpha)^2} I(y_1 \geq x_1^T \theta_\gamma) \right] dS_0(y_1|x_1) d\Psi(x_1), \quad (8)$$

we see that  $(\theta_\gamma^*, \alpha_\gamma^*)$  does not depend on  $\vartheta^*$ , thus  $(\theta_\gamma^*, \alpha_\gamma^*)$  is a global maximum. From (8)  $M(\theta_\gamma^*, \vartheta, \alpha_\gamma^*)$  trivially has the maximizer

$$\vartheta_\gamma^* = \int \left[ \frac{(y_1 - x_1^T \theta_\gamma^*)^2}{(1 + \alpha_\gamma^*)^2} I(y_1 < x_1^T \theta_\gamma^*) + \frac{(y_1 - x_1^T \theta_\gamma^*)^2}{(1 - \alpha_\gamma^*)^2} I(y_1 \geq x_1^T \theta_\gamma^*) \right] dS_0(y_1|x_1) d\Psi(x_1),$$

and from the likelihood equations we have that

$$\widehat{\vartheta}_\gamma = \frac{1}{n} \left( \sum_{i=1}^n \frac{(y_i - x_i^T \widehat{\theta}_\gamma)^2}{(1 + \widehat{\alpha}_\gamma)^2} I(y_i \leq x_i^T \widehat{\theta}_\gamma) + \frac{(y_i - x_i^T \widehat{\theta}_\gamma)^2}{(1 - \widehat{\alpha}_\gamma)^2} I(y_i > x_i^T \widehat{\theta}_\gamma) \right). \quad (9)$$

In order to simplify notation, let us define

$$\rho(y_i, x_i, \theta_\gamma, \alpha) = \frac{(y_i - x_i^T \theta_\gamma)^2}{(1 + \alpha)^2} I(y_i \leq x_i^T \theta_\gamma) + \frac{(y_i - x_i^T \theta_\gamma)^2}{(1 - \alpha)^2} I(y_i > x_i^T \theta_\gamma).$$

Then, by the triangle inequality

$$\left| \widehat{\vartheta}_\gamma - \vartheta_\gamma^* \right| \leq \left| \widehat{\vartheta}_\gamma - \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) \right| + \left| \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) - \vartheta_\gamma^* \right|.$$

For the second term it follows, by the law of large numbers, that

$$\left| \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) - \vartheta_\gamma^* \right| \xrightarrow{P} 0.$$

For the first term we have

$$\begin{aligned}
\left| \widehat{\vartheta}_\gamma - \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) \right| &= \left| M_n(\widehat{\theta}_\gamma, 1/2, \widehat{\alpha}_\gamma) - M_n(\theta_\gamma^*, 1/2, \alpha_\gamma^*) \right| \\
&\leq \left| M_n(\widehat{\theta}_\gamma, 1/2, \widehat{\alpha}_\gamma) - M(\widehat{\theta}_\gamma, 1/2, \widehat{\alpha}_\gamma) \right| \\
&\quad + \left| M(\widehat{\theta}_\gamma, 1/2, \widehat{\alpha}_\gamma) - M(\theta_\gamma^*, 1/2, \alpha_\gamma^*) \right| \\
&\quad + \left| M(\theta_\gamma^*, 1/2, \alpha_\gamma^*) - M_n(\theta_\gamma^*, 1/2, \alpha_\gamma^*) \right| \\
&\leq 2 \sup_{(\theta_\gamma, \alpha) \in \Gamma} |M_n(\theta_\gamma, 1/2, \alpha) - M(\theta_\gamma, 1/2, \alpha)| \\
&\quad + \left| M(\widehat{\theta}_\gamma, 1/2, \widehat{\alpha}) - M(\theta_\gamma^*, 1/2, \alpha_\gamma^*) \right|.
\end{aligned}$$

By using (6), the consistency of  $(\widehat{\theta}_\gamma, \widehat{\alpha}_\gamma)$ , and the continuous mapping theorem it follows that  $\left| \widehat{\vartheta}_\gamma - \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) \right| \xrightarrow{P} 0$ . Consequently,  $\widehat{\vartheta} \xrightarrow{P} \vartheta_\gamma^*$ , which completes the proof.

#### Two-piece Laplace errors ( $k = 2$ )

The proof strategy is analogous to that with  $k = 1$ . Denote  $M_n(\theta_\gamma, \vartheta, \alpha) = \frac{1}{n} \log L_2(\theta_\gamma, \vartheta, \alpha)$ . By the law of large numbers, we have that  $M_n(\theta_\gamma, \vartheta, \alpha) \xrightarrow{P} M(\theta_\gamma, \vartheta, \alpha)$ , for each  $(\theta_\gamma, \vartheta, \alpha) \in \Gamma$ . Moreover,

$$\begin{aligned}
|M(\theta_\gamma, \vartheta, \alpha)| &= \left| \mathbb{E} \left[ \log s_2(y_1 | x_1^T \theta_\gamma, \vartheta, \alpha) \right] \right| \leq \mathbb{E} \left[ \left| \log s_2(y_1 | x_1^T \theta_\gamma, \vartheta, \alpha) \right| \right] \\
&= \int \left| \log s_2(y_1 | x_1^T \theta_\gamma, \vartheta, \alpha) \right| dS_0(y_1 | x_1) d\Psi(x_1) \\
&= \int_{y < x^T \theta_\gamma} \left| \log \frac{1}{\sqrt{\vartheta}} f \left( \frac{y_1 - x_1^T \theta_\gamma}{\sqrt{\vartheta}(1 + \alpha)} \right) \right| dS_0(y_1 | x_1) d\Psi(x_1) \\
&\quad + \int_{y_1 \geq x_1^T \theta_\gamma} \left| \log \frac{1}{\sqrt{\vartheta}} f \left( \frac{y_1 - x_1^T \theta_\gamma}{\sqrt{\vartheta}(1 - \alpha)} \right) \right| dS_0(y_1 | x_1) d\Psi(x_1),
\end{aligned}$$

where  $f(z) = 0.5 \exp(-|z|)$ . For the first term in the last inequality we have, by integrating over the whole space and the triangle inequality, the following upper bound

$$\begin{aligned}
&\int \left| \log \frac{1}{\sqrt{\vartheta}} f \left( \frac{y_1 - x_1^T \theta_\gamma}{\sqrt{\vartheta}(1 + \alpha)} \right) \right| dS_0(y_1 | x_1) d\Psi(x_1) \\
&\leq |\log 2\sqrt{\vartheta}| + \int \frac{|y_1 - x_1^T \theta_\gamma|}{\sqrt{\vartheta}(1 + \alpha)} dS_0(y_1 | x_1) d\Psi(x_1) < \infty,
\end{aligned}$$



where the finiteness follows by assumption A4 with  $j = 1$ . An analogous result is obtained for the second term. Now, let  $\vartheta = \vartheta^*$  be an arbitrary fixed value for the (squared) scale parameter. From Proposition 2, it follows that for positive-definite  $X^T X$  (which is guaranteed by assumption A2, for  $n > n_0$ ) we have that  $M_n(\theta_\gamma, \vartheta^*, \alpha)$  is concave in  $(\theta, \alpha)$ , which by the convexity lemma in Pollard (1991) implies that

$$\sup_{(\theta_\gamma, \alpha) \in K} |M_n(\theta_\gamma, \vartheta^*, \alpha) - M(\theta_\gamma, \vartheta^*, \alpha)| \xrightarrow{P} 0, \quad (10)$$

for any compact set  $K \subseteq \Gamma$ , and also that  $M(\theta_\gamma, \vartheta^*, \alpha)$  is concave in  $(\theta_\gamma, \alpha)$  and thus has a unique maximum  $(\theta_\gamma^*, \alpha_\gamma^*)$ . That is, for a distance measure  $d(\cdot)$  and every  $\varepsilon > 0$  we have

$$\sup_{d((\theta_\gamma^*, \vartheta^*, \alpha_\gamma^*), (\theta_\gamma, \vartheta^*, \alpha)) \geq \varepsilon} M(\theta_\gamma, \vartheta^*, \alpha) < M(\theta_\gamma^*, \vartheta^*, \alpha_\gamma^*). \quad (11)$$

The consistency of  $(\widehat{\theta}_\gamma, \widehat{\alpha}_\gamma) \xrightarrow{P} (\theta_\gamma^*, \alpha_\gamma^*)$  follows directly from (10) and (11) together with Theorem 5.7 from van der Vaart (1998). To see that  $\widehat{\vartheta}_\gamma \xrightarrow{P} \vartheta_\gamma^*$ , note first that from

$$\begin{aligned} M(\theta_\gamma, \vartheta^*, \alpha) = & -\log(2\sqrt{\vartheta^*}) - \frac{1}{\sqrt{\vartheta^*}} \int \left[ \frac{|y_1 - x_1^T \theta_\gamma|}{1 + \alpha} I(y_1 < x_1^T \theta_\gamma) \right. \\ & \left. + \frac{|y_1 - x_1^T \theta_\gamma|}{1 - \alpha} I(y_1 \geq x_1^T \theta_\gamma) \right] dS_0(y_1|x_1) \Psi(x_1), \end{aligned} \quad (12)$$

we see that  $(\theta_\gamma^*, \alpha_\gamma^*)$  does not depend on  $\vartheta^*$ , thus  $(\theta_\gamma^*, \alpha_\gamma^*)$  is a global maximum. From (8)  $M(\theta_\gamma^*, \vartheta, \alpha_\gamma^*)$  trivially has the maximizer

$$\vartheta_\gamma^* = \left\{ \int \left[ \frac{|y_1 - x_1^T \theta_\gamma^*|}{1 + \alpha_\gamma^*} I(y_1 < x_1^T \theta_\gamma^*) + \frac{|y_1 - x_1^T \theta_\gamma^*|}{1 - \alpha_\gamma^*} I(y_1 \geq x_1^T \theta_\gamma^*) \right] dS_0(y_1|x_1) d\Psi(x_1) \right\}^2,$$

and from the likelihood equations we have that

$$\widehat{\vartheta}_\gamma = \left[ \frac{1}{n} \left( \sum_{i=1}^n \frac{|y_i - x_i^T \widehat{\theta}_\gamma|}{1 + \widehat{\alpha}_\gamma} I(y_i \leq x_i^T \widehat{\theta}_\gamma) + \frac{|y_i - x_i^T \widehat{\theta}_\gamma|}{1 - \widehat{\alpha}_\gamma} I(y_i > x_i^T \widehat{\theta}_\gamma) \right) \right]^2. \quad (13)$$

Let us define

$$\rho(y_i, x_i, \theta_\gamma, \alpha) = \frac{|y_i - x_i^T \theta_\gamma|}{1 + \alpha} I(y_i \leq x_i^T \theta_\gamma) + \frac{|y_i - x_i^T \theta_\gamma|}{1 - \alpha} I(y_i > x_i^T \theta_\gamma).$$

Then, by the triangle inequality

$$\left| \sqrt{\widehat{\vartheta}_\gamma} - \sqrt{\vartheta_\gamma^*} \right| \leq \left| \sqrt{\widehat{\vartheta}_\gamma} - \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) \right| + \left| \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) - \sqrt{\vartheta_\gamma^*} \right|.$$

For the second term in the right-hand side of the last equation, it follows, by the law of large numbers and the continuous mapping theorem, that

$$\left| \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) - \sqrt{\vartheta_\gamma^*} \right| \xrightarrow{P} 0.$$

For the first term we have

$$\begin{aligned} \left| \sqrt{\widehat{\vartheta}_\gamma} - \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) \right| &= \left| M_n(\widehat{\theta}_\gamma, 1, \widehat{\alpha}_\gamma) - M_n(\theta_\gamma^*, 1, \alpha_\gamma^*) \right| \\ &\leq \left| M_n(\widehat{\theta}_\gamma, 1, \widehat{\alpha}_\gamma) - M(\widehat{\theta}_\gamma, 1, \widehat{\alpha}_\gamma) \right| \\ &\quad + \left| M(\widehat{\theta}_\gamma, 1, \widehat{\alpha}_\gamma) - M(\theta_\gamma^*, 1, \alpha_\gamma^*) \right| \\ &\quad + \left| M(\theta_\gamma^*, 1, \alpha_\gamma^*) - M_n(\theta_\gamma^*, 1, \alpha_\gamma^*) \right| \\ &\leq 2 \sup_{(\theta_\gamma, \alpha) \in \Gamma} |M_n(\theta_\gamma, 1, \alpha) - M(\theta_\gamma, 1, \alpha)| \\ &\quad + \left| M(\widehat{\theta}_\gamma, 1, \widehat{\alpha}_\gamma) - M(\theta_\gamma^*, 1, \alpha_\gamma^*) \right|. \end{aligned}$$

By using (10), the consistency of  $(\widehat{\theta}_\gamma, \widehat{\alpha}_\gamma)$ , and the continuous mapping theorem it follows that  $\left| \sqrt{\widehat{\vartheta}_\gamma} - \frac{1}{n} \sum_{i=1}^n \rho(y_i, x_i, \theta_\gamma^*, \alpha_\gamma^*) \right| \xrightarrow{P} 0$ . Consequently,  $\widehat{\vartheta}_\gamma \xrightarrow{P} \vartheta_\gamma^*$ , which completes the proof.

## 2.4 Proof of Proposition 4

### Two-piece normal errors ( $k = 1$ )

The proof technique consists of showing first that  $\dot{m}_\eta(y_1, x_1)$  is dominated by an  $L^2$  function (square integrable),  $K(y_1, x_1)$ , for  $\eta$  in a neighborhood of  $\eta_\gamma^*$ . Then, we prove that the function  $Pm_\eta$  admits a second-order Taylor expansion at  $\eta_\gamma^*$  and that the matrix  $V_{\eta_\gamma^*}$  is nonsingular. Finally, we appeal

to the consistency result in Proposition 3 in order to apply Theorem 5.23 of van der Vaart (1998) to prove the asymptotic normality of  $\widehat{\eta}_\gamma$ .

We first note that under assumptions A1–A4, where A4 is assumed to be satisfied for  $j = 4$  throughout, Proposition 3 implies the existence and uniqueness of  $\eta_\gamma^*$ . The gradient of  $m_\eta(y_1, x_1)$ , which is given by (i) in Proposition 1 (with  $n = 1$ ), is bounded for all  $\eta \in \Gamma$  and for each  $(y_1, x_1)$ , due to the compactness of  $\Gamma$ . Now, a direct application of the Minkowski inequality implies that  $\|\dot{m}_\eta(y_1, x_1)\|$  is upper bounded by the sum of the absolute values of the entries of  $\dot{m}_\eta(y_1, x_1)$ . Let us now define  $K(y_1, x_1) = \sup_{\eta \in \mathcal{B}_{\eta_\gamma^*}} \|\dot{m}_\eta(y_1, x_1)\|$ , where  $\mathcal{B}_{\eta_\gamma^*} \subset \Gamma$  is any neighborhood of  $\eta_\gamma^*$ , whose projection over  $\theta$  coincides with  $\mathcal{B}_{\theta_\gamma^*}$ . Thus, from the expression of  $\dot{m}_\eta(y_1, x_1)$  together with assumption A4, it follows that

$$\int K(y_1, x_1)^2 dS_0(y_1|x_1) d\Psi(x_1) < \infty,$$

Then, by using the mean value theorem and the Cauchy-Schwartz inequality, it follows that for  $\eta_1, \eta_2 \in \mathcal{B}_{\eta_0}$ , with probability 1,

$$\begin{aligned} |m_{\eta_1}(y_1, x_1) - m_{\eta_2}(y_1, x_1)| &= |\dot{m}_{\eta_\star}(y_1, x_1)^T (\eta_1 - \eta_2)| \\ &\leq \|\dot{m}_{\eta_\star}(y_1, x_1)\| \cdot \|\eta_1 - \eta_2\| \\ &\leq K(y_1, x_1) \cdot \|\eta_1 - \eta_2\|, \end{aligned}$$

where  $\eta_\star = (1 - c)\eta_1 + c\eta_2$ , for some  $c \in (0, 1)$ .

Now, for each  $x_1$ :

$$\begin{aligned} Pm_{\eta|x_1} = \mathbb{E}[m_\eta|x_1] &= -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\vartheta) \\ &\quad - \frac{1}{2\vartheta(1+\alpha)^2} \int_{-\infty}^{x_1^T \theta_\gamma} (y_1 - x_1^T \theta_\gamma)^2 dS_0(y_1|x_1) \\ &\quad - \frac{1}{2\vartheta(1-\alpha)^2} \int_{x_1^T \theta_\gamma}^{\infty} (y_1 - x_1^T \theta_\gamma)^2 dS_0(y_1|x_1). \end{aligned}$$

Thus, the gradient of  $Pm_{\eta|x_1}$  is given by

$$\begin{aligned}\frac{\partial}{\partial\theta_\gamma}Pm_{\eta|x_1} &= -\frac{x_1}{\vartheta(1+\alpha)^2}I_1 + \frac{x_1}{\vartheta(1-\alpha)^2}I_2, \\ \frac{\partial}{\partial\vartheta}Pm_{\eta|x_1} &= -\frac{1}{2\vartheta} + \frac{I_3}{2\vartheta^2(1+\alpha)^2} + \frac{I_4}{2\vartheta^2(1-\alpha)^2}, \\ \frac{\partial}{\partial\alpha}Pm_{\eta|x_1} &= \frac{I_3}{\vartheta(1+\alpha)^3} - \frac{I_4}{\vartheta(1-\alpha)^3},\end{aligned}$$

Then, the second derivative matrix is given by

$$\begin{aligned}\frac{\partial^2}{\partial\theta_\gamma^2}Pm_{\eta|x_1} &= -\frac{x_1x_1^T[(1+\alpha)^2 - 4\alpha S_0(x_1^T\theta_\gamma|x_1)]}{\vartheta(1-\alpha^2)^2}, \\ \frac{\partial^2}{\partial\vartheta^2}Pm_{\eta|x_1} &= \frac{1}{2\vartheta^2} - \frac{I_3}{\vartheta^3(1+\alpha)^2} - \frac{I_4}{\vartheta^3(1-\alpha)^2}, \\ \frac{\partial^2}{\partial\alpha^2}Pm_{\eta|x_1} &= -\frac{3I_3}{\vartheta(1+\alpha)^4} - \frac{3I_4}{\vartheta(1-\alpha)^4}, \\ \frac{\partial^2}{\partial\vartheta\partial\theta_\gamma}Pm_{\eta|x_1} &= \frac{x_1}{\vartheta^2(1+\alpha)^2}I_1 - \frac{x_1}{\vartheta^2(1-\alpha)^2}I_2, \\ \frac{\partial^2}{\partial\alpha\partial\theta_\gamma}Pm_{\eta|x_1} &= \frac{2x_1}{\vartheta(1+\alpha)^3}I_1 + \frac{2x_1}{\vartheta(1-\alpha)^3}I_2, \\ \frac{\partial^2}{\partial\vartheta\partial\alpha}Pm_{\eta|x_1} &= -\frac{I_3}{\vartheta^2(1+\alpha)^3} + \frac{I_4}{\vartheta^2(1-\alpha)^3},\end{aligned}$$

where  $I_1 = \int_{-\infty}^{x_1^T\theta_\gamma} S_0(y_1|x_1)dy_1$ , and  $I_2 = \int_{x_1^T\theta_\gamma}^{\infty} [1 - S_0(y_1|x_1)] dy_1$ ,  $I_3 = \int_{-\infty}^{x_1^T\theta_\gamma} (y_1 - x_1^T\theta_\gamma)^2 dS_0(y_1|x_1)$ , and  $I_4 = \int_{x_1^T\theta_\gamma}^{\infty} (y_1 - x_1^T\theta_\gamma)^2 dS_0(y_1|x_1)$ . These entries are finite for all  $\eta \in \Gamma$  by assumption A4. Note that  $Pm_\eta = \mathbb{E}[Pm_{\eta|x_1}]$ , where the expectation is taken over  $x_1$ . Assumptions A1–A4 together with Proposition (3) imply that  $Pm_\eta$  is finite and that this expectation is concave and has a unique maximum at  $\eta_\gamma^*$ . From assumption A5,

$$\begin{aligned}\left.\frac{\partial}{\partial\theta_\gamma}Pm_\eta\right|_{\eta=\eta_\gamma^*} &= \mathbb{E}\left[\left.\frac{\partial}{\partial\theta_\gamma}Pm_{\eta|x_1}\right|_{\eta=\eta_\gamma^*}\right] = 0, \\ \left.\frac{\partial}{\partial\alpha}Pm_\eta\right|_{\eta=\eta_\gamma^*} &= \mathbb{E}\left[\left.\frac{\partial}{\partial\alpha}Pm_{\eta|x_1}\right|_{\eta=\eta_\gamma^*}\right] = 0,\end{aligned}$$

which in turn imply that  $\frac{\partial^2}{\partial\vartheta\partial\theta_\gamma}Pm_\eta = 0$  and  $\frac{\partial^2}{\partial\vartheta\partial\alpha}Pm_\eta = 0$  at  $\eta = \eta_\gamma^*$ . Thus, the matrix of

second derivatives evaluated at  $\eta_\gamma^*$  has the following structure:

$$V_\eta = \begin{pmatrix} \frac{\partial^2}{\partial \theta_\gamma^2} Pm_\eta & 0 & \frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta \\ 0 & \frac{\partial^2}{\partial \vartheta^2} Pm_\eta & 0 \\ \frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta & 0 & \frac{\partial^2}{\partial \alpha^2} Pm_\eta \end{pmatrix}.$$

Consequently, the determinant of this matrix is given by

$$\det V_\eta = \frac{\partial^2}{\partial \vartheta^2} Pm_\eta \times \det \begin{pmatrix} \frac{\partial^2}{\partial \theta_\gamma^2} Pm_\eta & \frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta \\ \frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta & \frac{\partial^2}{\partial \alpha^2} Pm_\eta \end{pmatrix}.$$

The determinant on the right-hand side of this expression, evaluated at  $\eta_\gamma^*$ , is non-zero since the  $Pm_\eta$  is concave with respect to  $(\theta_\gamma, \alpha)$ , as shown in Proposition 3. Moreover, the fact that the first derivative  $\frac{\partial}{\partial \vartheta} Pm_\eta = 0$  at  $\eta = \eta_\gamma^*$  together with the fact that  $\eta_\gamma^*$  is the unique maximizer implies that  $\frac{\partial^2}{\partial \vartheta^2} Pm_\eta \neq 0$ . Consequently, the matrix of second derivatives of  $Pm_\eta$  is nonsingular at  $\eta_\gamma^*$ . The asymptotic normality result follows by Theorem 5.23 from van der Vaart (1998).

#### Two-piece Laplace errors ( $k = 2$ )

First, we note that under assumptions A1–A4, where  $j = 2$  in A4 throughout, Proposition 3 implies the existence and uniqueness of  $\eta_\gamma^*$ . The gradient of  $m_\eta(y_1, x_1)$ , which is given by (i) in Proposition 2 (with  $n = 1$ ), is bounded for almost all  $\eta \in \Gamma$  and for each  $(y_1, x_1)$ , due to the compactness of  $\Gamma$ . Now, a direct application of the Minkowski inequality implies that  $\|\dot{m}_\eta(y_1, x_1)\|$  is upper bounded almost surely by the sum of the absolute values of the entries of  $\dot{m}_\eta(y_1, x_1)$ . Let us now define  $K(y_1, x_1) = \sup_{\eta \in \mathcal{B}_{\eta_\gamma^*}} \|\dot{m}_\eta(y_1, x_1)\|$ , where  $\mathcal{B}_{\eta_\gamma^*} \subset \Gamma$  is any neighborhood of  $\eta_\gamma^*$ , whose projection over  $\theta_\gamma$  coincides with  $\mathcal{B}_{\theta_\gamma^*}$ . Thus, from the expression of  $\dot{m}_\eta(y_1, x_1)$  together with assumption A4, it follows that

$$\int K(y_1, x_1)^2 dS_0(y_1|x_1) d\Psi(x_1) < \infty,$$

Then, by using the mean value theorem and the Cauchy-Schwartz inequality, it follows that for  $\eta_1, \eta_2 \in \mathcal{B}_{\eta_\gamma^*}$ , with probability 1,

$$\begin{aligned} |m_{\eta_1}(y_1, x_1) - m_{\eta_2}(y_1, x_1)| &= |\dot{m}_{\eta_\star}(y_1, x_1)^T(\eta_1 - \eta_2)| \\ &\leq \|\dot{m}_{\eta_\star}(y_1, x_1)\| \cdot \|\eta_1 - \eta_2\| \\ &\leq K(y_1, x_1) \cdot \|\eta_1 - \eta_2\|, \end{aligned}$$

where  $\eta_\star = (1 - c)\eta_1 + c\eta_2$ , for some  $c \in (0, 1)$ .

Now, for each  $x_1$ :

$$\begin{aligned} Pm_{\eta|x_1} = \mathbb{E}[m_\eta|x_1] &= -\log(2) - \frac{1}{2}\log(\vartheta) - \frac{1}{\sqrt{\vartheta}(1+\alpha)} \int_{-\infty}^{x_1^T \theta_\gamma} S_0(y_1|x_1) dy_1 \\ &\quad - \frac{1}{\sqrt{\vartheta}(1-\alpha)} \int_{x_1^T \theta_\gamma}^{\infty} 1 - S_0(y_1|x_1) dy_1. \end{aligned}$$

Then, the gradient of  $Pm_{\eta|x_1}$  is given by

$$\begin{aligned} \frac{\partial}{\partial \theta_\gamma} Pm_{\eta|x_1} &= -\frac{x_1 S_0(x_1^T \theta_\gamma|x_1)}{\sqrt{\vartheta}(1+\alpha)} + \frac{x_1 [1 - S_0(x_1^T \theta_\gamma|x_1)]}{\sqrt{\vartheta}(1-\alpha)}, \\ \frac{\partial}{\partial \vartheta} Pm_{\eta|x_1} &= -\frac{1}{2\vartheta} + \frac{I_1}{2\vartheta^{3/2}(1+\alpha)} + \frac{I_2}{2\vartheta^{3/2}(1-\alpha)}, \\ \frac{\partial}{\partial \alpha} Pm_{\eta|x_1} &= \frac{I_1}{\sqrt{\vartheta}(1+\alpha)^2} - \frac{I_2}{\sqrt{\vartheta}(1-\alpha)^2}, \end{aligned}$$

where  $I_1 = \int_{-\infty}^{x_1^T \theta_\gamma} S_0(y_1|x_1) dy$ , and  $I_2 = \int_{x_1^T \theta_\gamma}^{\infty} 1 - S_0(y_1|x_1) dy$ , which are finite by assumption A4.

Then, the second derivative matrix is given by

$$\begin{aligned}
\frac{\partial^2}{\partial \theta_\gamma^2} Pm_{\eta|x_1} &= -\frac{2x_1 x_1^T s_0(x_1^T \theta_\gamma | x_1)}{\sqrt{\vartheta}(1-\alpha^2)}, \\
\frac{\partial^2}{\partial \vartheta^2} Pm_{\eta|x_1} &= \frac{1}{2\vartheta^2} - \frac{3I_1}{4\vartheta^{5/2}(1+\alpha)} - \frac{3I_2}{4\vartheta^{5/2}(1-\alpha)}, \\
\frac{\partial^2}{\partial \alpha^2} Pm_{\eta|x_1} &= -\frac{2I_1}{\sqrt{\vartheta}(1+\alpha)^3} - \frac{2I_2}{\sqrt{\vartheta}(1-\alpha)^3}, \\
\frac{\partial^2}{\partial \vartheta \partial \theta_\gamma} Pm_{\eta|x_1} &= \frac{x_1 S_0(x_1^T \theta_\gamma | x_1)}{2\vartheta^{3/2}(1+\alpha)} - \frac{x_1 [1 - S_0(x_1^T \theta_\gamma | x_1)]}{2\vartheta^{3/2}(1-\alpha)}, \\
\frac{\partial^2}{\partial \alpha \partial \theta_\gamma} Pm_{\eta|x_1} &= \frac{x_1 S_0(x_1^T \theta_\gamma | x_1)}{\sqrt{\vartheta}(1+\alpha)^2} + \frac{x_1 [1 - S_0(x_1^T \theta_\gamma | x_1)]}{\sqrt{\vartheta}(1-\alpha)^2}, \\
\frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_{\eta|x_1} &= -\frac{I_1}{2\vartheta^{3/2}(1+\alpha)^2} + \frac{I_2}{2\vartheta^{3/2}(1-\alpha)^2}.
\end{aligned}$$

These entries are finite for all  $\eta \in \Gamma$  by assumption A4. Note that  $Pm_\eta = \mathbb{E}[Pm_{\eta|x_1}]$ , where the expectation is taken over  $x_1$ . Assumptions A1–A4 together with Proposition 3, imply that  $Pm_\eta$  is finite and that this expectation is concave and has a unique maximum at  $\eta_\gamma^*$ . From assumption A5,

$$\begin{aligned}
\left. \frac{\partial}{\partial \theta_\gamma} Pm_\eta \right|_{\eta=\eta_\gamma^*} &= \mathbb{E} \left[ \left. \frac{\partial}{\partial \theta_\gamma} Pm_{\eta|x_1} \right|_{\eta=\eta_\gamma^*} \right] = 0, \\
\left. \frac{\partial}{\partial \alpha} Pm_\eta \right|_{\eta=\eta_\gamma^*} &= \mathbb{E} \left[ \left. \frac{\partial}{\partial \alpha} Pm_{\eta|x_1} \right|_{\eta=\eta_\gamma^*} \right] = 0,
\end{aligned}$$

which in turn imply that  $\frac{\partial^2}{\partial \vartheta \partial \theta_\gamma} Pm_\eta = 0$  and  $\frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta = 0$  at  $\eta = \eta_\gamma^*$ . Thus, it follows that the matrix of second derivatives evaluated at  $\eta_\gamma^*$  has the structure:

$$V_\eta = \begin{pmatrix} \frac{\partial^2}{\partial \theta_\gamma^2} Pm_\eta & 0 & \frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta \\ 0 & \frac{\partial^2}{\partial \vartheta^2} Pm_\eta & 0 \\ \frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta & 0 & \frac{\partial^2}{\partial \alpha^2} Pm_\eta \end{pmatrix}.$$

Consequently, the determinant of this matrix is given by

$$\det V_\eta = \frac{\partial^2}{\partial \vartheta^2} Pm_\eta \times \det \begin{pmatrix} \frac{\partial^2}{\partial \theta_\gamma^2} Pm_\eta & \frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta \\ \frac{\partial^2}{\partial \vartheta \partial \alpha} Pm_\eta & \frac{\partial^2}{\partial \alpha^2} Pm_\eta \end{pmatrix}.$$

The determinant on the right-hand side of this expression, evaluated at  $\eta_\gamma^*$ , is non-zero since the  $Pm_\eta$  is concave with respect to  $(\theta_\gamma, \alpha)$ , as shown in Proposition 3. Moreover, the fact that the first derivative  $\frac{\partial}{\partial \vartheta} Pm_\eta = 0$  at  $\eta = \eta_\gamma^*$  together with the fact that  $\eta_\gamma^*$  is the unique maximizer implies that  $\frac{\partial^2}{\partial \vartheta^2} Pm_\eta \neq 0$ . Consequently, the matrix of second derivatives of  $Pm_\eta$  is nonsingular at  $\eta_\gamma^*$ . The asymptotic normality result follows by Theorem 5.23 from van der Vaart (1998).

## 2.5 Proof of Corollary 1

The result when  $\epsilon_i \sim L(0, \vartheta)$  follows directly from Pollard (1991) Theorem 1, hence it suffices to find the expression for  $f_0$  under each assumed residual distribution. The median for a general two-piece distribution with a mode at 0 is given by  $\sqrt{\vartheta}(1+\alpha)F^{-1}\left(\frac{1}{2(1+\alpha)}\right)$  if  $\alpha > 0$  and  $\sqrt{\vartheta}(1-\alpha)F^{-1}\left(\frac{(1-2\alpha)}{2(1-\alpha)}\right)$  if  $\alpha \leq 0$ , where  $F(\cdot)$  is the cdf of the standard underlying distribution with mode 0,  $\vartheta = 1$  (Arellano-Valle et al. (2005), Expression (9)).

When  $\epsilon_i \sim N(0, \vartheta)$  we have  $m = 0$  and hence  $f_0 = N(0; 0, \vartheta) = 1/(\sqrt{2\pi\vartheta})$ . When  $\epsilon_i \sim \text{AN}(0, \vartheta, \alpha)$  we have  $m = \sqrt{\vartheta}(1+\alpha)\Phi^{-1}(0.5/(1+\alpha))$  if  $\alpha > 0$  and  $m = \sqrt{\vartheta}(1-\alpha)\Phi^{-1}(0.5(1-2\alpha)/(1-\alpha))$  if  $\alpha < 0$ , where  $\Phi^{-1}(\cdot)$  is the inverse standard cdf, and hence  $f_0 = \exp\left\{-\frac{1}{2}\left(\Phi^{-1}\left(\frac{0.5}{1+|\alpha|}\right)\right)^2\right\} \frac{1}{\sqrt{2\pi\vartheta}}$ . For the Laplace and Asymmetric Laplace, we note that the inverse cdf of the standard Laplace distribution evaluated at a quantile  $q \in [0, 1]$  is  $F^{-1}(q) = \log(2q)$  if  $q < 0.5$  and  $F^{-1}(q) = -\log(2(1-q))$  if  $q \geq 0.5$ . When  $\epsilon_i \sim L(0, \vartheta)$  we have  $m = 0$  and  $f_0 = 1/(2\sqrt{\vartheta})$ . Finally, when  $\epsilon_i \sim \text{AL}(0, \vartheta, \alpha)$  we have  $m = -\sqrt{\vartheta}(1+\alpha)\log(1+\alpha)$  if  $\alpha > 0$  and  $m = \sqrt{\vartheta}(1-\alpha)\log(1-\alpha)$  if  $\alpha < 0$ , from which it follows that  $f_0 = \frac{1}{2\sqrt{\vartheta}} \exp\{-\log(1+|\alpha|)\} = \frac{1}{2\sqrt{\vartheta}(1+|\alpha|)}$ .

The results for the true Normal model follows by using classic asymptotic results on least square estimators (see *e.g.* Newey and Powell (1987))



## 2.6 Proof of Proposition 5

We provide the proof for the asymmetric Normal and asymmetric Laplace ( $\alpha \neq 0$ ), their symmetric counterparts follow as particular cases.  $\eta_\gamma = (\theta_\gamma, \vartheta_\gamma, \alpha_\gamma)$  denotes the parameter vector under model  $\gamma$ ,  $\hat{\eta}_\gamma$  the MLE and  $\tilde{\eta}_\gamma$  the posterior mode for a given observed  $(y, X)$ . Further,  $M_k(\eta_\gamma) = E(\log L_k(\eta_\gamma))$  where the expectation is with respect to the data-generating truth and  $\eta_{\gamma^*}^* = \arg \max_{\eta \in \Gamma_\gamma} M_k(\eta_\gamma)$  is the optimal parameter value under  $\gamma$ . We wish to characterize the asymptotic behaviour of the Laplace-approximated Bayes factors

$$\frac{\hat{p}(y | \gamma)}{\hat{p}(y | \gamma^*)} = e^{\log L_k(\tilde{\eta}_\gamma) - \log L_k(\tilde{\eta}_{\gamma^*})} \times \frac{p(\tilde{\eta}_\gamma | \gamma)}{p(\tilde{\eta}_{\gamma^*} | \gamma^*)} \times (2\pi)^{\frac{p_\gamma - p_{\gamma^*}}{2}} \times \frac{|H_k(\tilde{\eta}_{\gamma^*})|^{\frac{1}{2}}}{|H_k(\tilde{\eta}_\gamma)|^{\frac{1}{2}}}, \quad (14)$$

when  $(y, X)$  arise from the data-generating model in Condition A1, which may differ from the assumed model. The term  $(2\pi)^{\frac{p_\gamma - p_{\gamma^*}}{2}}$  is a constant since  $p_\gamma$  and  $p_{\gamma^*}$  are fixed. The expression for  $H_1$  is given by (24) and recall that for  $H_2$  we are taking the asymptotic covariance in (26). Hence

$$|H_2(\tilde{\eta}_\gamma)| = n^{p_\gamma} \left| \frac{1}{n} H_2(\tilde{\eta}_\gamma) \right| = n^{p_\gamma} \begin{vmatrix} -\frac{1}{n} X^T X \frac{1}{\tilde{\vartheta}_\gamma(1 - \tilde{\alpha}_\gamma^2)} & \frac{\bar{x}}{\sqrt{\tilde{\vartheta}_\gamma(1 - \tilde{\alpha}_\gamma^2)}} & 0 \\ \frac{\bar{x}}{\sqrt{\tilde{\vartheta}_\gamma(1 - \tilde{\alpha}_\gamma^2)}} & -\frac{1}{4\tilde{\vartheta}_\gamma^2} & 0 \\ 0 & 0 & -\frac{2}{1 - \tilde{\alpha}_\gamma^2} \end{vmatrix}.$$

The determinant converges in probability to a negative constant since  $\tilde{\eta}_\gamma \xrightarrow{P} \eta_{\gamma^*}^*$  by Proposition 3, together with the continuous mapping theorem and the asymptotic Hessian (the limiting  $-H_2$ ) being positive definite. An analogous argument applies to  $H_1$ , hence  $n^{\frac{p_\gamma - p_{\gamma^*}}{2}} |H_k(\tilde{\eta}_{\gamma^*})|^{\frac{1}{2}} / |H_k(\tilde{\eta}_\gamma)|^{\frac{1}{2}} \xrightarrow{P} \tilde{a}_3$  for some constant  $\tilde{a}_3 > 0$ . In other words,  $|H_k(\tilde{\eta}_{\gamma^*})|^{\frac{1}{2}} / |H_k(\tilde{\eta}_\gamma)|^{\frac{1}{2}} = O_p\left(n^{\frac{p_{\gamma^*} - p_\gamma}{2}}\right)$ .

The proof strategy is to first show that when  $M_k(\eta_{\gamma^*}^*) - M_k(\eta_\gamma^*) < 0$  (*i.e.*  $\gamma^* \not\subset \gamma$ ) the log-first term of the right hand in (14) behaves asymptotically in probability as  $-na_1$ , for some constant  $a_1 > 0$ , and the logarithm of the second term converges in probability to a constant  $a_2$ . Thus,

$$\frac{1}{n} \log \left( \frac{\hat{p}(y | \gamma)}{\hat{p}(y | \gamma^*)} \right) = -a_1(1 + o_p(1)) + \frac{1}{n} \left( a_2 + \frac{p_\gamma - p_{\gamma^*}}{2} (a_3 - \log(n)) \right) \xrightarrow{P} -a_1$$

where  $a_3 = \log(2\pi) - \log(\tilde{a}_3)$ , as we wish to prove. Subsequently we shall show that when  $M_k(\eta_{\gamma^*}^*) - M_k(\eta_\gamma^*) = 0$  (the case  $\gamma^* \subset \gamma$ ) the first term is essentially the likelihood ratio test statistic and is

$O_p(1)$ , whereas, analogously to the results in Johnson and Rossell (2010) and Rossell and Telesca (2017), the second term converges to a positive constant under local priors, but it is  $O_p(\tilde{b}_n)$  where  $\tilde{b}_n = n^{p_{\gamma^*} - p_\gamma}$  under the pMOM prior and  $\tilde{b}_n = e^{-c\sqrt{n}}$  for some  $c > 0$  under the peMOM prior. This gives

$$\frac{\widehat{p}(y | \gamma)}{\widehat{p}(y | \gamma^*)} = e^{O_p(1)} O_p(\tilde{b}_n) O_p\left(n^{\frac{p_{\gamma^*} - p_\gamma}{2}}\right) = O_p(b_n)$$

where  $b_n = n^{\frac{p_{\gamma^*} - p_\gamma}{2}}$  for local priors,  $b_n = n^{3(p_{\gamma^*} - p_\gamma)/2}$  for the pMOM prior and  $b_n = e^{-c\sqrt{n}} n^{\frac{p_{\gamma^*} - p_\gamma}{2}}$  for the peMOM prior, as we wish to prove.

Consider first the case when  $\gamma^* \not\subset \gamma$ , which implies  $M_k(\eta_\gamma^*) - M_k(\eta_{\gamma^*}^*) < 0$ . Then by continuity of  $p(\eta_\gamma | \gamma)$  we have that  $p(\tilde{\eta}_\gamma | \gamma) \xrightarrow{P} p(\eta_\gamma^* | \gamma) \geq 0$ , and analogously  $p(\tilde{\eta}_{\gamma^*} | \gamma^*) \xrightarrow{P} p(\eta_{\gamma^*}^* | \gamma^*) > 0$  (strict positivity is ensured by the assumption of prior positivity at  $\eta_{\gamma^*}^*$ ). Hence  $p(\tilde{\eta}_\gamma | \gamma)/p(\tilde{\eta}_{\gamma^*} | \gamma^*) \xrightarrow{P} a_2$  for some constant  $a_2 \geq 0$ . Note that  $a_2 = 0$  when  $\theta_\gamma^*$  contains some zeroes and hence a non-local prior would take the value  $p(\eta_\gamma^* | \gamma) = 0$ , but this gives even faster Bayes factor rates in favor of  $\gamma^*$ . Regarding  $\log L_k(\tilde{\eta}_\gamma) - \log L_k(\tilde{\eta}_{\gamma^*})$ , the law of large numbers and uniform convergence of  $\log L_k$  to its expected value shown in Proposition 3 give that

$$\frac{1}{n} (\log L_k(\tilde{\eta}_\gamma) - \log L_k(\tilde{\eta}_{\gamma^*})) \xrightarrow{P} (M_k(\eta_\gamma^*) - M_k(\eta_{\gamma^*}^*)) < 0, \quad (15)$$

hence the constant  $a_1$  defined above is  $a_1 = M_k(\eta_{\gamma^*}^*) - M_k(\eta_\gamma^*) > 0$ .

Next consider the case when  $\gamma^* \subset \gamma$ , which implies  $M_k(\eta_\gamma^*) - M_k(\eta_{\gamma^*}^*) = 0$ . Since  $\tilde{\eta}_\gamma \xrightarrow{P} \eta_\gamma^*$  by Proposition 3, we have that under a local prior

$$\frac{p(\tilde{\eta}_\gamma | \gamma)}{p(\tilde{\eta}_{\gamma^*} | \gamma^*)} \xrightarrow{P} \frac{p(\eta_\gamma | \gamma)}{p(\eta_{\gamma^*} | \gamma^*)} > 0. \quad (16)$$

Under a non-local prior we still have  $p(\eta_{\gamma^*} | \gamma^*) > 0$  but in contrast  $p(\eta_\gamma | \gamma) = 0$ . Thus, it is necessary to characterize the rate at which the latter term vanishes. Briefly, following the proof of Theorem 1 in Koenker and Bassett (1982), the fact that  $\log L_k$  converges uniformly to its expectation (see the proof of our Proposition 3) and consistency of  $\tilde{\eta}_\gamma \xrightarrow{P} \eta_\gamma^*$  give that  $\log L_k$  can be approximated by a quadratic function plus a term that is  $o_p(1)$ . Then, the argument leading to Rossell and Telesca (2017), Proposition 2(i), gives that  $\tilde{\theta}_{\gamma_j} - \hat{\theta}_{\gamma_j} = O_p(1/n)$  and thus

$\tilde{\theta}_{\gamma j} = O_p(n^{-1/2})$  under the pMOM prior  $p_M$ , whereas  $\tilde{\theta}_{\gamma j} = O_p(n^{-1/4})$  under the peMOM prior  $p_E$ . It follows that  $\pi_M(\tilde{\theta}_\gamma) = O_p(1) \prod_{\theta_{\gamma j}^* \neq 0} \tilde{\theta}_{\gamma j}^2 = O_p(n^{-(p_\gamma - p_{\gamma^*})})$ , and  $\pi_E(\tilde{\eta}) = O_p(1) \prod_{\theta_{\gamma j}^* \neq 0} e^{O_p(1)/\tilde{\theta}_{\gamma j}^2} = O_p(e^{-c\sqrt{n}})$  for some  $c > 0$ , as desired.

To conclude the proof, since  $\log L_k(\tilde{\eta}_\gamma) - \log L_k(\tilde{\eta}_{\gamma^*}) = \lambda(y) + o_p(1)$  where  $\lambda(y) = \log L_k(\hat{\eta}_\gamma) - \log L_k(\hat{\eta}_{\gamma^*})$  is the likelihood ratio (LR) statistic, it only remains to show that  $\lambda(y) = O_p(1)$ . The strategy is to see that  $\lambda(y) = \lambda(y; \vartheta_\gamma^*)(1 + o_p(1))$ , where  $\lambda(y; \vartheta_\gamma^*) = \log L_k(\hat{\theta}_\gamma, \vartheta_\gamma^*, \hat{\alpha}_\gamma) - \log L_k(\hat{\theta}_{\gamma^*}, \vartheta_\gamma^*, \hat{\alpha}_{\gamma^*})$  is the LR obtained by plugging in the oracle  $\vartheta_\gamma^* = \vartheta_{\gamma^*}^*$ , then use classical results to prove that  $\lambda(y; \vartheta_\gamma^*) = O_p(1)$ . Taking derivatives of the likelihoods (Expressions (3) and (4) in the main paper) shows that for  $k = 1$  the MLE must satisfy

$$\hat{\vartheta}_\gamma = \frac{1}{n} \left( \sum_{i \in A(\theta)} \frac{(y_i - x_i^T \hat{\theta}_\gamma)^2}{(1 + \hat{\alpha})^2} + \sum_{i \notin A(\theta)} \frac{(y_i - x_i^T \hat{\theta}_\gamma)^2}{(1 - \hat{\alpha})^2} \right) = \frac{1}{n} (y - X_\gamma \hat{\theta}_\gamma)^T W_{\hat{\theta}_\gamma, \hat{\alpha}}^2 (y - X_\gamma \hat{\theta}_\gamma),$$

whereas for  $k = 2$  it satisfies

$$\hat{\vartheta}_\gamma^{\frac{1}{2}} = \frac{1}{n} \left( \sum_{i \in A(\theta)} \frac{|y_i - x_i^T \hat{\theta}_\gamma|}{(1 + \hat{\alpha})} + \sum_{i \notin A(\theta)} \frac{|y_i - x_i^T \hat{\theta}_\gamma|}{(1 - \hat{\alpha})} \right) = \frac{1}{n} |W_{\hat{\theta}_\gamma, \hat{\alpha}}^{\frac{1}{2}} (y - X_\gamma \hat{\theta}_\gamma)|.$$

Plugging  $\hat{\vartheta}_\gamma$  into the likelihoods gives

$$\begin{aligned} \lambda(y) &= -\frac{n}{2} \log \left( \frac{\hat{\vartheta}_\gamma}{\hat{\vartheta}_{\gamma^*}} \right) = -\frac{n}{2} \log \left( 1 + \frac{\hat{\vartheta}_\gamma - \hat{\vartheta}_{\gamma^*}}{\hat{\vartheta}_{\gamma^*}} \right) = -\frac{n}{2} \frac{\hat{\vartheta}_\gamma - \hat{\vartheta}_{\gamma^*}}{\hat{\vartheta}_{\gamma^*}} (1 + o_p(1)) \\ &= -\frac{n}{2} \frac{\hat{\vartheta}_\gamma - \hat{\vartheta}_{\gamma^*}}{\vartheta_{\gamma^*}^*} (1 + o_p(1)) = \lambda(y; \vartheta_{\gamma^*}^*) (1 + o_p(1)) \end{aligned} \quad (17)$$

since by Proposition 3  $\hat{\vartheta}_{\gamma^*} \xrightarrow{P} \vartheta_{\gamma^*}^* > 0$  and  $(\hat{\vartheta}_\gamma - \hat{\vartheta}_{\gamma^*})/\hat{\vartheta}_{\gamma^*} \xrightarrow{P} 0$ .

Finally we show that  $\lambda(y; \vartheta_{\gamma^*}^*) = O_p(1)$ , which implies  $\lambda(y; \vartheta_\gamma^*)(1 + o_p(1)) = O_p(1)$  and completes the proof. For ease of notation when  $k = 1$  define  $Z_n(\gamma) = (y - X_\gamma \hat{\theta}_\gamma)^T W_{\hat{\theta}_\gamma, \hat{\alpha}}^2 (y - X_\gamma \hat{\theta}_\gamma)$  and  $Z(\gamma) = (y - X_\gamma \hat{\theta}_\gamma)^T W_{\hat{\theta}_\gamma, \hat{\alpha}}^2 (y - X_\gamma \hat{\theta}_\gamma)$ , and when  $k = 2$  let  $Z_n(\gamma) = |W_{\hat{\theta}_\gamma, \hat{\alpha}}^{\frac{1}{2}} (y - X_\gamma \hat{\theta}_\gamma)|$ ,  $Z(\gamma) = |W_{\hat{\theta}_\gamma, \hat{\alpha}}^{\frac{1}{2}} (y - X_\gamma \hat{\theta}_\gamma)|$ . Then by definition

$$\lambda(y; \vartheta_{\gamma^*}^*) = \frac{Z_n(\gamma^*) - Z_n(\gamma)}{2\vartheta_{\gamma^*}^*}. \quad (18)$$

Now, note that  $Z_n(\gamma) = Z(\gamma) + Z(\gamma)(Z_n(\gamma) - Z(\gamma))/Z(\gamma) = Z(\gamma)(1 + o_p(1))$ , since Proposition 3 gives that  $\frac{1}{n}Z_n(\gamma) \xrightarrow{P} \vartheta_\gamma$ ,  $\frac{1}{n}Z(\gamma) \xrightarrow{P} \vartheta_\gamma$  and hence  $(Z_n(\gamma) - Z(\gamma))/Z(\gamma) \xrightarrow{P} 0$ . Following the same argument  $Z_n(\gamma^*) = Z(\gamma^*)(1 + o_p(1))$ , hence

$$\lambda(y; \vartheta_\gamma^*) = \frac{Z(\gamma^*) - Z(\gamma) + o_p(Z(\gamma^*) - Z(\gamma))}{2\vartheta_\gamma^*} = \frac{Z(\gamma^*) - Z(\gamma)}{2\vartheta_\gamma^*}(1 + o_p(1)). \quad (19)$$

The term  $(Z(\gamma^*) - Z(\gamma))/\vartheta_\gamma^*$  is the LR test statistic for fixed  $(\vartheta_\gamma^*, \alpha_\gamma^*)$  comparing  $\gamma$  and  $\gamma^* \subset \gamma$ .

When  $k = 2$  this is a quantile regression LR test statistic, which Koenker and Bassett (1982) showed to be asymptotically  $\chi_{p_\gamma - p_{\gamma^*}}^2$  (after rescaling by a constant) precisely under our Conditions A2-A3. When  $k = 1$ ,  $(Z(\gamma^*) - Z(\gamma))/\vartheta_\gamma^*$  is the LR test statistic for a weighted least squares problem regressing  $\tilde{y} = W_{\theta^*, \alpha^*}y$  on  $\tilde{X} = W_{\theta^*, \alpha^*}X$ , which can be shown to be  $O_p(1)$  under the conditions in Proposition 4. Briefly, as usual for any  $\gamma$  the total sum of squares can be decomposed as  $\tilde{y}^T \tilde{y} = \hat{\theta}_\gamma^T \tilde{X}_\gamma^T \tilde{X}_\gamma \hat{\theta}_\gamma + (\tilde{y} - \tilde{X}_\gamma \hat{\theta}_\gamma)^T (\tilde{y} - \tilde{X}_\gamma \hat{\theta}_\gamma)$ , hence  $Z(\gamma^*) - Z(\gamma) =$

$$(\tilde{y} - \tilde{X}_{\gamma^*} \hat{\theta}_{\gamma^*})^T (\tilde{y} - \tilde{X}_{\gamma^*} \hat{\theta}_{\gamma^*}) - (\tilde{y} - \tilde{X}_\gamma \hat{\theta}_\gamma)^T (\tilde{y} - \tilde{X}_\gamma \hat{\theta}_\gamma) = \hat{\theta}_\gamma^T \tilde{X}_\gamma^T \tilde{X}_\gamma \hat{\theta}_\gamma - \hat{\theta}_{\gamma^*}^T \tilde{X}_{\gamma^*}^T \tilde{X}_{\gamma^*} \hat{\theta}_{\gamma^*}. \quad (20)$$

Without loss of generality let  $\tilde{X}_\gamma = (\tilde{X}_{\gamma^*}, \tilde{X}_{\gamma \setminus \gamma^*})$ , where  $\tilde{X}_{\gamma \setminus \gamma^*}$  are the columns in  $\tilde{X}_\gamma$  not contained in  $\tilde{X}_{\gamma^*}$ . Let  $R = (I - \tilde{X}_{\gamma^*}(\tilde{X}_{\gamma^*}^T \tilde{X}_{\gamma^*})^{-1} \tilde{X}_{\gamma^*}^T) \tilde{X}_{\gamma \setminus \gamma^*}$  be orthogonal to the projection of  $\tilde{X}_\gamma$  onto  $\tilde{X}_{\gamma^*}$ , then clearly  $\tilde{X}_{\gamma^*}^T R = 0$  and  $(\tilde{X}_{\gamma^*}, R)$  span the column space of  $\tilde{X}_\gamma$ . Hence  $\hat{\theta}_\gamma^T \tilde{X}_\gamma^T \tilde{X}_\gamma \hat{\theta}_\gamma = \hat{\theta}_{\gamma^*}^T \tilde{X}_{\gamma^*}^T \tilde{X}_{\gamma^*} \hat{\theta}_{\gamma^*} + \hat{\theta}_R^T R^T R \hat{\theta}_R$ , where  $\hat{\theta}_R = (R^T R)^{-1} R^T y$ , giving that  $Z(\gamma^*) - Z(\gamma) = \hat{\theta}_R^T R^T R \hat{\theta}_R$ . By Proposition 4,  $\sqrt{n} \hat{\theta}_R \xrightarrow{D} N(0, \vartheta_\gamma^* V)$  for a fixed positive-definite matrix  $V$ .

To conclude, our Conditions A3-A4 guarantee  $\frac{1}{n} R^T R \xrightarrow{P} \Sigma_R$  for some fixed  $\Sigma_R$  and by the continuous mapping theorem  $\sqrt{n} \Sigma_R^{\frac{1}{2}} \hat{\theta}_R \xrightarrow{D} N(0, \vartheta_\gamma^* \Sigma_R^{\frac{1}{2}} V \Sigma_R^{\frac{1}{2}})$ . Hence  $\frac{n}{\vartheta_\gamma^*} \hat{\theta}_R^T \Sigma_R \hat{\theta}_R \xrightarrow{D} Q$ , where  $Q = O_p(1)$  is a sum of re-scaled central chi-square random variables with 1 degree of freedom. By Slutsky's theorem  $\frac{Z(\gamma^*) - Z(\gamma)}{\vartheta_\gamma^*} = \frac{n}{\vartheta_\gamma^*} \hat{\theta}_R^T (\frac{1}{n} R^T R) \hat{\theta}_R \xrightarrow{D} Q$ , as we wished to prove.

## 2.7 Proof of Corollary 2

The proof runs analogous to Rossell and Telesca (2017), Proposition 3(ii). Briefly, the BMA estimate is  $E(\theta_i | y) =$

$$E(\theta_i | \gamma^*, y)p(\gamma^* | y) + \sum_{\gamma^* \subset \gamma} E(\theta_i | \gamma, y)p(\gamma | y) + \sum_{\gamma^* \not\subset \gamma} E(\theta_i | \gamma, y)p(\gamma | y). \quad (21)$$

Suppose that  $\theta_i^* \neq 0$ . From Proposition 4, the difference between the MLE under  $\gamma$  and  $\theta_i^*$  is  $O_p(1/\sqrt{n})$ , and it can be shown that the difference between a Laplace approximation to  $E(\theta_i | \gamma, y)$  and the MLE is  $O_p(1/\sqrt{n})$  hence  $E(\theta_i | \gamma, y) - \theta_i^* = O_p(1/\sqrt{n})$ . Since  $p(\gamma^* | y) \xrightarrow{P} 1$  by Proposition 5, we have that  $E(\theta_i | \gamma^*, y)p(\gamma^* | y) = \theta_i^* + O_p(1/\sqrt{n})$ . If  $\theta_i^* = 0$  then by definition  $E(\theta_i | \gamma^*, y)p(\gamma^* | y) = 0$ .

Consider the second term in (21) where  $\gamma^* \subset \gamma$ ,

$$p(\gamma | y) \leq 1/(1 + B_{\gamma^*, \gamma}p(\gamma^*)/p(\gamma)) < B_{\gamma, \gamma^*}p(\gamma)/p(\gamma^*) = O_p(b_n^{(k)})p(\gamma)/p(\gamma^*) \leq O_p(b_n^{(k)})r^+,$$

where  $B_{\gamma^*, \gamma}$  is the Bayes factor between  $\gamma^*$  and  $\gamma$ . From Proposition 5, we have that  $b_n^{(k)} = n^{-(p_\gamma - p_{\gamma^*})/2}$  for a local prior,  $b_n^{(k)} = n^{-3(p_\gamma - p_{\gamma^*})/2}$  for the pMOM prior, and  $b_n^{(k)} = e^{-c\sqrt{n}}$ , for some  $c > 0$ , for the peMOM and piMOM priors. Also,  $E(\theta_i | \gamma, y) = \theta_i^* + O_p(1/\sqrt{n})$ . Therefore, if  $\theta_i^* \neq 0$ , we have  $E(\theta_i | \gamma, y)p(\gamma | y) = O_p(b_n^{(k)})r^+$ . If  $\theta_i^* = 0$ , then  $E(\theta_i | \gamma, y)p(\gamma | y) = O_p(b_n^{(k)}/\sqrt{n})p(\gamma)/p(\gamma^*) \leq O_p(b_n^{(k)}/\sqrt{n})r^+$ . The case for  $\gamma^* \not\subset \gamma$  proceeds similarly by noting that by Proposition 5 we have  $B_{\gamma, \gamma^*}r^- = O_p(e^{-cn})r^- = O_p(b_n^{(k)})$  for some  $c > 0$ , since  $e^{-cn}r^- = O(b_n^{(k)})$  by assumption.

Combining the previous results it follows that, if  $\theta_i^* \neq 0$ , then

$$E(\theta_i | y) = \theta_i^* + O_p(1/\sqrt{n}) + O_p(b_n^{(k)})r^+ = \theta_i^* + O_p(1/\sqrt{n}), \quad (22)$$

since  $b_n^{(k)}r^+ = O_p(1/\sqrt{n})$  by the assumption that  $r^+$  does not increase with  $n$ . Conversely if  $\theta_i^* = 0$ , then

$$E(\theta_i | y) = O_p(b_n^{(k)}/\sqrt{n})r^+, \quad (23)$$

giving the desired result.

### 3. APPROXIMATIONS TO THE INTEGRATED LIKELIHOOD

For ease of notation, we drop the subindex  $k$  denoting the set of active variables and let  $\theta = (\theta_1, \dots, \theta_{|k|})$  be their coefficients. Both the Laplace and Importance Sampling approximations require maximizing and evaluating the hessian of  $h_l(\theta, \vartheta, \tilde{\alpha}) = \log L(\theta, \tilde{\vartheta}, \tilde{\alpha}) + \log p(\theta, \tilde{\vartheta}, \tilde{\alpha})$ , where  $L(\cdot)$  and  $p(\cdot)$  are the appropriate likelihood and prior density. Denote by  $g_l(\theta, \tilde{\vartheta}, \tilde{\alpha})$  the gradient of  $h_l(\cdot)$  and by  $H_l(\theta, \tilde{\vartheta}, \tilde{\alpha})$  its hessian, Algorithm 1 finds the posterior mode.

**Algorithm 1. Posterior mode via Newton-Raphson**

1. Initialize  $(\theta^{(0)}, \tilde{\vartheta}^{(0)}, \tilde{\alpha}^{(0)}) = (\hat{\theta}, \log(\hat{\vartheta}), \text{atanh}(\hat{\alpha}))$  where  $(\hat{\theta}, \hat{\vartheta}, \hat{\alpha})$  is the MLE given by Algorithm 1.
  1. Set  $t = 1$  and repeat Steps 2-3 until  $e$  is below some small tolerance (default  $10^{-5}$ ).

2. Update  $(\theta^{(t)}, \tilde{\vartheta}^{(t)}, \tilde{\alpha}^{(t)}) =$

$$(\theta^{(t-1)}, \tilde{\vartheta}^{(t-1)}, \tilde{\alpha}^{(t-1)}) - H_l^{-1}(\theta^{(t-1)}, \tilde{\vartheta}^{(t-1)}, \tilde{\alpha}^{(t-1)})g_l(\theta^{(t-1)}, \tilde{\vartheta}^{(t-1)}, \tilde{\alpha}^{(t-1)}).$$

3. Compute  $e = \|(\theta^{(t)}, \tilde{\vartheta}^{(t)}, \tilde{\alpha}^{(t)}) - (\theta^{(t-1)}, \tilde{\vartheta}^{(t-1)}, \tilde{\alpha}^{(t-1)})\|^\infty$  where  $\|\mathbf{z}\|^\infty$  is the largest element of  $\mathbf{z}$  in absolute value. Set  $t = t + 1$ .

As usual, in the event that  $(\theta^{(t)}, \tilde{\vartheta}^{(t)}, \tilde{\alpha}^{(t)})$  does not increase  $h_l(\cdot)$ , Step 2 can be adjusted by adding a constant  $\lambda$  to the diagonal of  $H_l(\cdot)$ , which for large  $\lambda$  gives the direction of the gradient and is guaranteed to decrease  $h_l(\cdot)$ . However, we observed that this is extremely rare in practice. Usually, the simple Newton step increases  $h_l(\cdot)$  at each iteration and converges to the maximum in a few iterations.

Both  $g_l(\cdot)$  and  $H_l(\cdot)$  are the sum of a term coming from the log-likelihood plus a term coming from the log-prior density. The exact expressions are given below separately.

As an alternative to Algorithm 1, we also provide Algorithm 2 based on Coordinate Descent (*i.e.* successive univariate optimization). Note that the Newton steps to update  $\theta_j$  and  $\alpha$  are in the direction of the gradient and are hence guaranteed to increase the objective function for small enough  $\lambda$ . Step 2 takes advantage of the fact that the maximizer with respect to  $\tilde{\vartheta}$  for fixed  $(\theta, \alpha)$  is available in closed form.

**Algorithm 2. Posterior mode via CDA**

1. Initialize  $\theta^{(0)}$  to the least squares estimate,  $\alpha^{(0)} = 0$ ,  $t = 0$ .

2. For the MOM prior set  $\tilde{\vartheta}^{(t)} = \log(s/(n + p + 3a_\vartheta))$ , where

$$s = \left( b_\vartheta + \theta^{(t)T} \theta^{(t)} + \sum_{i \in A(\theta)} \frac{(y_i - x_i^T \theta^{(t)})^2}{(1 + \alpha^{(t)})^2} + \sum_{i \notin A(\theta)} \frac{(y_i - x_i^T \theta^{(t)})^2}{(1 - \alpha^{(t)})^2} \right).$$

For eMOM and iMOM use a Newton-Raphson step.

3. For  $j = 1, \dots, p$

(a) Set  $\lambda = 1$  and  $\theta^* = \theta_j^{(t-1)} - \lambda g^*/h^*$ , where  $g^*$  and  $h^*$  are the first and second derivatives of  $f(\theta_j) = \log L_1(\theta_1^{(t-1)}, \dots, \theta_{j-1}^{(t-1)}, \theta_j, \theta_{j+1}^{(t)}, \dots, \theta_p^{(t)}, \vartheta^{(t)}, \alpha) + \log p(\theta_j | \vartheta)$  evaluated at  $\theta_j = \theta_j^{(t-1)}$ .

(b) If  $f(\theta^*) > f(\theta_j^{(t-1)})$  set  $\theta_j^{(t)} = \theta^*$ , else set  $\lambda = 0.5\lambda$  and repeat Step 3-(1).

4. Let  $\tilde{\alpha}^* = \tilde{\alpha}^{(t-1)} - \lambda g^*/h^*$ , where  $g^*$  and  $h^*$  are the first and second derivatives of  $f(\tilde{\alpha}) = \log L_1(\theta^{(t)}, \vartheta^{(t)}, \tilde{\alpha}) + \log p(\tilde{\alpha})$  at  $\tilde{\alpha} = \tilde{\alpha}^{(t-1)}$ . If  $f(\tilde{\alpha}^*) > f(\tilde{\alpha}^{(t-1)})$  set  $\alpha^{(t)} = \tilde{\alpha}^*$ , else set  $\lambda = 0.5\lambda$  and repeat Step 4.

5. Compute  $e = \max |(\theta^{(t)}, \tilde{\vartheta}^{(t)}, \tilde{\alpha}^{(t)}) - (\theta^{(t-1)}, \tilde{\vartheta}^{(t-1)}, \tilde{\alpha}^{(t-1)})|$ . If  $e < 10^{-5}$  stop, else set  $t = t + 1$  and go back to Step 1.

3.1 Derivatives of the log-likelihood

Two-piece Normal Under the re-parameterization  $\tilde{\vartheta} = \log(\vartheta)$ ,  $\tilde{\alpha} = \text{atanh}(\alpha)$  the two-piece Normal log-likelihood (3) has gradient

$$\begin{pmatrix} \frac{1}{\exp(\tilde{\vartheta})} X^T W (y - X\theta) \\ -\frac{n}{2} + \frac{1}{2 \exp(\tilde{\vartheta})} (y - X\theta)^T W (y - X\theta) \\ \frac{1}{2 \exp(\tilde{\vartheta})} (y - X\theta)^T W^* (y - X\theta) \end{pmatrix},$$

where as usual  $W = \text{diag}(w)$ ,  $w_i = [1 + \tanh(\tilde{\alpha})]^{-2}$  if  $i \in A(\theta)$  and  $w_i = [1 - \tanh(\tilde{\alpha})]^{-2}$  if  $i \notin A(\theta)$ , and  $W^* = \text{diag}(w^*)$  with  $w_i^* = -\frac{2\text{sech}^2(\tilde{\alpha})}{(\tanh(\tilde{\alpha})+1)^3}$  if  $i \in A(\theta)$  and  $w_i^* = \frac{2\text{sech}^2(\tilde{\alpha})}{(1-\tanh(\tilde{\alpha}))^3}$  if  $i \notin A(\theta)$ . Its Hessian is given by

$$-e^{-\tilde{\vartheta}} \begin{pmatrix} X^T W X & X^T W (y - X\theta) & X^T W^* (y - X\theta) \\ \frac{1}{2}(y - X\theta)^T W (y - X\theta) & -\frac{1}{2}(y - X\theta)^T W^* (y - X\theta) & \\ & \frac{1}{2}(y - X\theta)^T W^{**} (y - X\theta) & \end{pmatrix}, \quad (24)$$

where  $W^{**} = \text{diag}(w^{**})$ , with  $w_i^{**} = 2e^{-4\tilde{\alpha}}(e^{2\tilde{\alpha}} + 2)$  if  $i \in A(\theta)$  and  $w_i^{**} = 2e^{2\tilde{\alpha}} + 4e^{4\tilde{\alpha}}$  if  $i \notin A(\theta)$ .

Two-piece Laplace The asymmetric Laplace  $\log L_2(\theta, \tilde{\vartheta}, \tilde{\alpha})$ , where  $\tilde{\vartheta} = \log(\vartheta)$ ,  $\tilde{\alpha} = \text{atanh}(\alpha)$  has gradient

$$\begin{pmatrix} -e^{-\tilde{\vartheta}/2} X^T \bar{w} \\ -\frac{n}{2} + \frac{1}{2} e^{-\tilde{\vartheta}/2} w^T |y - X\theta| \\ e^{-\tilde{\vartheta}/2} |y - X\theta|^T \bar{w}^* \end{pmatrix},$$

and hessian

$$e^{-\tilde{\vartheta}/2} \times \begin{pmatrix} 0 & \frac{1}{2} X^T \bar{w} & X^T w^* \\ \frac{1}{2} \bar{w}^T X & -\frac{1}{4} w^T |y - X\theta| & -\frac{1}{2} |y - X\theta|^T \bar{w}^* \\ (X^T w^*)^T & -\frac{1}{2} |y - X\theta|^T \bar{w}^* & -2 |y - X\theta|^T w^* \end{pmatrix}, \quad (25)$$

where  $w_i = \bar{w}_i = (1 + \alpha)^{-1}$ ,  $w_i^* = \bar{w}_i^* = e^{-2\alpha}$  if  $i \in A(\theta)$ , and  $w_i = (1 - \alpha)^{-1}$ ,  $\bar{w}_i = -w_i$ ,  $w_i^* = e^{2\alpha}$ ,  $\bar{w}_i^* = -w_i^*$  if  $i \notin A(\theta)$ . Naturally, symmetric Laplace errors are the particular case  $\alpha = 0$  and give  $w_i = w_i^* = 1$ .

Expected two-piece Laplace log-likelihood We derive  $\bar{L}_2 = E(\log L_2(\eta))$ , where  $\eta = (\theta, \vartheta, \alpha)$  and its derivatives under the data-generating model  $y_i = x_i^T \theta_0 + \epsilon_i$ , for some  $\theta_0 \in \mathbb{R}^p$  where  $\epsilon_i$  are independent across  $i = 1, \dots, n$  and arise from an arbitrary probability density function  $s_0(y_i | x_i)$ . After some algebra and noting that  $\epsilon_i = y_i - x_i^T \theta_0$  gives



$$\begin{aligned}\bar{L}_2 &= \int \log L_2(\eta) s_0(\epsilon|x) d\epsilon = -n \log(2) - \frac{n}{2} \log(\vartheta) - \sum_{i=1}^n \frac{1}{\sqrt{\vartheta}(1+\alpha)} \int_{-\infty}^{x_i^T(\theta-\theta_0)} S_0(\epsilon_i) d\epsilon_i \\ &\quad - \sum_{i=1}^n \frac{1}{\sqrt{\vartheta}(1-\alpha)} \int_{x_i^T(\theta-\theta_0)}^{\infty} (1 - S_0(\epsilon_i)) d\epsilon_i,\end{aligned}$$

where  $S_0(\epsilon_i) = S_0(\epsilon_i|0)$  is the cumulative probability function associated to  $s_0(\epsilon_i) = s_0(\epsilon_i|0)$ , where 0 indicates a zero covariate vector. Then taking derivatives we obtain

$$\begin{aligned}\frac{\partial}{\partial \theta} \bar{L}_2 &= \sum_{i=1}^n \left[ -\frac{x_i S_0(x_i^T(\theta - \theta_0))}{\sqrt{\vartheta}(1+\alpha)} + \frac{x_i [1 - S_0(x_i^T(\theta - \theta_0))]}{\sqrt{\vartheta}(1-\alpha)} \right], \\ \frac{\partial}{\partial \vartheta} \bar{L}_2 &= \sum_{i=1}^n \left[ -\frac{1}{2\vartheta} + \frac{I_{i1}}{2\vartheta^{3/2}(1+\alpha)} + \frac{I_{i2}}{2\vartheta^{3/2}(1-\alpha)} \right], \\ \frac{\partial}{\partial \alpha} \bar{L}_2 &= \sum_{i=1}^n \left[ \frac{I_{i1}}{\sqrt{\vartheta}(1+\alpha)^2} - \frac{I_{i2}}{\sqrt{\vartheta}(1-\alpha)^2} \right],\end{aligned}$$

where  $I_{i1} = \int_{-\infty}^{x_i^T(\theta-\theta_0)} S_0(\epsilon_i) d\epsilon_i$ ,  $I_{i2} = \int_{x_i^T(\theta-\theta_0)}^{\infty} (1 - S_0(\epsilon_i)) d\epsilon_i$ . The second derivatives are

$$\begin{aligned}\frac{\partial^2}{\partial \theta^2} \bar{L}_2 &= -\sum_{i=1}^n \frac{2x_i x_i^T s_0(x_i^T(\theta - \theta_0))}{\sqrt{\vartheta}(1-\alpha^2)}, \\ \frac{\partial^2}{\partial \vartheta^2} \bar{L}_2 &= \sum_{i=1}^n \left[ \frac{1}{2\vartheta^2} - \frac{3I_{i1}}{4\vartheta^{5/2}(1+\alpha)} - \frac{3I_{i2}}{4\vartheta^{5/2}(1-\alpha)} \right], \\ \frac{\partial^2}{\partial \alpha^2} \bar{L}_2 &= -\sum_{i=1}^n \left[ \frac{2I_{i1}}{\sqrt{\vartheta}(1+\alpha)^3} - \frac{2I_{i2}}{\sqrt{\vartheta}(1-\alpha)^3} \right], \\ \frac{\partial^2}{\partial \vartheta \partial \theta} \bar{L}_2 &= \sum_{i=1}^n \left[ \frac{x_i S_0(x_i^T(\theta - \theta_0))}{2\vartheta^{3/2}(1+\alpha)} - \frac{x_i [1 - S_0(x_i^T(\theta - \theta_0))]}{2\vartheta^{3/2}(1-\alpha)} \right], \\ \frac{\partial^2}{\partial \alpha \partial \theta} \bar{L}_2 &= \sum_{i=1}^n \left[ \frac{x_i S_0(x_i^T(\theta - \theta_0))}{\sqrt{\vartheta}(1+\alpha)^2} + \frac{x_i [1 - S_0(x_i^T(\theta - \theta_0))]}{\sqrt{\vartheta}(1-\alpha)^2} \right], \\ \frac{\partial^2}{\partial \vartheta \partial \alpha} \bar{L}_2 &= -\sum_{i=1}^n \left[ \frac{I_{i1}}{2\vartheta^{3/2}(1+\alpha)^2} + \frac{I_{i2}}{2\vartheta^{3/2}(1-\alpha)^2} \right].\end{aligned}$$

Simple inspection reveals that  $(\partial/\partial \theta) \bar{L}_2 = 0$  implies  $(\partial^2/\partial \theta \partial \vartheta) \bar{L}_2 = 0$ , and likewise  $(\partial/\partial \alpha) \bar{L}_2 = 0$  implies  $(\partial^2/\partial \theta \partial \alpha) \bar{L}_2 = 0$ . Since the maximum likelihood estimator  $(\hat{\theta}, \hat{\vartheta}, \hat{\alpha})$  converges in prob-

ability to the maximizer of  $\bar{L}_2$ , these second derivatives evaluated at  $(\hat{\theta}, \hat{\vartheta}, \hat{\alpha})$  also converge in probability to 0.

We wish to find an asymptotic expression for the remaining second derivatives evaluated at  $(\hat{\theta}, \hat{\vartheta}, \hat{\alpha})$  when the data-generating truth is  $\epsilon_i \sim \text{AL}(x_i^T \theta_0, \vartheta_0, \alpha_0)$  for some  $(\theta_0, \vartheta_0, \alpha_0)$ . Given that  $(\hat{\theta}, \hat{\vartheta}, \hat{\alpha}) \xrightarrow{P} (\theta_0, \vartheta_0, \alpha_0)$ , the expressions above require evaluating the density of an asymmetric Laplace  $s_0(0) = 1/(2\sqrt{\vartheta_0})$  and its cumulative probability function  $S_0(0) = (1 + \alpha_0)/2$ . Similarly, direct integration gives  $I_{i1} = \sqrt{\vartheta_0}(1 + \alpha_0)^2/2$  and  $I_{i2} = \sqrt{\vartheta_0}(1 - \alpha_0)^2/2$ .

$$\begin{aligned}
\frac{\partial^2}{\partial \theta^2} \bar{L}_2 &\xrightarrow{P} -X^T X \frac{1}{\vartheta_0(1 - \alpha_0^2)}, \\
\frac{\partial^2}{\partial \vartheta^2} \bar{L}_2 &\xrightarrow{P} \frac{n}{2\vartheta_0^2} - \frac{3n(1 + \alpha_0)}{8\vartheta_0^2} - \frac{3(1 - \alpha_0)}{8\vartheta_0^2} = -\frac{n}{4\vartheta_0^2}, \\
\frac{\partial^2}{\partial \alpha^2} \bar{L}_2 &\xrightarrow{P} -\frac{n}{1 + \alpha_0} - \frac{n}{1 - \alpha_0} = -\frac{2n}{1 - \alpha_0^2}, \\
\frac{\partial^2}{\partial \alpha \partial \theta} \bar{L}_2 &\xrightarrow{P} \frac{n\bar{x}}{\sqrt{\vartheta_0}} \left( \frac{1}{2(1 + \alpha_0)} + \frac{1}{2(1 - \alpha_0)} \right) = \frac{n\bar{x}}{\sqrt{\vartheta_0}(1 - \alpha_0^2)}. \tag{26}
\end{aligned}$$

### 3.2 Derivatives of the log-prior density

The log-prior density is  $\log p(\theta, \tilde{\vartheta}) = \log p(\theta \mid \tilde{\vartheta}) + \log p(\tilde{\vartheta})$  when  $\tilde{\alpha} = 0$  under the assumed model and  $\log p(\theta, \tilde{\vartheta}, \tilde{\alpha}) = \log p(\theta, \tilde{\vartheta}) + \log p(\tilde{\alpha})$  when  $\tilde{\alpha} \neq 0$ , where  $p(\theta \mid \tilde{\vartheta})$  and  $p(\tilde{\alpha})$  are the pMOM, piMOM or peMOM priors and  $p(\tilde{\vartheta}) = \text{IG}(e^{\tilde{\vartheta}}; a_{\tilde{\vartheta}}/2, b_{\tilde{\vartheta}}/2)e^{\tilde{\vartheta}}$ . For ease of notation let  $\theta^{-a}$  be the vector with elements  $\theta_i^{-a}$  for  $i = 1, \dots, |k|$ .

pMOM prior Straightforward algebra gives

$$\nabla \log p_M(\theta, \tilde{\vartheta}, \tilde{\alpha}) = \begin{pmatrix} 2\theta^{-1} - \theta e^{-\tilde{\vartheta}}/g_\theta \\ -\frac{3|k|+a_\vartheta}{2} + (\theta^T \theta/g_\theta + b_\vartheta)e^{-\tilde{\vartheta}}/2 \\ 2\tilde{\alpha}^{-1} - \tilde{\alpha}g_\alpha^{-1} \end{pmatrix},$$

$$\nabla^2 \log p_M(\theta, \tilde{\vartheta}, \tilde{\alpha}) = \begin{pmatrix} \text{diag}(-2\theta^{-2} - e^{-\tilde{\vartheta}}/g_\theta) & \theta e^{-\tilde{\vartheta}}/g_\theta & 0 \\ \theta^T e^{-\tilde{\vartheta}}/g_\theta & -e^{-\tilde{\vartheta}}(\theta^T \theta/g_\theta + b_\vartheta)/2 & 0 \\ 0 & 0 & -2\tilde{\alpha}^{-2} - g_\alpha^{-1} \end{pmatrix},$$

piMOM prior We obtain

$$\nabla \log p_I(\theta, \tilde{\vartheta}, \tilde{\alpha}) = \begin{pmatrix} -2\theta^{-1} + 2g_\theta e^{\tilde{\vartheta}} \theta^{-3} \\ (|k| - a_\vartheta)/2 + b_\vartheta e^{-\tilde{\vartheta}}/2 - g_\theta e^{\tilde{\vartheta}} \sum_i \theta_i^{-2} \\ -2\tilde{\alpha}^{-1} - 2g_\alpha \tilde{\alpha}^{-3} \end{pmatrix},$$

$$\nabla^2 \log p_I(\theta, \tilde{\vartheta}, \tilde{\alpha}) = \begin{pmatrix} \text{diag}(2\theta^{-2} - 6g_\theta e^{\tilde{\vartheta}} \theta^{-4}) & 2g_\theta e^{\tilde{\vartheta}} \theta^{-3} & 0 \\ (-2g_\theta e^{\tilde{\vartheta}} \theta^{-3})^T & -b_\vartheta e^{-\tilde{\vartheta}}/2 - e^{\tilde{\vartheta}} g_\theta \sum_i \theta_i^{-2} & 0 \\ 0 & 0 & 2\tilde{\alpha}^{-2} + 6g_\alpha \tilde{\alpha}^{-4} \end{pmatrix}.$$

peMOM prior We obtain

$$\nabla \log p_E(\theta, \tilde{\vartheta}, \tilde{\alpha}) = \begin{pmatrix} 2g_\theta e^{\tilde{\vartheta}} \theta^{-3} - \theta e^{-\tilde{\vartheta}} g_\theta^{-1} \\ -(|k| + a_\vartheta)/2 + (b_\vartheta + \theta^T \theta/g_\theta)e^{-\tilde{\vartheta}}/2 - g_\theta e^{\tilde{\vartheta}} \sum_i \theta_i^{-2} \\ 2g_\alpha \tilde{\alpha}^{-3} - \tilde{\alpha}g_\alpha^{-1} \end{pmatrix},$$

and  $\nabla^2 \log p_E(\theta, \tilde{\vartheta}, \tilde{\alpha}) =$

$$\begin{pmatrix} \text{diag}(-6g_\theta e^{\tilde{\vartheta}} \theta^{-4} - e^{-\tilde{\vartheta}} g_\theta^{-1}) & 2g_\theta e^{\tilde{\vartheta}} \theta^{-3} + \theta e^{-\tilde{\vartheta}} g_\theta^{-1} & 0 \\ (2g_\theta e^{\tilde{\vartheta}} \theta^{-3} + \theta e^{-\tilde{\vartheta}} g_\theta^{-1})^T & -(b_\vartheta + \theta^T \theta/g_\theta)e^{-\tilde{\vartheta}}/2 - e^{\tilde{\vartheta}} g_\theta \sum_i \theta_i^{-2} & 0 \\ 0 & 0 & -6g_\alpha \tilde{\alpha}^{-4} - g_\alpha^{-1} \end{pmatrix}.$$

### 3.3 Quadratic approximation to asymmetric Laplace log-likelihood

The goal is to approximate the curvature of the one-dimensional function  $f(\lambda) = \log L_2(\theta_\lambda, \hat{\vartheta}, \hat{\alpha})$  around  $\lambda = 0$ , where  $\theta_\lambda = (\hat{\theta}_1, \dots, \hat{\theta}_{j-1}, \hat{\theta}_j + \lambda, \hat{\theta}_{j+1}, \dots, \hat{\theta}_p)$  is fixed to the maximum likelihood estimator except for the  $j^{\text{th}}$  regression parameter, which is a function of  $\lambda \in \mathbb{R}$ . Given that  $f(0)$  is known and that its derivative at  $\lambda = 0$  is 0 ( $\hat{\theta}$  is a maximum) we seek  $h_j^* < 0$  such that  $f(\lambda) - f(0) \approx 0.5h_j^*\lambda^2$ . Our strategy is to evaluate  $f(\lambda_k)$  on a grid  $\lambda_k$  for  $k = 1, \dots, K$  and use the least-squares estimate  $h_j^* = 2 \sum_{k=1}^K \lambda_k^2 (f(\lambda_k) - f(0)) / \sum_{k=1}^K \lambda_k^4$ , where the form of  $\log L_2$  gives the simple expression

$$f(\lambda_k) - f(0) = -\frac{1}{\sqrt{\hat{\vartheta}}} \sum_{i=1}^n |r_i - \lambda_k x_{ij}| \left( \frac{\mathbb{I}(r_i \leq \lambda_k x_{ij})}{1 + \hat{\alpha}} + \frac{\mathbb{I}(r_i > \lambda_k x_{ij})}{1 - \hat{\alpha}} \right),$$

and  $r_i = y_i - x_i^T \hat{\theta}$ . Once  $h_1^*, \dots, h_p^*$  have been obtained we let  $D = \text{diag}(h_1^*/\bar{h}_{11}, \dots, h_p^*/\bar{h}_{pp})$  where  $\bar{H} = (X^T X) / (\hat{\vartheta}(1 - \hat{\alpha}^2))$  is the asymptotic hessian under asymmetric Laplace errors, and we approximate the hessian of  $\log L_2(\theta, \hat{\vartheta}, \hat{\alpha})$  around  $\theta = \hat{\theta}$  with  $H^* = D^{\frac{1}{2}} \bar{H} D^{\frac{1}{2}}$ . The construction ensures that the diagonal elements in  $H^*$  are  $h_1^*, \dots, h_p^*$ , i.e. the quadratic approximation matches the actual curvature of  $\log L_2$  along each canonical axis. From Section 4 the correlation structure borrowed from  $\bar{H}$  remains asymptotically valid as long as the residuals are independent and identically distributed, however in our experience the approximation usually suffices for practical purposes even when these assumptions is violated.

The problem has been thus reduced to choosing the grid  $\lambda_1, \dots, \lambda_K$ . One naive option is to take the  $n$  points of non-differentiability  $\lambda = r_i/x_{ij}$ , however, by the nature of least squares, this strategy tends to approximate better  $f(\lambda)$  for large  $\lambda^2$  and we are interested in local approximations around  $\lambda = 0$ , further evaluating  $f(\lambda)$  at  $n$  points requires  $O(n^2)$  operations for each  $j = 1, \dots, p$  and is thus computationally costly. Instead we evaluate  $f(\lambda)$  only at the  $K = 2$  points given by the endpoints of the asymptotic 95% confidence interval  $\lambda = \{-1.96\bar{v}_j, 1.96\bar{v}_j\}$  where  $\bar{v}_j$  is the  $j^{\text{th}}$  diagonal element in  $\bar{H}^{-1}$ . This simple strategy ensures that the approximation holds locally around  $\lambda = 0$  in the sense of having non-negligible likelihood, requires only  $O(n)$  operations and we have observed to deliver reasonably accurate approximations in practice. Our approximation is similar in spirit to the rank-based score test inversion used to obtain confidence intervals in quantile

regression, which has been amply described to deliver fairly precise intervals, with the important difference that rank inversion requires an ordering of observations that scales poorly with  $p$  and  $n$ .

Supplementary Figure 1S shows an example with the likelihood  $L_2$  (scaled to  $(0, 1)$ ) and the two quadratic approximations based on the asymptotic covariance and its least-squares adjustment for an intercept-only model ( $p = 1$ ) and  $n = 200$ . When residuals were truly generated from an asymmetric Laplace (left panel) the two quadratic approximations were essentially identical, however under truly normally distributed residuals the asymptotic covariance over-estimated the curvature.

[Figure 1 about here.]

#### 4. SUPPLEMENTARY RESULTS

[Table 1 about here.]

[Table 2 about here.]

##### 4.1 Simulation study with identically distributed errors

[Table 3 about here.]

[Table 4 about here.]

[Figure 2 about here.]

[Figure 3 about here.]

[Table 5 about here.]

[Table 6 about here.]

[Figure 4 about here.]

[Figure 5 about here.]

[Figure 6 about here.]

[Figure 7 about here.]

We assessed the sensitivity of the results of the  $p = 6$  simulation study in Section 6.1 of the main paper to the prior on the asymmetry coefficient by setting  $g_\alpha$  such that  $P(|\alpha| > 0.1) = 0.99$ . Supplementary Table 3S summarizes the inference on the error distribution and Supplementary Figure 2S the marginal variable inclusion probabilities. The latter were virtually identical to those in Figure 2 obtained under  $g_\alpha$  such that  $P(|\alpha| > 0.2) = 0.99$ , showing that variable inclusion is robust to moderate changes in  $g_\alpha$ .

We also assessed the accuracy of the Laplace approximations to the integrated likelihood  $p(y | \gamma)$  by comparing the results with those obtained with the importance sampling estimates with  $B = 10,000$  draws described in Section 5 of the main paper. Supplementary Figure 3S displays the results for  $g_\alpha = 0.357$ . These are extremely similar to those based on Laplace approximation in Figure 2.

Supplementary Figure 6S shows analogous results for  $p = 100$ , with  $g_\alpha = 0.357$  and  $p(y | \gamma)$  estimated via Laplace approximations.

#### 4.2 Simulation study with non-identically distributed errors

[Table 7 about here.]

[Table 8 about here.]

[Figure 8 about here.]

[Table 9 about here.]

Supplementary Table 7S shows the mean average posterior probability assigned to the Normal, asymmetric Normal, Laplace and asymmetric Laplace models under the heteroskedastic simulation (Section 6.2, main manuscript).

Supplementary Figure 8S shows marginal variable inclusion probabilities under the hetero-asymmetric simulation.

Supplementary Table 9S reports true and false positives for our simulation study mimicking Grünwald and van Ommen (2014) described in Section 6.2 of the main manuscript.

### 4.3 DLD data

[Table 10 about here.]

[Table 11 about here.]

Supplementary Table 10S shows the six genes with largest marginal inclusion probabilities  $p(\gamma_j = 1 \mid y)$  when conditioning on Normal errors and when inferring the error distribution. The figures were similar for the four top genes, but the Normal model assigned somewhat higher probability to FBXL19 substantially lower probability to MTMR1.

### REFERENCES

- Arellano-Valle, R., Gómez, H., and Quintana, F. (2005), “Statistical inference for a general class of asymmetric distributions,” *Journal of Statistical Planning and Inference*, 128(2), 427–443.
- Grünwald, P., and van Ommen, T. (2014), “Inconsistency of Bayesian inference for misspecified linear models, and a proposal for repairing it,” *arXiv*, 1412.3730, 1–70.
- Johnson, V., and Rossell, D. (2010), “Prior Densities for Default Bayesian Hypothesis Tests,” *Journal of the Royal Statistical Society B*, 72(2), 143–170.
- Johnson, V., and Rossell, D. (2012), “Bayesian model selection in high-dimensional settings,” *Journal of the American Statistical Association*, 24(498), 649–660.
- Koenker, R., and Bassett, G. (1982), “Tests of linear hypotheses and L1 estimation,” *Econometrica*, 50(6), 1577–1584.
- Newey, W. K., and Powell, J. L. (1987), “Asymmetric Least Squares Estimation and Testing,” *Econometrica*, 55(4), 819–847.
- Pollard, D. (1991), “Asymptotics for least absolute deviation regression estimators,” *Econometric Theory*, 7(2), 186–199.
- Rockafellar, R. (2015), *Convex analysis*, Princeton: Princeton university press.
- Rossell, D., and Telesca, D. (2017), “Non-local priors for high-dimensional estimation,” *Journal of the American Statistical Association*, 112, 254–265.

van der Vaart, A. (1998), *Asymptotic statistics*, New York: Cambridge University Press.



List of Figures

1S Quadratic approximation to  $L_2$  (solid grey) with  $p = 1, n = 200$  from asymptotic covariance (dotted black) and least-squares adjustment (solid black). Left:  $\epsilon_i \sim \text{AL}(0, 2, -0.5)$ ; Right:  $\epsilon_i \sim N(0, 2)$ . . . . . 34

2S Sensitivity analysis with  $g_\alpha = 0.087$ .  $P(\theta_i \neq 0 | y)$  for  $p = 5, \vartheta = 2, \theta = (0.5, 1, 1.5, 0, 0)$ ,  $n = 100, \rho_{ij} = 0.5$ . Black circles show the mean. . . . . 35

3S Monte Carlo estimates ( $B = 10,000$ ) under  $g_\alpha = 0.357$ .  $P(\theta_i \neq 0 | y)$  for  $p = 5, \vartheta = 2, \theta = (0.5, 1, 1.5, 0, 0)$ ,  $n = 100, \rho_{ij} = 0.5$ . Black circles show the mean. . . . . 36

4S  $P(\theta_i \neq 0 | y)$  for  $p = 100, \vartheta = 1, \theta = (0, 0.5, 1, 1.5, 0, \dots, 0)$ ,  $n = 100, \rho_{ij} = 0.5$ . Black circles show the mean. . . . . 37

5S  $P(\theta_i \neq 0 | y)$  for  $p = 500, \vartheta = 1, \theta = (0, 0.5, 1, 1.5, 0, \dots, 0)$ ,  $n = 100, \rho_{ij} = 0.5$ . Black circles show the mean. . . . . 38

6S  $P(\theta_i \neq 0 | y)$  for  $p = 100, \vartheta = 2, \theta = (0, 0.5, 1, 1.5, 0, \dots, 0)$ ,  $n = 100, \rho_{ij} = 0.5$ . Black circles show the mean. . . . . 39

7S  $P(\theta_i \neq 0 | y)$  for  $p = 500, \vartheta = 2, \theta = (0, 0.5, 1, 1.5, 0, \dots, 0)$ ,  $n = 100, \rho_{ij} = 0.5$ . Black circles show the mean. . . . . 40

8S  $P(\theta_i \neq 0 | y)$  for simulation with constant  $\vartheta = 0$  and varying  $\tanh(\alpha_i) \sim N(\text{atanh}(\bar{\alpha}), 1/4^2)$ , where  $\bar{\alpha} = 0$  for Normal and Laplace and  $\bar{\alpha} = -0.5$  for ANormal and ALaplace.  $P(\theta_i \neq 0 | y)$  for  $p = 6, \theta = (0, 0.5, 1, 1.5, 0, 0)$ ,  $n = 100, \rho_{ij} = 0.5$ . Black circles show the mean. . . . . 41

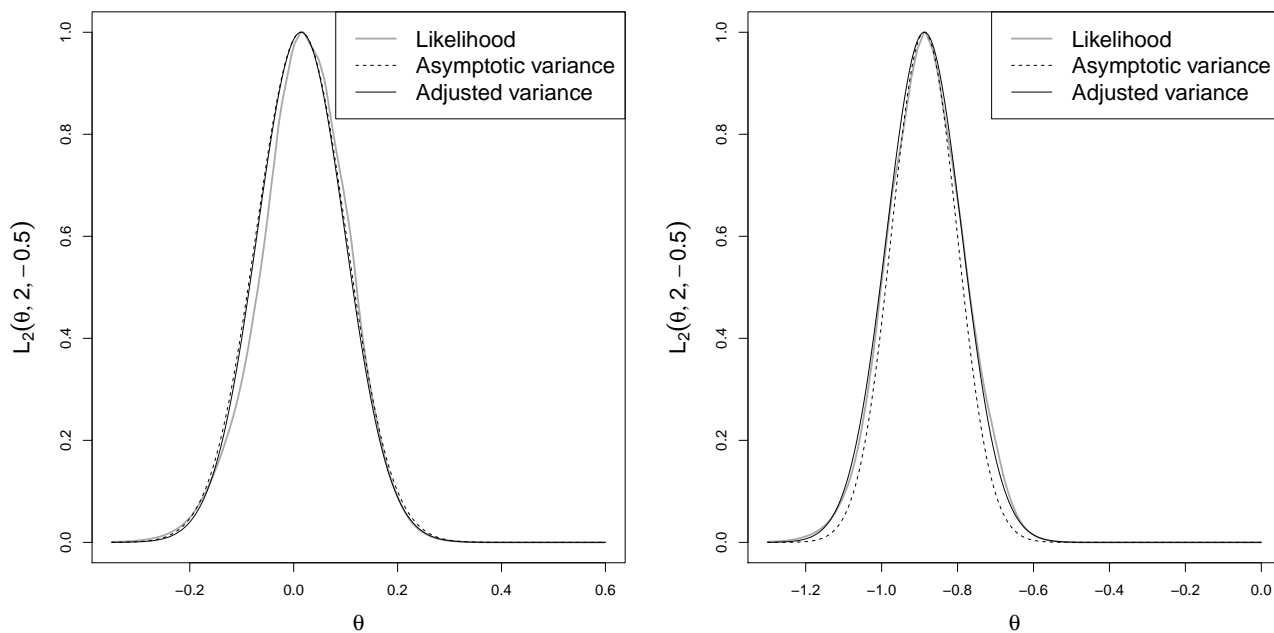


Figure 1S: Quadratic approximation to  $L_2$  (solid grey) with  $p = 1, n = 200$  from asymptotic covariance (dotted black) and least-squares adjustment (solid black). Left:  $\epsilon_i \sim \text{AL}(0, 2, -0.5)$ ; Right:  $\epsilon_i \sim N(0, 2)$ .

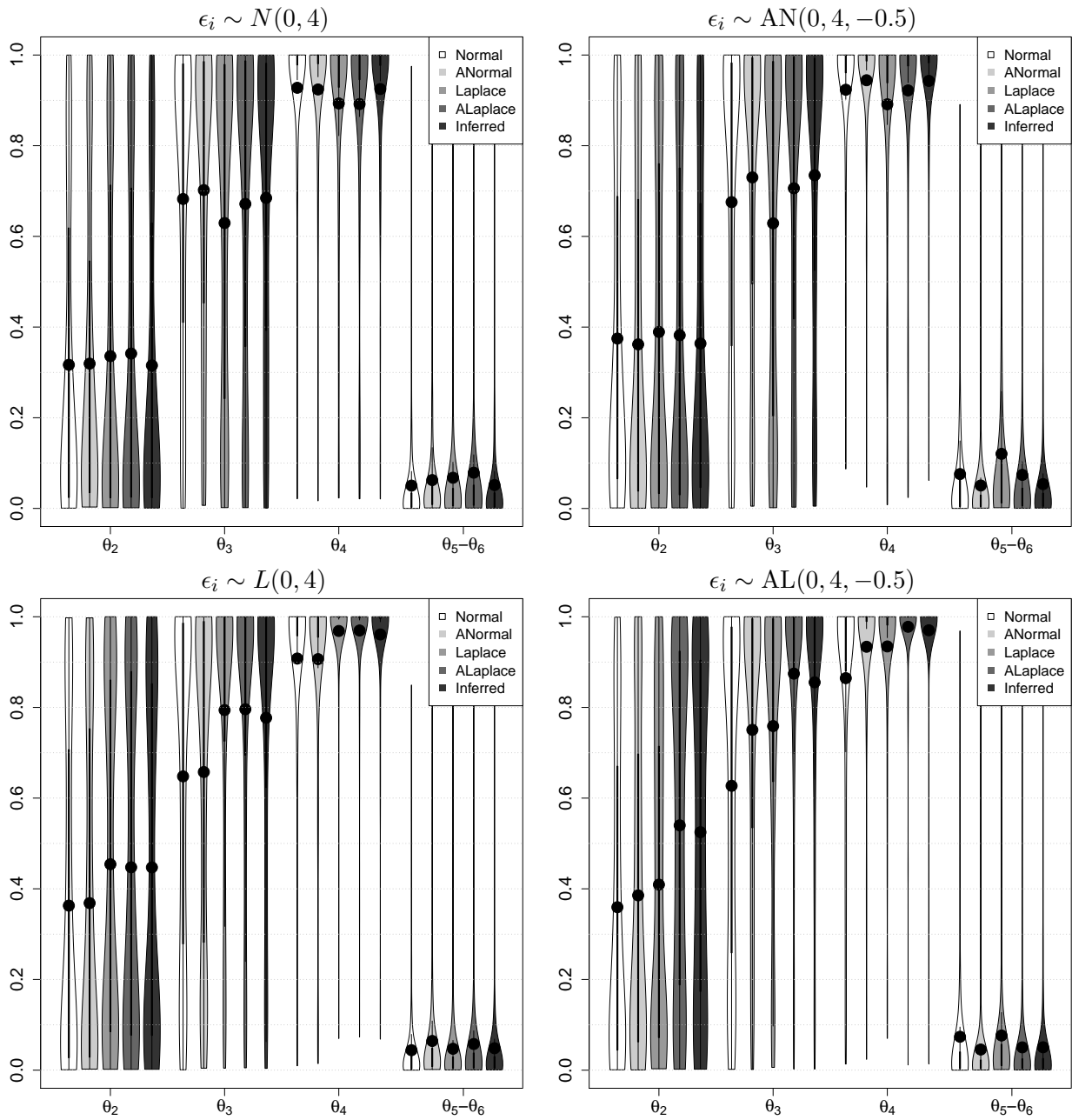


Figure 2S: Sensitivity analysis with  $g_\alpha = 0.087$ .  $P(\theta_i \neq 0 | y)$  for  $p = 5$ ,  $\vartheta = 2$ ,  $\theta = (0.5, 1, 1.5, 0, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Black circles show the mean.

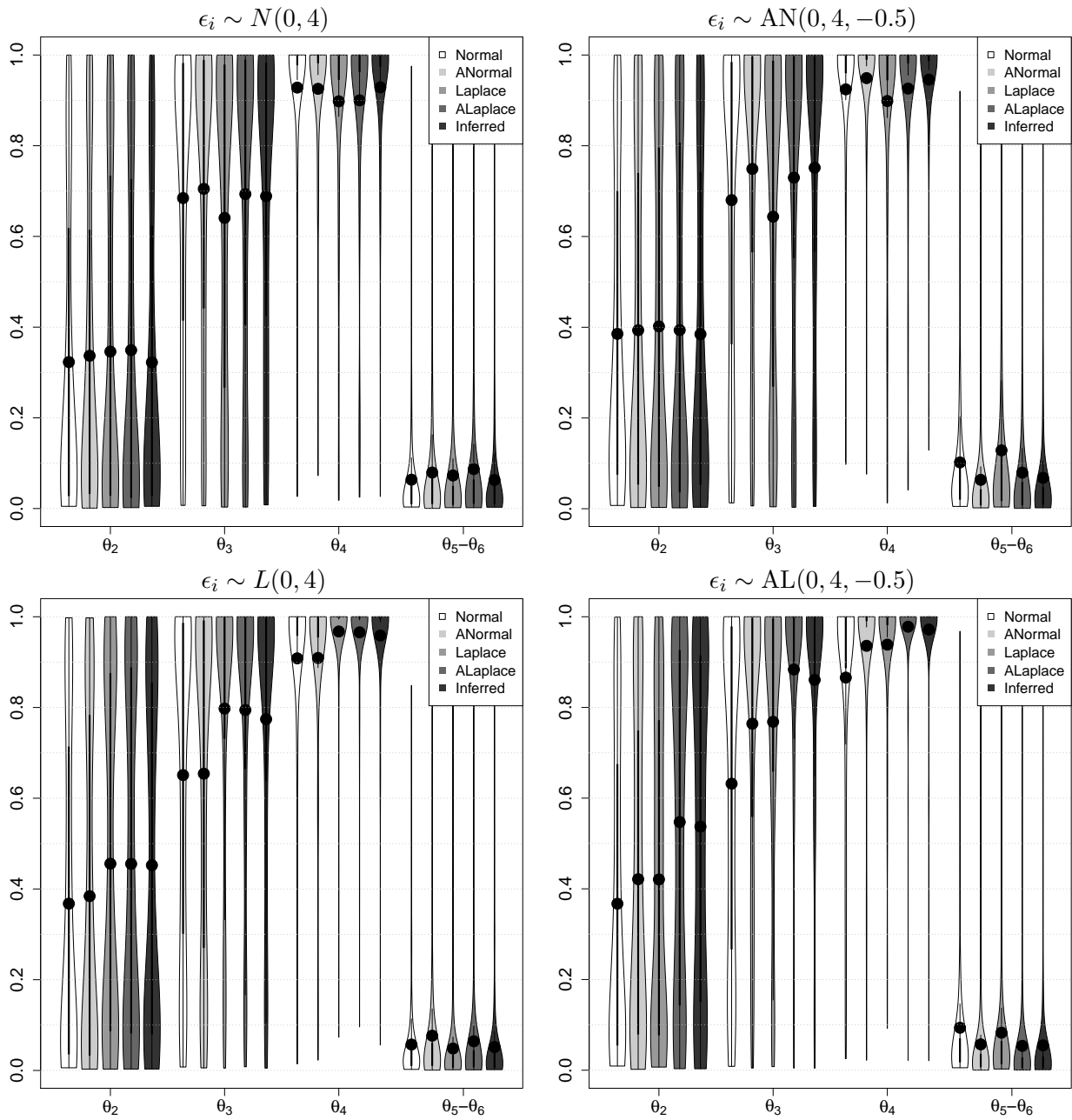


Figure 3S: Monte Carlo estimates ( $B = 10,000$ ) under  $g_\alpha = 0.357$ .  $P(\theta_i \neq 0 | y)$  for  $p = 5$ ,  $\vartheta = 2$ ,  $\theta = (0.5, 1, 1.5, 0, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Black circles show the mean.

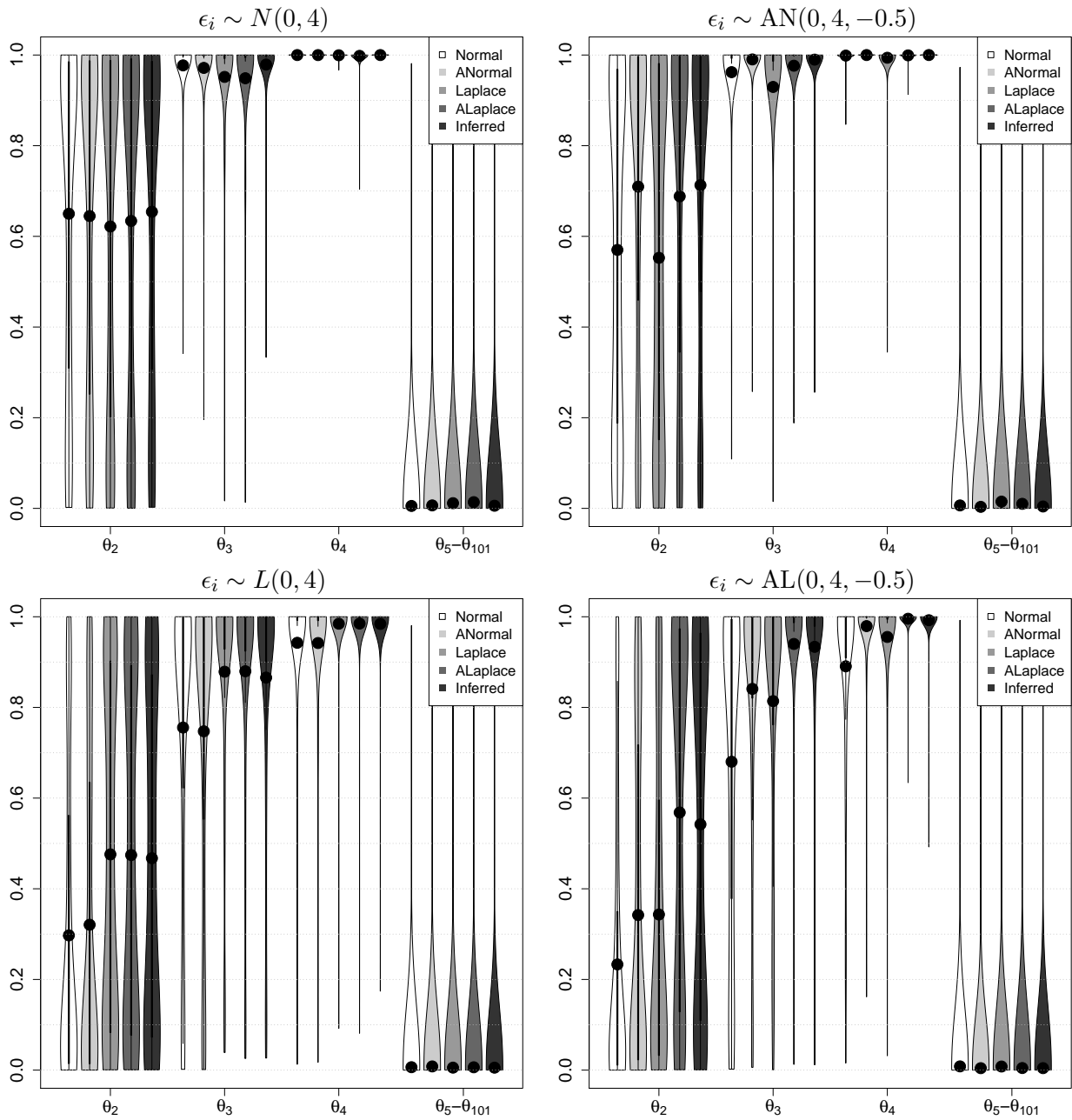


Figure 4S:  $P(\theta_i \neq 0 | y)$  for  $p = 100$ ,  $\vartheta = 1$ ,  $\theta = (0, 0.5, 1, 1.5, 0, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Black circles show the mean.

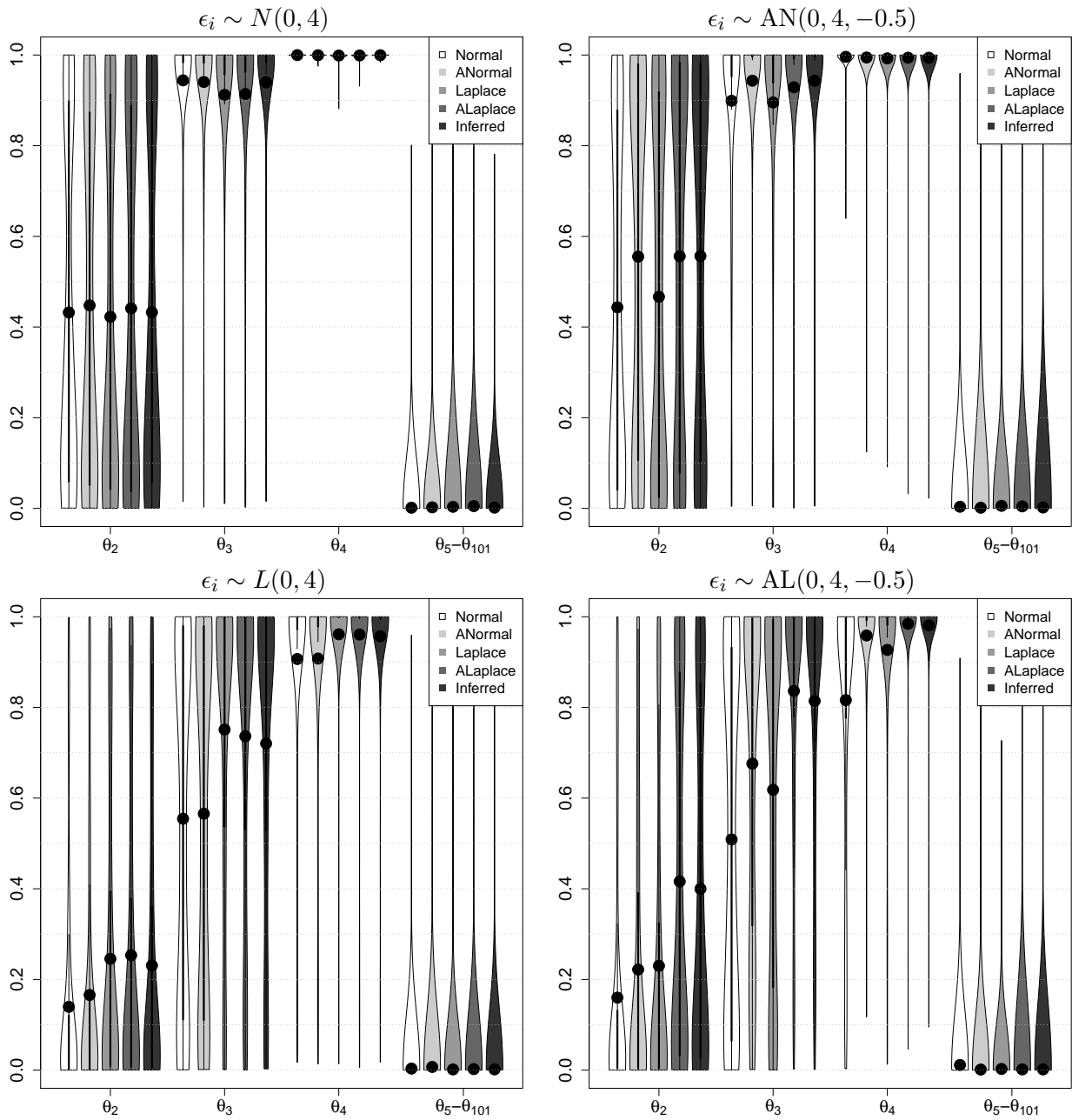


Figure 5S:  $P(\theta_i \neq 0 | y)$  for  $p = 500$ ,  $\vartheta = 1$ ,  $\theta = (0, 0.5, 1, 1.5, 0, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Black circles show the mean.

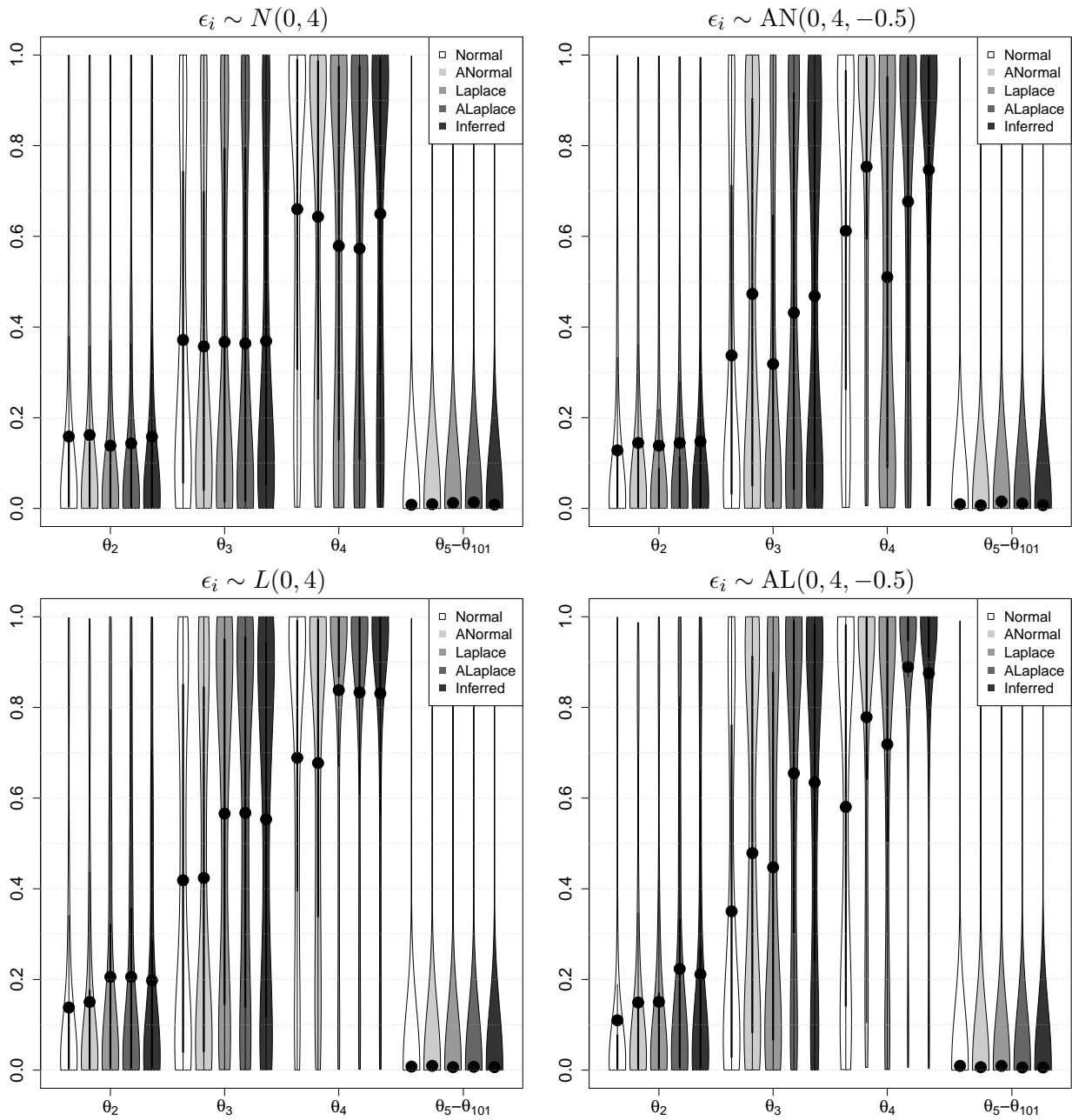


Figure 6S:  $P(\theta_i \neq 0 | y)$  for  $p = 100$ ,  $\vartheta = 2$ ,  $\theta = (0, 0.5, 1, 1.5, 0, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Black circles show the mean.

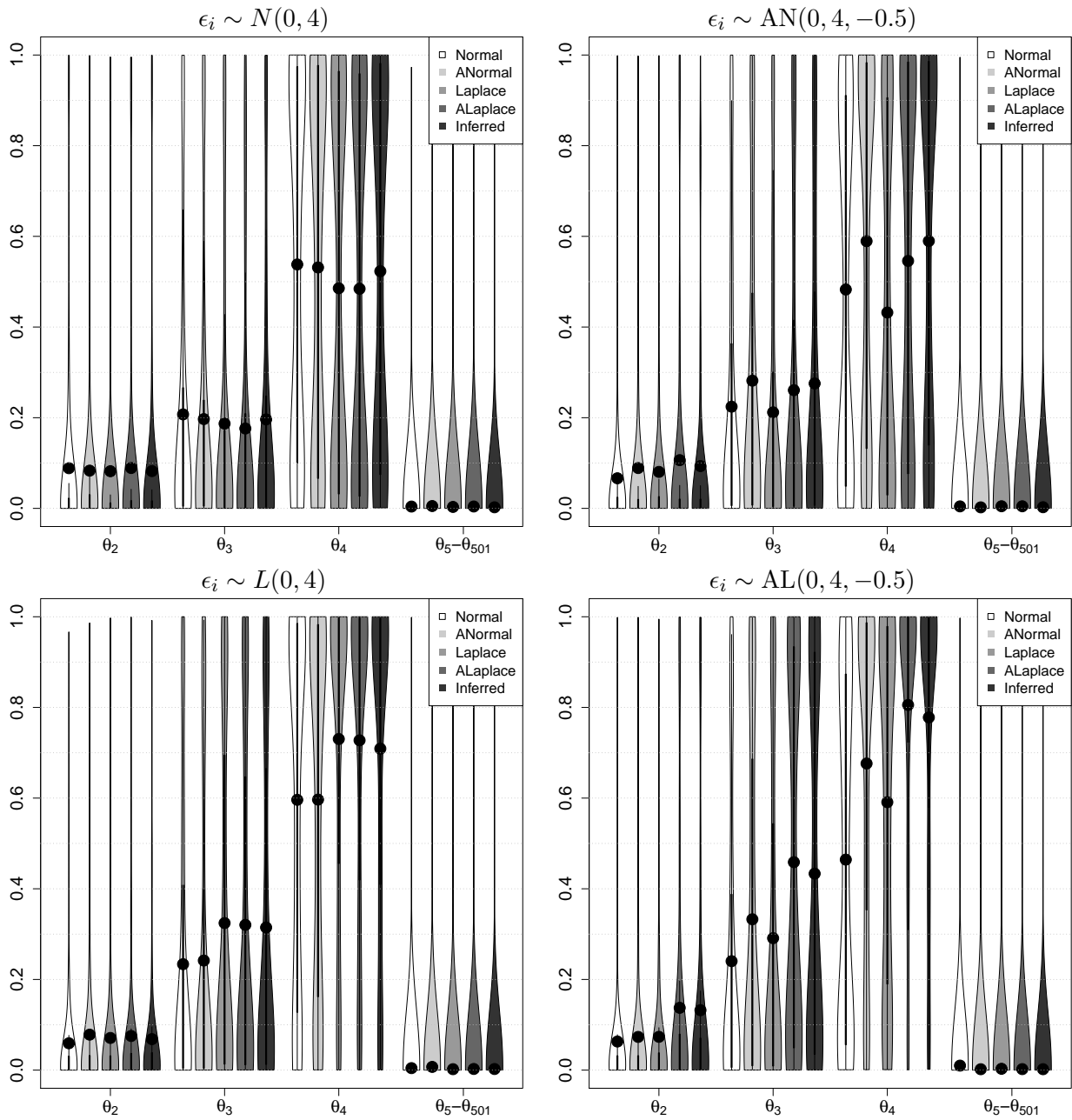


Figure 7S:  $P(\theta_i \neq 0 | y)$  for  $p = 500$ ,  $\vartheta = 2$ ,  $\theta = (0, 0.5, 1, 1.5, 0, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Black circles show the mean.



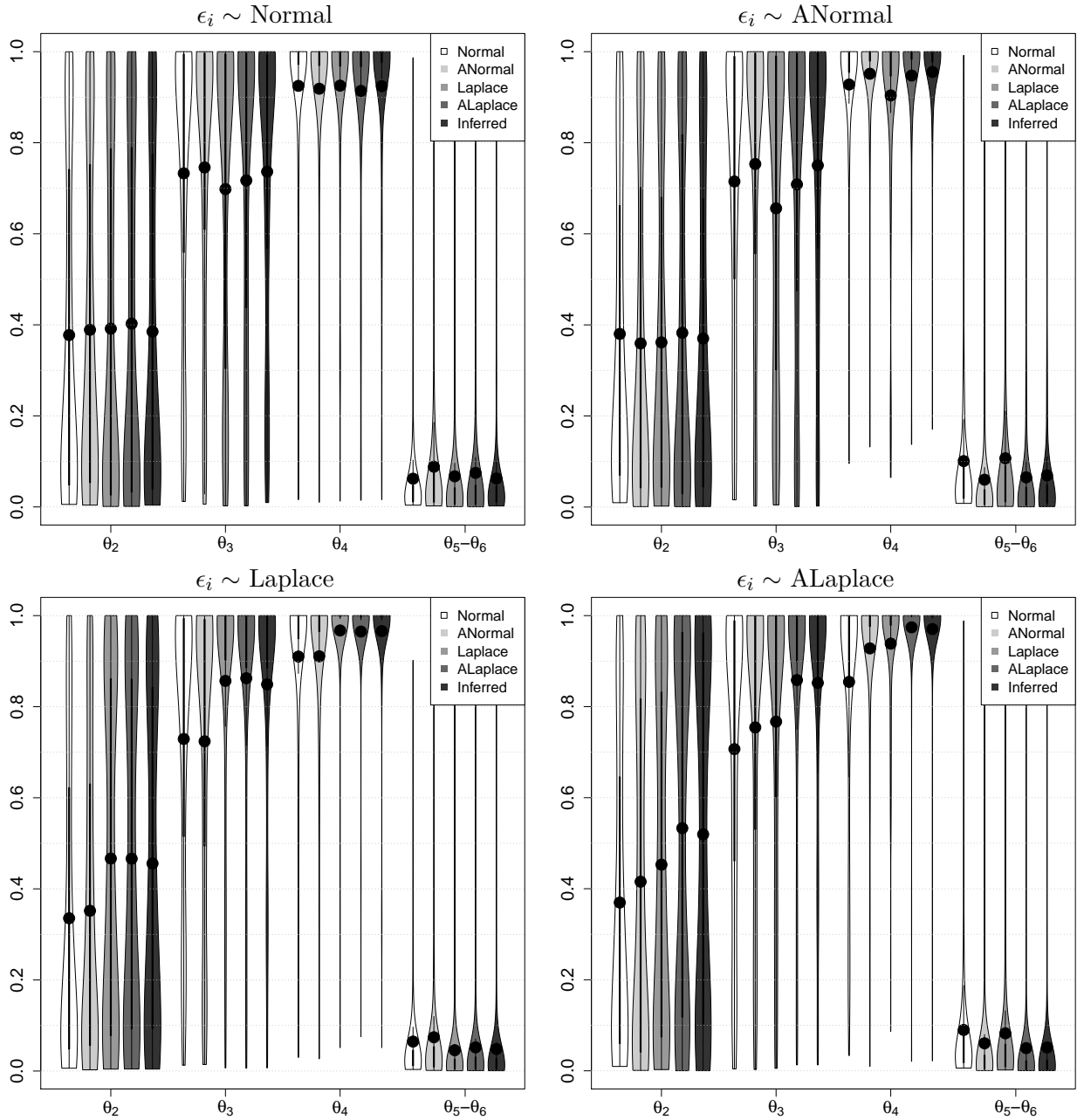


Figure 8S:  $P(\theta_i \neq 0 | y)$  for simulation with constant  $\vartheta = 0$  and varying  $\tanh(\alpha_i) \sim N(\text{atanh}(\bar{\alpha}), 1/4^2)$ , where  $\bar{\alpha} = 0$  for Normal and Laplace and  $\bar{\alpha} = -0.5$  for ANormal and ALaplace.  $P(\theta_i \neq 0 | y)$  for  $p = 6$ ,  $\theta = (0, 0.5, 1, 1.5, 0, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Black circles show the mean.

List of Tables

1S CPU time ( $10^{-4}$  seconds) on 3.4GHz Intel i7, 32Gb RAM, Windows 10.  $p = 6$ ,  $\vartheta = 4$ ,  $\theta = (0, 0.5, 0.75, 1, 0, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . . . . . 44

2S CPU time (seconds) on 8GB RAM Mac laptop with 1.6GHz Intel i5 processors running OS X 10.11.6  $p = 100$ ,  $\vartheta = 2$ ,  $\theta = (0, 0.5, 0.75, 1, 0, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . 45

3S Simulation study for  $p = 6$ . Posterior probability of the 4 error distributions under  $\vartheta = 2$ ,  $\theta = (0, 0.5, 1, 1.5, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . . . . . 46

4S Simulation study for  $p = 101, 501$ . Posterior probability of the 4 error distributions under  $g_\alpha = 0.357$ ,  $\theta = (0, 0.5, 1, 1.5, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Laplace approximation to  $p(y | \gamma)$  was used. . . . . 47

5S Simulation results under  $\vartheta = 1$ .  $\gamma_0$ : true predictors.  $\hat{\gamma}$ : selected variables. CC: number of correctly classified variables ( $\sum_{j=1}^p \mathbf{I}(\hat{\gamma}_j = \gamma_{0j})$ ). FP: number of false positives; TP: number of true positives. LASSO-LAD and LASSO-QR are equivalent when  $\alpha = 0$  . . . . . 48

6S Simulation results under  $\vartheta = 2$ .  $\gamma_0$ : true predictors.  $\hat{\gamma}$ : selected variables. CC: number of correctly classified variables ( $\sum_{j=1}^p \mathbf{I}(\hat{\gamma}_j = \gamma_{0j})$ ). FP: number of false positives; TP: number of true positives. LASSO-LAD and LASSO-QR are equivalent when  $\alpha = 0$  . . . . . 49

7S Inference on the error distribution under the  $p = 6$  simulation and heteroskedastic  $\vartheta_i \propto e^{x_i^T \theta}$  errors . . . . . 50

8S Average marginal  $P(\gamma_j = 1 | \mathbf{y})$  at multiple quantiles  $q = 0.05, 0.25, 0.5, 0.75, 0.95$  (i.e. conditioning on asymmetric Laplace errors with fixed  $\alpha = 2q - 1$ ) under the  $p = 6$  simulation and heteroskedastic  $\epsilon_i \sim N(0, \vartheta_i)$ ,  $\vartheta_i \propto e^{x_i^T \theta}$  errors. Simulation truth is  $\theta = (0, 0.5, 1, 1.5, 0, 0)$  . . . . . 51

9S Number of true and false positives in non-id example with 0.5 probability of degenerate  $(y_i, x_i) = (0, \dots, 0)$ .  $p = n = 50$ ,  $\theta^* = (0.1, 0.1, 0.1, 0.1, 0.1, 0, \dots, 0)$ ,  $\vartheta^* = 2$  . . . 52

10S Six genes with largest  $p(\gamma_j = 1 | y)$  in the DLD dataset under assumed normality and inferred error distribution. . . . . 53

11S DLD data. Top 5 models when conditioning on asymmetric Laplace residuals and  
fixed  $\alpha = -0.5, 0, 0.5$  . . . . . 54

Simulation truth $\epsilon_i \sim \text{AN}(0, 4, \alpha)$				
Fitted model	$\alpha = 0$	$\alpha = -0.25$	$\alpha = -0.5$	$\alpha = -0.75$
Normal	76.99	98.84	103.84	101.25
ANormal	92.08	86.10	102.60	115.64
Laplace	90.58	92.84	97.90	93.13
ALaplace	122.64	121.12	124.69	131.50
Simulation truth $\epsilon_i \sim \text{AL}(0, 4, \alpha)$				
Fitted model	$\alpha = 0$	$\alpha = -0.25$	$\alpha = -0.5$	$\alpha = -0.75$
Normal	76.62	96.05	99.85	97.78
ANormal	81.77	82.74	92.67	104.76
Laplace	90.29	93.42	92.88	91.25
ALaplace	117.30	113.69	115.08	122.79

Table 1S: CPU time ( $10^{-4}$  seconds) on 3.4GHz Intel i7, 32Gb RAM, Windows 10.  $p = 6$ ,  $\vartheta = 4$ ,  $\theta = (0, 0.5, 0.75, 1, 0, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ .

	Simulation truth			
	$N(0, 4)$	$AN(0, 4, -0.5)$	$L(0, 4)$	$AL(0, 4, -0.5)$
Normal	6.9	29.7	6.4	32.9
ANormal	52.9	21.9	41.3	22.2
Laplace	17.0	28.2	14.6	26.6
ALaplace	57.7	26.4	26.7	22.4
Inferred	6.1	22.6	13.5	23.0

Table 2S: CPU time (seconds) on 8GB RAM Mac laptop with 1.6GHz Intel i5 processors running OS X 10.11.6  $p = 100$ ,  $\vartheta = 2$ ,  $\theta = (0, 0.5, 0.75, 1, 0, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ .

Truth	Average $p(\gamma_{p+1}, \gamma_{p+2}   y)$			
	$\gamma_{p+1} = \gamma_{p+2} = 0$	$\gamma_{p+1} = 1, \gamma_{p+2} = 0$	$\gamma_{p+1} = 0, \gamma_{p+2} = 1$	$\gamma_{p+1} = \gamma_{p+2} = 0$
	$p = 6, g_\alpha = 0.357, \text{Laplace } p(\gamma   y)$			
$N(0, 2)$	0.91	0.02	0.06	0.00
$AN(0, 2, -0.5)$	0.11	0.81	0.01	0.06
$L(0, 2)$	0.14	0.00	0.84	0.02
$AL(0, 2, -0.5)$	0.02	0.12	0.01	0.85
	$p = 6, g_\alpha = 0.357, \text{Monte Carlo } p(\gamma   y)$			
$N(0, 2)$	0.91	0.02	0.06	0.00
$AN(0, 2, -0.5)$	0.11	0.81	0.01	0.07
$L(0, 2)$	0.12	0.01	0.85	0.02
$AL(0, 2, -0.5)$	0.02	0.12	0.01	0.85
	$p = 6, g_\alpha = 0.087, \text{Laplace } p(\gamma   y)$			
$N(0, 2)$	0.87	0.07	0.06	0.01
$AN(0, 2, -0.5)$	0.07	0.86	0.01	0.07
$L(0, 2)$	0.13	0.01	0.79	0.07
$AL(0, 2, -0.5)$	0.01	0.13	0.01	0.85

Table 3S: Simulation study for  $p = 6$ . Posterior probability of the 4 error distributions under  $\vartheta = 2$ ,  $\theta = (0, 0.5, 1, 1.5, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ .

Truth	Average $p(\gamma_{p+1}, \gamma_{p+2}   y)$			
	$\gamma_{p+1} = \gamma_{p+2} = 0$	$\gamma_{p+1} = 1, \gamma_{p+2} = 0$	$\gamma_{p+1} = 0, \gamma_{p+2} = 1$	$\gamma_{p+1} = \gamma_{p+2} = 0$
$p = 101, \vartheta = 1$				
$N(0, 2)$	0.91	0.01	0.08	0.00
$AN(0, 2, -0.5)$	0.03	0.86	0.00	0.11
$L(0, 2)$	0.15	0.01	0.83	0.02
$AL(0, 2, -0.5)$	0.00	0.13	0.01	0.86
$p = 101, \vartheta = 2$				
$N(0, 2)$	0.89	0.01	0.10	0.00
$AN(0, 2, -0.5)$	0.02	0.89	0.00	0.09
$L(0, 2)$	0.15	0.01	0.82	0.02
$AL(0, 2, -0.5)$	0.00	0.16	0.01	0.83
$p = 501, \vartheta = 1$				
$N(0, 2)$	0.85	0.00	0.14	0.00
$AN(0, 2, -0.5)$	0.01	0.85	0.01	0.14
$L(0, 2)$	0.18	0.00	0.80	0.02
$AL(0, 2, -0.5)$	0.00	0.15	0.00	0.84
$p = 501, \vartheta = 2$				
$N(0, 2)$	0.83	0.00	0.16	0.00
$AN(0, 2, -0.5)$	0.00	0.87	0.00	0.12
$L(0, 2)$	0.19	0.00	0.79	0.01
$AL(0, 2, -0.5)$	0.00	0.22	0.00	0.77

Table 4S: Simulation study for  $p = 101, 501$ . Posterior probability of the 4 error distributions under  $g_\alpha = 0.357$ ,  $\theta = (0, 0.5, 1, 1.5, \dots, 0)$ ,  $n = 100$ ,  $\rho_{ij} = 0.5$ . Laplace approximation to  $p(y | \gamma)$  was used.

	$p = 100$				$p = 500$			
	$p(\gamma_0   y)$	$p(\hat{\gamma} = \gamma_0)$	FP	TP	$p(\gamma_0   y)$	$p(\hat{\gamma} = \gamma_0)$	FP	TP
Truly $\epsilon \sim N(0, 1)$								
Normal	0.46	0.63	0.1	2.7	0.26	0.37	0.2	2.4
Two-piece Normal	0.43	0.63	0.2	2.7	0.24	0.38	0.3	2.4
Laplace	0.26	0.42	0.5	2.6	0.12	0.19	0.8	2.3
Two-piece Laplace	0.23	0.39	0.7	2.6	0.12	0.21	0.9	2.3
Inferred	0.45	0.62	0.2	2.7	0.25	0.37	0.2	2.4
LASSO-LS		0.00	12.4	3.0		0.00	20.4	2.9
LASSO-LAD		0.00	10.2	2.9		0.00	18.7	2.6
LASSO-QR		0.00	10.2	2.9		0.00	18.7	2.6
SCAD		0.07	4.2	2.9		0.01	7.3	2.8
Truly $\epsilon \sim AN(0, 1, -0.5)$								
Normal	0.38	0.55	0.2	2.6	0.21	0.34	0.5	2.4
Two-piece Normal	0.59	0.73	0.1	2.8	0.40	0.55	0.4	2.6
Laplace	0.20	0.35	0.7	2.5	0.07	0.14	1.2	2.4
Two-piece Laplace	0.33	0.48	0.5	2.7	0.18	0.32	1.1	2.5
Inferred	0.57	0.72	0.1	2.8	0.38	0.52	0.4	2.6
LASSO-LS		0.00	12.4	3.0		0.00	21.9	2.9
LASSO-LAD		0.00	9.8	2.8		0.00	18.1	2.6
LASSO-QR		0.00	9.0	2.9		0.00	15.1	2.7
SCAD		0.07	4.0	2.9		0.03	7.3	2.8
Truly $\epsilon \sim L(0, 1)$								
Normal	0.11	0.14	0.3	2.0	0.03	0.02	0.6	1.6
Two-piece Normal	0.11	0.15	0.3	2.1	0.04	0.04	1.1	1.7
Laplace	0.29	0.38	0.2	2.4	0.13	0.19	0.4	2.0
Two-piece Laplace	0.28	0.35	0.3	2.4	0.12	0.18	0.5	2.0
Inferred	0.28	0.38	0.2	2.4	0.12	0.18	0.4	2.0
LASSO-LS		0.00	11.3	2.8		0.00	21.4	2.5
LASSO-LAD		0.01	9.7	2.8		0.00	17.8	2.5
LASSO-QR		0.01	9.7	2.8		0.00	17.8	2.5
SCAD		0.02	5.0	2.7		0.00	9.0	2.4
Truly $\epsilon \sim AL(0, -0.5)$								
Normal	0.07	0.10	0.4	1.9	0.02	0.02	1.1	1.5
Two-piece Normal	0.21	0.27	0.2	2.2	0.11	0.15	0.3	2.0
Laplace	0.16	0.19	0.4	2.1	0.05	0.07	0.7	1.8
Two-piece Laplace	0.43	0.51	0.2	2.5	0.27	0.34	0.4	2.3
Inferred	0.41	0.48	0.2	2.5	0.25	0.33	0.4	2.2
LASSO-LS		0.00	11.6	2.8		0.00	20.1	2.5
LASSO-LAD		0.00	9.9	2.7		0.00	17.5	2.3
LASSO-QR		0.00	9.0	2.8		0.00	15.2	2.5
SCAD		0.01	5.2	2.6		0.01	9.4	2.3

Table 5S: Simulation results under  $\vartheta = 1$ .  $\gamma_0$ : true predictors.  $\hat{\gamma}$ : selected variables. CC: number of correctly classified variables ( $\sum_{j=1}^p \mathbf{I}(\hat{\gamma}_j = \gamma_{0j})$ ). FP: number of false positives; TP: number of true positives. LASSO-LAD and LASSO-QR are equivalent when  $\alpha = 0$



	$p = 100$				$p = 500$			
	$p(\gamma_0   y)$	$p(\hat{\gamma} = \gamma_0)$	FP	TP	$p(\gamma_0   y)$	$p(\hat{\gamma} = \gamma_0)$	FP	TP
Truly $\epsilon \sim N(0, 1)$								
Normal	0.01	0.01	0.4	1.2	0.00	0.00	0.8	0.9
Two-piece Normal	0.01	0.01	0.5	1.2	0.00	0.00	0.9	0.8
Laplace	0.00	0.00	0.7	1.1	0.00	0.00	1.0	0.8
Two-piece Laplace	0.00	0.01	0.8	1.1	0.00	0.00	1.1	0.8
Inferred	0.01	0.01	0.5	1.2	0.00	0.00	0.7	0.9
LASSO-LS		0.00	11.9	2.5		0.00	18.0	2.0
LASSO-LAD		0.00	8.9	2.0		0.00	15.6	1.4
LASSO-QR		0.00	8.9	2.0		0.00	15.6	1.4
SCAD		0.00	6.3	2.3		0.01	10.4	1.8
Truly $\epsilon \sim AN(0, 1, -0.5)$								
Normal	0.00	0.00	0.5	1.2	0.00	0.00	0.7	0.9
Two-piece Normal	0.01	0.01	0.4	1.4	0.00	0.01	0.7	1.1
Laplace	0.00	0.00	0.9	1.0	0.00	0.00	1.4	0.7
Two-piece Laplace	0.01	0.01	0.7	1.2	0.00	0.00	1.5	1.0
Inferred	0.01	0.01	0.4	1.4	0.00	0.01	0.9	1.0
LASSO-LS		0.00	11.0	2.4		0.00	19.4	1.9
LASSO-LAD		0.00	8.6	1.8		0.00	15.3	1.4
LASSO-QR		0.00	8.1	2.1		0.00	12.8	1.5
SCAD		0.00	6.1	2.1		0.00	10.1	1.8
Truly $\epsilon \sim L(0, 1)$								
Normal	0.01	0.01	0.4	1.3	0.00	0.00	0.8	0.9
Two-piece Normal	0.01	0.01	0.5	1.3	0.00	0.00	0.9	1.0
Laplace	0.05	0.06	0.4	1.7	0.01	0.01	0.7	1.2
Two-piece Laplace	0.05	0.07	0.4	1.7	0.01	0.01	0.8	1.2
Inferred	0.04	0.04	0.3	1.7	0.01	0.01	0.7	1.2
LASSO-LS		0.00	10.8	2.5		0.00	20.4	2.0
LASSO-LAD		0.01	9.3	2.5		0.00	17.1	2.0
LASSO-QR		0.01	9.3	2.5		0.00	17.1	2.0
SCAD		0.00	5.9	2.2		0.00	10.3	1.8
Truly $\epsilon \sim AL(0, -0.5)$								
Normal	0.00	0.00	0.5	1.1	0.00	0.00	0.9	0.8
Two-piece Normal	0.02	0.01	0.4	1.5	0.01	0.01	0.6	1.2
Laplace	0.02	0.01	0.6	1.3	0.00	0.01	0.8	1.0
Two-piece Laplace	0.09	0.12	0.3	1.9	0.04	0.05	0.7	1.5
Inferred	0.09	0.10	0.3	1.8	0.04	0.05	0.6	1.4
LASSO-LS		0.00	10.9	2.3		0.00	18.0	1.8
LASSO-LAD		0.00	9.4	2.3		0.00	15.6	1.7
LASSO-QR		0.00	8.3	2.5		0.00	14.0	2.0
SCAD		0.01	5.7	2.1		0.00	10.2	1.6

Table 6S: Simulation results under  $\vartheta = 2$ .  $\gamma_0$ : true predictors.  $\hat{\gamma}$ : selected variables. CC: number of correctly classified variables ( $\sum_{j=1}^p \mathbf{I}(\hat{\gamma}_j = \gamma_{0j})$ ). FP: number of false positives; TP: number of true positives. LASSO-LAD and LASSO-QR are equivalent when  $\alpha = 0$

Truth	Average $p(\gamma_{p+1}, \gamma_{p+2}   y)$			
	$\gamma_{p+1} = \gamma_{p+2} = 0$	$\gamma_{p+1} = 1, \gamma_{p+2} = 0$	$\gamma_{p+1} = 0, \gamma_{p+2} = 1$	$\gamma_{p+1} = \gamma_{p+2} = 0$
$N(0, \vartheta_i)$	0.000	0.000	0.914	0.086
$AN(0, \vartheta_i, -0.5)$	0.000	0.003	0.096	0.901
$L(0, \vartheta_i)$	0.000	0.000	0.906	0.094
$AL(0, \vartheta_i, -0.5)$	0.000	0.000	0.053	0.947

Table 7S: Inference on the error distribution under the  $p = 6$  simulation and heteroskedastic  $\vartheta_i \propto e^{x_i^T \theta}$  errors

	$P(\gamma_2 = 1   \mathbf{y})$	$P(\gamma_3 = 1   \mathbf{y})$	$P(\gamma_4 = 1   \mathbf{y})$	$P(\gamma_5 = 1   \mathbf{y})$	$P(\gamma_6 = 1   \mathbf{y})$
$q = 0.05$	0.425	0.834	0.961	0.017	0.015
$q = 0.25$	0.751	0.950	0.996	0.016	0.015
$q = 0.5$	0.796	0.970	0.999	0.020	0.016
$q = 0.75$	0.769	0.969	0.999	0.016	0.012
$q = 0.95$	0.473	0.912	0.987	0.016	0.016

Table 8S: Average marginal  $P(\gamma_j = 1 | \mathbf{y})$  at multiple quantiles  $q = 0.05, 0.25, 0.5, 0.75, 0.95$  (i.e. conditioning on asymmetric Laplace errors with fixed  $\alpha = 2q - 1$ ) under the  $p = 6$  simulation and heteroskedastic  $\epsilon_i \sim N(0, \vartheta_i)$ ,  $\vartheta_i \propto e^{x_i^T \theta}$  errors. Simulation truth is  $\theta = (0, 0.5, 1, 1.5, 0, 0)$

	TP	FP
Zellner, Normal errors	2.8	21.3
pMOM, Normal errors	3.0	12.0
pMOM, inferred errors	2.8	10.5
peMOM, Normal errors	1.9	2.9

Table 9S: Number of true and false positives in non-id example with 0.5 probability of degenerate  $(y_i, x_i) = (0, \dots, 0)$ .  $p = n = 50$ ,  $\theta^* = (0.1, 0.1, 0.1, 0.1, 0.1, 0, \dots, 0)$ ,  $\vartheta^* = 2$

Gene symbol	Normal	Inferred
C6orf226	1.000	1.000
ECH1	1.000	1.000
CSF2RA	1.000	1.000
RRP1B	0.944	0.999
FBXL19	0.993	0.658
MTMR1	0.183	0.467
SLC35B4	0.209	0.332
RAB3GAP2	0.007	0.040

Table 10S: Six genes with largest  $p(\gamma_j = 1 | y)$  in the DLD dataset under assumed normality and inferred error distribution.

$\alpha = -0.5$	
Model	$P(\gamma   y)$
C6orf226, ECH1, CSF2RA, FBXL19, RRP1B	0.384
SLC35B4, C6orf226, ECH1, CSF2RA, RRP1B	0.349
SLC35B4, C6orf226, MTMR1, ECH1, CSF2RA, RRP1B	0.127
C6orf226, MTMR1, ECH1, CSF2RA, FBXL19, RRP1B	0.049
C6orf226, MTMR1, RAB3GAP2, ECH1, CSF2RA, RRP1B	0.023
$\alpha = 0$	
Model	$P(\gamma   y)$
C6orf226, MTMR1, ECH1, CSF2RA, FBXL19, RRP1B	0.454
C6orf226, ECH1, CSF2RA, FBXL19, RRP1B	0.258
SLC35B4, C6orf226, MTMR1, ECH1, CSF2RA, RRP1B	0.108
SLC35B4, C6orf226, ECH1, CSF2RA, RRP1B	0.061
C6orf226, MTMR1, RAB3GAP2, ECH1, CSF2RA, RRP1B	0.016
$\alpha = 0.5$	
Model	$P(\gamma   y)$
C6orf226, ECH1, CSF2RA, FBXL19, RRP1B	0.399
SLC35B4, C6orf226, ECH1, CSF2RA, RRP1B	0.359
SLC35B4, C6orf226, MTMR1, ECH1, CSF2RA, RRP1B	0.120
C6orf226, MTMR1, ECH1, CSF2RA, FBXL19, RRP1B	0.051
SLC35B4, C6orf226, RAB3GAP2, ECH1, CSF2RA, RRP1B	0.008

Table 11S: DLD data. Top 5 models when conditioning on asymmetric Laplace residuals and fixed  $\alpha = -0.5, 0, 0.5$