

# Supplement to: Default Priors for the Intercept Parameter in Logistic Regressions by Boonstra, et al.

## S1 Properties of the Exponential Power (EP) Distribution

### Normalizing Constant

Let  $\alpha \sim \text{EP}_k(\sigma)$ , so that  $\pi(\alpha) \propto \exp\{-|\alpha/\sqrt{2}\sigma|^k\}$ ,  $\alpha \in (-\infty, \infty)$ . The normalizing constant is  $Z(\sigma, k) = 1/\left[\int_{-\infty}^{\infty} \exp\{-|x/\sqrt{2}\sigma|^k\}dx\right] = 0.5/\left[\int_0^{\infty} \exp\{-(x/\sqrt{2}\sigma)^k\}dx\right]$ . Define  $u = (x/k)^k$ , so that  $x = ku^{1/k}$  and  $dx = u^{1/k-1}du$ . Then,

$$\begin{aligned}\int_0^{\infty} \exp\{-(x/\sqrt{2}\sigma)^k\}dx &= \int_0^{\infty} u^{1/k-1} \exp\left\{-u(\sqrt{2}\sigma/k)^{-k}\right\} du \\ &= (\sqrt{2}\sigma/k)\Gamma(1/k).\end{aligned}$$

The last integral is the kernel of a gamma density, implying that  $U = (|\alpha|/k)^k \sim \text{Gamma}(1/k, (\sqrt{2}\sigma/k)^k)$ . Thus,  $Z(\sigma, k) = k/[\sqrt{8}\sigma\Gamma(1/k)]$ , and

$$\pi(\alpha|\sigma, k) = \frac{k \exp\{-|\alpha/\sqrt{2}\sigma|^k\}}{\sqrt{8}\sigma\Gamma(1/k)}$$

## Distribution Function

Using the above relationship and the symmetry of the ExpPow distribution, we have

$$\begin{aligned}\Pr(|\alpha| > x|\sigma, k) &= \Pr(U > (x/k)^k|\sigma, k) = \frac{\Gamma(1/k, (x/k)^k/(\sqrt{2}\sigma/k)^k)}{\Gamma(1/k)} \\ &= \frac{\Gamma(1/k, (x/\sqrt{2}\sigma)^k)}{\Gamma(1/k)}\end{aligned}$$

This can be calculated in R (R Core Team, 2016) with

```
> pgamma(q = (x/sqrt(2)/sigma)^k, shape = 1/k, lower.tail = FALSE)
```

## S2 Generating Scenarios in Simulation Study: Further Details

**Scenario 1**  $\mathbf{X}$  is multivariate Bernoulli with length  $p = 4$ , where  $\mathbf{X} = 1[\mathbf{X}^* > 0]$ ;  $\mathbf{X}^*$  is multivariate normal with length  $p$ ; each element has mean zero, unit variance, and pairwise correlation 0.50;  $\boldsymbol{\beta}^* = \{0.50, \dots, 0.50\}$ ;  $\alpha^* = -2.5$ ;  $\alpha^* + E(\mathbf{X})^\top \boldsymbol{\beta}^* = -1.5$ ;  $n \in \{50, 100, 200, 400\}$

**Scenario 2**  $\mathbf{X}$  is multivariate Bernoulli with length  $p = 25$ , where  $\mathbf{X} = 1[\mathbf{X}^* > 0]$ ;  $\mathbf{X}^*$  is multivariate normal with length  $p$ ; each element has mean  $\Phi^{-1}(0.25) = -0.67$ , unit variance, and pairwise correlation 0.15;  $\boldsymbol{\beta}^* = \{1.5, 0, \dots, 0\}$ ;  $\alpha^* = -2$ ;  $\alpha^* + E(\mathbf{X})^\top \boldsymbol{\beta}^* = -1.625$ ;  $n \in \{50, 100, 200\}$

**Scenario 3**  $\mathbf{X}$  is multivariate Bernoulli with length  $p = 25$ , where  $\mathbf{X} = 1[\mathbf{X}^* > 0]$ ;  $\mathbf{X}^*$  is multivariate normal with length  $p$ ; each element has mean  $\Phi^{-1}(0.25) = -0.67$ , unit variance, and pairwise correlation 0.15;  $\boldsymbol{\beta}^* = \{0.06, \dots, 0.06\}$ ;  $\alpha^* = -2$ ;  $\alpha^* + E(\mathbf{X})^\top \boldsymbol{\beta}^* = -1.625$ ;  $n \in \{50, 100, 200\}$

**Scenario 4**  $\mathbf{X}$  is multivariate Bernoulli with length  $p = 25$ , where  $\mathbf{X} = 1[\mathbf{X}^* > 0]$ ;  $\mathbf{X}^*$  is multivariate normal with length  $p$ ; each element has mean  $\Phi^{-1}(0.05) = -1.64$ , unit variance, and pairwise correlation 0.30;  $\boldsymbol{\beta}^* = \{\underbrace{3, \dots, 3}_{10}, \underbrace{0, \dots, 0}_{15}\}$ ;  $\alpha^* = -6.5$ ;  $\alpha^* + E(\mathbf{X})^\top \boldsymbol{\beta}^* = -5$ ;  $n \in \{100, 200, 400\}$

**Scenario 5**  $\mathbf{X}$  is multivariate Bernoulli with length  $p = 25$ , where  $\mathbf{X} = 1[\mathbf{X}^* > 0]$ ;  $\mathbf{X}^*$  is multivariate normal with length  $p$ ; each element has mean  $\Phi^{-1}(0.05) = -1.64$ ,

unit variance, and pairwise correlation 0.30;  $\beta^* = \{\underbrace{3, \dots, 3}_{10}, \underbrace{0, \dots, 0}_{15}\}$ ;  $\alpha^* = -4$ ;  $\alpha^* +$

$$E(\mathbf{X})^\top \beta^* = -2.5; n \in \{50, 100, 200\}$$

**Scenario 6**  $\mathbf{X}$  is multivariate normal with length  $p = 75$ ; each element has mean zero, unit variance, and pairwise correlation 0.30;  $\beta^* = \{2, 0, \dots, 0\}$ ;  $\alpha^* = \alpha^* + E(\mathbf{X})^\top \beta^* = -4$ ;  $n \in \{100, 200, 400\}$

**Scenario 7**  $\mathbf{X}$  is multivariate Bernoulli with length  $p = 75$ , where  $\mathbf{X} = 1[\mathbf{X}^* > 0]$ ;  $\mathbf{X}^*$  is multivariate normal with length  $p$ ; each element has mean  $\Phi^{-1}(0.25) = -0.67$ , unit variance, and pairwise correlation 0.30;  $\beta^* = \{2, 0, \dots, 0\}$ ;  $\alpha^* = -3.5$ ;  $\alpha^* + E(\mathbf{X})^\top \beta^* = -3$ ;  $n \in \{100, 200, 400\}$

**Scenario 8**  $\mathbf{X}$  is multivariate normal with length  $p = 150$ ; each element has mean zero, unit variance, and pairwise correlation 0.10;  $\beta^* = \{-0.5, -0.5, 0, \dots, 0\}$ ;  $\alpha^* = \alpha^* + E(\mathbf{X})^\top \beta^* = -3$ ;  $n \in \{100, 200, 400, 600\}$

**Scenario 9**  $\mathbf{X}$  is multivariate normal with length  $p = 150$ ; each element has mean zero, unit variance, and pairwise correlation 0.10;  $\beta^* = \{-1/150, \dots, -1/150\}$ ;  $\alpha^* = \alpha^* + E(\mathbf{X})^\top \beta^* = -3$ ;  $n \in \{100, 200, 400, 600\}$

### S3 Algorithm 2: Calculate $n_{\text{comp}}$ , $n_{\text{piv}}$ , or $n_{\text{over}}$

We assume that  $n_{\text{comp}}$  is to be calculated, and the algorithm is similar for the other separation statistics. For a given hyperplane of dimension  $p+1$ , say  $\mathbf{b}$ , which yields the one-dimensional linear predictor  $\{1, \mathbf{X}_i^\top\} \mathbf{b}$ ,  $i = 1, \dots, n$ , the minimum number of observations removed to induce separation is calculable in  $n \log(n)$  time by using a binary search. The minimum of these minima over all possible hyperplanes is precisely the  $n_{\text{comp}}$  statistic, but determining this is NP hard (Hoffgen et al., 1995; Christmann and Rousseeuw, 2001). The algorithm we used also reports a minimum of minima but searches over a subset of candidate hyperplanes that is intended to be likely to contain the best-separating hyperplane, i.e. correspond to  $n_{\text{comp}}$ . Thus, by construction, the value returned by our algorithm will always bound the true value of  $n_{\text{comp}}$  from above, but this may not be tight, i.e. we can only guarantee it is an upper-bound.

For a given dataset of  $n$  observations of  $\{Y_i, \mathbf{X}_i\}_i$  with  $\mathbf{X}_i$  being length- $p$ , the algorithm proceeds as follows.

1. Calculate the linear predictor from the full multivariable logistic regression of  $Y$  against all predictors  $\{\mathbf{X}_i\}_i$  simultaneously. Set  $n_{\text{comp}}^{\text{current}}$  equal to the number observations

needed to induce separation using this linear predictor. If this number is zero, proceed to step 5.

2. Now determine whether a different hyperplane constructed using a subset of observations can induce separation. The intuition is that hyperplanes constructed from subsamples that preferentially exclude the high-influence observations, i.e. those subsamples enriched with low-influence observations, may be able to better separate the data than those hyperplanes that are sensitive to high-influence observations. Calculate a length- $n$  vector of sampling weights  $\{w_i\}$ , where each  $w_i$  is equal to the square root of the reciprocal of the absolute difference between the corresponding observation’s linear predictor from the full multivariable logistic model and that calculated after leaving out the observation. This is used as a simple measure of that observation’s influence in estimation.
3. For  $k = 1, \dots, MAXk$ ,
  - (a) if  $\binom{n}{n-k} < NDRAWS$ , then, construct every possible subset of  $\binom{n}{n-k}$  observations.
  - (b) Otherwise, sample  $NDRAWS$  subsets of size  $n - k$  as follows. Identify the  $NDRAWS/2$  subsets of observations with the largest partial sums of weights  $w_i$  from step 3. For the remaining  $NDRAWS/2$  subsets, randomly sample observations in proportion to each weight  $w_i$ . Remove any duplicated sampled subsets from among the  $NDRAWS$  subsets.

For every subset of observations above, calculate the linear predictor for all  $n$  observations using the multivariable logistic regression of  $Y$  against all  $\mathbf{X}$ ’s fit to the subset of size  $n - k$ . Whenever the minimum number removed to induce separation falls below  $n_{\text{comp}}^{\text{current}}$ , set  $n_{\text{comp}}^{\text{current}}$  equal to this number. If  $n_{\text{comp}}^{\text{current}} = 0$ , step out of the loop and proceed to step 5.

4. Return the final value of  $n_{\text{comp}}^{\text{current}}$  as the estimate of  $n_{\text{comp}}$  and the corresponding hyperplane that resulted in this greatest separation.

We extended this algorithm to also approximate  $n_{\text{piv}}$  (the minimum number of observations necessary to induce pivotal separation among the remaining subsets, defined in Section 2.1) and  $n_{\text{over}}$  (the minimum number of observations necessary to remove to induce quasi-complete separation among the remaining subset (Christmann and Rousseeuw, 2001)). Applied to five datasets considered in Christmann and Rousseeuw (2001), and using  $MAXk = 8$  and  $NDRAWS = 1000$  (for small  $n$ ) or  $NDRAWS = 200$  (for large  $n$ ), Algorithm 2 was able to match or improve upon (i.e. identify a verifiable lower upperbound of  $n_{\text{comp}}$  or  $n_{\text{over}}$ ) an alternative algorithm considered in that paper (see Table 4 in the manuscript). We were not able obtain a sixth dataset (‘Hemophilia’) considered by those authors

In the simulation study and analysis in the manuscript, we used  $MAKk = 8$  and a larger value of  $NDRAWS = 5000$ .

## S4 Supplemental Tables and Figure

Table S1: Values of the extreme boundaries  $\text{logit}^{-1}(s_n)$  and the scale parameter  $\sigma_n$  of the EP distribution with  $k = 2$  or  $k = 4$  and the Logistic distribution, for varying sample sizes  $n$ , as calculated by **Algorithm 1** with  $\delta = 1$  and  $q = 0.01$ .

$n$	$\text{logit}^{-1}(s_n)$	EP( $k = 2$ )	EP( $k = 4$ )	Logistic
		$\sigma_n$	$\sigma_n$	$\sigma_n$
250	6.21	2.41	3.52	1.17
500	6.91	2.68	3.91	1.30
1000	7.60	2.95	4.30	1.44
2000	8.29	3.22	4.70	1.57
4000	8.99	3.49	5.09	1.70
8000	9.68	3.76	5.48	1.83

Table S2: Median AUCs ( $\times 100$ ), across 200 datasets, using the posterior mean of  $\beta$  for all combinations of six priors on the intercept parameter  $\alpha$  and two priors on  $\beta$ . Larger numbers are better. In each row, values in *italics* correspond to the largest median AUC, and values in **bold** had a larger AUC than the italicized value in at least 33% of the datasets. All metrics were calculated separately for each prior on  $\beta$ . The columns  $n_{\text{comp}}$  and  $n_{\text{piv}}$  give the median values of these statistics across the 200 simulated datasets.

Scenario	$p$	$n$	median		$\beta \sim \text{HS}(\bar{p}_{\text{eff}})$						$\beta \sim \text{Logis}(1)$					
			$n_{\text{comp.}}$	$n_{\text{piv.}}$	$t_3(10)$	$t_\infty(10)$	$\text{EP}_2(\sigma_n)$	$\text{EP}_4(\sigma_n)$	$\text{EP}_{10}(\sigma_n)$	$\text{Logis}(\sigma_n)$	$t_3(10)$	$t_\infty(10)$	$\text{EP}_2(\sigma_n)$	$\text{EP}_4(\sigma_n)$	$\text{EP}_{10}(\sigma_n)$	$\text{Logis}(\sigma_n)$
1	4	50	9	9	<b>66.3</b>	<b>66.4</b>	<b>66.4</b>	<b>66.3</b>	<b>66.4</b>	<i>66.5</i>	<b>65.5</b>	<i>65.5</i>	<b>65.5</b>	<b>65.5</b>	<b>65.5</b>	<b>65.5</b>
1	4	100	20	20	<b>66.9</b>	<i>67.1</i>	<b>67.0</b>	<b>67.0</b>	<b>67.1</b>	<b>67.1</b>	<i>66.6</i>	<b>66.5</b>	<b>66.5</b>	<b>66.5</b>	<b>66.5</b>	<b>66.5</b>
1	4	200	40	40	<b>67.7</b>	<b>67.7</b>	<b>67.7</b>	<i>67.8</i>	<b>67.7</b>	<b>67.7</b>	<b>67.4</b>	<i>67.4</i>	<b>67.4</b>	<b>67.4</b>	<b>67.4</b>	<b>67.4</b>
1	4	400	84	84	<i>67.9</i>	<b>67.9</b>	<b>67.9</b>	<b>67.9</b>	<b>67.9</b>	<b>67.9</b>	<b>67.8</b>	<b>67.8</b>	<b>67.8</b>	<i>67.8</i>	<b>67.8</b>	<b>67.8</b>
2	25	50	0	1	<b>56.5</b>	<b>56.5</b>	<b>57.8</b>	<b>57.5</b>	<b>57.3</b>	<i>58.0</i>	<b>56.6</b>	<b>56.5</b>	<b>56.6</b>	<b>56.5</b>	<b>56.4</b>	<i>56.7</i>
2	25	100	2	7.5	<b>62.7</b>	<b>63.1</b>	<i>63.1</i>	<b>62.9</b>	<b>62.8</b>	<b>63.1</b>	<b>58.5</b>	<b>58.5</b>	<b>58.5</b>	<b>58.5</b>	<b>58.4</b>	<i>58.6</i>
2	25	200	22	24.5	<b>65.3</b>	<b>65.3</b>	<b>65.3</b>	<i>65.3</i>	<b>65.3</b>	<b>65.3</b>	<b>61.2</b>	<i>61.3</i>	<b>61.3</b>	<b>61.2</b>	<b>61.2</b>	<b>61.3</b>
3	25	50	0	1	<b>53.0</b>	<b>52.9</b>	<i>53.0</i>	<b>52.9</b>	<b>53.0</b>	<b>53.0</b>	<b>52.8</b>	<i>52.8</i>	<b>52.8</b>	<b>52.7</b>	<b>52.8</b>	<b>52.8</b>
3	25	100	2	8	<b>52.6</b>	<b>52.6</b>	<b>52.5</b>	<b>52.6</b>	<b>52.6</b>	<i>52.6</i>	<b>52.5</b>	<b>52.5</b>	<b>52.6</b>	<i>52.6</i>	<b>52.5</b>	<b>52.5</b>
3	25	200	22	25	<b>53.7</b>	<b>53.7</b>	<b>53.8</b>	<i>53.8</i>	<b>53.8</b>	<b>53.7</b>	<b>52.8</b>	<b>52.7</b>	<b>52.8</b>	<b>52.8</b>	<i>52.8</i>	<b>52.8</b>
4	25	100	0	0	<b>80.0</b>	78.9	<b>82.2</b>	<i>82.4</i>	<b>82.2</b>	<b>82.1</b>	<i>91.7</i>	<b>91.7</b>	91.1	91.2	91.2	91.3
4	25	200	1	2	92.5	92.2	<b>93.6</b>	<i>93.6</i>	<b>93.6</b>	<b>93.5</b>	<b>94.8</b>	<i>94.8</i>	<b>94.7</b>	<b>94.7</b>	<b>94.8</b>	<b>94.8</b>
4	25	400	2	7	<b>96.1</b>	<i>96.1</i>	<b>96.1</b>	<b>96.1</b>	<b>96.1</b>	<b>96.1</b>	<b>96.3</b>	<i>96.3</i>	<b>96.3</b>	<b>96.3</b>	<b>96.3</b>	<b>96.3</b>
5	25	50	1	1	72.4	72.5	<i>73.9</i>	<b>73.7</b>	73.5	<b>73.8</b>	<b>83.9</b>	<b>83.9</b>	<b>84.1</b>	<b>84.1</b>	<i>84.3</i>	<b>84.0</b>
5	25	100	2	2	79.8	79.8	<i>80.6</i>	<b>80.5</b>	80.4	<b>80.6</b>	<b>88.0</b>	<b>87.9</b>	<b>88.0</b>	<i>88.2</i>	<b>88.0</b>	<b>88.0</b>
5	25	200	6	6	<b>90.0</b>	<b>90.1</b>	<b>90.1</b>	<i>90.1</i>	<b>90.0</b>	<b>90.0</b>	<b>90.8</b>	<i>90.8</i>	<b>90.8</b>	<b>90.8</b>	<b>90.8</b>	<b>90.8</b>
6	75	100	0	0	<b>85.6</b>	<b>85.6</b>	<b>86.2</b>	<b>86.2</b>	<b>86.2</b>	<i>86.3</i>	<i>74.5</i>	74.1	70.3	68.6	67.3	72.2
6	75	200	0	1	88.5	88.5	<b>88.7</b>	<b>88.7</b>	<b>88.7</b>	<i>88.7</i>	<i>78.5</i>	78.4	74.7	72.4	71.4	76.7
6	75	400	0	8	<b>88.7</b>	<b>88.7</b>	<i>88.8</i>	<b>88.7</b>	<b>88.7</b>	<b>88.8</b>	<i>79.0</i>	79.0	78.0	77.4	77.2	78.6
7	75	100	0	0	<b>60.3</b>	<b>60.1</b>	<b>63.6</b>	<b>63.8</b>	<b>63.5</b>	<i>64.1</i>	<i>57.7</i>	<b>57.7</b>	57.0	<b>56.5</b>	56.4	57.2
7	75	200	0	2	<b>67.7</b>	<b>67.5</b>	<i>69.5</i>	<b>69.4</b>	<b>69.1</b>	<b>69.4</b>	<b>58.5</b>	<i>58.6</i>	<b>57.8</b>	<b>57.6</b>	<b>57.3</b>	<b>58.2</b>
7	75	400	0.5	12	<b>71.7</b>	<b>71.9</b>	<b>72.0</b>	<b>71.9</b>	<b>72.0</b>	<i>72.1</i>	<b>60.2</b>	<i>60.4</i>	<b>60.2</b>	<b>59.9</b>	<b>59.9</b>	<b>60.2</b>
8	150	100	0	0	<b>55.7</b>	<b>55.7</b>	<i>56.5</i>	<b>56.5</b>	<b>56.4</b>	<b>56.4</b>	<b>55.7</b>	<i>55.7</i>	<b>55.1</b>	<b>55.0</b>	<b>55.0</b>	<b>55.5</b>
8	150	200	0	0	<b>58.4</b>	<b>58.4</b>	<i>58.6</i>	<b>58.5</b>	<b>58.3</b>	<b>58.4</b>	<i>56.3</i>	<b>56.1</b>	<b>55.7</b>	<b>54.9</b>	<b>54.6</b>	55.9
8	150	400	0	4	<i>61.3</i>	<b>61.2</b>	<b>61.0</b>	<b>61.1</b>	<b>61.1</b>	<b>61.2</b>	<i>55.6</i>	<b>55.6</b>	<b>55.0</b>	<b>54.6</b>	<b>54.4</b>	<b>55.5</b>
8	150	600	0	14	<b>63.4</b>	<b>63.6</b>	<b>63.5</b>	<b>63.4</b>	<b>63.4</b>	<i>63.6</i>	<i>56.1</i>	<b>56.0</b>	55.7	55.9	55.9	55.9
9	150	100	0	0	<b>54.7</b>	<b>55.2</b>	<b>55.6</b>	<b>55.8</b>	<i>56.1</i>	<b>56.0</b>	<b>55.0</b>	<b>55.1</b>	<i>55.2</i>	<b>55.1</b>	<b>55.0</b>	<b>55.0</b>
9	150	200	0	0	<b>56.0</b>	<b>55.9</b>	<b>56.2</b>	<b>55.9</b>	<b>55.9</b>	<i>56.3</i>	<i>55.0</i>	<b>54.9</b>	<b>54.9</b>	<b>54.3</b>	<b>54.0</b>	<b>54.8</b>
9	150	400	0	3	<b>56.2</b>	<b>56.1</b>	<i>56.7</i>	<b>56.3</b>	<b>56.2</b>	<b>56.6</b>	<b>54.2</b>	<b>54.3</b>	<b>54.2</b>	<b>53.7</b>	<b>53.8</b>	<i>54.3</i>
9	150	600	0	11	<b>56.9</b>	<b>57.0</b>	<b>56.9</b>	<b>56.9</b>	<b>56.9</b>	<i>57.1</i>	<i>54.2</i>	<b>54.1</b>	<b>54.0</b>	<b>54.1</b>	<b>53.9</b>	<b>54.1</b>
Avg. Rank (1-6)					3.77	3.80	3.29	3.36	3.48	3.30	3.05	3.13	3.67	3.77	3.91	3.45

## References

- Christmann, A. and Rousseeuw, P. (2001). Measuring overlap in binary regression. *Computational Statistics & Data Analysis* **37**, 65–75.
- Hoffgen, K., Simon, H., and Vanhorn, K. (1995). Robust trainability of single neurons. *Journal of Computer and System Sciences* **50**, 114–125.

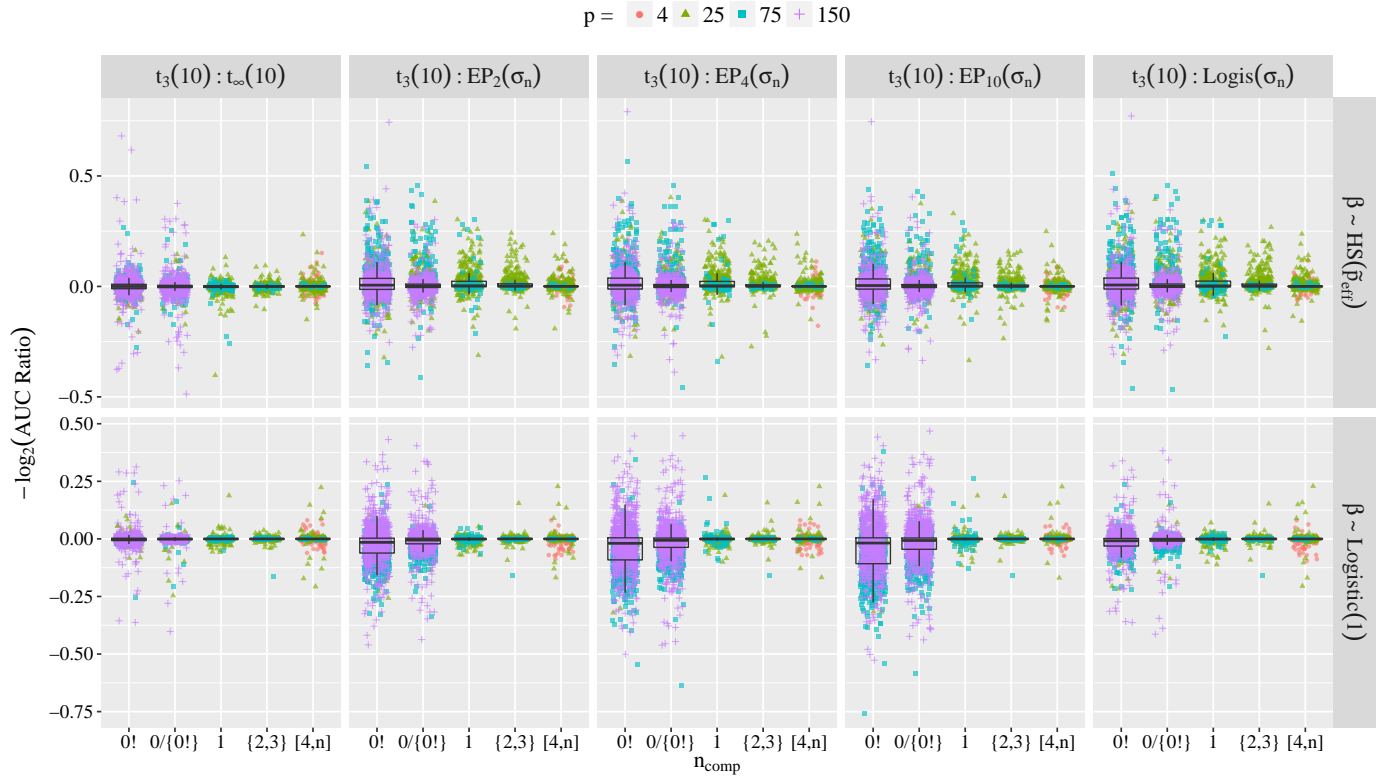


Figure S1: Comparison of  $t_3(10)$  prior on  $\alpha$  against five alternative priors on  $\alpha$  (columns) under two different priors on  $\beta$  (rows). Each point represents an individual simulated dataset. The  $y$ -axis gives area under the curve (AUC) ratios on the negative- $\log_2$  scale when using the posterior of mean of  $\beta$  to classify observations, and the  $x$ -axis defines groups based upon separation: “0!” indicates pivotal separation ( $n_{\text{piv}} = 0$ ); “0/{0!}” indicates complete but not pivotal separation ( $n_{\text{piv}} > n_{\text{comp}} = 0$ ); and the remaining categories correspond to value(s) of  $n_{\text{comp}}$ . **Positive values on the  $y$ -axis indicate that the given prior on  $\alpha$  yielded better classification of observations than a  $t_3(10)$  prior on  $\alpha$ .** Different plot characters are used to indicate  $p$ , the number of predictors. In total, each panel contains 6000 points (30 unique scenario-sample size configuration times 200 simulated datasets per configuration).

R Core Team (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.